

Hear-Your-Click: Interactive Video-to-Audio Generation via Object-aware Contrastive Audio-Visual Fine-tuning

Yingshan Liang
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
liangys23@mails.tsinghua.edu.cn

Keyu Fan
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
fky23@mails.tsinghua.edu.cn

Zhicheng Du
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
duzc24@mails.tsinghua.edu.cn

Yiran Wang
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
wangyr23@mails.tsinghua.edu.cn

Qingyang Shi
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
shiqy23@mails.tsinghua.edu.cn

Xinyu Zhang
Shenzhen International Graduate
School, Tsinghua University
Shenzhen, China
z-xy23@mails.tsinghua.edu.cn

Jiasheng Lu
Huawei Technologies Co., Ltd.
Shenzhen, China
lujiasheng2@huawei.com

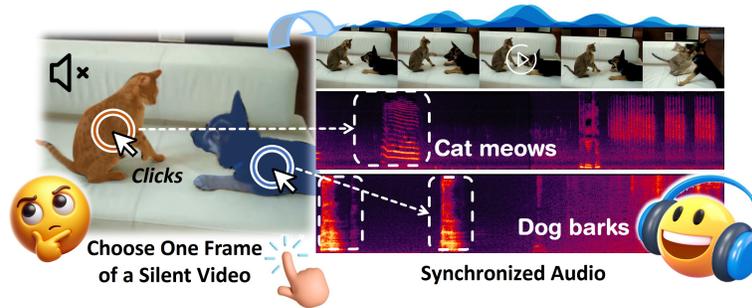


Figure 1. Overview of Hear-Your-Click.

Abstract

Video-to-audio (V2A) generation shows great potential in fields such as film production. Despite significant advances, current V2A methods, which rely on global video information, struggle with complex scenes and often fail to generate audio tailored to specific objects or regions in the videos. To address these limitations, we introduce Hear-Your-Click, an interactive V2A framework that enables users to generate sounds for specific objects in the videos by simply clicking on the frame. To achieve this, we propose Object-aware Contrastive Audio-Visual Fine-tuning (OCAV) with a Mask-guided Visual Encoder (MVE) to obtain object-level visual features aligned with corresponding audio segments. Furthermore, we tailor two data augmentation strategies: Random Video Stitching (RVS) and Mask-guided Loudness Modulation (MLM), aimed at enhancing the model’s sensitivity

to the segmented objects. To effectively measure the audio-visual correspondence, we design a new evaluation metric, the CAV score, for evaluation. Extensive experiments demonstrate that our framework offers more precise control and improved generation performance across various metrics. Project Page: <https://github.com/SynapGrid/Hear-Your-Click>

CCS Concepts: • Computing methodologies → Artificial intelligence; Computer vision representations; Natural language processing.

Keywords: Video-to-Audio Generation, Contrastive Learning, Fine-grained Control, Diffusion

1 Introduction

Video-to-Audio (V2A) generation has shown great potential for applications in diverse domains, including film production, social media, and accessibility services. Recent research

in the field has yielded significant advances, particularly with diffusion models demonstrating remarkable performance in cross-modal generation tasks [7, 19, 21, 30, 35].

Despite considerable advances, several unresolved challenges persist in the field of Video-to-Audio (V2A) generation. One of the primary challenges is the comprehension of intricate semantic information within videos. Videos often contain multiple objects and their corresponding actions, requiring V2A models to generate appropriate sound effects for each object while accounting for their interactions. This complexity poses a significant challenge, as most V2A models struggle to handle such rich semantic information effectively.

Another challenge is the lack of fine-grained control over the audio generation process. Existing methods predominantly rely on global video information, which is insufficient for customizing audio for specific objects or regions within the video according to user needs [20, 29, 44]. This limitation hinders the ability of the generated audio to meet the specific requirements of different scenarios, thereby restricting the practical applicability of V2A models in domains such as interactive media and film production. Introducing finer control mechanisms would not only enhance the accuracy of the generated audio but also improve the feasibility of V2A models in real-world applications.

To address these challenges, this work proposes an interactive V2A generation framework named Hear-Your-Click. As shown in Fig. 1, by allowing users to select specific objects within a video with a single click, Hear-Your-Click generates sounds that correspond to the selected regions, ensuring precise synchronization and alignment with the visual content. This approach provides users with greater control over the audio generation process, enabling more detailed and customizable interactions with the generated audio output. Furthermore, it addresses the suboptimal performance of global V2A approaches in handling multi-object videos by transferring control to the user, thereby enhancing the overall quality and applicability of the generated audio.

To enable interactive audio generation, we develop the VGG-AnimSeg dataset, extending VGGSound with object-specific video-mask-audio triplets created via multimodal aligners [42] and video object segmentation (VOS) [4]. We then introduce Object-aware Contrastive Audio-Visual Fine-tuning (OCAV), leveraging the Mask-guided Visual Encoder (MVE) to first extract object-level visual features and then align them with audio through contrastive learning. Two data augmentation strategies, Random Video Stitching (RVS) and Mask-Guided Loudness Modulation (MLM), are developed to improve model’s sensitivity to segmented objects. Finally, we train a Latent Diffusion Model (LDM) conditioned on the features extracted by MVE to generate audio. To better assess audio-visual correspondence, we introduce a novel metric, the CAV score, for our evaluation. Our comprehensive experimental evaluations demonstrate superior performance improvements across multiple metrics.

Our main contributions are summarized as below:

- To the best of our knowledge, Hear-Your-Click is the first interactive V2A framework that allows users to generate object-specific sounds via a simple click on the videos.
- To achieve interactive control, we propose Object-aware Contrastive Audio-Visual Fine-tuning (OCAV) alongside Mask-guided Visual Encoder, Random Video Stitching (RVS) and Mask-guided Loudness Modulation (MLM), collectively enhancing the model’s responsiveness to selected objects.
- We conduct thorough experiments to validate our approach on the VGG-AnimSeg dataset which is tailored specifically for our framework. The results across multiple metrics, including the CAV score we propose, showcase state-of-the-art performance in object-level audio-visual alignment.

