
Read, Watch and Scream!

Sound Generation from Text and Video

Yujin Jeong* Yunji Kim Sanghyuk Chun Jiyoung Lee†

NAVER AI Lab

Abstract

Multimodal generative models have shown impressive advances with the help of powerful diffusion models. Despite the progress, generating sound solely from text poses challenges in ensuring comprehensive scene depiction and temporal alignment. Meanwhile, video-to-sound generation limits the flexibility to prioritize sound synthesis for specific objects within the scene. To tackle these challenges, we propose a novel video-and-text-to-sound generation method, called ReWaS, where video serves as a conditional control for a text-to-audio generation model. Our method estimates the structural information of audio (namely, energy) from the video while receiving key content cues from a user prompt. We employ a well-performing text-to-sound model to consolidate the video control, which is much more efficient for training multimodal diffusion models with massive triplet-paired (audio-video-text) data. In addition, by separating the generative components of audio, it becomes a more flexible system that allows users to freely adjust the energy, surrounding environment, and primary sound source according to their preferences. Experimental results demonstrate that our method shows superiority in terms of quality, controllability, and training efficiency. Our demo is available at <https://naver-ai.github.io/rewas>.

1 Introduction

Recent generative models have developed dramatically, making content creation easier for people, including images, videos, and audio, based on text prompts. Especially, text-to-video generation models such as Make-a-Video [35] and Sora [1] show the impressive emergence of generative models in the video domain, showing remarkable utility in film and advertising. While we are fully immersed in the video content by watching and listening, unfortunately, these generated videos are silent. Generating the sound of a video is a challenging task requiring both a contextual and temporal understanding of the video. Figure 1 shows an actual example of when both text and video controls are required for generating realistic audio for the given video. In the video, a dog is growling while holding a toy in his mouth. A human can imagine the sound of the video; the dog growls lowly, and the growling sounds like the dog is biting something. When the person grips and pulls the toy, the dog will treat the human by growling louder. Finally, when the dog shakes his head, the growling will become louder. If the generated audio does not understand the visual information, it will be a random growling sound, and it will not sound like the dog bites something. If the audio is not controlled by the text, the generated audio might be only related to the dog, *e.g.*, a barking sound.

Table 1 shows the recent attempts to generate an audio sample from the given video or the given text context. As shown in the table, there are two major directions to generate an audio sample from the given video directly. For example, there have been studies of a sound effect (SFX) generation with

*Works done during an internship at NAVER AI Lab.

†Corresponding author: lee.j@navercorp.com

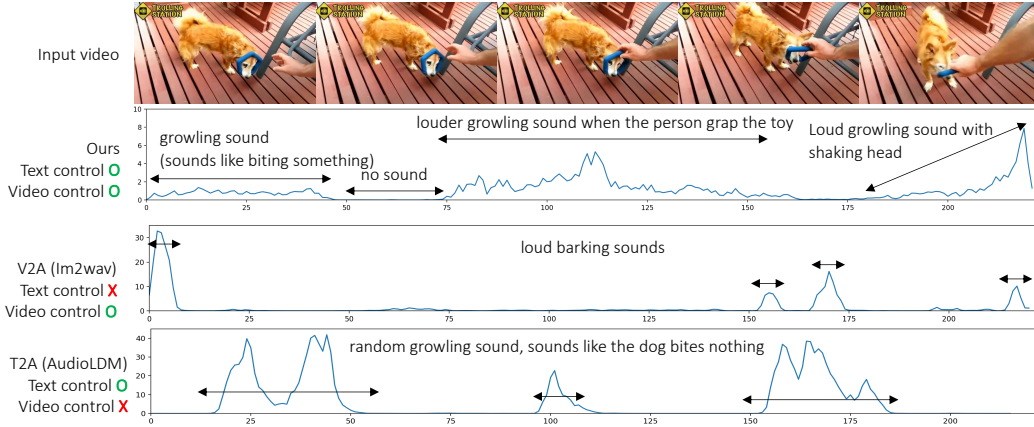


Figure 1: An example of audio generation requiring both text and video control. The text instruction “dog growling” is used for the text control. The video-to-audio (V2A) [34] or text-to-audio (T2A) [25] generation methods cannot understand the detailed semantics from texts (the dog is growling, not barking) or video (the dog is biting something, and the alignment), respectively.

short moments for video editing tasks [4, 7, 43], known as Foley. These methods are restricted to the pre-defined sound effect classes and can only control discrete information, such as onset. As another attempt, video-to-audio (V2A) generation methods have been proposed recently [27, 44, 18, 34]. They hardly generate sounds from multiple objects together but likely generate only simple sounds. Furthermore, both SFX and V2A methods cannot take text controls, which provide more rich control by users. Figure 1 shows the example when there is no text control; a V2A method just generates audio of “barking” rather than “growling” by focusing on the dog in the video. Our work aims to generate general sounds conditioned on both video and text for user control.

As another line of research, text-to-audio (T2A) generation has been actively studied [17, 25, 26, 10, 16]. Despite their diverse and high-quality audio generation quality, they lack a temporal understanding of video-only information. For example, like the example in Figure 1, text-only condition can make irrelevant audio to the video, namely, it might not be aligned with the given video, and it might not consider information only implied by visual information (*e.g.*, the dog holds something in the mouth). To tackle the problem, we may need more controllability to the T2A model, such as AudioLDM [25]. Recently, a few studies [41, 12] tried to control the pre-trained AudioLDM more precisely based on ControlNet [46]. Although they can control the pitch, temporal order, energy, or rhythm of the generated audio, their generation process needs expensive timestamp-wise annotations for each control feature. Namely, the existing T2A works cannot reflect the visual information hidden from the given text or pre-defined control features, and requires expensive timestamp-wise control features.

In this work, we propose a novel video-and-text-to-sound generation approach, named Read, Watch and Scream (ReWaS), by integrating video as a conditional control for a well-established text-to-audio model. Our method is based on a state-of-the-art text-to-audio generation method, AudioLDM [25]. While a text prompt specifies the subject, we additionally employ a control feature extracted from the video. More specifically, our method presents an energy adapter on AudioLDM motivated from ControlNet [46], an efficient structure control method for text-to-image generation. Since a video feature does not directly imply the structure of the audio, we estimate the temporal *energy* information, a basic audio structural information, from the video.

The energy operates as a time-varying control to complement the sound according to the dynamics of the given video. As shown in Figure 1, ReWaS successfully understands complex information from both text and video. Here, we define energy as the mean of frequency in each audio frame, which is related to both visual dynamics and semantics [20, 12]. It is relatively simple to estimate from a video rather than complex acoustic features (*e.g.*, MFCC, mel-spectrograms). Therefore, our energy control facilitates connecting video for T2A model. Lastly, our method offers a more flexible and effective solution thanks to the efficient training [46].

We compare our method and other state-of-the-art video-to-audio generation models [7, 43, 27, 18, 34] on two video-audio aligned datasets, VGGSound [2] and GreatestHits [28]. In the experiments,

Table 1: Comparison of audio generation methods. We consider four factors: Can it make a general sound? Can it take text or visual control? and is its training efficient?

Method	General sound?	Text control?	Visual control?	Efficient training?
Sound effect (SFX) generation [4, 7, 43]	✗	✗	✓ [†]	✗
Video-to-audio (V2A) [27, 44, 18, 34]	✓ [‡]	✗	✓	✗
Text-to-audio (T2A) [17, 25, 26, 10, 16]	✓	✓	✗	✗
T2A + Control [41, 12]	✓	✓	✗	✓
ReWaS (ours)	✓	✓	✓	✓

[†] Unable to adjust continuous sound variations (*i.e.*, energy). [‡] Hardly generate sounds of multiple subjects together.

