



GAMA: A Large Audio-Language Model with Advanced Audio Understanding and Complex Reasoning Abilities

Anonymous ACL submission

Abstract

Perceiving and understanding non-speech sounds and non-verbal speech is essential to making decisions that help us interact with our surroundings. In this paper, we propose GAMA, a novel **General-purpose Large Audio-Language Model (LALM)** with **Advanced Audio Understanding and Complex Reasoning Abilities**. We build GAMA by integrating an LLM with multiple types of audio representations, including features from a custom Audio Q-Former, a multi-layer aggregator that aggregates features from multiple layers of an audio encoder. We fine-tune GAMA on a large-scale audio-language dataset, which augments it with audio understanding capabilities. Next, we propose *CompA-R* (Instruction-Tuning for **Complex Audio Reasoning**), a synthetically generated instruction-tuning (IT) dataset with instructions that require the model to perform complex reasoning on the input audio. We instruction-tune GAMA with *CompA-R* to endow it with complex reasoning abilities, where we further add a soft prompt as input with high-level semantic evidence by leveraging event tags of the input audio. Finally, we also propose *CompA-R-test*, a human-labeled evaluation dataset for evaluating the capabilities of LALMs on open-ended audio question-answering that requires complex reasoning. Through automated and expert human evaluations, we show that GAMA outperforms all other LALMs in literature on diverse audio understanding tasks by margins of 1%-84%. Further, GAMA IT-ed on *CompA-R* proves to be superior in its complex reasoning and instruction following capabilities¹.

1 Introduction

Large Language Models (LLMs) possess impressive abilities to understand and reason about the world through language (Zhao et al., 2023). While spoken language understanding tasks, like automatic speech recognition, have had a long history

¹We will open-source code and data on paper acceptance

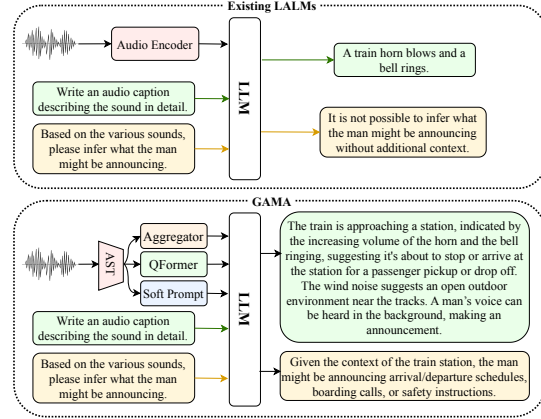


Figure 1: Comparison of existing LALMs (LTU by Gong et al. (2024) here) and GAMA. With improved audio understanding abilities (via diverse audio feature integration) and training on our proposed *CompA-R*, GAMA can provide more detailed captions of input audio and is also able to answer questions regarding it that demand complex reasoning.

of benefiting from language comprehension with (L)LMs (Watanabe et al., 2018; Hu et al., 2024), the ability to improve the perception and understanding of non-speech sounds and non-verbal speech through language has been less explored (from hereon we refer to these kinds of audios or sound as “audio” in the paper). Beyond visual and language perception, the ability to understand audio is unarguably important and necessary for autonomous agents to interact with the world.

Contrastive Language Audio Pre-training (CLAP) (Elizalde et al., 2023a) was one of the first Audio-Language Models (ALM) to improve audio understanding through a language interface. Following this, several attempts have been made to improve CLAP and its reasoning abilities (Ghosh et al., 2024b). On the other hand, Deshmukh et al. propose Pengi, a pre-trained decoder-only LLM coupled with an audio-encoder, that can solve all kinds of audio tasks by framing them as open-ended text-generation tasks. Similarly, Large Audio Language Models (LALMs) like

LTU (Gong et al., 2024) and SALMONN (Tang et al., 2024) follow a similar architecture and attempt to solve audio tasks by empowering the model with instruction following capabilities (Wei et al., 2022). Specifically, all audio tasks are first framed in instruction-response pairs. The model is then fine-tuned on these pairs to learn audio reasoning and, thereby, instruction following. As an emergent ability, these models also show remarkable capabilities in open-ended question answering by reasoning over the input audio. However, two significant problems still persist: (1) All these models employ simple connection modules between the audio encoder and the language decoder to enable the latter with audio understanding capabilities. This hinders comprehensive multimodal connection and alignment, thereby increasing the risk of hallucinations and leading to suboptimal performance (Liu et al., 2023a). (2) Complex reasoning with LALMs is still under-explored. While these models excel at audio event detection (in various forms like captioning, event classification, etc.) and information-seeking questions (e.g., close-ended audio questions like “How many birds are squawking?”), they fail to provide a faithful response for questions involving complex reasoning like “Identifying the context of laughter and its relationship with the automotive sounds in the recording. Draw a conclusion on the possible scenario occurring.”. We define complex reasoning for LALMs in Section 3.2 and show examples in Fig. 1 and Fig. 4.

Main Contributions. Our primary contributions are as follows:

- **A Novel LALM.** We introduce GAMA, an LALM with advanced audio understanding and complex reasoning abilities. To improve audio perception and understanding abilities, we propose integrating an LLM with multiple types of audio features that encode diverse aspects of information about the input audio. Specifically, we couple the output features from an Audio Q-Former and an Audio Spectrogram Transformer (AST) (Gong et al., 2021), where the AST is further equipped with an *aggregation module*. While the Audio Q-Former possesses impressive semantic generalization capabilities (Li et al., 2023), the AST possesses strong knowledge of surface-level audio properties. Additionally, inspired by the fact that different layers in audio models learn

audio information at different scales (Singla et al., 2022), the aggregation module aggregates the features from multiple layers of AST, which helps encode diverse knowledge. Both representations are passed through MLP layers that connect these features into the word embedding space before adding them as the prefix. As a result, GAMA possesses improved audio understanding capabilities by moving away from the simple coupling of audio encoders and linear layers commonly employed as connection modules to align the audio and textual modalities, which generally suffer from comprehensive multimodal alignment (Liu et al., 2023a). GAMA is first fine-tuned on a large-scale audio-language corpus, and the resulting model outperforms all other models on standard audio and music understanding benchmarks.

- **A Novel Instruction Tuning Dataset.** To endow an LALM with complex reasoning abilities, we propose `CompA-R`, a dataset synthetically generated with LLMs with multi-aspect information and human-written in-context examples. Specifically, we prompt GPT to synthesize an instruction-response pair by guiding it with various metadata related to the audio.
- **A Novel Evaluation Dataset.** To evaluate an LALM’s complex reasoning abilities, we develop `CompA-R-test`, a human-labeled benchmark. Specifically, `CompA-R-test` evaluates an LALM on open-ended AQA that demands complex reasoning over the audio. GAMA-IT (GAMA fine-tuned on `CompA-R`) shows significant improvements on `CompA-R-test` over all other baselines from literature.

2 Related Work

Large Multi-Modal and Audio-Language Models. Prior to the exploration of LLMs as efficient reasoners, encoder-based multi-modal language models, trained to learn a shared space between language and other modalities, have shown great promise. For example, CLAP, inspired by CLIP (Radford et al., 2021) in vision, showed state-of-the-art performance on audio-language tasks like retrieval, zero-shot classification, etc.

LLMs pre-trained at an incredible scale with the next token prediction objective implicitly com-

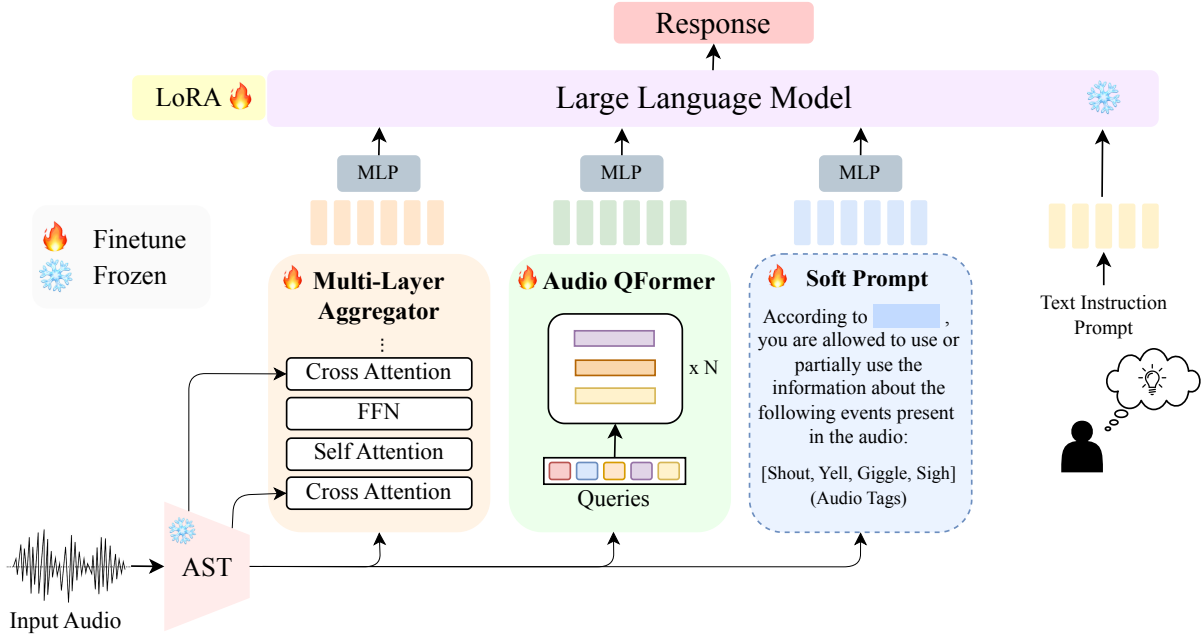


Figure 2: Illustration of GAMA. The GAMA architecture involves integrating various audio features with a text-only (pre-trained) LLM. The audio features are added as a prefix to the text instruction (by the user), and the LLM responds using a text output. We feed GAMA with 3 different types of audio features: ① The Multi-Layer Aggregator that takes as input features from the last and multiple other layers of the AST and calculates cross-attention to output a feature that encodes diverse surface features from audio. ② The pre-trained Audio Q-Former takes as input the last layer features of AST and outputs a feature that encodes the audio into a semantically rich and generalized space. ③ After fine-tuning, we instruction-tune this model on our proposed CompA-R and further feed the model with a soft prompt with audio tags from AST (with an added classification head). This additional information serves as high-level semantic knowledge to improve complex reasoning abilities.

press world knowledge in their parameters (Zhao et al., 2023). These models learn general-purpose representations, which can then be aligned with the desired response characteristics (Zhang et al., 2023). Instruction Tuning (IT), the process of fine-tuning an LLM with instruction-response pairs, has proved to be one of the most popular forms of alignment. Recent work shows that LLMs can also be instruction-tuned for multi-modal alignment. LLaVa (Liu et al., 2024), a pioneering work on multi-modal vision-language alignment, showed that fine-tuning an LLM on visual instruction-response pairs with additional vision features as prefix can endow the model with visual reasoning and understanding abilities. Several works following LLaVa improve aspects of LLMs and have achieved impressive performance on several vision-language tasks (Zhang et al., 2024). On the other hand, LALMs like LTU and SALMONN showed impressive performance on several audio-language tasks by reasoning over the audio. Though these models extensively evaluate several closed- and open-ended tasks, their ability to perform complex reasoning is largely under-explored.

Instruction Tuning and Complex Reasoning. IT-

based alignment has also shown significant improvements for LLMs on Natural Language Understanding tasks, unlocking impressive capabilities (Bubeck et al., 2023), suggesting that fine-tuning is key to building and improving LLM-based agents. Very recently, (Xu et al., 2024) and (Cui and Wang, 2024) show that well-curated IT data can improve various reasoning capabilities in LLMs, like logical, mathematical, complex reasoning, etc. More specifically, IT teaches LLMs better and more effective methods to reason about a problem, presented in the input instruction (like step-by-step reasoning (Kojima et al., 2022)).

3 Methodology

In the next sub-sections, we first describe the GAMA architecture and its components in detail, followed by fine-tuning GAMA on audio-language pairs, CompA-R creation, and instruction-tuning GAMA on CompA-R.

3.1 GAMA Architecture

Fig. 2 illustrates the architecture of GAMA. GAMA builds on the same base architecture proposed in prior works (Gong et al., 2024) but introduces sev-

213 eral novel components for improving audio percep- 263
214 tion. More specifically, we feed the pre-trained 264
215 LLM with features from multiple audio encoders, 265
216 including a pre-trained Audio-Q-Former and a pre- 266
217 trained AST that encode diverse audio knowledge. 267
218 Additionally, unlike prior work, we do not just use 268
219 the last layers of the AST but couple it with a multi- 269
220 layer aggregator that takes features from multiple 270
221 layers as input and outputs a feature that is aware of 271
222 various low-level and high-level properties of the 272
223 input audio. Finally, to endow the model with effec- 273
224 tive complex reasoning abilities, we employ AST 274
225 again to extract high-level semantic knowledge, i.e., 275
226 audio event tags, as supplementary information. 276

227 3.1.1 Audio Spectrogram Transformer (AST) 278

228 Audio Spectrogram Transformer (AST), was one 279
229 of the first attempts to model audio signals with 280
230 a pure Transformer network. We employ an AST 281
231 model fine-tuned on the AudioSet dataset. AST 282
232 has been employed as an audio encoder and a fea- 283
233 ture extractor in a wealth of prior works due to its 284
234 high informativeness (Gong et al., 2023, 2024). To 285
235 extract the last-layer features, we drop the audio 286
236 classification head and employ it only for event 287
237 classification for soft prompts. 288

238 3.1.2 Audio Q-Former 289

239 **Motivation.** Our primary goal is to integrate GAMA 290
240 with an audio encoder that possesses strong se- 291
241 mantic generalization capabilities for any input au- 292
242 dio. Prior work has extensively explored CLAP- 293
243 style training for learning audio-language encoders. 294
244 However, other methods and architectures have 295
245 rarely been explored. As a more powerful alter- 296
246 native, we explore the Q-Former architecture pro- 297
247 posed by (Li et al., 2023). We show in Table 1 298
248 that Q-Former outperforms other audio-language 299
249 models and achieves SOTA on audio understanding 300
250 benchmarks. 301

251 **Architecture.** The architecture of our Audio Q- 302
252 Former is based on the Querying Transformer pro- 303
253 posed in Li et al. (2023), which is initialized from 304
254 BERT (Devlin et al., 2018) and has Q querying 305
255 tokens. We employ AST as the audio encoder (in 306
256 place of the ViT-based vision encoder) and keep 307
257 the rest of the architecture the same. Similar to 308
258 the original implementation, we train the model in 309
259 two stages. For the first stage, we solve three tasks, 310
260 namely the Audio-Text Matching loss, the Audio- 311
261 Grounded Text Generation loss, and the Audio- 312
262 Text Contrastive Learning loss. For the second

stage, we employ LLaMa-2_{7B} as the language de- 263
264 coder and solve the language-modeling loss. For 265
266 training, we use 2.5M+ audio-caption pairs (de- 267
268 tailed in Section E.2). For architectural details, we 269
270 refer our readers to Li et al. (2023). 271

Training with Caption Augmentation. Addition- 272
273 ally, due to the lack of large-scale audio caption 274
275 pairs, we adopt a *caption-augmentation* method- 276
277 ology to augment the existing audios with diverse 278
279 additional captions. More specifically, we instruct 280
281 an LLM to generate k rewrites of the original cap- 282
283 tion. We employ two different prompts that rewrite 284
285 the input caption with two different objectives: 286

Prompts. For Prompt 1, our primary aim is that 287
288 the resultant rewrite should describe each acous- 289
290 tic event in the caption similarly but more vividly. 291
292 These augmentations help the model learn various 293
294 distinctive characteristics of the audio concepts cor- 295
296 responding to the acoustic events. For Prompt 1, 296
297 our primary aim is such that the resultant rewrite 298
299 should describe each acoustic event in the cap- 300
301 tion differently from the original caption. These 302
303 augmentations aid the model in understanding the 304
305 diverse linguistic expressions that can describe a 306
307 single audio concept. We show examples below: 308
309 (more examples in Table 13): 310

(1) **Original Caption:** Someone made a cool 311
312 vocal for a dubstep track.

