OBJECT-AWARE AUDIO-VISUAL SOUND GENERATION

Anonymous authors

Paper under double-blind review

Abstract

Generating accurate sounds for complex audio-visual scenes is challenging, especially when multiple objects and sound sources are present. In this paper, we introduce an *object-aware sound generation* model that aligns generated sounds with visual objects in a scene. By grounding sound generation in object-centric representations, our model learns to associate specific visual objects with their corresponding sounds. We fine-tune a conditional latent diffusion model with dot-product attention to improve sound-object alignment. At test time, users can compositionally generate sounds by selecting objects via segmentation masks. We theoretically validate our test-time object-grounding ability, ensuring that even subtle sounds can be represented. Quantitative and qualitative evaluations show that our model outperforms baselines, achieving better alignment between objects and their associated sounds.



Input Image

Audio generated to match a user-selected object

Figure 1: **Object-aware sound generation**. We generate sound aligned with specific visual objects in complex scenes. Users can select objects in the scene using segmentation masks, and the model generates audio corresponding to the selected objects. Here, we show a busy street with multiple sound sources (left). After training, our model generates object-specific audio (right), such as crowd noise for people, engine sounds for cars, and ambient wind for the sky. **Please refer to our supplement and project webpage to watch and listen to the results.**

1 INTRODUCTION

Generating the full sound texture (McDermott & Simoncelli, 2011) of real-world environments is a significant challenge in audio and audio-visual research. While early models have focused on synthesizing sound based on scene categories, text descriptions, and visual contexts (Kong et al., 2019; Yang et al., 2023; Van Den Doel et al., 2001), they often fail to represent specific sound sources in complex environments. In scenes such as a busy city street (Figure 1), where multiple distinct sound events (e.g., car engines, footsteps, crowd noise) co-occur, these models often produce incomplete soundscapes (Pijanowski et al., 2011), overlooking important audio events.

Existing approaches can be largely classified as vision-based or text-based. Vision-based models (Sheffer & Adi, 2023) attempt to synthesize sound by analyzing the entire visual scene, but in
environments with many overlapping sound sources, they tend to generate blended audio that misses
subtle yet important details, like footsteps. Text-based models (Liu et al., 2023) respond to detailed
prompts but face a similar challenge: certain sound events are either *forgotten* or *underrepresented* due to differences in the weight of each event in the latent space. For example, given a prompt de-

scribing both prominent and subtle sounds in a scene, the model might focus on only some of these
events, omitting others like footsteps, even though they were explicitly mentioned (Wu et al., 2023).
This occurs because the model assigns less importance to certain sounds, causing them to be ignored
or poorly generated. While some have attempted to manually reweight sound events in the latent
space (Xue et al., 2024), such interventions remain labor-intensive and impractical for large-scale
applications.

060 To overcome these limitations, we propose an *object-aware sound generation* model that grounds 061 sound generation in the visual domain. Inspired by object-centric learning (Greff et al., 2019), which 062 decomposes scenes into discrete objects, our model associates visual objects with their correspond-063 ing sound sources, ensuring that no sound events are overlooked. We build on an off-the-shelf con-064 ditional audio generation model (Liu et al., 2023), enhancing it with dot-product attention (Vaswani et al., 2017) to learn sound-object associations through self-supervision. This method overcomes 065 the problem of *forgetting* sound events, enabling the generation of various relevant sounds in com-066 plex scenes. To provide finer control and interactivity, we replace the attention with segmentation 067 masks (Kirillov et al., 2023) at test time, allowing users to select specific objects in a scene (e.g., 068 cars, groups of people) to generate the corresponding sounds within simple mouse clicks. This en-069 sures that even subtle sound events, like footsteps or distant conversations, are captured accurately by grounding sound generation in specific objects rather than relying on scene-wide analysis. 071

Through quantitative evaluations and human perceptual studies, we demonstrate that our model outperforms existing baselines, generating more complete and contextually relevant soundscapes. In addition, we provide qualitative results and theoretical analysis demonstrating that our object-grounding mechanism is functionally equivalent to segmentation masks. Through our evaluations, we show:

- Visual grounding from text provides supervision for learning compositional sound generation.
- Specifying different objects within a scene leads to predictable changes in the types of generated sounds.
- Our model learns to generate sound from in-the-wild visual data.
- 081

077

079

2 RELATED WORK

082 083

084 **Predicting sound from images and text.** Generating sounds from visual and textual inputs has 085 gained notable attention recently. Image-based methods focus on synthesizing sounds from visual cues such as physical interactions (Van Den Doel et al., 2001; Owens et al., 2016), human move-086 ments (Gan et al., 2020; Su et al., 2021; Ephrat & Peleg, 2017; Prajwal et al., 2020; Hu et al., 2021), 087 musical instrument performances (Koepke et al., 2020), and content from open-domain images and 880 videos (Zhou et al., 2018; Iashin & Rahtu, 2021; Sheffer & Adi, 2023; Luo et al., 2023). These 089 approaches typically generate audio that corresponds to the entire visual scene without isolating in-090 dividual sound sources, resulting in holistic sound generation. Text-based methods aim to produce 091 sounds from textual descriptions using generative models like GANs and diffusion models (Yang 092 et al., 2023; Kreuk et al., 2023; Liu et al., 2023; Huang et al., 2023b). However, when prompts con-093 tain multiple sound events, these methods often struggle to capture all the desired audio elements 094 (Wu et al., 2023), potentially missing some sounds. Unlike these models, our method distinguishes 095 itself by generating sounds compositionally and creating individual audio outputs for user-selected objects within images. This offers enhanced control and precision in sound generation. 096

Object discovery. Object-centric learning aims to represent visual scenes as compositions of dis-098 crete objects, enabling models to understand and manipulate individual entities within a scene. Unsupervised object discovery methods have been developed to decompose scenes into object repre-100 sentations without explicit annotations (Greff et al., 2019; Burgess et al., 2019). The Slot Attention 101 mechanism (Locatello et al., 2020) introduced a way to learn such representations by utilizing a set 102 of latent variables, or "slots," that iteratively attend to different parts of the input to capture indi-103 vidual objects. Subsequent works (Greff et al., 2019; Burgess et al., 2019) have sought to enhance 104 the stability and robustness of these models. In the audio-visual realm, prior studies have explored 105 object discovery (Arandjelovic & Zisserman, 2018; Rouditchenko et al., 2019; Afouras et al., 2020; Chen et al., 2021; Mo & Morgado, 2022; Hamilton et al., 2024) by leveraging the correspondence 106 between audio and visual modalities. However, these methods primarily focus on recognition and 107 localization tasks and do not address the generation of audio content based on visual inputs. In



Figure 2: **Model architecture.** We encode the reference spectrogram via a pre-trained latent encoder. An image and text prompt are processed by separate encoders, and their embeddings are fused using an attention mechanism to highlight relevant objects. We then feed these conditioned features and noisy latent into a latent diffusion model to generate the object-specific audio. Finally, the latent decoder reconstructs the spectrogram, and a pre-trained HiFi-GAN vocoder generates the final audio waveform. At test time, we replace the attention with a user-provided segmentation mask, and the latent encoder for the reference spectrogram is *not* used.

127 128 129

132

123

124

125

126

contrast, our model generates sounds corresponding to user-selected objects within visual frames,
 without requiring explicit object segmentations and representations during training.

133 Audio-visual learning. Many works have focused on audio-visual associations due to their inherent correspondence in videos. A line of works explores the semantic correspondence, identifying 134 which sounds and visuals are commonly associated with one another (Arandjelovic & Zisserman, 135 2017). This includes representation learning (Morgado et al., 2021; Huang et al., 2023a), source lo-136 calization (Chen et al., 2021; Harwath et al., 2018; Chen et al., 2023), audio stylization (Chen et al., 137 2022a; Li et al., 2024), as well as scene classification (Chen et al., 2020; Gemmeke et al., 2017; Du 138 et al., 2023a) and generation (Li et al., 2022b; Sung-Bin et al., 2023). Other studies leverage spa-139 tial correspondence between audio and visual streams (Owens & Efros, 2018; Korbar et al., 2018; 140 Patrick et al., 2021) to tackle tasks like source separation (Zhao et al., 2018; 2019; Ephrat et al., 141 2016; Gao et al., 2018; Li et al., 2020), Foley sound synthesis (Owens et al., 2016; Du et al., 2023b), 142 and audio spatialization (Gao & Grauman, 2019; Morgado et al., 2018; Yang et al., 2020). Inspired 143 by these works, we aim to generate sound from the user-selected objects within visual frames.

144 145

146 147

3 OBJECT-AWARE SOUND GENERATION

Our goal is to generate sound from user-selected objects within a scene in a compositional way. We cast this problem by learning the correlation between audio and its corresponding visual scene and then using this correlation to predict the sound from the activated region. To achieve this, we: (i) fine-tune an off-the-shelf conditional audio generation model for sound synthesis; (ii) train an audioguided visual object grounding model to isolate the desired object; (iii) theoretically demonstrate the equivalence between the segmentation mask and our grounding model.

