

MECAT: A MULTI-EXPERTS CONSTRUCTED BENCHMARK FOR FINE-GRAINED AUDIO UNDERSTANDING TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

While large audio-language models have advanced open-ended audio understanding, they still fall short of nuanced human-level comprehension. This gap persists largely because current benchmarks, limited by data annotations and evaluation metrics, fail to reliably distinguish between generic and highly detailed model outputs. To this end, this work introduces MECAT, a **Multi-Expert Constructed Benchmark for Fine-Grained Audio Understanding Tasks**. Generated via a pipeline that integrates analysis from specialized expert models with Chain-of-Thought large language model reasoning, MECAT provides multi-perspective, fine-grained captions and open-set question-answering pairs. The benchmark is complemented by a novel metric: DATE (**D**iscriminative-**E**nhanced **A**udio **T**ext **E**valuation). This metric penalizes generic terms and rewards detailed descriptions by combining single-sample semantic similarity with cross-sample discriminability. A comprehensive evaluation of state-of-the-art audio models is also presented, providing new insights into their current capabilities and limitations. The data and code will be made publicly available.

1 INTRODUCTION

The human auditory system is highly effective at processing complex acoustic scenes. It can distinguish subtle variations in sound, such as telling a dog’s playful bark from a defensive growl (Plack, 2023), and isolate target speech from noisy backgrounds, an ability known as the cocktail party effect.

A central goal of machine hearing is to replicate this auditory intelligence to interpret raw audio signals as semantically rich perception (Lyon, 2017). Early works in machine hearing focused on closed-ended tasks such as sound event classification and automatic speech recognition. Large language models (LLM) have spurred the development of large audio-language models (LALMs), which have driven a shift towards more general open-ended tasks like audio captioning and audio question answering (Chu et al., 2023; Du et al., 2023; Hu et al., 2024; Shu et al., 2023; Wang et al., 2023; Tang et al., 2024; Rubenstein et al., 2023; Chen et al., 2023; Huang et al., 2024).

Despite these advances, current LALMs still fall short of achieving the comprehensive understanding that characterizes human hearing (Sakshi et al., 2025). This work argues that despite ongoing improvements in model architectures and data, a crucial and often-overlooked bottleneck is the existing evaluation benchmark.

The first challenge lies with data annotations, which suffers from several limitations. To begin with, the annotations in current benchmarks are often overly simplistic, consisting of event-level captions that lack detail Mei et al. (2024); Kim et al. (2019); Drossos et al. (2020) or question-answering tasks confined to close-ended formats Lipping et al. (2022); Wang et al. (2025); Sakshi et al. (2025). Furthermore, they typically adopt a monolithic perspective, providing a single, global description that fails to account for the selective nature of human hearing. Compounding these issues is a “one-to-many” data redundancy problem, where the same audio clips, often from AudioSet Gemmeke et al. (2017), are reused across multiple benchmarks, limiting the assessment of model generalization.

The second challenge is rooted in evaluation metrics. Traditional lexical-matching metrics, on the one hand, penalizes semantically correct but lexically different descriptions. Embedding-based metrics, on the other hand, better align with human perception, they often fail to distinguish between generic, vague captions and highly detailed, accurate ones. Even the more recent LLM-as-judge method, while demonstrating strong discriminative ability, is often hindered by practical constraints such as high costs and slow inference speeds, as well as its high dependency on model selection and prompt design.

Thus, current benchmarks inadequately evaluate audio understanding, as they often reward generic captions (e.g., A dog is barking and people are talking) for distinct scenarios (e.g., an excited bark in a park vs. a defensive bark during an argument). This limits their ability to differentiate between models with true perceptual accuracy and those producing vague outputs.

To this end, we introduce MECAT, a Multi-Expert Constructed Benchmark for Fine-Grained Audio Understanding Tasks. By integrating analysis from a series of specialized audio-related experts models, including content-specific models (e.g., for speech, music, and sound events) and content-unrelated models (e.g., for audio quality, reverberation and intensity), followed by Chain-of-Thought (CoT) enhanced LLM reasoning (Guo et al., 2025), MECAT provide fine-grained captions alongside open-set question-answering pairs. The captions primarily focus on providing a comprehensive, multi-perspective description of the acoustic scene, while the QA pairs are designed to probe for specific details and higher-level contextual reasoning that descriptive tasks alone cannot fully assess. Furthermore, we introduce a novel metric *DATE* (Discriminative-Enhanced Audio Text Evaluation), which is designed to better quantify the detail and accuracy of model’s response. It uniquely combines a weighted single-sample semantic similarity that penalizes generic terms while emphasizing discriminative phrases, and a cross-sample discriminability score that explicitly rewards the model’s responses for exceeding general descriptions. This design enables DATE to robustly distinguish between superficial and context-rich model outputs.

2 RELATED WORKS

2.1 AUDIO CAPTIONING BENCHMARK

Audio captioning benchmarks have been pivotal in advancing audio understanding works (Wu et al., 2019; Kim et al., 2019; Drossos et al., 2020; Yuan et al., 2025; Manco et al., 2023; Liu et al., 2024a;b). Early dataset like AudioCaps (Kim et al., 2019) and Clotho (Drossos et al., 2020) primarily relied on manual annotation, where human annotators provide one or more captions for each audio clip. While foundational, these benchmarks face a critical limitation: the coarse-grained nature of their annotations. During the annotation process, human annotators often produce generic, events-level descriptions rather than capturing the nuanced acoustic details of a scene. This results in a gold standard that lacks the specificity needed to evaluate fine-grained understanding.

While newer methods using LLMs for automatic labeling, such as in AutoACD (Sun et al., 2024) and LPMusicCaps (Doh & Nam, 2023), have improved scalability, they did not solve the granularity problem. Caption quality suffers from coarse input metadata like titles and tags, perpetuating generic descriptions.

2.2 AUDIO QUESTION-ANSWERING BENCHMARK

Audio Question Answering (QA) presents a more targeted evaluation of a model’s audio understanding abilities (Lipping et al., 2022; Wang et al., 2025; Li et al., 2022; Sakshi et al., 2025; Ma et al., 2025). Datasets like ClothoAQA (Lipping et al., 2022) and MusicAVQA (Li et al., 2022) have been developed with manually crafted question-answer pairs. However, similar to captioning benchmarks, they suffer from limitations that hinder the assessment of detailed understanding.

The main issue is their reliance on close-ended answer formats designed for easier automatic scoring. For example, many questions in ClothoAQA are limited to “yes/no” answers (Lipping et al., 2022), while other benchmarks like MMAU (Sakshi et al., 2025) utilize a multiple-choice format. While convenient for evaluation, these formats prevent the assessment of a model’s ability to generate detailed, descriptive answers and may encourage models to learn shallow pattern matching rather than deep understanding.

2.3 EVALUATION METRICS FOR AUDIO CAPTION AND QA

The evaluation of open-ended audio caption and QA is critically dependent on the choice of metric. However, existing metrics fail to adequately assess the fine-grained descriptive capabilities of modern generative models.

Traditional metrics, such as BLEU (Papineni et al., 2002), CIDEr (Vedantam et al., 2015), and SPICE (Anderson et al., 2016), operate by measuring lexical overlap with reference texts. This reliance on n-gram matching unfairly penalizes novel yet accurate descriptions that do not share the exact wording of the references.

To overcome the limitations of lexical matching, embedding similarity-based metrics were introduced. Approaches like FENSE (Zhang et al., 2022) were specifically designed for audio captioning. However, our experiments found that they still struggle to effectively distinguish between a generic, vague response and a highly detailed and accurate one. More recently, LLM-as-judge has been adopted for evaluating open-ended generation (Wang et al., 2025; Zheng et al., 2023). These methods show strong correlation with human judgment and, as our experiments confirmed, possess a high sensitivity to response specificity. However, LLM-as-judge suffer from practical limitations, such as high computational costs, slow evaluation speeds, and a strong dependency on the choice of the LLM and the design of the prompt (Lee et al., 2024; Zheng et al., 2023).

3 MECAT BENCHMARK OVERVIEW

As illustrated in Figure 1, MECAT is a comprehensive benchmark for fine-grained audio understanding, distinguished by its unique data sources, broad domain coverage, and two core evaluation tasks: MECAT-Caption and MECAT-QA.

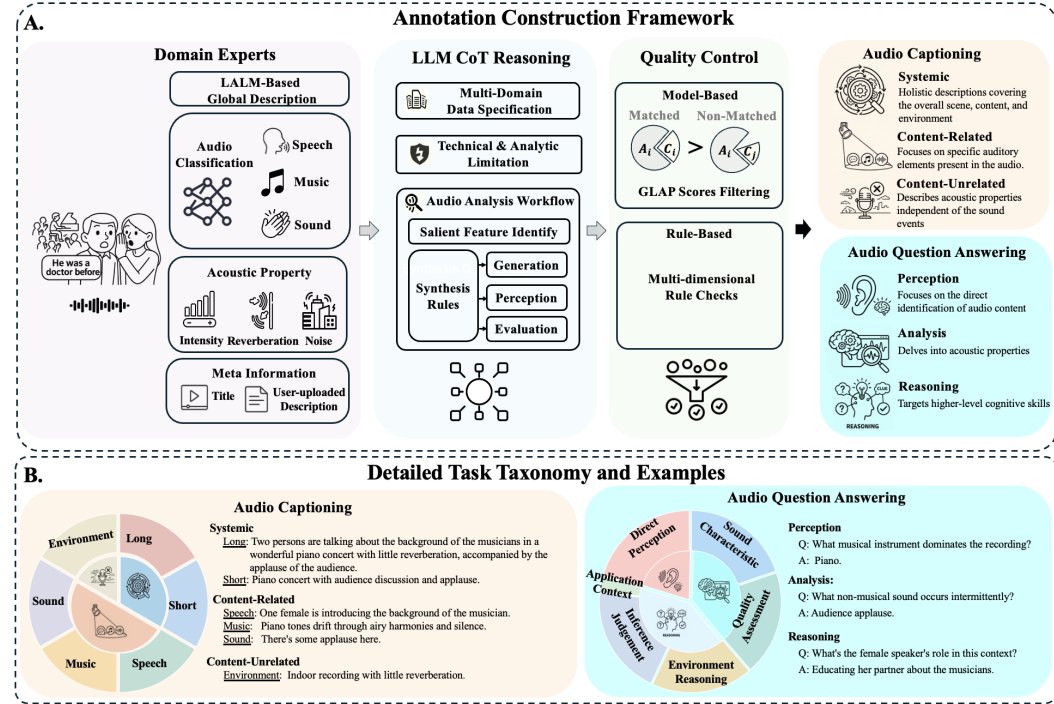


Figure 1: Overview of the MECAT Benchmark. (A) The proposed annotation construction pipeline. (B) Detailed task taxonomy and representative examples for Audio Captioning and Question Answering, showcasing the diversity of the dataset.

Table 1: Comparison of MECAT with Recent General Sound Evaluation Benchmark Datasets. † MP-LLM: Multiple Experts Models and LLM; ‡ Multi-Domain: This includes speech, music and sound-events (◊ denotes that domain were not elaborated in detail); § Extended Multi-Domain: This includes speech, music, sound-events, combinations thereof, and silence.

