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

Anonymous Author(s)

Affiliation Address email

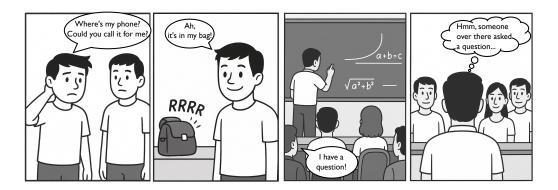


Figure 1: **Audio-Visual Spatial Reasoning.** (Left) A phone rings out of sight inside a bag; although the sound's semantic cue ("ring tone") is present, spatial reasoning is required to locate the true source among visually silent objects. (Right) In a classroom, several students share the same semantic 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

Many audio-visual learning methods have focused on aligning audio and visual information, either through semantic or temporal correspondence. However, most of these works have utilized monaural audio, which does not contain information about the spatial location of the sound source. In contrast, humans and other animals utilize binaural hearing to perceive this spatial information. Combining spatial sound and visual perception enables powerful high-level reasoning: for example, a person looking for their phone may hear the ringing sound coming from a backpack sitting on a table, and quickly infer that the missing phone is inside the backpack. In this paper, we investigate the problem of Audio-Visual **Spatial Reasoning.** We design a spatial audio-visual question answering dataset to cover scenarios where semantic correspondence between audio and visual signals is absent but spatial alignment exists, as well as cases with multiple audio-visual semantic correspondences that require spatial reasoning to disambiguate. We propose a model that learns spatial comprehension across the audio and vision modalities by connecting them with a large language model and experimentally demonstrate that spatial sound perception is an essential part of our task.

16 17

3

5

6

8

9

10

11

12

13

14

15

8 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) [7, 21, 33, 27, 30, 23, 29], source separation [3, 14, 16, 48, 46, 15, 45], and audio-visual synchronization [12, 6, 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 scenarios where spatial reasoning is necessary. For instance, understanding that a ringtone sound is emanating from a backpack requires spatial reasoning, as the backpack does not semantically match the sound. Another example is when a single sound (e.g., speech) could correspond to multiple visual objects (e.g., several students in a classroom), where spatial cues help pinpoint the actual source. These examples highlight the limitations of previous methods, emphasizing the need to address spatial reasoning beyond basic perception.

Previous studies in spatial audio reasoning have primarily focused on audio-only approaches, excluding visual information while incorporating language as a modality for spatial interpretation. [13] aligns audio and text embeddings for spatial tasks, while [52] leverages large language models for spatial audio question answering. While spatial audio itself provides rich information for spatial reasoning, integrating visual information into these tasks is a natural extension, as visual signals inherently convey spatial context. This combination not only enhances spatial perception and localization capabilities, but also enables more sophisticated spatial reasoning, such as handling scenarios involving sounding sources and nearby visual objects.

In this paper, we address the problem of **Audio-Visual Spatial Reasoning**, which involves understanding the spatial relationship between a sound and the visual context. This task goes beyond simply perceiving and localizing a sound source, as it requires reasoning about spatial cues to infer relationships and interactions between objects. To support research on this problem, we construct a large-scale dataset of 1 million question-answer pairs, specifically designed to serve as both the training and evaluation set for spatial audio-visual reasoning in diverse scenarios. The vision and spatial audio is rendered using SoundSpaces 2.0 [4], with source audio clips sampled from VGGSound[8]. 3D objects associated with these sounds are generated using Stable Diffusion 3[35] and InstantMesh [49], and then are placed within the virtual environments. This dataset serves as a comprehensive benchmark for spatially intricate settings, providing questions that assess spatial alignment between modalities, relative locations between sounding and non-sounding objects, and localization of sound sources among multiple visual objects of the same category as the query audio.

Furthermore, we propose a multi-modal framework, Hear You Are LLM, which leverages spatial audio and visual encoders to integrate spatial information. The model is trained to handle all the spatial reasoning tasks from our dataset, enabling it to address scenarios where semantic alignment alone is insufficient. We experimentally demonstrate that our proposed method effectively addresses the audio-visual spatial reasoning problem, outperforming existing baseline models including a state-of-the-art monaural sound source localization method [39, 40] and a large language model-based audio-visual model that lacks spatial understanding. These results highlight the importance of incorporating spatial audio-visual knowledge to achieve robust multi-modal reasoning. To summarize, our main contributions are as follows:

- We define a new task, audio-visual spatial reasoning, focusing on understanding spatial relationships between sound and visual context, going beyond basic semantic perception such as sound source localization (object detection based on audio queries) and audio-visual segmentation.
- We propose *Hear You Are LLM*, a multi-modal modeling framework that integrates spatial audio and visual encoders with a large language model to handle complex spatial reasoning tasks.