154 155

156

3.1 CONDITIONAL AUDIO GENERATION MODEL

157 **Conditional latent diffusion model.** We adopt a pre-trained conditional latent diffusion model (Liu et al., 2023) to generate audio conditioned on textual inputs. Building upon denoising diffusion probabilistic models (Ho et al., 2020) and latent diffusion models (Rombach et al., 2022), our model operates in a compressed latent space to improve computational efficiency. Specifically, given a text prompt t_q describing the desired sound and a noise vector $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the model iteratively denoises the latent variables over N steps to generate the corresponding audio. 162 Our model is trained to predict the added noise at each denoising step n, conditioned on the textual 163 input t_q . The training objective minimizes the difference between the predicted noise and the true 164 noise:

 $\mathcal{L}_{\theta} = \mathbb{E}_{\boldsymbol{z}_0, \boldsymbol{t}_q, \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \mathbf{I}), n} \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\boldsymbol{z}_n, n, \boldsymbol{t}_q) \|_2^2, \qquad (1)$

167 $\mathcal{L}_{\theta} = \mathbb{L}_{z_0, t_q, \epsilon \sim \mathcal{N}}(0, 1), n \| \mathcal{C} = \mathcal{C}_{\theta}(\mathcal{A}_n, n; \mathbf{c}_q) \|_2^2$, (1) 168 where z_0 is the latent representation of the ground truth audio, z_n is the noisy latent at step n, and 169 ϵ_{θ} is the denoising model parameterized by θ .

170 Mel-spectrograms compression. We compress mel-spectrograms into a lower-dimensional latent 171 space using a variational autoencoder (VAE) (Kingma & Welling, 2013). The VAE encodes the mel-172 spectrogram $a \in \mathbb{R}^{T \times F}$ into a latent representation $z \in \mathbb{R}^{T' \times F' \times d}$, where T' and F' are reduced 173 temporal and frequency dimensions, and d is the dimensionality of the latent embeddings.

Textual representation. We represent the textual input t_q using a pre-trained text encoder from CLAP (Elizalde et al., 2023), which maps the text into an embedding space $\mathcal{E}_t(t_q) \in \mathbb{R}^L$, where *L* denotes the embedding dimension. These text embeddings capture semantic information about the desired sound and are used to condition the diffusion model through cross-attention mechanisms (Vaswani et al., 2017).

Classifier-free guidance. We employ classifier-free guidance (CFG) (Ho & Salimans, 2022) to encourage the model to learn both conditional and unconditional denoising. During training, we randomly omit the conditioning input t_q with a 10% probability. At test time, we use a guidance scale $\lambda \ge 1$ to interpolate between the conditional and unconditional predictions:

184 185

189

195

196

208

213

214 215 $\tilde{\boldsymbol{\epsilon}}_{\theta}(\boldsymbol{z}_n, n, \boldsymbol{t}_q) = \lambda \cdot \boldsymbol{\epsilon} \theta(\boldsymbol{z}_n, n, \boldsymbol{t}_q) + (1 - \lambda) \cdot \boldsymbol{\epsilon}_{\theta}(\boldsymbol{z}_n, n, \varnothing) , \qquad (2)$

where $\epsilon_{\theta}(z_n, n, \emptyset)$ is the unconditional prediction. This approach enhances adherence to the conditioning text while maintaining diversity in the generated audio.

190Waveform reconstruction.After generating the latent representation of the audio, we reconstruct191the corresponding waveform.The decoder part of the VAE transforms the latent representation z_0 192back into a mel-spectrogram.Subsequently, a pre-trained HiFi-GAN neural vocoder (Kong et al.,1932020a) is used to synthesize the time-domain audio waveform from the mel-spectrogram, producing194

3.2 TEXT-GUIDED VISUAL OBJECT GROUNDING MODEL

197 Visual representation. To ground the visual objects corresponding to the desired sound, we ex-198 tract features from the input image using a pre-trained visual encoder. Specifically, we utilize CLIP 199 (Radford et al., 2021) to encode the image into a set of visual patches embeddings $\mathcal{E}_v(i_q) \in \mathbb{R}^{P \times L}$, 200 where i_q is the input image, P is the number of patches, and L denotes the embedding dimen-201 sion (matching that of the text embeddings). These embeddings capture both semantic and spatial 202 information of the visual scene.

Scaled dot-product attention. We employ scaled dot-product attention (Vaswani et al., 2017) to fuse the textual and visual inputs, allowing the model to focus on specific objects within the scene. Before computing the attention, the text embeddings $\mathcal{E}_t(t_q)$ and patch embeddings $\mathcal{E}_v(i_q)$ are linearly projected to obtain the query, key, and value matrices. Specifically, we compute:

$$\boldsymbol{Q} = \mathcal{E}_t(\boldsymbol{t}_q)\boldsymbol{W}^Q, \quad \boldsymbol{K} = \mathcal{E}_v(\boldsymbol{i}_q)\boldsymbol{W}^K, \quad \boldsymbol{V} = \mathcal{E}_v(\boldsymbol{i}_q)\boldsymbol{W}^V, \tag{3}$$

where W^Q , W^K , and W^V are learnable projection matrices.

We then computes the attention weights between the projected text and each projected image patch, grounding the text in the visual domain:

Attention
$$(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \operatorname{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^{\top}}{\sqrt{d_k}}\right)\boldsymbol{V}$$
, (4)

where d_k is the dimensionality of the key embedding.

After obtaining the attention output, we apply an MLP layer (Murtagh, 1991) to further refine the fused representations, which enables the model to attend to image regions corresponding to the text input. In this way, we integrate the images i_q with the diffusion process, allowing the model to learn to focus on the relevant regions in the image through self-supervision.

- Learnable positional encoding. To enhance the model's ability to localize objects within the
 image, we incorporate learnable positional encodings (Devlin, 2018) into the attention mechanism.
 These encodings are added to the key and value embeddings, providing spatial information about
 the image patches. By learning positional information, the model can better distinguish between
 objects in different locations, improving grounding precision.
- 226

220

Segmentation mask at test time. After training, we have the flexibility to substitute the attention weights derived from the scaled dot-product attention with segmentation masks generated by the segment anything model (SAM) (Kirillov et al., 2023). We rescale the raw outputs of SAM into a normalized mask $m_q \in \mathbb{R}^P$, matching the mean and variance of the attention weights. This allows us to generate the desired object's sound by focusing on the regions specified by the segmentation mask. Since SAM's masks can be obtained using either text prompts or point clicks, our model supports interactive and compositional sound generation, allowing users to intuitively select objects of interest and generate their associated sounds.

234 235

236

244 245

246 247

260 261

3.3 THEORETICAL ANALYSIS

One may notice that our training pipeline uses both text and image encoders, but the test-time computation involves only the image encoder, where the softmax attention weights are replaced by the segmentation masks. This indicates an *out-of-distribution* generalization ability, where our model trained on the softmax attention weights computed by CLAP & CLIP embeddings (Equation 4) is able to generalize well on the segmentation masks computed by SAM. We hypothesize that this ability is rooted in the alignment of contrastive losses and the dot-product attention mechanism. Recall that the InfoNCE loss (Oord et al., 2018) for the text encoder in contrastive learning is given by:

$$\mathcal{L}_{t}\left(\mathcal{E}_{t}, \mathcal{E}_{v}\right) = \mathbb{E}_{x^{T}, x_{1:N}^{I}} \left[-\log \frac{\exp\left(\langle \mathcal{E}_{v}(x^{T}), \mathcal{E}_{t}(x_{1}^{I}) \rangle / \tau\right)}{\sum_{j=1}^{N} \exp\left(\langle \mathcal{E}_{v}(x^{T}), \mathcal{E}_{t}(x_{j}^{I}) \rangle / \tau\right)} \right]$$
(5)

where (x^T, x_1^I) is the matching text-image pair, and x_2^I, \ldots, x_N^I are the negative image samples associated with x^T . Notice that if we substitute x^T with the text input t_q , $x_{1:N}^I$ with the image patches i_q , and x_1^I with the matching image patch (with the text input), then the loss in Equation 5 becomes the Maximum Likelihood Estimation (MLE) loss of the softmax attention weights in Equation 4 (under proper scaling in the exponents). Therefore, the encoders \mathcal{E}_v , \mathcal{E}_t are able to assign high attention weights to image patches that match with textual inputs, and low attention weights to irrelevant image patches, working effectively as the segmentation mask at test time. As such, the audio generation model is trained with the ability to focus only on the selected objects by segmentation masks.

In the following theorem, we formalize the above argument into a test-time error guarantee. We let f denote the composition of the trained MLP layers and the audio generation model that maps an attention a_q to an audio output s_q on query q, and v denote the value metric that maps a soundimage-mask tuple (s, i, m) to a real number $v(s, i, m) \in \mathbb{R}$. Our goal is to bound

$$\operatorname{err}_{\operatorname{test}} := \mathbb{E}_q[v(f^*(p_q V^*), i_q, p_q) - v(f(\boldsymbol{m}_q V), i_q, \boldsymbol{m}_q)]$$

i.e., the expected (over the randomness of test query q) distance between the optimal value $v(f^*(p_q V^*), i_q, p_q)$ and the value of the trained model $v(f(m_q V), i_q, m_q)$ at test time. Here, f* and V* are the ground-truth counterpart of f and value matrix, $p_q \in \Delta^P$ is the (normalized) ground-truth mask of query q such that $p_{q,k} = \frac{\mathbb{P}(t_q | i_{q,k})}{\sum_{l=1}^{P} \mathbb{P}(t_q | i_{q,l})}$ for patch index $k \in \{1, \ldots, P\}$, a_q represents the attention computed by Equation 4. Note that $f(m_q V)$, the audio output of the trained model, depends on the segmentation mask m_q instead of the ground-truth mask p_q or text input t_q .