ReWaS shows a superior audio generation performance quantitatively and qualitatively. For example, our method shows the best fidelity score (FID) and structure prediction (energy MAE) on both datasets. Furthermore, ReWaS outperforms V2A methods in human evaluation for three categories (audio quality, relevance to the video, and temporal alignment between audio and video) with a significant gap (almost +1 point for every category in 5-scale MOS). Also, as shown in the qualitative study, our method successfully deals with the temporal alignment of the visual information. For example, our method can capture the “short transition” of the skateboarding video when the boarder jumps into the air, and no skateboarding sound appears in the video.

2 Related work

2.1 Text-to-audio generation

Early work for conditional audio generation was built upon generative adversarial networks [24, 6], normalizing flows [21], and variational autoencoders [39]. Recently, several studies based on diffusion models have shown promising progress on a broad range of acoustic domains. DiffSound [45] employs a diffusion-based token decoder for the first time to transfer text features into mel-spectrogram tokens. Make-An-Audio [17], AudioLDM [25], AudioLDM2 [26], Tango [10] and Make-An-Audio2 [16] are well-founded in latent diffusion model (LDM) [32], demonstrating high-quality results with large scale training. A series of LDM predicts mel-spectrograms using a VQ-VAE decoder, and a pretrained vocoder generates raw waveforms from the generated mel-spectrograms. While these methods successfully generate high-quality audio samples for the given text prompt, they are only designed for taking text conditions, unable to understand visual semantics.

Meanwhile, there have been a few attempts based on ControlNet [46], an efficient training method for structure control for text-to-image generation. ControlNet utilizes hints (*e.g.*, Canny edge maps, scribbles, human pose, depth maps) to provide structural composition to the generated images. Inspired by this, MusicControlNet [41] showed control over melody, dynamics, and rhythm, while Guo et al. [12] built a FusionNet between each layer of the U-Net, enabling the fusion of control embeddings for temporal order, pitch, and energy controls. They have demonstrated that incorporating control signals into the pretrained audio generative models provides more explicit and fine-grained control over the generated audio, leading to performance improvement and adherence to the desired characteristics. However, designing these time-varying controls still requires skilled labor for users. To address this challenge, we generate energy control through a given video, which is a practical function for creating SFX, post-production for filmmaking, and utilizing AI-generated silent videos.

2.2 Video-to-audio generation

Existing video-to-audio (V2A) generation methods have focused on achieving two main characteristics: (i) audiovisual relevance and (ii) temporal synchronization. The first stream aims to represent general sound by leveraging datasets such as VGGSound [2] and AudioSet [9]. Given a set of video features, SpecVQGAN [18] learns a transformer to sample quantized representations (*i.e.*, codebook) based on visual features to decode spectrogram. Im2wav [34] utilizes rich semantic representations obtained from a pre-trained CLIP [30] as sequential visual conditioning for an audio language model, and applies classifier-free guidance [15] to steer the generation process. Recently, diffusion-based models have shown the stunning ability to generate high-quality audio [27, 44]. DiffFoley [27] is a diffusion-based video-to-audio generation model that learns temporal and semantic alignment features through contrastive learning. Although it achieves better alignment between visual and audio inputs with prior training, it necessitates tremendous training data, such as the utilization of both

VGGSound and AudioSet for alignment training. Seeing-and-hearing [44] is another diffusion-based model that optimizes the generation process using ImageBind [11] which learns joint embedding space for six modalities (image, text, audio, depth, thermal, and IMU). However, they often struggle to generate temporally aligned sounds at short times in the video (*e.g.*, dog barking, people laughing).

On the other hand, other research works [4, 43] have focused on creating simplistic SFX (*e.g.*, stick hits) using datasets like CountixAV [47] and GreatestHits [28], which provide fewer classes but more precisely temporal aligned data. CondFoleyGen [7] trains a transformer model to autoregressively predict a sequence of audio codes for a spectrogram VQGAN, conditioned on the given audiovisual example. Syncfusion [4] predicts a discrete onset label that denotes the beginning of a sound for repetitive actions. Recent SonicVisionLM [43] employs a large language model (LLM) to utilize text as an intermediate product that facilitates user interaction for personalized sound generation. They freeze Tango [10] and train ControlNet with timestamp estimated by a video for 23 SFX categories exclusively, where a sound event timestamp detection module is trained on a dataset comprising videos and sound timing data. While these approaches are limited to generating only pre-defined training SFX categories, our method generates sounds for various categories from the visual context.

3 Preliminary

3.1 Text-to-audio latent diffusion model

In this paper, we specifically utilize AudioLDM [25] which generates a latent of mel-spectrogram z computed by VAE [22]. The diffusion model ϵ_θ of AudioLDM is trained to predict the noise added to a given data by minimizing the objective function, $\mathcal{L}_{\text{diff}} = \mathbb{E}_{z_0, \epsilon, t} \|\epsilon - \epsilon_\theta(z_t, t, \mathbf{E}_a)\|_2^2$, where ϵ represents the noise added at time t , z_t is noisy latent induced via the forward process and \mathbf{E}_a denotes the embedding of the audio x obtained from the CLAP audio encoder $f_{\text{audio}}(\cdot)$ [42]. Here, the model is conditioned by \mathbf{E}_a using classifier free guidance [15].

In the sampling process, the generation starts from a noise z_T sampled from $\mathcal{N}(0, I)$ and the text embedding \mathbf{E}_y from the CLAP text encoder $f_{\text{text}}(\cdot)$. The reverse process conditioned on \mathbf{E}_y gradually generates the audio prior z_0 using the modified noise estimation $\hat{\epsilon}_\theta(z_t, t, \mathbf{E}_y) = (1+w)\epsilon_\theta(z_t, t, \mathbf{E}_y) - w\epsilon_\theta(z_t, t)$, where $w \in [0, 1]$ is a guidance weight to balance the trade-off of the audio condition \mathbf{E}_a . The VAE decoder decodes the sampled latent z to predict a mel-spectrogram. Finally, the decoded mel-spectrogram is converted to a raw audio sample using the HiFi-GAN vocoder [23].

Although AudioLDM enables text-conditional audio generation, it still lacks of understanding of visual contents and their temporal information. In this paper, we add a visual control to the pre-trained AudioLDM. Instead of directly using a visual feature to control, we extract more essential information from the given video, which will be discussed in Section 3.2.

3.2 Rethinking video-to-sound generation with temporal alignment

We hypothesize a video input can bring principal information missing in a text prompt, *i.e.*, temporal alignment for a visual scene. Prior sound generation methods from visual content were immersed in the overall audiovisual relevance [3, 34, 38]. Since then, recent works [7, 4, 43] have attempted to generate SFX for video by estimating the onset [4], or audio timestamp [43] from videos. However, these methods produce an unnatural sound in that discrete conditions such as onset or timestamp cannot serve continuous sound variations.

In this work, we consider *energy*, interpreted as “loudness” or “dynamics”, to produce a continuous condition. Figure 2 shows that energy is a continuous time-varying signal, including envelope components of sound such as peak, attack, sustain, and decay. Energy can be obtained cheaply and automatically by computing the frame-level magnitude of mel-spectrograms [31]. Moreover, we empirically observe that energy can also implicitly imply multiple semantics of the video. For example, Figure 3 shows the generated audio by ReWaS contains multiple semantics beyond the text control “playing darts”; it generates people talking sounds for the intervals having less power than the sound of a dart hitting the target. In

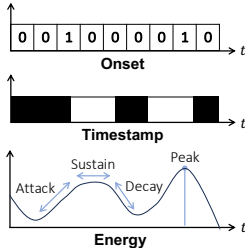


Figure 2: Limitation of timestamp annotations.

