

# Kling-Foley: Multimodal Diffusion Transformer for High-Quality Video-to-Audio Generation

Jun Wang\*, Xijuan Zeng\*, Chunyu Qiang\*, Ruilong Chen, Shiyao Wang, Le Wang, Wangjing Zhou, Pengfei Cai, Jiahui Zhao, Nan Li, Zihan Li, Yuzhe Liang, Xiaopeng Wang, Haorui Zheng, Ming Wen, Kang Yin, Yiran Wang, Nan Li, Feng Deng, Liang Dong, Chen Zhang, Di Zhang, Kun Gai

Kuaishou Technology, Beijing, China

{wangjun06,zengxijuan,qiangchunyu}@kuaishou.com

## Abstract

We propose Kling-Foley, a large-scale multimodal Video-to-Audio generation model that synthesizes high-quality audio synchronized with video content. In Kling-Foley, we introduce multimodal diffusion transformers to model the interactions between video, audio, and text modalities, and combine it with a visual semantic representation module and an audio-visual synchronization module to enhance alignment capabilities. Specifically, these modules align video conditions with latent audio elements at the frame level, thereby improving semantic alignment and audio-visual synchronization. Together with text conditions, this integrated approach enables precise generation of video-matching sound effects. In addition, we propose a universal latent audio codec that can achieve high-quality modeling in various scenarios such as sound effects, speech, singing, and music. We employ a stereo rendering method that imbues synthesized audio with a spatial presence. At the same time, in order to make up for the incomplete types and annotations of the open-source benchmark, we also open-source an industrial-level benchmark Kling-Audio-Eval. Our experiments show that Kling-Foley trained with the flow matching objective achieves new audio-visual SOTA performance among public models in terms of distribution matching, semantic alignment, temporal alignment and audio quality. Homepage is available at: <https://klingfoley.github.io/Kling-Foley/>

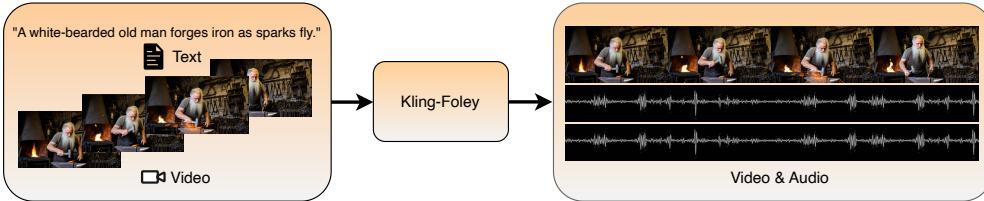


Figure 1: We propose Kling-Foley, a large-scale multimodal Video-to-Audio generation model. Taking an input video and an optional text prompt, the model synthesizes high-fidelity audio that is semantically aligned and temporally synchronized with the video content, encompassing elements such as sound effects and background music. Significantly, Kling-Foley can produce audio sequences of arbitrary duration, dynamically adapting to the length of the input video.

\*Equal contributions.

## 1 Introduction

Video generation has emerged as a focal point in generative AI research, with many models producing visually convincing results [1, 2, 3, 4]. However, these generated videos remain silent, requiring separate post-production audio dubbing to align with human perception and real-world expectations. Manual dubbing proves costly, inefficient, and heavily reliant on specialized expertise.

To automate this process, Text-to-Audio (TTA) models generate non-speech audio [5, 6, 7, 8, 9], such as sound effects and background music from text descriptions. The core objective of TTA involves translating natural language into high-fidelity, semantically aligned audio signals. However, TTA faces critical limitations: since it processes text alone, the generated audio often not temporally or semantically aligned with the video content. For instance, the footsteps may not match the shoe movements (temporal misalignment), or a "rain" description might yield light rain sounds while the video shows a downpour (semantic misalignment).

To resolve these issues, Video-to-Audio (V2A) models leverage both text descriptions and visual input[10, 11, 12, 13, 14]. Unlike TTA frameworks, V2A systems typically extract frame-level video features via pre-trained encoders to supplement textual inputs. While many V2A approaches extend pre-trained TTA models, their performance suffers due to the scarcity of high-quality, triple-modality (video-audio-text) datasets. Consequently, existing models often exhibit poor temporal/semantic alignment and low-fidelity audio output. MM-Audio[15], the current state-of-the-art (SOTA) model, addresses these weaknesses via a multimodal joint-training paradigm. Nevertheless, its reliance on limited open-source datasets restricts it to ambient sound and sound effects, excluding capabilities like background music generation. At the same time, it can only support fixed lengths, discarding those shorter than the fixed length, and splitting longer video clips, which is not conducive to training the model for video dubbing from videos with natural duration distributions.

To address these challenges, we propose Kling-Foley, a large-scale multimodal V2A generation model that generates high-quality audio synchronized with video content. Our method is based on multimodal diffusion transformer, using large-scale variable-length multimodal data, combining video semantic alignment, audio-visual synchronization alignment, and text conditions to achieve accurate joint control modeling of multi-scenario sound effects.

Critically, sound effect generation tasks lack an available multimodal benchmark that includes visual, auditory, and caption information. Most existing benchmarks support only partial modalities. For example, AudioSet [16] provides audio with category labels, while VGGSound [17] includes visual, audio, and label information. These datasets are not sufficient for a comprehensive evaluation of sound generation models. To address this gap, we construct a high-quality and reliable benchmark dataset named Kling-Audio-Eval, based on a self-developed labeling system and a rigorous manual annotation process. To the best of our knowledge, Kling-Audio-Eval is the first benchmark for sound effect generation that includes video, video caption, audio, audio caption, and sound event labels. The dataset contains 20,935 manually annotated samples and covers nine major sound event scenarios, such as traffic, human sounds, animal sounds, and more.

In addition, we evaluate model performance from multiple perspectives by conducting assessments across four dimensions: Distribution Matching, Audio Quality, Semantic Alignment, and Temporal Alignment. This multi-dimensional approach provides a comprehensive understanding of the model's capabilities.

Our key contributions are as follows:

- We introduce Kling-Foley, a novel V2A framework that generates high-fidelity audio perfectly synchronized with video content, achieving SOTA performance across audio quality, semantic alignment, and audiovisual synchronization metrics.
- We propose a visual semantic representation module and an audio-visual synchronization module that align video features with audio representations at each frame, enhancing semantic relevance and temporal coherence alongside text conditioning.
- We design a universal latent audio codec enabling high-fidelity modeling of various scenarios, including sound effect, speech, singing, and music.
- We release Kling-Audio-Eval, the first industry-scale multimodal benchmark with synchronized video, text descriptions, and audio featuring 20,935 manually annotated samples across

nine sound scenarios, addressing the limitations of existing datasets through comprehensive multimodal annotations and scenario coverage.

## 2 Related Work

### 2.1 Text-to-Audio

Early approaches to audio generation leveraged GANs [18], normalizing flows [19], and VAEs [20]. Recent models largely follow two architectural paradigms: Transformer-based and Diffusion-based.

Transformer-based models like AudioGen [5] use autoregressive transformer to predict discrete audio tokens, while MAGNET [21] introduces masked generative modeling for non-autoregressive generation.

Diffusion-based methods such as DiffSound [6] decode text into mel-spectrogram tokens via diffusion. Latent diffusion techniques [22] further drive models like AudioLDM [23], AudioLDM2 [7], Tango [8], and Make-An-Audio [24]. These models commonly rely on shared audio-text embedding spaces (e.g., CLAP [25]) and spectrogram autoencoders, with strategies such as instruction tuning [26] or pseudo prompt enhancement to improve data efficiency.

### 2.2 Video-to-Audio

V2A generation aims to synthesize sound from silent videos. Early work like [14] used SampleRNN [27] to directly generate waveforms from frames. Later methods, such as SpecVQGAN [28] and Im2Wav [29], use visual features (e.g., CLIP [30]) to condition transformer models. Diff-Foley [31] enhances temporal alignment via large-scale audio-visual pretraining and latent diffusion.

However, audio-visual datasets like VGGSound [32] are relatively small (550 hours), limiting model scalability. To address this, recent works extend pretrained TTA models for video input [13, 33, 10, 11, 12, 34]. For example, Seeing and Hearing [10] uses ImageBind [35] to convert video into text for AudioLDM, while V2A-Mapper [11] maps visual features to CLAP embeddings. To incorporate temporal dynamics, models like SonicVLM [12], ReWaS [13], and FoleyCrafter [36] integrate time-aware control modules.

Recent advancements in V2A generation have emphasized multimodal alignment and cross-modal contrastive learning to improve semantic relevance and temporal synchronization. VATT (Video-and-Text-to-Audio Transformer) [37] introduces a convolution-free Transformer architecture that jointly learns representations from video, audio, and text modalities. In parallel, VTA-LDM (Video-to-Audio Latent Diffusion Model) [38] employs a CLIP-based vision encoder to extract high-resolution visual features and integrates auxiliary embeddings (e.g., positional embeddings, textual prompts) to guide audio generation. V-AURA [39] introduces a cross-modal feature fusion strategy in an autoregressive framework, leveraging a high-framerate visual feature extractor (6x higher than prior work) to capture fine-grained motion details. For efficiency, FRIEREN [40] employs rectified flow matching (RFM) with reflow and one-step distillation. MMAudio [15] proposes a joint training paradigm combining audio-visual (VGGSound) and audio-text (WavCaps) data within a unified multimodal transformer.

