

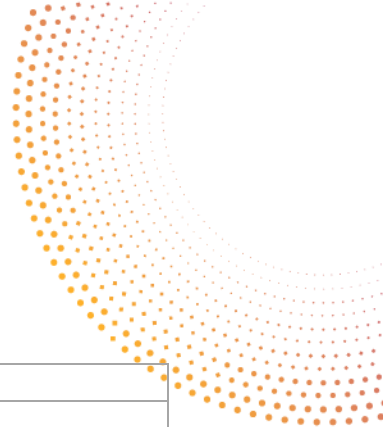


# D1.6

## Final Data Management Plan

<b>Project Title</b>	AI4Media - A European Excellence Centre for Media, Society and Democracy
<b>Contract No.</b>	951911
<b>Instrument</b>	Research and Innovation Action
<b>Thematic Priority</b>	H2020-EU.2.1.1. - INDUSTRIAL LEADERSHIP - Leadership in enabling and industrial technologies - Information and Communication Technologies (ICT) / ICT-48-2020 - Towards a vibrant European network of AI excellence centres
<b>Start of Project</b>	1 September 2020
<b>Duration</b>	48 months

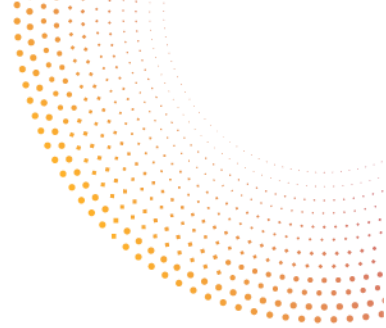




<b>Deliverable title</b>	Final Data Management Plan
<b>Deliverable number</b>	D1.6
<b>Deliverable version</b>	1.0
<b>Previous version(s)</b>	-
<b>Contractual date of delivery</b>	31 August 2024
<b>Actual date of delivery</b>	26 August 2024
<b>Deliverable filename</b>	AI4Media_D1.6_Final_Data_Management_Plan_final.docx
<b>Nature of deliverable</b>	ORDP: Open Research Data Pilot
<b>Dissemination level</b>	Public
<b>Number of pages</b>	289
<b>Work Package</b>	WP1
<b>Task(s)</b>	T1.1
<b>Partner responsible</b>	CERTH
<b>Author(s)</b>	Filareti Tsalakanidou (CERTH), Yiannis Kompatsiaris(CERTH), Vasileios Mezaris (CERTH), Symeon Papadopoulos (CERTH)
<b>Editor</b>	Filareti Tsalakanidou (CERTH)
<b>EC Project Officer</b>	Evangelia Markidou

<b>Abstract</b>	This deliverable presents the final Data Management Plan (DMP) of AI4Media and offers an update of the initial DMP (D1.2) that was delivered in M6. The final DMP summarises the strategy for data management within the AI4Media project and provides a detailed description of the datasets that have been collected, processed or generated within the project. It describes the handling of data during and after the project lifetime and discusses how they are curated and preserved. It also specifies which datasets will be openly accessible and how they will be shared, while also presenting the methodology and standards used to increase data interoperability. The report discusses 166 datasets: 61 research datasets created by project partners within the project, 81 research datasets of third-parties used in the project’s research activities, and 24 non-research datasets collected within the project.
<b>Keywords</b>	Artificial intelligence, data management, data management plan, data collection, research data, non-research data, FAIR data, metadata, open data, discoverability, accessibility, re-usability, interoperability, data security, ethical & legal aspects, open repositories





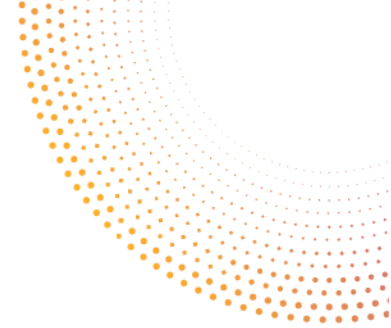
## Copyright

© Copyright 2024 AI4Media Consortium

This document may not be copied, reproduced, or modified in whole or in part for any purpose without written permission from the AI4Media Consortium. In addition to such written permission to copy, reproduce, or modify this document in whole or part, an acknowledgement of the authors of the document and all applicable portions of the copyright notice must be clearly referenced.

All rights reserved.





## Contributors

NAME	ORGANISATION
Filareti Tsalakanidou	CERTH
Yiannis Kompatsiaris	CERTH
Vasileios Mezaris	CERTH
Symeon Papadopoulos	CERTH
Evlampios Apostolidis	CERTH
Spiros Baxevanakis	CERTH
Noémie Krack	KUL
Lidia Dutkiewicz	KUL
Nicu Sebe	UNITN
Marco Formentini	UNITN
Lorenzo Seidenari	UNIFI
Ioannis Pitas	AUTH
Ioannis Patras	QMUL
Ioannis Maniadis Metaxas	QMUL
Anna Hansen	UvA
Matthew Barthet	UM
Henning Müller	HES-SO
Riccardo Fratti	HES-SO
Adrian Popescu	CEA
Sven Becker	FhG-IAIS
Milica Gerhardt	FhG-IDMT
Lucile Sassatelli	UCA
Lucia Vadicamo	CNR
Giuseppe Amato	CNR
Fabrizio Sebastiani	CNR
Adrian Tormos	BSC
Enrique Lopez	BSC
Dario Garcia Gasulla	BSC
Victor Bros	IDIAP
David Alonso del Barrio	IDIAP
Remi Mignot	IRCAM
Werner Bailer	JR
Danae Tsabouraki	ATC
Rasa Bocyte	NISV
Birgit Gray	DW
Chaja Libot	VRT
Alberto Messina	RAI
Fulvio Negro	RAI
Maurizio Montagnuolo	RAI
Roberto Iacoviello	RAI
Angelo Bruccoleri	RAI
Stefano Scotta	RAI
Lorenzo Canale	RAI



Samuel Almeida	F6S
Ellie Shtereva	F6S
Candela Bravo	LOBA
Christoffer Holmgard	MODL
Angel Spasov	IMG
Carmen Mac Williams	GAR

## Peer Reviews

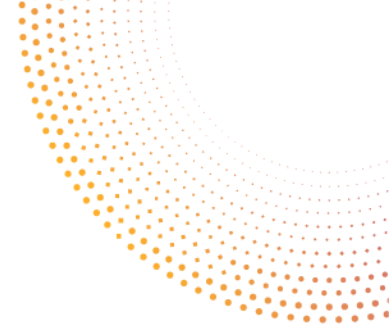
NAME	ORGANISATION
Henning Müller	HES-SO
Danae Tsabouraki	ATC

## Revision History

VERSION	DATE	REVIEWER	MODIFICATIONS
0.1	29/04/2024	Filareti Tsalakanidou, Yiannis Kompatsiaris	First draft sent to partners for contributions
0.2	11/07/2024	Filareti Tsalakanidou	Updated version including inputs from all partners (new datasets added in sections 4, 5, 6)
0.3	15/07/2024	Filareti Tsalakanidou	Updated version including additional datasets
0.4	15/07/2024	Filareti Tsalakanidou	Ready for internal review
0.5	23/08/2024	Henning Müller, Danae Tsabouraki	Internal review
0.6	26/08/2024	Filareti Tsalakanidou	Updated version based on internal review comments
1.0	26/08/2024	Filareti Tsalakanidou, Yiannis Kompatsiaris	Final version

The information and views set out in this report are those of the author(s) and do not necessarily reflect the official opinion of the European Union. Neither the European Union institutions and bodies nor any person acting on their behalf.

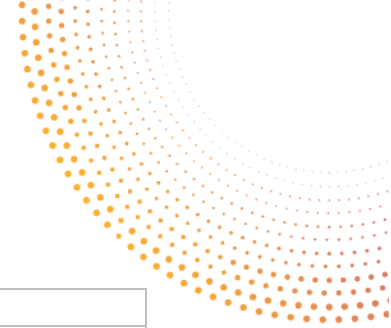




## Table of Acronyms and Abbreviations

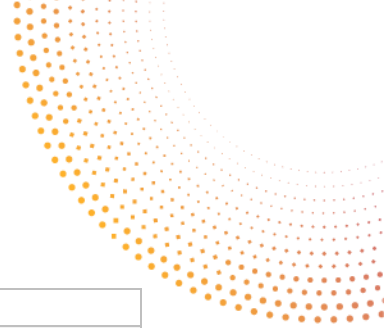
Acronym	Meaning
<b>1D, 2D, 3D, 4D</b>	One-, two-, three-, four-dimensional
<b>2FA</b>	2 Factor Authentication
<b>AI</b>	Artificial Intelligence
<b>AIDA</b>	International AI Doctoral Academy
<b>AIoD</b>	AI-on-Demand Platform
<b>API</b>	Application Programming Interface
<b>ASV</b>	Automatic Speaker Verification
<b>AWS</b>	Amazon Web Services
<b>BSD-2</b>	Berkeley Software Distribution 2
<b>CC</b>	Creative Commons
<b>CLEF</b>	Conference and Labs of the Evaluation Forum
<b>CoMA</b>	Convolutional Mesh Autoencoders
<b>CSV</b>	Comma-Separated Values file format
<b>DB</b>	Database
<b>DFDC</b>	DeepFake Detection Challenge
<b>DFDM</b>	DeepFakes from Different Models dataset
<b>DFEW</b>	Dynamic Facial Expression in-the-Wild dataset
<b>DMP</b>	Data Management Plan
<b>DoA</b>	Description of Action
<b>DOI</b>	Digital Object Identifier
<b>DPO</b>	Data Protection Officer
<b>DW</b>	Deutsche Welle
<b>EC</b>	European Commission
<b>EEA</b>	European Economic Area
<b>EEG</b>	ElectroEncephaloGraphy
<b>ENF</b>	Electrical Network Frequency
<b>EU</b>	European Union
<b>EULA</b>	End User License Agreement
<b>FAIR</b>	Findable, Accessible, Interoperable, Reusable
<b>FaVCi2D</b>	Face Verification with Challenging Imposters and Diversified Demographics dataset
<b>FFHQ</b>	Flickr-Faces-HQ
<b>FoR</b>	Fake-or-Real
<b>FSTP</b>	Financial Support to Third Parties
<b>GAN</b>	Generative Adversarial Networks
<b>GDPR</b>	General Data Protection Regulation
<b>HDF5</b>	Hierarchical Data Format 5
<b>HTTPS</b>	HyperText Transfer Protocol Secure
<b>IAM</b>	Identity and Access Management
<b>ID</b>	Identity
<b>ILSVRC2012</b>	Large Scale Visual Recognition Challenge 2012
<b>ISO</b>	International Organization for Standardization
<b>ISO/IEC</b>	International Organization for Standardization/ International





<b>Acronym</b>	<b>Meaning</b>
	Electrotechnical Commission
<b>IT</b>	Information Technology
<b>JAMS</b>	JSON Annotated Music Specification
<b>JSON</b>	JavaScript Object Notation
<b>JWT</b>	JSON Web Token
<b>KoDF</b>	Korean DeepFake (dataset)
<b>LLM</b>	Large Language Model
<b>LM</b>	Language Model
<b>MAD-TSC</b>	Multilingual Aligned Dataset for Target-dependent Sentiment Classification
<b>MED</b>	Multimedia Event Detection
<b>MRI</b>	Magnetic Resonance Imaging
<b>MSD</b>	Million Song Dataset
<b>N/A</b>	Non applicable
<b>NEC</b>	Non European Countries
<b>NEFER</b>	Neuromorphic Event-based Facial Expression Recognition dataset
<b>NGO</b>	Non-Governmental Organization
<b>NIR</b>	Near InfraRed
<b>NLP</b>	Natural Language Processing
<b>NN</b>	Neural Network
<b>OCR</b>	Optical Character Recognition
<b>ODSS</b>	Open Dataset of Synthetic Speech
<b>OTP</b>	One Time Password
<b>PII</b>	Personally Identifiable Information
<b>POPD</b>	Protection of Personal Data
<b>RDBMS</b>	Relational DataBase Management System
<b>RGB</b>	Red Green Blue
<b>SALAMI</b>	Structural Analysis of Large Amounts of Music Information
<b>SCC</b>	Standard Contractual Clause
<b>SME</b>	Small-Medium Enterprise
<b>SNR</b>	Signal-to-Noise Ratio
<b>SotA</b>	State of the Art
<b>SR</b>	Super Resolution
<b>SSH</b>	Social Sciences and Humanities
<b>SSL</b>	Secure Sockets Layer
<b>TF</b>	Tensorflow
<b>TSV</b>	Tab-Separated Values
<b>TTS</b>	Text to Speech
<b>UC</b>	Use Case
<b>UGC</b>	User Generated Content
<b>URL</b>	Uniform Resource Locator
<b>UTC</b>	Coordinated Universal Time
<b>UTF-8</b>	8-bit Unicode Transformation Format
<b>VPN</b>	Virtual Private Network
<b>WAV</b>	Waveform Audio file format
<b>WDF</b>	WildDeepFake dataset

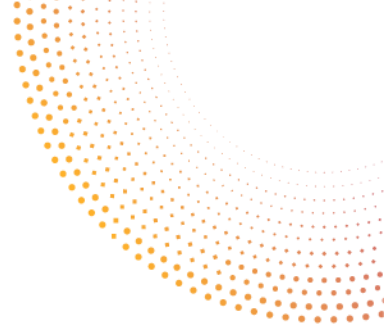




<b>Acronym</b>	<b>Meaning</b>
<b>WIPO</b>	World Intellectual Property Organization
<b>WNID</b>	WordNet ID
<b>WP</b>	Work Package
<b>XAI</b>	Explainable Artificial Intelligence
<b>XML</b>	eXtensible Markup Language



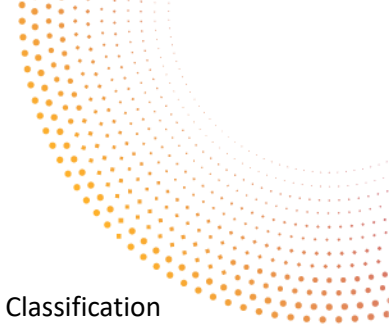




## Index of Contents

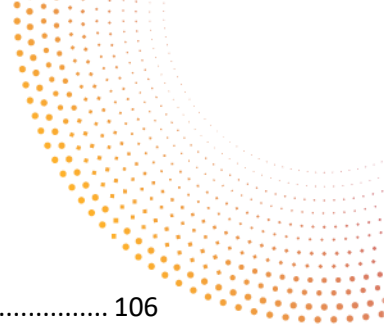
Table of Acronyms and Abbreviations .....	6
Index of Contents .....	9
Index of Tables .....	16
1. Executive Summary .....	17
2. Introduction.....	20
3. Data management methodology .....	21
3.1 Data summary .....	21
3.2 Making data findable, including provisions for metadata .....	25
3.3 Making data openly accessible.....	27
3.4 Making data interoperable.....	29
3.5 Increase data re-use (through clarifying licenses) .....	30
3.6 Allocation of resources.....	32
3.7 Data security.....	33
3.8 Ethical & legal aspects.....	35
3.9 Other issues.....	39
4. Data management plan for research datasets created within AI4Media.....	40
4.1 Datasets collected in the context of WP2 .....	45
4.1.1 Questionnaires for AI technology roadmap.....	45
4.1.2 Workshop data for AI technology impact and policy.....	46
4.2 Datasets collected in the context of WP3 .....	47
4.2.1 FaVCI2D image dataset for demographically diversified face verification .....	47
4.2.2 Mixamo-Kinetics dataset.....	49
4.2.3 100-Driver dataset for distracted driver classification.....	50
4.2.4 LeQua 2022 datasets.....	52
4.2.5 Product reviews for ordinal quantification dataset .....	53
4.2.6 Product reviews dataset.....	55
4.2.7 UCI and OpenML datasets for ordinal quantification .....	56
4.2.8 Cherenkov telescope data for ordinal quantification .....	57
4.3 Datasets collected in the context of WP4 .....	59
4.3.1 White matter multiple sclerosis lesion segmentation datasets.....	59
4.4 Datasets collected in the context of WP5 .....	60
4.4.1 ObyGaze12 dataset for detection of visual objectification in films .....	60





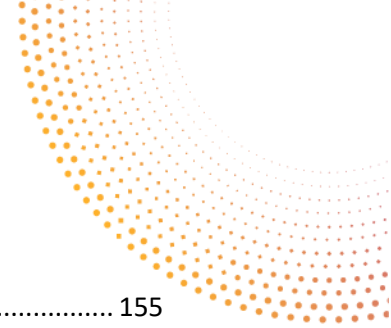
4.4.2	MAD-TSC Multilingual Aligned Dataset for Target-dependent Sentiment Classification	61
4.4.3	ÖWF Object Detection dataset .....	63
4.4.4	PeopleAtPlaces dataset.....	64
4.4.5	ToDY dataset for visual time of day and season classification.....	65
4.4.6	VISIONE Feature Repository.....	67
4.4.7	COCO, LVIS, Open Images V4 classes mapping .....	68
4.4.8	Bus violence dataset.....	69
4.4.9	Pest Sticky Traps dataset.....	71
4.4.10	Virtual World Fallen People dataset .....	72
4.4.11	SR_BVI-DVC super-resolution dataset.....	73
4.4.12	BSC4K super-resolution dataset.....	75
4.4.13	Night and Day Instance Segmented Park dataset .....	76
4.4.14	ArXiv abstracts for authorship analysis dataset .....	77
4.4.15	Florence4D facial expression dataset.....	78
4.4.16	Neuromorphic Event-based Facial Expression Recognition dataset.....	79
4.4.17	PEM360 dataset of 360° videos .....	81
4.4.18	VRT-Sum video summarization dataset .....	82
4.4.19	CA-SUM pretrained video summarization models.....	83
4.4.20	Audio phylogeny dataset.....	85
4.4.21	RAI CMM documentaries dataset .....	86
4.4.22	RAI CMM-ANTS newscasts dataset .....	87
4.4.23	RAI CMM mixed dataset.....	88
4.4.24	Cross lingual news dataset .....	89
4.4.25	YouTube RAI channel dataset .....	90
4.5	Datasets collected in the context of WP6 .....	91
4.5.1	GreekPolitics Twitter dataset .....	91
4.5.2	Covid-19 Twitter dataset.....	94
4.5.3	Twitter Text dataset .....	97
4.5.4	Twitter COVID-19 discussions topics dataset.....	99
4.5.5	ODSS Open Dataset of Synthetic Speech .....	101
4.5.6	CelebHQGaze image dataset for gaze estimation.....	102
4.5.7	European news about Covid vaccination dataset .....	103
4.5.8	Suisse Romande local news dataset.....	105





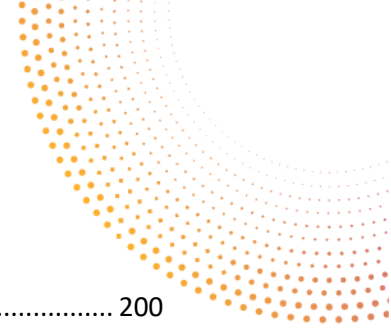
4.5.9	Lausanne news dataset .....	106
4.5.10	Suisse Allemande local news dataset.....	107
4.5.11	References on YouTube dataset .....	108
4.5.12	Political Barometer dataset.....	110
4.5.13	ElecDeb60To16-fallacy dataset.....	111
4.6	Datasets collected in the context of WP8 .....	112
4.6.1	Data from user research activities in Use Case 1 .....	112
4.6.2	Data from user research activities in Use Case 2 .....	114
4.6.3	Questionnaires for the collection of user requirements for Use Case 3.....	116
4.6.4	Questionnaires for the evaluation of Use Case 3.....	117
4.6.5	Data from user research activities in Use Case 4 .....	119
4.6.6	Questionnaires for the evaluations of Use Case 5-B.....	120
4.6.7	Data from user research activities in Use Case 7 .....	122
4.6.8	Truly Media dataset .....	123
4.6.9	Game glitches dataset for Use Case 5.....	125
4.6.10	Musical production for AI co-creation dataset .....	126
4.6.11	Current affairs transcripts dataset for Use Case 4 .....	127
4.6.12	User survey data for AI industrial needs (for T8.4) .....	129
5.	Data management plan for third-party research datasets used in AI4Media .....	132
5.1	Datasets used in the context of WP3 .....	135
5.1.1	Enron email dataset .....	135
5.1.2	Facebook wall dataset.....	137
5.1.3	Affect Game Annotation (AGAIN) dataset .....	138
5.1.4	IMDB movie reviews dataset .....	139
5.1.5	MHAD 2D pose dataset .....	141
5.1.6	20Newsgroups dataset.....	142
5.1.7	HP Amazon reviews dataset.....	143
5.1.8	JRCAcquis legislative text dataset .....	145
5.1.9	Kindle document dataset .....	146
5.1.10	OHSUMED MEDLINE document dataset .....	148
5.1.11	RCV1 Reuters stories dataset .....	149
5.1.12	RCV1RCV2 Reuters stories dataset .....	151
5.1.13	Reuters-21578 dataset .....	152
5.1.14	11 Tweet Sentiment datasets .....	154





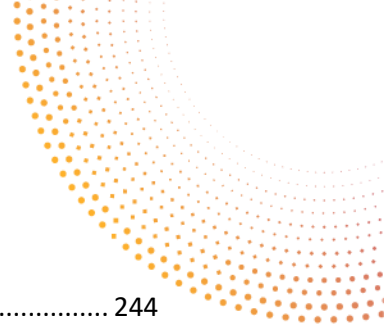
5.1.15	WipoGamma patent document dataset .....	155
5.1.16	WIDER FACE face detection dataset .....	157
5.1.17	Caltech 101 dataset.....	158
5.1.18	CUB-200-2021 image dataset.....	159
5.1.19	Describable Textures dataset.....	161
5.1.20	Food-101 dataset .....	162
5.1.21	MAMe dataset.....	163
5.1.22	MIT-ISR dataset .....	164
5.1.23	Oulu Knots dataset.....	165
5.1.24	Oxford Flower dataset.....	166
5.1.25	Oxford-IIIT Pet dataset .....	168
5.1.26	Stanford Dogs dataset.....	169
5.2	Datasets used in the context of WP4 .....	170
5.2.1	ImageNet-ILSVRC2012 image classification dataset .....	170
5.2.2	FFHQ dataset for GAN training.....	172
5.2.3	MNIST image dataset .....	173
5.2.4	Interestingness10k image +video dataset.....	174
5.2.5	MediaEval Memorability 2020 dataset .....	176
5.2.6	ImageCLEF DrawnUI 2021 dataset .....	177
5.2.7	FCVID event recognition dataset.....	178
5.2.8	YLI-MED event recognition dataset.....	180
5.3	Datasets used in the context of WP5 .....	181
5.3.1	SumMe video summarization dataset .....	181
5.3.2	TVSum video summarization dataset.....	183
5.3.3	RAI Monuments of Italy dataset .....	184
5.3.4	LVIS image dataset .....	185
5.3.5	CIFAR10/100 image dataset.....	187
5.3.6	STL-10 image dataset .....	188
5.3.7	CCNet text dataset for multilingual representation learning.....	189
5.3.8	360 Video Viewing Dataset in Head Mounted Virtual Reality .....	191
5.3.9	Predicting Head Movement in Panoramic Video Dataset.....	193
5.3.10	Gaze prediction in Dynamic 360° Immersive Videos Dataset .....	195
5.3.11	Your Attention is Unique Dataset .....	196
5.3.12	Dataset of Head and Eye Movements for 360° Videos .....	198





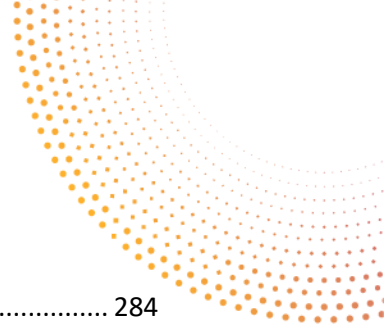
5.3.13	Dim-sim dataset for music similarity search .....	200
5.3.14	SPAM dataset for music segmentation .....	201
5.3.15	SALAMI dataset for music segmentation .....	203
5.3.16	Harmonix dataset for music segmentation.....	204
5.3.17	Free Music Archive dataset.....	205
5.3.18	LAKH MIDI music dataset .....	206
5.3.19	Piano Audio and MIDI music datasets.....	208
5.3.20	GiantSteps music datasets .....	209
5.3.21	MS COCO dataset .....	210
5.3.22	BVI-DVC dataset .....	212
5.3.23	Adience dataset.....	213
5.3.24	IMDB-Wiki dataset .....	214
5.4	Datasets used in the context of WP6 .....	215
5.4.1	Deepfake Detection Challenge dataset.....	215
5.4.2	FaceForensics++ dataset .....	217
5.4.3	Visual profile impact rating and ranking – ImageCLEFaware dataset.....	218
5.4.4	DEAP EEG dataset.....	220
5.4.5	SEED EEG dataset .....	221
5.4.6	SEED-IV EEG dataset.....	223
5.4.7	Clotho audio captioning dataset .....	224
5.4.8	ASVspoof2019 dataset .....	225
5.4.9	MOBIPHONE audio dataset.....	227
5.4.10	Fake-or-Real (FoR) audio dataset.....	228
5.4.11	ForenSynths dataset.....	230
5.4.12	SynthBuster dataset .....	231
5.4.13	Latent Diffusion Training-set.....	232
5.4.14	Diffusion datasets.....	233
5.4.15	FakeAVCeleb dataset.....	234
5.4.16	ForgeryNet dataset .....	236
5.4.17	Korean DeepFake dataset .....	237
5.4.18	WildDeepFake dataset .....	238
5.4.19	DeepFakes from Different Models (DFDM) dataset .....	240
5.4.20	DF-Platter dataset .....	241
5.4.21	MAFW dataset.....	242





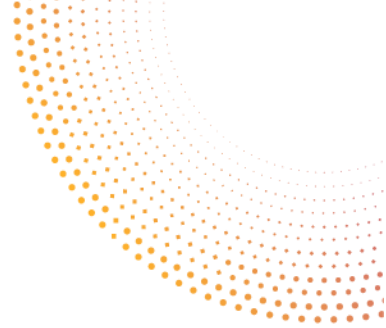
5.4.22	Dynamic Facial Expression in-the-Wild (DFEW) video dataset .....	244
5.4.23	FERV39k video dataset .....	245
6.	Data management plan for non-research datasets collected in AI4Media .....	248
6.1	Datasets collected in the context of WP1 .....	249
6.1.1	AI4Media consortium contact info dataset .....	249
6.2	Datasets collected in the context of WP2 .....	251
6.2.1	AI Media Observatory Dataset .....	251
6.2.2	Curated dataset of resources on how the media sector responds to content crawling for AI model training .....	253
6.3	Datasets collected in the context of WP7 .....	254
6.3.1	AI-Cafe mailing list members and AI-Cafe participants.....	254
6.3.2	Candidate AI assets dataset .....	255
6.4	Datasets collected in the context of WP8 .....	256
6.4.1	User data from Truly Media for Use case 1.....	256
6.5	Datasets collected in the context of WP9 .....	258
6.5.1	AIDA course offerings dataset.....	258
6.5.2	AIDA students dataset.....	259
6.5.3	AIDA mailing list dataset .....	261
6.5.4	AIDA AI educational resources dataset.....	262
6.5.5	AIDA lecturers dataset .....	263
6.5.6	AIDA AI Excellence Lecture Series dataset .....	265
6.5.7	AIDA curators dataset .....	266
6.5.8	AIDA website analytics dataset .....	267
6.5.9	AI4Media Junior Fellows exchange program dataset .....	268
6.6	Datasets collected in the context of WP10 .....	270
6.6.1	Competitive call application datasets .....	270
6.6.2	Sub-granted projects dataset.....	272
6.6.3	External experts and evaluators datasets .....	274
6.6.4	Participants of competitive call related events datasets.....	276
6.7	Datasets collected in the context of WP11 .....	278
6.7.1	AI4Media associate members contact info dataset.....	278
6.7.2	AI4Media newsletter subscribers dataset.....	279
6.7.3	AI4Media website messages dataset.....	281
6.7.4	AI4Media website analytics dataset .....	282





6.7.5	Dataset of registrants for AI4Media events .....	284
7.	Conclusions.....	286





## Index of Tables

Table 1: Template for the presentation of the data management plan for a specific dataset...	40
Table 2: Summary of research datasets created within AI4Media .....	42
Table 3: Summary of third-party research datasets used in AI4Media .....	132
Table 4: Summary of non-research datasets collected in AI4Media .....	248





## 1. Executive Summary

Various datasets have been used, collected or generated during the lifetime of the AI4Media project in order to pursue the project's research agenda and accomplish the project's objectives, as described in the Description of Action (DoA). A variety of data (videos, images, audio, text, social media, user profile data, questionnaires, contact info, etc.) have been used, collected or generated aiming to:

- define the user requirements and use case scenarios;
- develop cutting-edge Artificial Intelligence (AI) technologies for specific media-related fields as well as for human-centered and society-centered AI, in the context of the seven use cases;
- assess the effectiveness of the developed AI4Media technologies in a series of evaluation activities involving end-users;
- establish and support the International AI Doctoral Academy (AIDA); and
- establish the AI4Media Network by attracting and involving AI researchers and SMEs through open call procedures.

More specifically, we can distinguish among the following types of research data that were collected at different stages of the project:

- First, the leaders of the seven AI4Media use cases collected the **use case requirements**. This involved feedback from end users via questionnaires, interviews, focus groups and other similar methods, aiming to identify user needs with regard to the functionalities of the AI tools to be integrated in each use case. The objective of collecting such data was to define the user and use case requirements and guide the technical development of the AI4Media tools.
- In order to **develop novel AI methodologies and tools** to advance AI research for the media and also support the seven AI4Media use cases, multiple datasets (existing or new ones) were used or created by technical partners in the context of WP3, WP4, WP5 and WP6. To cover the needs of the different use cases (AI for social media, news, vision, games, social sciences and humanities, human co-creation, and automatic content organisation and moderation), heterogeneous datasets were used or generated, including **video, images, audio, text, social media posts, user activity data**, etc. To this end, several existing datasets were deployed, either owned by AI4Media technical and use-case partners or openly shared by third parties (usually academic or research institutions). In addition, several new datasets were also created in the context of the project, e.g., datasets including tweets for disinformation detection related to specific news events.
- Finally, as part of the **use cases evaluation**, data was collected from the end users of the proposed AI solutions, including both user activity data generated during the use of the various AI4Media tools as well as survey data (e.g., questionnaires) aiming to assess the impact and effectiveness of the proposed AI technologies and use case demonstrators.



In addition to the research data described above, non-research datasets that help us establish the operation of the International AI Doctoral Academy, the AI4Media Network and the two Open Calls were also collected. These mainly include:

- Personal data of lecturers, researchers and post-graduate students collected as part of registration processes, course attendance or instruction, and curation of educational resources in the context of the **International AI Doctoral Academy (AIDA)**.
- **Applicant data** collected by the platform used for submitting applications to get funding via the two AI4Media open calls. Data for the selected **sub-granted projects** and information of **internal experts and evaluators** involved in this process was also gathered.
- Information about AI4Media's **Associate Members** collected to facilitate their involvement in the project and also about researchers contributing to the **AI Media Observatory**. Also, contact information of people participating in AI4Media's events, workshops or other **dissemination & communication activities**.

The aforementioned (research and non-research) data requires a clear plan on how it is going to be managed, i.e., stored, accessed, shared, protected against unauthorized or improper use, etc. The main goals of AI4Media's Data Management Plan (DMP) are to:

1. Outline the datasets used/collected/generated in the project, including the context and procedures of the use/collection/generation, as well as the degree of privacy and confidentiality of the data;
2. Outline the procedures for FAIR (findable, accessible, interoperable, reusable) data;
3. Outline the measures that are foreseen for the adequate management of the data from an ethical and a security point of view.

The scope of the DMP is to describe the data management life cycle for all datasets to be used or collected in all Work Packages (WP) during the course of the 48 months of the AI4Media project. An initial version of the DMP (D1.2) was submitted at the beginning of the project, in M6, reflecting progress and expectations at that point in time. The initial DMP offered information on the general data management policy to be followed by the project using the "Template for HORIZON 2020 DATA MANAGEMENT PLAN (DMP)" (version 1.0, released on 13.10.2016 by the European Commission) and answered explicitly the questions listed there. Moreover, it listed the datasets that had already been collected or used by the partners at that point, as well as datasets that project partners foresaw to use or generate during the course of the project. For each dataset, explicit information was provided with regard to data use, data collection, data discoverability, data accessibility, data interoperability and metadata, data re-usability and open data, data security, and relevant ethical & legal aspects of data processing.

This final version of the DMP offers an update of the initial plan and extends the initial list of research and non-research datasets by adding several new datasets that were used, collected or generated in the project between M7 and M48. The initial version of the DMP discussed 82 datasets in total: 70 research datasets (19 created within the project by project partners and 51 third-party research datasets used by project partners) as well as 12 non-research datasets. **This final version of the DMP presents 166 datasets: 142 research datasets (61 created within**



**the project** by project partners and 81 third-party research datasets used by the partners) **and 24 non-research datasets.**

This deliverable was compiled through the collaborative work of the project coordinator and the consortium partners who are involved in data collection, production, and processing. Each consortium partner received a call-to-action to identify the datasets most relevant to their respective activities and deliverables.

## 2. Introduction

This deliverable presents the final Data Management Plan (DMP) of the AI4Media project and aims to provide a detailed description of the datasets that were collected, processed or generated during the course of the project. It describes the handling of research and non-research data during and after the project lifetime and discusses how they will be curated and preserved. It also provides information about which datasets are or will be openly accessible and how they are or will be shared, while also presenting standards used to increase data interoperability.

The DMP provides an analysis of the main elements of the data management policy that was adopted by the AI4Media consortium with regard to all the datasets that were used in or generated by the project. It reflects the consortium's comprehensive approach towards data management. The DMP is a living document which evolved during the lifetime of the project. An initial version (D1.2) was delivered in M6 while this final version (D1.6) is delivered at the end of the project, in M48, and reflects the following changes: i) an update of the general data management strategy of the project (section 3) and especially the data summary (section 3.1), data security (section 3.7), and ethical & legal aspects (section 3.8); ii) the collection, creation and use of new datasets that were not included or foreseen in the initial DMP; more specifically, D1.2 presented 82 (research and non-research) datasets while D1.6 presents 166 datasets in total (sections 4-6).

The remainder of the deliverable is structured as follows.

**Section 3** describes the methodology followed for drafting the DMP and provides a general overview of the project's data management policy.

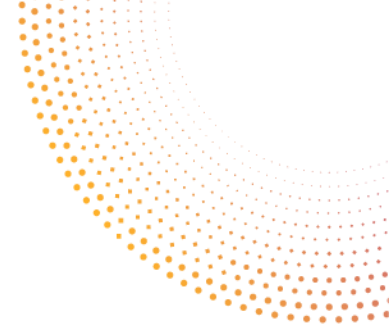
**Section 4** summarizes the management plan for the various research datasets created within AI4Media, following the aforementioned methodology. 61 research datasets are presented in total.

**Section 5** addresses the management of research datasets that are used within AI4Media but have been created by third parties, e.g. public or private research datasets used by the partners of WPs 3, 4, 5, and 6 to develop and test new AI methodologies, algorithms and tools. 81 third-party research datasets are presented in total.

**Section 6** describes how we handle non-research data collected within AI4Media, e.g. data collected from open calls, in the context of AIDA, from dissemination activities, etc. 24 non-research datasets are presented in total.

Finally, **Section 7** concludes the deliverable.





### 3. Data management methodology

The methodological approach that has been used for the compilation of both D1.2 and D1.6 follows the “*Template for HORIZON 2020 DATA MANAGEMENT PLAN (DMP)*”, version 1.0, released on 13.10.2016 by the EC. Taking into account the proposed methodology, the AI4Media DMP addresses the following points on a dataset-by-dataset basis:

- Data summary;
- FAIR data
  - Making data findable, including provisions for metadata;
  - Making data openly accessible;
  - Making data interoperable;
  - Increase data re-use;
- Allocation of resources;
- Data security;
- Ethical aspects;
- Other issues.

In the following subsections of section 3, we briefly present the kind of questions associated with each point in this list. Moreover, for each question we provide a **summary of the general strategy** adopted by the project consortium for handling different dataset categories. **Detailed answers to these questions on a dataset-by-dataset basis** (i.e. for each identified dataset individually) are provided in sections 4, 5 and 6.

#### 3.1 Data summary

The Data Summary addresses the following issues:

- Outline the purpose of the collected/ generated data and its relation to the objectives of the AI4Media project;
- Outline the types and formats of data already collected/ generated and/ or foreseen for generation at this stage of the project;
- Outline the reusability of the existing data;
- Outline the origin of the data;
- Outline the expected size of the data;
- Outline the data utility.

In this field, the data that will be generated or collected is described, including references to their origin (in cases where data is collected), nature, scale, to whom it could be useful, and whether it underpins a scientific publication. With regard to the individual questions, our generic DMP approach is summarized below (detailed answers for each dataset are given in sections 4, 5 and 6).

**What is the purpose of the data collection/generation and its relation to the objectives of the project?**

The main goal of AI4Media is to become a centre of excellence and a wide network of researchers across Europe and beyond, with a focus on delivering the next generation of core AI advances to serve the key sector of media, to make sure that the European values of ethical



and trustworthy AI are embedded in future AI deployments, and to re-imagine AI as a crucial beneficial enabling technology in the service of Society, Democracy and Media. This goal is achieved through six AI4Media pillars:

- The **European Media AI Observatory**, which sets and maintains a research and innovation agenda for media AI, while anticipating the social and economic disruptive potential of emerging technologies (WP2).
- An **intensive research and innovation plan** in core areas of Media AI where Europe has or can acquire a competitive advantage, generating technologies which will enrich the AIoD platform (WP3, 4, 5, 6).
- A **portfolio of use-cases**, aiming to provide direct application of the AI4Media technologies developed within WP3-WP6 to strengthen the competitiveness of European businesses in the broader media sector and the European society (WP8).
- A **targeted programme of cascade funding** to increase engagement of actors outside the consortium and build an ecosystem around the network, in turn benefiting from it and bringing innovation to the market (WP10).
- The **International AI Doctoral Academy (AIDA)**, which will nurture a new generation of PhD talent and ensure young skilled researchers remain in Europe (WP9).
- The **AI4Media Virtual Center of Excellence**, which will function as a portal and network nexus for all Media AI research and innovation activities in Europe (WP11).

To achieve the project's main goal and relevant objectives the following types of datasets are used, collected, or generated:

- **Use case requirements data** (in the form of questionnaires, interviews, focus groups, etc.) is collected by partners involved in use cases in the context of WP8 to identify user needs, use-case scenarios, and desired software functionalities. The objective of collecting such data is to define the user-based system requirements and guide the design and development of the various AI methodologies, tools, and demonstrators that will support the seven AI4Media use cases:
  - *AI for Social Media and against Disinformation*: AI technology for detecting disinformation in social media, aiming to support journalistic fact-checking and verification workflows in news organisations;
  - *AI for News*: Smart News Assistant AI solutions to help journalists ensure that published content is both relevant for its audience and a trustworthy source of information;
  - *AI in Vision*: AI algorithms and tools for high quality video production and video content automation in TV production;
  - *AI for Social Sciences and Humanities*: AI tools that allow researchers and journalists to sift, connect, and analyze various data and media collections in search of factual responses to broad societal research questions;
  - *AI for Games*: AI tools to advance game design and development process, focusing on i) improved music analysis and synthesis for games, and ii) game testing and quality assurance;
  - *AI for Human Co-creation*: AI-based audio generation methods to help music composers create music;
  - *AI for (re-)Organisation and Content Moderation*: AI-based advanced content moderation solutions for media companies and AI technologies to effectively



organize vast digital archives and collections of images and videos.

Details on these datasets can be found in section 4.6.4.5. Compliance with the legal obligations for processing personal data of users in user research activities are ensured. For more information about this, please have a look at Ethics deliverables D12.1 and D12.2 and the Ethics Management Plan (D1.3, D1.5).

- **Evaluation data** is collected from the end users of the proposed AI solutions in the context of the use cases and WP8, including i) user activity data automatically generated during the use of the various AI4Media tools, ii) data collected during the use cases implementation to assess the evolution of the users, and iii) user feedback data (e.g. from questionnaires, focus groups, etc.) aiming to assess the impact and effectiveness of the proposed AI technologies and also to guide industry partners in harmonising AI research with industrial needs. Details on these datasets can be found in section 4.6.
- **Media-related datasets** (existing or new ones) are used or collected by technical partners in the context of WPs 3, 4, 5, 6 in order to develop and test novel AI algorithms and methods for the media, with a focus on the seven AI4Media use cases. To develop algorithms and tools to address the different needs of each use case, heterogeneous datasets are used or generated, including video, images, audio, text, social media posts, user profiles, user/user group activity data, etc. Several existing datasets are also deployed and re-used, either owned by AI4Media technical and use-case partners or openly shared by third parties (usually academic or research institutions). Details on these datasets can be found in sections 5.1-5.4. In addition, new datasets were also created and used, e.g., datasets including tweets for disinformation detection or datasets including news items for local news analysis. Details on these datasets can be found in sections 4.2-4.5.
- **User survey data** is collected in the context of WP2 to analyse EU AI policy, create an AI roadmap, assess the social/economic/political impact from future advances in media AI technology, and draft relevant policy recommendations. Details on these datasets can be found in section 4.1.
- **Personal data** and other information collected from researchers contributing to the AI Media Observatory (WP2). Details on these datasets can be found in section 6.2.
- **Personal data (e.g., name, email, affiliation, etc.) of lecturers, researchers and students** is collected in the context of the AIDA (WP9) as part of their registration process and course attendance (for students), course or lecture delivery (for lecturers) and curation of educational resources (researchers/curators). Details on these datasets can be found in section 6.5.
- **Data of applicants** (individual researchers and organizations) that used the platform set up by F6S in the context of WP10 to submit applications to get funding from **AI4Media's two open calls**. Data for the selected sub-granted projects and information of internal experts and evaluators is also gathered. Details on these datasets can be found in section 6.6 .
- **Personal data of participants in AI4Media's dissemination and community building events** (WP11, WP7) were also collected where necessary to allow better organization of these events and better services to attendees. Data of associate members are also



collected in the context of WP11 via online questionnaire forms. Details on these datasets can be found in sections 6.7 and 6.3.

- **Personal data of members of the consortium** to facilitate project management and coordination. Details on these datasets can be found in section 6.1.

#### **What types and formats of data will the project generate/collect?**

As mentioned above, the project uses different types of data (video, images, audio, social media posts, system log data, questionnaires, applicant personal data, etc.) from various sources and also generated datasets including data of the aforementioned types. The data is available in various formats (e.g., json and csv files for social media data; mpeg and avi files for video; excel and doc files for questionnaire data; etc.) and diverse databases (MySQL, Mongo DB etc.). Again, detailed answers for each dataset are given in sections 4, 5 and 6.

#### **Will you re-use any existing data and how?**

The datasets presented in section 5 (media-related datasets used in WPs 3, 4, 5, 6) include existing public or private datasets available from project partners, third-party research/academic organizations, or third-party media companies. The data will be re-used in AI4Media to develop and test innovative AI algorithms, methodologies and software tools, aiming to support the needs of the seven pilots, as described above.

The data presented in section 4 is new research data generated by the project.

Beyond direct project purposes, all project data will be used for scientific publications (except from the personal data collected from participants or applicants involved in various AI4Media education, dissemination or open call activities presented in section 6; such data have no scientific value).

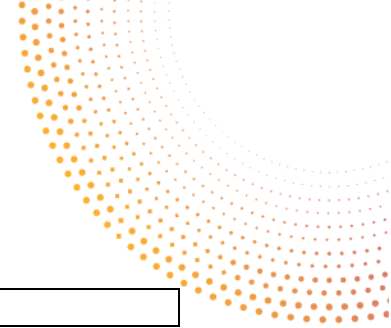
#### **What is the origin of the data?**

The data originates from various sources:

- Individual researchers that openly share their data in open repositories such as GitHub and Zenodo or via their webpages.
- Research and academic organizations that openly share data in open repositories or institutional repositories.
- Use case partners (media organizations) that share existing datasets with the technical partners of the consortium to help them train and test their algorithms and software.
- Social media APIs (e.g. Twitter, etc.).
- Web (e.g., online news articles & comments).
- Online forms filled in by applicants or participants in the context of AI4Media education, dissemination and open call activities.
- Questionnaires filled in by project partners (e.g., user requirements questionnaires), by end-users of AI4Media tools (e.g., evaluation questionnaires), etc.
- Use of AI4Media software tools by the end users during the AI4Media use case trials (this includes automatically collected user/software analytics).

#### **What is the expected size of the data?**





Dataset sizes are discussed on a dataset-by-dataset basis in sections 4, 5 and 6.

**To whom might it be useful ('data utility')?**

The aforementioned data is useful to project partners for identifying user and software requirements for the various use cases; designing, developing, and testing (and improving) the AI4Media methodologies, algorithms, and tools of WPs 3,4,5,6; and assessing the effectiveness of these tools in media use cases involving end users. Also, it is useful for the better design, organization, and execution of activities related to AIDA, AI Media Observatory, open calls for third-party funding, outreach and dissemination.

Media-related datasets and evaluation data may also be useful to researchers with a focus on the development of AI technologies for the media in general or specific aspects of media in particular (e.g., music for games, deepfake detection, or content moderation, just to name a few). The data can also be useful to social scientists that want to examine the impact of AI on media or the impact of media on the society (e.g., in the context of elections).

**3.2 Making data findable, including provisions for metadata**

This point addresses the following issues:

- Are the data produced and/ or used in the project discoverable and identifiable?
- What naming conventions are followed?
- Will search keywords be provided that optimize possibilities for re-use?
- Are clear version numbers provided?
- What metadata will be created?

In general, most of the data produced or used by the project is or will be identifiable and discoverable. With regard to the individual questions, our generic DMP approach is summarized below (again, detailed answers for each dataset are given in in sections 4, 5 and 6):

**Are the data produced and/ or used in the project discoverable and identifiable?**

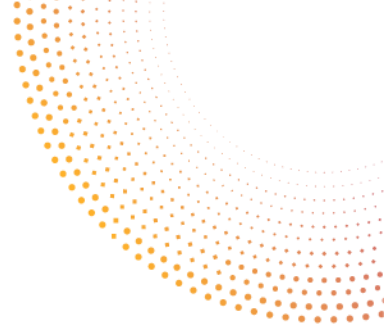
Publicly available third-party datasets that we will re-use to develop and test our AI technologies are already easily discoverable and identifiable from the original sources.

Datasets that are created within AI4Media and are made publicly available by partners are uploaded to open repositories like Zenodo or GitHub, thus making this data both easily discoverable and identifiable from the outside (since they are assigned a DOI). These datasets are also shared through the AIoD platform. Some datasets are shared through AI4Media partners' institutional (open) repositories.

With regard to datasets that will only be used internally in the project, some of these will be stored on partners' servers. This data will only be discoverable and identifiable from registered users, using simple queries with keywords.

Other internal data such as user requirement data and evaluation data are stored on user partners' servers and access is provided only to selected institutional users involved in the processing of this data.





### What naming conventions are followed?

A specific naming convention is used to identify the various AI4Media datasets:

*AI4Media\_Data\_<serial number of dataset>\_<WPno>\_<data type>\_<dataset title/ID>\_v<version no>*

- The <serial number of each dataset> is assigned manually in the order of presentation in this deliverable.
- The <WPno> reveals the WP in the context of which this data is collected or generated and processed.
- The <data type> takes one of the following values: *SOCIALMEDIA, VIDEO, AUDIO, EMAIL, TEXT, EEG, EMAIL, QUESTIONNAIRE, INTERVIEW, USER-RESEARCH, ACTIVITY-LOG, DEMOGRAPHIC<sup>1</sup>, SURVEY...*
- The <dataset title> is a descriptive title showing what kind of data is included in the dataset, e.g. “*Deepfake-Detection-Challenge-Dataset*”, “*UseCase1-UserReq2021*”.
- The <version no> is the dataset version. Different updated versions of the same dataset may be generated during the project lifetime.

For example, a dataset including questionnaires from the evaluation of the use case 3 demonstrator can be named *AI4Media\_Data\_47\_WP8\_QUESTIONNAIRE\_UseCase3-UserReqCollection2021\_v1*. For a third-party dataset which is named by its creators SumMe<sup>2</sup> and includes videos and relevant annotation data, we can use the following name: *AI4Media\_Data\_91\_WP5\_VIDEO\_SumMeGycli14\_v1*.

### Will search keywords be provided that optimize possibilities for re-use?

Keywords will be provided in the cases where this is applicable.

### Are clear version numbers provided?

For datasets that are made publicly available by the AI4Media partners in open repositories, versioning is supported by these repositories.

Versioning is also supported by the project wiki.

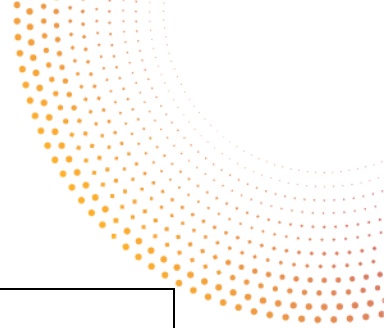
### What metadata will be created?

For datasets that are shared via open repositories, the metadata standards used by these repositories will be used.

<sup>1</sup> *DEMOGRAPHIC* refers to personal data of applicants, participants, attendees, etc. collected through online forms or docs in the context of our education, dissemination and open call activities.

<sup>2</sup> <https://gyglim.github.io/me/vsum/index.html#benchmark>





Metadata for data uploaded at the project wiki is also supported.

### 3.3 Making data openly accessible

This point addresses the following issues:

- Which data produced and/ or used in the project will be made openly available as the default?
- How will the data be made accessible (e.g. by deposition in a repository)?
- What methods or software tools are needed to access the data?
- Is documentation about the software needed to access the data included? Is it possible to include the relevant software (e.g. in open source code)?
- Where will the data and associated metadata, documentation and code be deposited? Have you explored appropriate arrangements with the identified repository?
- If there are restrictions on use, how will access be provided?
- Is there a need for a data access committee?
- Are there well-described conditions for access (i.e. a machine-readable license)?
- How will the identity of the person accessing the data be ascertained?

With regard to the individual questions about data accessibility, our generic DMP approach is summarized below (again, detailed answers for each dataset are given in sections 4, 5 and 6:

#### **Which data produced and/ or used in the project will be made openly available as the default?**

Many of the datasets used in the project is open data already made openly available by third parties (see section 5). More specifically,

- by individual researchers or research/academic organizations (e.g. benchmark video datasets, audio datasets, etc.);
- by media companies (e.g. archive of news articles from a news organization like Reuters etc.).

Since this data is already open, as a general policy, we will not re-share it. Sharing some of this data will be handled on a case-by-case basis, and will only be pursued in cases where the data license allows it and AI4Media researchers estimate that re-sharing of the data (in some new form) provides some additional benefit for the research or industrial community. In any case, we intend to provide open software tools shared on platforms like GitLab or GitHub that will allow other researchers to easily crawl and collect data from all the open data sources used in AI4Media.

With regard to datasets created within the project (see section 4), several of these datasets are already publicly shared on open repositories like Zenodo. They can also be accessed through links available on the [project website](#)<sup>3</sup>.

In addition to open data, there are also privately owned datasets. Most of these datasets are owned by the media companies and news organizations involved in AI4Media as well as by

<sup>3</sup> AI4Media open datasets: <https://www.ai4media.eu/open-datasets/>



some technical and academic partners and have been usually collected and created over a period of years or in the context of other projects or internal processes. Such data (e.g. RAI's video dataset of Italian monuments – see section 5.3.3) will be provided to the project for research purposes, but will not be shared openly, and will be only used internally by project partners. However, effort will be made to make some of this (or part of this) data openly available in cooperation with the data owners, wherever this is legally possible.

Data that is collected by the project in the form of questionnaires or other survey methods addressed to project partners (e.g. to collect user requirements, see section 4.5.5), end users (to evaluate the AI4Media technologies, see section 4.5.5), and other parties (e.g. applicant data requesting AI4Media funding, see section 6.6), will not be made openly accessible since they may contain personal or sensitive information. Where possible and in case there is added value from their sharing (mainly for evaluation data), data will be anonymized before being shared.

In all cases, the aforementioned data (whether public, private, or personal) are processed and analysed aiming to achieve the project objectives. Where appropriate, the analysis results are made open as part of public project deliverables or publications available in open repositories.

#### **How will the data be made accessible (e.g. by deposition in a repository)?**

Data to be openly shared will be deposited in open repositories like Zenodo but also GitHub or GitLab. The datasets will also be shared through the AIoD platform. Some datasets will also be shared on AI4Media partners' institutional repositories. Links to all open datasets are included in the project website in a [dedicated page](#).

Datasets that will be only used internally by project partners will be stored either on the project wiki on CERTH's servers or in the servers of project partners.

#### **What methods or software tools are needed to access the data?**

Different methods and software tools are required to access the data depending on the dataset as is explained in sections 4, 5 and 6 (e.g. web-browser, API).

*In some cases, to collect the data from the original sources, dedicated crawlers were developed by partners, which will be publicly shared in open repositories whenever possible. These can be openly used by other researchers to download the same data from the original data sources.*

#### **Is documentation about the software needed to access the data included? Is it possible to include the relevant software (e.g. in open source code)?**

Where this is applicable, the relevant software and its documentation will be included.

#### **Where will the data and associated metadata, documentation and code be deposited? Have you explored appropriate arrangements with the identified repository?**



Data to be openly shared will be deposited in open repositories like Zenodo as well as in AIO platform which can be accessed by registered users. These are widely used repositories adopting standard and simple procedures to allow data sharing by researchers. No need for appropriate arrangements is foreseen.

**If there are restrictions on use, how will access be provided?**

If such cases are identified, access could be provided either through use of consent and anonymisation or by regulating and restricting access to specific users.

**Is there a need for a data access committee?**

Not at this point.

**Are there well-described conditions for access (i.e. a machine-readable license)?**

Such licenses will be used for the data we plan to make openly available.

**How will the identity of the person accessing the data be ascertained?**

This will be dealt on a case-by-case basis. For the datasets we plan to share, open access will be granted. For the data that will be used only internally by project partners, access control procedures will be in place that define access rights and provide secure access with username/password credentials.

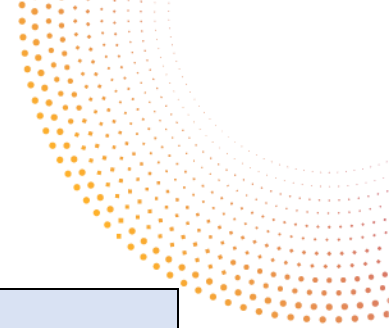
### 3.4 Making data interoperable

This point describes data interoperability specifying what data and metadata vocabularies, standards or methodologies are followed in order to facilitate interoperability. Moreover, it addresses whether a standard vocabulary is used for all data types present in the dataset in order to allow inter-disciplinary interoperability. Specifically, it addresses the following issues:

- Are the data produced in the project interoperable?
- What data and metadata vocabularies, standards or methodologies will you follow to make your data interoperable?
- Will you be using standard vocabularies for all data types present in your dataset, to allow inter-disciplinary interoperability?
- In case it is unavoidable that you use uncommon or generate project specific ontologies or vocabularies, will you provide mappings to more commonly used ontologies?

AI4Media uses data coming from diverse sources. To be able to easily integrate, analyse, and share these diverse types of data, mechanisms for data harmonization and integration will be adopted wherever possible aiming to ensure data interoperability. With regard to the individual questions about data interoperability, our generic DMP approach is summarized below (again, detailed answers for each dataset are given in sections 4 and 5):





**Are the data produced in the project interoperable?**

Effort will be made to make most of the data produced by the project interoperable. This will be pursued in the context of Task T7.1 *Publication of AI resources to the AI-on-demand platform*, which focuses on uploading AI4Media resources (including datasets) on the AIoD platform (see deliverables D7.1, D7.2 and D7.4).

**What data and metadata vocabularies, standards or methodologies will you follow to make your data interoperable?**

In order to ensure interoperability and maximum re-use of AI4Media data, project partners try to collect existing and new data in standardized formats, following well-known data representation models and metadata vocabularies.

Standard data vocabularies are adopted for different types of datasets (social media data, audiovisual data, user analytics, etc.) while a common approach is used for metadata creation based on OpenAIRE guidelines<sup>4</sup>. More information can be found in sections 4 and 5.

**Will you be using standard vocabularies for all data types present in your dataset, to allow inter-disciplinary interoperability?**

To facilitate the exchange of information and sharing of data, we try to rely on accepted and widely used standards whenever this is possible.

**In case it is unavoidable that you use uncommon or generate project specific ontologies or vocabularies, will you provide mappings to more commonly used ontologies?**

This is examined on a case-by-case basis.

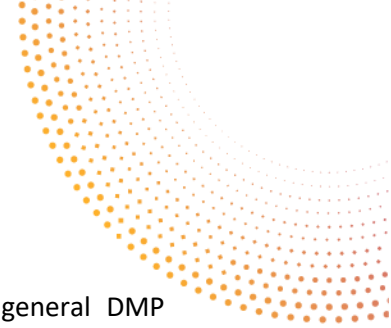
### 3.5 Increase data re-use (through clarifying licenses)

This point addresses the following issues:

- How will the data be licensed to permit the widest re-use possible?
- When will the data be made available for re-use? If an embargo is sought to give time to publish or seek patents, specify why and how long this will apply, bearing in mind that research data should be made available as soon as possible.
- Are the data produced and/ or used in the project useable by third parties, in particular after the end of the project? If the re-use of some data is restricted, explain why.
- How long is it intended that the data remains re-usable?
- Are data quality assurance processes described?

<sup>4</sup> OpenAIRE Guidelines for Data Archives: <https://guidelines.openaire.eu/en/latest/>





With regard to the individual questions about increasing data re-use, our general DMP approach is summarized below (again, detailed answers per dataset are given in sections 4, 5 and 6):

**How will the data be licensed to permit the widest re-use possible?**

This will be examined on a case-by-case basis depending on the dataset. Our general approach can be summarized as follows:

- In case of data coming from external open sources or in cases where the data comes with a license on its own, the data is shared under the same licence (if we decide to reshare it).
- For other cases, a CC-BY 4.0 (Creative Commons Attribution 4.0 International License) license will be selected wherever possible. This license allows open sharing but also allows keeping some control over the data (e.g. requires attribution).

**When will the data be made available for re-use? If an embargo is sought to give time to publish or seek patents, specify why and how long this will apply, bearing in mind that research data should be made available as soon as possible.**

This is examined on a case-by-case basis (see sections 4, 5 and 6). In general, effort is made for the data to be made available as soon as possible.

**Are the data produced and/ or used in the project useable by third parties, in particular after the end of the project? If the re-use of some data is restricted, explain why.**

This is examined on a case-by-case basis (see sections 4, 5 and 6). The datasets that we openly share will be re-usable after the end of the project through the AoD platform and open repositories like Zenodo and GitHub.

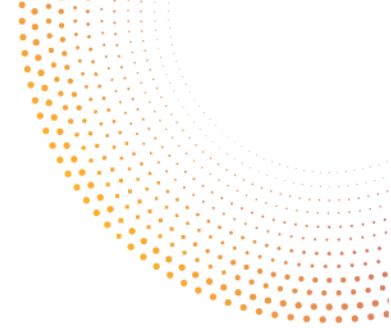
**How long is it intended that the data remains re-usable?**

This is examined on a case-by-case basis (see sections 4, 5 and 6). The datasets that we openly share will be re-usable at least for a few years after the end of the project.

**Are data quality assurance processes described?**

Automatic data cleaning techniques are employed during data collection. Data cleaning consists of identifying incomplete, incorrect, inaccurate, or inconsistent parts of the data and then replacing, modifying, or deleting such data. This is necessary for improving data quality and producing a clean, uniform, and consistent dataset for integration; the quality of the data reflects directly upon the quality and accuracy of the data analysis results. For datasets including questionnaire data, a manual quality control is performed by partners to ensure consistency of replies. These quality assurance processes, including data cleaning are discussed in sections 4, 5 and 6.





### 3.6 Allocation of resources

This point addresses the following issues:

- Estimate the costs for making the data FAIR and describe the method of covering these costs;
- Identify responsibilities for data management in the project;
- Describe costs and potential value of long term preservation.

With regard to the individual questions, our generic DMP approach is summarized below (again detailed answers for each dataset are given in sections 4, 5 and 6):

#### **Estimate the costs for making the data FAIR and describe the method of covering these costs.**

Costs for publications are covered by the project budget. Other costs for making the data FAIR will be covered by the individual partners that will share the data. Open data sharing will be achieved by depositing the data in open repositories like Zenodo, the AIO platform or partners' institutional open repositories where no costs for data sharing are foreseen. Resources to make the data interoperable are already foreseen in the DoA as part of the work performed in Task 7.1 "*Publication of AI resources to the AI-on-demand platform*".

#### **Identify responsibilities for data management in the project.**

A data manager role has been established in the project to ensure that data processing actions within AI4Media are in line with the law. CERTH has been appointed as the beneficiary responsible for data management and has cooperated with technical and pilot partners to draft a detailed data management plan that clearly identifies how each dataset used or created by the project will be handled. CERTH has appointed a Data Protection Officer (DPO) that is responsible for closely monitoring the execution of the data management plan and ensuring that project partners handle project datasets appropriately. The DPO works in the central administration of CERTH and has significant experience in data management and data protection as the DPO of many Horizon and H2020 projects coordinated by CERTH. The work to be done in T4.1 by KUL, the project's legal and ethical expert, on identifying the applicable legal frameworks helps with the process of data management. Moreover, ethics deliverables D12.1 and D12.2 and the Ethics Management Plan (D1.3, D1.5) clarify a lot of issues with regard to human participation, informed consent procedures, and protection of personal data.

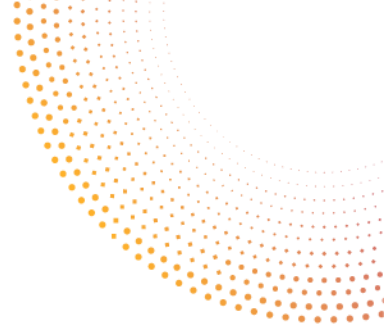
#### **Are the resources for long term preservation discussed (costs and potential value, who decides and how what data will be kept and for how long)?**

In principal, data used in or generated by the project will be preserved in partners' servers for at least 3 years beyond the lifetime of the project. After this period, the data will either be deleted or preserved for a longer period of time in dedicated repositories based on the internal agreements between partners.

Open data will continue to be available in open repositories like Zenodo after the project ends.







### 3.7 Data security

This point addresses data recovery, as well as secure storage and transfer of sensitive data. Specifically, this point addresses the following questions:

- Is the data safely stored in certified repositories for long-term preservation and curation?
- What provisions are in place for data security?

All software tools and data storage mechanisms developed within AI4Media are designed to safeguard collected data against unauthorized use and to comply with all national and EU regulations. Engineering best practices and state-of-the-art data security measures are incorporated as well as GDPR considerations, and respective guidelines and principles.

As explained above, AI4Media datasets are either be openly shared (by uploading them in open repositories) or shared internally among specific partners (e.g. stored in partners' servers). In addition, some datasets will be stored in third-party cloud servers. Below, we examine the data security strategy for the aforementioned data storage options.

#### **Open repositories**

Datasets to be openly shared, will be deposited in certified repositories like Zenodo that have in place strong mechanisms and protocols for data security and long-term data preservation. The same stands for the AIoD platform.

#### **AI4Media wiki**

The data is stored on the project wiki<sup>5</sup>, which is hosted on a dedicated web server in CERTH's premises. The wiki web site uses for its domain an SSL certificate enabling the SHA256RSA signature algorithm and forces all visits to use HTTPS to ensure the traffic is secure. The wiki is restricted only to registered users while registration is possible only by invitation. Access to the wiki requires username/password authentication. CERTH pays special attention to security and respects the privacy and confidentiality of the users' personal data by fully complying with the applicable national (Greek), European and International framework, and the European Union's General Data Protection Regulation (GDPR) 2016/679. The AI4Media wiki uses a file-based RDBMS to enhance security as no ports for separate DB-instances are open. Web server and file-based DB are running on a Linux encrypted partition, which conforms to the data-at-rest GDPR guidelines. Moreover, at a higher level in CERTH's data center, state-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimize the risk of data loss.

The data will be preserved in the wiki until the end of the project and for three years after that and will then be deleted.

#### **Partners' servers**

---

<sup>5</sup> <https://mklab.itl.gr/AI4Media/doku.php>



AI4Media partners have significant experience in data handling and protection both in the context of their institutional operation as well as in the context of their participation in other EU projects. To that end, the beneficiaries already have in place operational policies regarding potential ethics issues as well as privacy and security guidelines for data protection, adhering to national and EU regulations.

Ultimately, each partner is responsible for the data protection and security mechanisms in their own servers. In the following, we list some regulations and rules with regard to data protection, which can be followed by partners:

- Compliance with the internationally recognized and globally accepted and recently adopted standard, ISO/IEC 27701:2019 Security technique, which is an extension to the ISO/IEC 27001:2013 15 (Information technology - Security techniques - Information security management systems – Requirements).
- Compliance with the requirements of the ISO/IEC 27017:201516 (Information technology - Security techniques - Code of practice for information security controls based on ISO/IEC 27002 for cloud services), which provides controls and implementation guidance for cloud service providers as well as cloud service customers.
- Adherence to the internationally recognised and globally accepted standard, ISO 2701817 (Information technology - Security techniques - Code of practice for protection of personally identifiable information (PII) in public clouds actions as PII processors). The standard is designed for user privacy protection. The certification combines legal requirements for data processing with technical criteria for information security systems. The goal of ISO 27018 is to provide a set of uniform security controls to public cloud computing service providers who act as personal data processors. It implements measures to protect Personally Identifiable Information (PII).
- Adherence to the General Data Protection Regulation (2016/679/EU, GDPR).
- Appropriate strong access control mechanisms (at a server level, virtual private network (VPN) level, or Virtual Machine level) to provide only the necessary level of access to specific users.
- Robust encryption of personal data at-rest and in-transfer, but also of non-personal data wherever necessary or possible.
- (Pseudo-)anonymization of personal data according to the GDPR.
- Schedule of regular (daily or weekly) backups to enable rollback in case of significant hardware storage failure and thus minimize data loss.

The data will be preserved in partner servers until the end of the project and for at least three years after that and will then be deleted.

### **Third-party cloud servers**

In the framework of WP8, ATC will use Truly Media, a web-based platform for collaborative verification of User Generated Content (UGC), as the main demonstrator for Use Case 1. In order to facilitate the operability and functioning of Truly Media, ATC works with third-party service providers for hosting the service. The following third-party contractors are used for data storage in Truly Media:

- Amazon Web Services EMEA SARL Greek Branch, 59-61 Agiou Konstantinou St., Marousi, 15124, Athens, Greece: Storage for all multimedia files uploaded and hosted by the platform.





- Heroku, Inc., 1 Market St., Suite 300, San Francisco, CA 94105, United States: Hosting of backend processes.
- MongoDB Limited, Building Two, Number One Ballsbridge, Ballsbridge, Dublin 4, Ireland: Platform data storage.

It is noted that ATC has bound its data processors with data processing agreements concluded pursuant to Article 28 of the GDPR. In cases where third-party service providers are used, we have ensured that – even when these providers reside outside of the European Economic Area (EEA) – all data is stored in servers located in the EEA and there is no data transfer outside the EEA.

Appropriate and detailed security policies, rules, and technical measures are implemented to protect data that are used by the Truly Media platform and are stored on the platform from improper or unauthorized access, including use of firewalls where appropriate. Security measures also include 2FA (2 Factor Authentication) with OTP (One Time Password) for extra security during login, as well as Auth2.0 and JWT for authentication and authorisation. End-to-end encryption protects from man-in-the-middle attacks and data theft. All ATC employees and data processors, who have access to and are associated with the processing of personal data, are obliged to respect the confidentiality of the stored personal data. Moreover, ATC’s development team has received training from external auditors for security awareness and security best practices to avoid vulnerabilities in source code. External auditors have performed black-box penetration testing to ensure that the platform is fully secure. ATC’s Data Protection Officer ensures that all processes we follow are fully compliant with the GDPR provisions.

More details on the security measures adopted by each partner to protect the various datasets collected or generated by the project are provided in sections 4, 5 and 6.

### 3.8 Ethical & legal aspects

This point covers any ethical or legal issues that can have an impact on data sharing, including references to ethics deliverables and the ethics section (i.e. Section 5) of the DoA. Specifically, it addresses the following issues:

- Are there any ethical or legal issues that can have an impact on data sharing?
- Is informed consent for data sharing and long-term preservation included in questionnaires dealing with personal data?

When a dataset cannot be shared, the reasons for this will be outlined (e.g. ethical restrictions, rules governing privacy and personal data protection, intellectual property, and commercial sensitivity).

With regard to the individual questions, our generic DMP approach is summarized below (again, detailed answers for each dataset are given in sections 4, 5 and 6):

**Are there any ethical or legal issues that can have an impact on data sharing?**

Addressing legal and ethics challenges is an important part of the AI4Media work plan. As already indicated in section 5 “*Ethics and Security*” of the DoA, special attention has been paid to these issues since the very beginning of the project. A legal partner (KUL) is part of the consortium, and provided guidance and relevant expertise throughout the project, including



presentations on data protection and privacy and a background note on processing personal data for scientific research purposes. The background notes explored the ethical and legal consideration of dataset re-use. In addition, the organizational structure of AI4Media includes an external Ethics Advisory Board. This Board advises on ethics issues as well as on the data processing procedures adopted in AI4Media and offers expertise on the matter.

A dedicated task deals specifically with such issues: Task 1.3 *“Ethical issue management”*, led by KUL. T1.3 delivered an initial Ethics Management Plan (D1.3) in M12 and the final Ethics Management Plan (D1.5) in M48. The Ethics Management Plan includes the responses to the ethics requirements, as requested by the EC ethics review in relation to protection of personal data, human participants, misuse of research findings, and participation of non-EU countries. It also includes elements reporting on KUL’s online session focusing on ethics impact assessments in the context of technology development. The session provided an overview of various types of impact assessments, shared practical challenges, and offered tips for conducting them effectively.

Moreover, the following tasks provided analysis of relevant legal frameworks (GDPR, AI Act, Digital Services Act, Data Governance Act, Data Act) and offered clarifications with regard to AI ethics, privacy and data governance. Ethics considerations related to the use of data in Generative AI were also addressed in:

- Task T2.1 *“Analysis of the EU policy on AI and the forthcoming Commission’s legislative proposal on AI regulation”* as part of deliverable D2.1;
- Task 4.1 *“Legal and ethical frameworks for trusted AI”* as part of deliverables D4.3 and D4.4.

In addition to the aforementioned tasks, WP12 *“Ethics requirements”* also deals with specific ethics and legal requirements. These included the design of consent forms and information sheets for the collection of personal data from human participants (D12.1. - *H - Requirement No. 1*), requirements for the collection and protection of personal data (D12.2 - *POPD - Requirement No. 2*), risk assessment to prevent misuse of research findings including generated data (D12.3 - *M - Requirement No. 3*), and procedures to ensure data transfer to/from non-EU countries as required by national/EU legislation (D12.4 - *NEC - Requirement No. 4*).

Handling of ethics, legal, and privacy issues is one of the building blocks of AI4Media. Special focus has been given to privacy rights and data protection regime under the GDPR. In the framework of AI4Media research activities, the consortium processed personal data collected during the project or personal data and anonymized data collected during other previous professional experiences external to the project. In other words, partners of the consortium had in some case processed data already stored in partners’ datasets or data initially collected by third parties. Article 5(1)b GDPR stipulates that personal data shall be ‘collected for specified, explicit and legitimate purposes and not further processed in a manner that is incompatible with those purposes’. Article 5(1)b GDPR, as complemented by Recital 50 GDPR, allows further processing where the new purposes are compatible with the initial ones. Furthermore, it specifies that further processing should be considered to be compatible when is carried out for ‘archiving purposes in the public interest, scientific or historical research purposes or statistical purposes’ (emphasis added). When personal data are further used for compatible purposes, such as in the case of further processing of personal data for research



purposes, 'no legal basis separate from that which allowed the collection of the personal data is required'.