(1) **Rewritten Caption by Prompt 1:** A capti- 313
314 vating vocal performance ignites the dubstep 315
316 track, delivering a hypnotic and enthralling 317
318 sound that reverberates through the air.

(1) **Rewritten Caption by Prompt 2:** The dub- 319
320 step track features a slick, stylish vocal per- 321
322 formance that adds a layer of sophistication 323
324 to its heavy beats and basslines.

(2) **Original Caption:** Someone eating crisps 325
326 and talking.

(2) **Rewritten Caption by Prompt 1:** Crunchy 327
328 crisps mingle with the sound of a lively con- 329
330 versation, creating a cozy and intimate atmo- 331
332 sphere.

(2) **Rewritten Caption by Prompt 2:** The 333
334 crunch of crisps and the rustle of papers cre- 335
336 ate a cozy, intimate atmosphere, accompanied 337
338 by the gentle hum of a conversation. 339

340 During training, for each audio sample, we 341
342 choose the original caption with a probability $p =$ 343
344 0.4 or one of the rewritten versions (with a proba- 345
346 bility $1 - p$), where each rewritten caption has an 347
348 equal probability of selection. Both instructions 349
350

are provided in Appendix B. We employ LLaMa-2-13B (Touvron et al., 2023) with human-written in-context examples. We randomly sample 5 in-context examples from a collection of 50.

3.1.3 Multi-Layer Aggregator

Motivation. To extract additional details about the input audio, we devise a multi-layer aggregator that integrates multi-level hidden features of the pre-trained AST. Although AST has a global reception field in all layers, different layers learn auditory information at different scales (Singla et al., 2022), i.e., the middle layers encode more generic features (e.g., basic sounds, textures), while deeper layers capture high-level concepts (e.g., speech intonations, complex sound patterns). By aggregating these features, the multi-layer aggregator outputs features that encode a more holistic and fine-grained understanding of the audio. Thus, our multi-layer aggregator makes fine-grained auditory knowledge more likely to be learned while training.

Architecture. Our multi-layer aggregator is a transformer-style network consisting of two transformer layers for aggregating the hidden features of the audio encoder. Given the hidden features A_j and A_k from the middle layers in the audio encoder, the aggregation module uses two blocks to sequentially integrate the former two features with the last layer feature A_i . Each block \mathcal{B} is composed of self-attention, cross-attention, and Feed-forward network (FFN) arranged in a sequential manner. Finally, the output features \bar{A} is generated as follows,

$$\bar{A} = \mathcal{B}_2 (\mathcal{B}_1 (A_i; A_j); A_k) \quad (1)$$

$$\mathcal{B}(X; Y) = \text{FFN}(\text{Cross-Attn}(\text{Attn}(X), Y)). \quad (2)$$

In practice, we employ $j = 4$ and $k = 8$ from AST as our input to the multi-layer aggregator.

3.1.4 Soft Prompt

Motivation. Though models like AST and Audio Q-Former have shown much promise in audio tasks, a major problem still exists: real-world audio generally has multiple and overlapping acoustic events, and understanding all such events from model features proves to be inherently complex (Ghosh et al., 2024b). This eventually leads to sub-optimal performance for complex reasoning, where the explicit knowledge of *plausible* acoustic events in the audio can improve model responses. Thus, to improve fine-grained audio perception capabilities, we augment GAMA with high-level semantic understanding of the input audio. To do this, we employ

an off-the-shelf audio model to extract high-level semantic knowledge, i.e., audio event tags, as supplementary information. However, as audio event classification is not a solved problem, errors in tag predictions are inevitable. Thus, to mitigate the potential adverse effects of inaccurate predictions, we are inspired by prompt tuning to introduce a soft prompting technique that enables the model to utilize the embedded tags within the instructions adaptively.

Architecture. Fig. 2 shows an example of how we design our soft prompt together with an instruction. Specifically, we construct a fixed instruction template where we add the audio event tags along with the soft prompt, where the soft prompt is a trainable vector. In contrast to standard prompt tuning, where the model activations are generally steered towards completing the task for which the prompt is optimized, in our version the direction is specified by a tailored input sentence, “According to <hint>, you are allowed to use or partially use the following tags:”, and “<hint>” will be replaced by the soft prompt. This design allows us to select valuable information from tags adaptively rather than serving a specific task, as seen in standard prompt tuning methods. We only employ the soft prompt in the instruction tuning for complex reasoning step and not in the fine-tuning step. We provide a rationale in Appendix C.1.

3.1.5 Connection Module

We employ a multi-layer perceptron (MLP) to connect audio features into the word embedding space. All features are passed through separate MLP layers before being added as prefixes to word embeddings of the text instruction prompt.

3.2 CompA-R

Motivation. We define complex reasoning as the capability of an LALM to understand the input audio, every individual acoustic event in the audio, and reason the corresponding scene in which the audio might have occurred, such that it can infer nuanced relationships between them and its underlying contexts, thereby enabling it to draw sophisticated conclusions. We design CompA-R with the primary goal of endowing LALMs with complex reasoning abilities. We are motivated by the primary finding that current SOTA LALMs can only perform well in prompts that require describing the audio (e.g., *Describe the audio*) or reasoning-based prompts where identifying the

acoustic events present in the audio would suffice for a faithful response (e.g., *What type of video can this audio be used for dubbing?*). However, when posed with complex reasoning questions, these models often hallucinate or fail to provide a faithful response (see Fig. 4). Inspired by a wealth of prior work that shows how IT on well-curated datasets can align model behaviors for the execution of novel skills like reasoning and complex problem solving (Xu et al., 2024), we propose a systematic multi-stage pipeline to synthesize instruction-response pairs for CompA-R. CompA-R trains a model to engage in complex reasoning by querying it with instructions that cannot be directly inferred by identifying individual audio events and would require analyzing each event and its context in relation to other scene elements and world knowledge.

Synthesis Pipeline. We employ the AudioSet-strong subset to synthesize CompA-R. Our data synthesis pipeline consists of 3 stages: *i) Caption Generation.* To generate a caption that is aware of both the audio and the visual scene, we feed GPT-4 with multiple types of information about the audio and its corresponding video. These include a caption of the middle frame of the video generated using BLIP-2 (Li et al., 2023), objects in the frame identified using Grounding DINO (Liu et al., 2023c), image labels for the frame using the ImageNet (Deng et al., 2009) ontology obtained from CLIP, environment context using PlaceCNN (Zhou et al., 2017), caption of the audio obtained using RECAP (Ghosh et al., 2024a) and audio event tags using the AudioSet ontology obtained from AST. Finally, we prompt GPT-4 to aggregate these descriptions into a comprehensive caption. *ii) Dataset Synthesis.* We pass the generated caption together with the ground-truth acoustic event information and their corresponding time slices to GPT-4. We prompt GPT-4 with 3 human-written exemplars (which are randomly sampled from a pool of 50 exemplars) to synthesize an instruction-response pair. The exemplars and prompt are designed such that the synthesized instructions demand complex reasoning. We synthesize a total of 25000 instruction-response pairs. *iii) Human Verification.* We discard instructions due to unintended noise and hallucinations. We, the authors of this paper, manually verify a subset of CompA-R corresponding to 500 unique audios for creating the test set, i.e., CompA-R-test. The remainder of the synthesized dataset is used as the training set. We describe

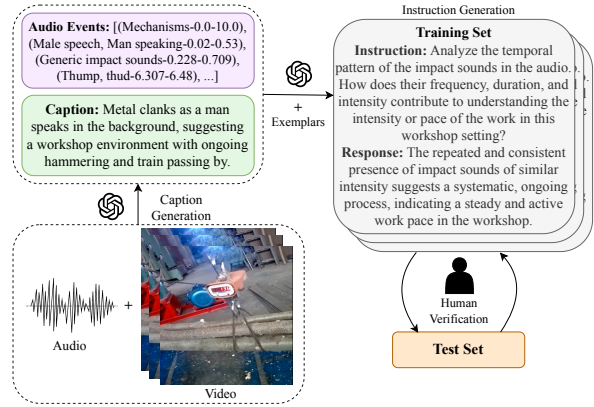


Figure 3: **Pipeline for synthesizing CompA-R.** For an audio in the AudioSet-strong dataset, we first use the audio and its corresponding video to generate a caption (described in Section 3.2). This caption is then fed into GPT-4 together with the ground-truth time slices for each event in the audio (available for AudioSet-strong). We then instruct GPT-4 to generate instruction-response pairs where the instruction is accompanied by human-written exemplars (prompt in Fig. 5). To generate the test set, we perform expert human verification for a subset of the generated dataset.

the process and annotation details further in Appendix G.1. This finally led to 200,234 unique pairs in training and 1,561 in testing.

3.3 Training

Fine-tuning. We fine-tune GAMA on the OpenAQA training set released by Gong et al. (2024). We use a fraction of all the instances due to the unavailability of the entire AudioSet and resource constraints. Dataset details are provided in Appendix H.1. Additionally, we augmented OpenAQA with 4 more datasets, including MusicCaps, MusicQA, NSynth, and Magna, to improve its music understanding capabilities. For fine-tuning, we follow the exact same 4-stage method proposed by Gong et al. (2024) where all parameters of all encoders are trainable, and we train only the LoRA modules of the LLM. We request our readers to refer to Gong et al. (2024) for more details.

Instruction Tuning on CompA-R. Post fine-tuning, we instruction-tune GAMA on CompA-R to endow it with complex reasoning abilities. Following common conventions (Liu et al., 2023b), we fine-tuned only the LoRA modules. We call the Instruction Tuned GAMA as GAMA-IT. Although fine-tuning on AQA also endows GAMA with instruction-following capabilities, CompA-R differs in the nature of training instances (thereby the capabilities it endows) and thus we differentiate with such a naming convention for ease of reading.

Model	ESC50# (Acc)	DCASE# (Mi-F1)	VS [†] (Acc)	TUT [†] (Acc)	BJO [†] (Acc)	VGG (Acc)	FSD (mAP)	NS _{ins} (ACC)	NS _{src} (ACC)	GTZAN [†] (ACC)	MSD [†] (ACC)	AudioSet (mAP)	Classif. Avg.	AudioCaps (SPICE)	Clotho (SPICE)	Cap. Avg.	ClothoAQA (ACC)
<i>Audio-Language encoder-based models. They are generalizable to unseen labels, but a pre-defined label set is required for inference.</i>																	
AudioCLIP	69.4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CLAP (Elizalde et al., 2023a)	82.6	30.0	48.4	29.6	47.5	24.0	30.2	22.7	16.4	25.0	44.0	5.8	29.4	-	-	-	-
CLAP (Wu* et al., 2023a)	89.1	31.3	47.1	35.6	48.0	26.3	30.8	25.2	18.9	26.3	46.9	6.2	36.0	-	-	-	-
CompA-CLAP	90.1	30.6	49.5	35.8	48.2	29.5	31.5	24.9	17.0	26.1	46.2	6.2	36.3	-	-	-	-
Audio Q-Former (w/o rw) (ours)	91.9	31.1	49.9	38.9	50.4	33.2	34.7	27.5	22.0	30.4	48.3	8.2	38.9	-	-	-	-
CLAP (w/rw) (ours)	90.7	30.9	50.7	36.2	53.4	30.7	37.1	24.1	17.4	27.6	48.1	6.1	37.7	-	-	-	-
Audio Q-Former (w/rw) (ours)	92.4	32.5	50.2	39.1	51.5	35.1	35.3	29.2	22.3	31.3	47.5	8.9	39.6	-	-	-	-
<i>Audio-Language generation-based models. They directly output label names and do not need a pre-defined label set is needed at inference.</i>																	
Qwen-Audio-Chat	71.7	32.4	74.2	16.9	50.8	17.5	39.8	30.2	41.3	41.6	69.1	13.4	41.1	14.7	9.8	12.3	32.3
LTU	81.7	37.5	53.3	19.9	67.8	50.3	43.9	28.0	41.8	9.9	74.2	18.3	42.4	16.9	11.7	15.8	25.1
SALMONN	16.4 [†]	18.0 [†]	16.9 [†]	7.8 [†]	25.0 [†]	23.3 [†]	22.1 [†]	16.2 [†]	33.7 [†]	10.1 [†]	28.8 [†]	13.4 [†]	17.9	8.3	7.6	8.0	23.1 [†]
Pengi	80.8 [†]	29.6 [†]	46.4 [†]	18.4 [†]	47.3 [†]	16.6 [†]	35.8	39.2	46.0	11.9	93.0	11.5	39.7	12.7	7.0	9.9	63.6
AudioGPT	41.3	20.9	35.8	14.9	21.6	5.6	18.8	40.9	15.6	11.9	28.5	12.7	22.4	6.9	6.2	6.6	33.4
GAMA (ours)	82.6	38.4	52.4	21.5	69.5	52.2	47.8	63.9	99.5	13.8	85.6	19.2	53.9	18.5	13.5	16.0	71.6
w/o AST & Aggregator	80.5	36.9	51.6	19.2	66.2	50.8	45.3	62.4	89.6	11.6	83.2	17.3	51.2	17.2	12.4	14.8	68.3
w/ Last Layer Features	81.3	37.6	50.2	20.4	68.2	51.7	45.8	62.6	92.3	11.2	81.5	18.1	51.7	17.7	12.8	15.3	69.5
w/o Audio Q-Former	79.7	37.4	51.3	20.2	68.0	51.6	46.4	60.1	90.4	11.6	79.8	18.4	51.2	16.9	11.9	14.4	61.2
w/ CLAP	81.8	38.4	52.2	21.6	69.1	52.0	47.5	58.8	99.5	12.4	77.9	19.0	52.5	17.2	13.1	15.1	66.4

Table 1: Comparison of GAMA with baselines on evaluation datasets described on close-ended general audio and music understanding benchmarks. GAMA outperforms most ALMs on most settings. [†] and # indicate zero-shot and weak zero-shot, respectively. **Note:** Qwen-Audio-Chat does not provide training details. We also mark baseline values which are zero-shot.

