

AudioLDM 2: Learning Holistic Audio Generation with Self-supervised Pretraining

Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong,
Yuping Wang, Wenwu Wang, Yuxuan Wang, Mark D. Plumbley

Abstract—Although audio generation shares commonalities across different types of audio, such as speech, music, and sound effects, designing models for each type requires careful consideration of specific objectives and biases that can significantly differ from those of other types. To bring us closer to a unified perspective of audio generation, this paper proposes a framework that utilizes the same learning method for speech, music, and sound effect generation. Our framework introduces a general representation of audio, called “language of audio” (LOA). Any audio can be translated into LOA based on AudioMAE, a self-supervised pre-trained representation learning model. In the generation process, we translate any modalities into LOA by using a GPT-2 model, and we perform self-supervised audio generation learning with a latent diffusion model conditioned on LOA. The proposed framework naturally brings advantages such as in-context learning abilities and reusable self-supervised pretrained AudioMAE and latent diffusion models. Experiments on the major benchmarks of text-to-audio, text-to-music, and text-to-speech demonstrate state-of-the-art or competitive performance against previous approaches. Our code, pretrained model, and demo are available at <https://audioldm.github.io/audioldm2>.

Index Terms—audio generation, diffusion model, self-supervised learning, speech synthesis, AIGC

I. INTRODUCTION

ARTIFICIAL intelligence generated content (AIGC) refers to any digital content such as images, videos, text, or audio that has been fully or partially created by an AI system without human involvement in the creative process [1]. Of particular interest is the ability of AI to produce audio content based on text, phonemes, or images [2]–[4]. AI-based audio generation has a wide potential in applications including synthesizing human or artificial voices for digital assistants [5], generating sound effects and background music for movies, and games [6], and automating the production of podcasts and audiobooks [7].

AI-based audio generation is often undertaken in separate sub-domains, such as the generation of speech [2], music [8], sound effects [4], and specific types of sounds such as footsteps and violin sounds [9], [10]. To address the specific challenges in each sub-domain, most previous works design task-specific inductive biases, which are predefined constraints that guide the learning process to a specific problem space. For

example, pitch and duration predictors are often used in speech synthesis to model the prosody of speech [2], [11], while MIDI representation [12] and domain-specific pre-trained modules are often used in music generation [8], [13].

Despite significant progress being made in developing specialized models for specific sub-domains of audio generation, the limitations of such specialization restrict the broader application of audio-generation models in complex auditory scenarios. The question of whether a unified approach can be developed to generate various types of audio signals still remains unanswered. In real-world cases, such as in movie scenes, different types of sound can occur simultaneously, requiring a more general approach to modelling audio generation. While there are works that address audio generation in a general domain, they mostly focus on generating audio with correct audio events with limited attention to detail. For example, previous text-to-audio generation research tends to generate unintelligible speech [4], [14], [15]. Moreover, while inductive biases have been useful in addressing the challenges of specific sub-domains, conclusions about a specific design drawn from one domain may not necessarily transfer to another. Recent advancements in addressing problems from a unified perspective have yielded substantial progress [16]–[19]. This trend highlights the potential of constructing a unified audio generation framework.

This paper presents a novel and versatile framework, called *AudioLDM 2*, that can generate any type of audio with flexible conditions, without the need for domain-specific inductive bias. The core idea is to introduce a new “language of audio” (LOA), which is a sequence of vectors that represent the semantic information of an audio clip. This approach allows us to translate human-understandable information into LOA and synthesize audio representation conditioned on LOA. This idea is similar to the use of onomatopoeia in [20] to describe environmental sounds. However, although onomatopoeia can effectively mimic certain sounds like animal noises or simple actions (e.g., “splash” for water), it can not encompass the full range of audio nuances. In theory, the “language of audio” should be able to represent both fine-grained acoustic information (e.g., “what does the speaker say”) and coarse-grained semantic information (e.g., “what is that sound”). Considering these requirements, we propose to leverage the features extracted by an audio masked autoencoder (AudioMAE) [21], an audio-generative self-supervised pretraining framework. An AudioMAE is pre-trained on diverse audio content, and its dual generative and reconstructive pre-training approach makes it potentially a strong option for representing

Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Wenwu Wang, and Mark D. Plumbley are with the Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey, Guilford, UK. Email: {haohe.liu, yi.yuan, xubo.liu, x.mei, w.wang, m.plumbley}@surrey.ac.uk.

Qiao Tian, Qiuqiang Kong, Yuping Wang and Yuxuan Wang: are with the Speech, Audio & Music Intelligence (SAMI) Group, ByteDance Inc. Email: {tianqiao.wave, kongqiuqiang, wangyuping, wangyuxuan.11}@bytedance.com.

audio in generative tasks.

Specifically, we utilize a GPT-2 language model [22] to translate conditioning information into AudioMAE features. The input conditions for the GPT-2 are flexible, including the representation of text, audio, image, video, and so on. We then use a latent diffusion model [23] to synthesize audio based on the AudioMAE features. The latent diffusion model can be optimized in a self-supervised manner, allowing for pre-training with large-scale unlabelled audio data. Our language-modelling approach with GPT-2 enables us to leverage recent advancements in language models [24], while alleviating challenges such as high inference computation costs and error accumulation that appeared in previous audio autoregressive models [8], [25]. The improvement is largely attributed to the shorter length of the LOA sequence. The continuous nature of LOA also potentially provides a richer representation power than the discrete tokens used in previous models [8], [13], [26].

Our experimental results demonstrate that *AudioLDM 2* achieves state-of-the-art (SoTA) performance on text-to-audio (TTA), and text-to-music (TTM) generation tasks, when evaluated on AudioCaps [27] and MusicCaps [8], respectively. On text-to-speech (TTS) generation tasks, *AudioLDM 2* achieves performance comparable with the SoTA by significantly outperforming a strong baseline FastSpeech2 [11]. In addition to using text conditions, we showcase the capability of utilizing visual modality conditions for audio generation, such as image-to-audio generation. Moreover, we explore some of the peripheral functionalities of *AudioLDM 2*, such as in-context learning for audio, music, and speech. In comparison to the original AudioLDM [4], *AudioLDM 2* is able to self-supervised pretraining on the latent diffusion model, and enjoy the benefit of auto-regressive modeling of LOA with GPT-2 model. Besides, while retaining the same ability, *AudioLDM 2* shows substantial advancements over AudioLDM in quality, versatility, and capacity to generate speech with intelligible content. Overall, our contributions are as follows:

- We propose a novel and versatile audio generation model that is capable of performing conditional generation of audio, music, and intelligible speech.
- The proposed method is based on a universal representation of audio, which enables large-scale self-supervised pretraining of the core latent diffusion model without audio annotation and helps to combine the advantages of both the auto-regressive and the latent diffusion model.
- Our experiments show *AudioLDM 2* achieves state-of-the-art (SoTA) performance in text-to-audio and text-to-music generation, while also delivering competitive results in text-to-speech generation, comparable to the current SoTA.

II. RELATED WORK

A. Conditional Audio Generation

Audio generation is an emerging topic that focuses on modelling the generation of general audio, including recent models such as AudioGen [3], AudioLDM [4], and Make-an-Audio [15]. AudioGen treats audio generation as a conditional language modelling task, while the other two works

approach this task by latent diffusion. Studies on image-to-audio and video-to-audio generation, such as Im2Wav [28] and SpecVQGAN [29], are also areas of interest to researchers. Additionally, there are audio generation approaches that do not rely on conditioning, such as AudioLM [26], which performs audio language modelling based on a neural codec. Even though audio generation usually includes the topic of speech generation, previous works on text-to-audio generation tend to generate unintelligible speech [3], [4], [14], [15].

The field of audio generation encompasses sub-domains such as text-to-speech (TTS) and text-to-music (TTM). The former focuses on generating speech signals from transcriptions, while the latter involves creating a music clip from a textual description. Cutting-edge TTS models like FastSpeech2 [11], GradTTS [30], and NaturalSpeech [2] have made significant strides, producing speech of such high quality that it is nearly indistinguishable from human speech. Various techniques have been introduced to address speech generation in TTS, such as the monotonic alignment algorithm [31], which aligns phoneme features with spectrogram features, and a prosody predictor [11], used to guide model training and enhance expressiveness. Recent advances in TTM are evident in models like MusicLM [8], Noise2Music [32], MusicGen [33], and MeLoDy [13]. Similar to AudioLDM, the MusicLM model aligns music and language embeddings through contrastive pretraining modules, which enables text-free model optimization, alleviating the scarcity of music-text pairs. MusicLM also includes a semantic modeling stage based on w2v-BERT [34] to enhance the model performance. MusicGen uses a language modeling approach for music generation, enhanced with a mechanism for conditioning the model with melodic features for improved controllability. Meanwhile, MeLoDy, a diffusion model guided by language modeling, achieves significant computational reduction in music generation compared to MusicLM.