• We construct *Hear You Are QA*, the first large-scale dataset specifically designed for audio-visual spatial reasoning, consisting of 1 million question-answer pairs across diverse spatial scenarios for training and evaluation. We will open source both the dataset and the training code.

76 2 Related Work

77

2.1 Audio-Visual Sound Source Localization

Audio-visual sound source localization is the task of detecting the object or area that corresponds to the query audio in the visual scene. Following the development of deep learning, Senocak et al. [37, 38] suggested a semantic alignment-based approach by proposing a cross-modal attention mechanism with contrastive learning. The field has advanced in the direction of better cross-modal alignment by leveraging negative-free self-supervised learning [42], intra-modality similarity learning [43], and the use of multiple positive learning [39], aligning with representation learning methods. However, these methods rely on monaural audio and are limited to audio-visual semantic correspondence without spatial understanding.

Different approaches have focused more on spatial audio for sound source localization. Anoopcherian 86 et al. [20] proposed a 3D sound source localization method trained on a dataset with four-channel 87 audio and multi-view visual scenes synthesized using SoundSpaces 2.0. Their approach localizes 88 sound within the visual scene, but the visual counterpart of the sound is not visible in their setting, as they only localize the area of the sound source. Shimada et al. [41] constructed an audio-visual 91 sound source localization and detection dataset in which audio-visual alignment is guaranteed. In their framework, the visual signal serves as an auxiliary modality to improve sound localization and detection. In contrast, we present an audio-visual scene that includes both sound-producing and silent 93 objects, allowing the model to learn a broader range of spatial reasoning tasks that require contextual 94 understanding beyond basic localization. 95

96 2.2 Spatial Audio Reasoning

Following recent advancements in audio understanding [18, 1, 24] and reasoning [19, 36], several 97 approaches have been proposed to address spatial audio reasoning. [52] synthesize the spatial sound 98 question answering dataset with the SoundSpaces 2.0 simulator and train a spatial audio encoder and 99 a large language model for spatial audio understanding and reasoning. This framework handles tasks 100 such as sound event detection, direction and distance estimation, and spatial reasoning, for example, 101 "What is the sound on the left side of the sound of the dog barking?" Another line of research explores 102 spatial audio reasoning through contrastive language-audio pretraining, with synthetic first-order 103 ambisonics [13]. However, these approaches do not incorporate the vision modality, which opens 104 another dimension for reasoning. 105

2.3 Audio-Visual LLMs

106

120

Inspired by the advancements of Large Language Models (LLMs), recent studies have extended these 107 models to Multimodal Large Language Models (MLLMs) to tackle a wider range of multimodal tasks. In the audio-visual domain, Grounding GPT [26] introduces multimodal grounding for audio, 109 image, and video data using LLMs. Meerkat [10] aligns audio-visual features using optimal transport 110 and attention consistency, and CAT [51] aggregates question-related clues in audio-visual scenarios. From a benchmarking standpoint, AVHBench, AVTRUSTBENCH, and AV-Odyssey Bench [11, 112 44, 17] provide comprehensive benchmarks targeting hallucination detection [44], reliability and robustness [11], and both foundational capabilities and high-level reasoning [17]. While recent studies have advanced multimodal learning, they primarily rely on monaural audio, limiting their ability to handle spatial reasoning. As spatial reasoning enables a broader range of tasks and more closely reflects real-world scenarios, it must be addressed to achieve comprehensive audio-visual 117 understanding. We propose a new dataset and model specifically designed for spatial reasoning in 118 audio-visual tasks. 119

3 Creation of Hear You Are QA Dataset

Our goal is to train a model to learn both semantic and spatial reasoning, for audio-visual inputs. To this end, we introduce the Hear You Are QA Dataset. Constructing large-scale audio-visual scene data



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.

Table 1: Spatial audio visual question types and base prompts/answers

Q1. Spatial Correspondence

Q: What is the sound class category? Where is the sound coming from? A: phone ringing; cupboard

O2-4. Relative Location

(**Distance**) **Q:** Is the sound source of the siren closer to the agent than it is to the cat?

(Direction) 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

(Angle) 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)

O: 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 Q: What is the sound class category? Is the sound source visible in the scene? A: fox; visible

with real-world spatial audio is time-consuming and challenging, requiring specialized equipment 123

such as ambisonic or dummy head microphones. To efficiently build a diverse dataset with various 124

objects and sound events, we adopt a simulation-based approach to generate both the scenes and 125

spatial audio. 126

Spatial Audio Simulator. We employ the SoundSpaces 2.0 simulator [4], which renders geometry-127

based acoustics, adding realistic reverberation for any source–receiver pair. Users can freely vary 128

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. 130

Scene meshes come from Matterport3D [2], a collection of 90 fully scanned buildings averaging 131

24.5 rooms across 2.61 floors and 517.34 m² of floor space. We use 72 scenes for training, 9 for 132

validation and 9 for testing. Given a source location, monaural signal, receiver position, and heading, 133

the observed signal is obtained by convolving the monaural signal with the environment's room 134

135 impulse response. We configure the receiver to record a binaural audio signal with the default Head

Related Transfer Function (HRTF) provided by SoundSpaces 2.0.