Theorem 3.1. Let $\epsilon_{sam} := \mathbb{E}_q[||\mathbf{m}_q - p_q||_{\ell_1}]$ denote the expected ℓ_1 error of the segmentation model. Let ϵ_f, ϵ_V denote the expected error of f and V under the pre-trained CLAP & CLIP embeddings respectively, and $\epsilon_{\text{contrast}}$ denote the expected contrastive loss of the encoders, more precisely,

$$\epsilon_f = \mathbb{E}_q[v(f^*(a_q), \boldsymbol{i}_q, p_q)] - \mathbb{E}[v(f(a_q), \boldsymbol{i}_q, p_q)], \ \epsilon_{\boldsymbol{V}} = \|\boldsymbol{V} - \boldsymbol{V}^*\|_{\infty}$$

276

277

278 279 280

281

282

283

284

285

286

287

288 289

$$\epsilon_{\text{contrast}} = \mathbb{E}_{q,d \sim p_q} \left[-\log \frac{\exp\left(\langle \mathcal{E}_v(\boldsymbol{t}_q), \mathcal{E}_t(i_{q,d}) \rangle_{\Sigma}\right)}{\sum_{k=1}^{P} \exp\left(\langle \mathcal{E}_v(\boldsymbol{t}_q), \mathcal{E}_t(i_{q,k}) \rangle_{\Sigma}\right)} \right] - \mathbb{E}_{q,d \sim p_q} \left[-\log p_{q,d} \right].$$

where $\langle \cdot, \cdot \rangle_{\Sigma}$ is the local inner product under $\Sigma := \mathbf{W}^{K}(\mathbf{W}^{Q})^{\top}/\sqrt{d_{k}}$. Suppose $\|\mathbf{V}^{*}\|_{\infty}, \|\mathbf{V}\|_{\infty} \leq B_{v}, v$ is L_{v} -Lipschitz, and f, f^{*} are L_{f} -Lipschitz, then we have

$$\operatorname{err}_{\operatorname{test}} \leq L_v \cdot \left(L_f \cdot \left(\epsilon_{\boldsymbol{V}} + B_v \cdot \left(\epsilon_{\operatorname{sam}} + 2\sqrt{2\epsilon_{\operatorname{contrast}}} \right) \right) + \epsilon_{\operatorname{sam}} \right) + \epsilon_f.$$

Due to space constraints, the proof is deferred to Appendix A.5. Theorem 3.1 implies that the testtime error can be upper bounded by the error of the pre-trained CLAP & CLIP encoders, the error of the segmentation model, and the error of the trained model under pre-trained encoders. Since the latter errors are usually small due to massive training and the regularity parameters L_v, L_f, B_v are commonly modest, our method can be guaranteed to achieve high accuracy. This explains why we are able to substitute the attention weights derived from the scaled dot-product attention with segmentation masks generated by the segmentation model at test time. Our theory is further corroborated by Section 4.3, where using dot-product attention weights achieves performance on par with using segmentation masks, while additive attention fails completely.

290 4 EXPERIMENTS 291

292 4.1 EXPERIMENT SETUP293

Dataset. We use the Sound-VECaps dataset (Yuan et al., 2024) as our primary data source. This 294 dataset is derived from AudioSet (Gemmeke et al., 2017), which consists of 4,616 hours of video 295 clips, each paired with corresponding labels and captions. To tailor the dataset for our task, we per-296 form several preprocessing steps: (i) employ Llama (Touvron et al., 2023) to rephrase the original 297 captions, ensuring they focus only on visible sounding objects for better consistency; (ii) exclude 298 clips containing voiceovers and music by applying keyword-based filters such as "speech" and "mu-299 sic"; and (iii) train and use an off-the-shelf audio-visual matching model to retain only those videos 300 with high correspondence scores. This reduces the dataset to 748 hours of video. Please see Ap-301 pendix A.2 for more details on the dataset refinement. 302

Model architecture. Building upon the AudioLDM (Liu et al., 2023), our model integrates image 303 inputs through a grounding model (Sec. 3.2). We employ the same VAE and HiFi-GAN vocoder, 304 which are trained on a combination of the AudioSet (Gemmeke et al., 2017), AudioCaps (Kim et al., 305 2019), BBC Sound Effects (Corporation, 2017), and Freesound (Fonseca et al., 2021) datasets. The 306 VAE is configured with a latent dimensionality d of 8 channels. For embedding extraction, we utilize 307 the "ViT-B/32" CLAP audio encoder (Elizalde et al., 2023) and the CLIP image encoder (Radford 308 et al., 2021). These embeddings are then incorporated into the U-Net-based diffusion model through 309 cross-attention (Vaswani et al., 2017). We implement a linear noise schedule consisting of N = 1000310 diffusion steps, from $\beta_1 = 0.0015$ to $\beta_N = 0.0195$. The DDIM sampling method (Song et al., 2020) 311 is used with 200 steps to facilitate efficient generation. At test time, we apply CFG with a guidance 312 scale λ set to 2, as defined in Equation 2.

313

Training configuration. To facilitate parallel training, each video's soundtrack is either truncated or zero-padded to achieve a fixed duration of 10 seconds and then converted to a 16 kHz sample rate in 32-bit floating-point PCM format. We apply a 512-point discrete Fourier transform with a frame length of 64 ms and a frame shift of 10 ms. For each video, a single visual frame is randomly chosen to serve as the input image. The model is then trained using the AdamW optimizer (Loshchilov & Hutter, 2017) with a batch size of 64, a learning rate of 10^{-4} , $\beta_1 = 0.95$, $\beta_2 = 0.999$, $\epsilon = 10^{-6}$, and a weight decay of 10^{-3} over 300 epochs.

320

Evaluation metrics. We use both quantitative and qualitative metrics (see Appendix A.3 for more evaluation details) to evaluate the performance of our model. For the objective evaluation, we employ several metrics, including Sound Event Accuracy (ACC), which leverages the PANNs model (Kong et al., 2020b) to predict and sample sound event logits based on the annotated labels and then

Method	$ACC (\uparrow)$	FAD (\downarrow)	$\mathbf{KL}\left(\downarrow\right)$	IS (\uparrow)	AVC (\uparrow)	$\left \mathbf{OVL}\left(\uparrow\right) \right.$	RET (†)	REI (\uparrow)	REO (\uparrow)
Ground Truth	/	/	/	/	0.962	$ 4.12 \pm 0.06$	4.02 ± 0.05	4.06 ± 0.07	/
AudioLDM 1	0.314	3.761	1.542	1.541	0.701	2.76 ± 0.03	3.08 ± 0.07	2.88 ± 0.02	2.12 ± 0.03
AudioLDM 2	0.502	2.981	1.141	1.785	0.747	2.97 ± 0.02	3.21 ± 0.04	3.06 ± 0.04	2.44 ± 0.02
Make-an-Audio	0.309	3.555	1.443	1.673	0.712	2.74 ± 0.08	3.06 ± 0.05	2.89 ± 0.05	2.08 ± 0.04
Im2Wav	0.499	3.602	1.526	1.872	0.798	2.88 ± 0.05	3.12 ± 0.04	3.01 ± 0.05	2.48 ± 0.06
SpecVQGAN	0.611	2.515	1.142	1.965	0.825	2.94 ± 0.04	3.26 ± 0.03	3.11 ± 0.06	2.51 ± 0.04
Diff-Foley	0.683	1.908	0.783	2.010	0.842	3.09 ± 0.06	3.43 ± 0.05	3.32 ± 0.03	2.52 ± 0.06
Ours	0.859	1.271	0.517	2.102	0.891	$\mid \textbf{ 3.31} \pm \textbf{0.04}$	$\textbf{3.62} \pm \textbf{0.05}$	$\textbf{3.48} \pm \textbf{0.04}$	$\textbf{3.74} \pm \textbf{0.07}$

Table 1: Quantitative comparison of our method and baselines across different metrics. The subjective OVL, RET, REI, and REO scores are presented with 95% confidence intervals.