AI4Media is a research project receiving funding from the EU's Horizon 2020 research and innovation programme. In the light of the information provided by other partners, in the AI4Media project the partners will indeed use data that got previously collected for an initial purpose, for research purposes. Such further processing is assumed compatible with the initial purpose within the meaning of the GDPR under two conditions: (i) the data were initially collected based on a lawful basis and (ii) that appropriate technical and organisational measures are in place to safeguard the rights of the data subjects (Art. 89 GDPR).

Where relevant and in respect with our research obligations, anonymisation and pseudonymisation techniques were used to safeguard the personal data processed. Where relevant, personal data were pseudonymised according to the current state of the art, and the additional information necessary for re-identifying the individual was kept separately (according to Art. 4(5) of the GDPR). Pseudonymisation is expressly mentioned as a measure within Art. 89 GDPR.

Collected personal data also include processing of special categories of data (so-called 'sensitive data'), such as data revealing political opinions. More specifically, in the context of WP6 one of the research objectives is enhancing opinion mining performance in politically charged texts (e.g. tweets), especially in the presence of implicitly expressed opinions, sarcasm and metaphors. According to Article 9 GDPR, the processing of such special categories of personal data is prohibited. However, according to Article 9(2)(j) GDPR this prohibition does not apply when the processing is necessary for scientific or research purposes in accordance with Article 89(1) GDPR. The safeguards of Article 89(1) GDPR include the use of specific 'technical and organizational measures to ensure data minimisation, such as pseudonymisation, or the anonymisation of data where possible'.

The AI4Media Consortium respected and fully complied with the GDPR provisions. Any processing of personal data in AI4Media is covered by the appropriate legal ground (e.g. informed consent or legitimate interest; for more, see all the personal data processing activities reported by partners in the final Ethics Management Plan (D1.5).

Moreover, for personal data processed in the context of the AI4Media use cases, end users were provided with informed consent forms and information sheets when personal data had been collected and processed for research purposes. Considering the ethical aspect along with the legal requirements on consent as set out under the GDPR, informed consent is a key condition for autonomous decision-making, which demonstrates the notion of control over one's personal data. By providing comprehensive information for the envisaged purposes, the aim was to sufficiently inform the person concerned about the use of his or her data in order to provide control over how the data is being managed. As indicated, the data involved has been pre-processed in a pseudonymised and confidential manner. Only anonymised (through aggregation) research results may be scientifically exchanged or disseminated.

Finally, in the context of AI4Media, transfer of data, including personal data, between EU and non-EU countries was also at stake (this issue is addressed in D12.4 – "*NEC - Requirement No. 4*"). The AI4Media consortium includes partners that are based in Switzerland (HES-SO, IDIAP) and the United Kingdom (QMUL, F6S). Thus, the data processing activities that involve these



partners include the import/export of personal data to these countries. As regards Switzerland, the European Commission has recognised this country as providing adequate protection with the Decision 2000/518/EC. According to Article 96 GDPR, international agreements involving the transfer of personal data to third countries which were concluded prior to May 2016 shall remain in force until amended, replaced or revoked. Thus, this decision is still valid. The effect of this decision is that personal data can flow from the EU to Switzerland without any further safeguard being necessary. In others words, transfers to Switzerland will be assimilated to intra-EU transmissions of data, as the European Commission’s website clearly explains.<sup>6</sup> In January 2024, the EC published its report on the first review of the functioning of the eleven adequacy decisions.<sup>7</sup> The report investigated if these adequacy decisions stood the test of time and address new developments and challenges. Based on the report findings, the EC concludes that Switzerland continues to provide an adequate level of protection for personal data transferred from the EU.

As regards the United Kingdom, after the Brexit, the European Commission has also concluded an adequacy decision with the UK.<sup>8</sup> “Personal data can now flow freely from the European Union to the United Kingdom where it benefits from an essentially equivalent level of protection to that guaranteed under EU law.”<sup>9</sup>

#### **Is informed consent for data sharing and long-term preservation included in questionnaires dealing with personal data?**

Consent has been identified as a lawful basis for the processing of data gathered during AI4Media use case trials. Accordingly, the end users that participate in use case activities like evaluation were provided with informed consent forms and information sheets (see D12.1 “*H - Requirement No. 1*” and D12.2 “*POPD - Requirement No. 2*”). Consent is a key tenet of the new data protection legislation (GDPR) and can only be obtained when providing the individual control over the processing of their personal data, as only lawfully obtained consent ensures that the processing is fair and transparent to the data subject.

AI4Media sought informed consent of all research participants. In order to obtain informed consent, the consortium provided prospective participants sufficient opportunity to consider whether or not to participate and this under circumstances that minimise the possibility of coercion or undue influence. Participants will be provided with informed consent forms and information sheets where personal data has been collected and processed for research purposes. Considering the ethical aspect along with the legal requirements on consent as set

<sup>6</sup> [https://ec.europa.eu/info/law/law-topic/data-protection/international-dimension-data-protection/adequacy-decisions\\_en](https://ec.europa.eu/info/law/law-topic/data-protection/international-dimension-data-protection/adequacy-decisions_en)

<sup>7</sup> Report on the first review of the functioning of the adequacy decisions adopted pursuant to Article 25(6) of Directive 95/46/EC, [https://commission.europa.eu/document/f62d70a4-39e3-4372-9d49-e59dc0fda3df\\_en](https://commission.europa.eu/document/f62d70a4-39e3-4372-9d49-e59dc0fda3df_en)

<sup>8</sup> Commission Implementing Decision (EU) 2021/1772 of 28 June 2021 pursuant to Regulation (EU) 2016/679 of the European Parliament and of the Council on the adequate protection of personal data by the United Kingdom (notified under document C(2021)4800) (Text with EEA relevance), <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32021D1772>

<sup>9</sup> [https://ec.europa.eu/commission/presscorner/detail/%20en/ip\\_21\\_3183](https://ec.europa.eu/commission/presscorner/detail/%20en/ip_21_3183)



out under the GDPR, informed consent is a key condition for autonomous decision-making, which demonstrates the notion of control over one's personal data. By providing comprehensive information for the envisaged purposes, the aim was to sufficiently inform the person concerned about the use of his or her data in order to provide control over how the data is being managed. As indicated, the data involved has been pre-processed in a pseudonymised and confidential manner. Only anonymised (through aggregation) research results may be scientifically exchanged or disseminated.

The collected data are used solely for the research purposes of AI4Media, are not transferred to any third Parties (as specified earlier) and will be deleted three years after the end of the project.

The consortium guarantees that all personal data collected during the project will be kept secure and unreachable by unauthorized persons. The data is handled with appropriate confidentiality and technical security, as required by law in the individual countries and EU laws and recommendations, mainly the General Data Protection Regulation (GDPR) (Regulation (EU) 2016/67<sup>10</sup> of the EU. All AI4Media partners have in place their own data privacy and security policies, which are compliant with EU regulations.

Before obtaining written consent, information concerning the data processing operations was handed to trial participants. The specific information requirements are laid down in Art. 13 and 14 GDPR. Accordingly, in order to provide information to the data subjects in a clear manner and to give the individual participants a genuine choice with regard to the envisaged data processing, the information sheets gave research participants information about, inter alia:

- Purposes of data collection, data processing and data analysis;
- Types of personal data processed;
- Transfer of their personal data between the use case leader and the relevant technical partner(s), involved in the trials;
- The rights they have as data subjects, and information on how to exercise them;
- The period for which the data will be stored.

The participation at the research was entirely voluntary and the participants had the right to withdraw from the research at any time without any adverse consequences.

### 3.9 Other issues

Other issues refer to other national/ funder/ sectoral/ departmental procedures for data management used in the project.

As mentioned above, all consortium partners (media organizations as well the research organizations and SMEs/industry that participate in the project) have in place their own data privacy and security policies, which are compliant with EU regulations and especially the GDPR.

---

<sup>10</sup> <https://www.eugdpr.org/eugdpr.org.html>





## 4. Data management plan for research datasets created within AI4Media

This section presents the research datasets that partners of the AI4Media consortium created during the project lifetime. An initial list of 19 datasets was presented in the initial DMP (D1.2) delivered in M6. This list has since been updated to include additional datasets created or collected between M7 and M48. **61 datasets in total are presented in this section.**

In the following sub-sections, we present the DMP plan for the different research datasets that were created within AI4Media, organized per WP. The DMP information for each dataset is provided in the form of a Table with specific fields, following the methodological approach described in section 3. The dataset presentation template is shown below (Table 1). It is the same template also used in D1.2. The field *DMP component* refers to the dataset reference name (i.e. dataset id).

Table 1: Template for the presentation of the data management plan for a specific dataset

DMP component	AI4Media_Data_DatasetNo_WPX_TypeofData_DatasetTitle_Version Partner: Short name of partner processing this data
Data Summary	<p><u>Purpose</u>: Short description of data + What is the purpose of data collection/generation (and its relation to project objectives) in the context of AI4Media?</p> <p><u>Type/format</u>: What is the type/format of the data?</p> <p><u>Re-use of existing data</u>: Are you re-using an existing dataset and how?</p> <p><u>Data origin</u>: What is the origin/source of the data?</p> <p><u>Expected size</u>: What is the expected data/dataset size?</p> <p><u>Data utility</u>: To whom will this data be useful and how? (inside the project and also to third parties, if applicable)</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Are the data produced and/or used in the project discoverable with metadata, identifiable and locatable by means of a standard identification mechanism (e.g. persistent and unique identifiers such as Digital Object Identifiers)?</p> <p><u>Search keywords</u>: Will search keywords be provided that optimize possibilities for re-use?</p> <p><u>Versioning</u>: Do you provide clear version numbers?</p> <p><u>Metadata creation</u>: Specify standards for metadata creation (if any). If there are no standards in your discipline, describe what type of metadata will be created and how.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Will data produced and/or used in the project be made openly available as the default? If certain datasets cannot be shared (or need to be shared under restrictions), explain why, clearly separating legal and contractual reasons from voluntary restrictions.</p> <p><u>How it will be accessible</u>: How will the data be made accessible (e.g. by deposition in an open repository)?</p> <p><u>Methods/software tools to access data</u>: What methods or software tools are needed to access the data? Also, is documentation about the software needed to access the data included? Is it possible to include the relevant software (e.g. in open source</p>





	<p>code)?</p> <p><u>Repository</u>: Where will the data and associated metadata, documentation and code be deposited? Preference should be given to certified repositories which support open access where possible</p> <p><u>Restrictions on access</u>: If there are restrictions on use, how will access be provided?</p>
Making data interoperable	<p><u>Interoperability</u>: Are the data produced in the project interoperable, that is allowing data exchange and re-use between researchers, institutions, organisations, countries, etc. (i.e. adhering to standards for formats, as much as possible compliant with available (open) software applications, and in particular facilitating re-combinations with different datasets from different origins)?</p> <p><u>Data and metadata vocabularies</u>: Specify what data and metadata vocabularies, standards or methodologies you will follow to facilitate interoperability</p> <p><u>Use of standard vocabularies</u>: Specify whether you will be using standard vocabulary for all data types present in your data set, to allow inter-disciplinary interoperability?</p> <p><u>Mappings to commonly used vocabularies</u>: In case it is unavoidable that you use uncommon or generate project specific ontologies or vocabularies, will you provide mappings to more commonly used ontologies?</p>
Increase data re-use	<p><u>Licence</u>: Specify how the data will be licensed to permit the widest reuse possible. E.g. Open Data License (Creative Commons CC Zero License, Creative Common Attribution License-CC-BY v4.0, etc.).</p> <p><u>Availability for re-use</u>: When will data be made available for re-use. If applicable, specify why and for what period a data embargo is needed</p> <p><u>Usable by third parties after end of project</u>: Specify whether the data produced and/or used in the project is useable by third parties, in particular after the end of the project? If the re-use of some data is restricted, explain why.</p> <p><u>Re-use timeframe</u>: Specify the length of time for which the data will remain re-usable</p> <p><u>Data quality assurance process</u>: Describe data quality assurance processes</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: Estimate the costs for making your data FAIR. Describe how you intend to cover these costs</p> <p><u>Costs for long-term preservation</u>: Describe costs and potential value of long term preservation</p>
Data security	<p><u>Security measures</u>: Security measures implemented for data protection (incl. controlled access, user authentication, firewalls, VPNs, encryption, back-ups, etc.)</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Are there any ethical or legal issues that can have an impact on data sharing?</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Is informed consent for data sharing and long term preservation included in questionnaires dealing with personal data? Discuss</p>
Other Issues	<p>Refer to other national/funder/sectorial/departmental procedures for data management that you may be using (if any)</p>

In the initial DMP, 19 research datasets were presented. Since then the number of datasets created within AI4Media has risen to 61 datasets, 42 more than in D1.2. **34 out of these 61 datasets are openly shared.** The remaining 27 datasets cannot be openly shared due to



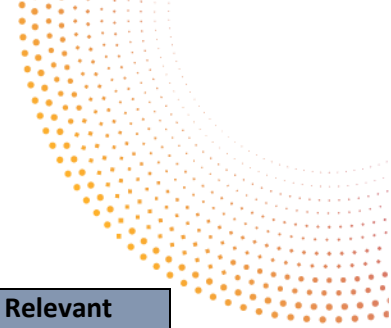
various ethical and legal aspects that prevent/prohibit sharing, examined on a dataset-by-dataset basis (e.g. some social media data cannot be shared due to the platform's terms of use, use case evaluation data cannot be shared because they usually include personal data of end users, datasets of news articles cannot be shared because the copyright is owned by the relevant media organisations, etc.)

Below, we provide a Table that briefly summarizes the 61 datasets created within AI4Media and offers a glance at the structure of this section and its subsections. New datasets that were not included in the initial DMP are indicated with yellow.

*Table 2: Summary of research datasets created within AI4Media*

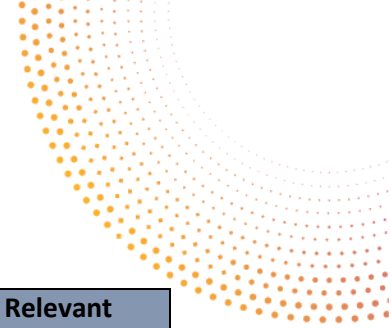
DMP component	WP	Short summary	Relevant sub-section
<b>Data collected in WP2 (European AI Vision, Policy and Common Research Agendas)</b>			4.1
AI4Media_Data_01_WP2_QUESTIONNAIRE_AI-tech-roadmap-2021_v1	WP2	Questionnaires for AI technology roadmap	4.1.1
AI4Media_Data_02_WP2_SURVEY_AI-impact-policy_v1	WP2	Workshop data for AI technology impact and policy recommendations	4.1.2
<b>Data collected in WP3 (New Learning Paradigms &amp; Distributed AI)</b>			4.2
AI4Media_Data_03_WP3_IMAGE_FaV CI2D_v1	WP3	FaVCI2D image dataset for demographically diversified face verification	4.2.1
AI4Media_Data_04_WP3_Video_Mixamo_v1	WP3	Mixamo-Kinetics dataset	4.2.2
AI4Media_Data_05_WP3_Image_100_Driver_v1	WP3	100-Driver dataset for distracted driver classification	4.2.3
AI4Media_Data_06_WP3_TEXT_LeQua_v1	WP3	LeQua 2022 datasets	4.2.4
AI4Media_Data_07_WP3_TEXT_ProductReviewsForOrdinalQuantification_v1	WP3	Product reviews for ordinal quantification dataset	4.2.5
AI4Media_Data_08_WP3_TEXT_ProductReviewsDataset_v1	WP3	Product reviews dataset	4.2.6
AI4Media_Data_09_WP3_UCIOpenMLdatasetsOrdinalQuantification_v2	WP3	UCI and OpenML datasets for ordinal quantification	4.2.7
AI4Media_Data_10_WP3_CherenkovTelescopeDataOrdinalQuantification_v2	WP3	Cherenkov telescope data for ordinal quantification	4.2.8
<b>Data collected in WP4 (Explainability, Robustness and Privacy in AI)</b>			4.3
AI4Media_Data_11_WP4_MRI_MultipleSclerosisLesionSegmentation_v1	WP4	White matter multiple sclerosis lesion segmentation datasets	4.3.1
<b>Data collected in WP5 (Content-centered AI)</b>			4.4
AI4Media_Data_12_WP5_Video_ObyGaze12_v1	WP5	ObyGaze12 dataset for detection of visual objectification in films	4.4.1
AI4Media_Data_13_WP5_TEXT_MAD-TSC_v1	WP5	MAD-TSC Multilingual Aligned Dataset for Target-dependent Sentiment Classification	4.4.2
AI4Media_Data_14_WP5_Video_OWFObjOD_v1	WP5	ÖWF Object Detection dataset	4.4.3
AI4Media_Data_15_WP5_Video_People@Places_v1	WP5	4PeopleAtPlaces dataset	4.4.4





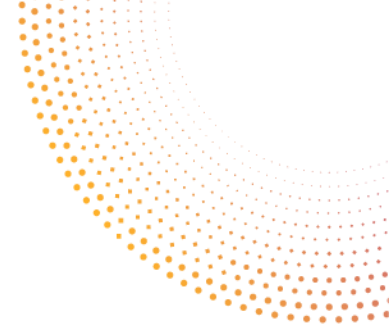
DMP component	WP	Short summary	Relevant sub-section
AI4Media_Data_16_WP5_Video_ToDY_v1	WP5	ToDY dataset for visual time of day and season classification	4.4.5
AI4Media_Data_17_WP5_ImageFeatures_VisioneFeatureRepository_v1	WP5	VISIONE Feature Repository	4.4.6
AI4Media_Data_18_WP5_TEXT_COCO-LVIS-Open ImagesV4-classes Mapping_v1	WP5	COCO, LVIS, Open Images V4 classes mapping	4.4.7
AI4Media_Data_19_WP5_Video_BusViolence_v1	WP5	Bus violence dataset	4.4.8
AI4Media_Data_20_WP5_Image_PestStickyTraps_v1	WP5	Pest Sticky Traps dataset	4.4.9
AI4Media_Data_21_WP5_IMAGE_VWFP_v1	WP5	Virtual World Fallen People dataset	4.4.10
AI4Media_Data_22_WP5_VIDEO_SR-BVI-DVC_v1	WP5	SR_BVI-DVC super-resolution dataset	4.4.11
AI4Media_Data_23_WP5_VIDEO_BSC4K_v1	WP5	BSC4K super-resolution dataset	4.4.12
AI4Media_Data_24_WP5_IMAGE_NDISPark_v1	WP5	Night and Day Instance Segmented Park dataset	4.4.13
AI4Media_Data_25_WP5_TEXT_ArxivAbstracts_v1	WP5	ArXiv abstracts for authorship analysis dataset	4.4.14
AI4Media_Data_26_WP5_4D_Florence4D_v1	WP5	Florence4D facial expression dataset	4.4.15
AI4Media_Data_27_WP5_EventRGB-NEFER_v1	WP5	Neuromorphic Event-based Facial Expression Recognition dataset	4.4.16
AI4Media_Data_28_WP5_Video_PEM360_v1	WP5	PEM360 dataset of 360° videos	4.4.17
AI4Media_Data_29_WP5_VIDEO-TEXT_VRT-Sum_v1	WP5	VRT-Sum video summarization dataset	4.4.18
AI4Media_Data_30_WP5_MODEL_CA-SUM-models_v1	WP5	CA-SUM pretrained video summarization models	4.4.19
AI4Media_Data_31_WP5_Audio_AudioPhylogeny_v1	WP5	Audio phylogeny dataset	4.4.20
AI4Media_Data_32_WP5_Video_CMM-Doc-dataset_v1	WP5	RAI CMM documentaries dataset	4.4.21
AI4Media_Data_33_WP5_Video_CMM-ANTS-dataset_v1	WP5	RAI CMM-ANTS newscasts dataset	4.4.22
AI4Media_Data_34_WP5_Video_CMM-dataset_v1	WP5	RAI CMM mixed dataset DMP component	4.4.23
AI4Media_Data_35_WP5_Text_Cross-lingual-news_v1	WP5	Cross lingual news dataset	4.4.24
AI4Media_Data_36_WP5_Misc_YouTubeDataset_v1	WP5	YouTube RAI channel dataset	4.4.25
<b>Data collected in WP6 (Human- and Society-centred AI)</b>			<b>4.5</b>
AI4Media_Data_37_WP6_SOCIALMEDIA_GreekTwitterPolitics_v1	WP6	Greek Politics Twitter dataset	4.5.1
AI4Media_Data_38_WP6_SOCIALMEDIA_Covid19Twitter_v1	WP6	Covid-19 Twitter dataset	4.5.2
AI4Media_Data_39_WP6_TEXT_TextTwitter_v1	WP6	Twitter text dataset	4.5.3





DMP component	WP	Short summary	Relevant sub-section
AI4Media_Data_40_WP6_TEXT_Covid19DiscussionTopicsTwitter_v1	WP6	Twitter COVID-19 discussions topics dataset	4.5.4
AI4Media_Data_41_WP6_AUDIO_ODSS_v1	WP6	ODSS Open Dataset of Synthetic Speech	4.5.5
AI4Media_Data_42_WP6_IMAGE_CelebHQGaze_v1	WP6	CelebHQGaze image dataset for gaze estimation	4.5.6
AI4Media_Data_43_WP6_TEXT_European-news-covid-vaccination_v1	WP6	European news about Covid vaccination dataset	4.5.7
AI4Media_Data_44_WP6_TEXT_SuisseRomandeLocalNews_v1	WP6	Suisse Romande local news dataset	4.5.8
AI4Media_Data_45_WP6_TEXT_LausanneNews_v1	WP6	Lausanne news dataset	4.5.9
AI4Media_Data_46_WP6_TEXT_SuisseAllemandeNews_v1	WP6	Suisse Allemande local news dataset	4.5.10
AI4Media_Data_47_WP6_TEXT_YouTubeReferences_v1	WP6	References on YouTube dataset	4.5.11
AI4Media_Data_48_WP6_TEXT_PoliticalBarometer_v1	WP6	Political Barometer dataset	4.5.12
AI4Media_Data_49_WP6_TEXT_ElecDeb60To16_v1	WP6	ElecDeb60To16-fallacy dataset	4.5.13
<b>Data collected in WP8 (Use cases &amp; demonstrators in media, society and politics)</b>			4.6
AI4Media_Data_50_WP8_USER-RESEARCH_UseCase1_DW_v1	WP8	Data from user research activities in Use Case 1	4.6.1
AI4Media_Data_51_WP8_USER-RESEARCH_UseCase2_VRT_v1	WP8	Data from user research activities in Use Case 2	4.6.2
AI4Media_Data_52_WP8_QUESTIONNAIRE_UseCase3-UserReqCollection2021_v1	WP8	Questionnaires for the collection of user requirements for Use Case 3	4.6.3
AI4Media_Data_53_WP8_QUESTIONNAIRE_UseCase3Evaluation_v1	WP8	Questionnaires for the evaluation of Use Case 3	4.6.4
AI4Media_Data_54_WP8_USER-RESEARCH_UseCase4_NISV_v1	WP8	Data from user research activities in Use Case 4	4.6.5
AI4Media_Data_55_WP8_QUESTIONNAIRE_UseCase5Evaluation_v1	WP8	Questionnaires for the evaluations of Use Case 5-B	4.6.6
AI4Media_Data_56_WP8_QUESTIONNAIRE_VIDEO_UseCase7-UserFeedbackData_v1	WP8	Data from user research activities in Use Case 7	4.6.7
AI4Media_Data_57_WP8_Text_TrulyMediaDataset-UseCase1-ATC_v1	WP8	TrulyMedia social media + web dataset	4.6.8
AI4Media_Data_58_WP8_Video_GameGlitches-UC5-MODL_v1	WP8	Game glitches dataset for Use Case 5	4.6.9
AI4Media_Data_59_WP8_Audio_RAW compositions_v1	WP8	Musical production for AI co-creation dataset	4.6.10
AI4Media_Data_60_WP8_TEXT_CurrentAffairsTranscripts-UC4-NISV_v1	WP8	Current affairs transcripts dataset for Use Case 4	4.6.11
AI4Media_Data_61_WP8_QUESTIONNAIRE_AI-IndustrialNeeds-T8.4_v1	WP8	User survey data for AI industrial needs (for T8.4)	4.6.12



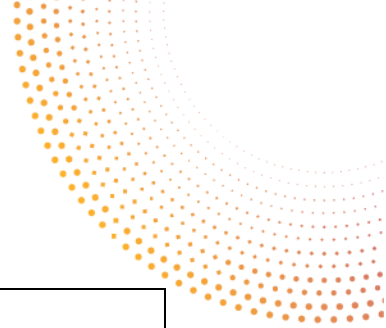


## 4.1 Datasets collected in the context of WP2

### 4.1.1 Questionnaires for AI technology roadmap

DMP component	AI4Media_Data_01_WP2_QUESTIONNAIRE_AI-tech-roadmap-2021_v1 Partner: CERTH
Data Summary	<p><b>Purpose:</b> A structured questionnaire was developed by CERTH to collect information on the evolving landscape of AI technologies for the media sector as part of an anonymous public survey. The questionnaire was addressed to both AI researchers working on multimedia AI but also to people working in the media industry or whose work is closely related to this industry (e.g. researchers studying the media, media regulators, people working in relevant NGOs, etc.). The survey aimed to collect their opinions on the benefits, risks, technological trends and challenges of AI use in the media industry as well as their experience on AI strategies and AI skills in media organisations, their insights on the most promising ways to facilitate AI adoption and knowledge transfer and, finally, their perceptions about ethical use of AI. The online survey was filled anonymously by 150 people. The collected questionnaires were then analyzed and the results of this analysis were included in the roadmap for AI technologies and applications in the media sector (D2.3) in the context of T2.3.</p> <p><b>Type/format:</b> Excel files containing questions and user responses.</p> <p><b>Re-use of existing data:</b> No.</p> <p><b>Data origin:</b> Questionnaires filled online by AI researchers, media professionals, media regulators, etc.</p> <p><b>Expected size:</b> A few MBs in total.</p> <p><b>Data utility:</b> It is useful to WP2 partners for the development of the roadmap for AI in the media sector. It could also be useful to other researchers that work on the same field (AI surveys) but also to policy makers dealing with AI-related issues.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> The survey replies are stored in CERTH’s servers. No DOI is assigned.</p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> We do not intend to make this data openly available since they are used as part of an internal exercise that facilitated the authoring of D2.3. However, a summary of this data is presented in D2.3, which is available on the project website. In case of a report or paper submitted for publication, all research findings will be integrated into the report or paper. Datasets will not be added to the publication.</p> <p><b>How it will be accessible:</b> Stored in CERTH servers, it is only internally accessible by CERTH members.</p> <p><b>Methods/software tools to access data:</b> Web-browser</p> <p><b>Repository:</b> CERTH’s internal servers.</p> <p><b>Restrictions on access:</b> Shared with selected project partners upon request. Access Control List: All partners (R), CERTH (W)</p>



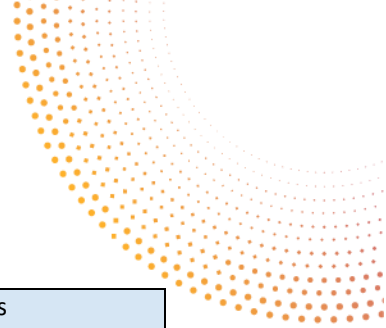


Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will not be shared.</p> <p><u>Availability for re-use</u>: This data is not expected to be re-used. It has been used by WP2 partners to develop the AI roadmap. It will remain on CERTH servers for three years after the end of the project.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data will be stored on a server in CERTH's premises. Access requires username/password authentication. CERTH fully complies with the applicable national, European and International framework, and the GDPR. Moreover, state-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimize the risk of data loss. The data will be preserved there for three years after the end of the project and will then be deleted.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

#### 4.1.2 Workshop data for AI technology impact and policy

DMP component	AI4Media_Data_02_WP2_SURVEY_AI- impact-policy_v1 Partner: UvA
Data Summary	<p><u>Purpose</u>: Six workshops have been conducted by UvA with other partners (KUL &amp; NISV) to collect information on the social and economic impact of AI for media technologies in the context of T2.4. The data collected was analyzed and the results of this analysis were used for drafting the white papers on the social, economic, and political impact of media AI Technologies (D2.5) as well as other shorter report formats for dissemination.</p> <p><u>Type/format</u>: Minutes of workshops (online and in person).</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Input/feedback from workshops from consortium partners, associate members and invited external participants.</p> <p><u>Expected size</u>: A few MBs.</p> <p><u>Data utility</u>: Only internal for WP2 research. The collected data was used for the</p>





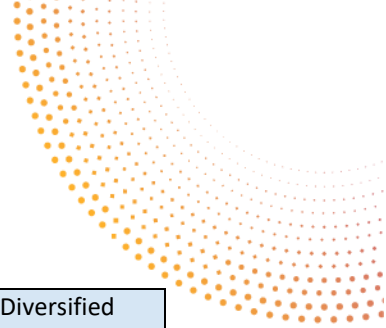
	authoring of D2.5, D2.4 and D2.6 as well as shorter reports and factsheets disseminated on the AI Media Observatory.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> We do not intend to make this data openly available. A summary of this data has been presented in D2.5 (and the methodology used) and relevant short reports, which are available on the project website. No personal information will be published.</p> <p><u>How it will be accessible:</u> Stored in UvA servers.</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> Shared among UvA, KUL and NISV</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<u>Security measures:</u> The data is stored on UvA's servers that use state-of-the-art IT security measures and company policies to mitigate the risk of illegitimate access while also complying with the GDPR.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> No.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	No

## 4.2 Datasets collected in the context of WP3

### 4.2.1 FaVCI2D image dataset for demographically diversified face verification

<b>DMP component</b>	AI4Media_Data_03_WP3_IMAGE_FaVCI2D_v1 Partner: CEA, UPB
----------------------	--

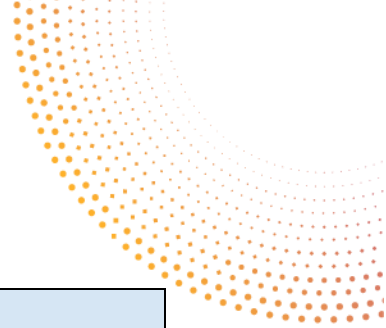




<p>Data Summary</p>	<p><b>Purpose:</b> This dataset (“Face Verification with Challenging Imposters and Diversified Demographics dataset”) includes face pairs for personalities included in Wikipedia. Focus is put on having a demographically diversified sample of persons (age, origin, gender, professions) and on the inclusion of difficult pairs of images. Its creation is necessary because most existing datasets are biased on at least one important demographic dimension. As a consequence, their use does not allow a fair evaluation of face verification algorithms. FAVCI2D will be used in WP3 (T3.3) to evaluate transferability of face analysis models and in WP4 (T4.6) as a contribution to improved benchmarking of AI systems.</p> <p><b>Type/format:</b> JPEG images and accompanying metadata in json format.</p> <p><b>Re-use of existing data:</b> No, the dataset is created within AI4Media.</p> <p><b>Data origin:</b> <a href="https://commons.wikimedia.org">https://commons.wikimedia.org</a> and Bing Image Search</p> <p><b>Expected size:</b> ~1 GB</p> <p><b>Data utility:</b> It is useful to WP3 partners to benchmark face verification algorithms in a setting which is close to realistic conditions.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><b>Is data discoverable:</b> Data is available at <a href="https://github.com/AIMultimediaLab/FaVCI2D-Face-Verification-with-Challenging-Imposters-and-Diversified-Demographics">https://github.com/AIMultimediaLab/FaVCI2D-Face-Verification-with-Challenging-Imposters-and-Diversified-Demographics</a></p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> GitHub supports versioning</p> <p><b>Metadata creation:</b> N/A</p>
<p>Making data openly accessible</p>	<p><b>Data openly accessible:</b> The data is openly accessible via GitHub at <a href="https://github.com/cea-list-lasti">https://github.com/cea-list-lasti</a></p> <p><b>How it will be accessible:</b> The data can be downloaded from an online archive after completing a form.</p> <p><b>Methods/software tools to access data:</b> N/A</p> <p><b>Repository:</b> GitHub</p> <p><b>Restrictions on access:</b> The user should accept the terms of use.</p>
<p>Making data interoperable</p>	<p><b>Interoperability:</b> The file structure makes the use of the dataset easy.</p> <p><b>Data and metadata vocabularies:</b> The dataset is structured around identities, with a JSON file including necessary attributes such as: wiki ID, name, age, origin, gender, professions and two pairs of images associated to it. The first pair includes two images of the person, while the second includes an image of the target person and an imposter image. The imposter image belongs to another identity which is visually similar to the target person.</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
<p>Increase data re-use</p>	<p><b>Licence:</b> The data is released under the <a href="#">FaVCI2D Terms of Use</a>, and the code is released under the CC license.</p> <p><b>Availability for re-use:</b> N/A</p>





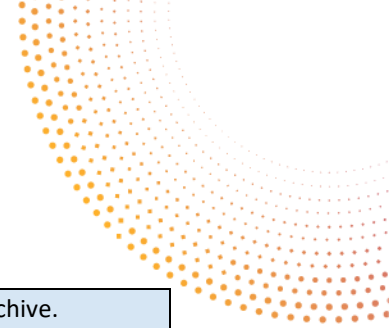


	<p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The full dataset (including images and non-anonymized metadata) will be hosted on UPB's servers. UPB fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The version of the dataset which is shared publicly includes data minimization, in compliance with art. 9 of GDPR.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.2.2 Mixamo-Kinetics dataset

DMP component	AI4Media_Data_04_WP3_Video_Mixamo_v1 Partner: UNITN
Data Summary	<p><u>Purpose</u>: The first large scale dataset for benchmarking domain adaptation methods for action recognition in the challenging task of transferring knowledge from the synthetic to the real domain. Mixamo-Kinetics is used in WP3 (T3.3) for video domain adaptation.</p> <p><u>Type/format</u>: mp4</p> <p><u>Re-use of existing data</u>: No, the dataset is created within AI4Media.</p> <p><u>Data origin</u>: The dataset comprises 36,195 videos, divided into 14 action categories and two domains, i.e., the source domain (synthetic videos generated from Mixamo) and the target domain (real videos from Kinetics (<a href="https://github.com/cvdfoundation/kinetics-dataset">https://github.com/cvdfoundation/kinetics-dataset</a>)).</p> <p><u>Expected size</u>: ~80 GB</p> <p><u>Data utility</u>: It is useful to AI4Media partners and other researchers that want to perform video domain adaptation.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available at <a href="https://drive.google.com/drive/folders/1vGeWBXO5hcFIEll-PLw6zbg-5QsauQG?usp=sharing">https://drive.google.com/drive/folders/1vGeWBXO5hcFIEll-PLw6zbg-5QsauQG?usp=sharing</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: GitHub supports versioning</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is openly accessible via GitHub at <a href="https://github.com/vturrisi/CO2A">https://github.com/vturrisi/CO2A</a></p>





	<p><u>How it will be accessible</u>: The data can be downloaded from an online archive.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: GitHub</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: The dataset comprises 36,195 videos, divided into 14 action categories and two domains, i.e., the source domain (synthetic videos from the Mixamo dataset) and the target domain (real videos from the Kinetics dataset). The source dataset (Mixamo) consists of 24,533 videos synthetically generated using the 3D characters from Mixamo. The dataset comprises videos depicting actions performed by 6 distinct avatars, with different backgrounds, camera positions and random 3D objects in the scene. Also, key-points are provided for each character following the scheme from the COCO dataset but without the keypoints for eyes and ears. Each frame is generated with a resolution of 512 by 512 and the mean length of the videos is 138 frames. The target dataset (Kinetics) was created considering 11,662 videos from 14 action categories extracted from the Kinetics dataset. The overlapping actions between the source and the target datasets are swing dancing, breakdancing, salsa dancing, throwing, capoeira, jogging, shouting, side kick, clapping, texting, golf putting, squat, punching and backflip.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data and the code are released under the CC license.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The full dataset is shared in a google drive repository.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The version of the dataset which is shared publicly includes data minimization, in compliance with art. 9 of GDPR.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.2.3 100-Driver dataset for distracted driver classification

<b>DMP component</b>	<b>AI4Media_Data_05_WP3_Image_100Driver_v1</b> <b>Partner: UNITN</b>
Data Summary	<u>Purpose</u> : 100-Driver is a large-scale, diverse posture-based distracted driver dataset. 100-Driver involves different types of variations that closely meet real-world applications, including changes in the vehicle, person, camera view, lighting, and



	<p>modality. It includes 4 settings for investigating practical problems of DDC, including the traditional setting without domain shift and 3 challenging settings (i.e., cross-modality, cross-view, and cross-vehicle) with domain shifts.</p> <p><u>Type/format</u>: mp4.</p> <p><u>Re-use of existing data</u>: No, the dataset is created within AI4Media.</p> <p><u>Data origin</u>: The dataset comprises more than 470K images taken by 4 cameras observing 100 drivers over 79 hours from 5 vehicles.</p> <p><u>Expected size</u>: ~60 GB</p> <p><u>Data utility</u>: It is useful to AI4Media partners to perform behavior understanding in different driving domains.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: Data is available at <a href="https://100-driver.github.io/">https://100-driver.github.io/</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The data is openly accessible at <a href="https://100-driver.github.io/">https://100-driver.github.io/</a></p> <p><u>How it will be accessible</u>: The data can be downloaded from an online archive after completing a licensing form.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: <a href="https://100-driver.github.io/">https://100-driver.github.io/</a></p> <p><u>Restrictions on access</u>: The user should accept the terms of use.</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: During data collection, we elaborately controled the diversity of raw data in terms of vehicles (5 vehicles, Mazda 3 axela, Lynk&amp;co 03, Toyota C-HR, Hyundai X25 and Ankaï A6), camera locations (4 Xiaomi-C1 cameras in front-left, front, front-right, and side-right, modalities (RGB and NIR), lighting conditions (from morning to afternoon, from spring to winter, and different weather conditions), drivers (100 participants), appearance variations (changing clothes, wearing mask, hat and sunglasses). The RGB modality has been captured in daytime while the Near Infrared (NIR) modality was collected in nighttime. To ensure the appearance variations, a part of the drivers (25% in daytime and 15% in nighttime) were recorded over multiple time periods, leading to huge appearance variations especially in clothes and lighting. Following the collection setting, we initially obtained 79.34 hours of video. The overall annotation process was conducted by 20 experts. To boost the efficiency of data annotation, we first grouped the data by drivers. In addition, we aligned the start and end times for the 4 cameras of the same driver, so that we could label each predefined class for all 4 cameras at once based on the timestamp. Each individual behaviour is labelled with behaviour class, modality type, driver ID, camera ID, vehicle ID and scene ID. Given the labelled video clips, we conducted a down sampling to generate more diverse data considering the high similarity between adjacent frames. We further removed outliers with very different content from the labelled class. We finally produce a total of 470,208 samples to form our 100-Driver dataset.</p>

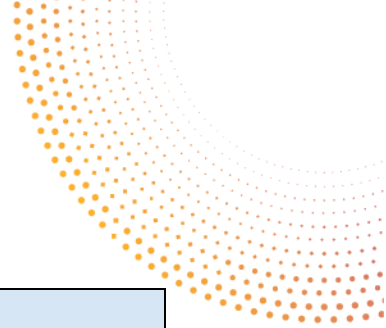


	<p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The data is protected by copyright. The user acquires no ownership, rights, or title of any kind in all or any parts of the dataset. Terms of use clarified at <a href="https://100-driver.github.io/img/100_Driver_Licensing.pdf">https://100-driver.github.io/img/100_Driver_Licensing.pdf</a></p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> Data already publicly shared for research purposes.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The full dataset is available on request.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> During collection, we equipped a safety officer giving relevant instructions according to the road condition. Each participant was informed of the risks involved in data collection and has signed a GDPR informed consent to allow the data to be publicly available for research study.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.2.4 LeQua 2022 datasets

DMP component	AI4Media_Data_06_WP3_TEXT_LeQua_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> The aim of LeQua 2022 (the 1st edition of the CLEF “Learning to Quantify” lab) is to allow the comparative evaluation of methods for “learning to quantify” in textual datasets, i.e., methods for training predictors of the relative frequencies of the classes of interest in sets of unlabelled textual documents. These predictors (called “quantifiers”) will be required to issue predictions for several such sets, some of them characterised by class frequencies radically different from the ones of the training set.</p> <p>LeQua 2022 offers two tasks (T1 and T2), each admitting two subtasks (A and B):</p> <p>T1A: This task is concerned with evaluating binary quantifiers;</p> <p>T1B: This task is concerned with evaluating single-label multi-class quantifiers;</p> <p>T2A: Like Task T1A, this task is concerned with evaluating binary quantifiers. Unlike in Task T1A, participants will be provided with the raw text of the documents;</p> <p>T2B: Like Task T1B, this task is concerned with evaluating single-label multi-class quantifiers; like in Task T2A, participants will be provided with the raw text of the documents.</p> <p><u>Type/format:</u> Raw text and already processed vector formats</p> <p><u>Re-use of existing data:</u> Yes, we use an existing dataset</p> <p><u>Data origin:</u> <a href="https://jmcauley.ucsd.edu/data/amazon/">https://jmcauley.ucsd.edu/data/amazon/</a></p>





	<p><u>Expected size</u>: 7.1GB</p> <p><u>Data utility</u>: Benchmark for learning to count tasks used in experiments in T3.7.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available online (<a href="https://zenodo.org/records/6546188">https://zenodo.org/records/6546188</a>) and indexed in Google</p> <p><u>Search keywords</u>: quantification, prevalence estimation</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: Available on Zenodo <a href="https://zenodo.org/records/6546188">https://zenodo.org/records/6546188</a></p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: Text data is in simple plain text format. Vector data is in csv format. No metadata.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use</u>: Data is already available online.</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

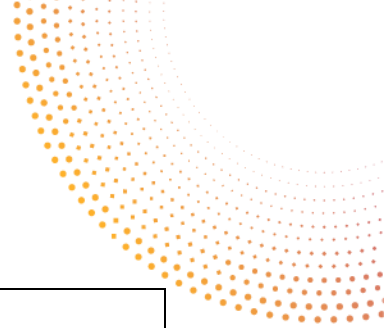
#### 4.2.5 Product reviews for ordinal quantification dataset

<b>DMP component</b>	<b>AI4Media_Data_07_WP3_TEXT_ProductReviewsForOrdinalQuantification_v1</b> <b>Partner: CNR</b>
Data Summary	<u>Purpose</u> : This data set comprises a labeled training set, validation samples, and testing samples for ordinal quantification. The goal of quantification is not to predict the class



	<p>label of each individual instance, but the distribution of labels in unlabeled sets of data. The data is extracted from the McAuley data set of product reviews in Amazon, where the goal is to predict the 5-star rating of each textual review. We have sampled this data according to three protocols that are designed for the evaluation of quantification methods. Used for experiments in T3.7.</p> <p><u>Type/format</u>: Text processed vector format by means of a Large Language Model (Roberta)</p> <p><u>Re-use of existing data</u>: Yes, we use an existing dataset</p> <p><u>Data origin</u>: <a href="https://jmcauley.ucsd.edu/data/amazon/">https://jmcauley.ucsd.edu/data/amazon/</a></p> <p><u>Expected size</u>: 41.9GB</p> <p><u>Data utility</u>: Benchmark for Ordinal Quantification tasks used in T3.7.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available online (<a href="https://zenodo.org/records/8405476">https://zenodo.org/records/8405476</a>) and indexed in Google.</p> <p><u>Search keywords</u>: ordinal quantification, prevalence estimation</p> <p><u>Versioning</u> yes</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: Available on Zenodo at <a href="https://zenodo.org/records/8405476">https://zenodo.org/records/8405476</a></p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: Text data is represented in vector format, with one vector per line, comma-separated values. No metadata.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use</u>: Data is already available online.</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>



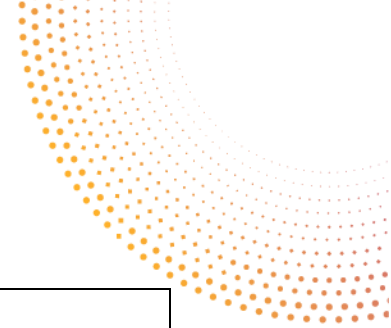


Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.2.6 Product reviews dataset

DMP component	AI4Media_Data_08_WP3_TEXT_ProductReviewsDataset_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> This data set comprises a labelled training set used in the experimentation of the paper "Binary Quantification and Dataset Shift: An Experimental Investigation". The data is extracted from the McAuley data set of product reviews on Amazon.</p> <p><u>Type/format:</u> Raw text partitioned in samples. A file per sample</p> <p><u>Re-use of existing data:</u> Yes, we use an existing dataset</p> <p><u>Data origin:</u> <a href="https://jmcauley.ucsd.edu/data/amazon/">https://jmcauley.ucsd.edu/data/amazon/</a></p> <p><u>Expected size:</u> 5.2GB</p> <p><u>Data utility:</u> Benchmark for quantification under dataset shift, used for experiments in T3.7</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is available online (<a href="https://zenodo.org/records/8421611">https://zenodo.org/records/8421611</a>) and indexed in Google.</p> <p><u>Search keywords:</u> quantification, dataset shift, prevalence estimation</p> <p><u>Versioning</u> yes</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Yes</p> <p><u>How it will be accessible:</u> Available on Zenodo at <a href="https://zenodo.org/records/8421611">https://zenodo.org/records/8421611</a></p> <p><u>Methods/software tools to access data:</u> Web browser</p> <p><u>Repository:</u> Zenodo</p> <p><u>Restrictions on access:</u> No</p>
Making data interoperable	<p><u>Interoperability:</u> The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies:</u> Text data is represented in vector format, with one vector per line, comma-separated values. No metadata.</p> <p><u>Use of standard vocabularies:</u> No</p> <p><u>Mappings to commonly used vocabularies:</u> No</p>
Increase data re-use	<p><u>Licence:</u> Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use:</u> Data is already available online.</p> <p><u>Usable by third parties after end of project:</u> Yes</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>





Allocation of resources	<u>Costs for making data FAIR:</u> N/A <u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

#### 4.2.7 UCI and OpenML datasets for ordinal quantification

DMP component	AI4Media_Data_09_WP3_UCIOpenMLdatasetsOrdinalQuantification_v2 Partner: CNR
Data Summary	<p><u>Purpose:</u> This dataset comprises four different labelled datasets, all targeted at the evaluation of ordinal quantification algorithms (i.e., algorithms that estimate the distribution of the class labels in a set of unseen data items, where there is a total order defined on the set of classes). Used for experiments in T3.7.</p> <p><u>Type/format:</u> The datasets are represented as csv files. Each of the four datasets consists of a set of samples to be used in the experimentation. Each sample is represented as a set of indices into the original UCI or OpenML dataset, so that the sample can be extracted unambiguously from the original dataset.</p> <p><u>Re-use of existing data:</u> Yes, we use existing UCI and OpenML datasets; what we do is providing scripts to extract, from each of them, sets of samples that are suitable for testing ordinal quantification algorithms.</p> <p><u>Data origin:</u> <a href="https://archive.ics.uci.edu/ml/index.php">https://archive.ics.uci.edu/ml/index.php</a> and <a href="https://www.openml.org">https://www.openml.org</a></p> <p><u>Expected size:</u> 25.6MB</p> <p><u>Data utility:</u> These datasets will be useful to anyone wishing to test the accuracy of ordinal quantification algorithms.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data is available online (<a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>) and indexed in Google.</p> <p><u>Search keywords:</u> Ordinal quantification, dataset shift, prevalence estimation</p> <p><u>Versioning:</u> Yes (the current version is v0.2.0)</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Yes</p> <p><u>How it will be accessible:</u> The data is and will always be accessible on Zenodo (<a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>)</p> <p><u>Methods/software tools to access data:</u> Via a web browser. Detailed instructions for correct use are given at <a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>.</p> <p><u>Repository:</u> Zenodo</p> <p><u>Restrictions on access:</u> No</p>



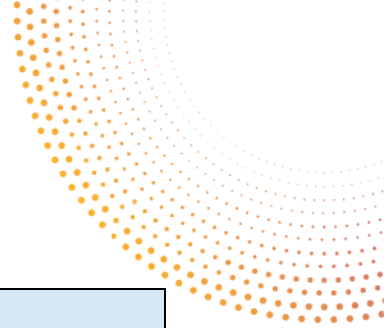


Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. Detailed instructions for correct use are given at <a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>.</p> <p><u>Data and metadata vocabularies</u>: Data is simply a set of pointers to individual data items in the original UCI or OpenML datasets, using comma-separated values. No metadata.</p> <p><u>Use of standard vocabularies</u>: No.</p> <p><u>Mappings to commonly used vocabularies</u>: No.</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International.</p> <p><u>Availability for re-use</u>: Data is already available online at <a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a></p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloadable from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No ethical or legal issues that can have an impact on data sharing and that we can think of.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No informed consent for data sharing and long term preservation included in questionnaires dealing with personal data, since no personal data are involved in these datasets.</p>
Other Issues	N/A

#### 4.2.8 Cherenkov telescope data for ordinal quantification

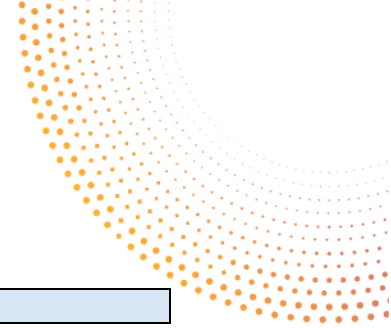
DMP component	AI4Media_Data_10_WP3_CherenkovTelescopeDataOrdinalQuantification_v2 Partner: CNR
Data Summary	<p><u>Purpose</u>: This dataset is targeted at the evaluation of ordinal quantification algorithms (i.e., algorithms that estimate the distribution of the class labels in a set of unseen data items, where there is a total order defined on the set of classes). Used for experiments in T3.7.</p> <p><u>Type/format</u>: The dataset is represented as csv files. It consists of a set of samples to be used in the experimentation. Each sample is represented as a set of indices into the original FACT database of telescope observations, so that the sample can be extracted unambiguously from the original dataset.</p> <p><u>Re-use of existing data</u>: Yes, we use the existing FACT dataset; what we do is providing scripts to extract from it sets of samples that are suitable for testing ordinal quantification algorithms.</p> <p><u>Data origin</u>: <a href="https://factdata.app.tu-dortmund.de/">https://factdata.app.tu-dortmund.de/</a></p>





	<p><u>Expected size</u>: 46.7MB</p> <p><u>Data utility</u>: While not a dataset pertaining to media, this dataset will be useful to anyone wishing to test the accuracy of ordinal quantification algorithms, which are indeed very important in the world of media.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is available online (<a href="https://zenodo.org/records/8172813">https://zenodo.org/records/8172813</a>) and indexed in Google.</p> <p><u>Search keywords</u>: Ordinal quantification, dataset shift, prevalence estimation</p> <p><u>Versioning</u>: Yes (the current version is v0.2.0)</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: The data is and will always be accessible on Zenodo (<a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>)</p> <p><u>Methods/software tools to access data</u>: Web browser. Detailed instructions for correct use are given at <a href="https://zenodo.org/records/8177302">https://zenodo.org/records/8177302</a>.</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. Detailed instructions for correct use are given at <a href="https://zenodo.org/records/8172813">https://zenodo.org/records/8172813</a>.</p> <p><u>Data and metadata vocabularies</u>: Data is simply a set of pointers to individual data items in the original FACT dataset, using comma-separated values. No metadata.</p> <p><u>Use of standard vocabularies</u>: No.</p> <p><u>Mappings to commonly used vocabularies</u>: No.</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International.</p> <p><u>Availability for re-use</u>: Data is already available online at <a href="https://zenodo.org/records/8172813">https://zenodo.org/records/8172813</a></p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloadable from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No ethical or legal issues that can have an impact on data sharing and that we can think of.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No informed consent for data sharing and long term preservation included in questionnaires dealing with personal data, since no personal data are involved in these datasets.</p>





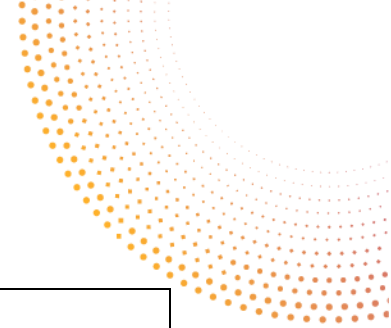
Other Issues	N/A
--------------	-----

### 4.3 Datasets collected in the context of WP4

#### 4.3.1 White matter multiple sclerosis lesion segmentation datasets

DMP component	AI4Media_Data_11_WP4_MRI_MultipleSclerosisLesionSegmentation_v1 Partner: HES-SO
Data Summary	<p><u>Purpose</u>: Segmentation of white matter lesions in Multiple Sclerosis patients, addressing variability in imaging conditions across multiple centers. Used in AI4Media WP4 to promote research in model robustness assessment, interpretability and uncertainty estimation on real-world complex data with a 3D structure.</p> <p><u>Type/format</u>: 3D MRI scans with T1-weighted and FLAIR contrasts, pre-processed to 1mm isovoxel space.</p> <p><u>Re-use of existing data</u>: Yes, the dataset reuses scans contributed by partners in another project</p> <p><u>Data origin</u>: MRI scans collected from multiple clinical centers: Rennes, Bordeaux, Lyon, Ljubljana, Best, and Lausanne.</p> <p><u>Expected size</u>: -</p> <p><u>Data utility</u>: Benchmark for uncertainty estimation and model robustness. Used for testing relevant techniques in WP4.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, available in Zenodo.</p> <p><u>Search keywords</u>: MRI data, Uncertainty Estimation, Robustness, Distributional Shift, Semantic Segmentation, NeuroImaging</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: Metadata follow the Zenodo standards and include information about the imaging centers, scan parameters, and pre-processing steps clearly defined in the reference paper.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Partially. The in-domain training and development data will be openly accessible in Zenodo at <a href="https://zenodo.org/records/7051658">https://zenodo.org/records/7051658</a> and <a href="https://zenodo.org/records/7051692">https://zenodo.org/records/7051692</a>. The Lausanne evaluation set will be restricted and accessible via a Grand-Challenge.</p> <p><u>How it will be accessible</u>: Data deposited in Zenodo. The Lausanne data will be accessible via Grand-Challenge.</p> <p><u>Methods/software tools to access data</u>: Standard MRI viewing software</p> <p><u>Repository</u>: Data in Zenodo at <a href="https://zenodo.org/records/7051658">https://zenodo.org/records/7051658</a> and <a href="https://zenodo.org/records/7051692">https://zenodo.org/records/7051692</a>; Software and models on github at <a href="https://github.com/Shifts-Project/shifts">https://github.com/Shifts-Project/shifts</a></p> <p><u>Restrictions on access</u>: Lausanne data restricted due to privacy and licensing agreements</p>
Making data interoperable	<p><u>Interoperability</u>: Data adheres to standard formats (NIFTI for MRI images)</p> <p><u>Data and metadata vocabularies</u>: Use of standard medical imaging vocabularies and</p>





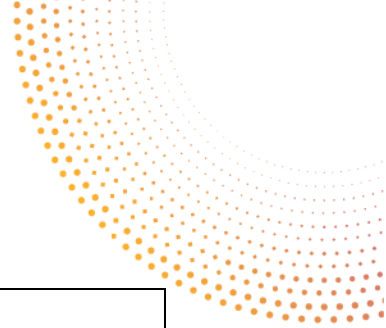
	<p>Zenodo metadata standards.</p> <p><u>Use of standard vocabularies</u>: Yes, standard medical imaging</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons CC BY NC SA 4.0. Data can be downloaded after signing <a href="#">OFSEP</a> data usage agreement.</p> <p><u>Availability for re-use</u>: Lausanne data will be accessible post-publication via controlled access</p> <p><u>Usable by third parties after end of project</u>: Yes, all data except for the restricted Lausanne dataset</p> <p><u>Re-use timeframe</u>: Indefinitely</p> <p><u>Data quality assurance process</u>: Expert validation of segmentation masks and consistency checks across centers</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : Data anonymization, controlled access, user authentication
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Privacy concerns for patient data</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes</p>
Other Issues	N/A

## 4.4 Datasets collected in the context of WP5

### 4.4.1 ObyGaze12 dataset for detection of visual objectification in films

DMP component	AI4Media_Data_12_WP5_Video_ObyGaze12_v1 Partner: UCA
Data Summary	<p><u>Purpose</u>: The ObyGaze12 dataset is made of 1,914 movie clips densely annotated by experts for objectification concepts identified in film studies and psychology. The purpose is to reveal and quantify the usage of complex temporal patterns operated in cinema to produce the cognitive perception of objectification.</p> <p><u>Type/format</u>: csv</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Films are annotated from scratch by experts</p> <p><u>Expected size</u>: 2,000 clips densely annotated with 8 visual concepts</p> <p><u>Data utility</u>: This data is useful to AI researchers to tackle a new AI task, that of detecting visual objectification of characters in videos, and to humanities researchers who can analyze quantitatively and qualitatively the data annotations.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, soon to be formatted with the Croissant ML format. Already online on <a href="#">Github</a>, freely accessible and usable.</p> <p><u>Search keywords</u>: Yes, through the Croissant ML format.</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: Croissant ML is going to be added to the formatting to have a</p>



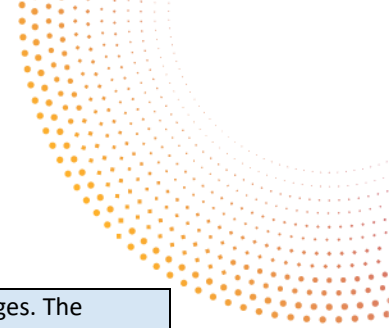


	machine-readable version of the dataset, including the metadata.
Making data openly accessible	<p><u>Data openly accessible</u>: This dataset produced and used in the project is made openly available at <a href="https://github.com/husky-helen/ObyGaze12">https://github.com/husky-helen/ObyGaze12</a></p> <p><u>How it will be accessible</u>: Deposition in an open repository</p> <p><u>Methods/software tools to access data</u>: The data is readable by custom code, and soon described in the Croissant ML format for easier access. As of May 2024, it is publicly hosted on <a href="#">Github</a> with code examples to read and use in an ML environment.</p> <p><u>Repository</u>: Currently on Github and soon on Zenodo</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: Will soon be with Croissant ML formatting</p> <p><u>Data and metadata vocabularies</u>: Will soon be with Croissant ML formatting</p> <p><u>Use of standard vocabularies</u>: Will soon be with Croissant ML formatting</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: CC-BY v4.0</p> <p><u>Availability for re-use</u>: Immediately</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Not determined</p> <p><u>Data quality assurance process</u>: The validity of the dense annotations has been evaluated with different forms of inter-annotator agreement (see publication)</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: Already covered by another project (engineers on data formatting)</p> <p><u>Costs for long-term preservation</u>: Zenodo is used, there are no additional costs as data is light to host</p>
Data security	<u>Security measures</u> : N/A
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: The data corresponds to annotations of films by experts. Experts are project partners. Films can be used for research as per the <a href="#">Directive (EU) 2019/790 of the European Parliament and of the Council of 17 April 2019 on copyright and related rights</a> translated into the <a href="#">French law in 2021</a>.</p>
Other Issues	Refer to other national/funder/sectorial/departmental procedures for data management that you may be using (if any)

#### 4.4.2 MAD-TSC Multilingual Aligned Dataset for Target-dependent Sentiment Classification

<b>DMP component</b>	<b>AI4Media_Data_13_WP5_TEXT_MAD-TSC_v1</b> <b>Partner: CEA</b>
Data Summary	<u>Purpose</u> : MAD-TSC (Multilingual Aligned Dataset for Target-dependent Sentiment Classification) enables sentiment classification in political news. Sentiment is annotated for individual entities occurring in sentences sampled from geographically and politically diversified news sources. The dataset includes 5,110 annotated entity mentions from 4,714 unique sentences. The sentences are aligned across eight





	<p>languages since all sentences have professional translations in all languages. The dataset facilitates the comparison of sentiment classification approaches across languages.</p> <p>MAD-TSC enables a fine-grained understanding of political news when combined with structured knowledge bases such as <a href="#">Wikidata</a> and <a href="#">ParlGov</a>. The predictions can be easily aggregated to assess the orientation of news sources, the public discussion of impactful topics, or the opinions expressed about demographic segments.</p> <p><u>Type/format</u>: The dataset includes textual data and associated annotations. It is distributed in JSON format.</p> <p><u>Re-use of existing data</u>: The sentences are sampled from news articles available via the VoxEurop website.</p> <p><u>Data origin</u>: <a href="https://voxeurop.eu/en/">https://voxeurop.eu/en/</a></p> <p><u>Expected size</u>: 5MB</p> <p><u>Data utility</u>: The dataset is useful to multiple stakeholders. NLP research groups can use it to train and evaluate fine-grained sentiment analysis methods in the eight languages included. Social and/or political scientists can use models pre-trained with the dataset to automatically analyze political texts shared on the Web. Media organizations, such as AI4Media industrial partners, can use the MAD-TSC to analyze the positioning of their own content and that of other news organizations toward different political tendencies and the potential biases toward demographic segments. CEA and VRT exemplified the utility of the dataset via an analysis of Dutch-language news sources from Belgium. Relevant to both WP5 and WP6.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: MAD-TSC has a DOI reference obtained via Zenodo (<a href="https://doi.org/10.5281/zenodo.7940057">https://doi.org/10.5281/zenodo.7940057</a>). The dataset is referenced in the associated publication (<a href="https://doi.org/10.18653/v1/2023.acl-long.461">https://doi.org/10.18653/v1/2023.acl-long.461</a>). It is also distributed on Github, freely accessible and usable.</p> <p><u>Search keywords</u>: Yes, the dataset is searchable with its title as keywords in any public search engine.</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: Metadata include the annotations associated with the entity mentions from the sentences. They are stored in JSON format using a usual three-class sentiment scale (-1=negative; 0 = neutral; 1=positive)</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: Yes, the dataset is available via <a href="#">Github</a> and <a href="#">Zenodo</a></p> <p><u>How it will be accessible</u>: <a href="https://github.com/EvanDufraisse/MAD_TSC">https://github.com/EvanDufraisse/MAD_TSC</a> and <a href="https://zenodo.org/records/7940057">https://zenodo.org/records/7940057</a></p> <p><u>Methods/software tools to access data</u>: The dataset is readable by custom code. It is hosted on Github, along with detailed usage examples.</p> <p><u>Repository</u>: Github and Zenodo</p> <p><u>Restrictions on access</u>: No restriction</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: The dataset is distributed following common practices in the sentiment classification sub-community.</p> <p><u>Data and metadata vocabularies</u>: Metadata are formatted according to common</p>



	<p>practices in the sentiment classification sub-community</p> <p><u>Use of standard vocabularies:</u> Standard vocabulary is used for formatting annotations.</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> MIT</p> <p><u>Availability for re-use:</u> Immediate</p> <p><u>Usable by third parties after end of project:</u> Yes</p> <p><u>Re-use timeframe:</u> Not determined</p> <p><u>Data quality assurance process:</u> Each entity mention was annotated by at least three annotators. Only mentions with sufficient inter-annotator agreement were included.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> The annotation costs were covered via AI4Media and another project.</p> <p><u>Costs for long-term preservation:</u> None. The dataset is hosted on Zenodo</p>
Data security	<p><u>Security measures:</u> Security measures implemented for data protection (incl. controlled access, user authentication, firewalls, VPNs, encryption, back-ups, etc.)</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> No</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The data corresponds to annotations of films by experts. Experts are project partners. Films can be used for research as per the <a href="#">Directive (EU) 2019/790 of the European Parliament and of the Council of 17 April 2019 on copyright and related rights</a> translated into the <a href="#">French law in 2021</a>.</p>
Other Issues	<p>Following the guidelines and requirements of the French national research agency (ANR)</p>

#### 4.4.3 ÖWF Object Detection dataset

DMP component	AI4Media_Data_14_WP5_Video_OWFOd_v1 Partner: JR
Data Summary	<p><u>Purpose:</u> This dataset contains object detection annotations and quality metadata of historic film with scientific and educational content.</p> <p><u>Type/format:</u> JSON and CSV files, and Python code to access the MPEG-4 video files</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> The content is from the collection "Österreichische Bundesinstitut für den Wissenschaftlichen Film (ÖWF)" of the Austrian Mediathek.</p> <p><u>Expected size:</u> 100MB</p> <p><u>Data utility:</u> Researchers/developers working on object detection.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes</p> <p><u>Search keywords:</u> No</p> <p><u>Versioning:</u> Yes</p> <p><u>Metadata creation:</u> Annotations are in JSON format, adhering to the widely used MS COCO format.</p>
Making data	<p><u>Data openly accessible:</u> Data available on Zenodo at</p>



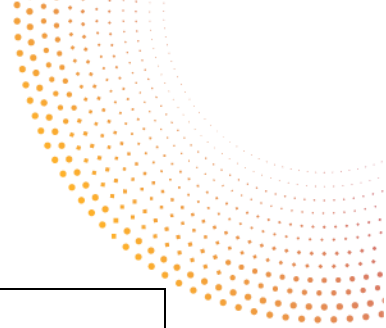
openly accessible	<p><a href="https://zenodo.org/records/8400030">https://zenodo.org/records/8400030</a></p> <p><u>How it will be accessible</u>: Downloadable from Zenodo</p> <p><u>Methods/software tools to access data</u>: Python code for working with the data is included</p> <p><u>Repository</u>: Data is available in Zenodo. Code is available at <a href="#">GitHub</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: Data Card is provided</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: MIT License</p> <p><u>Availability for re-use</u>: Available</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Annotations have been manually checked by two experts</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: 1,000€, covered by AI4Media</p> <p><u>Costs for long-term preservation</u>: None, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Persons are depicted, but the content has been cleared to be provided on the web</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

#### 4.4.4 PeopleAtPlaces dataset

DMP component	AI4Media_Data_15_WP5_Video_People@Places_v1 Partner: JR
Data Summary	<p><u>Purpose</u>: The PeopleAtPlaces dataset provides annotations per image for training classifiers for bustle (i.e., more or less populated) and cinematographic shot type (from extreme close-up to extreme long shot). The annotations are provided for images of the Places365-Standard training set, and have been generated using an automatic pipeline. 100 images per class are provided as validation set, and the annotations of these images have been manually checked and corrected.</p> <p><u>Type/format</u>: JSON and CSV files, and Python code to access the JPEG images</p> <p><u>Re-use of existing data</u>: Yes</p> <p><u>Data origin</u>: Places365-Standard dataset</p> <p><u>Expected size</u>: ~1GB</p> <p><u>Data utility</u>: Researchers/developers working video shot classification</p>







Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: No</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: Annotations are in JSON format (following human pose description formats) and CSV</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available in Zenodo at <a href="https://zenodo.org/records/8398916">https://zenodo.org/records/8398916</a></p> <p><u>How it will be accessible</u>: Downloadable from Zenodo</p> <p><u>Methods/software tools to access data</u>: Python code for reproducing the dataset is included</p> <p><u>Repository</u>: Data is available in Zenodo. Code is available at <a href="#">GitHub</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: No</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: CC-BY-4.0 License</p> <p><u>Availability for re-use</u>: Available</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Annotations have been manually checked</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: 1,000€ covered by a national research project</p> <p><u>Costs for long-term preservation</u>: None, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Persons are depicted, but the content has been cleared to be provided with the source dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

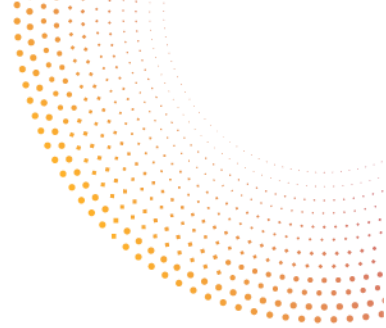
#### 4.4.5 ToDY dataset for visual time of day and season classification

<b>DMP component</b>	<b>AI4Media_Data_16_WP5_Video_ToDY_v1</b> <b>Partner: JR</b>
Data Summary	<p><u>Purpose</u>: The dataset provides training and validation data for classifying images by time of day and season (time of year). The images are taken from the Skyfinder dataset, containing webcam images along with timestamps and geolocation. The annotations have been automatically derived from the metadata, which is sufficiently precise for season. For time of day, the annotation have been manually checked and corrected. The dataset contains 2,790 training files and 311 validation files per class for</p>



	<p>season, and 986 training files and 110 validation files per class for time of day.</p> <p><u>Type/format</u>: CSV files, and Python code to access the MPEG-4 video files</p> <p><u>Re-use of existing data</u>: Yes</p> <p><u>Data origin</u>: Skyfinder dataset</p> <p><u>Expected size</u>: 500MB</p> <p><u>Data utility</u>: Researchers/developers working on visual classification of seasons and times of day</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: No</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: Annotations are in CSV format, partly mined from existing metadata, partly manually annotated</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Available on Zenodo at <a href="https://zenodo.org/records/8398861">https://zenodo.org/records/8398861</a></p> <p><u>How it will be accessible</u>: Downloadable from Zenodo</p> <p><u>Methods/software tools to access data</u>: Python code for recreating the dataset is provided</p> <p><u>Repository</u>: Data available on Zenod and code is available on GitHub at <a href="https://github.com/wbailer/ToDY">https://github.com/wbailer/ToDY</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: CC-BY-4.0 License</p> <p><u>Availability for re-use</u>: Available</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Annotations have been manually checked</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: 1,000€ covered by a national research project</p> <p><u>Costs for long-term preservation</u>: nNone, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: n/a</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

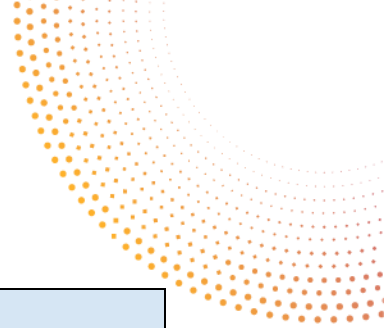




#### 4.4.6 VISIONE Feature Repository

DMP component	AI4Media_Data_17_WP5_ImageFeatures_VisioneFeatureRepository_v1 Partner: CNR
Data Summary	<p><b>Purpose:</b> Three datasets containing a diverse set of features extracted from the V3C1+V3C2, VBSLHE, and MVK video datasets, respectively. These features were extracted by CNR and employed in the VISIONE video retrieval system during the latest editions of the Video Browser Showdown (VBS) competition (<a href="https://www.videobrowsershowdown.org/">https://www.videobrowsershowdown.org/</a>).</p> <p>This repository comprises the following files related to V3C1+V3C2, VBSLHE, and MVK datasets:</p> <ul style="list-style-type: none"> <li>• Tab-separated files reporting the segmentation of each video into shots.</li> <li>• A Python script designed to extract the middle frame of each video segment</li> <li>• An archive of image features extracted using the <a href="#">ALADIN</a> model for all the segment's middle frames.</li> <li>• An archive of image features extracted using the <a href="#">CLIP ViT-H/14 - LAION-2B</a> model for all the segment's middle frames.</li> <li>• An archive of image features extracted using the <a href="#">CLIP ViT-L/14</a> model for all the segment's middle frames. .</li> <li>• An archive of image features extracted using the <a href="#">CLIP2Video</a> model for all the segment's middle frames.</li> <li>• An archive of objects detected using the <a href="#">Faster R-CNN+Inception ResNet (trained on the Open Images V4)</a> model for all the segment's middle frames.</li> <li>• An archive of objects detected using <a href="#">Mask R-CNN (trained on LVIS)</a> model for all the segment's middle frames. .</li> <li>• An archive of objects detected using <a href="#">VfNet (trained on COCO dataset)</a> model for all the segment's middle frames .</li> </ul> <p><b>Type/format:</b> CSV files, HDF5 files and a Python Script,</p> <p><b>Re-use of existing data:</b> Yes, we are reusing the original videos or keyframes from the three datasets: V3C1+V3C2, VBSLHE, and MVK. These were used solely for feature extraction and are not redistributed.</p> <p><b>Data origin:</b> V3C1+V3C2, VBSLHE and MVK datasets</p> <p><b>Expected size:</b> 90GB</p> <p><b>Data utility:</b> These datasets are essential for testing and enhancing video retrieval systems, particularly for research groups involved in the VBS competition. They also provide researchers with pre-extracted image features for more than 3M video keyframes, saving the time and effort required for feature extraction.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes</p> <p><b>Search keywords:</b> Yes</p> <p><b>Versioning:</b> Yes</p> <p><b>Metadata creation:</b> N/A.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> Yes, available in Zenodo</p> <p><b>How it will be accessible:</b> Downloadable from Zenodo</p> <p><b>Methods/software tools to access data:</b> The data can be accessed using any software</p>