### 2.3 Audio Representation

The masked-based self-supervised learning paradigm, validated in natural language processing (NLP) with BERT [41] and in computer vision with MAE [42], has become a mainstream approach for audio representation learning. Most models in this domain operate on audio spectrograms. Early works such as SSAST [43] pioneered the framework of applying masked modeling to spectrogram patches. This approach was further developed by models like Audio-MAE [44] and MaskSpec [45], which are heavily inspired by MAE and employ high-ratio random masking with an asymmetric encoder-decoder architecture to reconstruct original spectral features.

Current reconstruction-based audio representation learning methods can be divided into discrete and continuous types. Discrete representations, commonly produced via residual vector quantization, compress audio into discrete tokens and reconstruct them with a decoder. Models such as EnCodec [46], DAC [47] and AudioDec [48] follow this paradigm. Continuous representations do not rely on quantization. Instead, they learn smooth latent spaces through models like AudioLDM [49, 50],

Tango [51, 52], Make-An-Audio [53, 54], and DiffRhythm [55], which use variational encoders to capture audio features.

Contrastive learning has become a pivotal method in cross-modal representation learning, with vision models like CLIP [56], Florence [57], and ALIGN [58] effectively aligning image and text in a shared semantic space. In the audio domain, similar frameworks such as Wav2CLIP [59], AudioCLIP [60], and CLAP [61] have achieved strong performance by learning global semantic representations of audio.

### 3 Preliminary

#### 3.1 Conditioning Encoders

Building upon the original SD3 framework [62], we introduce several targeted improvements to enhance both semantic representation and temporal alignment across modalities for multimodal audio generation:

**Text Encoder.** We adopt the T5-Base model proposed by Raffel et al. [63] as the backbone for textual representation. T5 (Text-to-Text Transfer Transformer) introduces a unified framework that reformulates all NLP tasks into a text-to-text paradigm, allowing for a highly flexible and task-agnostic approach to model training and inference. Pretrained on the Colossal Clean Crawled Corpus (C4), T5-Base benefits from exposure to vast and diverse language data, enabling it to capture nuanced semantic relationships and contextual dependencies across a broad range of domains. In our system, the T5-Base functions as the primary text encoder, transforming natural language inputs—such as prompts, descriptions, or queries—into rich latent embeddings. These embeddings serve as a semantic anchor that guides downstream multimodal alignment and generation processes. Furthermore, the model’s ability to generalize well across tasks allows us to maintain high adaptability and performance with minimal task-specific tuning.

**Vision Encoder.** For the extraction of high-level visual semantics and effective multimodal grounding, we incorporate the ViT-bigG-14-QuickGELU model within the MetaCLIP framework, as introduced by Ma et al. [64]. ViT-bigG, a large Vision Transformer architecture, utilizes QuickGELU activations and benefits from extensive pretraining on large-scale, high-quality image-text pairs. The MetaCLIP training strategy enhances this foundation by employing an expert clustering mechanism to better capture the alignment between visual and linguistic modalities. This results in a visual encoder that excels at producing domain-robust, semantically meaningful image embeddings. These embeddings not only preserve fine-grained visual details but also remain aligned with their textual counterparts, making the model particularly suitable for tasks involving vision-language fusion, retrieval, and generation. In our multimodal pipeline, the visual encoder plays a critical role in grounding generated content in the visual domain, ensuring that the output remains coherent and contextually relevant.

**Align Encoder.** Temporal consistency across modalities is crucial in many real-world multimodal applications, such as video generation, dubbing, or speech-driven animation. To address this, we utilize the Synchformer model proposed by Iashin et al. [65], a transformer-based architecture explicitly designed for audio-visual synchronization. Synchformer leverages sparse synchronization cues—like lip movements, phoneme timing, and audio features - to infer fine-grained alignment across the temporal dimension. Unlike traditional alignment methods that may rely on dense supervision or heuristic rules, Synchformer employs self-attention and cross-modal transformer to model temporal dependencies efficiently. Its compact yet expressive design enables it to operate robustly in both constrained and unconstrained environments, ensuring that multimodal outputs (e.g., lip-synced avatars or synchronized video narration) maintain a high level of realism and temporal coherence. In our system, the align encoder acts as a mediator that refines and aligns latent representations from the text and vision encoders, ensuring consistent temporal flow and cross-modal harmony throughout the generative process.

#### 3.2 Multimodal Diffusion Transformer

To enable effective generation from arbitrary combinations of text, video, and audio inputs, we design a unified multimodal conditioning framework. This framework supports flexible pairwise or tri-modal combinations such as TTA, V2A, and Text-Video-to-Audio (TV2A) in a single model. At the core of

our design is a modular encoder architecture and a joint conditioning mechanism that harmonizes modality-specific representations into a temporally aligned and semantically coherent latent space.

We build upon the MM-DiT framework introduced in SD3 [62], extending it with enhanced temporal synchronization modules and a dynamic masking strategy. During training, each modality is encoded independently, and missing modalities are replaced with learnable placeholders. The encoded features are then projected into a shared embedding space, where audio and video tokens are augmented with RoPE-based temporal positional encodings [66] to facilitate alignment. To enable generation over sequences of variable durations, we further introduce learnable duration embeddings as timing-aware vectors, which are fused with global conditioning features. These fused embeddings serve as the foundation for downstream joint attention and generation. Additionally, inspired by FLUX[67], we introduce audio-only single-modality blocks by simply removing the data streams of the other two modalities (i.e., converting joint attention into self-attention).

### 3.3 Flow Matching

During training, we use a conditional flow matching objective [68] for generative modeling [69]. This modeling method utilizes conditions  $C$  such as text or video embedding after encoding to learn a time-dependent conditional velocity vector field function  $v_\theta(t, C, x)$  that describes the direction and speed of the input  $x$  at timestep  $t$ , the variable  $\theta$  represents the network parameters that need to be learned. This function guides the Gaussian noise latent variable  $x_0$  to fit and approximate the target latent audio variable  $x_1$  using an ordinary differential equation solver to numerically integrate over the time interval  $t \in [0, 1]$ . The function  $u$  represents the target conditional vector field, and  $p$  represents the conditional probability path.  $q$  denotes the distribution of training data.

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t,q,p} \|v_\theta(t, C, tx_1 + (1-t)x_0) - u(t, tx_1 + (1-t)x_0)\|^2, \quad (1)$$

During inference, we set the time step  $t$  to 0.05 and employ the Euler method to numerically integrate  $v_\theta(t, C, x)$ , reconstructing the final latent audio representation from the initial Gaussian noise.

### 3.4 Aligned RoPE Positional Embedding

Precise temporal alignment plays a critical role in effectively synchronizing audiovisual content. To incorporate temporal information into the attention layers, we employ Rotary Positional Embeddings (RoPE)[66], which are applied to the queries and keys within the visual and auditory branches prior to joint attention (see Figure 2). Since the textual modality does not possess an intrinsic temporal structure comparable to audio and video, it is excluded from this positional encoding step. Furthermore, because visual tokens are sampled at a lower temporal resolution than audio tokens, we proportionally scale the frequency components of the visual positional embeddings to match the higher temporal rate of the auditory modality. This adjustment facilitates alignment across modalities. While these modified embeddings mitigate temporal misalignment, they are insufficient on their own to guarantee robust synchronization. Therefore, we introduce a dedicated synchronization module to further improve temporal coherence.

## 4 Kling-Foley

### 4.1 Overview

Inspired by MMAudio[15] we propose Kling-Foley. The core of our approach is to model the interactions between video, audio, and text modalities. We adopt the MM-DiT block design of SD3[62] and introduce two new components for temporal alignment: aligned RoPE[66] position embeddings to adapt to sequences of different frame rates, and 1D convolutions and MLPs to capture local temporal structure. At the same time, a duration learnable module is added to control the production of variable-length audio from naturally distributed video duration features. In addition, we incorporate audio-specific unimodal blocks based on FLUX [67] to make the network deeper with the same parameters without sacrificing multimodal capabilities. This architecture allows the model to selectively focus on different modalities based on the input, supporting joint training of audio-visual and audio-text data.

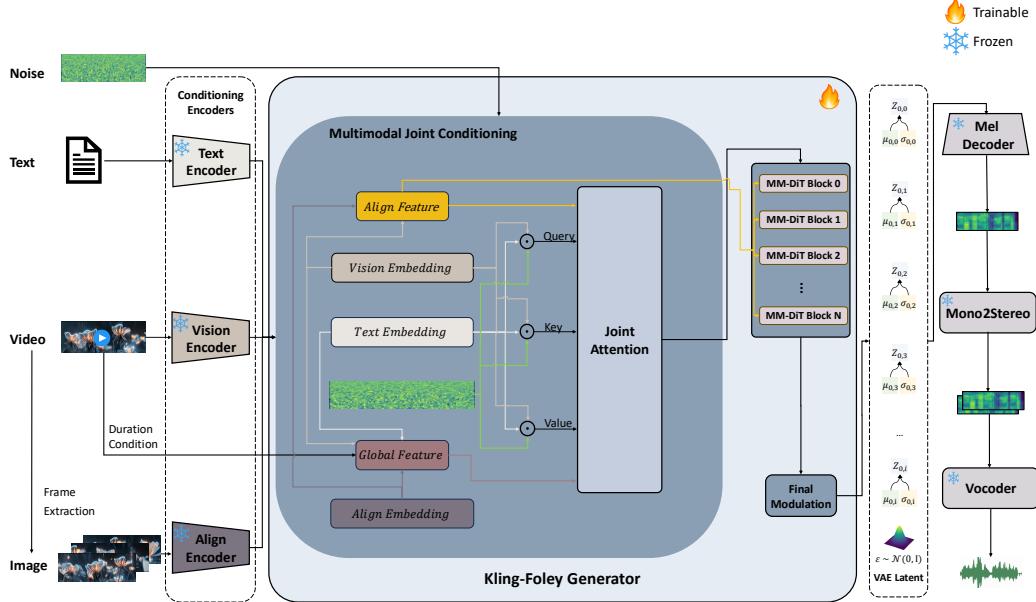


Figure 2: The core of Kling-Foley is a multimodal-controlled flowmatching model. Text, video, and temporally extracted video frames serve as conditional inputs. The multimodal features are then fused via a Multimodal Joint Conditioning module, which feeds into the MMDit Block for processing. This module predicts VAE latents, which a pretrained mel decoder subsequently reconstructs into a monaural mel-spectrogram. The monaural spectrogram is then converted to stereo spectrogram via a Mono2Stereo module. Finally, the stereo spectrogram is passed through a vocoder to generate the output waveform.