2 Related Work

2.1 Video-to-Audio Generation

Significant advancements in Video-to-Audio (V2A) generation have been fueled by innovations in generative models including Generative Adversarial Networks (GANs) [13] and diffusion models [34]. Various studies propose innovative solutions to V2A challenges. SpecVQGAN [20] trains a codebook on spectrograms to generate sounds. Im2Wav [36] employs CLIP embeddings [33] and Transformers [39] to generate audio from images. Diff-Foley [29] introduces Contrastive Audio-Visual Pre-Training (CAVP) to enhance audio-visual synchronization. Seeing and Hearing [44] introduces a multimodal latent aligner to bridge existing video and audio generation models. TiVA [41] and MaskVAT [32] focus on synchronization and propose the use of audio layout and sequence-to-sequence mask generative model, respectively. Other methods, such as SonicVisionLM [43], SVA [1] and FoleyCrafter [48], use text descriptions as a mediating modality between video and audio, allowing them to generate high-quality audio using large language models (LLMs) or pre-trained text-to-audio (T2A) models, but they also face the challenge of achieving accurate synchronization [48]. The above methods focus primarily on global video information, which can lead to missing local details and difficulties in processing complex multi-object videos. And they generally lack mechanisms for fine-grained control over the V2A generation process.

2.2 Audio-Visual Alignment

Audio-Visual Pre-Training aims to obtain joint representations for improved retrieval, classification and generation. Methods such as CAV-MAE [12], AV-MAE [9], MAViL [18] and CrossMAE [14] leverage the Masked Auto-Encoder (MAE) [17] to learn cross-modal correlations, thereby enhancing the performance of joint representations in audio-visual retrieval

and classification tasks. Inspired by CLIP [33], methods like Morgado [31], Diff-Foley [29] and SCAV [38] make use of contrastive learning to align audio and video features. Other methods, such as AudioCLIP [15] and ImageBind [11], extend multimodal alignment to achieve state-of-the-art results in zero-shot cross-modal tasks. These studies demonstrate the superiority of audio-visual pre-training, and thus an increasing number of works [29, 44] have applied these techniques to V2A generation tasks.

3 Hear-Your-Click

3.1 Task Formulation

Given a T -frame video $\mathcal{V} \in \mathbb{R}^{T \times H \times W \times 3}$ and a target object \mathcal{S} , our objective is to generate the corresponding Mel spectrogram of the audio \mathcal{A} , where $\mathcal{A} \in \mathbb{R}^{T' \times N}$, with T' representing the temporal length and N denoting the mel bins. To achieve this, we first obtain binary masks \mathcal{M} to delineate \mathcal{S} , where $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_T\}$ and $\mathcal{M}_t \in \{0, 1\}^{H \times W}$ refers to the t -th frame mask. In the inference stage, the target object \mathcal{S} is specified through the user’s clicks on video frames. In the training stage, \mathcal{S} is specified through textual prompts labeled by human. Leveraging the information of \mathcal{S} , we perform prompt-guided video segmentation to obtain \mathcal{M} . Then \mathcal{M} together with corresponding \mathcal{V} and \mathcal{A} form triplets $(\mathcal{V}, \mathcal{M}, \mathcal{A})$, which are used to supervise the model in learning the mappings between the input $(\mathcal{V}, \mathcal{M})$ and the output \mathcal{A} .

Method overview. First, we tailor a new dataset named VGG-AnimSeg based on the task defined above (Sec. 3.2). Then we propose Object-aware Contrastive Audio-Visual Fine-tuning (OCAV) to align video-mask-audio triplets (Sec. 3.3), including a Mask-guided Visual Encoder (MVE) (Sec. 3.3.1) designed to extract object-level visual features, and two data augmentation strategies: Random Video Stitching (RVS) (Sec. 3.3.4) and Mask-guided Loudness Modulation (MLM) (Sec. 3.3.3). Finally, we detail the training of a latent diffusion model conditioned on the visual features obtained through OCAV (Sec. 3.4), along with the construction of the interactive inference framework (Sec. 3.5).

3.2 VGG-AnimSeg Dataset

Given that existing audio-visual datasets fail to adequately support our objective of focusing on specific objects, we have constructed a new dataset based on the VGGSound dataset [2] to address this limitation. The VGGSound dataset contains over 200K 10-second video clips spanning more than 300 categories, each accompanied by human-labeled textual descriptions. To ensure distinct vocal subjects in the videos as well as cleaner audio tracks, we initially selected all animal-related videos from VGGSound. This selection process resulted in a subset containing 68 classes of textual descriptions.

To filter out noisy samples, we use multimodal joint embeddings to select the samples where both modalities (audio and video) exhibit a strong alignment with textual descriptions. Specifically, we first leverage the Contrastive Language-Audio Pretraining (CLAP) model [6] and Contrastive Language-Image Pre-Training (CLIP) model [33] to extract audio, image, and text embeddings. For each sample, we compute the cosine similarity between the audio-text and image-text embedding pairs. Based on these scores, we select 400 training samples and 40 test samples per textual description with the highest average similarity, resulting in a dataset of approximately 30,000 samples.

To further ensure precise correspondence between audio and visual content, we employ DEVA [4] to generate video binary masks via text-prompted segmentation using the text descriptions. This approach ensures a one-to-one correspondence between the audio tracks and their corresponding visual segments.

3.3 OCAV

Audio-visual alignment bridges semantic and temporal gaps between audio and visual modalities, with Contrastive Audio-Visual Pre-training (CAVP) [29] representing a notable advancement in achieving global alignment between these modalities. However, CAVP’s focus on global synchronization overlooks fine-grained, object-level details critical for precise audio-visual understanding. To address this limitation while retaining CAVP’s strengths, we propose the Object-aware Contrastive Audio-Visual Fine-tuning (OCAV) framework.

Given triplets $(\mathcal{V}, \mathcal{M}, \mathcal{A})$, the objective of OCAV is to extract object-level visual features from video-mask pairs $(\mathcal{V}, \mathcal{M})$, where target objects are determined by \mathcal{M} , while ensuring that these visual features are aligned with the features of the audio \mathcal{A} . The overview of OCAV is illustrated in Fig. 2.

3.3.1 Mask-guided Visual Encoder (MVE). Traditional video encoders, which process the entire video input, often fail to concentrate exclusively on the target objects indicated by the masks \mathcal{M} . Therefore, we propose Mask-guided Visual Encoder (MVE) to extract object-specific visual features from video-mask pairs efficiently. The MVE architecture integrates binary masks as supplementary inputs alongside the primary video stream, thereby enhancing the model’s responsiveness to the designated objects. This dual-input structure consists of two branches: a video-encoding branch f_v and a mask-encoding branch f_m .