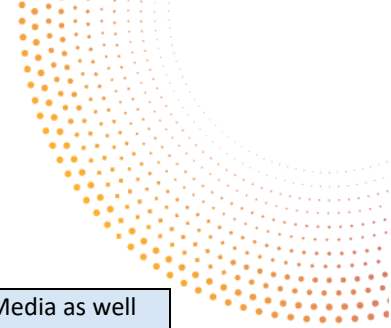


	<p>that supports HDF5 and CSV file formats, such as Python or Matlab.</p> <p><u>Repository</u>: The data has been deposited in three Zenodo repositories:  <a href="https://zenodo.org/records/8188570">https://zenodo.org/records/8188570</a> (features extracted from V3C1+V3C2)  <a href="https://zenodo.org/records/8355037">https://zenodo.org/records/8355037</a> (features extracted from MVK),  <a href="https://zenodo.org/records/10013329">https://zenodo.org/records/10013329</a> (Features extracted from VBSLHE)</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: Yes, we will use standard vocabularies for all data types in our dataset to allow inter-disciplinary interoperability. This includes using the HDF5 format for structured data storage, and standard CSV formatting for tabular data.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: undetermined</p> <p><u>Data quality assurance process</u>: checks for data integrity and completeness</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: none, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

#### 4.4.7 COCO, LVIS, Open Images V4 classes mapping

DMP component	AI4Media_Data_18_WP5_TEXT_COCO- LVIS-Open ImagesV4-classes Mapping_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: This repository contains a mapping between the classes of COCO, LVIS, and Open Images V4 datasets into a unique set of 1,460 classes. COCO contains 80 classes, LVIS contains 1,460 classes, Open Images V4 contains 601 classes.</p> <p>The mapping of these classes was obtained using a semi-automatic procedure in order to have a unique final list of 1,460 classes. Moreover, a hierarchy for each class, using wordnet, has been generated.</p> <p><u>Type/format</u>: CSV and TXT file</p> <p><u>Re-use of existing data</u>: Using the class names of COCO, LVIS, and Open Images V4 datasets</p> <p><u>Data origin</u>: COCO, LVIS, and Open Images V4 datasets</p> <p><u>Expected size</u>: 1 MB</p>



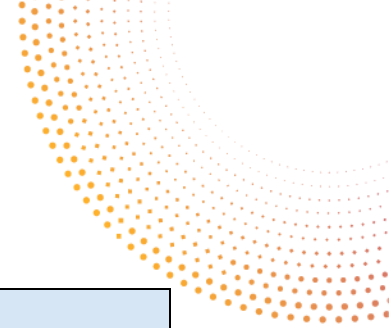


	<b>Data utility:</b> This data will be useful for researchers and developers in AI4Media as well as third parties. It can be used, for example, for systems that aim to utilize multiple object detectors trained on these datasets simultaneously by leveraging the unified class mapping and hierarchy /as done in the VISIONE video retrieval system.
Making data findable, incl. provisions for metadata	<b>Is data discoverable:</b> Yes, available on Zenodo <b>Search keywords:</b> Yes <b>Versioning:</b> Yes <b>Metadata creation:</b> N/A
Making data openly accessible	<b>Data openly accessible:</b> Yes, available in Zenodo: <a href="https://zenodo.org/records/7194300">https://zenodo.org/records/7194300</a> <b>How it will be accessible:</b> Downloadable from Zenodo <b>Methods/software tools to access data:</b> Any software supporting TXT and CSV files <b>Repository:</b> Zenodo <b>Restrictions on access:</b> No
Making data interoperable	<b>Interoperability:</b> Yes <b>Data and metadata vocabularies:</b> N/A <b>Use of standard vocabularies:</b> N/A <b>Mappings to commonly used vocabularies:</b> N/A
Increase data re-use	<b>Licence:</b> Creative Commons Attribution 4.0 International <b>Availability for re-use:</b> Yes <b>Usable by third parties after end of project:</b> Yes <b>Re-use timeframe:</b> Undetermined <b>Data quality assurance process:</b> Manually checks for integrity, completeness
Allocation of resources	<b>Costs for making data FAIR:</b> N/A <b>Costs for long-term preservation:</b> N/A
Data security	<b>Security measures:</b> Zenodo has appropriate security mechanisms
Ethical aspects	<b>Possible ethical and legal aspects preventing sharing:</b> No <b>Is informed consent for data sharing and long term preservation given:</b> N/A
Other Issues	N/A

#### 4.4.8 Bus violence dataset

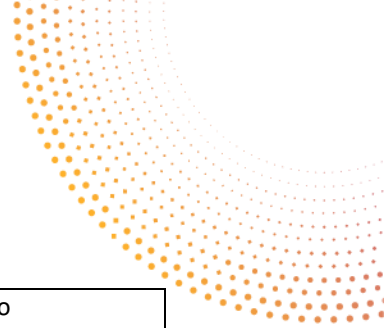
<b>DMP component</b>	<b>AI4Media_Data_19_WP5_Video_BusViolence_v1</b> <b>Partner:</b> CNR
Data Summary	<b>Purpose:</b> The Bus Violence dataset is a large-scale collection of videos depicting violent and non-violent situations in public transport environments. This benchmark was gathered from multiple cameras located inside a moving bus where several people simulated violent actions, such as stealing an object from another person, fighting between passengers, etc. It contains 1,400 video clips manually annotated as having or not violent scenes, making it one of the biggest benchmarks for video violence





	<p>detection in the literature.</p> <p>Specifically, videos are recorded from three cameras at 25 frames per second - two cameras located in the corners of the bus (with resolution 960x540 px) and one fisheye in the middle (1280x960 px). The clips have a minimum length of 16 frames and a maximum of 48 frames, capturing a very precise action (either violence or non-violence). The dataset is perfectly balanced, containing 700 videos of violence and 700 videos of non-violence.</p> <p><u>Type/format</u>: MP4 files</p> <p><u>Re-use of existing data</u>: No, this is original data.</p> <p><u>Data origin</u>: Surveillance cameras mounted on public transportation vehicles</p> <p><u>Expected size</u>: 500MB</p> <p><u>Data utility</u>: This dataset is useful to researchers working in AI and Computer Vision to train and test solutions for visual detection of violence scenes.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: Yes</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available on Zenodo <a href="https://zenodo.org/records/7044203">https://zenodo.org/records/7044203</a></p> <p><u>How it will be accessible</u>: From Zenodo</p> <p><u>Methods/software tools to access data</u>: The data can be accessed using any software that supports MP4 standard video files.</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: Yes, vocabulary is the standard binary one for violence detection (1 – video containing violence, 0 – video without violence).</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Checks for data integrity and completeness</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: none, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p>



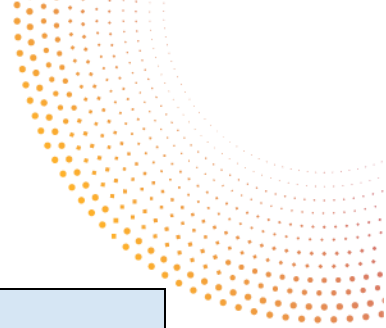


aspects	<u>Is informed consent for data sharing and long term preservation given:</u> No
Other Issues	N/A

#### 4.4.9 Pest Sticky Traps dataset

DMP component	AI4Media_Data_20_WP5_Image_PestStickyTraps_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> The Pest Sticky Traps dataset is a collection of yellow chromotropic sticky trap pictures specifically designed for training/testing deep learning models to automatically count insects and estimate pest populations.</p> <p>Images were manually annotated by some experts of the Department of Agriculture, Food and Environment of the University of Pisa (Italy) by putting a dot over the centroids of each identified insect. Specifically, we labeled insects as belonging to the category “whitefly” considering two different species, i.e., the sweet potato whitefly (<i>Bemisia tabaci</i>) (Gennadius) and the greenhouse whitefly (<i>Trialeurodes vaporariorum</i>) (Westwood).</p> <p><u>Type/format:</u> JPG and CSV files</p> <p><u>Re-use of existing data:</u> No, this is original data.</p> <p><u>Data origin:</u> Pictures from smartphone</p> <p><u>Expected size:</u> 150MB</p> <p><u>Data utility:</u> The dataset can be useful to researchers working in AI and Computer vision to develop and tests solutions for visual counting. In particular, density estimation of insects.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes</p> <p><u>Search keywords:</u> Yes</p> <p><u>Versioning:</u> Yes</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Yes, available in Zenodo: <a href="https://zenodo.org/records/7801239">https://zenodo.org/records/7801239</a></p> <p><u>How it will be accessible:</u> From Zenodo</p> <p><u>Methods/software tools to access data:</u> The data can be accessed using any software that supports JPG and JSON standard files.</p> <p><u>Repository:</u> Zenodo repository</p> <p><u>Restrictions on access:</u> No</p>
Making data interoperable	<p><u>Interoperability:</u> Yes</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use:</u> Yes</p>





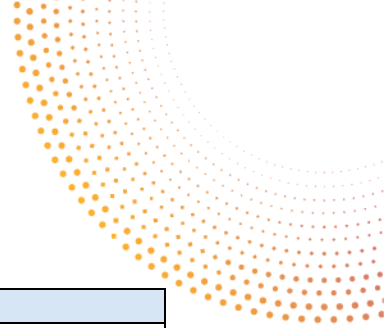
	<p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Checks for data integrity and completeness</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: None, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

#### 4.4.10 Virtual World Fallen People dataset

DMP component	AI4Media_Data_21_WP5_IMAGE_VWFP_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: A synthetic dataset for visual fallen people detection comprising images extracted from the highly photo-realistic video game Grand Theft Auto V developed by Rockstar North. Each image is labeled by the game engine, which provides bounding boxes and statuses (fallen or non-fallen) of people present in the scene.</p> <p>The dataset comprises 6,071 synthetic images depicting 7,456 fallen and 26,125 non-fallen pedestrian instances in various looks, camera positions, background scenes, lightning, and occlusion conditions.</p> <p>VWFP aims to improve fallen-people detection systems that often lack difficult-to-gather training data (e.g., people in high-risk environments) using synthetic data from virtual worlds.</p> <p><u>Type/format</u>: PNG images and annotations in CSV file.</p> <p><u>Re-use of existing data</u>: No, the dataset is created within AI4Media.</p> <p><u>Data origin</u>: Grand Theft Auto V developed by Rockstar North</p> <p><u>Expected size</u>: 11.5 GB</p> <p><u>Data utility</u>: It is useful to researchers working on AI and computer Vision to develop and test solutions to detect fallen people.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available at <a href="https://zenodo.org/records/6394684">https://zenodo.org/records/6394684</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Zenodo supports versioning</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is available on Zenodo at <a href="https://zenodo.org/records/6394684">https://zenodo.org/records/6394684</a></p> <p><u>How it will be accessible</u>: The data can be downloaded freely.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: Zenodo</p>







	<u>Restrictions on access</u> : None.
Making data interoperable	<u>Interoperability</u> : The file structure makes the use of the dataset easy. <u>Data and metadata vocabularies</u> : Annotations contain bounding box coordinates and class (fallen/non-fallen) for each person in the images. <u>Use of standard vocabularies</u> : N/A <u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<u>Licence</u> : GNU General Public License (GPL) v3 <u>Availability for re-use</u> : N/A <u>Usable by third parties after end of project</u> : Data already publicly shared. <u>Re-use timeframe</u> : N/A <u>Data quality assurance process</u> : N/A
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : Zenodo has appropriate security mechanisms
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : N/A <u>Is informed consent for data sharing and long term preservation given</u> : N/A
Other Issues	N/A

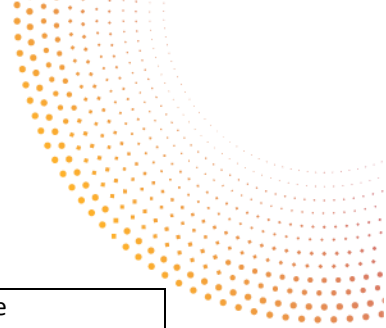
#### 4.4.11 SR\_BVI-DVC super-resolution dataset

DMP component	AI4Media_Data_22_WP5_VIDEO_SR-BVI-DVC_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: This dataset is comprised of 2,000 sequences of 64 frames each, all in 4K resolution. The source material for these sequences is 200 original clips from the BVI-DVC public dataset. Each of these clips has been upscaled using nine distinct super-resolution algorithms, encompassing both traditional and Deep Learning-based methods. The traditional methods employed include nearest-neighbor interpolation, bicubic interpolation, and bilinear interpolation. The Deep Learning-based methods utilized are Real-BasicVSR, RVRT, BasicVSR, SwinIR-Real, SwinIR-Classical, and Real-ESRGAN. The objective of this dataset is to create pairs of "real" and "fake" or upscaled video content in 4K resolution, featuring identical scenes. This dataset has been used for the training and evaluation of a Super-Resolution detection model specifically designed for 4K videos.</p> <p><u>Type/format</u>: Each video is represented by a sequence of 64 .png frames at 3840 x 2160 resolution</p> <p><u>Re-use of existing data</u>: The original 4K clips are being used as "real" or natively 4K content.</p> <p><u>Data origin</u>: The original dataset is a collection from different sources: Videvo Free Stock Video Footage set, IRIS32 Free 4K Footage set, Harmonics database, BVI-Texture database, MCML 4K video quality database, BVI-HFR database, SJTU 4K video database, LIVE-Netflix database, Mitch Martinez Free 4K Stock Footage set, Dareful Free 4K Stock Video data set, MCL-V database, MCL-JCV database, Netflix Chimera,</p>



	<p>TUM HD databases, Ultra Video Group-Tampere University database</p> <p><u>Expected size:</u> Around 2TB</p> <p><u>Data utility:</u> The dataset can be employed to train and evaluate Super-Resolution models, Super-Resolution detection models or image/video quality assessment models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No, because the dataset described is stored on BSC's internal corporate systems, although there are plans for future sharing.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> The README includes information about the original video sources, as well as details about upscaling algorithms and their parameters.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is not openly accessible, although there are plans for future sharing.</p> <p><u>How it will be accessible:</u> Via a third-party repository link, only shared via explicit request.</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> The BSC research team has access to the dataset.</p>
Making data interoperable	<p><u>Interoperability:</u> The data is preserved in two formats: individual .png files and compressed files for each upscaling method. This allows for flexibility in accessing and utilizing the data.</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The original database was compiled by the University of Bristol, Bristol, UK, comprising sequences originally generated by various sources. In the same manner, all intellectual property rights remain with the originators of each sequence. The test sequences shall only be used for academic research (no commercial use).</p> <p><u>Availability for re-use:</u> Yes, in the future, for academic research only.</p> <p><u>Usable by third parties after end of project:</u> Yes, in the future, for academic research only.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The data is stored in secure BSC machines, with user control access via unique credentials and VPN. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>





Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Source video license incompatibilities.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.4.12 BSC4K super-resolution dataset

DMP component	AI4Media_Data_23_WP5_VIDEO_BSC4K_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> This dataset contains paired video sequences at 1080p and 4K resolution recorded simultaneously. The motivation of the dataset is to overcome the challenges introduced by artificially downscaling high-resolution content to obtain the low-resolution counterparts, a common method in Super-Resolution tasks. The first version of the dataset contains 33 4K and 33 1080p videos, cut to 64 frames each, recorded indoor and outdoor with a single DSLR camera. The 1080p versions are upscaled to 4K to create pairs of original/synthetic data. The traditional methods employed include nearest-neighbor interpolation, bicubic interpolation, and bilinear interpolation. The Deep Learning-based methods utilized are Real-BasicVSR, RVRT, BasicVSR, SwinIR-Real, SwinIR-Classical, and Real-ESRGAN.</p> <p><u>Type/format:</u> Each video is represented by a sequence of 64 .png frames at 3840 x 2160 resolution or 1920 x 1080.</p> <p><u>Re-use of existing data:</u> All data is manually collected and original.</p> <p><u>Data origin:</u> All data is manually collected and original.</p> <p><u>Expected size:</u> Around 300GB</p> <p><u>Data utility:</u> The dataset can be employed to train and evaluate Super-Resolution models, Super-Resolution detection models or image/video quality assessment models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The dataset is not publicly available. There are no plans for data sharing.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> The dataset includes information and metadata about the 4K and 1080p video sources, as well as details about upscaling algorithms and their parameters.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is not openly accessible. There are no plans for data sharing.</p> <p><u>How it will be accessible:</u> N/A</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> The BSC research team has access to the dataset.</p>
Making data interoperable	<p><u>Interoperability:</u> The data is preserved in two formats: individual .png files and compressed files for each upscaling method. This allows for flexibility in accessing and</p>

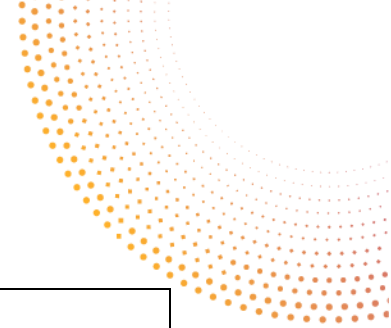


	<p>utilizing the data.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: TBD</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored in secure BSC machines, with user control access via unique credentials and VPN. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Every individual who is present in the video sequences has provided their explicit consent for their recording.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.4.13 Night and Day Instance Segmented Park dataset

DMP component	AI4Media_Data_24_WP5_IMAGE_NDISPark_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: A collection of images of parking lots for vehicle detection, segmentation, and counting. Each image is manually labeled with pixel-wise masks and bounding boxes localizing vehicle instances.</p> <p>The dataset includes about 250 images depicting several parking areas describing most of the problematic situations that we can find in a real scenario: seven different cameras capture the images under various weather conditions and viewing angles. Another challenging aspect is the presence of partial occlusion patterns in many scenes such as obstacles (trees, lampposts, other cars) and shadowed cars.</p> <p>The main peculiarity is that images are taken during the day and the night, showing utterly different lighting conditions.</p> <p><u>Type/format</u>: JPG and JSON files</p> <p><u>Re-use of existing data</u>: No, this is original data.</p> <p><u>Data origin</u>: Surveillance cameras</p> <p><u>Expected size</u>: 120MB</p> <p><u>Data utility</u>: This dataset can be used by researchers working on AI and Computer Vision to develop and test solutions for visual parking space analysis.</p>



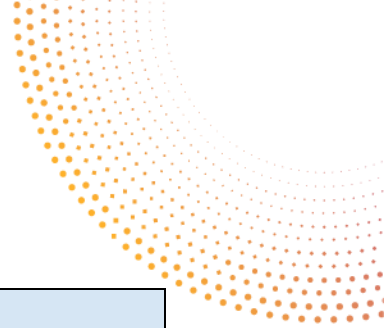


Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: Yes</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available on Zenodo: <a href="https://zenodo.org/records/6560823">https://zenodo.org/records/6560823</a></p> <p><u>How it will be accessible</u>: From Zenodo</p> <p><u>Methods/software tools to access data</u>: The data can be accessed using any software that supports JPG and JSON standard files.</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: Checks for data integrity and completeness</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: None, using public repositories</p>
Data security	<p><u>Security measures</u>: Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	N/A

#### 4.4.14 ArXiv abstracts for authorship analysis dataset

<b>DMP component</b>	<b>AI4Media_Data_25_WP5_TEXT_ArXivAbstracts_v1</b> <b>Partner: CNR</b>
Data Summary	<p><u>Purpose</u>: This data set comprises a labelled training set used in the experimentation of the paper "Binary Quantification and Dataset Shift: An Experimental Investigation" from CNR. The data is extracted from the McAuley data set of product reviews on Amazon.</p> <p><u>Type/format</u>: Raw text partitioned in samples. A file per sample</p> <p><u>Re-use of existing data</u>: This is sample from ArXiv data</p> <p><u>Data origin</u>: <a href="https://info.arxiv.org/help/bulk_data/index.html">https://info.arxiv.org/help/bulk_data/index.html</a></p>



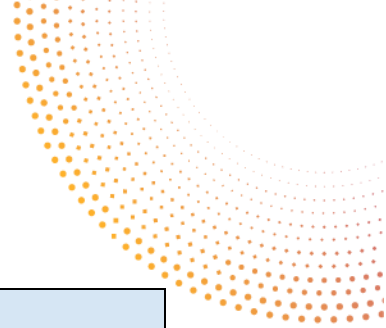


	<p><u>Expected size</u>: 1.3GB</p> <p><u>Data utility</u>: Benchmark for Authorship Analysis</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available online (<a href="https://zenodo.org/records/7404702">https://zenodo.org/records/7404702</a> ) and indexed in Google.</p> <p><u>Search keywords</u>: author analysis, text classification</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available in Zenodo: <a href="https://zenodo.org/records/7404702">https://zenodo.org/records/7404702</a></p> <p><u>How it will be accessible</u>: From Zenodo</p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: A single csv file with abstract of a paper and a unique author's identifier. No metadata.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International</p> <p><u>Availability for re-use</u>: Data is already available online.</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the Zenodo servers. Zenodo fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.4.15 Florence4D facial expression dataset

<b>DMP component</b>	<b>AI4Media_Data_26_WP5_4D_Florence4D_v1</b> <b>Partner: UNIFI</b>
Data Summary	<u>Purpose</u> : Collection of 3D sequences representing transitions between a diverse set of facial expressions. This dataset aims at filling the gap in the availability of 4D data.



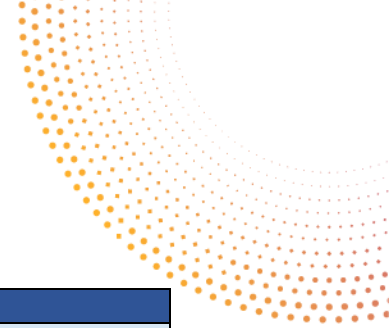


	<p><u>Type/format</u>: Sequences of point cloud data in OBJ format</p> <p><u>Re-use of existing data</u>: Part of the dataset contains CoMA meshes.</p> <p><u>Data origin</u>: CoMA, synthetic models from the web, real 3D Scans.</p> <p><u>Expected size</u>: 30 GB</p> <p><u>Data utility</u>: The dataset is useful for all research and industrial teams working on 3D representation learning and facial expression recognition.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, available in UNIFI website (<a href="https://www.micc.unifi.it/resources/datasets/florence-4d-facial-expression/">https://www.micc.unifi.it/resources/datasets/florence-4d-facial-expression/</a>)</p> <p><u>Search keywords</u>: Florence 4D Facial Expression Dataset</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Accessible from UNIFI website <a href="https://www.micc.unifi.it/resources/datasets/florence-4d-facial-expression/">https://www.micc.unifi.it/resources/datasets/florence-4d-facial-expression/</a> after completing a form <a href="https://forms.gle/AgLNSjjMZYXdUAK68">https://forms.gle/AgLNSjjMZYXdUAK68</a></p> <p><u>How it will be accessible</u>: From UNIFI website</p> <p><u>Methods/software tools to access data</u>: OBJ is an open format.</p> <p><u>Repository</u>: Google drive</p> <p><u>Restrictions on access</u>: Data accessible after form completion</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: OBJ format, annotation in plain text</p> <p><u>Use of standard vocabularies</u>: Expressions are classified according to Plutchik's wheel of emotions</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Common Attribution License-CC-BY v4.0</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Free for non-commercial purposes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Manual checking</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: N/A</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Written consents were collected from subjects for using their 3D face scans.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes.</p>
Other Issues	N/A

#### 4.4.16 Neuromorphic Event-based Facial Expression Recognition dataset

DMP	AI4Media_Data_27_WP5_EventRGB_NEFER_v1
-----	--

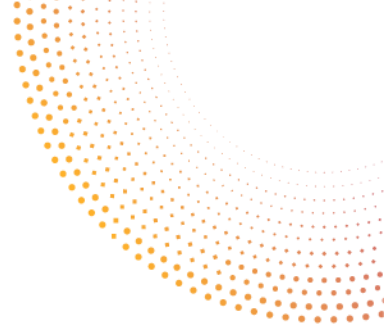




component	Partner: UNIFI
Data Summary	<p><u>Purpose</u>: Paired RGB and event videos representing human faces labeled with the respective emotions and also annotated with face bounding boxes and facial landmarks. Currently, the only event+RGB emotion dataset.</p> <p><u>Type/format</u>: RGB, Raw Event Data, Reconstructed Event Frames</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Acquisition</p> <p><u>Expected size</u>: 10 Gb</p> <p><u>Data utility</u>: The dataset is useful for all research and industrial teams working on neuromorphic representation learning and facial expression recognition.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, DOI: <a href="https://doi.org/10.1109/CVPRW59228.2023.00432">https://doi.org/10.1109/CVPRW59228.2023.00432</a></p> <p><u>Search keywords</u>: NEFER</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Accessible on GitHub <a href="https://github.com/miccunifi/NEFER">https://github.com/miccunifi/NEFER</a></p> <p><u>How it will be accessible</u>: From GitHub</p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: GitHub, Google drive</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes</p> <p><u>Data and metadata vocabularies</u>: RGB+CSV</p> <p><u>Use of standard vocabularies</u>: Expressions are classified according to Ekman's emotion classification</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Common Attribution License-CC-BY v4.0</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Free for non-commercial purposes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Manual checking</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: N/A</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Written consents were collected from subjects for using their 3D face scans</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes</p>
Other Issues	N/A







#### 4.4.17 PEM360 dataset of 360° videos

DMP component	AI4Media_Data_28_WP5_Video_PEM360_v1 Partner: UCA
Data Summary	<p><u>Purpose:</u> PEM360 includes user head movements and gaze recordings in 360° videos, along with self-reported emotional ratings of valence and arousal, and continuous physiological measurement of electrodermal activity and heart rate. Its purpose is to understand the connection between user attention, user emotions and immersive content.</p> <p><u>Type/format:</u> csv</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> The 360° video stimuli are selected from an existing publicly available dataset</p> <p><u>Expected size:</u> 34 users on 7 360° videos each, amounting to 238 traces of gaze motion and emotion ratings and physiological measures</p> <p><u>Data utility:</u> This data is useful to better understand non-motion factors (such as emotion or video content) that can non-deterministically impact models to predict user attention.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data is publicly available, but no metadata file can be automatically parsed.</p> <p><u>Search keywords:</u> PEM360</p> <p><u>Versioning:</u> No</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> This dataset produced and used in the project is made openly available at <a href="https://gitlab.com/PEM360/PEM360/">https://gitlab.com/PEM360/PEM360/</a></p> <p><u>How it will be accessible:</u> Deposition in an open repository</p> <p><u>Methods/software tools to access data:</u> The data is readable by custom code made available in the dataset repository</p> <p><u>Repository:</u> Gitlab</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> CC-BY</p> <p><u>Availability for re-use:</u> Immediately</p> <p><u>Usable by third parties after end of project:</u> Yes</p> <p><u>Re-use timeframe:</u> Not determined</p> <p><u>Data quality assurance process:</u> The validity of the collected user traces has been analyzed through EDA and known correlations have been retrieved.</p>



Allocation of resources	<p><u>Costs for making data FAIR:</u> Not engaged</p> <p><u>Costs for long-term preservation:</u> Gitlab is used, there are no additional costs as data is light to host</p>
Data security	<u>Security measures:</u> N/A
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> No</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> Yes, the user experiment carried out to collect the data has been approved by an IRB and informed consent has been formally collected (and user retribution financially).</p>
Other Issues	N/A

#### 4.4.18 VRT-Sum video summarization dataset

DMP component	AI4Media_Data_29_WP5_VIDEO-TEXT_VRT-Sum_v1 Partner: VRT, CERTH
Data Summary	<p><u>Purpose:</u> This dataset is composed of a set of videos of varying content (e.g. sports, news, interviews), the professionally-edited video summaries (one per video), and a short text that is used as a voice-over during the presentation of the video summary (one or two per video). It will be used in T5.1 for training and evaluation purposes, assisting the development of text-driven methods for video summarization,</p> <p><u>Type/format:</u> The video files are in MP4 format, and the text descriptions (in English and Dutch) are in WORD files.</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> The VRT partner of AI4Media</p> <p><u>Expected size:</u> ~12.6GB</p> <p><u>Data utility:</u> It will be useful to WP5 partners working on the development of video summarization methods, for training and evaluation purposes.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> Any metadata generated to assist training and evaluation of video summarization methods (e.g. deep feature vectors representing the visual content of video frames and the textual content of the description, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be created to facilitate their re-use.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No. The videos and the associated text descriptions are proprietary data of VRT that we do not have the right to make publicly available.</p> <p><u>How it will be accessible:</u> N/A</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> No. The videos and the associated text descriptions are proprietary data of VRT that we do not have the right to share.</p>



	<p><u>Data and metadata vocabularies:</u> Any metadata generated to assist training and evaluation of text-driven video summarization methods (e.g. deep feature vectors representing the visual content of video frames and the textual content of the description, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be also created to facilitate their re-use.</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The videos and the associated text descriptions are proprietary data of VRT that we do not have the right to grant a licence.</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> The videos and the associated text descriptions are proprietary data of VRT, that we do not have the right to make available for re-use by third parties.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset and any created metadata (e.g. deep feature vectors representing the visual content of video frames and the textual content of the description, or data about the shot-level structure of the videos) will be hosted on CERTH's servers. CERTH fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimise the risk of data loss.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	No

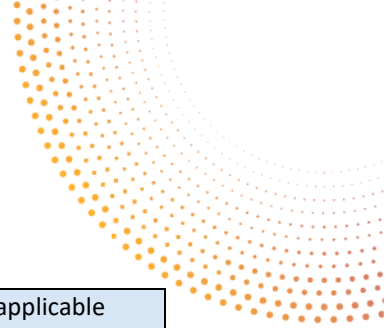
#### 4.4.19 CA-SUM pretrained video summarization models

<b>DMP component</b>	<b>AI4Media_Data_30_WP5_MODEL_CA-SUM-models_v1</b> <b>Partner: CERTH</b>
Data Summary	<p><u>Purpose:</u> This dataset contains pretrained models of the CA-SUM network architecture for video summarization, that is presented in CERTH's work titled "<a href="#">Summarizing Videos using Concentrated Attention and Considering the Uniqueness and Diversity of the Video Frames</a>", in Proc. ACM ICMR 2022. It will be used in T5.1 for comparison with other summarization approaches, and for assisting the implementation of the relevant use cases of WP8.</p> <p><u>Type/format:</u> The model files are in PT format.</p> <p><u>Re-use of existing data:</u> No</p>



	<p><u>Data origin</u>: CERTH</p> <p><u>Expected size</u>: ~195MB</p> <p><u>Data utility</u>: It will be useful to WP5 partners working on the development of video summarization methods, for comparison purposes. It will be useful also to WP8 partners involved in use cases that require the production of video summaries.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: CA-SUM pretrained models, video summarization</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available on Zenodo</p> <p><u>How it will be accessible</u>: The pretrained models are publicly available on Zenodo at <a href="https://zenodo.org/records/6562992">https://zenodo.org/records/6562992</a></p> <p><u>Methods/software tools to access data</u>: The pretrained models can be found using any web browser and the above listed search keywords</p> <p><u>Repository</u>: Zenodo at <a href="https://zenodo.org/records/6562992">https://zenodo.org/records/6562992</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Yes. The pretrained models can be used both on Windows and Linux machines.</p> <p><u>Data and metadata vocabularies</u>: A documentation of this dataset is available on Zenodo (<a href="https://zenodo.org/records/6562992">https://zenodo.org/records/6562992</a>). Details about the CA-SUM video summarization method are available on the relevant scientific publication (<a href="https://zenodo.org/records/6759007">https://zenodo.org/records/6759007</a>)</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is provided for academic, non-commercial use only.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes. The dataset is publicly available for use.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: No. The dataset is publicly available on the Zenodo platform that is cost-free.</p> <p><u>Costs for long-term preservation</u>: No. The dataset is publicly available on the Zenodo platform that is cost-free.</p>
Data security	<p><u>Security measures</u>: The dataset is hosted on <a href="https://zenodo.org/">Zenodo</a>, which is a EU open research repository. Zenodo is powered by <a href="https://www.cern.ch/en/infrastructure/data-centre">CERN Data Centre</a> and the <a href="https://www.invenio.org/">Invenio digital library framework</a> and is fully run on open source products all the way through. Zenodo servers are managed via <a href="https://openstack.org/">OpenStack</a> and <a href="https://puppet.com/">Puppet</a> configuration management system which ensures that the servers always have the latest security patches applied. The</p>





	dataset is hosted also on CERTH's servers. CERTH fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimise the risk of data loss.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	No

#### 4.4.20 Audio phylogeny dataset

DMP component	AI4Media_Data_31_WP5_Audio_AudioPhylogeny_v1 Partner: FhG-IDMT
Data Summary	<p><u>Purpose:</u> The IDMT Audio Phylogeny Dataset contains audio phylogeny trees for evaluation of audio phylogeny algorithms. It includes two different sets of phylogeny trees with 60 trees each, where every tree contains 20 nodes (audio files). The main difference between these two sets is in the set of transformations T.</p> <p><u>Type/format:</u> audio files (.wav) and phylogeny trees metadata as text (.txt)</p> <p><u>Re-use of existing data:</u> Yes. 3 existing audio files are used as roots for the processing.</p> <p><u>Data origin:</u> VCTK corpus version 0.92. Under Creative Commons License: Attribution 4.0 International</p> <p><u>Expected size:</u> 1.5 GB</p> <p><u>Data utility:</u> This dataset is useful in the context of T5.6 for the evaluation of algorithms for audio phylogeny analysis.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable. The dataset is hosted on Zenodo: <a href="https://zenodo.org/records/8135331">https://zenodo.org/records/8135331</a></p> <p><u>Search keywords:</u> Audio phylogeny dataset IDMT</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> For every audio phylogeny tree generated metadata is provided describing the process of generations by indicating parent-child couples and transformations applied.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The dataset is available on Zenodo at <a href="https://zenodo.org/records/8135331">https://zenodo.org/records/8135331</a></p> <p><u>How it will be accessible:</u> Downloadable from Zenodo</p> <p><u>Methods/software tools to access data:</u> Web-browser to download the data as zip file</p> <p><u>Repository:</u> Zenodo: <a href="https://zenodo.org/records/8135331">https://zenodo.org/records/8135331</a></p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p>



	<p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> Creative Commons Attribution 4.0 International License.</p> <p><u>Availability for re-use:</u> Yes</p> <p><u>Usable by third parties after end of project:</u> Yes, it's an open dataset</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> Zenodo has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> One half of the audio data originates from an open dataset of speech samples whose speakers have given a permission for data sharing and reuse. The other half of audio data is AI generated music with no ethical or licensing concerns.</p>
Other Issues	N/A

#### 4.4.21 RAI CMM documentaries dataset

DMP component	AI4Media_Data_32_WP5_Video_CMM-Doc-dataset_v1 Partner: RAI
Data Summary	<p><u>Purpose:</u> The CMM Doc dataset is a collection of selected documentaries with their segmentation from the Rai archive. The documentaries are used to experiment with the video segmentation techniques considered within the project and to compare the quality of these segmentations with the ones present in the archive (ground truth).</p> <p><u>Type/format:</u> mp4, csv, wav</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> Rai Archive</p> <p><u>Expected size:</u> 12GB, around 15 hours of audio/video content plus the csv files describing the segmentations.</p> <p><u>Data utility:</u> It is useful to evaluate the performance of different video segmentation approaches, comparing each of them with a ground truth.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No, data is stored on RAI's private servers.</p> <p><u>How it will be accessible:</u> N/A</p>



	<u>Methods/software tools to access data:</u> N/A <u>Repository:</u> N/A <u>Restrictions on access:</u> N/A
Making data interoperable	<u>Interoperability:</u> N/A <u>Data and metadata vocabularies:</u> N/A <u>Use of standard vocabularies:</u> N/A <u>Mappings to commonly used vocabularies:</u> N/A
Increase data re-use	<u>Licence:</u> N/A <u>Availability for re-use:</u> N/A <u>Usable by third parties after end of project:</u> N/A <u>Re-use timeframe:</u> N/A <u>Data quality assurance process:</u> N/A
Allocation of resources	<u>Costs for making data FAIR:</u> N/A <u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> N/A
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> Copyright owned by RAI <u>Is informed consent for data sharing and long-term preservation given:</u> N/A
Other Issues	N/A

#### 4.4.22 RAI CMM-ANTS newscasts dataset

DMP component	AI4Media_Data_33_WP5_Video_CMM-ANTS-dataset_v1 Partner: RAI
Data Summary	<p><u>Purpose:</u> The CMM-ANTS dataset is a collection of selected newscasts with their segmentation from the Rai archive. The programmes are used to experiment with the video segmentation techniques considered within the project and to compare the quality of these segmentations with the ones present in the archive (ground truth).</p> <p><u>Type/format:</u> mp4, csv, wav</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> Rai Archive</p> <p><u>Expected size:</u> 5GB, around 15 hours of audio/video content plus the csv files describing the segmentations.</p> <p><u>Data utility:</u> It is useful to evaluate the performance of different video segmentation approaches, comparing each of them with a ground truth.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p>



	<u>Metadata creation</u> : N/A
Making data openly accessible	<u>Data openly accessible</u> : No, data stored on RAI's private servers <u>How it will be accessible</u> : N/A <u>Methods/software tools to access data</u> : N/A <u>Repository</u> : N/A <u>Restrictions on access</u> : N/A
Making data interoperable	<u>Interoperability</u> : N/A <u>Data and metadata vocabularies</u> : N/A <u>Use of standard vocabularies</u> : N/A <u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<u>Licence</u> : N/A <u>Availability for re-use</u> : N/A <u>Usable by third parties after end of project</u> : N/A <u>Re-use timeframe</u> : N/A <u>Data quality assurance process</u> : N/A
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : N/A
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : N/A <u>Is informed consent for data sharing and long-term preservation given</u> : N/A
Other Issues	N/A

#### 4.4.23 RAI CMM mixed dataset

DMP component	AI4Media_Data_34_WP5_Video_CMM-dataset_v1 Partner: RAI
Data Summary	<p><u>Purpose</u>: CMMdataset is a collection of selected content of different kinds with their segmentation (if available) from the Rai archive. These contents are used to experiment with the video segmentation techniques considered within the project and to compare the quality of these segmentations with the ones present in the archive (ground truth).</p> <p><u>Type/format</u>: mp4, wav, csv</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Rai Archive</p> <p><u>Expected size</u>: 70GB, around 90 hours of audio/video content plus the csv files describing the segmentations.</p> <p><u>Data utility</u>: It is useful to evaluate the performance of different video segmentation approaches, comparing each of them with a ground truth.</p>



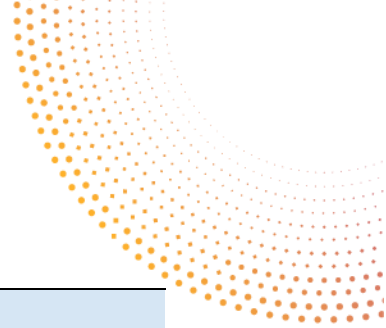


Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No, data stored on RAI's private servers</p> <p><u>How it will be accessible:</u> N/A</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> N/A</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long-term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.4.24 Cross lingual news dataset

<b>DMP component</b>	<b>AI4Media_Data_35_WP5_Text_Cross-lingual-news_v1</b> <b>Partner: RAI</b>
Data Summary	<p><u>Purpose:</u> The cross lingual news dataset is a collection of news items published by international online newspapers and press agencies. It is used by RAI and CNR to train and test CNR's Generalized Funneling component.</p> <p><u>Type/format:</u> json</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> News items from RSS feeds and press agencies</p> <p><u>Expected size:</u> 20 MB</p> <p><u>Data utility:</u> It is useful to WP5 partners to evaluate and benchmark text classification</p>





	models.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, this is a private dataset.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : N/A
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.4.25 YouTube RAI channel dataset

DMP component	AI4Media_Data_36_WP5_Misc_YouTubeDataset_v1 Partner: RAI
Data Summary	<p><u>Purpose</u>: A collection of URLs to RAI's YouTube channel content items + chapterisation of data.</p> <p><u>Type/format</u>: json</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: RAI YouTube channel</p> <p><u>Expected size</u>: 1.3 MB</p>



	<b>Data utility:</b> It is useful to WP5 partners to evaluate and benchmark media editorial segmentation.
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes</p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> Yes, provided by hosting platform</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> Yes</p> <p><b>How it will be accessible:</b> <a href="https://huggingface.co/datasets/raicrits/YouTube_RAI_dataset">https://huggingface.co/datasets/raicrits/YouTube_RAI_dataset</a></p> <p><b>Methods/software tools to access data:</b> git</p> <p><b>Repository:</b> <a href="https://huggingface.co/datasets/raicrits/YouTube_RAI_dataset">https://huggingface.co/datasets/raicrits/YouTube_RAI_dataset</a></p> <p><b>Restrictions on access:</b> N/A</p>
Making data interoperable	<p><b>Interoperability:</b> N/A</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> Apache 2.0</p> <p><b>Availability for re-use:</b> Yes</p> <p><b>Usable by third parties after end of project:</b> Yes</p> <p><b>Re-use timeframe:</b> N/A</p> <p><b>Data quality assurance process:</b> N/A</p>
Allocation of resources	<p><b>Costs for making data FAIR:</b> N/A</p> <p><b>Costs for long-term preservation:</b> N/A</p>
Data security	<b>Security measures:</b> N/A
Ethical aspects	<p><b>Possible ethical and legal aspects preventing sharing:</b> N/A</p> <p><b>Is informed consent for data sharing and long term preservation given:</b> N/A</p>
Other Issues	N/A

## 4.5 Datasets collected in the context of WP6

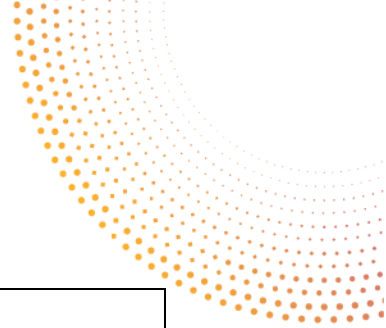
### 4.5.1 GreekPolitics Twitter dataset

<b>DMP component</b>	<b>AI4Media_Data_37_WP6_SOCIALMEDIA_GreekPolitics_v1</b> <b>Partner:</b> AUTH
Data Summary	<b>Purpose:</b> The AUTH GreekPolitics Dataset contains 2,578 tweet IDs from Twitter posts with politically charged content in the Greek language, spanning the period January 2014 – March 2021. Manually annotated ground-truth labels along 4 sentiment dimensions are provided for each tweet: polarity ('1' = positive, '0' = neutral, '-1' =



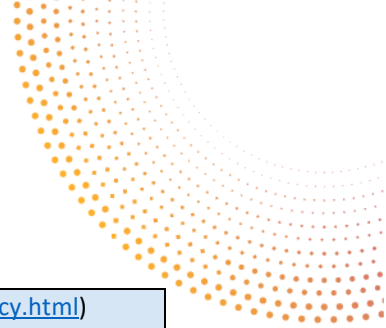
	<p>negative), figurativeness ('1' = figurative, '0' = literal), aggressiveness ('1' = aggressive, '0' = non-aggressive) and bias ('1' = partisan, '0' = non-partisan). This data is used in Task 6.4 "<i>AI for Healthier Political Debate</i>" to investigate political public opinion mining and/or monitoring. The dataset is used by AUTH for training and testing the new algorithms developed within T6.4.</p> <p><u>Type/format</u>: txt &amp; csv files</p> <p><u>Re-use of existing data</u>: We use data mined from Twitter using the official Twitter API.</p> <p><u>Data origin</u>: GreekPolitics Twitter posts were collected based on specific query hashtags related to the Greek political scene, using the official Twitter API. These hashtags are mainly related to the names of the various political parties and popular politicians represented in the Greek parliament over the past decade, while variants of them (in both the Greek and the Latin alphabet) were also exploited.</p> <p><u>Expected size</u>: ~ 10 GB</p> <p><u>Data utility</u>: The data is useful to WP6 partners that process Greek-language tweets. The data will also be useful to other social media researchers, but also to social or political scientists that study public opinion in Greece.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: The dataset is stored at an internal AUTH server and it is discoverable via AUTH's webpage <a href="https://aiia.csd.auth.gr/auth-greekpolitics-dataset/">https://aiia.csd.auth.gr/auth-greekpolitics-dataset/</a>.</p> <p>The data comprising the dataset (i.e. the tweets) are discoverable via the tweet ID; anyone can access the associated tweet text through the Twitter API. There is a chance that a subset of the GreekPolitics tweets may have been deleted, in the time interval that has passed since the dataset was originally constructed by AUTH.</p> <p><u>Search keywords</u>: AUTH GreekPolitics Dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: In order to access the GreekPolitics Dataset (<a href="https://aiia.csd.auth.gr/auth-greekpolitics-dataset/">https://aiia.csd.auth.gr/auth-greekpolitics-dataset/</a>), interested parties should complete and sign a license agreement: <a href="http://aiia.csd.auth.gr/wp-content/uploads/2023/10/AUTH_GreekPolitics_Dataset_License.pdf">http://aiia.csd.auth.gr/wp-content/uploads/2023/10/AUTH_GreekPolitics_Dataset_License.pdf</a>.</p> <p><u>How it will be accessible</u>: After completing the license agreement, interested parties receive FTP credentials for downloading the dataset from AUTH's servers.</p> <p><u>Methods/software tools to access data</u>: Download via FTP</p> <p><u>Repository</u>: AUTH servers</p> <p><u>Restrictions on access</u>: Only available upon signing a license agreement.</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: Using Twitter's standard data representation model.</p> <p><u>Data and metadata vocabularies</u>: The following metadata vocabulary is used:</p> <p style="padding-left: 40px;">created_at: UTC time when Tweet was created.</p> <p style="padding-left: 40px;">id: Unique identifier for Tweet.</p>





	<p>id_str: The string representation of Tweet unique identifier.</p> <p>text: UTF-8 text of the status update.</p> <p>source: Utility used to post the Tweet.</p> <p>truncated: Indicates whether the value of the text parameter was truncated.</p> <p>geo: The geo object of the status.</p> <p>coordinates: The coordinates of the status.</p> <p>in_reply_to_status_id: Original Tweet's ID if the Tweet is a reply.</p> <p>in_reply_to_status_id_str: Original Tweet's ID if the Tweet is a reply.</p> <p>in_reply_to_user_id: Tweet's author ID If the Tweet is a reply.</p> <p>in_reply_to_user_id_str: Tweet's author ID If the Tweet is a reply.</p> <p>in_reply_to_screen_name: Screen name of original Tweet's author if the Tweet is a reply.</p> <p>user: The user who posted this Tweet.</p> <p>coordinates: Geographic location of Tweet.</p> <p>place: Tweet is associated with a place.</p> <p>quoted_status_id: Tweet ID of the quoted Tweet.</p> <p>quoted_status_id_str: Tweet ID of the quoted Tweet.</p> <p>is_quote_status: Quoted Tweet indicator.</p> <p>quoted_status: Tweet object of the original Tweet that was quoted.</p> <p>retweeted_status: Original Tweet that was retweeted.</p> <p>quote_count: Number of quotes.</p> <p>reply_count: Number of replies.</p> <p>retweet_count: Number of retweets.</p> <p>favorite_count: Number of likes.</p> <p>geo: Contains place details in GeoJSON format.</p> <p>place_type: Specified the particular type of information represented by this place information, such as a city name, or a point of interest.</p> <p><u>Use of standard vocabularies:</u> Standard vocabularies as used in the Twitter APIs.</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> Data is shared after interested parties should and sign a license agreement: <a href="http://aiia.csd.auth.gr/wp-content/uploads/2023/10/AUTH_GreekPolitics_Dataset_License.pdf">http://aiia.csd.auth.gr/wp-content/uploads/2023/10/AUTH_GreekPolitics_Dataset_License.pdf</a></p> <p>Twitter data is used by complying to Twitter's Developer Agreement and Policy</p>



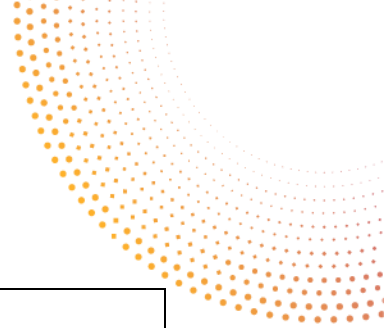


	<p><a href="https://developer.twitter.com/en/developer-terms/agreement-and-policy.html">https://developer.twitter.com/en/developer-terms/agreement-and-policy.html</a></p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> Yes</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Data quality assurance is ensured through a data pre-processing step, including data cleaning and harmonization as part of data crawling.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The data is stored at an internal AUTH server. Appropriate security measures are implemented and compliance with GDPR is considered.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Twitter data cannot be redistributed directly because of the terms of service of Twitter APIs. We only share the Tweet IDs.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	No

#### 4.5.2 Covid-19 Twitter dataset

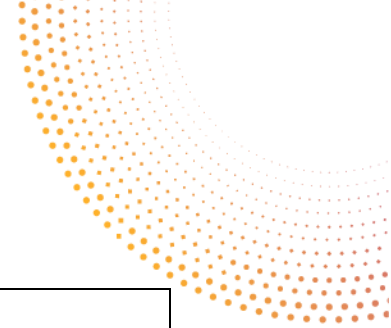
DMP component	AI4Media_Data_38_WP6_SOCIALMEDIA_Covid19Twitter_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> The dataset includes COVID-19 related tweets. The tweets were collected in real time using Twitter’s stream API and a set of appropriate keywords and hashtags. The collection spans the period from March 2020 to March 2021. This data is used in Task 6.4 “AI for Healthier Political Debate” mainly to research the properties of COVID-19 related discussions and healthy discussions on social media in general. The dataset is used to train and test the new algorithms developed within T6.4.</p> <p><u>Type/format:</u> MongoDB database composed of JSON files (one file per tweet)</p> <p><u>Re-use of existing data:</u> The original data was stored by Twitter. This dataset was originally collected for several research projects including AI4Media.</p> <p><u>Data origin:</u> Twitter API</p> <p><u>Expected size:</u> Compressed: ~2 TB, uncompressed ~ 6.5 TB</p> <p><u>Data utility:</u> The data will be used by T6.4 partners for research on topics including, but not limited to discussion healthiness estimation, argumentation mining and analysis, deep fake video analysis, long text analysis and sentiment analysis. Dataset is available to other partners within AI4Media to conduct research in other topics as well.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The tweets within the dataset are discoverable via twitter.com and via twitter API. The dataset is stored on the BSC servers and will be available to AI4Media partners on demand. The dataset will not be made discoverable outside of the consortium.</p>





	<p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> The data is in the raw format provided by the Twitter Streaming API, with no modifications nor metadata added.</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible:</u> The dataset will not be openly accessible. Redistribution of data (Twitter posts) is against the terms of service of Twitter APIs.</p> <p>In case of a report or paper submitted for publication, all research findings will be integrated into the report or paper. Datasets will not be added to the publication.</p> <p><u>How it will be accessible:</u> N/A</p> <p><u>Methods/software tools to access data:</u> The data is accessed via MongoDB API, either directly via secure connections or from code (for instance, via pymongo package for Python)</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> The access is provided to AI4Media partners on demand.</p>
<p>Making data interoperable</p>	<p><u>Interoperability:</u> Using Twitter’s standard data representation model.</p> <p><u>Data and metadata vocabularies:</u> The following metadata vocabulary is used:</p> <ul style="list-style-type: none"> <li><i>_id:</i> MongoDB unique identifier for the Tweet.</li> <li><i>created_at:</i> UTC time when Tweet was created.</li> <li><i>timestamp_ms:</i> Unix timestamp in miliseconds.</li> <li><i>id:</i> Unique identifier for Tweet.</li> <li><i>id_str:</i> The string representation of Tweet unique identifier.</li> <li><i>text:</i> UTF-8 text of the status update.</li> <li><i>source:</i> Utility used to post the Tweet.</li> <li><i>truncated:</i> Indicates whether the value of the text parameter was truncated.</li> <li><i>in_reply_to_status_id:</i> Original Tweet’s ID if the Tweet is a reply</li> <li><i>in_reply_to_status_id_str:</i> Original Tweet’s ID if the Tweet is a reply</li> <li><i>in_reply_to_user_id:</i> Tweet’s author ID If the Tweet is a reply</li> <li><i>in_reply_to_user_id_str:</i> Tweet’s author ID If the Tweet is a reply</li> <li><i>in_reply_to_screen_name:</i> Screen name of original Tweet’s author if the Tweet is a reply</li> <li><i>user:</i> The user who posted this Tweet.</li> <li><i>coordinates:</i> Geographic location of Tweet</li> <li><i>place:</i> Tweet is associated with a place</li> <li><i>quoted_status_id:</i> Tweet ID of the quoted Tweet.</li> <li><i>quoted_status_id_str:</i> Tweet ID of the quoted Tweet.</li> </ul>

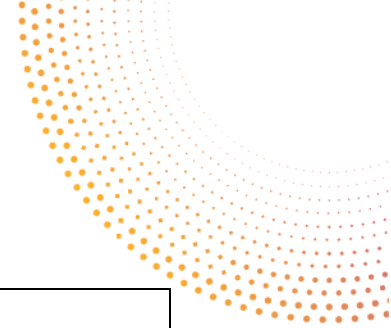




	<p><i>is_quote_status</i>: Quoted Tweet indicator</p> <p><i>quoted_status</i>: Tweet object of the original Tweet that was quoted.</p> <p><i>retweeted_status</i>: Original Tweet that was retweeted.</p> <p><i>quote_count</i>: Number of quotes</p> <p><i>reply_count</i>: Number of replies</p> <p><i>retweet_count</i>: Number of retweets</p> <p><i>favorite_count</i>: Number of likes</p> <p><i>entities</i>: Entities which have been parsed out of the text of the Tweet</p> <p><i>extended_entities</i>: Entities which have been parsed out of the text of the Tweet if there is media content</p> <p><i>favorited</i>: Indicates whether this Tweet has been liked by the authenticating user</p> <p><i>retweeted</i>: Indicates whether this Tweet has been Retweeted by the authenticating user</p> <p><i>possibly_sensitive</i>: URL contained in the Tweet may contain content or media identified as sensitive</p> <p><i>filter_level</i>: The maximum value of the filter_level parameter to still stream this Tweet</p> <p><i>lang</i>: Machine-detected language of the Tweet text</p> <p><i>matching_rules</i>: Provides the id and tag associated with the rule that matched the Tweet</p> <p><i>geo</i>: [deprecated] Coordinates in [lat, long] format</p> <p><u>Use of standard vocabularies</u>: Standard vocabularies as used in the Twitter APIs, plus 2 additional database fields.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will not be shared. Twitter data is used by complying to Twitter’s Developer Agreement and Policy (<a href="https://developer.twitter.com/en/developer-terms/agreement-and-policy.html">https://developer.twitter.com/en/developer-terms/agreement-and-policy.html</a>)</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Twitter data is collected as originally provided by Twitter API. Additional post-processing of data is left for the researchers.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data will be stored in BSC servers. Access will be allowed only via ad hoc credentials created for each user.</p>
Ethical	<p><u>Possible ethical and legal aspects preventing sharing</u>: Twitter data cannot be</p>





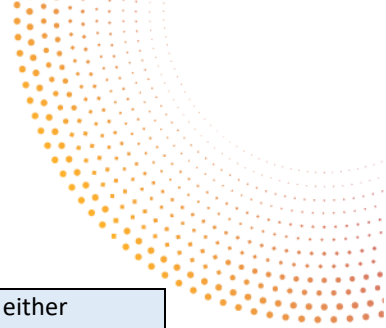


aspects	redistributed because of the terms of service of Twitter APIs. <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 4.5.3 Twitter Text dataset

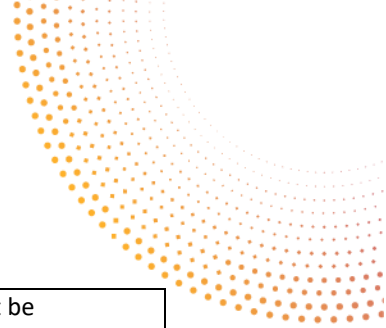
DMP component	AI4Media_Data_39_WP6_TEXT_TextTwitter_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> The dataset contains reduced versions of COVID-19 related tweets (mostly the text field and attachments) from the COVID-19 Twitter dataset described in section 4.5.2. It uses the Solr database functionality to provide efficient text search throughout the whole collection. The tweets were collected in real time using Twitter’s stream API and a set of appropriate keywords and hashtags. Twets were collected from March 2020 to March 2021. These tweets were than parsed, relevant fields processed and extracted, and then added to the Solr DB. This data is used in task 6.4 “AI for Healthier Political Debate” mainly to research the properties of COVID-19 related discussions and healthy discussions on social media in general. The dataset is used to train and test the new algorithms developed within T6.4.</p> <p><u>Type/format:</u> Solr database</p> <p><u>Re-use of existing data:</u> The <a href="#">original tweets were stored by Twitter</a>. This dataset is based on the dataset described in section 4.5.2, originally collected by BSC.</p> <p><u>Data origin:</u> Twitter API</p> <p><u>Expected size:</u> ~ 160 GB</p> <p><u>Data utility:</u> The data will be used by T6.4 partners for research on topics including, but not limited to discussion healthiness estimation, argumentation mining and analysis, deep fake video analysis, long text analysis and sentiment analysis. Dataset will be available to other partners within AI4Media to conduct research in other topics as well.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The tweets within the dataset are discoverable via twitter.com and via twitter API. The dataset is stored on the BSC servers and will be available to AI4Media partners on demand. The dataset will not be made discoverable outside of the consortium.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> The data is in the raw format provided by the Twitter Streaming API, with no modifications nor metadata added.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The dataset will not be openly accessible. Redistribution of data (Twitter posts) is against the terms of service of Twitter APIs.</p> <p>In case of a report or paper submitted for publication, all research findings will be integrated into the report or paper. Datasets will not be added to the publication.</p> <p><u>How it will be accessible:</u> N/A</p>





	<p><u>Methods/software tools to access data</u>: The data is accessed via Solr API, either directly via secure connections or from code (for instance, via pymongo package for Python)</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: The access is provided to AI4Media partners on demand.</p>
Making data interoperable	<p><u>Interoperability</u>: If the same data schemas are used, different datasets can be merged.</p> <p><u>Data and metadata vocabularies</u>: The following metadata vocabulary is used:</p> <p><i>id</i>: Unique Tweet id as provided by Twitter</p> <p><i>timestamp_ms</i>: Unix timestamp (in milliseconds) of when the Tweet was created_at</p> <p><i>text</i>: Text of the Tweet. If the Tweet is a retweet or quote, text of the original tweet is appended to this field</p> <p><i>lang</i>: Machine-detected language of the Tweet text (provided by Twitter)</p> <p><i>attachment_photos</i>: JSON list of links to the photos attached to the Tweet</p> <p><i>attachment_videos</i>: JSON list of links and metadata of the videos attached to the Tweet, for versions of different quality</p> <p><i>attachment_texts</i>: Tesseract character recognition algorithms was applied to photos attached to the Tweet. If more than 40 symbols is extracted, the extracted text is contained in this field (and is fully searchable)</p> <p><i>attachment_gifs</i>: JSON list of links and metadata of the GIFs attached to the Tweet, as videos of different quality</p> <p><i>attachment_urls</i>: If the Tweet contained a URL link to anything except a Tweet, it is contained here.</p> <p><u>Use of standard vocabularies</u>: When applicable, Twitter API field names were used.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will not be shared. Twitter data is used by complying to Twitter's Developer Agreement and Policy (<a href="https://developer.twitter.com/en/developer-terms/agreement-and-policy.html">https://developer.twitter.com/en/developer-terms/agreement-and-policy.html</a>)</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Twitter data is collected as originally provided by Twitter API. Text extraction from photos is performed via Tesseract OCR, a top-end character recognition and extraction module.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data will be stored in BSC servers. Access will be allowed only via ad hoc credentials created for each user.</p>



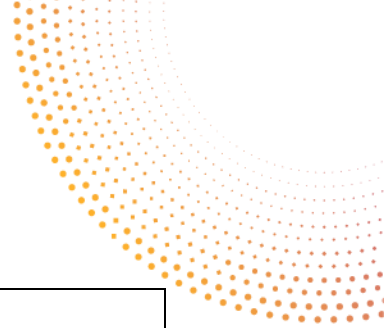


Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Twitter data cannot be redistributed because of the terms of service of Twitter APIs.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.5.4 Twitter COVID-19 discussions topics dataset

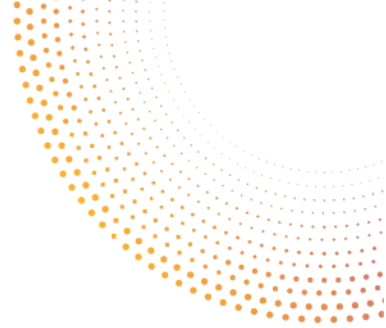
DMP component	AI4Media_Data_40_WP6_TEXT_Covid19DiscussionTopicsTwitter_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> Dataset of tweet IDs and timestamps in CSV format, organized by COVID-19-related discussion topics they belong to. It contains a total of 3,423,260 entries between 10 topics and covers 2 time periods: March to August 2020 and August 2020 to March 2021. The tweet IDs correspond to the tweets included in the COVID-19 Twitter dataset described in section 4.5.2 The topics include:</p> <ol style="list-style-type: none"> <li>1. Azithromycin and covid</li> <li>2. Bleach and covid</li> <li>3. 5G and covid</li> <li>4. Covid is artificial</li> <li>5. Amrmy mobilizations during covid</li> <li>6. Invermictin and covid</li> <li>7. Moderna vaccine</li> <li>8. Astrazeneca vaccine</li> <li>9. Blood types and covid susceptibility</li> <li>10. H1N1 and covid</li> </ol> <p>All timestamps are in UNIX seconds.</p> <p>This data is used in task 6.4 “AI for Healthier Political Debate” mainly to research the properties of COVID-19 related discussions and healthy discussions on social media in general. The dataset is used to train and test the new algorithms developed within T6.4.</p> <p><u>Type/format:</u> CSV</p> <p><u>Re-use of existing data:</u> The original tweets were stored by Twitter. This dataset includes tweet IDs and timestamps and is based on the dataset described in section 4.5.2, originally collected by BSC.</p> <p><u>Data origin:</u> Twitter API</p> <p><u>Expected size:</u> ~ 100 MB</p> <p><u>Data utility:</u> The data is used by T6.4 partners for research on topics including, but not limited to discussion healthiness estimation, argumentation mining and analysis, deep fake video analysis, long text analysis and sentiment analysis. Dataset will be available to other researchers to perform research in similar topics.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The dataset is discoverable through the AI4Media website (<a href="https://www.ai4media.eu/open-datasets/">https://www.ai4media.eu/open-datasets/</a>), which includes a link to the data repository (<a href="https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8">https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8</a>)</p> <p>The original tweets are discoverable via twitter.com and via twitter API using the</p>





	<p>tweet IDs included in the dataset</p> <p><u>Search keywords:</u> AI4Media Twitter Discussions Topics dataset</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Dataset available on the cloud at <a href="https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8">https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8</a></p> <p><u>How it will be accessible:</u> Access via web-browser - direct download of data</p> <p><u>Methods/software tools to access data:</u> Web browser</p> <p><u>Repository:</u> <a href="https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8">https://docs.hpai.cloud/s/A6EwRCrswm3B5Q8</a></p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> Use of Twitter IDs and UNIX timestamps</p> <p><u>Data and metadata vocabularies:</u> The following metadata vocabulary is used:</p> <p style="padding-left: 40px;"><i>id: Tweet ID</i></p> <p style="padding-left: 40px;"><i>timestamp_ms:</i> Unix timestamp (in milliseconds) of when the Tweet was created at</p> <p><u>Use of standard vocabularies:</u> Use of Twitter IDs and UNIX timestamps</p> <p><u>Mappings to commonly used vocabularies:</u> Use of Twitter IDs and UNIX timestamps</p>
Increase data re-use	<p><u>Licence:</u> MIT license</p> <p><u>Availability for re-use:</u> Yes</p> <p><u>Usable by third parties after end of project:</u> Yes</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Yes</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> Data stored on BSC-HPAI's repo on <a href="#">Nextcloud Hub</a>. Nextcloud features a host of innovative security technologies from brute force protection to advanced server side and integrated end-to-end, client side encryption with enterprise-grade key handling and a wide range of security hardenings.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Twitter data cannot be redistributed because of the terms of service of Twitter APIs. Thus we won't share tweetIDs</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

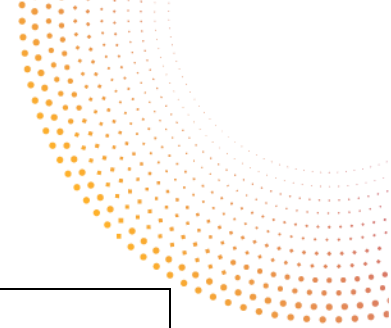




#### 4.5.5 ODSS Open Dataset of Synthetic Speech

DMP component	AI4Media_Data_41_WP6_AUDIO_ODSS_v1 Partner: FhG-IDMT, CERTH
Data Summary	<p><u>Purpose</u>: ODSS is a multilingual, multispeaker dataset of synthetic and natural speech, designed to foster research and benchmarking of novel studies on synthetic speech detection. ODSS comprises audio utterances generated from text by state-of-the-art synthesis methods, paired with their corresponding natural counterparts. The synthetic audio data includes several languages, with an equal representation of genders. This dataset is required to learn the statistical distribution of features derived from synthetic speech data – to be distinguished, e.g., from pristine speech data.</p> <p><u>Type/format</u>: Audio (wav and csv files)</p> <p><u>Re-use of existing data</u>: Yes, re-utterances and transcripts of speech corpora of pristine speech.</p> <p><u>Data origin</u>: The data was generated using Text-to-Speech synthesis methods.</p> <p><u>Expected size</u>: 2.4GB</p> <p><u>Data utility</u>: Useful in the context of T6.2 for developing and testing new algorithms for content manipulation/synthesis detection.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, data is hosted on Zenodo and has a DOI reference (<a href="https://doi.org/10.5281/zenodo.8370668">https://doi.org/10.5281/zenodo.8370668</a>)</p> <p><u>Search keywords</u>: ODSS, Open Dataset of Synthetic Speech</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: Metadata is provided in CSV format, which has the advantage of being both machine-readable and human-readable.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, data available on Zenodo at <a href="https://zenodo.org/records/8370669">https://zenodo.org/records/8370669</a></p> <p><u>How it will be accessible</u>: Downloadable via zenodo.org.</p> <p><u>Methods/software tools to access data</u>: Any software supporting WAV and CSV files.</p> <p><u>Repository</u>: Zenodo (<a href="https://zenodo.org/records/8370669">https://zenodo.org/records/8370669</a>)</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Audio files are present in WAV format, which is largely interoperable</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: CC-BY-SA v4.0</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: Undetermined</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of	<p><u>Costs for making data FAIR</u>: N/A</p>



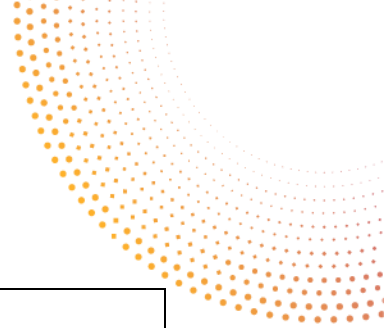


resources	<u>Costs for long-term preservation:</u> None
Data security	<u>Security measures:</u> N/A
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Requirements of the CC-BY-SA 4.0 licence need to be fulfilled.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> All recordings used to train the neural networks involved in the data generation, as well as all pristine recordings included in the dataset were acquired in controlled conditions and with the informed consent of all the speakers involved, in compliance with the GDPR.</p>
Other Issues	N/A

#### 4.5.6 CelebHQGaze image dataset for gaze estimation

DMP component	AI4Media_Data_42_WP6_IMAGE_CelebHQGaze_v1 Partner: UNITN
Data Summary	<p><u>Purpose:</u> Common gaze-correction methods usually require annotating training data with precise gaze, and head pose information. Solving this problem using an unsupervised method remains an open problem, especially for high-resolution face images in the wild, which are not easy to annotate with gaze and head pose labels. To address this issue, we created a high-resolution CelebHQGaze dataset of faces with corrected gazes. CelebHQGaze consists of 29,255 high resolution celebrity images that are collected from CelebAHQ. It consists of 21,005 face images with the eyes staring at the camera and 8,250 face images with eyes looking somewhere else. Similarly to CelebGaze, we extract facial landmarks and generate the mask. All images are cropped to 512 × 512, and the mask size is fixed to 46×80. Useful to WP6 for gaze correction applications.</p> <p><u>Type/format:</u> JPEG images and accompanying metadata in json format.</p> <p><u>Re-use of existing data:</u> No, the dataset is created within AI4Media.</p> <p><u>Data origin:</u> <a href="#">CelebA dataset</a></p> <p><u>Expected size:</u> ~3 GB</p> <p><u>Data utility:</u> It is useful to WP6 partners (especially T6.2) to correct the gaze of portraits.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes, it is discoverable via GitHub and the project website</p> <p><u>Search keywords:</u> CelebHQGaze</p> <p><u>Versioning:</u> GitHub supports versioning</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is openly accessible via GitHub at <a href="https://github.com/zhangqianhui/GazeAnimationV2">https://github.com/zhangqianhui/GazeAnimationV2</a></p> <p><u>How it will be accessible:</u> The data can be downloaded from an online archive after completing a form.</p> <p><u>Methods/software tools to access data:</u> Web browser</p> <p><u>Repository:</u> GitHub</p> <p><u>Restrictions on access:</u> The user should accept the terms of use.</p>



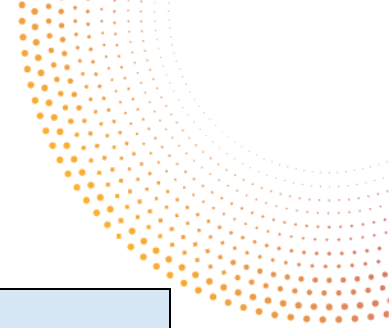


Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: The repository consists of two datasets: CelebGaze and CelebHQGaze. CelebGaze consists of 25,283 celebrity images, most of which have been collected from CelebA and a minority from the Internet. Specifically, there are 21,832 face images with the eyes staring at the camera and 3,451 face images with the eyes looking somewhere else. We cropped all the images to 256 x 256 and computed the eye mask region using Dlib. Specifically, we used Dlib to extract 68 facial landmarks, and we computed the mean of 6 points near the eye region, which is the center point of the mask. The size of the mask is fixed to 30 x 50. CelebHQGaze consists of 29,255 high resolution celebrity images that are collected from CelebAHQ. It consists of 21,005 face images with the eyes staring at the camera and 8,250 face images with eyes looking somewhere else. Similarly to CelebGaze, we extract facial landmarks and generate the mask. All images are cropped to 512 x 512, and the mask size is fixed to 46 x 80.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data and the code are released under the CC license.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The full dataset is shared in a google drive repository.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.5.7 European news about Covid vaccination dataset

DMP component	AI4Media_Data_43_WP6_TEXT_European-news-covid-vaccination_v1 Partner: IDIAP
Data Summary	<p><u>Purpose</u>: The dataset consists of 50,000+ full newspaper articles in five countries (France, Italy, Spain, Switzerland, and UK), and four languages (French, Spanish, Italian, and English) over a period of 22 months, as well as a translated English version for the full dataset. The common theme of the articles is about Covid vaccination and the objective of this dataset was to study the local component and the peculiarities of each country addressing the common theme of vaccination. From this original <a href="#">dataset</a>, a subsample of articles focused on the No-Vax movement (1,786 articles) has been used to carry out two additional lines of research, one concerning the identification of <a href="#">frames</a> in these articles and the other concerning the study of quality articles on the subject of <a href="#">disinformation</a>. The dataset is used in T6.5 for the analysis of hyper-local news.</p>

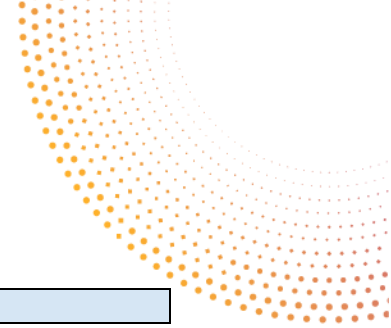




	<p><u>Type/format</u>: csv file</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: All the data is collected from the 19 webpages of the newspapers.</p> <p><u>Expected size</u>: 460 MB</p> <p><u>Data utility</u>: The dataset contains articles from European newspapers with a great tradition and reputation, with the content in their original language, as well as in English. This dataset can be of interest for research lines such as disinformation detection since this content could be labelled as truthful, as well as for political analysis of how information is presented by different ideologies or countries.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on IDIAP's servers, and can only be accessed by authorized users.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: We have not obtained permission from the newspapers to share the data. We have obtained authorization only for data extraction and use for research purposes.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>