Sound Sources. Previous spatial audio datasets include either a limited number of class cate-137

gories [41] or classes that are not guaranteed to be visually observable [52, 20]. To construct a 138

large-scale audio-visual dataset, we adopt VGGSound [8], which contains 200,000 in-the-wild 10-139

second YouTube clips, each annotated with one of 309 audio event classes. However, some of these 140

classes correspond to events that typically occur outdoors or are difficult to associate with a single

141 visual object (e.g., "Airplane Flyby", "People Marching"). To enhance the visual reliability and 142

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 146 datasets, we generate our own 3D objects to be placed within the Matterport3D environments, either 147 as sounding objects or as distractor objects. Specifically, we first select 150 class categories from 148 VGGSound and 40 from ImageNet, and generate 2D images for each category using Stable Diffusion 149 3. After manually filtering out low-quality or unrealistic generations, we select 40 visually plausible 150 images per category. These 2D images are then lifted into 3D object meshes using the method 151 from [49]. For each sounding object category, we reserve 32 images for training, 4 for validation and 152 4 for testing. 153

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 [5], 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-40 to paraphrase and expand each base question into multiple human-like variations.

4 Method

156

157

168

171

172

173

174

175

176

177

178

179

180

181

182

183

186

187

189

191

192

193

195

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, $\mathbf{v} \in \mathbb{R}^{N_v \times C_v}$, where N_v is the number of visual tokens and C_v is the feature dimension of each token. We preserve the full spatial layout of patch tokens without pooling. The audio encoder f_a takes the input spectrogram of a and produces a set of audio tokens, $\mathbf{a} \in \mathbb{R}^{N_a \times C_a}$, where N_a is the number of audio tokens and C_a is the corresponding feature dimension. Each modality-specific encoder is followed by a

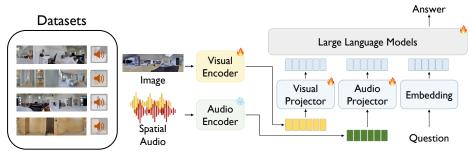


Figure 3: **The pipeline of our framework:** feature extraction, projection, and multimodal reasoning. We extract spatial audio and visual features using pre-trained encoders, project them into a shared embedding space, and integrate the embeddings with the question embedding to generate the answer.

projector that maps the extracted features into the hidden dimension of the language model. The visual projector attends to the spatial visual features to generate N_V projected tokens, and the audio projector similarly produces N_A tokens from the audio features. These projected tokens are then passed to the large language model for multi-modal reasoning.

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

Warm Start of the Encoders. To ensure the effectiveness of each modality-specific representation, the audio and visual encoders, along with their respective projectors, are pretrained in a unimodal setting using a large language model. We utilize the panorama image and binaural audio from our dataset and construct two types of auxiliary questions for each modality: classification and localization tasks. For the visual encoder, the classification task involves identifying visual objects at specific coordinates, phrased as "What visual objects did you detect at ({azimuth}, {elevation}), {distance} meters?", and the localization task asks for the predicted azimuth, elevation, and distance to a specified object class, stated as "What are the predicted azimuth and elevation angles, and the distance to the {class category}?". The audio encoder is trained with analogous tasks: the classification task asks "What sound did you detect?", while the localization task prompts for spatial coordinates of the sound source with the question "What are the predicted azimuth and elevation angles, and the distance to the sound source?". The visual encoder adopts a progressive training scheme, first focusing on classification to learn semantic representations and then incorporating spatial grounding through a combined classification and localization task. The audio encoder is trained on both tasks jointly from the beginning.

5 Experiments

5.1 Implementation Details

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

Audio Encoder f_a . We use the pretrained Spatial-AST binaural audio encoder from [52]. The model takes binaural audio spectrograms as input and generates a sequence of audio tokens that preserve spatial acoustic cues. The encoder was pretrained using the same audio event classification and localization tasks proposed in [52]. This encoder is kept frozen throughout the entire training process.

Modality-specific Projectors and Large Language Model We adopt the Q-Former architecture as the projector for both modalities. The audio-side projector is based on the implementation and pretrained weights from BAT [52], while the visual-side projector is adapted from BLIP-2 [25], using only the first two attention layers and their corresponding pretrained weights. The number of query

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 70.27	46.37
Ours	R+B+Q	52.69	77.61	61.67	75.44	73.21		64.27

tokens is set to $N_1 = 64$ for audio and $N_2 = 128$ for vision. All projector parameters are fully trainable. We adopt Qwen2-7B-Instruct [50] as our LLM backbone.

