

# Audio-Guided Visual Animation

Lin Zhang<sup>1</sup>

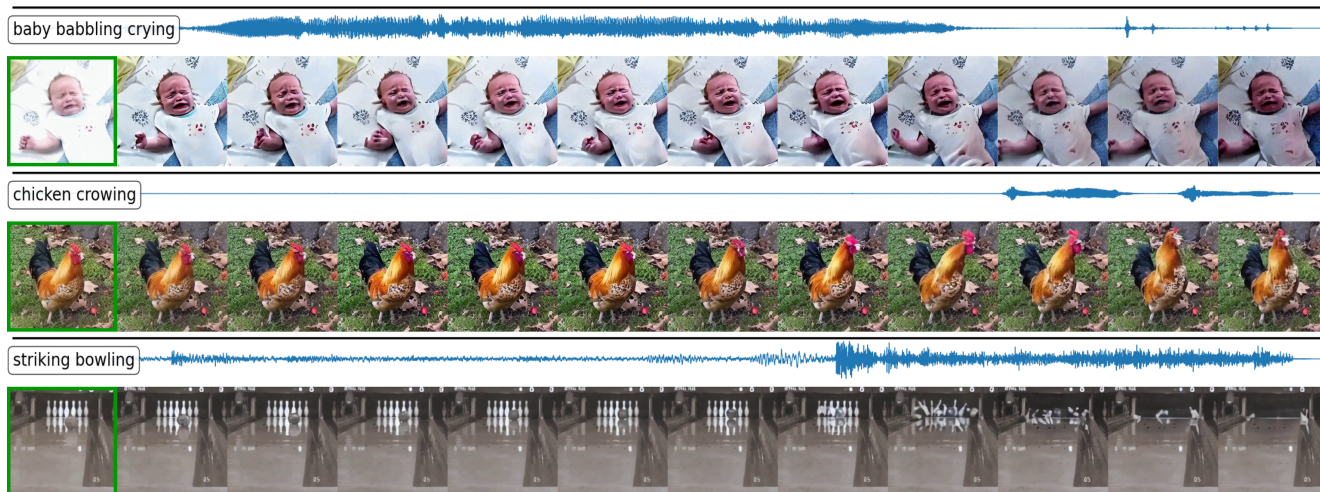
Shentong Mo<sup>2</sup>

Yijing Zhang<sup>1</sup>

Pedro Morgado<sup>1</sup>

University of Wisconsin Madison<sup>1</sup>

Carnegie Mellon University<sup>2</sup>



**Figure 1.** Given an audio and an image (green box), we produce animations beyond image stylization with complex but natural dynamics, synchronized with input audio at each frame. Results are produced by our A2VD model trained on the proposed AVSync15 dataset.

## Abstract

Contemporary visual generation methods often fall short in effectively control along temporal dimension. In response, we introduce Audio-Guided Visual Animation (AGVA), a task aimed at generating image animations that are temporally synchronized with audio cues. To address the absence of datasets tailored for this task, we present AVSync15, the first benchmark with highly synchronized audio-visual dynamics. Curated from the extensive but noisy VGGSound dataset, AVSync15 consists of 15 diverse audio-visual categories ranging from animal sounds, human actions, musical instruments, to triggered events. In addition, we introduce A2VD, a diffusion model capable of producing semantically aligned and temporally synchronized image animations from audio. We provide thorough evaluations to validate AVSync15 as a reliable dataset for synchronous video dynamics generation task and the superior performance of A2VD. Moreover, we explore various potentials of the trained A2VD in a range of audio-guided visual generation applications, bringing in new vision for controllable visual generation.

## 1. Introduction

Generative modeling has witnessed remarkable progress in recent years, largely due to the development of stronger and robust architectures such as diffusion [13, 27, 30] models. Conditional generation and, in particular, text-to-image generation [26, 27], given its immense application potential and the availability of high-quality training data [28], has been the focal point. Nevertheless, the success of text-to-image generation has also spurred exploration of generation in other modalities, including text-to-video [7, 16, 29, 35], text-to-audio [14], audio-to-image [10, 31], among others.

While text conditioning has been well investigated, the potential of audio has been largely overlooked. The temporal dimension of audio signals offers a unique advantage over text for video generation. While text conditioning provides direct control over global semantics, audio conditioning can provide additional fine-grained control at each frame of the video generation process. However, current audio-conditioned generation models, including those used in audio-to-image [10, 19, 31] and audio-to-text [17] generation, often encode audio into a single global semantic

feature, neglecting its temporal aspect. Even in prior audio-to-video generation works [19], the focus has been primarily on semantic correlation, with temporal synchronization between audio and video motion remaining largely unexplored. Although there exists some recent work [15, 18] attempting to generate audio-synchronized video, they mainly studied monotonous audio classes [19] such as weather and environment sound, where the audio cues can be simply connected with visual effects by changing image texture and styles. The complex visual dynamics triggering the sound, such as object motion and interaction, are completely ignored.

In this work, we bridge this gap by tackling a more challenging generation task – Audio-Guided Visual Animation (AGVA). The goal of AGVA is to animate images, generating a video sequence with motion dynamics semantically-aligned and temporally-synchronized with an input audio. Aiming at synchronizing the dynamic aspects of audio with appropriate visual changes in the animation, it necessitates a sophisticated use of the audio’s temporal structure. AGVA not only expands the scope of conditional generation, but also introduces a novel dimension of fine-grained control beyond text tokens for multi-modal content generation.