## 4.2 Variable Duration Control

To support variable-length audio-visual generation and enhance temporal control, we introduce discrete duration embeddings as part of the global conditioning mechanism [70]. Specifically, two scalar properties are computed per training clip: the start time seconds-start and the total duration seconds-total of the original video-audio sequence. These values are embedded into learnable per-second embeddings and concatenated with global textual and visual features. The resulting timing-aware global condition is fused with the flow timestep embedding via a shallow MLP and applied to all transformer layers using Adaptive Layer Normalization [71]. Each adaLN layer modulates token-wise activations by scaling and shifting normalized features based on the global conditioning vector.

## 4.3 Joint Attention Mechanism

In our architecture, we enable cross-modal communication through a joint attention strategy (see Figure 2). Inspired by prior work[72], we integrate the query, key, and value matrices across the textual, visual, and auditory modalities into a unified attention computation. Specifically, these modality-specific components are concatenated and passed through a shared scaled dot-product attention module[73]. After the attention operation, the resulting unified output is segmented back into the original three modalities based on the initial token groupings. While this joint mechanism facilitates rich cross-modal interaction, we emphasize that, on its own, it does not ensure temporal synchrony across streams such as audio and video.

## 4.4 Flexible Pairwise Training

To effectively support multimodal generation under arbitrary combinations of available inputs (text, video, and audio), we adopt a conditional training strategy that reflects the modular structure outlined in Algorithm 1. Each modality is first encoded independently: video inputs are processed by

---

**Algorithm 1** Multimodal Conditional Training Strategy

---

**Input:** Audio  $x$ , Video  $V$  (optional), Text  $T$  (optional)

**Hyperparameters:** total steps  $S$ , frame rates, number of transformer blocks  $N_1, N_2$

**Initialize:** Learned empty embeddings  $e_v, e_t$  for missing V/T

**for** each training step  $s = 1$  to  $S$  **do**

**Stage 1: Modality-Specific Encoding**

**if** Video  $V$  available **then**

Extract sync features  $F_{\text{sync}} \leftarrow \text{SyncFormer}(V)$

Extract video features  $F_v \leftarrow \text{MetaCLIP}_{\text{visual}}(V)$

**else**

$F_{\text{sync}} \leftarrow e_v, F_v \leftarrow e_v$

**end if**

**if** Text  $T$  available **then**

$F_t \leftarrow \text{T5}(T)$

**else**

$F_t \leftarrow e_t$

**end if**

Encode audio  $x$  to latent representation  $x_0$  with diffusion noise

**Stage 2: Conditional Synchronization Module**

Project  $F_{\text{sync}}$  and upsample to frame-aligned sync feature

Generate align feature  $a_f$  and global feature  $g_f$

**Stage 3: Transformer-based Joint Processing**

Project and align  $F_v, F_t, x_0$

Inject positional embeddings (RoPE) to visual/audio queries and keys

**for**  $i = 1$  to  $N_1$  **do**

Apply multimodal transformer block with joint attention over  $F_v, F_t, x_0$  and conditions  $a_f, g_f$

**end for**

**for**  $i = 1$  to  $N_2$  **do**

Apply single-modal transformer block to refine audio flow path

**end for**

**Stage 4: Output Flow Prediction**

Apply adaptive LayerNorm and 1D-Conv to predict audio flow  $w$

Use  $w$  in reverse diffusion to reconstruct waveform

**Stage 5: Loss Computation and Optimization**

Compute losses:  $\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t,q(x_1),p_t(x|x_1)} \|v_t(x) - u_t(x|x_1)\|^2$

Backpropagate and update model

**end for**

---

MetaCLIP and SyncFormer to extract semantic and synchronization features, while textual inputs are encoded via T5. Missing modalities are substituted with learned placeholder embeddings ( $e_v, e_t$ ), ensuring a consistent representation space.

The synchronization features are projected and upsampled to produce two types of conditioning tokens: align feature and global feature, the latter incorporating learnable embeddings of start time and total duration. These global tokens are fused with diffusion timestep embeddings via a shallow MLP and modulate each transformer layer through Adaptive LayerNorm [71].

In the joint transformer stage, modality-specific features are projected into a shared latent space. RoPE embeddings [66] are added to audio and visual tokens—rescaled for temporal alignment—to encode temporal structure. Joint attention enables cross-modal interaction, while audio-only transformer blocks [67], applied after joint fusion, provide efficient unimodal refinement, benefiting tasks like audio continuation and TTA.

## 4.5 Latent Audio Codec

The latent audio codec extends our prior VQ-CTAP framework[74, 75], inheriting its core components while introducing key modifications to optimize audio reconstruction. As illustrated in Figure 3, the core of the latent audio codec is a Mel-VAE composed of three primary components: a mel encoder, a mel decoder, and a discriminator. The audio encoder processes an input waveform sampled at 44.1 kHz, generating embeddings at a rate of 43 Hz (equivalent to 1024 times downsampling relative to

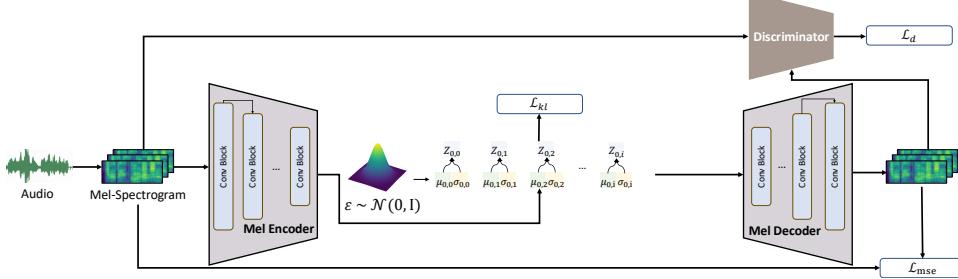


Figure 3: The main body of latent audio codec is a Mel-VAE, which jointly trains a mel encoder, a mel decoder, and a discriminator. The VAE structure enables the model to learn a continuous and complete distribution of latent spaces, significantly enhancing its audio representation capabilities.

the input sampling rate). Critically, the VAE structure enables the model to learn a continuous and complete distribution of latent spaces, significantly enhancing audio representation capabilities.

#### 4.5.1 Structure

As illustrated in Figure 3,  $A_{in}$  denotes the input batch of target audio data,  $A_{in} \in \mathbb{R}^{B \times T_s \times D_s}$ , where  $B$  is the batch size,  $T_s$  is the number of time frames, and  $D_s$  is the number of spectral components (mel-spectrogram bands). The  $S_{in}$  is passed through the audio encoder:  $A = \text{AudioEncoder}(A_{in})$

where  $A \in \mathbb{R}^{B \times T_s / 2 \times d}$  are the target audio representations. The audio encoder compress the length of the audio representations by a factor of 2.

Following encoding, the model utilizes the VAE structure[76] to parameterize the latent distribution.

In Equation (5),  $D_{KL}$  refers to the KL loss,  $\mathcal{N}(\cdot)$  represents a Gaussian distribution, and  $(\hat{\mu}, \hat{\sigma})$  denotes the (mean, variance) of the audio representation latent space distribution. Random operations in the network cannot be processed by backpropagation, "reparameterization trick" is introduced to VAE:  $\mathbf{z} = \hat{\mu} + \hat{\sigma} \odot \phi; \phi \sim \mathcal{N}(0, I)$

To address the KL collapse problem[77], a margin  $\Delta$  is introduced to limit the minimum value of the KL loss as shown:

$$\mathcal{L}_{kl} = \max(0, D_{KL}[\mathcal{N}(\hat{\mu}, \hat{\sigma}^2) || \mathcal{N}(0, I)] - \Delta) \quad (2)$$

To enable reconstruction capability using the pre-trained latent representations, the latent variable  $\mathbf{z}$  serves as input to the audio decoder for predicting the mel-spectrogram:  $A_d = \text{AudioDecoder}(\mathbf{z})$

A mean squared error (MSE) loss is used to compare the predicted mel-spectrogram  $A_d$  with the ground-truth mel-spectrogram  $A_{in}$ :

$$\mathcal{L}_{mse} = MSE(A_{in}, A_d) \quad (3)$$

The discriminator module is pivotal for adversarial training, enforcing the generator to produce high-fidelity audio spectrograms indistinguishable from real data. Its design integrates multi-scale feature analysis, gradient regularization, and dynamic loss weighting, addressing spectral artifacts and training instability common in audio synthesis.