Models	CompA-R-test (GPT-4/Human)				OpenAQA				Dense Captioning		
	Clarity	Correctness	Engagement	Avg.	Clarity	Correctness	Engagement	Avg.	AudioCaps	Clotho	Avg.
Qwen-Audio-Chat	3.5 / 3.4	3.3 / 3.4	3.6 / 3.2	3.5 / 3.5	3.6	3.6	3.5	3.6	3.8	3.6	3.7
LTU	3.5 / 4.0	3.2 / 3.3	3.4 / 3.5	3.4 / 3.6	3.5	3.7	3.5	3.6	3.5	3.6	3.5
SALMONN	2.6 / 2.8	2.4 / 2.3	2.0 / 2.2	2.3 / 2.4	2.4	2.5	2.7	2.5	2.8	3.1	2.9
Pengi	1.8 / 1.6	1.5 / 1.4	1.3 / 1.2	1.5 / 1.4	1.7	1.5	1.4	1.5	2.6	2.8	2.7
AudioGPT	1.3 / 1.4	1.6 / 1.5	1.4 / 1.7	1.4 / 1.5	1.6	1.5	1.5	1.5	2.7	2.9	2.8
LTU w/ CompA-R	3.5 / 4.0	3.2 / 3.3	3.4 / 3.5	3.6 / 3.6	3.5	3.7	3.5	3.6	3.7	3.8	3.8
GAMA-IT (ours)	4.3 / 4.5	3.9 / 4.1	3.9 / 4.3	4.0 / 4.3	4.0	4.2	3.8	4.0	4.3	4.1	4.2
w/o Soft Prompt	4.1 / 4.2	3.7 / 3.8	3.6 / 3.4	3.8 / 3.8	3.9	3.8	3.7	3.8	4.1	3.9	4.0
w/o Aggregator	4.0 / 4.2	3.5 / 3.5	3.6 / 3.5	3.7 / 3.7	3.7	3.7	3.5	3.6	3.7	3.8	3.8
w/o Audio Q-Former	3.8 / 3.7	3.4 / 3.6	3.5 / 3.3	3.6 / 3.5	3.4	3.9	3.5	3.6	3.7	3.5	3.6

Table 2: Comparison of GAMA with other baselines on open-ended AQA on OpenAQA, complex open-ended AQA on CompA-R-test and Dense Captioning on 500 instances from AudioCaps and Clotho.

3.4 Experimental Setup

Hyper-parameters. For the fine-tuning stage, we follow the exact same hyper-parameter setup proposed by Gong et al. (2024). However, we scale down our batch sizes to 4, 2, 2, and 2 (due to compute constraints) with an effective batch size of 256 in all stages. For Instruction Tuning, we employ a batch size of 2, an effective batch size of 256, and a learning rate of 1e-4. For both training and evaluation, we sampled audio at 16kHz.

Baselines. We compare GAMA with *i*) generation-based LALMs: LTU, Qwen-Audio, SALMONN, Pengi and AudioGPT. We only employ the original checkpoints open-sourced by the authors and do not re-train the models due to compute constraints (except LTU, which we retrain on our version of OpenAQA, the same batch size as GAMA, and with LLaMa-2 as the LLM). We do not compare with Audio Flamingo (Kong et al., 2024) as the checkpoint was not available at the time of writing the paper, and we are constrained by compute for training it from scratch. *ii*) audio-language encoders: CLAP by Wu* et al. (2023b) and Elizalde et al. (2023b), CompA-CLAP (Ghosh et al., 2024b), AudioCLIP (Guzhov et al., 2021) and Audio Q-Former. For dense captioning and close- and open-ended AQA, we evaluate using GAMA-IT. For all

other tasks, we evaluate using the only fine-tuned version of GAMA (rationale in Appendix C).

Evaluation Datasets and Metrics. Evaluation metrics used for all evaluation datasets are mentioned in Table 2 and detailed statistics about each dataset is mentioned in Section H.2. For classification, zero-shot evaluation refers to datasets GAMA that have never been seen during training; weak zero-shot evaluation refers to datasets GAMA that have not been seen in training but are sourced from the same project as part of the training data, and seen datasets refer to datasets GAMA has been trained on. Similar to Deshmukh et al. (2023); Gong et al. (2024), we first caption the audio and retrieve the most similar label using SentenceBERT. We employ either accuracy (Acc), Micro-F1 (Mi-F1), or Mean Average Precision (mAP) for classification evaluation. For captioning, we also propose dense captioning, which evaluates a model for its capability to identify every event in the audio and the context of its occurrence with respect to other events in the audio (more in Section 4). For evaluation, we randomly select a subset of 500 samples from AudioCaps and Clotho. We also employ human evaluation for OpenAQA, CompA-R-test, and dense captioning. For human evaluation, we ask human annotators to score the caption on a scale

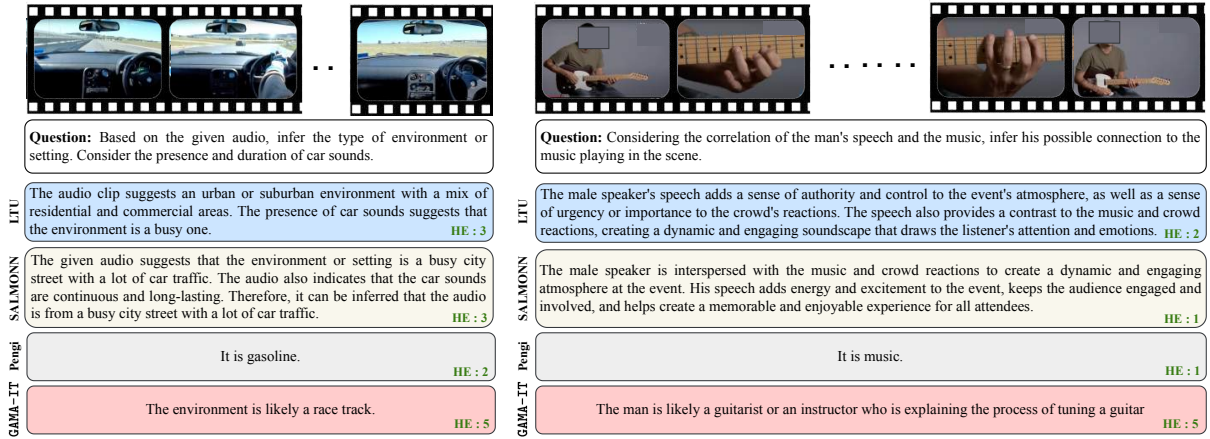


Figure 4: Qualitative comparison of GAMA with other baselines on instances from *CompA-R-test*. Both instances challenge an LLM with a question about the input audio that requires advanced understanding and complex reasoning regarding the audio and its individual events. GAMA is able to provide faithful and accurate responses through improved audio perception and reasoning capabilities. HE refers to scores assigned by human evaluators. Note that the video is only provided for illustration purposes and not provided as input to the LLM. More examples here: <https://gamaaudio.github.io/gamaaudio/>

of 1-5 and report the score averaged across the 3. More details on recruitment and background of annotators can be found in Appendix D. Finally, due to human evaluation being prohibitively expensive, we also propose an automated evaluation methodology for complex open-ended AQA on *CompA-R-test*. We evaluate model responses using text-only GPT-4, where we provide it with the audio caption generated in Section 3.2 and the gold-standard audio event with timestamps (prompt in Appendix B).

4 Results and Analysis

Quantitative Results. Table 1 compares GAMA with other baselines on classification and captioning tasks. For zero-shot classification evaluation on VocalSound (VS) (Gong et al., 2022), TUT 2017 (TUT) (Mesaros et al., 2018), Beijing Opera (BJO) (Tian et al., 2014), GTZAN (GTZ) (Park et al., 2022) and Medley-solos-DB (MDB) (Lostanlen et al., 2018), GAMA outperforms our baselines by 2%-67%. For weak zero-shot evaluation on ESC-50 (Piczak, 2015) and DCASE2017 Task 4 (DCASE) (Mesaros et al., 2017), GAMA outperforms our baselines by 1%-66%. Finally, for in-domain evaluation on VGGSound (VGG) (Chen et al., 2020), FSD50K (FSD) (Fonseca et al., 2021), AudioSet (AS) (Gemmeke et al., 2017) and NSynth (NS) (Engel et al., 2017) GAMA outperforms our baselines by 1%-84%. GAMA sees the steepest drop in performance when the AST and Aggregator are removed (i.e., only Audio Q-Former is employed).

Table 2 compares GAMA with other baselines on AQA (open-ended and complex open-ended)

and dense captioning. GAMA outperforms all our baselines on all settings. GAMA shows absolute improvement of 4% - 50% on OpenAQA, 8% - 58% on *CompA-R-test* and 8% - 30% on Dense Captioning. Similar to the tasks in Table 1, performance on benchmarks suffers the most when without the Audio Q-Former (when only the AST and Aggregator are employed). Audio Q-Former proves to especially effective (over employing CLAP) in AQA.

Qualitative Results. Fig. 4 compares GAMA-IT against other LLMs from literature with instances from *CompA-R-test*. All models compared by default possess audio chat or open-ended AQA capabilities. GAMA-IT is able to provide more faithful responses that are both correct and preferred more by humans. We provide additional comparisons in Figs. 8, 9, 10, 11, 12, and our demo page: (where we also show comparisons of dense captioning).

5 Conclusion

In this paper, we propose GAMA, an LLM with improved audio perception abilities. We integrate an LLM with multiple types of audio representations, which are responsible for providing diverse knowledge about the input audio. GAMA fine-tuned on a mixture of open-source datasets outperforms prior audio-language models by significant margins on 16 datasets spanning 4 tasks. Next, we propose *CompA-R*, an instruction-tuning dataset that we synthesize using a robust pipeline for endowing an LLM with complex reasoning abilities. GAMA IT-ed on *CompA-R* outperforms baselines on complex open-ended AQA and dense captioning.

592 Limitations and Future Work

593 GAMA and our experimental setup have several lim-
594 itations, including:

- 595 • For the scope of our experiments, we do not
596 evaluate and compare music understanding
597 extensively. We do not do this as we do not
598 train GAMA on diverse and large-scale music
599 datasets. We also acknowledge that it is possi-
600 ble to employ the GAMA architecture for com-
601 prehensive music understanding if trained on
602 large-scale music understanding datasets. As
603 part of future work, we plan to release a music-
604 only version of GAMA, similar to Gardner et al.
605 (2024).
- 606 • We do not employ larger LLMs, for exam-
607 ple, the 13B versions of the LLaMA family,
608 similar to Tang et al. (2024) and Gong et al.
609 (2024), due to compute constraints.
- 610 • The audio-encoder(s) in GAMA have more pa-
611 rameters than in our baselines. However, we
612 also acknowledge that this adds to only a frac-
613 tion of the total parameter count of the LLM.

614 References

615 2018. [A dump of BBC’s sound effects library.](#)
616 This dump was created using the script found at
617 <https://github.com/FThompson/BBCSoundDownloader>.
618 Identifier: BBCSoundEffectsComplete.

619 2023. [SoundBible - Free Sound Clips, Sound Bites, and](#)
620 [Sound Effects.](#) Accessed: 25 September 2023.

621 Andrea Agostinelli, Timo I. Denk, Zalán Borsos,
622 Jesse Engel, Mauro Verzetti, Antoine Caillon,
623 Qingqing Huang, Aren Jansen, Adam Roberts, Marco
624 Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Chris-
625 tian Frank. 2023. [Musilm: Generating music from](#)
626 [text.](#)

627 Sébastien Bubeck, Varun Chandrasekaran, Ronen El-
628 dan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
629 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund-
630 berg, et al. 2023. Sparks of artificial general intelli-
631 gence: Early experiments with gpt-4. *arXiv preprint*
632 *arXiv:2303.12712*.

633 Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew
634 Zisserman. 2020. [Vggsound: A large-scale audio-](#)
635 [visual dataset.](#)

636 Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shil-
637 iang Zhang, Zhijie Yan, Chang Zhou, and Jingren
638 Zhou. 2023. Qwen-audio: Advancing universal
639 audio understanding via unified large-scale audio-
640 language models. *arXiv preprint arXiv:2311.07919*.