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

Baselines. Since no existing method directly addresses our proposed task, we introduce three baselines adapted from related domains. The first two baselines are audio-visual sound source localization approaches. Specifically, we adopt the framework proposed in [39, 40], which has demonstrated strong performance on synthetic benchmarks and exhibits robustness with multiple visual objects. [31] learns audio-driven embeddings compatible with the text encoder of CLIP[34] and leverages the CLIP-based segmentation network [28] to achieve tight localization results. Although they do not handle language understanding, we evaluate them using cross-modal retrieval and localization metrics. Implementation details are provided in the supplementary material. The third baseline is the VideoLLaMA2[9], multi-modal large language model (MLLM), the closest prior work to ours in terms of multimodal reasoning. For a fair comparison, we replace its original vision and audio encoders with the same encoders used in our method, Spatial AST[52] and SigLIP2 NaFLEX [47], and fine-tune the model on our proposed dataset using the same LLM backbone. Notably, the baseline uses monaural audio input, whereas our method leverages binaural cues. Since the sound source localization approaches are not designed for reasoning tasks (e.g., Q2, Q3, Q4, Q5, Q6, Q9), we evaluate them only on tasks that do not require language processing. The metrics in Table 2 cover classification and direction of arrival (DoA). Q1 (aligned) and Q1 (non-matching) indicate sound source localization task where the source is semantically aligned and non-aligned with the audio, respectively.

5.2 Main Results

We present our results in Table 2, showing that only our model effectively addresses spatial reasoning scenarios. For sound classification tasks (Q1, Q7, Q8), sound source localization approaches outperform the Question Only setting, which serves as a random baseline. VideoLLaMA2 shows comparable performance to our model, particularly in Q1 (aligned) and Q7 (DoA), where semantic cues are sufficient for localization due to the presence of a single matching visual object with audio. Monaural audio is sufficient to localize the sound source, allowing baseline models to perform consistently without spatial audio cues. However, in Q1 (non-matching) and Q8 (DoA), spatial reasoning is essential for different reasons. In Q1 (non-matching), the visual object at the sound source is semantically unrelated to the audio, requiring spatial cues to correctly associate the sound with the aligned object. In Q8 (DoA), multiple objects share the same sound category, making it necessary to differentiate between them using spatial cues. In both cases, baseline models perform significantly worse. VideoLLaMA2, which shares the same architecture as ours but lacks binaural audio, achieves approximately 50% accuracy in Q8 (DoA), indicating its inability to distinguish between visually similar objects that semantically match the audio. Since all baseline models use only monaural audio, they lack spatial information, making spatial reasoning impossible.

Table 3: **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.

		Trained	and tested o	Trained on R+B+Q, tested on		Random Chance		
metric	R+B+Q	R+M+Q	B+Q	B+Q M+Q		R+M+Q	B+Q	Q
Q1								
sound accuracy ↑	52.69	51.01	52.53	51.40	27.28	54.03	46.86	3.50
coming-from accuracy ↑	69.64	64.10	26.40	26.40	56.22	61.92	23.39	2.72
Q2 (Yes or No) ↑	84.74	83.77	55.63	50.87	85.28	83.55	54.11	50.11
Q3								
3-field accuracy ↑	69.73	66.52	32.40	18.67	74.46	66.42	24.57	18.56
Avg. distance error (m) ↓	0.39	0.41	1.20	1.31	0.36	0.47	1.37	1.34
Q4								
DoA accuracy ↑	65.68	59.03	12.86	12.43	58.06	56.14	11.38	9.80
Avg. DoA error (°) ↓	15.41	20.21	81.18	87.38	18.59	23.55	86.49	85.48
Avg. distance error (m) ↓	0.38	0.47	1.10	1.21	0.38	0.51	1.32	1.21
Q2-invisible audio ↑	72.46	70.40	57.14	48.00	73.03	70.51	52.91	50.63
O3-invisible audio								
3-field accuracy ↑	59.52	47.29	34.14	18.45	41.64	45.56	25.49	18.22
Avg. distance error (m) ↓	0.75	0.98	1.20	1.33	1.02	1.12	1.39	1.38
O4-invisible audio								
DoA accuracy ↑	41.18	16.71	11.18	11.76	13.53	16.47	11.29	9.88
Avg. DoA error (°) \	39.81	69.25	80.51	84.56	77.15	75.39	85.24	84.81
Avg. distance error (m) ↓	0.71	1.08	1.13	1.21	1.16	1.04	1.32	1.23
<u>Q5</u>								
class accuracy ↑	72.43	74.26	25.79	25.63	74.87	72.82	22.18	2.78
sounding accuracy ↑	75.60	64.54	59.48	37.72	36.63	65.93	75.93	41.36
Q6								
class accuracy ↑	81.06	81.61	51.78	50.47	83.78	80.72	42.72	3.72
sounding accuracy ↑	72.33	52.33	59.33	38.67	31.94	49.28	75.67	41.67
07								
class accuracy ↑	75.44	70.88	51.64	53.62	37.35	73.53	51.68	2.56
DoA accuracy ↑	73.21	68.57	47.30	7.80	37.52	64.04	48.38	7.89
Avg. DoA error (°) \	14.75	22.41	33.02	88.31	56.66	24.55	35.25	90.92
Avg. distance error (m) \	0.30	0.33	0.50	0.53	0.44	0.36	0.79	0.53
08								
class accuracy ↑	70.27	75.33	48.42	48.02	69.89	71.90	32.51	0.78
DoA accuracy ↑	64.27	46.37	47.69	8.46	43.72	39.76	49.41	7.61
Avg. DoA error (°) ↓	23.78	50.80	32.32	89.93	51.90	52.46	32.45	89.40
Avg. distance error (m) ↓	0.36	0.44	0.48	0.51	0.42	0.46	0.85	0.52
O9								
sound accuracy ↑	54.00	51.14	51.14	52.20	27.17	55.57	47.25	2.81
visiblity accuracy †	75.22	72.94	38.99	39.79	33.31	76.35	49.42	42.31