Despite the potential applications, AGVA presents several challenges. The first challenge relates to a lack of quality training datasets to learn video dynamics synchronization, which requires strong correlations at audio-visual content at each moment in the video. In other words, sound sources should be easily located in the scene and their visual motions should be clearly associated with the corresponding sound (temporally synchronized and semantically appropriate). In addition to the requirements for audio-visual synchronization, the video content should also be of high quality for generation. However, existing audio-visual datasets are either too noisy [5, 6, 9], containing a large number of unassociated audio-visual pairs [24], or are overwhelmed by ambient sound categories that lack meaningful motion cues [19].

To address this, we curated a high-quality dataset for AGVA, denoted AVSync15, from the noisy VGGSound dataset [5]. Due to the overwhelming noise of real-world videos as in VGGSound, direct manual curation would be laborious and inefficient. We thus leveraged a two-step data cleaning pipeline with automatic and manual curation steps. In the first step, we use a variety of signal processing techniques and foundation models to automatically filter videos according to several metrics, from raw pixel differences to high-level audio-visual synchronization. Then, to ensure the high quality of training data, we further conduct manual curation in the second step, resulting into a final dataset with 1500 sounded videos uniformly spreaded over 15 diverse categories, from animal sounds to triggered events.

The second challenge lies in the development of effective AGVA models capable of generating natural and highly syn-

chronized video motions. The closest work to ours are AADiff [18] and TPoS[15], which however have been discussed above to focus on environment sound categories lacking motion cues and merely stylize images along the temporal dimension. AADiff even simplified audio features into temporal amplitudes and a global semantic feature, thus generating visual effect re-weighted by audio amplitude at each frame. Audio-video synchronization in the world however is much more complex. For instance, dog barking involves not only "opening mouth" motion synchronized with the "barking" timestamp, but also subtle details such as the dog’s change of pose at the other timestamps, e.g., raising head.

To address this challenge, we propose a novel architecture named Audio-to-Video Latent Diffusion model (A2VD), which builds upon a pre-trained latent diffusion model with text conditioning [27], and modify it for more effective synchronization. First, to enable precise semantic and synchronized audio control at each timestep, we leverage the pre-trained ImageBind [10] model to encode audio into time-aware semantic tokens, and fuse them into image latent features via cross attention. To capture complex video motions, we incorporate temporal attention layers with learnable positional embeddings to the diffusion model. Finally, to encourage faithful animation of the provided image, we introduce temporal convolutions and attention layers to always lookup on the input image, *i.e.*, first-frame conditioning.

With the carefully designed model and dataset, we are able to obtain a well-trained model specialized for AGVA, and produce animations with visually pleasing and audio-synchronized motions (Fig. 1). We provide thorough experiments to validate the effectiveness of the proposed dataset, AVSync15, and the architecture design of A2VD. We also demonstrate how to deploy A2VD for a variety of audio-guided applications, including editing and replacing the input audio. Code and dataset will be open-sourced.

## 2. Related Work

**Conditional visual generation** Many conditional visual generation models based on diffusion process [13, 30] have emerged recently. Benefiting from more efficient architecture, large-scale training data [28], and aligned semantic space [25], Latent Diffusion Model [27] has achieved great success to generate semantically aligned images conditioned on text. This inspired researchers to explore various visual generation tasks, such as text-to-video [3, 16, 35, 36], audio-to-image [10, 31], and audio-to-video [15, 19]. While some work adopted a training-free strategy [16, 18, 35] or trained from scratch, the others augmented the architecture of a pretrained text-to-image model by carefully adding some trainable layers to learn task-specific information [3, 15].

In this work, we extend this trend by developing a model for the Audio-Guided Visual Animation task, augmenting

pre-trained StableDiffusion models with temporal layers and a synchronized audio conditioning mechanism.

**Audio-to-Video generation** Audio, like text, has been widely used as a semantic signal in visual generation tasks [10, 31]. However, these approaches often overlook the temporal aspect inherent in audio. Traditionally, this temporal aspect has been leveraged to generate synchronized talking faces [23, 37–39], but synchronized visual content generation across a broader range of classes has been relatively unexplored. Recent advances include AADiff [18], which re-weights the word-image cross-attention map in LDM at each timestep using audio amplitude to produce visually pleasing results. Similarly, TPoS [15] learns segmented audio semantic features to fuse with LDM, aiming to create audio-synchronized video content. These methods, however, primarily focus on monotonous sound classes, as demonstrated in the Landscapes dataset [19]. They are limited to modifying appearance and style within an image without capturing the natural dynamics of video content.

Addressing this limitation, our work introduced AVSync15, a high-quality dataset specifically designed for the Audio-Guided Visual Animation task. AVSync15 stands out from previous efforts by focusing on synchronization cues between audio and visual dynamics. This allows for object-centric animation generation, moving beyond mere visual effect animation. We further propose a model, A2VD, to facilitate AGVA task by training on AVSync15.

### 3. Audio-Guided Visual Animation

In this work, we introduce Audio-Guided Visual Animation (AGVA), where the goal is to generate a video conditioned on an audio clip and a single image. Formally, given a  $T$  second audio clip  $\mathbf{a}$  and an image  $\mathbf{x}_1$ , the goal of AGVA is to generate a sequence of  $T \times r - 1$  frames ( $\mathbf{x}_2, \dots, \mathbf{x}_{T \times r}$ ) constituting the video animation, where  $r$  is the desired frame rate. Despite the simple formulation, our AGVA task is challenging in that an effective generated video sequence must be 1) composed of high-quality generated frames, 2) semantically aligned with the given image  $\mathbf{x}_1$  and audio  $\mathbf{a}$ , 3) temporally coherent to model the natural motion dynamics, and (4) the motion of frames is well synchronized with given audio  $\mathbf{a}$ . Previous works on audio-reactive video generation [15, 18] mainly study visual effects animation on monotonous classes, thus hardly satisfy 3 and 4 especially when requiring generating object actions.

