

Revisiting Audio-language Pretraining for Learning General-purpose Audio Representation

Anonymous ACL submission

Abstract

Audio-language pretraining (ALP) holds promise for learning general-purpose audio representation, yet remains underexplored. Crucially, there is no consensus on whether audio-language models can build effective general-purpose audio encoders, nor a systematic understanding of how pretraining objectives behave across diverse tasks and scales. We identify three key barriers: limited scale of audio-text corpora, insufficient caption diversity, and lack of systematic exploration and evaluation. To fill this gap, we present the first principled empirical study of ALP. We first introduce CaptionStew, a 10.7M caption dataset aggregating open-source audio-text corpora across multiple domains and captioning focuses. We then conduct the first comprehensive evaluation comparing contrastive and captioning objectives for learning audio representation across speech, music, and environmental sound tasks. Our results not only demonstrate that ALP yields competitive, transferable representations, but reveal critical trade-offs: contrastive learning offers superior data efficiency, while captioning exhibits better scalability. Furthermore, we find that supervised initialization provides diminishing returns at scale, challenging common practices. By grounding these claims in empirical evidence, we establish a viable pathway toward general-purpose audio representation learning, guiding future research.

1 Introduction

Representation learning has long been central to audio processing¹. Current approaches are predominated by supervised learning (Kong et al., 2020; Chen et al., 2022a; Desplanques et al., 2020) and self-supervised learning (Chen et al., 2023; Baevski et al., 2020; Hsu et al., 2021; Li et al., 2024), which

have consistently enhanced performance across various speech and audio benchmarks (Yang et al., 2021; Turian et al., 2022; Yuan et al., 2023). Despite these successes, most existing methods remain optimized for narrow task scopes rather than general-purpose use; models excelling at environmental sound classification, for example, often fail to capture speaker identity or paralinguistic attributes, and vice versa (Turian et al., 2022). Thus, learning audio representations that transfer robustly across diverse audio processing tasks remains an actively pursued and unresolved challenge.

A promising alternative is audio-language pretraining (ALP) (Elizalde et al., 2023; Wu et al., 2023), which grounds audio perception in natural language descriptions (captions). In this framework, text serves as a flexible semantic scaffold, enabling supervision spanning multiple levels of granularity, from coarse event categories (e.g., “dog barking,” “applause”) to fine-grained acoustic attributes (e.g., speaking style or musical structure), offering a unified path toward general audio understanding (Sakshi et al., 2025; Huang et al., 2025; Yang et al., 2024b; Su et al., 2025).

The success of vision-language pretraining underscores this promise. Models like CLIP (Radford et al., 2021) and AIM-v2 (Fini et al., 2025) not only power vision-language alignments but also produce representations that benefit various vision tasks (Liu et al., 2023; Minderer et al., 2022; Crowson et al., 2022). In contrast, audio-language models (ALMs) have not yet seen similar adoption. Existing models (Elizalde et al., 2023; Wu et al., 2023; Mei et al., 2024; Bai et al., 2025) remain largely confined to retrieval tasks, leaving the community without a systematic understanding of whether ALP can support general-purpose audio representation learning. Fundamental questions remain unanswered: how do different pretraining objectives behave and scale, and how does transfer performance vary across heterogeneous audio

¹In this work, audio processing refers to audio understanding, speech analysis and music understanding, while excluding automatic speech recognition

081 tasks such as speaker identification and audio event
082 classification? The absence of empirical evidence
083 regarding these questions has hindered progress
084 and led to uncertainty of design choices.

085 We identify three key challenges that have con-
086 strained progress. **First**, large-scale, web-mined
087 image–text corpora (Schuhmann et al., 2022; Gadre
088 et al., 2023) contain billions of pairs, but no compa-
089 rable resource exists for audio. Current audio cap-
090 tion datasets barely exceed one million pairs (Bai
091 et al., 2025; Mei et al., 2024; Kim et al., 2019;
092 Drossos et al., 2020), fundamentally limiting the
093 scaling potential of ALMs. **Second**, widely used
094 audio caption datasets focus predominantly on de-
095 scribing *what* is present in the audio, with limited
096 coverage of the rich range of acoustic attributes that
097 characterize different audio signals. For instance,
098 captions rarely characterize speaker characteristics,
099 musical attributes, or environmental acoustics. This
100 imbalanced focus limits the model’s ability to learn
101 representations that capture the full range of audio
102 semantics. **Third**, prior ALP works have primar-
103 ily focused on contrastive learning and audio–text
104 retrieval benchmarks. Systematic studies on al-
105 ternative pretraining objectives (captioning) and
106 comprehensive evaluations across a wide suite of
107 audio understanding tasks remain scarce, limiting
108 our understanding of what drives effective ALP.

109 In this work, we revisit ALP with the goal of re-
110 assessing its viability for learning general-purpose
111 audio representation. Rather than proposing a new
112 model architecture, we provide **a foundational em-
113 pirical study that fills the critical knowledge gap**
114 described above, establishing a rigorous baseline to
115 guide future research in accordance with scientific
116 best practices. We begin by aggregating diverse
117 open-source audio caption datasets into a unified re-
118 source, **CaptionStew**, enabling analysis at substan-
119 tially larger scales and with greater caption diver-
120 sity than prior work. Using this testbed, we conduct
121 the first comprehensive evaluation of ALP across
122 diverse downstream tasks and evaluation protocols,
123 showing that it yields competitive and transferable
124 representations across speech, music, and environ-
125 mental audio domains. Through a controlled com-
126 parison between contrastive and captioning objec-
127 tives, we reveal a consistent trade-off: contrastive
128 learning exhibits superior data efficiency, while cap-
129 tioning demonstrates better scalability. We further
130 analyze key training factors—data scaling and su-
131 pervised initialization—showing that not all tasks
132 benefit uniformly from increased data, and that

133 the gains from supervised initialization diminish
134 at larger scales and for tasks beyond audio event
135 classification, challenging common practices in the
136 field. Finally, we discuss how limited lexical di-
137 versity in existing caption datasets might constrain
138 performance scaling on certain attributes, suggest-
139 ing potential directions for improvement.

140 Taken together, our study reveals actionable in-
141 sights that were previously undocumented for au-
142 dio community and occasionally contradict trends
143 from other modalities. They establish ALP as a
144 practical and competitive approach for learning
145 general-purpose audio representations and high-
146 light key factors for future progress. To facilitate
147 further research, we will release data, training and
148 evaluation code, and pretrained models.

149 2 Related Work

150 **Audio Representation Learning.** Supervised
151 models trained on labeled datasets have been fun-
152 damental to the field, including audio event classi-
153 fiers (Kong et al., 2020; Gong et al., 2021; Chen
154 et al., 2022a; Dinkel et al., 2024), speech recog-
155 nition systems (Radford et al., 2023) and speaker
156 recognition models (Snyder et al., 2018; Desplan-
157 ques et al., 2020). These approaches remain widely
158 adopted due to their strong performance on spe-
159 cific target tasks. Self-supervised learning methods
160 have also emerged, demonstrating benefits across
161 speech (Baevski et al., 2020; Hsu et al., 2021; Chen
162 et al., 2022b), audio (Huang et al., 2022; Chen et al.,
163 2023; Li and Li, 2022), and music (Li et al., 2024;
164 Zhu et al., 2025) without requiring labeled data.

165 **Audio–Language Pretraining.** ALP has emerged
166 as a promising approach for learning cross-modal
167 representations. Most existing work focuses on
168 contrastive objectives (Elizalde et al., 2023; Wu
169 et al., 2023, 2022), with recent extensions exploring
170 combinations with other objectives (Xu et al., 2023;
171 Zhu et al., 2024; Niizumi et al., 2025). The field
172 has also witnessed evolution in datasets, transition-
173 ing from human-annotated ones (Kim et al., 2019;
174 Drossos et al., 2020; Agostinelli et al., 2023) to
175 recent LLM-augmented ones (Mei et al., 2024; Bai
176 et al., 2025; Chen et al., 2025; Sun et al.), alongside
177 domain-specific resources covering speaker char-
178 acteristics (Diwan et al., 2025) and fine-grained
179 musical attributes (Roy et al., 2025).

180 **Universal Audio Understanding.** The evalua-
181 tion of audio understanding has evolved from task-
182 specific benchmarks (Yang et al., 2021; Turian

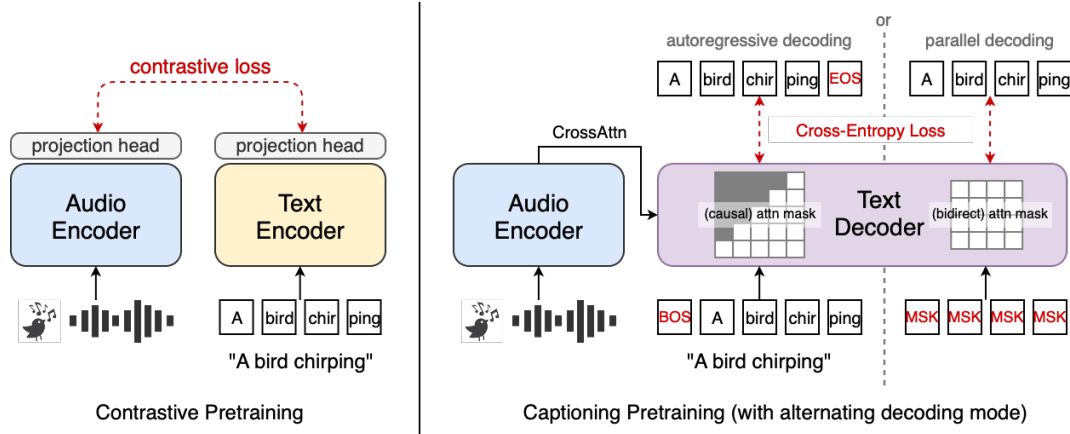


Figure 1: Audio-language pretraining objective studied in this work: contrastive and captioning.

et al., 2022; Yuan et al., 2023) toward more complex evaluation framework. Recent developments have emphasized LLM-based audio understanding systems (Ghosh et al., 2024; Gong et al., 2024; Dinkel et al., 2025; Goel et al., 2025; Chu et al., 2024; Tang et al., 2024) that can handle natural language queries and complex reasoning tasks. This shift has driven the development of corresponding evaluation benchmarks that assess models’ abilities across diverse audio understanding scenarios (Sakshi et al., 2025; Yang et al., 2024b; Huang et al., 2025; Ma et al., 2025). Our work contributes to this trend by providing the first comprehensive evaluation of ALP across discriminative tasks, audio-language alignment, and open-form question answering, bridging the gap between representation learning and universal audio understanding.