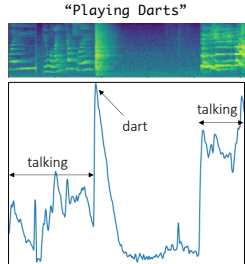


Figure 3: Energy can imply multiple semantics.

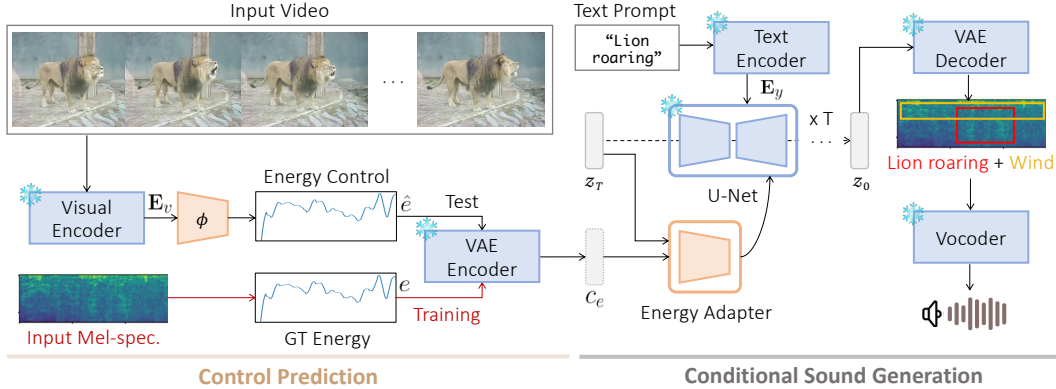


Figure 4: Overall architecture of ReWaS. Our model predicts energy control from a given video, and generates sound with text prompt and control condition. Red lines are used in training only.

other words, the energy information can imply the existence of different objects in the video as well as the main object.

4 Method

In this work, we introduce a novel sound generation method conditioned on both text and video, to generate a waveform that is temporally well aligned with the visual input. As shown in Figure 4, our model consists of two main parts: (i) *control prediction module*, which intermediately predicts energy control from the video. (Section 4.1) (ii) *conditional sound generation module*, which uses the energy control signal as a condition in the diffusion process to generate corresponding audio outputs (Section 4.2), which are both temporally and semantically aligned with text and video.

4.1 Control prediction from video

Energy control. ReWaS is based on a T2A generation method, specifically AudioLDM that uses CLAP embedding space for text and audio alignment. A naïve approach using video as a condition is to align latent space between audio-video-text. Luo et al. [27] attempted to align tri-modal embeddings in a unified space by large-scale contrastive learning prior to training diffusion models. To more efficiently overcome this challenge, we design an energy control as an intermediate bridge from video to audio. We speculate that energy control brings three advantages: First, the power of audio is intuitively correlated to visual dynamics and semantics [20, 37]. With the natural fact that people can imagine the power of sound from the size of the instance or distance to the object, we regard audio energy as a visually correlated signal that can be certainly obtained from video. Second, as shown in previous works [31, 13], energy plays as a structural condition for audio generation. Therefore, it is well-suited to parameter-efficient fine-tuning methods such as ControlNet [46]. Finally, using temporal acoustic signals for generating audio needs a skilled user to annotate the pitch, melody, or rhythm for every timestamp. It makes the audio generation phase impractical and difficult for the public to control. Meanwhile, energy is highly related to physical interactions implicated in visual signals; thus, it can be easily estimated from the video. Our approach does not require timestamp-wise fine-grained user control, but automatically estimating energy structure from the given video.

Video embedding. To predict the energy control from video input, we use the feature extracted from a pretrained SynchFormer [19] video encoder. We empirically observe that the image encoder (e.g., CLIP [30]) is limited to V2A generation, especially from a temporal alignment perspective. We finally take video embedding $\mathbf{E}_v \in \mathbb{R}^{S \times C}$, where S is the number of segments and C is the dimension of latent. The implementation details for this process are described in Appendix A.1.

Training energy control from video. Similar to Ren et al. [31], we calculate the energy from the mel-spectrogram by averaging the frequency bins and further smoothing the time-sequential energy information. We first transform the raw waveform to the mel-spectrogram, $\text{mel} \in \mathbb{R}^{D \times W}$, where D represents the number of mel-frequency bins, and W is the width of the spectrogram following AudioLDM [25]. However, we empirically observe that the computed energy fluctuates a lot for

each temporal window, which hinders stable training. We resolve the issue by taking a smoothing operator. The energy of audio $e \in \mathbb{R}^W$ is defined as $e_a = \text{Smoothing} \left(\frac{1}{D} \sum_{d=1}^D \text{mel}_{w,d} \right)$. We use the second-order Savitzky-Golay filter [40] with a window length of 9 for smoothing.

We estimate \hat{e} by using a shallow projection module ϕ from the video encoder output (See Figure 4 “Control Prediction”). For efficient training, we resize e_a by taking the nearest-neighbor interpolation to have the same number of segments S as the visual representations. We also can apply the same resize method to video embeddings at inference time. Now, we train our energy control prediction module ϕ by minimizing the following loss function $\mathcal{L}_e = \|\phi(\mathbf{E}_v) - \text{Resize}(e)\|_2^2$.

The output \hat{e} of the projection module is used for energy control at inference time. We train ϕ separately to diffusion models for training efficiency. In addition, our energy estimation module is not specialized for generation models, thus our energy control can be utilized in other ways.

4.2 Conditional sound generation

Adding control signal to diffusion model. To reflect the energy control signal, we train the energy adapter following the framework of ControlNet [46]. The weights of the energy adapter are initialized from pretrained parameters of diffusion models, and connected to AudioLDM with zero convolution layers. Compared to training audiovisual alignment into the latent space in diffusion model [27, 44], our adapter takes the benefit of robust fine-tuning speed (*e.g.*, [27] uses 8 A100 GPUs for 140 hours for feature alignment and LDM tuning, whereas ReWaS uses 4 V100 GPUs for total 33 hours). To add the control feature for z_t , the energy control e_a is duplicated by the number of mel-filterbanks, and transferred to the VAE encoder for the purpose of encoding, followed by a fully-connected layer and SiLU activation [8]. This latent control feature c_e is added to the z_0 , where z_0 is an audio prior obtained from the VAE encoder. Thus, given a text embedding \mathbf{E}_y and latent control feature c_e , we train energy adapter by optimizing the following objective: $\mathcal{L}_c = \mathbb{E}_{z_0, t, \mathbf{E}_y, c_e, \epsilon \sim \mathcal{N}(0,1)} \|\epsilon - \epsilon_\theta(z_t, t, \mathbf{E}_y, c_e)\|_2^2$. During training, we randomly drop \mathbf{E}_y with the probability 0.3 for better controls. We denote that \mathcal{L}_c and \mathcal{L}_e are optimized separately.

Sound generation. We employ DDIM sampler [36] to generate sound from the noise. The reverse sampling process is conditioned on both text prompt and video. We replace e to $\hat{e} = \phi(\mathbf{E}_v)$ at inference. Once the mel spectrogram is generated by the VAE decoder, it can be transformed into a raw waveform using a pretrained neural vocoder [23] as explained in Section 3.1.