In the processing of the video input for f_v , we first compare feeding masked videos with original videos, and observe that masked videos lead to more stable and clearer audio as they block out background interference. This is denoted as:

$$\mathbf{x}_{mv} = \text{norm}(f_v(\mathcal{V} \odot \mathcal{M})) \quad (1)$$

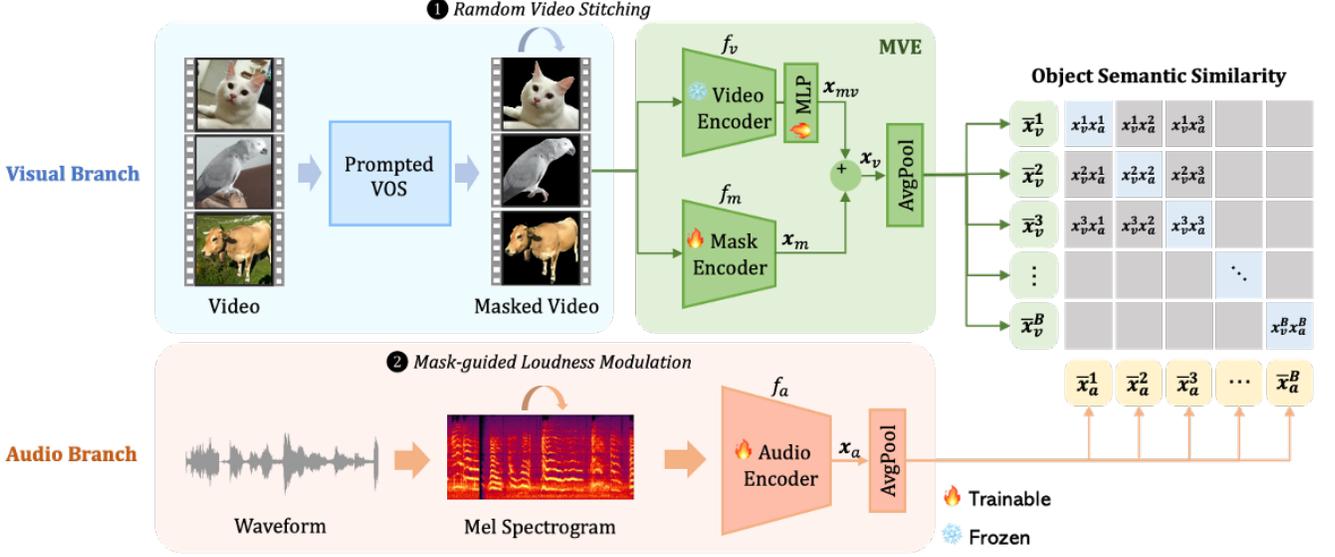


Figure 2. Overview of OCAV. Object-level visual features extracted by MVE are aligned with corresponding audio features. Training videos and audio tracks are augmented through (1) Random Video Stitching and (2) Mask-guided Loudness Modulation.

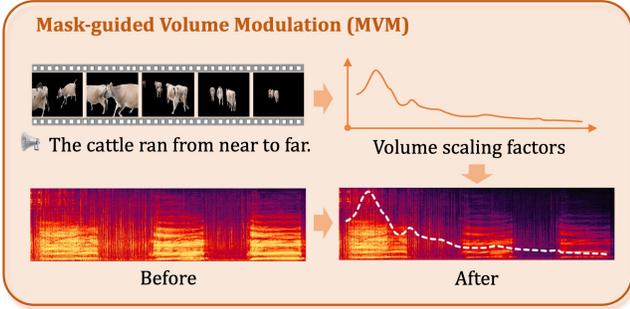


Figure 3. Enhance the correlation between object distance and volume changes.

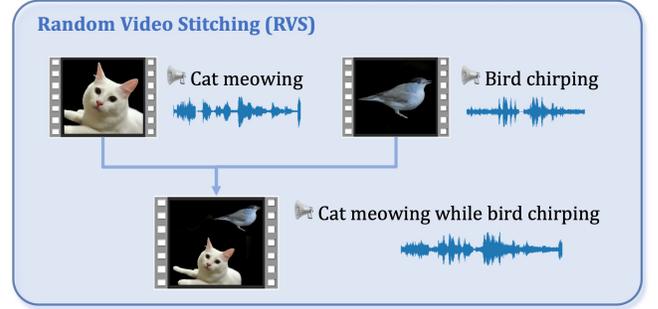


Figure 4. Improve the model's ability to handle multi-object scenes.

where $\mathbf{x}_{mv} \in \mathbb{R}^{T \times d}$, d represents the dimension of the features, and $norm(\cdot)$ represents the Euclidean norm along the second dimension to unify the feature scale.

Regarding the mask input, we design f_m with a convolutional backbone to capture the temporal characteristics of the masks. These extracted features are then fused with \mathbf{x}_{mv} to obtain the ultimate visual features \mathbf{x}_v :

$$\mathbf{x}_m = norm(f_m(\mathcal{M})) \quad (2)$$

$$\mathbf{x}_v = norm(\mathbf{x}_{mv} + \mathbf{x}_m) \quad (3)$$

where $\mathbf{x}_m, \mathbf{x}_v \in \mathbb{R}^{T \times d}$.

3.3.2 Contrastive Learning. In parallel with MVE, we adopt a convolutional encoder f_a to extract temporal audio features:

$$\mathbf{x}_a = norm(f_a(\mathcal{A})) \quad (4)$$

where $\mathbf{x}_a \in \mathbb{R}^{T \times d}$. We then compute the average of \mathbf{x}_v and \mathbf{x}_a along the time axis to obtain $\bar{\mathbf{x}}_v, \bar{\mathbf{x}}_a \in \mathbb{R}^d$.