332

333

336 compute the mean accuracy across the dataset. We also measure the semantic alignment between 337 the output and target using three established metrics: (i) Fréchet Audio Distance (FAD) (Kilgour 338 et al., 2019), which quantifies how close the generated audio is to the real audio in latent space; 339 (ii) Kullback-Leibler Divergence (KL), which assesses the alignment of distributions between the 340 generated and target audio; and (iii) the Inception Score (IS) (Salimans et al., 2016), which evaluates the diversity of the generated audio. Additionally, the Audio-Visual Correspondence (AVC) (Arand-341 jelovic & Zisserman, 2017) is used to measure the semantic coherence between the input image and 342 the resulting audio, indicating how well the sounds match the visual context. We report this using 343 the average cosine similarity of features extracted by OpenL3 (Cramer et al., 2019). 344

345 For subjective evaluation, we conduct a human study to assess the quality and relevance of the gen-346 erated audio. We present both the holistic samples and the object-selected samples. Each participant is provided with an input image, along with the corresponding generated audio, and is asked to rate 347 each sample on a scale from 1 to 5 based on several criteria: (i) Overall Quality (OVL), which eval-348 uates the general quality of the audio; (ii) Relevance to the Text Prompt (RET), which assesses how 349 well the audio matches any associated text description; (iii) Relevance to the Input Image (REI), 350 which judges the alignment between the audio and the visual content; and Relevance to the Selected 351 Object (REO), which focuses on how well the generated audio aligns with a specific object in the 352 visual scene. 353

Baselines. We compare our method with several baseline models, each of which is adapted for our task:

- AudioLDM 1 & 2 (Liu et al., 2023; 2024): These models are originally designed for text-to-audio generation, but we modify them by swapping their text embeddings with image embeddings. We fine-tune these models on our dataset for a fair comparison.
 Molue on Audia (Usang et al., 2023b): Male on Audia supports either text or image memory for a fair comparison.
 - Make-an-Audio (Huang et al., 2023b): Make-an-Audio supports either text or image prompts for sound generation. We extract its image-based branch and fine-tune it on our dataset.
 - Im2Wav (Sheffer & Adi, 2023): Im2Wav is an image-guided open-domain audio generation model that operates auto-regressively. Since the original model generates only 4 seconds of audio, we retrain it on our dataset to adapt it to our task.
 - **SpecVQGAN** (Iashin & Rahtu, 2021): SpecVQGAN is a two-stream VQGAN model (Esser et al., 2021) designed for video-to-audio generation. We modify it by randomly sampling a single frame from video data and fine-tune it for our task.
 - **Diff-Foley** (Luo et al., 2023): Diff-Foley is a diffusion model that generates sound semantically and temporally aligned with the video. Similar to SpecVQGAN, we fine-tune it on our dataset using randomly sampled video frames.
- 368 369 370

371

360

361

362

363

364

366

367

4.2 COMPARISON TO BASELINES

Quantitative results. Table 1 compares our approach against the baselines on the Sound-VECaps dataset. Our model outperforms the baselines across different metrics, highlighting its ability to produce high-quality audio. In particular, our method achieves the best ACC metrics, indicating its capacity to generate sound closely linked to the visual objects in the scene. Diff-Foley shows competitive performance among the baselines, likely due to its contrastive representations, which map visual and audio features to a shared latent space, improving audio-visual consistency. Although Im2Wav and SpecVQGAN achieve reasonable AVC scores, they struggle with FAD and KL, indi-

³³⁴ 335



Figure 3: Qualitative model comparison. We show sound generation results for our method and the baselines, each of which is conditioned on an image, text, or segmentation mask.

392 cating they fall short in generating high-quality sounds. Similarly, AudioLDM and Make-an-Audio 393 show relatively lower accuracy and semantic alignment, which could be due to their original design 394 for the text-to-audio task rather than the image-guided one. Notably, our model significantly sur-395 passes Diff-Foley in terms of FAD and KL, suggesting that it can generate audio that is not only 396 realistic but also semantically linked to the visual inputs. These results indicate the advantage of our 397 method in leveraging visual cues for more contextually relevant sound generation.

398 For subjective evaluation, we randomly select 100 generated samples from the test set, with 50 of 399 them manually processed to create segmentation masks for specific objects within a scene. These 400 samples are then rated by 50 participants. Our model receives the highest average ratings across all 401 subjective measures, with a particularly notable lead in REO, suggesting that it generates sounds 402 aligned with the objects in the image. Interestingly, we observe that all the baselines achieve rel-403 atively close scores for REO, which demonstrates that our method is particularly good at linking 404 audio to object-level visual cues, a feature that is less evident in the baselines. Moreover, partici-405 pants consistently rated the OVL, RET, and REI of our model higher, further validating the objective metrics and highlighting its improved contextual alignment. 406

Qualitative results. Figure 3 compares our method with the baselines on the Sound-VECaps 408 dataset. In the first example, where both a dog and a goose are present, all baselines only gen-409 erate dog growls, missing the goose honks. Our method, however, captures both sounds, illustrating 410 its object-aware capability. Similarly, in the second and third examples, involving a car with distant 411 chatter and a train with people talking, the baselines produce either one of the sound events but not 412 all simultaneously. By contrast, our model successfully generates the complete soundscape. The 413 final example presents a small jet in the background with the crowd cheering. Vision-based models 414 fail to detect the jet due to the jet's small size in the image, generating only the crowd and wind 415 noises, while text-based models struggle to combine multiple sounds. Our approach accurately captures all relevant sounds, highlighting its ability to generate accurate sounds aligned with complex 416 visual scenes. For a more direct experience, please view the results video in the supplement and on 417 the project webpage. 418

419

407

389

390 391

4.3 ABLATION STUDY AND ANALYSIS 420

421

423

426

430

Table 2 summarizes the ablation experiments. We explore the following model variations: (i) freez-422 ing the latent diffusion weights rather than fine-tuning them; (ii) replacing single-head attention with multi-head attention; (iii) substituting text-image attention with audio-image attention; (iv) altering 424 the attention mechanism from dot-product to additive attention; and (v) using text-image attention 425 instead of segmentation masks during inference. We also show additional results in Appendix A.4.

427 Effect of freezing diffusion weights. We test the impact of freezing the latent diffusion model 428 weights instead of fine-tuning them during training. We observe that freezing the weights degrades 429 the performance, which suggests that fine-tuning is required to achieve more coherent audio.

Impact of attention head. We compare our single-head attention mechanism with the multi-head 431 counterpart (Vaswani et al., 2017). The multi-head approach enhances the alignment between tex-

Method	ACC (†)	FAD (\downarrow)	$\mathbf{KL}\left(\downarrow\right)$	IS (†)	AVC (†)
(i) Frozen Diffusion	0.692	1.543	1.047	1.943	0.733
(ii) Multi-Head Attention	0.415	2.238	1.903	2.115	0.887
(iii) Audio-Image Attention	0.634	1.761	1.232	1.731	0.692
(iv) Additive Attention	0.103	15.747	7.425	1.343	0.137
(v) Text-Image Attention	0.856	1.270	0.520	2.097	0.890
Ours	0.859	1.271	0.517	2.102	0.891

Table 2: Quantitative ablation studies on the Sound-VECaps dataset.

tual inputs and the generated audio, leading to a stronger correspondence between text descriptions and sound outputs. However, this improvement reduces controllability when specifying specific audio characteristics based on the segmentation mask. We conjecture that this limitation arises because each head in the multi-head attention focuses on different regions of the input (Voita et al., 2019; Hamilton et al., 2024). While this strategy increases text-audio alignment, the lack of a clear definition for each head's specific scope reduces the interpretability of the final results. This likely contributes to the masking results deviating from expectations.

Choice of attention modality. We assess the effectiveness of text-image attention compared to audio-image attention. The audio-image attention variant shows a decline in performance, which could be attributed to the inherent limitations of the CLAP model in representing overlapping audios. This limitation probably introduces noise, thereby weakening the model's ability to form audio-visual associations essential for sound generation.

Evaluation of attention scoring mechanism. We investigate the role of the attention scoring
function by replacing dot-product attention with the additive one (Bahdanau, 2014). The additive
attention variant collapses significantly, indicating that segmentation masks are not a suitable replacement for this attention. Explained by the theory in Section 3.3, this could be because addition
operations are not compatible with the contrastive losses used by CLAP & CLIP and segmentation
masks generated by SAM, which disrupts our grounding model.

Role of segmentation masks during inference. We compare the standard text-image attention mechanism to the proposed segmentation masks at test time. The results show that text-image attention achieves performance on par with the segmentation mask approach (ours). This suggests that both methods provide similar levels of spatial and semantic guidance for audio generation. This finding also supports the theory discussed in Section 3.3.

465 466 467

439 440 441

448

466 4.4 CROSS-DATASET EVALUATION

Visualization between grounding and masking. In 468 Figure 4, we visualize the comparison between the atten-469 tion maps generated by our model and the segmentation 470 masks produced by SAM. For this, we use images from 471 Places (Zhou et al., 2017) and text prompts derived from 472 BLIP (Li et al., 2022a). To visualize the attention maps, 473 we apply bilinear interpolation to match the resolution of 474 the segmentation masks. Our results show a strong align-475 ment between our model's attention maps and the seg-476 mentation masks, providing empirical support for the the-477 oretical analysis in Section 3.3 and the findings of the ablation study in Section 4.3. While the segmentation masks 478 represent a form of "hard" attention, directly highlight-479 ing specific regions, our model produces "soft" attention 480 maps that provide a probabilistic focus on the relevant 481 areas. This similarity indicates that, through training, our 482 model effectively learns to capture object-specific regions 483 similar to those identified by segmentation, achieving the 484 desired grounding in a flexible manner. Furthermore, this



Figure 4: **Visualization results**. We visualize the difference between attention maps and segmentation masks using images from Places (Zhou et al., 2017) and text prompts from BLIP (Li et al., 2022a).

⁴⁸⁵ observation suggests that attention maps can be replaced with segmentation masks at test time.