5 Experiments

5.1 Experimental settings

Datasets. To ensure a fair comparison with existing baselines, we train both the control prediction module and the adapter in the conditional sound generation module on VGGSound [2]. VGGSound is a large-scale audiovisual dataset containing approximately 200k video clips of 10-second duration, accompanied by corresponding audio tracks. The dataset covers 309 classes of general sounds, roughly categorizing them into acoustic events, music, and people. The videos are sourced from YouTube, providing a diverse and realistic corpus. For our experiments, we leverage a subset of 160k videos from VGGSound due to the availability of public videos at the time of training. Since the VGGSound includes plentiful general sound examples, our method trained on the VGGSound enables general-purpose sound generation for real-world scenarios. We split the train data list into training and validation subsets following SpecVQGAN [18]. To evaluate temporal alignment accuracy, we evaluate ReWaS on Greatest Hits [28] test set including the videos of hitting a drumstick with materials. Since Greatest Hits samples have a distinct audio property compared to the other audio samples, we fine-tune our method on the Greatest Hits training samples. To the best of our knowledge, it is the first proposed evaluation protocol to evaluate the model on both of these datasets.

Evaluation protocol. In all experiments, we generate 5 seconds of audio corresponding to the input video clips and text prompt. We use sound class categories as the text prompt. We observe that the existing baselines used different numbers of generated samples and audio sequence time, which could lead to inconsistencies in the evaluation. For a fair comparison, we re-implement all baseline methods to generate one audio sample with 5 seconds from one video in VGGSound dataset, consistent with ReWaS. On the other hand, 2 seconds of generation results are used in Greatest Hits.

Table 2: Performance comparison on VGGSound [2] with reproduced five seconds audio samples.

Model	FID↓	IS↑	MKL↓	CLAP↑	Energy MAE↓	Training Params.↓
SpecVQGAN [18]	38.88	7.41	7.22	12.34	2.27	379M
Im2wav [34]	33.10	10.72	5.90	25.70	3.31	365M
Diff-Foley [27]	<u>30.53</u>	17.03	6.61	<u>38.11</u>	2.08	859M
ReWaS (Ours)	29.09	<u>15.96</u>	<u>6.37</u>	41.70	1.29	204M

Table 3: Performance comparison for sound generation on Greatest Hits [28]. All baselines are evaluated on re-produced results. † denotes the models trained on Greatest Hits. We use material types as text prompts, while CondFoleyGen uses both reference audio and video as inputs.

Model	FID↓	IS↑	MKL↓	Acc↑	AP↑	Energy MAE↓
SpecVQGAN [18]	42.72	1.57	6.13	18.09	62.78	4.29
Im2wav [34]	85.64	1.64	<u>4.83</u>	14.36	60.44	3.77
Diff-Foley [27]	86.40	1.74	5.07	23.94	<u>65.61</u>	3.40
CondFoleyGen [7]†	40.89	1.63	7.08	23.94	60.24	<u>2.94</u>
ReWaS (Onset)†	<u>39.96</u>	<u>3.09</u>	7.02	21.81	65.88	3.64
ReWaS (Energy)†	36.70	4.85	3.92	19.15	63.28	2.93

Baselines. We compare ReWaS against existing video-to-sound generation approaches in priority, SpecVQGAN [18], Im2wav [34], and Diff-Foley [27], which are trained on the VGGSound and AudioSet datasets. For a fair comparison, we take the following steps: We first generate the full-length audio by SpecVQGAN and Diff-Foley (10s and 8s, respectively) for the test videos. Then, we extract the 5-second clip corresponding to the same video frames used in our method. Since Im2wav is designed to generate sound with a fixed length of 4 seconds, we first generate the initial 4 seconds and extend it by generating an additional 1 second, resulting in a 5-second audio clip. In the temporal alignment evaluation, we consider CondFoleyGen [7] as a main baseline, which is trained on the Greatest Hits dataset.

Evaluation metrics. To evaluate perceptual quality and diversity, we employ Fréchet distance (FID) [14] and inception score (IS) [33]. IS assesses sample quality and diversity, while FID measures distribution-level similarity. We also report the mean of KL divergence (MKL) [18] for paired samples to evaluate the relevance to the condition. Following the implementation of SpecVQGAN [18], we use a Melception classifier to report those scores. We also measure the alignment between the generated audio and sound categories with CLAP score [17] in VGGSound. In the Greatest Hits experiment, we report onset accuracy (Acc) and average precision (AP), following the evaluation protocol introduced by CondFoleyGen [7]. The onset of sound events is a discrete signal obtained by the thresholding of the amplitude gradient. Therefore, it is challenging to detect onsets in natural sound like VGGSound benchmark. To address this issue, we report the mean absolute error (MAE), following the approach in [12], to evaluate the difference energy signal from real and generated sounds for the first time in the sound generation task conditioned on video. Although these evaluation metrics can evaluate different properties of the generated audio, all of them measure the difference between the generated audio and the “ground truth” audio corresponding to the original video. However, a sound is multimodal information; one video can sound differently (*e.g.*, the voice of a human can vary). Consequently, existing quantitative evaluation metrics have challenges in measuring whether the generated audio is truly suited to the given video. To tackle the issue, we conduct a user study to evaluate the quality of the generated audio samples. More details are explained in the next subsection.

5.2 Results

Quantitative results. Table 2 shows the quantitative comparisons on the VGGSound benchmark. We note that category classes are used as text prompts in the VGGSound. Compared to GAN-based SpecVQGAN and language model-based Im2wav, the models based on diffusion (Diff-Foley and ReWaS) show high-fidelity sounds (*i.e.*, better FID, IS and MKL). We train 22M parameters for video projection to audio conditional control, and 182M parameters for fine-tuning the AudioLDM with our energy adapter. Especially, while we use only a quarter of training parameters compared to Diff-Foley [27], our method outperforms Diff-Foley on all metrics except IS score. CLAP scores

illustrate the importance of text prompts for semantic alignment. Also, our method surprisingly outperforms in terms of MAE. This result demonstrates the accuracy of our control prediction module, and generated outputs from ReWaS are most temporally closer to the real content.

In addition, we evaluate how the generated audio and the condition video are temporally aligned on Greatest Hits. The dataset distribution of Greatest Hits highly differs from the general audio samples; hence, we fine-tune ReWaS on the Greatest Hits training samples. Table 3 shows the results. ReWaS achieves the best audio generation quality (FID, IS, MKL), and reasonable onset accuracy and AP, although ReWaS is not specially designed for Foley like Diff-Foley and CondFoleyGen. Finally, our method shows the best energy prediction score which forms the overall sound structure.

We also implement ReWaS with the onset signal, which is a discrete sequence to represent 0 or 1 (details can be found in Appendix A.1 and A.2). Compared to ReWaS with energy, our onset model shows more accurate results for temporal alignment metrics. It demonstrates that although the onset signal is advantageous in predicting the exact starting point of the sound, the energy is advantageous in terms of the overall audio quality and the homogeneity of the sound (*i.e.*, energy MAE).

User study. The quantitative results are limited to measuring how the generated audio sounds realistic and aligned to the given video. To complement the quantitative evaluations, we conduct a human evaluation study to assess the subjective quality of the generated audio concerning the input video. We ask the human evaluators to evaluate the quality of the audio samples generated by SpecVQGAN [18], Im2wav [34], Diff-Foley [27], and ReWaS.

Table 4: Human evaluation of V2A methods for comparison of audio quality, audiovisual relevance, and temporal alignment with 5-scale MOS.