In the training stage, f_v and f_a are initialized with their pre-trained weights from Diff-Foley [29], whereas f_m is initialized randomly. We freeze most of f_v to retain its prior knowledge, while keeping its final multi-layer perceptron (MLP) block trainable. Both f_m and f_a are kept fully trainable. For each batch, we randomly sample time-synchronized clips from $(\mathcal{V}, \mathcal{M}, \mathcal{A})$ and extract feature pairs $\{(\bar{\mathbf{x}}_v^i, \bar{\mathbf{x}}_a^i)\}_{i=1}^B$, where B is batch size. We use the following contrastive objective to supervise the model training [29, 33, 37]:

$$\mathcal{L}_{\text{contrast}} = \frac{1}{B} \sum_{i=1}^B \left\{ -\frac{1}{2} \log \frac{\exp(\phi(\bar{\mathbf{x}}_v^i, \bar{\mathbf{x}}_a^i)/\tau)}{\sum_j \exp(\phi(\bar{\mathbf{x}}_v^i, \bar{\mathbf{x}}_a^j)/\tau)} - \frac{1}{2} \log \frac{\exp(\phi(\bar{\mathbf{x}}_v^j, \bar{\mathbf{x}}_a^i)/\tau)}{\sum_j \exp(\phi(\bar{\mathbf{x}}_v^j, \bar{\mathbf{x}}_a^i)/\tau)} \right\} \quad (5)$$

where τ is the temperature parameter and ϕ denotes cosine similarity function:

$$\phi(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \quad (6)$$

3.3.3 Mask-guided Loudness Modulation (MLM). Temporal inconsistencies in the dataset, such as sounds continuing after their sources move out of the frame, create ambiguities for models. These ambiguities make it difficult for models to determine whether sounds should continue or stop once the sound sources are no longer visible. Imperfect audio-video synchronization, including off-screen sounds and background noise, complicates audio-video alignment and introduces extraneous interference to joint representation learning.

Our objective is to mitigate these issues by ensuring that the cessation of sound effects precisely corresponds to the departure of the target object from the camera view during inference. Additionally, we aim to adjust the audio loudness in accordance with the relative distances of the objects from the camera. To accomplish these objectives, we incorporate Mask-guided Loudness Modulation (MLM) into our training methodology. Through strengthening the correlation between object masks and audio loudness, MLM facilitates a more accurate learning of the mapping relationship between visual masks and auditory components.

As illustrated in Fig. 2, given a set of binary masks M , we first compute the ratio of unmasked pixels to the total number of pixels for each frame:

$$\lambda_k = \frac{\sum_{\substack{0 \leq i < W \\ 0 \leq j < H}} \mathcal{M}_k(i, j)}{H \times W}, k = 1, 2, \dots, T \quad (7)$$

Then we normalize these values to ensure the maximum value is 1:

$$\lambda'_k = \frac{\lambda_k}{\max(\lambda_1, \lambda_2, \dots, \lambda_T)}, k = 1, 2, \dots, T \quad (8)$$

Next, we apply linear interpolation to resample these values from the length T to match the length of 1-D audio signals, thereby obtaining a new set of loudness scaling factors $\Lambda = \{\lambda'_1, \lambda'_2, \dots, \lambda'_T\}$. Finally, Λ is element-wise multiplied with audio signals to modulate the loudness based on the presence and extent of the unmasked region over time.

3.3.4 Random Video Stitching (RVS). Video augmentation methods are extensively applied in various video processing tasks, including classification, action recognition and object detection [3, 5, 45, 47]. These techniques facilitate more robust generalization and improved performance of deep learning models by generating diverse training data. Moreover, studies such as LeMDA [27] and MixGen [16] offer novel insights into multimodal data augmentation. However, such methods have been underutilized in the domain of audio-visual alignment.

In this work, we propose Random Video Stitching (RVS) during OCAV to enhance the model’s ability to identify and understand individual objects. Specifically, for each video

in the dataset, another video is randomly selected, and each frame is stitched either horizontally or vertically with the corresponding frame from the original video. Simultaneously, the audio tracks from both videos are overlapped. By incorporating RVS, we aim to improve the model’s capacity to handle complex, multi-object scenes and enhance its overall performance and generalization.

Theoretically, RVS can significantly increase training samples by combining different video segments, but too many augmented samples can lead to overfitting and reduced performance on new data. Therefore, we apply RVS to the 1,126 training samples from the VGG-AnimSeg-1k dataset and integrate the augmented data into the final training set, ensuring a balance between original and augmented samples.

3.4 Latent Diffusion Model

We utilize Latent Diffusion Models (LDMs) [34] as our generative model due to their proven efficacy in cross-modal generation tasks. To mitigate the issues of failing to generate audio for the target object or generating incorrect audio, we propose augmenting the MVE visual features \mathbf{x}_v with CLIP features \mathbf{x}_v^* of the masked frames, as they provide rich, high-level semantic information that helps disambiguate between different objects. This is denoted as:

$$\mathbf{x}_c = \mathbf{x}_v + \mathbf{x}_v^* \quad (9)$$

where $\mathbf{x}_c, \mathbf{x}_v^* \in \mathbb{R}^{T \times d}$, and \mathbf{x}_v^* is extracted from the T frames of a video. Then we train our model conditioned on \mathbf{x}_c .

Given a Mel spectrogram $\mathcal{A} \in \mathbb{R}^{T' \times N}$, we first map it into a latent space through $\mathbf{z}_0 = E(\mathcal{A})$, where $\mathbf{z}_0 \in \mathbb{R}^{C \times T' \times N^*}$ and $E(\cdot)$ is a pre-trained latent encoder. In the latent space, LDM employs a forward diffusion process to gradually add noise to \mathbf{z}_0 , transforming it to a Gaussian distribution \mathcal{N} :

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I}) \quad (10)$$

$$q(\mathbf{z}_t | \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t; \sqrt{\bar{\alpha}_t} \mathbf{z}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (11)$$

where $\bar{\alpha}_t = \prod_{i=1}^t (1 - \beta_i)$, β_t is the noise schedule parameter that controls the amount of noise added at each time step. Then LDM is trained to predict the noise added at each step of the forward diffusion process:

$$\mathcal{L}_{LDM} = \mathbb{E}_{\mathbf{z}_0, t, \epsilon} \|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, \mathbf{x}_c)\|_2^2 \quad (12)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and ϵ_θ is predicted by the model. After training, LDM follows the reverse process to generate new samples. Starting from a random noise latent $\mathbf{z}_T \sim \mathcal{N}(0, \mathbf{I})$, the model iteratively denoises to reconstruct a meaningful latent:

$$p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}_c) = \mathcal{N}(\mathbf{z}_{t-1}; \mu_\theta(\mathbf{z}_t, t, \mathbf{x}_c), \sigma_t^2 \mathbf{I}) \quad (13)$$

$$p_\theta(\mathbf{z}_{0:T} | \mathbf{x}_c) = p(\mathbf{z}_T) \prod_{t=1}^T p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}_c) \quad (14)$$