Wanyun Cui and Qianle Wang. 2024. [Ada-instruct:](#)
641 [Adapting instruction generators for complex reason-](#)
642 [ing.](#) 643

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
644 and Li Fei-Fei. 2009. Imagenet: A large-scale hier-
645 archical image database. In *2009 IEEE conference*
646 *on computer vision and pattern recognition*, pages
647 248–255. Ieee. 648

Soham Deshmukh, Benjamin Elizalde, Rita Singh, and
649 Huaming Wang. 2023. [Pengi: An audio language](#)
650 [model for audio tasks.](#) 651

Soham Deshmukh, Benjamin Elizalde, and Huaming
652 Wang. 2022. Audio retrieval with wavtext5k and clap
653 training. *arXiv preprint arXiv:2209.14275*. 654

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
655 Kristina Toutanova. 2018. Bert: Pre-training of deep
656 bidirectional transformers for language understand-
657 ing. *arXiv preprint arXiv:1810.04805*. 658

Konstantinos Drossos, Samuel Lipping, and Tuomas
659 Virtanen. 2020. Clotho: An audio captioning dataset.
660 In *ICASSP 2020-2020 IEEE International Confer-*
661 *ence on Acoustics, Speech and Signal Processing*
662 *(ICASSP)*, pages 736–740. IEEE. 663

Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Is-
664 mail, and Huaming Wang. 2023a. Clap learning
665 audio concepts from natural language supervision.
666 In *ICASSP 2023-2023 IEEE International Confer-*
667 *ence on Acoustics, Speech and Signal Processing*
668 *(ICASSP)*, pages 1–5. IEEE. 669

Benjamin Elizalde, Soham Deshmukh, Mahmoud Al
670 Ismail, and Huaming Wang. 2023b. [Clap learning](#)
671 [audio concepts from natural language supervision.](#)
672 In *ICASSP 2023 - 2023 IEEE International Confer-*
673 *ence on Acoustics, Speech and Signal Processing*
674 *(ICASSP)*, pages 1–5. 675

Jesse Engel, Cinjon Resnick, Adam Roberts, Sander
676 Dieleman, Mohammad Norouzi, Douglas Eck, and
677 Karen Simonyan. 2017. Neural audio synthesis of
678 musical notes with wavenet autoencoders. In *Inter-*
679 *national Conference on Machine Learning*, pages
680 1068–1077. PMLR. 681

Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic
682 Font, and Xavier Serra. 2021. Fsd50k: an open
683 dataset of human-labeled sound events. *IEEE/ACM*
684 *Transactions on Audio, Speech, and Language Pro-*
685 *cessing*, 30:829–852. 686

Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic
687 Font, and Xavier Serra. 2022. [Fsd50k: An open](#)
688 [dataset of human-labeled sound events.](#) 689

Joshua P Gardner, Simon Durand, Daniel Stoller, and
690 Rachel M Bittner. 2024. [LLark: A multimodal foun-](#)
691 [dation model for music.](#) 692

693	Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman,	Rongjie Huang, Mingze Li, Dongchao Yang, Jia-	746
694	Aren Jansen, Wade Lawrence, R Channing Moore,	tong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu,	747
695	Manoj Plakal, and Marvin Ritter. 2017. Audio set:	Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. 2024.	748
696	An ontology and human-labeled dataset for audio	Audiogpt: Understanding and generating speech, mus-	749
697	events. In <i>2017 IEEE international conference on</i>	ic, sound, and talking head. In <i>Proceedings of</i>	750
698	<i>acoustics, speech and signal processing (ICASSP)</i> ,	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	751
699	pages 776–780. IEEE.	ume 38, pages 23802–23804.	752
700	Sreyan Ghosh, Sonal Kumar, Chandra Kiran	Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee,	753
701	Reddy Evuru, Ramani Duraiswami, and Dinesh	and Gunhee Kim. 2019. Audiocaps: Generating cap-	754
702	Manocha. 2024a. Recap: Retrieval-augmented	tions for audios in the wild. In <i>Proceedings of the</i>	755
703	audio captioning . In <i>ICASSP 2024 - 2024 IEEE</i>	<i>2019 Conference of the North American Chapter of</i>	756
704	<i>International Conference on Acoustics, Speech and</i>	<i>the Association for Computational Linguistics: Hu-</i>	757
705	<i>Signal Processing (ICASSP)</i> , pages 1161–1165.	<i>man Language Technologies, Volume 1 (Long and</i>	758
706	Sreyan Ghosh, Ashish Seth, Sonal Kumar, Utkarsh	<i>Short Papers)</i> , pages 119–132.	759
707	Tyagi, Chandra Kiran Reddy Evuru, Ramaneswaran	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	760
708	S, S Sakshi, Oriol Nieto, Ramani Duraiswami, and	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	761
709	Dinesh Manocha. 2024b. Compa: Addressing the	guage models are zero-shot reasoners. <i>Advances in</i>	762
710	gap in compositional reasoning in audio-language	<i>neural information processing systems</i> , 35:22199–	763
711	models . In <i>The Twelfth International Conference on</i>	22213.	764
712	<i>Learning Representations</i> .	Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping,	765
713	Yuan Gong, Yu-An Chung, and James Glass. 2021.	Rafael Valle, and Bryan Catanzaro. 2024. Audio	766
714	Ast: Audio spectrogram transformer. <i>arXiv preprint</i>	flamingo: A novel audio language model with few-	767
715	<i>arXiv:2104.01778</i> .	shot learning and dialogue abilities .	768
716	Yuan Gong, Hongyin Luo, Alexander H. Liu, Leonid	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	769
717	Karlinisky, and James R. Glass. 2024. Listen, think,	2023. Blip-2: Bootstrapping language-image pre-	770
718	and understand . In <i>The Twelfth International Confer-</i>	training with frozen image encoders and large lan-	771
719	<i>ence on Learning Representations</i> .	guage models. In <i>International conference on ma-</i>	772
720	Yuan Gong, Andrew Rouditchenko, Alexander H. Liu,	<i>chine learning</i> , pages 19730–19742. PMLR.	773
721	David Harwath, Leonid Karlinisky, Hilde Kuehne,	Samuel Lipping, Parthasaarathy Sudarsanam, Konstanti-	774
722	and James R. Glass. 2023. Contrastive audio-visual	nos Drossos, and Tuomas Virtanen. 2022. Clotho-	775
723	masked autoencoder . In <i>The Eleventh International</i>	qa: A crowdsourced dataset for audio question an-	776
724	<i>Conference on Learning Representations</i> .	swering. In <i>2022 30th European Signal Processing</i>	777
725	Yuan Gong, Jin Yu, and James Glass. 2022. Vocal-	<i>Conference (EUSIPCO)</i> , pages 1140–1144. IEEE.	778
726	sound: A dataset for improving human vocal sounds	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	779
727	recognition. In <i>ICASSP 2022-2022 IEEE Interna-</i>	Lee. 2023a. Improved baselines with visual instruc-	780
728	<i>tional Conference on Acoustics, Speech and Signal</i>	tion tuning. <i>arXiv preprint arXiv:2310.03744</i> .	781
729	<i>Processing (ICASSP)</i> , pages 151–155. IEEE.	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	782
730	Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang	Lee. 2023b. Visual instruction tuning .	783
731	Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	784
732	Dawn Song. 2023. The false promise of imitating	Lee. 2024. Visual instruction tuning. <i>Advances in</i>	785
733	proprietary llms. <i>arXiv preprint arXiv:2305.15717</i> .	<i>neural information processing systems</i> , 36.	786
734	Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas	Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao	787
735	Dengel. 2022. Audioclip: Extending clip to image,	Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang	788
736	text and audio. In <i>ICASSP 2022</i> .	Su, Jun Zhu, et al. 2023c. Grounding dino: Marrying	789
737	Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas	dino with grounded pre-training for open-set object	790
738	Dengel. 2021. Audioclip: Extending clip to image,	detection. <i>arXiv preprint arXiv:2303.05499</i> .	791
739	text and audio .	Vincent Lostanlen, Carmine-Emanuele Cella, Rachel	792
740	Yuchen Hu, CHEN CHEN, Chao-Han Huck Yang,	Bittner, and Slim Essid. 2018. Medley-solos-db: a	793
741	Ruizhe Li, Chao Zhang, Pin-Yu Chen, and Ensiong	crosscollection dataset for musical instrument recog-	794
742	Chng. 2024. Large language models are efficient	nition. <i>Zenodo</i> .	795
743	learners of noise-robust speech recognition . In <i>The</i>	Vincent Lostanlen, Carmine-Emanuele Cella, Rachel	796
744	<i>Twelfth International Conference on Learning Repre-</i>	Bittner, and Slim Essid. 2019. Medley-solos-DB: a	797
745	<i>sentations</i> .	cross-collection dataset for musical instrument recog-	798
		nition .	799

800	Annamaria Mesaros, Toni Heittola, Aleksandr Diment,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	856
801	Benjamin Elizalde, Ankit Shah, Emmanuel Vincent,	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	857
802	Bhiksha Raj, and Tuomas Virtanen. 2017. Dcase	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	858
803	2017 challenge setup: Tasks, datasets and baseline	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	859
804	system. In <i>DCASE 2017-Workshop on Detection and</i>	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	860
805	<i>Classification of Acoustic Scenes and Events</i> .	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	861
806	Annamaria Mesaros, Toni Heittola, and Tuomas	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	862
807	Virtanen. 2018. A multi-device dataset for ur-	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	863
808	ban acoustic scene classification. <i>arXiv preprint</i>	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	864
809	<i>arXiv:1807.09840</i> .	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	865
810	Irene Martin Morato and Annamaria Mesaros. 2021.	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	866
811	Macs - multi-annotator captioned soundscapes .	Melanie Kambadur, Sharan Narang, Aurelien Rod-	867
812	Junwoo Park, Youngwoo Cho, Gyuhyeon Sim, Hojoon	riguez, Robert Stojnic, Sergey Edunov, and Thomas	868
813	Lee, and Jaegul Choo. 2022. Enemy spotted: in-	Scialom. 2023. Llama 2: Open foundation and fine-	869
814	game gun sound dataset for gunshot classification and	tuned chat models .	870
815	localization. In <i>2022 IEEE Conference on Games</i>	George Tzanetakis, Georg Essl, and Perry Cook. 2001.	871
816	(<i>CoG</i>), pages 56–63. IEEE.	Automatic musical genre classification of audio sig-	872
817	Karol J Piczak. 2015. Esc: Dataset for environmental	nals .	873
818	sound classification. In <i>Proceedings of the 23rd ACM</i>	Shinji Watanabe, Takaaki Hori, Shigeaki Karita, Tomoki	874
819	<i>international conference on Multimedia</i> , pages 1015–	Hayashi, Jiro Nishitoba, Yuya Unno, Nelson En-	875
820	1018.	rique Yalta Soplin, Jahn Heymann, Matthew Wiesner,	876
821	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Nanxin Chen, et al. 2018. Espnet: End-to-end speech	877
822	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	processing toolkit. <i>arXiv preprint arXiv:1804.00015</i> .	878
823	try, Amanda Askell, Pamela Mishkin, Jack Clark,	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	879
824	et al. 2021. Learning transferable visual models from	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.	880
825	natural language supervision. In <i>International confer-</i>	Dai, and Quoc V Le. 2022. Finetuned language mod-	881
826	<i>ence on machine learning</i> , pages 8748–8763. PMLR.	els are zero-shot learners . In <i>International Confer-</i>	882
827	Yaman Kumar Singla, Jui Shah, Changyou Chen, and	ence on Learning Representations .	883
828	Rajiv Ratn Shah. 2022. What do audio transform-	Yusong Wu*, Ke Chen*, Tianyu Zhang*, Yuchen Hui*,	884
829	ers hear? probing their representations for language	Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023a.	885
830	delivery & structure. In <i>2022 IEEE International</i>	Large-scale contrastive language-audio pretraining	886
831	<i>Conference on Data Mining Workshops (ICDMW)</i> ,	with feature fusion and keyword-to-caption augmen-	887
832	pages 910–925. IEEE.	tation. In <i>IEEE International Conference on Acous-</i>	888
833	Sonniss Limited. 2022. Sonniss Game Audio . Regis-	<i>tics, Speech and Signal Processing, ICASSP</i> .	889
834	tered in England, UK. Company number: 09377364.	Yusong Wu*, Ke Chen*, Tianyu Zhang*, Yuchen Hui*,	890
835	Accessed: 25 September 2023.	Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023b.	891
836	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao	Large-scale contrastive language-audio pretraining	892
837	Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao	with feature fusion and keyword-to-caption augmen-	893
838	Zhang. 2024. SALMONN: Towards generic hearing	tation. In <i>IEEE International Conference on Acous-</i>	894
839	abilities for large language models . In <i>The Twelfth</i>	<i>tics, Speech and Signal Processing, ICASSP</i> .	895
840	<i>International Conference on Learning Representa-</i>	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,	896
841	<i>tions</i> .	Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei	897
842	Mi Tian, Ajay Srinivasamurthy, Mark Sandler, and	Lin, and Daxin Jiang. 2024. WizardLM: Empow-	898
843	Xavier Serra. 2014. A study of instrument-wise on-	ering large pre-trained language models to follow	899
844	set detection in beijing opera percussion ensembles.	complex instructions . In <i>The Twelfth International</i>	900
845	In <i>2014 IEEE International Conference on Acoustics,</i>	<i>Conference on Learning Representations</i> .	901
846	<i>speech and signal processing (icassp)</i> , pages 2159–	Duzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong,	902
847	2163. IEEE.	Dan Su, Chenhui Chu, and Dong Yu. 2024. Mm-	903
848	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	llms: Recent advances in multimodal large language	904
849	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	models. <i>arXiv preprint arXiv:2401.13601</i> .	905
850	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,	906
851	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-	907
852	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	wei Zhang, Fei Wu, et al. 2023. Instruction tuning	908
853	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	for large language models: A survey. <i>arXiv preprint</i>	909
854	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	<i>arXiv:2308.10792</i> .	910
855	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan		

911	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	formal framework, we leave this as part of future	959
912	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	research.	960
913	Zhang, Junjie Zhang, Zican Dong, et al. 2023. A		
914	survey of large language models. <i>arXiv preprint</i>		
915	<i>arXiv:2303.18223</i> .		
916	Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude		962
917	Oliva, and Antonio Torralba. 2017. Places: A 10		963
918	million image database for scene recognition. <i>IEEE</i>		964
919	<i>Transactions on Pattern Analysis and Machine Intel-</i>		965
920	<i>ligence</i> .		966
921			967
922			968
923			969
924			970
925			971
926			972
927			973
928			974
929			975
930			976
931			977
932			978
933			979
934			980
935			981
936			982
937			983
938			984
939			985
940			986
941			987
942			988
943			989
944			990
945			991
946			992
947			993
948			994
949			995
950			996
951			997
952			998
953			999
954			1000
955			1001
956			1002
957			1003
958			1004
			1005
			1006
			1007

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

A Additional Results

B Prompts employed for LLMs

Fig. 5 illustrates the prompt employed for synthesizing CompA-R. Fig. 6 illustrates the prompt employed for evaluating model responses on CompA-R. For dense captioning, we just prompt the model: *Write an audio caption describing the sound in detail*.

C GAMA-IT vs GAMA and Evaluation Choices.