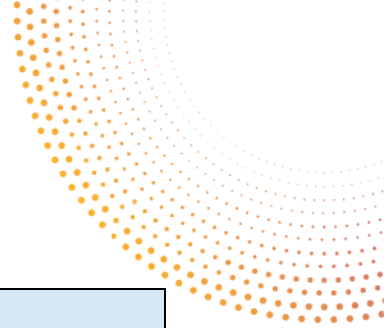


Other Issues	N/A
--------------	-----

#### 4.5.8 Suisse Romande local news dataset

DMP component	AI4Media_Data_44_WP6_TEXT_SuisseRomandeLocalNews_v1 Partner: IDIAP
Data Summary	<p><b>Purpose:</b> The dataset contains 130,155 articles sourced from the websites of three Swiss francophone newspapers: Arc Info, La Cote, and Le Nouvelliste, spanning the time period from 01/01/2015 to 30/06/2022. The dataset, compiled from the temporary data feeds provided by the press agency, consists of the articles in their entirety including the title, headline, and content, along with metadata for each article. The collected articles are primarily in French language and are categorically sorted by topics and region, everything encoded in the JSON format. The dataset is used in T6.5 for the analysis of hyper-local news.</p> <p><b>Type/format:</b> JSON</p> <p><b>Re-use of existing data:</b> No</p> <p><b>Data origin:</b> Local newspapers websites (Arc Info, La Cote, and Le Nouvelliste)</p> <p><b>Expected size:</b> 1 GB</p> <p><b>Data utility:</b> Used in T6.5 for the analysis of hyper-local news. Useful to researchers studying local news press.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes, with a unique identifier: <a href="https://zenodo.org/records/10256911">https://zenodo.org/records/10256911</a></p> <p><b>Search keywords:</b> Suisse Romande local news</p> <p><b>Versioning:</b> Only one release version</p> <p><b>Metadata creation:</b> The metadata are the ones available on the websites associated to their respective articles</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The dataset is accessible via Zenodo upon request (the authors' names are anonymized).</p> <p><b>How it will be accessible:</b> From Zenodo at <a href="https://zenodo.org/records/10256911">https://zenodo.org/records/10256911</a></p> <p><b>Methods/software tools to access data:</b> Only need for a json reader, the documentation explains the process in Python.</p> <p><b>Repository:</b> Available on Zenodo</p> <p><b>Restrictions on access:</b> Access upon request</p>
Making data interoperable	<p><b>Interoperability:</b> The identifier and the metadata provide unique sets of attributes for each article, which can be identified this way.</p> <p><b>Data and metadata vocabularies:</b> Date and newspapers.</p> <p><b>Use of standard vocabularies:</b> No use of standard vocabularies</p> <p><b>Mappings to commonly used vocabularies:</b> The mapping is explained in the documentation.</p>
Increase data re-use	<p><b>Licence:</b> Worldwide distribution is restricted to academic institutions / non profit</p>





	<p>research institutions.</p> <p><u>Availability for re-use</u>: Yes, no embargo.</p> <p><u>Usable by third parties after end of project</u>: Available for research purposes</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality assessment + back-up system</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: Marginal (small dataset without updates)</p>
Data security	<p><u>Security measures</u>: Accessible only upon request. Dataset stored in IDIAP's servers and all appropriate security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The press agency authorized the use for research purposes only.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: The consent was given by the owner of the rights to the articles</p>
Other Issues	N/A

#### 4.5.9 Lausanne news dataset

DMP component	AI4Media_Data_45_WP6_TEXT_LausanneNews_v1 Partner: IDIAP
Data Summary	<p><u>Purpose</u>: Dataset of 2,666 articles scraped from the newspaper Lausanne Cité, all published in French, from the time period between January 1, 2019 and November 30, 2022. The dataset is used in T6.5 for the analysis of hyper-local news.</p> <p><u>Type/format</u>: pickle</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: All the data is collected from the webpage of the newspaper Lausanne Cité</p> <p><u>Expected size</u>: 8.6 MB</p> <p><u>Data utility</u>: Used in T6.5 for the analysis of hyper-local news. Useful to researchers studying local news press.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p>

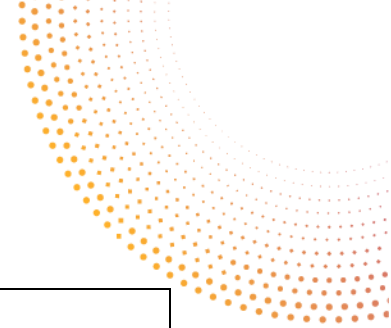


	<u>Restrictions on access</u> : N/A
Making data interoperable	<u>Interoperability</u> : N/A <u>Data and metadata vocabularies</u> : N/A <u>Use of standard vocabularies</u> : N/A <u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<u>Licence</u> : N/A <u>Availability for re-use</u> : N/A <u>Usable by third parties after end of project</u> : N/A <u>Re-use timeframe</u> : N/A <u>Data quality assurance process</u> : N/A
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : The data is stored on IDIAP's servers, and can only be accessed by authorized users.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : We have not obtained permission from the newspaper to share the data. We have obtained authorization only for data extraction and use for research purposes by IDIAP. <u>Is informed consent for data sharing and long term preservation given</u> : N/A
Other Issues	N/A

#### 4.5.10 Suisse Allemande local news dataset

DMP component	AI4Media_Data_46_WP6_TEXT_SuisseAllemandeNews_v1 Partner: IDIAP
Data Summary	<p><u>Purpose</u>: The dataset consists of around 14,500 full articles from 2 local newspapers in the canton of Zurich, Tagblatt der Stadt Zürich and Winterthurer Zeitung, with all the content in German. These are articles between the years 2012 and 2022. The dataset is used in T6.5 for the analysis of hyper-local news.</p> <p><u>Type/format</u>: csv file</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: All the data is collected from the webpages of the newspapers Tagblatt der Stadt Zürich and Winterthurer Zeitung.</p> <p><u>Expected size</u>: Few MBs</p> <p><u>Data utility</u>: Used in T6.5 for the analysis of hyper-local news. Useful to researchers studying local news press.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No</p>



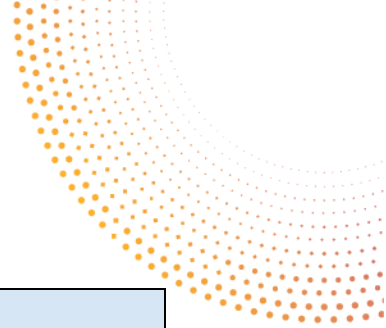


	<u>Metadata creation</u> : No
Making data openly accessible	<p><u>Data openly accessible</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The data is stored on IDIAP's servers, and can only be accessed by authorized users.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: We have not obtained permission from the newspapers to share the data. We have obtained authorization only for data extraction and use for research purposes by IDIAP.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.5.11 References on YouTube dataset

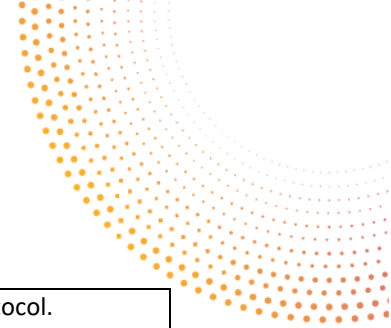
DMP component	AI4Media_Data_47_WP6_TEXT_YouTubeReferences_v1 Partner: IDIAP
Data Summary	<p><u>Purpose</u>: The dataset describes how YouTube knowledge communications videos put and use references to specify external information sources that are used to create the videos. 44 English-speaking YouTube channels that specialize in knowledge communication were manually chosen, and one video per channel was chosen for the analysis. For each video, the referencing scheme (<i>in-video</i> and <i>bibliography</i>) and general context about the channel/creators (visual style, upload date, the number of subscribers, topic, details of the creator group, running time of the video) were summarized and categorized. Additionally, a sample of 129 references was collected to verify whether the resources were actually available or not. The dataset is used in T6.5 for the analysis of referencing practices in YouTube.</p>





	<p><u>Type/format</u>: csv file</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: The channel list was collected based on past literature that listed a selection of knowledge communication channels on YouTube and extended by using the recommendation of the YouTube platform. The video and reference list items were arbitrarily selected from each channel.</p> <p><u>Expected size</u>: 1 MB</p> <p><u>Data utility</u>: To the best of our knowledge, this dataset is the first systematic observation of referencing practices in YouTube. By linking the accessibility and availability of external references included in online videos, we can extend the research questions to automatic validation of video content with regard to references, or towards the fundamental debate about information accessibility and information quality.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is not openly accessible. There are no plans for data sharing.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on IDIAP's servers, and can only be accessed by authorized users.</p>
Ethical	<p><u>Possible ethical and legal aspects preventing sharing</u>: The data is not publicly available</p>



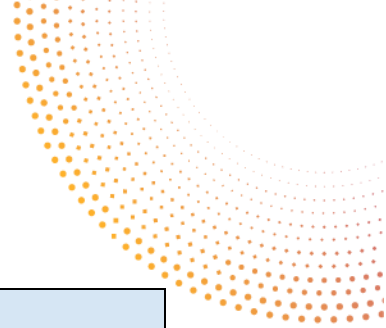


aspects	as that was not requested in the Institutional Ethical approval of the protocol. <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

#### 4.5.12 Political Barometer dataset

DMP component	AI4Media_Data_48_WP6_TEXT_PoliticalBarometer_v1 Partner: AUTH
Data Summary	<p><u>Purpose:</u> This Dataset contains text from Twitter posts with politically charged content in Greek language to be used in the Political Barometer software (<a href="https://icarus.csd.auth.gr/political-barometer/">https://icarus.csd.auth.gr/political-barometer/</a>) to predict the outcome of the Greek elections of June 2024. Due to Twitter's change in API access policy and limitations, this dataset remains private. Twitter IDs and URLs weren't stored in a consistent manner and format. The dataset is used for the research activities of T6.4, focusing on online political debate analysis.</p> <p><u>Type/format:</u> Text</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> Tweets in Greek language collected from Twitter, spanning the period June 2022 - March 2024.</p> <p><u>Expected size:</u> 0.25 GB</p> <p><u>Data utility:</u> The data is useful for text and sentiment analysis and analysis of online political debate.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> No</p> <p><u>Metadata creation:</u> No</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No</p> <p><u>How it will be accessible:</u> N/A</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> No</p>



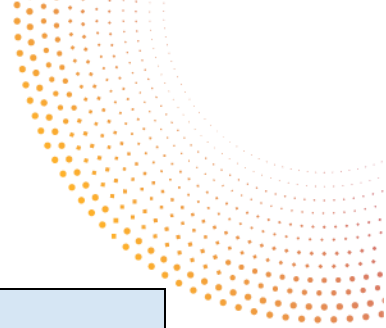


	<p><u>Availability for re-use</u>: No</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data is stored on the Intranet of the AIIA Laboratory (AUTH). Appropriate security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Due to Twitter’s change in API access policy and limitations, this dataset remains cannot be shared. Vaild identifiers to make data eligible for sharing weren’t store in a consistent manner.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: N/A</p>
Other Issues	N/A

#### 4.5.13 ElecDeb60To16-fallacy dataset

DMP component	AI4Media_Data_49_WP6_TEXT_ElecDeb60To16_v1 Partner: UCA
Data Summary	<p><u>Purpose</u>: The ElecDeb60To16-fallacy dataset collects televised debates of the presidential election campaigns in the U.S. from 1960 to 2016. It contains political debates annotated with argument components (evidence, claim), relations (support, attack), 6 fallacy types, and 14 sub-fallacy types. This data is meant to be used in Task 6.4 “AI for Healthier Political Debate” to aid in the construction of text machine-learning models that classify fallacious arguments.</p> <p><u>Type/format</u>: csv files</p> <p><u>Re-use of existing data</u>: We extended the ElecDeb60To16 dataset by adding the annotations for fallacious arguments.</p> <p><u>Data origin</u>: The political debates' text was collected from the website of the Commission on Presidential Debates, which provided transcripts of the debates broadcasted on TV and held among the leading candidates for the presidential and vice presidential nominations in the US.</p> <p><u>Expected size</u>: 25 MB</p> <p><u>Data utility</u>: Useful in the context of T6.4 for political debate analysis. Also, useful to any researcher working on falacious argument analysis and debate analysis.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is publicly available and discoverable from the project website and GitHub, but no metadata file can be automatically parsed.</p> <p><u>Search keywords</u>: ElecDeb60To16-fallacy dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data	<p><u>Data openly accessible</u>: The dataset is publicly available on GitHub at</p>





openly accessible	<p><a href="https://github.com/pierpaologoffredo/IJCAI2022/tree/main/dataset">https://github.com/pierpaologoffredo/IJCAI2022/tree/main/dataset</a></p> <p><u>How it will be accessible:</u> Downloadable via GitHub</p> <p><u>Methods/software tools to access data:</u> Given that the data's format is csv, any spreadsheet software or text editor can be used to access it.</p> <p><u>Repository:</u> GitHub (<a href="https://github.com/pierpaologoffredo/IJCAI2022">https://github.com/pierpaologoffredo/IJCAI2022</a>)</p> <p><u>Restrictions on access:</u> No</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International</p> <p><u>Availability for re-use:</u> The data is available for re-use</p> <p><u>Usable by third parties after end of project:</u> The dataset is publicly available for non-commercial purposes</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Inter-annotator agreement analysis</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> Data is publicly available on GitHub. GitHub has appropriate security mechanisms</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> No</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

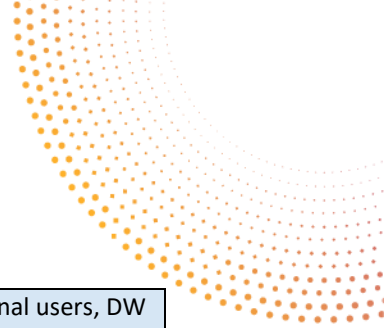
## 4.6 Datasets collected in the context of WP8

### 4.6.1 Data from user research activities in Use Case 1

DMP component	AI4Media_Data_50_WP8_USER-RESEARCH_UseCase1_DW_v1 Partner: DW
Data Summary	<p><u>Purpose:</u> In order to realise Use Case 1 in WP8, partner DW has designed various research activities with <i>professional</i> users, who work at DW or fact-checking organisations (journalists, verification experts or managers). This research is required for the definition of user requirements (D8.1), further user requirements gathered for the first White Paper (D8.3) as well as the user evaluation of new functions integrated into the Truly Media demonstrator and related trustworthy AI elements. The following, mainly qualitative research methods, have been conducted with small, selected target groups: tailored questionnaire-based surveys, interviews, focus groups or demo/evaluation sessions with single users. User research of this kind identifies the users' needs and their opinion regarding new functions/tools provided by the AI4Media project (aspects such as usefulness, usability, business and journalistic value,</p>

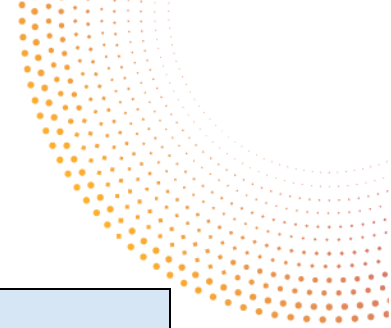






	<p>or performance). By preparing or conducting this research with professional users, DW has generated a dataset containing very limited <i>Participant Personal Data</i> and <i>User Responses related to surveys with questionnaires</i> that were conducted anonymously. Any data will be stored internally on DW's systems/servers and will not be shared with external organisations or AI4Media project partners. The limited Participant Personal Data included information such as Name, Company Email and Company. The data will be used for internal analysis purposes by DW and the production of aggregated results summaries in Deliverables D8.1 (M12), D8.3 (M24) and D8.6 (M48).</p> <p><u>Type/format</u>: Documents containing Participant Personal Data and documents containing questions/user responses from user research activities.</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Internal research to identify potential user research participants, questionnaires or other research methods (focus groups, demo sessions or interviews).</p> <p><u>Expected size</u>: A few KBs per document - a few MBs in total.</p> <p><u>Data utility</u>: This data has been used in the context of WP8 to define the user requirements and to evaluate the AI-enhanced demonstrator (see above). Analysed, aggregated results have been used by use case and technical partners to develop/improve functionalities, the demonstrators and the plans for commercial exploitation.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: No, because the dataset is stored on DW's internal corporate systems and there are no plans for data sharing. The deliverables containing aggregated, analysed results are classified as confidential (D8.1, D8.3 and D8.6) and therefore only available to the consortium.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The data will not be openly accessible since it contains personal information.</p> <p><u>How it will be accessible</u>: It is planned that the data will only be accessible by project partner DW.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
<p>Increase data re-use</p>	<p><u>Licence</u>: The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p>





	<p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset described will be stored on internal corporate systems/servers of DW. Data handling and protection follows standard DW operations and national/EU guidelines. IT security measures mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The data will not be shared.</p> <p><u>Is informed consent for data sharing and long-term preservation given:</u> Questionnaire-based surveys were conducted anonymously where users did not provide their personal data such as Name, Company, or email. They were informed about the purpose of this research, the anonymous approach and how responses are used/reported.</p>
Other Issues	No

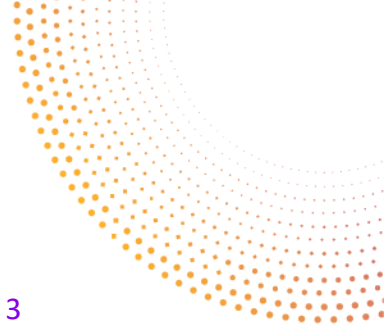
#### 4.6.2 Data from user research activities in Use Case 2

DMP component	AI4Media_Data_51_WP8_USER-RESEARCH_UseCase2_VRT_v1 Partner: VRT
Data Summary	<p><u>Purpose:</u> In order to realise Use Case 2, VRT has conducted qualitative research activities with journalists, who work at VRT and other European media organisations. This research was required for the definition of the user requirements (D8.1), further user requirements gathered for the first White Paper (D8.3) as well as the user evaluation of new AI functionalities as part of the Smart News Assistant (D8.3, D8.6). Mainly qualitative research methods have been conducted such as interviews, focus groups or demo/evaluation sessions with individual users. By preparing or conducting this research with professional users, VRT has generated a dataset containing Online meeting recordings and User Responses related to surveys. Questionnaires were conducted anonymously. Any data will be stored internally on VRT's systems/servers and will not be shared. Recordings are removed after the project. The limited Participant Personal Data is likely to include information such as Name, Company Email, Company, Position, and – if necessary – Area of Work and Additional Contact Details. The data will be used for internal analysis purposes by VRT and the production of aggregated results summaries in WP8 deliverables.</p> <p><u>Type/format:</u> Documents containing Participant Personal Data and documents containing questions/user responses from user research activities. Video recordings of demo sessions and interviews.</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Internal research to identify potential user research participants, questionnaires or other research methods (focus groups, demo sessions or interviews).</p> <p><u>Expected size:</u> A few KBs per document - a few MBs in total.</p> <p><u>Data utility:</u> This data has been used in the context of WP8 to define the final set of user requirements for UC2 and to evaluate the AI functionalities in the UC2</p>



	demonstrators. Analysed, aggregated results have been used by use case and technical partners to develop/improve functionalities, the demonstrators and the plans for commercial exploitation.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, because the dataset described is stored on VRT's internal corporate systems and there are no plans for data sharing. The WP8 deliverables containing aggregated, analysed results are classified as confidential and therefore only available to the consortium.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data will not be openly accessible since it contains personal information.</p> <p><u>How it will be accessible</u>: The data will only be accessible by VRT.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The dataset described will be stored on internal corporate systems/servers of VRT. Data handling and protection follows standard VRT operations and national/EU guidelines. IT security measures mitigate most of the risk of illegitimate access.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The data will not be shared since it contains personal information of end users.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: Questionnaire-based surveys were conducted anonymously where users did not provide their personal data such as Name, Company, or email. They were informed data was only used for research analysis purposes.</p>
Other Issues	No

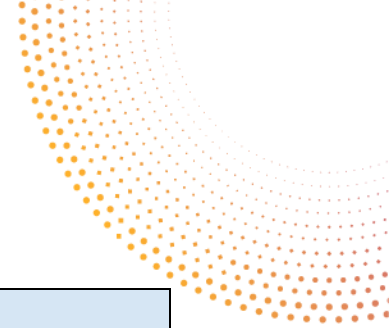




#### 4.6.3 Questionnaires for the collection of user requirements for Use Case 3

DMP component	AI4Media_Data_52_WP8_QUESTIONNAIRE_UseCase3-UserReqCollection2021_V1 Partner: RAI
Data Summary	<p><u>Purpose:</u> A structured questionnaire has been developed by RAI to collect user opinions about end user requirements proposed by UC3. The questionnaire includes questions about each presented feature/epic/user story; its goal is to understand if and how each feature can be useful in a production workflow and to highlight a possible development priority. The questionnaire was filled by RAI journalist/editorial staff people/archivists and by people with the same profiles belonging to other broadcasters. The collected questionnaires have been analysed and the results of this analysis have been used to drive UC3 features development. The results and related analyses are presented in D8.1.</p> <p><u>Type/format:</u> Text documents containing questions and user responses.</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> Questionnaires filled by end-users in the context of use case 3 during the year 1 of the project.</p> <p><u>Expected size:</u> A few KBs per questionnaire. A few MBs in total.</p> <p><u>Data utility:</u> This data has been used in the context of WP8 to evaluate requirements proposed by UC3. The consortium used requirements evaluation results to define a development plan for use case 3 in the project lifetime.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The user requirement evaluation questionnaires for use case 3 are stored on RAI servers. A summary of the requirements evaluation results is available in D8.1.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Raw data is not openly accessible since it contains personal information. It is protected and shared only among project partners and a summary of this data has been presented in D8.1. In case of further reports or papers submitted for publication, all findings will be integrated into the reports or papers. Datasets will not be added to the publication.</p> <p><u>How it will be accessible:</u> The dataset is stored on RAI servers; the questionnaires are only accessible by project partners.</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>



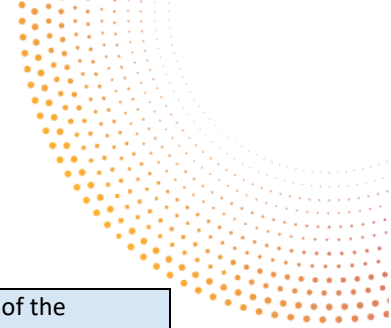


Increase data re-use	<p><u>Licence</u>: The data is not licensed since it can only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality has been assured by asking end users to fill out the questionnaire in their own languages. Feedback collected has been translated to English, in order to ensure accurate data collection and analysis.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data is stored on local repositories inside RAI's premises, subject to the same security measures already used for IT infrastructure in RAI. These include network isolation from external internet accesses, firewalling, account-based access control management to the storage where the data copies are located. RAI fully complies with the applicable national, European data security frameworks, and the GDPR.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The dataset will not be shared since it may contain personal information of end-users.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: An Informed Consent Form was prepared for the participation in the requirements collection activity to manage treatment of personal data. The questionnaires include a Privacy Notice that specifies that the treatment of the data is confidential, complies with GDPR and is carried out exclusively for analytical and statistical purposes.</p>
Other Issues	No

#### 4.6.4 Questionnaires for the evaluation of Use Case 3

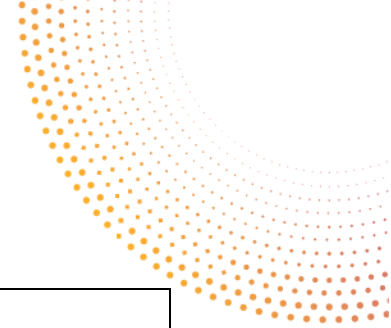
DMP component	AI4Media_Data_53_WP8_QUESTIONNAIRE_UseCase3Evaluation_V1 Partner: RAI
Data Summary	<p><u>Purpose</u>: Structured questionnaires have been developed by RAI for the evaluation of the developed AI services/tools in the context of UC3 pilot trials. The questionnaires include questions that cover issues such as usefulness, usability, visualisation &amp; interaction, learnability, encountered problems and future expectations, etc. as well as user demographics. This dataset includes questionnaires filled by the end users to assess the tools developed for all evaluation phases. Data collected through the questionnaires is used exclusively for analytical and statistical purposes. The evaluation results are presented in relevant WP8 deliverables.</p> <p><u>Type/format</u>: Text documents containing questions and user responses.</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Questionnaires filled by end-users in the context of use case 3 during the different pilot phases.</p> <p><u>Expected size</u>: A few KBs per questionnaire. A few MBs in total.</p> <p><u>Data utility</u>: This data is used in the context of WP8 to evaluate different versions of AI4Media services/tools for UC3. The evaluation results of each pilot phase have been</p>





	used by the technical partners to improve and extend the functionalities of the developed services/tools during the following development phases. The final evaluation results will help partners to improve these services/tools as part of further commercial exploitation.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The evaluation questionnaires (raw data) are stored on RAI servers. A summary of the evaluation results have been presented in relevant WP8 deliverables.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The raw data is not openly accessible since it contains personal information. It is protected and shared only among project partners. However, a summary of the evaluation results has been presented in relevant WP8 deliverables. In case of further reports or papers submitted for publication, all research findings will be integrated into the reports or papers. Datasets will not be added to the publication.</p> <p><u>How it will be accessible</u>: The dataset is stored on RAI servers; the questionnaires are only accessible by project partners.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Dataset is not licensed since it can only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality has been assured by asking end users to fill out the evaluation questionnaire in their own languages. Feedback collected has been translated to English, in order to ensure accurate data collection and analysis.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data is stored on local repositories inside RAI's premises, subject to the same security measures already used for IT infrastructure in RAI. These include network isolation from external internet accesses, firewalling, account-based access control management to the storage where the data copies are located. RAI fully complies with the applicable national, European data security frameworks, and the GDPR.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The dataset will not be shared</p>





	<p>since it may contain personal information of end-users.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> An Informed Consent Form has been prepared for the participation of end-users in the evaluation activities to manage treatment of personal data. The questionnaires include a Privacy Notice that specifies that the treatment of the data is confidential, complies with GDPR and is carried out exclusively for analytical and statistical purposes.</p>
Other Issues	No

#### 4.6.5 Data from user research activities in Use Case 4

DMP component	AI4Media_Data_54_WP8_USER-RESEARCH_UseCase4_NISV_v1 Partner: NISV
Data Summary	<p><u>Purpose:</u> In order to realise Use Case 4 in WP8, partner NISV performed various qualitative research activities with researchers from social sciences and humanities (SSH). This research was required to define the user requirements brought forward in the project, as well as evaluate the CLARIAL Media Suite demonstrator. The following, mainly qualitative research methods were executed with small, selected target groups: tailored questionnaires, interviews, focus groups or demo/evaluation sessions with individual users or small groups.</p> <p>By conducting this research with professional users, NISV generated a dataset containing <i>Participant Personal Data</i> and <i>User Responses</i>. This data is stored internally on NISV's systems/servers and it is not planned to share this data with external organisations or AI4Media project partners. Participant Personal Data include information such as name, company email, position, and area of work/expertise. The data will be used for internal analysis purposes by NISV and the production of aggregated results summaries in WP8 deliverables.</p> <p><u>Type/format:</u> Documents containing Participant Personal Data and documents containing questions/user responses from user research activities.</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Internal research to identify potential user research participants, questionnaires or other research methods (focus groups, demo sessions or interviews).</p> <p><u>Expected size:</u> A few KBs per document - a few MBs in total.</p> <p><u>Data utility:</u> This data was used in the context of WP8 to define the final set of user requirements and to evaluate the demonstrator for Use Case 4. Analysed, aggregated results were shared with technical partners to develop/improve functionalities, the demonstrators and the plans for commercial exploitation. Data presented in the deliverables has been anonymised.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No, because the dataset described will be stored on NISV's internal corporate systems and there are no plans for data sharing. The WP8 deliverables containing aggregated, analysed results are classified as confidential and therefore only available to the consortium.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>



Making data openly accessible	<p><u>Data openly accessible</u>: The data will not be openly accessible since it contains personal information.</p> <p><u>How it will be accessible</u>: The data will only be accessible by project partner NISV.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it is only used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u> N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset described is stored on internal corporate systems/servers of NISV. Data handling and protection follows standard NISV operations and national/EU guidelines. IT security measures mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The data will not be shared since it contains personal information of end users.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: An <i>Informed Consent Form</i> or similar methods to obtain consent was prepared for potential participants in the user research activities described above. This includes information about the purpose of the research, the level of anonymisation of personal information provided, how responses are used/reported and data treatment/compliance.</p>
Other Issues	No

#### 4.6.6 Questionnaires for the evaluations of Use Case 5-B

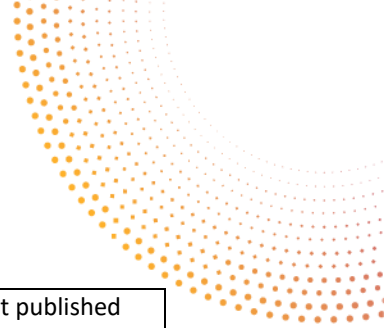
<b>DMP component</b>	<b>AI4Media_Data_55_WP8_QUESTIONNAIRE_UseCase5Evaluation_v1</b> <b>Partner: IRCAM</b>
Data Summary	<p><u>Purpose</u>: For the evaluation of UC5-B (music for games), online questionnaires have been prepared by IRCAM, and filled by the participants of the UC5 tests. The questions are about: the usability and readability of the frontend, the usefulness of the prototype, efficiency of the algorithms, and the relevancy of the approach. The data were first collected in an pseudonymised via Google Forms (<a href="https://docs.google.com/forms">https://docs.google.com/forms</a>), then downloaded on a private directory of the experimentator (at IRCAM) for analysis and statistics. This data was then summarised</p>





	<p>to be published in WP8 deliverables.</p> <p><u>Type/format</u>: Google forms, and Microsoft Excel files.</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Google forms filled by participants of the UC5-B demonstrator evaluation, and Excel documents filled by the experimenter of IRCAM during the interviews.</p> <p><u>Expected size</u>: few MBs.</p> <p><u>Data utility</u>: These data were used to improve the UC5-B demonstrator at each phase of its development along the project lifetime.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The raw (pseudonymised) data are stored locally on a private directory at IRCAM and are not shared publicly. A summary and statistics are published in relevant WP8 deliverables.</p> <p><u>How it will be accessible</u>: The raw datasets are not accessible, their summary and statistics are published in relevant WP8 deliverables.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: No, internal use only</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: The forms were filled by participants during the online tests, then an interview with the experimenter of IRCAM made it possible to confirm the replies, and to obtain additional details, comments and opinions.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The raw datasets are stored locally at IRCAM on a private directory of the experimenter. They are secured by the Computer and Network System of IRCAM. The questionnaires have been done in a pseudonymised way in order to find the response of each participant during the interview, but the association between the pseudo name and true identity of the participant are not stored.</p>



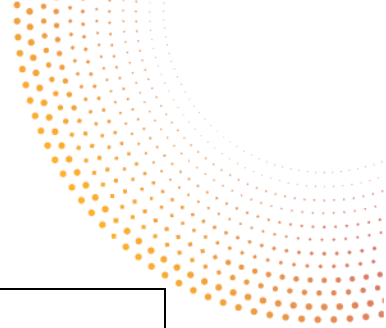


Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The datasets are not published since the contain personal data.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The questionnaires began after an Informed Consent question about the use the collected data. All the evaluation participants validated this question.</p>
Other Issues	No

#### 4.6.7 Data from user research activities in Use Case 7

DMP component	AI4Media_Data_56_WP8_QUESTIONNAIRE_VIDEO_UseCase7-UserFeedbackData_v1 Partner: IMG
Data Summary	<p><u>Purpose:</u> IMG performed interviews with users of the demonstrators in Use Case 7. They also filmed test users in the Living Lab trials they conducted (content managers and team members) using the demonstrators. A structured questionnaire was developed by IMG to collect user opinions about the demonstrators in UC7. Along with the questionnaire, images, videos and heatmaps of how users interact with the demonstrators were captured. The collected feedback has been analysed, and the results of this analysis were used to drive UC7 demonstrators' further improvement.</p> <p><u>Type/format:</u> Text documents containing questions and user responses, along with image and video, as well as heatmap, of the user's interaction with the demonstrators.</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> Questionnaires filled by end-users in the context of UC7. Images, videos and heatmaps of users' interaction with the two demonstrators.</p> <p><u>Expected size:</u> Several MBs per user. Several GBs in total.</p> <p><u>Data utility:</u> The data from the tests on UC7 demonstrators with real users from media companies, in Living Lab trials, were used to validate their performance and help tailor them further to meet user requirements.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The user feedback data for UC7 was securely stored on IMG servers. Summaries of the feedback data and relevant analyses were included in WP8 deliverables.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The raw data will not be openly accessible since it contains personal information. A summary of this data will only be presented in WP8 deliverables. In case of a report or paper submitted for publication, all findings will be integrated into the report or paper. Datasets will not be added to the publication.</p> <p><u>How it will be accessible:</u> The dataset is securely stored on IMG servers; the data will be only accessible by IMG.</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>



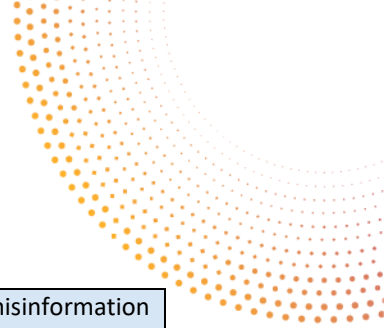


Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: The data collection process will be supervised by IMG.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data are stored on local repositories inside IMG's premises, subject to the same security measures already used for IT infrastructure in IMG. These include network isolation from external internet accesses, firewalling, and access control management to the storage where the data copies are located. IMG fully complies with the applicable national, European data security frameworks, and the GDPR.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The dataset will not be shared since it may contain personal information of end-users.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Consent forms were provided to people that took part in the trials and shared this data with IMG for the purpose of improving our technology and software demonstrators. A Privacy Notice specified that the treatment of the data is confidential, complies with GDPR and is carried out exclusively for analytical and statistical purposes.</p>
Other Issues	No

#### 4.6.8 Truly Media dataset

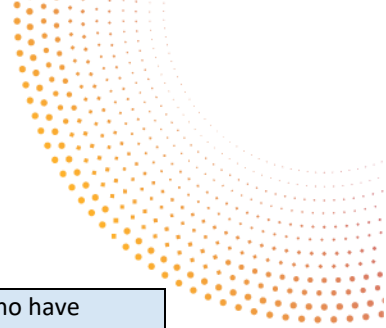
DMP component	AI4Media_Data_57_WP8_Text_TrulyMediaDataset-UseCase1-ATC_v1 Partner: ATC
Data Summary	<p><u>Purpose</u>: Truly Media is used in the context of Use Case 1 as the base demonstrator. The dataset was created for testing purposes. It contains a selection of social media posts from different social media platforms and other web content (news articles, blog posts) depending on the interests and searches made by Truly Media users. The dataset contains social media posts and their content, along with all the accompanying data, such as user name, date and time of post, profile picture, likes/retweets etc. The extent of the accompanying data depends on the API specifications of each social media platform (e.g. Twitter/X provides number of likes for a post, while Facebook does not). No private social media data will be included in the dataset. All social media data will be collected through APIs. The dataset may also contain notes and annotations made to the social media posts and other news items by Truly Media users, as well as manually inserted information to Truly Media by its users, such as the location of an event or a post, other related URLs, similar images found online, etc.</p> <p>AI4Media AI tools will use the relevant data in order to extract useful information and</p>





	<p>produce insights to be used by use case end-users for detecting misinformation attempts as part of UC1.</p> <p><u>Type/format</u>: JSON files</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Social media (Twitter, Facebook, VKontakte, Reddit, YouTube) and web</p> <p><u>Expected size</u>: A few MBs.</p> <p><u>Data utility</u>: The dataset will be used by AI4Media partners that develop tools for Use Case 1 for training purposes.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: That data will be discoverable only by AI4Media partners that will develop tools for Use Case 1 and that will take part in evaluation activities. The dataset is accessible only through Truly Media's environment and is stored on Truly Media's servers.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The dataset will not be shared outside the AI4Media consortium.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: No</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use</u>: No</p> <p><u>Usable by third parties after end of project</u>: No</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset is stored by ATC on third-party cloud servers. Appropriate and detailed security policies, rules, and technical measures are implemented to protect data that are used by the Truly Media platform and stored on the platform from improper or unauthorized access, including use of firewalls where appropriate. Security measures also include 2FA (2 Factor Authentication) with OTP (One Time Password) for extra security during login, as well as Auth2.0 and JWT for authentication and authorisation. End-to-end encryption protects from man-in-the-</p>



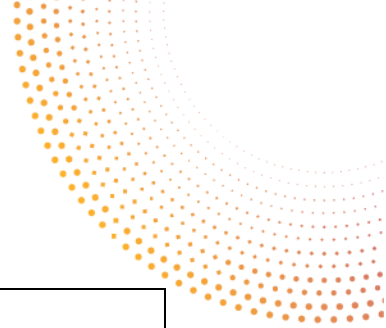


	<p>middle attacks and data theft. All ATC employees and data processors, who have access to and are associated with the processing of personal data, are obliged to respect the confidentiality of the stored personal data. Moreover, ATC's development team has received training from external auditors for security awareness and security best practices to avoid vulnerabilities in source code. The external auditors have performed black-box penetration testing to ensure that the platform is fully secure. ATC's Data Protection Officer ensures that all processes followed are fully compliant with the GDPR provisions.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The dataset will not be shared outside AI4Media since it contains personal data. Only AI4Media partners that develop tools for UC1 and ATC employees working on the project will have access to this dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 4.6.9 Game glitches dataset for Use Case 5

DMP component	AI4Media_Data_58_WP8_Video_GameGlitches-UC5-MODL_v1 Partner: MODL
Data Summary	<p><u>Purpose:</u> This dataset will contain several short videos of games' glitches and corresponding game logs, generated using the modl.ai platform. This dataset will be used in the context of Use Case 5 (feature 5A) to train an AI capable of automatically recognizing different kinds of game glitches.</p> <p><u>Type/format:</u> movie data (.mp4), Game log files (.txt)</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Videos of game glitches created internally in the company (MODL).</p> <p><u>Expected size:</u> ~4 GBs</p> <p><u>Data utility:</u> Training data that will be used to improve the technology in the context of Use Case 5A.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes, via Github and project website</p> <p><u>Search keywords:</u> MODL, AI4Media, Game glitches dataset</p> <p><u>Versioning:</u> Yes</p> <p><u>Metadata creation:</u> No</p>
Making data openly accessible	<p><u>Data openly accessible:</u> A non-personally identifiable version of the data set has been made generally available on GitHub at <a href="https://github.com/modl-ai/ai4media">https://github.com/modl-ai/ai4media</a>.</p> <p><u>How it will be accessible:</u> Via GitHub</p> <p><u>Methods/software tools to access data:</u> Github.com / git</p> <p><u>Repository:</u> GitHub: <a href="https://github.com/modl-ai/ai4media">https://github.com/modl-ai/ai4media</a></p> <p><u>Restrictions on access:</u> No</p>
Making data interoperable	<p><u>Interoperability:</u> No</p> <p><u>Data and metadata vocabularies:</u> No</p>



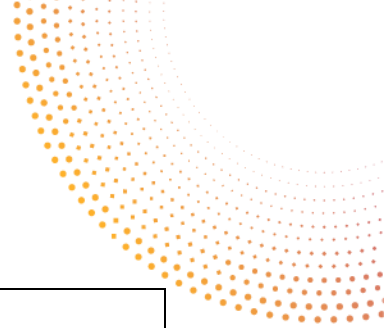


	<p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: Open use, no liability, warranty, or trademark usage.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes</p> <p><u>Re-use timeframe</u>: At least 1 year</p> <p><u>Data quality assurance process</u>: No</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset is entirely synthetic and openly available. It contains no identifying information. Stored in GitHub.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

#### 4.6.10 Musical production for AI co-creation dataset

DMP component	AI4Media_Data_59_WP8_Audio_RAWcompositions_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: New AI4Media dataset generated in audio format from the Demonstrator of AI co-creation in WP8. The Demonstrator for UC6 aims to produce a set of tools to assist artists - especially music composers - to produce in an easy and accessible manner novel creations with the help of an AI engine. This approach requires the set up of a platform to load, train and experiment with generative models. These models are trained from an initial dataset, in our case principally formed by a collection of RAW audio files. These audio files may be decorated with a set of labels to extend the description of each track. Once trained, these models are used by the artist to generate new collections of audio material that can be used for further training of new models, to be used as an inspiration material for novel human compositions, or directly as an audio used in novel content. The input dataset used for training may be selected from an existing collection of audio material, but the output material has to be presented in a consistent way for exploration and analysis by the artist and developers. Along with the platform for audio generation and model selection, this audio dataset is the main validation of the demonstrator.</p> <p><u>Type/format</u>: Raw audio data, uncompressed music</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Generated audio from a ML application.</p> <p><u>Expected size</u>: Around 2GB.</p> <p><u>Data utility</u>: Useful to creators using audio to improve the composition of new music. Outside the project participants, general audience of music compositions.</p>
Making data	<p><u>Is data discoverable</u>: Data is not uniquely identifiable, but can be labeled according to</p>



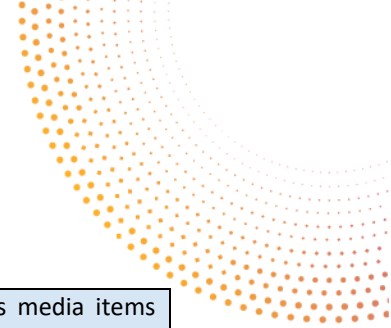


findable, incl. provisions for metadata	<p>some label set.</p> <p><u>Search keywords</u>: No.</p> <p><u>Versioning</u>: Yes, version numbers are assigned according to the training conditions.</p> <p><u>Metadata creation</u>: Labels referred to the generation of creation, and conditions of the generative process.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Data will be openly accessible when required. Access will be granted through BSC data storage facilities, through an open repository like Zenodo, or through a dedicated channel for audio file sharing, like soundcloud. Special cases for sharing content include if an author/creator uses his own data with the generative model; then, data may be partially restricted.</p> <p><u>How it will be accessible</u>: Online repository.</p> <p><u>Methods/software tools to access data</u>: Raw audio is made accessible online.</p> <p><u>Repository</u>: Online repository – TBD.</p> <p><u>Restrictions on access</u>: Access control for restricted.</p>
Making data interoperable	<p><u>Interoperability</u>: No.</p> <p><u>Data and metadata vocabularies</u>: NA</p> <p><u>Use of standard vocabularies</u>: NA</p> <p><u>Mappings to commonly used vocabularies</u>: NA</p>
Increase data re-use	<p><u>Licence</u>: Data will be shared under a Creative Commons CC Zero License.</p> <p><u>Availability for re-use</u>: Data available without embargo.</p> <p><u>Usable by third parties after end of project</u>: Data available for further analysis.</p> <p><u>Re-use timeframe</u>: Data will be in audio format, with no restriction on use.</p> <p><u>Data quality assurance process</u>: Data is in audio format, but quality of these tracks cannot objectively be assessed.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: Audio data in our case is not subject to FAIR criteria.</p> <p><u>Costs for long-term preservation</u>: Server hosting, without additional maintenance.</p>
Data security	<p><u>Security measures</u>: Controlled-access for restricted parts of the dataset, with periodic backup copies.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: We don't foresee any ethical implications for sharing our dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Not needed.</p>
Other Issues	No

#### 4.6.11 Current affairs transcripts dataset for Use Case 4

<b>DMP component</b>	<b>AI4Media_Data_60_WP8_TEXT_CurrentAffairsTranscripts-UC4-NISV_v1</b> <b>Partner: NISV</b>
Data Summary	<u>Purpose</u> : Within the context of Use Case 4, NISV and IDIAP developed and validated a framework that leverages LLMs for automated frame detection in TV transcripts. For this purpose, NISV created a dataset of transcripts of two current affairs programmes





	<p>on Dutch television. Specifically, it contains a selection of 2,000 news media items from broadcasts of news programs EenVan-daag (1,000 items) and Nieuwsuur (1,000 items).</p> <p>EenVandaag 2 (OneToday) is a daily evening program broadcasted on Dutch public television channel NPO1. EenVandaag has the format of a news programme with current issues and background information behind the news. The programme is about 30 minutes long and deals with various news topics during an episode. It has multiple presenters introducing various news items, and also interview experts live in the studio. Nieuwsuur 3 (News Hour) is also an evening programme and is broadcasted on NPO2, a Dutch public television channel. The broadcasts are between 30 and 45 minutes long, and also have the format of a news programme with current issues and background information behind the news. This programme also has multiple presenters and live interviews with experts.</p> <p>We chose these two current newscasts as they provide a good overview of Dutch news items on a daily basis. These shows also provided a large corpus of items over the years, with little change in the show format. Due to this, the data is very consistent.</p> <p><u>Type/format</u>: CSV file containing transcripts.</p> <p><u>Re-use of existing data</u>: The transcripts are extracted from television programmes in the NISV archive.</p> <p><u>Data origin</u>: Transcripts created within AI4Media by NISV</p> <p><u>Expected size</u>: A few KBs in total</p> <p><u>Data utility</u>: For analysis of the framing detection framework, the spoken words in the video recordings of these programmes were transcribed. This was done with the open-source, Kaldi automatic transcriber that NISV uses on audiovisual materials in its archive. This software can automatically transcribe Dutch spoken language into text. This pre-processing step resulted in a dataset of 2,000 texts covering news between 2014 and 2018, varying in length, with an average number of 499 words for Nieuwsuur and 664 words for EenVandaag. The data was then translated from Dutch to English using the Deepl API.</p> <p>The resulting dataset was classified by a human annotator and GPT-3.5 using a prompt composed by the definition of different frame types.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: No. The dataset is stored on NISV's internal server.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: It is not possible to make this dataset openly accessible since the Dutch public broadcaster owns copyright for the transcripts and has not given a permission to make this dataset accessible. NISV is in ongoing negotiations with the broadcaster to make exceptions for this dataset and others created in similar research contexts.</p> <p><u>How it will be accessible</u>: The data is stored on NISV's internal server. Based on a research agreement, temporary access to the dataset was given to IDIAP for the purposes of the experiment.</p>



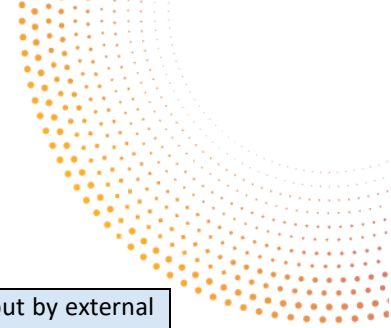


	<p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since copyright restrictions prevent reuse.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u> N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset is stored on internal corporate systems/servers of NISV. Data handling and protection follows standard NISV operations and national/EU guidelines. IT security measures mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: It is not possible to make this dataset openly accessible since the Dutch public broadcaster owns copyright for the transcripts and has not given a permission to make this dataset accessible.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

#### 4.6.12 User survey data for AI industrial needs (for T8.4)

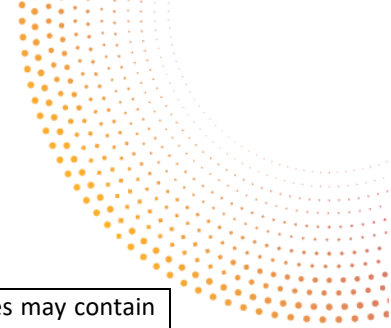
DMP component	AI4Media_Data_61_WP8_QUESTIONNAIRE_AI-IndustrialNeeds-T8.4_v1 Partner: ATC
Data Summary	<p><u>Purpose</u>: Structured questionnaires were developed by WP8 partners in the framework of T8.4 “Harmonising AI research with industrial needs” for the collection of input by external stakeholders. The questionnaires include questions that aim to gather input and insights from industry stakeholders regarding the AI needs of media organisations. This dataset includes questionnaires filled by the industry stakeholders. Data collected through the questionnaires is used exclusively for analytical and statistical purposes to guide T8.4 partners in harmonising AI research with the needs of media organisations.</p> <p><u>Type/format</u>: Text documents containing questions and responses.</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Survey participants</p> <p><u>Expected size</u>: A few KBs per questionnaire. A few MBs in total</p> <p><u>Data utility</u>: This data is used in the context of WP8, and more specifically T8.4</p>





	<p>“Harmonising AI research with industrial needs” for the collection of input by external stakeholders. Data collected through the questionnaires is used exclusively for analytical and statistical purposes to guide T8.4 partners in harmonising AI research with the needs of media organisations.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The questionnaires (raw data) are stored on the servers of the partners that participate in the relevant T8.4 activities. A summary of the results is included in D8.3 and D8.6, which are confidential and will be accessible only by consortium partners and the EC.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The raw data will not be openly accessible since it contains personal information. It is protected and shared only among selected project partners. A summary of the results will be available through confidential deliverables D8.3 and D8.6. In case of a report or paper submitted for publication, all research findings will be integrated into the report or paper. Datasets will not be added to the publication.</p> <p><u>How it will be accessible:</u> Stored in the servers of T8.4 partners that participate in the related activities, the questionnaires are only accessible by WP8 project partners that participate in the related T8.4 activities.</p> <p><u>Methods/software tools to access data:</u> N/A</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Data quality is assured by asking participants to fill out the questionnaire in their own languages, and in the case this is not possible by ensuring that survey participants clearly understand all the questions. Feedback collected has been translated to English, in order to ensure accurate data collection and analysis.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The questionnaires are stored on the partners’ servers. Appropriate security mechanisms enforced by each partner include use of firewalls and restricted access to the folders only to partner employees that work on AI4Media.</p>





Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The questionnaires may contain personal information of end-users.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> Informed Consent Forms were prepared for the participation in T8.4 validation activities. The questionnaires include a Privacy Notice that specifies that the treatment of the data is confidential, complies with GDPR, and is carried out exclusively for analytical and statistical purposes (see D12.1).</p>
Other Issues	No





## 5. Data management plan for third-party research datasets used in AI4Media

This section presents existing third-party research datasets that were used by AI4Media partners to develop and test new AI algorithms, methodologies, and tools in the context of WPs 3, 4, 5 and 6. Most of these datasets are already publicly shared by their third party owners while some are private datasets that have been provided to AI4Media partners for research purposes exclusively.

In the initial DMP, 51 third-party research datasets were identified. Since then the number of third-party datasets used within AI4Media has risen to **81 datasets** in total. In the following subsections, we present the DMP plan for these datasets using the same template as in Section 4 (see Table 1).

The following Table briefly summarizes the 81 datasets presented in this section and offers a glance at the structure of the section and its subsections. New datasets that were not included in the initial DMP are indicated with yellow.

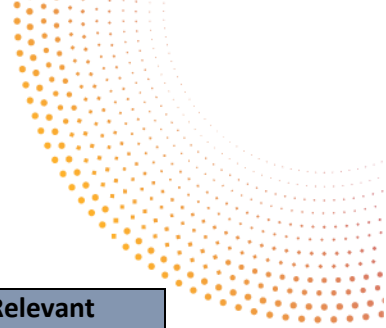
*Table 3: Summary of third-party research datasets used in AI4Media*

DMP component	WP	Short summary	Relevant sub-section
<b>Data used in WP3 (New Learning Paradigms &amp; Distributed AI)</b>			5.1
AI4Media_Data_62_WP3_EMAIL_Enron_v1	WP3	Enron email dataset	5.1.1
AI4Media_Data_63_WP3_SOCIALMEDIA_FacebookWall_v1	WP3	Facebook Wall dataset	5.1.2
AI4Media_Data_64_WP3_IMAGE_AFFECT_GAME_ANNOTATION_AGAIN_v1	WP3	Affect Game Annotation (AGAIN) dataset	5.1.3
AI4Media_Data_65_WP3_TEXT_IMDB_Reviews_v1	WP3	IMDB movie reviews dataset	5.1.4
AI4Media_Data_66_WP3_TEXT_MHAD_2D_Pose_v1	WP3	MHAD 2D pose dataset	5.1.5
AI4Media_Data_67_WP3-5_TEXT_20Newsgroups_v1	WP3,5	20Newsgroups dataset	5.1.6
AI4Media_Data_68_WP3-5_TEXT_HPAmazonReviews_v1	WP3,5	HP Amazon reviews dataset	5.1.7
AI4Media_Data_69_WP3-5_TEXT_JRCAcquis_v1	WP3,5	JRCAcquis legislative text dataset	5.1.8
AI4Media_Data_70_WP3-5_TEXT_Kindle_v1	WP3,5	Kindle document dataset	5.1.9
AI4Media_Data_71_WP3-5_TEXT_OHSUMED_v1	WP3,5	OHSUMED MEDLINE document dataset	5.1.10
AI4Media_Data_72_WP3-5_TEXT_RCV1-Reuters_v1	WP3,5	RCV1 Reuters stories dataset	5.1.11
AI4Media_Data_73_WP3-5_TEXT_RCV1RCV2-Reuters_v1	WP3,5	RCV1RCV2 Reuters stories dataset	5.1.12
AI4Media_Data_74_WP3-5_TEXT_Reuters-21578_v1	WP3,5	Reuters-21578 dataset	5.1.13
AI4Media_Data_75_WP3-	WP3,5	11 Tweet Sentiment datasets	5.1.14



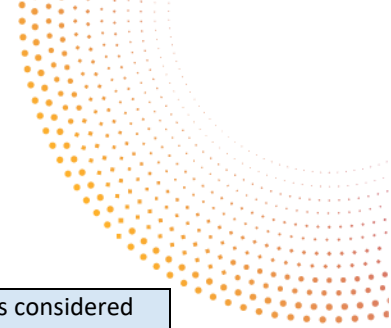
DMP component	WP	Short summary	Relevant sub-section
5_TEXT_11TweetSentiment_v1			
AI4Media_Data_76_WP3-5_TEXT_WipoGamma_v1	WP3	WipoGamma patent document dataset	5.1.15
AI4Media_Data_77_WP3_Image_WIDER_FACE_v1	WP3	WIDER FACE face detection dataset	5.1.16
AI4Media_Data_78_WP3_IMAGE_Caltech101_v1	WP3	Caltech 101 dataset	5.1.17
AI4Media_Data_79_WP3_IMAGE_CUB-200-2011_v1	WP3	CUB-200-2011 image dataset	5.1.18
AI4Media_Data_80_WP3_IMAGE_Describable-Textures_v1	WP3	Describable Textures dataset	5.1.19
AI4Media_Data_81_WP3_IMAGE_Food101_v1	WP3	Food-101 dataset	5.1.20
AI4Media_Data_82_WP3_IMAGE_MAMe_v1	WP3	MAMe dataset	5.1.21
AI4Media_Data_83_WP3_IMAGE_MIT-ISR_v1	WP3	MIT-ISR dataset	5.1.22
AI4Media_Data_84_WP3_IMAGE_OuluKnots_v1	WP3	Oulu Knots dataset	5.1.23
AI4Media_Data_85_WP3_IMAGE_OxfordFlowers_v1	WP3	Oxford Flower dataset	5.1.24
AI4Media_Data_86_WP3_IMAGE_OxfordIIITPet_v1	WP3	Oxford-IIIT Pet dataset	5.1.25
AI4Media_Data_87_WP3_IMAGE_StanfordDogs_v1	WP3	Stanford Dogs dataset	5.1.26
<b>Data used in WP4 (Explainability, Robustness and Privacy in AI)</b>			5.2
AI4Media_Data_88_WP4_IMAGE_ImageNet_01	WP4	ImageNet-ILSVRC2012 image classification dataset	5.2.1
AI4Media_Data_89_WP4_IMAGE_FFHQ_v1	WP4	FFHQ dataset for GAN training	5.2.2
AI4Media_Data_90_WP4_IMAGE_MNIST_v1	WP4	MNIST image dataset	5.2.3
AI4Media_Data_91_WP4_IMAGE_Video_DeepLearning10k	WP4	Interestingness10k image +video dataset	5.2.4
AI4Media_Data_92_WP4_VIDEO_Memorability2020	WP4	MediaEval Memorability 2020 dataset	5.2.5
AI4Media_Data_93_WP4_IMAGE_DrawnUI2021	WP4	ImageCLEF DrawnUI 2021 dataset	5.2.6
AI4Media_Data_94_WP4_Video_FCVID_v1	WP4	FCVID event recognition dataset	5.2.7
AI4Media_Data_95_WP4_Video_YLI-MED_v1	WP4	YLI-MED event recognition dataset	5.2.8
<b>Data used in WP5 (Content-centered AI)</b>			5.3
AI4Media_Data_96_WP5_VIDEO_SumMeGycli14_v1	WP5	SumMe video summarization dataset	5.3.1
AI4Media_Data_97_WP5_VIDEO_TVSumSong15_v1	WP5	TVSum video summarization dataset	5.3.2
AI4Media_Data_98_WP5_VIDEO_MonumentsOfItaly_v1	WP5	RAI Monuments of Italy dataset	5.3.3

DMP component	WP	Short summary	Relevant sub-section
AI4Media_Data_99_WP5_IMAGE_LV IS_v1	WP5	LVIS image dataset	5.3.4
AI4Media_Data_100_WP5_IMAGE_CIFAR_v1	WP5	CIFAR10/100 image dataset	5.3.5
AI4Media_Data_101_WP5_IMAGE_STL10_v1	WP5	STL-10 image dataset	5.3.6
AI4Media_Data_102_WP5_TEXT_CCNET_v1	WP5	CCNet text dataset for multilingual representation learning	5.3.7
AI4Media_Data_103_WP5_VIDEO_360VideoViewingHeadMountedVR_v1	WP5	360 Video Viewing Dataset in Head Mounted Virtual Reality	5.3.8
AI4Media_Data_104_WP5_VIDEO_HeadMovementinPanoramicVideo_v1	WP5	Predicting Head Movement in Panoramic Video Dataset	5.3.9
AI4Media_Data_105_WP5_VIDEO_GazePredictionDynamic360ImmersiveVideos_v1	WP5	Gaze prediction in Dynamic 360° Immersive Videos Dataset	5.3.10
AI4Media_Data_106_WP5_VIDEO_YourAttentionIsUnique_v1	WP5	Your Attention is Unique Dataset	5.3.11
AI4Media_Data_107_WP5_VIDEO_HeadEyeMovements360Videos_v1	WP5	Dataset of Head and Eye Movements for 360° Videos	5.3.12
AI4Media_Data_108_WP5_Audio_dim-sim_music_v1	WP5	Dim-sim dataset for music similarity search	5.3.13
AI4Media_Data_109_WP5_Audio_spam_music_v1	WP5	SPAM dataset for music segmentation	5.3.14
AI4Media_Data_110_WP5_Audio_salami_music_v1	WP5	SALAMI dataset for music segmentation	5.3.15
AI4Media_Data_111_WP5_Audio_harmonix_music_v1	WP5	Harmonix dataset for music segmentation	5.3.16
AI4Media_Data_112_WP5_Audio_FreeMusicArchive_v1	WP5	Free Music Archive dataset	5.3.17
AI4Media_Data_113_WP5_Audio_LAKH-MIDI_v1	WP5	LAKH MIDI music dataset	5.3.18
AI4Media_Data_114_WP5_Audio_MIDI_Piano_v1	WP5	Piano Audio and MIDI music datasets	5.3.19
AI4Media_Data_115_WP5_Audio_GiantSteps_v1	WP5	GiantSteps music datasets	5.3.20
AI4Media_Data_116_WP5_IMAGE_MSCOCO_v1	WP5	MS COCO dataset	5.3.21
AI4Media_Data_117_WP5_VIDEO_BVI-DVC_v1	WP5	BVI-DVC dataset	5.3.22
AI4Media_Data_118_WP5_Image_Adience_v1	WP5	Adience dataset	5.3.23
AI4Media_Data_119_WP5_Image_IMDB-Wiki_v1	WP5	IMDB-Wiki dataset	5.3.24
<b>Data used in WP6 (Human- and Society-centred AI)</b>			<b>5.4</b>
AI4Media_Data_120_WP6_VIDEO_Deepfake-Detection-Challenge-Dataset_v1	WP6	Deepfake Detection Challenge video dataset	5.4.1
AI4Media_Data_121_WP6_VIDEO_FaceForensics+_v1	WP6, WP4	FaceForensics++ video dataset	5.4.2
AI4Media_Data_122_WP4_IMAGE_ImageCLEFaware-dataset_v1	WP6	Visual profile impact rating and ranking – ImageCLEFaware dataset	5.4.3



DMP component	WP	Short summary	Relevant sub-section
AI4Media_Data_123_WP6_EEG_DEA P_v1	WP6	DEAP EEG dataset	5.4.4
AI4Media_Data_124_WP6_EEG_SEE D_v1	WP6	SEED EEG dataset	

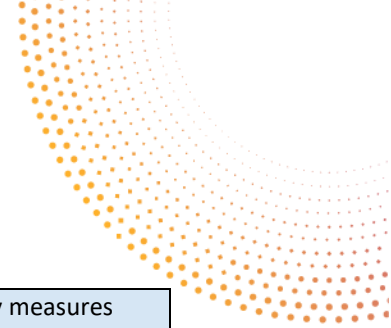




	<p>CERTH in T3.5 to simulate decentralized settings, where each employee is considered to own a separate device. Results involving this dataset will be reported in D3.2 and D3.4.</p> <p><u>Type/format</u>: Csv file containing timestamp, an anonymized email sender and anonymized email receiver</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: Emails of Enron employees obtained by FERC and available at <a href="http://networkrepository.com/ia-enron-email-dynamic.php">http://networkrepository.com/ia-enron-email-dynamic.php</a></p> <p><u>Expected size</u>: 4.2 MB</p> <p><u>Data utility</u>: It is useful to WP3 partners to benchmark graph mining algorithms, such as node ranking and graph neural networks on link prediction tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted in the Network Repository.</p> <p><u>Search keywords</u>: Enron email dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="http://networkrepository.com/ia-enron-email-dynamic.php">http://networkrepository.com/ia-enron-email-dynamic.php</a></p> <p><u>How it will be accessible</u>: Shared through a third-party repository link</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file</p> <p><u>Repository</u>: Network Repository (<a href="http://networkrepository.com">http://networkrepository.com</a>)</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under an attribution licence (<a href="http://networkrepository.com/policy.php">http://networkrepository.com/policy.php</a>)</p> <p><u>Availability for re-use</u>: The loading and pre-processing mechanism developed for using the dataset in experiments will be made publicly available to ensure reproducibility of research.</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: This dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection</p>







	frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.2 Facebook wall dataset

DMP component	AI4Media_Data_63_WP3_SOCIALMEDIA_FacebookWall_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: This dataset captures a Facebook friendship graph where nodes are users and edges between the users represent wall post events. It will be used in T3.5 to simulate decentralized settings, where each user is considered to own a separate device. Results involving this dataset will be reported in D3.2 and D3.4.</p> <p><u>Type/format</u>: Csv file containing timestamp, an anonymized message sender and anonymized message receiver</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: Facebook - <a href="http://networkrepository.com/ia-facebook-wall-wosn-dir.php">http://networkrepository.com/ia-facebook-wall-wosn-dir.php</a></p> <p><u>Expected size</u>: 6.7 MB</p> <p><u>Data utility</u>: It is useful to WP3 partners to benchmark graph mining algorithms, such as node ranking and graph neural networks on link prediction tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted in the Network Repository.</p> <p><u>Search keywords</u>: facebook-wall-wosn</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="http://networkrepository.com/ia-facebook-wall-wosn-dir.php">http://networkrepository.com/ia-facebook-wall-wosn-dir.php</a></p> <p><u>How it will be accessible</u>: Shared through a third-party repository link</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file</p> <p><u>Repository</u>: Network Repository (<a href="http://networkrepository.com">http://networkrepository.com</a>)</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data	<u>Licence</u> : The dataset is already publicly available under an attribution licence



re-use	<p><a href="http://networkrepository.com/policy.php">http://networkrepository.com/policy.php</a></p> <p><u>Availability for re-use:</u> The loading and pre-processing mechanism developed for using the dataset in experiments will be made publicly available to ensure reproducibility of research.</p> <p><u>Usable by third parties after end of project:</u> This is an open dataset.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.1.3 Affect Game Annotation (AGAIN) dataset

DMP component	AI4Media_Data_64_WP3_IMAGE_AGAIN_v1 Partner: UM
Data Summary	<p><u>Purpose:</u> AGAIN is a large-scale affective corpus that features over \$1,100\$ in-game videos (with corresponding gameplay data) from nine different games, which are annotated for arousal from 124 participants in a first-person continuous fashion. Even though AGAIN is created for the purpose of investigating the generality of affective computing across dissimilar tasks, affect modelling can be studied within each of its 9 specific interactive games. AGAIN will likely be used in WP3 to evaluate affect-driven quality diversity algorithms.</p> <p><u>Type/format:</u> Annotated in-game video footage images and accompanying telemetry-and metadata.</p> <p><u>Re-use of existing data:</u> Yes, the dataset is created within the TAMED Marie Curie project.</p> <p><u>Data origin:</u> Annotated by MTurk workers, based on original artefacts created for the TAMED Marie Cure project: <a href="https://again.institutedigitalgames.com/">https://again.institutedigitalgames.com/</a></p> <p><u>Expected size:</u> 42.4 GB</p> <p><u>Data utility:</u> It will be useful to AI4Media partners that investigate affect detection in relation to media.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is made available on the UM IDG Google Drive storage and the following site: <a href="https://again.institutedigitalgames.com">https://again.institutedigitalgames.com</a></p> <p><u>Search keywords:</u> dataset, videogame, affect, arousal, video repository, affective</p>

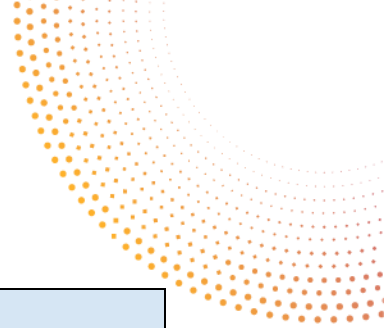


	<p>computing, affective dataset</p> <p><u>Versioning</u>: Google Drive supports versioning.</p> <p><u>Metadata creation</u>: Metadata written in a readable and searchable JSON file.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is openly accessible via <a href="https://again.institutedigitalgames.com">https://again.institutedigitalgames.com</a></p> <p><u>How it will be accessible</u>: The data can be downloaded from an online archive after completing a form.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: <a href="https://drive.google.com/drive/u/1/folders/1f4e00A0JH6FE8n5v7ozjJztObPcdLzoJ">https://drive.google.com/drive/u/1/folders/1f4e00A0JH6FE8n5v7ozjJztObPcdLzoJ</a></p> <p><u>Restrictions on access</u>: The user should accept the terms of use.</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: datapackage.json includes schemas and field descriptions of the dataset.</p> <p><u>Use of standard vocabularies</u>: The dataset metadata follows Data Package specification.</p> <p><u>Mappings to commonly used vocabularies</u>: The Kaggle API implements the same data specifications for datasets.</p>
Increase data re-use	<p><u>Licence</u>: The data is released under MIT License.</p> <p><u>Availability for re-use</u>: The data is available openly for reuse for research purposes.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: The data has been sorted and cleaned for further use. Unusable data is removed from the dataset, and both raw unprocessed and preprocessed ready-to-use packages are available from the remaining data.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The full dataset (including images and non-anonymized metadata) will be hosted on UM's servers. UM fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes</p>
Other Issues	N/A

#### 5.1.4 IMDB movie reviews dataset

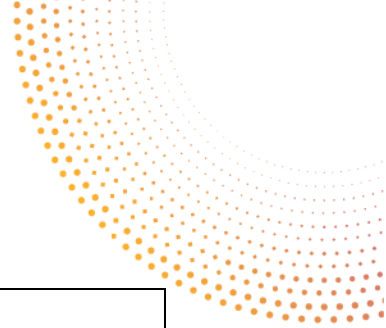
<b>DMP component</b>	<b>AI4Media_Data_65_WP3_TEXT_IMDBReviews_v1</b> <b>Partner: IDF</b>
Data Summary	<u>Purpose</u> : This dataset includes movie reviews extracted from the IMDB website. It consists in 25,000 training and 25,000 test reviews annotated as positive or negative.





	<p>This data will be used as benchmark for T3.5 by IDF.</p> <p><u>Type/format</u>: Raw text and already processed bag of words formats</p> <p><u>Re-use of existing data</u>: Yes, we use an existing dataset</p> <p><u>Data origin</u>: <a href="https://ai.stanford.edu/~amaas/data/sentiment/">https://ai.stanford.edu/~amaas/data/sentiment/</a></p> <p><u>Expected size</u>: 250MB</p> <p><u>Data utility</u>: As benchmark for T3.5 for IDF.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is available online (<a href="https://ai.stanford.edu/~amaas/data/sentiment/">https://ai.stanford.edu/~amaas/data/sentiment/</a>) and indexed in Google. It is also directly available in some frameworks such as tensorflow.</p> <p><u>Search keywords</u>: imdb, sentiment analysis, stanford</p> <p><u>Versioning</u> N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Data has already been made publicly available by its owners. We will not re-share it.</p> <p><u>How it will be accessible</u>: Accessible for original source: <a href="https://ai.stanford.edu/~amaas/data/sentiment/">https://ai.stanford.edu/~amaas/data/sentiment/</a></p> <p><u>Methods/software tools to access data</u>: Download the dataset. It is also available directly from some framework such as tensorflow</p> <p><u>Repository</u>: <a href="https://ai.stanford.edu/~amaas/data/sentiment/">https://ai.stanford.edu/~amaas/data/sentiment/</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. <u>Data and metadata vocabularies</u>: Data will not be altered as it is easy to read.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: When using this dataset, please cite <i>Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. (2011). Learning Word Vectors for Sentiment Analysis. The 49th Annual Meeting of the Association for Computational Linguistics (ACL 2011).</i></p> <p><u>Availability for re-use</u>: Data is already available online.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on IDF servers provided that this is not forbidden by the owner of the dataset. IDF fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>



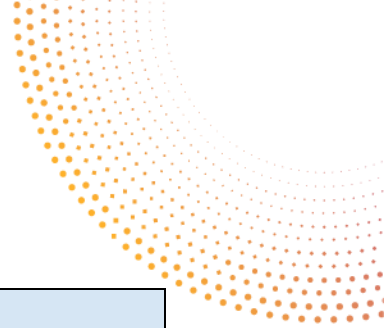


Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.1.5 MHAD 2D pose dataset

DMP component	AI4Media_Data_66_WP3_TEXT_MHAD2DPose_v1 Partner: IDF
Data Summary	<p><u>Purpose:</u> This is a human action recognition dataset consisting of 2D pose estimations extracted from videos. The 2D pose estimation was performed on a subset of the Berkeley Multimodal Human Action Database (MHAD) dataset. Six actions (Jumping in place, Jumping jacks, Punching (boxing), Waving two hands, Waving one hand, Clapping hands) are considered for a total of 1438 videos recorded at 22Hz. This data will be used as benchmark for T3.5 by IDF.</p> <p><u>Type/format:</u> Text files</p> <p><u>Re-use of existing data:</u> Yes</p> <p><u>Data origin:</u> Pose estimation performed on a subset of the MHAD dataset, available at <a href="https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input">https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input</a></p> <p><u>Expected size:</u> 250MB</p> <p><u>Data utility:</u> As benchmark for T3.5 for IDF.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is available online at Github <a href="https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input">https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input</a> and indexed in Google.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Data is already available online, shared by its owners. We will not re-share it.</p> <p><u>How it will be accessible:</u> Can be downloaded from GitHub at <a href="https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input">https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input</a></p> <p><u>Methods/software tools to access data:</u> Download the dataset.</p> <p><u>Repository:</u> GitHub</p> <p><u>Restrictions on access:</u> Access based on license on <a href="https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input">https://github.com/stuarteiffert/RNN-for-Human-Activity-Recognition-using-2D-Pose-Input</a></p>
Making data interoperable	<p><u>Interoperability:</u> The file structure makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies:</u> Data will not be altered as it is easy to read.</p> <p><u>Use of standard vocabularies:</u> No</p> <p><u>Mappings to commonly used vocabularies:</u> No</p>





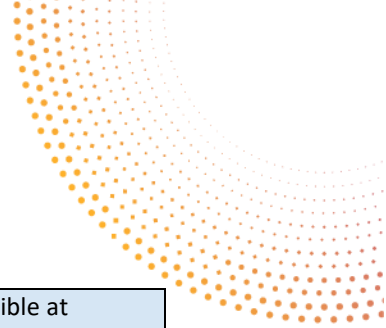
Increase data re-use	<p><u>Licence</u>: The data is already openly shared under a BSD-2 license.</p> <p><u>Availability for re-use</u>: Data is already available.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on IDF servers provided that this is not forbidden by the owner of the dataset. IDF fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes.</p>
Other Issues	N/A

### 5.1.6 20Newsgroups dataset

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_67_WP3-5_TEXT_20Newsgroups_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: The 20Newsgroups dataset consists of about 20,000 messages posted in the early '90s on 20 different Usenet discussion groups. The dataset is used in AI4media (and has been used elsewhere since the '90s) as a benchmark (training set + test set) for testing text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format</u>: Raw text</p> <p><u>Re-use of existing data</u>: Yes, this dataset has been in the public domain since the '90s.</p> <p><u>Data origin</u>: The data consist of messages posted in the early '90s on Usenet discussion groups.</p> <p><u>Expected size</u>: Approximately 20,000 documents.</p> <p><u>Data utility</u>: It is, and it has been for decades, useful to text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data have been made available by its creator on <a href="http://qwone.com/~jason/20Newsgroups/">http://qwone.com/~jason/20Newsgroups/</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords</u>: No search keywords provided.</p> <p><u>Versioning</u>: There is only one version which has been made available by its creator.</p> <p><u>Metadata creation</u>: No metadata were attached to this dataset by its creator, aside from the set of classes that is to be used for labelling the documents.</p>





Making data openly accessible	<p><u>Data openly accessible</u>: The creator of the dataset makes it openly accessible at <a href="http://qwone.com/~jason/20Newsgroups/">http://qwone.com/~jason/20Newsgroups/</a>. We will not reshare the data.</p> <p><u>How it will be accessible</u>: The dataset has been accessible from its creator's home page ever since the '90s.</p> <p><u>Methods/software tools to access data</u>: The only software tool needed to access the data is a web browser.</p> <p><u>Repository</u>: <a href="http://qwone.com/~jason/20Newsgroups/">http://qwone.com/~jason/20Newsgroups/</a>.</p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Already publicly available without license at <a href="http://qwone.com/~jason/20Newsgroups/">http://qwone.com/~jason/20Newsgroups/</a>.</p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.7 HP Amazon reviews dataset

This dataset will also be used in the context of WP5.