Task	Labeling	Dataset	Test Size	Domain	Source
Caption	Manual	AudioCaps (Kim et al., 2019)	~1.0k	Multi-Domain ^{‡,◊}	AudioSet
		Clotho (Drossos et al., 2020)	~1.0k	Multi-Domain ^{‡,◊}	Clotho
		SongDescriber (Manco et al., 2023)	0.7k	Music	MTG-Jamendo
	LLM	AudioCaps-Enhanced (Yuan et al., 2025)	0.9k	Multi-Domain ^{‡,◊}	AudioSet
		AutoACD (Sun et al., 2024)	1.0k	Multi-Domain ^{‡,◊}	AudioSet
		LPMusicCaps-MSD (Doh & Nam, 2023)	35k	Music	Song Dataset
QA	Manual	LPMusicCaps-MTT (Doh & Nam, 2023)	5k	Music	MagnaTagATune
		MECAT-Caption (Ours)	20k	Extended Multi-Domain [§]	ACAV100M
		ClothoQA (Lipping et al., 2022)	2k	Multi-Domain ^{‡,◊}	Clotho
	LLM	WavCaps-QA (Wang et al., 2025)	0.3k	Multi-Domain ^{‡,◊}	AudioSet and 2 others
		MusicAVQA (Li et al., 2022)	6k	Music	YouTube
		Audiocaps-QA (Wang et al., 2025)	0.3k	Multi-Domain ^{‡,◊}	AudioSet
	MP-LLM [†]	MMAU (Sakshi et al., 2025)	10k	Multi-Domain [‡]	AudioSet and 12 others
		EvalSIFT (Pandey et al., 2025)	30k	Speech	Open-source ASR
		MECAT-QA (Ours)	20k	Extended Multi-Domain [§]	ACAV100M

3.1 DATASET DESCRIPTION

To ensure data source novelty, MECAT is constructed from a carefully selected subset of ACAV100M (Lee et al., 2021). This approach contrasts with benchmarks, such as AudioCaps (Kim et al., 2019), Clotho (Drossos et al., 2020), and WavCaps-QA (Wang et al., 2025), which predominantly draw from a limited pool of sources such as AudioSet (Gemmeke et al., 2017) and Clotho (Drossos et al., 2020) (see Table 1). The dataset comprises approximately 20,000 Creative Commons-licensed audio clips, each with a maximum duration of 10 seconds **which is sufficient to contain one or a few salient acoustic events, while still allowing us to attach dense supervision for fine-grained, clip-local understanding.**

Based on this unique data foundation, MECAT encompasses eight distinct audio domains designed to comprehensively represent real-world acoustic scenarios. These categories include four *Pure* domains: silence (000), speech (S00), sound events (00A), and music (0M0), as well as all four possible combinations of *Mixed* domains that reflect the complexity of natural auditory environments (SM0, S0A, 0MA, and SMA). This extended multi-domain coverage, with its distribution detailed in Figure 2, enables a nuanced evaluation of models on complex acoustic scenes, such as those that combine piano music with spoken discussion and audience applause.

3.2 TASKS DEFINITION

As illustrated in Figure 1-B, the MECAT-Caption task delivers multi-perspective annotations for comprehensive evaluation. Each audio clip is annotated with a rich set of captions organized into three categories, which together comprise six distinct sub-categories. The first category, *Systemic Captions*, consists of two sub-categories: a concise short caption focused on primary audio content and a detailed long caption encompassing contextual details and event interactions. The second category, *Content-Specific Captions*, includes three sub-categories for the independent analysis of speech, music, and sound events. Crucially, to assess model performance across different levels of acoustic complexity, the evaluation for each content type is performed on corresponding pure domains (e.g., pure speech - S00) and all mixed domains. Notably, these captions also explicitly

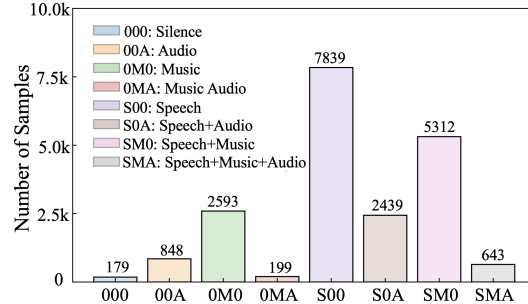


Figure 2: Distribution of audio samples across extended multi-domains in the MECAT, including speech, music, audio, combinations thereof, and silence.

state when a corresponding domain is absent. The final category is a single *Content-Unrelated Caption* that focuses exclusively on acoustic characteristics like audio quality and reverberation. For each of these six sub-categories, three synonymous reference captions are provided, yielding a total of 18 reference captions per clip and creating a significantly richer vocabulary than existing datasets (see Appendix B for more details).

The final score $\text{Score}_{\text{Cap}}$ for the MECAT-Caption task is calculated as a weighted average of the three main categories:

$$\text{Score}_{\text{Cap}} = 0.4 \cdot S_{\text{Systemic}} + 0.4 \cdot S_{\text{Content-Specific}} + 0.2 \cdot S_{\text{Content-Unrelated}},$$

where the category scores are themselves weighted sums of their sub-categories:

$$S_{\text{Systemic}} = 0.8 \cdot S_{\text{Long}} + 0.2 \cdot S_{\text{Short}},$$

$$S_{\text{Content-Specific}} = 0.6 \cdot S_{\text{Speech}} + 0.3 \cdot S_{\text{Music}} + 0.1 \cdot S_{\text{Sound}}.$$

The score for each content type (S_{Speech} , S_{Music} , S_{Sound}) is calculated as the unweighted mean of its performance on the corresponding pure domains (e.g., S00, 0M0, 00A) and all mixed domains. All coefficients are set heuristically to reflect the relative importance of overall scene, content detail, and acoustic context. Within the Content-Specific category, the 0.6/0.3/0.1 weights roughly follow the relative prevalence of Speech, Music, and Sound in ACAV100M.

Complementing the captioning task, MECAT-QA facilitates evaluation through targeted, probing questions. Each audio clip is paired with five question-answer pairs that span different cognitive skills, resulting in over 100,000 QA pairs in total. These pairs are organized into three cognitive categories. The first, *Perception*, focuses on the direct identification of audio content through its Direct Perception (DP) sub-category. The second, *Analysis*, delves into acoustic properties via two sub-categories: Sound Characteristics (SC), for examining properties like pitch, and Quality Assessment (QAS), for evaluating technical fidelity. The final and most complex category, *Reasoning*, targets higher-level cognitive skills through three sub-categories: Environment Reasoning (ER), requiring acoustic scene inference; Inference & Judgement (IJ), involving logical deductions; and Application Context (AC), testing the understanding of practical scenarios.

The scoring for MECAT-QA is designed to ensure equal contribution from each cognitive skill. The overall score is the unweighted arithmetic mean of the scores from all six individual sub-categories:

$$\text{Score}_{\text{QA}} = (S_{\text{DP}} + S_{\text{SC}} + S_{\text{QAS}} + S_{\text{ER}} + S_{\text{IJ}} + S_{\text{AC}})/6.$$

4 ANNOTATION CONSTRUCTION

This section details the MECAT annotation construction pipeline. As illustrated in Figure 1, the process starts with a audio classification stage identifying the domain of each audio clip. Based on the resulting domains, the clip is then processed by a series of specialized expert models.

The structured outputs from these experts are subsequently synthesized using LLM CoT reasoning to generate fine-grained captions and open-set QA pairs. The pipeline concludes with a rigorous quality control stage to ensure the reliability of all final annotations. The complete list of the used models is available in Appendix C.

4.1 DOMAIN EXPERTS

For each audio clip, we first use Audio Flamingo 2 (Ghosh et al., 2025) to generate a global, event-level summarization in natural language. Furthermore, we apply a series of domain expert models for more detailed analysis.

Audio Classification For each audio clip, we use CED-Base (Dinkel et al., 2024a) to predict AudioSet (Gemmeke et al., 2017) labels for every 2-second, non-overlapping interval. This process results in a sequence of multi-label predictions for each clip. Based on the CED prediction, we categorize each clip into one of eight distinct domains: 000, 00A, 0M0, S00, SM0, 0MA, S0A, SMA, as detailed in Dataset Description Section.

Speech-focused Analysis For speech-domain clips (S00, S0A, SM0, SMA), we employ a speech-focused analysis pipeline (Figure 3-A). The pipeline consists of automatic speech recognition, language identification, and speaker diarization. Using the temporal boundaries from diarization, we extract each speaker’s attributes, including gender, age, emotion, and English accent. The probabilities of these results are also utilized for subsequent LLM reasoning.

Music-focused Analysis For music-domain clips (0M0, SM0, 0MA, SMA), a music-focused analysis pipeline is employed (Figure 3-B). It consists of LALM-based global description of music content (Audio Flamingo 2 (Ghosh et al., 2025)), musical attribute analysis, and music separation. Musical attribute analysis provides a series of perceptual and technical attributes such as emotions and tempo. The music separation module isolates vocal tracks from the instrumental background, which are then routed to the speech analysis pipeline.

Sound Events-focused Analysis For audios in 00A, we directly utilize the events labels predicted by the CED-Base model during the audio classification stage.

Acoustic Properties Analysis To extract fundamental signal characteristics and assess the recording environment, we apply a universal acoustic property analysis pipeline to all audio clips (Figure 3-C). The analysis content includes signal intensity, speech quality assessment, and reverberation. Signal intensity is quantified via Root Mean Square (RMS). For audio quality, we conduct both DNSMOS (Reddy, 2021) and NISQA2 (Mittag et al., 2021) assessments to measure signal distortion, background noise, and perceptual quality. We also characterize the acoustic environment by estimating the reverberation time of the recording space.

4.2 LLM CoT REASONING

Our pipeline employs a Chain-of-Thought (CoT) guided LLM (Deepseek-R1: Guo et al., 2025) to synthesize a set of rich annotations. The model is instructed to reason over the outputs from all preceding analyses and the metadata. This reasoning process weighs evidence from various sources to resolve inconsistencies and identifying salient features. The final output consists of captions and corresponding QA pairs, where each item is annotated with a confidence level. The complete prompt is shown in Appendix D.

4.3 QUALITY CONTROL

The model-based filtering use GLAP (Dinkel et al., 2025) to compute the cosine similarity between audio clip and its systemic long caption embeddings. A sample is kept only if the similarity of its correct audio-caption pair exceeds its average similarity with a set of 6 other randomly selected captions by an empirically set threshold of 6.

We further apply rule-based filtering including LLM confidence thresholding, domain consistency between audio classification and LLM output, and hallucination removal (Barański et al., 2025).