GAMA is first fine-tuned on OpenAQA and then instruction-tuned on CompA-R for complex reasoning. We call the instruction-tuned version GAMA-IT. We do not evaluate GAMA-IT on general tasks like classification and vanilla captioning². GAMA-IT is aligned to generate detailed descriptions as part of the complex reasoning stage, and we found a lack of metrics and methods that can faithfully evaluate such descriptions for classification or captioning. For example, the retrieval-based classification evaluation method, employed extensively in prior work, including ours, uses a Sentence-BERT to retrieve the label closest to the description for classification evaluation. During our preliminary analysis, we found that Sentence-BERT, which just performs retrieval using semantic matching, is unable to faithfully retrieve the correct label despite the caption mentioning the label as an audio event. We further investigated CLAP as our retrieval model for evaluation and found that it suffers from the same limitations. We attribute this to the detailed and dense nature of the descriptions and the fact that these models only focus on high-level semantic meaning for retrieval. Our initial experiments show that LLM prompting serves as a feasible alternative for automatic evaluation (beyond human evaluation) using such dense descriptions, but due to the lack of resources and a

²**Note:** Both depend on the description of the input audio generated by the model

formal framework, we leave this as part of future research.

C.1 Soft Prompts

We employ the soft prompt only in the instruction tuning stage for learning complex reasoning and not in the fine-tuning step. We do this for 2 reasons: (i) Fine-tuned GAMA is only expected to solve generic audio tasks like classification, captioning, etc. Thus, we hypothesize that such high-level semantic cues are not necessary for effective and optimal performance. (ii) Since fine-tuning is done on a large-scale dataset and acoustic event classification is far from accurate, our soft prompt method might add unwanted noise to the training process, thereby leading to sub-optimal performance. On the contrary, our instruction-tuning stage, which is done on relatively low-resource data and is only responsible for aligning a model for complex reasoning, is robust to inaccurate audio tags due to our soft-prompting methodology.

D Additional Details: Human Study

Note. Our institution’s Institutional Review Board (IRB) has granted approval for both human studies presented in the paper.

Background and Recruitment for Dense Captioning and CompA-R-test Evaluation. We recruit 3 professionals for human evaluation of dense captioning and CompA-R-test evaluation. All these 3 professionals come with at least a Ph.D. in Engineering or Sciences and were asked to use headphones to first analyze the audio and then judge the response quality. The authors of this paper gave these annotators 5 examples of responses and the corresponding judgments. The work was done voluntarily and not paid. We refrain from recruiting crowd raters as prior research has noticed discrepancies in evaluation by them (Gudibande et al., 2023). More precisely, they have been shown to possess a tendency to rate an answer with a high score only by visualizing the style of answering and not the exact factual information making up the response.

All 3 human annotators score the response between 1-5 and we report score averaged across the 3. Prior to evaluation all annotators were given at least 10 examples from the authors of the paper of generations and their corresponding scores. For evaluation, only the audio was provided to them with a software that could play the audio and has

fields to input the scores.

Background and Recruitment for OpenQA. Since the size of OpenQA is relatively larger than CompA-R-test, we perform evaluation on Amazon Mechanical Turk similar to Gong et al. (2024). Evaluation was done with a total of 267 unique human evaluators and each generation was scored by 2 evaluators. The same software was used for evaluation as CompA-R-test.

E Additional Details: Audio Q-Former

E.1 Audio Q-Former Training Details

Pre-training Hyper-parameter. For Stage 1 of training, we employ a training batch size of 192, an initial learning rate of 1e-4, a minimum learning rate of 1e-5, and a warm-up learning rate of 1e-6. We do cosine decay as the learning rate scheduling technique. We do warmup for 5000 steps. Stage 1 was pre-trained on 8 A6000 GPUs for 100 epochs. For Stage 2 of training, we keep the exact same settings as Stage 1 but change the batch size to 128. **Fine-tuning.** For zero-shot audio classification evaluation, we fine-tune the Audio Q-Former after Stage 1 pre-training on the same corpus presented in Table 3 and using the same Stage 1 objective. The only difference in the fine-tuning step is that we train the AST model, which is otherwise kept frozen in the pre-training stage.

Fine-tuning Hyper-parameter. For fine-tuning, we again use the same hyper-parameter setting as Stage 1 pre-training but use a batch size of 64.

E.2 Training Dataset Details

Table 3 provides dataset statistics of all individual datasets used for training Audio Q-Former. We employ ≈ 2.2 M audio-caption pairs for training with no speech-transcription pairs.

Dataset	#Audio-Caption Pairs
Audio Set (Gemmeke et al., 2017) ³	1591364
Free Sound (Fonseca et al., 2022) ⁴	259020
VGGSound (Chen et al., 2020) ⁵	185161
AudioSet Strong (CompA Version) (Ghosh et al., 2024b) ⁶	108311
MACS (Morato and Mesaros, 2021) ⁷	14400
BBC (BBC, 2018) ⁸	31201
AudioCaps (Kim et al., 2019) ⁹	48649
Clotho (Drossos et al., 2020) ¹⁰	18735
SONISS (Sonmiss Limited, 2022) ¹¹	1602
Musical Instrument (Agostinelli et al., 2023) ¹²	7990
SoundBible (sou, 2023) ¹³	1232
WavText5K (Deshmukh et al., 2022) ¹⁴	4347
MusicCaps (Agostinelli et al., 2023) ¹⁵	2645
GTZAN (Tzanetakis et al., 2001) ¹⁶	6014
Medley-solos (Lostanlen et al., 2019) ¹⁷	732

Table 3: List of open-source datasets used for collating our final dataset for training ReCLAP with ≈ 2.2 M audio-caption pairs. All datasets are free to use for research purposes.

E.3 Augmentation Examples

Table 11 illustrates prompt augmentations for two categories from each dataset. Table 12 illustrates caption augmentations for training Audio Q-Former.

E.4 Hyper-parameter Tuning

E.4.1 Number of the custom prompts N

In this subsection, we show the effect of the number of custom prompts N on the final zero-shot audio classification performance. Table 4 compares performance across $N=\{1,2,3,4,5\}$. As we see, the optimal performance is achieved at $N=2$, and model performance decreases with an increase in N . This decline is hypothesized to be due to the introduction of more noise into the process with each additional caption.

N	1	2	3	4	5
Score	38.1	39.6	<u>39.0</u>	39.4	36.5

Table 4: Impact of N on ZSAC with ReCLAP.

E.4.2 Probability of choosing rewritten captions

In this subsection, we show the effect of probability p on the final zero-shot audio classification performance.

p	0.2	0.4	0.6	0.8
Score	35.1	39.6	<u>38.4</u>	37.0

Table 5: Impact of p on ZSAC with ReCLAP.

¹<https://research.google.com/audioset/download.html>
²<https://huggingface.co/datasets/cvssp/WavCaps>
³<https://www.robots.ox.ac.uk/vgg/data/vggsound/>
⁴<https://research.google.com/audioset/download.html>
⁵<https://zenodo.org/records/5114771>
⁶<https://sound-effects.bbcrewind.co.uk/>
⁷<https://research.google.com/audioset/download.html>
⁸<https://zenodo.org/records/4783391>
⁹<https://labs.freesound.org/datasets/>
¹⁰<https://www.kaggle.com/datasets/soumendraprasad/musical-instruments-sound-dataset>
¹¹<https://soundbible.com/>
¹²<https://github.com/microsoft/WavText5K>
¹³https://github.com/seungheondoh/music_caps_dl
¹⁴<https://www.kaggle.com/datasets/andradaolteanu/gtzan-dataset-music-genre-classification>
¹⁵<https://zenodo.org/records/1344103>
¹⁶<https://zenodo.org/records/1344103>

1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113

F Baseline Details

AudioCLIP. (Guzhov et al., 2022) AudioCLIP is an extension of the CLIP model that can handle audio in addition to text and images by incorporating the ESResNeXt audio model in the CLIP framework. It was trained on the AudioSet dataset, which contains millions of audio clips with corresponding labels.

CLAP. (Elizalde et al., 2023a) CLAP (Contrastive Language-Audio Pre-training), similar to CLIP, is an audio-language model trained with contrastive learning between audio data and their corresponding natural language descriptions. Representations are obtained from audio encoders and text encoders. Wu* et al. (2023b) further extend this using a feature fusion mechanism and keyword-to-caption augmentation into the model design to further enable the model to process audio inputs of variable lengths and enhance performance.

CompA-CLAP. (Elizalde et al., 2023a) CompA-CLAP, an extension to CLAP, is trained on completely open-sourced datasets and further fine-tuned using specific algorithms and datasets to improve compositional reasoning.

Pengi. (Deshmukh et al., 2023) Pengi was one of the first efforts to achieve general-purpose audio understanding through free-form language generation with transfer learning. Precisely, Pengi integrates an audio encoder with a decoder-only pre-trained language model (LM) where the audio features serve as a prefixes for the LM during response generation. Following this, similar to our evaluation strategy, they prompt the model to caption the input audio and calculate the similarity between the caption and the ground-truth audio label for zero-shot classification.

LTU. (Gong et al., 2024) As a concurrent work to Pengi, took a step forward and showed that substituting the pre-trained language model with an LLM can induce an LALM with reasoning capabilities. Precisely, they achieved this by integrating an audio encoder to LLaMA (Touvron et al., 2023) and fine-tuning the model on close-ended and open-ended instruction-tuning datasets. Finally, beyond just close-ended tasks, they also evaluate their models on open-ended reasoning tasks and show superior performance compared to baselines.

AudioGPT. (Huang et al., 2024) Different from Pengi and LTU, AudioGPT differs in how the audio models and LLMs are integrated for complet-

ing audio tasks. More specifically, different from end-to-end training and alignment, they integrate a closed-source model (ChatGPT) with a pre-trained audio model, already capable of completing the required task, using a modality-transfer transformer τ . The integration or interaction between the two models is accomplished using the prompts. Additionally, AudioGPT is capable of solving more tasks, which include human verbal speech, beyond just non-verbal speech like Pengi and LTU.

SALMONN. (Tang et al., 2024) SALMONN follows a similar architecture to LTU and Pengi and does prefix conditioning with an LLM. However, in addition to an audio encoder, they also integrate a speech encoder for speech or verbal audio understanding. Precisely, the audio and speech features are concatenated before feeding them as prefixes to the LLM. SALMONN shows unique reasoning capabilities over speech inputs overlaid with non-verbal audio.

Qwen-Audio. (Chu et al., 2023) Qwen follows a similar architecture to LTU, Pengi, and SALMONN, i.e., adding audio features as prefix to the model, and additionally employs a novel multi-task learning formulation for pre-training. More specifically, they append specific tags to specific parts of the instruction-response text pairs and train the model on diverse speech, non-speech, and music tasks. Post-pre-training, similar to GAMA, employs an instruction-tuning stage for alignment. The resultant model, Qwen-Audio-Chat, is able to respond to respond to diverse queries about the input speech and audio.

G Additional Details: CompA-R

G.1 Annotation and Annotator Details

As mentioned earlier, CompA-R was cleaned and CompA-R-test was verified by the paper authors themselves. To preserve anonymity, we briefly provide some details about the authors. All authors of the paper are either enrolled in or have graduated from a graduate degree (MS and/or Ph.D.). All authors have at least 2 years of professional research experience at a academic or industry lab. Their research experience spans across speech, audio and language processing. This provides them with adequate knowledge to faithfully complete the process.

For CompA-R-test verification, after at least 3 authors verified the test set, with proper rationales (which they were also asked to provide) the lead

1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163

author cross-verified all instances. The verification was done manually on local laptops and no kind of application was used which was made specifically for this. More details will be provided on camera-ready.

H Additional Details: General

H.1 GAMA Training Dataset Details

Table 6 shows statistics of all datasets used for fine-tuning and instruction-tuning GAMA. Table 7 shows statistics of CompA-R, which is sourced entirely from the AudioSet-Strong dataset.

Dataset	# Audio Samples	# QA Pairs
AudioSet-Strong	102K	636K
AudioSet	500K	441K
VGGSound	184K	336K
FSD50K	41K	82K
AudioCaps	46K	90K
FreeSound	91K	91K
Clotho	5K	32K
Sound Bible	1.2K	12K
NSynth(Instrument+Source)	301K	602K
Clotho AQA	1.5K	4.2K
MusicCaps	5.5K	2.8K
MusicQA	13.1K	118K
Magna	51.7K	51.7K
Sum (Closed-Ended)	1,217K	2,555K
AudioSet-Strong (Open-Ended)	91K	901K
AudioSet-20K	19K	184K
VGGSound (Open-Ended)	184K	907K
FSD50K (Open-Ended)	41K	403K
AudioCaps (Open-Ended)	46K	478K
Freesound (Open-Ended)	91K	791K
Clotho (Open-Ended)	5K	89K
Sound Bible (Open-Ended)	1.2K	10K
Sum (Open-Ended)	453K	3,764K
Total	1,670K	6,319K

Table 6: The statistics of the OpenAQA dataset.

Dataset	# Audio Samples	# QA Pairs
AudioSet-Strong	62613	200234
Total	62613	200234

Table 7: The statistics of the CompA-R dataset.

H.2 GAMA Evaluation Dataset Details

Table 8 shows statistics of all datasets used for evaluating GAMA. Table 10 shows statistics of CompA-R-test, which is sourced entirely from the AudioSet-Strong dataset.

¹<https://www.kaggle.com/datasets/modaresimr/sound-event-detection-audioset-strong>

²<https://zenodo.org/records/4060432>

³<https://www.tensorflow.org/datasets/catalog/nsynth>

⁴<https://zenodo.org/records/6473207>

Dataset	# Instances
AudioSet-Strong ¹⁸	102K
AudioSet	500K
VGGSound	184K
FSD50K ¹⁹	41K
AudioCaps	46K
FreeSound	91K
Clotho	5K
Sound Bible	1.2K
NSynth _{instrument} ²⁰	4K
NSynth _{source} ²¹	4K
Clotho AQA ²²	1.3K
GTZAN	3K
Medley-solos-DB	12.2K

Table 8: The statistics of the datasets used for evaluation of GAMA.