The generator aims to deceive the discriminator by maximizing the probability of generated spectrograms being classified as real. This is formalized as a non-saturating loss:

$$\mathcal{L}_g = -\mathbb{E}[D(A_d)] \quad (4)$$

Where  $D(\cdot)$  represents the discriminator's output logit.

The discriminator learns to distinguish real and generated spectrograms via a weighted binary objective:

$$\mathcal{L}_d = \underbrace{\gamma_{step} \cdot \mathcal{L}_{adv}}_{\text{adversarial loss}} + \underbrace{\lambda_{R1} \cdot \mathcal{R}_1}_{\text{gradient penalty}} \quad (5)$$

$$\mathcal{L}_{adv} = \mathbb{E}[\max(0, 1 - A_{in})] + \mathbb{E}[\max(0, 1 + D(A_d))] \quad (6)$$

$$\mathcal{R}_1 = \mathbb{E}_x[|\nabla_x \text{AudioDecoder}(\mathbf{z})|_2^2] \quad (7)$$

Where  $\gamma_{step}$  and  $\lambda_{R1}$  represent the dynamic weighting.

In the experiment, we employ a mel-spectrogram encoder structurally similar to that used in Make-An-Audio2 [54]. This encoder comprises 32 stacked 1D-convolutional layers. The mel decoder mirrors this architecture with transposed convolutions for mel-spectrogram reconstruction.

#### 4.5.2 Multi-Stage Stepping Optimization Strategy

A stepping optimization strategy is designed to ensure effective model convergence by gradually injecting and adjusting the influence of various loss components, as shown in Algorithm 2. The training process involves the following losses:  $\mathcal{L}_{kl}$ ,  $\mathcal{L}_{mse}$ ,  $\mathcal{L}_d$ , and  $\mathcal{L}_g$ . The variable  $step$  represents the current training step. Initially, the model is trained using  $\mathcal{L}_{mse}$ . When the  $step$  exceeds the specified starting step for  $\mathcal{L}_{kl}$ ,  $\mathcal{L}_{kl}$  is added to the training process. The weight for  $\mathcal{L}_{kl}$  increases gradually as the training progresses. Once the  $step$  surpasses the specified ending step, the weight for  $\mathcal{L}_{kl}$  is fixed at  $kl\_upper$ . Similarly, when the  $step$  exceeds the specified starting step for  $\mathcal{L}_d$  &  $\mathcal{L}_g$ , These losses are incorporated into the training process. The weight for  $\mathcal{L}_d$  &  $\mathcal{L}_g$  also increases gradually during training. Once the  $step$  exceeds the specified ending step, the weight for  $\mathcal{L}_d$  &  $\mathcal{L}_g$  is fixed at  $gan\_upper$ . In the final training stage, we freeze both the audio encoder and VAE, training exclusively the audio decoder and discriminator. This focused optimization addresses potential error propagation by fine-tuning the vocoder using mel-spectrograms generated from the audio decoder's output. Algorithm 2 outlines the step-wise inclusion of different losses and their corresponding weight adjustments. This optimization strategy aims to facilitate effective model convergence by gradually introducing and adjusting the influence of various loss components throughout the training process.

---

**Algorithm 2** Multi-Stage Stepping Optimization Strategy

---

```

Initialize hyperparameters:
stage1end, stage2end, stage3end {Phase transition steps}
klupper, ganupper {Max loss weights}
for each training step do
    if step ≤ stage1end then
        Stage 1: Base Reconstruction
        Train full model with  $\mathcal{L}_{mse}$  {Pure MSE focus}
    else if step ≤ stage2end then
        Stage 2: Regularization Enhancement
         $\gamma \leftarrow kl_{upper} \cdot \min(1, \frac{step - stage1_{end}}{stage2_{end} - stage1_{end}})$ 
         $\mathcal{L}_{total} \leftarrow \mathcal{L}_{mse} + \gamma \mathcal{L}_{kl}$ 
        Train full model with  $\mathcal{L}_{total}$ 
    else if step ≤ stage3end then
        Stage 3: Adversarial Introduction
         $\delta \leftarrow gan_{upper} \cdot \min(1, \frac{step - stage2_{end}}{stage3_{end} - stage2_{end}})$ 
         $\mathcal{L}_{total} \leftarrow \mathcal{L}_{mse} + kl_{upper} \mathcal{L}_{kl} + \delta (\mathcal{L}_g + \mathcal{L}_d)$ 
        Train full model with  $\mathcal{L}_{total}$ 
    else
        Stage 4: Decoder Refinement
         $\mathcal{L}_{total} \leftarrow \mathcal{L}_{mse} + kl_{upper} \mathcal{L}_{kl} + gan_{upper} (\mathcal{L}_g + \mathcal{L}_d)$ 
        Freeze audio encoder and VAE parameters
        Train only audio decoder and discriminator with  $\mathcal{L}_{total}$ 
    end if
end for

```

---

## 4.6 Mono-to-Stereo

The process utilizes a Mono2Stereo module to convert the monaural mel-spectrogram into dual-channel mel-spectrograms. Critically, this module only predicts the ratio of the left and right mel-spectrograms relative to the monaural mel-spectrogram. This targeted prediction significantly reduces data dependency and enhances training stability. Finally, these left and right mel-spectrograms are processed by a vocoder to generate corresponding waveforms for each channel, which are concatenated along the channel dimension to produce the final stereo audio.

## 5 Data

### 5.1 Overview

To train a multimodal generative model capable of synthesizing diverse and realistic sound effects, it is essential to construct a large-scale training dataset that is broad in coverage and tightly aligned across modalities. Current research in sound effect generation faces two critical limitations. First, most existing datasets are relatively small, typically containing only tens of thousands of audio samples, which are insufficient to support the training of large-scale generative models that require high data diversity. Second, the majority of these datasets are incomplete in modality structure—lacking alignment among audio, video, and natural language—which significantly limits the model’s ability to utilize conditional inputs effectively. For instance, the VGG-Sound dataset [17] contains audio-video pairs but only includes coarse class labels, without natural language descriptions of the sound content. On the other hand, datasets such as AudioCaps [78], Clotho [79], and WavCaps [80] focus primarily on audio-caption alignment, yet do not include accompanying video streams, making them less suitable for training generation models conditioned on both vision and language.

To address these challenges, we construct a new large-scale multimodal sound effects dataset from scratch, consisting of over 100 million samples. Each sample contains a raw video segment, a corresponding monaural audio clip, and a structured textual description of the audio. The three modalities are tightly aligned and sourced from real-world online video content.

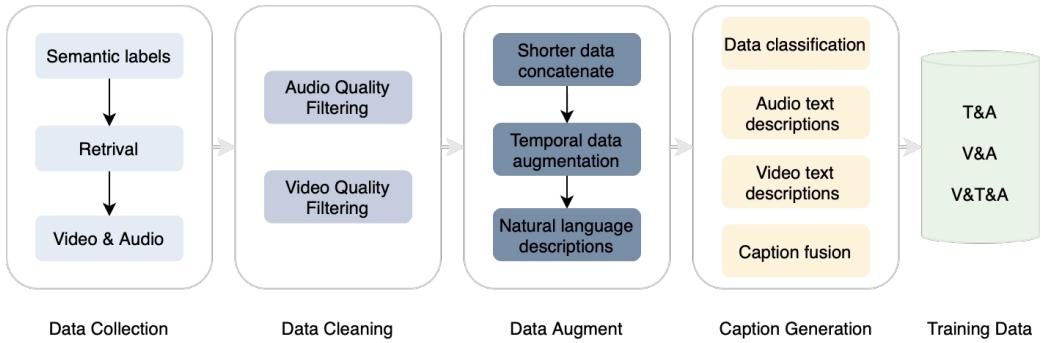


Figure 4: Audio and video data undergo preprocessing and quality filtering to obtain high-quality single-event audio and video segments. Subsequently, synthetic multi-event audio samples are generated through temporal augmentation, and large models are used to generate and extract keywords and classification captions for audio and video. Finally, various caption information is combined to produce the final training captions.

### 5.2 Data Construction

In this work, we constructed three types of paired data: text-audio, video-audio, and video-text-audio. Our overall data processing workflow is shown in Figure 4.

