

Hear you are: Teaching LLMs Spatial Reasoning with Vision and Spatial Sound

Anonymous CVPR submission

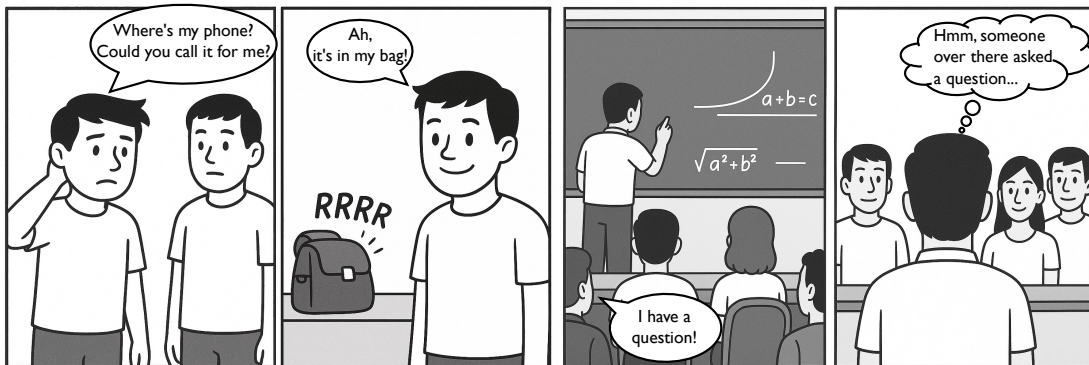


Figure 1. **Audio-Visual Spatial Reasoning.** (Left) A phone rings out of sight inside a bag; although the sound’s semantic cue (“ring tone”) is present, spatial reasoning is required to locate the true source among visually silent objects. (Right) In a classroom, several students share the same semantic cue (“speech”), so the teacher must rely on spatial audio to identify which student asked the question. These examples illustrate that accurate audio-visual understanding demands not only semantic alignment but also spatial comprehension.

Abstract

001 *Many audio-visual learning methods have focused on align-*
002 *ing audio and visual information, either through semantic*
003 *or temporal correspondence. However, most of these works*
004 *have utilized monaural audio, which does not contain in-*
005 *formation about the spatial location of the sound source. In*
006 *contrast, humans and other animals utilize binaural hear-*
007 *ing to perceive this spatial information. Combining spatial*
008 *sound and visual perception enables powerful high-level*
009 *reasoning: for example, a person looking for their phone*
010 *may hear the ringing sound coming from a backpack sitting*
011 *on a table, and quickly infer that the missing phone is in-*
012 *side the backpack. In this paper, we investigate the problem*
013 *of **Audio-Visual Spatial Reasoning**. We design a spatial*
014 *audio-visual question answering dataset to cover scenar-*
015 *ios where semantic correspondence between audio and vi-*
016 *sual signals is absent but spatial alignment exists, as well*
017 *as cases with multiple audio-visual semantic correspondences*
018 *that require spatial reasoning to disambiguate. We propose*
019 *a model that learns spatial comprehension across the audio*
020 *and vision modalities by connecting them with a large lan-*
021 *guage model and experimentally demonstrate that spatial*
022 *sound perception is an essential part of our task.*

1. Introduction

We live in a world full of sights and sounds, naturally associating what we hear with what we see. Several cues help us connect the two, such as the visual appearance and audible characteristics of an object, the synchronization between an action or event and its corresponding sound, and the direction from which the sound arrives, through binaural hearing. We rely on these audio-visual cues to locate a missing mobile device, or to know when an emergency vehicle is approaching as we are driving. This natural ability to connect auditory and visual information has motivated advancements in audio-visual machine learning, such as sound source localization (object detection based on audio queries) [8, 21, 23, 27, 29, 30, 33], source separation [3, 14–16, 45, 46, 48], and audio-visual synchronization [7, 12, 32]. However, most of these studies, which commonly use monaural audio, focus on the semantic correspondence between a sound and the visual appearance of the object that made the sound, or the audio-visual temporal alignment between an event and the sound it creates. These past approaches often overlook spatial cues that provide information about where a sound is coming from.

Binaural audio becomes essential when semantic matching is ambiguous or misleading. Figure 1 illustrates two