In this paper, we propose a unified framework for audio generation, which encompasses a breadth of topics including, but not limited to, speech, sound effect, and music generation.

B. Diffusion Models

Diffusion models [35], [36] have demonstrated high sample quality in a variety of tasks including image generation [37]–[39], image restoration [40], speech generation [41]–[43], and video generation [44], [45]. In the realm of speech or audio synthesis, these models have been explored for both mel-spectrogram generation [30], [46] and waveform generation [47]–[49]. However, the iterative nature of generation in a high-dimensional data space often results in slow training and inference speeds. One solution involves the use of diffusion models in a more restricted latent space, a strategy exemplified in image generation [23]. This idea has been adopted in various audio generation works, including AudioLDM [4], Make-An-Audio [15], and TANGO [50]. These works utilize latent diffusion models trained on a continuous latent space. On the other hand, there are also studies that explore diffusion in the discrete latent space. For instance, DiffSound [14] employs a discrete autoencoder to mitigate redundancy [51], [52] in the

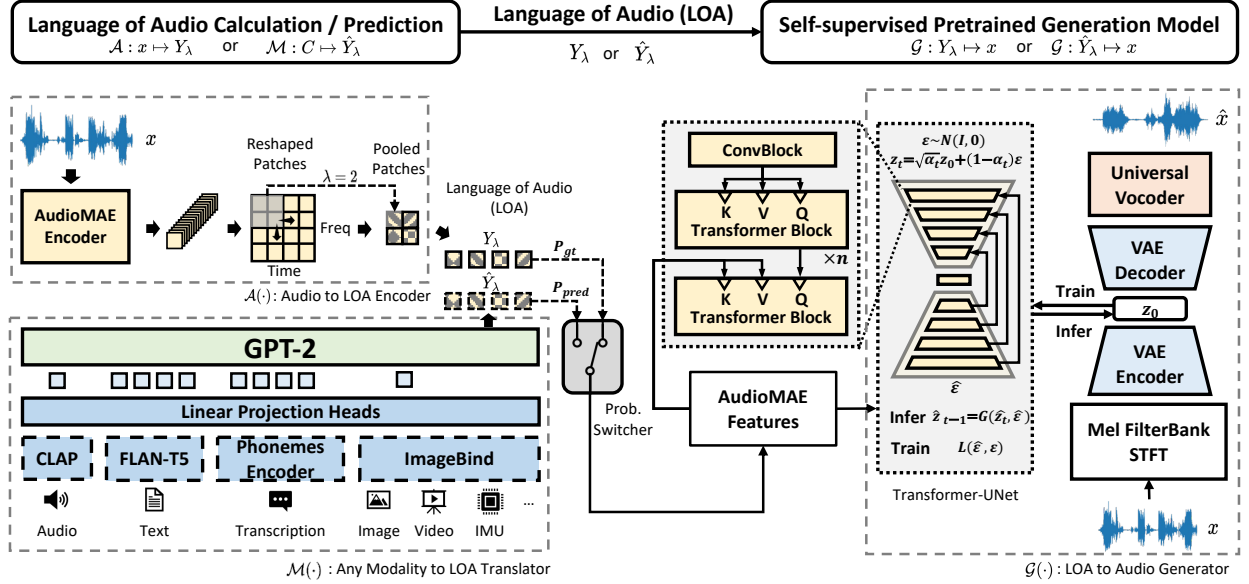


Fig. 1. The overview of the *AudioLDM 2* architecture. The AudioMAE feature is a proxy that bridges the *conditioning information to LOA translation* stage (modelled by GPT-2) and the *LOA to audio generation* stage (modelled by the latent diffusion model). The probabilistic switcher controls the probability of the latent diffusion model using the ground truth AudioMAE (P_{gt}) and the GPT-2 generated AudioMAE feature (P_{pred}) as the condition. Both the AudioMAE and latent diffusion models are self-supervised pre-trained with audio data.

audio waveform and create a compressed representation of mel-spectrograms. DiffSound utilizes text-conditional discrete diffusion models to generate discrete tokens.

III. AUDIOLDM 2

A. Overview

Let $x \in \mathbb{R}^{L_s}$ represent an audio signal, where L_s is the length of the audio samples in x . An audio generation process can be denoted as $\mathcal{H} : C \mapsto x$, where C is the conditioning information and \mathcal{H} is the conditional audio generation system. The condition C can be flexible, including text [4], image [53], video [29], and functional magnetic resonance imaging (fMRI) [54].

The direct generation of x from C is usually challenging [55]. Motivated by regeneration learning [56], we propose to utilize an intermediate feature Y , as an abstraction of x , to bridge the gap between C and x , as introduced in Section III-B1. We call the feature Y the language of audio (LOA). The LOA feature is calculated by $Y = \mathcal{A}(x)$ in which \mathcal{A} performs audio to LOA encoding either by hand-crafted rules or self-supervised representation learning [56]. As illustrated in Figure 1, with the intermediate representation Y , the overall audio generation process can be denoted as

$$\mathcal{H}_0 = \mathcal{G} \circ \mathcal{M} : C \mapsto \hat{Y} \mapsto x, \quad (1)$$

where \hat{Y} is the estimation of the ground truth LOA. As denoted in (1), the audio generation process of *AudioLDM 2* includes the following two steps:

(i) *Conditioning information to LOA translation*: The function $\mathcal{M} : C \mapsto \hat{Y}$ aims to produce the LOA Y based on C , which could be the conditional information from any modality, such as image and text. As a potentially better representation

of C , the generated \hat{Y} will be used in later stages as the conditioning information for audio generation.

(ii) *LOA to audio generation*: Followed by \mathcal{M} , function \mathcal{G} accepts an LOA estimation \hat{Y} as input condition and estimates the audio data x . During the training process, when the training data x is available, the ground truth Y will be also available using $\mathcal{A}(\cdot)$, allowing the optimization of \mathcal{G} in a self-supervised manner. Specifically, instead of using the LOA estimation $\hat{Y} = \mathcal{M}(C)$, we condition the generation of x based on the $Y = \mathcal{A}(x)$, which can be formulated as

$$\mathcal{H}_1 = \mathcal{G} \circ \mathcal{A} : x \mapsto Y \mapsto \hat{x}. \quad (2)$$

We introduce the detail of $\mathcal{A}(\cdot)$ in Section III-B. Since the process \mathcal{H}_1 only involves x as the training data, Equation (2) means model \mathcal{G} can be optimized in a self-supervised manner without any audio annotation. This self-supervised scheme can alleviate the scarcity of the audio data labels [4] and provide a robust backbone for the overall generation system.

The following sections provide a detailed introduction to *AudioLDM 2*. In Section III-B, we discuss the audio representations employed in *AudioLDM 2*, including the AudioMAE and VAE features. These features also serve as the generation targets for the two stages within *AudioLDM 2*. Section III-C introduces the auto-regressive modeling of the AudioMAE feature with GPT-2. In Section III-D, we elucidate the process of generating audio waveforms via the latent diffusion model, which applies a VAE for feature compression and generates audio conditioned on the LOA. The LOA here can be based on either ground truth or GPT-2-generated data, which corresponds to self-supervised training and joint training with GPT-2 (Section III-D3), respectively.

B. Audio Representation Learning

1) Semantic Representation Learning with the AudioMAE:

To accurately represent diverse types of audio, encompassing speech, music, and sound effects, the LOA Y should effectively capture both the semantic and the acoustic details of audio signals. Therefore, we propose to use a self-supervised pretrained AudioMAE [21] as the representation extraction module for function \mathcal{A} for its generality and high accuracy on the downstream audio classification task [21].

The audio masked autoencoder (AudioMAE) is an audio self-supervised pre-training model, which learns representations from unlabeled audio data without relying on manually labeled annotations. An AudioMAE consists of an encoder and a decoder, both realized with an architecture similar to the vision transformer (ViT) [57]. During self-supervised pre-training, input patches to the encoder, which are usually mel spectrograms, are randomly masked and the decoder learns to reconstruct the masked patches [21]. Compared with other audio self-supervised pretraining models, AudioMAE has the following two advantages:

(i) *The AudioMAE has been verified to work well in the general audio domain.* For example, an AudioMAE can be effectively pre-trained on AudioSet [58], with state-of-the-art performance on the downstream audio classification tasks. In comparison, typical audio self-supervised models focus on a specific domain, such as the MERT [59] on music and the HuBERT [60] on speech.

(ii) *AudioMAE features are potentially better for generative tasks than other discriminative pre-training methods.* Building upon the contrastive loss or next token prediction classification loss as the learning objective, previous systems such as wav2vec [61] and BYOL-A [62] utilize a discriminative approach during pre-training. In comparison, AudioMAE focuses on a generative process by learning the reconstruction of the masked patches.