**Data Collection** The generative capacity of sound synthesis models is largely determined by the variety of sound sources and the range of semantic labels present in the training data. To ensure broad coverage, we build our label set based on the hierarchical structure defined in the AudioSet [16] ontology, selecting categories from its top three levels. This ontology provides a clear semantic

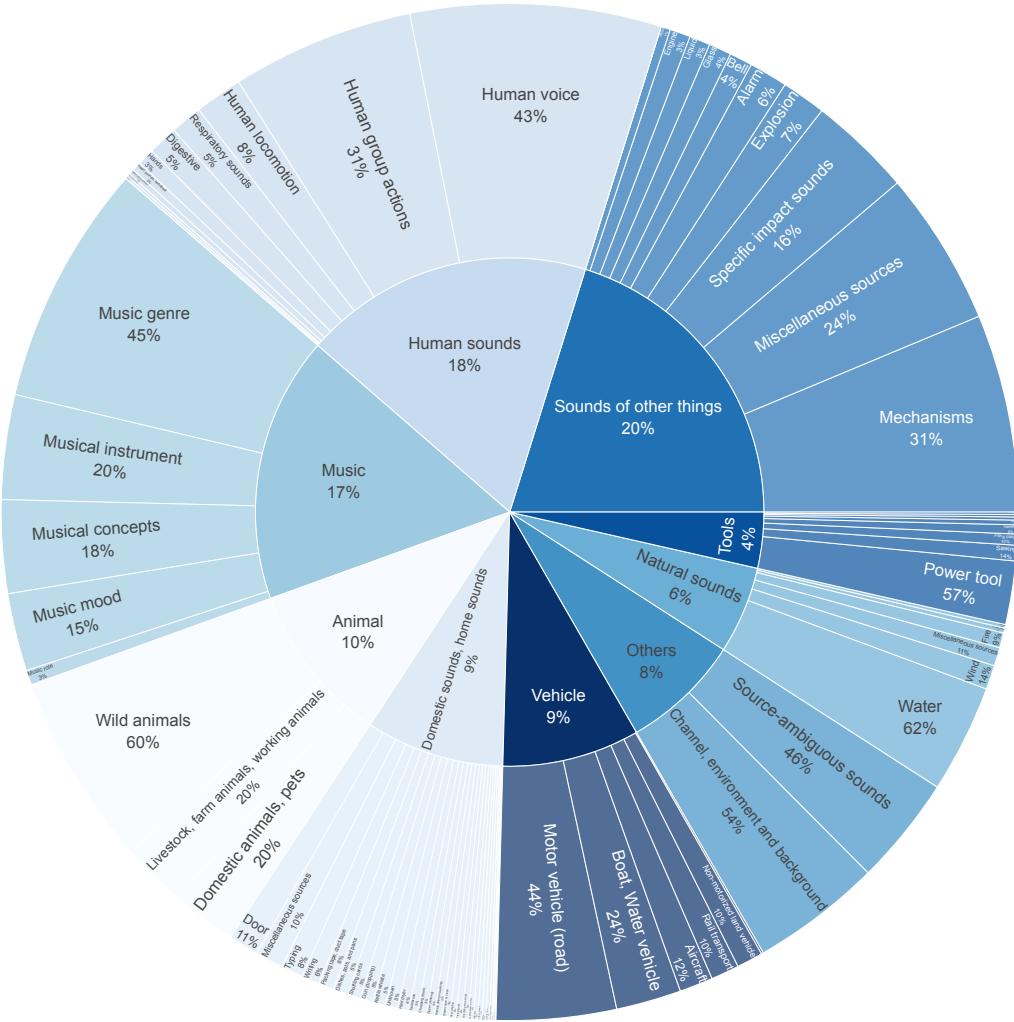


Figure 5: Category distribution of sound events in the training set. The broad coverage of real-world acoustic events ensures the diversity and generalizability required for training open-domain sound generation models.

hierarchy, serving as a principled foundation for constructing a systematic keyword vocabulary for data mining.

Using this label set, we derive a keyword bank to guide the large-scale retrieval process. We query video platforms using these keywords and filter candidate videos and channels based on metadata such as titles, descriptions, and tags to ensure semantic relevance. To further enhance long-tail coverage and content diversity, we supplement the collected videos with manually curated resources and samples from existing open-source datasets. The resulting raw multimodal data serves as the foundation for our dataset and is subsequently passed through a dedicated cleaning pipeline to ensure quality and alignment across modalities.

**Data Cleaning** We only retain data with video resolution above 720P and a small proportion of subtitles, and we uniformly convert the audio to WAV format with a 44k sample rate, 16-bit depth, and stereo channels. For audio, we perform quality filtering based on SNR, MOS score, clipping ratio, and audio bandwidth. We use VAD to select audio data with a silence ratio of less than 0.2. We employ the CLAP model to calculate the consistency between audio and text labels, retaining data with high consistency. Finally, we split longer videos and audio into 10-second segments.

**Data Augment** For shorter data, we concatenate shorter videos and audio to obtain data with a duration of 10 seconds, which enables the model to better response to dynamic visual input. To enhance the temporal alignment ability of our model, we introduce temporal data augmentation [54] by concatenating single-event video and audio clips according to different temporal rules to generate multi-event data. The text captions corresponding to the temporally augmented video and audio are obtained by merging the original single-event data text captions. Additionally, we extract key information such as sound sources, objects, scenes, emotions, gender, actions, and modifiers from the audio. Using a large model combined with the extracted keywords, we transform the unstructured original text descriptions into semantically complete natural language descriptions.

**Caption Extraction** Video and audio can typically yield textual descriptions containing different information. To obtain text captions that are as accurate, detailed, and complete as possible, we utilize both video and audio to derive the final text caption. First, we use audio classification model to classify the video and audio, retaining data and corresponding category labels for four categories: sound effects, music, speech, and singing [81]. For different categories of data, we employ corresponding audio understanding large models to extract audio text descriptions from the audio, while also extracting video text descriptions from the video. Subsequently, we input the audio descriptions, video descriptions, and the enhanced natural language text descriptions into a large model to obtain the final fused text caption [2].

**Training Data** As illustrated in Figure 5, we visualize the distribution of high-level sound categories in our training set. Our training data contains textaudio, videoaudio, and videotextaudio three types of paired data. Our dataset spans a wide variety of real-world acoustic scenarios, including natural environments, human activities, animal sounds, mechanical operations, and transportation, providing a solid foundation for learning diverse generative patterns and improving the realism and controllability of synthesized audio.

### 5.3 Benchmark Dataset

Several audio-visual datasets have been proposed to support sound-related tasks, as shown in Table 1. AudioSet [16] is one of the largest general-purpose audio datasets, but its heavy reliance on human annotation leads to high construction costs. VGGSound [17] improves scalability through audio-visual alignment, making it more practical for sound generation evaluation. EPIC-SOUNDS [82] focuses on audio-driven actions, offering precise temporal boundaries and fine-grained labels.

However, a common limitation of these datasets is the lack of textual descriptions (captions) for both audio and video modalities, which hinders comprehensive evaluation in caption-aware or text-conditioned generation scenarios. To address this issue, a common solution is to incorporate additional audio-text test sets, such as Clotho [79] and AudioCaps [78]. While Clotho provides high-quality captions, its scale is limited. AudioCaps offers more annotated samples, but only covers audio modality and lacks video context. WavCaps [80], though large-scale, is weakly labeled and unsuitable for evaluation.

Table 1: Comparison of our dataset with existing datasets.

Dataset	#Test clips	Length	Video	Video Caption	Audio	Audio Caption	#Class
AudioSet [16]	18K	50h	Yes	No	Yes	No	527
VGGSound [17]	15K	41.7h	Yes	No	Yes	No	309
Epic Sounds [82]	10K	13.9h	Yes	No	Yes	No	44
Clotho [79]	1K	6h	No	No	Yes	Yes	-
AudioCaps [78]	1K	2.7h	No	No	Yes	Yes	-
<b>Kling-Audio-Eval</b>	<b>21K</b>	<b>58.6h</b>	Yes	Yes	Yes	Yes	1919

Despite these efforts, there remains no benchmark that jointly supports vision, audio, and language modalities for evaluating sound effect generation. **To fill this gap, we introduce Kling-Audio-Eval — the first high-quality multimodal benchmark that combines video, video captions, audio, audio captions, and sound event labels.** Our dataset is constructed through a carefully designed

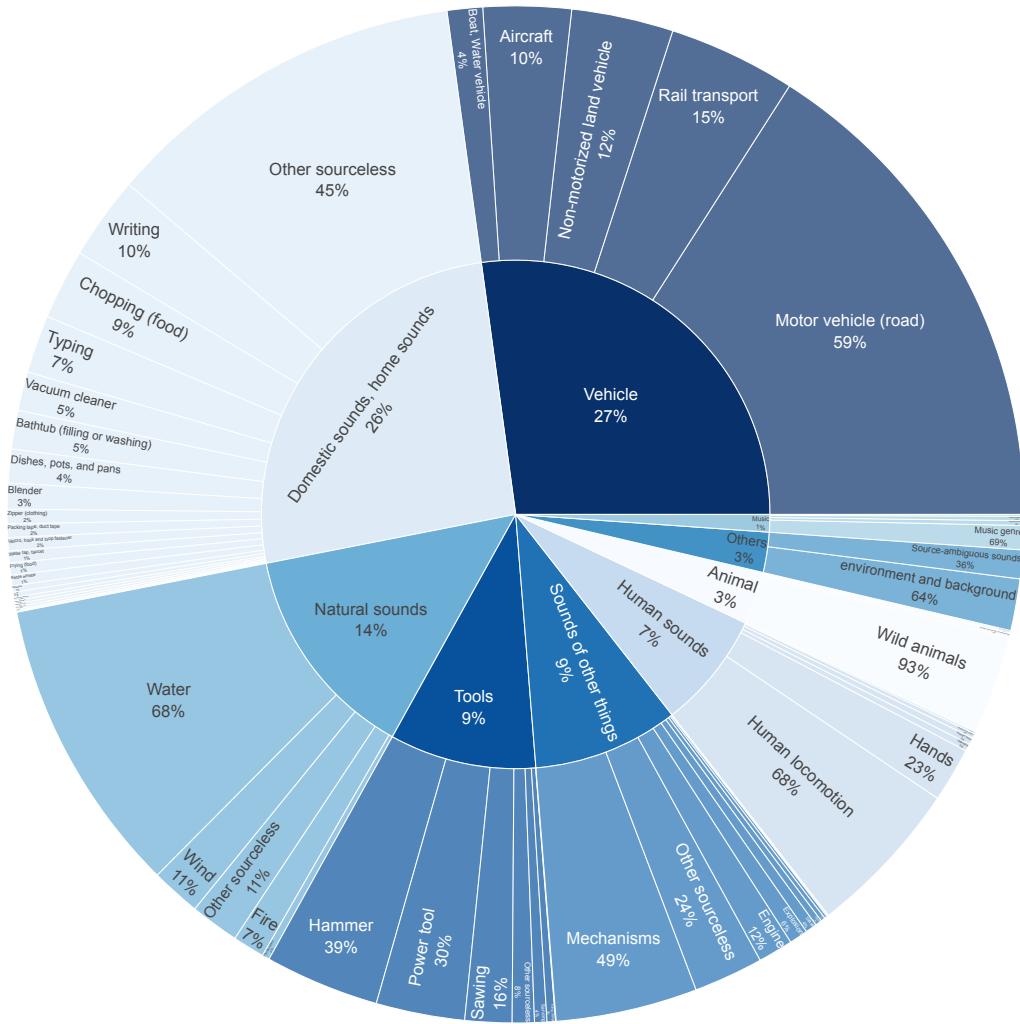


Figure 6: Distribution of various categories in the benchmark dataset.

taxonomy and extensive human annotation, enabling robust evaluation across multiple modalities and generation tasks.