Model	Audio Quality \uparrow	Relevance \uparrow	Temporal Alignment \uparrow
SpecVQGAN [18]	2.76	2.50	2.64
Im2wav [34]	2.97	3.18	3.01
Diff-Foley [27]	2.89	2.97	2.98
ReWaS (Ours)	3.70	4.04	3.68

We use three evaluation criteria: audio quality, relevance between audio and video, and temporal alignment. Detailed user instructions are depicted in Appendix A.3. We use a five-point Likert scale to measure mean opinion score (MOS), where an ideal video clip with its ideal audio receives a rating of 5 across all criteria. We recruit human annotators via two separate channels: Amazon Mechanical Turk (AMT) and local hiring. For AMT, we recruit 50 annotators for each criterion, and each annotator evaluates five generated samples for each method (*i.e.*, each annotator evaluates 20 generated audios). For locally hired annotators, we ask them to evaluate 20 generated samples for each method and criterion. There were 23 local annotators, and each of them evaluated 240 samples.

Table 4 shows that ReWaS achieves the best in all categories. This subjective human evaluation result is consistent with our quantitative and qualitative findings, further validating the effectiveness of our approach in generating high-quality, relevant, and temporally synchronized audio for the given video.

Qualitative results. Figure 5 shows qualitative results in baselines and ReWaS. Given the skateboarding video, SpecVQGAN and Diff-Foley fail to generate the sound of skate wheels rolling on the floor. Although Im2wav generates that sound, it cannot capture a short transition.

We also demonstrate the effectiveness of the text prompt in Figure 6 with CLAP similarity. In the input video shown in Fig. 6 (a), rain streaks are barely visible, while we want to hear the sound of rain. ReWaS can emphasize the desired sound with the help of a text prompt ‘*raining*’.

Meanwhile, videos in the wild may include redundant frames as shown in Fig. 6 (b). In this case, V2A methods also struggle to generate corresponding sound. However, ReWaS can effectively calibrate the semantics by user prompt. More examples can be found on our demo page.

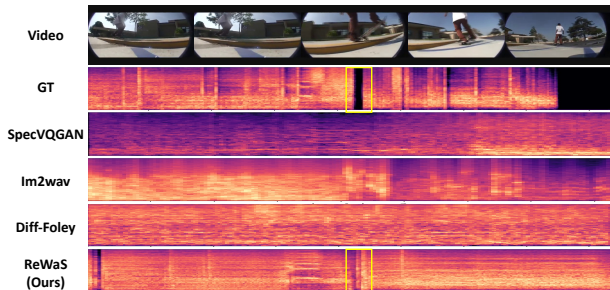


Figure 5: ReWaS successfully generates corresponding audio compared to reference (ground-truth) audio. Surprisingly, when the skateboarder jumps into the air, only ReWaS succeeded in detecting short transition moments (yellow box). The input text prompt is “skateboarding”.

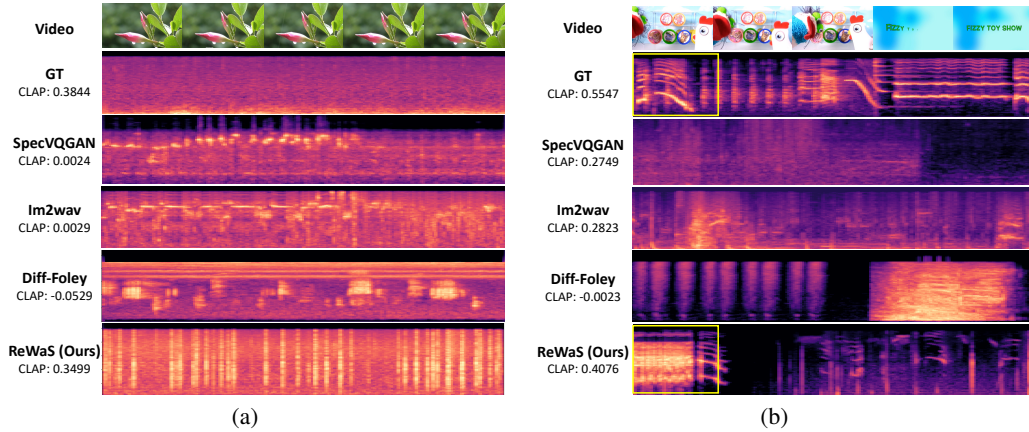


Figure 6: Effectiveness of text prompt. Videos in the real world are sometimes noisy. For example, when videos (a) are hard to distinguish the semantics or (b) contain redundant frames, text prompts used in ReWaS calibrate the results. Text prompt in (a) is “raining”, and (b) is “chicken clucking”.

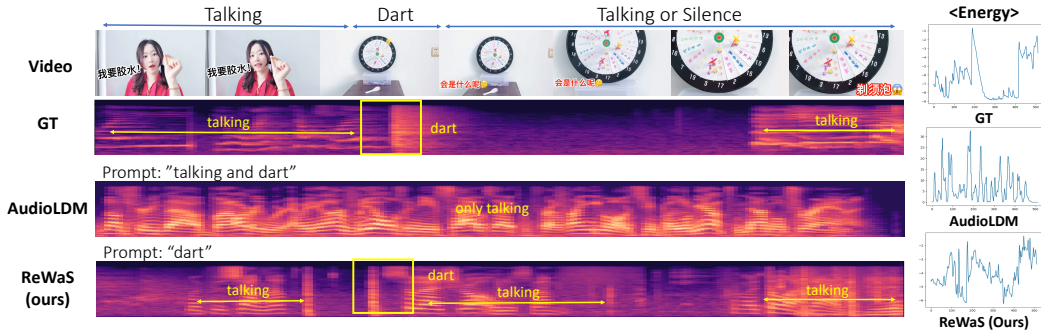


Figure 7: Effectiveness of video input. In ReWaS, multiple semantics in the video are transferred through energy control. Even if there is missing information in the text prompt, the energy control complements this.

5.3 Discussion and limitation

The impact of the quality of the energy control. ReWaS estimates the audio energy from the video. However, the estimated energy would propagate an estimation error, which may cause a significant performance drop. To verify the robustness of the energy prediction module, we compare the control by our video-to-energy prediction module and the energy directly extracted from the ground truth audio. Table 5 shows that although we use the estimated energy, the quality of the generated audio is very similar to the audio samples controlled by the ground truth audio energy. It supports the idea that energy information is highly related to visual information and is easy to estimate solely using video.

Effectiveness of video input. Figure 7 illustrates examples including multiple semantics. For instance, when a person talking and playing a dart game in an input video, the original AudioLDM generates only the sound of talking, ignoring ‘dart’ prompt. Additionally, aligning generated sound with video is challenging in AudioLDM. In comparison, ReWaS not only generates both the sound of talking and dart but also aligns the sound with the frames. Since the energy encompasses not just volumes but also the power of acoustical features (frequency), this supports the notion that multiple semantics can also be conveyed by our energy condition. Furthermore, the result demonstrates the limitation of T2A methods for automatic Foley synthesis, because they cannot watch a video. Another example is in Figure A.5.

Table 5: Impact of the energy control’s quality. Numbers are measured on VGGSound.

Control	FAD↓	IS↑	KL↓	CLAP↑	MAE ↓
Text & the ground-truth audio energy (upper bound)	29.06	18.21	6.12	44.97	0.95
Text & the estimated energy from the video	29.09	15.96	6.37	41.70	1.29

Limitation. Although our approach successfully leverages the text and video control simultaneously, our method shares the limitation of AudioLDM, namely, hallucination in generated samples. For example, for a given “basketball bounce” video, ReWaS generates a squeaking sound, even if the player is standing still. This problem might be mitigated if we can use a better AudioLDM model.

6 Conclusion

This paper proposes ReWaS, a novel video-and-text-to-sound generation framework. Our main contribution is that audio structural condition, namely energy, is inferred from video to efficiently and effectively input visual condition to the robust T2A model. Therefore, ReWaS can generate complex sounds in the real world without the need for a difficult control design. Quantitative results on VGGSound and Greatest Hits datasets, subjective human study, and qualitative results consistently support that ReWaS can generate natural, temporally well-aligned, and relevant audio for the given video by employing text and video as control.