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

042

043

044

045

046

| | | |
|-----|---|-----|
| 047 | scenarios where spatial reasoning is necessary. For in- | 100 |
| 048 | stance, understanding that a ringtone sound is emanating | 101 |
| 049 | from a backpack requires spatial reasoning, as the backpack | 102 |
| 050 | does not semantically match the sound. Another example is | 103 |
| 051 | when a single sound (e.g., speech) could correspond to mul- | 104 |
| 052 | multiple visual objects (e.g., several students in a classroom), | 105 |
| 053 | where spatial cues help pinpoint the actual source. These | |
| 054 | examples highlight the limitations of previous methods, em- | |
| 055 | phasizing the need to address spatial reasoning beyond basic | |
| 056 | perception. | |
| 057 | Previous studies in spatial audio reasoning have primarily | |
| 058 | focused on audio-only approaches, excluding visual in- | |
| 059 | formation while incorporating language as a modality for | |
| 060 | spatial interpretation. [13] aligns audio and text embed- | |
| 061 | dings for spatial tasks, while [52] leverages large language | |
| 062 | models for spatial audio question answering. While spatial | |
| 063 | audio itself provides rich information for spatial reasoning, | |
| 064 | integrating visual information into these tasks is a natural | |
| 065 | extension, as visual signals inherently convey spatial con- | |
| 066 | text. This combination not only enhances spatial perception | |
| 067 | and localization capabilities, but also enables more sophis- | |
| 068 | ticated spatial reasoning, such as handling scenarios involv- | |
| 069 | ing sounding sources and nearby visual objects. | |
| 070 | In this paper, we address the problem of Audio-Visual | |
| 071 | Spatial Reasoning , which involves understanding the spa- | |
| 072 | tial relationship between a sound and the visual context. | |
| 073 | This task goes beyond simply perceiving and localizing a | |
| 074 | sound source, as it requires reasoning about spatial cues | |
| 075 | to infer relationships and interactions between objects. To | |
| 076 | support research on this problem, we construct a large- | |
| 077 | scale dataset of 1 million question-answer pairs, specifi- | |
| 078 | cally designed to serve as both the training and evalua- | |
| 079 | tion set for spatial audio-visual reasoning in diverse sce- | |
| 080 | narios. The vision and spatial audio is rendered using | |
| 081 | SoundSpaces 2.0 [5], with source audio clips sampled | |
| 082 | from VGGSound[6]. 3D objects associated with these | |
| 083 | sounds are generated using Stable Diffusion 3[35] and In- | |
| 084 | stantMesh [49], and then are placed within the virtual en- | |
| 085 | vironments. This dataset serves as a comprehensive bench- | |
| 086 | mark for spatially intricate settings, providing questions that | |
| 087 | assess spatial alignment between modalities, relative loca- | |
| 088 | tions between sounding and non-sounding objects, and lo- | |
| 089 | calization of sound sources among multiple visual objects | |
| 090 | of the same category as the query audio. | |
| 091 | Furthermore, we propose a multi-modal framework, | |
| 092 | Hear You Are LLM, which leverages spatial audio and vi- | |
| 093 | sual encoders to integrate spatial information. The model is | |
| 094 | trained to handle all the spatial reasoning tasks from our | |
| 095 | dataset, enabling it to address scenarios where semantic | |
| 096 | alignment alone is insufficient. We experimentally demon- | |
| 097 | strate that our proposed method effectively addresses the | |
| 098 | audio-visual spatial reasoning problem, outperforming ex- | |
| 099 | isting baseline models including a state-of-the-art monau- | |
| | ral sound source localization method [39, 40] and a large | 100 |
| | language model-based audio-visual model that lacks spatial | 101 |
| | understanding. These results highlight the importance of | 102 |
| | incorporating spatial audio-visual knowledge to achieve ro- | 103 |
| | burst multi-modal reasoning. To summarize, our main con- | 104 |
| | tributions are as follows: | 105 |
| | • We define a new task, audio-visual spatial reasoning, | 106 |
| | focusing on understanding spatial relationships between | 107 |
| | sound and visual context, going beyond basic semantic | 108 |
| | perception such as sound source localization (object de- | 109 |
| | tection based on audio queries) and audio-visual segmen- | 110 |
| | tation. | 111 |
| | • We propose <i>Hear You Are LLM</i> , a multi-modal model- | 112 |
| | ing framework that integrates spatial audio and visual en- | 113 |
| | coders with a large language model to handle complex | 114 |
| | spatial reasoning tasks. | 115 |
| | • We construct <i>Hear You Are QA</i> , the first large-scale | 116 |
| | dataset specifically designed for audio-visual spatial rea- | 117 |
| | soning, consisting of 1 million question-answer pairs | 118 |
| | across diverse spatial scenarios for training and evalua- | 119 |
| | tion. We will open source both the dataset and the training | 120 |
| | code. | 121 |
| | 2. Related Works | 122 |
| | 2.1. Audio-Visual Sound Source Localization | 123 |
| | Audio-visual sound source localization is the task of detect- | 124 |
| | ing the object or area that corresponds to the query audio in | 125 |
| | the visual scene. Following the development of deep learn- | 126 |
| | ing, Senocak <i>et al.</i> [37, 38] suggested a semantic alignment- | 127 |
| | based approach by proposing a cross-modal attention mech- | 128 |
| | anism with contrastive learning. The field has advanced in | 129 |
| | the direction of better cross-modal alignment by leveraging | 130 |
| | negative-free self-supervised learning [42], intra-modality | 131 |
| | similarity learning [43], and the use of multiple positive | 132 |
| | learning [39], aligning with representation learning meth- | 133 |
| | ods. However, these methods rely on monaural audio and | 134 |
| | are limited to audio-visual semantic correspondence with- | 135 |
| | out spatial understanding. | 136 |
| | Different approaches have focused more on spatial au- | 137 |
| | dio for sound source localization. He <i>et al.</i> [20] proposed | 138 |
| | a 3D sound source localization method trained on a dataset | 139 |
| | with four-channel audio and multi-view visual scenes syn- | 140 |
| | thesized using SoundSpaces 2.0. Their approach localizes | 141 |
| | sound within the visual scene, but the visual counterpart of | 142 |
| | the sound is not visible in their setting, as they only lo- | 143 |
| | calize the area of the sound source. Shimada <i>et al.</i> [41] | 144 |
| | constructed an audio-visual sound source localization and | 145 |
| | detection dataset in which audio-visual alignment is guar- | 146 |
| | anteed. In their framework, the visual signal serves as an | 147 |
| | auxiliary modality to improve sound localization and de- | 148 |
| | tection. In contrast, we present an audio-visual scene that | 149 |
| | includes both sound-producing and silent objects, allowing | 150 |

151 the model to learn a broader range of spatial reasoning tasks
152 that require contextual understanding beyond basic local-
153 ization.

154 2.2. Spatial Audio Reasoning

155 Following recent advancements in audio understanding [1,
156 18, 24] and reasoning [19, 36], several approaches have
157 been proposed to address spatial audio reasoning. [52] syn-
158 thesize the spatial sound question answering dataset with
159 the SoundSpaces 2.0 simulator and train a spatial audio en-
160 coder and a large language model for spatial audio under-
161 standing and reasoning. This framework handles tasks such
162 as sound event detection, direction and distance estimation,
163 and spatial reasoning, for example, “What is the sound on
164 the left side of the sound of the dog barking?” Another
165 line of research explores spatial audio reasoning through
166 contrastive language-audio pretraining, with synthetic first-
167 order ambisonics [13]. However, these approaches do not
168 incorporate the vision modality, which opens another di-
169 mension for reasoning.

170 2.3. Audio-Visual LLMs

171 Inspired by the advancements of Large Language Mod-
172 els (LLMs), recent studies have extended these mod-
173 els to Multimodal Large Language Models (MLLMs) to
174 tackle a wider range of multimodal tasks. In the audio-
175 visual domain, GroundingGPT [26] introduces multimodal
176 grounding for audio, image, and video data using LLMs.
177 Meerkat [10] aligns audio-visual features using optimal
178 transport and attention consistency, and CAT [51] ag-
179 gregates question-related clues in audio-visual scenarios.
180 From a benchmarking standpoint, AVHBench, AVTRUST-
181 BENCH, and AV-Odyssey Bench [11, 17, 44] provide
182 comprehensive benchmarks targeting hallucination detec-
183 tion [44], reliability and robustness [11], and both foun-
184 dational capabilities and high-level reasoning [17]. While
185 recent studies have advanced multimodal learning, they pri-
186 marily rely on monaural audio, limiting their ability to han-
187 dle spatial reasoning. As spatial reasoning enables a broader
188 range of tasks and more closely reflects real-world scenar-
189 ios, it must be addressed to achieve comprehensive audio-
190 visual understanding. We propose a new dataset and model
191 specifically designed for spatial reasoning in audio-visual
192 tasks.

193 3. Creation of Hear You Are QA Dataset

194 Our goal is to train a model to learn both semantic and spa-
195 tial reasoning, for audio-visual inputs. To this end, we in-
196 troduce the Hear You Are QA Dataset. Constructing large-
197 scale audio-visual scene data with real-world spatial audio
198 is time-consuming and challenging, requiring specialized
199 equipment such as ambisonic or dummy head microphones.
200 To efficiently build a diverse dataset with various objects

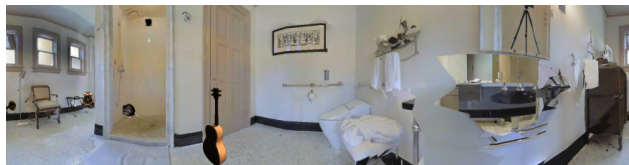


Figure 2. **Image sample from Hear You Are QA dataset.** The dataset consists of diverse indoor scenes captured in 360° panoramic views, featuring various object arrangements and providing a comprehensive range of spatial contexts for analysis.

and sound events, we adopt a simulation-based approach to
generate both the scenes and spatial audio.