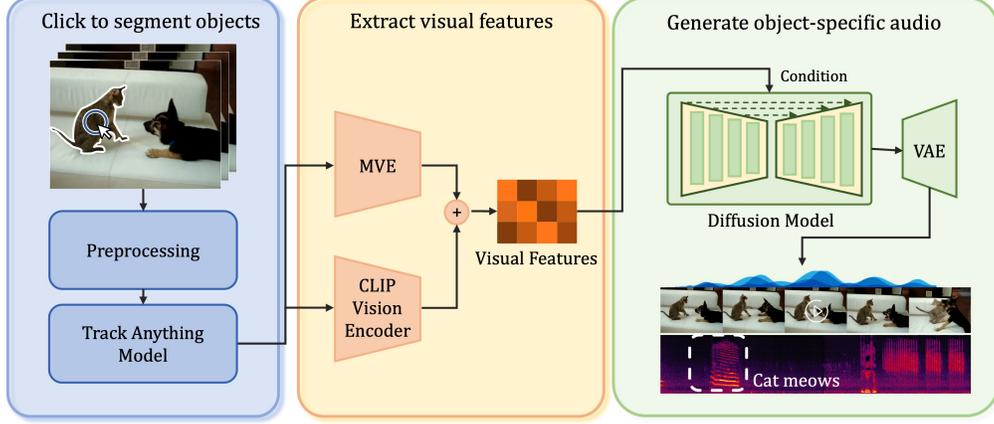


Figure 5. Overview of interactive inference process. Users can upload a silent video and select a frame to refine the mask via clicks based on Segment Anything Model (SAM). The Track Anything Model (TAM) then propagates these masks for MVE and CLIP to extract visual features. Finally, a trained LDM generates the final audio conditioned on these features.

where μ_θ and σ_t^2 are defined as follows:

$$\mu_\theta(z_t, t, \mathbf{x}_c) = \frac{1}{\sqrt{\alpha_t}} \left(z_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(z_t, t, \mathbf{x}_c) \right) \quad (15)$$

$$\sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} (1 - \alpha_t) \quad (16)$$

z_0 is then mapped back to the original data space using the decoder D : $\hat{\mathcal{A}} = D(z_0)$.

3.5 Iterative Inference

There is a growing trend in research towards interactive content generation [10, 24, 30, 40] because it engages users directly in the process, allowing for real-time feedback and customization. This makes the system more adaptable to specific user needs and enhances the overall user experience, providing greater enjoyment. Inspired by these studies, we develop an interactive interface for Hear-Your-Click, utilizing the promptable segmentation capabilities of the Segment Anything Model (SAM) [22] and the Track Anything Model (TAM) [46]. This interface allows users to upload a silent video and then select a single frame to specify target regions via clicks. Upon selection, SAM generates corresponding masks in real-time based on user clicks. Users can iteratively refine these masks to ensure high precision in delineating the target object. Once the mask is finalized, TAM propagates it throughout the entire video sequence using semi-supervised video object segmentation, generating the comprehensive video mask \mathcal{M} . Next, we extract MVE features \mathbf{x}_v and CLIP features \mathbf{x}_c^* from $(\mathcal{V}, \mathcal{M})$, forming the condition embedding \mathbf{x}_c . Finally, we sample the final audio using trained LDM conditioned on \mathbf{x}_c . Fig. 5 provides a detailed overview of the inference pipeline, from user interaction to audio generation.

4 Experiments

4.1 Experimental Setup

4.1.1 Dataset. Our models are trained and evaluated on the VGG-AnimSeg dataset, which consists of 27,200 training samples and 2,720 testing samples. When data augmentation via Random Visual Sampling (RVS) is applied, an additional 6,800 augmented training samples are incorporated, resulting in a total of 34,000 samples for model training.

4.1.2 Evaluation Metrics. For quantitative assessment, we use Frechet Distance (FD), Frechet Audio Distance (FAD), Inception Score (IS), Kullback-Leibler Divergence (KL) and Kernel Inception Distance (KID) following other audio generation studies [20, 25, 26, 29, 44]. FD, FAD, KL and KID are employed to quantify the similarity between the generated audio and the ground truth, while IS is utilized to assess the diversity and quality of the generated audio samples. To better evaluate the correspondence between video and audio elements, we introduce a new metric, the CAV score. The score leverages the model from C-MCR [42], which integrates CLIP [33] and CLAP [6] to provide audio-image contrastive representations. For each audio-video pair, we obtain per-frame image embeddings and an audio embedding, and then calculate the similarity between the average image embedding and the audio embedding. A higher CAV score indicates a better match between the objects in the generated audio and the original video.

4.1.3 Implementation Details. For the VOS setting, the IoU threshold and NMS threshold of SAM are set to 0.88 and 0.8. TAM is configured to use 15 voting frames. For data preprocessing, we resample 10-second video clips to 4 fps and resize each frame to 224×224 pixels. For audio samples, we compute Mel spectrograms with a hop size of 250 and 128 mel bins during OCAV training, and a hop size