References

- [1] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 1
- [2] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 2, 3, 6, 7, 15
- [3] Lele Chen, Sudhanshu Srivastava, Zhiyao Duan, and Chenliang Xu. Deep cross-modal audio-visual generation. In *Thematic Workshops of ACM MM*, 2017. 4
- [4] Marco Comunità, Riccardo F Gramaccioni, Emilian Postolache, Emanuele Rodolà, Danilo Comminiello, and Joshua D Reiss. Syncfusion: Multimodal onset-synchronized video-to-audio foley synthesis. In *ICASSP*, 2024. 2, 3, 4, 13
- [5] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9268–9277, 2019. 13
- [6] Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. In *AAAI*, 2018. 3
- [7] Yuexi Du, Ziyang Chen, Justin Salamon, Bryan Russell, and Andrew Owens. Conditional generation of audio from video via foley analogies. In *CVPR*, 2023. 2, 3, 4, 7
- [8] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 2018. 6
- [9] Jort F Gemmeke, Daniel PW 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 *ICASSP*, 2017. 3
- [10] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction tuned llm and latent diffusion model. *arXiv preprint arXiv:2304.13731*, 2023. 2, 3, 4
- [11] 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. 4
- [12] Zhifang Guo, Jianguo Mao, Rui Tao, Long Yan, Kazushige Ouchi, Hong Liu, and Xiangdong Wang. Audio generation with multiple conditional diffusion model. In *AAAI*, 2024. 2, 3, 7
- [13] Zhifang Guo, Jianguo Mao, Rui Tao, Long Yan, Kazushige Ouchi, Hong Liu, and Xiangdong Wang. Audio generation with multiple conditional diffusion model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18153–18161, 2024. 5
- [14] 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 *NeurIPS*, 2017. 7

- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3, 4
- [16] Jiawei Huang, Yi Ren, Rongjie Huang, Dongchao Yang, Zhenhui Ye, Chen Zhang, Jinglin Liu, Xiang Yin, Zejun Ma, and Zhou Zhao. Make-an-audio 2: Temporal-enhanced text-to-audio generation. *arXiv preprint arXiv:2305.18474*, 2023. 2, 3
- [17] 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. 2, 3, 7, 14
- [18] Vladimir Iashin and Esa Rahtu. Taming visually guided sound generation. *arXiv preprint arXiv:2110.08791*, 2021. 2, 3, 6, 7, 8
- [19] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Synchformer: Efficient synchronization from sparse cues. *arXiv preprint arXiv:2401.16423*, 2024. 5, 13
- [20] 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 *ICCV*, 2023. 2, 5
- [21] Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-tts: A generative flow for text-to-speech via monotonic alignment search. In *NeurIPS*, 2020. 3
- [22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 4
- [23] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. In *NeurIPS*, 2020. 4, 6
- [24] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. In *ICLR*, 2023. 3
- [25] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023. 2, 3, 4, 5, 14
- [26] Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D. Plumbley. AudioLDM 2: Learning holistic audio generation with self-supervised pretraining. *arXiv preprint arXiv:2308.05734*, 2023. 2, 3
- [27] Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. In *NeurIPS*, 2024. 2, 3, 5, 6, 7, 8
- [28] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *CVPR*, 2016. 2, 4, 6, 7
- [29] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra, Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, and Joao F Henriques. Keeping your eye on the ball: Trajectory attention in video transformers. In *NeurIPS*, 2021. 13
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 3, 5
- [31] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. FastSpeech 2: Fast and high-quality end-to-end text to speech. In *ICLR*, 2020. 4, 5
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3
- [33] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *NeurIPS*, 2016. 7
- [34] Roy Sheffer and Yossi Adi. I hear your true colors: Image guided audio generation. In *ICASSP*, 2023. 2, 3, 4, 7, 8
- [35] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. In *ICLR*, 2022. 1
- [36] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2020. 6

- [37] 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 *CVPR*, 2023. 5
- [38] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. In *NeurIPS*, 2024. 4
- [39] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *NeurIPS*, 2017. 3
- [40] Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 2020. 6
- [41] Shih-Lun Wu, Chris Donahue, Shinji Watanabe, and Nicholas J Bryan. Music controlnet: Multiple time-varying controls for music generation. *arXiv preprint arXiv:2311.07069*, 2023. 2, 3
- [42] Yusong Wu*, Ke Chen*, Tianyu Zhang*, Yuchen Hui*, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP*, 2023. 4
- [43] Zhifeng Xie, Shengye Yu, Mengtian Li, Qile He, Chaofeng Chen, and Yu-Gang Jiang. Sonicvisionlm: Playing sound with vision language models. *arXiv preprint arXiv:2401.04394*, 2024. 2, 3, 4
- [44] Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners. In *CVPR*, 2024. 2, 3, 4, 6
- [45] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023. 3
- [46] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 2, 3, 5, 6
- [47] Yunhua Zhang, Ling Shao, and Cees GM Snoek. Repetitive activity counting by sight and sound. In *CVPR*, 2021. 4

A Appendix

A.1 Data preprocessing and feature extraction

During training, we randomly extract 5-second segments from the VGGSound dataset and 2-second segments from the Greatest Hits dataset. However, during the testing phase, we extract video clips ranging from 2 to 7 seconds in duration for the VGGSound dataset, and from 0 to 2 seconds for the Greatest Hits dataset. Video frames are uniformly sampled at 25 fps. Since ReWaS generates audio based on 5-second videos, we duplicated frames from the Greatest Hits dataset to match the length of these 5-second videos. Subsequently, we trimmed the generated audio to a duration of 2 seconds.

We employ SynchFormer [19] trained on VGGSound for the sparse synchronized setting as a video encoder. The video encoder employed in SynchFormer is based on Motionformer [29] pre-trained on Something-Something v2, and fine-tuned on VGGSound and AudioSet. Therefore, the video encoder is strong enough to encode motion dynamics and semantics. We freeze the parameters in the video encoder, and solely train a projection module to estimate energy control. We extract a video feature in the short video clip (0.64 sec). Thus we use a total of 112 length visual embeddings for a 5s video. We note that, for a fair comparison, RGB frames are only used in all methods including ReWaS.

Audios of all videos used in our experiments are resampled to 16kHz sampling rate. We follow the default setting of AudioLDM to compute the mel-spectrogram. Specifically, we use 64-bin mel-spectrograms with 1024 window length. While f_{\min} and f_{\max} are 0 and 8000 respectively, the hop size is 160 and the FFT size is 1024.

To train ReWaS with the onset signal, we follow mostly the same settings as Synchfusion [4]. Specifically, frames are sampled at 15 fps, and resized to 112x112 without image augmentation, resulting in a total of 30 frames per 2 seconds clip. We resample audios to a 16 kHz sampling rate. Using the open-source library Librosa³, we detect the onset signal and interpolate it to 30 frames using nearest-neighbor interpolation to match the frame length.

A.2 Architecture and training details

Energy signal. To encode a video feature into 1-dimensional energy, a projection module ϕ consists of a linear layer, two transformer blocks, and MLPs consisting of four FC layers. We use 768 hidden dimensions for the first linear layer and transformer blocks, and the four FC layers' output dimensions are 128, 64, 16, and 1. The total parameter of ϕ is 22M. We choose AudioLDM-M⁴, and the number of training parameters for fine-tuning AudioLDM with our energy adapter is 182M. ReWaS is optimized by AdamW and the learning rate is fixed to 3e-5 during training. We train ReWaS with 4 V100 GPUs for 33 hours on VGGSound, and 1 hour on GreatestHits respectively.