#### 3.1. AVSync15: A High-Quality Audio-Visual Dataset for Synchronized Video Generation

We start by observing that existing audio-visual datasets are either too challenging [5, 6, 9] for audio-to-video generation tasks due to noises like rapid scene changes/camera motion, missing/static frames, and out-of-scene audio, or

only contain ambient and style classes [19, 20] like weather.

Thus, to facilitate research on AGVA, we assembled a high-quality dataset specifically designed for audio-guided video generation, ensuring a close synchronization between audio and visuals. In broad terms, the selection of videos for our dataset was based on the following criteria. 1) *High Correlation*: Every significant visual change in the video should be closely associated with audio at each timestamp, and vice versa. 2) *Dynamic Content*: We sought content rich in temporal changes, excluding ambient or monotonous classes like running fans or rain. 3) *Quality and Relevance*: Both video and audio needed to be clean, stable, and representative of their respective categories.

**Preliminary curation** We constructed our dataset from VGGSound [5], a large-scale dataset with 309 diverse audio classes. Similar to VGGSoundSync [6], we began by narrowing down the videos to 149 classes with potentially clear audio-visual synchronization cues, removing ambient classes without video synchronization events, such as *hair dryer drying*. We refer to this intermediate dataset as VGGSS. From VGGSS, we further deployed a sequence of automatic cleaning steps and a final manual selection stage to identify appropriate videos. We summarize curation procedures below and provide every details in the Suppl.

**Automatic curation** First, we utilize PySceneDetect [1] to split videos with sharp scene changes to different scenes. These scenes are still likely to contain both high-quality and low-quality short sub-clips. To maximize usage, we split each scene into 3-second clips with 0.5-second strides, and filtered out unsuitable clips based on the following metrics:

Raw Pixel Difference We calculate average pixel differences between consecutive frames and remove clips with static frames/small or excessive motion/large value.

Image Semantics Difference To complement the above metric in the semantic space with potential zoom-in/out static frames and semantic transitions, we calculate a similar score as above by encoding images into CLIP features.

Waveform Amplitude We exclude clips whose maximum waveform amplitude is low, indicating weak audio cues.

CLIP Semantic Alignment With pre-trained ImageBind [10], we compute average Image-Audio and Image-Text CLIP alignment scores [25] (cosine similarity of CLIP features) in a video, removing clips with low scores to ensure cross-modal semantic alignment.

Audio-Video Synchronization To measure audio-visual synchronization, we follow VGGSoundSync [6] to contrastively train an audio-visual synchronization classifier on VGGSS, ending up with a comparable 40.85% test accuracy. The model outputs an unbounded AVSync score  $\phi_{ij}$  for each input audio-video pair  $(\mathbf{a}_i, \mathbf{v}_j)$ . During training, these scores are computed for the synchronized pair  $(\mathbf{a}_i, \mathbf{v}_i)$  and multiple temporally shifted pairs from the



same instance. Contrastive loss is then applied on these shifted pairs to maximize the synchronization probability  $p_{sync,i} = \frac{\exp(\phi_{ii})}{\sum_j \exp(\phi_{ij})}$  to distinguish the synchronized pair from shifted ones. With  $p_{sync}$  as a synchronization indicator, we removed the low-scoring clips.

We empirically determined the thresholds of each metric by prioritizing quality, acknowledging that some acceptable clips might be discarded to maintain a high-quality final dataset. After automatic curation, we merged consecutive 3-second clips sampled from the same video, and further removed categories with less than 100 examples to address category imbalances, resulting in a dataset with 76 categories and 39,902 examples. We refer to this dataset as AVSync-AC (Audio-Visual Synchronization with Automatic Curation).

**Manual curation** Manual intervention is still essential to ensure quality. To this end, we selected a diverse set of 15 categories with clear audio-visual cues from AVSync-AC for further manual refinement, including categories ranging from animals and human actions to triggered tools and musical instruments. Manual curation once again sought to identify appropriate videos for AVGA using the criteria above: high correlation, dynamic content, and quality and relevance. We also extracted sub-clips with minimal duration of 2-seconds from examples when necessary.

**Dataset comparison** The final dataset, AVSync15, contains 100 videos per category, each 2 to 10 seconds long. We allocated 90 videos for training and 10 for testing in each category. We provide an overview of AVSync15 in Fig. 2. To verify the effectiveness of our curation pipeline, we randomly sample 3 splits with 1500 data on the selected 15 categories from VGGSS and AVSync-AC, and quantitatively compare them with AVSync15 in Fig. 3. We also compare AVSync15 with existing audio-visual datasets in Suppl., highlighting its attributes for AGVA tasks.

### 3.2. AGVA Evaluation Metrics

AGVA is a multi-faceted generation task, requiring carefully designed and diverse metrics for comprehensively evaluation. Therefore, in our proposed benchmark, we first evaluate the generated video from the following conventional metrics:

**Visual Quality** Following previous works [4, 8], we use Fréchet Inception Distance (FID) [11] to measure image quality of image frames and Fréchet Video Distance (FVD) [33] to evaluate the quality of the generated videos.

**Semantic Alignment** To assess semantic alignment, we reuse the CLIP alignment scores mentioned in *automatic curation*, i.e., IA-Align and IT-Align.

**Human Evaluation** While automated metrics are useful for quantitative evaluation, they are not always aligned with human perception. A human evaluation study is thus used to assess the quality of the generated videos. Specifically, we ask human raters to compare videos generated by

multiple models and select the best according to Image Quality, Temporal Frame Consistency, and Audio-Visual Synchronization. We provide more details on the human evaluation study in the supplementary material.