For an input audio signal x , AudioMAE first calculates the log mel spectrogram $X \in \mathbb{R}^{T \times F}$, where T represents the time steps of the mel spectrogram, and F denotes the mel bins. The mel spectrogram X is then treated as an image and split into patches each of size $P \times P$, serving as the input for the AudioMAE encoder. The patch size P is typically designed to be a common factor of T and F . Patch splitting and embedding are performed using a convolutional neural network with a kernel size of P , a stride of P , and D output channels. This yields an output shape of $T' \times F' \times D$, where D is the AudioMAE embedding dimension, $T' = T/P$, and $F' = F/P$. The resulting output feature of the AudioMAE encoder, $E \in \mathbb{R}^{T' \times F' \times D}$, has the same shape as the input and is usually treated as the feature for downstream tasks after pretraining [21].

2) *AudioMAE Feature Post Processing:* As shown in Figure 1, once the AudioMAE features E are computed, we introduce an additional pooling step to aggregate E into Y_λ , where $\lambda \in I^+$ represents a hyper-parameter used in the post-processing pooling step. This pooling step aims to reduce the sequence length, facilitating easier estimation in the function \mathcal{M} . Specifically, we perform a two-dimensional average-max pooling [51] on the first two dimensions of $E \in \mathbb{R}^{T' \times F' \times D}$, in

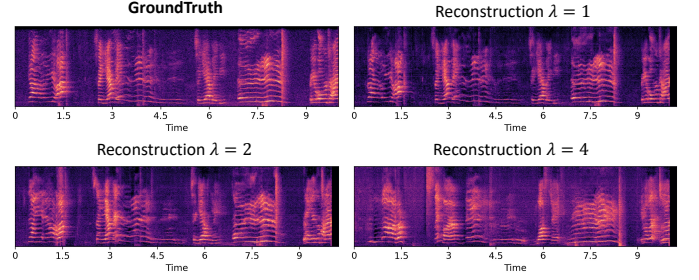


Fig. 2. The influence of λ on audio reconstruction from LOA Y_λ . The reconstruction closely resembles the ground truth when $\lambda = 1$, suggesting that $Y_{\lambda=1}$ retains sufficient audio details. However, with $\lambda = 2$ or 4, the reconstruction diverges slightly from the original audio, indicating that while the post-processed AudioMAE feature may not include all details, it nonetheless accurately preserves semantic content.

which the pooling kernel size and stride have the same value $\lambda \in I^+$. The two-dimensional pooling operation can help to preserve the time-frequency relationship in the output. The final output after pooling, Y_λ , is reshaped into an embedding sequence with shape $L_\lambda \times D$, in which $L_\lambda = T'F'/\lambda^2$. To facilitate implementation, The value of λ is chosen in a way that L_λ is always a positive integer. We demonstrate the effect of different choices of λ in Figure 2. In the remaining sections of this paper, if λ is not specified, we'll refer to Y_λ simply as Y .

3) *Acoustic Representation Learning with VAE:* We use a VAE for feature compression and for learning an audio representation z , which has a significantly smaller dimension than x [4]. The VAE we used in this work is a convolutional architecture that consists of encoders with down-sampling and decoders with up-sampling following the architecture described in [4]. The forward pass of the VAE can be formulated as $\mathcal{V} : X \mapsto z \mapsto \hat{X}$, where X is the mel-spectrogram of x and \hat{X} is the reconstruction of x . The reconstruction \hat{X} can be converted to the audio waveform \hat{x} using a pre-trained HiFiGAN vocoder [63]. Following AudioLDM [4], we calculate a reconstruction loss and a discriminative loss based on X and \hat{X} to optimize the parameters of the VAE. We also calculate the KL divergence between z and a standard Gaussian ($\mu = 0$, $\sigma^2 = 1$) as a loss function to limit the variance of the VAE latent space.

4) *Comparison between AudioMAE and VAE:* Since both AudioMAE and VAE are based on autoencoders for representation learning, one might wonder why we use a VAE for representation learning instead of directly modeling the AudioMAE latent space. Part of the reason is that AudioMAE does not primarily focus on reconstruction quality, and its latent space compression ratio is not as high as that of the VAE. On the other hand, the VAE exhibits good reconstruction ability and a higher compression level than AudioMAE, making VAE more suitable for mel-spectrogram compression. Furthermore, as shown in Figure 3, we visualize the latent representation of AudioMAE and VAE on the ESC-50 [64] dataset using tSNE [65]. The visualization demonstrates that the latent representation of AudioMAE can group similar audio at a closer region in the latent space, whereas the representation of VAE exhibits more overlap between different

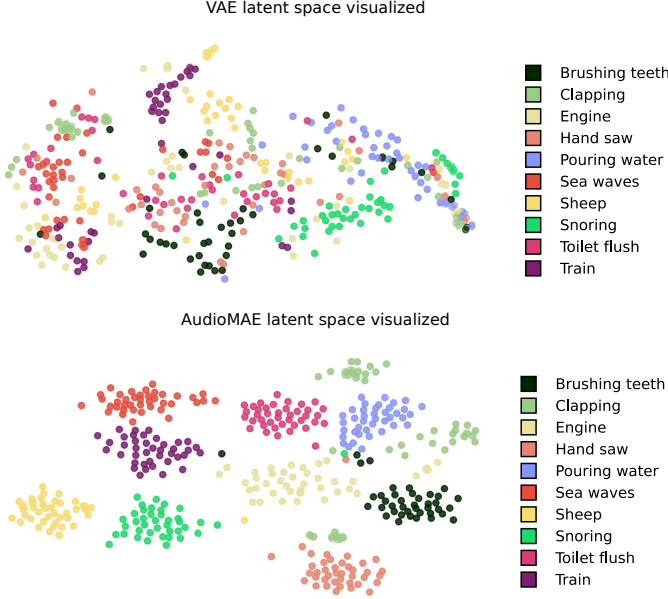


Fig. 3. Visualization of the latent space based on tSNE and ten randomly selected classes in the ESC50 [64] dataset. Each point in the figure represents an audio clip. The AudioMAE feature space tends to group similar audio clips together, indicating more semantic structure than in the VAE feature.

audio classes. This indicates that the representations for the AudioMAE and VAE are distinct. AudioMAE contains more information on the semantic side, while VAE representation is less semantically structured.

C. Conditioning Information to LOA Translation with GPT-2

This subsection introduces the design of the function \mathcal{M} . As introduced in Section III-A, the input to the model $\mathcal{G} : Y \mapsto x$ can be calculated using the AudioMAE. However, during inference, when we perform audio generation with the condition C , the ground truth LOA $Y = \mathcal{A}(x)$ is not available. Therefore, we need another model that can generate \hat{Y} given C , denoted by $\mathcal{M}_\theta : C \rightarrow \hat{Y}$, where θ represents trainable parameters.

Specifically, we treat the generation of Y as a language modelling task and choose the GPT-2 (Generative Pre-trained Transformer 2) [22] model as the backbone. GPT-2 is based on a transformer architecture and was originally trained on 8 million documents for a total of 40 GB of text using an unsupervised learning approach [22]. GPT-2 has been used in a variety of natural language processing tasks, such as text completion, question answering, and language translation [66], [67]. Initialized with pre-trained weights, we finetune the GPT-2 model based on teacher forcing [68], so that during model training, \hat{y}_l will be generated based on both the condition C and the ground truth sequence y_1, \dots, y_{l-1} , where y_l is the l -th vector in LOA sequence Y . Specifically, the GPT-2 model \mathcal{M}_θ is trained to maximize the likelihood of a sequence $Pr(y_1, y_2, \dots, y_L | C)$, which can be interpreted into the following optimization objective:

$$\arg\max_\theta \mathbb{E}_C [Pr(y_1 | C; \theta) \prod_{l=2}^L Pr(y_l | y_1, \dots, y_{l-1}, C; \theta)], \quad (3)$$

where \mathbb{E}_C represents the expectation operator with respect to the variable C . We calculate the mean squared error loss [55] between y_l and $\hat{y}_l = \mathcal{M}_\theta(y_1, \dots, y_{l-1}, C)$ to optimize Equation (3). We directly optimize the regression of continuous vectors y_l , without discretizing the AudioMAE feature space and estimating the token index. The condition C in Equation (3) can encompass a flexible range of data representations, including audio representations, text embeddings, phoneme embeddings, or visual clues. We adopt the mixture of experts [69] approach and use multiple encoders as feature extractors to calculate C . Given K systems as the feature extraction modules, the shape of the output from the k -th system $C_k, k \in \{1, \dots, K\}$ is $L_k \times D_k$, in which L_k is the sequence length of the k -th system and D_k is the dimension of the feature. We apply a linear transformation layer after the output of each feature extraction module to unify the embedding dimension to D_0 for easier process of the GPT-2 model. For modules that extract global features from the input without sequential information, such as CLAP [70] or ImageBind [18], we have $L_k = 1$. The final condition $C = [C_1, \dots, C_K]$ is a concatenation of C_k along the sequence length dimension. The final condition C has a shape of $L \times D_0$, where $L = \sum_{k=1}^K L_k$. We introduce several condition modules we used in this paper as follows.