3 Audio-language Pretraining

ALP learns audio representations by establishing correspondence between audio signals and captions. The core concept is to leverage text as structured semantic supervision, enabling models to capture diverse information across speech, music, and environmental sounds within a unified framework. ALMs typically employ a two-tower architecture: an audio encoder f_a that maps raw audio signals into contextual representations, and a text component f_t whose design depends on the training objective. As shown in Figure 1, we explore two complementary paradigms that differ fundamentally in how they establish audio-text correspondence: contrastive and captioning objective. These approaches represent discriminative and generative perspectives of ALP, respectively.

3.1 Contrastive Objective

Contrastive objective is proven to be a robust representation learning method (Chen et al., 2020b; Radford et al., 2021; Baevski et al., 2020) and have been a dominant approach for ALP (Elizalde et al., 2023; Wu et al., 2023, 2022). This approach aligns audio and text representations in a shared embedding space by maximizing similarity between paired samples while minimizing similarity between mismatched pairs. Given a batch of paired samples $\{(a_i, t_i)\}_{i=1}^N$, the audio encoder produces frame- (or patch-) level representations that are pooled and projected to audio embeddings \mathbf{z}_i^a , while the text encoder f_t generates corresponding text embeddings \mathbf{z}_i^t . The symmetric InfoNCE loss (Oord et al., 2018) is applied to optimize both modalities:

$$\mathcal{L}_{\text{con}} = \frac{-1}{2N} \sum_{i=1}^N \left[\log \frac{\exp(\text{sim}(\mathbf{z}_i^a, \mathbf{z}_i^t)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{z}_i^a, \mathbf{z}_j^a)/\tau)} + \log \frac{\exp(\text{sim}(\mathbf{z}_i^t, \mathbf{z}_i^a)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathbf{z}_i^t, \mathbf{z}_j^a)/\tau)} \right], \quad (1)$$

where $\text{sim}(\cdot, \cdot)$ denotes cosine similarity and τ is a learnable temperature parameter. This objective encourages paired audio-text samples to be close in embedding space, encouraging semantic organization where similar content is grouped together.

3.2 Captioning Objective

Captioning objective takes a generative approach to audio-language alignment, learning representations by generating textual descriptions from audio. We argue that captioning presents a promising yet underexplored alternative for ALP. Theoretically, the cross-attention mechanism provides

frame-level supervision on the audio representation, offering denser learning signals than the utterance-level alignment used in contrastive learning. Also, since captioning models the joint distribution over all caption tokens, it is inherently more sensitive to fine-grained attributes, relations, and word order, enabling richer relational grounding (Yuksekgonul et al., 2023; Hsieh et al., 2023; Tschannen et al., 2023). Moreover, caption-based supervision is increasingly relevant given recent efforts toward general audio understanding systems (Dinkel et al., 2025; Goel et al., 2025)

Given an audio signal a_i , the encoder f_a produces contextual representations \mathbf{Z}_i^a , which are fed into a transformer decoder g_t through cross-attention. Inspired by CapPa (Tschannen et al., 2023), we alternate between two decoding modes—autoregressive and parallel prediction—to enhance audio encoder representation learning. In the autoregressive decoding, the decoder generates caption tokens (y_1, \dots, y_T) sequentially, with each token conditioned on the audio representation and previously generated tokens. Training follows the teacher-forcing approach with a cross-entropy loss:

$$\mathcal{L}_{\text{cap}} = - \sum_{t=1}^T \log p_{\theta}(y_t | y_{<t}, \mathbf{Z}_i^a), \quad (2)$$

In parallel prediction, we replace the decoder input tokens with [MASK] tokens and remove the causal attention mask, forcing simultaneous prediction of all tokens based solely on audio features:

$$\mathcal{L}_{\text{par}} = - \sum_{t=1}^T \log p_{\theta}(y_t | \mathbf{Z}_i^a), \quad (3)$$

This mode eliminates reliance on prior autoregressive context and forces each token prediction to depend solely on the audio representation, thereby strengthening encoder supervision. In a preliminary experiment, we observe that incorporating the parallel mode yields stronger representations than using a purely autoregressive decoder. We adopt mixed training where a random fraction of each minibatch uses standard autoregression while the remainder use parallel decoding.

4 CaptionStew Dataset

To investigate the potential of audio–language pre-training for general-purpose representation learning, we collect a large-scale and diverse audio caption dataset that addresses key limitations in existing corpora. Audio signals inherently encode

information across multiple dimensions—timbre, pitch, rhythm, semantic events, emotional tone, and acoustic environment—each amenable to different linguistic descriptions. However, existing large-scale audio caption datasets typically rely on a single caption-generation pipeline (Appendix A.3), where all captions are produced through the same procedure—either human annotation following uniform guidelines or LLM-based synthesis—and consequently share a homogeneous linguistic style. This uniformity offers consistency and scalability but introduces systematic stylistic biases and restricts linguistic diversity. Moreover, single-pipeline captions tend to exhibit limited syntactic variation and a narrow descriptive focus on only a subset of audio characteristics, often overlooking complementary acoustic attributes.

To fully leverage text as a flexible semantic scaffold for diverse audio representation learning, we embrace caption diversity across sources, styles, and descriptive granularities. Rather than creating captions through a single pipeline, we aggregate existing open-source corpora (Kim et al., 2019; Drossos et al., 2020; Agostinelli et al., 2023; Mei et al., 2024; Chen et al., 2025; Bai et al., 2025; Diwan et al., 2025; Roy et al., 2025). These datasets span multiple audio domains—general sound events, expressive speech, and musical performance—and employ fundamentally different caption creation methodologies. This aggregation yields captions that describe complementary audio aspects with varying granularity, from coarse event categories to fine-grained acoustic attributes. Please refer to Appendix A.3 for detail and examples of each source dataset. The resulting dataset, **CaptionStew** (denoted by CS10M), contains 9.3 million audio samples paired with 10.7 million captions, spanning 37,290 hours across speech, music, and environmental domains. Compared to existing collections, CaptionStew achieves both greater scale and broader coverage. This not only facilitates the learning of general-purpose audio representations but also provides a standardized, reproducible testbed for rigorous empirical study. Table 1 presents a comparison with existing audio caption datasets.

5 Experimental Setup

5.1 Implementation Details

We pretrain all models on CaptionStew. The audio encoder uses a Zipformer-M architecture (Yao

Table 1: Comparison of publicly available audio caption datasets. The number of audio-text pairs (#pair) and number of unique words (#vocab) are shown here.

Audio Caption Dataset	#pair	#vocab
<i>Human-annotated</i>		
AudioCaps (Kim et al., 2019)	46K	4,844
Clotho (Drossos et al., 2020)	5K	4,366
MusicCaps (Agostinelli et al., 2023)	5K	3,730
<i>LLM-augmented</i>		
WavCaps (Mei et al., 2024)	403K	18,372
AudioSetCaps (Bai et al., 2025)	1.9M	21,783
FusionAudio (Chen et al., 2025)	1.2M	18,403
AutoACD (Sun et al.)	1.5M	20,491
CaptionStew (Ours)	10.7M	56,586

Table 2: Datasets used for evaluating linear probing, audio-language task and open-form question answering performance (separated by lines). All metrics are higher the better. †reported with AIR-Bench (Yang et al., 2024b).

Evaluation Dataset	Task	Metrics
FSD-50k	Multi-label audio event classification	mAP
VggSound	Single-label audio event classification	accuracy
VoxCeleb2	Speaker identification	accuracy
CREMA	Speech emotion recognition	accuracy
MagnaTagATune	Music tagging	mAP
NSynth	Musical instrument classification	accuracy
AS-strong	Sound event detection	PSDS1
AudioCaps	Text-to-audio retrieval	Recall@1
ParaSpeechCaps	Audio captioning	RougeL
MusicCaps		
ClothoAQA	Open-formed question answering	Score†
ParaLMQA		
MusicQA		

et al., 2024), chosen for its efficiency on long sequences and fast convergence. For contrastive pre-training, the text encoder follows BERT-base architecture (Devlin et al., 2019). For captioning pretraining, the text decoder adopts the BART-base decoder architecture (Lewis et al., 2020). We use twice as many encoder layers (12) as decoder layers (6) to ensure comparable training speed across objectives. We experiment with two scenarios: training from scratch (*-scratch*) or initialized from pre-trained checkpoints (*-init*), following prior works in ALP (Wu et al., 2023; Mei et al., 2024; Bai et al., 2025). Please refer to Appendix A.1 for the full implementation details.

5.2 Evaluation Protocols and Datasets

We evaluate pretrained audio encoders across three protocols assessing discriminative capabilities, audio-language alignment, and open-formed question answering. All experiments probe frozen representations from the audio encoder’s final layer to ensure fair comparison. Table 2 and Appendix A.4 details the datasets and task metrics.

Linear Probing trains simple linear classifier on frozen representations. We evaluate across a diverse set of tasks across audio domains, including audio event classification (AEC) (Fonseca et al., 2021; Chen et al., 2020a), sound event detection (SED) (Hershey et al., 2021), speaker identification (SID) (Chung et al., 2018), speech emotion recognition (SER) (Cao et al., 2014), music tagging (MTAG) (Law et al., 2010) and musical instrument classification (INST) (Engel et al., 2017).

Audio-language Alignments follow the LiT protocol (Zhai et al., 2022), adapting either pre-trained text encoder (Liu et al., 2019) or text de-

coder (Lewis et al., 2020) to align with frozen audio representations for performing retrieval and captioning tasks. We evaluate on audio-caption datasets spanning diverse domains: AudioCaps (AC) (Kim et al., 2019) for general sound event descriptions; ParaSpeechCaps (PSC) (Diwan et al., 2025) for speaking-style and acoustic-environment descriptions; and MusicCaps (MC) (Agostinelli et al., 2023) for fine-grained musical descriptions. **Open-formed Question Answering.** Acknowledging the trend of combining audio encoders with large language models (LLMs) for general audio understanding (Ghosh et al., 2024; Gong et al., 2024), we connects frozen audio encoders to a LLM (Qwen2.5-7B-Instruct Yang et al. (2024a)) through lightweight adaptors. We train only the adaptor on multiple audio QA datasets that span distinct domains: sound event understanding (Liping et al., 2022), speaker-related and paralinguistic understanding (Huo et al., 2025), and music understanding (Liu et al., 2024). Evaluation is conducted on the corresponding tracks (sound, speaker-related, music; see Appendix A.4) of AIR-Bench (Yang et al., 2024b).