5.3 Ablation Studies

Table 3 shows that both image (R: RGB) and binaural audio (B) inputs are crucial for spatial reasoning. It compares R+B+Q, R+M+Q (M: monaural), B+Q, M+Q, and R+Q (Q: question), highlighting that binaural audio provides spatial cues while monaural lacks directional information. The following is an analysis of the performance for each question type.

Question 1 involves sound and visual object classification, with half of the samples containing a non-matching visual object at the sound source. Both R+B+Q and R+M+Q show similar sound classification accuracy (52.69% and 51.01%), suggesting comparable semantic cues from monaural and binaural audio. However, in coming-from accuracy, R+B+Q (69.64%) outperforms R+M+Q (64.10%), highlighting the spatial advantage of binaural audio.

Questions 2, 3, and 4 assess distance and relative location between the sound source and visual objects, requiring spatial reasoning across modalities. For visible audio, R+M+Q achieves 66.52% in Q3 and 59.03% in Q4, performing similarly to R+B+Q (69.73% and 65.68%). When the sound source is invisible, R+B+Q shows a clear advantage, outperforming R+M+Q in Q3 (59.52% vs. 47.29%) and Q4 (41.18% vs. 16.71%). This highlights the role of binaural audio in capturing spatial cues that monaural audio with visual input cannot provide.

Questions 5 and 6 both involve identifying the sound-producing object but differ in complexity

based on the number of visual objects that match the sound. In Q5, with only one matching object, visual context alone provides sufficient spatial information for localization. R+M+Q leverages visual cues effectively, achieving a sounding accuracy of 64.54%. With no visual ambiguity, the model can reliably associate the sound with the correct object using spatial information from the visual signal. In Q6, two visually similar objects match the sound, introducing ambiguity. R+M+Q's performance drops to 52.33%, as visual context alone is no longer sufficient to distinguish between the two objects, leading to random guessing. In contrast, B+Q and R+B+Q maintain consistent performance across both questions. In Q5, they achieve 59.48% and 75.60%, respectively, and in Q6, their performance remains stable at 59.33% and 72.33%. This stability is due to binaural audio, which provides explicit spatial cues, enabling the model to localize the sound source based solely on directional information, unaffected by visual similarity. These results indicate that when there is only one matching object (Q5), R+M+Q can effectively use visual spatial information. However, when multiple visually similar objects are present (Q6), spatial audio cues become essential, allowing B+Q and R+B+Q to maintain stable performance regardless of visual similarity. These results highlight the importance of binaural audio in resolving ambiguity in complex visual scenes.

Questions 7 and 8 both involve sound classification and localization but differ in the number of visual objects that correspond to the audio, with two in Q8 and one in Q7. In Q8, two visually similar objects correspond to the audio, making it difficult for the model to distinguish between them using visual information alone. R+M+Q and B+Q show similar DoA accuracy (46.37% and 47.69%), but their Avg. DoA errors differ, with R+M+Q at 50.80° and B+Q at 32.32°. R+M+Q relies on visual context for spatial cues, but semantic ambiguity between the two objects complicates localization, leading to random selection and higher error. In contrast, B+Q, using binaural audio, focuses solely on directional information, perceiving only one sound source without considering object-level ambiguity, resulting in a lower error. R+B+Q achieves the lowest error (23.78°) by combining spatial audio and visual inputs. In Q7, the audio corresponds to a single object, eliminating semantic ambiguity. In this case, the performance of R+M+Q and B+Q reverses from Q8. R+M+Q records a lower error (22.41°) than B+Q (33.02°), indicating that when only one object is present, visual spatial information can effectively guide localization without semantic confusion. These results support the findings in Q5 and Q6, emphasizing the role of spatial audio in disambiguating visually similar objects.

