
ACAV-1M: Data Curation and Benchmarking for Audio-Visual Representation Learning (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 In this supplementary material, we provide the following material:

- 2 • dataset documentation and intended uses in Section 1,
- 3 • dataset website in Section 2,
- 4 • croissant metadata in Section 3,
- 5 • author statement in Section 4,
- 6 • hosting, licensing, and maintenance plan in Section 5,
- 7 • algorithm for our data curation pipeline in Section 6,
- 8 • addition implementation and datasets details in Section 7.

9 1 Dataset Documentation & Intended Uses

10 The *ACAV-1M* dataset is designed to facilitate research in audio-visual representation learning. It
11 includes synchronized audio and visual data curated from various sources to ensure a diverse
12 and comprehensive collection for training and evaluating machine learning models. The dataset
13 documentation follows the *datasheets for datasets* framework, providing detailed information on the
14 dataset's composition, collection process, and intended uses.

15 Composition: The dataset consists of 1 million audio-visual pairs, including video clips with corre-
16 sponding audio tracks from various domains such as user-generated content.

17 Collection Process: Data was collected using automated scripts and manual curation to ensure quality
18 and relevance. Metadata includes source URLs, timestamps, and content descriptions.

19 Intended Uses: The dataset is intended for developing and benchmarking models in audio-visual
20 representation learning, including tasks like video classification, audio-visual synchronization, and
21 cross-modal retrieval.

22 Ethical Considerations: We ensured that the dataset adheres to ethical guidelines, including the
23 exclusion of sensitive or inappropriate content and respect for copyright and privacy concerns.

24 This document is based on *Datasheets for Datasets* by Gebru *et al.* [1].

MOTIVATION

25 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a
26 specific gap that needed to be filled? Please provide a description.

27 The *ACAV-1M* was created to fill a significant gap in multimodal learning where audio and visual
28 data are integrated systematically. It aims to enhance robust models that leverage both modalities
29 for improved understanding and interaction, designed specifically for tasks like audio-visual
30 classification, localization, retrieval, and segmentation.

31

32 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**
33 **company, institution, organization)?**

34 The dataset was created by a collaborative effort involving researchers from various academic
35 institutions specializing in machine learning and computer vision, under the coordination of a leading
36 university's computer science department.

37

38 **What support was needed to make this dataset?** (e.g. who funded the creation of the dataset? If
39 there is an associated grant, provide the name of the grantor and the grant name and number, or if it
40 was supported by a company or government agency, give those details.)

41 No. The creation of the *ACAV-1M* was not supported by any grants from several research funding
42 agencies. However, the dataset development received technical support and infrastructure from the
43 host university.

44

45 **Any other comments?**

46 No.

47

COMPOSITION

48 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
49 **countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and
50 interactions between them; nodes and edges)? Please provide a description.

51 The instances in the *ACAV-1M* represent synchronized audio-visual clips from diverse settings,
52 including music performances, public speeches, and everyday activities, ensuring a wide range of
53 scenarios for robust multimodal learning.

54

55 **How many instances are there in total (of each type, if appropriate)?**

56 The dataset comprises approximately 100,000 video clips, each paired with corresponding audio
57 tracks that have been meticulously synchronized and annotated.

58

59 **Does the dataset contain all possible instances or is it a sample (not necessarily random)**
60 **of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the
61 sample representative of the larger set (e.g., geographic coverage)? If so, please describe how
62 this representativeness was validated/verified. If it is not representative of the larger set, please
63 describe why not (e.g., to cover a more diverse range of instances, because instances were withheld
64 or unavailable).

65 The dataset is a curation subset of the original *ACAV* [2] dataset with 100 million samples.

66

67 **What data does each instance consist of?** "Raw" data (e.g., unprocessed text or images) or
68 features? In either case, please provide a description.

69 Each instance consists of "Raw" video and audio data. Additional metadata include synchronization
70 points, annotations for source localization, and labels for classification and segmentation tasks.

71

72 **Is there a label or target associated with each instance?** If so, please provide a description.

73 Yes, each instance includes captions associated with the video and audio. For various tasks, we
74 include labels for each instance like classification (audio-visual context), segmentation masks, and

75 localization coordinates.

76

77 **Is any information missing from individual instances?** If so, please provide a description,
78 explaining why this information is missing (e.g., because it was unavailable). This does not include
79 intentionally removed information, but might include, e.g., redacted text.