5.3 Baseline Methods

Recognizing the broad adoption of pretrained audio event classifiers in transfer learning (Alonso-Jiménez et al., 2023; Cappellazzo et al., 2024), audio-language modeling (Elizalde et al., 2023; Wu et al., 2023) and general audio understanding (Gong et al., 2024; Ghosh et al., 2024; Dinkel et al., 2025), we select our pretrained Zipformer-based audio event classifier (denoted by Zipformer-AEC, described in Appendix A.1) as the primary baseline. We also compare against rep-

Table 3: Evaluation results across tasks and protocols. [†] numbers quoted from other papers with consistent evaluation setup. [‡]state-of-the-art results on each task without any training constraints (e.g. full-finetuning) (see Appendix A.5). ^{††}no available prior work. ^{‡‡}results of speaker emotion recognition, gender recognition, and age prediction in AIR-Bench (Yang et al., 2024b), respectively.

(a) Linear Probing (with mean pooling)									
Method	Model Initialization	Audio-lang. Pretraining	linear probing						
			AEC FSD50k	AEC VggSound	SID VoxCeleb2	SER CREMA	MTAG MagnaTagATune	INST NSynth	SED AS-Strong
<i>Existing SSL Models</i>									
BEATs (Chen et al., 2023)	SSL	–	0.565 [†]	–	–	–	0.400 [†]	75.90 [†]	0.034 [†]
Wav2vec 2.0 (Baevski et al., 2020)	SSL	–	0.342 [†]	–	<u>51.60</u>	56.10	0.317 [†]	40.20 [†]	–
MERT (Li et al., 2024)	SSL	–	–	–	–	–	0.402 [†]	72.60 [†]	–
<i>Our Supervised Baselines</i>									
Zipformer-AEC (Yao et al., 2024)	AudioSet SL	–	0.656	<u>56.46</u>	18.84	67.14	0.407	67.19	<u>0.216</u>
<i>Our Audio-lang. Pretrained</i>									
Contrastive- <i>scratch</i>	–	CS10M	0.625	50.87	46.67	67.71	0.406	67.30	0.132
Captioning- <i>scratch</i>	–	CS10M	0.580	47.79	33.43	63.60	0.401	63.10	0.124
Contrastive- <i>init</i>	AudioSet SL	CS10M	0.664	54.70	38.17	68.84	0.406	69.38	0.187
Captioning- <i>init</i>	AudioSet SL	CS10M	0.652	53.13	26.23	65.86	<u>0.410</u>	67.16	0.145
SOTA [‡]			0.655	59.50	96.20	– ^{††}	0.414	79.20	0.374

(b) Audio-language Alignment / Open-form QA											
Method	Captioning			Retrieval			Open-formed QA				
	AC	PSC	MC	AC	PSC	MC	Sound	Speaker-related ^{‡‡}	Music		
<i>Our Supervised Baselines</i>											
Zipformer-AEC (Yao et al., 2024)	46.7	45.5	<u>22.9</u>	40.5	49.2	24.6	7.01	36.5 / 46.2 / 37.2	5.61		
<i>Our Audio-lang. Pretrained</i>											
Contrastive- <i>scratch</i>	46.6	46.3	22.1	39.3	63.2	27.4	6.65	37.9 / 81.3 / 63.4	5.86		
Captioning- <i>scratch</i>	46.7	46.5	22.9	36.9	60.2	23.0	6.69	44.2 / 65.4 / 69.0	5.97		
Contrastive- <i>init</i>	47.2	46.2	22.5	42.8	60.6	29.4	6.73	35.1 / 67.3 / 64.5	5.63		
Captioning- <i>init</i>	47.2	45.9	22.6	42.2	55	28.2	7.06	32.4 / 49.5 / 45.6	5.50		
SOTA [‡]			52.2	– ^{††}	26.2	44.4	– ^{††}	– ^{††}	6.99	60.0 / 82.5 / 62.4	6.79

representative self-supervised learning (SSL) models, each specialized for particular audio domains: BEATs (Chen et al., 2023) for environmental sound (or general audio); Wav2vec 2.0 (Baevski et al., 2020) for speech signal; and MERT for music pieces. Together, these baselines provide a broad comparative context for studying ALP toward general-purpose audio representation.

6 Experiment Results

We present our evaluation results in Table 3. Our analysis reveals key insights about objective design, representation quality, and the role of initialization. **Contrastive vs. Captioning Objectives.** The two pretraining paradigms exhibit complementary strengths across evaluation protocols. On linear probing tasks, contrastive learning consistently outperforms captioning, particularly excelling at audio event classification and speaker identification. However, it is worth noting that this gap narrows substantially when the classifier learns

to aggregate information across frames through multi-head attention pooling (Appendix A.5). This observation reflects the objectives’ inherent designs: contrastive learning explicitly optimizes for linearly separable clip-level representations, while captioning relies on cross-attention mechanisms over frame-level representations for text sequence generation. This finding aligns with recent work highlighting how downstream module choices significantly impact the assessment of audio representation quality (Zaiem et al., 2023). For language-involved tasks, both objectives demonstrate competitive performance, with captioning showing slight advantages in open-form question answering across multiple domains. This suggests captioning’s potential for language-involved audio understanding tasks.

Impact of Supervised Initialization. Initializing from supervised pretraining (AS SL) provides substantial benefits across most tasks, with notable improvements on audio event classification, sound event detection and audio-text retrieval. The gains

are particularly pronounced for contrastive objectives, suggesting that supervised pretraining provides useful inductive biases for contrastive learning. However, these benefits diminish (or disappear entirely) when the attributes required for downstream tasks diverge from AudioSet’s ontology. On speaker identification and music tagging, scratch-trained models often match or exceed initialized variants, indicating that AudioSet’s focus on distinguishing between sound categories may bias representations toward event-level semantics rather than the speaker characteristics (voice timbre, speaking style) or musical structure (chord, rhythm) essential for these tasks. These findings challenge common initialization practices for ALP and suggest the need for tailored pretraining strategies when targeting general-purpose audio representation learning.

Competitive Performance Across Domains. Our audio-language representations achieve strong transferability across diverse audio domains. Compared to supervised baselines (Zipformer-AEC), our overall best-performing model (Contrastive-init) demonstrates superior performance on speaker identification, music understanding and audio-text retrieval while maintaining competitiveness on audio-event classification. Against domain-specialized SSL methods (BEATs, Wav2vec 2.0, MERT), our approach consistently shows competitive performance. This cross-domain performance validates our hypothesis that diverse caption aggregation enables broadly transferable representations, establishing ALP as a viable path toward learning general-purpose audio representation.

6.1 Data-Scaling Experiments

To understand the scalability of audio–language pretraining, we conduct controlled experiments using CaptionStew subsets at 400K, 1M, 4M, and 10M (whole corpus) audio-text pairs. Figure 2 reveals distinct scaling patterns across objectives and evaluation protocols.

Scaling Patterns. Most tasks demonstrate consistent performance improvements with increased data scale, validating the potential of large-scale ALP. However, notable exceptions emerge that reveal fundamental limitations of current approaches. Sound event detection, particularly for models initialized with AudioSet pretraining, exhibits a reverse scaling trend where performance degrades with more caption data. This suggests a potential conflict between natural language supervision—which typically describes audio characteristics and

attributes—and temporal localization tasks requiring precise event boundaries. Additionally, emotion recognition and instrument classification show weaker scaling gains compared to other tasks, likely reflecting limited caption diversity for these specific attributes in existing corpora, which we will discuss in Sec. 6.2.

Contrastive vs. Captioning Scaling. Contrastive learning consistently outperforms captioning at varying data scales, particularly under less data and on discriminative tasks such as audio event classification. However, captioning demonstrates slightly better scaling properties, with distinct patterns emerging across task categories. or language-involved tasks—especially captioning and question answering—captioning matches or surpasses contrastive learning at our current 10M-pair scale. On linear probing benchmarks, the gap remains substantial, with scaling trends suggesting captioning would require hundreds of millions of pairs to achieve parity with contrastive methods.

Impact of Initialization at Scale. AudioSet initialization provides immediate performance gains but introduces diminishing returns at larger scales. Both contrastive learning and captioning show decreasing benefits from initialization as data scale increases, with scratch and initialized models achieving matched performance at larger scales on some tasks. This suggests that pretrained initialization effectively bootstraps learning at small scales but may constrain the model’s ability to adapt to the broader semantic space covered by large-scale caption data, potentially due to mismatch between AudioSet’s ontology and diverse audio descriptions.

Overall, these findings reveal complementary behaviors: contrastive pretraining achieves superior data efficiency at current scales, while captioning shows better scalability, especially for language-involved tasks. Importantly, the diminishing returns of initialization at scale indicate that large-scale caption data can provide sufficient semantic supervision independent of domain-specific pretraining, challenging current practices of ALP and opening possibilities for learning general-purpose representations from diverse captions alone.

6.2 Dataset Analysis

To understand the linguistic characteristics of CaptionStew, we analyze caption diversity across constituent datasets through visualization and quantitative methods. Figure 4 provides compelling evidence of our aggregation strategy’s success through

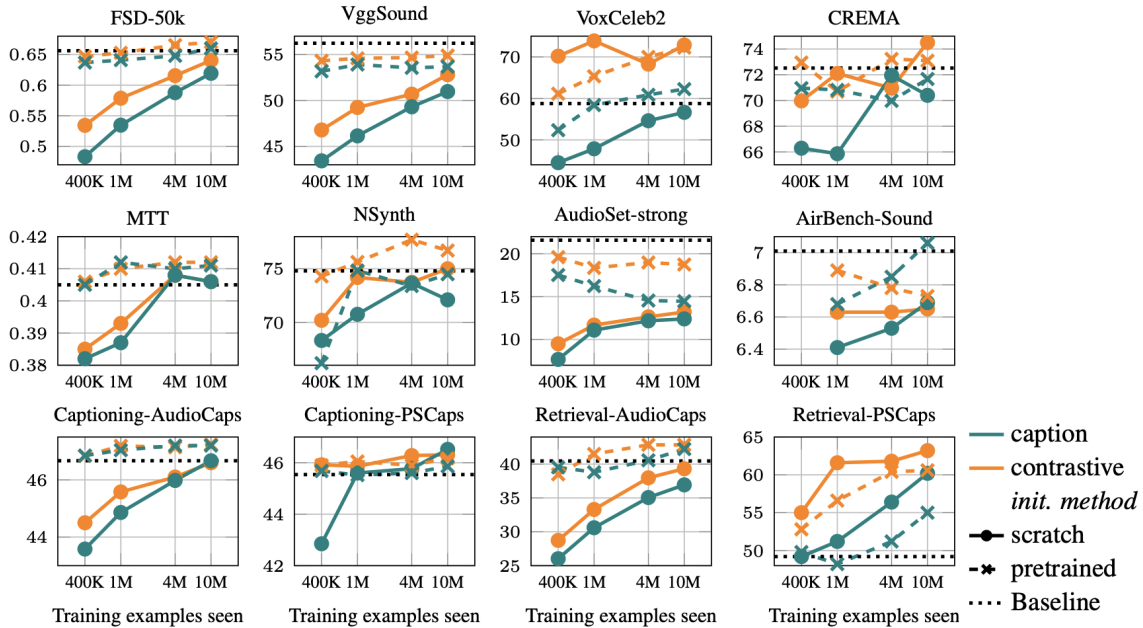


Figure 2: Data scaling behavior of contrastive vs. captioning objectives across representative tasks.

t-SNE visualization (Maaten and Hinton, 2008) of sentence embeddings (Reimers and Gurevych, 2019) from sampled captions, revealing distinct clustering patterns by source that demonstrate complementary linguistic perspectives: AudioSetCaps and WavCaps overlap in audio event descriptions and aligns with human annotated dataset, while JamendoMaxCaps creates a distinct cluster focused on music-specific terminology, and ParaSpeechCaps forms a separate cluster emphasizing speaking styles and paralinguistic attributes. These minimal overlaps confirm that each dataset contributes distinct caption styles and descriptive focuses.