CLAP or contrastive language and audio pretraining [70], is a system that learns a joint audio-text embedding space, in which paired audio and language data have a closer distance in the latent space. CLAP has been successfully applied as a conditioning module to audio generation such as AudioLDM [4]. In this study, we employ a pre-trained CLAP¹ text encoder as the default conditioning module for extracting text embeddings as conditions. However, in scenarios where text captions (e.g., “A man is speaking happily with background static noise”) are unavailable, such as for text-to-speech tasks, we use the CLAP audio encoder as the conditioning module instead of using CLAP text encoder, in the same way as [4].

FLAN-T5. The CLAP model, as a module that calculates global-level conditions, has been found to have issues in capturing the temporal information in the text data [71]. To allow for this, we use another pretrained text encoder to capture the semantic information of the textual input, which might contain useful details such as temporal orders. Specifically, we utilize FLAN-T5 [72], which is an enhanced version of the text-to-text transfer transformer (T5) model [73] based on the finetuning on a mixture of tasks².

Phoneme Encoder is a widely adopted module in text-to-speech research for extracting helpful information regarding phonemes [2], [11], which are the smallest units of sound in a language that can distinguish one word from another [74]. In this work, we follow the structure introduced in NaturalSpeech [2] to build a phoneme encoder, in the form of a stack of transformer encoder layers. We preprocess the textual input into phonemes using the open-source tool Espeak phonemizers³ and append a stop token after each phoneme

¹<https://github.com/LAION-AI/CLAP>

²<https://huggingface.co/google/flan-t5-large>

³<https://github.com/espeak-ng/espeak-ng>

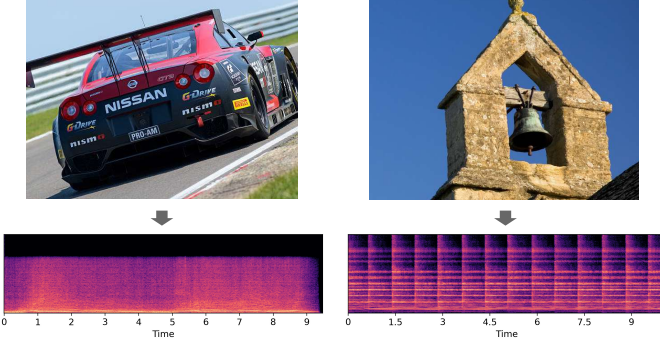


Fig. 4. Examples of image-to-audio generation.

sequence to mark the end of the sequence for the transformer model.

ImageBind [18] is a system with a similar idea to the CLAP model, but has more modalities aligned in a single embedding space, including image, text, video, audio, depth map, thermal map, and inertial measurement units (IMUs). This means that, once the conditional model is trained with one modality as a condition (e.g., images), other modalities can be used as conditions as well. Since the output of all pretrained encoders from ImageBind⁴ are aligned and have the same shape, we utilize an audio encoder branch to calculate an embedding as a condition during training and switch to the image encoder branch to calculate image embedding as a condition during inference.

Except for the phoneme encoder, which does not have a readily available pre-trained weights, the parameters of all other pre-trained feature extraction models are kept frozen during the experiment.

D. LOA to Audio Generation with Latent Diffusion Model

We model the process $\mathcal{G} : Y \mapsto x$ with a latent diffusion model (LDM) [23], which is a variant of the denoising diffusion probabilistic models (DDPM) [35]. In contrast to DDPM, which directly models the training data, the LDM learns the reverse diffusion process in a variational autoencoder (VAE)-based compressed latent space [75], which can reduce the computational cost. Similar ideas have been adapted to audio generation, such as AudioLDM [4].

1) *Latent Diffusion Model*: We follow the formulation in [35] to implement the LDM. Given a VAE representation z , the forward transition is defined as a T steps Markov process in a way that does not include trainable parameters. Given the data z_{t-1} at diffusion step $t-1$, the data distribution of z_t at step $t \in 2, \dots, T$ can be formulated as

$$q(z_t|z_{t-1}) = \sqrt{1 - \beta_t}z_{t-1} + \sqrt{\beta_t}\epsilon_t, \quad (4)$$

in which the noise schedule hyper-parameter $\beta_t \in [0, 1]$ determines how quickly the noise is blended into the data. By recursive substitution of $q(z_t|z_{t-1})$ in Equation (4) [35], we can derive the distribution of z_t given z_0 as

$$q(z_t|z_0) = \sqrt{\alpha_t}z_0 + \sqrt{1 - \alpha_t}\epsilon_t, \quad (5)$$

where $\alpha_t = \prod_{s=1}^t 1 - \beta_s$ and $\epsilon_t \sim N(0, I)$. At the final step $t = T$, the distribution of z_t will be close to a standard Gaussian distribution [35].

The LDM learns a backward transition from the prior distribution $N(0, I)$ to the data distribution z . The reverse process models the conditional distribution $Pr(z_{0:T}|Y; \phi) = Pr(z_0|z_1, Y; \phi) \prod_{t=2}^T Pr(z_{t-1}|z_t, Y; \phi) \cdot Pr(z_T)$, in which Y is the LOA as the condition signal and the ϕ denotes the parameter of the model for learning the reverse diffusion. If we marginalize $z_{1:T}$ we can derive the lower bound of $\log[Pr(z_0|Y; \phi)]$ based on the evidence lower bound (ELBO) and Bayes' rule [35]:

$$\log[Pr(z_0|Y; \phi)] \geq \log[Pr(z_0|z_1, Y; \phi)] - \sum_{t=2}^T KL[Pr(z_{t-1}|z_t, Y; \phi) || q(z_{t-1}|z_t, z_0)], \quad (6)$$

where $KL(\cdot)$ is the function for calculating KL divergence, and $q(z_{t-1}|z_t, z_0)$ is the target conditional diffusion distribution that has a closed-form solution given z_0 and z_t [35]. Following [35], we can derive the loss function that maximizes the lower bound of Equation (6) as:

$$\argmin_{\phi} [\mathbb{E}_{z_0, Y, t \sim \{1, \dots, T\}} ||\mathcal{G}(\sqrt{\alpha_t}z_0 + \sqrt{1 - \alpha_t}\epsilon_t, t, Y; \phi) - \epsilon_t||]. \quad (7)$$

As shown in Figure 1, we propose to use a Transformer-UNet (T-UNet) architecture as the function \mathcal{G} in Equation (7), which is an improved version of the UNet used in AudioLDM [4]. Similar to the UNet used in other works [76], the T-UNet architecture consists of a series of encoders with downsampling and a series of decoders with upsampling, and there are skip connections between encoders and decoders at the same scale. To enhance the modelling capacity of the T-UNet, we insert multiple transformer blocks after the convolution operation in each encoder and decoder block. Specifically, we have $n_{\text{trans}} + 1$ transformer blocks, in which the first n_{trans} transformer blocks are a stack of self-attention layers [77] and feed-forward networks. To incorporate the condition information Y from the ground truth LOA or \hat{Y} from $\mathcal{M}(\cdot)$ (Section III-C), as shown in Figure 1, the last transformer block changes the self-attention layer to cross-attention, which accepts the LOA as key and value and fuses with the feature from the previous transformer block as the query. Except for text-to-speech generation, we add an extra cross-attention layer in the transformer block to accept the text embedding from FLAN-T5 [72] as an extra condition to enhance the audio-text relationship learning.

2) *Classifier-free Guidance*: For diffusion models, controllable generation can be achieved by introducing guidance at each sampling step. Classifier-free guidance [78], [79] (CFG) has been the state-of-the-art technique for guiding diffusion models. During training, we randomly discard our condition Y in Equation (7) with a fixed probability (e.g., 10%) to train both the conditional LDMs $\mathcal{G}(z_t, t, Y; \phi)$ and

⁴<https://github.com/facebookresearch/ImageBind>

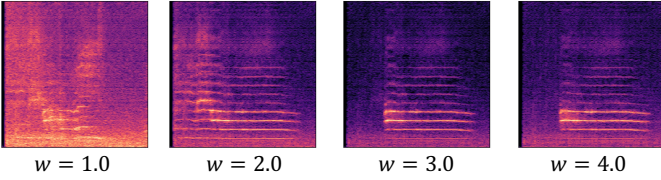


Fig. 5. The samples generated with different classifier-free guidance scales. The text prompt is “A cat is meowing”.

the unconditional LDMs $\mathcal{G}(z_t, t, \phi)$. For generation, we use LOA \hat{Y} or Y as the condition and perform sampling with a modified noise estimation $\mathcal{G}'(z_t, t, Y; \phi)$:

$$\mathcal{G}'(z_t, t, Y; \phi) = w\mathcal{G}(z_t, t, \phi) + (1 - w)\mathcal{G}(z_t, t, Y; \phi), \quad (8)$$

where w determines the guidance scale.