80 No.

81

82 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social
83 network links)?** If so, please describe how these relationships are made explicit.

84 No.

85

86 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so,
87 please provide a description of these splits, explaining the rationale behind them.

88 No.

89

90 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a
91 description.

92 No.

93

94 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,
95 websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there
96 guarantees that they will exist, and remain constant, over time; b) are there official archival versions
97 of the complete dataset (i.e., including the external resources as they existed at the time the dataset
98 was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external
99 resources that might apply to a future user? Please provide descriptions of all external resources and
100 any restrictions associated with them, as well as links or other access points, as appropriate.

101 Yes. The dataset is a curation subset of the original ACAV [2] dataset with 100 million samples.

102

103 **Does the dataset contain data that might be considered confidential (e.g., data that is protected
104 by legal privilege or by doctor-patient confidentiality, data that includes the content of
105 individuals' non-public communications)?** If so, please provide a description.

106 No.

107

108 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,
109 or might otherwise cause anxiety?** If so, please describe why.

110 No.

111

112 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

113 No.

114

115 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how
116 these subpopulations are identified and provide a description of their respective distributions within
117 the dataset.

118 No.

119

120 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or
121 indirectly (i.e., in combination with other data) from the dataset?** If so, please describe how.

122 No.

123

124 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that**
125 **reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or**
126 **union memberships, or locations; financial or health data; biometric or genetic data; forms of**
127 **government identification, such as social security numbers; criminal history)?** If so, please
128 provide a description.

129 No.

131 **Any other comments?**

132 No.

133

COLLECTION

134 **How was the data associated with each instance acquired?** Was the data directly observable (e.g.,
135 raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived
136 from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If data was
137 reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If
138 so, please describe how.

139 The data was acquired through a combination of public domain resources and contributions from
140 collaborating institutions, where scenarios were staged and recorded under controlled conditions to
141 ensure quality and diversity.

142

143 **Over what timeframe was the data collected?** Does this timeframe match the creation timeframe
144 of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please
145 describe the timeframe in which the data associated with the instances was created. Finally, list when
146 the dataset was first published.

147 Data collection spanned over half one year, culminating in the dataset's release in 2024. The temporal
148 alignment of collection and creation ensured the relevance and recency of the data.

149

150 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
151 **sensor, manual human curation, software program, software API)?** How were these mechanisms
152 or procedures validated?

153 We have alignment filtering mechanisms to curate our dataset from the original ACAV [2] dataset.

154

155 **What was the resource cost of collecting the data?** (e.g. what were the required computational
156 resources, and the associated financial costs, and energy consumption - estimate the carbon footprint.
157 See Strubell *et al.*[3] for approaches in this area.)

158 We use four A100 GPUs to curate data and train our models.

159

160 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
161 **probabilistic with specific sampling probabilities)?**

162 We used alignment filtering mechanisms.

163

164 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
165 **and how were they compensated (e.g., how much were crowdworkers paid)?**

166 Authors are involved in the data curation process.

167

168 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so,
169 please provide a description of these review processes, including the outcomes, as well as a link or
170 other access point to any supporting documentation.

171 No.

172

173 **Does the dataset relate to people?** If not, you may skip the remainder of the questions in this
174 section.

175 No.

176

177 **Did you collect the data from the individuals in question directly, or obtain it via third parties
178 or other sources (e.g., websites)?**

179 No.

180

181 **Were the individuals in question notified about the data collection?** If so, please describe (or
182 show with screenshots or other information) how notice was provided, and provide a link or other
183 access point to, or otherwise reproduce, the exact language of the notification itself.

184 No.

185

186 **Did the individuals in question consent to the collection and use of their data?** If so, please
187 describe (or show with screenshots or other information) how consent was requested and provided,
188 and provide a link or other access point to, or otherwise reproduce, the exact language to which the
189 individuals consented.

190 No.

191

192 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke
193 their consent in the future or for certain uses?** If so, please provide a description, as well as a link
194 or other access point to the mechanism (if appropriate)

195 No.

196

197 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data
198 protection impact analysis) been conducted?** If so, please provide a description of this analysis,
199 including the outcomes, as well as a link or other access point to any supporting documentation.

200 No.

201

202 **Any other comments?**

203 No.

204

PREPROCESSING / CLEANING / LABELING

205 **Was any preprocessing/cleaning/labeling of the data done(e.g., discretization or bucketing,
206 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing
207 of missing values)?** If so, please provide a description. If not, you may skip the remainder of the
208 questions in this section.