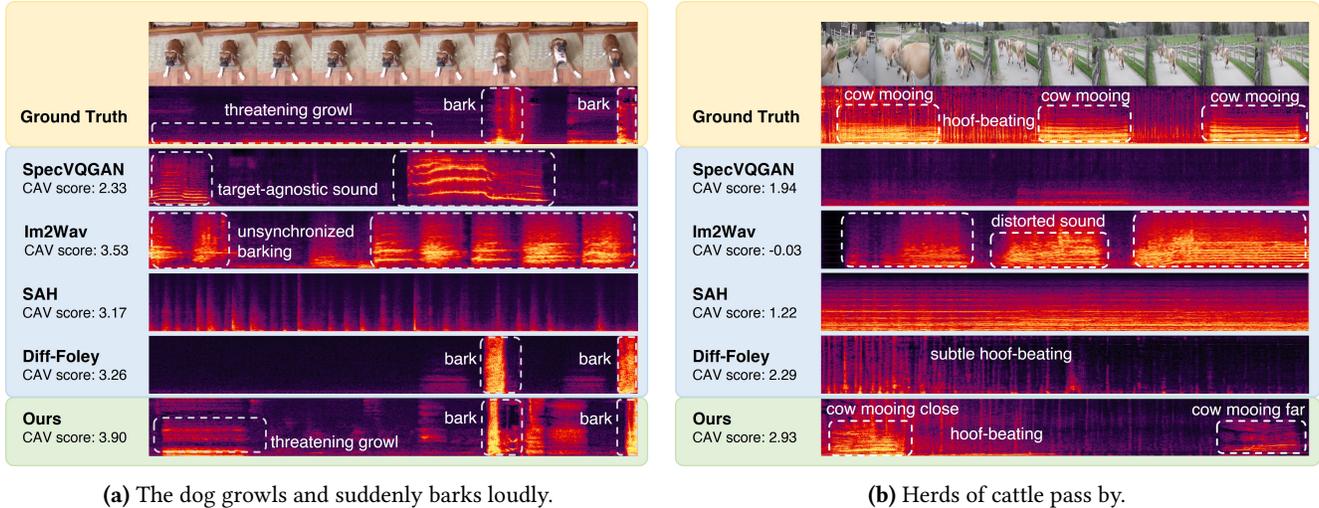


Figure 6. Visualization of the generated samples. (a) Our method successfully reproduces similar sounds with the ground truth while preserving synchronization, demonstrating great performance in capturing transient motion. (b) As the cows run away, the volume of the audio generated by our method diminishes, while preserving finer details.

of 256 and 128 mel bins during LDM training. Considering the limited size of our dataset, which is insufficient for training the model from scratch, we construct our framework based on Diff-Foley and fine-tune our model using its pre-trained weights. Specifically, the f_v and f_m in MVE utilize the SlowOnly [8] architecture, known for its proficiency in detecting temporal action changes within videos. f_a leverages an audio encoder provided by PANNs [23]. The architecture of LDM, including the latent encoder and decoder, is inherited from Stable Diffusion-V1.4 (SD-V1.4) [34]. During the OCAV training process, we randomly sample 4-second segments from the original audio and videos, denoted as $\mathcal{A} \in \mathbb{R}^{256 \times 128}$ and $\mathcal{V} \in \mathbb{R}^{16 \times 224 \times 224 \times 3}$, yielding corresponding features $\mathbf{x}_a \in \mathbb{R}^{16 \times 512}$, $\mathbf{x}_v \in \mathbb{R}^{16 \times 512}$. For LDM training, we first extract visual features through MVE from the 10-second videos, resulting in $\mathbf{x}_v \in \mathbb{R}^{40 \times 512}$. Following the Im2Wav [36] method, we extract per-frame CLIP embeddings, resulting in $\mathbf{x}_v^* \in \mathbb{R}^{40 \times 512}$, and add them to \mathbf{x}_v to form the condition embedding \mathbf{x}_c for the LDM. At the inference stage, we set the Classifier-Free Guidance Scale to 4.5, the Classifier-Guidance Scale to 50, and use the DPM-Solver [28] sampler with 50 inference steps.

4.2 Comparative Experiments

4.2.1 Qualitative Results. We conduct a comparative experiment of our method with several recently published, open-source V2A approaches, including SpecVQGAN [20], Im2Wav [36], Seeing and Hearing (SAH) [44] and Diff-Foley [29]. The qualitative evaluation focuses on three main aspects: the capacity to produce coherent audio, the relevance of the generated audio to the target objects, and the temporal alignment between audio and visual content, particularly

focusing on the movements within it. As shown in Fig. 6, our method excels in generating audio most similar to the original for given video inputs, accurately capturing nuances such as dog barks and cow moos, while maintaining synchronization with their movements. Other methods fall short, either producing visually accurate but poorly synchronized or low-quality audio (Im2Wav), or generating irrelevant sounds (SAH). SpecVQGAN, despite its high spectral detail, fails to correctly identify video objects, resulting in inaccurate audio generation. While Diff-Foley succeeds in crafting synchronized audio, it occasionally struggles to reproduce the exact sounds associated with specific objects. Overall, our approach demonstrates a superior balance between audio quality, synchronization, and object-specific sound accuracy.

Method	FD↓	FAD↓	IS↑	KL↓	KID↓	CAV↑
SpecVQGAN	75.55	5.51	3.79	3.53	0.021	1.50
Im2Wav	63.35	11.94	3.27	3.13	0.021	2.27
SAH	98.41	11.30	4.46	5.09	0.039	1.76
Diff-Foley	59.00	6.60	5.67	3.95	0.015	2.22
Ours ¹	49.48	4.04	5.20	3.01	0.012	2.55
Ours ²	48.78	5.02	4.49	2.82	0.010	2.67

Table 1. Results of quantitative comparison experiments. Our¹ represents inference conditioned on MVE features, and Our² represents inference conditioned on MVE+CLIP features.

4.2.2 Quantative Results. For each video from VGG-AnimSeg-4k test set, we generate an 8.2-second audio clip for evaluation. Tab. 1 demonstrates our method presents superior

performance compared to other methods in terms of FD, FAD, KL, and KID, indicating that our results closely resemble the ground truth. Our excellent CAV score further highlights our method’s exceptional object alignment between the generated audio and input videos.

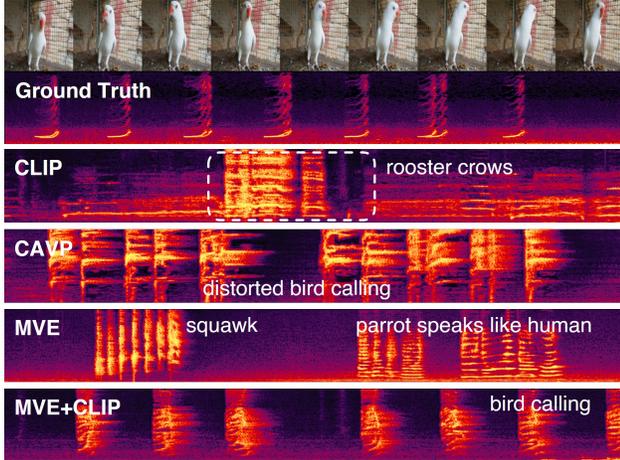


Figure 7. Comparison of different visual features. Our method (MVE and MVE+CLIP) performs better in capturing the semantic information and temporal conditions in the video, while others generate target-agnostic sounds or low-fidelity audio.