Question 9 involves sound classification and localization while also requiring the model to determine whether the object is visually present at the sound source. This task demands both audio and visual semantic understanding. Both multi-modal settings (R+B+Q, R+M+Q) successfully address this question.

Modality Setting Cross-Evaluation. To assess the impact of vision signals and binaural audio during training, we evaluate the model trained on R+B+Q under R+M+Q and B+Q settings. While Q7 and Q8 show minimal change, Q5 and Q6 exhibit noticeable gaps in sounding accuracy. This might come from Q5 and Q6 only requiring yes/no responses given a location, without the detailed localization required in Q7 and Q8. Consequently, the model in the B+Q setting may not effectively leverage spatial reasoning for these tasks. However, with visual signals, the model gains implicit spatial cues that align audio locations with the visual scene, potentially enhancing spatial audio understanding. Thus, the presence of visual information may be beneficial even for learning spatial audio cues.

6 Conclusion

We introduce a new task, audio-visual spatial reasoning, along with the *Hear You Are LLM* and QA dataset. Unlike prior work that focuses on semantic or temporal alignment, our approach emphasizes spatial reasoning by integrating binaural audio and visual inputs. We build a large-scale dataset covering diverse spatial scenarios and propose a multimodal framework combining spatial encoders with a large language model. Experiments show that monaural audio with vision or unimodal binaural methods lack the capacity for spatial reasoning. These results underscore the importance of spatial reasoning in robust multimodal understanding and set a new benchmark in audio-visual learning.

7 Limitations and Future Directions

While our framework effectively addresses spatial reasoning by integrating binaural audio and visual context, several real-world scenarios remain unaddressed. Specifically, our approach does not consider moving sound sources, actions associated with visual objects, or occluded objects positioned behind walls or in separate rooms. These aspects are critical for capturing dynamic spatial interactions. Future work will focus on extending the dataset to incorporate these complexities, enabling more comprehensive audio-visual reasoning in realistic settings.

2 References

- 1353 [1] A. Baade, P. Peng, and D. Harwath. Mae-ast: Masked autoencoding audio spectrogram transformer. In *INTERSPEECH*, 2022.
- [2] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, et al. Matterport3D: Learning from RGB-D data in indoor environments. *In 3DV*, 2017.
- 357 [3] M. Chatterjee, J. Le Roux, N. Ahuja, and A. Cherian. Visual scene graphs for audio source separation. In *ICCV*, 2021.
- [4] C. Chen, C. Schissler, S. Garg, P. Kobernik, A. Clegg, P. Calamia, D. Batra, P. Robinson, and K. Grauman. Soundspaces 2.0: A simulation platform for visual-acoustic learning. *NeurIPS*, 2022.
- [5] C. Chen, W. Sun, D. Harwath, and K. Grauman. Learning audio-visual dereverberation. arXiv preprint
 arXiv:2106.07732, 2021.
- [6] H. Chen, W. Xie, T. Afouras, A. Nagrani, A. Vedaldi, and A. Zisserman. Audio-visual synchronisation in the wild. In *BMVC*, 2021.
- [7] H. Chen, W. Xie, T. Afouras, A. Nagrani, A. Vedaldi, and A. Zisserman. Localizing visual sounds the hard
 way. In CVPR, 2021.
- [8] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020.
- [9] Z. Cheng, S. Leng, H. Zhang, Y. Xin, X. Li, G. Chen, Y. Zhu, W. Zhang, Z. Luo, D. Zhao, and L. Bing.
 Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. arXiv preprint
 arXiv:2406.07476, 2024.
- 372 [10] S. Chowdhury, S. Nag, S. Dasgupta, J. Chen, M. Elhoseiny, R. Gao, and D. Manocha. Meerkat: Audio-373 visual large language model for grounding in space and time. In *ECCV*, 2024.
- 111 S. Chowdhury, S. Nag, S. Dasgupta, Y. Wang, M. Elhoseiny, R. Gao, and D. Manocha. Avtrustbench:
 Assessing and enhancing reliability and robustness in audio-visual llms. *arXiv preprint arXiv:2501.02135*,
 2025.
- 377 [12] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In ACCV, 2017.
- 138 [13] B. Devnani, S. Seto, Z. Aldeneh, A. Toso, E. Menyaylenko, B.-J. Theobald, J. Sheaffer, and M. Sarabia. Learning spatially-aware language and audio embeddings. In *NeurIPS*, 2024.
- [14] C. Gan, D. Huang, H. Zhao, J. B. Tenenbaum, and A. Torralba. Music gesture for visual sound separation.
 In CVPR, 2020.
- 382 [15] R. Gao and K. Grauman. Co-separating sounds of visual objects. In ICCV, 2019.
- 183 [16] R. Gao and K. Grauman. Visualvoice: Audio-visual speech separation with cross-modal consistency. In CVPR, 2021.
- 185 [17] K. Gong, K. Feng, B. Li, Y. Wang, M. Cheng, S. Yang, J. Han, B. Wang, Y. Bai, Z. Yang, et al. Avodyssey bench: Can your multimodal llms really understand audio-visual information? arXiv preprint arXiv:2412.02611, 2024.
- 388 [18] Y. Gong, Y.-A. Chung, and J. Glass. Ast: Audio spectrogram transformer. In INTERSPEECH, 2021.
- 189 [19] Y. Gong, H. Luo, A. H. Liu, L. Karlinsky, and J. R. Glass. Listen, think, and understand. In ICLR, 2024.
- 390 [20] Y. He, S. Shin, A. Cherian, N. Trigoni, and A. Markham. Soundloc3d: Invisible 3d sound source localization and classification using a multimodal rgb-d acoustic camera. In *WACV*, 2025.
- [21] D. Hu, R. Qian, M. Jiang, X. Tan, S. Wen, E. Ding, W. Lin, and D. Dou. Discriminative sounding objects
 localization via self-supervised audiovisual matching. *NeurIPS*, 2020.
- E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 396 [23] X. Hu, Z. Chen, and A. Owens. Mix and localize: Localizing sound sources in mixtures. In CVPR, 2022.
- 397 [24] P.-Y. Huang, H. Xu, J. Li, A. Baevski, M. Auli, W. Galuba, F. Metze, and C. Feichtenhofer. Masked autoencoders that listen. In *NeurIPS*, 2022.