209 Yes. We use multimodal LLM to .

210

211 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support
212 unanticipated future uses)?** If so, please provide a link or other access point to the “raw” data.

213 No.

214

215 **Is the software used to preprocess/clean/label the instances available?** If so, please provide a
216 link or other access point.

217 No.

218

219 **Any other comments?**

220 No.

221

USES

222 **Has the dataset been used for any tasks already?** If so, please provide a description.

223 Yes, *ACAV-1M* has been employed in several benchmarking tasks within the research group, including
224 preliminary studies on audio-visual perception tasks.

225

226 **Is there a repository that links to any or all papers or systems that use the dataset?** If so,
227 please provide a link or other access point.

228 No.

229

230 **What (other) tasks could the dataset be used for?**

231 Beyond the current uses, the dataset holds potential for tasks in automated content generation,
232 assistive technologies, and advanced surveillance systems.

233

234 **Is there anything about the composition of the dataset or the way it was collected and
235 preprocessed/cleaned/labeled that might impact future uses?** For example, is there anything that
236 a future user might need to know to avoid uses that could result in unfair treatment of individuals or
237 groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms,
238 legal risks) If so, please provide a description. Is there anything a future user could do to mitigate
239 these undesirable harms?

240 No.

241

242 **Are there tasks for which the dataset should not be used?** If so, please provide a description.

243 No.

244

245 **Any other comments?**

246 No.

247

DISTRIBUTION

248 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,
249 organization) on behalf of which the dataset was created?** If so, please provide a description.

250 No.

251

252 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does the
253 dataset have a digital object identifier (DOI)?

254 The dataset is available via a website page and can be accessed through the dataset page, which
255 ensures controlled and ethical usage aligned with academic standards.

256

257 **When will the dataset be distributed?**

258 The dataset will be available upon publication.

259

260 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
261 **and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and
262 provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU,
263 as well as any fees associated with these restrictions.

264 No.

266 **Have any third parties imposed IP-based or other restrictions on the data associated with**
267 **the instances?** If so, please describe these restrictions, and provide a link or other access point
268 to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these
269 restrictions.

270 No.

272 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**
273 **instances?** If so, please describe these restrictions, and provide a link or other access point to, or
274 otherwise reproduce, any supporting documentation.

275 No.

277 **Any other comments?**

278 No.

MAINTENANCE

280 **Who is supporting/hosting/maintaining the dataset?**

281 The dataset is maintained by the authors, with plans for ongoing updates and expansions based on
282 community feedback and technological advancements.

284 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

285 The owner of the dataset can be contacted by email.

287 **Is there an erratum?** If so, please provide a link or other access point.

288 No.

290 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
291 **instances)?** If so, please describe how often, by whom, and how updates will be communicated to
292 users (e.g., mailing list, GitHub)?

293 Yes, the dataset is scheduled for regular reviews and updates to address any errors, introduce new
294 instances, and phase out obsolete data, with all changes communicated through the dataset's official
295 repository.

297 **If the dataset relates to people, are there applicable limits on the retention of the data**
298 **associated with the instances (e.g., were individuals in question told that their data would be**
299 **retained for a fixed period of time and then deleted)?** If so, please describe these limits and
300 explain how they will be enforced.

301 No.

303 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please
304 describe how. If not, please describe how its obsolescence will be communicated to users.

305 Yes. It will be maintained on the dataset website.

307 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
308 **them to do so?** If so, please provide a description. Will these contributions be validated/verified? If
309 so, please describe how. If not, why not? Is there a process for communicating/distributing these
310 contributions to other users? If so, please provide a description.

311 Yes. We will open the opportunity for other researchers to augment the dataset for additional
312 benchmarks.

313

314 **Any other comments?**

315 No.

316

317 **2 Dataset Website**

318 The dataset and its documentation can be accessed at the following URL:

319 <https://acav1m.github.io>

320 This website provides an overview of the dataset, download links, and additional resources such
321 as example code, tutorials, and a forum for community discussions. Users can explore the dataset
322 through an interactive interface, which includes search and filter options to facilitate easy access to
323 specific subsets of the data.

324 **3 Croissant Metadata**

325 The Croissant metadata for the ACAV-1M dataset is available at:

326 <https://acav1m.github.io>

327 This metadata record documents the dataset’s structure, including descriptions of the files, their
328 formats, and the fields within each record. The metadata adheres to the Croissant format, ensuring
329 interoperability and ease of use with ML tools and platforms.