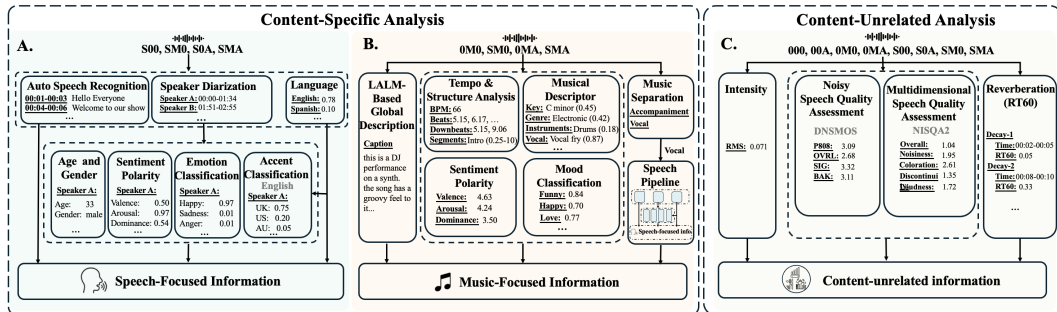


Figure 3: Domain Experts for Speech, Music, and Acoustic Properties.

5 METRIC DESIGN

Existing evaluation metrics demonstrate significant limitations when evaluating fine-grained, detailed descriptions. To address this, we propose DATE, a metric built on Sentence-BERT (Reimers & Gurevych, 2019) that improves semantic assessment by combining single-sample semantic similarity and cross-sample discriminability score.

Single-Sample Semantic Similarity We apply [embedding-level](#) term frequency-inverse document frequency (TF-IDF) weighting to token embeddings from Sentence-BERT to emphasize tokens that are frequent within a single sample but rare across the dataset ([details in Appendix E](#)). The weighted embedding vector \mathbf{v}_T for a given sentence T is computed as:

$$\mathbf{v}_T = \sum_{t \in T} (\text{TF}_{\text{emb}}(t, T) \cdot \text{IDF}_{\text{emb}}(t)) \cdot E(t),$$

where t is a token in T . The term $\text{TF}_{\text{emb}}(t, T)$, $\text{IDF}_{\text{emb}}(t)$, and $E(t)$ are the frequency, inverse document frequency, and the Sentence-BERT embedding, respectively. The single-sample semantic similarity, $S_{\text{sim},i}$, is the cosine similarity between the weighted embeddings of the candidate and reference text. $S_{\text{sim},i} = (\mathbf{v}_{\text{cand}} \cdot \mathbf{v}_{\text{ref}}) / (\|\mathbf{v}_{\text{cand}}\| \|\mathbf{v}_{\text{ref}}\|)$.

Cross-Sample Discriminability An ideal description should be clearly distinguishable from descriptions of other audio samples. We construct a cross-sample similarity matrix, \mathcal{M} , where each element $M_{i,j}$ is the score between the reference description for audio i and the candidate description for audio j . For each sample i , we rank the correctly matched score $M_{i,i}$ against all candidate scores $\{M_{i,j}\}_{j=1}^N$. Denoting this rank as r_i , the discriminability score is defined as:

$$S_{\text{dis},i} = 1 - r_i / N.$$

This rewards candidates that rank highly for their correct reference, approaching 1 for top ranks and 0 for bottom ranks.

DATE To ensure a balanced evaluation for both descriptive accuracy and uniqueness, the DATE score for each sample DATE_i is defined as the harmonic mean of its semantic similarity ($S_{\text{sim},i}$) and discriminability ($S_{\text{dis},i}$):

$$\text{DATE}_i = \frac{2 \cdot S_{\text{sim},i} \cdot S_{\text{dis},i}}{S_{\text{sim},i} + S_{\text{dis},i}} \in [0, 1].$$

The DATE score of a dataset is computed as $\text{DATE} = \frac{1}{N} \sum_{i=1}^N \text{DATE}_i$.

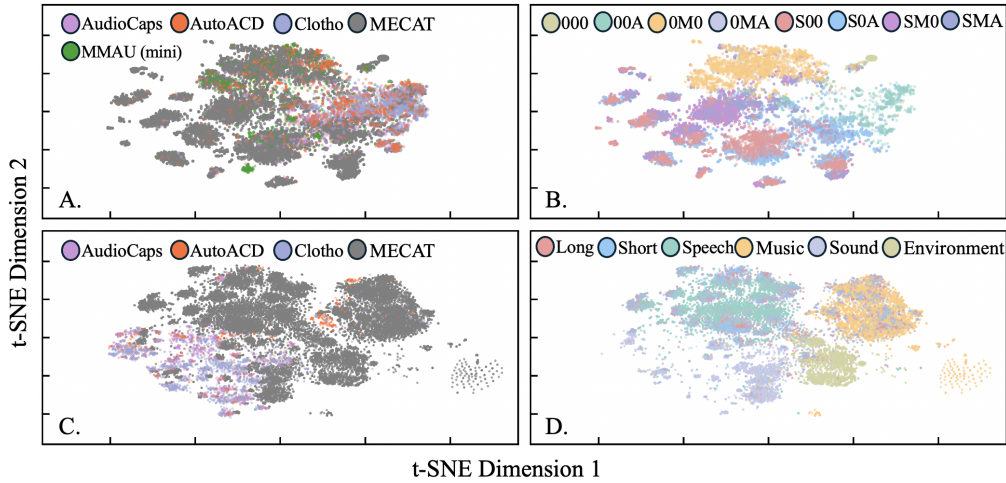


Figure 4: t-SNE plots of MECAT audio embeddings compared to other benchmarks (A), further clustered by domain (B). Caption embeddings are visualized in C and clustered by categories in D. Audio embeddings and captions embeddings are extracted from Dasheng-Base (Dinkel et al., 2024b) and Sentence-BERT (Reimers & Gurevych, 2019) respectively.

6 RESULT AND DISCUSSION

This section presents an analysis of MECAT’s data diversity, the analysis of the DATE metric, and a comprehensive evaluation of state-of-the-art models on MECAT.

6.1 DATA ANALYSIS

Distribution and Diversity The t-SNE analysis in Figure 4 highlights MECAT’s superior coverage. Regarding *audio* (Figure 4-A/B), MECAT spans the full feature space with distinct internal clusters for pure domains, whereas existing benchmarks remain densely clustered around sound-event regions. Regarding *captions* (Figure 4-C/D), MECAT exhibits significantly richer semantic diversity driven by distinct content categories (speech, music, events) rather than simple length variations.

Quality Validation Annotation trustworthiness was further verified via a human preference A/B test ($N = 150$). As detailed in Table F.1, generic and incorrect baselines were significantly outperformed by MECAT captions ($> 94\%$ win rates), while statistical parity with human references was achieved (56.9% win rate). Consequently, the reliability of the automated pipeline against hallucinations or vagueness is confirmed (comprehensive details in Appendix F).

6.2 METRIC ANALYSIS

This section validates our proposed metric, DATE, against the strong baseline FENSE, using an LLM-as-judge method as the upper-bound reference (prompts and reliability analysis in Appendices G and H).

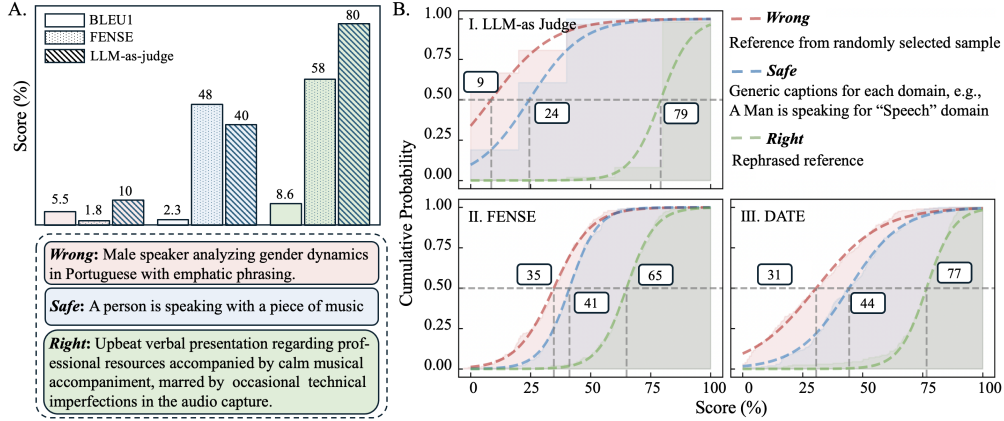


Figure 5: Metric Analysis. (A) Case study of existing metrics. Reference: “An animated woman’s voice shares information about learning materials while melodic instruments play quietly underneath, with persistent low-quality artifacts in the recording.” (B) Cumulative Distribution Functions (CDF) of LLM-as-judge, FENSE, and DATE on Caption (left) and QA (right). Larger distances between CDF curves indicate better discriminative ability of the metric.

Qualitative Analysis The case study in Figure 5-A exposes critical flaws in existing metrics. Lexical-based metrics like BLEU-1 are semantically unreliable, assigning higher scores to Wrong captions than Safe ones. While FENSE improves upon this, it struggles to distinguish high-quality (Right) from vague (Safe) captions, showing a negligible score gap ($\Delta \approx 10$). In contrast, DATE aligns with the clear separation observed in LLM-as-judge scores. DATE demonstrates a clear advantage over FENSE, evidenced by significantly larger median score spans for both Right vs. Wrong (DATE: 46 vs. FENSE: 30) and Right vs. Safe (DATE: 33 vs. FENSE: 24). **Quantitative Analysis** The Cumulative Distribution Function (CDF) curves in Figure 5-B further quantify discriminative power, where larger inter-curve distances indicate superior performance.

Alignment with Human Judgment To assess the alignment between DATE and human perception, we utilized the same 150 A/B caption pairs described in Appendix F. Captions preferred by human evaluators received substantially higher DATE scores (Mean: 90.9) compared to non-preferred

ones (Mean: 49.3). This significant margin indicates that DATE is strongly correlated with human preferences regarding accuracy and detail, more information are detailed in Table F.2.

6.3 MODEL PERFORMANCE ON MECAT

6.3.1 OVERALL PERFORMANCE

An extensive collection of publicly available models was evaluated, including 15 models for Captioning (13 LALMs, 2 traditional baselines) and 13 LALMs for QA. The evaluated models are strictly categorized into two primary types: traditional Caption-Only models (which constitute the non-LALM baselines for the Captioning task, e.g., EnClap, Pengi) and Large Audio-Language Models (LALMs). The LALMs are further stratified into four architectural subcategories: i) *Audio-focused LALMs* (e.g., Kimi-Audio), ii) *Omni LALMs* (e.g., Qwen-Omni), iii) *Multimodal LALMs*, and iv) the *Gemini series*. Detailed specifications regarding the model architectures and the corresponding prompts utilized in this study are provided in Appendix I.

Performance on MECAT-Caption. As shown in Table 2, LALMs demonstrate a substantial advantage over traditional baselines, with scores ranging from 29.4 (Pengi) to 51.6 (Gemini 2.5 Flash), attributed to superior instruction-following capabilities. In terms of domain stability, while Speech tasks remain robust, Music and Sound tasks suffer significant degradation (10% \sim 25%) when transitioning from Pure to Mixed domains, likely due to the speech-centric bias of current architectures. Furthermore, the universal suboptimal performance on Content-Unrelated tasks highlights an urgent need to capture intrinsic sound properties beyond high-level event recognition.