486 Compositional sound 487 generation. We ask 488 whether our model will 489 generate object-specific 490 sounds by isolating individual objects within 491 a scene. As shown in 492 Figure 5, we use the same 493 image for each scene, 494 separating different objects 495 (cars, people, seagulls, 496 etc.) to generate corre-497 sponding audio outputs. 498 The results illustrate that 499 our model successfully 500



Figure 5: **Compositional sound generation**. Our model generates objectspecific sounds in the city (left) and beach (right) scenes, and composes a complete soundscape when multiple objects are selected.

learns to generate distinct sounds for each object, such as car engines or footsteps, reflecting
 their unique sound textures. Furthermore, when multiple objects are selected together, the model
 compositionally generates the entire soundscape that represents the scene property. This capability
 highlights our model's strength in decomposing and synthesizing audio-visual elements for sound
 generation.

adaptation Sound to 506 visual texture changes. 507 We explore whether our 508 method can generate 509 soundscapes that adapt to 510 changes in visual textures, 511 inspired by audio-visual 512 video editing (Lee et al., 513 2023). Starting with images from the Places (Zhou 514



Figure 6: Generating soundscapes from visual texture changes. We generate different soundscapes by manipulating the visual textures of the same scene, such as changing weather (left) or materials (right).

et al., 2017) and Greatest Hits (Owens et al., 2016) datasets, we apply an off-the-shelf image
translation model (Park et al., 2020; Li et al., 2022b) to create paired scenes (e.g., sunny-rainy,
water-grass), and then overlay full-image segmentation masks on top. As illustrated in Figure 6,
our model generates context-appropriate soundscapes. For instance, it generates rain sounds for
dark skies, wind sounds for clear skies, water splashing for watery surfaces, and grass crunching for
grassy areas. This demonstrates that our model successfully captures variations in visual textures to
generate corresponding audio.

521 522 523

524

505

5 CONCLUSION

525 In this paper, we proposed an *object-aware sound generation* model, focusing on aligning generated sounds with specific visual objects in complex scenes. To achieve this, we developed a diffusion 526 model grounded in object-centric representations, enhancing the association between objects and 527 their corresponding sounds. Our theoretical analysis demonstrates that the object-grounding mech-528 anism is functionally equivalent to segmentation masks. Quantitative and qualitative evaluations 529 show that our model surpasses baselines in sound-object alignment, enabling cross-dataset gener-530 alization and compositional sound generation. We hope this work not only advances controllable 531 sound generation but also inspires further exploration into the relationships between objects and 532 soundscapes. 533

Limitations and broader impacts. Our model shows promising results in generating object specific sounds from images but has certain limitations. First, since our model relies on static images, it may struggle to produce non-stationary audio synchronizing with dynamic events, such as impact
 sounds (Figure 6). Additionally, it may lack precise control over the type of sound produced for an
 object, leading to potential ambiguity. For example, a car might be associated with various sounds, such as siren or engine noise (Figure 3). Lastly, while useful for content creation like filmmaking, our model also poses a potential risk, as it could be exploited to create misleading videos.

540 ETHICS STATEMENT 541

This paper introduces an *object-aware sound generation* model. It is trained on publicly available datasets, such as AudioSet and Sound-VECaps, which do not contain personally identifiable information. We have taken steps to ensure compliance with data usage policies, and our model does not involve human subjects or raise privacy concerns. We believe our work poses minimal ethical risks, as it focuses on enhancing sound-object alignment in a controlled research environment. However, we encourage responsible use of our model, particularly when applied to real-world scenarios.

548 549

556

557

561 562

563

564 565

566

567

568

569

573

577

578

579

584

585

586

REPRODUCIBILITY STATEMENT

To facilitate the reproducibility of our results, we provide detailed information in multiple sections of this paper and its appendix. A comprehensive description of the dataset is presented in Section 4.1 of the main paper, with additional data refinement details included in Appendix A.2. The key training configurations, including hyperparameters, are outlined in Section 4.1. Our proposed method is illustrated in Section 3, and the source code has been made available in the supplement for reference.

References

- Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised
 learning of audio-visual objects from video. In *Computer Vision–ECCV 2020: 16th Euro- pean Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pp. 208–224.
 Springer, 2020.
 - Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE international conference on computer vision*, pp. 609–617, 2017.
 - Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *Proceedings of the European* conference on computer vision (ECCV), pp. 435–451, 2018.
 - Dzmitry Bahdanau. Neural machine translation by jointly learning to align and translate. *arXiv* preprint arXiv:1409.0473, 2014.
- 570 Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt
 571 Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and represen572 tation. *arXiv preprint arXiv:1901.11390*, 2019.
- Changan Chen, Ruohan Gao, Paul Calamia, and Kristen Grauman. Visual acoustic matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18858–18868, 2022a.
 - Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audiovisual dataset. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 721–725. IEEE, 2020.
- Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
 - Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*, 2022b.
- Ziyang Chen, Shengyi Qian, and Andrew Owens. Sound localization from motion: Jointly learning sound direction and camera rotation. *arXiv preprint arXiv:2303.11329*, 2023.
- 589 British Broadcasting Corporation. BBC Sound Effects, 2017. Available: https://
 590 sound-effects.bbcrewind.co.uk/search.
 591
- Aurora Linh Cramer, Ho-Hsiang Wu, Justin Salamon, and Juan Pablo Bello. Look, listen, and learn more: Design choices for deep audio embeddings. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3852–3856. IEEE, 2019.

594 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. 595 arXiv preprint arXiv:1810.04805, 2018. 596 Chenzhuang Du, Jiaye Teng, Tingle Li, Yichen Liu, Tianyuan Yuan, Yue Wang, Yang Yuan, and 597 Hang Zhao. On uni-modal feature learning in supervised multi-modal learning. In International 598 Conference on Machine Learning, pp. 8632–8656. PMLR, 2023a. 600 Yuexi Du, Ziyang Chen, Justin Salamon, Bryan Russell, and Andrew Owens. Conditional gener-601 ation of audio from video via foley analogies. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2426–2436, 2023b. 602 603 Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning 604 audio concepts from natural language supervision. In ICASSP 2023-2023 IEEE International 605 Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023. 606 Ariel Ephrat and Shmuel Peleg. Vid2speech: speech reconstruction from silent video. In 2017 IEEE 607 International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5095–5099. 608 IEEE, 2017. 609 610 Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T 611 Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent 612 audio-visual model for speech separation. ACM Transactions on Graphics (TOG), 37(4), 2016. 613 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image 614 synthesis. In Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-615 tion, pp. 12873-12883, 2021. 616 617 Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic Font, and Xavier Serra. Fsd50k: an open dataset of human-labeled sound events. IEEE/ACM Transactions on Audio, Speech, and Language 618 Processing, 30:829–852, 2021. 619 620 Chuang Gan, Deng Huang, Peihao Chen, Joshua B Tenenbaum, and Antonio Torralba. Foley mu-621 sic: Learning to generate music from videos. In Computer Vision-ECCV 2020: 16th European 622 Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16, pp. 758–775. Springer, 623 2020. 624 Ruohan Gao and Kristen Grauman. 2.5 d visual sound. In Proceedings of the IEEE/CVF Conference 625 on Computer Vision and Pattern Recognition, pp. 324-333, 2019. 626 627 Ruohan Gao, Rogerio Feris, and Kristen Grauman. Learning to separate object sounds by watching 628 unlabeled video. In Proceedings of the European Conference on Computer Vision (ECCV), pp. 35-53, 2018. 629 630 Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing 631 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for 632 audio events. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing 633 (ICASSP), pp. 776-780. IEEE, 2017. 634 Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel 635 Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation 636 learning with iterative variational inference. In International conference on machine learning, pp. 637 2424-2433. PMLR, 2019. 638 639 Mark Hamilton, Andrew Zisserman, John R Hershey, and William T Freeman. Separating the" 640 chirp" from the" chat": Self-supervised visual grounding of sound and language. In Proceedings 641 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13117–13127, 2024. 642 643 David Harwath, Adria Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. 644 Jointly discovering visual objects and spoken words from raw sensory input. In Proceedings of 645 the European conference on computer vision (ECCV), pp. 649–665, 2018. 646 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint 647

arXiv:2207.12598, 2022.