Quantitative analysis reveals both the benefits and limitations (Table 10). CaptionStew achieves substantial vocabulary expansion (56,586 unique words vs. 4,060-27,906 for individual datasets) However, this growth doesn't yield proportional lexical diversity. The Distinct-n (Li et al., 2015) of CaptionStew remain low, falling short of image caption dataset (Changpinyo et al., 2021) and text corpora (Merity et al., 2016). This constraint stems from datasets with limited linguistic variation, particularly JamendoMaxCaps and ParaSpeechCaps with extremely low Distinct-n scores.

These findings highlight that simply combining datasets doesn't guarantee improved linguistic diversity, revealing broader limitations in current ALP approaches. Also, the constrained diversity in certain aspect may partially explain weaker scaling behavior observed for certain tasks, as models

encounter repetitive linguistic patterns despite increased data volume, aligning with vision-language findings on caption diversity's importance for representation quality (Santurkar et al., 2023; Chan et al., 2022). This analysis motivates developing enhanced aggregation pipeline and more diverse caption generation methods to better capture the full spectrum of information in audio signals, thereby fully realizing the potential of large-scale ALP.

7 Conclusions

We revisited audio–language pretraining with the goal of establishing a rigorous baseline for general-purpose audio representation learning. By aggregating and harmonizing diverse datasets into CaptionStew, we addressed the data scarcity issues that have hindered the field and enabled a rigorous comparison of training objectives and data scales. Our comprehensive evaluation yielded several actionable insights: (1) audio–language pretraining produces competitive representations across speech, music, and environmental sounds; (2) contrastive and captioning objectives exhibit complementary strengths regarding efficiency and scalability; and (3) standard supervised initializations may be unnecessary or even detrimental at scale. Finally, our analysis highlighted the limited lexical diversity in current caption datasets as a key frontier for future improvement. We hope these empirical foundations will accelerate the development of future general-purpose audio representation learning.

618 Limitation

619 While this work provides valuable empirical in- 669
620 sights for audio-language pretraining, we acknowl- 670
621 edge several important limitations that present op- 671
622 portunities for future research. 672

623 **Dataset Construction and Quality.** CaptionStew 673
624 aggregates captions from multiple sources with 674
625 varying generation methodologies, including LLM- 675
626 synthesized descriptions that may introduce sys- 676
627 tematic biases or artifacts. We do not perform 677
628 extensive quality control or human verification 678
629 across the aggregated corpus, which could impact 679
630 model training. Additionally, our dataset analysis
631 reveals that simple aggregation does not guaran-
632 tee improved linguistic diversity—CaptionStew’s
633 lexical diversity metrics remain lower than mature
634 image-text corpora. However, our design choice
635 prioritizes semantic diversity over linguistic vari-
636 ety, as evidenced by the t-SNE clustering analysis
637 showing distinct descriptive focuses across con-
638 stituent datasets. While more sophisticated cura-
639 tion strategies could improve quality, our goal was
640 to establish whether diverse caption aggregation
641 can benefit audio representation learning, which
642 our results support despite these limitations.

643 **Limited Technical Novelty.** Our work primarily 681
644 combines existing techniques—contrastive learn- 682
645 ing, captioning objectives, and dataset aggrega- 683
646 tion—rather than introducing fundamentally new 684
647 methods. The mixed autoregressive/parallel train- 685
648 ing approach is adapted from vision-language 686
649 work (CapPa), and our architectural choices fol- 687
650 low standard practices. We acknowledge that 688
651 the technical contributions are largely empirical 689
652 rather than methodological. However, this aligns 690
653 with our primary goal of systematically evaluating 691
654 audio-language pretraining’s potential for general- 692
655 purpose representation learning. The field currently 693
656 lacks comprehensive comparative studies across ob- 694
657 jectives, evaluation protocols, and training factors. 695
658 Our systematic analysis reveals important insights 696
659 about scaling behaviors and initialization effects 697
660 that have practical implications for practitioners, 698
661 even if the underlying techniques are not novel. 699

662 **Limited Model and Data Scalability.** Our experi- 700
663 ments are constrained to 10M audio-text pairs and 701
664 relatively modest model sizes compared to state- 702
665 of-the-art vision-language systems that leverage 703
666 billions of samples and much larger architectures. 704
667 This scale limitation may not fully reflect the po- 705
668 tential of audio-language pretraining, particularly

670 for the captioning objective which our results sug- 671
672 gest benefits from larger-scale training. Addition- 673
674 ally, we do not explore recent advances in large 675
676 language model integration or more sophisticated 676
677 architectural designs that could improve perfor- 677
678 mance. These constraints stem from computational 678
679 resource limitations and our focus on controlled 679
680 comparisons rather than pushing absolute perfor-
681 mance boundaries. Future work with larger scales
682 may reveal different scaling dynamics and stronger
683 evidence for general-purpose capabilities. 684

680 Ethical Considerations

681 This work investigates audio–language pretraining 681
682 (ALP) as a framework for learning general-purpose 682
683 audio representations through large-scale empirical 683
684 analysis. While our study is methodological in na- 684
685 ture and does not directly deploy end-user systems, 685
686 it raises several ethical considerations related to 686
687 data sources and potential misuse. 687

688 **Data sourcing and privacy.** CaptionStew is con- 688
689 structed by aggregating existing open-source au- 689
690 dio–text datasets. These datasets are collected 690
691 under diverse licenses and data collection prac- 691
692 tices, and may include audio containing human 692
693 speech or environmental recordings. We rely on 693
694 the original dataset providers’ compliance with con- 694
695 sent, anonymization, and licensing requirements, 695
696 and we do not introduce new data collection or 696
697 re-identification procedures. Nevertheless, large- 697
698 scale aggregation may amplify latent biases or ar- 698
699 tifacts present in individual sources, including de- 699
700 mographic imbalance, recording context skew, or 700
701 stylistic biases introduced by caption generation 701
702 process. Users of the dataset and pretrained models 702
703 should be aware of these limitations. 703

704 **Potential misuse.** General-purpose audio repre- 704
705 sentations can benefit applications such as acces- 705
706 sibility tools, audio search, and audio understand- 706
707 ing, but they may also lower barriers to harmful 707
708 uses, including surveillance, profiling, or misuse 708
709 of speaker-related attributes. Our released models 709
710 are intended for research purposes, and we do not 710
711 claim suitability for high-stakes or safety-critical 711
712 scenario without safeguards and validation. 712

713 In summary, this work aims to clarify the ca- 713
714 pabilities and limitations of audio–language pre- 714
715 training through transparent empirical study. We 715
716 hope to support more responsible development and 716
717 evaluation of general-purpose audio representation 717
718 learning. 718

References

Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. 2016. Youtube8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*.

Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, and 1 others. 2023. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*.

Pablo Alonso-Jiménez, Xavier Serra, and Dmitry Bogdanov. 2023. Efficient supervised training of audio transformers for music representation learning. In *Ismir 2023 Hybrid Conference*.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670*.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Jisheng Bai, Haohe Liu, Mou Wang, Dongyuan Shi, Wenwu Wang, Mark D Plumbley, Woon-Seng Gan, and Jianfeng Chen. 2025. Audiosetcaps: An enriched audio-caption dataset using automated generation pipeline with large audio and language models. *IEEE Transactions on Audio, Speech and Language Processing*.

Shikhar Bharadwaj, Samuele Cornell, Kwanghee Choi, Satoru Fukayama, Hye-jin Shim, Soham Deshmukh, and Shinji Watanabe. 2025. Openbeats: A fully open-source general-purpose audio encoder. *arXiv preprint arXiv:2507.14129*.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335–359.

Houwei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini Verma. 2014. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE transactions on affective computing*, 5(4):377–390.

Umberto Cappellazzo, Daniele Falavigna, Alessio Brutti, and Mirco Ravanelli. 2024. Parameter-efficient transfer learning of audio spectrogram transformers. In *2024 IEEE 34th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE.

David M Chan, Austin Myers, Sudheendra Vijayanarasimhan, David A Ross, Bryan Seybold, and John F Canny. 2022. What’s in a caption? dataset-specific linguistic diversity and its effect on visual description models and metrics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4740–4749. 775
776
777
778
779
780
781

Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*. 782
783
784
785

Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020a. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 721–725. IEEE. 786
787
788
789
790

Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2022a. Htsat: A hierarchical token-semantic audio transformer for sound classification and detection. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 646–650. IEEE. 791
792
793
794
795
796
797

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, and 1 others. 2022b. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518. 798
799
800
801
802
803
804

Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, Wanxiang Che, Xi-angzhan Yu, and Furu Wei. 2023. Beats: Audio pre-training with acoustic tokenizers. In *International Conference on Machine Learning*, pages 5178–5193. PMLR. 805
806
807
808
809
810

Shunian Chen, Xinyuan Xie, Zheshu Chen, Liyan Zhao, Owen Lee, Zhan Su, Qilin Sun, and Benyou Wang. 2025. Fusionaudio-1.2 m: Towards fine-grained audio captioning with multimodal contextual fusion. *arXiv preprint arXiv:2506.01111*. 811
812
813
814
815

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020b. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmlR. 816
817
818
819
820

Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, and 1 others. 2024. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*. 821
822
823
824

Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. 2018. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*. 825
826
827

Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. 2022. Vqgan-clip: Open domain 828
829
830