Performance on MECAT-QA. The hierarchy in the QA task (Table 3) largely mirrors captioning, yet with a notable shift: the gap between proprietary and open-weight models narrows significantly. Qwen3-Omni (52.3) slightly outperforms both Gemini-2.5-Flash (52.1) and Pro (51.5). At a granular level, we observe a distinct *capability dichotomy*: models excel in Direct Perception and content-based Inference, but degrade significantly on tasks requiring the analysis of intrinsic acoustic properties, such as Quality Assessment and Environment Reasoning. For instance, even top models fail to exceed a score of 40 in Quality Assessment, confirming that current LALMs prioritize high-level semantics over nuanced acoustic interpretation.

6.3.2 IN-DEPTH ANALYSIS: LALM BOTTLENECKS

While standard metrics rank the models, they do not fully reveal the underlying behavioral tendencies. We utilize the Captioning task as a representative case study to investigate two critical bottlenecks: lack of discriminability and robustness against hallucination.

Table 2: Model performance (DATE %) on MECAT-Caption. **Bold** indicates the best performance, and underline indicates the second best. [†] indicates that the model or its previous version (Audio Flamingo 2) was explicitly used in the data construction process.

Type	Model	Systemic		Content-Specific						Content Unrelated	Score _{Cap}
		Long	Short	Speech		Music		Sound			
				Pure	Mixed	Pure	Mixed	Pure	Mixed	Env	
Caption-Only	Pengi	43.5	46.8	27.2	29.5	29.3	13.1	42.8	14.6	7.1	29.4
	EnClap	48.6	53.1	30.2	31.8	17.9	15.9	48.8	15.2	6.8	31.9
LALM	Phi-4-Multimodal-Instruct	42.4	44.0	26.9	31.3	14.9	24.0	28.5	18.1	13.1	30.0
	Kimi-Audio-7B-Instruct	49.5	54.2	30.0	31.3	27.7	16.9	43.1	16.2	7.0	32.8
	Baichuan-Audio-Instruct	42.6	36.5	46.0	40.4	21.3	20.7	44.8	17.7	15.1	33.7
	Audio Flamingo 2 [†]	48.6	49.7	30.5	34.3	28.8	25.6	41.2	18.5	17.5	35.3
	Baichuan-Omni	47.0	50.9	43.5	41.7	35.2	13.7	34.3	19.7	11.3	35.6
	Mimo-Audio-Instruct	56.5	56.9	45.8	44.9	35.8	19.4	46.8	21.0	9.8	40.1
	Audio Flamingo 3 [†]	52.5	51.5	49.3	48.8	40.4	24.8	50.6	21.9	11.5	40.4
	Qwen3-Omni	47.9	43.7	50.2	48.7	51.2	26.8	49.0	19.5	18.2	40.4
	Step-audio-2-mini	55.6	58.7	44.2	43.6	35.3	32.0	42.8	18.9	16.1	41.5
	Qwen2.5-Omni 3B	56.4	55.2	42.5	41.3	46.6	29.7	52.9	23.9	19.4	42.5
	Qwen2.5-Omni 7B	61.1	56.5	39.9	40.9	32.1	30.9	50.7	23.8	17.9	42.6
	Gemini-2.5-Flash	65.6	63.9	57.5	57.5	52.9	41.0	52.2	28.3	22.1	51.6
	Gemini-2.5-Pro	62.3	62.4	56.6	57.5	53.6	38.7	53.4	29.9	24.0	50.6

Table 3: Model Performance (DATE %) on MECAT-QA. **Bold** indicates the best performance, and underline indicates the second best. [†] indicates that the model or its previous version (Audio Flamingo 2) was explicitly used in the data construction process.

Model	Perception	Analysis		Reasoning			Score _{QA}
	Direct Perception	Sound Characteristics	Quality Assessment	Environment Reasoning	Inference & Judgment	Application Context	
Kimi-Audio-7B-Instruct	45.6	39.2	18.7	34.6	48.9	41.2	38.0
Baichuan-Audio-Instruct	40.7	45.2	31.0	35.1	49.0	46.9	41.3
Audio Flamingo 2 [†]	45.1	46.3	34.9	37.5	44.0	42.4	41.7
Baichuan-Omni	43.6	44.7	33.7	39.9	49.3	49.1	43.4
Phi-4-Multimodal-Instruct	48.4	46.3	34.7	40.2	49.3	48.7	44.6
Mimo-Audio-Instruct	<u>59.3</u>	49.3	24.9	39.1	52.7	46.2	45.2
Step-Audio-2-mini	57.7	54.3	37.2	39.2	48.9	48.0	47.6
Audio Flamingo 3 [†]	53.8	50.2	36.0	43.0	54.5	49.6	47.8
Qwen2.5-Omni 3B	55.7	53.2	38.6	41.1	51.8	50.8	48.5
Qwen2.5-Omni 7B	57.8	52.9	<u>39.1</u>	44.0	53.2	50.8	49.6
Qwen3-Omni	61.7	<u>54.6</u>	39.3	45.0	56.9	56.1	52.3
Gemini-2.5-Flash	56.3	55.3	37.7	46.8	58.6	58.0	<u>52.1</u>
Gemini-2.5-Pro	55.5	54.4	37.7	47.6	<u>57.3</u>	<u>56.6</u>	51.5

The Critical Role of Discriminability. A granular analysis of the Pure Speech subset using DATE (Appendix J) underscores the interplay between Semantic Similarity and Discriminability (note that for other subtasks, discriminability could be derived from the Similarity and DATE scores in Table 2 and Appendix K). The Gemini series dominates by synergizing top-tier similarity with exceptional discriminability (e.g., 78.4 vs. runner-up 64.7). Notably, Qwen3-Omni illustrates the vital role of discriminative power: despite moderate semantic similarity, its strong discriminability effectively compensates for this gap, securing a top-tier ranking. This reinforces that for high-quality captioning, the capacity to distinguish unique acoustic features is just as critical as aligning with ground-truth semantics.

Hallucination in Silent Segments. Our qualitative evaluation of silent segments (Appendix L) reveals a significant robustness issue. While models like Gemini and Audio Flamingo 2 correctly identify silence, many LALMs fail to inhibit generation when valid acoustic cues are absent, resulting in hallucinations of specific but unrelated text (e.g., *I’m gonna be a daddy*” or *Thank you*”). This tendency to over-generate exposes a fundamental vulnerability, underscoring the necessity for improved rejection mechanisms in future architectures.

7 CONCLUSION AND FUTURE WORK

In this work, we introduced MECAT, a Multi-Experts Constructed Benchmark leveraged by Chain-of-Thought reasoning to advance fine-grained audio understanding in Captioning and QA tasks. Complementing this, we proposed DATE, a novel metric tailored to penalize vague terminology and incentivize detailed, discriminative descriptions.

Our comprehensive evaluation reveals distinct patterns in the current landscape of Large Audio-Language Models (LALMs). A primary finding is the strong bias toward speech content; most models exhibit a diminished capacity to perceive and reason about non-speech elements, such as background music, acoustic events, and audio quality. Furthermore, models struggle to generate discriminative content, often defaulting to safe, generic captions that fail to distinguish unique audio features. Perhaps most critically, our qualitative analysis uncovers a significant robustness issue: when valid acoustic cues are absent (i.e., in silent segments), many models fail to inhibit generation, resulting in hallucinated speech content.

Future work should address specific constraints identified in this study. Regarding data, expanding the benchmark’s coverage to include extremely long and streaming audio contexts is necessary to address current limitations in duration distributions. DATE is designed for settings with diverse audio inputs and rich, heterogeneous annotations. However, in scenarios where the audio clips and reference labels are highly homogeneous (e.g., different levels of white noise that are not differentiated in the reference captions) the DATE score will degenerate. Therefore, the refinement of metrics to capture subtle semantic discrepancies is a priority.

8 REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our research, we provide the following details. The data and code associated with this benchmark will be made publicly available upon publication. A comprehensive description of our data processing steps is provided in Section 3. All models used in this study, along with their citations, are listed in Appendix C. Furthermore, the specific prompts used for caption and QA generation, scoring, and LALM evaluation are provided in Appendices D, E, and G, respectively.

REFERENCES

- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *Proceedings of the European conference on computer vision*, pp. 382–398. Springer, 2016.
- Mateusz Barański, Jan Jasiński, Julitta Bartolewska, Stanisław Kacprzak, Marcin Witkowski, and Konrad Kowalczyk. Investigation of whisper asr hallucinations induced by non-speech audio. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2025.
- Hervé Bredin. pyannote. audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe. In *Proceedings of the 24th Interspeech Conference (interspeech)*, pp. 1983–1987. ISCA, 2023.
- Felix Burkhardt, Johannes Wagner, Hagen Wierstorf, Florian Eyben, and Björn Schuller. Speech-based age and gender prediction with transformers. In *Speech Communication; 15th ITG Conference*, pp. 46–50. VDE, 2023.
- Feilong Chen, Minglun Han, Haozhi Zhao, Qingyang Zhang, Jing Shi, Shuang Xu, and Bo Xu. X-LLM: Bootstrapping advanced large language models by treating multi-modalities as foreign languages. *arXiv preprint arXiv:2305.04160*, 2023.
- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-Audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Heinrich Dinkel, Yongqing Wang, Zhiyong Yan, Junbo Zhang, and Yujun Wang. Ced: Consistent ensemble distillation for audio tagging. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 291–295. IEEE, 2024a.
- Heinrich Dinkel, Zhiyong Yan, Yongqing Wang, Junbo Zhang, Yujun Wang, and Bin Wang. Scaling up masked audio encoder learning for general audio classification. In *Proceedings of the 25th Interspeech Conference (interspeech)*, pp. 547–551, 2024b.
- Heinrich Dinkel, Zhiyong Yan, Tianzi Wang, Yongqing Wang, Xingwei Sun, Yadong Niu, Jizhong Liu, Gang Li, Junbo Zhang, and Jian Luan. Glap: General contrastive audio-text pretraining across domains and languages, 2025.
- Seunghoon Doh and Juhan Nam. Lp-musiccaps: Llm-based pseudo music captioning. In *Proceedings of the 24th International Society for Music Information Retrieval Conference*. International Society for Music Information Retrieval Conference, 2023.

- Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: an audio captioning dataset. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 736–740. IEEE, 2020.
- Zhihao Du, Jiaming Wang, Qian Chen, Yunfei Chu, Zhifu Gao, Zerui Li, Kai Hu, Xiaohuan Zhou, Jin Xu, Ziyang Ma, et al. LauraGPT: Listen, attend, understand, and regenerate audio with GPT. *arXiv preprint arXiv:2310.04673*, 2023.
- Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 776–780. IEEE, 2017.
- Sreyan Ghosh, Zhifeng Kong, Sonal Kumar, S Sakshi, Jaehyeon Kim, Wei Ping, Rafael Valle, Dinesh Manocha, and Bryan Catanzaro. Audio flamingo 2: An audio-language model with long-audio understanding and expert reasoning abilities. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, pp. 1–48, 2025. URL <https://openreview.net/forum?id=xWu5qpDK6U>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Scott H. Hawley. Shaart: Speech and hearing in audio and real time. <https://github.com/drscotthawley/SHAART>, 2023.
- Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Lingwei Meng, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, et al. Wavllm: Towards robust and adaptive speech large language model. In *Proceedings of the Findings of the Association for Computational Linguistics (EMNLP)*, pp. 4552–4572, 2024.
- Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audiogpt: Understanding and generating speech, music, sound, and talking head. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 23802–23804, 2024.
- Jaeyong Kang and Dorien Herremans. Towards unified music emotion recognition across dimensional and categorical models, 2025. URL <https://arxiv.org/abs/2502.03979>.
- Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *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)*, pp. 119–132, 2019.
- Taejun Kim and Juhan Nam. All-in-one metrical and functional structure analysis with neighborhood attentions on demixed audio. In *Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2023.
- KimiTeam, Ding Ding, Zeqian Ju, Yichong Leng, Songxiang Liu, Tong Liu, Zeyu Shang, Kai Shen, Wei Song, Xu Tan, Heyi Tang, Zhengtao Wang, Chu Wei, Yifei Xin, Xinran Xu, Jianwei Yu, Yutao Zhang, Xinyu Zhou, Y. Charles, Jun Chen, Yanru Chen, Yulun Du, Weiran He, Zhenxing Hu, Guokun Lai, Qingcheng Li, Yangyang Liu, Weidong Sun, Jianzhou Wang, Yuzhi Wang, Yuefeng Wu, Yuxin Wu, Dongchao Yang, Hao Yang, Ying Yang, Zhilin Yang, Aoxiong Yin, Ruibin Yuan, Yutong Zhang, and Zaida Zhou. Kimi-audio technical report, 2025. URL <https://arxiv.org/abs/2504.18425>.
- Qiuqiang Kong, Yin Cao, Haohe Liu, Keunwoo Choi, and Yuxuan Wang. Decoupling magnitude and phase estimation with deep resunet for music source separation. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference (ISMIR)*, pp. 342–349. Citeseer, 2021.

- 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/CVF International Conference on Computer Vision*, pp. 10274–10284, 2021.
- Yebin Lee, Imseong Park, and Myungjoo Kang. Fleur: An explainable reference-free evaluation metric for image captioning using a large multimodal model. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3732–3746, 2024.
- Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer questions in dynamic audio-visual scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19108–19118, 2022.
- Tianpeng Li, Jun Liu, Tao Zhang, Yuanbo Fang, Da Pan, Mingrui Wang, Zheng Liang, Zehuan Li, Mingan Lin, Guosheng Dong, et al. Baichuan-audio: A unified framework for end-to-end speech interaction. *arXiv preprint arXiv:2502.17239*, 2025.
- Yadong Li, Haoze Sun, Mingan Lin, Tianpeng Li, Guosheng Dong, Tao Zhang, Bowen Ding, Wei Song, Zhenglin Cheng, Yuqi Huo, et al. Baichuan-omni technical report. *arXiv preprint arXiv:2410.08565*, 2024.
- Samuel Lipping, Parthasaarathy Sudarsanam, Konstantinos Drossos, and Tuomas Virtanen. Clotho-aqa: A crowdsourced dataset for audio question answering. In *Proceedings of the 30th European Signal Processing Conference (EUSIPCO)*, pp. 1140–1144. IEEE, 2022.
- Jizhong Liu, Gang Li, Junbo Zhang, Heinrich Dinkel, Yongqing Wang, Zhiyong Yan, Yujun Wang, and Bin Wang. Enhancing automated audio captioning via large language models with optimized audio encoding. In *Proceedings of the 25th Interspeech Conference (interspeech)*, pp. 1135–1139, 2024a.
- Jizhong Liu, Gang Li, Junbo Zhang, Chenyu Liu, Heinrich Dinkel, Yongqing Wang, Zhiyong Yan, Yujun Wang, and Bin Wang. Leveraging ced encoder and large language models for automated audio captioning. *Proceedings of the DCASE Challenge*, pp. 1–4, 2024b.
- Richard F Lyon. *Human and machine hearing*. Cambridge University Press, 2017.
- Ziyang Ma, Zhisheng Zheng, Jiaxin Ye, Jinchao Li, Zhifu Gao, Shiliang Zhang, and Xie Chen. Emotion2vec: Self-supervised pre-training for speech emotion representation. *arXiv preprint arXiv:2312.15185*, 2023.
- Ziyang Ma, Yinghao Ma, Yanqiao Zhu, Chen Yang, Yi-Wen Chao, Ruiyang Xu, Wenxi Chen, Yuanzhe Chen, Zhuo Chen, Jian Cong, et al. Mmar: A challenging benchmark for deep reasoning in speech, audio, music, and their mix. *arXiv preprint arXiv:2505.13032*, 2025.
- Ilaria Manco, Benno Weck, Seungheon Doh, Minz Won, Yixiao Zhang, Dmitry Bogdanov, Yusong Wu, Ke Chen, Philip Tovstogan, Emmanouil Benetos, et al. The song describer dataset: a corpus of audio captions for music-and-language evaluation. *arXiv preprint arXiv:2311.10057*, 2023.
- Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:3339–3354, 2024.
- Gabriel Mittag, Babak Naderi, Assmaa Chehadi, and Sebastian Möller. Nisqa: A deep cnn-self-attention model for multidimensional speech quality prediction with crowdsourced datasets. In *Proceedings of the 22nd Interspeech Conference (interspeech)*, pp. 2127–2131, 2021.
- Prabhat Pandey, Rupak Vignesh Swaminathan, KV Girish, Arunasish Sen, Jian Xie, Grant P Strimel, and Andreas Schwarz. Sift-50m: A large-scale multilingual dataset for speech instruction fine-tuning. *arXiv preprint arXiv:2504.09081*, 2025.

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Christopher J Plack. *The sense of hearing*. Routledge, 2023.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, pp. 28492–28518, 2023.
- Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, Ju-Chieh Chou, Sung-Lin Yeh, Szu-Wei Fu, Chien-Feng Liao, Elena Rastorgueva, François Grondin, William Aris, Hwidong Na, Yan Gao, Renato De Mori, and Yoshua Bengio. SpeechBrain: A general-purpose speech toolkit, 2021. arXiv:2106.04624.
- Chandan KA et al. Reddy. Dnsmos: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6493–6497. IEEE, 2021.
- Chandan KA et al. Reddy. Dnsmos p. 835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 886–890. IEEE, 2022.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- Paul K. Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quiry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. AudioPaLM: A large language model that can speak and listen. *arXiv preprint arXiv:2306.12925*, 2023.
- S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Raman Duraiswami, Sreyan Ghosh, and Dinesh Manocha. Mmau: A massive multi-task audio understanding and reasoning benchmark. In *Proceedings of the 13th International Conference on Learning Representations (ICLR)*, pp. 1–36, 2025.
- Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang, Ruihua Zhang, Daochen Shi, Qiqi Xiang, and Yemin Shi. LLASM: Large language and speech model. *arXiv preprint arXiv:2308.15930*, 2023.
- Luoyi Sun, Xuenan Xu, Mengyue Wu, and Weidi Xie. Auto-acd: A large-scale dataset for audio-language representation learning. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 5025–5034, 2024.
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. In *Proceedings of the 20th International Conference on Learning Representations (ICLR)*, pp. 1–23, 2024.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575, 2015.
- Johannes Wagner, Andreas Triantafyllopoulos, Hagen Wierstorf, Maximilian Schmitt, Felix Burkhardt, Florian Eyben, and Björn W Schuller. Dawn of the transformer era in speech emotion recognition: closing the valence gap. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10745–10759, 2023.

- Bin Wang, Xunlong Zou, Geyu Lin, Shuo Sun, Zhuohan Liu, Wenyu Zhang, Zhengyuan Liu, Aiti Aw, and Nancy Chen. Audiobench: A universal benchmark for audio large language models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 4297–4316, 2025.
- Chen Wang, Minpeng Liao, Zhongqiang Huang, Jinliang Lu, Junhong Wu, Yuchen Liu, Chengqing Zong, and Jiajun Zhang. BLSP: Bootstrapping language-speech pre-training via behavior alignment of continuation writing. *arXiv preprint arXiv:2309.00916*, 2023.
- Boyong Wu, Chao Yan, Chen Hu, Cheng Yi, Chengli Feng, Fei Tian, Feiyu Shen, Gang Yu, Haoyang Zhang, Jingbei Li, et al. Step-audio 2 technical report. *arXiv preprint arXiv:2507.16632*, 2025.
- Mengyue Wu, Heinrich Dinkel, and Kai Yu. Audio caption: Listen and tell. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 830–834. IEEE, 2019.
- LLM-Core-Team Xiaomi. Mimo-audio: Audio language models are few-shot learners, 2025. URL <https://github.com/XiaomiMiMo/MiMo-Audio>. GitHub repository.
- Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, et al. Qwen2. 5-omni technical report. *arXiv preprint arXiv:2503.20215*, 2025a.
- Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting He, Xinfa Zhu, et al. Qwen3-omni technical report. *arXiv preprint arXiv:2509.17765*, 2025b.
- LI Yizhi, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenchua Lin, Anton Ragni, Emmanouil Benetos, et al. Mert: Acoustic music understanding model with large-scale self-supervised training. In *Proceedings of the 12th International Conference on Learning Representations (ICLR)*, pp. 1–24, 2023.
- Yi Yuan, Dongya Jia, Xiaobin Zhuang, Yuanzhe Chen, Zhuo Chen, Yuping Wang, Yuxuan Wang, Xubo Liu, Xiyuan Kang, Mark D Plumbley, et al. Sound-vecaps: Improving audio generation with visually enhanced captions. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2025.
- Tanghaoran Zhang, Yue Yu, Xinjun Mao, Yao Lu, Zhixing Li, and Huaimin Wang. Fense: A feature-based ensemble modeling approach to cross-project just-in-time defect prediction. *Empirical Software Engineering*, 27(7):162, 2022.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- Juan Zuluaga-Gomez, Sara Ahmed, Danielius Visockas, and Cem Subakan. Commonaccent: Exploring large acoustic pretrained models for accent classification based on common voice. In *Proceedings of the 24th Interspeech Conference (interspeech)*, pp. 5291–5295. ISCA, 2023.

A LLM USAGE

In the preparation of this manuscript, we utilized Gemini 2.5 Pro model (accessed in July 2025) primarily for proofreading and grammatical corrections. The tool was used to improve the clarity and readability of the text. All authors have reviewed and edited the final manuscript and take full responsibility for its content.

B VOCABULARY SIZE OF AUDIO CAPTIONING TESTSET

This section introduces vocabulary size comparison to demonstrate the lexical diversity of MECAT-Caption. The following table indicates that the vocabulary size of MECAT-Caption contains about 4-17 times more words than the existing dataset.