648 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in 649 neural information processing systems, 33:6840–6851, 2020. 650 Chenxu Hu, Qiao Tian, Tingle Li, Wang Yuping, Yuxuan Wang, and Hang Zhao. Neural dubber: 651 Dubbing for videos according to scripts. Advances in neural information processing systems, 34: 652 16582–16595, 2021. 653 654 Po-Yao Huang, Vasu Sharma, Hu Xu, Chaitanya Ryali, Haoqi Fan, Yanghao Li, Shang-Wen Li, 655 Gargi Ghosh, Jitendra Malik, and Christoph Feichtenhofer. Mavil: Masked audio-video learners, 656 2023a. 657 Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin 658 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced 659 diffusion models. In International Conference on Machine Learning (ICML), 2023b. 660 661 Vladimir Iashin and Esa Rahtu. Taming visually guided sound generation. In The British Machine 662 Vision Conference (BMVC), 2021. 663 Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: 664 A reference-free metric for evaluating music enhancement algorithms. In INTERSPEECH, pp. 665 2350-2354, 2019. 666 Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating 667 captions for audios in the wild. In Proceedings of the 2019 Conference of the North American 668 Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol-669 ume 1 (Long and Short Papers), pp. 119–132, 2019. 670 671 Tae Kyun Kim. T test as a parametric statistic. Korean journal of anesthesiology, 68(6):540-546, 672 2015. 673 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint 674 arXiv:1312.6114, 2013. 675 676 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete 677 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In Proceed-678 ings of the IEEE/CVF International Conference on Computer Vision, pp. 4015–4026, 2023. 679 A Sophia Koepke, Olivia Wiles, Yael Moses, and Andrew Zisserman. Sight to sound: An end-to-end 680 approach for visual piano transcription. In ICASSP 2020-2020 IEEE International Conference on 681 Acoustics, Speech and Signal Processing (ICASSP), pp. 1838–1842. IEEE, 2020. 682 Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for 683 efficient and high fidelity speech synthesis. Advances in Neural Information Processing Systems, 684 33:17022–17033, 2020a. 685 686 Qiuqiang Kong, Yong Xu, Turab Iqbal, Yin Cao, Wenwu Wang, and Mark D Plumbley. Acoustic 687 scene generation with conditional samplernn. In ICASSP 2019-2019 IEEE International Confer-688 ence on Acoustics, Speech and Signal Processing (ICASSP), pp. 925–929. IEEE, 2019. 689 Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. Panns: 690 Large-scale pretrained audio neural networks for audio pattern recognition. IEEE/ACM Transac-691 tions on Audio, Speech, and Language Processing, 28:2880-2894, 2020b. 692 693 Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models 694 from self-supervised synchronization. In Proceedings of the Advances in Neural Information Processing Systems, 2018. 696 Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi 697 Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. In Inter-698 national Conference on Learning Representations (ICLR), 2023. 699 Seung Hyun Lee, Sieun Kim, Innfarn Yoo, Feng Yang, Donghyeon Cho, Youngseo Kim, Huiwen 700 Chang, Jinkyu Kim, and Sangpil Kim. Soundini: Sound-guided diffusion for natural video edit-701 ing. arXiv preprint arXiv:2304.06818, 2023.

702 703 704	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre- training for unified vision-language understanding and generation. In <i>International conference on</i> <i>machine learning</i> , pp. 12888–12900. PMLR, 2022a.
705 706 707	Tingle Li, Qingjian Lin, Yuanyuan Bao, and Ming Li. Atss-net: Target speaker separation via attention-based neural network. In <i>Interspeech</i> , pp. 1411–1415, 2020.
708 709	Tingle Li, Yichen Liu, Andrew Owens, and Hang Zhao. Learning visual styles from audio-visual associations. In <i>European Conference on Computer Vision</i> , pp. 235–252. Springer, 2022b.
710 711 712 713	Tingle Li, Renhao Wang, Po-Yao Huang, Andrew Owens, and Gopala Anumanchipalli. Self- supervised audio-visual soundscape stylization. In <i>Proceedings of the European Conference on</i> <i>Computer Vision</i> , 2024.
714 715 716	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. In <i>International Conference on Machine Learning (ICML)</i> , 2023.
717 718 719 720	Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 2024.
721 722 723 724	Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. <i>Advances in neural information processing systems</i> , 33:11525–11538, 2020.
725 726	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. <i>arXiv preprint arXiv:1711.05101</i> , 2017.
727 728 729	Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. <i>arXiv preprint arXiv:2306.17203</i> , 2023.
730 731	Josh H McDermott and Eero P Simoncelli. Sound texture perception via statistics of the auditory periphery: evidence from sound synthesis. <i>Neuron</i> , 71(5):926–940, 2011.
732 733	Mary L McHugh. Interrater reliability: the kappa statistic. <i>Biochemia medica</i> , 22(3):276–282, 2012.
734 735	Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In <i>European Conference</i> on <i>Computer Vision</i> , pp. 218–234. Springer, 2022.
736 737	Pedro Morgado, Nuno Vasconcelos, Timothy Langlois, and Oliver Wang. Self-supervised generation of spatial audio for 360 video. In <i>Advances in Neural Information Processing Systems</i> , 2018.
739 740 741	Pedro Morgado, Nuno Vasconcelos, and Ishan Misra. Audio-visual instance discrimination with cross-modal agreement. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 12475–12486, 2021.
742 743	Fionn Murtagh. Multilayer perceptrons for classification and regression. <i>Neurocomputing</i> , 2(5-6): 183–197, 1991.
744 745 746	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predic- tive coding. <i>arXiv preprint arXiv:1807.03748</i> , 2018.
747 748 749	Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 631–648, 2018.
750 751 752 753	Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 2405–2413, 2016.
754 755	Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16</i> , pp. 319–345. Springer, 2020.

756 Mandela Patrick, Po-Yao Huang, Ishan Misra, Florian Metze, Andrea Vedaldi, Yuki M Asano, and João F Henriques. Space-time crop & attend: Improving cross-modal video representation learn-758 ing. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10560– 759 10572, 2021. 760 Bryan C Pijanowski, Luis J Villanueva-Rivera, Sarah L Dumyahn, Almo Farina, Bernie L Krause, 761 Brian M Napoletano, Stuart H Gage, and Nadia Pieretti. Soundscape ecology: the science of 762 sound in the landscape. BioScience, 61(3):203-216, 2011. 763 764 KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. Learning individual 765 speaking styles for accurate lip to speech synthesis. In Proceedings of the IEEE/CVF Conference 766 on Computer Vision and Pattern Recognition, pp. 13796–13805, 2020. 767 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 768 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 769 models from natural language supervision. In International conference on machine learning, pp. 770 8748-8763. PMLR, 2021. 771 772 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-773 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-774 ence on computer vision and pattern recognition, pp. 10684–10695, 2022. 775 Andrew Rouditchenko, Hang Zhao, Chuang Gan, Josh McDermott, and Antonio Torralba. Self-776 supervised audio-visual co-segmentation. In ICASSP 2019-2019 IEEE International Conference 777 on Acoustics, Speech and Signal Processing (ICASSP), pp. 2357–2361. IEEE, 2019. 778 779 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 780 Improved techniques for training gans. Advances in neural information processing systems, 29, 2016. 781 782 Roy Sheffer and Yossi Adi. I hear your true colors: Image guided audio generation. In ICASSP 2023-783 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 784 1-5. IEEE, 2023. 785 786 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv 787 preprint arXiv:2010.02502, 2020. 788 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-789 hanced transformer with rotary position embedding. Neurocomputing, 568:127063, 2024. 790 791 Kun Su, Xiulong Liu, and Eli Shlizerman. How does it sound? Advances in Neural Information 792 Processing Systems, 34:29258–29273, 2021. 793 Kim Sung-Bin, Arda Senocak, Hyunwoo Ha, Andrew Owens, and Tae-Hyun Oh. Sound to visual 794 scene generation by audio-to-visual latent alignment. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6430-6440, 2023. 796 797 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-798 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-799 tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023. 800 Kees Van Den Doel, Paul G Kry, and Dinesh K Pai. Foleyautomatic: physically-based sound ef-801 fects for interactive simulation and animation. In Proceedings of the 28th annual conference on 802 *Computer graphics and interactive techniques*, pp. 537–544, 2001. 803 804 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 805 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa-806 tion processing systems, 30, 2017. 807 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head 808 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. arXiv preprint 809 arXiv:1905.09418, 2019.

810 811 812 813	Ho-Hsiang Wu, Oriol Nieto, Juan Pablo Bello, and Justin Salamon. Audio-text models do not yet leverage natural language. In <i>ICASSP 2023-2023 IEEE International Conference on Acoustics</i> , <i>Speech and Signal Processing (ICASSP)</i> , pp. 1–5. IEEE, 2023.
814 815	Jinlong Xue, Yayue Deng, Yingming Gao, and Ya Li. Auffusion: Leveraging the power of diffusion and large language models for text-to-audio generation. <i>arXiv preprint arXiv:2401.01044</i> , 2024.
816 817 818	Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 2023.
819 820 821 822	Karren Yang, Bryan Russell, and Justin Salamon. Telling left from right: Learning spatial correspondence of sight and sound. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9932–9941, 2020.
823 824 825	Yi Yuan, Dongya Jia, Xiaobin Zhuang, Yuanzhe Chen, Zhengxi Liu, Zhuo Chen, Yuping Wang, Yuxuan Wang, Xubo Liu, Mark D Plumbley, et al. Improving audio generation with visual enhanced caption. <i>arXiv preprint arXiv:2407.04416</i> , 2024.
826 827 828 829	Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 570–586, 2018.
830 831	Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 1735–1744, 2019.
832 833 834	Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 40(6):1452–1464, 2017.
835 836 837 838	Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L Berg. Visual to sound: Gener- ating natural sound for videos in the wild. In <i>Proceedings of the IEEE conference on computer</i> <i>vision and pattern recognition</i> , pp. 3550–3558, 2018.
839 840 841	
842 843 844	
845 846	
847 848 849	
850 851 852	
853 854	
855 856 857	
858 859	
860 861 862	
863	

864 **RESULTS VIDEO** A.1 865

866

867

868

870 871

872

873

874 875

876 877

878

883

913

Our results video provided in the supplement, as well as on the project webpage, showcases our model's ability to generate accurate sound textures based on the mask prompts. Specifically, the video demonstrates:

- Our model can compositionally generate object-specific sounds within complex scenes.
- Despite being trained on the Sound-VECaps dataset (Yuan et al., 2024), our model can be successfully applied to out-of-domain visual scenes, including those from the Places dataset (Zhou et al., 2017), the Greatest Hits dataset (Owens et al., 2016), and even random web images.
- Our model can capture variations in visual textures to generate corresponding audio.