330 **4 Author Statement**

331 We, the authors of the ACAV-1M dataset, bear full responsibility for any violations of rights and
332 confirm that all data included in the dataset complies with the relevant licenses and ethical guidelines.
333 The dataset is released under the Creative Commons Attribution 4.0 International License (CC BY
334 4.0), which allows for sharing, adaptation, and use of the data with appropriate credit given to the
335 original authors.

336 **5 Hosting, Licensing, & Maintenance Plan**

337 The dataset is hosted on a dedicated server managed by our institution, ensuring reliable access and
338 download speeds. We also provide mirror links through major cloud storage providers to ensure
339 redundancy and availability. The dataset is licensed under the Creative Commons Attribution 4.0
340 International License.

341 **6 Pseudo Algorithm for ACAV-1M Data Curation and Alignment**

342 Algorithm 1 is a pseudo-algorithm that encapsulates the data curation and alignment filtering processes
343 described for the *ACAV-1M* dataset. This algorithm is structured to provide a clear, step-by-step
344 procedure that reflects the robust methodologies used in preparing the dataset. This algorithm also
345 provides a structured approach to processing and aligning the data within the *ACAV-1M* ensuring
346 that each component (video, audio, and textual caption) is effectively synchronized and semantically

Algorithm 1 Pseudo Algorithm for Data Curation and Alignment Filtering

```
1: Input: Raw video and audio data
2: Output: Curated dataset with aligned audio-visual captions

3: Data Curation Process:
4: for each video clip in dataset do
5:   Extract raw audio and video streams
6:   Use VideoLLaVA [4] to generate sentence-level descriptions from video
7:   Condense descriptions using a general LLM [5] into a single comprehensive caption
8:   Attach the comprehensive caption to the corresponding video clip
9: end for

10: Alignment Filtering Process:
11: for each item in curated dataset do
12:   Language Alignment:
13:   Calculate normalized cosine similarity between text captions and audio-visual content
14:   if similarity < 0.5 then
15:     Flag for review or reprocessing
16:   end if
17:   Instance Alignment:
18:   Assess synchronization between audio and visual streams using ImageBind [6]
19:   if similarity < 0.5 then
20:     Adjust synchronization parameters and re-align
21:   end if
22:   Temporal Alignment:
23:   Check for alignment within a temporal window of 1 second per segment
24:   if average alignment threshold < 0.5 then
25:     Refine temporal synchronization parameters
26:   end if
27: end for

28: Return finalized dataset with validated and aligned captions
```

347 coherent. The algorithm is designed to be part of a larger document or paper, offering clarity on the
348 methods and steps taken to curate and align data within the dataset.

349 7 Implementation & Dataset Details

350 In this section, we provide more implementation and dataset details.

351 **Audio-visual classification.** For linear probing, we follow the prior work [7, 8] and extract frozen
352 audio-visual representations from our *ACAV-1M* pre-trained audio-visual masked autoencoder. Then
353 we attach a linear layer as a head to the frozen features for training with the audio-visual classes.
354 During training, we only fine-tune the linear head to evaluate the quality of pre-trained features.
355 The models are trained for 50 epochs using the Adam optimizer [9] with a learning rate of $1e-4$
356 and a batch size of 128. For fine-tuning, we use the same optimizer and batch size settings, but all
357 parameters are learnable.

358 **Audio-Visual Source Localization.** For sound source localization, we train all baselines [10, 11, 12]
359 using the same backbone (*i.e.*, ViT-Base) for audio/visual encoder with different proposed objectives
360 in their original papers. The final localization map is generated through bilinear interpolation of the
361 similarity map between audio/visual features from the last self-attention layer. The models are trained
362 for 30 epochs using the Adam optimizer [9] with a learning rate of $1e-4$ and a batch size of 128.

363 **Audio-Visual Retrieval.** The retrieval task processes video frames sampled at 8 fps and utilizes
364 combined low-level visual features from ResNet-152 [13] and 3D ResNet models [14], both pre-
365 trained on respective large-scale datasets. Audio features are extracted using VGGish [15], pre-trained

on AudioSet [16]. The complete model, integrating these features, is trained using Adam to optimize retrieval effectiveness across 40 epochs.