Dataset	# Vocab
AudioCaps	5,581
AudioCaps-Enhanced	1,260
AutoACD	3,517
Clotho	1,852
StrongDescriber	2,726
LPMusicCaps-MTT	1,666
MECAT-Caption	22,595

C DEPLOYED ACOUSTIC MODELS IN PROCESSING PIPELINE

This section introduces the acoustic models deployed in our processing pipeline. These models are categorized into Content-Specific models (including Speech, Music, and Sound analysis) and Content-Unrelated models (Environment analysis), each designed to handle different aspects of audio understanding tasks.

Category	Subcategory	Model	Analysis Task
Content Specific	Speech	Speechbrain-ECAPA(Ravanelli et al., 2021)	Language Recognition
		Whisper Large v2 (Radford et al., 2023)	Auto Speech Recognition
		Pyannote-SD 3.1 (Bredin, 2023)	Speaker Diarization
		Emotion2Vec (Ma et al., 2023)	Speaker Emotion Recognition
		Audeering-DSER (Wagner et al., 2023)	Dimensional Speaker Emotion Recognition
		Audeering-AGR (Burkhardt et al., 2023)	Age and Gender Recognition
		CommonAccent (Zuluaga-Gomez et al., 2023)	English Accent Recognition
	Music	Music Structure Analyzer (Kim & Nam, 2023)	Tempo & Structure
		Music2Emo (Kang & Herremans, 2025)	Emotion (Sentiment Polarity and Mood)
		MERT (Yizhi et al., 2023)	Musical Descriptor
		ByteSep (Kong et al., 2021)	Music Separation
Content Unrelated	Sound	Audio Flamingo 2 (Ghosh et al., 2025)	AudioLLM
		CED (Dinkel et al., 2024a)	Sound Event Recognition
	Environment	DNSMOS (Reddy, 2021; 2022)	Noisy Speech Quality Assessment
		NISQA V2.0 (Mittag et al., 2021)	Multidimensional Speech Quality Assessment
		SHAART (Hawley, 2023)	Reverberation

D LLM AUDIO ANALYSIS SYNTHESIS PROMPTS

Act as an expert audio analysis synthesizer to process multi-model JSON outputs through this workflow

Step 1: Multi-Domain Data Specifications

1.1 Multi-Domain Input Integration

- a) Speech: Speech recognition, speech emotion, speaker diarization and so on
- b) Music: Structure analysis, technical descriptors, emotion and so on
- c) Sound: Event detection timestamps, classifications
- d) Environment: Acoustic characteristics, interference markers
- e) Meta-info: Title and description of original video where audio clip was extracted

1.2 Data Integrity Challenges

- a) Missing fields
- b) Contradictory model outputs
- c) Confidence score variances

Step 2: Technical and Analytical Limitations

2.1 Model and System Constraints

- a) No speech recognition ability in audio captioning models (e.g., audio-flamingo variants)
- b) Accuracy disparities across the analyzed domains
- c) Potential conflicting information between models

2.2 Audio Content Heterogeneity

- a) Hybrid audio types (e.g., speech, music, sound-event, environment)
- b) Variable audio properties (e.g., clip lengths or quality)
- c) Reliable topic or domain, but absent or non-relevant details in Meta-info

Step 3: Audio Analysis Workflow

3.1 Salient Feature Identification

3.1.1 Identify dominant characteristics of this audio:

What makes this specific audio clip unique according to the analysis? Examples include:

- a) Specific spoken phrases
- b) Dominant musical styles or moods
- c) Significant sound events
- d) The overall acoustic scene
- e) Notable quality issues
- f) Complex interplay of elements

3.1.2 Supporting Evidence Extraction:

Gather the key details describing these salient features from the relevant JSON fields

3.2 Synthesis Rules

3.2.1 Generation Rules:

- a) Critically weigh evidence from different fields, considering inaccuracies or conflicts and accounting for domain-specific limitations
- b) Prioritize information most reliable or central to the audio’s character based on overall data patterns
- c) Carefully identify conflicting information between fields and avoid mentioning conflicting aspects in the final caption. Focus only on consistent and unopposed information. Do not invent details not present in the data
- d) Crucial Constraint:
 - The final generated text must strictly describe only the analyzed content of the audio segment itself
 - It must not refer to the topic, title, description, or inferred subject matter from the overall video metadata
 - Avoid phrases like ”in a clip from a video about...” or similar references to the source video’s topic
 - Prohibit using parentheses to provide detailed explanation in any output, e.g., Moderate tempo (88 BPM)

3.2.2 Perspective Rules:

ALL answers must be created from the perspective of someone who ONLY LISTENED to the audio without any technical/model references or quantitative metrics (e.g., BPM, MOS, etc.)

3.2.3 Evaluation Rules:

Assign a confidence level (High or Low) based on the following aspects:

- a) Consistency: Are the different analyses in the JSON generally consistent or contradictory? High consistency increases confidence
- b) Completeness: Is key information present? (Fewer gaps = higher confidence)
- c) Clarity: How clearly does the consistent data point to the audio’s nature? (Less ambiguity in reliable data = higher confidence)
- d) Metadata Context Usefulness: How relevant and useful was the overall video metadata in confirming or contextualizing findings from the clip’s direct analysis?

3.3 Caption Development Framework

3.3.1 Systematic Caption

- a) Short (< 15 words):
 - Protocol: Primary domain characteristics + Most prevalent characteristic from cross-model correlation
 - Example: Blues guitar performance at live concert with audience reactions
- b) Long (1-2 sentences):
 - Protocol: Primary domain + significant secondary elements + notable quality factors
 - Example: A live concert recording featuring guitar with crowd cheers, despite occasional microphone static

3.3.2 Content-Focused Caption

- a) Speech:
 - Protocol: ASR content + paralinguistic context

- Example: Two speakers discussing jazz history, with piano accompaniment
 - b) Music:
 - Protocol: Technical descriptors + performance context
 - Example: Upbeat electronic track with distant traffic noise
 - c) Sound:
 - Protocol: Event taxonomy + spatial relationships
 - Example: Office environment with printer hum and keyboard typing, mild echo present
- ### 3.3.3 Content-Unrelated Caption
- a) Environment:
 - Protocol: Acoustic properties + interference profile
 - Example: Studio recording with noticeable background interference
- ### 3.3.4 Caption Variants
- a) Lexical substitution (WordNet-based synonyms)
 - b) Structural reordering (active/passive voice)
 - c) Descriptive equivalence ('crowd cheers' → 'audience applause')
- ### 3.3.5 Null Handling
- When no domain-specific elements are detected:
- a) Use explicit 'None' declaration in content field
 - b) Generate null statement variants (e.g., 'No discernible speech content', 'Musical elements appear absent')
- ## 3.4 QUESTION-ANSWERING DESIGN
- ### 3.4.1 Content Categories
- Include questions across:
- a) Direct Perception (sound type, volume, duration)
 - b) Sound Characteristics (timbre, rhythm, frequency characteristics)
 - c) Environmental Perception (recording setting, echo, background noise)
 - d) Quality Assessment (clarity, interference factors)
 - e) Inference and Judgment (sound source, generation method, object properties)
 - f) Application Context (use cases, semantic meaning)
- ### 3.4.2 Difficulty Levels
- Include a mix of:
- a) Basic: Direct descriptive questions (e.g., 'What sound is heard?')
 - b) Intermediate: Analytical questions (e.g., 'What are the characteristics of this sound?')
 - c) Advanced: Inferential questions (e.g., 'In what environment was this recorded?')
 - d) Complex: Comprehensive judgment questions (e.g., 'Based on the sound, what is the most likely material?')
- ### 3.4.3 Question Distribution
- Basic (25%) — Intermediate (35%) — Advanced (25%) — Complex (15%)
- ### 3.4.4 Question Variety
- Include:

- a) Ensure questions cover all listed categories
- b) Avoid repetitive question patterns or formats
- c) Include both yes/no questions and open-ended questions
- d) Include some questions about what is NOT present in the audio
- e) Include some comparative questions (e.g., 'Does this sound more like X or Y?')

3.4.5 Cognitive Levels

Include:

- a) Include questions requiring simple recognition
- b) Include questions requiring analysis of components
- c) Include questions requiring synthesis of information
- d) Include questions requiring evaluation or judgment

Step 4: Structured Output Specification (JSON Format)

Confidence: High/Low

Possible Conflicts: None or list of conflicting fields

Reasoning: 2-3 line evaluation considering model consensus and data quality

Short-Caption: Single-sentence essence

Short-Caption-Variants-1: Paraphrased version 1

Short-Caption-Variants-2: Paraphrased version 2

Main-Caption: Integrated summary

Main-Caption-Variants-1: Paraphrased version 1

Main-Caption-Variants-2: Paraphrased version 2

Speech-Captions: Speech-focused analysis or NONE

Speech-Caption-Variants-1: Paraphrased version 1

Speech-Caption-Variants-2: Paraphrased version 2

Music-Captions: Music-focused analysis or NONE

Music-Caption-Variants-1: Paraphrased version 1

Music-Caption-Variants-2: Paraphrased version 2

Sound-Captions: Sound-focused analysis or NONE

Sound-Caption-Variants-1: Paraphrased version 1

Sound-Caption-Variants-2: Paraphrased version 2

Environment-Caption: Environment-focused analysis

Environment-Caption-Variants-1: Paraphrased version 1

Environment-Caption-Variants-2: Paraphrased version 2

QA-Pair-1-id: 1 or None

QA-Pair-1-difficulty: basic, intermediate, advanced, or complex

QA-Pair-1-category: direct perception, sound characteristics, environmental perception, quality assessment, inference judgment, or application context

QA-Pair-1-question: question content

QA-Pair-1-answer: answer content

1134 QA-Pair-2-id: 2 or None
1135
1136 QA-Pair-2-difficulty: basic, intermediate, advanced, or complex
1137 QA-Pair-2-category: direct perception,sound characteristics, environmental perception,quality
1138 assessment, inference judgment,or application context
1139 QA-Pair-2-question: question content
1140
1141 QA-Pair-2-answer: answer content
1142 // ... 3 more QA pairs following the same pattern
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

E EMBEDDING-LEVEL TF-IDF CALCULATION FOR DATE

E.1 RATIONALE FOR EMBEDDING-LEVEL TF-IDF

Classic term frequency-inverse document frequency (TF-IDF) relies on discrete, hard-count token occurrences to calculate term frequency (TF). However, in natural language, the semantic relationship between words (e.g., "dog" and "canine") is lost when treating them as independent tokens.

The Discriminability based Audio Task Evaluation (DATE) metric enhances the representational power of sentence embeddings by incorporating an Embedding-level TF-IDF weighting scheme. This method leverages the semantic information encoded in word embeddings to compute a soft, non-integer Term Frequency, thereby improving the quality of the resulting sentence representation for downstream similarity and discrimination calculations.

E.2 METHOD OF CALCULATION

The Embedding-level TF-IDF weight for a word w in a sentence s is calculated as the product of its semantic-aware Term Frequency (TF_{emb}) and its Inverse Document Frequency (IDF_{emb}):

$$\text{TF-IDF}_{\text{emb}}(w, s) = \text{TF}_{\text{emb}}(w, s) \times \text{IDF}_{\text{emb}}(w)$$

The calculation proceeds in three main steps:

Semantic-aware Term Frequency (TF_{emb}) Instead of counting a word's exact occurrences (which would result in an integer count), TF_{emb} is calculated by measuring the average similarity of the word's embedding to the embeddings of all other words within the same sentence. This accounts for semantic context and relatedness.