3) *Joint Finetuning*: We perform joint finetuning with the GPT-2 and latent diffusion models based on Equation (1), (7), and (3). As demonstrated by Table V, we found that joint finetuning significantly enhances the overall performance of the *AudioLDM 2* system. As depicted in Figure 1, the *probabilistic switcher* controls the source of the conditioning signal during the joint training process. During training, the switcher dynamically chooses between ground truth AudioMAE features and GPT-generated AudioMAE features, with probabilities set to P_{gt} and P_{pred} , respectively.

IV. EXPERIMENT SETUP

A. Dataset

The datasets used in this work include AudioSet [58], WavCaps [80], AudioCaps (AC) [27], VGGSound [81], Free Music Archive (FMA) [82], Million Song Dataset (MSD) [83], LJSpeech (LJS) [84], and GigaSpeech [85]. AudioSet is the largest audio classification dataset at the time of writing, with around two million ten-seconds of audio and 527 different classes. WavCaps is a dataset with ChatGPT-assisted weakly-labeled audio captions. WavCaps contains 403,050 audio clips with an average duration of 68 seconds. AudioCaps is a subset of AudioSet with handcrafted captions, containing about 46,000 ten-second audio clips. VGGSound is a large-scale single-label audio-visual dataset, which contains over 200,000 videos. We only utilize the audio data and the labels in the VGGSound. FMA is a large music dataset without captions, containing 106,574 music tracks from 16,341 artists and 14,854 albums. For the Million Song Dataset, we only utilize the labelled subset proposed in [86], which contains around 510,000 music tracks with metadata such as tags, titles, and artist names. LJSpeech is a single-speaker speech dataset with 13,100 short audio clips and detailed transcriptions. GigaSpeech is a multi-speaker large-scale English speech recognition corpus with around 10,000 hours of audio labeled with transcriptions. The test and development set of GigaSpeech are not included during training. All the audio data used in this work are resampled to 16 kHz for easier comparison with previous works [4], [15]. We use only the audio data with paired text labels to train the GPT-2 model by optimizing Equation (3). We train the latent diffusion model with all the

audio data regardless of annotation by optimizing the objective in Equation (6) in a self-supervised manner.

B. Evaluation Metrics

We mainly focus on the text-to-audio generation task to evaluate the effectiveness of *AudioLDM 2*. We follow the evaluation protocol of AudioGen [3], which calculates both objective metrics such as Frechet Audio Distance (FAD), Kullback-Leibler Divergence (KL), and subjective metrics including Overall Impression (OVL) and Audio and Text Relation (REL). FAD is a reference-free audio quality measure that is calculated based on the distribution distance between the feature of the target and generated audios, extracted from the VGGish [87] model. KL divergence measures the similarity between the generated and target audio with the label calculated by the audio tagging model, Patch-out Transformer [88], in the same way as AudioGen [3]. We use a similar evaluation protocol for text-to-music generation. For the text-to-speech task, we utilize the commonly used mean opinion score (MOS) for evaluation [74].

C. Subjective Evaluation

We use Amazon Mechanical Turk⁵, a crowd-sourced platform, to perform the subjective evaluation on metrics including OVL, REL, and MOS. The instructions on how to perform evaluation are clearly illustrated for the raters with examples. To ensure the credibility of the evaluation result, we set requirements for the crowd-source worker with a minimum average approval rate of 60% and with at least 50 approvals in the record. Each audio clip is evaluated by at least 10 different raters. All three subjective metrics have a Likert scale [89] between one and five, where a larger number indicates better performance. Study raters received payment at or above the US minimum wage. We average the scores among all raters and samples as the final score for a system.

D. Model Architecture Details

We perform the experiment with two sizes of the latent diffusion model, *AudioLDM 2* and *AudioLDM 2-Large*, with transformer layer numbers $n_{\text{trans}} = 2$ and $n_{\text{trans}} = 6$ (Section III-D), respectively. We use a pre-trained AudioMAE⁶ with a patch size of 16×16 and no overlapping, resulting in a 768-dimension feature sequence with length 512 for every ten seconds of mel spectrogram. In a similar way to the idea introduced in [90], on calculating the LOA Y , we gather the output of the last 16 transformer layers from the AudioMAE encoder and perform averaging as the final Y . The GPT-2 model we employ has an embedding dimension of 768 with 12 layers of transformers. For joint fine-tuning, we set the probability of using ground truth LOA Y and LOA estimation \hat{Y} as $P_{\text{gt}} = 0.25$, and $P_{\text{pred}} = 0.75$, respectively.

For the generation of audio and music, we combine the text embeddings from the CLAP text encoder and FLAN-T5 as conditioning and designate $Y_{\lambda=8}$ as the target sequence

⁵<https://requester.mturk.com/>

⁶<https://github.com/facebookresearch/AudioMAE>

TABLE I
PERFORMANCE COMPARISON ON THE AUDIOCAPS EVALUATION SET. *AudioLDM 2* OUTPERFORMS PREVIOUS APPROACHES BY A LARGE MARGIN ON BOTH SUBJECTIVE AND OBJECTIVE EVALUATION.

Model	Duration (h)	Param	FAD↓	KL↓	CLAP (%)↑	OVL ↑	REL ↑
GroundTruth	-	-	-	-	25.1	4.04	4.08
AudioGen-Large	6824	1 B	1.82	1.69	-	-	-
Make-an-Audio	3000	453 M	2.66	1.61	-	-	-
AudioLDM-Large-FT	9031	739 M	1.96	1.59	-	-	-
AudioLDM-M	9031	416 M	4.53	1.99	14.1	3.61	3.55
Make-an-Audio 2	3700	937 M	2.05	1.27	17.3	3.68	3.62
TANGO	145	866 M	1.73	1.27	17.6	3.75	3.72
<i>AudioLDM 2-AC</i>	145	346 M	1.67	1.01	24.9	3.88	3.90
<i>AudioLDM 2-AC-Large</i>	145	712 M	1.42	0.98	24.3	3.89	3.87

for GPT. The conditioning modules for speech generation are configured differently, primarily due to the need to better preserve the fine-grained phoneme information in speech signals through a smaller λ value. Thus, for speech generation, we incorporate both CLAP and the phoneme encoder, designating $Y_{\lambda=1}$ as the target sequence to retain more details. For the speech data, since there are no available audio captions (different from transcriptions), we adopt a similar approach as AudioLDM [4] to utilize the CLAP audio encoder to compute the embedding as a condition during model training, and employ the CLAP text encoder during inference. This method also facilitates prompt-based speaker control, as demonstrated in Figure 7.

E. Training and Inference Setup

The latent diffusion model and the GPT-2 model are initially trained separately. We randomly choose $\lambda \in \{1, 2, 4, 8\}$ during pre-training of the latent diffusion model to enhance the model robustness under conditions Y_{λ} with different λ . We train the latent diffusion model and finetune the GPT-2 model on eight NVIDIA A100 80GB GPUs. The setting of the latent diffusion model is mostly the same as that of AudioLDM [4], except that we change the default classifier-free guidance scale during the Denoising Diffusion Implicit Models (DDIM) [91] sampling to 3.5. For both GPT-2 finetuning and the latent diffusion model, we utilize the AdamW [92] optimizer with a learning rate of 10^{-4} and 10000 steps of linear warming up without decay.

V. RESULT

We evaluated our proposed system on three primary audio generation tasks: text-to-audio, text-to-music, and text-to-speech. The three basic systems were trained on three different datasets: AudioCaps (general audio), MSD (music), and LJSpeech (speech), and are denoted as *AudioLDM 2-AC*, *AudioLDM 2-MSD*, and *AudioLDM 2-LJS*, respectively. The model *AudioLDM 2-Full* represents a version capable of performing both audio and music generation simultaneously, with training data scaled up to 29510 hours. In contrast with AudioLDM [4], we do not perform additional model finetuning on AudioCaps for model trained with the full-scale datasets. Models with the suffix *Large* indicate larger-sized model variants, such as *AudioLDM 2-Full-Large*.

A. Text-to-Audio Generation

We compare the performance of our proposed model with several state-of-the-art systems, including AudioGen-Large [3], Make-an-Audio [15], AudioLDM [4], Make-an-Audio 2 [93], and TANGO [50]. To generate the samples for subjective evaluation, we adopt AudioLDM-M, an AudioLDM with 652M parameters, from HuggingFace⁷ and run with 100 reverse diffusion steps. The result of Make-an-Audio 2 is provided by the author [93]. We use the pre-trained TANGO model open-sourced on GitHub⁸ to reproduce their result.