4.3 Ablation Studies

4.3.1 MVE Ablation. To validate the effectiveness of the visual features extracted by MVE, we compare them with other visual features. We fine-tune the LDM using different visual features extracted from the training set of VGG-AnimSeg-4k, and then generate 10 audio samples for each test video for evaluation. As shown in Tab. 2, our method achieves best in multiple metrics. Although CAVP features achieve the highest IS score, they perform worse in other metric. Since IS is originally designed to measure the quality of generated images, it does not guarantee the quality of the generated audio. Additionally, adding CLIP features to MVE features (referred to as MVE+CLIP) further enhances performance. This improvement is attributed to the rich and high-level semantic information provided by CLIP features, which may be ignored in convolutional backbone. Fig. 7 shows a qualitative comparison of a representative example.

4.3.2 MLM and RVS Ablation. We conduct an ablation study to validate the effectiveness of RVS and MLM. As shown in Tab. 3, MLM significantly improves the scores of FD, FAD, KL, KID, and CAV. The qualitative comparison in Fig. 8 demonstrates that the model trained with MLM generates audio with more pronounced changes in loudness corresponding to the distance of the target object. However, when RVS is added to our base model, there is a slight decrease

Visual Features	FD↓	FAD↓	IS↑	KL↓	KID↓	CAV↑
CLIP	42.78	5.73	5.39	3.12	0.015	2.55
CAVP	47.38	6.73	6.65	4.14	0.016	1.95
MVE	38.89	4.03	5.79	3.17	0.012	3.06
MVE+CLIP	35.41	4.90	5.93	2.90	0.011	2.69

Table 2. Investigate the effect of conditioning on different visual features. The results showcase the superiority of MVE features, and MVE+CLIP further enhances performance.

in performance. This may be due to the resizing process in RVS, which can distort the aspect ratio of the original images. Despite this, Fig. 9 demonstrates that RVS enables the model to handle multi-object scenes more effectively, highlighting its importance in complex scenarios.

MLM	RVS	FD↓	FAD↓	IS↑	KL↓	KID↓	CAV↑
		41.29	5.22	5.86	3.36	0.012	2.28
✓		34.67	4.81	5.53	3.07	0.010	2.47
	✓	50.19	4.87	5.56	3.62	0.018	2.30
✓	✓	40.11	4.81	5.37	3.16	0.014	2.35

Table 3. Validation of MLM and RVS.

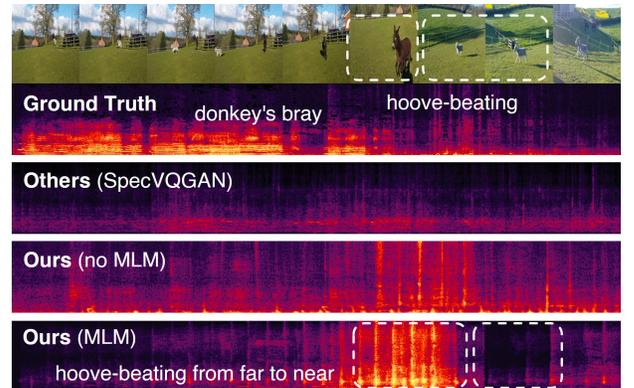


Figure 8. An ablation study of MLM. When the donkey runs closer, our results show a clearer increase in loudness, and when it moves away, the sound diminishes.

5 Conclusion

In this work, we introduce Hear-Your-Click, an interactive V2A framework. To support this framework, we develop the VGG-AnimSeg dataset and propose Object-aware Contrastive Audio-Visual Fine-tuning, which includes a Mask-Guided Visual Encoder and two data augmentation strategies. For a more precise evaluation of audio-visual alignment, we introduce the CAV score alongside traditional metrics.

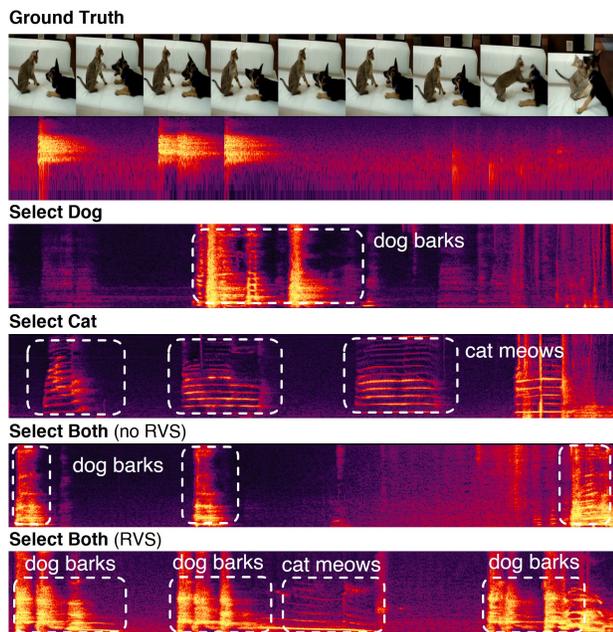


Figure 9. An ablation study of RVS. RVS enables the model to generate corresponding sounds when selecting multiple objects.

Our experimental results demonstrate superior performance across various objective indicators and qualitative assessments. Additionally, we have created an interactive V2A interface for Hear-Your-Click, a feature not found in other V2A methods. These results highlight the effectiveness of our approach in capturing object-specific local details within videos, leading to more accurate audio generation. We anticipate that our work will inspire further advancements in V2A research and facilitate its broader application.