Dataset	Evaluation Metric
<i>Classification (zero-shot)</i>	
VocalSound (VS) (Gong et al., 2022)	Acc.
TUT 2017 (TUT) (Mesaros et al., 2018)	Acc.
Beijing Opera (BJO) (Tian et al., 2014)	Acc.
GTZAN (GTZ) (Park et al., 2022)	Acc.
Medley-solos-DB (MDB) (Lostanlen et al., 2018)	Acc.
<i>Classification (weak zero-shot)</i>	
DCASE2017 Task 4 (DCASE) (Mesaros et al., 2017)	Mi-F1
ESC-50 (Piczak, 2015)	Acc.
<i>Classification (seen)</i>	
VGGSound (VGG) (Chen et al., 2020)	Acc.
FSD50K (FSD) (Fonseca et al., 2021)	mAP
AudioSet (AS) (Gemmeke et al., 2017)	mAP
NSynth (NS) (Engel et al., 2017)	Acc.
<i>Captioning (vanilla & dense)</i>	
AudioCaps (Kim et al., 2019)	SPICE & Human
Clotho (Drossos et al., 2020)	SPICE & Human
<i>AQA (close-ended)</i>	
Clotho AQA (Lipping et al., 2022)	Acc.
<i>AQA (open-ended)</i>	
OpenAQA (Gong et al., 2024)	Human
<i>AQA (complex open-ended)</i>	
CompA-R-test (ours)	GPT-4 & Human

Table 9: List of evaluation datasets and their corresponding evaluation metrics for GAMA.

H.3 Other Details

Model Parameters: GAMA has a total of $\approx 7B$ parameters. Out of this, LLaMA-2-7B has 32 transformer-encoder layers and $\approx 6.7B$ parameters, the Audio Q-Former has $\approx 280M$ parameters, and our LoRA modules introduce 4.2 M learnable parameters for fine-tuning. The AST used in our experiments (audio-encoder of CAV-MAE (Gong et al., 2023)) has $\approx 85M$ parameters with 12 transformer-encoder layers, 768-hidden-state, and 12 attention-heads.

Compute Infrastructure: All our experiments are conducted on four NVIDIA A6000 GPUs. Training GAMA required four days of continuous training. Training GAMA-IT requires 4 hours of training. Pre-training Audio Q-Former requires 7 days each for stages 1 and 2.

Dataset	# Audio Samples	# QA Pairs
CompA-R-test	500	1561
Total	500	1561

Table 10: The statistics of the CompA-R-test dataset.

Implementation Software and Packages: We implement all our models in PyTorch ²³ and use the HuggingFace ²⁴ implementations of T5_{large} and the original implementation of HTSAT_{tiny} ²⁵.

For our baselines, we use the original GitHub repository provided by the authors: LAION-CLAP ²⁶, CompA-CLAP ²⁷, CLAP ²⁸, Wav2CLIP ²⁹, AudioCLIP ³⁰, MMT ³¹, ML-ACT ³², Pengi ³³, LTU ³⁴, AudioGPT ³⁵, SALMONN ³⁶, Qwen-Audio ³⁷.

Potential Risks. GAMA might encode biases from the pre-trained LLM or during its fine-tuning stage. Additionally, Audio Q-Former used as a backbone for audio-to-text/music generation, might generate synthetic audio that is misused.

²³<https://pytorch.org/>

²⁴<https://huggingface.co/>

²⁵<https://github.com/RetroCirce/HTS-Audio-Transformer>

²⁶<https://github.com/LAION-AI/CLAP/tree/main>

²⁷<https://github.com/Sreyan88/CompA>

²⁸<https://github.com/microsoft/CLAP>

²⁹<https://github.com/descriptinc/lyrebird-wav2clip>

³⁰<https://github.com/AndreyGuzhov/AudioCLIP>

³¹<https://github.com/akoepke/audio-retrieval-benchmark>

³²<https://github.com/akoepke/audio-retrieval-benchmark>

³³<https://github.com/microsoft/pengi>

³⁴<https://github.com/YuanGongND/ltu>

³⁵<https://github.com/aigc-audio/audiogpt>

³⁶<https://github.com/bytedance/salmonn>

³⁷<https://github.com/QwenLM/Qwen-Audio>


```

# Prompt 1

I will provide you with 2 different types of information about a 10-second audio clip:

1. A list where each comma-separated element indicates the individual events occurring in the audio at various time segments. For example, '(Speech-0.0-0.64)' would mean human speech between 0.0 second to 0.64 second.
2. A caption of the audio describing in a brief and abstract manner the scene in which the audio takes place.

I want you to act as a Prompt Generator. According to the event information and the caption, design some instructions and corresponding responses. The instruction should be designed in a way such that it can be answered only from the audio without the caption and any other detail provided. The instruction should involve one or more hops of complex knowledge and complex reasoning based on the scene created by the audio and the corresponding caption. Ensure that the knowledge and reasoning chains in the instructions are precise and sufficiently challenging, to the extent that only well-educated people and experts in the respective field can provide adequate responses.

The instructions must meet the following conditions:
1. Do NOT use phrases like 'according to the caption' in both the questions and answers; you should ask and answer as if you were observing the image by yourself.
2. The questions and answers should be as diverse as possible.
3. Please don't ask some simple questions about the intensity of the audio or the gender speaking the utterance; your questions must involve some knowledge.
4. Your instructions should not be answered directly based on the image and your instructions. Instead, it requires the test-taker to carefully observe the image and have a deep knowledge of the content within the image in order to answer correctly.
5. If a question cannot be answered, please do not ask.

Come up with 3 diverse instructions for the knowledge topics above with different language styles and accurate answers. The instructions should contain interrogative sentences and declarative sentences. The answers should be less than 30 words.

Output format, which is a list of jsons:

[['Instruction': instruction example, 'Answer': answer example, 'Knowledge topic': The specific knowledge topic], ['Instruction': instruction example, 'Answer': answer example, 'Knowledge topic': The specific knowledge topic], ...]
Here are some examples of inputs and outputs:

Input list of audio events: ['(Speech-0.0-0.64)', '(Mechanisms-0.0-10.0)', '(Dog-0.221-0.547)', '(Dog-0.803-0.966)', '(Generic impact sounds-0.885-1.129)', '(Tick-0.99-1.083)', '(Dog-1.432-1.665)', '(Speech-1.537-4.901)', '(Dog-1.921-2.119)', '(Dog-2.456-3.202)', '(Dog-3.434-3.597)', '(Dog-4.016-4.121)', '(Dog-4.936-5.39)', '(Generic impact sounds-5.204-5.611)', '(Dog-5.774-5.972)', '(Speech-5.984-6.787)', '(Tick-6.508-6.636)', '(Dog-6.717-8.266)', '(Generic impact sounds-7.649-8.277)', '(Laughter-8.347-9.488)', '(Dog-9.767-10.0)']
Caption: A baby cries while a woman laughs, creating a joyful and lively atmosphere in a domestic setting.

Output list of jsons: [['Instruction': 'Analyze the sounds in the audio and determine the most likely cause of the laughter heard towards the end of the recording. Consider the potential interactions between the different sound sources and their temporal overlaps.', 'Answer': 'The laughter likely results from the playful interaction between the dogs and the baby, as indicated by the overlapping sounds of dogs and the baby's presence.', 'Knowledge topic': 'Human and Animal Behavior Interpretation'], ['Instruction': 'From the given audio, infer the type of domestic setting depicted in the scene. Base your inference on the variety and sequence of sounds, particularly focusing on the interaction between the human speaking, the dog barking, and other background noises that may be there.', 'Answer': 'The setting is likely a home with an active family environment, evidenced by the continuous presence of dogs, speech, and everyday household sounds.', 'Knowledge topic': 'Environmental Acoustics and Domestic Soundscapes'], ['Instruction': 'Considering the duration and placement of speech and laughter in the audio, infer the possible emotional dynamics between the speakers. How do these elements interact to shape the scene's atmosphere?', 'Answer': 'The scene likely transitions from a more chaotic or lively mood and finally to a more joyful and relaxed atmosphere.']]

Input list of audio events: ['(Insect-0.0-0.724)', '(Mechanisms-0.0-9.777)', '(Female speech, woman speaking-0.737-1.434)', '(Bird vocalization, bird call, bird song-1.243-1.775)', '(Insect-2.376-3.182)', '(Female speech, woman speaking-3.386-3.509)', '(Insect-4.397-5.23)', '(Dog-7.906-8.78)', '(Surface contact-8.603-9.654)']
Caption: 'Birds chirp in the distance as a dog barks, creating a lively atmosphere in a peaceful outdoor setting.'
Output list of jsons: [['Instruction': 'What time of day this scene is likely set in?', 'Answer': 'The concurrent presence of insect and bird sounds suggests a natural, outdoor environment, possibly during early morning or evening when such wildlife is typically active.', 'Knowledge topic': 'Environmental Sound Analysis and Wildlife Behavior'], ['Instruction': 'Analyze the presence and timing of the dog's barking in the latter part of the audio. Considering the preceding sounds and infer the dog's possible reaction or behavior in this context.', 'Answer': 'The dog's barking following the peaceful nature sounds and speech could indicate a response to a new stimulus, possibly a visitor or an animal in the area.', 'Knowledge topic': 'Animal Behavior Analysis in Diverse Sound Environments'], ['Instruction': 'Deduce the woman's likely activity or purpose in this setting.', 'Answer': 'The woman might be engaging in an outdoor activity like gardening or bird-watching.', 'Knowledge topic': 'Human activity recognition through scene analysis']]

Input list of audio events: ['(Music-0.0-10.0)', '(Male singing-0.0-10.0)', '(Male speech, man speaking-0.354-1.364)', '(Male speech, man speaking-7.674-10.0)', '(Crowd-7.681-10.0)']
Caption: 'A basketball bounces while music plays, and a man speaks in an indoor stage environment.'
Output list of jsons: [['Instruction': 'Considering the presence of crowd sounds towards the end of the audio, deduce the nature of the event taking place. How do the elements of music, singing, and speech suggest the type of event and audience involvement?', 'Answer': 'The event seems to be a live performance or concert, with the crowd's reaction indicating an engaged and responsive audience, typical in such settings.', 'Knowledge topic': 'Event Atmosphere Analysis'], ['Instruction': 'Given the continuous presence of music and male singing throughout the audio, analyze the role of the man's speech in shaping the atmosphere of the scene. How does his speech, interspersed with music and singing, contribute to the overall environment?', 'Answer': 'The man's speech likely serves as commentary or narration, adding a personal or interactive element to the musical performance, enhancing the audience's engagement.', 'Knowledge topic': 'Music and Speech Dynamics'], ['Instruction': 'Identify the genre of music being played and explain how it complements the atmosphere of the indoor stage environment.', 'Answer': 'The genre is likely upbeat or energetic, enhancing the lively ambiance of a sports or performance event in an indoor setting.', 'Knowledge topic': 'Music Genre Detection and Scene Analysis']]

Input list of audio events: (timestamp events)
Caption: (caption)
Output list of jsons:

```

Figure 5: Prompts/Instructions used for caption augmentation with LLaMa-7B. Prompts are indexed according to the description in Section 3.2.

Complex AQA Evaluation Prompt

Please act as an impartial judge and evaluate the quality of the response provided with respect to the details provided. You will rate the quality of the response on multiple aspects, such as Helpfulness, Clarity, Factuality, Depth and Engagement. The response has been provided by an AI agent for a query related to an input audio, which the agent can perceive. You will be provided with 4 kinds of information for evaluating the response:

1. A list where each comma-separated element indicates the individual events occurring in the audio at various time segments. For example, '(Speech-0.0-0.64)' would mean human speech between 0.0 second to 0.64 second.
2. A caption of the audio describing in a brief and abstract manner the scene in which the audio takes place.
3. The question asked to the AI agent related to the audio.
4. An answer provided by an expert judge which you can consider as a reference.
5. An answer by the AI agent.

##Query: {query}

Evaluate

Aspects

- Helpfulness: Rate the response based on how well it addresses the users query about the audio and provides a relevant answer. A score of 5 indicates the answer fully aids the user, while a 1 suggests it offers little to no help.
- Clarity: Rate the response based on how well-structured it is, with ideas presented in a clear and coherent manner. A high score of 5 means the answer is clear and logically structured, while a 1 suggests a disjointed or confusing reply.
- Correctness: Evaluate the correctness or accuracy of the response provided with respect to the information provided to you. A perfect 5 indicates the response is entirely correct and accurate, while a 1 suggests it has significant errors or has not provided an answer to the question asked at all.
- Depth: Determine the level of detail and thoroughness in the response. A score of 5 means the answer delves deeply into the aspects of the input image for answering the question, while a 1 indicates it barely scratches the surface.

Format

Given the query and the extra information about the audio provided (the caption and comma-separated list of individual individual events), please rate the quality of the output by scoring it from 1 to 5, individually on **each aspect**. You are allowed to use all 3 information provided to you about the audio, in any way you want, to judge the response.

Now, please output your scores in the following json format by filling in the placeholders in [].

```
{ 'helpfulness': { 'reason': '[your rationale]', 'score': '[score from 1 to 5]' }, 'clarity': { 'reason': '[your rationale]', 'score': '[score from 1 to 5]' }, 'correctness': { 'reason': '[your rationale]', 'score': '[score from 1 to 5]' }, 'depth': { 'reason': '[your rationale]', 'score': '[score from 1 to 5]' }, 'engagement': { 'reason': '[your rationale]', 'score': '[score from 1 to 5]' } }
```

Only return the json and nothing else."

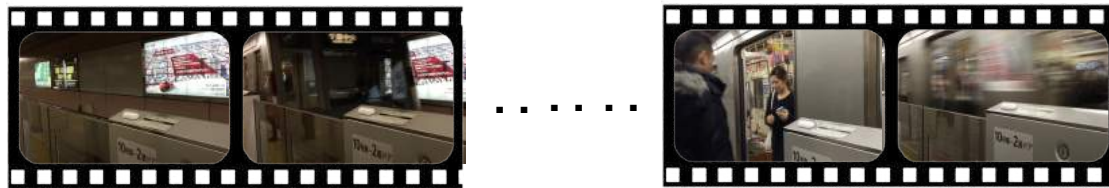
Figure 6: Prompts/Instructions used for evaluating the answers of models on CompA-R benchmark with GPT-4 as a judge.

Caption Evaluation Prompt

I want you to act as a Caption Evaluator. I will provide you with an audio caption generated by an AI agent. The agent was asked to generate a dense and detailed caption of the audio. To evaluate the caption, I will provide you with 2 different types of information about the 10-second audio clip:

1. A list where each comma-separated element indicates the individual events occurring in the audio at various time segments. For example, '(Speech-0.0-0.64)' would mean human speech between 0.0 second to 0.64 second.
2. A scene caption of the audio describing in a brief and abstract manner the scene in which the audio takes place. Using these two pieces of information, assign a score of 1-10 to the caption, where 1 is the lowest score and 10 is the highest score. Your evaluation should be based on the detailedness, correctness, and bluntness of the caption. Return a JSON with a single key 'score', where the value of the key is the score. Here are the details: {query}

Figure 7: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.