Onset signal. We re-implement the video onset detection module from Synchfusion [4] as a projection module ϕ , consisting of a ResNet(2+1)D-18 network followed by a few fully-connected layers. The module generates binary labels using a sigmoid function and applies a threshold of 0.5. Finally, it produces onset signals (0 or 1) for each video frame. To train ϕ , we adopt the same training settings as Synchfusion. Specifically, we train the video onset detection module for 100 epochs with a batch size of 16, using the AdamW optimizer with a weight decay of 10^{-3} and a learning rate of 10^{-4} . Due to the imbalance of onset sound in the data, we utilize a weighted binary cross-entropy loss [5] with a weight of 10. The onset adapter fine-tuned using AudioLDM-M has the same architecture as the energy adapter, and is trained on Greatest Hits.

A.3 User study

Figure A.1, Figure A.2, and Figure A.3 show the user instructions used in our human evaluation. Before launching Amazon MTurk (AMT), we first conducted an in-lab study with 23 participants; each participant evaluated 20 audio samples for each method and each criterion, namely, they evaluated 240 ($20 \times 4 \times 3$) generated audio samples. Based on the observation from the in-lab study, we have set the compensation level for each HIT to \$0.45 so that a worker can earn \$15 per hour. At the same time, we observed that a number of participants had trouble keeping focus on the evaluation

³Module available at <https://librosa.org/doc/>

⁴weights in <https://github.com/haoheliu/AudioLDM>

with 240 samples (each sample takes five seconds). To prevent the low-quality responses from MTurk annotators, we split each evaluation Human Intelligence Task (HIT) on a smaller scale. Each AMT annotator evaluates five audio samples for each method and one additional ground truth audio to prevent random guessing. We published 50 HITs for each criterion, and 150 responses were collected. Finally, we observe that many AMT annotators consistently score high for all questions (*e.g.*, 4 or 5). To ignore noisy responses, we omit responses having an average score larger than 4.0 for 21 questions. 55 responses were omitted after this filtering process.

Instruction 1

How natural is this audio recording?

Please focus on examining the audio quality and naturalness (noise, timbre, sound clarity, and high-frequency details).

1. Listen to the sample (Click ****Play**** button to listen audio samples)
2. Select an option
 - Excellent: 5 (Completely natural audio)
 - Good: 4 (Mostly natural audio)
 - Fair: 3 (Equally natural and unnatural audio)
 - Poor: 2 (Mostly unnatural audio)
 - Bad: 1 (Completely unnatural audio)

Figure A.1: User instruction for audio quality (naturalness) test.

A.4 T2A framework

We conduct an ablation study on the T2A model to demonstrate the flexibility of our approach. A key advantage of our method is its modular design, allowing seamless substitution of any T2A diffusion framework. To validate this, we replace the AudioLDM-M [25] backbone with Make-An-Audio [17]. The video encoder used in Make-An-Audio is re-trained to predict the appropriate energy scale of mel-spectrogram, which is configured with 80 frequency bins and a hop size of 256 samples, different from

Instruction 2

How much is the sound related to the object or material in video?

Please focus on examining the relevance between video and audio, not considering the quality and temporal alignment (*i.e.* sound timing).

1. Watch the sample (Click ****Play**** button to watch video samples)
2. Select an option
 - Excellent: 5 (Completely relevant audio)
 - Good: 4 (Mostly relevant audio)
 - Fair: 3 (Equally relevant and irrelevant audio)
 - Poor: 2 (Mostly irrelevant audio)
 - Bad: 1 (Completely irrelevant audio)

Figure A.2: User instruction for video-audio relevance test

Instruction 3

How much is the sound temporally aligned to the movements of objects or material in video?

Please focus on examining the temporal alignment between video and audio, not considering audio quality and naturalness.

1. Watch the sample (Click ****Play**** button to watch video samples)
2. Select an option
 - Excellent: 5 (Completely aligned audio)
 - Good: 4 (Mostly aligned audio)
 - Fair: 3 (Equally aligned and non-aligned audio)
 - Poor: 2 (Mostly non-aligned audio)
 - Bad: 1 (Completely non-aligned audio)

Figure A.3: User instruction for temporal alignment test.

Table A.1: Performance comparison on VGGSound [2] with different T2A framework.

Model	FID↓	IS↑	MKL↓	CLAP↑	Energy MAE↓	Training Params.↓
ReWaS (Make-An-Audio)	34.55	6.35	4.76	42.37	0.75	74.3M
ReWaS (AudioLDM)	29.09	15.96	6.37	41.70	1.29	204M

the AudioLDM-M configuration. Make-An-Audio is notable for its parameter efficiency, requiring significantly fewer parameters than AudioLDM. This reduction in model complexity translates to substantially shorter training times, with the entire model converging in less than one day. Table A.1 presents a quantitative comparison between the two backbones on VGGSound.

Interestingly, ReWaS built upon Make-An-Audio achieves performance comparable to its AudioLDM-M-based counterpart across most metrics. This underscores the robustness of our framework, demonstrating its ability to maintain high performance even when integrated with a more compact backbone. However, it is worth noting that the Make-An-Audio variant exhibits a slight degradation in overall audio quality such as FID and IS scores. This discrepancy can be attributed to the inherent limitations of the Make-An-Audio architecture, which is computationally efficient but often fails to generate high-quality audio over large-scale T2A frameworks. In other words, the development of T2A models could advance our model.

A.5 More qualitative results

We illustrate estimated energy from video in Figure A.4. The results show the correlation between our energy control generated from video and GT energy obtained from reference audio.

We provide more examples to show the effectiveness of ReWaS qualitatively. Figure A.5 is another qualitative example that demonstrates the ReWaS can generate multiple semantics without providing the context in textual prompts. In Figure A.6, ReWaS successfully generates "silent" audio sounds with visual input. In contrast, other baseline models fail to produce the corresponding sound (e.g., alarm clock ringing) due to misaligned visual and sound contexts, often generating unintended sounds or remaining silent when they should produce sound. This suggests that baseline models may resort to generating random sounds when faced with a change in frame scene since the visual context and sound context are misaligned but they also generate the sound where they have to be silent. Figure A.7 shows the effectiveness of ReWaS in aligning sound with corresponding frames, particularly in scenarios where man-made video-audio data lacks correlation due to the sound being added later to the frames. ReWaS exhibits superior performance in generating suitable audio sounds with better temporal synchronization accurately capturing small movements of objects (e.g., lip synchronization). For example, the energy of the growling sound becomes bigger as the lion opens

its mouth. Lastly, Figure A.8 illustrates ReWaS’s ability to generate temporally synchronized sounds, further highlighting its effectiveness.

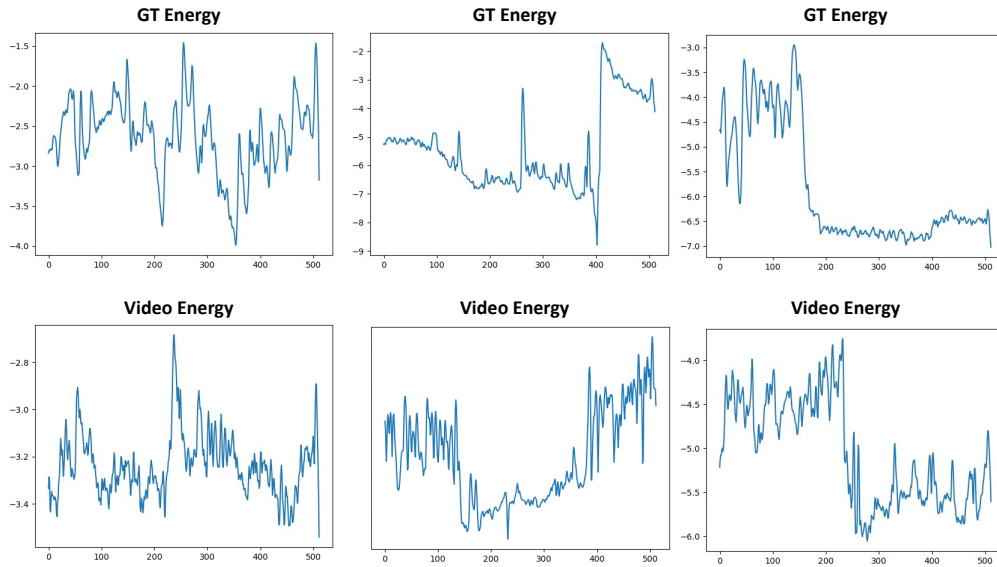


Figure A.4: Examples of energy controls from input videos.