As shown in Table I, our proposed *AudioLDM 2-AC* significantly outperforms the previous systems across all three objective metrics. The previous best-performing system, TANGO, achieves a CLAP score of 17.6, while our proposed system surpasses it with a substantially higher CLAP score of 24.9. *AudioLDM 2-Large* also attains the best KL divergence score of 0.98, considerably improving upon the previous SoTA of 1.27. For the FAD score, our model reaches 1.42, establishing a new SoTA for text-to-audio generation. Our subjective evaluation results are mostly consistent with the objective metrics, confirming the effectiveness of *AudioLDM 2-AC*, which achieves an OVL of 3.88 and a REL of 3.90, surpassing AudioLDM and the previous SoTA TANGO by a significant margin. The difference between *AudioLDM 2-AC* and the GroundTruth, which are real audios from the AudioCaps dataset [27], is merely 0.16 and 0.18 for OVL and REL, respectively, demonstrating the strong performance of our proposed system. The AudioLDM-M we used is not finetuned on the AudioCaps dataset, which may explain its degraded performance compared with the metric score reported in [4].

To investigate the scalability of *AudioLDM 2* in terms of model size and dataset scale, we further trained *AudioLDM 2* on a much larger dataset containing 29,510 hours of data using two different model sizes. The results are shown in Table II. The FAD score generally shows improvement after scaling up the model size, while the KL divergence and CLAP scores do not exhibit clear improvements, indicating that scaling the model size might be more beneficial for enhancing audio quality than audio-text relations. Despite the

⁷<https://huggingface.co/spaces/haoheliu/audioldm-text-to-audio-generation>

⁸<https://github.com/declare-lab/tango>

TABLE II

PERFORMANCE COMPARISON ON THE AUDIOCAPS EVALUATION SET ON DIFFERENT TRAINING DATA SCALES. THE MODELS WITH FULL-SCALE TRAINING DATA ARE NOT FINETUNED ON AUDIOCAPS.

Model	FAD↓	KL↓	CLAP (%)↑	OVL ↑	REL ↑
<i>AudioLDM 2-AC</i>	1.67	1.01	24.9	3.88	3.90
<i>AudioLDM 2-Full</i>	2.13	1.42	19.4	3.76	3.81
<i>AudioLDM 2-Full</i>	1.78	1.60	19.1	3.83	3.77
<i>AudioLDM 2-AC-Large</i>	1.42	0.98	24.3	3.89	3.87
<i>AudioLDM 2-Full-Large</i>	1.86	1.64	18.2	3.79	3.80

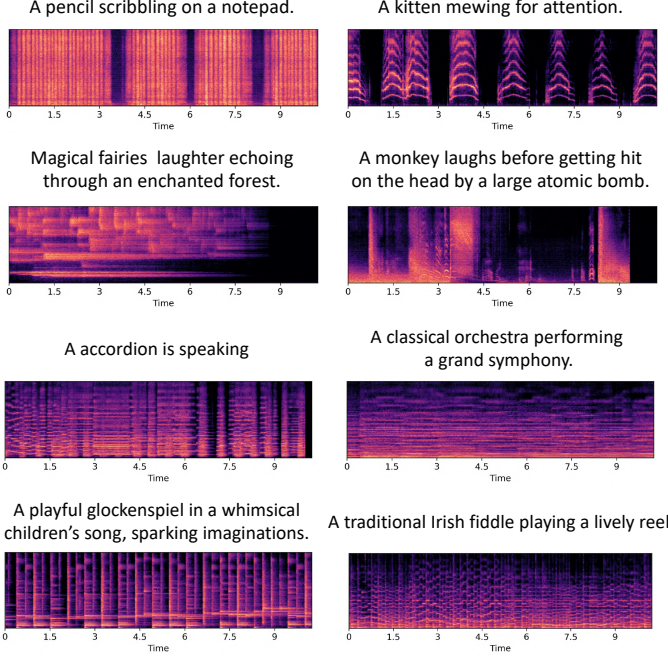


Fig. 6. Examples for text-to-audio generation.

significant increase in training data, we did not observe significant improvements in the objective evaluation metrics. On the contrary, all three metrics showed degraded performance after training on more data. This is potentially because our test set has a limited distribution, while the large-scale training data covers a much wider distribution. The mismatch between the training and test data distributions results in poorer objective scores.

Nevertheless, when compared with the AudioLDM-M (FAD 4.53) in Table I, which is also a large-scale pre-trained text-to-audio model without finetuning on AudioCaps, *AudioLDM 2* with full-scale training data achieves significantly better performance (FAD 1.42 ~2.13), showing a substantial improvement over AudioLDM-M.

B. Text-to-Music Generation

In this section, we compare our proposed model with other text-to-music generation models, including MusicGen [33], MusicLM [8], MeLoDy [13], Mousai [94], AudioLDM [4], and Riffusion [6]. The output of AudioLDM is obtained in the same way as Table I. MusicGen is reproduced using the official Github repository⁹.

⁹<https://github.com/facebookresearch/audiocraft>

TABLE III

PERFORMANCE COMPARISON ON THE MUSICCAPS EVALUATION SET. THE SUPERScript [†] INDICATES RESULTS REPRODUCED USING PUBLICLY AVAILABLE IMPLEMENTATIONS. THE OPEN-SOURCE VERSION OF MUSICGEN-MEDIUM EXCLUDES VOCAL SOUNDS, RESULTING IN SLIGHTLY INFERIOR PERFORMANCE COMPARED TO THE ORIGINAL REPORT [33]. ALL GENERATED AUDIO CLIPS WERE RESAMPLED TO 16kHz PRIOR TO EVALUATION.

Model	FAD↓	KL↓	CLAP (%)↑	OVL↑	REL↑
GroundTruth	-	-	25.3	3.82	4.26
Riffusion	14.80	2.06	19.0	-	-
Mousai	7.50	1.59	-	-	-
MeLoDy	5.41	-	-	-	-
MusicLM	4.00	-	-	-	-
MusicGen-Medium	3.4	1.23	32.0	-	-
MusicGen-Medium [†]	4.89	1.35	29.1	3.37	3.38
AudioLDM-M [†]	3.20	1.29	36.0	3.03	3.25
<i>AudioLDM 2-MSD</i>	4.47	1.32	29.4	3.41	3.30
<i>AudioLDM 2-Full</i>	3.13	1.20	30.1	3.34	3.54

As shown in Table III, our proposed method significantly outperforms these strong baselines. For instance, *AudioLDM 2-Full* outperforms MusicGen by 36%, 11%, and 3.4% on FAD, KL and CLAP scores, respectively. The *AudioLDM 2-MSD* model, which is only trained on music data, does not achieve better performance on objective metrics than the more general *AudioLDM 2-Full*. This result suggests that learning audio generation from a general perspective can benefit the performance in specialised domains as well, demonstrating the advantages of our proposed general framework. The general model *AudioLDM 2-Full* achieves a significantly higher 3.54 REL score than the other systems, indicating better textual understanding ability. The AudioLDM-M model achieves a significantly higher CLAP score than the remaining systems, which suggests that AudioLDM may benefit from being directly conditioned by the same CLAP model during training. The high performance of AudioLDM may also stem from the diversity of audio training data, which also includes music and sound effects, which further supports the benefits of training a general-purpose model. However, the subjective evaluation in Table III indicates that the subjective performance of AudioLDM-M is not as good as suggested by the objective metrics.

C. Text-to-Speech Generation

We compare our proposed model with the widely-adopted FastSpeech2¹⁰ model on the LJSpeech test set. To study the upper bound of our system, we add a setting called GT-AudioMAE that utilizes the ground truth LOA Y to the function \mathcal{G} for audio generations. Our proposed *AudioLDM 2-LJS* is trained on the LJSpeech training split. To further explore the potential of our system, we pre-train the GPT-2 model in function \mathcal{M} on the GigaSpeech dataset before finetuning on LJSpeech. This version is denoted as *AudioLDM 2-LJS-Pretrained*.

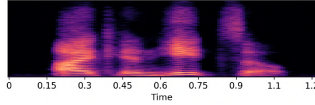
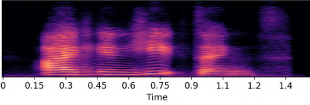
As shown in Table IV, with the pre-trained GPT-2 model, *AudioLDM 2-LJS-Pretrained* achieves a MOS of 4.00, significantly outperforming FastSpeech2. Our subjective evaluation

¹⁰<https://huggingface.co/facebook/fast-speech2-en-ljspeech>

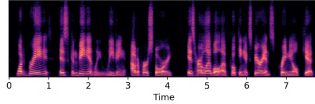
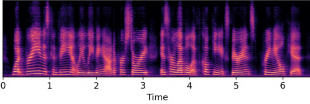
TABLE IV
TEXT-TO-SPEECH PERFORMANCE EVALUATED ON THE LJSPEECH TEST SET.