Question: Based on the combination of subway sounds and male speech, suggest the actions the man might be performing at this time.


LTU The man may be giving a lecture or presentation in the subway car while people are moving around him, possibly taking notes. **HE : 3**

SALMONN It is difficult to determine the exact actions the man might be performing based solely on the combination of subway sounds and male speech. However, it is possible that the man is giving a public announcement or warning to the passengers on the subway. **HE : 1**

Pengi It is a man. **HE : 1**

GAMA-IT The man is likely giving instructions to passengers or announcing a station stop. He may also be explaining the train's route or destination. **HE : 5**

Figure 8: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.



Question: Given the variety of sounds including mechanisms, generic impact sounds, and water splashes, determine the most probable activity occurring during the audio clip. Consider the presence of baby laughter in your analysis.

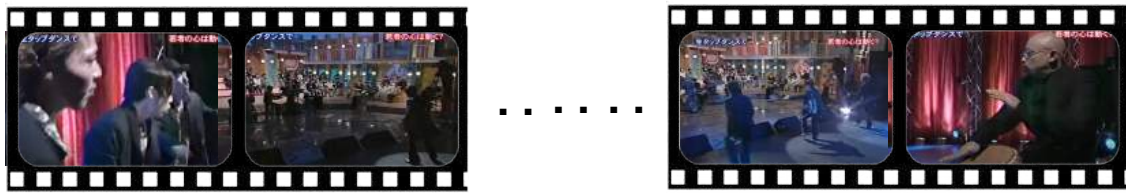
LTU The most probable activity occurring during the audio clip is a family or group of people having fun together, possibly playing with babies or toddlers. The mechanical sounds could be from a playground or other outdoor equipment. **HE : 2**

SALMONN Based on the variety of sounds including mechanisms, generic impact sounds, and water splashes, as well as the presence of baby laughter, the most probable activity occurring during the audio clip is likely a bath or a water-related activity. **HE : 5**

Pengi It is a combo. **HE : 1**

GAMA-IT The woman is likely bathing her child, as indicated by the sound of a bathtub filling up with water. **HE : 5**

Figure 9: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.



Question: Inferring from the nature and arrangement of sounds, what can be said about the skill level of the tap dancer?

LTU

The skill level of the tap dancer cannot be inferred from the audio clip.

HE : 1

SALMONN

The skill level of the tap dancer can be inferred from the nature and arrangement of sounds. If the sounds are complex and intricate, it suggests that the dancer is highly skilled. On the other hand, if the sounds are simple and repetitive, it suggests that the dancer may not be as skilled.

HE : 3

Pengi

It is a big crowd.

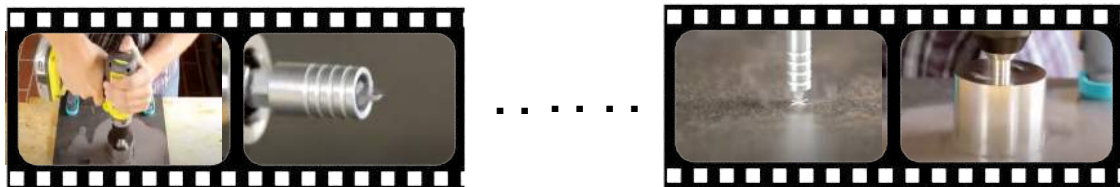
HE : 1

GAMA-IT

The tap dancing is likely performed by a professional or an expert. The rhythm and tempo are consistent throughout the entire duration of the clip.

HE : 5

Figure 10: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.



Question: Determine the type of power tool being used in this woodworking process based on the available sounds. Consider the mechanisms sounds.

LTU

Based on the available sounds, it is difficult to determine the exact type of power tool being used. However, the sharp and metallic sound suggests that a circular saw or a hand drill might be present in the scene.

HE : 1

SALMONN

Based on the available sounds, it is likely that a circular saw is being used in this woodworking process. The circular saw is a power tool that uses a spinning circular blade to cut wood, and the sound of the blade cutting through the wood can be heard in the audio. Additionally, the sound of the motor and other mechanisms of the circular saw can also be heard in the audio.

HE : 2

Pengi

It is a cameraman.

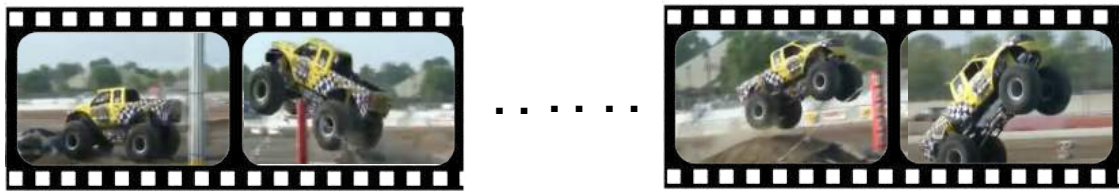
HE : 1

GAMA-IT

The sound is a drill, which is usually used for woodwork.

HE : 5

Figure 11: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.



Question: Analyze the audio and infer the likely engine type based on the sound of acceleration and revving.

LTU

Based on the sound of acceleration and revving, it is likely that the engine type is a high-performance or sports car engine. HE : 2

SALMONN

The audio contains the sound of an engine revving and accelerating. Based on this, it is likely that the engine type is a high-performance gasoline engine, such as a V8 or a turbocharged four-cylinder. The sound of the engine revving and accelerating is a common characteristic of high-performance engines, which are designed to produce a lot of power and torque. The sound of the engine can vary depending on the specific make and model of the vehicle, but the general characteristics of a high-performance engine can be identified by the sound of the engine revving and accelerating. HE : 2

Pengi

It is diesel. HE : 4

GAMA-IT

The truck is probably a diesel, as indicated by the heavy engine sounds and the absence of a typical gasoline engine's characteristic high-pitched whine at idle. HE : 5

Figure 12: Prompts/Instructions used for evaluating the captions generated by models Clotho and AudioCaps with GPT-4 as a judge.

Dataset	Category	Acoustic Property	Captions
AudioSet	Baby cry, infant cry	1: "a high-pitched, piercing wail" 2: "a persistent, rhythmic sobbing"	1. A tiny voice pierced the stillness of the night, demanding attention with its <i>high-pitched piercing wail</i> . 1. In the midst of the bustling market, a <i>high-pitched wail</i> rose above the chatter, signaling a tiny dissenter among us. 2. The <i>persistent, rhythmic sobbing</i> of a newborn punctuates the night's silence. 2. Amidst the rustling leaves, a baby's <i>rhythmic sobs</i> weave an unexpected lullaby.
	Stomach rumble	1: "a low, gurgling growl" 2: "a grumbling undercurrent"	1. As she leaned in for the kiss, a <i>low, gurgling growl</i> betrayed her pre-date fasting. 1. In the midst of the exam, a <i>low, gurgling growl</i> echoed from her stomach, punctuating the tension. 2. As the classroom fell into a hushed anticipation of the next lecture, a <i>grumbling undercurrent</i> betrayed the student's skipped breakfast. 2. Amidst the solemnity of the library, a <i>grumbling undercurrent</i> served as a comical reminder of lunchtime's approach.
ESCC50	siren	1: "Continuous, shrill alarm" 2: "Ear-splitting, cyclic blare"	1. The night was pierced by the <i>continuous, shrill alarm</i> of the siren, signaling an urgency that couldn't be ignored. 1. The relentless, <i>shrill alarm</i> of the siren wove through the corridors of the hospital, prompting swift movements from the staff. 2. The <i>ear-splitting, cyclic blare</i> cut through the night as the ambulance raced down the street. 2. The <i>ear-splitting, cyclic blare</i> of the air raid siren cast a shadow of dread over the city.
	sheep	1: "bleating", 2: "baaing"	1. In the tranquil meadow, a chorus of woolly creatures serenaded the dawn with their gentle <i>bleating</i> . 1. The shepherd smiled as the flock's <i>bleating</i> echoed through the valley, signaling a return to the fold. 2. A chorus of baaing accompanied the farmer as he made his morning rounds in the misty fields. 2. Under the starry sky, the gentle baaing of the flock blended with the whispers of the night.
NSynth	flute	1: "Melodious and silvery, carrying a light, airy tune that seems to float on the breeze." 2: "Clear and resonant, with a pure tone that sings above the orchestra like a bird in the morning sky." 1: "Thumping, providing a rhythmic pulse that can drive the beat of music." 2: "Booming, with a powerful, pervasive quality that can reverberate through a space."	1. The flute's melody weaved through the garden, mingling with the rustle of leaves. 1. Amidst the hustle of the market, the silvery flute tune danced above the clamor, a ribbon of tranquility. 2. The flute's crystalline melody soared above the orchestra, a lark greeting the dawn. 2. Amidst the rustle of the forest, the flute's song danced through the leaves, pure and high.
	bass	1: "Thumping, providing a rhythmic pulse that can drive the beat of music." 2: "Booming, with a powerful, pervasive quality that can reverberate through a space."	1. The bass pulsed through the dance floor, a heartbeat synchronizing every dancer's move. 1. Amidst the serene silence of the night, the bass from the distant festival throbbed like a gentle earthquake. 2. The bass pulsed like a second heartbeat, filling the room with its unyielding presence. 2. As the bassline dropped, it seemed to command the very air, a force unseen yet unforgettable.
FSD50K	Slam	1: "an abrupt, resonant boom that startles anyone nearby" 2: "a sharp, impactful smack as two hard surfaces collide with force"	1. The mailbox lid clapped shut, a resonant signal marking the departure of the day's correspondence. 1. The oven door's heavy thud resonated in the kitchen, a prelude to the aroma of freshly baked bread. 2. The kitchen was filled with the aroma of spices and the sharp smack of dough being forcefully thrown onto the countertop. 2. In the crisp morning air, the sharp smack of the newspaper hitting the doorstep announced the arrival of daily news.
	Dishes, pots, and pans	1: "Clanging and clattering" 2: "Metallic clinking and clunking"	1. A symphony of clanging and clattering announces the busy bustle of a restaurant kitchen in full swing. 1. The rhythmic clanging and clattering of pots and pans punctuate the air as grandma orchestrates her holiday feast. 2. The metallic clinking and clunking heralded the start of the dinner rush in the bustling restaurant kitchen. 2. A symphony of metallic clinking and clunking rose from the sink as grandma washed up after the family feast.
TUT Urban	bus	1: "a deep, rumbling engine", "2": "the low, steady hum of the diesel motor"	1. The city pulse beats with a <i>deep, rumbling engine</i> , heralding the arrival of the morning commute. 1. A gentle giant purrs in the stillness of dawn, its <i>deep, rumbling engine</i> announcing the start of a journey. 2. Market stalls buzz with life, their vibrant colors and smells underscored by the bus's <i>diesel hum</i> rolling down the avenue. 2. Leaves rustle in the autumn breeze, a natural chorus to the bus's <i>diesel motor humming</i> along the cobblestone path.
	residential area	1: "The symphony of children's laughter and chatter fills the air, punctuated by the occasional bark of a dog and the hum of lawn mowers in the distance." 2: "A serene hush blankets the neighborhood, broken occasionally by the soft whoosh of passing cars and the rustle of leaves stirred by a gentle breeze."	1. The neighborhood comes alive with the melody of playful banter and the sporadic chorus of canines. 1. Amidst the gentle drone of distant lawn mowers, the air vibrates with juvenile mirth and convivial exchanges. 2. The neighborhood rests under a tranquil silence, punctuated now and then by the whisper of tires on asphalt and the soft dance of leaves in the wind. 2. Calmness envelops the streets, save for the faint hum of vehicles gliding by and the tender shuffling of foliage in the zephyr's caress.
Urban-Sound 8K	air conditioner	1: "a steady humming" 2: "a low, monotonous droning"	1. The room filled with the steady humming of the air conditioner as they focused intently on their chess match. 1. A steady humming enveloped the library, where pages turned almost in rhythm with the air conditioning's constant song. 2. The air conditioner's low, monotonous droning became the unlikely lullaby for a midsummer's nap. 2. Amid the quiet study hall, the air conditioner's low, monotonous droning was a steady companion to the students' focused brows.
	gun shot	1: "A loud, sharp crack that echoes through the air." 2: "A thunderous boom that startles and reverberates."	1. The night's silence shattered with a loud, sharp crack echoing through the air. 1. A burst of sudden, sharp noise split the tranquil afternoon, reverberating off the canyon walls. 2. A thunderous boom startles a flock of birds into the sky, their wings flapping frantically against the silence that had just been. 2. The night's silence was shattered by a boom, reverberating through the alleyways and causing stray cats to scurry.
VGG Sound	mouse squeaking	1: "a high-pitched, sharp chirp" 2: "a soft, repetitive squeal"	1. In the moonlit barn, a tiny silhouette pauses to release its high-pitched, sharp chirp, disturbing the stillness of the hay-strewn loft. 2. Amidst the rustling leaves, a diminutive creature contributes its sharp chirp to the dusk chorus, a minuscule soloist in nature's vast orchestra. 3. A soft, repetitive squeal punctuated the silence of the old attic. 4. The cheese plate on the kitchen counter became the stage for a soft, repetitive squeal.
	typing on typewriter	1: "a rhythmic series of sharp clicks" 2: "a steady clatter of keys striking paper"	1. Fingers dance across keys, a rhythmic series of sharp clicks punctuating the silence of the library. 1. In the attic, a story unfolds to the staccato beat of a rhythmic series of sharp clicks. 2. Each steady clatter of keys striking paper weaves a tapestry of words, painting stories on the blank canvas. 2. In the dimly lit corner of the library, the rhythmic dance of metallic hammers against the page composes a silent symphony.

Table 11: Examples of prompt augmentations.