Based on existing sound libraries and internal statistics, we selected the 1,000 most frequent third-level labels (Section 5.1). These were chosen to ensure full coverage of all first- and second-level categories, providing broad content representation. Following a strict data cleaning process (Section 5.2), we then selected 30,000 samples from these categories, each with pre-generated captions and sound event labels, for further human verification and annotation.

The manual annotation process covers four main aspects: audio captions, video captions, sound event labels, and audio-visual quality assessment, aiming to ensure consistency, accuracy, and practical usability. The specific annotation guidelines are as follows:

- **Caption Correction:** Review the pre-generated audio and video captions, and revise any errors or omissions using concise and clear language.
  - **Modal Independence:** Ensure that audio and video captions are annotated independently. For example, audio captions should not rely on visual information, and vice versa.
  - **Label Verification:** Check whether the assigned first- and second-level labels match the actual content. If not, select the correct labels from the predefined taxonomy.

- Valid Sample Selection: Only retain audio-visual clips that meet the following criteria: (i) Foreground audio must not contain human voices. (ii) Sound effects must originate from tasks or objects, while speech, singing, and musical sounds must be produced by visible individuals in the video. (iii) Video duration must be at least 5 seconds, and the sound effect must last at least 2 seconds. (iv) Only ambient off-screen sounds are permitted, such as birdsong in a forest setting; (v) Videos are allowed to include watermarks, logos, or subtitles; (vi) Background music must not contain vocals.

In total, we collected 20,935 high-quality samples to form the final test set, and the original 1,000 third-level labels were further refined into 1,919. The category distribution is illustrated in Figure 6.

## 6 Experiments and Results

### 6.1 Experimental Settings

We benchmark against the following methods: V2A-Mapper [11] maps CLIP visual embeddings into the CLAP audio-text space to enable AudioLDM-based generation, while FoleyCrafter [36] incorporates time-conditional estimators for improved temporal coherence. VATT [37] and VTALDM [38] adopt transformer-based architectures for joint video, audio, and text representation learning. V-AURA [39] employs a high-framerate visual encoder ( $6\times$  higher than prior work) in an autoregressive fusion framework to enhance motion sensitivity. To improve efficiency, FRIEREN [40] introduces rectified flow matching with reflow and one-step distillation. MMAudio [15] unifies audio-visual and audio-text data training via a multimodal Transformer trained on VGGSound and WavCaps. MMAudio’s results are derived from the audio results generated by inference using the open source model. Finally, ReWaS [13] addresses cross-modal gaps by introducing energy-based constraints and hand-crafted audio features.

### 6.2 Training Strategy

**Learning Rate Scheduling.** We adopt a smooth inverse decay schedule with exponential warmup [83, 84, 85]. The learning rate at step  $t$  is given by:

$$\text{LR}(t) = \max \left( \text{LR}_{\text{final}}, \text{LR}_{\text{base}} \cdot \left( 1 + \frac{t}{\gamma_{\text{inv}}} \right)^{-p} \right) \cdot (1 - w^{t+1}), \quad (8)$$

where  $\gamma_{\text{inv}}$  controls the decay speed,  $p$  determines the curvature of the decay, and  $w \in [0, 1]$  enables a smooth exponential increase in the early training stage. The  $\max(\cdot)$  operation ensures that the learning rate does not fall below a specified  $\text{LR}_{\text{final}}$ . InverseLR scheduler offers a continuous, smooth decay that aligns more closely with the training dynamics. This not only accelerates convergence but also enhances stability across diverse and complex multimodal training scenarios.

**Scaling Strategy.** We investigate the impact of model scaling on multimodal audio-language learning by progressively increasing model capacity from 1.5B to 6B parameters. Inspired by empirical scaling laws [86, 87], we scale the model along three dimensions—depth  $d$ , hidden dimension  $d_{\text{hidden}}$ , and number of attention heads  $h$ , while preserving a consistent architectural ratio  $h/d \in [20, 100]$  [88].

Formally, we follow the parameterization:

$$d_{\text{hidden}} = 64 \cdot h, \quad h = \alpha \cdot d \quad (9)$$

where  $\alpha$  is a scaling factor controlling the heads-to-depth ratio. This ensures balanced compute utilization across attention and MLP blocks.

Our base model uses 17 heads and depth 23 (1.5B), while the 3B and 6B variants increase the head count to 23 and 32 respectively, keeping depth consistent at 27, thereby scaling width and total representational capacity.

### 6.3 Objective Metrics

Following the setup in MMAudio, we evaluate both our model and selected baselines across four dimensions on the VGGSound [17] test set (15,220 samples). Results are presented in the table below. The objective metrics of V2A used are as follows:

Table 2: Results of V2A on VGGSound

<b>Method</b>	<b>Distribution matching</b>		<b>Semantic align.</b>	<b>Temporal align.</b>	<b>Audio quality</b>	
	<b>FDPANNs</b> ↓	<b>KLPANNs</b> ↓	<b>IB-score</b> ↑	<b>DeSync</b> ↓	<b>SDR</b> ↑	<b>MCD</b> ↓
ReWaS[13]	17.54	2.87	14.82	1.06	-	-
VTA-LDM[38]	14.49	2.23	24.73	1.26	-	-
V-AURA[39]	14.80	2.42	27.64	0.65	-	-
VATT[37]	10.63	<b>1.48</b>	25.00	1.20	-	-
Frieren[40]	11.45	2.73	22.78	0.85	-	-
FoleyCrafter[36]	16.24	2.30	25.68	1.23	-	-
V2A-Mapper[11]	8.40	2.69	22.58	1.23	-	-
MMAudio[15]	<b>6.29</b>	1.77	29.26	0.45	-3.09	2.84
Kling-Foley	7.60	1.86	<b>30.75</b>	<b>0.43</b>	<b>-2.41</b>	<b>2.60</b>

- **FD (Fréchet Distance):** This metric assesses the similarity between the distributions of generated and ground-truth audio features. It is computed on feature embeddings from pre-trained audio classifier PANNs[89]. A lower FD value signifies that the generated feature distribution is closer to the ground-truth distribution, indicating higher fidelity.
- **KL (Kullback-Leibler Divergence):** KL Divergence measures the difference between the probability distributions of audio events in the generated set versus the ground-truth set. It is calculated using the output predictions from pre-trained classifier PANNs[89]. A lower score is better, indicating that the generated audio has a similar event distribution to the reference audio.
- **ImageBind Score (IB-score)[90]:** This metric evaluates the semantic coherence between a video and its generated audio. It calculates the cosine similarity of features extracted from both the video and audio modalities using the unified ImageBind model. A higher score reflects better cross-modal semantic consistency.
- **DeSync Score[65]:** This metric quantifies audio-video synchronization by predicting the temporal misalignment between the visual stream and the generated audio. It employs Synchformer[65] to output the predicted offset in seconds. A lower absolute value indicates better synchronization.

To comprehensively evaluate the capabilities of our latent audio codec, we conduct a direct comparison with MMAudio[15]. The comparison is performed across four distinct tasks: sound effect, music, speech, and singing. For each task, we test on 500 out-of-domain audio samples, resulting in a total of 2,000 test cases to ensure a robust evaluation. Specifically, for the Sound effect, music, and speech scenarios, we utilize the evaluation set from the Codec-SUPERB @ SLT 2024 challenge\*, while the Sing task is evaluated on a proprietary, internally constructed dataset.

The objective metrics of latent audio codec used are as follows:

- **PESQ[91] (Perceptual Evaluation of Speech Quality):** This metric rates the perceptual quality of speech on a scale from -0.5 to 4.5. It is designed to model subjective quality scores, making it a strong indicator of human perception. A higher score is better. We report this metric for the Speech and Sing tasks.
- **SI-SDR[92] (Scale-Invariant Signal-to-Distortion Ratio):** SI-SDR measures the fidelity of the waveform in the time domain, independent of the overall signal amplitude. It provides a robust assessment of signal reconstruction. A higher value is better.
- **SDR[93] (Signal-to-Distortion Ratio):** This metric quantifies the ratio of the original signal’s power to that of the reconstruction error, serving as a fundamental measure of distortion. A higher value indicates better signal integrity.
- **LSD (Log-Spectral Distance):** LSD gauges the discrepancy in frequency content by calculating the error between the log-power spectra of the generated and reference audio. A lower value signifies a more accurate spectral envelope.
- **MCD (Mel-Cepstral Distortion):** MCD measures the distance between Mel-Frequency Cepstral Coefficients, providing a crucial evaluation of timbral texture and vocal naturalness, which are highly relevant to human hearing. A lower value is better.