831	image generation and editing with natural language	of large vision encoders. In <i>Proceedings of the Com-</i>	887
832	guidance. In <i>European conference on computer vi-</i>	puter Vision and Pattern Recognition Conference,	888
833	sion, pages 88–105. Springer.	pages 9641–9654.	889
834	Brecht Desplanques, Jenthe Thienpondt, and Kris	Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic	890
835	Demuynck. 2020. Ecapa-tdnn: Emphasized	Font, and Xavier Serra. 2021. Fsd50k: an open	891
836	channel attention, propagation and aggregation in	dataset of human-labeled sound events. <i>IEEE/ACM</i>	892
837	tdnn based speaker verification. <i>arXiv preprint</i>	<i>Transactions on Audio, Speech, and Language Pro-</i>	893
838	<i>arXiv:2005.07143</i> .	<i>cessing</i> , 30:829–852.	894
839	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang,	895
840	Kristina Toutanova. 2019. Bert: Pre-training of deep	Jonathan Hayase, Georgios Smyrnis, Thao Nguyen,	896
841	bidirectional transformers for language understand-	Ryan Marten, Mitchell Wortsman, Dhruva Ghosh,	897
842	ing. In <i>Proceedings of the 2019 conference of the</i>	Jieyu Zhang, and 1 others. 2023. Datacomp: In	898
843	<i>North American chapter of the association for com-</i>	search of the next generation of multimodal datasets.	899
844	<i>putational linguistics: human language technologies,</i>	<i>Advances in Neural Information Processing Systems</i> ,	900
845	<i>volume 1 (long and short papers)</i> , pages 4171–4186.	36:27092–27112.	901
846	Heinrich Dinkel, Gang Li, Jizhong Liu, Jian	Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman,	902
847	Luan, Yadong Niu, Xingwei Sun, Tianzi Wang,	Aren Jansen, Wade Lawrence, R Channing Moore,	903
848	Qiyang Xiao, Junbo Zhang, and Jiahao Zhou.	Manoj Plakal, and Marvin Ritter. 2017. Audio set:	904
849	2025. Midashenglm: Efficient audio understand-	An ontology and human-labeled dataset for audio	905
850	ing with general audio captions. <i>arXiv preprint</i>	events. In <i>2017 IEEE international conference on</i>	906
851	<i>arXiv:2508.03983</i> .	<i>acoustics, speech and signal processing (ICASSP)</i> ,	907
852	Heinrich Dinkel, Zhiyong Yan, Yongqing Wang, Junbo	pages 776–780. IEEE.	908
853	Zhang, Yujun Wang, and Bin Wang. 2024. Scaling	Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Ki-	909
854	up masked audio encoder learning for general audio	ran Reddy Evuru, Utkarsh Tyagi, S Sakshi, Oriol	910
855	classification. In <i>Proc. Interspeech 2024</i> , pages 547–	Nieto, Ramani Duraiswami, and Dinesh Manocha.	911
856	551.	2024. Gama: A large audio-language model with ad-	912
857	Anuj Diwan, Zhisheng Zheng, David Harwath, and Eu-	vanced audio understanding and complex reasoning	913
858	nsul Choi. 2025. Scaling rich style-prompted text-to-	abilities. In <i>Proceedings of the 2024 Conference on</i>	914
859	speech datasets. <i>arXiv preprint arXiv:2503.04713</i> .	<i>Empirical Methods in Natural Language Processing</i> ,	915
860	Konstantinos Drossos, Samuel Lipping, and Tuomas	pages 6288–6313.	916
861	Virtanen. 2020. Clotho: An audio captioning dataset.	Arushi Goel, Sreyan Ghosh, Jaehyeon Kim, Sonal Ku-	917
862	In <i>ICASSP 2020-2020 IEEE International Confer-</i>	mar, Zhifeng Kong, Sang-gil Lee, Chao-Han Huck	918
863	<i>ence on Acoustics, Speech and Signal Processing</i>	Yang, Ramani Duraiswami, Dinesh Manocha, Rafael	919
864	<i>(ICASSP)</i> , pages 736–740. IEEE.	Valle, and 1 others. 2025. Audio flamingo 3: Advanc-	920
865	Janek Ebberts, Reinhold Haeb-Umbach, and Romain	ing audio intelligence with fully open large audio	921
866	Serizel. 2022. Threshold independent evaluation of	language models. <i>arXiv preprint arXiv:2507.08128</i> .	922
867	sound event detection scores. In <i>ICASSP 2022-2022</i>	Yuan Gong, Yu-An Chung, and James Glass. 2021. Ast:	923
868	<i>IEEE International Conference on Acoustics, Speech</i>	Audio spectrogram transformer. In <i>Proc. Interspeech</i>	924
869	<i>and Signal Processing (ICASSP)</i> , pages 1021–1025.	2021, pages 571–575.	925
870	IEEE.	Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid	926
871	Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Is-	Karlinsky, and James Glass. 2024. Listen, think, and	927
872	mail, and Huaming Wang. 2023. Clap learning	understand. In <i>International Conference on Learning</i>	928
873	audio concepts from natural language supervision.	<i>Representations</i> .	929
874	In <i>ICASSP 2023-2023 IEEE International Confer-</i>	Yuan Gong, Andrew Rouditchenko, Alexander H	930
875	<i>ence on Acoustics, Speech and Signal Processing</i>	Liu, David Harwath, Leonid Karlinsky, Hilde	931
876	<i>(ICASSP)</i> , pages 1–5. IEEE.	Kuehne, and James Glass. 2022. Contrastive	932
877	Jesse Engel, Cinjon Resnick, Adam Roberts, Sander	audio-visual masked autoencoder. <i>arXiv preprint</i>	933
878	Dieleman, Mohammad Norouzi, Douglas Eck, and	<i>arXiv:2210.07839</i> .	934
879	Karen Simonyan. 2017. Neural audio synthesis of	Haorui He, Zengqiang Shang, Chaoren Wang, Xuyuan	935
880	musical notes with wavenet autoencoders. In <i>In-</i>	Li, Yicheng Gu, Hua Hua, Liwei Liu, Chen Yang,	936
881	<i>ternational conference on machine learning</i> , pages	Jiaqi Li, Peiyang Shi, and 1 others. 2024. Emilia: An	937
882	1068–1077. PMLR.	extensive, multilingual, and diverse speech dataset for	938
883	Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter,	large-scale speech generation. In <i>2024 IEEE Spoken</i>	939
884	Michal Klein, David Haldimann, Sai Aitharaju, Vic-	<i>Language Technology Workshop (SLT)</i> , pages 885–	940
885	tor G Turrisi da Costa, Louis Béthune, Zhe Gan, and 1	890. IEEE.	941
886	others. 2025. Multimodal autoregressive pre-training		

942	Shawn Hershey, Daniel PW Ellis, Eduardo Fonseca, Aren Jansen, Caroline Liu, R Channing Moore, and Manoj Plakal. 2021. The benefit of temporally-strong labels in audio event classification. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 366–370. IEEE.	999
943		1000
944		1001
945		1002
946		
947		1003
948		1004
949	Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. 2023. Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 31096–31116. Curran Associates, Inc.	1005
950		1006
951		1007
952		1008
953		1009
954		1010
955	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. <i>IEEE/ACM transactions on audio, speech, and language processing</i> , 29:3451–3460.	1011
956		1012
957		1013
958		1014
959		
960		1015
961	Chien-yu Huang, Wei-Chih Chen, Shu-wen Yang, Andy T Liu, Chen-An Li, Yu-Xiang Lin, Wei-Cheng Tseng, Anuj Diwan, Yi-Jen Shih, Jiatong Shi, and 1 others. 2025. Dynamic-superb phase-2: A collaboratively expanding benchmark for measuring the capabilities of spoken language models with 180 tasks. In <i>The Thirteenth International Conference on Learning Representations</i> .	1016
962		1017
963		
964		1018
965		1019
966		1020
967		1021
968		1022
969	Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metzger, and Christoph Feichtenhofer. 2022. Masked autoencoders that listen. <i>Advances in Neural Information Processing Systems</i> , 35:28708–28720.	1023
970		1024
971		1025
972		1026
973		
974	Mingyue Huo, Wei-Cheng Tseng, Yiwen Shao, Hao Zhang, and Dong Yu. 2025. Auden-voice: General-purpose voice encoder for speech and language understanding. <i>arXiv preprint arXiv:2511.15145</i> .	1027
975		1028
976		1029
977		1030
978	Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. 2019. Audiocaps: Generating captions for audios in the wild. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 119–132.	1031
979		1032
980		1033
981		1034
982		1035
983		
984		1036
985	Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. 2020. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 28:2880–2894.	1037
986		1038
987		1039
988		1040
989		1041
990		1042
991	Luca A Lanzendörfer, Constantin Pinkl, Nathanaël Perraudin, and Roger Wattenhofer. 2025. Bootstrapping language-audio pre-training for music captioning. In <i>ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	1043
992		1044
993		1045
994		1046
995		1047
996		
997	Edith Law, Kris West, M Mandel, M Bay, and JS Downie. 2010. Evaluation of algorithms using	1048
998		1049
		1050
		1051
		1052
	games: the case of music annotation. In <i>Proceedings of the 11th International Society for Music Information Retrieval Conference (ISMIR)</i> . Utrecht, the Netherlands.	
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880.	
	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. <i>arXiv preprint arXiv:1510.03055</i> .	
	Xian Li and Xiaofei Li. 2022. Atst: Audio representation learning with teacher-student transformer. <i>arXiv preprint arXiv:2204.12076</i> .	
	Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, and 1 others. 2024. Mert: Acoustic music understanding model with large-scale self-supervised training. In <i>ICLR</i> .	
	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
	Samuel Lipping, Parthasaarathy Sudarsanam, Konstantinos Drossos, and Tuomas Virtanen. 2022. Clotho-aqa: A crowdsourced dataset for audio question answering. In <i>2022 30th European Signal Processing Conference (EUSIPCO)</i> , pages 366–370. IEEE.	
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36:34892–34916.	
	Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. 2024. Music understanding llama: Advancing text-to-music generation with question answering and captioning. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 286–290. IEEE.	
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	
	Steven R Livingstone and Frank A Russo. 2018. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. <i>PLoS one</i> , 13(5):e0196391.	