Audio-Visual Video Parsing. Following the data pre-processing in previous work [17], we sample video frames at 8 fps from the 10-second videos with 10 non-overlapping snippets of 1 second. For low-level visual features, we concatenate 2D and 3D visual features extracted by ResNet-152 [13] pre-trained on ImageNet [18] and 3D ResNet [14] pre-trained on Kinetics-400 [19]. We utilize VGGish [15] pre-trained on AudioSet [16] to extract the audio features. The model is trained with Adam [9] optimizer with $\beta_1=0.9$, $\beta_2=0.999$ and with an initial learning rate of $3e-4$. We train the model with a batch size of 16 for 40 epochs. Note that each video includes at least 1s audio or visual event, and 7202 video clips are annotated with more than one event category. We use 10,000 video clips with only video-level event labels for training. Following the official splits [17] of validation and test sets, we develop and test the model on the remaining 1879 videos with the segment-level annotations, *i.e.*, the speech event for audio starts at 1s and ends at 5s.

Audio-Visual Scene-Aware Dialog. In the audio-visual scene-aware dialog task, our model employs an advanced dialog generation framework that integrates audio and visual information to produce contextually relevant conversations. The dialog system utilizes a Transformer-based architecture, which processes inputs from both modalities through separate encoders before merging them in a fusion layer. This approach allows the model to understand the context provided by both the audio and visual data streams effectively. The model is optimized using the Adam optimizer with a learning rate of $1e-4$ and a batch size of 64. Training is conducted for up to 30 epochs, with early stopping based on performance on a validation set to prevent overfitting.

Audio-Visual Question-Answering. For the AVQA task, our implementation focuses on integrating spatial and temporal grounding techniques to accurately answer questions based on the video and audio content. The system employs a dual-stream encoder that separately processes visual and audio inputs. The encoded features are then combined using a co-attention mechanism that aligns audio and visual elements relevant to the question context. This integration allows the model to focus on specific segments of audio and video that are crucial for answering the given question. The model is trained using the Adam [9] optimizer with an initial learning rate of $3e-4$, reduced by a factor of 0.1 upon plateauing of validation loss. The system is trained for 40 epochs with a batch size of 32.

Audio-Visual Segmentation. For segmentation, we follow the prior work [20], and apply an upsampling decoder on features from the last self-attention layer to generate the final segmentation mask. We use the binary cross entropy (BCE) loss between the prediction and ground-truth masks for training. The models are trained for 20 epochs using the Adam optimizer [9] with a learning rate of $1e-4$ and a batch size of 128.

Audio-Visual Source Separation. For sound source separation, we follow the previous method [21, 22] and attach an audio U-Net decoder to our pre-trained audio-visual encoders for separating sounds from the mixture. The decoder depth for self-attention layers is 8, and the decoder receives the representations of the audio mixture and the visual embeddings. We also apply multiple transposed convolutions and an output head to predict a time-frequency separation mask. This separation mask is then used to multiply the input mixture STFT to separate the audio. Similarly to [21], the target masks refer to the time-frequency bins where the source is the most dominant component in the mixture. The sound source separation is achieved by optimizing a binary cross-entropy loss over these binary targets. The model is trained for 20 epochs using the Adam optimizer [9] with a learning rate of $1e-4$ and a batch size of 128.

Dataset Details. We evaluated our method using several prominent audio-visual datasets:

- **Flick-SoundNet [23]:** a dataset consisting of natural soundscapes with associated Flickr images with 4,500 audio-visual pairs for training and testing the model on 250 audio-visual pairs of sounding objects and extended 250 non-sounding objects;

- 414 • **VGG-Instruments [12]**: contains video clips of musical instrument performances, with 32k
415 video clips of 10s lengths from 36 musical instrument classes, a subset of VGG-Sound [24],
416 and each video only has one single instrument class;
- 417 • **MUSIC [21]**: consists of 448 untrimmed YouTube music videos of solos and duets from 11
418 instrument categories;
- 419 • **VGG-Music [22]**: a dataset that features a collection of music videos with annotations
420 related to the genre and instruments present;
- 421 • **VGGSound [24]**: a comprehensive dataset that includes a wide variety of sound categories
422 and corresponding visual scenes, which contains categories, such as animals, instruments,
423 vehicles, people, etc;
- 424 • **AudioSet [16]**: a collection of 2,084,320 human-labeled 10-second sound clips drawn from
425 YouTube videos with 632 audio event classes;
- 426 • **AVSBench [20]**: a benchmark for testing audio-visual synchronization and alignment
427 in diverse settings, including 4,932 videos (in total 10,852 frames) from 23 categories,
428 including instruments, humans, animals, etc.
- 429 • **MSR-VTT [25]**: A large-scale video description dataset that includes 10,000 video clips,
430 each paired with 20 human-annotated captions, useful for tasks involving video understand-
431 ing and retrieval.
- 432 • **LLP [17]**: The Look, Listen, and Parse (LLP) dataset contains densely labeled video
433 segments that are used to train and evaluate models on tasks requiring fine-grained temporal
434 understanding of video content.
- 435 • **MUSIC-AVQA [26]**: A dataset specifically curated for audio-visual question answering
436 in 11,849 YouTube video clips of 10 seconds long from 25 different event categories. It
437 combines visual and audio clues to answer complex queries about the content and context
438 of musical pieces.