Furthermore, when evaluating a generated video  $\hat{v}$  from audio  $\mathbf{a}$ , we cannot generate its shifted videos, which poses challenges to compute aforementioned synchronization probability  $p_{sync}$  as a metric as in *Automatic Curation*. This can be solved by contrasting on groundtruth video  $v$  instead, i.e.,  $p_{sync} = \frac{2 \times \exp(\phi(\hat{v}, \mathbf{a}))}{\exp(\phi(\hat{v}, \mathbf{a})) + \exp(\phi(v, \mathbf{a}))}$ , where the multiplier 2 is used to normalize  $p_{sync}$  into range  $[0, 1]$ . However, such a perception metric is still inaccurate for generated results. This is because during conventional training [6], the synchronization classifier only contrasts on audio-video pairs sampled from the same instances, thus is implicitly conditioned on the semantically aligned audio-visual content as  $p_{sync|align}$ . In practice, when the generated visual content drift too much from groundtruth frame, the metric can fail. We thus propose Aligned Synchronization (**AlignSync**), a more robust metric that more faithfully measures audio-visual synchronization. We first measure semantic alignment probability between generated  $\hat{v}$  and condition audio  $\mathbf{a}$  by normalizing Image-Audio Alignment score into  $[0, 1]$ , i.e.,  $\frac{IA-Align+1}{2}$ . We then recover the implicit semantics condition of  $p_{sync}$  by multiplying it with the semantic alignment probability, i.e.,

$$AlignSync = p_{sync} \cdot \frac{IA-Align + 1}{2} \quad (1)$$

The metric then is used to evaluate synchronization between generated results and input audio, with our pretrained synchronization classifier and ImageBind. We provide empirical justifications for AlignSync in supplementary material.

## 4. Audio-to-Video Generation

### 4.1. Preliminary: Text-to-Image Latent Diffusion

Text-to-image latent diffusion model (LDM [27]) encode images  $x$  into a lower-dimensional latent space  $z = \mathcal{E}(x)$  using a pre-trained perceptual auto-encoder, and learn the conditional distribution  $p(z|\tau)$  of latent space given a CLIP-encoded text prompt condition  $\tau$ . Specifically, LDM models the conditional distribution by learning to gradually denoise latents  $z^k$  at each diffusion step  $k$ , which are obtained by corrupting the image latent  $z$  by normally distributed noise  $\epsilon$  over  $k$  time steps. A denoising UNet architecture parameterized by  $\theta$  is deployed to estimate the added noise  $\epsilon$  by minimizing the following objective

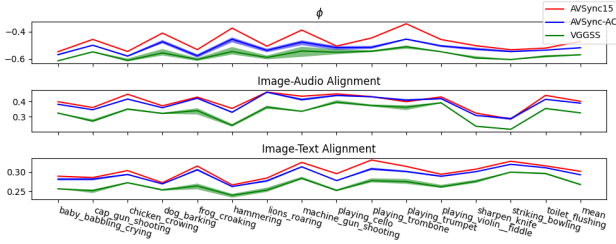
$$\mathcal{L}_{LDM} = \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), k} [\|\epsilon - \epsilon_{\theta}(z^k, k, \tau)\|_2^2] \quad (2)$$

During inference, LDMs start from a random gaussian noise map  $z^K$ , and iterate over  $K$  reverse diffusion steps [13, 30], gradually predicting and removing residual noise  $z^{k-1} = z^k - \epsilon_{\theta}(z^k, k, \tau)$ , until the image latent is found  $z^0$ . LDMs then decode the latent into image





**Figure 2.** Overview of 15 categories in AVSync15. Left to right: *baby babbling crying, dog barking, lions roaring, chicken crowing, frog croaking, playing cello, playing trombone, playing trumpet, playing violin, cap gun shooting, machine gun shooting, hammering, shapen knife, striking bowling, toilet flushing.*



**Figure 3.** Category-wise comparison of AVSync score  $\phi$ , IA-Align, and IT-Align scores on AVSync15 and equivalently sized 3 splits of VGGSS and AVSync-AC.

$x^0 = \mathcal{D}(z^0)$  using the pre-trained decoder  $\mathcal{D}$ . We will refer to the images by their latent representation  $z$  (rather than  $x$ ), for simplicity, throughout the rest of this paper.

## 4.2. Audio-to-Video Latent Diffusion

Given the impressive performance of latent diffusion models, we seek to adapt them for audio-guided video generation. However, most of current approaches [10] use audio primarily for its global semantics as opposed to temporal synchronization. We thus propose an Audio-to-Video Latent Diffusion model (A2VD), which starts from a pretrained image LDM and incorporates synchronized audio control and trainable temporal layers for improved video consistency.

Given an image  $z_1$ , a text prompt (*i.e.*, the category name)  $y$  which is encoded by CLIP into  $\tau$ , and a  $T$ -seconds audio signal  $\mathbf{a}$ , our model generates sequences of  $r \times T - 1$  future frames  $\{z_t\}_{t=2}^{rT}$  depicting plausible evolutions of the image  $z_1$  over time synchronized with the provided audio, through iterative denoising. Given a video dataset with synchronized audio-visual content, the denoising UNet,  $\epsilon_\theta(z_{2:rT}^k, k; z_1, \mathbf{a}, \tau)$ , can be easily trained by randomly sampling sequences of frames, using the first frame as the input image  $z_1$  (together with the corresponding audio and text prompt conditioning), and using the remaining frames as targets for  $z_{2:rT}^0$ . The overall architecture, illustrating how A2VD incorporates the various conditioning signals,  $z_1, \mathbf{a}, \tau$ , is shown in Fig. 4, and discussed in detail below.