\*[https://github.com/voidful/Codec-SUPERB/tree/SLT\\_Challenge?tab=readme-ov-file](https://github.com/voidful/Codec-SUPERB/tree/SLT_Challenge?tab=readme-ov-file)

Table 3: Results of Latent Audio Codec

Task Type	Model	PESQ $\uparrow$	SISDR $\uparrow$	SDR $\uparrow$	LSD $\downarrow$	MCD $\downarrow$	Mel Loss $\downarrow$	STFT Loss $\downarrow$
Sound Effect	GT	-	-24.84	-2.53	0.61	0.76	0.41	0.86
	MMAudio	-	-29.57	-3.29	<b>0.70</b>	1.45	0.81	1.18
	Kling-Foley	-	<b>-29.47</b>	<b>-2.41</b>	<b>0.70</b>	<b>1.35</b>	<b>0.78</b>	<b>1.14</b>
Music	GT	-	-24.99	-2.50	3.07	5.76	0.81	1.60
	MMAudio	-	-30.42	-2.54	3.07	5.72	<b>1.13</b>	<b>1.83</b>
	Kling-Foley	-	<b>-29.85</b>	<b>-2.16</b>	<b>3.02</b>	<b>5.44</b>	1.20	1.95
Singing	GT	4.02	-15.04	-2.15	0.69	0.29	0.33	0.78
	MMAudio	2.80	<b>-22.84</b>	-3.90	0.83	1.08	0.76	1.14
	Kling-Foley	<b>2.88</b>	-23.55	<b>-2.83</b>	<b>0.81</b>	<b>0.70</b>	<b>0.60</b>	<b>1.00</b>
Speech	GT	3.72	-8.30	-0.32	2.51	3.18	0.54	1.10
	MMAudio	3.10	-26.64	-2.62	2.51	3.09	<b>0.85</b>	<b>1.32</b>
	Kling-Foley	<b>3.27</b>	<b>-26.32</b>	<b>-2.27</b>	<b>2.48</b>	<b>2.91</b>	0.89	1.37

- **Mel Loss & STFT Loss:** These metrics directly quantify the reconstruction error at the spectral level by calculating the L1 distance between the predicted and ground-truth spectrograms (Mel and STFT, respectively). They reflect the model’s ability to accurately reproduce the underlying spectral structure. Lower values are better.

#### 6.4 Inference Optimize

Inference utilizes static computation graph technology such as *torch.compile*’s JIT-compiling to achieve acceleration. We need to maintain a fixed input shape, and when encountering shorter inputs, we apply padding. Next, the original Conv1d kernel is decomposed into kernel’s size individual linear layers, converting the convolution into multiple small matrix multiplications (GEMM). This enables the use of highly optimized basic linear algebra subprograms libraries like cuBLAS.

#### 6.5 Results

**Video-to-Audio.** The results of latent audio codec are presented in Table 2. Distribution matching pertains to the similarity in distribution between the generated audio and the real audio within the feature space. FDPANNs and KLPANNs are utilized as metrics (where lower values are more favorable). For Kling-Foley, FDPANNs registers at 7.60, trailing only MMAudio. The KLPANNs results of VATT is notably superior to those of other models. However, the KLPANN results of Kling-foley value is 1.86, which is also better than most baseline models, for instance, ReWaS has a KLPANNs of 2.87. In terms of distribution matching indicators, the MMAudio model also achieved good results, mainly because the video dataset is basically derived from the VGGSound training datasets, so the generated audio has a distribution similar to the original audio.

Semantic alignment gauges the semantic consistency between the generated audio and the video content, with the IB-score (higher scores are preferred) serving as the indicator. Kling-Foley attains an IB-score of 30.75, exceeding MMAudio’s 29.26 and V-AURA’s 27.64. This indicates that our model has the most robust semantic understanding capability. The reason is that we employ Metaclip which offers enhanced visual semantic understanding. We also use the T5-base model, and the model is known for its superior text semantic understanding.

Temporal alignment assesses the synchronization of audio and video events, with DeSync (lower values are better) as the primary metric. Kling-Foley achieves the result of 0.43, marginally better than MMAudio and significantly better than other models. For example, VATT has a DeSync of 1.20. This can be credited to our implementation of a more refined temporal alignment module.

**Latent Audio Codec.** The results of the latent audio codec are presented in Table 3, where our model is benchmarked against the MMAudio baseline across four tasks. The data reveals a consistent trend: our proposed model demonstrates superior or highly competitive performance across all scenarios. A key finding is our model’s consistent advantage in metrics that are crucial for perceptual quality and signal fidelity, such as PESQ, SDR, and especially MCD. This suggests that our approach more effectively generates audio that is not only spectrally accurate but also perceptually closer to the ground-truth reference.

Specifically, in the sound effect, singing, and speech tasks, our model surpasses the baseline in nearly all metrics, with particularly significant gains in perceptual quality (e.g., PESQ of 3.27 vs. 3.10 for Speech) and timbral accuracy (e.g., MCD of 0.70 vs. 1.08 for Singing). For Music generation, while MMAudio is competitive in direct spectral reconstruction losses (Mel/STFT Loss), our model achieves better performance in the more perceptually relevant SDR and MCD metrics, indicating a more faithful synthesis of complex musical textures. Collectively, these results validate the effectiveness and robustness of our proposed model for high-quality, multi-domain audio generation.

## 7 Conclusion

In this work, we propose Kling-Foley, a novel multimodal framework for high-quality V2A generation that achieves synchronization between audio outputs and visual content. Furthermore, we address the critical gap in evaluation resources by releasing Kling-Audio-Eval, the first industry-scale multimodal benchmark featuring 20,935 manually annotated samples across nine sound scenarios, providing comprehensive video-audio-text annotations for rigorous assessment. Extensive experiments validate that Kling-Foley achieves SOTA performance across most metrics: audio quality, distribution matching, semantic alignment, and audio-visual synchronization. Future work will focus on extending the framework to support longer video sequences, and enhancing cross-modal alignment for complex auditory scenes.

## 8 Limitations, Ethics and Safety

Kling-Foley provides an industrial-grade solution through multimodal alignment and stereo rendering technology. It replaces the traditional Foley sound effect artist’s manual annotation process and significantly reduces the time and economic cost of video dubbing. The cross-modal audio sequence representation supports the unified modeling of mixed scenes of sound effects, voice, and music, and is suitable for interactive media such as games and virtual live broadcasts. Combined with text conditions, it accurately controls the semantics of sound effects (such as "glass shattering + from far to near") and provides refined sound effect generation capabilities. It supports stereo spatial rendering of targets moving in a specific direction (such as vehicles passing by and bird trajectories) to enhance the immersion of film and television works.

**Limitations** Despite the leading performance of the model, there are still the following technical bottlenecks. Insufficient modeling of complex physical processes, the generation fidelity of multi-object interactive sound effects (such as layered voices in crowd conversations and chain reactions of object collisions) is low, and acoustic logic errors are prone to occur. Challenges of long-term dependency: due to the limited modeling ability of stream matching training for long-range time relationships, video clips longer than 20 seconds may experience audio and video synchronization drift. Although the industrial-grade Benchmark makes up for the lack of annotation, the quality of sound effects in niche scenes (such as cultural relic restoration and surgical operations) fluctuates due to insufficient training samples.

**Ethics** Emphasize that Kling-Foley is an auxiliary tool for sound effects artists, not a substitute. The sound effects generated by the model must be manually reviewed before commercial use to avoid the disruptive impact of technology on the traditional Foley industry. Use data enhancement technology to balance sound effect samples from different cultural backgrounds to avoid the model’s bias towards specific timbres or scenes (such as associating explosion sounds with male narration by default). Open text input to users to allow manual adjustment of sound effect styles and reduce the impact of built-in bias in the model on creation.

**Safety** Given the harmful social impact that the abuse of sound effects may have, we have implemented a number of security procedures in related products to prevent abuse throughout the development and potential deployment of the model. We have also implemented a multi-level watermarking scheme that is mandatory at all levels of content creation, such as embedding invisible watermarks in the spectrum of generated sound effects, supporting the tracing of content sources through self-developed tools, and synchronously generating visual watermarks to prevent the abuse of deep fake videos. In addition, it is prohibited to generate sound effects that may cause public panic (such as large-scale explosions, sirens), unless compliance approval is obtained for specific scenarios such as film and television production.