- [25] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image
 encoders and large language models. In *International conference on machine learning*. PMLR, 2023.
- 401 [26] Z. Li, Q. Xu, D. Zhang, H. Song, Y. Cai, Q. Qi, R. Zhou, J. Pan, Z. Li, V. Tu, et al. Groundinggpt: Language enhanced multi-modal grounding model. In *ACL*, 2024.
- 403 [27] J. Liu, C. Ju, W. Xie, and Y. Zhang. Exploiting transformation invariance and equivariance for self-404 supervised sound localisation. In *ACM MM*, 2022.
- 405 [28] T. Lüddecke and A. Ecker. Image segmentation using text and image prompts. In CVPR, 2022.
- 406 [29] S. Mo and P. Morgado. A closer look at weakly-supervised audio-visual source localization. In *NeurIPS*, 407 2022.
- 408 [30] S. Mo and Y. Tian. Audio-visual grouping network for sound localization from mixtures. In CVPR, 2023.
- 409 [31] S. Park, A. Senocak, and J. S. Chung. Can clip help sound source localization? In WACV, 2024.
- 410 [32] S. J. Park, M. Kim, J. Hong, J. Choi, and Y. M. Ro. Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory. In *AAAI*, 2022.
- 412 [33] R. Qian, D. Hu, H. Dinkel, M. Wu, N. Xu, and W. Lin. Multiple sound sources localization from coarse to 413 fine. In *ECCV*, 2020.
- 414 [34] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
 415 J. Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- 416 [35] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- 418 [36] S. Sakshi, U. Tyagi, S. Kumar, A. Seth, R. Selvakumar, O. Nieto, R. Duraiswami, S. Ghosh, and D. Manocha. Mmau: A massive multi-task audio understanding and reasoning benchmark. In *ICLR*, 2024.
- 420 [37] A. Senocak, T.-H. Oh, J. Kim, M.-H. Yang, and I. S. Kweon. Learning to localize sound source in visual scenes. In *CVPR*, 2018.
- 422 [38] A. Senocak, T.-H. Oh, J. Kim, M.-H. Yang, and I. S. Kweon. Learning to localize sound sources in visual scenes: Analysis and applications. *IEEE TPAMI*, 2021.
- 424 [39] A. Senocak, H. Ryu, J. Kim, T.-H. Oh, H. Pfister, and J. S. Chung. Sound source localization is all about cross-modal alignment. In *ICCV*, 2023.
- 426 [40] A. Senocak, H. Ryu, J. Kim, T.-H. Oh, H. Pfister, and J. S. Chung. Aligning sight and sound: Advanced sound source localization through audio-visual alignment. *arXiv* preprint arXiv:2407.13676, 2024.
- 428 [41] K. Shimada, A. Politis, P. Sudarsanam, D. A. Krause, K. Uchida, S. Adavanne, A. Hakala, Y. Koyama, N. Takahashi, S. Takahashi, et al. Starss23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events. In *NeurIPS*, 2023.
- 431 [42] Z. Song, Y. Wang, J. Fan, T. Tan, and Z. Zhang. Self-supervised predictive learning: A negative-free 432 method for sound source localization in visual scenes. In *CVPR*, 2022.
- 433 [43] W. Sun, J. Zhang, J. Wang, Z. Liu, Y. Zhong, T. Feng, Y. Guo, Y. Zhang, and N. Barnes. Learning audio-visual source localization via false negative aware contrastive learning. In *CVPR*, 2023.
- 435 [44] K. Sung-Bin, O. Hyun-Bin, J. Lee, A. Senocak, J. S. Chung, and T.-H. Oh. Avhbench: A cross-modal hallucination benchmark for audio-visual large language models. *arXiv* preprint arXiv:2410.18325, 2024.
- 437 [45] R. Tan, A. Ray, A. Burns, B. A. Plummer, J. Salamon, O. Nieto, B. Russell, and K. Saenko. Language-438 guided audio-visual source separation via trimodal consistency. In *CVPR*, 2023.
- 439 [46] Y. Tian, D. Hu, and C. Xu. Cyclic co-learning of sounding object visual grounding and sound separation.
 440 In *CVPR*, 2021.
- [47] M. Tschannen, A. Gritsenko, X. Wang, M. F. Naeem, I. Alabdulmohsin, N. Parthasarathy, T. Evans,
 L. Beyer, Y. Xia, B. Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved
 semantic understanding, localization, and dense features. arXiv preprint arXiv:2502.14786, 2025.
- 444 [48] E. Tzinis, S. Wisdom, A. Jansen, S. Hershey, T. Remez, D. P. Ellis, and J. R. Hershey. Into the wild with audioscope: Unsupervised audio-visual separation of on-screen sounds. In *ICLR*, 2020.