**First-frame conditioning** While LDMs do not inherently support image conditioning, previous studies have introduced image inversion methods for this purpose. Image inversion however can be both inaccurate [30] and time-consuming [21]. To circumvent this, we directly input  $z$  into the UNet model as the known latents of the first frame, irrespective of the diffusion step  $k$ . For all subsequent frames, we adhere to the original LDM design, using independently

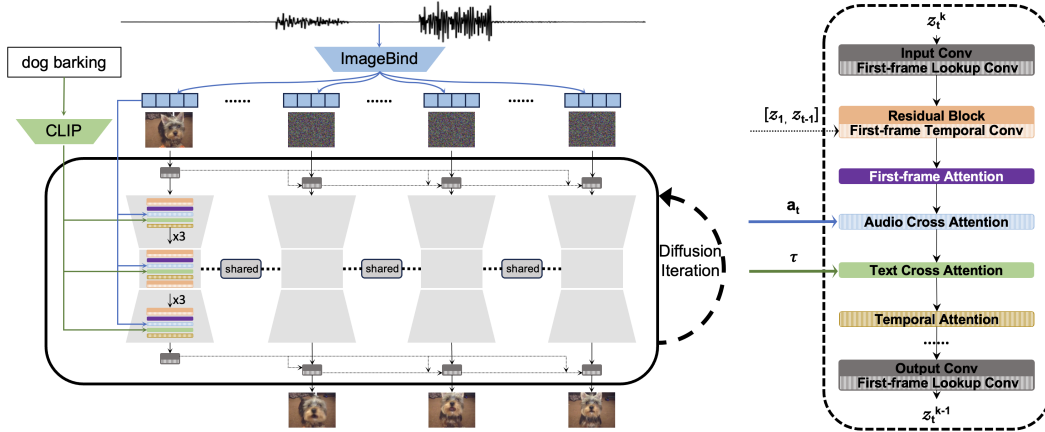
sampled noised latents  $(z_2^k, \dots, z_{rT}^k)$  as initial inputs.

**Temporal convolutions and temporal attentions** To generate temporally consistent videos, we incorporate temporal convolutions and attention layers into the UNet backbone. Similar to R(2+1)D [32], we append a 1D temporal convolution layer with a kernel size of 3 after each 2D conv layer. We also introduce a unique first-frame lookup to better adhere to the starting image (see next paragraph). In addition to temporal convolutions, we also include temporal attention layers [3] with learnable temporal positional encodings [34] to effectively model long-range visual dependencies. Each frame index  $t$  is converted into a sinusoidal positional embedding, which, after a learnable linear projection, is added to the corresponding frame’s latents. Each frame’s local representation,  $z_{hwt}$ , is then updated by attending to all frame latents at the same position, including the base frame  $(z_{hw1}, z_{hw2}, \dots, z_{hw(rT)})$ , using a learnable self-attention.

**First-frame lookups and first-frame attention** To ensure that video generation adheres to the input image, we adjusted the receptive field of all temporal convolutions to always encompass  $z_1$ , thereby preventing it from being overlooked when generating distant future frames. Specifically, the receptive field of each frame  $z_t$  include frames  $(z_1, z_{t-1}, z_t)$  as opposed to  $(z_{t-1}, z_t, z_{t+1})$ . This first-frame lookup mechanism is applied to three components in the UNet, *i.e.*, the input/output conv layer and all residual blocks.

In addition to first-frame lookups, we further modify the spatial self-attention layers of the backbone UNet model, adopting the first-frame conditioning strategy proposed in [16] for enhanced image animation. Specifically, in these layers, the representation of each frame,  $z_t$ , are updated by attending to the base image’s representation,  $z_1$ , rather than the frame  $z_t$  itself. Consistent with [16], we keep the first-frame attention layers frozen during this process.

**Audio conditioning** To facilitate audio-guided generation, we employ a pretrained ImageBind audio encoder [10] to encode the audio. This encoder computes a global audio token,  $\mathbf{a}^g$ , encapsulating global semantics, and  $F_a \times T_a$  patch tokens,  $\mathbf{a}_{f,t}$ , across a grid of  $F_a$  frequency bands and  $T_a$  time steps, providing local synchronization cues. We achieve frame-specific audio guidance by dividing the patch tokens temporally into  $rT$  segments, corresponding to the base frame  $z_1$  and the following  $(rT - 1)$  frames to be generated, and appending the global token to each segment. The resulting sequence of audio tokens,  $\mathbf{a}_t$ , for the first given frame ( $\mathbf{a}_1$ ) and for generated frames ( $\{\mathbf{a}_t\}_{t=2}^{rT}$ ), are then



**Figure 4.** A2VD overview. *Left:* We use ImageBind to encode audio into semantically aware time-dependent feature tokens  $\{\mathbf{a}_t\}_1^T$ , and CLIP to encode audio category into prompt condition  $\tau$ . During inference, the model receives first frame and subsequent frame noise latents, and iteratively refines the subsequent noises via diffusion. First-frame lookup convolutions at input/output conv layers and intermediate residual layers (hidden for ease of visualization), first-frame attention, audio cross attention, and temporal attention layers are introduced to learn synchronized visual motion. *Right:* Anatomy of modules processing each frame. Trainable layers are marked with vertical stripes. Different frames share the same UNet.

input into the UNet model and fused with the representations of the corresponding frame,  $\mathbf{z}_t$ , via cross-attention [34].

**Text conditioning** We follow original LDMs [27] by feeding in the audio category as prompt to every frame latent via frozen text cross attention layer.