A.2 DATASET REFINEMENT

879 We use the Sound-VECaps dataset (Yuan et al., 2024), derived from AudioSet (Gemmeke et al., 880 2017), as the primary source for this task. The original dataset comprises 4,616 hours of video clips, each paired with corresponding labels and captions. To adapt this dataset for our use, we apply the 882 following refinement steps.

884 Audio-visual matching. To ensure strong 885 correspondence between audio and visual 886 inputs, we train an audio-visual matching model (Figure 8), which consists of a 6-887 layer non-causal transformer with a rotary positional embedding mechanism (Su et al., 889 2024). Visual embeddings are extracted using 890 the ViT-B/16 Transformer module from CLIP 891 (Radford et al., 2021), while audio embed-892 dings are generated using the BEATs model 893 (Chen et al., 2022b). Both embeddings are 894 then passed through a 3-layer MLP to match a 895 768-dimensional space. The model is trained in a self-supervised manner (Owens & Efros, 896 2018; Korbar et al., 2018), treating audio-897 visual pairs from the same temporal instance 898 as matches and those from different videos as 899 mismatches, which allows the model to learn 900 audio-visual correspondences without human annotations. 901

902 For training efficiency, the videos are standardized to 8 frames per second, with each 903 frame resized to 224x224 pixels. During the 904 evaluation, our model achieves an accuracy of 905 91% for matching scenarios and 85% for non-906 matching scenarios on a set of 100 matched 907 and 100 mismatched samples, indicating its 908 effectiveness in capturing audio-visual align-909 ment. We use this model to score each clip 910 in the Sound-VECaps dataset, with results 911 shown in Figure 7. A threshold of 0.6 is then 912 applied to filter the dataset.



Figure 7: Distribution of matching scores. We present the scores for audio-visual pairs in the Sound-VECaps dataset.



Figure 8: Architecture of the audio-visual matching model. We train a model to quantify the correspondence between a video and its corresponding soundtrack.

914 **Caption rephrasing.** To ensure captions focus exclusively on visible sounding objects, we utilize 915 Llama (Touvron et al., 2023) with a tailored prompt (Figure 9). Given the video and audio captions, our prompt instructs the model to generate a single sentence highlighting the common features 916 between the audio and visual content. The prompt emphasizes including only events present in both 917 modalities, while excluding modality-specific details such as overly specific visual features. The

18	Role-System:
19	You are a helpful assistant for identifying audio-visual events and generating sentences. Your task is to identify the overlapping or
0	common features between a 10-second audio and the corresponding visual description, and help the user to generate a single sentence
,	of caption that represents this intersection.
	The caption feature is a sentence generated by an audio-caption model: {enclap_caption} .
,	The label feature is several audio events that happened in the audio: {audio_label} .
	Lastly, the user is given several sentences which are the image description of the scene for each second, connected by "and then".
	Please identify all the audio events and visual elements based on all three features and try to conclude in one single sentence to describe
	this scene with the shared audio-visual events or actions that present sound and sight together.
	Please emphasize time features to present the order of each event, such as "and then", "followed by", "after" for order; "and", "while" etc.,
	for parallel events.
	Intersection Focus:
	 Based on the first capiton reactine, you might need to change or acter any wrong additio event, improve the sentence with more reactines, such as the workbursthe amotion of any penalty the description of the acrond the acrond the acrond the sentence with more reactines.
	such as the weather, the entrough of any people, the description of the data and so on.
	 Neep only the reduces that are common between the adult and visual descriptions. If an event of element is mendored in both the audio and the visual description include it in the final capition.
	 Omit any feature or detail that is present in only one modality. This includes removing overly specific visual details, such as the color
	shape any text or label, name and what pende are writing and so on that do not align with the audio description and vice yersa
	Please ensure that the final caption accurately reflects the common elements of the audio-visual scene, maintaining the order of
	occurrence, and capturing the shared background, foreground, and context.
	Role-User:
	The descriptions of the frames are: {frame_caption}
	Figure 0: Prompt for I long . We extract common features between the audio and visual contion using I lama
	Figure 9. I found to Liama. We extract common features between the autor and visual capiton using Liama
	ensuring the resulting caption locuses on events present in both modalities while avoiding overly specific de-

937 938 model is guided to capture the order and parallel occurrence of events using temporal markers like 939

"and then," "followed by," and "while." This process enhances the consistency between audio and visual descriptions.

Audio filtering. We filter out clips containing human vocalizations (e.g., singing, talking), 942 voiceovers, and music using a sound event detection model (Kong et al., 2020b) and the meta-943 data from AudioSet. This step ensures that the remaining audio data largely consists of ambient and 944 context-specific sounds that are more likely to align with the visual content. 945

946 After applying these refinement steps, the resulting data is reduced to 748 hours of video clips that 947 exhibit high audio-visual correspondence.

948 949

950

958

935

936

940

941

tails.

ADDITIONAL EVALUATION DETAILS A.3

951 ACC. We use the PANNs model (Kong et al., 2020b) to compute ACC for each audio clip, lever-952 aging annotations provided by AudioSet. First, we process each audio clip through the pre-trained 953 PANNs model to obtain the logit values for all possible sound event classes. Using the AudioSet 954 annotations, we then sample the logits corresponding to the annotated labels for each clip. Since 955 these logits are the softmax outputs, they represent the model's confidence for each event, allowing us to interpret them as accuracy scores for the labeled events. We then compute the mean of these 956 sampled logits across all clips in the dataset to obtain the final ACC score. 957

FAD, KL, and IS. We measure FAD, KL, and IS using the AudioLDM-Eval toolbox¹. The refer-959 ence and generated audio files are organized into separate folders, and the toolbox is run in paired 960 mode. 961

962 AVC. We measure AVC using a two-stream network Arandjelovic & Zisserman (2017). One 963 stream extracts audio features, while the other extracts visual features. We use OpenL3 Cramer 964 et al. (2019) to obtain these features and compute the cosine similarity for each image-audio pair. 965 Specifically, we employ the "env" content type model with a 512-dimensional linear spectrogram 966 representation. 967

968 **Human evaluation.** We conducted a human evaluation to assess the quality and relevance of the 969 generated audio using Amazon Mechanical Turk. The interface for this study is shown in Figure 10. 970 Each participant was presented with an input image and the corresponding generated audio, then 971

¹https://github.com/haoheliu/audioldm_eval



Figure 10: **Human evaluation interface.** We show the interface used for the subjective evaluation of generated audio samples. Participants are presented with input text, an image, and a corresponding audio sample, and are instructed to rate the audio on four criteria. All ratings must be completed before advancing to the next sample.

Scale	ACC (\uparrow)	FAD (\downarrow)	$\mathrm{KL}\left(\downarrow\right)$	IS (†)	AVC (†)
$\lambda = 1.0$	0.413	2.021	0.914	1.336	0.674
$\lambda = 1.5$	0.657	1.558	0.762	1.617	0.751
$\lambda = 2.0$	0.859	1.271	0.517	2.102	0.891
$\lambda = 2.5$	0.807	1.440	0.589	2.012	0.853
$\lambda = 3.0$	0.796	1.482	0.576	2.023	0.841

Table 3: Quantitative results under different CFG scales.

rated each sample on a scale from 1 to 5 based on the following criteria: (i) Overall Quality (OVL),
assessing the general audio quality; (ii) Relevance to Input Text (RET), measuring the alignment
of the audio with the associated text description; (iii) Relevance to Input Image (REI), evaluating
how well the audio corresponds to the visual content; and (iv) Relevance to Selected Object (REO),
focusing on the alignment of the audio with a specific object in the image.

We randomly selected 100 samples for evaluation, each rated by 50 unique participants to ensure reliability. The samples included both holistic and object-specific audio. To control for random responses, we incorporated a set of noise-only samples. Consistently low scores for these control samples confirmed the reliability of participants. Additionally, we ensured that each participant spent at least 90 seconds evaluating each sample to guarantee thoughtful assessment.

To further validate our results, we computed the inter-rater reliability using Cohen's kappa (McHugh, 2012), which indicated a substantial agreement among raters ($\kappa = 0.78$). Furthermore, we conducted a statistical significance test (paired t-test) (Kim, 2015) between our model and baselines for each criterion, confirming that the improvements reported are statistically significant (p < 0.01). The final scores presented in the main paper are the mean ratings across all participants.

1017

1019

990

991

1000 1001

1018 A.4 ADDITIONAL RESULTS

Different CFG scales. We evaluate our model's performance across CFG scales ranging from 1.0 to 3.0. As shown in Table 3, there is a consistent improvement in metrics as λ increases from 1.0 to 2.0, reaching peak performance at $\lambda = 2.0$. However, further increasing λ beyond 2.0 results in a gradual decline across most metrics.