Spatial Audio Simulator. We employ the SoundSpaces 2.0 simulator [5], which renders geometry-based acoustics, adding realistic reverberation for any source–receiver pair. Users can freely vary wall materials, object properties, and microphone-array geometry, letting us create a rich, controllable dataset while retaining exact ground-truth parameters, e.g., every source’s 3D position and orientation. Scene meshes come from Matterport3D [2], a collection of 90 fully scanned buildings averaging 24.5 rooms across 2.61 floors and 517.34 m² of floor space. We use 72 scenes for training, 9 for validation and 9 for testing. Given a source location, monaural signal, receiver position, and heading, the observed signal is obtained by convolving the monaural signal with the environment’s room impulse response. We configure the receiver to record a binaural audio signal with the default Head Related Transfer Function (HRTF) provided by SoundSpaces2.0.

Sound Sources. Previous spatial audio datasets include either a limited number of class categories [41] or classes that are not guaranteed to be visually observable [20, 52]. To construct a large-scale audio-visual dataset, we adopt VGGSound [6], which contains 200,000 in-the-wild 10-second YouTube clips, each annotated with one of 309 audio event classes. However, some of these classes correspond to events that typically occur outdoors or are difficult to associate with a single visual object (e.g., “Airplane Flyby”, “People Marching”). To enhance the visual reliability and realism of our dataset, we manually exclude categories typically occur outdoors, or are visually ambiguous. We follow the original testing splits provided by VGGSound, and create a validation set of the same size as the testing set by sampling clips from the VGGSound training split.

Visual Objects. Due to the limited number of sound-emitting categories in existing 3D object datasets, we generate our own 3D objects to be placed within the Matterport3D environments, either as sounding objects or as distractor objects. Specifically, we first select 150 class categories from VGGSound and 40 from ImageNet, and generate 2D images for each category using Stable Diffusion 3. After manually filtering out low-quality or unrealistic gen-

Table 1. Spatial audio visual question types and base templates.

| |
|---|
| Q1. Spatial Correspondence |
| Q: What is the sound class category? Where is the sound coming from? |
| A: phone ringing; cupboard |
| Q2-4. Relative Location (Distance, Direction, Angle) |
| Q: Is the sound source of the siren closer to the agent than it is to the cat? |
| A: Yes |
| Q: Can you estimate the distance from the accordion sound to the dog, and the relative location of the accordion from the dog? |
| A: right; behind; upper; 2.3 m |
| Q: Can you estimate the distance from the accordion sound to the dog, and the angle between the agent’s gaze directions toward the accordion and the dog? |
| A: 30; 10; 2.3 m |
| Q5. Spatial & Semantic Correspondence (One visual object semantically matches the audio) |
| Q: What is the object in the scene located at $(-30, -12)$, 2.549 m? Is it making a sound? |
| A: bird squawking; making sound |
| Q6. Spatial & Semantic Correspondence (Multiple visual objects semantically match the audio) |
| Q: What is the object in the scene located at $(150, -14)$, 1.735 m? Is it making a sound? |
| A: canary calling; making sound |
| Q7. Spatial & Semantic Correspondence (One visual object semantically matches the audio) |
| Q: Given multiple visual objects, which one is making a sound, and where is it located? |
| A: bird squawking; -30 ; -12 ; 2.549 m |
| Q8. Spatial & Semantic Correspondence (Multiple visual objects semantically match the audio) |
| Q: Could you determine the sound class category, and which object of that category in the scene is making the sound? |
| A: canary calling; 150; -14 ; 1.735 m |
| Q9. Semantic Co-occurrence |
| Q: What is the sound class category? Is the sound source visible in the scene? |
| A: cat; not visible |

erations, we select 40 visually plausible images per category. These 2D images are then lifted into 3D object meshes using the method from [49]. For each sounding object category, we reserve 32 images for training, 4 for validation and 4 for testing.

Audio-Visual Scene Construction. Each audio-visual scene consists of a 360° panoramic image as Figure 2 and corresponding binaural audio. We stitch 18 images, each with a horizontal FoV of 20 degrees as in [4], to form a 360° view. The final image resolution is set to 224×812 , and the center of the image is aligned with the front-facing direction of the observing agent in SoundSpaces 2.0.

We inject the aforementioned sound source and 3D objects into random locations within the scene, excluding placements where objects are occluded by walls or located in a different room. Each scene includes one sound source. The sound source, depending on the question scenario, is assigned to either a semantically matching object from a VGGSound category, a random object from a different category (VGGSound or ImageNet), or a random empty location within the scene.

One potential concern is that rendering artifacts, such as visible seams between injected objects and the original

scene, could serve as shortcuts for the model. To mitigate this and increase the visual complexity of the scene, we randomly insert up to three random objects sampled from categories distinct from the main visual objects in the scene.

Crafting Questions. We manually defined nine different “base” questions that require spatial audio-visual understanding, summarized in Table 1. When filling a question template, we use handcrafted rules to automatically populate the missing fields in the question and answer using the scene construction parameters. The questions cover four main categories: spatial correspondence (Q1), relative location (Q2, Q3, Q4), spatial and semantic correspondence (Q5, Q6, Q7, Q8), and semantic co-occurrence (Q9).

Spatial Correspondence questions aim to evaluate whether the model can correctly associate an audio signal with its spatially aligned visual source. To assess the model’s robustness, we include counterfactual examples in which semantically mismatched visual objects and sounds (e.g., a piano and dog barking) are placed at the same location. This setting discourages reliance on semantic priors and encourages the model to learn true spatial correspondence between audio and visual modalities without hallucination. **Relative Location** questions assess the model’s ability to understand the spatial relationship between audio and visual information. These include determining whether a sound source is located to the left, right, front, or behind the agent, as well as reasoning about vertical position (e.g., above or below), angular direction, and relative distance with respect to a visual reference. **Spatial and Semantic Correspondence** questions evaluate whether the model can jointly associate the correct object class (semantic) and its location (spatial) based on the audio signal. **Semantic Co-occurrence** questions focus on learning spatial audio understanding regardless of whether the corresponding visual object is explicitly visible, encouraging the model not to solely rely on an object’s appearance. To diversify the question set and improve naturalness, we utilize ChatGPT-4o to paraphrase and expand each base question into multiple human-like variations.

4. Method

Our aim is to construct a model that can answer the questions in our proposed dataset by leveraging both visual and spatial audio inputs. To this end, we design and train a multi-modal large language model with both visual and binaural audio inputs. The overall architecture is illustrated in Figure 3.