- **Word Embeddings:** Word-level embeddings are generated for all tokens in the corpus with the Sentence-Bert.
- **Word-to-Word Similarity Matrix (WordSim):** A full WordSim matrix is computed by taking the dot product of the L_2 -normalized embeddings (E) of all unique non-padding words in the batch.

$$\text{WordSim} = \text{Normalize}(E) \cdot \text{Normalize}(E)^T$$

- **Term Frequency Calculation:** The semantic Term Frequency for a word w_i in a sentence s is calculated by summing the squared similarity scores with all other words w_j in that sentence (excluding special tokens like [CLS] and [SEP]):

$$\text{TF}_{\text{emb}}(w_i, s) = \sum_{w_j \in s, j \neq i} \text{WordSim}(i, j)^2$$

- **Result:** This calculation yields a non-integer TF_{emb} value, where words semantically central to the sentence receive a higher score.

Embedding-aware Inverse Document Frequency (IDF_{emb}) The Inverse Document Frequency (IDF) component measures a word's uniqueness across the entire document corpus. In the Embedding-level approach, the document frequency (DF) is calculated based on the similarity of a unique word's embedding to the embeddings of all tokens across all documents.

- **Word-to-Document Similarity Matrix (Word2DocSim):** A Word2DocSim matrix is calculated by taking the dot product between the unique word embeddings and the embeddings of all tokens (word pieces) in the corpus.
- **Document Frequency Accumulation:** The document frequency (DF) for a word w is accumulated across the entire corpus by summing the squared Word2DocSim values. This accumulation is performed sentence-by-sentence based on token availability.
- **IDF Calculation:** The final IDF_{emb} is calculated using a log-normalization formula, where N is the corpus size:

$$\text{IDF}_{\text{emb}}(w) = \log \left(\frac{N + 1}{\text{DF}(w) + 1} \right) + 1$$

1242 **Final Weighting and Normalization** The TF_{emb} is multiplied by the IDF_{emb} to get the final raw
1243 TF-IDF weight for each non-special token. These weights are then normalized to ensure stability.
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

F HUMAN VALIDATION OF DATA INTEGRITY AND METRIC ALIGNMENT

This section provides comprehensive details on the human evaluation studies conducted to validate the integrity of the MECAT dataset and to assess the alignment of the proposed DATE metric with human judgments.

F.1 ANNOTATION QUALITY VALIDATION VIA HUMAN PREFERENCE

A Human Preference A/B test ($N = 150$ caption pairs spanning all domains) was conducted to verify the quality and trustworthiness of the MECAT pipeline-generated references. Evaluators were instructed to select the better caption based on the stringent criterion: “accuracy first, then level of detail.”

MECAT references (A) were compared against three opponent types (B) to probe specific quality aspects: Safe Captions (generic descriptions, testing discriminability); Wrong Captions (factually incorrect references, testing accuracy); and Human References (expert-written ground truth, establishing the quality ceiling).

The results are presented in Table F.1. MECAT references were strongly favored over both Safe and Wrong captions ($> 94\%$ win rates), confirming the pipeline’s effectiveness in mitigating vague or factually incorrect outputs. Crucially, the quality of MECAT references was found to be statistically on par with human-written references (56.9% win rate, as the 95% confidence interval spans 50%).

Table F.1: Human Preference A/B Validation Results ($N = 150$ pairs). Win rates indicate preference for the MECAT Reference (A).

Opponent Type (B)	Size	Reference (A) Win Rate	95% CI (Wilson)
Overall	150	82.7%	[75.8%, 87.9%]
Safe Captions	52	94.2%	[84.4%, 98.0%]
Wrong Captions	47	97.9%	[88.9%, 99.6%]
Human References	51	56.9%	[43.3%, 69.5%]

F.2 ALIGNMENT OF DATE WITH HUMAN JUDGMENTS

The alignment of the DATE metric with human preferences was assessed using the 150 A/B caption pairs detailed in Appendix F.1, where listener choices established the human “gold standard”.

For this validation, the human-preferred caption and the non-preferred caption from each A/B pair were separately treated as a candidate hypothesis. The DATE score was then computed for each candidate against the audio segment’s original ground truth reference.

A comparison of mean DATE scores, presented in Table F.2, reveals a substantial difference: captions selected by human evaluators received substantially higher mean DATE scores (90.9) compared to non-preferred captions (49.3). This large margin ($\Delta \approx 41.6$) robustly demonstrates that DATE is highly correlated with human preference regarding accuracy and detail, thereby validating its utility as a fine-grained evaluation tool.

Table F.2: Alignment between DATE Score and Human Preference ($N = 150$ A/B pairs).

Caption Group	DATE
Human-Preferred Captions	90.9
Non-Preferred Captions	49.3

G LLM-AS-JUDGE PROMPTS IN EVALUATION

This section provides the prompt template required for LLM-as-Judge method. The evaluation *tasks* primarily include audio captioning and audio question-answering. In the template, the *description*, *subtask*, and *scoring_aspects* parameters can be referenced from the corresponding columns in the task table above, while *ref_texts* represents the samples to be evaluated.

G.1 Evaluation Prompt Template

You are tasked with evaluating if a set of candidate $\{\text{tasks}\}$ responses accurately addresses the same audio as a reference set of answers. You will focus on the $\{\text{description}\}$ for the subtask ' $\{\text{subtask}\}$ '.

Evaluation Steps:

- a) First, carefully compare the candidate answers with the reference answers
- b) Assess the accuracy and precision of how the audio characteristics are captured in the responses, then provide a 0-10 fine-grained score:
 - 10 = perfect match with the reference content
 - 0 = completely wrong
- c) Provide detailed scoring reasoning, explaining why you gave this score

Scoring Aspects: $\{\text{scoring_aspects}\}$

Score rubric (0-10 Scale):

- points 9-10: Excellent - Highly accurate, comprehensive, well-expressed
- points 8: Very Good - Accurate with minor gaps, clear expression
- points 7: Good - Mostly accurate, some missing details
- points 6: Acceptable - Basic accuracy, meets minimum HIGH standard
- points 4-5: Below Standard - Some correct elements but major issues
- points 2-3: Poor - Limited accuracy, significant problems
- points 0-1: Very Poor - Major errors or completely incorrect

You need to evaluate the following sample: $\{\text{ref_texts}\}$

Please return JSON-formatted evaluation results for the sample.

Return format (strict JSON array):

sample_id: sample ID

subtask: subtask_name

fine_score: <numerical value 0-10>

reasoning: detailed scoring rationale, including comparative analysis with reference answers

Table G.2: Category, subcategory, descriptions and scoring aspects of captioning evaluation with LLM-as-Judge method

Category	Subcategory	Description & Scoring Aspects
Systemic	Short	Description: quality of short audio descriptions Scoring Aspects: a) accuracy of core content capture (most important) b) conciseness and completeness of expression c) semantic consistency with reference descriptions
		Description: quality of detailed audio descriptions Scoring Aspects: a) comprehensiveness and richness of description details b) accuracy of detailed descriptions c) logical structure and expression coherence
		Description: accuracy of speech content recognition Scoring Aspects: a) accuracy rate of speech content recognition b) accurate description of speaker characteristics (gender, accent, etc.) c) description of speech quality and environment
	Music	Description: quality of music content description Scoring Aspects: a) accuracy of music type, style, and rhythm identification b) identification of instruments and musical elements c) description of musical emotion and atmosphere
		Description: accuracy of sound event identification Scoring Aspects: a) accurate identification and classification of sound sources b) description of sound occurrence timing and duration c) description of sound intensity, pitch and other characteristics
	Environment	Description: accuracy of environment and recording quality description Scoring Aspects: a) identification of recording environment (indoor/outdoor, space size, etc.) b) assessment of audio technical quality (distortion, noise, etc.) c) description of environmental atmosphere and background characteristics

Table G.3: Category, subcategory, descriptions and scoring aspects of question-answering evaluation with LLM-as-Judge method

Category	Subcategory	Description & Scoring Aspects
Perception	Direct Perception	Description: accuracy of direct audio content identification
		Scoring Aspects:
		a) correct identification of primary audio elements (most important) b) accurate detection of presence/absence of specific sounds c) precise recognition of obvious audio features and events
Analysis	Sound Characteristics	Description: quality of sound property analysis
		Scoring Aspects:
		a) accurate description of sound attributes (pitch, volume, timbre, etc.) b) correct identification of sound sources and their properties c) precise characterization of audio dynamics and patterns
	Quality Assessment	Description: accuracy of audio quality evaluation
		Scoring Aspects:
		a) correct assessment of technical audio quality (clarity, distortion, etc.) b) accurate evaluation of recording conditions and fidelity c) appropriate judgment of audio production quality
	Environment Reasoning	Description: quality of environmental context inference
		Scoring Aspects:
		a) accurate inference of recording location and setting b) correct identification of spatial and acoustic properties c) logical deduction of environmental factors affecting audio
Reasoning	Inference Judgment	Description: accuracy of complex audio analysis and reasoning
		Scoring Aspects:
		a) correct interpretation of implicit audio information b) accurate temporal reasoning and sequence understanding c) logical inference of causality and relationships between audio elements
	Application Context	Description: relevance and appropriateness of contextual understanding
		Scoring Aspects:
		a) accurate understanding of audio's intended purpose or context b) appropriate application of domain-specific knowledge c) correct interpretation of cultural, social, or professional context

H VALIDATION OF LLM-AS-JUDGE AS A REFERENCE METRIC

To validate the effectiveness of LLM-as-Judge as a reference metric, we assessed its performance on three distinct sets of responses with varying quality levels: Right (detailed and accurate rephrasings of the ground-truth reference), Safe (generic, vague descriptions, e.g., "A man is speaking" for all speech-only audio), and Wrong (factually incorrect references randomly selected from other samples). As shown in following table, our analysis confirms that LLM-as-Judge method serves as a reliable evaluator. It successfully distinguishes between the quality tiers, with mean scores consistently following the expected Right > Safe > Wrong order for both captioning and QA tasks. Furthermore, its inter-rater reliability, measured by Fleiss' Kappa (κ), is substantial for QA ($\kappa = 0.73$) and moderate for captioning ($\kappa = 0.43$). However, the significant practical limitations of LLM-as-Judge method—including high computational cost, slow speed, and sensitivity to prompt engineering—motivate our development of the DATE metric as an efficient and scalable alternative.