**Classifier-free audio guidance** Classifier-free guidance [12] is a technique used in generative models to amplify the influence of the input prompt on the generated output, without the need for a separate classifier to validate the output. We extend this concept to amplify audio guidance for improved synchronization. Specifically, we first compute a null audio embedding,  $\mathbf{a}_\emptyset$ , by encoding an all-zero waveform. During training, we randomly replace  $\mathbf{a}$  with  $\mathbf{a}_\emptyset$  with a 20% probability, thus, training the model for both audio-conditioned and unconditioned generation. Then, during inference, we can enhance the effect of audio guidance by scaling the latents generated from the unconditional generation to the conditioned generation with a factor  $\eta$

$$\mathbf{z}_{2:T}^{k-1} = (1 - \eta) \cdot \epsilon_\theta(\mathbf{z}_{2:T}^k, k; \mathbf{z}_1, \mathbf{a}_\emptyset, \tau) + \eta \cdot \epsilon_\theta(\mathbf{z}_{2:T}^k, k; \mathbf{z}_1, \mathbf{a}, \tau) \quad (3)$$

In practice, we made classifier-free audio guidance optional and did not use prompt classifier-free guidance as in LDM by always feeding in audio category as condition  $\tau$ .

## 5. Experiments

### 5.1. Implementation

**Dataset** Most of existing audio-visual datasets [5, 6, 9] are in poor quality for visual generation tasks due to camera motion, cluttered background, and non-centric objects. The other high-quality datasets [19] contain mostly ambient sounds, such as whether and color style, without desired

natural video dynamics and audio-visual synchronization motion cues. One previously collected dataset potential for AGVA is *The Greatest Hits* [22], which captures the unique audio-visual responses of various objects and materials (such as dirt, water, or a desk) when struck by a stick. It is of high quality however with limited diversity, as all videos contain unitary motion of hitting, similar to the *hammering* category in our AVSync15 dataset. Thus, we conducted our main experiments on the proposed AVSync15 with 15 categories, 1350 training videos, and 150 testing videos, and used *The Greatest Hits* for further evaluation of the proposed model. We also verified the effectiveness of the proposed dataset curation pipeline by evaluating models trained on equivalently sized uncurated versions of the dataset (VGGSS), and with automatic curation alone (AVSync-AC).

**Baselines** To verify the effectiveness of audio input, we first provide two baselines without audio input. One is the recently introduced image-text-to-video baseline VideoCrafter [7] pretrained on the large-scale video dataset WebVid10M [2] with 10M high-quality videos. The second is I2VD, our proposed model trained without audio input. Secondly, we re-implemented AADiff [18], an image editing based training-free synchronized audio-to-video generation model. We do not compare to TPoS [15] since its architecture is challenging for faithful re-implementation.

**Training and evaluation** We adopted publicly-available StableDiffusion-V1.5 [27] and ImageBind [10] trained on subsets of LAION2B-en [28] and AudioSet [9] respectively. We only trained the temporal layers with stripes in Fig. 4. In all experiments, we use Adam optimizer with a batch size of 64, a constant learning rate of 0.0001, and weight decay of 0.01. Models were trained for 37000 iterations with 256x256

**Table 1. (a):** Overview of Audio-Guided Visual Animation performance on AVSync15. User study on the right side shows the number of votes to compare 4 models on 3 metrics: Image Quality, Frame Consistency, and Audio-Video Synchronization. **(b):** Performance on TheGreatestHits. **(c):** Effect of first-frame conditioning on models trained and evaluated on AVSync15.

Model	FID↓	FVD↓	IA-Align↑	IT-Align↑	AlignSync↑	User Study		
						IQ	FC	Sync
VideoCrafter [7]	11.20	0.120	36.9 ± 9.8	29.9 ± 3.4	59.08 ± 10.06	38	20	12
AADiff [18]	16.53	0.172	34.6 ± 11.3	29.1 ± 4.2	61.15 ± 9.68	37	4	5
I2VD	11.40	0.064	38.3 ± 9.6	30.4 ± 2.8	60.74 ± 9.22	62	90	91
A2VD w/o $\eta$	11.19	0.070	38.6 ± 9.6	30.5 ± 2.8	62.24 ± 8.57	-	-	-
A2VD w/ $\eta = 1$	11.49	0.076	38.4 ± 9.6	30.3 ± 2.9	62.01 ± 9.00	-	-	-
A2VD w/ $\eta = 4$	11.13	0.075	38.6 ± 9.6	30.4 ± 2.8	63.06 ± 8.51	163	186	192
A2VD w/ $\eta = 8$	11.18	0.091	38.1 ± 9.6	30.3 ± 2.8	63.31 ± 8.56	-	-	-
A2VD w/ $\eta = 12$	11.40	0.118	37.2 ± 9.6	29.9 ± 2.9	63.24 ± 8.58	-	-	-

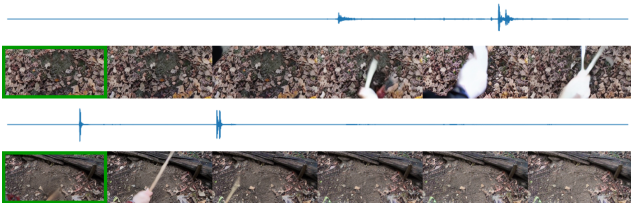
  

Model	FID↓	FVD↓	AlignSync↑
I2VD	8.38	0.068	49.60 ± 7.38
A2VD w/o $\eta$	8.20	0.065	50.54 ± 7.38
A2VD w/ $\eta = 1$	8.28	0.051	50.92 ± 7.26
A2VD w/ $\eta = 4$	8.07	0.040	51.51 ± 7.25

FF-Lookups	FF-Attn	FID↓	FVD↓	AlignSync↑
✗	✗	11.36	0.071	61.83 ± 9.07
✓	✗	11.16	0.068	61.97 ± 8.92
✓	✓	11.19	0.070	62.24 ± 8.57

(a)