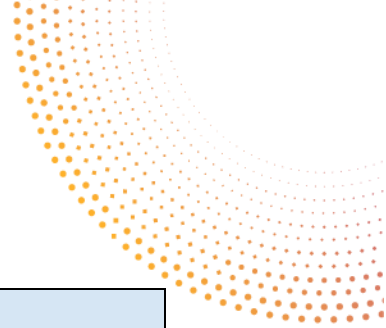
<b>DMP component</b>	<b>AI4Media_Data_68_WP3-5_TEXT_HPAmazonReviews_v1</b> <b>Partner: CNR</b>
Data Summary	<p><u>Purpose</u>: The HP dataset consists of 27,932 documents. The dataset is used in AI4media as a benchmark (training set of 9,533 documents + test set of 18,399) for testing text quantification methods, e.g., in the context of T3.7 and T5.4. Every document in the dataset is a plain-text tokenized review, with associated a sentiment</p>



	<p>label: “positive” or “negative”.</p> <p><u>Type/format</u>: Raw text, with a binary label associated to every document.</p> <p><u>Re-use of existing data</u>: Yes. The dataset was built by CNR prior to the start of the project.</p> <p><u>Data origin</u>: The dataset consists of product reviews for the Amazon Kindle e-book reader, collected from public reviews published on the Amazon.com website.</p> <p><u>Expected size</u>: 25,421 documents.</p> <p><u>Data utility</u>: Quantification is an emerging research topic in the field of aggregated data analysis. Sentiment quantification on text is of special interest for its usefulness in text mining application on social data.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data have been made available by publishing it on Zenodo (<a href="https://doi.org/10.5281/zenodo.4117827">https://doi.org/10.5281/zenodo.4117827</a>). In every paper we write, we cite the resource and the URL from where the dataset can be downloaded.</p> <p><u>Search keywords</u>: No search keywords provided.</p> <p><u>Versioning</u>: There is only one version, which has been made available by its creator.</p> <p><u>Metadata creation</u>: The dataset has no metadata attached.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already publicly shared at <a href="https://doi.org/10.5281/zenodo.4117827">https://doi.org/10.5281/zenodo.4117827</a>. We will not reshare the data.</p> <p><u>How it will be accessible</u>: The data are downloadable from the Zenodo website.</p> <p><u>Methods/software tools to access data</u>: The dataset consists of two text files. A Web browser is required to download them.</p> <p><u>Repository</u>: The dataset is published on Zenodo with DOI 10.5281/zenodo.4117827.</p> <p><u>Restrictions on access</u>: There are no restrictions on the access and use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: No data or metadata vocabularies are used, since content is plain text and no metadata are available.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared on Zenodo under a Creative Commons Attribution 4.0 International license.</p> <p><u>Availability for re-use</u>: Data are already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: Yes.</p> <p><u>Re-use timeframe</u>: Unlimited.</p> <p><u>Data quality assurance process</u>: The data has not been subjected to a data quality assurance process.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: No</p> <p><u>Costs for long-term preservation</u>: No</p>
Data security	<p><u>Security measures</u>: Data is released with CC licence, no security measures for data</p>







	protection have been implemented.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> As the creators of this dataset, CNR did not foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.1.8 JRCAcquis legislative text dataset

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_69_WP3-5_TEXT_JRCAcquis_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> JRC-Acquis is a collection of legislative texts of European Union law written between the 1950s and 2006, and classified according to a set of 6,000 classes which describe what the text is about; the data are multilingual, i.e., each news story is written in one of 22 official European languages. The dataset is used in AI4media (and has been used elsewhere since the mid 2000's) as a benchmark (training set + test set) for testing multilingual text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format:</u> Raw text</p> <p><u>Re-use of existing data:</u> Yes, this dataset has been in the public domain since the mid 2000's.</p> <p><u>Data origin:</u> The stories are legislative texts of European Union law.</p> <p><u>Expected size:</u> About 111,700 documents.</p> <p><u>Data utility:</u> It is, and it has been for many years, useful to multilingual text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their multilingual text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data have been made available by its creators from the EU website at <a href="https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis">https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords:</u> JRC Acquis</p> <p><u>Versioning:</u> There is only one version (3.0) which has been made available by its creators.</p> <p><u>Metadata creation:</u> No metadata were attached to this dataset by its creators, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The creators of the dataset make it openly accessible at <a href="https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis">https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis</a>. We will not reshare the data.</p> <p><u>How it will be accessible:</u> The dataset has been accessible from <a href="https://ec.europa.eu/">https://ec.europa.eu/</a> from the beginning</p> <p><u>Methods/software tools to access data:</u> The only software tool needed to access the data is a web browser.</p>



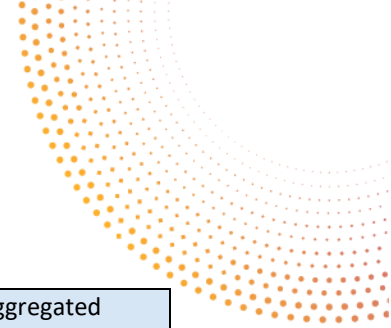
	<p><u>Repository</u>: The data is deposited in the UCI ML repository.</p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Already publicly available dataset. License and terms of use defined at <a href="https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis">https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis</a>.</p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.9 Kindle document dataset

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_70_WP3-5_TEXT_Kindle_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: The Kindle dataset consists of 25,421 documents. The dataset is used in AI4media as a benchmark (training set of 21,591 documents + test set of 3,821) for testing text quantification methods, e.g., in the context of T3.7 and T5.4. Every document in the dataset is a plain-text tokenized review, with associated a sentiment label: “positive” or “negative”.</p> <p><u>Type/format</u>: Raw text, with a binary label associated to every document.</p> <p><u>Re-use of existing data</u>: The dataset was built by CNR prior to the start of the project.</p> <p><u>Data origin</u>: The dataset consists of product reviews for the Amazon Kindle e-book reader, collected from public reviews published on the Amazon.com website.</p> <p><u>Expected size</u>: 25,421 documents.</p>





	<p><b>Data utility:</b> Quantification is an emerging research topic in the field of aggregated data analysis. Sentiment quantification on text is of special interest for its usefulness in text mining application on social data.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> The data have been made available by publishing it on Zenodo (<a href="https://doi.org/10.5281/zenodo.4117827">https://doi.org/10.5281/zenodo.4117827</a>). In every paper we write, we cite the resource and the URL from where the dataset can be downloaded.</p> <p><b>Search keywords:</b> No search keywords provided.</p> <p><b>Versioning:</b> There is only one version, which has been made available by its creator.</p> <p><b>Metadata creation:</b> The dataset has no metadata attached.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> Yes, the data is already openly accessible at <a href="https://doi.org/10.5281/zenodo.4117827">https://doi.org/10.5281/zenodo.4117827</a>. We will not reshare the data.</p> <p><b>How it will be accessible:</b> The data are downloadable from the Zenodo website.</p> <p><b>Methods/software tools to access data:</b> The dataset consists of two text files. A Web browser is required to download them.</p> <p><b>Repository:</b> The dataset is published on Zenodo with DOI 10.5281/zenodo.4117827.</p> <p><b>Restrictions on access:</b> There are no restrictions on the access and use of this dataset.</p>
Making data interoperable	<p><b>Interoperability:</b> The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><b>Data and metadata vocabularies:</b> No data or metadata vocabularies are used, since content is plain text and no metadata are available.</p> <p><b>Use of standard vocabularies:</b> No.</p> <p><b>Mappings to commonly used vocabularies:</b> No.</p>
Increase data re-use	<p><b>Licence:</b> The data is already openly shared in Zenodo under a Creative Commons Attribution 4.0 International license.</p> <p><b>Availability for re-use:</b> Data are already available for re-use.</p> <p><b>Usable by third parties after end of project:</b> Yes.</p> <p><b>Re-use timeframe:</b> Unlimited.</p> <p><b>Data quality assurance process:</b> The data has not been subjected to a data quality assurance process.</p>
Allocation of resources	<p><b>Costs for making data FAIR:</b> No.</p> <p><b>Costs for long-term preservation:</b> No.</p>
Data security	<p><b>Security measures:</b> Data is released with CC licence, no security measures for data protection have been implemented.</p>
Ethical aspects	<p><b>Possible ethical and legal aspects preventing sharing:</b> As the creators of this dataset, CNR did not foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><b>Is informed consent for data sharing and long term preservation given:</b> N/A</p>
Other Issues	N/A



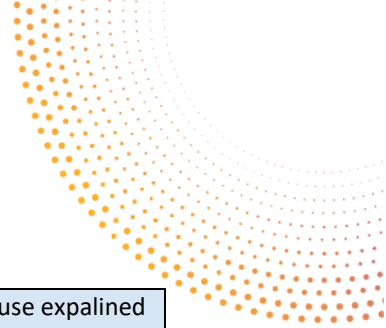


### 5.1.10 OHSUMED MEDLINE document dataset

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_71_WP3-5_TEXT_OHSUMED_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: The OHSUMED dataset consists of a set of about 348,000 MEDLINE documents spanning the years from 1987 to 1991, and classified according to a set of classes representing disease, which describe what the document is about. The dataset is used in AI4media (and has been used elsewhere since the mid '90s) as a benchmark (training set + test set) for testing text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format</u>: Raw text</p> <p><u>Re-use of existing data</u>: Yes, this dataset has been in the public domain since the mid '90s.</p> <p><u>Data origin</u>: The documents consist of title+abstract+metadata from scientific articles available from the MEDLINE service. Ohio State University (OSU) released this dataset to the public in the mid '90s.</p> <p><u>Expected size</u>: About 348,000 documents.</p> <p><u>Data utility</u>: It is, and it has been for decades, useful to text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data have been made available by their creator at <a href="https://dmice.ohsu.edu/hersh/ohsumed/index1.html">https://dmice.ohsu.edu/hersh/ohsumed/index1.html</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords</u>: No search keywords provided.</p> <p><u>Versioning</u>: There is only one version which has been made available by its creator.</p> <p><u>Metadata creation</u>: No metadata were attached to this dataset by its creator, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The creator of the dataset makes it openly accessible at <a href="https://dmice.ohsu.edu/hersh/ohsumed/index1.html">https://dmice.ohsu.edu/hersh/ohsumed/index1.html</a>. We will not reshare the data.</p> <p><u>How it will be accessible</u>: The dataset has been accessible from its creator's home page ever since the mid '90s.</p> <p><u>Methods/software tools to access data</u>: The only software tool needed to access the data is a web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>





Increase data re-use	<p><u>Licence</u>: The data is already openly shared by the creator under terms of use explained at <a href="https://dmice.ohsu.edu/hersh/ohsumed/index1.html">https://dmice.ohsu.edu/hersh/ohsumed/index1.html</a>.</p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.11 RCV1 Reuters stories dataset

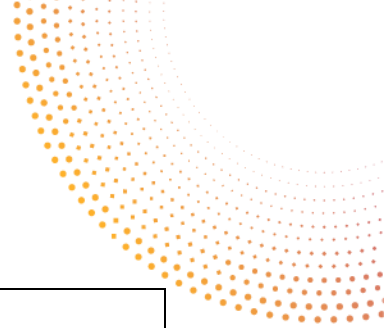
This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_72_WP3-5_TEXT_RCV1-Reuters_v1 Partner: CNR
Data Summary	<p><u>Purpose</u>: The RCV1-v2 dataset consists of about 800,000 news stories written by Reuters journalists, and classified according to a set of 101 classes related to economics, which describe what the news story is about. The dataset is used in AI4media (and has been used elsewhere since the mid 2000's) as a benchmark (training set + test set) for testing text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format</u>: Raw text</p> <p><u>Re-use of existing data</u>: Yes, this dataset has been in the public domain since the mid 2000's.</p> <p><u>Data origin</u>: The stories were written by Reuters journalists. Reuters released this dataset to the public in the mid 2000's.</p> <p><u>Expected size</u>: about 800,000 documents.</p> <p><u>Data utility</u>: It is, and it has been for decades, useful to text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data have been made available by its creator on his home page at <a href="http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm">http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm</a>, from where it can be obtained by signing an agreement; however, the data in preprocessed (matrix) form can also be simply downloaded, without signing any</p>



	<p>agreement. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords</u>: No search keywords provided.</p> <p><u>Versioning</u>: There is only one version which has been made available by its creator.</p> <p><u>Metadata creation</u>: No metadata were attached to this dataset by its creator, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The creator of the dataset makes it openly accessible at <a href="http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm">http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm</a>. We will not reshare the data.</p> <p><u>How it will be accessible</u>: The dataset has been accessible from its creator's home page ever since the mid 2000's.</p> <p><u>Methods/software tools to access data</u>: The only software tool needed to access the data is a web browser.</p> <p><u>Repository</u>: Data available at <a href="http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm">http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm</a></p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared by their creator without a license. However, those who want to acquire the data are encouraged to sign an agreement: <a href="http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm">http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_REA_DME.htm</a></p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing. The data are made available in preprocessed (matrix) form, a form from which the original text cannot be reconstructed. Reuters personnel have stated that</p>





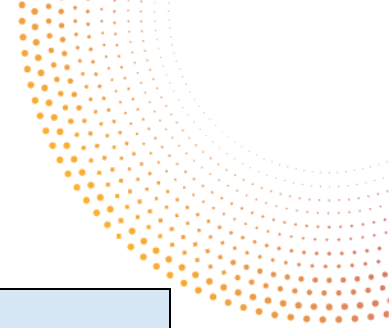
	distributing term/document matrices is not a violation of the Agreement. <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 5.1.12 RCV1RCV2 Reuters stories dataset

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_73_WP3-5_TEXT_RCV1RCV2-Reuters_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> The RCV1 RCV2 dataset consists of about 111,700 news stories written by Reuters journalists, and classified according to a set of 101 classes related to economics, which describe what the news story is about; the data are multilingual, i.e., each news story is written in one of 5 different languages. The dataset is used in AI4media (and has been used elsewhere since year 2000) as a benchmark (training set + test set) for testing multilingual text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format:</u> Raw text</p> <p><u>Re-use of existing data:</u> Yes, this dataset has been in the public domain since the mid 2000's.</p> <p><u>Data origin:</u> The stories were written by Reuters journalists. Reuters released this dataset to the public in the mid 2000's.</p> <p><u>Expected size:</u> about 111,700 documents.</p> <p><u>Data utility:</u> It is, and it has been for decades, useful to multilingual text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their multilingual text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data have been made available by its creators from the UCI Machine Learning Repository at <a href="https://archive.ics.uci.edu/ml/datasets/Reuters+RCV1+RCV2+Multilingual%2C+Multi+ew+Text+Categorization+Test+collection">https://archive.ics.uci.edu/ml/datasets/Reuters+RCV1+RCV2+Multilingual%2C+Multi+ew+Text+Categorization+Test+collection</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords:</u> No search keywords provided.</p> <p><u>Versioning:</u> There is only one version which has been made available by its creators.</p> <p><u>Metadata creation:</u> No metadata were attached to this dataset by its creators, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The dataset is already openly accessible at <a href="https://archive.ics.uci.edu/ml/datasets/Reuters+RCV1+RCV2+Multilingual%2C+Multi+ew+Text+Categorization+Test+collection">https://archive.ics.uci.edu/ml/datasets/Reuters+RCV1+RCV2+Multilingual%2C+Multi+ew+Text+Categorization+Test+collection</a> by its creator. We will not reshare the data.</p> <p><u>How it will be accessible:</u> The dataset has been accessible from the UCI ML repository ever since 2013.</p> <p><u>Methods/software tools to access data:</u> The only software tool needed to access the data is a web browser.</p>





	<p><u>Repository</u>: The data is deposited in the UCI ML repository.</p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared by their creator without a license. Users of the dataset should acknowledge its use, by referring to: <i>M.-R. Amini, N. Usunier, C. Goutte. Learning from Multiple Partially Observed Views - an Application to Multilingual Text Categorization. Advances in Neural Information Processing Systems 22, p. 28-36, 2009</i></p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing. The data are made available in preprocessed (matrix) form, a form from which the original text cannot be reconstructed. Reuters personnel have stated that distributing term/document matrices is not a violation of the Agreement.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

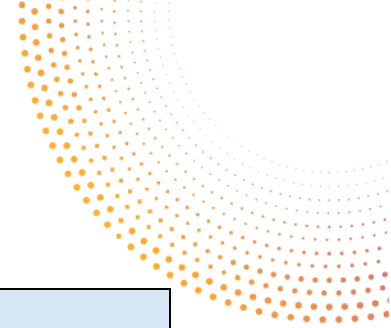
### 5.1.13 Reuters-21578 dataset

This dataset will also be used in the context of WP5.

<b>DMP component</b>	<b>AI4Media_Data_74_WP3-5_TEXT_Reuters-21578_v1</b> <b>Partner: CNR</b>
Data Summary	<p><u>Purpose</u>: The Reuters-21578 dataset consists of 12,904 news stories written by Reuters journalists in the late '90s, and classified according to a set of 115 classes related to economics (e.g., "acquisitions", "interest rates"), which describe what the news story is about. The dataset is used in AI4media (and has been used elsewhere since the '90s) as a benchmark (training set + test set) for testing text classification systems, e.g., in the context of T3.7 and T5.4.</p>

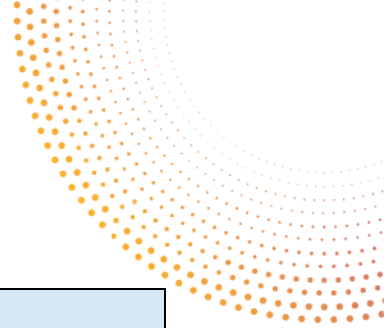






	<p><u>Type/format</u>: Raw text</p> <p><u>Re-use of existing data</u>: Yes, this dataset has been in the public domain since the '90s.</p> <p><u>Data origin</u>: The stories were written by Reuters journalists. Reuters released this dataset to the public in the mid '90s.</p> <p><u>Expected size</u>: 12,904 documents.</p> <p><u>Data utility</u>: It is, and it has been for decades, useful to text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their text classification systems.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: The data have been made available by its creator on his home page and on the UCI machine learning repository at <a href="https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection">https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><u>Search keywords</u>: No search keywords provided.</p> <p><u>Versioning</u>: There is only one version which has been made available by its creator.</p> <p><u>Metadata creation</u>: No metadata were attached to this dataset by its creator, aside from the set of classes that is to be used for labelling the documents.</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The dataset is already openly accessible at <a href="https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection">https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection</a> by its creator. We will not reshare the data.</p> <p><u>How it will be accessible</u>: The dataset has been accessible from its creator's home page ever since the '90s.</p> <p><u>Methods/software tools to access data</u>: The only software tool needed to access the data is a web browser.</p> <p><u>Repository</u>: The data is deposited in the UCI machine learning repository at <a href="https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection">https://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection</a></p> <p><u>Restrictions on access</u>: There are no restrictions on the use of this dataset.</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
<p>Increase data re-use</p>	<p><u>Licence</u>: The copyright for the text of newswire articles and Reuters annotations in the Reuters-21578 collection resides with Reuters Ltd. Reuters Ltd. and Carnegie Group, Inc. have agreed to allow the free distribution of this data *for research purposes only*.</p> <p>Users of the dataset should acknowledge its use, refer to the data set by the name "Reuters-21578, Distribution 1.0", and inform of the current location of the dataset.</p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p>





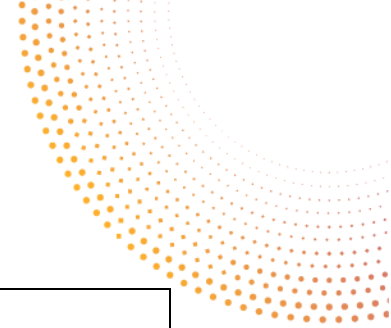
	<p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 5.1.14 11 Tweet Sentiment datasets

This dataset will also be used in the context of WP5.

DMP component	AI4Media_Data_75_WP3-5_TEXT_11TweetSentiment_v1 Partner: CNR
Data Summary	<p><u>Purpose:</u> The 11 Tweet Sentiment Datasets is a set of 11 datasets (called GASP, HCS, OMD, Sanders, SemEval2013, SemEval2014,, SemEval2015, SemEval2016, SST, WA, WB, respectively), all of a similar nature, often used all together for experimentation purposes, consisting of tweets classified according to the Positive, Neutral, Negative, sentiment-based classes. The datasets are used in AI4media (and have been used elsewhere in the last ten years) as benchmarks (training sets + test sets) for testing sentiment classification systems or sentiment quantification systems, e.g., in the context of T3.7 and T5.4.</p> <p><u>Type/format:</u> Vectors of features extracted from text</p> <p><u>Re-use of existing data:</u> Yes, these datasets have been in the public domain for 5 years or more.</p> <p><u>Data origin:</u> The data consist of posts crawled from Twitter by several authors; they are in the form of feature vectors so as to comply with the Twitter terms of use.</p> <p><u>Expected size:</u> Approximately 20,000 tweets altogether.</p> <p><u>Data utility:</u> These datasets are useful to sentiment classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their sentiment classification systems.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The datasets are available, among other places, on Zenodo at <a href="https://zenodo.org/record/4255764">https://zenodo.org/record/4255764</a>. In every paper we write, we indicate the URL from where the datasets can be downloaded. This URL can be located by just typing "tweet sentiment classification datasets Zenodo" into any web search engine.</p> <p><u>Search keywords:</u> No search keywords provided.</p>



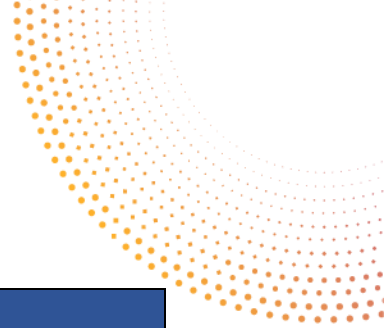


	<p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: No metadata are attached to these datasets, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The datasets are openly accessible on Zenodo at <a href="https://zenodo.org/record/4255764">https://zenodo.org/record/4255764</a>. We will not reshare the datasets.</p> <p><u>How it will be accessible</u>: The datasets are openly accessible on Zenodo, an open-access repository.</p> <p><u>Methods/software tools to access data</u>: The only software tool needed to access the data is a web browser.</p> <p><u>Repository</u>: The data are deposited in the Zenodo repository at <a href="https://zenodo.org/record/4255764">https://zenodo.org/record/4255764</a></p> <p><u>Restrictions on access</u>: There are no restrictions on the use of these datasets.</p>
Making data interoperable	<p><u>Interoperability</u>: The datasets consist of interoperable data, for the simple fact that they consist of raw text.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared on Zenodo under a <a href="#">Creative Commons Attribution 4.0 International</a> license.</p> <p><u>Availability for re-use</u>: Data are already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: Without limits.</p> <p><u>Re-use timeframe</u>: Perpetual</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The authors who made the datasets available did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing, since the tweets that the datasets consist of are made available in the form of feature vectors only, which means that the tweets in their original form cannot be recovered.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.15 WipoGamma patent document dataset

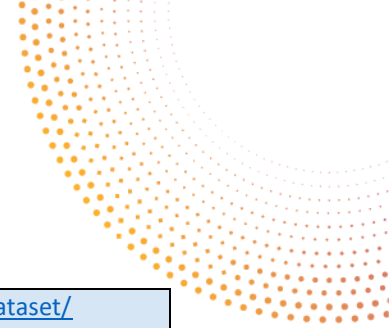
This dataset will also be used in the context of WP5.





DMP component	AI4Media_Data_76_WP3-5_TEXT_WipoGamma_v1 Partner: CNR
Data Summary	<p><b>Purpose:</b> The WipoGamma dataset consists of about 1,100,000 patent documents made available by the World Intellectual Property Organization (WIPO), and classified according to classes representing sectors and subsectors of technology, which describe what the patent document is about. The dataset is used in AI4media (and has been used elsewhere) as a benchmark (training set + test set) for testing text classification systems, e.g., in the context of T3.7 and T5.4.</p> <p><b>Type/format:</b> Raw text</p> <p><b>Re-use of existing data:</b> Yes, this dataset has been in the public domain for years.</p> <p><b>Data origin:</b> The dataset consists of patent documents made available by the World Intellectual Property Organization.</p> <p><b>Expected size:</b> About 1,100,000 documents.</p> <p><b>Data utility:</b> It is, and it has been for years, useful to text classification researchers (see e.g., T3.7 and T5.4) wishing to benchmark their text classification systems.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> The data have been made available by WIPO on their website, at <a href="https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/">https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/</a>. In every paper we write, we indicate the URL from where the dataset can be downloaded. This URL can be located by just typing the dataset name into any web search engine.</p> <p><b>Search keywords:</b> No search keywords provided.</p> <p><b>Versioning:</b> There is only one version which has been made available by its creators.</p> <p><b>Metadata creation:</b> No metadata were attached to this dataset by its creators, aside from the set of classes that is to be used for labelling the documents.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The dataset is already openly accessible at <a href="https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/">https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/</a> by WIPO. We will not reshare the data.</p> <p><b>How it will be accessible:</b> The dataset has been accessible from the WIPO website after filling in a registration form: <a href="https://www.wipo.int/classifications/ipc/en/forms/index.html">https://www.wipo.int/classifications/ipc/en/forms/index.html</a></p> <p><b>Methods/software tools to access data:</b> The only software tool needed to access the data is a web browser.</p> <p><b>Repository:</b> <a href="https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/">https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/</a></p> <p><b>Restrictions on access:</b> There are no restrictions on the use of this dataset.</p>
Making data interoperable	<p><b>Interoperability:</b> The dataset consists of interoperable data, for the simple fact that it consists of raw text.</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> The data is already openly shared by WIPO under the terms of use defined at</p>



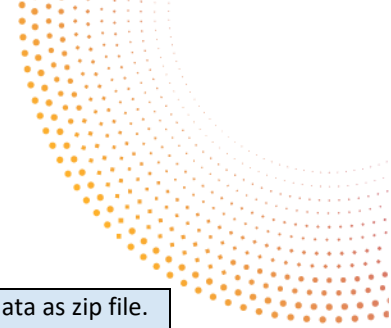


	<p><a href="https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/">https://www.wipo.int/classifications/ipc/en/ITsupport/Categorization/dataset/</a></p> <p><u>Availability for re-use</u>: Data is already available for re-use.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CNR servers. CNR fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The creator of the dataset did not, in all evidence, foresee any ethical or legal issues that can have an impact on data sharing.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.16 WIDER FACE face detection dataset

DMP component	AI4Media_Data_77_WP3_Image_WIDER_FACE_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: WIDER FACE is a public face detection benchmark dataset, of which images are selected from the publicly available WIDER dataset. The dataset features a high degree of variability in scale, pose and occlusion in the depicted faces. It is used by CERTH to evaluate the face detection models in WP3.</p> <p><u>Type/format</u>: jpg</p> <p><u>Re-use of existing data</u>: Yes</p> <p><u>Data origin</u>: The dataset contains more than 32K images available at <a href="http://shuoyang1213.me/WIDERFACE/">http://shuoyang1213.me/WIDERFACE/</a>.</p> <p><u>Expected size</u>: 1.8 GB</p> <p><u>Data utility</u>: It is useful to WP3 partners to evaluate and benchmark face detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable and publicly available.</p> <p><u>Search keywords</u>: wider face dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="http://shuoyang1213.me/WIDERFACE/">http://shuoyang1213.me/WIDERFACE/</a>.</p> <p><u>How it will be accessible</u>: Shared through a third-party repository link.</p>



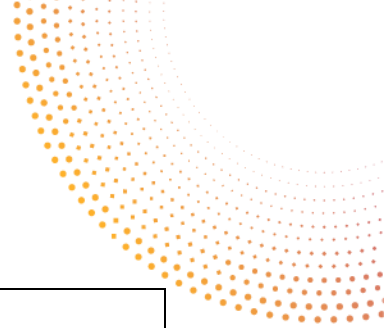


	<p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: <a href="http://shuoyang1213.me/">http://shuoyang1213.me/</a></p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under Creative Common License.</p> <p><u>Availability for re-use</u>: The dataset is already publicly available.</p> <p><u>Usable by third parties after end of project</u>: Yes.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: This dataset is downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.17 Caltech 101 dataset

DMP component	AI4Media_Data_78_WP3_IMAGE_Caltech101_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: This public dataset contains images of 101 different types of objects, with boundary annotations. Each image has a class label associated. It is used for training models for image classification or object segmentation tasks.</p> <p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The dataset is a set of images collected via Google Image search.</p> <p><u>Expected size</u>: 151MB total</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification or object segmentation tasks.</p>
Making data	<u>Is data discoverable</u> : The data can be acquired through the website



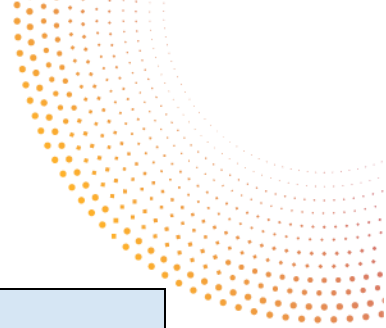


findable, incl. provisions for metadata	<p><a href="https://data.caltech.edu/records/mzrjq-6wc02">https://data.caltech.edu/records/mzrjq-6wc02</a>.</p> <p><u>Search keywords</u>: caltech 101</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://data.caltech.edu/records/mzrjq-6wc02">https://data.caltech.edu/records/mzrjq-6wc02</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: <a href="https://data.caltech.edu/records/mzrjq-6wc02">https://data.caltech.edu/records/mzrjq-6wc02</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution 4.0 International. Allows re-distribution and re-use of the work on the condition that the creator is appropriately credited.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.18 CUB-200-2021 image dataset

<b>DMP component</b>	<b>AI4Media_Data_79_WP3_IMAGE_CUB-200-2011_v1</b> <b>Partner: BSC</b>
Data Summary	<p><u>Purpose</u>: This public dataset contains 11,788 images of 200 bird species, each associated with a Wikipedia article and organized by taxonomic classification. Annotations in the form of image locations are included for several body parts of the birds.</p>

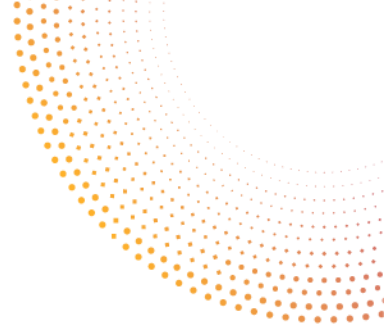




	<p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The dataset is a collection of images collected via Flickr Image search.</p> <p><u>Expected size</u>: Around 2GB total</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification or object detection/location tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data can be acquired at the website <a href="https://www.vision.caltech.edu/datasets/cub_200_2011/">https://www.vision.caltech.edu/datasets/cub_200_2011/</a>.</p> <p><u>Search keywords</u>: cub200, cub-200-2011</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://www.vision.caltech.edu/datasets/cub_200_2011/">https://www.vision.caltech.edu/datasets/cub_200_2011/</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: <a href="https://www.vision.caltech.edu/datasets/cub_200_2011/">https://www.vision.caltech.edu/datasets/cub_200_2011/</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: No license found.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A



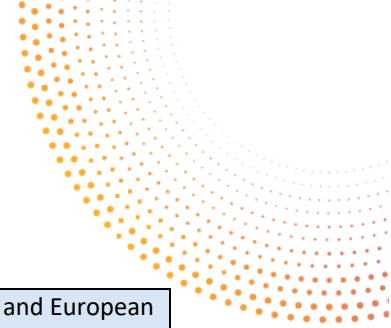




### 5.1.19 Describable Textures dataset

DMP component	AI4Media_Data_80_WP3_IMAGE_Describable-Textures_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: The Describable Textures Dataset contains 5,640 textured images, organized according to a list of 47 terms, inspired from human perception. The images contain at least 90% of the surface representing the category attribute. Each image has a class label associated.</p> <p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The dataset is a collection of images collected via Google and Flickr Image searches.</p> <p><u>Expected size</u>: Around 1.1GB total</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data can be acquired at the website <a href="https://www.robots.ox.ac.uk/~vgg/data/dtd/">https://www.robots.ox.ac.uk/~vgg/data/dtd/</a>.</p> <p><u>Search keywords</u>: describable textures dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://www.robots.ox.ac.uk/~vgg/data/dtd/">https://www.robots.ox.ac.uk/~vgg/data/dtd/</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as tar.gz file.</p> <p><u>Repository</u>: <a href="https://www.robots.ox.ac.uk/~vgg/data/dtd/">https://www.robots.ox.ac.uk/~vgg/data/dtd/</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Data only available for research purposes.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will</p>





	be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 5.1.20 Food-101 dataset

DMP component	AI4Media_Data_81_WP3_IMAGE_Food101_v1 Partner: BSC
Data Summary	<u>Purpose:</u> The Food-101 dataset contains 101,000 images of food dishes, separated in 101 food categories. Each image has a class label associated. <u>Type/format:</u> Each image is stored as a .jpg file. <u>Re-use of existing data:</u> Yes. <u>Data origin:</u> The images that form the dataset come from individual posts on Foodspotting, a (now defunct) social network for food picture sharing. <u>Expected size:</u> 8.6GB total <u>Data utility:</u> It is useful to WP3 partners working in image classification tasks.
Making data findable, incl. provisions for metadata	<u>Is data discoverable:</u> The data can be acquired at the website <a href="https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/">https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/</a> . <u>Search keywords:</u> food 101 dataset <u>Versioning:</u> N/A <u>Metadata creation:</u> N/A
Making data openly accessible	<u>Data openly accessible:</u> The data is accessible from <a href="https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/">https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/</a> . It is a public dataset. <u>How it will be accessible:</u> The dataset is accessible via direct download. <u>Methods/software tools to access data:</u> Web-browser to download the data as tar.gz file. <u>Repository:</u> <a href="https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/">https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/</a> <u>Restrictions on access:</u> None.
Making data interoperable	<u>Interoperability:</u> N/A <u>Data and metadata vocabularies:</u> N/A <u>Use of standard vocabularies:</u> N/A <u>Mappings to commonly used vocabularies:</u> N/A
Increase data re-use	<u>Licence:</u> Data available for research purposes. Other uses may be negotiated with the original owners of each picture. <u>Availability for re-use:</u> N/A <u>Usable by third parties after end of project:</u> The data is openly available and is subject

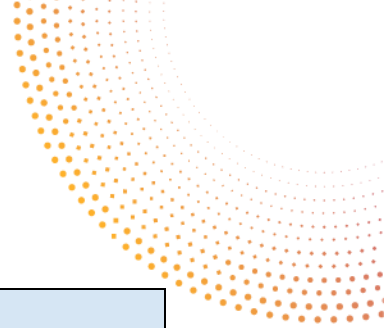


	<p>to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.21 MAMe dataset

DMP component	AI4Media_Data_82_WP3_IMAGE_MAMe_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: The dataset contains images of pieces of art from three museums, organized in 29 classes representing different mediums and techniques the art pieces are made on/with. Each image has a class label associated.</p> <p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The dataset is a collection of images published by the Metropolitan Museum of Art of New York, the Los Angeles County Museum of Art and the Cleveland Museum of Art. All images are under a CC0 license.</p> <p><u>Expected size</u>: The low-resolution version of the dataset (all images resized to 256x256px) takes 765MB.</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data can be acquired at the website <a href="https://hpa.bsc.es/MAMe-dataset/">https://hpa.bsc.es/MAMe-dataset/</a>.</p> <p><u>Search keywords</u>: mame dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://hpa.bsc.es/MAMe-dataset/">https://hpa.bsc.es/MAMe-dataset/</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as a zip file.</p> <p><u>Repository</u>: <a href="https://hpa.bsc.es/MAMe-dataset/">https://hpa.bsc.es/MAMe-dataset/</a></p>



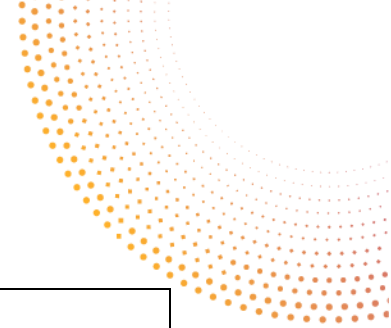


	<u>Restrictions on access</u> : None.
Making data interoperable	<u>Interoperability</u> : N/A <u>Data and metadata vocabularies</u> : N/A <u>Use of standard vocabularies</u> : N/A <u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<u>Licence</u> : CC BY-NC-ND 4.0 <u>Availability for re-use</u> : N/A <u>Usable by third parties after end of project</u> : The data is openly available and is subject to licenses from creators. <u>Re-use timeframe</u> : N/A <u>Data quality assurance process</u> : N/A
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : N/A <u>Is informed consent for data sharing and long term preservation given</u> : N/A
Other Issues	N/A

### 5.1.22 MIT-ISR dataset

DMP component	AI4Media_Data_83_WP3_IMAGE_MIT-ISR_v1 Partner: BSC
Data Summary	<u>Purpose</u> : The dataset contains 15,620 images representing 67 different types of indoor scenes. Each image has a class label associated. A subset of the images are also segmented and annotated with the objects that they contain. <u>Type/format</u> : Each image is stored as a .jpg file. <u>Re-use of existing data</u> : Yes. <u>Data origin</u> : The dataset is a collection of images obtained via Google, AltaVista and Flickr image searches, and from the LabelMe dataset. <u>Expected size</u> : 6.5GB total <u>Data utility</u> : It is useful to WP3 partners working in image classification tasks and object segmentation/recognition tasks.
Making data findable, incl. provisions for metadata	<u>Is data discoverable</u> : The data can be acquired at the website <a href="https://web.mit.edu/torralba/www/indoor.html">https://web.mit.edu/torralba/www/indoor.html</a> . <u>Search keywords</u> : mit indoor scenes dataset <u>Versioning</u> : N/A



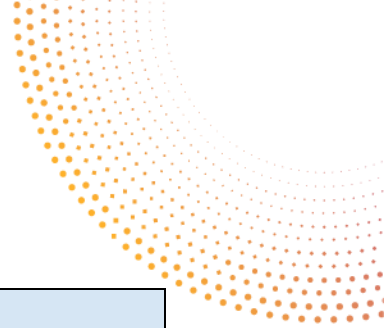


	<u>Metadata creation</u> : N/A
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://web.mit.edu/torralba/www/indoor.html">https://web.mit.edu/torralba/www/indoor.html</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as a tar file.</p> <p><u>Repository</u>: <a href="https://web.mit.edu/torralba/www/indoor.html">https://web.mit.edu/torralba/www/indoor.html</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Research purposes only.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.1.23 Oulu Knots dataset

<b>DMP component</b>	<b>AI4Media_Data_84_WP3_IMAGE_OuluKnots_v1</b> <b>Partner: BSC</b>
Data Summary	<p><u>Purpose</u>: The Oulu Knots dataset contains images depicting wood and various defects on it. Each image has a class label associated.</p> <p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The images in the dataset were produced by members of the University of Oulu.</p>



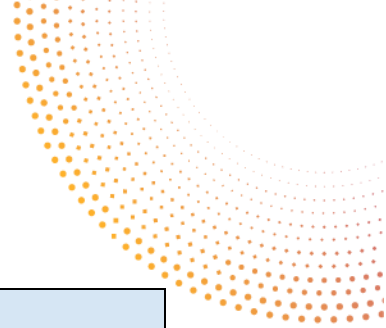


	<p><u>Expected size</u>: 11MB</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data appears to not be reachable currently.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: This is a public dataset, although it appears to not be reachable currently.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: No license found.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is not available.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset was downloaded from the original source and was stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 5.1.24 Oxford Flower dataset

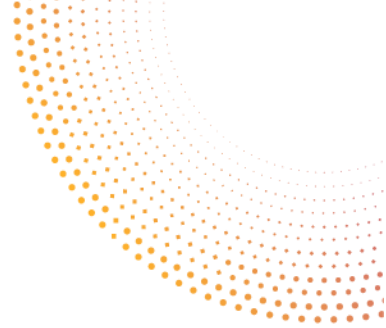
<b>DMP component</b>	<b>AI4Media_Data_85_WP3_IMAGE_OxfordFlowers_v1</b> <b>Partner: BSC</b>
Data Summary	<u>Purpose</u> : The Oxford Flower dataset contains images depicting 102 types of flowers. Each image has a class label and a segmentation annotation.





	<p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The images in the dataset are a mix of own production from the Visual Geometry Group at University of Oxford and images collected from the web.</p> <p><u>Expected size</u>: Around 0.7GB total</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification tasks and object segmentation tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data can be acquired at the website <a href="https://www.robots.ox.ac.uk/~vgg/data/flowers/102/">https://www.robots.ox.ac.uk/~vgg/data/flowers/102/</a>.</p> <p><u>Search keywords</u>: flowers 102, oxford flowers dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://www.robots.ox.ac.uk/~vgg/data/flowers/102/">https://www.robots.ox.ac.uk/~vgg/data/flowers/102/</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as a tgz file.</p> <p><u>Repository</u>: <a href="https://www.robots.ox.ac.uk/~vgg/data/flowers/102/">https://www.robots.ox.ac.uk/~vgg/data/flowers/102/</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: No license found.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A



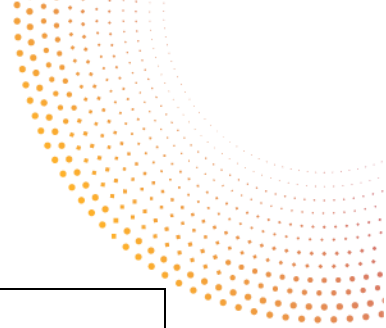


### 5.1.25 Oxford-IIIT Pet dataset

DMP component	AI4Media_Data_86_WP3_IMAGE_OxfordIIITPet_v1 Partner: BSC
Data Summary	<p><u>Purpose</u>: The dataset contains images depicting 37 breeds of cats and dogs. Each image has a class label, bounding box of the animal's head and a segmentation annotation.</p> <p><u>Type/format</u>: Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The images in the dataset have been collected from Catster, Dogster, Flickr and Google image search.</p> <p><u>Expected size</u>: Around 1.3GB total</p> <p><u>Data utility</u>: It is useful to WP3 partners working in image classification tasks and object segmentation tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data can be acquired at the website <a href="https://www.robots.ox.ac.uk/~vgg/data/pets/">https://www.robots.ox.ac.uk/~vgg/data/pets/</a>.</p> <p><u>Search keywords</u>: oxford iiit pet</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is accessible from <a href="https://www.robots.ox.ac.uk/~vgg/data/pets/">https://www.robots.ox.ac.uk/~vgg/data/pets/</a>. It is a public dataset.</p> <p><u>How it will be accessible</u>: The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as a tar.gz file.</p> <p><u>Repository</u>: <a href="https://www.robots.ox.ac.uk/~vgg/data/pets/">https://www.robots.ox.ac.uk/~vgg/data/pets/</a></p> <p><u>Restrictions on access</u>: None.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Creative Commons Attribution-ShareAlike 4.0 International License</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p>







	<u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 5.1.26 Stanford Dogs dataset

DMP component	AI4Media_Data_87_WP3_IMAGE_StanfordDogs_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> The dataset contains images depicting 120 breeds of dogs. Each image has a class label and bounding box.</p> <p><u>Type/format:</u> Each image is stored as a .jpg file.</p> <p><u>Re-use of existing data:</u> Yes.</p> <p><u>Data origin:</u> The images in the dataset come from the ImageNet1K dataset.</p> <p><u>Expected size:</u> Around 1.2GB total</p> <p><u>Data utility:</u> It is useful to WP3 partners working in image classification tasks and object detection tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data can be acquired at the website <a href="http://vision.stanford.edu/aditya86/ImageNetDogs/">http://vision.stanford.edu/aditya86/ImageNetDogs/</a>.</p> <p><u>Search keywords:</u> oxford iiit pet</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is accessible from <a href="http://vision.stanford.edu/aditya86/ImageNetDogs/">http://vision.stanford.edu/aditya86/ImageNetDogs/</a>. It is a public dataset.</p> <p><u>How it will be accessible:</u> The dataset is accessible via direct download.</p> <p><u>Methods/software tools to access data:</u> Web-browser to download the data as a tar.gz file.</p> <p><a href="http://vision.stanford.edu/aditya86/ImageNetDogs/">http://vision.stanford.edu/aditya86/ImageNetDogs/</a></p> <p><u>Restrictions on access:</u> None.</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<u>Licence:</u> No license found.



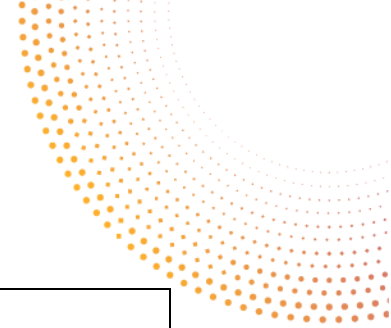
	<p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	<p>Refer to other national/funder/sectorial/departmental procedures for data management that you may be using (if any)</p>

## 5.2 Datasets used in the context of WP4

### 5.2.1 ImageNet-ILSVRC2012 image classification dataset

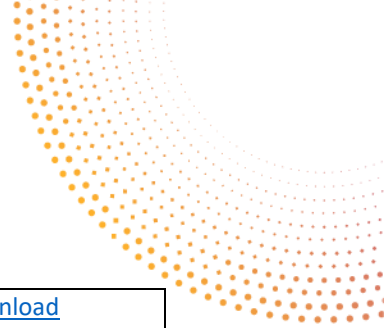
DMP component	AI4Media_Data_88_WP4_IMAGE_Imagenet_01 Partner: CERTH
Data Summary	<p><u>Purpose</u>: ImageNet is an image dataset organized according to the WordNet hierarchy: images annotated with concept labels. ImageNet is among the most popular large-scale image dataset for image semantic concept classification tasks. In this version of the dataset, we use a subset of ImageNet, the so called Large Scale Visual Recognition Challenge 2012 (ILSVRC2012) dataset. ILSVRC2012 contains approximately 1.3 million images, associated object bounding boxes, and 1,000 semantic concept categories. It will be used by CERTH in T4.3 for evaluating the explainable AI methods developed in this task.</p> <p><u>Type/format</u>: jpeg</p> <p><u>Re-use of existing data</u>: Yes</p> <p><u>Data origin</u>: ImageNet project site: <a href="http://www.image-net.org/index">http://www.image-net.org/index</a></p> <p><u>Expected size</u>: Train partition: 137 GB; Validation partition: 6.28 GB; Test partition: 12.7 GB.</p> <p><u>Data utility</u>: It is useful in the context of T4.3 for evaluating the XAI methods developed in this task. In general, this dataset is also useful for any researcher that wants to train deep learning models for image classification/localization using a large-scale image dataset.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted on the ImageNet project site (<a href="http://image-net.org/download">http://image-net.org/download</a>). It is discoverable by googling "ImageNet Dataset".</p> <p><u>Search keywords</u>: N/A</p>





	<p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No. The data is accessible through registration on the ImageNet project site (<a href="http://image-net.org/download">http://image-net.org/download</a>) under certain restrictions. The data will not be re-shared by AI4Media partners.</p> <p><u>How it will be accessible</u>: The data is hosted on ImageNet project site and requires registration to download. The images are provided for non-commercial research and/or educational purposes under certain conditions and terms. Details are provided on: <a href="http://image-net.org/download">http://image-net.org/download</a></p> <p><u>Methods/software tools to access data</u>: Creation of an ImageNet account as described on: <a href="http://image-net.org/signup.php?next=download-images">http://image-net.org/signup.php?next=download-images</a></p> <p><u>Repository</u>: The data repository (data, metadata, documentation, processing code) is hosted on the ImageNet project site.</p> <p><u>Restrictions on access</u>: ImageNet does not own the copyright of the images. ImageNet provides the images for non-commercial research and/or educational purposes under certain conditions and terms. The users have to sign up for an ImageNet Account.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable and widely used in the research community.</p> <p><u>Data and metadata vocabularies</u>: Images are in jpeg format. All synsets are assigned to an integer ID between 1 and 1000 (ILSVRC2012_ID). Moreover, its synset has a WordNet ID (WNID) used to uniquely identify a synset in ImageNet or WordNet</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: A mapping to WordNet is already provided.</p>
Increase data re-use	<p><u>Licence</u>: A registration is required to download the dataset. ImageNet does not own the copyright of the images. However, it provides the images for non-commercial research and/or educational purposes under certain conditions and terms. Details are provided on: <a href="http://image-net.org/download">http://image-net.org/download</a></p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The images of the dataset are provided for non-commercial research and/or educational purposes under certain</p>



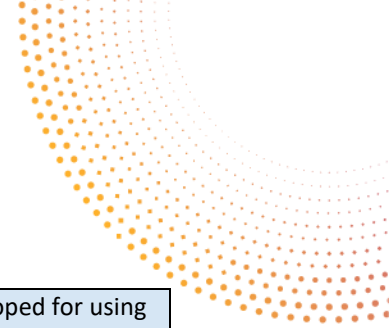


	<p>conditions and terms. Details are provided on: <a href="http://image-net.org/download">http://image-net.org/download</a></p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A (Note that ImageNet does not own the copyright of the images)</p>
Other Issues	N/A

## 5.2.2 FFHQ dataset for GAN training

DMP component	AI4Media_Data_89_WP4_IMAGE_FFHQ_v1 Partner: CEA
Data Summary	<p><u>Purpose:</u> Flickr-Faces-HQ (FFHQ) is a high-quality image dataset of human faces, originally created as a benchmark for generative adversarial networks (GAN). We will use the dataset stored as multi-resolution TF records. It will be used in T4.3 to evaluate the generative models developed in the task. Results involving this dataset will be reported in D4.1, D4.4 and D4.6.</p> <p><u>Type/format:</u> JSON file containing metadata and 70k images stored as multi-resolution TF records.</p> <p><u>Re-use of existing data:</u> Yes, were reusing an existing dataset</p> <p><u>Data origin:</u> <a href="https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html">https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html</a></p> <p><u>Expected size:</u> 575 MB</p> <p><u>Data utility:</u> It is useful to WP4 partners to benchmark generative models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable. The dataset is hosted on the NVIDIA website (actually Google drive of NVIDIA).</p> <p><u>Search keywords:</u> FFHQ, Flickr-Faces-HQ</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is already openly accessible at <a href="https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html">https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html</a>. We will not reshare the data.</p> <p><u>How it will be accessible:</u> Already shared through a third-party repository link</p> <p><u>Methods/software tools to access data:</u> python scripts to download data are provided</p> <p><u>Repository:</u> Network Repository (<a href="http://networkrepository.com">http://networkrepository.com</a>)</p> <p><u>Restrictions on access:</u> None</p>
Making data interoperable	<p><u>Interoperability:</u> The data is interoperable.</p> <p><u>Data and metadata vocabularies:</u> The metadata schema can be found at <a href="https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html#articleHeader4">https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html#articleHeader4</a></p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The dataset is publicly available under an attribution licence (<a href="https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html">https://reposhub.com/python/deep-learning/NVlabs-ffhq-dataset.html</a>)</p>





	<p><u>Availability for re-use</u>: The loading and pre-processing mechanism developed for using the dataset in experiments will be made publicly available to ensure reproducibility of research.</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on CEA's servers. CEA fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The dataset was made available from FlickrR images with appropriate licence. It will be kept only during the length of the T4.3 and will not be re-shared since it is available on NVIDIA website.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.2.3 MNIST image dataset

DMP component	AI4Media_Data_90_WP4_IMAGE_MNIST_v1 Partner: IBM
Data Summary	<p><u>Purpose</u>: The MNIST database of handwritten digits has a training set of 60,000 examples, and a test set of 10,000 examples.</p> <p><u>Type/format</u>: Images</p> <p><u>Re-use of existing data</u>: Yes, were reusing an existing dataset</p> <p><u>Data origin</u>: <a href="https://tensorflow.google.cn/datasets/catalog/mnist?hl=en">https://tensorflow.google.cn/datasets/catalog/mnist?hl=en</a></p> <p><u>Expected size</u>: 21 MB</p> <p><u>Data utility</u>: It is useful to WP4 partners to benchmark generative AI models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted on the Google Tensorflow website at <a href="https://tensorflow.google.cn/datasets/catalog/mnist?hl=en">https://tensorflow.google.cn/datasets/catalog/mnist?hl=en</a></p> <p><u>Search keywords</u>: MNIST</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="https://tensorflow.google.cn/datasets/catalog/mnist?hl=en">https://tensorflow.google.cn/datasets/catalog/mnist?hl=en</a> , <a href="http://yann.lecun.com/exdb/mnist/">http://yann.lecun.com/exdb/mnist/</a> and <a href="https://keras.io/api/datasets/mnist/">https://keras.io/api/datasets/mnist/</a></p> <p><u>How it will be accessible</u>: Already shared through a third-party repository link</p>

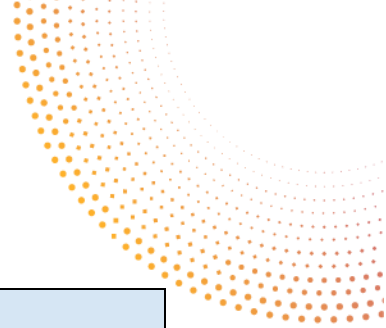


	<p><u>Methods/software tools to access data</u>: python scripts to download data are provided</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data simply consists of images so it is interoperable with any program that can read images</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u> N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly available under an attribution licence <a href="#">Creative Commons Attribution-Share Alike 3.0 licence</a>.</p> <p><u>Availability for re-use</u>: The loading and pre-processing mechanism developed for using the dataset is available to the public through the tensorflow library.</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: N/A (The dataset is already publicly available)</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 5.2.4 Interestingness10k image +video dataset

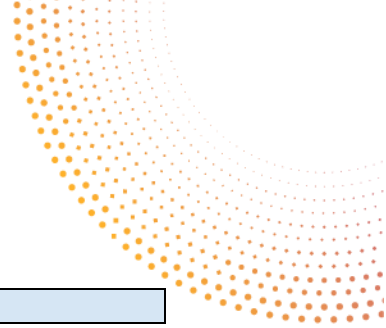
DMP component	AI4Media_Data_91_WP4_IMAGE_VIDEO_Interestingness10k Partner: UPB
Data Summary	<p><u>Purpose</u>: Interestingness10k is the most comprehensive collection of image and video information annotated for training and evaluating algorithms for visual interestingness prediction. It is used for T4.6 'Benchmarking of AI Systems' and for T6.6 'Measuring and Predicting User Perception of Social Media' and for T6.3 'Hybrid, privacy-enhanced recommendation'. Results involving this data are to be reported to the task corresponding deliverables.</p> <p><u>Type/format</u>: 9,831 images, 4 hours of video, interestingness scores determined based on more than 1M pair-wise annotations of 800 trusted annotators, some pre-computed multi-modal descriptors, and 192 system output results as baselines.</p> <p><u>Re-use of existing data</u>: No, it was newly generated. Part of the data usage and algorithms analysis was carried out within the project. The data was generated outside the project.</p> <p><u>Data origin</u>: IDF dataset: <a href="https://www.interdigital.com/data_sets/interestingness-dataset">https://www.interdigital.com/data_sets/interestingness-dataset</a></p>





	<p><u>Expected size</u>: 20 GB</p> <p><u>Data utility</u>: It is useful to WP4 and WP6 partners.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted on the Interdigital website.</p> <p><u>Search keywords</u>: Interestingness10k, Predicting Media Interestingness.</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: The data comes with metadata.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="https://www.interdigital.com/data_sets/interestingness-dataset">https://www.interdigital.com/data_sets/interestingness-dataset</a>. We will not reshare the data.</p> <p><u>How it will be accessible</u>: Already shared through a third-party repository link.</p> <p><u>Methods/software tools to access data</u>: Confirmation of data usage agreement via email is required.</p> <p><u>Repository</u>: Interdigital.</p> <p><u>Restrictions on access</u>: Access is made via email request.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is interoperable.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under a Creative Commons license that allows redistribution.</p> <p><u>Availability for re-use</u>: All the details of the data are available on the website. A detailed article describing the data, usage and many baseline systems is also available <i>M.G. Constantin, L.-D. Ștefan, B. Ionescu, Q.-K.-N. Duong, C.-H. Demarty, M. Sjöberg "Visual Interestingness Prediction: A Benchmark Framework and Literature Review", International Journal of Computer Vision</i></p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on UPB servers. UPB fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: We didn't identify any as materials are already publicly available Creative Commons content. The annotations are not recording any personal information of the user.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A.</p>





Other Issues	N/A
--------------	-----

### 5.2.5 MediaEval Memorability 2020 dataset

DMP component	AI4Media_Data_92_WP4_VIDEO_Memorability2020 Partner: UPB
Data Summary	<p><b>Purpose:</b> The Predicting Media Memorability 2020 dataset is a collection of videos annotated for their long- and short-term memorability impact on users. All materials are under Creative Commons licenses that allow redistribution. It is used for T4.6 'Benchmarking of AI Systems' and for T6.6 'Measuring and Predicting User Perception of Social Media' and for T6.3 'Hybrid, privacy-enhanced recommendation'. Results involving this data are to be reported to the task corresponding deliverables.</p> <p><b>Type/format:</b> 1,500 short videos, long- (72 hours) and short-(24 hours) term memorability scores, some pre-computed multi-modal descriptors.</p> <p><b>Re-use of existing data:</b> The video data is recovered from the TRECVID 2019 Video-to-Text dataset. The annotations were newly created. Part of the data usage and algorithms analysis was carried out within the project. The data was generated outside the project.</p> <p><b>Data origin:</b> Videos recovered from the TRECVID 2019 Video-to-Text dataset.</p> <p><b>Expected size:</b> 1.8 GB</p> <p><b>Data utility:</b> It is useful to WP4 and WP6 partners.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Data is discoverable.</p> <p><b>Search keywords:</b> Predicting Media Memorability.</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> The data comes with metadata.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The data is open for distribution. It has not yet been published publicly.</p> <p><b>How it will be accessible:</b> The data can be obtained via request to owners.</p> <p><b>Methods/software tools to access data:</b> N/A.</p> <p><b>Repository:</b> N/A</p> <p><b>Restrictions on access:</b> The data is provided by the authors.</p>
Making data interoperable	<p><b>Interoperability:</b> The data is interoperable.</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> The dataset is publicly available under a Creative Commons license that allows redistribution.</p> <p><b>Availability for re-use:</b> All the details of the data are available on the <a href="https://multimediaeval.github.io/editions/2020/tasks/memorability/">https://multimediaeval.github.io/editions/2020/tasks/memorability/</a> website. A detailed article describing the data, usage and baseline systems is also available <i>Alba García Seco De Herrera, Rukiye Savran Kiziltepe, Jon Chamberlain, Mihai Gabriel</i></p>



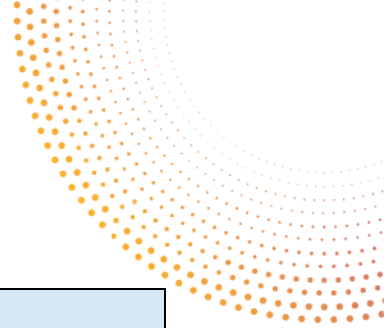


	<p><i>Constantin, Claire-Hélène Demarty, Faiyaz Doctor, Bogdan Ionescu and Alan F. Smeaton, Overview of MediaEval 2020 Predicting Media Memorability task: What Makes a Video Memorable? MediaEval Workshop 2021.</i></p> <p><u>Usable by third parties after end of project:</u> This is an open dataset.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on UPB servers. UPB fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> We didn't identify any as materials are already publicly available Creative Commons content. The annotations are not recording any user information.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A.</p>
Other Issues	N/A

### 5.2.6 ImageCLEF DrawnUI 2021 dataset

DMP component	AI4Media_Data_93_WP4_IMAGE_drawnUI2021 Partner: UPB
Data Summary	<p><u>Purpose:</u> The ImageCLEF drawnUI 2021 dataset contains hand drawn images of website user interface units and real website screenshots which are annotated for their user interface components. It serves for training systems capable of automatically identifying a website template from a drawing or a image of it. It is used for T4.6 'Benchmarking of AI Systems'. Results involving this data are to be reported to the task corresponding deliverables.</p> <p><u>Type/format:</u> 4,291 hand drawn images and 9,630 screenshot images, manual labelling of the positions of UI bounding boxes.</p> <p><u>Re-use of existing data:</u> The data and annotations were newly created. Part of the data usage and algorithms analysis was carried out within the project. The data was generated outside the project.</p> <p><u>Data origin:</u> Images from ImageCLEF drawnUI 2021 dataset.</p> <p><u>Expected size:</u> 9 GB</p> <p><u>Data utility:</u> It is useful to WP4 partners.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable.</p> <p><u>Search keywords:</u> ImageCLEFdrawnUI.</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> The data comes with metadata.</p>
Making data	<p><u>Data openly accessible:</u> The data is not open yet. The data set is so far owned by</p>



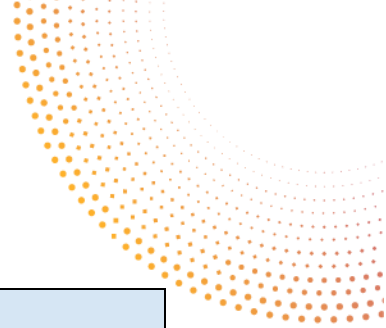


openly accessible	<p>teleportHQ. We plan to release it publicly.</p> <p><u>How it will be accessible</u>: The data can be obtained via request to authors.</p> <p><u>Methods/software tools to access data</u>: N/A.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: The data is provided by the authors.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is interoperable.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is owned by company teleportHQ, Romania.</p> <p><u>Availability for re-use</u>: All the details of the data are available on the <a href="https://www.imageclef.org/2021/drawnui">https://www.imageclef.org/2021/drawnui</a> website.</p> <p><u>Usable by third parties after end of project</u>: The data set is so far owned by teleportHQ. We plan to release it publicly.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on UPB servers. UPB fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: We didn't identify any as materials are either generated by the authors or recovered from public sources. The annotations are not recording any user information.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A.</p>
Other Issues	N/A

### 5.2.7 FCVID event recognition dataset

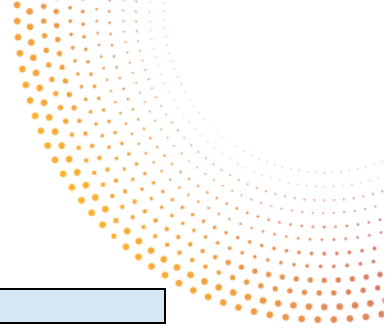
<b>DMP component</b>	<b>AI4Media_Data_94_WP4_Video_FCVID_v1</b> <b>Partner: CERTH</b>
Data Summary	<p><u>Purpose</u>: Fudan-Columbia Video Dataset (FCVID) is a public video event recognition dataset, containing 91,223 Web videos annotated manually according to 239 categories. It is used by CERTH to evaluate the proposed video event recognition models.</p> <p><u>Type/format</u>: .AVI/.FLV</p> <p><u>Re-use of existing data</u>: Yes.</p> <p><u>Data origin</u>: The dataset can be obtained by sending a request email to the authors at</p>





	<p><a href="https://fvl.fudan.edu.cn/dataset/fcvid/">https://fvl.fudan.edu.cn/dataset/fcvid/</a>.</p> <p><u>Expected size</u>: 2TB</p> <p><u>Data utility</u>: It is useful to WP4 partners to evaluate and benchmark video event recognition models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable and publicly available after sending a request email to the authors.</p> <p><u>Search keywords</u>: fcvid dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible upon email request at <a href="https://fvl.fudan.edu.cn/dataset/fcvid/">https://fvl.fudan.edu.cn/dataset/fcvid/</a>.</p> <p><u>How it will be accessible</u>: Shared through a third-party repository link after the email request has been approved by the authors.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: <a href="https://fvl.fudan.edu.cn/">https://fvl.fudan.edu.cn/</a></p> <p><u>Restrictions on access</u>: The dataset can be obtained after sending an email request to the authors.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: proprietary, for research purposes only</p> <p><u>Availability for re-use</u>: The dataset is already publicly available after sending a request email to the authors.</p> <p><u>Usable by third parties after end of project</u>: Yes.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: This dataset is downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Data should not be shared without the approval of the authors.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>





Other Issues	N/A
--------------	-----

### 5.2.8 YLI-MED event recognition dataset

DMP component	AI4Media_Data_95_WP4_Video_YLI-MED_v1 Partner: CERTH
Data Summary	<p><u>Purpose:</u> The YLI-MED corpus is an index of videos from the YFCC100M specialized for Multimedia Event Detection (MED) research. The videos are categorized as depicting one of 10 target events, or no target event. It is used by CERTH to evaluate the proposed event recognition models.</p> <p><u>Type/format:</u> mp4.</p> <p><u>Re-use of existing data:</u> Yes.</p> <p><u>Data origin:</u> The dataset contains more than 50K videos available at <a href="https://multimediacommons.wordpress.com/yli-multimedia-event-detection-subcorpus/">https://multimediacommons.wordpress.com/yli-multimedia-event-detection-subcorpus/</a>.</p> <p><u>Expected size:</u> 90 GB</p> <p><u>Data utility:</u> It is useful to WP4 partners to evaluate and benchmark event recognition models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable and publicly available.</p> <p><u>Search keywords:</u> YLI MED dataset</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is already openly accessible at <a href="https://multimediacommons.wordpress.com/yli-multimedia-event-detection-subcorpus/">https://multimediacommons.wordpress.com/yli-multimedia-event-detection-subcorpus/</a>.</p> <p><u>How it will be accessible:</u> Shared through a Amazon S3 data bucket</p> <p><u>Methods/software tools to access data:</u> Web-browser to download the data manually or to access the AWS CLI.</p> <p><u>Repository:</u> <a href="https://multimediacommons.wordpress.com/">https://multimediacommons.wordpress.com/</a>.</p> <p><u>Restrictions on access:</u> None.</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The dataset is publicly available under Creative Common License.</p> <p><u>Availability for re-use:</u> The dataset is already publicly available.</p> <p><u>Usable by third parties after end of project:</u> Yes.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>



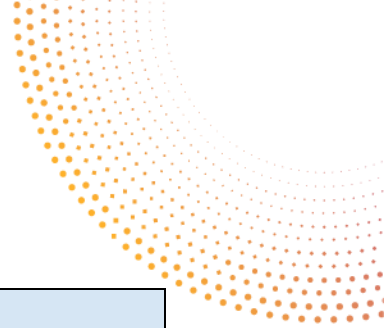
Allocation of resources	<u>Costs for making data FAIR:</u> N/A <u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> This dataset is downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

## 5.3 Datasets used in the context of WP5

### 5.3.1 SumMe video summarization dataset

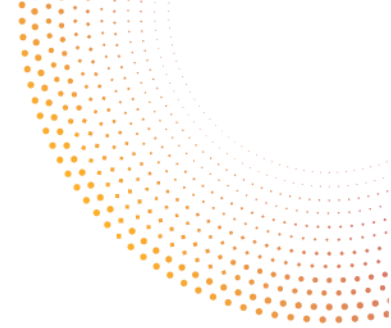
DMP component	AI4Media_Data_96_WP5_VIDEO_SumMeGycli14_v1 Partner: CERTH
Data Summary	<p><u>Purpose:</u> This dataset is composed of a set of videos from various genres (e.g. holidays, sports) and the associated ground-truth data that indicate the preferences of multiple human annotators with respect to the optimal visual summary for each video. It will be used in T5.1 for training and evaluating purposes, assisting the development of deep-learning-based architectures for video summarization.</p> <p><u>Type/format:</u> Video files in MP4 and WEBM format; MAT files with the ground-truth annotations; TXT and XLS files with information about the statistics of each video category; Matlab and Python scripts for evaluating the performance of a video summarization algorithm.</p> <p><u>Re-use of existing data:</u> Yes.</p> <p><u>Data origin:</u> This dataset was introduced in an ECCV 2014 paper titled “Creating Summaries from User Videos”, and was made publicly available through the following link: <a href="https://gyglim.github.io/me/vsum/index.html#benchmark">https://gyglim.github.io/me/vsum/index.html#benchmark</a></p> <p><u>Expected size:</u> ~2.5GB</p> <p><u>Data utility:</u> It will be useful to WP5 partners working on the development of video summarization methods, for training and evaluation purposes.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data are already publicly available at: <a href="https://gyglim.github.io/me/vsum/index.html#benchmark">https://gyglim.github.io/me/vsum/index.html#benchmark</a></p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> Any metadata generated to assist training and evaluation of video summarization methods (e.g. deep feature vectors representing the visual content of video frames, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be also created to facilitate their re-use.</p>





Making data openly accessible	<p><u>Data openly accessible</u>: The data are already openly accessible at: <a href="https://gyglim.github.io/me/vsum/index.html#benchmark">https://gyglim.github.io/me/vsum/index.html#benchmark</a></p> <p><u>How it will be accessible</u>: The original data are available through a third-party repository link. Any created metadata and their associated documentation will be made publicly available through a GitHub repository.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: <a href="https://gyglim.github.io/me/vsum/index.html#benchmark">https://gyglim.github.io/me/vsum/index.html#benchmark</a> Any generated metadata will be deposited at a publicly accessible GitHub repository.</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: Any metadata generated to assist training and evaluation of video summarization methods (e.g. deep feature vectors representing the visual content of video frames, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be also created to facilitate their re-use.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under an attribution non-commercial licence (<a href="https://gyglim.github.io/me/vsum/index.html#benchmark">https://gyglim.github.io/me/vsum/index.html#benchmark</a>)</p> <p><u>Availability for re-use</u>: The data are already publicly-available for re-use. Any generated metadata will be made permanently publicly-available for re-use as soon as they are complete and appropriately documented.</p> <p><u>Usable by third parties after end of project</u>: The dataset is already available for use by third parties.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the original dataset and the associated ECCV2014 paper should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A