Audio and Visual Encoders with Projector. Given an image v and its corresponding audio a , our backbone networks extract features from each modality. The vision encoder f_v processes a panoramic image frame and outputs a sequence of spatially aligned visual tokens, $\mathbf{v} \in \mathbb{R}^{N_v \times C_v}$,

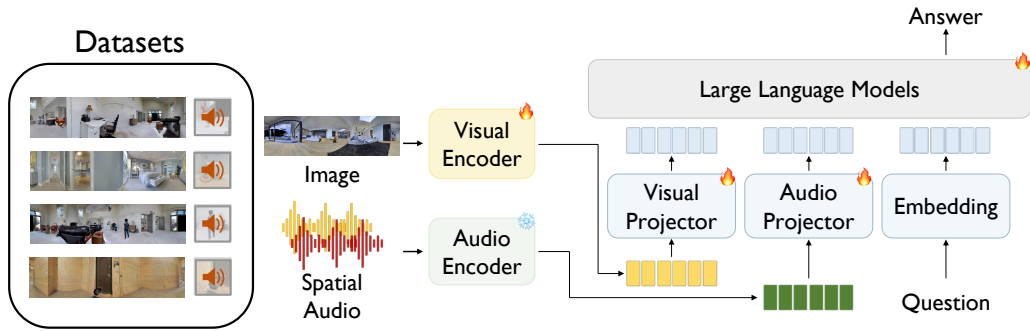


Figure 3. **The pipeline of our framework:** feature extraction, projection, and multimodal reasoning. We extract spatial audio and visual features using pre-trained encoders, project them into a shared embedding space, and integrate the embeddings with the question embedding to generate the answer.

317 where N_v is the number of visual tokens and C_v is the feature
 318 dimension of each token. We preserve the full spatial
 319 layout of patch tokens without pooling. The audio encoder
 320 f_a takes the input spectrogram of a and produces a set of
 321 audio tokens, $\mathbf{a} \in \mathbb{R}^{N_a \times C_a}$, where N_a is the number of audio
 322 tokens and C_a is the corresponding feature dimension. Each
 323 modality-specific encoder is followed by a projector that maps
 324 the extracted features into the hidden dimension of the language
 325 model. The visual projector attends to the spatial visual features
 326 to generate N_V projected tokens, and the audio projector
 327 similarly produces N_A tokens from the audio features. These
 328 projected tokens are then passed to the large language model for
 329 multi-modal reasoning.

330 **Large Language Model.** To bridge the audio and visual
 331 encoders, we utilize a large language model that takes as
 332 input the projected audio and image tokens along with the
 333 embedded question text. During fine-tuning, the model is
 334 optimized to generate the correct answer based on the given
 335 question and the corresponding multimodal inputs. Training
 336 is performed using the standard language modeling objective
 337 function that maximizes the likelihood of the target
 338 sequence using a cross-entropy loss applied at each token
 339 position.

340 **Warm Start of the Encoders.** To ensure the effectiveness
 341 of each modality-specific representation, the audio and vi-
 342 sual encoders, along with their respective projectors, are
 343 pretrained in a unimodal setting using a large language
 344 model. We utilize the panorama image and binaural audio
 345 from our dataset and construct two types of auxiliary
 346 questions for each modality: classification and localization
 347 tasks. For the visual encoder, the classification task involves
 348 identifying visual objects at specific coordinates, phrased
 349 as “What visual objects did you detect at ($\{\text{azimuth}\}$,
 350 $\{\text{elevation}\}$), $\{\text{distance}\}$ meters?”, and the local-
 351 ization task asks for the predicted azimuth, elevation, and

distance to a specified object class, stated as “What are the
 352 predicted azimuth and elevation angles, and the distance to
 353 the $\{\text{class category}\}$?”. The audio encoder is trained
 354 with analogous tasks: the classification task asks “What
 355 sound did you detect?”, while the localization task prompts
 356 for spatial coordinates of the sound source with the ques-
 357 tion “What are the predicted azimuth and elevation angles,
 358 and the distance to the sound source?”. The visual encoder
 359 adopts a progressive training scheme, first focusing on clas-
 360 sification to learn semantic representations and then incor-
 361 porating spatial grounding through a combined classifica-
 362 tion and localization task. The audio encoder is trained on
 363 both tasks jointly from the beginning. 364

365 5. Experiments

366 5.1. Implementation Details

367 **Image Encoder f_v .** We use a SigLIP2 [47] vision en-
 368 coder with the NaFLEX setting, which supports flexible im-
 369 age resolutions and aspect ratios. The encoder processes a
 370 panoramic image and outputs a sequence of patch tokens.
 371 We apply LoRA [22] to fine-tune the patch embedding and
 372 attention layers of the encoder during both the uni-modal
 373 training and the audio-visual end-to-end training.

374 **Audio Encoder f_a .** We use the pretrained Spatial-AST bin-
 375 aural audio encoder from [52]. The model takes binaural
 376 audio spectrograms as input and generates a sequence of
 377 audio tokens that preserve spatial acoustic cues. The en-
 378 coder was pretrained using the same audio event classifica-
 379 tion and localization tasks proposed in [52]. This encoder
 380 is kept frozen throughout the entire training process.

381 **Modality-specific Projectors and Large Language
 382 Model.** We adopt the Q-Former architecture as the projec-
 383 tor for both modalities. The audio-side projector is based on

Table 2. **Evaluation of baseline models on sound source localization that requires spatial understanding.** R, B, M, Q refer to RGB Image, Binaural Audio, Monaural Audio, and Question (Text) in this table.

| Method | Modality | Q1 (class) | Q1 (aligned) | Q1 (non-matching) | Q7 (class) | Q7 (DoA) | Q8 (class) | Q8 (DoA) |
|-----------------|----------|--------------|--------------|-------------------|--------------|--------------|--------------|--------------|
| Question Only | Q | 3.50 | 3.00 | 2.44 | 2.56 | 7.89 | 0.78 | 7.61 |
| ISSL [39, 40] | R+M | 26.97 | 28.83 | 12.94 | 28.46 | 23.18 | 26.94 | 21.0 |
| ACL-SSL [31] | R+M | 40.56 | 32.83 | 10.61 | 40.41 | 30.68 | 41.11 | 24.33 |
| VideoLLaMA2 [9] | R+M+Q | 51.01 | 77.44 | 50.75 | 70.88 | 68.57 | 75.33 | 46.37 |
| Ours | R+B+Q | 52.69 | 77.61 | 61.67 | 75.44 | 73.21 | 70.27 | 64.27 |

Table 3. **Uni-modal performance of Audio and Vision.** Detection accuracy is denoted as Det., mean angular error as Ang., the proportion of samples with angular error greater than 30° as Ang. > 30, and mean distance error in meters as Dist.

| Modality | Det. | Ang. (°) | Ang. > 30 | Dist. (m) |
|----------|-------|----------|-----------|-----------|
| Audio | 0.575 | 38.01 | 0.289 | 0.476 |
| Vision | 0.633 | 26.89 | 0.161 | 0.332 |

384 the implementation and pretrained weights from BAT [52],
 385 while the visual-side projector is adapted from BLIP-2 [25],
 386 using only the first two attention layers and their corre-
 387 sponding pretrained weights. The number of query tokens
 388 is set to $N_1 = 64$ for audio and $N_2 = 128$ for vision. All
 389 projector parameters are fully trainable. We adopt Qwen2-
 390 7B-Instruct [50] as our LLM backbone.