- 1024
- **Different thresholds of audio-visual matching.** We test our model's performance across different audio-visual matching thresholds, varying from 0.4 to 0.8 (Figure 7). The same held-out test set is

		Threshold	ACC (↑)	FAD (↓)	KL (↓)	IS (†)	AVC (↑)	-
		0.4	0.521	1 874	0.888	1 432	0.696	_
		0.5	0.743	1.536	0.691	1.625	0.774	
		0.6	0.859	1.271	0.517	2.102	0.891	
		0.7	0.845	1.387	0.612	1.987	0.882	
		0.8	0.812	1.501	0.664	2.005	0.879	-
	Table	• 4· Quantit	ative result	s under diff	erent audi	o-visual i	matching s	cores
	Tuble	, i. Quantit	utive result	s under um	erent addi	o visuuri	indicining 5	0103.
					771 (1)	TG (4)		
		Method	ACC (↑)	FAD (↓)	$\mathbf{KL}(\downarrow)$	IS (↑)	AVC (↑)	
		w/o PE w/ PE	0.787 0.859	1.493 1.271	0.674 0.517	1.913 2.102	0.779 0.891	
	Table 5: 0	Comparisor	n of model	performanc	e with and	l without	positional	encoding.
used	to assess the m	etrics, wit	h results	presented	in Table	4. We e	empiricall	y find that the model
achie	ves optimal perf	formance a	at a thresh	old of 0.6.				
Effec	t of positional	encoding.	We ass	ess the im	pact of po	ositional	encoding	g (PE) on our model's
perfo	rmance. As sho	wn in Tab	le 5, remo	ving posit	ional enc	oding le	ads to a si	ignificant degradation
acros	s all metrics, his	ghlighting	its import	tance in the	e model's	overall	performa	nce.
		0					-	
• ~		Tree						
A.5	PROOF OF	THEORI	EM 3.1					
				_				
Proof	E For notation s	simplicity,	let $u_q \in \Delta$	Δ^P denote	the softm	nax atten	tion weig	ht computed on query
q sucl	h that $u_{al} = \overline{\Sigma}$	$\exp(\langle \mathcal{E}_v(\boldsymbol{t}_q) \rangle)$	$(\mathcal{E}_t(i_{q,l}))$	$\frac{1}{2}$. We f	first state	the follo	owing lem	ima.
1	\sum	$\mathcal{E}_{k=1}^{I} \exp(\langle \mathcal{E}_{v} \rangle)$	$(\boldsymbol{t}_q), \mathcal{E}_t(\boldsymbol{i}_{q,k})$	$_{z})\rangle _{\Sigma})$			0	
Low	no A 5 1 II J	on the area	a aar diti -	na in The		wa kari -		
Lemi	na A.5.1. Unde	er ine sam	e conditio	ns in Theo	orem 3.1,	we nave		
			$\mathbb{E}_{a}[u]$	$ p_{\alpha} _{\ell_{1}}$	$<\sqrt{2\epsilon_{cc}}$	ntrast		
			$-q \ln \alpha q$	[PY 11]	v = ~co	mast		
Proot	For notation	simplicity.	let $u_{\sigma} \in$	Δ^P denot	e the atte	ention m	ask comn	buted on query a such
that	$\exp(\langle \mathcal{E}_i \rangle)$	$v(t_q), \mathcal{E}_t(i_{q,l})$	$ \rangle_{\Sigma})$ N	Intice that			Joinp	
mat U	$p_{q,l} = \frac{1}{\sum_{k=1}^{P} \exp(k)}$	$(\langle \overline{\mathcal{E}_v(\boldsymbol{t}_q)}, \mathcal{E}_t(\boldsymbol{t}_q), \mathcal{E}_t(\boldsymbol{t}_q))$	$\overline{m{i}_{q,k}) angle_{\Sigma}}$. N	where that				
		г				-	1	
	с II.		ex	$\operatorname{p}\left(\langle \mathcal{E}_v(\boldsymbol{t}_q) \right)$	$,\mathcal{E}_t(oldsymbol{i}_{q,d})$	$\rangle_{\Sigma})$	न्त	[]
	$\epsilon_{\text{contrast}} = \mathbb{E}_{\epsilon}$	$q, d \sim p_q \mid -$	$\frac{\log \overline{\nabla^P}}{\nabla^P}$	ovn (/s	$(+) \epsilon (2)$	<u>,)_)</u>	$ - \mathbb{E}_{q,d\sim}$	$p_q \left[-\log p_{q,d}\right]$
		L	$\sum_{k=1}^{k=1}$	$\lim_{v \to v} \left(\langle \mathcal{E}_v \rangle \right)$	$(\boldsymbol{\iota}_q), \boldsymbol{\mathcal{E}}_t(\boldsymbol{\imath})$	q,k) Σ)		
	п э	Γ,	$p_{q,d}$					
	$=\mathbb{E}_{0}$	$q, d \sim p_q \left \log \right $	$\frac{1}{u_{a,d}}$					
	יתו	ן ה ([a, a]					
	$=\mathbb{E}_{0}$	$_q [D_{\mathrm{KL}}(p_q$	$[d, u_{q,d})]$					
where	Dur denotes t	the KL die	tance Ry	Pinsker's	inequality	v and Co	auchy_Sch	warz inequality
where	$\mathcal{D}_{\mathrm{KL}}$ denotes (uie KL uis	ance. by	I IIISKEI S	incquailt	y and Ca	ucity-Sell	iwarz mequanty,
		ϵ	$_{\rm contrast} =$	$\mathbb{E}_{a}\left[D_{\mathrm{KL}}\right]$	$p_{a,d}, u_{a,d}$	()]		
			001101 0.00	1 -	± 4,∞) ~4,u			
			\geq	$\frac{1}{2} \cdot \mathbb{E}_q \left[\ p \right]$	$p_{q,d} - u_q$	$_{,d}\ _{\ell_1}^2$		
				2 - ⁻	-			
			\geq	$\frac{1}{2} \cdot (\mathbb{E}_q [\parallel])$	$p_{q,d} - u_d$	$_{q,d}\ _{\ell_1}])^2$		
				2				

It follows that

 $\mathbb{E}_q[\|u_q - p_q\|_{\ell_1}] \le \sqrt{2\epsilon_{\text{contrast}}}.$

Returning to the proof of Theorem 3.1, let $s_q := f(a_q) = f(u_q V)$ denote the audio output on query q by the trained model. We decompose $\operatorname{err}_{\text{test}}$ by err_{test} $=\underbrace{\mathbb{E}_q[v(f^*(p_q\boldsymbol{V}^*),\boldsymbol{i}_q,p_q)]}_{A} - \underbrace{\mathbb{E}_q[v(f^*(u_q\boldsymbol{V}^*),\boldsymbol{i}_q,p_q)]}_{B} + \underbrace{\mathbb{E}_q[v(f^*(u_q\boldsymbol{V}^*),\boldsymbol{i}_q,p_q)]}_{B} - \underbrace{\mathbb{E}_q[v(f^*(a_q),\boldsymbol{i}_q,p_q)]}_{B}$ $+\underbrace{\mathbb{E}_{q}[v(f^{*}(a_{q}), \boldsymbol{i}_{q}, p_{q})] - \mathbb{E}_{q}[v(f(a_{q}), \boldsymbol{i}_{q}, p_{q})]}_{C} + \underbrace{\mathbb{E}_{q}[v(f(a_{q}), \boldsymbol{i}_{q}, p_{q})] - \mathbb{E}_{q}[v(f(a_{q}), \boldsymbol{i}_{q}, m_{q})]}_{D}}_{D}$ + $\underbrace{\mathbb{E}_q[v(f(a_q), i_q, m_q)] - \mathbb{E}_q[v(f(m_q V), i_q, m_q)]}_{\Gamma}$. By Lemma A.5.1 and $\|V^*\|_{\infty} \leq B_v$, we have $A < \mathbb{E}_{q}[L_{v} \cdot L_{f} \cdot B_{v} \cdot \|u_{q} - p_{q}\|_{\ell_{1}}]$ $\leq L_v \cdot L_f \cdot B_v \cdot \sqrt{2\epsilon_{\text{contrast}}}.$ Since $\|V^* - V\|_{\infty} \le \epsilon_v$ and $\|u_q\|_1 = 1$, we have $B = \mathbb{E}_{q}[v(f^{*}(u_{q}\boldsymbol{V}^{*}), \boldsymbol{i}_{q}, p_{q})] - \mathbb{E}_{q}[v(f^{*}(u_{q}\boldsymbol{V}), \boldsymbol{i}_{q}, p_{q})]$ $\leq L_v \cdot L_f \cdot \epsilon_V.$ By definition, $C \leq \epsilon_f$. Using the definition $\epsilon_{sam} = \mathbb{E}_q[||\boldsymbol{m}_q - p_q||_{\ell_1}]$, we have $D \leq \mathbb{E}_{q}[L_{v} \cdot \|\boldsymbol{m}_{q} - p_{q}\|_{\ell_{1}}]$ $\leq L_v \cdot \epsilon_{\rm sam}.$ and using $\|V\|_{\infty} \leq B_v$ with Lemma A.5.1, $E \leq \mathbb{E}_{q}[L_{v} \cdot L_{f} \cdot B_{v} \cdot \|\boldsymbol{m}_{q} - u_{q}\|_{\ell_{1}}]$ $\leq L_v \cdot L_f \cdot B_v \cdot (\epsilon_{\text{sam}} + \sqrt{2\epsilon_{\text{contrast}}}).$ Combining, we have $\operatorname{err}_{\operatorname{test}} \leq L_v \cdot \left(L_f \cdot \left(\epsilon_{\mathbf{V}} + B_v \cdot \left(\epsilon_{\operatorname{sam}} + 2\sqrt{2\epsilon_{\operatorname{contrast}}} \right) \right) + \epsilon_{\operatorname{sam}} \right) + \epsilon_f.$ This completes the proof.