Original Caption	Augmented caption
A man speaks followed by the sound of shuffling cards in a small room.	<ol style="list-style-type: none"> 1. A deep, resonant voice fills the small room, accompanied by the soft shuffle of cards as they change hands, creating an intimate and deliberate atmosphere. 2. The sound of a man's voice echoes through the small space, punctuated by the subtle rustle of cards as they are shuffled and arranged, invoking a sense of purposeful deliberation. 3. A deep voice speaks, followed by the subtle shuffle of cards, creating an intimate and anticipatory atmosphere in the small room. 4. The gentle rustle of cards breaks the silence, punctuated by a man's voice, evoking a sense of anticipation and private reflection in the cozy space.
A person strums an acoustic guitar, creating melodic music with the sound of a bell ringing in the background.	<ol style="list-style-type: none"> 1. Soothing melodies flow from the acoustic guitar, harmonizing with the soft chime of a distant bell, crafting a peaceful ambiance. 2. The acoustic guitar's strings vibrate with grace, weaving a melodic tapestry that intertwines with the gentle ring of a bell, transporting the listener to a serene realm. 3. The gentle strumming of an acoustic guitar weaves a melodic tapestry, intertwined with the soft chime of a background bell, creating a soothing and harmonious atmosphere. 4. The rhythmic plucking of an acoustic guitar crafts a lively and uplifting melody, complemented by the delicate ringing of a background bell, transporting the listener to a serene and joyful realm.
Dogs bark while people talk in the background, creating a lively atmosphere in a field.	<ol style="list-style-type: none"> 1. Lively chatter and joyful barks fill the air, capturing the playful spirit of a sunny day in a field. 2. The rhythmic sounds of dogs barking and people talking blend together, creating a vibrant and lively ambiance in the open field. 3. The chatter of people and the joyful barks of dogs fill the air, creating a vibrant and lively atmosphere in the field. 4. The sound of playful dogs and lively conversation fills the field, evoking a sense of happiness and energy.
A man's voice is heard speaking over a radio as a vehicle passes by in the background.	<ol style="list-style-type: none"> 1. A clear, crisp voice pierces the airwaves, intertwining with the distant hum of a vehicle, creating an engaging audio experience. 2. The man's voice on the radio blends seamlessly with the subtle rumble of a passing vehicle, forming a captivating auditory tapestry. 3. A voiceover speaks over a radio, complemented by the distant hum of a vehicle passing by, creating a dynamic and engaging audio experience. 4. A man's voice broadcasts over the radio, intertwining with the subtle rumble of a vehicle in the background, forming a captivating audio landscape.
A woman speaks while a bird chirps in the background, creating a tranquil atmosphere in a natural setting.	<ol style="list-style-type: none"> 1. A gentle voice echoes through the forest, harmonizing with the chirping of birds, creating a soothing ambiance. 2. The sound of a gentle voice blends seamlessly with the melodic chirping of birds, transporting the listener to a serene natural setting. 3. The woman's gentle voice blends with the soothing chirps of a bird, creating a serene ambiance reminiscent of a peaceful afternoon in nature. 4. The woman's words are accompanied by the melodic chirping of a bird, transporting the listener to a calming and picturesque outdoor setting.
Water rushes as people talk in the background near a hot spring, creating a serene ambiance.	<ol style="list-style-type: none"> 1. Soothing waters create a peaceful ambiance, punctuated by the gentle chatter of people nearby, as if they are harmonizing with the soothing sounds of the hot spring. 2. The calm trickle of water creates an intimate atmosphere, with the soft murmur of voices in the background adding a sense of connection and tranquility to the space. 3. A soothing, babbling sound fills the air as people converse near a steaming hot spring, creating a tranquil atmosphere. 4. The gentle gurgling of water intertwines with the chatter of people in the background, crafting a peaceful and relaxing ambiance.
Soft music plays in the background as a speech is heard faintly, creating a calm and peaceful atmosphere.	<ol style="list-style-type: none"> 1. A soothing melody floats in the background, complementing the faint speech, creating a tranquil ambiance. 2. The soft strains of music blend with the subtle speech, fostering a sense of serenity and calmness in the atmosphere. 3. Soothing tunes fill the air, complemented by a gentle speech, creating an atmosphere of tranquility and serenity. 4. Mellow music and soft speech blend together, crafting a calming environment that soothes the senses.
A car engine revs up and then slows down, creating a vroom sound, as the vehicle accelerates in the audio.	<ol style="list-style-type: none"> 1. The car's engine purrs and then decelerates, emitting a smooth and powerful vroom sound as it shifts gears, creating a dynamic and energizing atmosphere. 2. The vehicle's engine roars to life, producing a bold and intense vroom sound as it speeds up, then gradually slows down, immersing the listener in a thrilling and exhilarating experience. 3. The car's engine purrs powerfully, then decelerates, creating a smooth and steady vroom sound as the vehicle gains speed. 4. The car's engine roars to life, building momentum with a series of sharp vroom sounds before shifting gears and slowing down.
Background music plays softly as the theme music gradually fades in, creating a melodic ambiance in an arena/performance setting.	<ol style="list-style-type: none"> 1. The arena comes alive with a subtle, soothing melody that gradually builds in intensity, creating an electrifying ambiance. 2. The soft strains of background music fill the air, setting the tone for an exhilarating performance in a vibrant arena setting. 3. Soft, melodic strains fill the air as the theme music subtly builds, establishing a harmonious ambiance in the arena. 4. The arena comes alive with a gentle, orchestral tune that gradually gains momentum, creating an uplifting and energetic atmosphere.

Table 12: Examples of caption augmentations.

Instruction-Response Pairs	AudioSet ID	Caption	Timestamp Events
<p>Instruction:Analyze the audio to understand the potential emotional state or mood of the man. How does the progression from typing to speech to chewing reflect his transition through different phases of work or activity? Output:The man initially seems engaged and focused during the typing and speaking portion, which might then transition into relaxation during the break, suggested by the chewing sound.</p>	YCecEf0abd4Y	A man speaks while typing on a keyboard in a small room, followed by the sound of chewing.	'(Generic impact sounds-0.0-1.037)', '(Background noise-0.0-10.0)', '(Generic impact sounds-1.191-1.421)', '(Generic impact sounds-2.01-2.202)', '(Generic impact sounds-2.343-2.574)', '(Male speech, man speaking-2.727-3.393)', '(Generic impact sounds-3.163-3.406)', '(Generic impact sounds-3.585-3.905)', '(Generic impact sounds-4.136-4.379)', '(Breathing-4.405-4.917)', '(Generic impact sounds-4.93-5.288)', '(Generic impact sounds-5.442-5.608)', '(Generic impact sounds-5.736-6.12)', '(Generic impact sounds-6.274-6.569)', '(Breathing-6.825-7.26)', '(Generic impact sounds-6.863-7.042)', '(Male speech, man speaking-7.81-8.873)', '(Generic impact sounds-8.041-8.348)', '(Breathing-9.001-9.36)', '(Human sounds-9.014-9.181)', '(Generic impact sounds-9.309-9.565)', '(Scrape-9.449-10.0)'
<p>Instruction:Considering the diverse array of sounds present in the audio, from insects to birds to a dog, infer the type of ecosystem this outdoor setting might represent. What does the combination of these sounds tell us about the biodiversity and potential human impact in this area? Output:The ecosystem is likely a suburban or rural area with a mix of wildlife and human habitation, indicated by the variety of animal sounds and intermittent female speech.</p>	YcQiEI7HLGJg	Birds chirp in the distance as a dog barks, creating a lively atmosphere in a peaceful outdoor setting.	'(Insect-0.0-0.724)', '(Mechanisms-0.0-9.777)', '(Female speech, woman speaking-0.737-1.434)', '(Bird vocalization, bird call, bird song-1.243-1.775)', '(Insect-2.376-3.182)', '(Female speech, woman speaking-3.386-3.509)', '(Insect-4.397-5.23)', '(Dog-7.906-8.78)', '(Surface contact-8.603-9.654)'
<p>Instruction:Given the presence of mechanisms throughout the audio and the interspersed generic impact sounds, infer the type of machinery that is likely operating in the background and its commonality in an office environment. Output:The continuous mechanism sound suggests a printer or copier, which are common in office settings.</p>	YXQ2XAXx7mKs	A printer hums while people converse in the background, creating a typical office ambiance.	'(Generic impact sounds-0.0-0.622)', '(Mechanisms-0.0-10.0)', '(Generic impact sounds-0.815-1.227)', '(Generic impact sounds-1.632-2.134)', '(Child speech, kid speaking-3.591-6.684)', '(Squeal-7.385-7.612)', '(Child speech, kid speaking-8.437-10.0)'

Table 13: Examples of CompA-R

Instruction-Response Pairs	AudioSet ID	Caption	Timestamp Events
Instruction: From the sequencing and overlapping of different sound events, infer the likely cause and process of the dripping sounds heard intermittently throughout the audio. Output: The dripping sounds may be a result of water overflowing from a filled sink or bath, supported by the earlier sounds of splashing and liquid gurgling.	YCU9A5xL3TVc	Water splashes and gurgles as it drips inside a small room, creating a soothing ambiance reminiscent of a tranquil bathroom.	'(Background noise-0.0-10.0)', '(Generic impact sounds-0.083-0.331)', '(Splash, splatter-0.67-1.174)', '(Liquid-1.385-1.956)', '(Splash, splatter-2.325-3.138)', '(Liquid-3.085-4.131)', '(Liquid-4.372-5.5)', '(Drip-4.949-5.047)', '(Drip-5.279-5.458)', '(Generic impact sounds-8.819-9.142)', '(Drip-9.511-9.649)'
Instruction: Assess the style and elements of the rapping and music in this audio clip. Based on this, what genre of music might the DJ be playing on stage? Output: Given the presence of rapping and electronic music, the DJ is likely playing Hip Hop or Electronic Dance Music (EDM).	YRjUZjMPP-nA	Electronic music plays as a whoosh sound follows, creating a lively atmosphere for the DJ performing on stage.	'(Rapping-0.0-0.376)', '(Music-0.0-10.0)', '(Rapping-0.685-1.663)', '(Rapping-2.295-2.837)', '(Sound effect-2.423-5.222)', '(Sound effect-7.427-10.0)'

Table 14: Examples of CompA-R

Instruction-Response Pairs	AudioSet ID	Caption	Timestamp Events
Instruction: Identify the role of the crowd's continuous conversation during the woman's speech. How does it contribute to the atmosphere of the scene? Output: The crowd's constant chatter indicates an informal or relaxed environment, suggesting that while the woman's speech is central, other conversations are also ongoing in the background.	Y6fRYeClf5U4	A woman delivers a speech while a crowd of people engage in conversation in an urban setting.	'(Crowd-0.0-10.0)', '(Wind-0.008-10.0)', '(Female speech, woman speaking-0.074-1.65)', '(Female speech, woman speaking-2.879-5.427)', '(Female speech, woman speaking-5.604-6.083)', '(Female speech, woman speaking-6.9-10.0)'
Instruction: Considering the presence of bird sounds and chainsaw noise, identify the probable location where this sound sequence is occurring. Output: The location is likely an outdoor area, possibly in a woodland or suburban setting where birds can be heard alongside human activity and power tools.	YbkG4M4TiXZg	A man speaks while a power tool revs up, indicating a speech event followed by the sound of a chainsaw.	'(Male speech, man speaking-0.0-0.268)', '(Chainsaw-0.0-10.0)', '(Male speech, man speaking-1.772-4.425)', '(Male speech, man speaking-5.008-8.118)', '(Bird vocalization, bird call, bird song-5.362-7.512)', '(Bird vocalization, bird call, bird song-8.244-8.709)', '(Bird vocalization, bird call, bird song-8.937-9.283)', '(Male speech, man speaking-9.661-10.0)'

Table 15: Examples of CompA-R-test

Instruction-Response Pairs	AudioSet ID	Caption	Timestamp Events
<p>Instruction:Analyze the frequency and consistency of the ticking sound. What is the likely source of this sound, and what implication might it have on the depicted setting? Output:The steady ticking likely comes from a clock, which coupled with the soft music, suggests a relaxed, cozy domestic setting, perhaps aimed at unwinding or relaxation.</p>	YCoBAR5Mbjys	The clock ticks steadily as soft music plays in the background, creating a calming atmosphere in a cozy living room.	'(Mechanisms-0.0-10.0)', '(Alarm clock-0.008-10.0)', '(Tick-0.386-0.583)', '(Tick-1.071-1.22)', '(Tick-1.764-1.906)', '(Tick-2.465-2.638)', '(Tick-3.197-3.331)', '(Tick-3.772-3.976)', '(Tick-4.346-4.48)', '(Tick-4.646-4.787)', '(Tick-5.087-5.22)', '(Tick-5.669-5.795)', '(Tick-6.031-6.15)', '(Tick-6.37-6.528)', '(Tick-6.724-6.795)', '(Tick-6.969-7.118)', '(Tick-7.386-7.614)', '(Tick-8.134-8.354)', '(Tick-8.882-9.094)', '(Tick-9.315-9.425)', '(Tick-9.575-9.685)'
<p>Instruction:Identify the type of vocal music that is being depicted in the audio based on the presence of singing and beatboxing. Output:This audio resembles A Capella, where voices impersonate the sounds of instruments, including rhythms often mimicked through beatboxing.</p>	Y6SvDRiIG2NY	A group of people sing and harmonize, creating vocal music with occasional beatboxing, in a room with a piano.	'(Male singing-0.0-6.594)', '(Music-0.0-10.0)', '(Mechanisms-0.0-10.0)', '(Breathing-7.064-8.314)', '(Breathing-8.911-10.0)', '(Male singing-9.713-10.0)'
<p>Instruction:Based on the audio, ascertain the possible relationship between the gunfire sounds, artillery fire, and music. How does the sequencing and manner of these sounds contribute to the atmosphere of the scene? Output:The gunfire and artillery sounds likely serve as a ceremonial display, with the music adding to the grandeur and solemnity of a military parade.</p>	YbJvOp4gmHBg	Gunshots and artillery fire echo through the air as music plays during a military parade at a raceway.	'(Music-0.0-10.0)', '(Generic impact sounds-0.166-0.307)', '(Artillery fire-0.32-0.704)', '(Generic impact sounds-0.781-0.948)', '(Generic impact sounds-1.063-1.165)', '(Generic impact sounds-1.524-1.677)', '(Generic impact sounds-2.625-2.881)', '(Artillery fire-3.035-3.521)', '(Generic impact sounds-3.611-3.777)', '(Generic impact sounds-4.213-4.43)', '(Generic impact sounds-5.096-5.262)', '(Artillery fire-5.288-5.762)', '(Generic impact sounds-5.89-6.095)', '(Generic impact sounds-6.479-6.812)', '(Generic impact sounds-6.94-7.106)', '(Artillery fire-7.222-7.606)', '(Generic impact sounds-8.207-8.425)', '(Artillery fire-8.476-8.988)', '(Generic impact sounds-9.206-9.385)', '(Generic impact sounds-9.654-9.795)'

Table 16: Examples of CompA-R-test