- 446 [49] J. Xu, W. Cheng, Y. Gao, X. Wang, S. Gao, and Y. Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv preprint arXiv:2404.07191*, 2024.
- 448 [50] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, et al. Qwen2 technical report, 2024. *arXiv preprint arXiv:2407.10671*, 2024.
- 450 [51] Q. Ye, Z. Yu, R. Shao, X. Xie, P. Torr, and X. Cao. Cat: Enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. In *ECCV*, 2024.
- 452 [52] Z. Zheng, P. Peng, Z. Ma, X. Chen, E. Choi, and D. Harwath. Bat: Learning to reason about spatial sounds with large language models. In *ICLR*, 2024.

4 NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We address audio-visual spatial reasoning by introducing a dataset that emphasizes spatial alignment over semantic correspondence and propose a model that integrates spatial sound cues with visual perception for enhanced multimodal reasoning.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The limitations are addressed in Sec. 7.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
 only tested on a few datasets or with a few runs. In general, empirical results often
 depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA].

506

507

508

509

510

511

512

513

516

517

518

519

520

521

522

523

525

526

527

528

529

530

531

532

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551 552

553

554

556

557

558

559

Justification: This work focuses on empirical evaluations rather than theoretical formulations or proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented
 by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide all implementation details in Sec. 5.1 and Appendix A.1, and the data and code will be released upon paper acceptance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: After acceptance, we will publish the code, data, and model for public use.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide all training and test details in Sec. 5.1 and Appendix A.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We did not include error bars as we were unable to run sufficient experiments due to resource limitations. The model was tuned using the validation set and the final performance was reported on the test set.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
 - The assumptions made should be given (e.g., Normally distributed errors).
 - It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

612

613

614

615 616

617

618

619

620

623

624

625

626

627 628

629

630 631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

648

649

650

651

653

654

655

656

657

658

659

660

661

663

Justification: We provide the regarding information in Sec. 5.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Yes, our work aligns with the ethical guidelines set by NeurIPS.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We provide the regarding information in Appendix A.3.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: Our study itself does not present significant risks. We employ publicly accessible diffusion models known to have certain risks, and we refer readers to their model cards for safeguard information.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We provide the regarding information in Appendix A.4.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739 740

741

742

745

746

747

748

749

750

751

752

753

754

755

756

757

758 759

760

761

762

763

764

765

767

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We provide the detail about the dataset in Appendix A.2.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The study does not include crowdsourcing or human subject research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA].

Justification: The study does not include crowdsourcing or human subject research.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: We provide the regarding information in Appendix A.5.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.