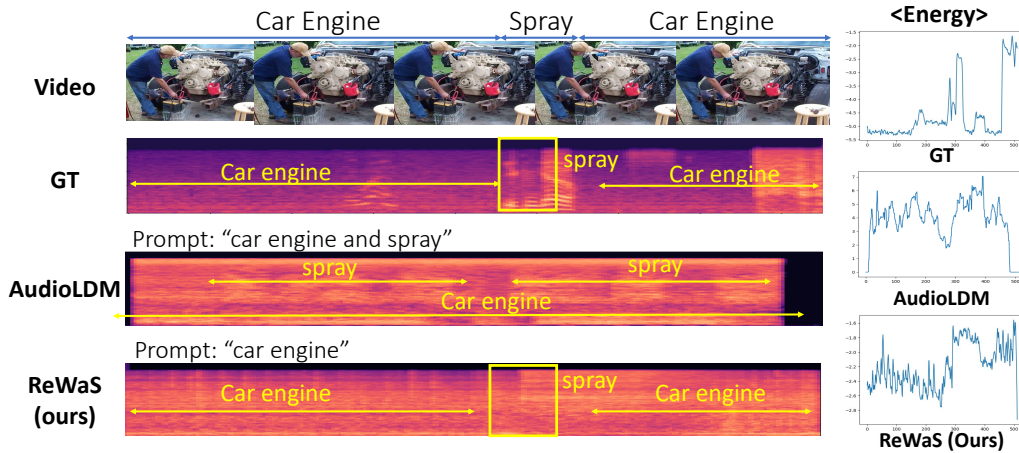


Figure A.5: Additional example of demonstrating the effectiveness of video input. ReWaS can make the sound missing in the text prompt (‘spray’), while simultaneously aligning multiple sounds with various semantics to the video.

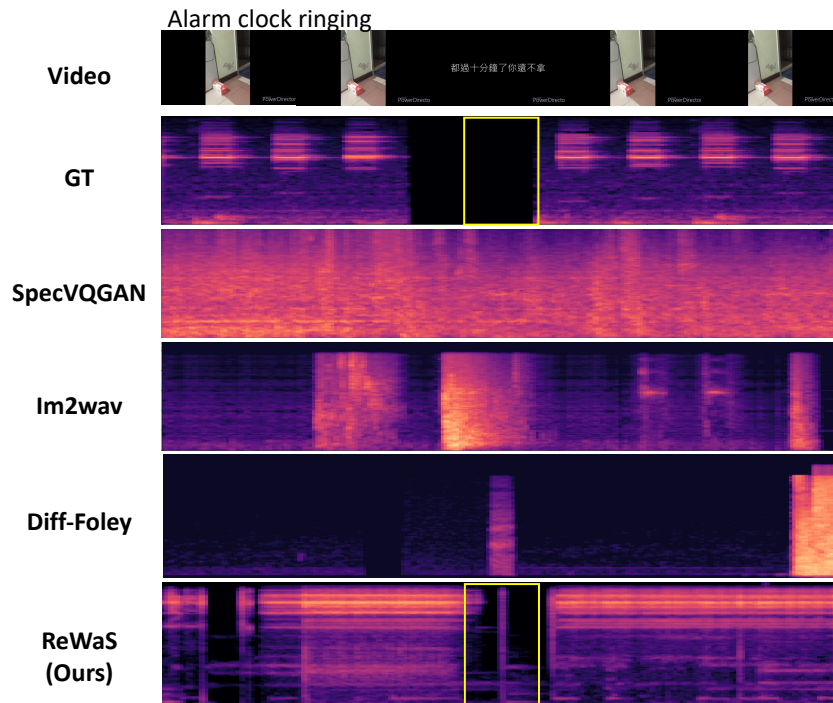


Figure A.6: Example of audio sound from misaligned visual input. ReWaS can make the desired sound and the silent moment like ground-truth sound.

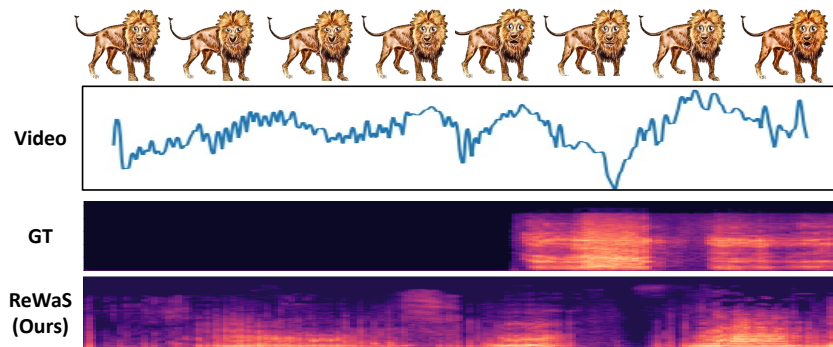


Figure A.7: Example of audio with improved synchronization, capturing small movements (e.g., a lion's lip synchronization).

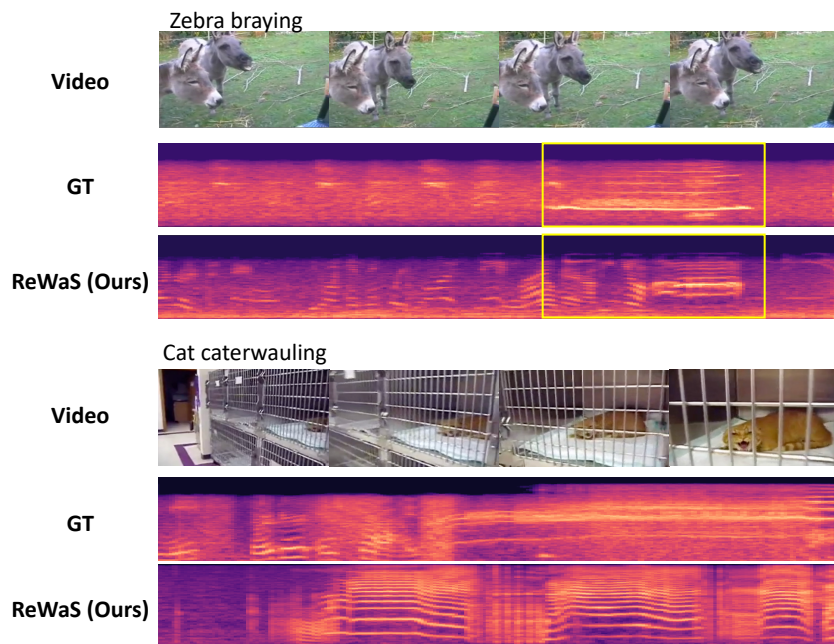


Figure A.8: Examples of generated audio sounds demonstrating the capability of temporal synchronization.