391 **Training Setup and Input Preprocessing.** Inputs to our
 392 model consist of a single 224×812 panoramic image and
 393 a 10-second audio binaural waveform sampled at 32 kHz.
 394 We preprocess the image input following [47] and the audio
 395 input following [52]. Our full model is trained for 3 epochs
 396 on 8 A5000 GPUs with an effective batch size of 128, using
 397 a LoRA rank of 16 for the image encoder and LLM back-
 398 bone. The training takes three days. Additional training
 399 details are provided in the supplementary material.

400 **Warm-start Performance.** Table 3 shows the uni-modal
 401 performance of the audio and visual encoders after warm-
 402 start pretraining. Both modalities achieve solid individ-
 403 ual results, demonstrating that each encoder learns effective
 404 modality-specific representations. These results serve as a
 405 reference for the subsequent multi-modal experiments.

406 **Baselines.** Since no existing method directly addresses
 407 our proposed task, we introduce three baselines adapted
 408 from related domains. The first two baselines are audio-
 409 visual sound source localization approaches. Specifically,
 410 we adopt the framework proposed in [39, 40], which has
 411 demonstrated strong performance on synthetic benchmarks
 412 and exhibits robustness with multiple visual objects. [31]
 413 learns audio-driven embeddings compatible with the text
 414 encoder of CLIP[34] and leverages the CLIP-based segmen-

415 tion network [28] to achieve tight localization results. Al-
 416 though they do not handle language understanding, we evalu-
 417 ate them using cross-modal retrieval and localization met-
 418 rics. Implementation details are provided in the supplemen-
 419 tary material. The third baseline is the VideoLLaMA2[9],
 420 multi-modal large language model (MLLM), the closest
 421 prior work to ours in terms of multimodal reasoning. For
 422 a fair comparison, we replace its original vision and audi-
 423 o encoders with the same encoders used in our method,
 424 Spatial AST[52] and SigLIP2 NaFLEX [47], and fine-tune
 425 the model on our proposed dataset using the same LLM
 426 backbone. Notably, the baseline uses monaural audio in-
 427 put, whereas our method leverages binaural cues. Since the
 428 sound source localization approaches are not designed for
 429 reasoning tasks (e.g., Q2, Q3, Q4, Q5, Q6, Q9), we evaluate
 430 them only on tasks that do not require language processing.
 431 The metrics in Table 2 cover classification and direction of
 432 arrival (DoA). Q1 (aligned) and Q1 (non-matching) indicate
 433 sound source localization task where the source is semanti-
 434 cally aligned and non-aligned with the audio, respectively.

5.2. Main Results 435

436 We present our results in Table 2, showing that only our
 437 model effectively addresses spatial reasoning scenarios. For
 438 sound classification tasks (Q1, Q7, Q8), sound source local-
 439 ization approaches outperform the Question Only setting,
 440 which serves as a random baseline. VideoLLaMA2 shows
 441 comparable performance to our model, particularly in Q1
 442 (aligned) and Q7 (DoA), where semantic cues are sufficient
 443 for localization due to the presence of a single matching
 444 visual object with audio. Monaural audio is sufficient to
 445 localize the sound source, allowing baseline models to per-
 446 form consistently without spatial audio cues. However, in
 447 Q1 (non-matching) and Q8 (DoA), spatial reasoning is es-
 448 sential for different reasons. In Q1 (non-matching), the vi-
 449 sual object at the sound source is semantically unrelated to
 450 the audio, requiring spatial cues to correctly associate the
 451 sound with the aligned object. In Q8 (DoA), multiple ob-
 452 jects share the same sound category, making it necessary
 453 to differentiate between them using spatial cues. In both
 454 cases, baseline models perform significantly worse. Vide-
 455 oLLaMA2, which shares the same architecture as ours but

Table 4. **Ablation study on modality settings for audio-visual spatial reasoning tasks.** R, B, M, Q refer to RGB Image, Binaural Audio, Monaural Audio, and Question (Text) in this table.

| metric | Trained and tested on | | | | | Trained on R+B+Q, tested on | | Random Chance | Oracle |
|-----------------------------|-----------------------|--------------|-------|-------|--------------|-----------------------------|--------------|---------------|--------|
| | R+B+Q | R+M+Q | B+Q | M+Q | R+Q | R+M+Q | B+Q | Q | O+Q |
| <i>Q1</i> | | | | | | | | | |
| sound accuracy ↑ | 52.69 | 51.01 | 52.53 | 51.40 | 27.28 | 54.03 | 46.86 | 3.50 | 98.08 |
| coming-from accuracy ↑ | 69.64 | 64.10 | 26.40 | 26.40 | 56.22 | 61.92 | 23.39 | 2.72 | 95.61 |
| <i>Q2 (Yes or No) ↑</i> | | | | | | | | | |
| | 84.74 | 83.77 | 55.63 | 50.87 | 85.28 | 83.55 | 54.11 | 50.11 | 94.70 |
| <i>Q3</i> | | | | | | | | | |
| 3-field accuracy ↑ | 69.73 | 66.52 | 32.40 | 18.67 | 74.46 | 66.42 | 24.57 | 18.56 | 85.21 |
| Avg. distance error (m) ↓ | 0.39 | 0.41 | 1.20 | 1.31 | 0.36 | 0.47 | 1.37 | 1.34 | 0.11 |
| <i>Q4</i> | | | | | | | | | |
| DoA accuracy ↑ | 65.68 | 59.03 | 12.86 | 12.43 | 58.06 | 56.14 | 11.38 | 9.80 | 86.63 |
| Avg. DoA error (°) ↓ | 15.41 | 20.21 | 81.18 | 87.38 | 18.59 | 23.55 | 86.49 | 85.48 | 4.17 |
| Avg. distance error (m) ↓ | 0.38 | 0.47 | 1.10 | 1.21 | 0.38 | 0.51 | 1.32 | 1.21 | 0.16 |
| <i>Q2-invisible audio ↑</i> | | | | | | | | | |
| | 72.46 | 70.40 | 57.14 | 48.00 | 73.03 | 70.51 | 52.91 | 50.63 | 94.63 |
| <i>Q3-invisible audio</i> | | | | | | | | | |
| 3-field accuracy ↑ | 59.52 | 47.29 | 34.14 | 18.45 | 41.64 | 45.56 | 25.49 | 18.22 | 84.31 |
| Avg. distance error (m) ↓ | 0.75 | 0.98 | 1.20 | 1.33 | 1.02 | 1.12 | 1.39 | 1.38 | 0.10 |
| <i>Q4-invisible audio</i> | | | | | | | | | |
| DoA accuracy ↑ | 41.18 | 16.71 | 11.18 | 11.76 | 13.53 | 16.47 | 11.29 | 9.88 | 86.24 |
| Avg. DoA error (°) ↓ | 39.81 | 69.25 | 80.51 | 84.56 | 77.15 | 75.39 | 85.24 | 84.81 | 4.27 |
| Avg. distance error (m) ↓ | 0.71 | 1.08 | 1.13 | 1.21 | 1.16 | 1.04 | 1.32 | 1.23 | 0.15 |
| <i>Q5</i> | | | | | | | | | |
| class accuracy ↑ | 72.43 | 74.26 | 25.79 | 25.63 | 74.87 | 72.82 | 22.18 | 2.78 | 97.50 |
| sounding accuracy ↑ | 75.60 | 64.54 | 59.48 | 37.72 | 36.63 | 65.93 | 75.93 | 41.36 | 100 |
| <i>Q6</i> | | | | | | | | | |
| class accuracy ↑ | 81.06 | 81.61 | 51.78 | 50.47 | 83.78 | 80.72 | 42.72 | 3.72 | 95.78 |
| sounding accuracy ↑ | 72.33 | 52.33 | 59.33 | 38.67 | 31.94 | 49.28 | 75.67 | 41.67 | 100 |
| <i>Q7</i> | | | | | | | | | |
| class accuracy ↑ | 75.44 | 70.88 | 51.64 | 53.62 | 37.35 | 73.53 | 51.68 | 2.56 | 93.83 |
| DoA accuracy ↑ | 73.21 | 68.57 | 47.30 | 7.80 | 37.52 | 64.04 | 48.38 | 7.89 | 92.05 |
| Avg. DoA error (°) ↓ | 14.75 | 22.41 | 33.02 | 88.31 | 56.66 | 24.55 | 35.25 | 90.92 | 4.05 |
| Avg. distance error (m) ↓ | 0.30 | 0.33 | 0.50 | 0.53 | 0.44 | 0.36 | 0.79 | 0.53 | 0.11 |
| <i>Q8</i> | | | | | | | | | |
| class accuracy ↑ | 70.27 | 75.33 | 48.42 | 48.02 | 69.89 | 71.90 | 32.51 | 0.78 | 95.15 |
| DoA accuracy ↑ | 64.27 | 46.37 | 47.69 | 8.46 | 43.72 | 39.76 | 49.41 | 7.61 | 90.67 |
| Avg. DoA error (°) ↓ | 23.78 | 50.80 | 32.32 | 89.93 | 51.90 | 52.46 | 32.45 | 89.40 | 3.86 |
| Avg. distance error (m) ↓ | 0.36 | 0.44 | 0.48 | 0.51 | 0.42 | 0.46 | 0.85 | 0.52 | 0.12 |
| <i>Q9</i> | | | | | | | | | |
| sound accuracy ↑ | 54.00 | 51.14 | 51.14 | 52.20 | 27.17 | 55.57 | 47.25 | 2.81 | 98.03 |
| visibility accuracy ↑ | 75.22 | 72.94 | 38.99 | 39.79 | 33.31 | 76.35 | 49.42 | 42.31 | 100 |