1053	Ziyang Ma, Yinghao Ma, Yanqiao Zhu, Chen Yang,	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	1109
1054	Yi-Wen Chao, Ruiyang Xu, Wenxi Chen, Yuanzhe	man, Christine McLeavey, and Ilya Sutskever. 2023.	1110
1055	Chen, Zhuo Chen, Jian Cong, and 1 others. 2025.	Robust speech recognition via large-scale weak super-	1111
1056	Mmar: A challenging benchmark for deep reasoning	vision. In <i>Proceedings of the 40th International Con-</i>	1112
1057	in speech, audio, music, and their mix. <i>arXiv preprint</i>	<i>ference on Machine Learning</i> , pages 28492–28518.	1113
1058	<i>arXiv:2505.13032</i> .		
1059	Laurens van der Maaten and Geoffrey Hinton. 2008.	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	1114
1060	Visualizing data using t-sne. <i>Journal of machine</i>	Sentence embeddings using siamese bert-networks.	1115
1061	<i>learning research</i> , 9(Nov):2579–2605.	<i>arXiv preprint arXiv:1908.10084</i> .	1116
1062	Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang	Julius Richter, Yi-Chiao Wu, Steven Krenn, Simon	1117
1063	Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley,	Welker, Bunlong Lay, Shinji Watanabe, Alexander	1118
1064	Yuexian Zou, and Wenwu Wang. 2024. Wavcaps:	Richard, and Timo Gerkmann. 2024. Ears: An	1119
1065	A chatgpt-assisted weakly-labelled audio captioning	anechoic fullband speech dataset benchmarked for	1120
1066	dataset for audio-language multimodal research.	speech enhancement and dereverberation. In <i>Proc.</i>	1121
1067	<i>IEEE/ACM Transactions on Audio, Speech, and Lan-</i>	<i>Interspeech 2024</i> , pages 4873–4877.	1122
1068	<i>guage Processing</i> , 32:3339–3354.		
1069	Stephen Merity, Caiming Xiong, James Bradbury, and	Abhinaba Roy, Renhang Liu, Tongyu Lu, and Dorien	1123
1070	Richard Socher. 2016. <i>Pointer sentinel mixture mod-</i>	Herremans. 2025. Jamendomaxcaps: A large scale	1124
1071	<i>els</i> . <i>Preprint</i> , arXiv:1609.07843.	music-caption dataset with imputed metadata. <i>arXiv</i>	1125
1072	Matthias Minderer, Alexey Gritsenko, Austin Stone,	<i>preprint arXiv:2502.07461</i> .	1126
1073	Maxim Neumann, Dirk Weissenborn, Alexey Doso-	S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth,	1127
1074	vitskiy, Aravindh Mahendran, Anurag Arnab,	Ramaneswaran Selvakumar, Oriol Nieto, Ramani Du-	1128
1075	Mostafa Dehghani, Zhuoran Shen, and 1 others. 2022.	raiswami, Sreyan Ghosh, and Dinesh Manocha. 2025.	1129
1076	Simple open-vocabulary object detection. In <i>Euro-</i>	Mmau: A massive multi-task audio understanding	1130
1077	<i>pean conference on computer vision</i> , pages 728–755.	and reasoning benchmark. In <i>The Thirteenth Interna-</i>	1131
1078	Springer.	<i>tional Conference on Learning Representations</i> .	1132
1079	Arsha Nagrani, Joon Son Chung, Weidi Xie, and	Justin Salamon, Christopher Jacoby, and Juan Pablo	1133
1080	Andrew Zisserman. 2020. Voxceleb: Large-scale	Bello. 2014. A dataset and taxonomy for urban sound	1134
1081	speaker verification in the wild. <i>Computer Speech &</i>	research. In <i>Proceedings of the 22nd ACM interna-</i>	1135
1082	<i>Language</i> , 60:101027.	<i>tional conference on Multimedia</i> , pages 1041–1044.	1136
1083	Tu Anh Nguyen, Wei-Ning Hsu, Antony D’Avirro,	Shibani Santurkar, Yann Dubois, Rohan Taori, Percy	1137
1084	Bowen Shi, Itai Gat, Maryam Fazel-Zarani, Tal Re-	Liang, and Tatsunori Hashimoto. 2023. Is a caption	1138
1085	mez, Jade Copet, Gabriel Synnaeve, Michael Hassid,	worth a thousand images? a study on representation	1139
1086	and 1 others. 2023. Espresso: A benchmark and	learning. In <i>ICLR</i> .	1140
1087	analysis of discrete expressive speech resynthesis. In	Christoph Schuhmann, Romain Beaumont, Richard	1141
1088	<i>Proc. Interspeech 2023</i> , pages 4823–4827.	Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti,	1142
1089	Daisuke Niizumi, Daiki Takeuchi, Masahiro Yasuda,	Theo Coombes, Aarush Katta, Clayton Mullis,	1143
1090	Binh Thien Nguyen, Yasunori Ohishi, and Noboru	Mitchell Wortsman, and 1 others. 2022. Laion-5b:	1144
1091	Harada. 2025. M2d2: Exploring general-purpose	An open large-scale dataset for training next genera-	1145
1092	audio-language representations beyond clap. <i>arXiv</i>	tion image-text models. <i>Advances in neural informa-</i>	1146
1093	<i>preprint arXiv:2503.22104</i> .	<i>tion processing systems</i> , 35:25278–25294.	1147
1094	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018.	David Snyder, Daniel Garcia-Romero, Gregory Sell,	1148
1095	Representation learning with contrastive predictive	Daniel Povey, and Sanjeev Khudanpur. 2018. X-	1149
1096	coding. <i>arXiv preprint arXiv:1807.03748</i> .	vectors: Robust dnn embeddings for speaker recog-	1150
1097	Soujanya Poria, Devamanyu Hazarika, Navonil Ma-	nition. In <i>2018 IEEE international conference on</i>	1151
1098	ajumder, Gautam Naik, Erik Cambria, and Rada Mi-	<i>acoustics, speech and signal processing (ICASSP)</i> ,	1152
1099	halcea. 2018. Meld: A multimodal multi-party	pages 5329–5333. IEEE.	1153
1100	dataset for emotion recognition in conversations.	Yi Su, Jisheng Bai, Qisheng Xu, Kele Xu, and Yong	1154
1101	<i>arXiv preprint arXiv:1810.02508</i> .	Dou. 2025. Audio-language models for audio-centric	1155
1102	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	tasks: A survey. <i>arXiv preprint arXiv:2501.15177</i> .	1156
1103	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	Luoyi Sun, Xuenan Xu, Mengyue Wu, and Weidi Xie.	1157
1104	try, Amanda Askell, Pamela Mishkin, Jack Clark, and	Auto-acd: A large-scale dataset for audio-language	1158
1105	1 others. 2021. Learning transferable visual models	representation learning. In <i>Proceedings of the 32nd</i>	1159
1106	from natural language supervision. In <i>International</i>	<i>ACM International Conference on Multimedia</i> , pages	1160
1107	<i>conference on machine learning</i> , pages 8748–8763.	5025–5034.	1161
1108	PmLR.		

1162	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang,	1219
1163	Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao	Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y Lin,	1220
1164	Zhang. 2024. Salmonn: Towards generic hearing	Andy T Liu, Jiatong Shi, Xuankai Chang, Guan-Ting	1221
1165	abilities for large language models. In <i>The Twelfth</i>	Lin, and 1 others. 2021. Superb: Speech processing	1222
1166	<i>International Conference on Learning Representa-</i>	universal performance benchmark. In <i>Proc. Inter-</i>	1223
1167	<i>tions</i> .	<i>speech 2021</i> , pages 1194–1198.	1224
1168	Michael Tschannen, Manoj Kumar, Andreas Steiner,	Zengwei Yao, Liyong Guo, Xiaoyu Yang, Wei Kang,	1225
1169	Xiaohua Zhai, Neil Houlsby, and Lucas Beyer. 2023.	Fangjun Kuang, Yifan Yang, Zengrui Jin, Long Lin,	1226
1170	Image captioners are scalable vision learners too. <i>Ad-</i>	and Daniel Povey. 2024. Zipformer: A faster and	1227
1171	<i>vances in Neural Information Processing Systems</i> ,	better encoder for automatic speech recognition. In	1228
1172	36:46830–46855.	<i>The Twelfth International Conference on Learning</i>	1229
1173	Joseph Turian, Jordie Shier, Humair Raj Khan, Bhiksha	<i>Representations</i> .	1230
1174	Raj, Björn W Schuller, Christian J Steinmetz, Colin	Ruibin Yuan, Yinghao Ma, Yizhi Li, Ge Zhang, Xingran	1231
1175	Malloy, George Tzanetakis, Gissel Velarde, Kirk Mc-	Chen, Hanzhi Yin, Yiqi Liu, Jiawen Huang, Zeyue	1232
1176	Nally, and 1 others. 2022. Hear: Holistic evaluation	Tian, Binyue Deng, and 1 others. 2023. Marble:	1233
1177	of audio representations. In <i>NeurIPS 2021 Competi-</i>	Music audio representation benchmark for universal	1234
1178	<i>tions and Demonstrations Track</i> , pages 125–145.	evaluation. <i>Advances in Neural Information Process-</i>	1235
1179	PMLR.	<i>ing Systems</i> , 36:39626–39647.	1236
1180	Luyu Wang, Pauline Luc, Yan Wu, Adria Recasens,	Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri,	1237
1181	Lucas Smaira, Andrew Brock, Andrew Jaegle, Jean-	Dan Jurafsky, and James Zou. 2023. When and why	1238
1182	Baptiste Alayrac, Sander Dieleman, Joao Carreira,	vision-language models behave like bags-of-words,	1239
1183	and 1 others. 2022. Towards learning universal audio	and what to do about it? In <i>The Eleventh Interna-</i>	1240
1184	representations. In <i>ICASSP 2022-2022 IEEE Inter-</i>	<i>national Conference on Learning Representations</i> .	1241
1185	<i>national Conference on Acoustics, Speech and Signal</i>	Salah Zaiem, Youcef Kemiche, Titouan Parcollet, Slim	1242
1186	<i>Processing (ICASSP)</i> , pages 4593–4597. IEEE.	Essid, and Mirco Ravanelli. 2023. Speech self-	1243
1187	Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and	supervised representation benchmarking: Are we	1244
1188	Juan Pablo Bello. 2022. Wav2clip: Learning robust	doing it right? In <i>Interspeech 2023</i> , pages 2873–	1245
1189	audio representations from clip. In <i>ICASSP 2022-</i>	2877.	1246
1190	<i>2022 IEEE International Conference on Acoustics,</i>	Piotr Żelasko, Daniel Povey, Jan Trmal, Sanjeev Khu-	1247
1191	<i>Speech and Signal Processing (ICASSP)</i> , pages 4563–	danpur, and 1 others. 2021. Lhotse: a speech data	1248
1192	4567. IEEE.	representation library for the modern deep learning	1249
1193	Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Tay-	ecosystem. <i>arXiv preprint arXiv:2110.12561</i> .	1250
1194	lor Berg-Kirkpatrick, and Shlomo Dubnov. 2023.	Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas	1251
1195	Large-scale contrastive language-audio pretraining	Steiner, Daniel Keysers, Alexander Kolesnikov, and	1252
1196	with feature fusion and keyword-to-caption augmen-	Lucas Beyer. 2022. Lit: Zero-shot transfer with	1253
1197	tation. In <i>ICASSP 2023-2023 IEEE International</i>	locked-image text tuning. In <i>Proceedings of the</i>	1254
1198	<i>Conference on Acoustics, Speech and Signal Process-</i>	<i>IEEE/CVF conference on computer vision and pat-</i>	1255
1199	<i>ing (ICASSP)</i> , pages 1–5. IEEE.	<i>tern recognition</i> , pages 18123–18133.	1256
1200	Xuenan Xu, Zhiling Zhang, Zelin Zhou, Pingyue Zhang,	Ge Zhu, Jordan Darefsky, and Zhiyao Duan. 2024. Ca-	1257
1201	Zeyu Xie, Mengyue Wu, and Kenny Q Zhu. 2023.	cophony: An improved contrastive audio-text model.	1258
1202	Blat: Bootstrapping language-audio pre-training	<i>IEEE/ACM Transactions on Audio, Speech, and Lan-</i>	1259
1203	based on audioset tag-guided synthetic data. In <i>Pro-</i>	<i>guage Processing</i> .	1260
1204	<i>ceedings of the 31st ACM International Conference</i>	Haina Zhu, Yizhi Zhou, Hangting Chen, Jianwei Yu,	1261
1205	<i>on Multimedia</i> , pages 2756–2764.	Ziyang Ma, Rongzhi Gu, Yi Luo, Wei Tan, and Xie	1262
1206	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	Chen. 2025. Muq: Self-supervised music represen-	1263
1207	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	tation learning with mel residual vector quantization.	1264
1208	Gao, Chengen Huang, Chenxu Lv, and 1 others.	<i>arXiv preprint arXiv:2501.01108</i> .	1265
1209	2024a. qwen2.5 technical report. <i>arXiv preprint</i>	A Appendix	1266
1210	<i>arXiv:2412.15115</i> .	A.1 Full Implementation details	1267
1211	Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue	We pretrain all models on CaptionStew. Training	1268
1212	Jiang, Xiaohuan Zhou, Yichong Leng, Yuanjun Lv,	data preparation is performed with the Lhotse (Že-	1269
1213	Zhou Zhao, Chang Zhou, and 1 others. 2024b. Air-	lasko et al., 2021) toolkit. All audio is resampled to	1270
1214	bench: Benchmarking large audio-language models	16 kHz and converted into 80-dimensional log-Mel	1271
1215	via generative comprehension. In <i>Proceedings of the</i>	filterbank features using a 25 ms window length	1272
1216	<i>62nd Annual Meeting of the Association for Compu-</i>		
1217	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages		
1218	1979–1998.		