Model	Mean Opinion Score \uparrow
GroundTruth	4.63 ± 0.08
GT-AudioMAE	4.14 ± 0.13
FastSpeech2	3.78 ± 0.15
<i>AudioLDM 2-LJS</i>	3.65 ± 0.21
<i>AudioLDM 2-LJS-Pretrained</i>	4.00 ± 0.13

Text: I can heat things up.



Text: What green is conveniently leaving out of her story is her level of cooking experience pre-meal kit.



Speaker prompt: A young girl is speaking

Speaker prompt: A young male reporter is speaking

Fig. 7. Examples of speaker-prompted text-to-speech generation. We use speaker prompts to describe the characteristics of the speaker and provide the model with the text transcription.

shows *AudioLDM 2-LJS-Pretrained* exhibits greater fluctuations in emotion, punctuation, and tone. This demonstrates the benefits of pretraining on diverse datasets like GigaSpeech before finetuning on smaller corpora. Without pretraining, our proposed model still achieves a competitive MOS (Mean Opinion Score) of 3.65, which is comparable with the 3.78 MOS of our baseline FastSpeech2.

D. Audio In-context Learning

In-context learning refers to the ability of the auto-regressive model, such as the one used in *AudioLDM 2*, to incorporate and utilize contextual information during the generation process [95]. As demonstrated in Figure 8, in text-to-speech generation, when a prompt speech is provided as context, *AudioLDM 2* can synthesize the remaining speech signal with a consistent speaker style. For text-to-audio and text-to-music tasks, the text typically only has a loose constraint in the generation process. The constraint refers to the flexible influence of the text during generation, allowing for creative interpretation. By comparison, as shown in Figure 8, the incorporation of audio context enhances controllability and aligns the results with both the text description and the audio content.

E. Ablation Studies

In order to validate our design choices of *AudioLDM 2*, we conducted a series of ablation studies on the text-to-audio generation task on the AudioCaps dataset. The results are shown in Table V. When the joint finetuning process between the GPT-2 model and the latent diffusion model was disabled (a), thereby only optimizing them separately, all three evaluation

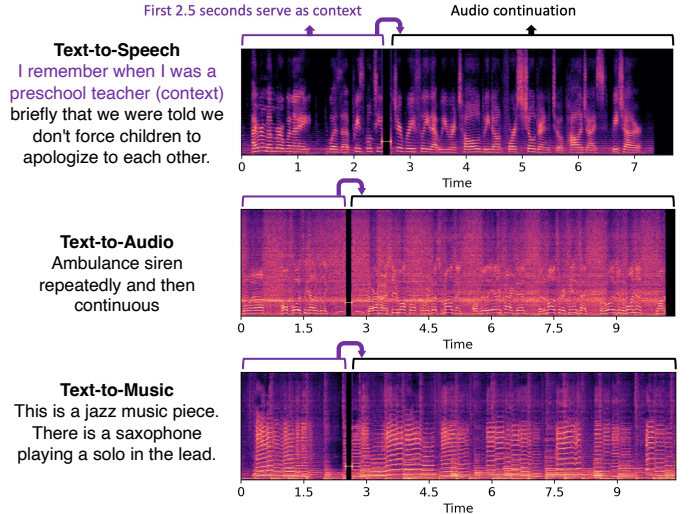


Fig. 8. In-context learning ability of *AudioLDM 2*. The left column shows the ground truth audio, where the leading 2.5 seconds are used as context for audio generation. The continuation of the audio context is shown in the right column. We manually insert a 0.15 seconds beep sound before the continuation for better demonstration.

TABLE V
ABLATION STUDIES ON THE AUDIOCAPS DATASET.

Setting	FAD \downarrow	KL \downarrow	CLAP (%) \uparrow
<i>AudioLDM 2</i>	1.67	1.01	24.9
a. w/o Joint finetuning	2.24	1.07	23.4
b. w/o CLAP embedding (GPT)	2.48	1.07	24.5
c. w/o FLAN-T5 embedding (GPT)	2.73	1.05	25.0
d. w/o FLAN-T5 crossattn (T-UNet)	1.38	1.30	21.1

metrics exhibited a marked deterioration, suggesting joint finetuning is helpful for the GPT-2 model to better cooperate with the LDM model. The GPT-2 model accepts inputs from both the CLAP and FLAN-T5 modules for text-to-audio generation. The removal of either module resulted in a degradation of the evaluation metrics (b-c). However, when only the CLAP module was used as an input (c), the CLAP score was improved. This improvement is likely due to the conditioning directly matching the evaluation metric. The removal of the cross-attention mechanism in the T-UNet model (d), which accepts the FLAN-T5 embeddings, led to a significant degradation in both the KL divergence and CLAP scores. However, it improved the FAD score, from 1.67 to 1.38. These results indicate that while AudioMAE conditioning alone can achieve better FAD, the use of FLAN-T5 conditioning provides additional language semantic information that assists the learning of the audio and text relationships. A similar effect is observed in Table II after removing the FLAN-T5 embeddings.

VI. CONCLUSION AND FUTURE WORKS

In this paper, we have presented *AudioLDM 2* for audio generation, achieving state-of-the-art or comparative performance on text-to-audio, text-to-music, and text-to-speech generation tasks. As a universal audio representation, the language of audio (LOA) we proposed enables self-supervised pre-training of the latent diffusion model, providing a robust foundation for the audio generation task. We further demonstrate the

versatility of our proposed method by performing audio in-context learning. *AudioLDM 2* opens doors for future works on audio generation from a unified perspective. Future work will focus on enabling the multi-task learning of the GPT-2 model to generate audio, music, and speech simultaneously with a single model.

ACKNOWLEDGMENTS

This research was partly supported by the British Broadcasting Corporation Research and Development (BBC R&D), Engineering and Physical Sciences Research Council (EPSRC) Grant EP/T019751/1 “AI for Sound”, and a PhD scholarship from the Centre for Vision, Speech and Signal Processing (CVSSP), Faculty of Engineering and Physical Science (FEPS), University of Surrey. For the purpose of open access, the authors have applied a Creative Commons Attribution (CC BY) license to any Author Accepted Manuscript version arising.