456 lacks binaural audio, achieves approximately 50% accuracy
457 in Q8 (DoA), indicating its inability to distinguish between
458 visually similar objects that semantically match the audio.
459 Since all baseline models use only monaural audio, they
460 lack spatial information, making spatial reasoning impos-
461 sible.

5.3. Ablation Studies

Table 4 shows that both image (R: RGB) and binaural au-
463 dio (B) inputs are crucial for spatial reasoning. It compares
464 R+B+Q, R+M+Q (M: monaural), B+Q, M+Q, and R+Q (Q:
465 question), highlighting that binaural audio provides spatial
466 cues while monaural lacks directional information. The fol-
467

| | | |
|-----|--|-----|
| 468 | lowing is an analysis of the performance for each question | 521 |
| 469 | type. Oracle performance assumes ideal audio and visual | 522 |
| 470 | encoders using metadata. | 523 |
| 471 | Question 1 involves sound and visual object classification, | 524 |
| 472 | with half of the samples containing a non-matching visual | 525 |
| 473 | object at the sound source. Both R+B+Q and R+M+Q show | 526 |
| 474 | similar sound classification accuracy (52.69% and 51.01%), | 527 |
| 475 | suggesting comparable semantic cues from monaural and | 528 |
| 476 | binaural audio. However, in coming-from accuracy, | 529 |
| 477 | R+B+Q (69.64%) outperforms R+M+Q (64.10%), high- | 530 |
| 478 | lighting the spatial advantage of binaural audio. | 531 |
| 479 | Questions 2, 3, and 4 assess distance and relative location | 532 |
| 480 | between the sound source and visual objects, requiring spa- | 533 |
| 481 | tial reasoning across modalities. For visible audio, R+M+Q | 534 |
| 482 | achieves 66.52% in Q3 and 59.03% in Q4, performing sim- | 535 |
| 483 | ilarly to R+B+Q (69.73% and 65.68%). When the sound | 536 |
| 484 | source is invisible, R+B+Q shows a clear advantage, out- | 537 |
| 485 | performing R+M+Q in Q3 (59.52% vs. 47.29%) and Q4 | 538 |
| 486 | (41.18% vs. 16.71%). This highlights the role of binaural | 539 |
| 487 | audio in capturing spatial cues that monaural audio with vi- | 540 |
| 488 | sual input cannot provide. | 541 |
| 489 | Questions 5 and 6 both involve identifying the sound- | 542 |
| 490 | producing object but differ in complexity based on the num- | 543 |
| 491 | ber of visual objects that match the sound. In Q5, with only | 544 |
| 492 | one matching object, visual context alone provides suffi- | 545 |
| 493 | cient spatial information for localization. R+M+Q lever- | 546 |
| 494 | ages visual cues effectively, achieving a sounding accuracy | 547 |
| 495 | of 64.54%. With no visual ambiguity, the model can reli- | 548 |
| 496 | ably associate the sound with the correct object using spa- | 549 |
| 497 | tial information from the visual signal. In Q6, two visu- | 550 |
| 498 | ally similar objects match the sound, introducing ambiguity. | 551 |
| 499 | R+M+Q's performance drops to 52.33%, as visual context | 552 |
| 500 | alone is no longer sufficient to distinguish between the two | 553 |
| 501 | objects, leading to random guessing. In contrast, B+Q and | 554 |
| 502 | R+B+Q maintain consistent performance across both ques- | 555 |
| 503 | tions. In Q5, they achieve 59.48% and 75.60%, respectively, | 556 |
| 504 | and in Q6, their performance remains stable at 59.33% and | 557 |
| 505 | 72.33%. This stability is due to binaural audio, which pro- | 558 |
| 506 | vides explicit spatial cues, enabling the model to localize | 559 |
| 507 | the sound source based solely on directional information, | |
| 508 | unaffected by visual similarity. These results indicate that | |
| 509 | when there is only one matching object (Q5), R+M+Q can | |
| 510 | effectively use visual spatial information. However, when | |
| 511 | multiple visually similar objects are present (Q6), spatial | |
| 512 | audio cues become essential, allowing B+Q and R+B+Q to | |
| 513 | maintain stable performance regardless of visual similarity. | |
| 514 | These results highlight the importance of binaural audio in | |
| 515 | resolving ambiguity in complex visual scenes. | |
| 516 | Questions 7 and 8 both involve sound classification and | |
| 517 | localization but differ in the number of visual objects that | |
| 518 | correspond to the audio, with two in Q8 and one in Q7. | |
| 519 | In Q8, two visually similar objects correspond to the au- | |
| 520 | dio, making it difficult for the model to distinguish between | |
| | them using visual information alone. R+M+Q and B+Q | 521 |
| | show similar DoA accuracy (46.37% and 47.69%), but their | 522 |
| | Avg. DoA errors differ, with R+M+Q at 50.80° and B+Q at | 523 |
| | 32.32°. R+M+Q relies on visual context for spatial cues, but | 524 |
| | semantic ambiguity between the two objects complicates | 525 |
| | localization, leading to random selection and higher error. | 526 |
| | In contrast, B+Q, using binaural audio, focuses solely on | 527 |
| | directional information, perceiving only one sound source | 528 |
| | without considering object-level ambiguity, resulting in a | 529 |
| | lower error. R+B+Q achieves the lowest error (23.78°) by | 530 |
| | combining spatial audio and visual inputs. In Q7, the audio | 531 |
| | corresponds to a single object, eliminating semantic ambi- | 532 |
| | guity. In this case, the performance of R+M+Q and B+Q | 533 |
| | reverses from Q8. R+M+Q records a lower error (22.41°) | 534 |
| | than B+Q (33.02°), indicating that when only one object | 535 |
| | is present, visual spatial information can effectively guide | 536 |
| | localization without semantic confusion. These results sup- | 537 |
| | port the findings in Q5 and Q6, emphasizing the role of spa- | 538 |
| | tial audio in disambiguating visually similar objects. | 539 |
| | Question 9 involves sound classification and localization | 540 |
| | while also requiring the model to determine whether the | 541 |
| | object is visually present at the sound source. This task de- | 542 |
| | mands both audio and visual semantic understanding. Both | 543 |
| | multi-modal settings (R+B+Q, R+M+Q) successfully ad- | 544 |
| | dress this question. | 545 |
| | Modality Setting Cross-Evaluation. To assess the im- | 546 |
| | pact of vision signals and binaural audio during training, | 547 |
| | we evaluate the model trained on R+B+Q under R+M+Q | 548 |