### 5.3.2 TVSum video summarization dataset

DMP component	AI4Media_Data_97_WP5_VIDEO_TVSumSong15_v1 Partner: CERTH
Data Summary	<p><b>Purpose:</b> This dataset is composed of a set of videos from 10 categories of the TRECVID MED dataset - including news, how-to's, user-generated-content and documentaries - and the associated ground-truth data that indicate the opinion of multiple human annotators with respect to the importance of video frames and fragments. It will be used in T5.1 for training and evaluating purposes, assisting the development of deep-learning-based architectures for video summarization.</p> <p><b>Type/format:</b> Video files in MP4 format; video thumbnails in JPG format; MAT and TSV files with the ground-truth annotations and video-level metadata; Matlab scripts for evaluating the performance of a video summarization algorithm.</p> <p><b>Re-use of existing data:</b> Yes</p> <p><b>Data origin:</b> This dataset was introduced in a CVPR 2015 paper titled "TVSum: Summarizing web videos using titles", and was made publicly available through the following links: <a href="https://github.com/yalesong/tvsum">https://github.com/yalesong/tvsum</a> and <a href="http://people.csail.mit.edu/yalesong/tvsum/">http://people.csail.mit.edu/yalesong/tvsum/</a></p> <p><b>Expected size:</b> ~650MB</p> <p><b>Data utility:</b> It will be useful to WP5 partners working on the development of video summarization methods, for training and evaluation purposes.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Data are already publicly-available at: <a href="http://people.csail.mit.edu/yalesong/tvsum/">http://people.csail.mit.edu/yalesong/tvsum/</a></p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> Any metadata generated to assist training and evaluation of video summarization methods (e.g. deep feature vectors representing the visual content of video frames, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be also created to facilitate their re-use.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The data are already openly accessible at: <a href="http://people.csail.mit.edu/yalesong/tvsum/">http://people.csail.mit.edu/yalesong/tvsum/</a></p> <p><b>How it will be accessible:</b> The original data are shared through a third-party repository link. Any created metadata and their associated documentation will be made publicly-available through a GitHub repo.</p> <p><b>Methods/software tools to access data:</b> Web-browser to download the data as tgz file.</p> <p><b>Repository:</b> Original data are already deposited at: <a href="http://people.csail.mit.edu/yalesong/tvsum/">http://people.csail.mit.edu/yalesong/tvsum/</a>. Any generated metadata will be deposited at a publicly-accessible GitHub repo.</p> <p><b>Restrictions on access:</b> None</p>
Making data interoperable	<p><b>Interoperability:</b> N/A</p> <p><b>Data and metadata vocabularies:</b> Any metadata generated to assist training and evaluation of video summarization methods (e.g. deep feature vectors representing</p>



	<p>the visual content of video frames, or data about the shot-level structure of the videos) will be stored in HDF5 files; a documentation of these metadata will be also created to facilitate their re-use.</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The videos of this dataset, collected from YouTube, come with a Creative Commons CC-BY (v3.0) license. (<a href="https://github.com/yalesong/tvsum#overview">https://github.com/yalesong/tvsum#overview</a>)</p> <p><u>Availability for re-use:</u> The data are already publicly-available for re-use. Any generated metadata will be made permanently publicly-available for re-use as soon as they are complete and appropriately documented.</p> <p><u>Usable by third parties after end of project:</u> The dataset is already available for use by third parties.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Links to the original dataset and the associated CVPR2015 paper should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.3.3 RAI Monuments of Italy dataset

DMP component	AI4Media_Data_98_WP5_VIDEO_MonumentsOfItaly_v1 Partner: RAI
Data Summary	<p><u>Purpose:</u> A collection of videos depicting various Italian monuments. This dataset has been useful as a reference for developments in WP5 about landmark recognition. The processing included indexing of images and storage of resulting features in a database for search/match.</p> <p><u>Type/format:</u> Videos in MP4 or WMV format</p> <p><u>Re-use of existing data:</u> All videos are coming from RAI Archives.</p> <p><u>Data origin:</u> RAI Archives, mainly news production.</p> <p><u>Expected size:</u> 4 GB</p> <p><u>Data utility:</u> The dataset can be used to test landmark recognition algorithms.</p>





Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: Videos are organised in folders with self-descriptive names.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, the dataset will not be made openly accessible. It is stored in a private repository and can only be accessed by partners based on bilateral agreements with RAI.</p> <p><u>How it will be accessible</u>: Bilateral agreement with RAI.</p> <p><u>Methods/software tools to access data</u>: File transfer.</p> <p><u>Repository</u>: Private RAI repository.</p> <p><u>Restrictions on access</u>: Direct access is not allowed.</p>
Making data interoperable	<p><u>Interoperability</u>: Video encoding schemes used in the datasets are widely accepted by all common software tools.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: RAI specific license.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: To be defined.</p> <p><u>Re-use timeframe</u>: To be defined.</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Data is stored on local repositories inside RAI's premises, subject to the same security measures already used for IT infrastructure in RAI. These include network isolation from external internet accesses, firewalling, account-based access control management to the storage where the data copies are located. RAI fully complies with the applicable national, European data security frameworks, and the GDPR.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Possible rights issues related to material included in news items for which license of usage has expired.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	No

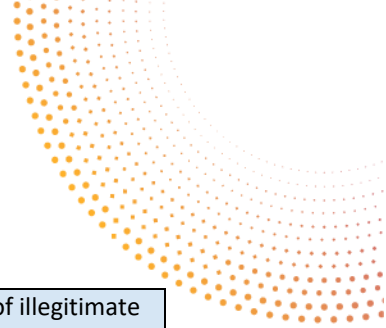
### 5.3.4 LVIS image dataset

<b>DMP component</b>	<b>AI4Media_Data_99_WP5_IMAGE_LVIS_v1</b> <b>Partner: JRC</b>
Data Summary	<u>Purpose</u> : LVIS is a dataset for large vocabulary instance segmentation, with > 1,200 categories, a large number of rare categories (long-tail), and high-quality instance



	<p>segmentation mask. The dataset will be used by JR within Task 5.3 (Learning with Scarce Data), specifically for training the few-shot object detection algorithms, which will be researched and developed within this task.</p> <p><u>Type/format</u>: Images, segmentation mask + json annotations</p> <p><u>Re-use of existing data</u>: Yes, we are reusing an existing dataset</p> <p><u>Data origin</u>: <a href="https://www.lvisdataset.org/dataset">https://www.lvisdataset.org/dataset</a></p> <p><u>Expected size</u>: 30 GB</p> <p><u>Date utility</u>: It is useful for various tasks, but especially for T5.3 (learning from scarce data) for the few-shot instance segmentation.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted in the Network Repository <a href="https://www.lvisdataset.org/dataset">https://www.lvisdataset.org/dataset</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="https://www.lvisdataset.org/dataset">https://www.lvisdataset.org/dataset</a> by its owners thus, we will not re-share it.</p> <p><u>How it will be accessible</u>: Shared through a third-party repository link</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file</p> <p><u>Repository</u>: Network Repository (<a href="https://www.lvisdataset.org/dataset">https://www.lvisdataset.org/dataset</a>)</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: LVIS has annotations for instance segmentations in a commonly used format similar to MS COCO. The annotations are stored using JSON. An API is provided to access and manipulate annotations.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The LVIS annotations along with this website are licensed under a Creative Commons Attribution 4.0 License. All LVIS dataset images come from the COCO dataset; please see <a href="https://cocodataset.org/#termsofuse">https://cocodataset.org/#termsofuse</a> for their terms of use.</p> <p><u>Availability for re-use</u>: Yes.</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on JR's servers. JR fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company</p>





	policies (firewalls, right-based file system, etc.) mitigate most of the risk of illegitimate access.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.5 CIFAR10/100 image dataset

DMP component	AI4Media_Data_100_WP5_IMAGE_CIFAR_v1 Partner: QMUL
Data Summary	<p><u>Purpose</u>: The CIFAR-10 and CIFAR-100 are labeled subsets of the 80 million tiny images dataset. CIFAR-10 consists of 60,000 labeled images, split between 10 including objects and animals. CIFAR-100 consists of 60,000 labeled images split between 100 “fine” classes and 20 “coarse” superclasses, including classes related to people. The CIFAR-10 and CIFAR-100 datasets are among the most popular image datasets for classification tasks. The datasets will be used to by QMUL to evaluate novel representation learning techniques developed for T5.3.</p> <p><u>Type/format</u>: Numpy array</p> <p><u>Re-use of existing data</u>: Yes</p> <p><u>Data origin</u>: <a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a></p> <p><u>Expected size</u>: ~160MB for each of the two datasets</p> <p><u>Data utility</u>: The datasets will be used for WP5 T5.3 to evaluate the developed methods and contrast their performance with existing works.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The dataset is discoverable from its website: <a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly available at <a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a>. We will not re-share it.</p> <p><u>How it will be accessible</u>: From original source.</p> <p><u>Methods/software tools to access data</u>: No specialized software is required to access the data. They can be downloaded directly from the source</p> <p><u>Repository</u>: <a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a></p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data are interoperable.</p> <p><u>Data and metadata vocabularies</u>: The data are provided in numpy format. The labels are also provided as numpy files, where each sample’s label is represented by an integer (0-9 for CIFAR-10, 0-99 for CIFAR-100). The dataset also includes a list that maps semantic labels to their numerical representation (e.g. label 0 in CIFAR-10</p>

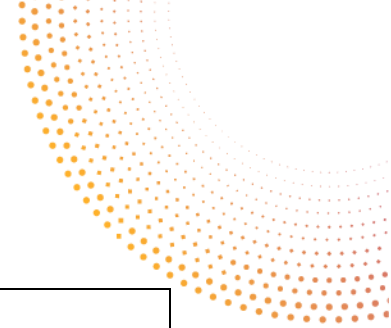


	<p>corresponds to the class “airplane”.</p> <p><u>Use of standard vocabularies:</u> The classes of the dataset relate to real world objects and entities and use common words as labels to identify them.</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> Already publicly shared.</p> <p><u>Availability for re-use:</u> The data are publicly available with no restrictions as to who can acquire it.</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> After downloading, the data are stored in the servers of Queen Mary University of London. Access requires username/password authentication. The security measures taken prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> N/A</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.3.6 STL-10 image dataset

<b>DMP component</b>	<b>AI4Media_Data_101_WP5_IMAGE_STL10_v1</b> <b>Partner: QMUL</b>
Data Summary	<p><u>Purpose:</u> STL-10 is a CIFAR-10 inspired dataset whose samples are drawn from the ImageNet dataset. It consists of 13,000 labeled samples that are equally split between 10 labels, the same as those of ImageNet. 5,000 of those samples belong to the training set and 8,000 to the test set. Furthermore, the dataset includes 100,000 unlabeled samples. Notably, the datasets will be used by QMUL to evaluate novel representation learning techniques developed for T5.3.</p> <p><u>Type/format:</u> Binary files</p> <p><u>Re-use of existing data:</u> Yes</p> <p><u>Data origin:</u> <a href="https://cs.stanford.edu/~acoates/stl10/">https://cs.stanford.edu/~acoates/stl10/</a></p> <p><u>Expected size:</u> 2.5 GB</p> <p><u>Data utility:</u> The datasets will be used for WP5 T5.3 to evaluate the developed methods and contrast their performance with existing works.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The dataset is discoverable from its website: <a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a>. No registration is required and no restrictions are present as to who can access the data</p> <p><u>Search keywords:</u> N/A</p>



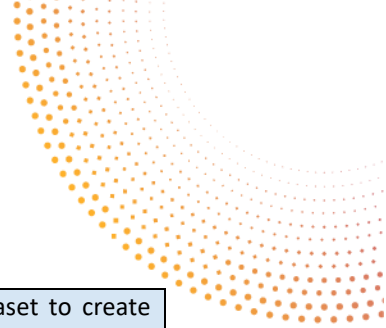


	<p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly available at <a href="https://cs.stanford.edu/~acoates/stl10/">https://cs.stanford.edu/~acoates/stl10/</a>. We will not re-share it.</p> <p><u>How it will be accessible</u>: From original source.</p> <p><u>Methods/software tools to access data</u>: No specialized software is required to access the data. They can be downloaded directly from the source. They can be read via python code that the source provides.</p> <p><u>Repository</u>: <a href="https://cs.stanford.edu/~acoates/stl10/">https://cs.stanford.edu/~acoates/stl10/</a></p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data are interoperable.</p> <p><u>Data and metadata vocabularies</u>: The data are provided in binary format. The labels are also provided as numpy files, where each sample's label is represented by an integer (range 1-10). The dataset also includes a txt file that maps semantic labels to their numerical representation. Finally, a separate txt file proposes specific data splits to be used for multi-fold validation purposes</p> <p><u>Use of standard vocabularies</u>: The classes of the dataset relate to real world objects and entities and use common words as labels to identify them.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Already publicly shared.</p> <p><u>Availability for re-use</u>: The data are publicly available with no restrictions as to who can acquire it.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the data are stored in the servers of Queen Mary University of London. Access requires username/password authentication. The security measures taken prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.7 CCNet text dataset for multilingual representation learning

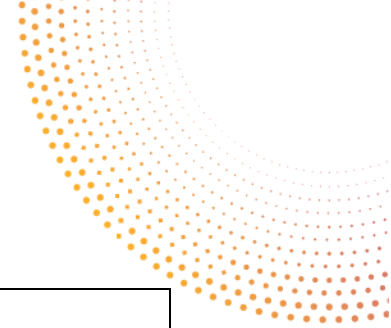
<b>DMP component</b>	AI4Media_Data_102_WP5_TEXT_CCNET_v1 Partner: CEA
Data	<u>Purpose</u> : CCNet is a large text dataset composed of high-quality monolingual datasets





Summary	<p>from Common Crawl for a variety of languages. CEA will use this dataset to create languages models in different languages. These models will be used to train information extraction models (e.g. opinion mining or named entity recognition). Results involving this dataset will be reported in D5.4.</p> <p><u>Type/format</u>: compressed JSON files (one per line)</p> <p><u>Re-use of existing data</u>: Yes, we are reusing an existing dataset</p> <p><u>Data origin</u>: <a href="https://github.com/facebookresearch/cc_net">https://github.com/facebookresearch/cc_net</a></p> <p><u>Expected size</u>: Unknown</p> <p><u>Data utility</u>: Useful for training language models. Will be used in WP5 for developing information extraction models in different languages.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable on the original source website.</p> <p><u>Search keywords</u>: CCNET text dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible at <a href="https://github.com/facebookresearch/cc_net">https://github.com/facebookresearch/cc_net</a>. The data will not be re-shared.</p> <p><u>How it will be accessible</u>: Shared through a third-party repository link</p> <p><u>Methods/software tools to access data</u>: python scripts to download data are provided.</p> <p><u>Repository</u>: Network Repository (<a href="http://networkrepository.com">http://networkrepository.com</a>)</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data is shared as JSON files.</p> <p><u>Data and metadata vocabularies</u>: The JSON files follow the following format (see <a href="https://github.com/facebookresearch/cc_net">https://github.com/facebookresearch/cc_net</a>):</p> <ul style="list-style-type: none"> <li>• url: webpage URL (part of CC)</li> <li>• date_download: date of download (part of CC)</li> <li>• digest: sha1 digest of the webpage (part of CC)</li> <li>• length: number of chars</li> <li>• nlines: number of lines</li> <li>• source_domain: web domain of the webpage</li> <li>• title: page title (part of CC)</li> <li>• raw_content: webpage content after deduplication</li> <li>• original_nlines: number of lines before deduplication</li> <li>• original_length: number of chars before deduplication</li> <li>• language: language detected by FastText LID</li> <li>• language_score: language score</li> <li>• perplexity: perplexity of a LM trained on Wikipedia</li> </ul> <p><u>Use of standard vocabularies</u>: N/A</p>





	<u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<p><u>Licence</u>: The data is shared through an MIT license.</p> <p><u>Availability for re-use</u>: The loading and pre-processing mechanism developed for using the dataset in experiments will be made publicly available to ensure reproducibility of research.</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The dataset will be downloaded from the original source and will be stored on CEA's servers. CEA fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A The dataset is already open. We will not reshare it.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.8 360 Video Viewing Dataset in Head Mounted Virtual Reality

DMP component	AI4Media_Data_103_WP5_VIDEO_360VideoViewingHeadMountedVR_v1 Partner: UCA
Data Summary	<p><u>Purpose</u>: The 360° Video Viewing Dataset in Head-Mounted Virtual Reality consists of the content (10 equirectangular videos) and sensory data (50 subjects) of 360-degree videos to Head Mounted Display (HMD).</p> <p>The dataset contains both content data (such as image saliency maps and motion maps derived from 360° videos) and sensor data (such as viewer head positions and orientations derived from HMD sensors).</p> <p>The content and sensor data are aligned using timestamps in the log files. The dataset was used by their creators to optimize 360° video streaming applications by the prediction of fixation points in 360° Videos.</p> <p>The dataset will be used to develop and test new algorithms for the prediction of head motion in 360° videos in WP5.</p> <p><u>Type/format</u>: videos are compressed using H-264 in MP4 container, while sensor (text) data is stored as comma separated values files in ASCII</p> <p><u>Re-use of existing data</u>: Yes, we are re-using an existing dataset from Lo, W. C., Fan, C. L., Lee, J., Huang, C. Y., Chen, K. T., &amp; Hsu, C. H. (2017, June). 360 video viewing dataset in head-mounted virtual reality. In Proceedings of the 8th ACM on Multimedia Systems Conference (pp. 211-216).</p>



	<p><u>Data origin</u>: Lo, W. C., Fan, C. L., Lee, J., Huang, C. Y., Chen, K. T., &amp; Hsu, C. H. (2017, June). 360 video viewing dataset in head-mounted virtual reality. In Proceedings of the 8th ACM on Multimedia Systems Conference (pp. 211-216). <a href="https://nmsl.cs.nthu.edu.tw/360-video-project/">https://nmsl.cs.nthu.edu.tw/360-video-project/</a></p> <p><u>Expected size</u>: ~1GB</p> <p><u>Data utility</u>: This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most. This dataset however can be leveraged in various novel applications in a much broader scope. For example, it can be used by researchers, engineers, and hobbyists to either optimize existing 360° video streaming applications (rate-distortion optimization, saliency prediction) and novel applications (crowd-driven camera movements, head motion prediction).</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is under the Networking and Multimedia Systems Lab of the Department of Computer Science at National Tsing Hua University (NTHU).</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on the website of the National Tsing Hua University <a href="https://nmsl.cs.nthu.edu.tw/360-video-project/">https://nmsl.cs.nthu.edu.tw/360-video-project/</a>.</p> <p><u>How it will be accessible</u>: The data is accessible directly by following a Google Drive link: <a href="https://drive.google.com/file/d/1s_EEUjUTa_N5u94Nuwir9gl3pwDEY_7_/view">https://drive.google.com/file/d/1s_EEUjUTa_N5u94Nuwir9gl3pwDEY_7_/view</a></p> <p><u>Methods/software tools to access data</u>: The data can be directly downloaded.</p> <p><u>Repository</u>: Web server of the National Tsing Hua University.</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Videos for saliency maps and motion maps are in mp4 format and information about sensory data (head orientation or viewed tiles) are text files in csv format. There is also a readme file in txt format explaining the different folders in the dataset.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To use the dataset, you should cite the work of <i>Lo, W. C., Fan, C. L., Lee, J., Huang, C. Y., Chen, K. T., &amp; Hsu, C. H. (2017, June). 360 video viewing dataset in head-mounted virtual reality. In Proceedings of the 8th ACM on Multimedia Systems Conference (pp. 211-216).</i></p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p>





	<u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on UCA servers. UCA fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A. <u>Is informed consent for data sharing and long term preservation given:</u> N/A (Note that actors in the dataset provided their consent for the creation of these traces and the users cannot be identified, and therefore no personal data is used).
Other Issues	N/A

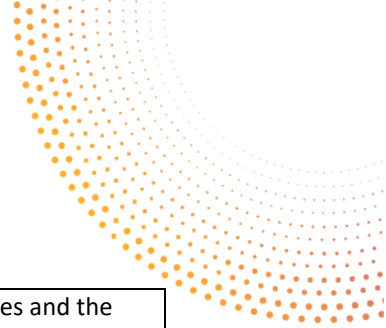
### 5.3.9 Predicting Head Movement in Panoramic Video Dataset

DMP component	AI4Media_Data_104_WP5_VIDEO_HeadMovementinPanoramicVideo_v1 Partner: UCA
Data Summary	<p><u>Purpose:</u> The dataset for the Head Movement Prediction in Panoramic Video contains both head movement and eye movement data of 56 subjects on 76 panoramic videos. Of 360-degree videos to Head Mounted Display (HMD).</p> <p>The dataset contains both content data (the 360° videos) and sensor data (the head positions of subjects exploring each video).</p> <p>The sensor data is stored in a Matlab file, this file includes 76 cells, corresponding to the head motion data of all 76 videos. Each cell records (longitude, latitude) of head motion pairs for 58 subjects in the columns of the Matlab file.</p> <p>The dataset will be used in WP5 to develop and test new algorithms for the prediction of head motion in 360° videos.</p> <p><u>Type/format:</u> videos are compressed using H-264 in MP4 container, while sensor data is stored as a MATLAB file.</p> <p><u>Re-use of existing data:</u> Yes, we are re-using an existing dataset from Xu, M., Song, Y., Wang, J., Qiao, M., Huo, L., &amp; Wang, Z. (2018). Predicting head movement in panoramic video: A deep reinforcement learning approach. IEEE transactions on pattern analysis and machine intelligence, 41(11), 2693-2708.</p> <p><u>Data origin:</u> M., Song, Y., Wang, J., Qiao, M., Huo, L., &amp; Wang, Z. (2018). Predicting head movement in panoramic video: A deep reinforcement learning approach. IEEE transactions on pattern analysis and machine intelligence, 41(11), 2693-2708.</p> <p><u>Expected size:</u> ~4GB</p> <p><u>Data utility:</u> This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most. This dataset however can be leveraged in various novel applications in a much broader scope. For example, it can be used by researchers, engineers, and hobbyists to either optimize existing 360° video streaming applications (rate-distortion optimization, saliency prediction) and novel applications (crowd-driven camera movements, head-eye relationship in video exploration).</p>
Making data findable, incl.	<u>Is data discoverable:</u> Yes, the data is hosted in the github repository of the authors



provisions for metadata	<p><a href="https://github.com/YuhangSong/DHP">https://github.com/YuhangSong/DHP</a>.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on the repository of the authors <a href="https://github.com/YuhangSong/DHP">https://github.com/YuhangSong/DHP</a>.</p> <p><u>How it will be accessible</u>: The data is accessible directly by following a Dropbox link: <a href="https://www.dropbox.com/s/swenk8b33vs6151/PVS-HM.tar.gz?dl=0">https://www.dropbox.com/s/swenk8b33vs6151/PVS-HM.tar.gz?dl=0</a></p> <p><u>Methods/software tools to access data</u>: The data can be directly downloaded.</p> <p><u>Repository</u>: Github</p> <p><u>Restrictions on access</u>: There is no restriction to access the dataset, it is hosted in the repository of the authors at <a href="https://github.com/YuhangSong/DHP">https://github.com/YuhangSong/DHP</a>, with a direct link to download in Dropbox <a href="https://www.dropbox.com/s/swenk8b33vs6151/PVS-HM.tar.gz?dl=0">https://www.dropbox.com/s/swenk8b33vs6151/PVS-HM.tar.gz?dl=0</a>. However, the authors recommend contacting them so that they can grant permission to the file.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Equirectangular videos are in mp4 format and information about sensory data (head orientation) is stored as a MATLAB file with an entry per video, each with a matrix with a column per user and alternating longitude and latitudes in the rows.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To use the dataset, you should cite the work of M., Song, Y., Wang, J., Qiao, M., Huo, L., &amp; Wang, Z. (2018). Predicting head movement in panoramic video: A deep reinforcement learning approach. IEEE transactions on pattern analysis and machine intelligence, 41(11), 2693-2708.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The authors recommend contacting them by mail to ask for permission to access the file, before a password was needed to download the dataset, but now the data is publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on UCA servers. UCA fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that</p>





	actors in the dataset provided their consent for the creation of these traces and the users cannot be identified, and therefore no personal data is used.)
Other Issues	N/A

### 5.3.10 Gaze prediction in Dynamic 360° Immersive Videos Dataset

DMP component	AI4Media_Data_105_WP5_VIDEO_GazePredictionDynamic360ImmersiveVideos_v1 Partner: UCA
Data Summary	<p><b>Purpose:</b> The dataset for Gaze Prediction in VR Videos contains 208 videos captured in dynamic scenes, and the traces of head and gaze position of 45 subjects, the traces of at least 31 subjects are recorded per video. The sensor data is stored in a csv files with a folder per user and a text file per trace.</p> <p>The dataset will be used in WP5 to develop and test new algorithms for the prediction of head motion in 360° videos.</p> <p><b>Type/format:</b> videos are compressed using H-264 in MP4 container, while sensor data is stored as text files.</p> <p><b>Re-use of existing data:</b> Yes, we are re-using an existing dataset from Xu, Y., Dong, Y., Wu, J., Sun, Z., Shi, Z., Yu, J., &amp; Gao, S. (2018). Gaze prediction in dynamic 360 immersive videos. In proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 5333-5342).</p> <p><b>Data origin:</b> Xu, Y., Dong, Y., Wu, J., Sun, Z., Shi, Z., Yu, J., &amp; Gao, S. (2018). Gaze prediction in dynamic 360 immersive videos. In proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 5333-5342). <a href="https://github.com/xuyanyu-shh/VR-EyeTracking">https://github.com/xuyanyu-shh/VR-EyeTracking</a>.</p> <p><b>Expected size:</b> ~4GB</p> <p><b>Data utility:</b> This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most. This dataset however can be leveraged in various novel applications in a much broader scope. For example, it can be used by researchers, engineers, and hobbyists to either optimize existing 360° video streaming applications (rate-distortion optimization, saliency prediction) and novel applications (crowd-driven camera movements, head-eye relationship in video exploration).</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes, the data is hosted in the github repository of the authors <a href="https://github.com/xuyanyu-shh/VR-EyeTracking">https://github.com/xuyanyu-shh/VR-EyeTracking</a>.</p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The data is already openly accessible on the repository of the authors <a href="https://github.com/xuyanyu-shh/VR-EyeTracking">https://github.com/xuyanyu-shh/VR-EyeTracking</a>.</p> <p><b>How it will be accessible:</b> The data is accessible by following a Baidu link: <a href="https://pan.baidu.com/s/1RKTZoeiLjKidW3S1E0I_yw">https://pan.baidu.com/s/1RKTZoeiLjKidW3S1E0I_yw</a> with the code "olxt" given in the github repository.</p> <p><b>Methods/software tools to access data:</b> The data can be downloaded from the baidu</p>

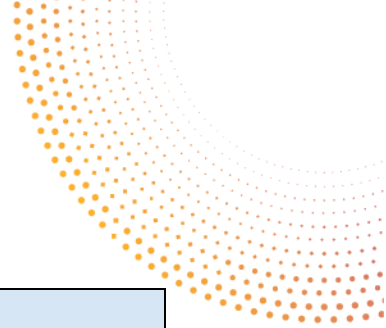


	<p>website using the baidu-pan-downloader <a href="https://github.com/dotennin/baidu-pan-downloader">https://github.com/dotennin/baidu-pan-downloader</a>.</p> <p><u>Repository</u>: Github</p> <p><u>Restrictions on access</u>: There is a password to access the Baidu link, the password "olxt" can be found in the github repository.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Equirectangular videos are in mp4 format and information about sensory data (head and gaze position) is stored as text files with a folder per user and a text file per video.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To use the dataset, you should cite the work of Xu, Y., Dong, Y., Wu, J., Sun, Z., Shi, Z., Yu, J., &amp; Gao, S. (2018). Gaze prediction in dynamic 360 immersive videos. In proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 5333-5342).</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data is already publicly available.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on UCA servers. UCA fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset provided their consent for the creation of these traces and the users cannot be identified, and therefore no personal data is used.)</p>
Other Issues	N/A

### 5.3.11 Your Attention is Unique Dataset

<b>DMP component</b>	<b>AI4Media_Data_106_WP5_VIDEO_YourAttentionIsUnique_v1</b> <b>Partner: UCA</b>
Data Summary	<p><u>Purpose</u>: The dataset for the PanoSalnet model consists on the post-processing of a publicly available dataset. The dataset includes 18 videos viewed by 48 users, from which 9 videos are selected. The dataset also includes the model weights of a Python Caffe implementation.</p> <p>The sensor data and the saliency data are stored in a Python dictionary with an entry per video, each with a list with an entry per user. This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to</p>





	<p>watch the most (WP5).</p> <p><u>Type/format</u>: All the data is contained within a python dictionary, numpy arrays are used to store the values of the saliency maps and the head motion traces.</p> <p><u>Re-use of existing data</u>: Yes, we are re-using an existing dataset from Nguyen, A., Yan, Z., &amp; Nahrstedt, K. (2018, October). Your attention is unique: Detecting 360-degree video saliency in head-mounted display for head movement prediction. In Proceedings of the 26th ACM international conference on Multimedia (pp. 1190-1198).</p> <p><u>Data origin</u>: Nguyen, A., Yan, Z., &amp; Nahrstedt, K. (2018, October). Your attention is unique: Detecting 360-degree video saliency in head-mounted display for head movement prediction. In Proceedings of the 26th ACM international conference on Multimedia (pp. 1190-1198). <a href="https://github.com/phananh1010/PanoSalNet">https://github.com/phananh1010/PanoSalNet</a>.</p> <p><u>Expected size</u>: ~1GB</p> <p><u>Data utility</u>: This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most. This dataset however can be leveraged in various novel applications in a much broader scope. For example, it can be used by researchers, engineers, and hobbyists to either optimize existing 360° video streaming applications (rate-distortion optimization, saliency prediction) and novel applications (crowd-driven camera movements).</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted in the github repository of the authors <a href="https://github.com/phananh1010/PanoSalNet">https://github.com/phananh1010/PanoSalNet</a>.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on the repository of the authors <a href="https://github.com/phananh1010/PanoSalNet">https://github.com/phananh1010/PanoSalNet</a>.</p> <p><u>How it will be accessible</u>: The data is accessible by following a Dropbox link: <a href="https://www.dropbox.com/s/smiplkqqlv0npsm/panosalnet_iter_800.caffemodel?dl=0">https://www.dropbox.com/s/smiplkqqlv0npsm/panosalnet_iter_800.caffemodel?dl=0</a> given in the github repository.</p> <p><u>Methods/software tools to access data</u>: N/A.</p> <p><u>Repository</u>: Github.</p> <p><u>Restrictions on access</u>: N/A.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Equirectangular saliency maps and sensory data (head position) are in numpy array format inside a Python dictionary the weights of the model are given in a caffemodel format.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To use the dataset, you should cite the work of Nguyen, A., Yan, Z., &amp; Nahrstedt, K. (2018, October). Your attention is unique: Detecting 360-degree video saliency in head-mounted display for head movement prediction. In Proceedings of the 26th ACM international conference on Multimedia (pp. 1190-1198).</p>

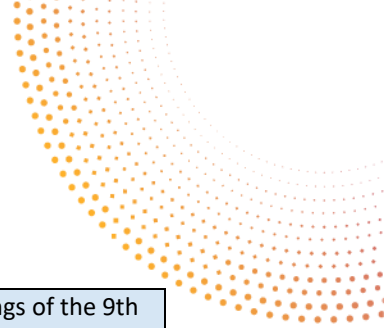


	<p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data is already publicly available.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on UCA servers. UCA fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset provided their consent for the creation of these traces and the users cannot be identified, and therefore no personal data is used.)</p>
Other Issues	N/A

### 5.3.12 Dataset of Head and Eye Movements for 360° Videos

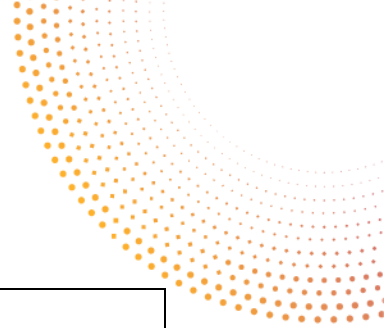
DMP component	AI4Media_Data_107_WP5_VIDEO_HeadEyeMovements360Videos_v1 Partner: UCA
Data Summary	<p><u>Purpose</u>: The dataset of Head and Eye Movements for 360° Videos is composed of 19 videos in 4K resolution in equirectangular format. It contains the exploration of 57 observers on all 19 videos for a duration of 20 seconds.</p> <p>Visual attention data is organized in folders according to their data type: whether if they come from Head-only movements or head and eye movements. For both cases, saliency maps are stored in a separate folder. It contains one saliency map per stimulus as a compressed binary file. The scanpaths are stored in a directory as CSV text files for each video. These CSV files contain all identified fixations for one video, ordered temporally for each observer one after the other in the file. The first data column reports fixation indexes for each participant, this value is incremented with each new fixation until reaching the end of an observer's trial, after which indexing starts over at 0 for the next observer. Next two columns are gaze positions in longitudes and latitudes normalized between 0 and 1 and then two columns contain the head positions in the same format.</p> <p>This dataset is useful in the context of WP5 to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most.</p> <p><u>Type/format</u>: videos in mp4, scanpaths in csv files and saliency maps in binary files.</p> <p><u>Re-use of existing data</u>: Yes, we are re-using an existing dataset from David, E. J., Gutiérrez, J., Coutrot, A., Da Silva, M. P., &amp; Callet, P. L. (2018, June). A dataset of head and eye movements for 360 videos. In Proceedings of the 9th ACM Multimedia Systems Conference (pp. 432-437) this dataset was used in the Salient360°! ICME'18 Grand Challenge.</p> <p><u>Data origin</u>: David, E. J., Gutiérrez, J., Coutrot, A., Da Silva, M. P., &amp; Callet, P. L. (2018,</p>





	<p>June). A dataset of head and eye movements for 360 videos. In Proceedings of the 9th ACM Multimedia Systems Conference (pp. 432-437). <a href="https://salient360.ls2n.fr/datasets/training-dataset/">https://salient360.ls2n.fr/datasets/training-dataset/</a>.</p> <p><u>Expected size</u>: ~1GB</p> <p><u>Data utility</u>: This dataset is useful in our context to train deep neural networks to predict which parts of 360° videos attract viewers to watch the most. This dataset however can be leveraged in various novel applications in a much broader scope. For example, it can be used by researchers, engineers, and hobbyists to either optimize existing 360° video streaming applications (rate-distortion optimization, saliency prediction) and novel applications for coding, transmitting, and rendering 360° content.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted in the website of the salient360! challenge <a href="https://salient360.ls2n.fr/datasets/training-dataset/">https://salient360.ls2n.fr/datasets/training-dataset/</a>.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on the repository of the Salient360! Challenge <a href="https://salient360.ls2n.fr/datasets/training-dataset/">https://salient360.ls2n.fr/datasets/training-dataset/</a>.</p> <p><u>How it will be accessible</u>: The data is accessible by a ftp link: <a href="ftp://gdchallenge18:1piN4nte5@ftp.ivc.polytech.univ-nantes.fr">ftp://gdchallenge18:1piN4nte5@ftp.ivc.polytech.univ-nantes.fr</a> given in the repository of the University of Nantes.</p> <p><u>Methods/software tools to access data</u>: The data could be downloaded using any software to download ftp files.</p> <p><u>Repository</u>: <a href="https://salient360.ls2n.fr/datasets/training-dataset/">https://salient360.ls2n.fr/datasets/training-dataset/</a>.</p> <p><u>Restrictions on access</u>: N/A.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Equirectangular videos are in mp4 format, saliency maps extracted from the videos are in binary format and sensory data (head orientation and eye motion) is stored in text files in csv format. There is also a readme file in txt format explaining the different folders in the dataset and its usage.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To use the dataset, you should cite the work of David, E. J., Gutiérrez, J., Coutrot, A., Da Silva, M. P., &amp; Callet, P. L. (2018, June). A dataset of head and eye movements for 360 videos. In Proceedings of the 9th ACM Multimedia Systems Conference (pp. 432-437).</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data is already publicly available.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of	<p><u>Costs for making data FAIR</u>: N/A</p>





resources	<u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on UCA servers. UCA fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A.  <u>Is informed consent for data sharing and long term preservation given:</u> N/A (Note that actors in the dataset provided their consent for the creation of these traces and they received monetary compensation for their time, the users cannot be identified in the dataset, and therefore no personal data is used).
Other Issues	N/A

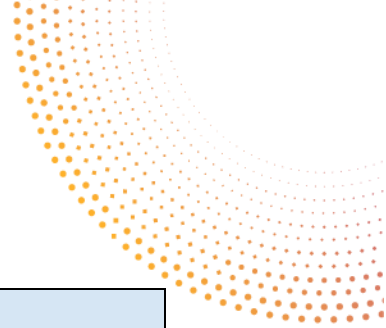
### 5.3.13 Dim-sim dataset for music similarity search

DMP component	AI4Media_Data_108_WP5_Audio_dim-sim_v1 Partner: FhG-IDMT
Data Summary	<p><u>Purpose:</u> The dim-sim dataset is a collection of user-annotated triplet ratings for music similarity. The dataset can be used to train and to evaluate algorithms for music similarity search and music recommendation. All annotations relate to audio files from the Million Song Dataset (MSD).</p> <p>4,000 3-second triplets were randomly sampled from the MSD and were each annotated by 5-12 annotators w.r.t. to which song being more similar to the anchor song. In total, the dataset includes 39,400 human annotations. Furthermore, a subset of cleaned annotations with higher agreement is additionally provided.</p> <p><u>Type/format:</u> Triplet annotations (csv / json) with audio file names (MSD) and similarity ratings<sup>11</sup></p> <p><u>Re-use of existing data:</u> Yes</p> <p><u>Data origin:</u> Million Song Dataset (MSD) <a href="http://millionsongdataset.com/">http://millionsongdataset.com/</a></p> <p><u>Expected size:</u> &lt;1MB</p> <p><u>Data utility:</u> The dataset is useful in the context of T5.6 for the evaluation of algorithms for disentangled music similarity search.</p>
	<p><u>Is data discoverable:</u> Data is discoverable. The dataset is hosted on Zenodo: <a href="https://zenodo.org/record/3889149#.XuovcxMzbyV">https://zenodo.org/record/3889149#.XuovcxMzbyV</a></p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The dataset is openly accessible.</p> <p><u>How it will be accessible:</u> Shared through a third-party repository link: <a href="https://zenodo.org/record/3889149#.XuovcxMzbyV">https://zenodo.org/record/3889149#.XuovcxMzbyV</a></p> <p><u>Methods/software tools to access data:</u> Web-browser to download the data as zip file.</p>

<sup>11</sup> As documented here: <https://jongpillee.github.io/multi-dim-music-sim/>







	<p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under a non-commercial attribution licence (Creative Commons Attribution Non Commercial 4.0 International)</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on FhG-IDMT's servers. FhG-IDMT fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

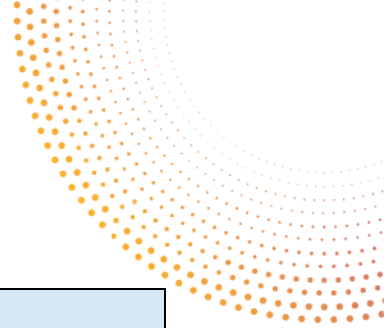
### 5.3.14 SPAM dataset for music segmentation

DMP component	AI4Media_Data_109_WP5_Audio_spam_music_v1 Partner: FhG-IDMT
Data Summary	<p><u>Purpose</u>: The SPAM dataset includes structural annotations for 50 tracks from 5 human annotators each<sup>12</sup>. These annotations include the segment boundary times as well as segment labels indicating similar segments such as chorus or vers. As the original audio files are copyright-protected, the dataset instead includes 5 different types of (pre-computed) audio features relating to timbre, tonality/harmony, and rhythm.</p> <p>Type/format: Annotations (json), Pre-computed audio features (hosted separately at CCRMA Stanford)<sup>13</sup>, audio file metadata (tsv), python scripts for dataset parsing</p> <p><u>Re-use of existing data</u>: Yes</p>

<sup>12</sup> <https://github.com/uriniето/msaf-data/tree/master/SPAM>

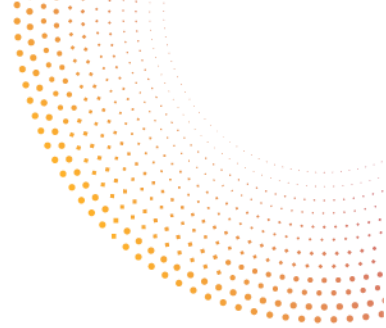
<sup>13</sup> <https://ccrma.stanford.edu/~uriniето/SPAM/SPAM-features.tgz>





	<p><u>Data origin</u>: 50 copyright-protected audio files of various origins</p> <p><u>Expected size</u>: ~1.2 GB</p> <p><u>Data utility</u>: The dataset is useful in the context of T5.6 for the evaluation of algorithms for music structure analysis.</p>
	<p><u>Is data discoverable</u>: Data is discoverable. The dataset is hosted on Github and (linked from there) on a website hosted by CCRMA: <a href="https://github.com/uriniето/msaf-data/tree/master/SPAM">https://github.com/uriniето/msaf-data/tree/master/SPAM</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The dataset is openly accessible.</p> <p><u>How it will be accessible</u>: Shared through a third-party repository link</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the dataset as zip file.</p> <p><u>Repository</u>: Github</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under a not-specified licence.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on FhG-IDMT's servers. FhG-IDMT fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A



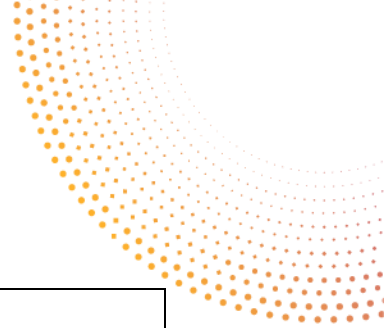


### 5.3.15 SALAMI dataset for music segmentation

DMP component	AI4Media_Data_110_WP5_Audio_salami_music_v1 Partner: FhG-IDMT
Data Summary	<p><b>Purpose:</b> The SALAMI (Structural Analysis of Large Amounts of Music Information) dataset includes structural annotations for around 1300 licence free music recordings from one or multiple annotators.</p> <p><b>Type/format:</b> Annotations (txt) + Metadata</p> <p><b>Re-use of existing data:</b> Yes</p> <p><b>Data origin:</b> 50 copyright-protected audio files of various origins: <a href="https://github.com/DDMAL/salami-data-public">https://github.com/DDMAL/salami-data-public</a></p> <p><b>Expected size:</b> ~1.2 GB</p> <p><b>Data utility:</b> The dataset is useful in the context of T5.6 for the evaluation of algorithms for music structure analysis.</p>
	<p><b>Is data discoverable:</b> Data is discoverable. The dataset is hosted on Github (<a href="https://github.com/DDMAL/salami-data-public">https://github.com/DDMAL/salami-data-public</a>) (a later version was hosted on the facility of McGill University: <a href="https://ddmal.music.mcgill.ca/research/SALAMI/annotation/">https://ddmal.music.mcgill.ca/research/SALAMI/annotation/</a> )</p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The dataset is openly accessible.</p> <p><b>How it will be accessible:</b> Shared through a third-party repository link: <a href="https://github.com/DDMAL/salami-data-public">https://github.com/DDMAL/salami-data-public</a></p> <p><b>Methods/software tools to access data:</b> Web-browser to download the dataset, audio files can be partly accessed via matching Youtube-links<sup>14</sup>.</p> <p><b>Repository:</b> Github</p> <p><b>Restrictions on access:</b> N/A</p>
Making data interoperable	<p><b>Interoperability:</b> N/A</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> The dataset is publicly available under a Creative Commons 0 licence.</p> <p><b>Availability for re-use:</b> Yes</p> <p><b>Usable by third parties after end of project:</b> This is an open dataset.</p> <p><b>Re-use timeframe:</b> N/A</p> <p><b>Data quality assurance process:</b> N/A</p>

<sup>14</sup> <https://github.com/jblsmith/matching-salami>





Allocation of resources	<u>Costs for making data FAIR:</u> N/A <u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on FhG-IDMT's servers. FhG-IDMT fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> Links to the dataset should be shared instead of raw data. <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 5.3.16 Harmonix dataset for music segmentation

DMP component	AI4Media_Data_111_WP5_Audio_harmonix_music_v1 Partner: FhG-IDMT
Data Summary	<u>Purpose:</u> The Harmonix dataset includes metrical (beats & downbeats) and structural annotations for 912 pop tracks.  Type/format: Annotations (json), Pre-computed audio features (mel-spectrograms) hosted on a dropbox account  <u>Re-use of existing data:</u> Yes  <u>Data origin:</u> 912 copyright-protected Western pop music recordings at <a href="https://github.com/uriniето/harmonixset">https://github.com/uriniето/harmonixset</a>  <u>Expected size:</u> ~1.3 GB  <u>Data utility:</u> The dataset is useful in the context of T5.6 for the evaluation of algorithms for music structure analysis.
	<u>Is data discoverable:</u> Data is discoverable. The dataset is hosted on Github ( <a href="https://github.com/uriniето/harmonixset">https://github.com/uriniето/harmonixset</a> ) and (linked from there) on an Dropbox account ( <a href="https://www.dropbox.com/s/zxnqlx0hxx0lsyc/Harmonix_melspecs.tgz?dl=0">https://www.dropbox.com/s/zxnqlx0hxx0lsyc/Harmonix_melspecs.tgz?dl=0</a> ). Audio files can be retrieved via matching Youtube-URLs.  <u>Search keywords:</u> N/A  <u>Versioning:</u> N/A  <u>Metadata creation:</u> N/A
Making data openly accessible	<u>Data openly accessible:</u> The dataset is openly accessible.  <u>How it will be accessible:</u> Shared through a third-party repository link <a href="https://github.com/uriniето/harmonixset">https://github.com/uriniето/harmonixset</a>  <u>Methods/software tools to access data:</u> Web-browser to download the dataset as zip file.  <u>Repository:</u> Github  <u>Restrictions on access:</u> N/A
Making data	<u>Interoperability:</u> N/A

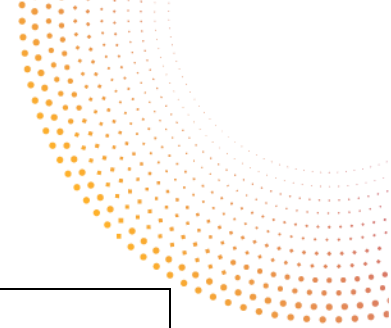


interoperable	<p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is publicly available under the MIT Licence.</p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: This is an open dataset.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on FhG-IDMT's servers. FhG-IDMT fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.17 Free Music Archive dataset

DMP component	AI4Media_Data_112_WP5_Audio_FMA_v1 Partner: IRCAM
Data Summary	<p><u>Purpose</u>: The Free Music Archive (FMA) song dataset is a collection of 106,574 musical recordings, from 16,341 artists. It comes with musical genre annotations, artist and album names, and other metadata. This dataset has mainly an interest for non-supervised learning, and training of GANs for the sound synthesis of musical mixes.</p> <p><u>Type/format</u>: MP3 files for the audio recordings, and CSV files for the metadata.</p> <p><u>Re-use of existing data</u>: Yes, we reuse a public dataset from EPFL.</p> <p><u>Data origin</u>: <a href="http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis">http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis</a></p> <p><u>Expected size</u>: 879 GB</p> <p><u>Data utility</u>: This dataset will be used by IRCAM for the sound synthesis and audio analysis (T5.2, T5.6, and Use Case 5).</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable in the original source: <a href="http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis">http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis</a></p> <p><u>Search keywords</u>: FMA music</p> <p><u>Versioning</u>: N/A</p>



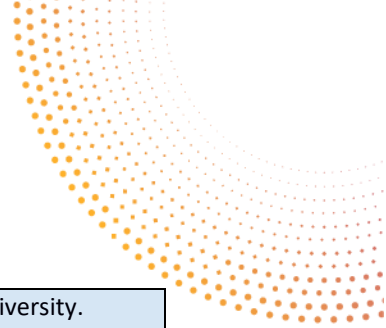


	<u>Metadata creation</u> : The metadata are written in the standard CSV files.
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already publicly shared at <a href="http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis">http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis</a>. We will not re-share the data.</p> <p><u>How it will be accessible</u>: Can be downloaded from <a href="http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis">http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis</a>.</p> <p><u>Methods/software tools to access data</u>: Web-browser</p> <p><u>Repository</u>: <a href="http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis">http://archive.ics.uci.edu/ml/datasets/FMA%3A+A+Dataset+For+Music+Analysis</a>.</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The data and metadata formats are standard, which makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies</u>: A documentation and a publication are accessible on the repository website for an explanation of the content.</p> <p><u>Use of standard vocabularies</u>: No</p> <p><u>Mappings to commonly used vocabularies</u>: No</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared under a Creative Commons license.</p> <p><u>Availability for re-use</u>: Already publicly shared</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : After downloading, the dataset is stored in IRCAM's servers, for which access requires username/password authentication. Security measures prevent illegitimate access (firewalls and rights-based-file system).
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.18 LAKH MIDI music dataset

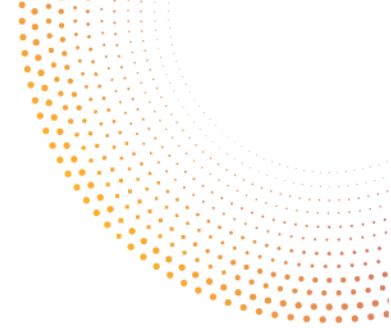
<b>DMP component</b>	<b>AI4Media_Data_113_WP5_Audio_LAKH-MIDI_v1</b> <b>Partner: IRCAM</b>
Data Summary	<p><u>Purpose</u>: The Lakh MIDI dataset is a collection of digital scores of 176,581 songs, with the MIDI format and other annotations. It will be used in WP5 for the sound synthesis of full musical mixes, and the learning of score augmentation.</p> <p><u>Type/format</u>: MIDI files.</p>





	<p><u>Re-use of existing data</u>: Yes, we re-use a public dataset from Columbia University.</p> <p><u>Data origin</u>: <a href="https://colinraffel.com/projects/lmd/">https://colinraffel.com/projects/lmd/</a></p> <p><u>Expected size</u>: 7.6 GB</p> <p><u>Data utility</u>: This dataset will be used by IRCAM for sound synthesis, MIDI score augmentation, and possibly for automatic transcription and time alignment (T5.2, T5.6, and Use Case 5).</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable in the original source: <a href="https://colinraffel.com/projects/lmd/">https://colinraffel.com/projects/lmd/</a></p> <p><u>Search keywords</u>: LAKH MIDI music</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: The Lakh dataset is made of MIDI files without external metadata.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already publicly shared at <a href="https://colinraffel.com/projects/lmd/">https://colinraffel.com/projects/lmd/</a>. We will not re-share the data.</p> <p><u>How it will be accessible</u>: Can be downloaded from <a href="https://colinraffel.com/projects/lmd/">https://colinraffel.com/projects/lmd/</a></p> <p><u>Methods/software tools to access data</u>: Web-browser</p> <p><u>Repository</u>: <a href="https://colinraffel.com/projects/lmd/">https://colinraffel.com/projects/lmd/</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: The MIDI format is standard in computer music. It is well documented and many libraries and software are available to read it.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is already shared under a CC-BY 4.0 license.</p> <p><u>Availability for re-use</u>: Already publicly shared</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the dataset is stored in IRCAM's servers, for which access requires username/password authentication. Security measures prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A





### 5.3.19 Piano Audio and MIDI music datasets

DMP component	AI4Media_Data_114_WP5_Audio_MIDI_Piano_v1 Partner: IRCAM
Data Summary	<p><b>Purpose:</b> The ENST-MAPS and the MAESTRO datasets are two collections of real piano recordings: (1) individual notes and (2) full musical piano pieces, coming with the corresponding digital scores. The recordings and the MIDI scores have been produced using a mechanic piano (YAMAHA disklavier). These data are useful for piano sound synthesis and analyses, such as automatic transcription.</p> <p><b>Type/format:</b> WAV files for the audio recordings, and MIDI files for the scores.</p> <p><b>Re-use of existing data:</b> Yes, we re-use public datasets from Telecom-Paris and Google AI.</p> <p><b>Data origin:</b> ENST-MAPS available on the Telecom-Paris website: <a href="http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/">http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/</a>, and MAESTRO available with the Magenta Framework of Google AI: <a href="https://magenta.tensorflow.org/datasets/maestro">https://magenta.tensorflow.org/datasets/maestro</a></p> <p><b>Expected size:</b> ENST-MAPS: 32 GB, and MAESTRO: 120GB.</p> <p><b>Data utility:</b> These datasets will be used by IRCAM for sound synthesis, MIDI score augmentation, and for piano synthesis and analysis (T5.2, T5.6, and Use Case 5).</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Data are discoverable in the original sources: <a href="http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/">http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/</a> and <a href="https://magenta.tensorflow.org/datasets/maestro">https://magenta.tensorflow.org/datasets/maestro</a>.</p> <p><b>Search keywords:</b> ENST MAPS piano, MAESTRO</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> The metadata are in the standard MIDI format for computer music.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The data are already publicly shared. We will not re-share the data.</p> <p><b>How it will be accessible:</b> Can be downloaded from <a href="http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/">http://www.tsi.telecom-paristech.fr/aao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/</a> and <a href="https://magenta.tensorflow.org/datasets/maestro">https://magenta.tensorflow.org/datasets/maestro</a>.</p> <p><b>Methods/software tools to access data:</b> Web-browser</p> <p><b>Repository:</b> France Telecom repository, Magenta Framework of Google AI</p> <p><b>Restrictions on access:</b> No</p>
Making data interoperable	<p><b>Interoperability:</b> The data and metadata formats are standard, which makes the use of the dataset easy.</p> <p><b>Data and metadata vocabularies:</b> Documentations and publications are accessible on the repository websites for an explanation of the content.</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>



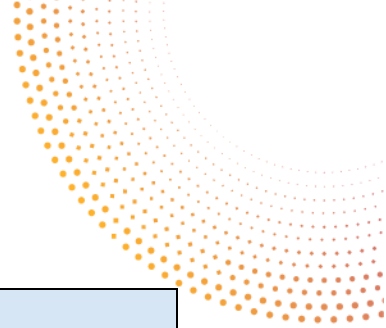


Increase data re-use	<p><u>Licence</u>: For ENST-MAPS, a subset is under Creative Commons License, and the three other subsets have no given License. MAESTRO is fully under a CC BY-NC-SA 4.0 (Non-Commercial Creative Commons).</p> <p><u>Availability for re-use</u>: Already publicly shared</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the dataset is stored in IRCAM's servers, for which access requires username/password authentication. Security measures prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: No</p>
Other Issues	No

### 5.3.20 GiantSteps music datasets

DMP component	AI4Media_Data_115_WP5_Audio_GiantSteps_v1 Partner: IRCAM
Data Summary	<p><u>Purpose</u>: The GiantSteps datasets are two collections of musical recordings annotated in tempo (664 songs) and in harmonic key (604 songs) for research purposes. These datasets are useful for training and testing tempo and key recognition methods in WP5.</p> <p><u>Type/format</u>: MP3 files for the audio recordings, and JAMS files for the annotations.</p> <p><u>Re-use of existing data</u>: Yes, we re-use a public dataset from Universitat Pompeu Fabra.</p> <p><u>Data origin</u>: Key dataset: <a href="https://github.com/GiantSteps/giantsteps-key-dataset">https://github.com/GiantSteps/giantsteps-key-dataset</a>, and Tempo dataset: <a href="https://github.com/GiantSteps/giantsteps-tempo-datase">https://github.com/GiantSteps/giantsteps-tempo-datase</a></p> <p><u>Expected size</u>: 1.8 GB</p> <p><u>Data utility</u>: This dataset will be used by IRCAM evaluate the recognition of tempo and harmonic key (T5.6, and Use Case 5).</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable in the original source: <a href="https://github.com/GiantSteps/">https://github.com/GiantSteps/</a></p> <p><u>Search keywords</u>: GiantSteps music</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: The metadata are in JAMS format (JSON Annotated Music Specification), a format for reproducible MIR research. It is documented on the following webpage: <a href="https://github.com/marl/jams">https://github.com/marl/jams</a></p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already publicly shared at <a href="https://github.com/GiantSteps/">https://github.com/GiantSteps/</a>. We will not re-share the data.</p> <p><u>How it will be accessible</u>: Can be downloaded from <a href="https://github.com/GiantSteps/">https://github.com/GiantSteps/</a></p>



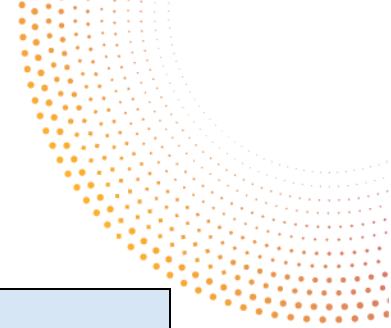


	<p><u>Methods/software tools to access data:</u> Web-browser</p> <p><u>Repository:</u> GitHub</p> <p><u>Restrictions on access:</u> No</p>
Making data interoperable	<p><u>Interoperability:</u> The data and metadata formats are standard, which makes the use of the dataset easy.</p> <p><u>Data and metadata vocabularies:</u> Documentations and publications are accessible on the repository websites for an explanation of the content.</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The datasets are already publicly shared in GitHub.</p> <p><u>Availability for re-use:</u> Already publicly shared</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> After downloading, the dataset is stored in IRCAM's servers, for which access requires username/password authentication. Security measures prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> No</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

### 5.3.21 MS COCO dataset

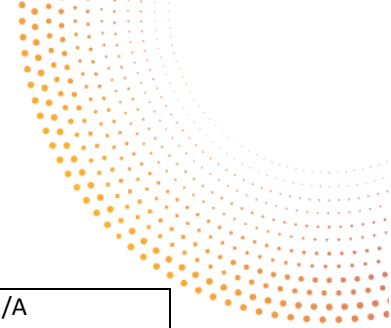
DMP component	AI4Media_Data_116_WP5_IMAGE_MSCOCO_v1 Partner: QMUL
Data Summary	<p><u>Purpose:</u> MS COCO is a dataset for object detection, object segmentation and image captioning. It contains over 330,000 images, more than 200,000 of which are labeled, featuring 80 object categories. It is used in WP5 to evaluate self-supervised methods via object detection fine-tuning.</p> <p><u>Type/format:</u> Image (jpeg)</p> <p><u>Re-use of existing data:</u> Yes</p> <p><u>Data origin:</u> <a href="https://cocodataset.org/">https://cocodataset.org/</a></p> <p><u>Expected size:</u> Train partition: 18GB, validation partition: 1 GB, Unlabeled partition: 19GB</p> <p><u>Data utility:</u> It is useful in the context of T5.3 for evaluating self-supervised methods in dense-prediction tasks. In general, this dataset is also useful for any researcher that wants to train deep learning models for object detection and segmentation or image</p>





	captioning.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data are discoverable in the MS COCO website at <a href="https://cocodataset.org/">https://cocodataset.org/</a></p> <p><u>Search keywords</u>: MS COCO dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is available at the MS COCO website at <a href="https://cocodataset.org/">https://cocodataset.org/</a></p> <p><u>How it will be accessible</u>: The data can be downloaded from the MS COCO website.</p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: The data are hosted at the MS COCO website at <a href="https://cocodataset.org/">https://cocodataset.org/</a></p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: Each image in the dataset is associated with a unique ID and contains metadata such as the file name, license, and URL. Annotations are provided in json format and include bounding boxes and segmentation masks with object categories, keypoints and captions.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset's annotations belong to the COCO Consortium and are licensed under a Creative Commons Attribution 4.0 License. The COCO Consortium does not own the copyright of the images. Use of the images must abide by the Flickr Terms of Use. The users of the images accept full responsibility for the use of the dataset, including but not limited to the use of any copies of copyrighted images that they may create from the dataset.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: Estimate the costs for making your data FAIR. Describe how you intend to cover these costs</p> <p><u>Costs for long-term preservation</u>: Describe costs and potential value of long term preservation</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded from the original source and will be stored on QMUL's servers. QMUL fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies (firewalls, right-based file system, etc.) mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Links to the dataset should be shared instead of raw data.</p>



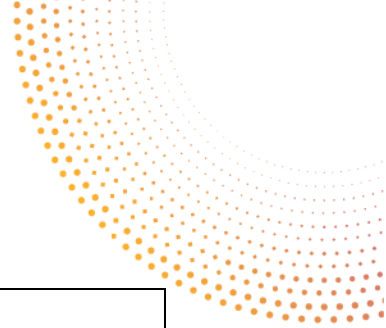


	<u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A

### 5.3.22 BVI-DVC dataset

DMP component	AI4Media_Data_117_WP5_VIDEO_BVI-DVC_v1 Partner: BSC
Data Summary	<p><u>Purpose:</u> The dataset contains 200 video clips at various resolutions, including 4K. These clips are used to generate synthetic real and fake video pairs to train and evaluate Video Super-Resolution detection models.</p> <p><u>Type/format:</u> Each video is represented by a sequence of 64 .png frames</p> <p><u>Re-use of existing data:</u> The original 4K clips from BVI-DVC are being used as “real” content. The lower resolution video frames are used to generate the “fake” content, by upscaling them with different Super-Resolution algorithms.</p> <p><u>Data origin:</u> The dataset is a collection from different sources: Videvo Free Stock Video Footage set, IRIS32 Free 4K Footage set, Harmonics database, BVI-Texture database, MCML 4K video quality database, BVI-HFR database, SJTU 4K video database, LIVE-Netflix database, Mitch Martinez Free 4K Stock Footage set, Dareful Free 4K Stock Video data set, MCL-V database, MCL-JCV database, Netflix Chimera, TUM HD databases, Ultra Video Group-Tampere University database</p> <p><u>Expected size:</u> 300GB combined for the 200 4K clips</p> <p><u>Data utility:</u> The dataset was originally created to train deep video compression algorithms. It can also be used to train and evaluate Super-Resolution models, as well as image/video quality assessment methods or other similar tasks.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data can be acquired from the website <a href="https://fan-aaron-zhang.github.io/BVI-DVC">https://fan-aaron-zhang.github.io/BVI-DVC</a>, after filling a registration form.</p> <p><u>Search keywords:</u> bvi-dvc dataset, deep video compression dataset</p> <p><u>Versioning:</u> There is one version of the dataset available</p> <p><u>Metadata creation:</u> Video metadata is included in the README, including resolution and scanning format, chroma sampling, and bit depth</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is accessible from <a href="https://fan-aaron-zhang.github.io/BVI-DVC">https://fan-aaron-zhang.github.io/BVI-DVC</a>, but requires access granted after filling a form.</p> <p><u>How it will be accessible:</u> The dataset is accessible via MS OneDrive.</p> <p><u>Methods/software tools to access data:</u> The dataset is accessible via MS OneDrive.</p> <p><u>Repository:</u> Data and relevant information are stored at <a href="https://fan-aaron-zhang.github.io/BVI-DVC">https://fan-aaron-zhang.github.io/BVI-DVC</a>.</p> <p><u>Restrictions on access:</u> Access is granted after filling a form with personal information.</p>
Making data interoperable	<p><u>Interoperability:</u> The data is interoperable, as long as original licenses and use restrictions are followed.</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p>



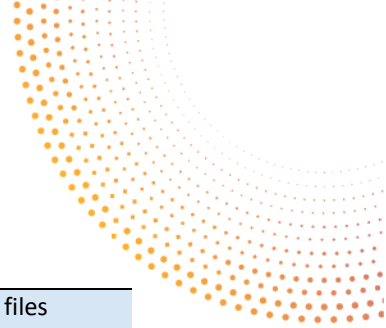


	<u>Mappings to commonly used vocabularies</u> : N/A
Increase data re-use	<p><u>License</u>: This database has been compiled by the University of Bristol, Bristol, UK, comprising sequences originally generated by various sources. All intellectual property rights remain with the originators of each sequence. The test sequences from source (15) shall only be used for academic research (no commercial use).</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: The data is openly available and is subject to licenses from creators.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The dataset will be downloaded from the original source and will be stored on BSC servers. BSC fully complies with the applicable national and European data protection framework. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.23 Adience dataset

DMP component	AI4Media_Data_118_WP5_Image_Adience_v1 Partner: RAI
Data Summary	<p><u>Purpose</u>: Adience is a dataset for gender and age classification tasks. It is used by RAI to fine tune the gender classification component.</p> <p><u>Type/format</u>: jpg</p> <p><u>Re-use of existing data</u>: yes</p> <p><u>Data origin</u>: The dataset contains about 27K images available at <a href="https://talhassner.github.io/home/projects/Adience/Adience-data.html">https://talhassner.github.io/home/projects/Adience/Adience-data.html</a></p> <p><u>Expected size</u>: 1.34 GB</p> <p><u>Data utility</u>: It is useful to WP5 partners to evaluate and benchmark face analysis models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes</p> <p><u>Search keywords</u>: adience dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: Third party website</p>



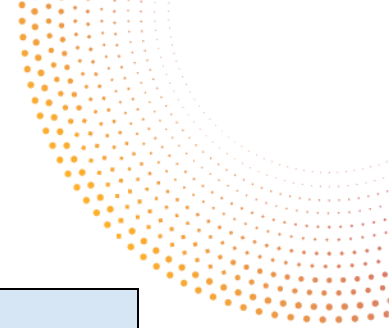


	<p><u>Methods/software tools to access data</u>: Web browser to download tar.gz files</p> <p><u>Repository</u>: <a href="https://talhassner.github.io/home/projects/Adience/Adience-data.html">https://talhassner.github.io/home/projects/Adience/Adience-data.html</a></p> <p><u>Restrictions on access</u>: Name and email required</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Academic research purpose</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: N/A</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 5.3.24 IMDB-Wiki dataset

DMP component	AI4Media_Data_119_WP5_Image_IMDB-Wiki_v1 Partner: RAI
Data Summary	<p><u>Purpose</u>: IMDB-Wiki is a publicly available dataset of face images of celebrities from IMDB and Wikipedia with age and gender labels. It is used by RAI to train the gender classification component used in UC1.</p> <p><u>Type/format</u>: jpg</p> <p><u>Re-use of existing data</u>: yes</p> <p><u>Data origin</u>: The dataset contains more than 523K images available at <a href="https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/">https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/</a></p> <p><u>Expected size</u>: 270 GB</p> <p><u>Data utility</u>: It is useful to WP5 partners to evaluate and benchmark face analysis models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable and publicly available.</p> <p><u>Search keywords</u>: imdb wiki dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>





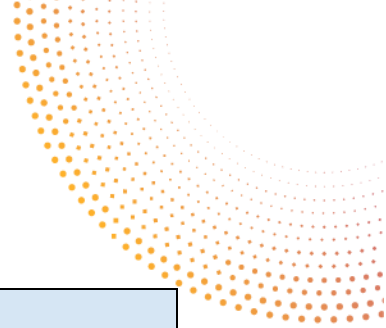
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: Third party website</p> <p><u>Methods/software tools to access data</u>: Web browser to download tar and tar.gz files</p> <p><u>Repository</u>: <a href="https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/">https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/</a></p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Academic research purpose</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: N/A</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: N/A</p>
Other Issues	N/A

## 5.4 Datasets used in the context of WP6

### 5.4.1 Deepfake Detection Challenge dataset

DMP component	AI4Media_Data_120_WP6_VIDEO_Deepfake-Detection-Challenge-Dataset_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: The Deepfake Detection Challenge dataset (DFDC) consists of more than 124k videos. The DFDC has enabled experts from around the world to come together, benchmark their deepfake detection models, try new approaches, and learn from each other's work. The dataset contains real videos and videos that have been manipulated with eight facial modification algorithms.</p> <p>This full dataset was used by participants during a Kaggle competition to create new and better models to detect manipulated media. The dataset was created by Facebook with paid actors who entered into an agreement to the use and manipulation of their likenesses in the creation of the dataset.</p> <p>The dataset will be used by CERTH in the context of T6.2 to develop and test new algorithms for the detection of deep fake videos in the web, focusing on facial</p>

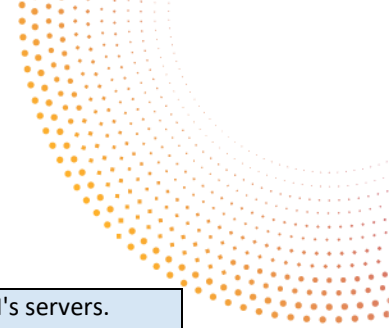




	<p>manipulation.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we are re-using an existing dataset from Kaggle.com</p> <p><u>Data origin</u>: Kaggle.com/ AWS</p> <p><u>Expected size</u>: ~500GB</p> <p><u>Data utility</u>: It is useful in the context of T6.2 for the detection of synthetic content in videos including faces. In general, this dataset is useful for any researcher that wants to train deep learning models for facial manipulation detection using a large-scale dataset.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted on Kaggle or AWS platforms and is discoverable by googling “Deepfake Detection Challenge dataset”. See <a href="https://ai.facebook.com/datasets/dfdc/">https://ai.facebook.com/datasets/dfdc/</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on Kaggle or AWS platforms. The data will not be re-shared by AI4Media partners.</p> <p><u>How it will be accessible</u>: The data is hosted on Kaggle or AWS platforms. Details on: <a href="https://ai.facebook.com/datasets/dfdc/">https://ai.facebook.com/datasets/dfdc/</a></p> <p><u>Methods/software tools to access data</u>: Creation of Kaggle account or AWS account with an IAM user and Access Keys.</p> <p><u>Repository</u>: Kaggle or AWS platforms</p> <p><u>Restrictions on access</u>: The user should accept licence agreement first.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is already interoperable.</p> <p><u>Data and metadata vocabularies</u>: Videos are in mp3 format and are accompanied with a metadata file that contains information about the authenticity of a particular video. Also, for manipulated videos there is information regarding the original video that was used to produce the manipulated video.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: To download the dataset from Kaggle the user had to agree to the challenge rules <a href="https://www.kaggle.com/c/deepfake-detection-challenge/rules">https://www.kaggle.com/c/deepfake-detection-challenge/rules</a>. The guidelines and licence for the DFDC dataset are listed in section 7. COMPETITION DATA.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>







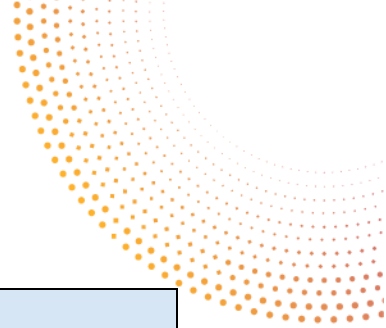
Data security	<u>Security measures</u> : The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : Licence prevents the user from re-sharing the dataset.  <u>Is informed consent for data sharing and long term preservation given</u> : N/A (Note that actors in the dataset provided their consent for the creation of these videos)
Other Issues	N/A

## 5.4.2 FaceForensics++ dataset

This dataset is also used in WP4.