References

- [1] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Dauméé III, and Kate Crawford. Datasheets for Datasets. *arXiv preprint arXiv:1803.09010*, 2018. 1
- [2] Sangho Lee, Jiwan Chung, Youngjae Yu, Gunhee Kim, Thomas Breuel, Gal Chechik, and Yale Song. Acav100m: Automatic curation of large-scale datasets for audio-visual video representation learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2, 3, 4
- [3] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and Policy Considerations for Deep Learning in NLP. *arXiv preprint arXiv:1906.02243*, 2019. 4
- [4] Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023. 9
- [5] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 9
- [6] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 9
- [7] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. 9
- [8] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *Proceedings of Advances In Neural Information Processing Systems (NeurIPS)*, 2022. 9
- [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 9, 10
- [10] Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In *Proceedings of European Conference on Computer Vision (ECCV)*, page 218–234, 2022. 9
- [11] Shentong Mo and Pedro Morgado. A closer look at weakly-supervised audio-visual source localization. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 9
- [12] Xixi Hu, Ziyang Chen, and Andrew Owens. Mix and localize: Localizing sound sources in mixtures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10483–10492, 2022. 9, 11
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 9, 10
- [14] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6450–6459, 2018. 9, 10
- [15] Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss, and Kevin W. Wilson. CNN architectures for large-scale audio classification. In *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017. 9, 10

- 484 [16] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Chan-
 485 ning Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled
 486 dataset for audio events. In *Proceedings of 2017 IEEE International Conference on Acoustics,
 487 Speech and Signal Processing (ICASSP)*, pages 776–780, 2017. [10](#), [11](#)
- 488 [17] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-
 489 supervised audio-visual video parsing. In *Proceedings of European Conference on Computer
 490 Vision (ECCV)*, page 436–454, 2020. [10](#), [11](#)
- 491 [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia. Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-
 492 Scale Hierarchical Image Database. In *Proceedings of IEEE/CVF Conference on Computer
 493 Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. [10](#)
- 494 [19] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the
 495 kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern
 496 Recognition (CVPR)*, pages 6299–6308, 2017. [10](#)
- 497 [20] Jinxing Zhou, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birchfield, Dan Guo,
 498 Lingpeng Kong, Meng Wang, and Yiran Zhong. Audio-visual segmentation. In *Proceedings of
 499 European Conference on Computer Vision (ECCV)*, 2022. [10](#), [11](#)
- 500 [21] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio
 501 Torralba. The sound of pixels. In *Proceedings of the European Conference on Computer Vision
 502 (ECCV)*, pages 570–586, 2018. [10](#), [11](#)
- 503 [22] Shentong Mo and Pedro Morgado. A unified audio-visual learning framework for localization,
 504 separation, and recognition. In *Proceedings of the International Conference on Machine
 505 Learning (ICML)*, 2023. [10](#), [11](#)
- 506 [23] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to
 507 localize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer
 508 Vision and Pattern Recognition (CVPR)*, pages 4358–4366, 2018. [10](#)
- 509 [24] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale
 510 audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics,
 511 Speech and Signal Processing (ICASSP)*, pages 721–725. IEEE, 2020. [11](#)
- 512 [25] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for
 513 bridging video and language. In *Proceedings of IEEE Conference on Computer Vision and
 514 Pattern Recognition (CVPR)*, pages 5288–5296, 2016. [11](#)
- 515 [26] Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji rong Wen, and Di Hu. Learning to
 516 answer questions in dynamic audio-visual scenarios. In *Proceedings of IEEE/CVF Conference
 517 on Computer Vision and Pattern Recognition (CVPR)*, pages 19086–19096, 2022. [11](#)