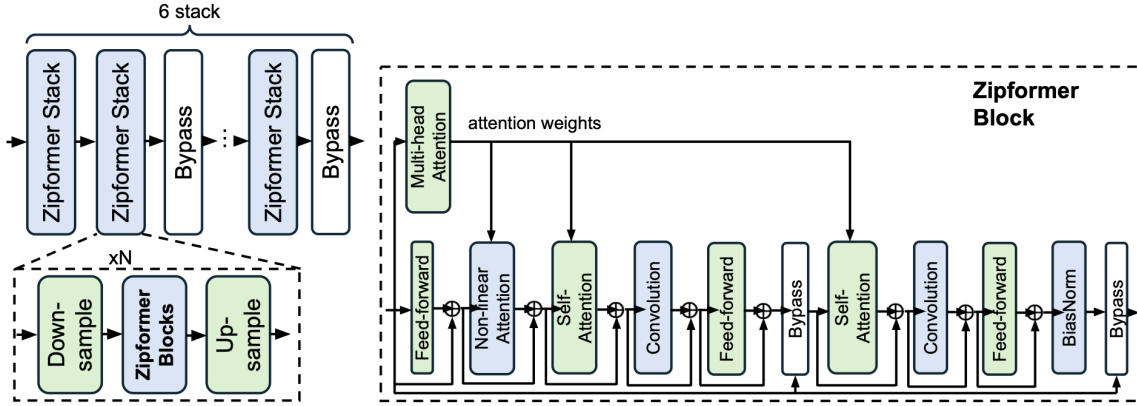


Figure 3: Model diagram of Zipformer.

Table 4: Zipformer performance across audio domains when trained from scratch on individual datasets, demonstrating cross-domain efficacy as a general audio encoder.

AudioSet (mAP)	VggSound (acc)	VoxCeleb2 (acc)	CREMA (acc)	MagnaTagATune (mAP)	NSynth-Instrument (acc)
0.46	54.2	84.8	65.4	0.38	78.8

and 10 ms hop size. Text is tokenized with a 50k-vocabulary BPE tokenizer (Lewis et al., 2020).

The audio encoder uses a Zipformer-M architecture (see Appendix A.2), chosen for its efficiency on long sequences and fast convergence. For contrastive pretraining, the text encoder follows BERT-base architecture (12 layers 768 hidden dimensions) (Devlin et al., 2019). For captioning pretraining, the text decoder adopts the BART-base decoder architecture (6 layers, 768 hidden dimensions) (Lewis et al., 2020). For the decoding mode, the ratio between autoregressive and parallel decoding is 0.25:0.75. It is worth noting that we use twice as many encoder layers as decoder layers to ensure comparable training speed across objectives.

Following prior works in audio-language pretraining (Elizalde et al., 2023; Wu et al., 2023; Mei et al., 2024; Bai et al., 2025), we experiment with two scenarios: training from scratch (denoted by *-scratch*) or initialized from pretrained checkpoints (denoted by *-init*). The audio encoder initializes from a Zipformer-based audio event classifier (Zipformer-AEC) trained on AudioSet (Gemmeke et al., 2017) with an mAP of 0.46, while text components use corresponding publicly available checkpoints. All models are trained on 8 Tesla V100 GPUs with an effective batch size of 640 seconds of audio per GPU. Training runs for 600k steps from scratch (14 days wall-clock time) or 200k steps if initialized from pretrained checkpoint.

A.2 Zipformer Model

Zipformer (Figure 3) employs a U-Net-inspired design with six Transformer stages that process sequences at multiple temporal resolutions. The stages operate at progressively decreasing then increasing frame rates (50, 25, 12.5, 6.25, 12.5, and 25 Hz), with residual and upsampling connections between stages to capture both fine-grained and long-range temporal patterns. We implement the original {2,2,3,4,3,2} block configuration, where each number indicates the blocks per stage. After processing through all stages, outputs are fused at 25 Hz to produce frame-level embeddings. The model incorporates several architectural improvements: BiasNorm for gradient stability, Swoosh activation functions for better convergence, and ScaledAdam optimizer. The resulting embeddings are 768-dimensional and used consistently across all downstream evaluation tasks.

Although Zipformer was originally designed for automatic speech recognition, we conducted preliminary experiments to validate its effectiveness as a general audio encoder across diverse domains. As in Table 4, our initial studies confirmed that Zipformer achieves competitive performance on environmental sound classification, music understanding, and speaker-related tasks, demonstrating its suitability as a unified backbone for multi-domain audio representation learning. This cross-domain efficacy makes it an appropriate choice for our experiments.

Table 5: Details of public-available datasets contribute to proposed CaptionStew dataset. We summarize their size, domain coverage, audio sources, captioning style, and generation pipelines.

Dataset	#audio/#cap	Domain	Audio source	Caption style	Caption generation pipeline
AudioCaps (Kim et al., 2019)	46k/46k	general (environmental, human/animal sounds)	AudioSet (Gemmeke et al., 2017)	Human-annotated, short description	crowdsourced
Clotho (Drossos et al., 2020)	5k/25k	environmental sounds	FreeSound	Human-annotated, short description	crowdsourced
MusicCaps (Agostinelli et al., 2023)	3k/3k	music	AudioSet	Expert musician-written, multi-sentence, fine-grained description	expert curation
WavCaps (Mei et al., 2024)	400k/400k	general (environmental, human/animal sounds)	AudioSet BBC Sound Effect FreeSound SoundBible	LLM-refined captions	three-stage pipeline: web-crawled raw descriptions → ChatGPT rewrite → filtering
AudioSetCaps (Bai et al., 2025)	1.9M/1.9M 4.0M/4.0M 182k/182k	general (environmental, human/animal sounds)	AudioSet YouTube8M (Abu-El-Haija et al., 2016) VggSound (Chen et al., 2020a)	LLM-generated, detailed, multi-sentence description	three-stage pipeline: LALM attribute extraction → LLM captioning → CLAP-based filtering
FusionAudio (Chen et al., 2025)	1.2M/1.2M	general (environmental, human/animal sounds)	AudioSet	LLM-augmented, multi-sentence, visual-enhanced description	multimodal context fusion (audio, visual, metadata) + LLM captioning
JamendoMaxCap (Roy et al., 2025)	360k/1.8M	music	Jamendo Platform	LLM-augmented, multi-sentence, fine-grained music description	retrieval-based metadata imputation + LLM captioning
ParaSpeechCaps (Diwan et al., 2025)	116k/116k (base) 924k/924k (scaled)	expressive speech	VoxCeleb1 (Nagrani et al., 2020) VoxCeleb2 (Chung et al., 2018) EARS (Richter et al., 2024) Expresso (Nguyen et al., 2023) Emilia (He et al., 2024)	Human-annotated/LLM-augmented, speaking-style description	crowdsourced / retrieval-based metadata imputation + LLM captioning

Table 6: Example caption sampled from each sourced dataset.

Dataset	Example Caption
AudioCaps	"Distant traffic sounds followed by a car passing closely."
Clotho	"Something is being sanded or dragged, manipulated, scraped."
MusicCaps	"This is an advertisement jingle music piece. It is an instrumental piece. The main theme is being played by the piano while there is a synth string sound in the melodic background. There is an emotional, heart-touching atmosphere. This piece could be used in the soundtrack of a drama movie during scenes of tragedy. It could also work well as an advertisement jingle where there is an attempted appeal to emotion."
WavCaps	"Music is playing while people are walking and crickets are chirping."
AudioSetCaps	"A choir performs a folk music piece, utilizing only their voices as instruments. The harmonious and uplifting sounds create an engaging and captivating listening experience."
FusionAudio	"A full choir is singing with powerful harmonized vocals"
JamendoMaxCaps	"The music is instrumental with a dominant piano sound, falling under the genres of ambient, classical, and contemporary. It carries a mood that is nostalgic and romantic, played in a 4/4 time signature at a tempo of 81.1 bpm. The piano piece evokes a sense of tranquility, making it suitable for scenarios depicting love scenes or peaceful moments in movies."
ParaSpeechCaps	"A male speaker delivers his words quickly with a medium-pitched voice. His speech exhibits a flowing rhythm and is recorded in an environment that is balanced in clarity. There is a subtle nasal quality to his speech, suggesting an American accent."