DMP component	AI4Media_Data_121_WP6_VIDEO_FaceForensics++_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: FaceForensics++ is a forensics dataset consisting of 1,000 original video sequences that have been manipulated with five automated face manipulation methods: Deepfakes, Face2Face, FaceSwap, FaceShifter and NeuralTextures. The data has been sourced from 977 YouTube videos and all videos contain a trackable mostly frontal face without occlusions, which enables automated tampering methods to generate realistic forgeries. Dataset is available in 3 different video qualities.</p> <p>The dataset is used by CERTH in the context of T6.2 to develop and test new algorithms for the detection of deep fake videos on the web, focusing on facial manipulation, and in the context of T4.3 to facilitate the development and evaluation of explainable deepfake detection methods.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we are re-using existing data. Original videos are taken from YouTube.</p> <p><u>Data origin</u>: youtube.com</p> <p><u>Expected size</u>: ~400GB (all video qualities)</p> <p><u>Data utility</u>: It is useful in the context of T6.2 for the detection of synthetic content in videos including faces. It is also useful in the context of T4.3 for producing explanations about the outcomes of deepfake detection. In general, this dataset is useful for any researcher that wants to train deep learning models for facial manipulation detection using a large-scale dataset.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable in the authors' GitHub repo <a href="https://github.com/ondyari/FaceForensics">https://github.com/ondyari/FaceForensics</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: GitHub supports versioning</p> <p><u>Metadata creation</u>: N/A</p>
Making data	<u>Data openly accessible</u> : The data is already openly accessible on GitHub at



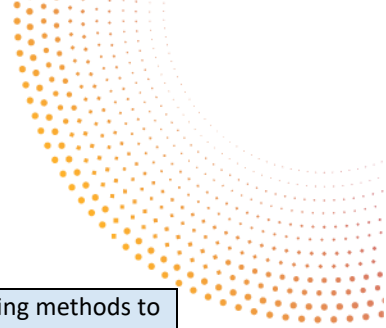


openly accessible	<p><a href="https://github.com/ondyari/FaceForensics">https://github.com/ondyari/FaceForensics</a> We will not re-share the data.</p> <p><u>How it will be accessible:</u> The data can be downloaded from the original source after filling an online form:  <a href="https://docs.google.com/forms/d/e/1FAIpQLSdRRR3L5zAv6tQ_CKxmK4W96tAab_pfBu2EKAgQbeDVhmXagg/viewform">https://docs.google.com/forms/d/e/1FAIpQLSdRRR3L5zAv6tQ_CKxmK4W96tAab_pfBu2EKAgQbeDVhmXagg/viewform</a></p> <p><u>Methods/software tools to access data:</u> Data owner provides a download script</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> The user should accept the terms of use:  <a href="http://kaldir.vc.in.tum.de/faceforensics_tos.pdf">http://kaldir.vc.in.tum.de/faceforensics_tos.pdf</a></p>
Making data interoperable	<p><u>Interoperability:</u> The file structure makes the use of the dataset easy. Original videos are in a separate folder from manipulated. Each manipulation method also appears in a separate folder.</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The data is released under the <a href="#">FaceForensics Terms of Use</a>, and the code is released under the MIT license.</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> Data already publicly shared.</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The user may provide research associates and colleagues with access to the database if they first agree to be bound by the terms and conditions.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A (Original data was mined from youtube and there is no consent from the subject appearing in the videos)</p>
Other Issues	N/A

### 5.4.3 Visual profile impact rating and ranking – ImageCLEFaware dataset

<b>DMP component</b>	AI4Media_Data_122_WP4_IMAGE_ImageCLEFaware-dataset_v1 Partner: CEA
----------------------	---





<p>Data Summary</p>	<p><b>Purpose:</b> This dataset is designed to evaluate the ability of machine learning methods to compute automatic ratings of visual user profiles made of photos in impactful real-life situations. Photos which compose the profiles have been sampled from the YFCC100M dataset. With regard to object detections, visual concept scores were crowdsourced in an experiment carried out by partner CEA. Object detectors were trained using a model trained with a combination of the publicly available MS-COCO, ImageNet and OpenImages datasets.</p> <p>The dataset will be used by partners CEA and UPB in the context of T6.7 to develop and test new algorithms for visual profile rating. A minimized version of the dataset (excluding photographs and with anonymized visual concepts) will be released publicly as part of the ImageCLEF 2021 evaluation campaign.</p> <p><b>Type/format:</b> Images (JPEG), metadata (JSON)</p> <p><b>Re-use of existing data:</b> Yes, we are re-using existing data. Original photos are from the YFCC100M dataset, which is itself collected from Flickr.</p> <p><b>Data origin:</b> flickr.com</p> <p><b>Expected size:</b> ~6MB</p> <p><b>Data utility:</b> It is useful in the context of T6.2 for providing users with feedback about potentially serious real-life effects of personal data sharing.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><b>Is data discoverable:</b> Data will be made available as a subproject of the CEA's github account: <a href="https://github.com/cea-list-lasti">https://github.com/cea-list-lasti</a></p> <p><b>Search keywords:</b> N/A</p> <p><b>Versioning:</b> GitHub supports versioning</p> <p><b>Metadata creation:</b> N/A</p>
<p>Making data openly accessible</p>	<p><b>Data openly accessible:</b> The data will be openly accessible via GitHub at <a href="https://github.com/cea-list-lasti">https://github.com/cea-list-lasti</a></p> <p><b>How it will be accessible:</b> The data can be downloaded from an online archive after completing a form.</p> <p><b>Methods/software tools to access data:</b> N/A</p> <p><b>Repository:</b> N/A</p> <p><b>Restrictions on access:</b> The user should accept the terms of use.</p>
<p>Making data interoperable</p>	<p><b>Interoperability:</b> The file structure makes the use of the dataset easy. Anonymized image detections are provided for train/val/test users in separate files. Anonymized visual concept ratings are provided per situation.</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
<p>Increase data re-use</p>	<p><b>Licence:</b> The data is released under the <a href="#">ImageCLEFaware Terms of Use</a>, and the code is released under the CC license.</p> <p><b>Availability for re-use:</b> N/A</p> <p><b>Usable by third parties after end of project:</b> Data already publicly shared.</p>



	<p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The full dataset (including images and non-anonymized metadata) will be hosted on CEA's servers. CEA fully complies with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The version of the dataset which is shared publicly includes data minimization, in compliance with art. 9 of GDPR.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 5.4.4 DEAP EEG dataset

DMP component	AI4Media_Data_123_WP6_EEG_DEAP_v1 Partner: QMUL
Data Summary	<p><u>Purpose:</u> DEAP is a dataset of human physiological signal recordings (EEG) and facial video recordings, originally created for emotion recognition purposes. QMUL will use the dataset stored in Python (.npy) and video (.avi) file format. It will be used by QMUL in T6.6 to evaluate the models developed in the task.</p> <p><u>Type/format:</u> CSV file containing metadata, 880 videos stored in .avi format.</p> <p><u>Re-use of existing data:</u> Yes, we will reuse an existing dataset</p> <p><u>Data origin:</u> <a href="https://www.eecs.qmul.ac.uk/mmv/datasets/deap/index.html">https://www.eecs.qmul.ac.uk/mmv/datasets/deap/index.html</a></p> <p><u>Expected size:</u> 12 GB</p> <p><u>Data utility:</u> It is useful to WP6 partners to evaluate EEG-based emotion recognition models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The dataset is discoverable from its website: <a href="https://www.eecs.qmul.ac.uk/mmv/datasets/deap/index.html">https://www.eecs.qmul.ac.uk/mmv/datasets/deap/index.html</a>. The dataset is stored on the servers of Queen Mary University of London.</p> <p><u>Search keywords:</u> DEAP</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The dataset is not openly accessible. Access is allowed only to users that have been given credentials after signing an End User License Agreement (EULA) form.</p> <p><u>How it will be accessible:</u> The EULA form has been sent and access has been granted, together with credentials to download it.</p> <p><u>Methods/software tools to access data:</u> The dataset is downloaded from an internet</p>



	<p>browser, without use of any other tool</p> <p><u>Repository</u>: Network Repository (<a href="http://networkrepository.com">http://networkrepository.com</a>)</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data are interoperable.</p> <p><u>Data and metadata vocabularies</u>: Vocabularies used in the dataset are clearly defined in the dataset description: <a href="https://www.eecs.qmul.ac.uk/mmv/datasets/deap/readme.html">https://www.eecs.qmul.ac.uk/mmv/datasets/deap/readme.html</a></p> <p><u>Use of standard vocabularies</u>: As defined above.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly available under an End User License Agreement: <a href="https://www.eecs.qmul.ac.uk/mmv/datasets/deap/doc/eula.pdf">https://www.eecs.qmul.ac.uk/mmv/datasets/deap/doc/eula.pdf</a></p> <p><u>Availability for re-use</u>: The dataset is available for re-use, only under the terms of the EULA form.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the data are stored in the servers of Queen Mary University of London. Access requires username/password authentication. The security measures taken prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The EULA terms clearly prevent sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Participants provided their consent for the creation of the dataset)</p>
Other Issues	N/A

#### 5.4.5 SEED EEG dataset

DMP component	AI4Media_Data_124_WP6_EEG_SEED_v1 Partner: QMUL
Data Summary	<p><u>Purpose</u>: SEED is a dataset of human physiological signal recordings (EEG), originally created for emotion recognition purposes. We will use the dataset stored in Matlab (.mat) file format. It will be used in T6.6 to evaluate the models developed in the task.</p> <p><u>Type/format</u>: Excel (.xls) files containing metadata, EEG signals and EEG features stored in .mat format.</p> <p><u>Re-use of existing data</u>: Yes, we will reuse an existing dataset</p> <p><u>Data origin</u>: <a href="http://bcmi.sjtu.edu.cn/~seed/seed.html">http://bcmi.sjtu.edu.cn/~seed/seed.html</a></p> <p><u>Expected size</u>: 10 GB</p>



	<p><u>Data utility</u>: It is useful to WP6 partners to evaluate EEG-based emotion recognition models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The dataset is discoverable from its website: <a href="http://bcmi.sjtu.edu.cn/~seed/seed.html">http://bcmi.sjtu.edu.cn/~seed/seed.html</a>. Access to the dataset is handled by Shanghai Jiao Tong University.</p> <p><u>Search keywords</u>: SEED</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The dataset is not openly accessible. Access is allowed only to users that have been given credentials after signing an End User License Agreement (EULA) form.</p> <p><u>How it will be accessible</u>: The EULA form has been sent and access has been granted, together with credentials to download it.</p> <p><u>Methods/software tools to access data</u>: The dataset is downloaded from an internet browser, without use of any other tool</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: None</p>
Making data interoperable	<p><u>Interoperability</u>: The data are interoperable.</p> <p><u>Data and metadata vocabularies</u>: Vocabularies used in the dataset are clearly defined in the dataset description: <a href="http://bcmi.sjtu.edu.cn/~seed/seed.html">http://bcmi.sjtu.edu.cn/~seed/seed.html</a></p> <p><u>Use of standard vocabularies</u>: As defined above.</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly available under an End User License Agreement: <a href="http://bcmi.sjtu.edu.cn/~seed/resource/license/license">http://bcmi.sjtu.edu.cn/~seed/resource/license/license</a>.</p> <p><u>Availability for re-use</u>: The dataset is available for re-use, only under the terms of the EULA form.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the data are stored in the servers of Queen Mary University of London. Access requires username/password authentication. The security measures taken prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The EULA terms clearly prevent sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Participants provided their consent for the creation of the dataset)</p>
Other Issues	N/A

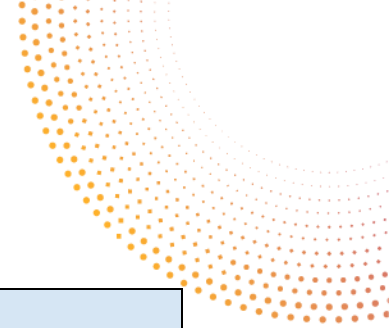




#### 5.4.6 SEED-IV EEG dataset

DMP component	AI4Media_Data_125_WP6_EEG_SEED-IV_v1 Partner: QMUL
Data Summary	<p><b>Purpose:</b> SEED-IV is a dataset of human physiological signal recordings (EEG), originally created for emotion recognition purposes. We will use the dataset stored in Matlab (.mat) file format. It will be used in T6.6 by QMUL to evaluate the models developed in the task.</p> <p><b>Type/format:</b> Excel (.xls) files containing metadata, EEG signals and EEG features stored in .mat format.</p> <p><b>Re-use of existing data:</b> Yes, we will reuse an existing dataset</p> <p><b>Data origin:</b> <a href="http://bcmi.sjtu.edu.cn/~seed/seed-iv.html">http://bcmi.sjtu.edu.cn/~seed/seed-iv.html</a></p> <p><b>Expected size:</b> 7 GB</p> <p><b>Data utility:</b> It is useful to WP6 partners to evaluate EEG-based emotion recognition models.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> <a href="http://bcmi.sjtu.edu.cn/~seed/seed-iv.html">http://bcmi.sjtu.edu.cn/~seed/seed-iv.html</a>. Access to the dataset is handled by Shanghai Jiao Tong University.</p> <p><b>Search keywords:</b> SEED-IV</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The dataset is not openly accessible. Access is allowed only to users that have been given credentials after signing an End User License Agreement (EULA) form.</p> <p><b>How it will be accessible:</b> The EULA form has been sent and access has been granted, together with credentials to download it.</p> <p><b>Methods/software tools to access data:</b> The dataset is downloaded from an internet browser, without use of any other tool.</p> <p><b>Repository:</b> N/A</p> <p><b>Restrictions on access:</b> None</p>
Making data interoperable	<p><b>Interoperability:</b> The data are interoperable.</p> <p><b>Data and metadata vocabularies:</b> Vocabularies used in the dataset are clearly defined in the dataset description: <a href="http://bcmi.sjtu.edu.cn/~seed/seed-iv.html">http://bcmi.sjtu.edu.cn/~seed/seed-iv.html</a></p> <p><b>Use of standard vocabularies:</b> As defined above.</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> The dataset is already publicly available under an End User License Agreement: <a href="http://bcmi.sjtu.edu.cn/~seed/resource/license/license-SEED-IV.pdf">http://bcmi.sjtu.edu.cn/~seed/resource/license/license-SEED-IV.pdf</a></p> <p><b>Availability for re-use:</b> The dataset is available for re-use, only under the terms of the EULA form.</p> <p><b>Usable by third parties after end of project:</b> N/A</p>





	<p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> After downloading, the data are stored in the servers of Queen Mary University of London. Access requires username/password authentication. The security measures taken prevent illegitimate access (firewalls and rights-based-file system).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The EULA terms clearly prevent sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A (Participants provided their consent for the creation of the dataset)</p>
Other Issues	N/A

#### 5.4.7 Clotho audio captioning dataset

<b>DMP component</b>	<p><b>AI4Media_Data_126_WP6_Audio_Clotho_v1</b></p> <p><b>Partner: CERTH</b></p>
Data Summary	<p><u>Purpose:</u> The Clotho dataset consists of audio samples of 15 to 30 seconds duration, in WAVE format, and each audio sample has five captions of eight to 20 words length. There is a total of 4,981 audio samples in the dataset with 24,905 captions.</p> <p>The dataset will be used by CERTH in the context of T6.2 to develop and test new algorithms for automated audio captioning, where general audio content is described using free text. The final system will use an input audio signal and it will output the textual description (i.e., the caption) of that signal.</p> <p><u>Type/format:</u> Audio (wav)</p> <p><u>Re-use of existing data:</u> Yes, we are reusing an existing dataset.</p> <p><u>Data origin:</u> <a href="https://zenodo.org/record/3490684">https://zenodo.org/record/3490684</a></p> <p><u>Expected size:</u> 6 GB</p> <p><u>Data utility:</u> It is useful in the context of T6.2 for automated audio captioning. In general, this dataset is useful for any researcher that wants to train deep learning models for audio signal processing by using a dataset captured in the wild.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable in the authors' GitHub repo <a href="https://github.com/audio-captioning/clotho-dataset">https://github.com/audio-captioning/clotho-dataset</a>.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> GitHub supports versioning</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is already openly accessible on GitHub at <a href="https://github.com/audio-captioning/clotho-dataset">https://github.com/audio-captioning/clotho-dataset</a>. We will not re-share the data.</p> <p><u>How it will be accessible:</u> It can be accessed via running the '.sh' scripts provided by the authors on the aforementioned GitHub link, or it can be downloaded directly from</p>



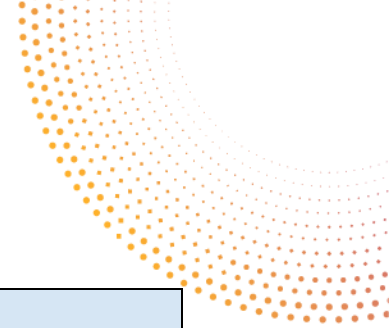


	<p>Zenodo.</p> <p><u>Methods/software tools to access data</u>: Data owner provides a download script.</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. The authors have provided all the necessary files for training, validation and testing, separately, along with the ground truth for each file.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is released under the <a href="#">Tampere University, Finland Terms of Use</a>.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The user may provide research associates and colleagues with access to the database if they first agree to be bound by the terms and conditions. The data will not be reshared by CERTH.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Original data was mined from Zenodo and there is no consent from the subject that is heard in the recordings).</p>
Other Issues	No

#### 5.4.8 ASVspoo2019 dataset

<b>DMP component</b>	<b>AI4Media_Data_127_WP6_Audio_ASVspoo_v1</b> <b>Partner: CERTH</b>
Data Summary	<p><u>Purpose</u>: The ASVspoo 2019 dataset consists of audio samples from 107 speakers (46 males and 61 females). It was released by the University of Edinburgh in collaboration with Google. The database encompasses two partitions for the assessment of logical access and physical access scenarios. The original waveform format is PCM and compressed using FLAC. No telephone or mobile codec was used.</p> <p>The dataset will be used by CERTH in the context of T6.2 to develop and test new</p>





	<p>algorithms for deepfake detection, focusing on manipulated speech.</p> <p><u>Type/format</u>: Audio (flac)</p> <p><u>Re-use of existing data</u>: Yes, we use an existing dataset</p> <p><u>Data origin</u>: Audio dataset created by Uni Edinburgh and Google <a href="https://datashare.ed.ac.uk/handle/10283/3336">https://datashare.ed.ac.uk/handle/10283/3336</a></p> <p><u>Expected size</u>: ~ 25 GB</p> <p><u>Data utility</u>: It is useful in the context of T6.2 for the detection of fake audios to spoof automatic speaker verification (ASV) systems. The dataset is suited not only to study of ASV replay spoofing and countermeasures, but also the study of fake audio detection in the case of e.g., smart home devices that facilitate the security of online banking and payment solutions.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data is discoverable in the University of Edinburgh DataShare repo: <a href="https://datashare.ed.ac.uk/handle/10283/3336">https://datashare.ed.ac.uk/handle/10283/3336</a></p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: The DataShare repository supports versioning (e.g., 2013, 2015, 2017 and 2019)</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data is already openly accessible on the DataShare repo of the University of Edinburgh. We will not re-share the data.</p> <p><u>How it will be accessible</u>: It can be accessed via downloading the necessary parts from the DataShare repo</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: DataShare repo of the University of Edinburgh</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. The authors have provided the evaluation plan for the database in the following link: <a href="https://www.asvspoof.org/asvspoof2019/asvspoof2019_evaluation_plan.pdf">https://www.asvspoof.org/asvspoof2019/asvspoof2019_evaluation_plan.pdf</a></p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is released under the Open Data Commons Attribution License</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers.</p>



	CERTH fully complies with the applicable national and European data protection frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The user may provide research associates and colleagues with access to the database if they first agree to be bound by the terms and conditions</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A(Original data was mined from the University of Edinburgh and there is no consent from the subject heard in the recordings).</p>
Other Issues	No

#### 5.4.9 MOBIPHONE audio dataset

DMP component	AI4Media_Data_128_WP6_Audio_MOBIPHONE_v1 Partner: FhG-IDMT
Data Summary	<p><u>Purpose:</u> MOBIPHONE is a forensics dataset consisting of audio recordings acquired by 21 mobile phones of various models from 7 different brands, collected by recording 10 utterances, uttered by 12 male speakers and another 12 female speakers, randomly chosen from the TIMIT database.</p> <p>The dataset is the only publicly available dataset for microphone classification at present and has been used by several publications in this domain for benchmarking and comparison with other state-of-the-art algorithms.</p> <p><u>Type/format:</u> Audio (wav)</p> <p><u>Re-use of existing data:</u> Yes, re-use of utterances from the TIMIT speech corpus.</p> <p><u>Data origin:</u> The data was created by the authors themselves, by replaying a subset of TIMIT utterances with a high-end loudspeaker, and recording the outcome with several devices at once</p> <p><u>Expected size:</u> ~941MB</p> <p><u>Data utility:</u> The dataset is useful in the context of T6.2 for developing and testing new algorithms for microphone classification and device identification, and for any research focused on manipulation detection on the basis of changes in the recording device.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Data is discoverable by reading the authors' paper on <a href="https://ieeexplore.ieee.org/document/6900732">https://ieeexplore.ieee.org/document/6900732</a>.</p> <p>A Google search for "microphone classification dataset" or "MOBIPHONE dataset" is not sufficient to discover the dataset.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> The data is already openly accessible on Dropbox at <a href="https://www.dropbox.com/sh/9n7fy7moi825bgk/WFLBKxUitV">https://www.dropbox.com/sh/9n7fy7moi825bgk/WFLBKxUitV</a>. We will not re-share the data.</p>



	<p><u>How it will be accessible</u>: Already accessible at the original dropbox location.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. Audios are separated in folders, with each folder corresponding to a single microphone. Splitting the data for training, validation and testing is left to the users</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data is not released under any license. Further communications with the authors are necessary to clear and clarify usage rights.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset is stored on Dropbox servers, which have mechanisms in place for minimizing the risk of data loss.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: From a legal perspective, the lack of license prevents usage of all sorts, including sharing and re-hosting. No ethical aspects are present which may prevent re-distribution. We plan on contacting the original authors for clearing any license issues, before making use of the data.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: All recordings in the dataset have been acquired in controlled conditions and with the informed consent of all the speakers involved, in compliance with the GDPR).</p>
Other Issues	N/A

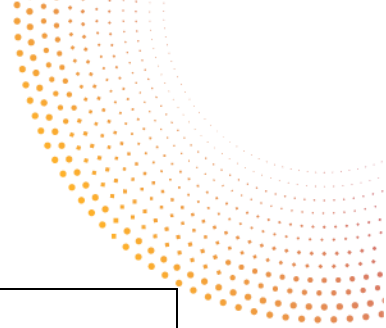
#### 5.4.10 Fake-or-Real (FoR) audio dataset

<b>DMP component</b>	<b>AI4Media_Data_129_WP6_Audio_Fake-or-Real_v1</b> <b>Partner: FhG-IDMT</b>
Data Summary	<p><u>Purpose</u>: The Fake-or-Real (FoR) dataset is a collection of utterances from real humans and computer generated speech. The dataset can be used to train classifiers for synthetic speech detection.</p> <p>The dataset aggregates data from the latest TTS solutions (such as Deep Voice 3 and Google Wavenet TTS) as well as a variety of real human speech, including pre-existing speech datasets, and the authors' own speech recordings.</p> <p>The data has been normalized in terms of speakers' genre, but is not clear how many of</p>



	<p>the natural voices have a synthetic counterpart for avoiding inherent biases. Nevertheless, it's the only database of this kind ever release to the research community.</p> <p><u>Type/format</u>: Audio (wav)</p> <p><u>Re-use of existing data</u>: The datasets is re-using recordings from the Arctic Dataset (<a href="http://festvox.org/cmu_arctic/">http://festvox.org/cmu_arctic/</a>), LJSpeech Dataset (<a href="https://keithito.com/LJ-Speech-Dataset/">https://keithito.com/LJ-Speech-Dataset/</a>), and VoxForge Dataset (<a href="http://www.voxforge.org">http://www.voxforge.org</a>).</p> <p><u>Data origin</u>: Previous datasets, own recordings from the authors</p> <p><u>Expected size</u>: ~22GB (four different variants)</p> <p><u>Data utility</u>: The dataset is useful in the context of T6.2 for the detection of synthetic speech content in audio files. The authors also included a variant of the dataset acquired by simulating a re-recording, to let the researcher community also address the problem of washed-up synthetic recordings</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: Data is discoverable by reading the authors' paper on <a href="https://ieeexplore.ieee.org/document/8906599">https://ieeexplore.ieee.org/document/8906599</a>.</p> <p>A Google search for "synthetic speech dataset" is sufficient to discover the authors' paper, and thus the dataset itself.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The data is already openly accessible on the authors' institutional page <a href="https://bil.eecs.yorku.ca/datasets/">https://bil.eecs.yorku.ca/datasets/</a>. We do not plan to re-share the data.</p> <p><u>How it will be accessible</u>: The data can be downloaded from the original source.</p> <p><u>Methods/software tools to access data</u>: Web-browser to download the data as zip file.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: The file structure makes the use of the dataset easy. The files are split into disjoint training, validation and testing set. Moreover, for each split, original audios are in a separate folder from synthetic ones.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
<p>Increase data re-use</p>	<p><u>Licence</u>: The data is released under the <a href="#">GNU General Public License (V3)</a>.</p> <p><u>Availability for re-use</u>: Yes, if the license terms and conditions are fulfilled.</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>



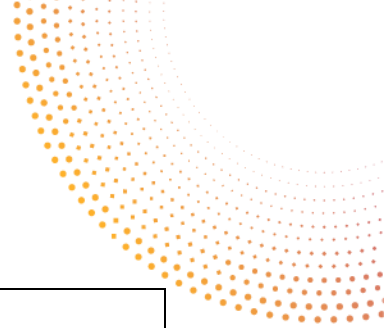


Allocation of resources	<u>Costs for making data FAIR:</u> N/A <u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on FHG-IDMT's servers. FHG-IDMT fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> N/A <u>Is informed consent for data sharing and long term preservation given:</u> All recordings in the dataset have been acquired in controlled conditions and with the informed consent of all the speakers involved, in compliance with the GDPR).
Other Issues	N/A

#### 5.4.11 ForenSynths dataset

DMP component	AI4Media_Data_130_WP6_Image_ForenSynths_v1 Partner: CERTH
Data Summary	<u>Purpose:</u> This dataset consists of synthesized images from 11 models (GAN and Deepfakes). The ProGAN-generated images are used for training purposes while all 11 models are used for evaluation purposes. <u>Type/format:</u> Images (png) <u>Re-use of existing data:</u> We re-use the dataset introduced in <a href="https://arxiv.org/abs/1912.11035">https://arxiv.org/abs/1912.11035</a> <u>Data origin:</u> <a href="https://github.com/peterwang512/CNNDetection">https://github.com/peterwang512/CNNDetection</a> <u>Expected size:</u> 80GB <u>Data utility:</u> CERTH uses this dataset to train and evaluate synthetic image detectors in T6.2. Useful for researchers working on sythetic image detection.
Making data findable, incl. provisions for metadata	<u>Is data discoverable:</u> Data discoverable through this link <a href="https://github.com/peterwang512/CNNDetection">https://github.com/peterwang512/CNNDetection</a> . <u>Search keywords:</u> ForenSynths dataset <u>Versioning:</u> N/A <u>Metadata creation:</u> N/A
Making data openly accessible	<u>Data openly accessible:</u> The data are already publicly available in GitHub at <a href="https://github.com/peterwang512/CNNDetection">https://github.com/peterwang512/CNNDetection</a> <u>How it will be accessible:</u> The repository provides Google Drive links to download the data. <u>Methods/software tools to access data:</u> Web browser. <u>Repository:</u> GitHub ( <a href="https://github.com/peterwang512/CNNDetection">https://github.com/peterwang512/CNNDetection</a> ) <u>Restrictions on access:</u> No
Making data interoperable	<u>Interoperability:</u> Standard formats are considered in this dataset, the data are easy to download and use.



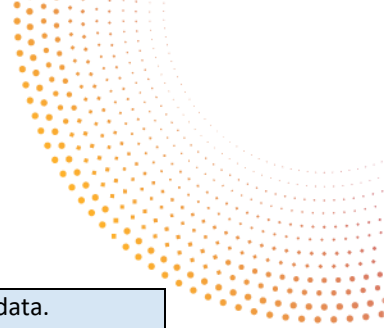


	<p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly shared in GitHub.</p> <p><u>Availability for re-use</u>: Already publicly shared.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the dataset is stored in CERTH's password protected servers.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 5.4.12 SynthBuster dataset

DMP component	AI4Media_Data_131_WP6_IMAGE_SynthBuster_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: This dataset contains AI-generated images from 9 different models, namely DALL-E 2, DALL-E 3, Adobe Firefly, Midjourney v5, Stable Diffusion 1.3, Stable Diffusion 1.4, Stable Diffusion 2, Stable Diffusion XL, and Glide. It is combined with real images from RAISE (<a href="https://dl.acm.org/doi/10.1145/2713168.2713194">https://dl.acm.org/doi/10.1145/2713168.2713194</a>) in order to evaluate binary synthetic image detection classifiers.</p> <p><u>Type/format</u>: Images (png)</p> <p><u>Re-use of existing data</u>: We re-use the dataset introduced in <a href="https://ieeexplore.ieee.org/document/10334046">https://ieeexplore.ieee.org/document/10334046</a></p> <p><u>Data origin</u>: <a href="https://zenodo.org/records/10066460">https://zenodo.org/records/10066460</a></p> <p><u>Expected size</u>: 12GB</p> <p><u>Data utility</u>: CERTH uses this dataset to train and evaluate synthetic image detectors in T6.2. Useful for researchers working on sythetic image detection.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data discoverable through this link <a href="https://zenodo.org/records/10066460">https://zenodo.org/records/10066460</a>.</p> <p><u>Search keywords</u>: Synthbuster dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly	<p><u>Data openly accessible</u>: The data are already publicly available in <a href="https://zenodo.org/records/10066460">https://zenodo.org/records/10066460</a></p>





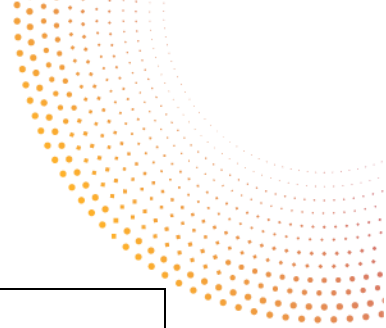
accessible	<p><u>How it will be accessible</u>: The repository provides a link to download the data.</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: Zenodo</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly shared in Zenodo.</p> <p><u>Availability for re-use</u>: Already publicly shared.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the dataset is stored in CERTH's password protected servers.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 5.4.13 Latent Diffusion Training-set

DMP component	AI4Media_Data_132_WP6_Image_LatentDiffusion_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: This dataset contains AI-generated images from the Latent Diffusion model. We use it in order to train and evaluate binary synthetic image detection classifiers.</p> <p><u>Type/format</u>: Images (png)</p> <p><u>Re-use of existing data</u>: We re-use the dataset introduced in <a href="https://ieeexplore.ieee.org/document/10095167">https://ieeexplore.ieee.org/document/10095167</a></p> <p><u>Data origin</u>: <a href="https://github.com/grip-unina/DMImageDetection">https://github.com/grip-unina/DMImageDetection</a></p> <p><u>Expected size</u>: 21GB</p> <p><u>Data utility</u>: CERTH uses this dataset to train and evaluate synthetic image detectors in T6.2. Useful for researchers working on sythetic image detection.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Data discoverable through this link <a href="https://github.com/grip-unina/DMImageDetection">https://github.com/grip-unina/DMImageDetection</a>.</p> <p><u>Search keywords</u>: On the detection of synthetic images generated by diffusion models</p>





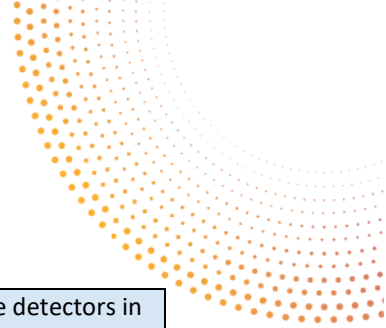


	<p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data are already publicly available in <a href="https://github.com/grip-unina/DMImageDetection">https://github.com/grip-unina/DMImageDetection</a></p> <p><u>How it will be accessible</u>: The repository provides a link to download the data.</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: GitHub</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset is already publicly shared in GitHub.</p> <p><u>Availability for re-use</u>: Already publicly shared.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: After downloading, the dataset is stored in CERTH's password protected servers.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: No</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

#### 5.4.14 Diffusion datasets

DMP component	AI4Media_Data_133_WP6_Image_DiffusionDatasets_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: This dataset contains AI-generated images from Guided, Latent Diffusion, Glide, and DALL-E models used in order to evaluate binary synthetic image detection classifiers.</p> <p><u>Type/format</u>: Images (png)</p> <p><u>Re-use of existing data</u>: We re-use the dataset introduced in <a href="https://arxiv.org/abs/2302.10174">https://arxiv.org/abs/2302.10174</a></p> <p><u>Data origin</u>: <a href="https://github.com/WisconsinAI/Vision/UniversalFakeDetect">https://github.com/WisconsinAI/Vision/UniversalFakeDetect</a></p> <p><u>Expected size</u>: 895MB</p>





	<b>Data utility:</b> CERTH uses this dataset to train and evaluate synthetic image detectors in T6.2. Useful for researchers working on sythetic image detection.
Making data findable, incl. provisions for metadata	<b>Is data discoverable:</b> Data discoverable through this link <a href="https://github.com/WisconsinAIVision/UniversalFakeDetect">https://github.com/WisconsinAIVision/UniversalFakeDetect</a> . <b>Search keywords:</b> Towards Universal Fake Image Detectors that Generalize Across Generative Models <b>Versioning:</b> N/A <b>Metadata creation:</b> N/A
Making data openly accessible	<b>Data openly accessible:</b> The data are already publicly available in <a href="https://github.com/WisconsinAIVision/UniversalFakeDetect">https://github.com/WisconsinAIVision/UniversalFakeDetect</a> <b>How it will be accessible:</b> The repository provides a link to download the data. <b>Methods/software tools to access data:</b> Web browser. <b>Repository:</b> GitHub <b>Restrictions on access:</b> No
Making data interoperable	<b>Interoperability:</b> Standard formats are considered in this dataset, the data are easy to download and use. <b>Data and metadata vocabularies:</b> N/A <b>Use of standard vocabularies:</b> N/A <b>Mappings to commonly used vocabularies:</b> N/A
Increase data re-use	<b>Licence:</b> The dataset is already publicly shared in GitHub. <b>Availability for re-use:</b> Already publicly shared. <b>Usable by third parties after end of project:</b> N/A <b>Re-use timeframe:</b> N/A <b>Data quality assurance process:</b> N/A
Allocation of resources	<b>Costs for making data FAIR:</b> N/A <b>Costs for long-term preservation:</b> N/A
Data security	<b>Security measures:</b> After downloading, the dataset is stored in CERTH's password protected servers.
Ethical aspects	<b>Possible ethical and legal aspects preventing sharing:</b> No <b>Is informed consent for data sharing and long term preservation given:</b> N/A
Other Issues	N/A

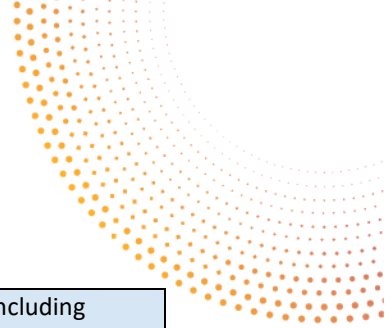
#### 5.4.15 FakeAVCeleb dataset

<b>DMP component</b>	<b>AI4Media_Data_134_WP6_VIDEO_FakeAVCeleb_v1</b> <b>Partner:</b> CERTH
<b>Data Summary</b>	<b>Purpose:</b> The Fake Audio Visual Celeb (FakeAVCeleb) dataset consists of 20,000 videos, of those only 500 are real and the remaining 19,500 are manipulated. The manipulation methods contained in the dataset are Faceswap, FSGAN, Wav2Lip and



	<p>SV2TTS. We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we re-use already published data.</p> <p><u>Data origin</u>: <a href="https://sites.google.com/view/fakeavcelebdash-lab/?pli=1">https://sites.google.com/view/fakeavcelebdash-lab/?pli=1</a></p> <p><u>Expected size</u>: ~163GB</p> <p><u>Data utility</u>: Used in T6.2 for training and testing video deepfake detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted here <a href="https://sites.google.com/view/fakeavcelebdash-lab/?pli=1">https://sites.google.com/view/fakeavcelebdash-lab/?pli=1</a></p> <p><u>Search keywords</u>: FakeAVCeleb</p> <p><u>Versioning</u>: Yes</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, at <a href="https://sites.google.com/view/fakeavcelebdash-lab/?pli=1">https://sites.google.com/view/fakeavcelebdash-lab/?pli=1</a></p> <p><u>How it will be accessible</u>: The data is hosted here <a href="https://sites.google.com/view/fakeavcelebdash-lab/?pli=1">https://sites.google.com/view/fakeavcelebdash-lab/?pli=1</a>.</p> <p><u>Methods/software tools to access data</u>: Browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: The user should accept a licence agreement first: <a href="https://docs.google.com/forms/d/e/1FAIpQLSfPDD3oV0auqmmWEgCSaTEQ6CGpFeB-ozQJ35x-B_0Xjd93bw/viewform">https://docs.google.com/forms/d/e/1FAIpQLSfPDD3oV0auqmmWEgCSaTEQ6CGpFeB-ozQJ35x-B_0Xjd93bw/viewform</a></p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: <a href="https://docs.google.com/forms/d/e/1FAIpQLSfPDD3oV0auqmmWEgCSaTEQ6CGpFeB-ozQJ35x-B_0Xjd93bw/viewform">https://docs.google.com/forms/d/e/1FAIpQLSfPDD3oV0auqmmWEgCSaTEQ6CGpFeB-ozQJ35x-B_0Xjd93bw/viewform</a></p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures</p>



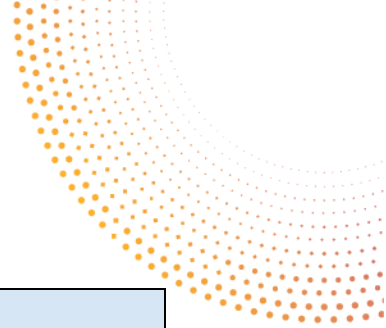


	and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset provided their consent for the creation of these videos)</p>
Other Issues	N/A

#### 5.4.16 ForgeryNet dataset

DMP component	AI4Media_Data_135_WP6_VIDEO_ForgeryNet_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: The ForgeryNet dataset consists of 221,247 videos, of those only 99,630 are real and the remaining 121,617 are manipulated with 10 total manipulation methods.</p> <p>We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we re-use published data.</p> <p><u>Data origin</u>: <a href="https://arxiv.org/pdf/2103.05630.pdf">https://arxiv.org/pdf/2103.05630.pdf</a></p> <p><u>Expected size</u>: 500GB</p> <p><u>Data utility</u>: Used in T6.2 for training and testing video deepfake detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is hosted here <a href="https://yinanhe.github.io/projects/forgerynet.html#download">https://yinanhe.github.io/projects/forgerynet.html#download</a></p> <p><u>Search keywords</u>: ForgeryNet dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: The data is hosted here <a href="https://yinanhe.github.io/projects/forgerynet.html#download">https://yinanhe.github.io/projects/forgerynet.html#download</a></p> <p><u>Methods/software tools to access data</u>: Browser.</p> <p><u>Repository</u>: Google drive</p> <p><u>Restrictions on access</u>: The user should accept a licence agreement first.</p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<u>Licence</u> : N/A





	<p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset provided their consent for the creation of these videos)</p>
Other Issues	N/A

#### 5.4.17 Korean DeepFake dataset

DMP component	AI4Media_Data_136_WP6_VIDEO_KoDF_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: The Korean DeepFake (KoDF) dataset consists of 237,942 videos, of those, 62,166 are real and the remaining 175,776 are manipulated with 6 manipulation methods.</p> <p>We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we re-use published data.</p> <p><u>Data origin</u>: <a href="https://deepbrainai-research.github.io/kodf/">https://deepbrainai-research.github.io/kodf/</a></p> <p><u>Expected size</u>: 4000GB (4TB)</p> <p><u>Data utility</u>: Used in T6.2 for training and testing video deepfake detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, The data can be found at <a href="https://deepbrainai-research.github.io/kodf/">https://deepbrainai-research.github.io/kodf/</a></p> <p><u>Search keywords</u>: Korean DeepFake dataset, KoDF</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, at <a href="https://deepbrainai-research.github.io/kodf/">https://deepbrainai-research.github.io/kodf/</a></p> <p><u>How it will be accessible</u>: The data can be found at <a href="https://deepbrainai-research.github.io/kodf/">https://deepbrainai-research.github.io/kodf/</a></p>

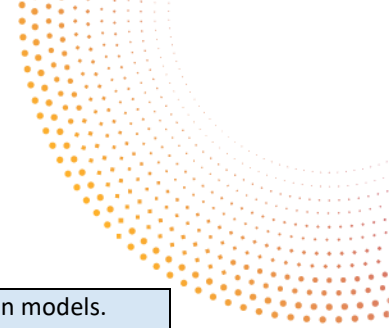


	<p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: <a href="https://deepbrainai-research.github.io/kodf/">https://deepbrainai-research.github.io/kodf/</a></p> <p><u>Restrictions on access</u>: The user should accept a <a href="#">licence agreement</a> first.</p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The user should accept this <a href="#">licence agreement</a> first.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset provided their consent for the creation of these videos)</p>
Other Issues	N/A

#### 5.4.18 WildDeepFake dataset

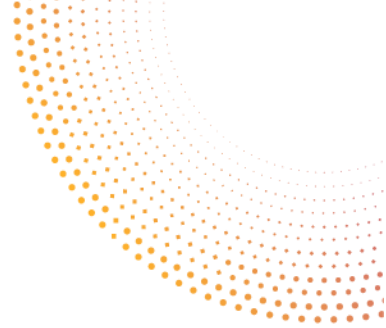
<b>DMP component</b>	<b>AI4Media_Data_137_WP6_VIDEO_WildDeepFake_v1</b> <b>Partner: CERTH</b>
Data Summary	<p><u>Purpose</u>: The WildDeepFake (WDF) dataset consists of 7,314 face sequences extracted from 707 deepfake videos; of those, 3,805 are real and the remaining 3,509 are fake. All of them were sourced from the internet.</p> <p>We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we re-use published data.</p> <p><u>Data origin</u>: <a href="https://github.com/deepfakeinthewild/deepfake-in-the-wild">https://github.com/deepfakeinthewild/deepfake-in-the-wild</a></p> <p><u>Expected size</u>: 100GB</p>





	<u>Data utility</u> : Used in T6.2 for training and testing video deepfake detection models.
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data can be found on GitHub at <a href="https://github.com/deepfakeinthewild/deepfake-in-the-wild">https://github.com/deepfakeinthewild/deepfake-in-the-wild</a></p> <p><u>Search keywords</u>: WildDeepFake Dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available on GitHub: <a href="https://github.com/deepfakeinthewild/deepfake-in-the-wild">https://github.com/deepfakeinthewild/deepfake-in-the-wild</a></p> <p><u>How it will be accessible</u>: Can be downloaded from GitHub</p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: GitHub: <a href="https://github.com/deepfakeinthewild/deepfake-in-the-wild">https://github.com/deepfakeinthewild/deepfake-in-the-wild</a></p> <p><u>Restrictions on access</u>: The user should accept a licence agreement first: <a href="https://docs.google.com/forms/d/e/1FAIpQLSfN-CrnxDz_Furv0KzcNdO_Nzf_3Gpy4s-P4qRjKBJuD2CaEA/viewform">https://docs.google.com/forms/d/e/1FAIpQLSfN-CrnxDz_Furv0KzcNdO_Nzf_3Gpy4s-P4qRjKBJuD2CaEA/viewform</a></p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Licence prevents the user from re-sharing the dataset: <a href="https://docs.google.com/forms/d/e/1FAIpQLSfN-CrnxDz_Furv0KzcNdO_Nzf_3Gpy4s-P4qRjKBJuD2CaEA/viewform">https://docs.google.com/forms/d/e/1FAIpQLSfN-CrnxDz_Furv0KzcNdO_Nzf_3Gpy4s-P4qRjKBJuD2CaEA/viewform</a></p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset have NOT provided their consent for the usage of these videos)</p>
Other Issues	N/A



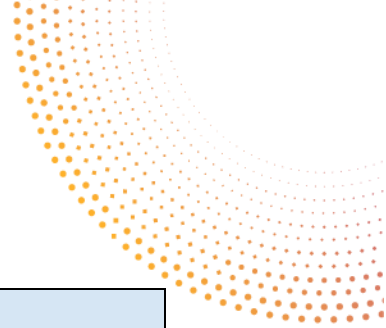


#### 5.4.19 DeepFakes from Different Models (DFDM) dataset

DMP component	AI4Media_Data_138_WP6_VIDEO_DFDM_v1 Partner: CERTH
Data Summary	<p><u>Purpose</u>: The DeepFakes From Different Models (DFDM) dataset consists of 6,450 fake videos sourced from YouTube and manipulated with 5 manipulation methods.</p> <p>We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format</u>: Videos (mp4)</p> <p><u>Re-use of existing data</u>: Yes, we re-use published data.</p> <p><u>Data origin</u>: <a href="https://github.com/shanface33/Deepfake_Model_Attribution">https://github.com/shanface33/Deepfake_Model_Attribution</a></p> <p><u>Expected size</u>: 183GB</p> <p><u>Data utility</u>: Used in T6.2 for training and testing video deepfake detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data can be found in GitHub at <a href="https://github.com/shanface33/Deepfake_Model_Attribution">https://github.com/shanface33/Deepfake_Model_Attribution</a></p> <p><u>Search keywords</u>: DeepFakes From Different Models, DFDM dataset, Model Attribution of Face-swap Deepfake Videos</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes, available on GitHub at <a href="https://github.com/shanface33/Deepfake_Model_Attribution">https://github.com/shanface33/Deepfake_Model_Attribution</a></p> <p><u>How it will be accessible</u>: From GitHub</p> <p><u>Methods/software tools to access data</u>: Browser</p> <p><u>Repository</u>: GitHub</p> <p><u>Restrictions on access</u>: The DFDM database is released only for academic research. Researchers from educational institute are allowed to use this database freely for noncommercial purpose.</p>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The DFDM database is released only for academic research. Researchers from educational institute are allowed to use this database freely for noncommercial purpose. Users should sign a form: <a href="https://docs.google.com/forms/d/e/1FAIpQLSeM-1pJ13RyPVgF0bGRQtLiupwWDvALD6rKa_Oa8sllulqtSA/viewform?vc=0&amp;c=0&amp;w=1&amp;flr=0">https://docs.google.com/forms/d/e/1FAIpQLSeM-1pJ13RyPVgF0bGRQtLiupwWDvALD6rKa_Oa8sllulqtSA/viewform?vc=0&amp;c=0&amp;w=1&amp;flr=0</a></p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p>







	<p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A (Note that actors in the dataset have NOT provided their consent for the usage of these videos)</p>
Other Issues	N/A

#### 5.4.20 DF-Platter dataset

DMP component	AI4Media_Data_139_WP6_VIDEO_DF_Platter_v1 Partner: CERTH
Data Summary	<p><u>Purpose:</u> The DF-Platter dataset consists of 133,260 videos sourced from YouTube and manipulated with 3 manipulation methods.</p> <p>We use this dataset for training and testing deepfake detection models in T6.2.</p> <p><u>Type/format:</u> Videos (mp4)</p> <p><u>Re-use of existing data:</u> Yes, we re-use published data.</p> <p><u>Data origin:</u> <a href="https://iab-rubic.org/df-platter-database#">https://iab-rubic.org/df-platter-database#</a></p> <p><u>Expected size:</u> ~500GB</p> <p><u>Data utility:</u> Used in T6.2 for training and testing video deepfake detection models.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes, the data can be found at <a href="https://iab-rubic.org/df-platter-database#">https://iab-rubic.org/df-platter-database#</a></p> <p><u>Search keywords:</u> DF-Platter dataset</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Yes, available at <a href="https://drive.google.com/drive/folders/1GeR-a2LfcMkcY6Qzpv2TP8utLtYFBmTs">https://drive.google.com/drive/folders/1GeR-a2LfcMkcY6Qzpv2TP8utLtYFBmTs</a></p> <p><u>How it will be accessible:</u> Download as a compressed file</p> <p><u>Methods/software tools to access data:</u> Browser</p> <p><u>Repository:</u> Google drive</p> <p><u>Restrictions on access:</u> The user should accept a licence agreement first: <a href="https://iab-">https://iab-</a></p>



	<a href="https://rubric.org/old/images/pdf/license/LICENSE_AGREE_DB_df-platter.pdf">rubric.org/old/images/pdf/license/LICENSE_AGREE_DB_df-platter.pdf</a>
Making data interoperable	<p><u>Interoperability</u>: Standard formats are considered in this dataset, the data are easy to download and use.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Licence prevents the user from re-sharing the dataset: <a href="https://rubric.org/old/images/pdf/license/LICENSE_AGREE_DB_df-platter.pdf">https://rubric.org/old/images/pdf/license/LICENSE_AGREE_DB_df-platter.pdf</a></p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Data already publicly shared.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset will be downloaded and stored on CERTH's servers. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate most of the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Licence prevents the user from re-sharing the dataset.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A (Note that actors in the dataset have NOT provided their consent for the usage of these videos)</p>
Other Issues	N/A

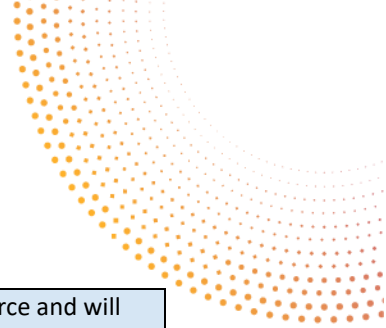
#### 5.4.21 MAFW dataset

DMP component	AI4Media_Data_140_WP6_VIDEO_MAFW_v1 Partner: QMUL
Data Summary	<p><u>Purpose</u>: MAFW is a large-scale, multi-modal, compound affective database for dynamic facial expression recognition in the wild. Clips in this database come from China, Japan, Korea, Europe, America and India, and cover various themes, e.g., variety, family, science fiction, suspense, love, comedy, and interviews, encompassing a wide range of human emotions. Each clip has been independently labeled 11 times by 11 well-trained annotators. MAFW database has enormous diversities, large quantities, and rich annotations, including: 10,045 number of video clips from movies, TV dramas, and short videos, a 11-dimensional expression distribution vector for each video clip, three kinds of annotations: (1) single expression label; (2) multiple expression label; (3) bilingual emotional descriptive text, two subsets: single-expression set, including 11 classes of single emotions; multiple-expression set, including 32 classes of multiple emotions, three automatic annotations: the frame-level 68 facial landmarks, bounding boxes of face regions, and gender, four</p>



	<p>benchmarks: uni-modal single expression classification, multi-modal single expression classification, uni-modal compound expression classification, and multi-modal compound expression classification. Used in T6.4 for facial expression recognition in video.</p> <p><u>Type/format</u>: mp4, png</p> <p><u>Re-use of existing data</u>: We are reusing an existing dataset</p> <p><u>Data origin</u>: <a href="https://mafw-database.github.io/MAFW/">https://mafw-database.github.io/MAFW/</a></p> <p><u>Expected size</u>: 35 GB</p> <p><u>Data utility</u>: It is useful in the context of T6.6 for evaluating video-based facial expression recognition models. In general, it can be used for evaluating models with spatial-temporal feature learning.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is discoverable in the MAFW website and GitHub: <a href="https://mafw-database.github.io/MAFW/">https://mafw-database.github.io/MAFW/</a></p> <p><u>Search keywords</u>: MAFW dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, the data are accessible for academic purposes upon contacting the team responsible for developing the dataset.</p> <p><u>How it will be accessible</u>: Users must contact the team responsible for developing the dataset after completing appropriate documentation including the terms and conditions for the dataset's use.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: The data are hosted at the MAFW website: <a href="https://mafw-database.github.io/MAFW/">https://mafw-database.github.io/MAFW/</a></p> <p><u>Restrictions on access</u>: The team responsible for the dataset must be contacted to allow access to the dataset.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: User must sign a license agreement: <a href="https://github.com/MAFW-database/MAFW/blob/main/academics/maf-w-academics-final.pdf">https://github.com/MAFW-database/MAFW/blob/main/academics/maf-w-academics-final.pdf</a></p> <p><u>Availability for re-use</u>: Yes</p> <p><u>Usable by third parties after end of project</u>: Yes, this is an open dataset</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>





Data security	<u>Security measures</u> : The dataset will be downloaded from the original source and will be stored on QMUL's servers. QMUL fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies (firewalls, right-based file system, etc.) mitigate most of the risk of illegitimate access.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : Links to the dataset should be shared instead of raw data.  <u>Is informed consent for data sharing and long term preservation given</u> : N/A
Other Issues	N/A

#### 5.4.22 Dynamic Facial Expression in-the-Wild (DFEW) video dataset

DMP component	AI4Media_Data_141_WP6_VIDEO_DFEW_v1 Partner: QMUL
Data Summary	<p><u>Purpose</u>: Dynamic Facial Expression in-the-Wild (DFEW) is a large-scale facial expression database with 16,372 very challenging video clips taken from movies. Clips in the DFEW database are of various challenging interferences, such as extreme illumination, occlusions, and capricious pose changes. Based on the crowdsourcing annotations, we hired 12 expert annotators, and each clip has been independently labeled ten times by them. DFEW database has enormous diversities, large quantities, and rich annotations, including: 16372 number of very challenging video clips from movies, a 7-dimensional expression distribution vector for each video clip, single-labeled expression annotation for classic seven discrete emotions, baseline classifier outputs based on single-labeled annotation.</p> <p><u>Type/format</u>: mp4, png</p> <p><u>Re-use of existing data</u>: We are reusing an existing dataset</p> <p><u>Data origin</u>: <a href="https://dfew-dataset.github.io/index.html">https://dfew-dataset.github.io/index.html</a></p> <p><u>Expected size</u>: 25.79 GB</p> <p><u>Data utility</u>: It is useful in the context of T6.6 for evaluating video-based facial expression recognition models. In general, it can be used for evaluating models with spatial-temporal feature learning.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data is discoverable in the DFEW website: <a href="https://dfew-dataset.github.io/index.html">https://dfew-dataset.github.io/index.html</a></p> <p><u>Search keywords</u>: DFEW dataset</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, the data are accessible for academic purposes upon contacting the team responsible for developing the dataset.</p> <p><u>How it will be accessible</u>: The data can be downloaded from the DFEW project website, after contacting the team responsible for the dataset and being given the password: <a href="https://dfew-dataset.github.io/download.html">https://dfew-dataset.github.io/download.html</a></p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: The data are hosted at the DFEW website: <a href="https://dfew-">https://dfew-</a></p>

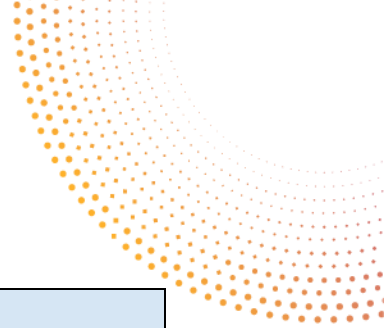


	<p><a href="https://dataset.github.io/download.html">dataset.github.io/download.html</a></p> <p><u>Restrictions on access:</u> The team responsible for the dataset must be contacted to allow access to the dataset.</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> Each video is individually annotated by ten annotators under professional guidance and assigned to one of the seven basic expressions (i.e., happiness, sadness, neutral, anger, surprise, disgust, and fear).</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> DFEW database is available for non-commercial research purposes only. User must sign a license agreement: <a href="https://dfew-dataset.github.io/download.html">https://dfew-dataset.github.io/download.html</a></p> <p><u>Availability for re-use:</u> Yes</p> <p><u>Usable by third parties after end of project:</u> Yes, this is an open dataset</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The dataset will be downloaded from the original source and will be stored on QMUL's servers. QMUL fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies (firewalls, right-based file system, etc.) mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Links to the dataset should be shared instead of raw data.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	N/A

#### 5.4.23 FERV39k video dataset

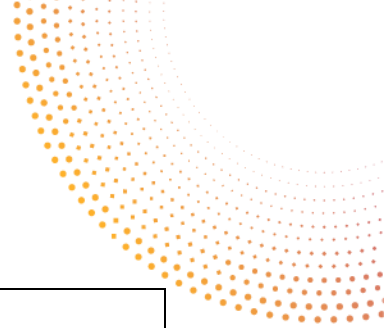
<b>DMP component</b>	<b>AI4Media_Data_142_WP6_VIDEO_FERV39k_v1</b> <b>Partner: QMUL</b>
Data Summary	<p><u>Purpose:</u> The FERV39k dataset is a large-scale multi-scene dataset designed for facial expression recognition in the wild. It contains 38,935 video clips labelled with 7 classic expressions across 22 fine-grained scenes in 4 isolated scenarios. The dataset uncovers several new challenges: 1) difficulty and confusion of 7 basic expression classes; 2) discrepancy across 4 scenarios; 3) unsatisfactory cross-scenario performance; 4) long-tail distribution of expressions and duration. The scenarios and scenes were designed following four reasons: 1) Plenty of video sources and samples; 2) Expandability of 22 fine-grained scenes; 3) Large variations and limited overlapping; 4) Distinct associations with scene context.</p> <p><u>Type/format:</u> mp4</p>





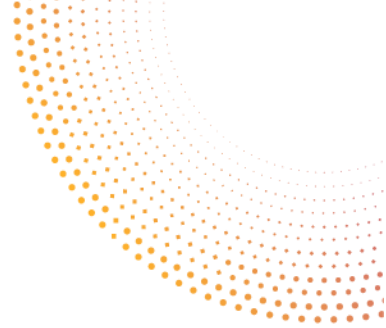
	<p><b>Re-use of existing data:</b> We are reusing an existing dataset</p> <p><b>Data origin:</b> <a href="https://wangyanckxx.github.io/Proj_CVPR2022_FERV39k.html">https://wangyanckxx.github.io/Proj_CVPR2022_FERV39k.html</a></p> <p><b>Expected size:</b> 13.45 GB</p> <p><b>Data utility:</b> It is useful in the context of T6.6 for evaluating video-based facial expression recognition models. In general, it can be used for evaluating models with spatial-temporal feature learning.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes, the data are discoverable in the FERV39k website: <a href="https://wangyanckxx.github.io/Proj_CVPR2022_FERV39k.html">https://wangyanckxx.github.io/Proj_CVPR2022_FERV39k.html</a></p> <p><b>Search keywords:</b> FERV39k dataset</p> <p><b>Versioning:</b> N/A</p> <p><b>Metadata creation:</b> N/A</p>
Making data openly accessible	<p><b>Data openly accessible:</b> No, the data are accessible for academic purposes upon contacting the team responsible for developing the dataset.</p> <p><b>How it will be accessible:</b> The data can be downloaded from Baidu or Google Drivers, after contacting the team responsible for the dataset and being given the password and link: <a href="https://github.com/wangyanckxx/FERV39k">https://github.com/wangyanckxx/FERV39k</a></p> <p><b>Methods/software tools to access data:</b> N/A</p> <p><b>Repository:</b> The data are hosted at the Baidu or Google Drivers.</p> <p><b>Restrictions on access:</b> The team responsible for the dataset must be contacted to allow access to the dataset: <a href="https://github.com/wangyanckxx/FERV39k">https://github.com/wangyanckxx/FERV39k</a></p>
Making data interoperable	<p><b>Interoperability:</b> N/A</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> The annotations of FERV39k are copyright and published under the Creative Commons Attribution-NonCommercial 4.0 International License.</p> <p><b>Availability for re-use:</b> Yes</p> <p><b>Usable by third parties after end of project:</b> This is an open dataset.</p> <p><b>Re-use timeframe:</b> N/A</p> <p><b>Data quality assurance process:</b> N/A</p>
Allocation of resources	<p><b>Costs for making data FAIR:</b> N/A</p> <p><b>Costs for long-term preservation:</b> N/A</p>
Data security	<p><b>Security measures:</b> The dataset will be downloaded from the original source and will be stored on QMUL's servers. QMUL fully complies with the applicable national and European data protection frameworks. State-of-the-art IT security measures and company policies (firewalls, right-based file system, etc.) mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><b>Possible ethical and legal aspects preventing sharing:</b> Links to the dataset should be</p>





	shared instead of raw data. <u>Is informed consent for data sharing and long term preservation given:</u> N/A
Other Issues	N/A





## 6. Data management plan for non-research datasets collected in AI4Media

This section presents non-research datasets collected by AI4Media partners to support the activities of WP1 (management), WP2 (AI Media Observatory), WP7 (contributions to AIoD), WP8 (use cases), WP9 (AIDA), WP10 (open calls) and WP11 (dissemination). **24 non-research datasets** have been identified in the final DMP (12 datasets had been identified in the previous version).

In the following, we present the DMP plan for each of these datasets using the template presented in Section 4 (see Table 1). The Table below briefly summarizes the 24 datasets presented in this section and offers a glance at the structure of the section and its subsections. New datasets that were not included in the initial DMP are indicated with yellow.

Table 4: Summary of non-research datasets collected in AI4Media

DMP component	WP	Short summary	Relevant sub-section
<b>Data collected in WP1 (Project Management)</b>			6.1
AI4Media_Data_143_WP1_Text_AI4MediaConsortiumContactInfo_v1	WP1	AI4Media consortium contact info dataset	6.1.1
<b>Data collected in WP2 (European AI vision, policy and common research)</b>			6.2
AI4Media_Data_144_WP2_TEXT-IMAGE_ObservatoryData_v1	WP2	AI Media Observatory Dataset	6.2.1
AI4Media_Data_145_WP2_Text_MediaResponseToDataCrawls_v1	WP2	Curated dataset of resources on how the media sector responds to content crawling for AI model training	6.2.2
<b>Data collected in WP7 (Integration with AI-on-Demand platform)</b>			6.3
AI4Media_Data_146_WP7_Text_AICafe_v1	WP7	AI-Cafe mailing list members and AI-Cafe participants	6.3.1
AI4Media_Data_147_WP7_Text_AI-Assets_v1	WP7	Candidate AI assets dataset	6.3.2
<b>Data collected in WP8 (Use cases &amp; demonstrators in media, society and politics)</b>			6.4
AI4Media_Data_148_WP8_UserData-TrulyMedia-UC1-ATC_v1	WP8	User data from Truly Media for Use case 1	6.4.1
<b>Data collected in WP9 (Doctoral Academy and exchange programme)</b>			6.5
AI4Media_Data_149_WP9_Text_AIDA CourseOfferings_v1	WP9	AIDA course offerings dataset	6.5.1
AI4Media_Data_150_WP9_Text_AIDA Students_v1	WP9	AIDA students dataset	6.5.2
AI4Media_Data_151_WP9_Text_AIDA MailingList_v1	WP9	AIDA mailing list dataset	6.5.3
AI4Media_Data_152_WP9_Text_AIDA AIEducationalResources_v1	WP9	AIDA AI educational resources dataset	6.5.4
AI4Media_Data_153_WP9_Text_AIDA Lecturers_v1	WP9	AIDA lecturers dataset	6.5.5
AI4Media_Data_154_WP9_Text_AIExcellenceLectures_v1	WP9	AIDA AI Excellence Lecture Series dataset	6.5.6
AI4Media_Data_155_WP9_Text_AIDA Curators_v1	WP9	AIDA curators dataset	6.5.7





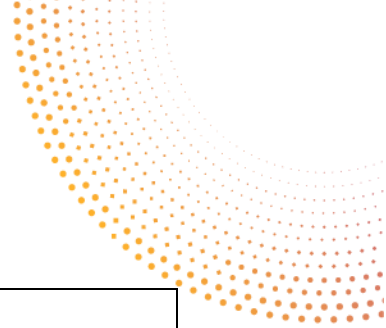
AI4Media_Data_156_WP9_TEXT_AIDAWebsiteGoogleAnalytics_v1	WP9	AIDA website analytics dataset	6.5.8
AI4Media_Data_157_WP9_Text_JuniorFellowsExchange_v1	WP9	AI4Media Junior Fellows exchange program dataset	6.5.9
<b>Data collected in WP10 (Community Outreach and Growth)</b>			6.6
AI4Media_Data_158_WP10_Competitive_call_application_datasets_v1	WP10	Competitive call application datasets	6.6.1
AI4Media_Data_159_WP10_sub-granted_projects_dataset_v1	WP10	Sub-granted projects datasets	6.6.2
AI4Media_Data_160_WP10_sub-granted_projects_dataset_v1	WP10	External experts and evaluators datasets	6.6.3
AI4Media_Data_161_WP10_ParticipantsCompetitiveCallsEvents_v1	WP10, WP11	Participants of competitive call related events datasets	6.6.4
<b>Data collected in WP11 (Communication, dissemination, exploitation and sustainability)</b>			6.7
AI4Media_Data_162_WP11_Text_AI4MediaAssociateMembersContactInfo_v1	WP11	AI4Media associate members contact info dataset	6.7.1
AI4Media_Data_163_WP11_Text_NewsletterSubscribers_v1	WP11	AI4Media newsletter subscribers dataset	6.7.2
AI4Media_Data_164_WP11_Text_WebsiteMessages_v1	WP11	AI4Media website messages dataset	6.7.3
AI4Media_Data_165_WP11_TEXT_AI4MediaWebsiteGoogleAnalytics_v1	WP9	AI4Media website analytics dataset	6.7.4
AI4Media_Data_166_WP11_Text_RegistrantsstoEvents_v1	WP11	Dataset of registrants for AI4Media events	6.7.5

## 6.1 Datasets collected in the context of WP1

### 6.1.1 AI4Media consortium contact info dataset

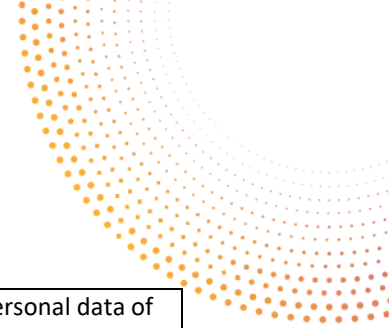
DMP component	AI4Media_Data_143_WP1_Text_AI4MediaConsortiumContactInfo_v1 Partner: CERTH
Data Summary	<p><b>Purpose:</b> This dataset contains business contact information from AI4Media consortium members, including names, affiliation, emails, office phone numbers, Skype Ids, office postal addresses, wiki user names, etc. The collected data are necessary for the communication among project partners and are collected in the context of WP1-Management.</p> <p><b>Type/format:</b> text</p> <p><b>Re-use of existing data:</b> No.</p> <p><b>Data origin:</b> Data provided by AI4Media partners to CERTH in emails or excel files.</p> <p><b>Expected size:</b> &lt; 1MB</p> <p><b>Data utility:</b> This data is necessary for facilitating communication among consortium members.</p>
Making data findable, incl. provisions for	<p><b>Is data discoverable:</b> No, data are not discoverable from outside. The data is stored on the project wiki and is only discoverable by consortium members with a wiki account.</p>





metadata	<p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: N/A.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No. The data will only be shared internally in AI4Media since it contains personal information.</p> <p><u>How it will be accessible</u>: Restricted access. The data is accessible only by wiki users. Information about wiki usernames and passwords can only be accessed by the wiki administrator.</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Data will not be shared.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on the project wiki, which is on a dedicated web server hosted in CERTH's premises. The wiki web site uses for its domain an SSL certificate enabling the SHA256RSA signature algorithm and forces all visits to use HTTPS to ensure the traffic is secure. The wiki is restricted only to registered users while registration is possible only by invitation. Access requires username/password authentication. CERTH fully complies with the applicable national, European and International framework, and the GDPR. The wiki uses a file-based RDBMS to enhance security. Web server and file-based DB are running on a Linux encrypted partition, which conforms to the data-at-rest GDPR guidelines. Moreover, state-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimize the risk of data loss. The data will be preserved there for three years after the end of the project and will then be deleted.</p>





Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Dataset includes personal data of consortium members and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> N/A</p>
Other Issues	No

## 6.2 Datasets collected in the context of WP2

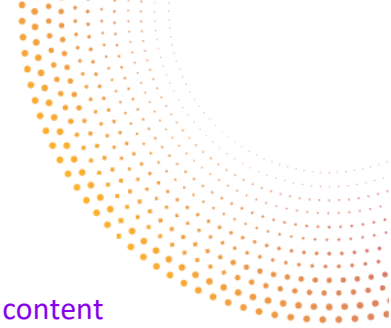
### 6.2.1 AI Media Observatory Dataset

DMP component	AI4Media_Data_144_WP2_TEXT-IMAGE_ObservatoryExpertsRepo_v1 Partner: LOBA
Data Summary	<p><u>Purpose:</u> The data was collected in order to build an <a href="#">AI Media expert repository</a> for the <a href="#">AI Media Observatory</a>. The AI Media Expert is an expert directory where visitors of the Observatory can easily search, find and contact a relevant expert within the field of media and AI - whether it be a technical, legal or social expert. The directory aims to help civil society and media professionals who are seeking knowledge on a specific topic on AI in media.</p> <p><u>Type/format:</u> The data is in textual format, with some elements such as photos being in image format (e.g., JPEG, PNG).</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> The data originates from applications submitted by experts via a <a href="#">form</a> on the Media AI Observatory website.</p> <p><u>Expected size:</u> The expected size per expert application is relatively small, with the resume limited to 50 words max and the description limited to 150 characters. Depending on the number of expert submissions, the total size is estimated to be a few MBs.</p> <p><u>Data utility:</u> The visible results of the datasets (meaning the webpage where the experts are featured) aims to help civil society and media professionals who are seeking knowledge on a specific topic on AI in media. However, the dataset is held internally and is not shared, but the information is freely accessible on the page: <a href="https://www.ai4media.eu/find-an-ai-media-expert/">https://www.ai4media.eu/find-an-ai-media-expert/</a></p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data is featured in a publicly accessible webpage on the AI Media Observatory. There are no barrier to access.</p> <p><u>Search keywords:</u> You can find the Expert Repository page on the Observatory by typing in a search engine: AI Media Expert AI4Media.</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> Metadata includes the expert's name, position, and keywords related to their expertise. This metadata is created and managed within the WordPress system used by the website.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Not applicable, the data are visible online but the dataset will not be shared for GDPR reasons.</p> <p><u>How it will be accessible:</u> Online on the AI Media Observatory:</p>



	<p><a href="https://www.ai4media.eu/find-an-ai-media-expert/">https://www.ai4media.eu/find-an-ai-media-expert/</a></p> <p><u>Methods/software tools to access data:</u> The data can be accessed using standard web browsers to access to the AI Media Observatory. No special software tools are needed beyond those used for website access.</p> <p><u>Repository:</u> The data is stored in the WordPress system of the AI Media Observatory website.</p> <p><u>Restrictions on access:</u> Access will be granted to data subjects who make a specific request to access their personal data. Otherwise, request will not be granted, the data are already visible online in the Observatory.</p>
Making data interoperable	<p><u>Interoperability:</u> Anyone interested can have a look at the webpage and experts featured in the repository.</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> NA. The data cannot be re-used otherwise another consent is needed from the individuals featured in the expert repository.</p> <p><u>Usable by third parties after end of project:</u> The data will remain visible on the observatory after the end of the project. Unless expert no longer want to be part of the repository, their data will be removed.</p> <p><u>Re-use timeframe:</u> The data will remain available during the period the website is maintained, for up to 5 years after the completion of the project.</p> <p><u>Data quality assurance process:</u> Data are updated as soon as a new form is received or when someone wishes to update or delete their information. They can reach out the observatory through the online contact form or through our email address: <a href="mailto:observatory@ai4media.eu">observatory@ai4media.eu</a>.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> The costs for long-term preservation include maintaining the website and server hosting.</p>
Data security	<p><u>Security measures:</u> Security measures include controlled access to the backend of the website, user authentication, firewalls, encryption, and regular back-ups.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The data contain personal data.</p> <p><u>Is informed consent for data sharing and long-term preservation given:</u> The data will not be shared. Participants are free to update or ask the deletion of their information at any time without any justification. The purpose for the Observatory is to be a long lasting AI4Media milestone. However, if the Observatory is no longer sustained, the website will no longer be supported after (31 August 2029) and the data will be deleted.</p>
Other Issues	N/A

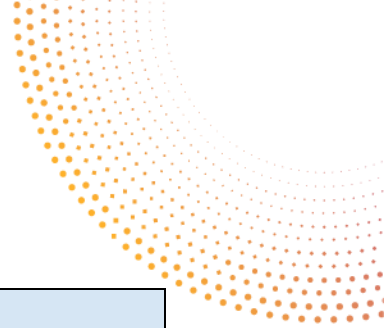




## 6.2.2 Curated dataset of resources on how the media sector responds to content crawling for AI model training

DMP component	AI4Media_Data_145_WP2_Text_MediaResponsetoDataCrawls_v1 Partner: NISV, KUL
Data Summary	<p><u>Purpose</u>: A curated repository of third-party resources (e.g. news articles, blogs) that monitors how media organisations are responding to content scraping for AI model training.</p> <p><u>Type/format</u>: An online spreadsheet with a description of third-party links, organised into topical categories (anti-scraping statement, legal action, licensing deals, new techniques &amp; methods, opinion pieces and analysis - general analysis of the topic).</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: The majority of resources were curated by AI4Media partners. Additionally, an open invitation via an online survey allows crowdsourced contributions from anyone interested in the topic.</p> <p><u>Expected size</u>: 40-80 entries.</p> <p><u>Data utility</u>: A curated repository of third-party resources (e.g. news articles, blogs) that monitors how media organisations are responding to content scraping for AI model training – will be used for the AI Media Observatory.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No</p> <p><u>Search keywords</u>: No</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: Yes</p> <p><u>How it will be accessible</u>: Publicly accessible link to the Google spreadsheet with the data: <a href="https://docs.google.com/spreadsheets/d/1UZF-H-SfzpeaZ7OjEmbJFippDESyRilfX3I5I_US7uE/edit">https://docs.google.com/spreadsheets/d/1UZF-H-SfzpeaZ7OjEmbJFippDESyRilfX3I5I_US7uE/edit</a></p> <p><u>Methods/software tools to access data</u>: Through a web browser</p> <p><u>Repository</u>: Google drive</p> <p><u>Restrictions on access</u>: No</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: Anyone interested in this dataset is free to download it, perform analysis based on it or combine it with other datasets.</p> <p><u>Usable by third parties after end of project</u>: NISV and KUL are committed to maintain the repository online for as long as deemed relevant. If at a certain point after the end of the project a decision will be made to stop maintaining the repository, a final</p>