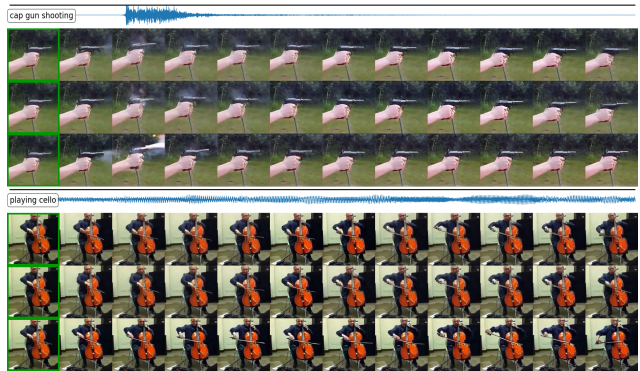
**Figure 5.** Results of A2VD with  $\eta = 4.0$  on The Greatest Hits. We use prompt *hitting with a stick* for all videos.

image size on AVSync15 and 16000 iterations with 128x256 on The Greatest Hits, both with 2-second input audio and 12-frame video in 6 FPS. We use the metrics in Section 3.2 for evaluation, *i.e.*, FID, FVD, IA-Align, IT-Align, AlignSync, and user study. We put more evaluation details in Suppl.

## 5.2. Main results

Table 1a compares the models on AVSync15. Since model receives the first frame as condition, which is a shortcut to generate animations, all reported FID/FVD scores are relatively low. As can be seen, A2VD outperforms all baseline methods in terms of AlignSync, while maintaining high visual quality and semantic alignment. In contrast, VideoCrafter and I2VD, without audio conditioning, score low on synchronization. VideoCrafter also struggles to generate high-quality, semantically aligned content, as evidenced by its performance on FVD, IA-Align, and IT-Align. Upon examining sample outputs from VideoCrafter, we observed that this could be attributed to its inability to accurately replicate the input image, resulting in semantic drift from given image. On the other hand, by directly reweighting the attention to prompt tokens using audio amplitude, AADiff achieves a relatively high AlignSync score. However, the simplistic use of audio amplitude for adjusting animation compromised visual quality and semantic alignment, resulting into frequent flickering and artifacts. AADiff also failed in categories such as playing violin, hammering, and striking bowling, where the dynamics are far more complex than stylization

(c)



**Figure 6.** Audio amplitude versus classifier-free audio guidance. We visualize videos generated with *top*: original audio with  $\eta = 1$ ; *mid*: 100× amplified audio with  $\eta = 1$ ; *bottom*: original audio with  $\eta = 8$ .

and weather. In contrast, A2VD was able to generate frame sequences with natural video dynamics and more aligned with the input audio. Generated samples are shown in Fig. 1. We provide more qualitative comparisons in Suppl.

**Human evaluation** Our user study asked participants to compare the videos generated by four different method and select the best one for each of the three criteria: visual quality, temporal consistency, and audio-visual synchronization. In total, we collected 900 responses, uniformly distributed among all classes, from 15 participants, and reported the number of votes obtained by each method in Table 1a. As shown, A2VD generated the best image animations on all three criteria, with an especially large margin on AlignSync.

## 5.3. Ablation studies

**Audio guidance** We explored the effect of classifier-free audio guidance factor  $\eta$ . By increasing  $\eta$  from 1.0 to 8.0, we observed that the generated frames have clearer visual effects indicative of the audio input, yielding generated videos that appear better synchronized. This observation was validated quantitatively in Table 1a and qualitatively in Fig. 6.



Fig. 6 also compares increased audio-guidance factors with videos generated with increased audio amplitudes. Prior audio→visual generation works [15, 18, 31] claimed that louder audio often leads to visual effects that are better aligned with the audio condition. To examine the influence of audio amplitude, we generated videos for two classes, *cap gun shooting* and *playing cello*, with the audio amplitude increased by a factor of 100. Our findings, illustrated in Figure 6, show that an extreme increase in audio amplitude does not distort the generated frames but slightly intensifies the visual effects, as seen by the presence of smoke. In comparison, audio guidance was much more effective in enhancing not only visual effects but also object dynamics, as indicated by more exaggerated hand moving when playing cello.

**Audio conditioning** We evaluated our model’s ability to incorporate audio cues by training an audio-unconditioned variant, I2VD. The quantitative results, presented in Table 1a and Table 1b for the AVSync15 and The Greatest Hits respectively, demonstrate that our proposed architecture effectively enhanced audio-visual synchronization without compromising visual quality or semantic alignment. The introduction of classifier-free audio guidance allows for a trade-off between visual quality and higher synchronization.

**Effect of first-frame conditioning** We also evaluated the effect of the proposed first-frame conditioning mechanism on model performance. As shown in Table 1c, the model trained with first-frame conditioning achieved better performance especially on AlignSync, suggesting that first-frame conditioning (as also observed in prior work on text-conditioned image animation) can help the model to better generate videos consistent with the original image.

**Effect of data curation** To assess the effectiveness of our data cleaning pipeline, we randomly sampled subsets from VGGSS and AVSync-AC, ensuring equal training data scale and balanced category distributions, *i.e.*, 90 training videos for each of the 15 categories. We trained A2VD on these subsets using the same training strategy, and reported their performance on Table 2. The inferior FID, FVD, and semantic alignment scores of VGGSS highlight the unsuitability of uncured data sources for video generation tasks. On the other hand, automatic curation, as deployed in AVSync-AC, enabled training generation models with visual quality and semantic alignment scores comparable to those from AVSync15. With manual cleaning, synchronization is further improved. We provide qualitative comparisons in Suppl.

#### 5.4. More applications

Although our model is trained for AGVA task on AVSync15, we find the trained model can be easily extended for more applications, including animating images on the internet, audio-to-video generation without image condition, generating into distant future guided by long audio, etc. Here we only discuss one of them and leave the others in Suppl.