A.3 Sourced Datasets for CaptionStew

CaptionStew aggregates eight open-source audio caption datasets to address data scarcity and limited diversity in current audio-language pretraining. The constituent datasets span environmental sounds, music, and expressive speech, with fundamentally different captioning approaches—from crowdsourced human annotation to expert curation to various LLM-based generation pipelines. Table 5 and Table 6 detail each dataset’s characteristics and provide example captions that illustrate the diverse descriptive styles, ranging from concise

event descriptions to detailed multi-sentence narratives with fine-grained acoustic and contextual information. During aggregation, we filter audio samples longer than one minute for computational efficiency and remove samples that overlap with common audio understanding benchmarks (Kim et al., 2019; Drossos et al., 2020; Kim et al., 2019; Agostinelli et al., 2023; Fonseca et al., 2021; Chen et al., 2020a; Salamon et al., 2014) to prevent data leakage. This approach preserves the unique characteristics of each source while creating a unified corpus that captures broader semantic coverage than individual datasets.

Table 7: Details of the dataset used for assessing audio representation. [†]evaluate by GPT-4 in AIR-Bench. [‡]synthesized with public available speech datasets (Ardila et al., 2019; Busso et al., 2008; Cao et al., 2014; Livingstone and Russo, 2018; Poria et al., 2018) with fixed question template.

Evaluation Dataset	Task	#samples	#class	train	eval	Metrics
FSD-50k	Multi-label audio event classification	37,168 / 10,231	200	✓	✓	mAP
VggSound	Single-label audio event classification	183,730 / 15,446	309	✓	✓	accuracy
VoxCeleb2	Speaker identification	1,092,009 / 36,693	5,994	✓	✓	accuracy
CREMA-D	Speech emotion recognition	6,030 / 706	6	✓	✓	accuracy
MagnaTagATune	Music tagging	19,425 / 4,856	50	✓	✓	mAP
NSynth	Musical instrument classification	289,205 / 4,096	11	✓	✓	accuracy
AudioSet-strong	Sound event detection	103,463 / 16,996	456	✓	✓	PSDS1
AudioCaps	Text-to-audio retrieval	49,838 / 975	–	✓	✓	Recall@1 RougeL
ParaSpeechCaps	Audio captioning	116,516 / 500	–	✓	✓	
MusicCaps		2,663 / 500	–	✓	✓	
ClothoQA		7,044	–	✓	×	Score [†]
ParaLMQA [‡]		160,000	–	✓	×	
MusicQA		70,011	–	✓	×	
AIRBench-chat-sound	Open-formed question answering	400	–	×	✓	
AIRBench-foundation-emotion		1,000	–	×	✓	
AIRBench-foundation-gender		1,000	–	×	✓	
AIRBench-foundation-age		1,000	–	×	✓	
AIRBench-chat-sound		400	–	×	✓	

A.4 Evaluation Datasets

Table 7 details the evaluation datasets and their metrics used for assessing audio representation quality across our three evaluation protocols: linear probing (Fonseca et al., 2021; Chen et al., 2020a; Chung et al., 2018; Cao et al., 2014; Law et al., 2010; Engel et al., 2017; Hershey et al., 2021; Ebberts et al., 2022), audio-language alignment (Kim et al., 2019; Diwan et al., 2025; Agostinelli et al., 2023; Lin, 2004) and open-form question answering (Lipping et al., 2022; Liu et al., 2024; Huo et al., 2025; Yang et al., 2024b).

A.5 Main Results (cont.)

Table 9 presents linear probing results when using multi-head attention pooling instead of mean pooling. With learned attention pooling, the performance gap between contrastive and captioning objectives narrows substantially, particularly evident on speaker identification where captioning-scratch achieves 72.86% compared to 46.67% with mean pooling (Table 3). This demonstrates that captioning models benefit significantly from adaptive pooling mechanisms, while contrastive learning’s explicit optimization for clip-level representations shows less sensitivity to pooling strategy. These results underscore the critical importance of appropriate downstream module selection when evaluating different pretraining paradigms, as the choice of pooling mechanism can dramatically influence conclusions about objective effectiveness. The improved performance across all methods with

attention pooling also suggests that frame-level representations from both objectives contain rich information that can be better exploited through learned aggregation. SOTA results and SSL baseline results in Table 3 and Table 8 are quoted collectively from Niizumi et al. (2025); Turian et al. (2022); Li and Li (2022); Wang et al. (2022); Bharadwaj et al. (2025); Gong et al. (2022); Lanzendörfer et al. (2025); Bai et al. (2025); Yang et al. (2024b).

A.6 Additional Results

Aside from learning representations, we also compare against state-of-the-art audio-text retrieval models to assess our approach’s performance on the specific task it was designed for. Table 9 presents retrieval results for our best-performing model (Contrastive-init) against state-of-the-art audio-text retrieval model (Bai et al., 2025). Our model achieving comparable or superior results on benchmarks in various audio domains, with particularly strong performance on speech and music retrieval. The results indicate that our general-purpose audio-language pretraining approach can compete with specialized retrieval models while offering broader applicability across diverse usage scenarios.

A.7 The Use of Large Language Model

The authors used large language models to assist with writing refinement and grammatical corrections during the drafting process. All technical content, experimental design, analysis, and conclusions remain the authors’ original contributions.

Table 8: Linear probing results when using multi-head attention pooling.

Method	Model Initialization	Audio-language Pretraining	linear probing					
			AEC FSD50k	AEC VggSound	SID VoxCeleb2	SER CREMA	MTAG MagnaTagATune	INST NSynth
<i>Our Supervised Baselines</i>								
Zipformer-AEC (Yao et al., 2024)	AS SL	–	0.656	<u>56.23</u>	58.76	72.52	0.405	67.19
<i>Our Audio-language Pretrained Models</i>								
Contrastive- <i>scratch</i>	–	CS10M	0.640	52.81	72.86	74.50	0.406	75.00
Captioning- <i>scratch</i>	–	CS10M	0.619	50.97	56.64	70.40	0.406	72.10
Contrastive- <i>init</i>	AS SL	CS10M	0.670	54.89	72.24	73.09	0.412	76.70
Captioning- <i>init</i>	AS SL	CS10M	0.660	53.68	62.24	71.67	0.411	74.49
SOTA [†]			0.655	59.50	96.20	–	0.414	79.20

Table 9: audio-text retrieval of the best performing model (Contrastive-init) against state-of-the-art audio-text retrieval model. [†]reproduce by ourselves.

Model	Text-to-audio			Audio-to-text		
	AudioCaps	ParaSpeechCaps	MusicCaps	AudioCaps	ParaSpeechCaps	MusicCaps
AudioSetCaps [†]	49.7 / 79.2	0.8 / 2.5	13.4 / 30.6	45.9 / 80.8	0.2 / 3.8	12.0 / 29.0
Contrastive-init (ours)	44.4 / 79.0	29.6 / 61.6	22.4 / 53.0	47.2 / 78.8	27.0 / 57.4	26.0 / 56.2

Table 10: Comparison of lexical statistics and diversity across audio caption datasets and text corpora. We report vocabulary size (#vocab), average sentence length (avg. sent), and Distinct-n.

Source	#vocab	avg. sent	Distinct-n			
			1	2	3	4
AudioCaps	5,572	8.46	0.011	0.113	0.309	0.519
WavCaps	18,372	7.77	0.026	0.184	0.420	0.646
AudioSetCaps	21,061	28.22	0.006	0.082	0.249	0.450
FusionAudio	18,403	13.81	0.009	0.111	0.322	0.546
JamendoMaxCaps	27,906	63.29	0.002	0.026	0.079	0.153
ParaSpeechCaps	4,060	28.50	0.001	0.015	0.051	0.112
CaptionStew(Ours)	56,586	32.23	0.006	0.080	0.231	0.401
CC12M	366,175	17.03	0.046	0.486	0.813	0.927
WikiText-103	531,346	74.29	0.031	0.365	0.757	0.930

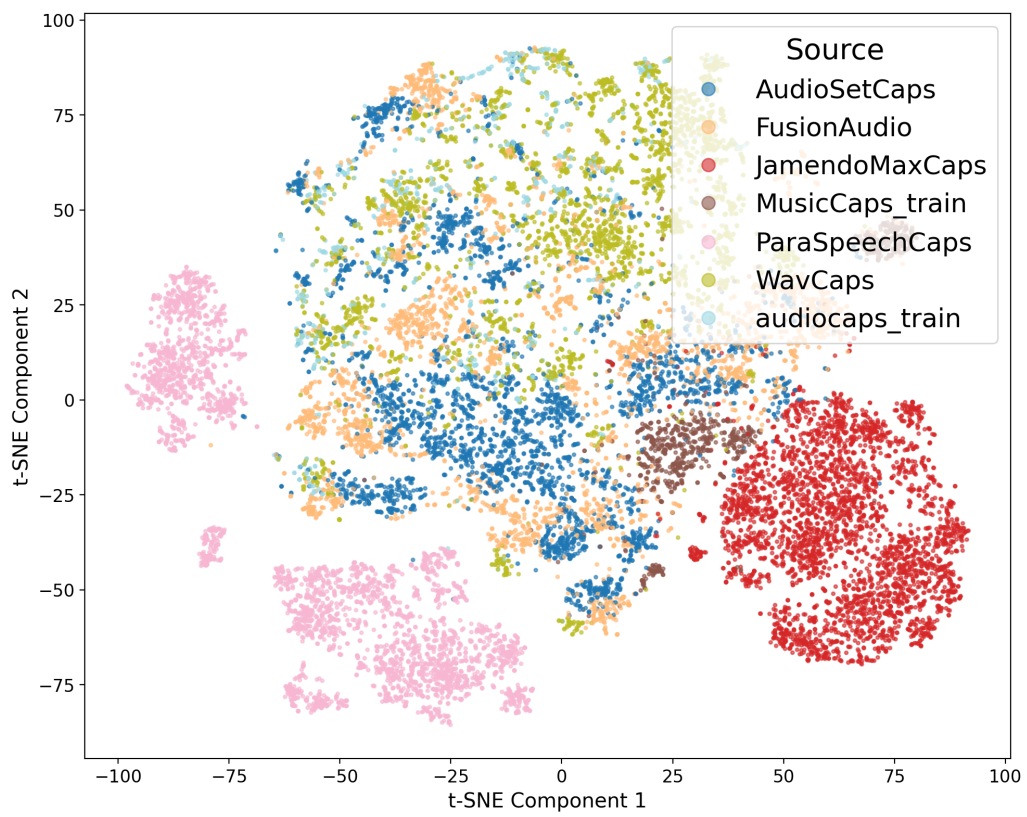


Figure 4: t-SNE visualization of sentence embedding of captions grouped by source.