Type	Mean		Fleiss' Kappa (κ)	
	Caption	QA	Caption	QA
Right	0.78	0.97	0.68	0.74
Safe	0.24	-	0.17	-
Wrong	0.13	0.12	0.45	0.72
Overall	-	-	0.43	0.73

I TASK-SPECIFIC PROMPTS FOR LALM IN MECAT TASKS

This section details the prompt strategies employed for Large Audio-Language Models (LALMs) during the MECAT-Caption evaluation. For specific Audio-focused LALMs that require specialized instruction formats—namely Audio-Flamingo 2 (Ghosh et al., 2025) and Kimi-Audio (KimiTeam et al., 2025)—the exact prompt templates are provided in Table C.1. For the remaining LALMs, the prompts are standardized as shown in Table C.2. This category encompasses a diverse range of architectures, including other Audio-focused models such as Mimo-Audio (Xiaomi, 2025), Step-Audio-2-mini (Wu et al., 2025), and Baichuan-Audio (Li et al., 2025); Omni LALMs represented by the Qwen-Omni series (Xu et al., 2025a;b) and Baichuan-Omni (Li et al., 2024); Multimodal LALMs specifically Phi-4-Multimodal (Abouelenin et al., 2025); and the state-of-the-art Gemini series (Comanici et al., 2025). Regarding system configurations, the system prompt was explicitly set to “*You are a helpful assistant*” for the **Qwen2.5-Omni** models, while the default system prompts were retained for all other architectures.

Regarding the MECAT-QA task, the prompt for each sample consists solely of the corresponding question, without additional task-specific templates.

Table I.1: Prompts for Audio-Flamingo2 and Kimi-Audio models in caption task

Category	Subcategory	Prompt
Systematic	Short	Provide a caption for this audio within 15 words
	Long	Provide a caption for this audio within 1-2 sentences
Content-Specific	Speech	Provide a caption for the speech content in this audio
	Music	Provide a caption for the music content in this audio
	Sound	Provide a caption for general sound excluding speech and music
Content-Unrelated	Environment	Provide a caption for quality or acoustic environment for this audio

Table I.2: Prompts for remaining models in caption task

Category	Subcategory	Prompt
Systematic	Short	Listen to the audio and provide a caption for this audio within 15 words
	Long	Listen to this audio and provide a caption for this audio within 1-2 sentences
Content-Specific	Speech	Listen to the audio and provide a caption describing the speech content in this audio
	Music	Listen to the audio and provide a caption for the music content in this audio
	Sound	Listen to the audio and provide a general sound excluding speech and music
Content-Unrelated	Environment	Listen to this audio and provide a caption for quality or acoustic environment for this audio

J FINE-GRAINED ANALYSIS OF SPEECH CAPTIONING

This section presents a detailed performance of all models on the pure speech subset of MECAT-Caption task. Table J.1 presents the results across three dimensions: similarity, discriminability, and DATE.

Table J.1: Model performance and rankings on the **Pure Speech** subset of MECAT-Caption, evaluated via Similarity, Discriminability, and DATE metrics. **Bold** and underlined values denote the best and second-best results, respectively, with rankings in gray. [†] indicates that the model or its previous version (Audio Flamingo 2) was explicitly used in the data construction process.

Type	Model	Similarity		Discriminability		DATE	
		Score	Rank	Score	Rank	Score	Rank
Caption-Only	Pengi	26.6	13	27.8	14	27.2	14
	EnClap	28.7	11	31.9	13	30.2	12
LALM	Phi-4-Multimodal-Instruct	26.6	13	27.2	15	26.9	15
	Kimi-Audio-7B-Instruct	25.6	15	36.2	11	30.0	13
	Baichuan-Audio-Instruct	37.2	7	60.3	4	46.0	5
	Audio Flamingo 2 [†]	28.5	12	32.8	12	30.5	11
	Audio Flamingo 3	46.6	1	52.3	7	49.3	4
	Mimo-Audio-Instruct	42.5	4	49.7	8	45.8	6
	Step-Audio-2-mini	36.6	8	55.8	6	44.2	7
	Baichuan-Omni	34.9	10	57.7	5	43.5	8
	Qwen2.5-Omni 3B	37.3	6	49.4	9	42.5	9
	Qwen2.5-Omni 7B	35.3	9	45.9	10	39.9	10
	Qwen3-Omni	41.0	5	64.7	3	50.2	3
	Gemini-2.5-Flash	<u>45.8</u>	2	<u>77.2</u>	2	57.5	1
	Gemini-2.5-Pro	44.3	3	78.4	1	<u>56.6</u>	2

K MODEL PERFORMANCE OF SIMILARITY ON MECAT TASKS

This section presents the complete Similarity scores for all models evaluated on MECAT, serving as a comparative reference for the DATE metrics reported in the main text (see Tables 2 and 3).

Table K.1: Model performance (Similarity %) on MECAT-Caption. **Bold** indicates the best performance, and underline indicates the second best. [†] indicates that the model or its previous version (Audio Flamingo 2) was explicitly used in the data construction process.

Type	Model	Systemic		Content-Specific						Content	Score _{Cap}
		Long	Short	Speech		Music		Sound		Unrelated	
				Pure	Mixed	Pure	Mixed	Pure	Mixed	Env	
Caption-Only	Pengi	37.5	41.0	26.6	29.2	39.6	11.8	35.4	16.2	17.8	29.5
	EnClap	40.5	45.0	28.7	29.5	39.3	15.0	41.2	17.3	17.9	31.6
LALM	Phi-4-Multimodal-Instruct	45.4	40.3	26.6	31.7	41.5	26.2	29.5	25.7	37.3	37.4
	Kimi-Audio-7B-Instruct	40.8	45.7	25.6	27.1	39.5	16.2	35.8	19.4	16.7	30.8
	Baichuan-Audio-Instruct	33.0	28.2	37.2	35.0	36.4	24.7	45.0	29.9	47.1	36.1
	Audio Flamingo 2 [†]	43.8	43.3	28.5	33.7	43.1	30.3	41.0	24.7	45.4	39.4
	Baichuan-Omni	39.2	42.5	34.9	35.4	41.0	13.2	40.0	32.3	29.4	35.0
	Mimo-Audio-Instruct	49.9	49.4	42.5	43.5	47.5	19.9	44.5	27.6	27.2	41.2
	Audio Flamingo 3 [†]	49.6	49.6	46.6	47.5	50.6	26.4	44.6	28.3	31.7	43.5
	Qwen3-Omni	38.2	33.6	34.1	34.5	49.0	34.1	41.4	20.8	40.2	37.4
	Step-Audio-2-mini	44.1	47.8	36.6	37.3	45.9	36.0	36.4	24.9	41.4	41.2
	Qwen2.5-Omni 3B	48.3	45.3	37.3	37.5	50.7	34.7	46.6	34.1	47.8	44.1
	Qwen2.5-Omni 7B	<u>52.7</u>	46.2	35.3	37.5	39.2	33.1	45.2	32.1	41.0	43.4
	Gemini-2.5-Flash	56.1	53.5	45.8	<u>46.6</u>	59.1	44.3	50.7	36.4	48.9	51.0
	Gemini-2.5-Pro	50.8	<u>49.9</u>	44.3	45.7	<u>58.5</u>	44.6	<u>49.6</u>	<u>35.0</u>	51.9	<u>49.3</u>

Table K.2: Model Performance (Similarity %) on MECAT-QA. **Bold** indicates the best performance, and underline indicates the second best. [†] indicates that the model or its previous version (Audio Flamingo 2) was explicitly used in the data construction process.

Model	Perception	Analysis		Reasoning			Score _{QA}
	Direct Perception	Sound Characteristics	Quality Assessment	Environment Reasoning	Inference & Judgment	Application Context	
Kimi-Audio-7B-Instruct	37.5	32.5	19.2	37.5	38.8	33.8	33.2
Baichuan-Audio-Instruct	35.2	36.6	36.0	38.1	39.5	39.6	37.5
Audio Flamingo 2 [†]	39.1	39.0	37.4	41.3	35.5	35.8	38.0
Baichuan-Omni	36.8	36.1	35.4	39.1	38.5	39.4	37.6
Phi-4-Multimodal-Instruct	41.2	37.6	36.6	40.3	39.0	40.1	39.1
Mimo-Audio-Instruct	<u>50.9</u>	40.5	27.0	40.7	41.9	38.5	39.9
Step-Audio-2-mini	48.6	44.6	39.1	38.2	38.7	39.3	41.4
Audio Flamingo 3 [†]	46.0	41.4	38.6	43.5	43.2	40.9	42.3
Qwen2.5-Omni 3B	47.2	43.8	39.7	43.2	41.0	41.9	42.8
Qwen2.5-Omni 7B	49.7	43.8	<u>40.5</u>	44.1	42.5	41.9	43.8
Qwen3-Omni	52.3	44.8	41.2	45.2	44.7	45.2	<u>45.6</u>
Gemini-2.5-Flash	47.9	46.1	39.7	<u>46.2</u>	47.1	47.9	45.8
Gemini-2.5-Pro	47.4	<u>45.2</u>	39.0	46.9	<u>45.7</u>	<u>46.3</u>	45.1

L ANALYSIS OF SPEECH HALLUCINATION IN SILENT SEGMENTS

This section qualitatively evaluates the grounding capabilities of various models when presented with audio segments containing no discernible speech. We specifically analyze the "speech" output field for two representative silent cases (C1 and C2) to detect potential hallucinations in low-signal environments.

Table L.1: Speech caption generated by different models for two silence segments in MECAT.

Type	Model	Output
Caption-Only	Pengi	C1: a person is playing a guitar. C2: someone is opening and closing a door.
	EnClap	C1: A low frequency of a person talking in the distance C2: A helicopter flying in the distance followed by a man speaking
LALM	Phi-4-Multimodal-Instruct	C1: Audio caption: A rumble is heard. C2: Audio caption: A man is speaking with music in the background.
	Kimi-Audio-7B-Instruct	C1: a machine is running C2: a person is playing a drum set with a bass drum and a snare drum.
	Baichuan-Audio-Instruct	C1: The speaker is talking about how they are going through their day... [Truncated] C2: The speaker is asking for help with their website... [Truncated]
	Audio Flamingo 2	C1: there is no speech content in this audio. C2: the audio contains a single, continuous sound without any speech or music.
	Baichuan-Omni	C1: A woman is speaking with background noise. C2: Someone is speaking with authority, likely giving instructions or delivering an important message.
	Mimo-Audio-Instruct	C1: A male voice, speaking in a neutral tone... says 'Thank you' in English... C2: The audio contains a sequence of sound effects... with no speech.
	Audio Flamingo 3	C1: A male voice says 'Thank you' amidst the sound of a waterfall. C2: A female voice says 'Thank you' in a neutral tone.
	Qwen3-Omni	C1: A person is speaking. C2: The audio contains only a single, sustained, low-pitched electronic tone.
	Step-audio-2-mini	C1: The speech content is "Oh no, I'm sorry." C2: There is no speech in this audio.
	Qwen2.5-Omni 3B	C1: The audio contains a speech saying 'I'm gonna be a daddy'. C2: The audio contains a speech segment where a male voice says 'you' in a neutral tone.
	Qwen2.5-Omni 7B	C1: The audio contains a speech segment where someone is saying 'I'm going to go ahead and do that.' C2: The audio contains a speech segment in which the speaker says 'you'.
	Gemini-2.5-Flash	C1: No discernible speech is present in this audio. C2: No speech detected.
	Gemini-2.5-Pro	C1: There is no speech in this audio. C2: There is no speech in this audio.