## References

- [1] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanyvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- [2] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024.
- [3] Kuaishou. Kling ai. <https://klingai.kuaishou.com/>, 2024.
- [4] 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.
- [5] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre D’efossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. *ArXiv*, abs/2209.15352, 2022.
- [6] 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*, 31:1720–1733, 2022.
- [7] Haohe Liu, Qiao Tian, Yiitan 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. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2871–2883, 2023.
- [8] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction-tuned llm and latent diffusion model. *ArXiv*, abs/2304.13731, 2023.
- [9] Jia-Bin 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*, abs/2305.18474, 2023.
- [10] Yazhou Xing, Yin-Yin He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7151–7161, 2024.
- [11] Heng Wang, Jianbo Ma, Santiago Pascual, Richard Cartwright, and Weidong (Tom) Cai. V2a-mapper: A lightweight solution for vision-to-audio generation by connecting foundation models. *ArXiv*, abs/2308.09300, 2023.
- [12] Zhifeng Xie, Shengye Yu, Qile He, and Mengtian Li. Sonic visionlm: Playing sound with vision language models. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26856–26865, 2024.
- [13] Yujin Jeong, Yunji Kim, Sanghyuk Chun, and Jiyoung Lee. Read, watch and scream! sound generation from text and video. In *AAAI Conference on Artificial Intelligence*, 2024.
- [14] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L. Berg. Visual to sound: Generating natural sound for videos in the wild. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3550–3558, 2017.
- [15] Ho Kei Cheng, Masato Ishii, Akio Hayakawa, Takashi Shibuya, Alexander Schwing, and Yuki Mitsufuji. Taming multimodal joint training for high-quality video-to-audio synthesis. *arXiv preprint arXiv:2412.15322*, 2024.
- [16] 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 *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE, 2017.

- [17] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 721–725. IEEE, 2020.
- [18] 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 Conference on Artificial Intelligence*, 2017.
- [19] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *ArXiv*, abs/2010.05646, 2020.
- [20] Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *Neural Information Processing Systems*, 2017.
- [21] Alon Ziv, Itai Gat, Gaël Le Lan, Tal Remez, Felix Kreuk, Alexandre D’efossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. Masked audio generation using a single non-autoregressive transformer. *ArXiv*, abs/2401.04577, 2024.
- [22] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021.
- [23] Haohe Liu, Zehua Chen, Yitian Yuan, Xinhao Mei, Xubo Liu, Danilo P. Mandic, Wenwu Wang, and Mark D. Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. In *International Conference on Machine Learning*, 2023.
- [24] Rongjie Huang, Jia-Bin Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiaoyue Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. *ArXiv*, abs/2301.12661, 2023.
- [25] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.
- [26] Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416, 2022.
- [27] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose M. R. Sotelo, Aaron C. Courville, and Yoshua Bengio. Samplernn: An unconditional end-to-end neural audio generation model. *ArXiv*, abs/1612.07837, 2016.
- [28] Vladimir E. Iashin and Esa Rahtu. Taming visually guided sound generation. *ArXiv*, abs/2110.08791, 2021.
- [29] Roy Sheffer and Yossi Adi. I hear your true colors: Image guided audio generation. *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2022.
- [30] 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 *International Conference on Machine Learning*, 2021.
- [31] Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. *ArXiv*, abs/2306.17203, 2023.
- [32] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 721–725, 2020.

- [33] Shentong Mo, Jing Shi, and Yapeng Tian. Text-to-audio generation synchronized with videos. *ArXiv*, abs/2403.07938, 2024.
- [34] Yiming Zhang, Yicheng Gu, Yanhong Zeng, Zhenning Xing, Yuancheng Wang, Zhizheng Wu, and Kai Chen. Foleydrafter: Bring silent videos to life with lifelike and synchronized sounds. *ArXiv*, abs/2407.01494, 2024.
- [35] 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 *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15180–15190, 2023.
- [36] Yiming Zhang, Yicheng Gu, Yanhong Zeng, Zhenning Xing, Yuancheng Wang, Zhizheng Wu, and Kai Chen. Foleydrafter: Bring silent videos to life with lifelike and synchronized sounds. *arXiv preprint arXiv:2407.01494*, 2024.
- [37] Hassan Akbari, Linagzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. *ArXiv*, abs/2104.11178, 2021.
- [38] Manjie Xu, Chenxing Li, Yong Ren, Rilin Chen, Yu Gu, Weihan Liang, and Dong Yu. Video-to-audio generation with hidden alignment. *ArXiv*, abs/2407.07464, 2024.
- [39] Ilpo Viertola, Vladimir E. Iashin, and Esa Rahtu. Temporally aligned audio for video with autoregression. *ArXiv*, abs/2409.13689, 2024.
- [40] Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jia-Bin Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. Frieren: Efficient video-to-audio generation network with rectified flow matching. In *Neural Information Processing Systems*, 2024.
- [41] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- [42] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022.
- [43] Yuan Gong, Cheng-I Lai, Yu-An Chung, and James Glass. Ssast: Self-supervised audio spectrogram transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10699–10709, 2022.
- [44] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Information Processing Systems*, 35:28708–28720, 2022.
- [45] Dading Chong, Helin Wang, Peilin Zhou, and Qingcheng Zeng. Masked spectrogram prediction for self-supervised audio pre-training. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [46] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio compression. *Transactions on Machine Learning Research*.
- [47] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved rvqgan. *Advances in Neural Information Processing Systems*, 36:27980–27993, 2023.
- [48] Yi-Chiao Wu, Israel D Gebru, Dejan Marković, and Alexander Richard. Audiodec: An open-source streaming high-fidelity neural audio codec. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [49] 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.

- [50] Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. AudioLDM 2: Learning holistic audio generation with self-supervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [51] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction guided latent diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3590–3598, 2023.
- [52] Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. Tango 2: Aligning diffusion-based text-to-audio generations through direct preference optimization. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 564–572, 2024.
- [53] 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 *International Conference on Machine Learning*, pages 13916–13932. PMLR, 2023.
- [54] 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.
- [55] Ziqian Ning, Huakang Chen, Yuepeng Jiang, Chunbo Hao, Guobin Ma, Shuai Wang, Jixun Yao, and Lei Xie. DiffRhythm: Blazingly fast and embarrassingly simple end-to-end full-length song generation with latent diffusion. *arXiv preprint arXiv:2503.01183*, 2025.
- [56] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [57] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- [58] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021.
- [59] Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning robust audio representations from clip. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4563–4567. IEEE, 2022.
- [60] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 976–980. IEEE, 2022.
- [61] 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 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [62] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- [63] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

- [64] Jiawei Ma, Po-Yao Huang, Saining Xie, Shang-Wen Li, Luke Zettlemoyer, Shih-Fu Chang, Wen-Tau Yih, and Hu Xu. Mode: Clip data experts via clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26354–26363, 2024.
- [65] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Synchformer: Efficient synchronization from sparse cues. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5325–5329. IEEE, 2024.
- [66] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [67] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- [68] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [69] Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- [70] Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent audio diffusion. In *Forty-first International Conference on Machine Learning*, 2024.
- [71] Zizheng Liu, Kexin Wang, Xiyang Huang, et al. Tango: Text-to-audio generation with hierarchical diffusion. *arXiv preprint arXiv:2307.05474*, 2023.
- [72] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- [73] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [74] Chunyu Qiang, Wang Geng, Yi Zhao, Ruibo Fu, Tao Wang, Cheng Gong, Tianrui Wang, Qiuyu Liu, Jiangyan Yi, Zhengqi Wen, et al. Vq-ctap: Cross-modal fine-grained sequence representation learning for speech processing. *IEEE Transactions on Audio, Speech and Language Processing*, 2025.
- [75] Chunyu Qiang, Hao Li, Yixin Tian, Ruibo Fu, Tao Wang, Longbiao Wang, and Jianwu Dang. Learning speech representation from contrastive token-acoustic pretraining. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10196–10200. IEEE, 2024.
- [76] Chunyu Qiang, Peng Yang, Hao Che, Xiaorui Wang, and Zhongyuan Wang. Style-label-free: Cross-speaker style transfer by quantized vae and speaker-wise normalization in speech synthesis. In *2022 13th International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 61–65, 2022.
- [77] Chunyu Qiang, Peng Yang, Hao Che, Ying Zhang, Xiaorui Wang, and Zhongyuan Wang. Improving prosody for cross-speaker style transfer by semi-supervised style extractor and hierarchical modeling in speech synthesis. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [78] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 119–132, 2019.
- [79] Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: An audio captioning dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 736–740. IEEE, 2020.

- [80] Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [81] Abul Ehtesham, Saket Kumar, Aditi Singh, and Tala Talaei Khoei. Movie gen: Swot analysis of meta’s generative ai foundation model for transforming media generation, advertising, and entertainment industries. In *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 00189–00195. IEEE, 2025.
- [82] Jaesung Huh, Jacob Chalk, Evangelos Kazakos, Dima Damen, and Andrew Zisserman. Epic-sounds: A large-scale dataset of actions that sound. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [83] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al. Attention is all you need. In *NeurIPS*, 2017.
- [84] Zeghidour Chen et al. Audiolum: a language modeling approach to audio generation. *arXiv preprint arXiv:2209.03143*, 2022.
- [85] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *ICML*, 2021.
- [86] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- [87] Tom B Brown et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [88] Chenglin Yang, Celong Liu, Xueqing Deng, Dongwon Kim, Xing Mei, Xiaohui Shen, and Liang-Chieh Chen. 1.58-bit flux. *arXiv preprint arXiv:2412.18653*, 2024.
- [89] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:2880–2894, 2020.
- [90] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15180–15190, 2023.
- [91] Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs. In *2001 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 01CH37221)*, volume 2, pages 749–752. IEEE, 2001.
- [92] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey. Sdr-half-baked or well done? In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 626–630. IEEE, 2019.
- [93] Emmanuel Vincent, Rémi Gribonval, and Cédric Févotte. Performance measurement in blind audio source separation. *IEEE transactions on audio, speech, and language processing*, 14(4):1462–1469, 2006.