| | and B+Q settings. While Q7 and Q8 show minimal change, | 549 |
| | Q5 and Q6 exhibit noticeable gaps in sounding accuracy. | 550 |
| | This might come from Q5 and Q6 only requiring yes/no | 551 |
| | responses given a location, without the detailed localiza- | 552 |
| | tion required in Q7 and Q8. Consequently, the model in | 553 |
| | the B+Q setting may not effectively leverage spatial reason- | 554 |
| | ing for these tasks. However, with visual signals, the model | 555 |
| | gains implicit spatial cues that align audio locations with | 556 |
| | the visual scene, potentially enhancing spatial audio under- | 557 |
| | standing. Thus, the presence of visual information may be | 558 |
| | beneficial even for learning spatial audio cues. | 559 |
| | 6. Conclusion | 560 |
| | We introduce a new task, audio-visual spatial reasoning, | 561 |
| | along with the <i>Hear You Are LLM</i> and QA dataset. Un- | 562 |
| | like prior work that focuses on semantic or temporal align- | 563 |
| | ment, our approach emphasizes spatial reasoning by inte- | 564 |
| | grating binaural audio and visual inputs. We build a large- | 565 |
| | scale dataset covering diverse spatial scenarios and propose | 566 |
| | a multimodal framework combining spatial encoders with | 567 |
| | a large language model. Experiments show that monaural | 568 |
| | audio with vision or unimodal binaural methods lack the | 569 |
| | capacity for spatial reasoning. These results underscore the | 570 |
| | importance of spatial reasoning in robust multimodal un- | 571 |
| | derstanding and set a new benchmark in audio-visual learn- | 572 |
| | ing. | 573 |

574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629

References

- [1] Alan Baade, Puyuan Peng, and David Harwath. Mae-ast: Masked autoencoding audio spectrogram transformer. In *INTERSPEECH*, 2022. 3
- [2] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3D: Learning from RGB-D data in indoor environments. In *3DV*, 2017. 3
- [3] Moitreyia Chatterjee, Jonathan Le Roux, Narendra Ahuja, and Anoop Cherian. Visual scene graphs for audio source separation. In *ICCV*, 2021. 1
- [4] Changan Chen, Wei Sun, David Harwath, and Kristen Grauman. Learning audio-visual dereverberation. *arXiv preprint arXiv:2106.07732*, 2021. 4
- [5] Changan Chen, Carl Schissler, Sanchit Garg, Philip Kobernik, Alexander Clegg, Paul Calamia, Dhruv Batra, Philip Robinson, and Kristen Grauman. Soundspaces 2.0: A simulation platform for visual-acoustic learning. *NeurIPS*, 2022. 2, 3
- [6] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 2, 3
- [7] H Chen, W Xie, T Afouras, A Nagrani, A Vedaldi, and A Zisserman. Audio-visual synchronisation in the wild. In *BMVC*, 2021. 1
- [8] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *CVPR*, 2021. 1
- [9] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 6
- [10] Sanjoy Chowdhury, Sayan Nag, Subhrajyoti Dasgupta, Jun Chen, Mohamed Elhoseiny, Ruohan Gao, and Dinesh Manocha. Meerkat: Audio-visual large language model for grounding in space and time. In *ECCV*, 2024. 3
- [11] Sanjoy Chowdhury, Sayan Nag, Subhrajyoti Dasgupta, Yaoting Wang, Mohamed Elhoseiny, Ruohan Gao, and Dinesh Manocha. Avtrustbench: Assessing and enhancing reliability and robustness in audio-visual llms. *arXiv preprint arXiv:2501.02135*, 2025. 3
- [12] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *ACCV*, 2017. 1
- [13] Bhavika Devnani, Skyler Seto, Zakaria Aldeneh, Alessandro Toso, Elena Menyaylenko, Barry-John Theobald, Jonathan Sheaffer, and Miguel Sarabia. Learning spatially-aware language and audio embeddings. In *NeurIPS*, 2024. 2, 3
- [14] Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. Music gesture for visual sound separation. In *CVPR*, 2020. 1
- [15] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019.
- [16] Ruohan Gao and Kristen Grauman. Visualvoice: Audio-visual speech separation with cross-modal consistency. In *CVPR*, 2021. 1
- [17] Kaixiong Gong, Kaituo Feng, Bohao Li, Yibing Wang, Mofan Cheng, Shijia Yang, Jiaming Han, Benyou Wang, Yutong Bai, Zhuoran Yang, and Xiangyu Yue. Av-odyssey bench: Can your multimodal llms really understand audio-visual information? *arXiv preprint arXiv:2412.02611*, 2024. 3
- [18] Yuan Gong, Yu-An Chung, and James Glass. Ast: Audio spectrogram transformer. In *INTERSPEECH*, 2021. 3
- [19] Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid Karlinsky, and James R Glass. Listen, think, and understand. In *ICLR*, 2024. 3
- [20] Yuhang He, Sangyun Shin, Anoop Cherian, Niki Trigoni, and Andrew Markham. Soundloc3d: Invisible 3d sound source localization and classification using a multimodal rgb-d acoustic camera. In *WACV*, 2025. 2, 3
- [21] Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou. Discriminative sounding objects localization via self-supervised audiovisual matching. *NeurIPS*, 2020. 1
- [22] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 5
- [23] Xixi Hu, Ziyang Chen, and Andrew Owens. Mix and localize: Localizing sound sources in mixtures. In *CVPR*, 2022. 1
- [24] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *NeurIPS*, 2022. 3
- [25] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Bliip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*. PMLR, 2023. 6
- [26] Zhaowei Li, Qi Xu, Dong Zhang, Hang Song, YiQing Cai, Qi Qi, Ran Zhou, Junting Pan, Zefeng Li, Vu Tu, et al. Groundingpt: Language enhanced multi-modal grounding model. In *ACL*, 2024. 3
- [27] Jinxiang Liu, Chen Ju, Weidi Xie, and Ya Zhang. Exploiting transformation invariance and equivariance for self-supervised sound localisation. In *ACM MM*, 2022. 1
- [28] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *CVPR*, 2022. 6
- [29] Shentong Mo and Pedro Morgado. A closer look at weakly-supervised audio-visual source localization. In *NeurIPS*, 2022. 1
- [30] Shentong Mo and Yapeng Tian. Audio-visual grouping network for sound localization from mixtures. In *CVPR*, 2023. 1
- [31] Sooyoung Park, Arda Senocak, and Joon Son Chung. Can clip help sound source localization? In *WACV*, 2024. 6
- [32] Se Jin Park, Minsu Kim, Joanna Hong, Jeongsoo Choi, and Yong Man Ro. Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory. In *AAAI*, 2022. 1
- [33] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, 2020. 1