References

- [1] Gehui Chen, Guan'an Wang, Xiaowen Huang, and Jitao Sang. 2024. Semantically consistent Video-to-Audio Generation using Multimodal Language Large Model. *arXiv preprint arXiv:2404.16305* (2024).
- [2] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 721–725.
- [3] Lihui Chen, Gemine Vivone, Zihao Nie, Jocelyn Chanussot, and Xiaomin Yang. 2023. Spatial data augmentation: Improving the generalization of neural networks for pansharpening. *IEEE Transactions on Geoscience and Remote Sensing* 61 (2023), 1–11.
- [4] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. 2023. Tracking anything with decoupled video segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1316–1326.
- [5] Jinhao Duan, Quanfu Fan, Hao Cheng, Xiaoshuang Shi, and Kaidi Xu. 2023. Improve video representation with temporal adversarial augmentation. *arXiv preprint arXiv:2304.14601* (2023).
- [6] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. 2023. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [7] Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. 2024. Fast timing-conditioned latent audio diffusion. *arXiv preprint arXiv:2402.04825* (2024).
- [8] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. 2019. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*. 6202–6211.
- [9] Mariana-Iuliana Georgescu, Eduardo Fonseca, Radu Tudor Ionescu, Mario Lucic, Cordelia Schmid, and Anurag Arnab. 2023. Audiovisual masked autoencoders. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 16144–16154.
- [10] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. 2023. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373* (2023).
- [11] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15180–15190.
- [12] Yuan Gong, Andrew Rouditchenko, Alexander H Liu, David Harwath, Leonid Karlinsky, Hilde Kuehne, and James Glass. 2022. Contrastive audio-visual masked autoencoder. *arXiv preprint arXiv:2210.07839* (2022).
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [14] Yuxin Guo, Siyang Sun, Shuailei Ma, Kecheng Zheng, Xiaoyi Bao, Shijie Ma, Wei Zou, and Yun Zheng. 2024. CrossMAE: Cross-Modality Masked Autoencoders for Region-Aware Audio-Visual Pre-Training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 26721–26731.
- [15] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. 2022. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 976–980.
- [16] Xiaoshuai Hao, Yi Zhu, Srikanth Appalaraju, Aston Zhang, Wanqian Zhang, Bo Li, and Mu Li. 2023. Mixgen: A new multi-modal data augmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 379–389.
- [17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16000–16009.
- [18] Po-Yao Huang, Vasu Sharma, Hu Xu, Chaitanya Ryal, Yanghao Li, Shang-Wen Li, Gargi Ghosh, Jitendra Malik, Christoph Feichtenhofer, et al. 2024. Mavil: Masked audio-video learners. *Advances in Neural Information Processing Systems* 36 (2024).
- [19] Qingqing Huang, Daniel S Park, Tao Wang, Timo I Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Frank, et al. 2023. Noise2music: Text-conditioned music generation with diffusion models. *arXiv preprint arXiv:2302.03917* (2023).
- [20] Vladimir Iashin and Esa Rahtu. 2021. Taming visually guided sound generation. *arXiv preprint arXiv:2110.08791* (2021).
- [21] Xuan Ju, Xian Liu, Xintao Wang, Yuxuan Bian, Ying Shan, and Qiang Xu. 2024. Brushnet: A plug-and-play image inpainting model with decomposed dual-branch diffusion. *arXiv preprint arXiv:2403.06976* (2024).
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4015–4026.
- [23] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. 2020. Panns: Large-scale pretrained audio neural

- networks for audio pattern recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28 (2020), 2880–2894.
- [24] Zhengqi Li, Richard Tucker, Noah Snively, and Aleksander Holynski. 2024. Generative image dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 24142–24153.
- [25] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. 2023. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503* (2023).
- [26] Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D. Plumbley. 2024. AudioLDM 2: Learning Holistic Audio Generation With Self-Supervised Pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 32 (2024), 2871–2883. <https://doi.org/10.1109/TASLP.2024.3399607>
- [27] Zichang Liu, Zhiqiang Tang, Xingjian Shi, Aston Zhang, Mu Li, Anshumali Shrivastava, and Andrew Gordon Wilson. 2022. Learning multimodal data augmentation in feature space. *arXiv preprint arXiv:2212.14453* (2022).
- [28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095* (2022).
- [29] Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. 2024. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. *Advances in Neural Information Processing Systems* 36 (2024).
- [30] Yue Ma, Yingqing He, Hongfa Wang, Andong Wang, Chenyang Qi, Chengfei Cai, Xiu Li, Zhifeng Li, Heung-Yeung Shum, Wei Liu, et al. 2024. Follow-your-click: Open-domain regional image animation via short prompts. *arXiv preprint arXiv:2403.08268* (2024).
- [31] Pedro Morgado, Ishan Misra, and Nuno Vasconcelos. 2021. Robust audio-visual instance discrimination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12934–12945.
- [32] Santiago Pascual, Chungshin Yeh, Ioannis Tsiamas, and Joan Serra. 2024. Masked Generative Video-to-Audio Transformers with Enhanced Synchronicity. *arXiv preprint arXiv:2407.10387* (2024).
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [35] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 22500–22510.
- [36] Roy Sheffer and Yossi Adi. 2023. I hear your true colors: Image guided audio generation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [37] Yuhan Shen, Huiyu Wang, Xitong Yang, Matt Feiszli, Ehsan Elhamifar, Lorenzo Torresani, and Effrosyni Mavroudi. 2024. Learning to Segment Referred Objects from Narrated Egocentric Videos. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2024), 14510–14520. <https://api.semanticscholar.org/CorpusID:272724136>
- [38] Ioannis Tsiamas, Santiago Pascual, Chungshin Yeh, and Joan Serra. 2024. Sequential contrastive audio-visual learning. *arXiv preprint arXiv:2407.05782* (2024).
- [39] A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems* (2017).
- [40] Xudong Wang, Trevor Darrell, Sai Saketh Rambhatla, Rohit Girdhar, and Ishan Misra. 2024. Instancediffusion: Instance-level control for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6232–6242.
- [41] Xihua Wang, Yuyue Wang, Yihan Wu, Ruihua Song, Xu Tan, Zehua Chen, Hongteng Xu, and Guodong Sui. 2024. Tiva: Time-aligned video-to-audio generation. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 573–582.
- [42] Zehan Wang, Yang Zhao, Haifeng Huang, Jiageng Liu, Aoxiong Yin, Li Tang, Linjun Li, Yongqi Wang, Ziang Zhang, and Zhou Zhao. 2023. Connecting multi-modal contrastive representations. *Advances in Neural Information Processing Systems* 36 (2023), 22099–22114.
- [43] Zhifeng Xie, Shengye Yu, Qile He, and Mengtian Li. 2024. Sonicvisionlm: Playing sound with vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 26866–26875.
- [44] Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and Qifeng Chen. 2024. Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7151–7161.
- [45] Haoran Xu, Jie Zhou, Mengduo Yang, and Jiaze Li. 2024. Shortform ugc video quality assessment based on multi-level video fusion with rank-aware. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, Vol. 7.
- [46] Jinyu Yang, Mingqi Gao, Zhe Li, Shang Gao, Fangjing Wang, and Feng Zheng. 2023. Track anything: Segment anything meets videos. *arXiv preprint arXiv:2304.11968* (2023).
- [47] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*. 6023–6032.
- [48] Yiming Zhang, Yicheng Gu, Yanhong Zeng, Zhenning Xing, Yuancheng Wang, Zhizheng Wu, and Kai Chen. 2024. Foleyrafter: Bring silent videos to life with lifelike and synchronized sounds. *arXiv preprint arXiv:2407.01494* (2024).