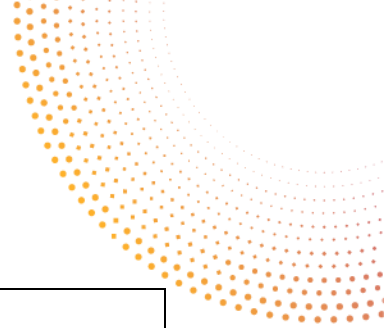
	<p>version of it will be uploaded to the AI4Media Zenodo account.</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Publicly crowdsourced submissions are reviewed by NISV and KUL teams.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<u>Security measures</u> : N/A
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

## 6.3 Datasets collected in the context of WP7

### 6.3.1 AI-Cafe mailing list members and AI-Cafe participants

DMP component	AI4Media_Data_146_WP7_Text_AI-Cafe_v1 Partner: GAR
Data Summary	<p><u>Purpose</u>: GAR organises and moderates regular AI-Cafe sessions since 2019 in the AI4EU Cafe and then continued the AI-Cafes for the AI4Media Project. GAR uses the tool Gotowebinar for the Cafes. In 2021, all former AI4EU Cafe participants were asked if they want to receive future invitations from AI4Media for the upcoming AI-Cafes. Only 10% answered in written form, they would like to receive further email invitations. These participants get regular invitations by email to the Cafe and they are part of a confidential email list.</p> <p><u>Type/format</u>: Text</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: The email list includes emails from participants who registered to the AI-Café.</p> <p><u>Expected size</u>: ~400 emails</p> <p><u>Data utility</u>: It is useful for the AI- Cafe organisation in the AI4Media project.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No</p> <p><u>Search keywords</u>: No</p> <p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No. The mailing list contains personal information.</p> <p><u>How it will be accessible</u>: N/A</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>



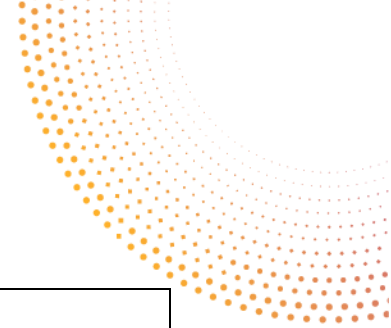


Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: We use the the following metadata fields: First name, Last name, Email</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: In the case the AI-Cafe continues after the AI4Media project ends, then again each list participant will receive an email, if they want to receive further invitations to future Cafes.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The list is stored in GAR internal servers and only selected GAR employees have access to it.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: It contains personal information.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	N/A

### 6.3.2 Candidate AI assets dataset

DMP component	AI4Media_Data_147_WP7_Text_AI-Assets_v1 Partner: FHG-IAIS
Data Summary	<p><u>Purpose</u>: This dataset is a list of project results from AI4Media consortium members, including contact information (names, affiliation, emails). The collected data is required to contact project partners to publish their project results on the AI-on-Demand platform.</p> <p><u>Type/format</u>: Excel</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: Data provided by AI4Media partners in work package reports, deliverables and on the Resources Library on the AI4Media website.</p> <p><u>Expected size</u>: &lt; 1MB</p> <p><u>Data utility</u>: This data is necessary to keep track of the publication of candidate AI Assets on the AioD platform.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, data is not discoverable from outside. The data is stored on the FhG internal storage and is only discoverable for selected employees working in the project.</p> <p><u>Search keywords</u>: N/A</p>





	<p><u>Versioning</u>: No</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No.</p> <p><u>How it will be accessible</u>: Restricted access. The data is accessible only by selected employees.</p> <p><u>Methods/software tools to access data</u>: Microsoft Excel.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: Data will not be shared.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on the FhG internal storage and is only discoverable for selected employees in the project. Access to the storage is protected by FhG user access management. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Dataset includes personal data of consortium members and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

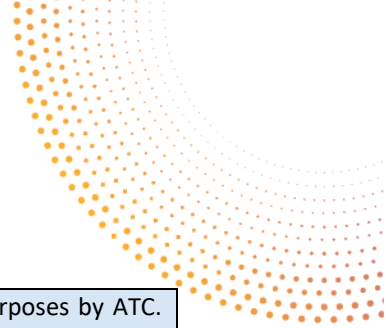
## 6.4 Datasets collected in the context of WP8

### 6.4.1 User data from Truly Media for Use case 1

<b>DMP component</b>	<b>AI4Media_Data_148_WP8_UserData-TrulyMedia-UC1-ATC_v1</b> <b>Partner: ATC</b>
Data Summary	<p><u>Purpose</u>: In order to realise Use Case 1 in WP8, <a href="#">Truly Media</a>, a web-based platform for collaborative verification co-owned by ATC and DW, will be used as the basis of the main demonstrator in UC1. In order for the test users, as well as for the project partners developing AI tools and components for UC1, to access Truly Media, ATC will collect and</p>

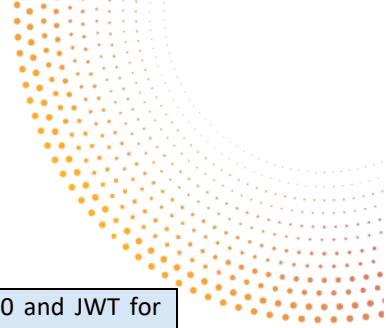






	<p>store personal data. The data will be used for registration and login purposes by ATC. ATC uses Twitter login that enables users to register and login to the platform with their Twitter credentials. Collected data include: Name; Email; Organization; Department; Role; Expertise; Office phone number; Twitter profile image; and Twitter handle.</p> <p><u>Type/format</u>: Data stored in JSON format.</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Truly Media users registering and logging in the platform.</p> <p><u>Expected size</u>: A few MBs.</p> <p><u>Data utility</u>: This data will be used in the context of WP8 to allow test users and project partners developing AI tools and components for UC1 to access Truly Media.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, because the dataset described will be stored on Truly Media's databases and there are no plans for data sharing.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The data will not be openly accessible since it contains personal information.</p> <p><u>How it will be accessible</u>: It is planned that the data will only be accessible by project partner ATC.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The data will not be licensed since it will only be used internally.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The dataset described will be stored by ATC on third-party cloud servers. Appropriate and detailed security policies, rules, and technical measures are implemented to protect data are used by the Truly Media platform and stored on the platform from improper or unauthorized access, including use of firewalls where appropriate. Security measures also include 2FA (2 Factor Authentication) with OTP</p>





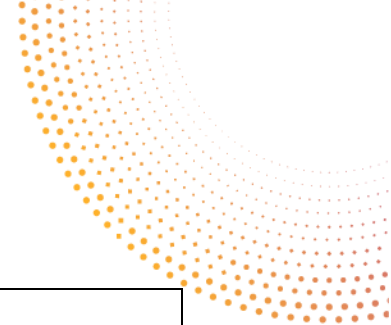
	(One Time Password) for extra security during login, as well as Auth2.0 and JWT for authentication and authorisation. End-to-end encryption protects from man-in-the-middle attacks and data theft. All ATC employees and data processors, who have access to and are associated with the processing of personal data, are obliged to respect the confidentiality of the stored personal data. Moreover, ATC's development team has received training from external auditors for security awareness and security best practices to avoid vulnerabilities in source code. External auditors have performed black-box penetration testing to ensure that the platform is fully secure. ATC's Data Protection Officer ensures that all processes followed are fully compliant with the GDPR provisions.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The data will not be shared since it contains personal information of end users.</p> <p><u>Is informed consent for data sharing and long-term preservation given:</u> Informed consent is given implicitly by users when completing the relevant information on the registration form and when authorising Truly Media to use their Twitter account for login purposes.</p>
Other Issues	No

## 6.5 Datasets collected in the context of WP9

### 6.5.1 AIDA course offerings dataset

DMP component	AI4Media_Data_149_WP9_Text_AIDACourseOfferings_v1 Partner: AUTH
Data Summary	<p><u>Purpose:</u> This dataset contains information about course offerings from collaborating AIDA full and international members, for the purpose of advertising them to AIDA students. Such information includes non-personal data (course title, date, course offer affiliation), minimal personal data about the lecturers (Full Name, lecturer affiliation) and also electronic business contact information (e-mail). The collected data are necessary for the execution of WP9 Tasks T9.1 and T9.2.</p> <p><u>Type/format:</u> text</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Lecturers will the details in web applications or provide them in shared excel files and are responsible for maintenance updates. All data is moderated by AUTH to resolve inconsistencies. At any time, lectures may alter or remove their personal data and course offerings that appear on i-aida.org website, using the web applications or contacting the website moderators (AUTH). Lecturers will maintain the responsibility/option of adding, editing, and deleting course offerings, having full access to their own content at all times.</p> <p><u>Expected size:</u> Few KBs</p> <p><u>Data utility:</u> The data is useful to AIDA registered students to apply for AIDA course offerings.</p>
Making data findable, incl.	<u>Is data discoverable:</u> Yes, the data will be discoverable in the web, though search engines.





provisions for metadata	<p><u>Search keywords</u>: Through web metadata.</p> <p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: Course title, offer type (web/short/semester), course affiliation and lecturer full names and affiliation.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The dataset is publicly available on <a href="http://www.i-aida.org">www.i-aida.org</a> website.</p> <p><u>How it will be accessible</u>: Though web.</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/AT</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality assurance is ensured through moderation from AUTH. AUTH will consistently check if the appearing information is in the correct form, and will request the respective lectures for updates, if necessary.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored at an internal AUTH server in an encrypted format.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: Yes. A form will be requested to be signed for every lecturer involved in AIDA program.</p>
Other Issues	No

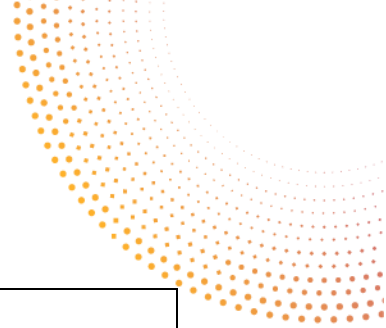
### 6.5.2 AIDA students dataset

<b>DMP component</b>	<b>AI4Media_Data_150_WP9_Text_AIDAStudents_v1</b> <b>Partner: AUTH</b>
Data Summary	<u>Purpose</u> : This dataset contains minimal personal identity information about AIDA PhD students (full name, e-mail, affiliation, supervisor id, gender) and details about their



	<p>progress in the form of courses attended, grades, ECTS collected within the AIDA program. The collected data are necessary for the execution of WP9 Tasks T9.1 and T9.2, for book-keeping and administrative purposes (e.g., provide AIDA certificates).</p> <p><u>Type/format</u>: text</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Students enter their personal details through a secure registration login process. Their supervisors are notified and validate this information. Details about the procedure are provided future i-aida.org website updates.</p> <p><u>Expected size</u>: Few KBs</p> <p><u>Data utility</u>: The data is useful for offering AIDA services to students (course attendance certificates etc.)</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: No.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, because it includes personal information of students.</p> <p><u>How it will be accessible</u>: It will be restricted</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: Personal registered student data will only be accessible by the students themselves, registered lecturers that offer the respective courses that the students attended, and AIDA website moderators (in encrypted/pseudo-anonymized form).</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality assurance will be ensured through moderation from AUTH.</p>
Allocation of	<p><u>Costs for making data FAIR</u>: N/A</p>



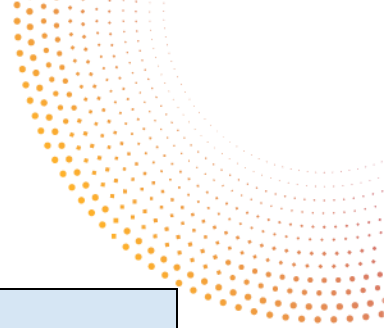


resources	<u>Costs for long-term preservation:</u> N/A
Data security	<u>Security measures:</u> The data is stored at an internal AUTH server in an encrypted form. Proper data security measures are implemented.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing:</u> The data will not be shared since it contains some personal information about AIDA students.  <u>Is informed consent for data sharing and long-term preservation given:</u> Yes. A form will be requested to be signed for every student involved in the AIDA program.
Other Issues	No

### 6.5.3 AIDA mailing list dataset

DMP component	AI4Media_Data_151_WP9_Text_AIDAMailingList_v1 Partner: AUTH
Data Summary	<u>Purpose:</u> This dataset contains business contact information from students and lecturers and interested researches in AIDA activities from the general AI community. The collected data are necessary for the execution of WP9 Task T9.3.  <u>Type/format:</u> text  <u>Re-use of existing data:</u> No.  <u>Data origin:</u> Mailing list participants will register by following the following instructions: <a href="https://lists.auth.gr/sympa/info/aida">https://lists.auth.gr/sympa/info/aida</a>  <u>Expected size:</u> Few KBs  <u>Data utility:</u> This mailing list will advertise AIDA activities to the general public.
Making data findable, incl. provisions for metadata	<u>Is data discoverable:</u> No personal data is retrievable except by the list moderators.  <u>Search keywords:</u> N/A  <u>Versioning:</u> Yes, though scheduled website backups.  <u>Metadata creation:</u> N/A.
Making data openly accessible	<u>Data openly accessible:</u> No  <u>How it will be accessible:</u> It is restricted.  <u>Methods/software tools to access data:</u> Web browser.  <u>Repository:</u> N/A  <u>Restrictions on access:</u> N/A
Making data interoperable	<u>Interoperability:</u> N/A  <u>Data and metadata vocabularies:</u> N/A  <u>Use of standard vocabularies:</u> N/A  <u>Mappings to commonly used vocabularies:</u> N/A



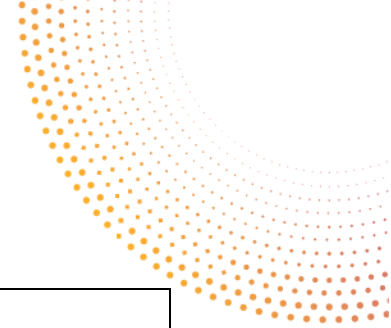


Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u> N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Personal data are stored in secure AUTH servers according to internal institutional procedures. Proper data security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Dataset includes personal data (email addresses) and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

#### 6.5.4 AIDA AI educational resources dataset

DMP component	AI4Media_Data_152_WP9_Text_AIDAAIEducationalResources_v1 Partner: AUTH
Data Summary	<p><u>Purpose</u>: This dataset contains information about resources provided by collaborating AIDA full and international members, for the purpose of advertising them to all AI enthusiasts. Such information includes non-personal data (resource title, resource description), and minimal personal data about the contributors (Full Name, contributor affiliation). The collected data are necessary for the execution of WP9 Tasks T9.1 and T9.2.</p> <p><u>Type/format</u>: text</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Lecturers enter the details in the AIDA website and are responsible for maintenance updates. AUTH moderates all data to resolve inconsistencies. At any time, lectures may alter or remove their personal data and resources that appear on i-aida.org website, using the web applications or contacting the website moderators (AUTH). Lecturers will maintain the responsibility/option of adding, editing, and deleting resources, having full access to their own content at all times.</p> <p><u>Expected size</u>: Few KBs</p> <p><u>Data utility</u>: The data will be useful to everyone to discover high-quality AI educational materials.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: Yes, the data will be discoverable on the web, through search engines.</p> <p><u>Search keywords</u>: Through web metadata.</p>



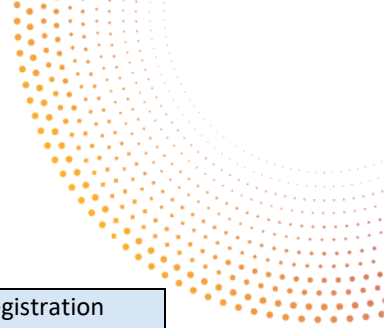


	<p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: Resource title, Resource author, Resource Description and Resource type</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The dataset is publicly available on <a href="http://www.i-aida.org">www.i-aida.org</a> website.</p> <p><u>How it will be accessible</u>: Though web.</p> <p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality assurance will be ensured though moderation from AUTH. AUTH will consistently check if the appearing information is in the correct form, and will request the respective lectures for updates, if necessary.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored at an internal AUTH server in an encrypted format. Proper data security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: Yes. A form will be requested to be signed for every lecturer involved in the AIDA program.</p>
Other Issues	N/A

### 6.5.5 AIDA lecturers dataset

<b>DMP component</b>	<b>AI4Media_Data_153_WP9_Text_AIDAlecturers_v1</b> <b>Partner: AUTH</b>
Data Summary	<p><u>Purpose</u>: This dataset will contain minimal personal identity information about AIDA Lecturers (full name, e-mail, affiliation). The collected data are necessary for the execution of WP9 Tasks T9.1 and T9.2, for book-keeping and administrative purposes.</p> <p><u>Type/format</u>: text</p> <p><u>Re-use of existing data</u>: No</p>

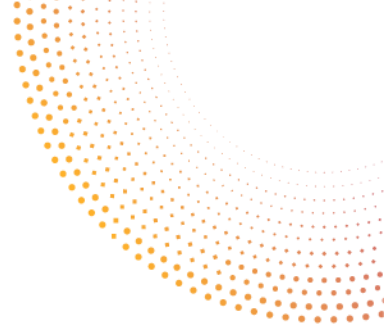




	<p><u>Data origin</u>: Lecturers will enter their personal details through a secure registration login process. Details about the procedure are analytically provided on the i-aida.org website.</p> <p><u>Expected size</u>: Few KBs</p> <p><u>Data utility</u>: The data will be useful for the AIDA website administrators (for book-keeping and administrative purposes)</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, except by the AIDA website administrators.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: No</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, because it includes the personal information of lecturers.</p> <p><u>How it will be accessible</u>: It will be restricted</p> <p><u>Methods/software tools to access data</u>: Web browser</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: Personal registered lecturer data will only be accessible by the lecturers themselves and website moderators (in encrypted/pseudo-anonymized form).</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality assurance will be ensured through moderation from AUTH.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored at an internal AUTH server in an encrypted form. Proper data security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: The data will not be shared since it contains some personal information about AIDA lectures.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: Yes. A form will be requested to be signed for every lecturer involved in the AIDA program.</p>
Other Issues	No



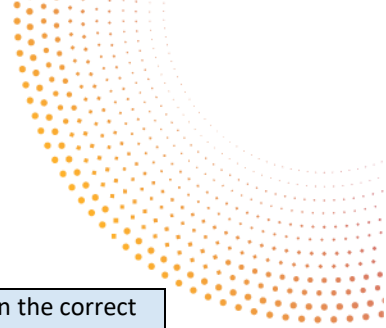




### 6.5.6 AIDA AI Excellence Lecture Series dataset

DMP component	AI4Media_Data_154_WP9_Text_AIExcellenceLectures_v1 Partner: AUTH
Data Summary	<p><b>Purpose:</b> This dataset contains information about AIDA’s AI Excellence Lecture Series provided by collaborating AIDA full and international members, for the purpose of advertising them to AI enthusiasts. Such information includes non-personal data (AI Excellence Lecture title, AI Excellence Lecture Abstract), and minimal personal data about the contributors (lecturer’s full name and short bio). The collected data are necessary for the execution of WP9 Task T9.3.</p> <p><b>Type/format:</b> text</p> <p><b>Re-use of existing data:</b> No</p> <p><b>Data origin:</b> Lecturers enter the details in the AIDA web application and are responsible for maintenance updates. All data is moderated by AUTH to resolve inconsistencies. At any time, lecturers may alter or remove their personal data and resources that appear on i-aida.org website, using the web applications or contacting the website moderators (AUTH).</p> <p><b>Expected size:</b> Few KBs</p> <p><b>Data utility:</b> The data will be useful to everyone to discover high-quality lecture materials.</p>
Making data findable, incl. provisions for metadata	<p><b>Is data discoverable:</b> Yes, the data will be discoverable in the web, though search engines.</p> <p><b>Search keywords:</b> Through web metadata.</p> <p><b>Versioning:</b> Yes, through scheduled website backups.</p> <p><b>Metadata creation:</b> Excellence Lecture title, Excellence Lecture Abstract, and lecturer full names and short bio.</p>
Making data openly accessible	<p><b>Data openly accessible:</b> The dataset is publicly available on <a href="http://www.i-aida.org">www.i-aida.org</a> website.</p> <p><b>How it will be accessible:</b> Though web.</p> <p><b>Methods/software tools to access data:</b> Web browser.</p> <p><b>Repository:</b> N/A</p> <p><b>Restrictions on access:</b> N/A</p>
Making data interoperable	<p><b>Interoperability:</b> N/A</p> <p><b>Data and metadata vocabularies:</b> N/A</p> <p><b>Use of standard vocabularies:</b> N/A</p> <p><b>Mappings to commonly used vocabularies:</b> N/A</p>
Increase data re-use	<p><b>Licence:</b> N/A</p> <p><b>Availability for re-use:</b> N/A</p> <p><b>Usable by third parties after end of project:</b> N/A</p> <p><b>Re-use timeframe:</b> N/A</p> <p><b>Data quality assurance process:</b> Data quality assurance is ensured though moderation</p>



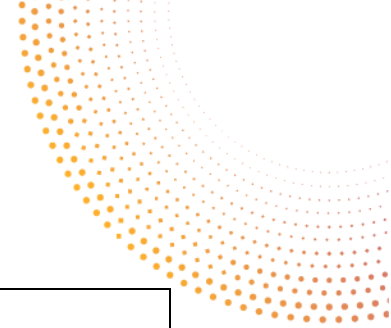


	from AUTH. AUTH will consistently check if the appearing information is in the correct form, and will request the respective lectures for updates, if necessary.
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : The data will be stored at an internal AUTH server in an encrypted format. Proper data security measures are implemented.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : N/A <u>Is informed consent for data sharing and long term preservation given</u> : Yes. A form will be requested to be signed for every lecturer involved in AIDA program.
Other Issues	No

### 6.5.7 AIDA curators dataset

DMP component	AI4Media_Data_155_WP9_Text_AIDACurators_v1 Partner: AUTH
Data Summary	<u>Purpose</u> : This dataset contains minimal personal identity information about AIDA Curators (full name, e-mail, affiliation, specific expertise, AI Educational Taxonomy). The collected data are necessary for the execution of WP9 Tasks T9.1 and T9.2, for book-keeping and administrative purposes. <u>Type/format</u> : text <u>Re-use of existing data</u> : No <u>Data origin</u> : Curators provide their personal details to system administrators. Details about the procedure are provided in the i-aida.org website. <u>Expected size</u> : Few KBs <u>Data utility</u> : The data will be useful for offering and managing AIDA Resources.
Making data findable, incl. provisions for metadata	<u>Is data discoverable</u> : No <u>Search keywords</u> : N/A <u>Versioning</u> : Yes, through scheduled website backups. <u>Metadata creation</u> : No
Making data openly accessible	<u>Data openly accessible</u> : No, because it includes the personal information of AIDA curators. <u>How it will be accessible</u> : It will be restricted <u>Methods/software tools to access data</u> : Web browser <u>Repository</u> : N/A <u>Restrictions on access</u> : Personal registered curator data will only be accessible by the curators themselves and AIDA website moderators (in encrypted/pseudo-anonymized form).
Making data interoperable	<u>Interoperability</u> : N/A <u>Data and metadata vocabularies</u> : N/A



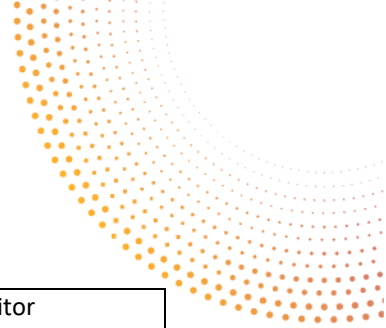


	<p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Data quality assurance is ensured through moderation from AUTH.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> The data will be stored at an internal AUTH server in an encrypted format. Proper data security measures are implemented.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The data will not be shared since it contains some personal information about AIDA curators.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> Yes. A form will be requested to be signed for every curator involved in AIDA program.</p>
Other Issues	No

### 6.5.8 AIDA website analytics dataset

DMP component	AI4Media_Data_156_WP9_TEXT_AIDAWebsiteGoogleAnalytics_v1 Partner: AUTH
Data Summary	<p><u>Purpose:</u> The dataset comprises data and statistics collected from Google Analytics related to the AIDA website. This includes metrics such as visitor numbers, page views, session duration, bounce rates, traffic sources, and user demographics. The purpose of collecting this data is to analyse the performance and reach of the AIDA website to understand user engagement, and to inform strategies for improving communication and dissemination activities. In addition to Google Analytics, additional statistics related to students and course attendance is tracked through WordPress.</p> <p><u>Type/format:</u> The data is in numerical and categorical format, typically exported as CSV or JSON files.</p> <p><u>Re-use of existing data:</u> No</p> <p><u>Data origin:</u> The data originates from Google Analytics tracking implemented on the AIDA website and WordPress feature.</p> <p><u>Expected size:</u> The expected size of the dataset is relatively small, estimated to be a few megabytes, depending on the volume of web traffic and the duration of data collection.</p> <p><u>Data utility:</u> This data is useful for monitoring and improving the website's performance and assess the impact of the content on the community.</p>
Making data findable, incl. provisions for	<p><u>Is data discoverable:</u> The data is identifiable and locatable through standard identification mechanisms.</p>



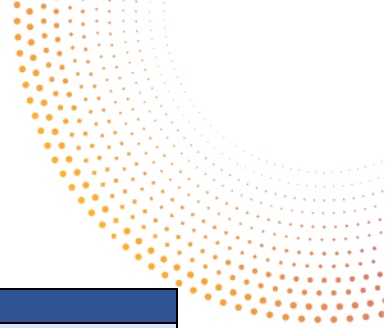


metadata	<p><u>Search keywords</u>: Search keywords such as "AIDA website analytics," "visitor statistics," and "web traffic data".</p> <p><u>Versioning</u>: N/A</p> <p><u>Metadata creation</u>: Metadata includes details such as the date range of data collection, types of metrics collected, and any filtering or processing applied. This is created following standard web analytics metadata practices.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No, the Google Analytics data will not be made openly available due to privacy and security concerns.</p> <p><u>How it will be accessible</u>: The data is mainly accessed by AUTH for monitoring purposes.</p> <p><u>Methods/software tools to access data</u>: The data can be accessed using Google Analytics and WordPress dashboards and reports, which are available through standard web browsers.</p> <p><u>Repository</u>: The data will be stored within the Google Analytics platform and WordPress and accessible through secure login.</p> <p><u>Restrictions on access</u>: Access will be restricted to authorized project members to ensure data privacy and security.</p>
Making data interoperable	<p><u>Interoperability</u>: The data is interoperable within the Google Analytics platform, allowing data exchange and re-use within the project team.</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A.</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: Data quality is assured through regular checks and validation processes, ensuring accuracy and reliability.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: Long-term preservation costs include maintaining access to the Google Analytics account.</p>
Data security	<p><u>Security measures</u>: Security measures include controlled access to the Google Analytics account, user authentication, firewalls, encryption, and regular backups.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: N/A</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: N/A</p>
Other Issues	N/A

### 6.5.9 AI4Media Junior Fellows exchange program dataset

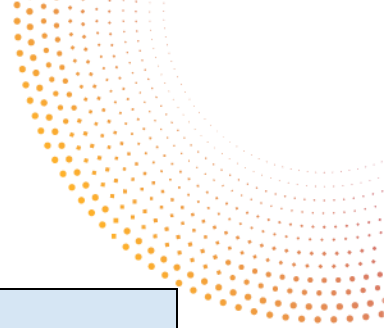
DMP	AI4Media_Data_157_WP9_Text_JuniorFellowsExchange_v1
-----	---





component	Partner: CERTH
Data Summary	<p><u>Purpose:</u> This dataset contains information of host and sender organizations involved in the AI4Media exchange program for Junior Fellows. This information includes: organization, organization type, country, exchange topics and interests, type of mobility, availability period, contact person name and email, logo/photo, short bios of researchers, organization profiles, etc. The collected data is necessary for the implementation of the Junior Fellows exchange program of AI4Media and it is collected in the context of WP9.</p> <p><u>Type/format:</u> Text</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Data provided by Host and Sender institutions (AI4Media partners but also other organizations involved in AI research, e.g. partners in other ICT-48 projects) when submitting their profiles through an online form in the project website in order to participate in AI4Media’s Junior Fellows exchange program.</p> <p><u>Expected size:</u> &lt;10 MBs</p> <p><u>Data utility:</u> This data is necessary for the operation of AI4Media’s Junior Fellows exchange program.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> Yes, data is discoverable from outside. The data will be displayed in a dedicated page in the project website.</p> <p><u>Search keywords:</u> Visitors of the website will be able to search for preferred profiles using keywords and filters.</p> <p><u>Versioning:</u> Yes, through scheduled website backups.</p> <p><u>Metadata creation:</u> N/A</p>
Making data openly accessible	<p><u>Data openly accessible:</u> Yes, the data will be openly accessible via the project website.</p> <p><u>How it will be accessible:</u> Displayed in a dedicated page in the project website.</p> <p><u>Methods/software tools to access data:</u> Web browser.</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>
Making data interoperable	<p><u>Interoperability:</u> N/A</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p>





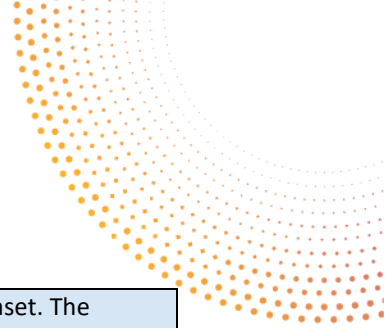
	<u>Data quality assurance process</u> : N/A
Allocation of resources	<u>Costs for making data FAIR</u> : N/A <u>Costs for long-term preservation</u> : N/A
Data security	<u>Security measures</u> : The data is stored on CERTH servers. Access requires username/password authentication. CERTH fully complies with the applicable national and European data protection frameworks & guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : Dataset includes personal data (name, email, short bio) from Host and Sender institutions that create their profiles to be publicly displayed in the website. Consent is provided for making this data openly accessible in the project website. <u>Is informed consent for data sharing and long term preservation given</u> : Yes.
Other Issues	No

## 6.6 Datasets collected in the context of WP10

### 6.6.1 Competitive call application datasets

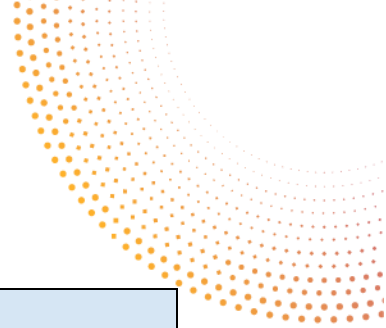
<b>DMP component</b>	<b>AI4Media_Data_158_WP10_Competitive_call_application_datasets_v1</b> <b>Partner: F6S</b>
Data Summary	<p><u>Purpose</u>: The competitive call applications dataset includes data collected during the open call application process. It contains:</p> <ul style="list-style-type: none"> <li>• The applications to the competitive calls;</li> <li>• The evaluation results;</li> <li>• Communications with applicants.</li> </ul> <p><u>Type/format</u>: Several types of data: documents, emails, texts etc.</p> <p><u>Re-use of existing data</u>: No, this is an original dataset to be created in the context of WP10 as a result of the two open call procedures.</p> <p><u>Data origin</u>: The data was collected from applicants who submitted a proposal to AI4Media - Open Call #1 or Open Call #2 via the F6S portal. Data was collected and consolidated by the leader of the evaluation process. Communication channels with applicants were defined in the “Guidelines for Applicants”; relevant messages were collected by the participating consortium members. All the information related to competitive calls was exported to relevant repositories to enable the long-term storage and access to data by consortium members and auditors. Once the competitive calls closed, all data was exported from the F6S portal to the project repository.</p> <p><u>Expected size</u>: Several GB</p> <p><u>Data utility</u>: Data collected was used by consortium members and external evaluators to perform the evaluation of the submitted applications and to decide whether they should be selected to participate in AI4Media. Information on applications selected to</p>





	<p>participate in the project was used to create the sub-granted project dataset. The dataset was also used to generate statistics and reports about the AI4Media project as requested by the Grant Agreement, which was then aggregated into anonymised data that do not compromise personal details of applicants nor any other confidential information about their projects.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: The data is stored in a relevant project platform. It will be discoverable only by authenticated users with metadata, identifiable by participant name and/or organisation and in some cases, be indexable/findable using a persistent and unique actor key.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, the platform supports versioning.</p> <p><u>Metadata creation</u>: A part of the data, such as applicant and evaluation data, is organised as structured data. Metadata used includes:</p> <ul style="list-style-type: none"> <li>• How data was created;</li> <li>• Time and date of creation or modification;</li> <li>• Source of data;</li> <li>• Who created the data.</li> </ul>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The raw data is not openly accessible due to AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information. However, anonymised or aggregated data from the open call participation is public.</p> <p><u>How it will be accessible</u>: The collected raw data is stored in a relevant project platform, only accessible by those partners with direct management of the open calls and respective information. Data including anonymised or aggregated data from the open call participation has been made public on the project website and public deliverables.</p> <p><u>Methods/software tools to access data</u>: Download files via web browser.</p> <p><u>Repository</u>: The collected data is stored on a relevant project platform. A dataset including anonymised or aggregated data from the open call participation has been made public on the project website.</p> <p><u>Restrictions on access</u>: Given the confidential nature of the information, access is restricted to those required for managing the data.</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: N/A.</p> <p><u>Data and metadata vocabularies</u>: Data is compiled using standard spreadsheet formats (Open Office, Excel, Google Sheets, etc.), standard text editors (Open Office, Word, Google format) and standard encoding and formatting to reduce the likelihood of incompatibility.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
<p>Increase data re-use</p>	<p><u>Licence</u>: The dataset will not be made openly available due to AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information.</p> <p><u>Availability for re-use</u>: No</p>





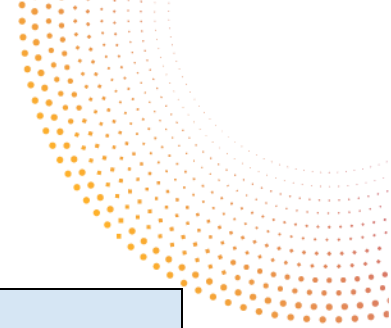
	<p><u>Usable by third parties after end of project:</u> No</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Completeness and conformity of the competitive calls exports is ensured by the leader of the evaluation process and the person responsible for performing the eligibility check as they will have to check all documents. Completeness and conformity of the evaluation data, including communications with applicants, is ensured by the evaluation committee as they will have to check all process documents to approve final ranking and selection.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> There are no additional costs to make data FAIR in the project, as the costs to operate each platform used in the project are already integrated into the project costs.</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> AI4Media will maintain protection of personal data and compliance with Data Regulations as per national and European legislation regarding the protection of personal data. Procedures have been kept in place for applicable technical means to avoid the loss, misuse, alteration, access by unauthorised persons and/or theft of the data provided to this entity. Notwithstanding, security measures (particularly for Internet accessible data) are not impregnable. To mitigate risk of unauthorised access, access controls will be applied to data sources. As an example, applications to the AI4Media project and forms have been made accessible to a limited number of team members using user-level security and permissions. AI4Media participants (data owners) will be able to exercise their right to be forgotten.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The raw data will not be made openly available due AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information. Only anonymised or aggregated data from the open call participation will be made public.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The applicants provided their informed consent for the collection and processing of their personal data in the context of AI4Media. The data will not be shared outside the consortium.</p>
Other Issues	N/A

### 6.6.2 Sub-granted projects dataset

DMP component	AI4Media_Data_159_WP10_sub-granted_projects_dataset_v1 Partner: F6S
Data Summary	<p><u>Purpose:</u> In each AI4Media competitive open call (Open Call #1 and Open Call #2), several projects were selected and invited to sign a sub-grant contract. The seed for this dataset is the information provided by applicants in the application process and by the selection committee in the evaluation report. The dataset was extended with all information needed to run the AI4Media programme, such as deliverables submitted by sub-grantees, evaluation reports, payment requests as will be defined in the Guidelines for Applicants. To summarize, the dataset includes (but is not limited to): applications, application evaluation, contracts, deliverables submitted, deliverables' evaluations, payment requests, proof of payments, amendments and copies or summaries of messages in any form, whose content may have an impact in the programme outcome.</p>







	<p><u>Type/format</u>: Several types of data: documents, emails, etc.</p> <p><u>Re-use of existing data</u>: No, this is an original dataset created in the context of WP10 as a result of the open call procedure.</p> <p><u>Data origin</u>: The data is collected from applicants, selection committee and project partners in the context of the WP10 open calls. The dataset is created during the implementation process as information becomes available.</p> <p><u>Expected size</u>: Several GB</p> <p><u>Data utility</u>: The data was useful to WP10 partners to run the open call procedures and monitor the smooth execution of the funded projects as defined in the Guidelines for Applicants. The dataset was also used to generate statistics and reports about the AI4Media project as requested by the Grant Agreement.</p>
<p>Making data findable, incl. provisions for metadata</p>	<p><u>Is data discoverable</u>: The data will be stored in a relevant project platform. It will be discoverable only by authenticated users with metadata, identifiable by participant name and/or organisation and in some cases, be indexable/findable using a persistent and unique actor key.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, the platform supports versioning.</p> <p><u>Metadata creation</u>: A part of the data, such as applicant and evaluation data, will be organised as structured data. Metadata used includes:</p> <ul style="list-style-type: none"> <li>• How data was created;</li> <li>• Time and date of creation or modification;</li> <li>• Source of data;</li> <li>• Who created data.</li> </ul>
<p>Making data openly accessible</p>	<p><u>Data openly accessible</u>: The raw data will not be openly accessible due to AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information. However, anonymised or aggregated data from the open call participation is public.</p> <p><u>How it will be accessible</u>: The collected raw data is stored in a relevant project platform, only accessible to those partners responsible with direct management of the open calls and respective information. Data including anonymised or aggregated info from the open call participation has been made public on the project website and project deliverables.</p> <p><u>Methods/software tools to access data</u>: Download files via web browser.</p> <p><u>Repository</u>: The collected data is stored on a relevant project platform. A dataset including anonymised or aggregated data from the open call participation has been made public on the project website.</p> <p><u>Restrictions on access</u>: Given the confidential nature of the information, access is restricted to only those required for managing the data.</p>
<p>Making data interoperable</p>	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: Data is compiled using standard spreadsheet formats (Open Office, Excel, Google Sheets, etc.), standard text editors (Open Office, Word, Google format) and standard encoding and formatting to reduce the likelihood of</p>



	<p>incompatibility.</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> The dataset will not be made openly available due AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information.</p> <p><u>Availability for re-use:</u> No</p> <p><u>Usable by third parties after end of project:</u> No</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> The Guidelines for Applicants will define the process to collect and verify the data with multiple checkpoints guaranteeing the quality of data.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> There are no additional costs to make data FAIR in the project, as the costs to operate each platform used in the project are already integrated into the project costs.</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> AI4Media will maintain protection of personal data and compliance with Data Regulations as per national and European legislation regarding the protection of personal data. Procedures have been kept in place for applicable technical means to avoid the loss, misuse, alteration, access by unauthorised persons and/or theft of the data provided to this entity. Notwithstanding, security measures (particularly for Internet accessible data) are not impregnable. To mitigate risk of unauthorised access, access controls will be applied to data sources. As an example, applications to the AI4Media project and forms have been made accessible to a limited number of team members using user-level security and permissions. AI4Media participants (data owners) will be able to exercise their right to be forgotten.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The raw data will not be made openly available due AI4Media's commitments to its applicants and sub-grantees in relation to personal information and business private information. Only anonymised or aggregated data from the open call participation will be made public.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The applicants provided their informed consent for the collection and processing of their personal data in the context of AI4Media. The data will not be shared outside the consortium.</p>
Other Issues	N/A

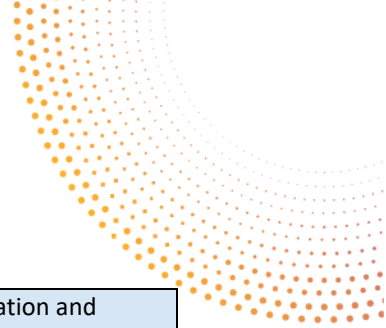
### 6.6.3 External experts and evaluators datasets

<b>DMP component</b>	<b>AI4Media_Data_160_WP10_external_experts_evaluators_datasets_v1</b> <b>Partner: F6S</b>
Data Summary	<p><u>Purpose:</u> As defined in the Guidelines for Applicants for the WP10 open calls, external evaluators were involved in multiple evaluation tasks during the AI4Media programme. This dataset is refers to the contractual relationship between experts and AI4Media. All data related with evaluation tasks is stored in the datasets "Competitive call applications" and "Sub-granted projects". The dataset includes (but is not limited to):</p>



	<p>contracts, payments request, proof of payments, evaluation reports and other.</p> <p><u>Type/format</u>: Several types of data: documents, emails, texts etc.</p> <p><u>Re-use of existing data</u>: No, this is an original dataset created in the context of WP10 as a result of the open call procedure.</p> <p><u>Data origin</u>: External experts' availability was collected through a form on the F6S portal by submitting their expression of interest. Contracts and related documents for selected evaluators, such as declarations of honour, receipts and other expenses information will be submitted either in a form in the F6S portal or by email.</p> <p><u>Expected size</u>: Several GB</p> <p><u>Data utility</u>: Data collected was used to manage the AI4Media programme as will be defined in the Guidelines for Applicants, including the evaluation of the submitted deliverables by external evaluators.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data is stored on a relevant project platform. It is discoverable only by authenticated users with metadata, identifiable by participant name and/or organisation and in some cases, be indexable/findable using a persistent and unique actor key.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, the platform supports versioning.</p> <p><u>Metadata creation</u>: A part of the data, such as applicant and evaluation data, will be organised as structured data. Metadata used includes:</p> <ul style="list-style-type: none"> <li>• How data was created;</li> <li>• Time and date of creation or modification;</li> <li>• Source of data;</li> <li>• Who created data.</li> </ul>
Making data openly accessible	<p><u>Data openly accessible</u>: The raw data is not openly accessible due to AI4Media's commitments to its applicants, sub-grantees and evaluators in relation to personal information and business private information. No data related to evaluators is public.</p> <p><u>How it will be accessible</u>: Data regarding evaluators will not be made accessible outside of those responsible for recruiting/ managing the evaluators</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: Access will be limited to those required for managing evaluators' data.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: Data is compiled using standard spreadsheet formats (Open Office, Excel, Google Sheets, etc.), standard text editors (Open Office, Word, Google format) and standard encoding and formatting to reduce the likelihood of incompatibility.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The dataset will not be made openly available due AI4Media's commitments to</p>





	<p>its applicants, sub-grantees and evaluators in relation to personal information and business private information.</p> <p><u>Availability for re-use:</u> No</p> <p><u>Usable by third parties after end of project:</u> No</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> The persons acting as project coordinator and project treasurer, as defined in the Grant Agreement, ensure that all documentation complies with legal requirements. The leader of the evaluation task ensures that the assigned tasks have been completed as contractually agreed and evidence is stored in relevant datasets.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> There are no additional costs to make data FAIR in the project, as the costs to operate each platform used in the project are already integrated into the project costs.</p> <p><u>Costs for long-term preservation:</u> N/A</p>
Data security	<p><u>Security measures:</u> AI4Media will maintain protection of personal data and compliance with Data Regulations as per national and European legislation regarding the protection of personal data. Procedures have been kept in place for applicable technical means to avoid the loss, misuse, alteration, access by unauthorised persons and/or theft of the data provided to this entity. Notwithstanding, security measures (particularly for Internet accessible data) are not impregnable. To mitigate risk of unauthorised access, access controls will be applied to data sources. As an example, applications to the AI4Media project and forms have been made accessible to a limited number of team members using user-level security and permissions. AI4Media participants (data owners) will be able to exercise their right to be forgotten.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The raw data will not be made openly available due AI4Media's commitments to its applicants, sub-grantees and evaluators in relation to personal information and business private information. Only anonymised or aggregated data will be made public.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The applicants provided their informed consent for the collection and processing of their personal data in the context of AI4Media. The data will not be shared outside the consortium.</p>
Other Issues	N/A

#### 6.6.4 Participants of competitive call related events datasets

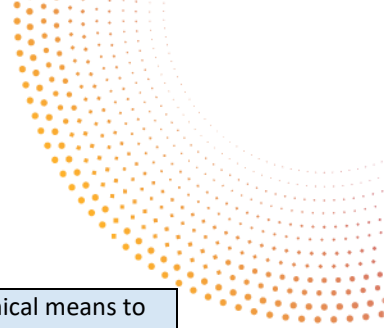
This dataset was collected in the context of both WP10 and WP11.

<b>DMP component</b>	<b>AI4Media_Data_161_WP10_ParticipantsCompetitiveCallsEvents_v1 Partner: F6S</b>
Data Summary	<p><u>Purpose:</u> During the organisation and delivery of the two competitive calls (AI4Media - Open Call #1 and Open Call #2) several information-related events were organised. The data included in this dataset is related to the participants of these events. Data collected - varying from event to event - included for example name, organisation, e-mail, country, industry.</p> <p><u>Type/format:</u> Several types of data: documents, emails, texts etc.</p>



	<p><u>Re-use of existing data</u>: No, this is an original dataset created in the context of WP10 as a result of the open call procedure.</p> <p><u>Data origin</u>: Participants' information was collected through event pages on the F6S portal or the registration functionalities of other event platforms.</p> <p><u>Expected size</u>: Several GB</p> <p><u>Data utility</u>: Data collected was used to identify potential applicants of the competitive calls or those with general interest in the competitive calls results.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: The data will be stored in a relevant project platform. It will be discoverable only by authenticated users with metadata, identifiable by participant name and/or organisation and in some cases, be indexable/findable using a persistent and unique actor key.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, the platform supports versioning.</p> <p><u>Metadata creation</u>: N/A</p>
Making data openly accessible	<p><u>Data openly accessible</u>: The raw data is not be openly accessible due to AI4Media's commitments in the framework of the competitive calls and respect to GDPR. Only anonymised or aggregated data will be made public and openly accessible.</p> <p><u>How it will be accessible</u>: Participation statistics without information about individuals will be included in AI4Media reports.</p> <p><u>Methods/software tools to access data</u>: N/A</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: Access will be limited to those required for managing the events.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A.</p> <p><u>Data and metadata vocabularies</u>: Data is compiled using standard spreadsheet formats (Open Office, Excel, Google Sheets, etc.), standard text editors (Open Office, Word, Google format) and standard encoding and formatting to reduce the likelihood of incompatibility.</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: The raw data is not openly accessible due to AI4Media's commitments in the framework of the competitive calls and respect to GDPR. Only anonymised or aggregated data will be made public and openly accessible.</p> <p><u>Availability for re-use</u>: No</p> <p><u>Usable by third parties after end of project</u>: No</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A as it will not be reused.</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: AI4Media will maintain protection of personal data and compliance with Data Regulations as per national and European legislation regarding the protection</p>





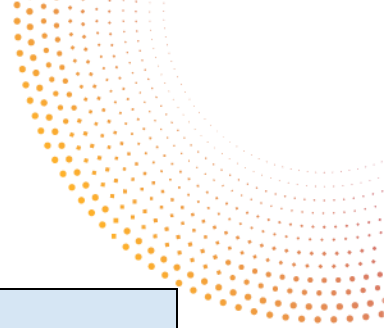
	of personal data. Procedures have been kept in place for applicable technical means to avoid the loss, misuse, alteration, access by unauthorised persons and/or theft of the data provided to this entity. Notwithstanding, security measures (particularly for Internet accessible data) are not impregnable. To mitigate risk of unauthorised access, access controls will be applied to data sources. AI4Media participants (data owners) will be able to exercise their right to be forgotten.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> The dataset includes personal and business information and therefore can not be shared. Only anonymised or aggregated data will be made public.</p> <p><u>Is informed consent for data sharing and long term preservation given:</u> The event participants provide their informed consent for the collection and processing of their personal data in the context of AI4Media. The data will not be shared outside the consortium.</p>
Other Issues	N/A

## 6.7 Datasets collected in the context of WP11

### 6.7.1 AI4Media associate members contact info dataset

DMP component	AI4Media_Data_162_WP1_Text_AI4MediaAssociateMembersContactInfo_v1 Partner: CERTH
Data Summary	<p><u>Purpose:</u> This dataset contains business contact information from AI4Media associate members, including name, affiliation and email of contact person. The collected data are necessary for the communication of AI4Media partners with the Associate members and are collected in the context of WP11.</p> <p><u>Type/format:</u> text</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Data provided by AI4Media Associate members to CERTH as part of an application form template (<a href="https://ai4media.eu/associate-members/">https://ai4media.eu/associate-members/</a>).</p> <p><u>Expected size:</u> &lt; 1MB</p> <p><u>Data utility:</u> This data is necessary for facilitating communication with the Associate members.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No, data are not discoverable from outside. The data is stored on the project wiki and is only discoverable by consortium members with a wiki account.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> Yes, through scheduled website backups.</p> <p><u>Metadata creation:</u> N/A.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No. The data will only be shared internally in AI4Media since it contains personal information.</p> <p><u>How it will be accessible:</u> Restricted access. The data is accessible only by wiki users.</p>



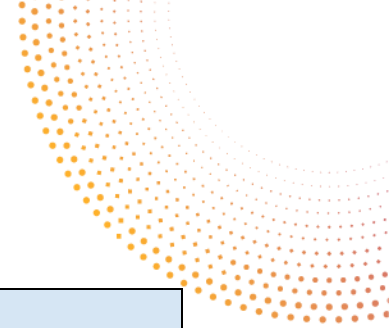


	<p><u>Methods/software tools to access data</u>: Web browser.</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u> N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on CERTH and LOBA servers and access is restricted only to authenticated users. CERTH and LOBA fully comply with the applicable national, European and International framework, and the GDPR. State-of-the-art IT security measures and company policies mitigate most of the risk of illegitimate access. Firewalls (to prevent illegitimate access from outside) and a rights-based-file system (to prevent illegitimate access from inside) are the countermeasures against this risk. Regular rolling daily backups are scheduled to minimize the risk of data loss. The data will be preserved there for three years after the end of the project and will then be deleted.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Dataset includes personal data (name and email) of associate members and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long term preservation given</u>: N/A</p>
Other Issues	No

### 6.7.2 AI4Media newsletter subscribers dataset

<b>DMP component</b>	<b>AI4Media_Data_163_WP11_Text_NewsletterSubscribers_v1</b> <b>Partner: LOBA</b>
Data Summary	<p><u>Purpose</u>: This dataset contains contact information including name and email, from people subscribing to the AI4Media newsletter. The collected data is necessary for distributing the newsletters to the subscribers and disseminating the project in the context of WP11.</p> <p><u>Type/format</u>: text</p>

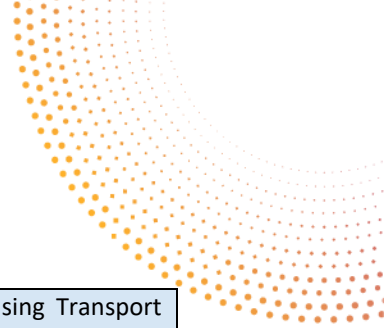




	<p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Data provided from people interested in subscribing to the projects' newsletter, through the subscription form available in the website (<a href="https://ai4media.eu/newsletters/">https://ai4media.eu/newsletters/</a>)</p> <p><u>Expected size</u>: &lt; 1MB</p> <p><u>Data utility</u>: This data is necessary for facilitating communication with the subscribers.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, data are not discoverable from outside. The data is stored on Zoho Campaigns, the email marketing software used by LOBA for managing and distribution of newsletters. The data is only accessed by LOBA.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: Yes, through scheduled website backups.</p> <p><u>Metadata creation</u>: N/A.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No. The data will only be shared internally in AI4Media since it contains personal information.</p> <p><u>How it will be accessible</u>: Restricted access. The data is accessible only by LOBA, who can share this information with the coordinator if requested.</p> <p><u>Methods/software tools to access data</u>: Website's back office (word press).</p> <p><u>Repository</u>: N/A</p> <p><u>Restrictions on access</u>: N/A</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A. The data will not be shared.</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on Zoho Campaigns, an email marketing software used for distribution of newsletters and email marketing. Zoho Campaigns fully complies with GDPR, from data collection and processing to managing data subject rights. The software handles and processes data, to ensure the additional level of security that GDPR encourages. Data at rest is encrypted using industry-standard AES-</p>





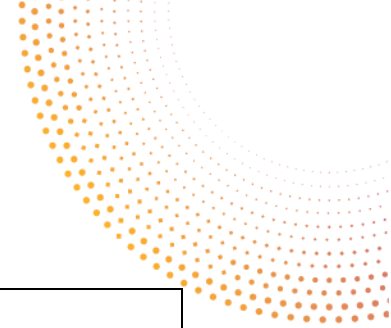


	256. All customer data is encrypted in transit over public networks using Transport Layer Security (TLS) 1.2/1.3 with Perfect Forward Secrecy (PFS) to protect it from unauthorized disclosure or modification.
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing:</u> Dataset includes personal data (name and email) of subscribers and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long-term preservation given:</u> N/A (Consent is given by the subscribers to only use their personal data for sending the newsletters)</p>
Other Issues	No

### 6.7.3 AI4Media website messages dataset

DMP component	AI4Media_Data_164_WP11_Text_WebsiteMessages_v1 Partner: CERTH, LOBA
Data Summary	<p><u>Purpose:</u> This dataset contains contact information (including name and email and messages) from people sending a message to AI4Media through the website's contact form. The collected data is necessary for replying to messages received through the website in the context of WP11.</p> <p><u>Type/format:</u> text</p> <p><u>Re-use of existing data:</u> No.</p> <p><u>Data origin:</u> Data provided from people sending messages through the contact form available in the website (<a href="https://ai4media.eu/contact/">https://ai4media.eu/contact/</a>)</p> <p><u>Expected size:</u> &lt; 1MB</p> <p><u>Data utility:</u> This data is necessary for facilitating communication with the stakeholders.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> No, data are not discoverable from the outside. The data is forwarded to the project's email list (<a href="mailto:info@ai4media.eu">info@ai4media.eu</a>) which only includes the emails from CERTH (project coordinator) and from LOBA (dissemination leader). The data is stored on CERTH's servers.</p> <p><u>Search keywords:</u> N/A</p> <p><u>Versioning:</u> Yes, through scheduled website backups.</p> <p><u>Metadata creation:</u> N/A.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No. The data will only be shared internally in AI4Media since it contains personal information.</p> <p><u>How it will be accessible:</u> Restricted access. The data is accessible only by CERTH and LOBA.</p> <p><u>Methods/software tools to access data:</u> Website back office (word press)</p> <p><u>Repository:</u> N/A</p> <p><u>Restrictions on access:</u> N/A</p>



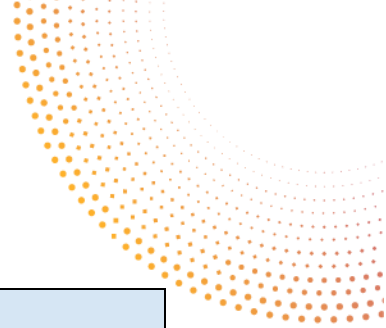


Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: The data is stored on CERTH's servers. Access requires username/password authentication. CERTH fully complies with the applicable national and European data protection frameworks &amp; guidelines, including the GDPR. State-of-the-art IT security measures and institutional policies mitigate the risk of illegitimate access, including firewalls (to prevent illegitimate access from outside) and a rights-based-file access system (to prevent illegitimate access from inside).</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Dataset includes personal data (name and email) of people sending messages to the project and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: N/A</p>
Other Issues	No

#### 6.7.4 AI4Media website analytics dataset

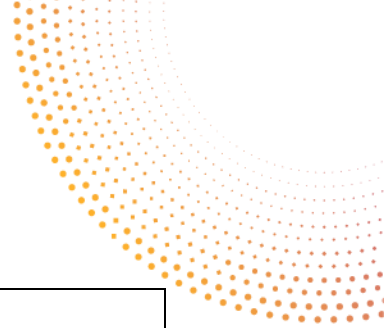
DMP component	AI4Media_Data_165_WP11_TEXT_AI4MediaWebsiteGoogleAnalytics_v1 Partner: LOBA
Data Summary	<p><u>Purpose</u>: The dataset comprises data and statistics collected from Google Analytics related to the AI4Media website. This includes metrics such as visitor numbers, page views, session duration, bounce rates, traffic sources, and user demographics. The purpose of collecting this data is to analyse the performance and reach of the AI4Media website, and particular web pages for example the AI Media Observatory, to understand user engagement, and to inform strategies for improving communication and dissemination activities. This aligns with the project objectives of enhancing visibility and impact of AI4Media's research and outputs.</p> <p><u>Type/format</u>: The data is in numerical and categorical format, typically exported as CSV or JSON files.</p> <p><u>Re-use of existing data</u>: No</p> <p><u>Data origin</u>: The data originates from Google Analytics tracking implemented on the</p>





	<p>AI4Media website.</p> <p><u>Expected size:</u> The expected size of the dataset is relatively small, estimated to be a few megabytes, depending on the volume of web traffic and the duration of data collection.</p> <p><u>Data utility:</u> This data is useful for monitoring and improving the website's performance and assess the impact of our content on the community, including the use and benefits of the AI Media Observatory.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable:</u> The data identifiable and locatable through standard identification mechanisms.</p> <p><u>Search keywords:</u> Search keywords such as "AI4Media website analytics," "visitor statistics," and "web traffic data"</p> <p><u>Versioning:</u> N/A</p> <p><u>Metadata creation:</u> Metadata includes details such as the date range of data collection, types of metrics collected, and any filtering or processing applied. This is created following standard web analytics metadata practices.</p>
Making data openly accessible	<p><u>Data openly accessible:</u> No, the Google Analytics data will not be made openly available due to privacy and security concerns.</p> <p><u>How it will be accessible:</u> The data is mainly accessed by LOBA for monitoring purposes, but access can be granted to any other AI4Media partner with interest in the information.</p> <p><u>Methods/software tools to access data:</u> The data can be accessed using Google Analytics dashboards and reports, which are available through standard web browsers.</p> <p><u>Repository:</u> The data will be stored within the Google Analytics platform and accessible through secure login.</p> <p><u>Restrictions on access:</u> Access will be restricted to authorized project members to ensure data privacy and security.</p>
Making data interoperable	<p><u>Interoperability:</u> The data is interoperable within the Google Analytics platform, allowing data exchange and re-use within the project team.</p> <p><u>Data and metadata vocabularies:</u> N/A</p> <p><u>Use of standard vocabularies:</u> N/A</p> <p><u>Mappings to commonly used vocabularies:</u> N/A</p>
Increase data re-use	<p><u>Licence:</u> N/A</p> <p><u>Availability for re-use:</u> N/A</p> <p><u>Usable by third parties after end of project:</u> N/A</p> <p><u>Re-use timeframe:</u> N/A</p> <p><u>Data quality assurance process:</u> Data quality is assured through regular checks and validation processes, ensuring accuracy and reliability.</p>
Allocation of resources	<p><u>Costs for making data FAIR:</u> N/A</p> <p><u>Costs for long-term preservation:</u> Long-term preservation costs include maintaining</p>





	access to the Google Analytics account.
Data security	<u>Security measures</u> : Security measures include controlled access to the Google Analytics account, user authentication, firewalls, encryption, and regular backups.
Ethical aspects	<u>Possible ethical and legal aspects preventing sharing</u> : N/A <u>Is informed consent for data sharing and long-term preservation given</u> : N/A
Other Issues	N/A

### 6.7.5 Dataset of registrants for AI4Media events

DMP component	AI4Media_Data_166_WP11_Text_RegistrantstoEvents_v1 Partner: LOBA, CERTH
Data Summary	<p><u>Purpose</u>: This dataset contains contact information (including name and email and messages) about individuals who register to attend various AI4Media events. The collected data is necessary for managing AI4Media events.</p> <p><u>Type/format</u>: The dataset is structured as a relational database or stored in a tabular format like CSV or Excel spreadsheets. Each record represents a registrant and includes fields such as name, email address, affiliation, role, event registered for, registration date, dietary restrictions and any specific preferences or interests related to the event topics.</p> <p><u>Re-use of existing data</u>: No.</p> <p><u>Data origin</u>: Data is primarily collected directly from registrants through online registration forms</p> <p><u>Expected size</u>: &lt; 1MB</p> <p><u>Data utility</u>: This data is used to facilitate event organisation, communication with registrants, personalisation of the event experience, and post-event follow-ups. It also aids in statistical analysis of participant demographics and interests to improve future event planning and targeting.</p>
Making data findable, incl. provisions for metadata	<p><u>Is data discoverable</u>: No, data is not discoverable from outside. The data is stored on registration systems used for managing event registrations. The data is only accessed by the event organisers.</p> <p><u>Search keywords</u>: N/A</p> <p><u>Versioning</u>: No.</p> <p><u>Metadata creation</u>: N/A.</p>
Making data openly accessible	<p><u>Data openly accessible</u>: No. The data will only be shared internally in AI4Media since it contains personal information.</p> <p><u>How it will be accessible</u>: Restricted access. The data is accessible only by event organisers.</p> <p><u>Methods/software tools to access data</u>: Registration systems used for managing event registrations e.g. eventbrite, google forms, pretix, etc.</p>



	<p><u>Repository</u>: Yes, on the registration systems.</p> <p><u>Restrictions on access</u>: Yes, only event organisers.</p>
Making data interoperable	<p><u>Interoperability</u>: N/A</p> <p><u>Data and metadata vocabularies</u>: N/A</p> <p><u>Use of standard vocabularies</u>: N/A</p> <p><u>Mappings to commonly used vocabularies</u>: N/A</p>
Increase data re-use	<p><u>Licence</u>: N/A</p> <p><u>Availability for re-use</u>: N/A</p> <p><u>Usable by third parties after end of project</u>: N/A</p> <p><u>Re-use timeframe</u>: N/A</p> <p><u>Data quality assurance process</u>: N/A</p>
Allocation of resources	<p><u>Costs for making data FAIR</u>: N/A</p> <p><u>Costs for long-term preservation</u>: N/A</p>
Data security	<p><u>Security measures</u>: Security practices include data encryption, access controls via authentication and authorisation, and regular security audits.</p>
Ethical aspects	<p><u>Possible ethical and legal aspects preventing sharing</u>: Dataset includes personal data (name and email, etc.) of registrants and will thus not be shared.</p> <p><u>Is informed consent for data sharing and long-term preservation given</u>: N/A (Consent is given by the registrants to only use their personal data for event purposes)</p>
Other Issues	No



## 7. Conclusions

This final version of the DMP identifies all the datasets managed by the AI4Media consortium during the project lifetime, organized per work package. The initial version of the DMP discussed 70 research datasets (19 created within the project by project partners and 51 third-party research datasets used by project partners) as well as 12 non-research datasets. In the present document, we present **142 research datasets in total**, including 81 pre-existing ones and **61 newly-created within the project** in the context of WP2, WP3, WP4, WP5, WP6, and WP8. Most of the datasets are openly available, to the benefit of the broader scientific community. In addition, we describe how we manage and protect non-research data, including **24 non-research datasets** collected in the context of WP1, WP2, WP7, WP8, WP9, WP10 and WP11.

The datasets collected and used in the project include: **media-related datasets** (including videos, audio files, social media posts, user profiles, etc.) that were employed for the design and development of AI methodologies, algorithms, and tools in the context of WP3, WP4, WP5 and WP6, aiming to support the seven use cases; **questionnaire data** including questionnaires collected from project partners, associate members and end-users aiming to identify use case requirements, provide guidance to technical partners for the development of the AI tools but also assess the effectiveness and impact of the developed tools and demonstrators; **user activity data and software analytics** automatically collected during the use cases evaluation, with the purpose of evaluating and improving the various AI tools integrated in the demonstrators; **personal data of members of the research, academic, student and business communities** participating in AI4Media educational, outreach and dissemination activities, applying for funding through the AI4Media open calls, or contributing to AI Media Observatory.

With regard to making the research data FAIR, our approach may be summarized as follows:

- **Findable:** The datasets that are made publicly available are uploaded in open repositories like Zenodo or GitHub, thus making this data both easily discoverable and identifiable from the outside. Datasets that are only used internally by project partners are stored on partners' servers. In both cases, the datasets are internally discoverable and identifiable using simple queries with keywords or filters.
- **Accessible:** In the context of the project, we try to make publicly available the research datasets created by consortium members in the context of WPs 2, 3, 4, 5, 6, and 8, wherever possible. The data is shared on open repositories like Zenodo and GitHub or institutional (open) repositories of partners. More specifically, **34 out of the 61 datasets created within AI4Media are already openly shared** with the community and are listed in the project website<sup>15</sup>. The remaining 27 datasets are not openly shared due to various ethical and legal aspects that prevent/prohibit sharing (e.g. some social media data cannot be shared due to social media platforms' terms of use, use case evaluation data cannot be shared because they usually include personal data of end users, datasets of news articles

---

<sup>15</sup> <https://www.ai4media.eu/open-datasets/>



cannot be shared because the copyright is owned by the relevant media outlets, datasets collected from media companies' archives cannot be shared due to copyright etc.).

Many of the data used in the project for research and development of new AI methods and tools (in WPs 3, 4, 5, 6) is already **open data**, made openly available by research organizations or media companies. Since it is already open, as a general policy, we do not re-share them. Sharing some of this data is handled on a case-by-case basis, and will only be pursued in cases where the data license allows it and re-sharing of the data (in some new form or after some processing) provides some additional benefit to the research community. In any case, we try to provide open software tools that will allow other researchers to easily crawl and collect data from all open data sources.

In addition to open data, there are also **privately owned datasets**. Such data has been provided to the project for research purposes and will be only used internally by project partners. Effort is made to make some of this data available in cooperation with the data owners, as long as there are no legal or ethical issues for their sharing.

**Data from surveys** addressed to project partners, associate members, and end users of AI tools in the context of the seven use cases will not be made openly accessible (at least most of them) since they may contain sensitive or personal information. Where possible and in case there is added value from their sharing, data is fully anonymized before being shared. The collected data (whether public, private, or personal) is analysed and analysis results are published as part of public project deliverables or publications that are available in open repositories like Zenodo.

- **Interoperable:** Effort is dedicated to making the data interoperable, wherever possible. The data and metadata vocabularies adopted for each data source have been also discussed.
- **Reusable:** Effort is made to increase the re-use of the data that we plan to make open, through clarifying licenses. Licensing is examined on a case-by-case basis depending on the dataset. In case of data coming from external sources or in cases where the data comes with a license of its own, the data will be re-shared (where necessary) under the same licence. For other datasets, a CC-BY 4.0 (Creative Commons Attribution 4.0 International License) license will be selected.

Datasets created within AI4Media are either openly shared (by uploading them in open repositories) or shared internally among specific partners (usually stored on the institutional servers of AI4Media partners). Datasets to be openly shared, are mostly deposited in certified repositories like Zenodo and GitHub that have in place strong mechanisms and protocols for **data security** and long-term data preservation. Similar mechanisms are in place in both the project wiki and the partners' institutional servers to ensure data protection.

Finally, addressing **legal and ethics challenges** is an important part of the AI4Media work plan. A legal partner (KUL) forms part of the consortium and dedicated tasks (T1.3 - *Ethical issue management* and Task 4.1 - *Legal and ethical frameworks for trusted AI*) and work packages (WP12 – *Ethics Management*) deal specifically with such issues. Moreover, an Ethical Advisory Board provides guidance on such issues. The aim is to identify the relevant EU legal and ethical frameworks and relevant requirements and provide guidance to partners on issues of data



collection and data privacy and protection. To conform to the GDPR, consent forms and information sheets have been provided to users whose personal data were collected and processed in the context of the AI4Media use cases or other project activities.

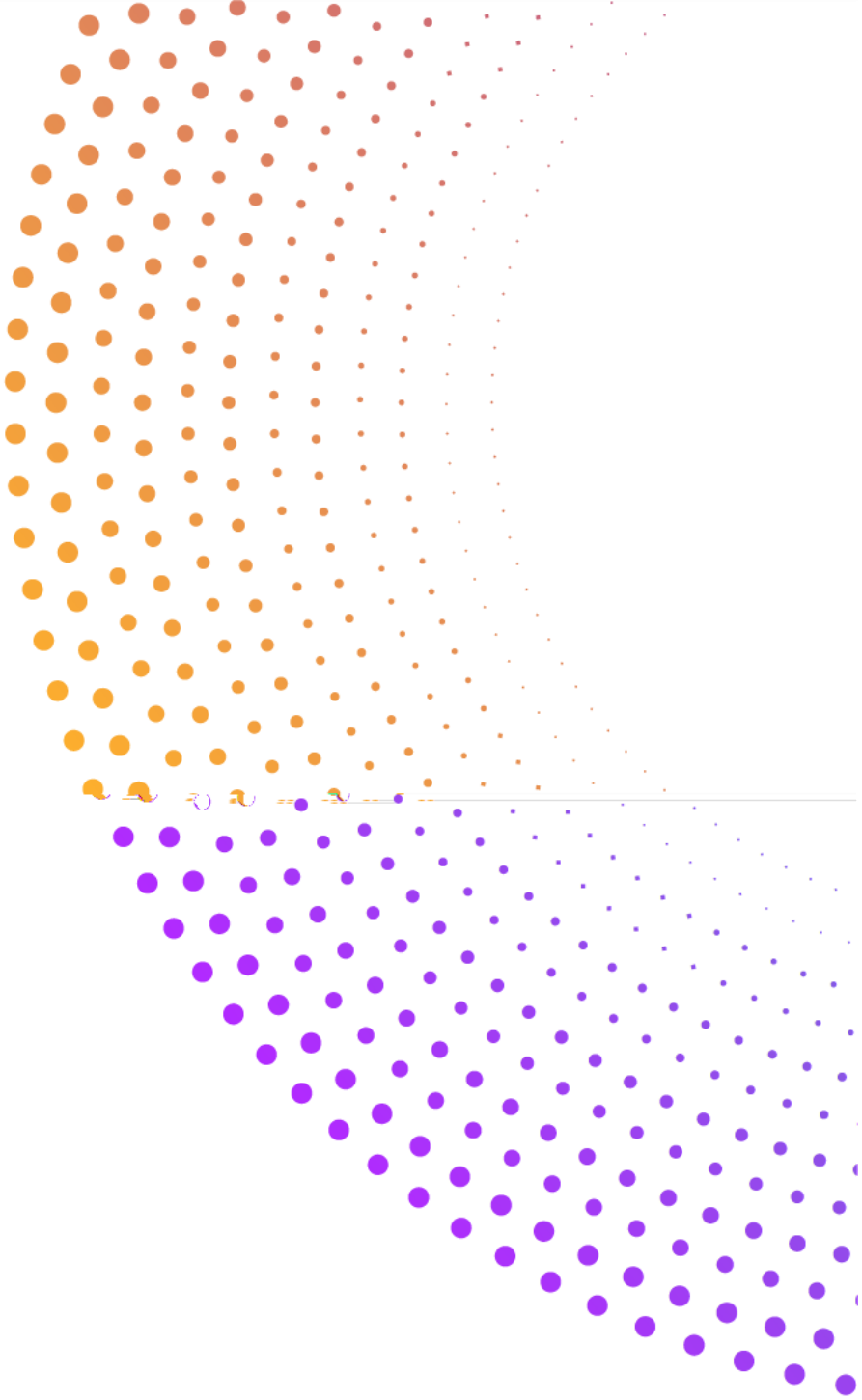
The DMP has been a **living document** collaboratively updated throughout the lifetime of the project.





# AI4media

ARTIFICIAL INTELLIGENCE FOR  
THE MEDIA AND SOCIETY



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 951911

[info@ai4media.eu](mailto:info@ai4media.eu)

[www.ai4media.eu](http://www.ai4media.eu)