| | | |
|-----|--|------------|
| 687 | [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In <i>ICML</i> , 2021. | |
| 688 | | |
| 689 | | |
| 690 | | |
| 691 | | |
| 692 | 6 | |
| 693 | [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In <i>CVPR</i> , 2022. | |
| 694 | | 2 |
| 695 | [36] S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. Mmau: A massive multi-task audio understanding and reasoning benchmark. In <i>ICLR</i> , 2024. | |
| 696 | | 3 |
| 697 | [37] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In <i>CVPR</i> , 2018. | |
| 698 | | 2 |
| 699 | [38] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound sources in visual scenes: Analysis and applications. <i>IEEE TPAMI</i> , 2021. | |
| 700 | | 2 |
| 701 | [39] Arda Senocak, Hyeonggon Ryu, Junsik Kim, Tae-Hyun Oh, Hanspeter Pfister, and Joon Son Chung. Sound source localization is all about cross-modal alignment. In <i>ICCV</i> , 2023. | |
| 702 | | 2, 6 |
| 703 | [40] Arda Senocak, Hyeonggon Ryu, Junsik Kim, Tae-Hyun Oh, Hanspeter Pfister, and Joon Son Chung. Aligning sight and sound: Advanced sound source localization through audio-visual alignment. <i>arXiv preprint arXiv:2407.13676</i> , 2024. | |
| 704 | | 2, 6 |
| 705 | [41] Kazuki Shimada, Archontis Politis, Parthasaarathy Sudarsanam, Daniel Krause, Kengo Uchida, Sharath Adavanne, Aapo Hakala, Yuichiro Koyama, Naoya Takahashi, Shusuke Takahashi, Tuomas Virtanen, and Yuki Mitsufuji. Starss23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events. In <i>NeurIPS</i> , 2023. | |
| 706 | | 2, 3 |
| 707 | [42] Zengjie Song, Yuxi Wang, Junsong Fan, Tieniu Tan, and Zhaoxiang Zhang. Self-supervised predictive learning: A negative-free method for sound source localization in visual scenes. In <i>CVPR</i> , 2022. | |
| 708 | | 2 |
| 709 | [43] Weixuan Sun, Jiayi Zhang, Jianyuan Wang, Zheyuan Liu, Yiran Zhong, Tianpeng Feng, Yandong Guo, Yanhao Zhang, and Nick Barnes. Learning audio-visual source localization via false negative aware contrastive learning. In <i>CVPR</i> , 2023. | |
| 710 | | 2 |
| 711 | [44] Kim Sung-Bin, Oh Hyun-Bin, JungMok Lee, Arda Senocak, Joon Son Chung, and Tae-Hyun Oh. Avhbench: A cross-modal hallucination benchmark for audio-visual large language models. <i>arXiv preprint arXiv:2410.18325</i> , 2024. | |
| 712 | | 3 |
| 713 | [45] Reuben Tan, Arijit Ray, Andrea Burns, Bryan A Plummer, Justin Salamon, Oriol Nieto, Bryan Russell, and Kate Saenko. Language-guided audio-visual source separation via trimodal consistency. In <i>CVPR</i> , 2023. | |
| 714 | | 1 |
| 715 | [46] Yapeng Tian, Di Hu, and Chenliang Xu. Cyclic co-learning of sounding object visual grounding and sound separation. In <i>CVPR</i> , 2021. | |
| 716 | | 1 |
| 717 | [47] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. <i>arXiv preprint arXiv:2502.14786</i> , 2025. | |
| 718 | | 5, 6 |
| 719 | [48] Efthymios Tzinis, Scott Wisdom, Aren Jansen, Shawn Hershey, Tal Remez, Daniel PW Ellis, and John R Hershey. Into the wild with audioscope: Unsupervised audio-visual separation of on-screen sounds. In <i>ICLR</i> , 2020. | |
| 720 | | 1 |
| 721 | [49] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. <i>arXiv preprint arXiv:2404.07191</i> , 2024. | |
| 722 | | 2, 4 |
| 723 | [50] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report, 2024. <i>arXiv preprint arXiv:2407.10671</i> , 2024. | |
| 724 | | 6 |
| 725 | [51] Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: Enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. In <i>ECCV</i> , 2024. | |
| 726 | | 3 |
| 727 | [52] Zhisheng Zheng, Puyuan Peng, Ziyang Ma, Xie Chen, Eun-sol Choi, and David Harwath. Bat: Learning to reason about spatial sounds with large language models. In <i>ICML</i> , 2024. | |
| 728 | | 2, 3, 5, 6 |
| 729 | | 744 |
| 730 | | 745 |
| 731 | | 746 |
| 732 | | 747 |
| 733 | | 748 |
| 734 | | 749 |
| 735 | | 750 |
| 736 | | 751 |
| 737 | | 752 |
| 738 | | 753 |
| 739 | | 754 |
| 740 | | 755 |
| 741 | | 756 |
| 742 | | 757 |
| 743 | | 758 |
| | | 759 |
| | | 760 |
| | | 761 |
| | | 762 |
| | | 763 |
| | | 764 |
| | | 765 |
| | | 766 |
| | | 767 |
| | | 768 |
| | | 769 |
| | | 770 |
| | | 771 |