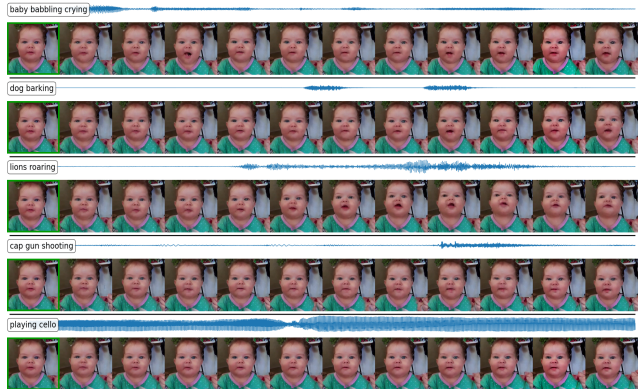


Figure 7. Baby animations controlled by different audios.

Table 2. Effect of training data curation steps on model performance. Models are evaluated on AVSync15 test set.

dataset	Automatic Curation	Manual Curation	FID↓	FVD↓	IA-Align↑	AlignSync↑
VGGSS	✗	✗	12.96	0.143	30.0 ± 12.0	58.95 ± 9.12
AVSync-AC	✓	✗	11.42	0.071	38.2 ± 9.6	61.54 ± 8.81
AVSync15	✓	✓	11.19	0.070	38.6 ± 9.6	62.24 ± 8.57

**Image animation with (un)related audio** To better understand how the model behaves when the audio condition is unrelated to the depicted scene, we animated a variety of images with audio from different categories. Fig. 7 shows an example where a baby’s face is animated to the sound of several audio signals. As can be seen, the baby’s mouth opens and his expression changes in sync with a variety of audio conditions. For example, lions typically roar for longer duration and with clear visual cues like raising head which has been well transferred to the baby in synchronization, without distorting image content. However, the baby’s face remained unaltered when the audio condition was unrelated, such as with the sound of a gun shooting and playing cello. This interesting behavior can be leveraged for a variety of audio-guided animation objectives in the real world.

## 6. Conclusion

In this paper, we propose to solve Audio-Guided Visual Animation (AGVA) task, with an emphasis on learning synchronization between provided audio and generated video dynamics. Lacking an appropriate dataset for the task, we adopted two-stage data cleaning pipeline to curate a clean and high-quality dataset AVSync15, as well as associate evaluation benchmark. We further proposed an Audio-to-Video Latent Diffusion model for training. As such, we can generate highly synchronized video with consistent motion dynamics from images and audios. We hope our research can inspire future work on conditioned generation.

## References

- [1] Pyscenedetect. <https://www.scenedetect.com/>. 3
- [2] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021. 6
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 5
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019. 4
- [5] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 2, 3, 6
- [6] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Audio-visual synchronization in the wild. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2021. 2, 3, 4, 6
- [7] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter1: Open diffusion models for high-quality video generation, 2023. 1, 6, 7
- [8] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. 4
- [9] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP 2017*, 2017. 2, 3, 6
- [10] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023. 1, 2, 3, 5, 6
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017. 4
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop on Deep Generative Models and Downstream Applications*, 2022. 6
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1, 2, 4
- [14] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. In *ICML*, 2023. 1
- [15] Yujin Jeong, Wonjeong Ryoo, Seunghyun Lee, Dabin Seo, Wonmin Byeon, Sangpil Kim, and Jinkyu Kim. The power of sound (tpos): Audio reactive video generation with stable diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 2, 3, 6, 8
- [16] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. In *ICCV*, 2023. 1, 2, 5
- [17] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *NAACL-HLT*, 2019. 1
- [18] Seungwoo Lee, Chaerin Kong, Donghyeon Jeon, and Nojun Kwak. Aadiff: Audio-aligned video synthesis with text-to-image diffusion. In *CVPR Workshop on Content Generation*, 2023. 2, 3, 6, 7, 8
- [19] Seung Hyun Lee, Gyeongrok Oh, Wonmin Byeon, Jihyun Bae, Chanyoung Kim, Won Jeong Ryoo, Sang Ho Yoon, Jinkyu Kim, and Sangpil Kim. Sound-guided semantic video generation. In *ECCV*, 2022. 1, 2, 3, 6
- [20] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++: Music conditioned 3d dance generation. In *ICCV*, 2021. 3
- [21] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *CVPR*, 2023. 5
- [22] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H. Adelson, and William T. Freeman. Visually indicated sounds. In *CVPR*, 2016. 6
- [23] Se Jin Park, Minsu Kim, Joanna Hong, Jeongsoo Choi, and Yong Man Ro. Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2022. 3
- [24] Nuno Vasconcelos Pedro Morgado, Ishan Misra. Robust audio-visual instance discrimination. In *Computer Vision and Pattern Recognition (CVPR), IEEE/CVF Conf. on*, 2021. 2
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2, 3
- [26] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. In *arXiv*, 2022. 1
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 2, 4, 6
- [28] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022. 1, 2, 6
- [29] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data, 2022. 1

- [30] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. [1](#), [2](#), [4](#), [5](#)
- [31] Kim Sung-Bin, Arda Senocak, Hyunwoo Ha, Andrew Owens, and Tae-Hyun Oh. Sound to visual scene generation by audio-to-visual latent alignment. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. [1](#), [2](#), [3](#), [8](#)
- [32] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *CVPR*, 2018. [5](#)
- [33] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. In *arXiv*, 2019. [4](#)
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. [5](#), [6](#)
- [35] Wen Wang, kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video editing using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. [1](#), [2](#)
- [36] Ruiqi Wu, Liangyu Chen, Tong Yang, Chunle Guo, Chongyi Li, and Xiangyu Zhang. Lamp: Learn a motion pattern by few-shot tuning a text-to-image diffusion model. *arXiv preprint arXiv:2310.10769*, 2023. [2](#)
- [37] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, Jinzheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. In *ICLR*, 2023. [3](#)
- [38] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [39] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [3](#)