REFERENCES

- [1] Y. Cao, S. Li, Y. Liu, Z. Yan, Y. Dai, P. S. Yu, and L. Sun, “A comprehensive survey of AI-generated content: A history of generative AI from GAN to ChatGPT,” *arXiv preprint:2303.04226*, 2023.
- [2] X. Tan, J. Chen, H. Liu, J. Cong, C. Zhang, Y. Liu, X. Wang, Y. Leng, Y. Yi, L. He *et al.*, “NaturalSpeech: End-to-end text to speech synthesis with human-level quality,” *arXiv preprint:2205.04421*, 2022.
- [3] F. Kreuk, G. Synnaeve, A. Polyak, U. Singer, A. Défossez, J. Copet, D. Parikh, Y. Taigman, and Y. Adi, “AudioGen: Textually guided audio generation,” *International Conference on Learning Representations*, 2022.
- [4] H. Liu, Z. Chen, Y. Yuan, X. Mei, X. Liu, D. Mandic, W. Wang, and M. D. Plumbley, “AudioLDM: Text-to-audio generation with latent diffusion models,” *International Conference on Machine Learning*, 2023.
- [5] D. Zhang, S. Li, X. Zhang, J. Zhan, P. Wang, Y. Zhou, and X. Qiu, “SpeechGPT: Empowering large language models with intrinsic cross-modal conversational abilities,” *arXiv preprint:2305.11000*, 2023.
- [6] S. Forsgren and H. Martiros, “Riffusion: Stable diffusion for real-time music generation, 2022,” URL <https://riffusion.com/about>, vol. 6, 2022.
- [7] X. Liu, Z. Zhu, H. Liu, Y. Yuan, M. Cui, Q. Huang, J. Liang, Y. Cao, Q. Kong, M. D. Plumbley *et al.*, “WavJourney: Compositional audio creation with large language models,” *arXiv preprint:2307.14335*, 2023.
- [8] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzett, A. Caillon, Q. Huang, A. Jansen, A. Roberts, M. Tagliasacchi *et al.*, “MusicLM: Generating music from text,” *arXiv preprint:2301.11325*, 2023.
- [9] R. Bresin, A. de Witt, S. Papetti, M. Civolani, and F. Fontana, “Expressive sonification of footsteps sounds,” *Proceedings of Interactive Sonification Workshop*, 2010.
- [10] J. Engel, L. Hantrakul, C. Gu, and A. Roberts, “DDSP: Differentiable digital signal processing,” *International Conference on Learning Representations*, 2020.
- [11] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, “FastSpeech 2: Fast and high-quality end-to-end text to speech,” in *International Conference on Learning Representations*, 2021.
- [12] D. Herremans and E. Chew, “Morpheus: Automatic music generation with recurrent pattern constraints and tension profiles,” in *Proceedings of IEEE TENCON*, 2016, pp. 282–285.
- [13] M. W. Lam, Q. Tian, T. Li, Z. Yin, S. Feng, M. Tu, Y. Ji, R. Xia, M. Ma, X. Song *et al.*, “Efficient neural music generation,” *arXiv preprint:2305.15719*, 2023.
- [14] D. Yang, J. Yu, H. Wang, W. Wang, C. Weng, Y. Zou, and D. Yu, “DiffSound: Discrete diffusion model for text-to-sound generation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 1720–1733, 2023.
- [15] R. Huang, J. Huang, D. Yang, Y. Ren, L. Liu, M. Li, Z. Ye, J. Liu, X. Yin, and Z. Zhao, “Make-An-Audio: Text-to-audio generation with prompt-enhanced diffusion models,” *International Conference on Machine Learning*, 2023.
- [16] H. Liu, Q. Kong, Q. Tian, Y. Zhao, D. Wang, C. Huang, and Y. Wang, “VoiceFixer: Toward general speech restoration with neural vocoder,” *arXiv preprint:2109.13731*, 2021.
- [17] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, “Data2Vec: A general framework for self-supervised learning in speech, vision and language,” in *International Conference on Machine Learning*, 2022, pp. 1298–1312.
- [18] R. Girdhar, A. El-Nouby, Z. Liu, M. Singh, K. V. Alwala, A. Joulin, and I. Misra, “ImageBind: One embedding space to bind them all,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 15 180–15 190.
- [19] Q. Kong, K. Chen, H. Liu, X. Du, T. Berg-Kirkpatrick, S. Dubnov, and M. D. Plumbley, “Universal source separation with weakly labelled data,” *arXiv preprint:2305.07447*, 2023.
- [20] Y. Okamoto, K. Imoto, S. Takamichi, R. Yamanishi, T. Fukumori, Y. Yamashita *et al.*, “Onoma-to-wave: Environmental sound synthesis from onomatopoeic words,” *Transactions on Signal and Information Processing*, vol. 11, no. 1, 2022.
- [21] H. Xu, J. Li, A. Baevski, M. Auli, W. Galuba, F. Metze, C. Feichtenhofer *et al.*, “Masked autoencoders that listen,” *Advances in Neural Information Processing Systems*, 2022.
- [22] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” 2019.
- [23] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.
- [24] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, “A survey of large language models,” *arXiv preprint:2303.18223*, 2023.
- [25] N. Zeghidour, A. Luebs, A. Omran, J. Skoglund, and M. Tagliasacchi, “SoundStream: An end-to-end neural audio codec,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 495–507, 2021.
- [26] Z. Borsos, R. Marinier, D. Vincent, E. Kharitonov, O. Pietquin, M. Sharifi, D. Roblek, O. Teboul, D. Grangier, M. Tagliasacchi *et al.*, “AudioLM: A language modeling approach to audio generation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 42, pp. 2523–2544, 2023.
- [27] C. D. Kim, B. Kim, H. Lee, and G. 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*, 2019, pp. 119–132.
- [28] R. Sheffer and Y. Adi, “I hear your true colors: Image guided audio generation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
- [29] V. Iashin and E. Rahtu, “Taming visually guided sound generation,” in *British Machine Vision Conference*, 2021.
- [30] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, “Grad-TTS: A diffusion probabilistic model for text-to-speech,” in *International Conference on Machine Learning*, 2021, pp. 8599–8608.
- [31] J. Kim, S. Kim, J. Kong, and S. Yoon, “Glow-TTS: A generative flow for text-to-speech via monotonic alignment search,” *Advances in Neural Information Processing Systems*, pp. 8067–8077, 2020.
- [32] Q. Huang, D. S. Park, T. Wang, T. I. Denk, A. Ly, N. Chen, Z. Zhang, Z. Zhang, J. Yu, C. Frank *et al.*, “Noise2Music: Text-conditioned music generation with diffusion models,” *arXiv preprint:2302.03917*, 2023.
- [33] J. Copet, F. Kreuk, I. Gat, T. Remez, D. Kant, G. Synnaeve, Y. Adi, and A. Défossez, “Simple and controllable music generation,” *arXiv preprint:2306.05284*, 2023.
- [34] Y.-A. Chung, Y. Zhang, W. Han, C.-C. Chiu, J. Qin, R. Pang, and Y. Wu, “W2V-Bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training,” in *IEEE Automatic Speech Recognition and Understanding Workshop*. IEEE, 2021, pp. 244–250.
- [35] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 6840–6851.
- [36] Y. Song, J. Sohl-Dickstein, D. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” in *International Conference on Learning Representations*, 2021.
- [37] P. Dhariwal and A. Nichol, “Diffusion models beat GANs on image synthesis,” in *Advances in Neural Information Processing Systems*, 2021.
- [38] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, “Hierarchical text-conditional image generation with CLIP latents,” *arXiv preprint:2204.06125*, 2022.
- [39] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi, “Photorealistic text-to-

- image diffusion models with deep language understanding,” *arXiv preprint:2205.11487*, 2022.
- [40] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, “Image super-resolution via iterative refinement,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4713–4726, 2022.
 - [41] N. Chen, Y. Zhang, H. Zen, R. Weiss, M. Norouzi, and W. Chan, “WaveGrad: Estimating gradients for waveform generation,” in *International Conference on Learning Representations*, 2021.
 - [42] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, “DiffWave: A versatile diffusion model for audio synthesis,” in *International Conference on Learning Representations*, 2021.
 - [43] Y. Leng, Z. Chen, J. Guo, H. Liu, J. Chen, X. Tan, D. Mandic, L. He, X.-Y. Li, T. Qin *et al.*, “BinauralGrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis,” *Advances in Neural Information Processing Systems*, 2022.
 - [44] U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni *et al.*, “Make-a-video: Text-to-video generation without text-video data,” in *International Conference on Learning Representations*, 2022.
 - [45] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet, and T. Salimans, “Imagen video: High definition video generation with diffusion models,” *arXiv preprint:2210.02303*, 2022.
 - [46] Z. Chen, Y. Wu, Y. Leng, J. Chen, H. Liu, X. Tan, Y. Cui, K. Wang, L. He, S. Zhao, J. Bian, and D. Mandic, “ResGrad: Residual denoising diffusion probabilistic models for text to speech,” *arXiv preprint:2212.14518*, 2022.
 - [47] M. Lam, J. Wang, R. Huang, D. Su, and D. Yu, “Bilateral denoising diffusion models,” in *International Conference on Learning Representations*, 2022.
 - [48] S. Lee, H. Kim, C. Shin, X. Tan, C. Liu, Q. Meng, T. Qin, W. Chen, S. Yoon, and T. Liu, “Priorgrad: Improving conditional denoising diffusion models with data-driven adaptive prior,” in *International Conference on Learning Representations*, 2022.
 - [49] Z. Chen, X. Tan, K. Wang, S. Pan, D. Mandic, L. He, and S. Zhao, “Infergrad: Improving diffusion models for vocoder by considering inference in training,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2022.
 - [50] D. Ghosal, N. Majumder, A. Mehri, and S. Poria, “Text-to-audio generation using instruction-tuned LLM and latent diffusion model,” *arXiv preprint:2304.13731*, 2023.
 - [51] X. Liu, H. Liu, Q. Kong, X. Mei, M. D. Plumbley, and W. Wang, “Simple pooling front-ends for efficient audio classification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
 - [52] H. Liu, X. Liu, Q. Kong, W. Wang, and M. D. Plumbley, “Learning the spectrogram temporal resolution for audio classification,” *arXiv preprint:2210.01719*, 2022.
 - [53] R. Sheffer and Y. Adi, “I hear your true colors: Image guided audio generation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
 - [54] T. I. Denk, Y. Takagi, T. Matsuyama, A. Agostinelli, T. Nakai, C. Frank, and S. Nishimoto, “Brain2Music: Reconstructing music from human brain activity,” *arXiv preprint:2307.11078*, 2023.
 - [55] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, “WaveNet: A generative model for raw audio,” in *ISCA Speech Synthesis Workshop*, 2016, pp. 125–125.
 - [56] X. Tan, T. Qin, J. Bian, T.-Y. Liu, and Y. Bengio, “Regeneration learning: A learning paradigm for data generation,” *arXiv preprint:2301.08846*, 2023.
 - [57] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2020.
 - [58] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, “AudioSet: An ontology and human-labeled dataset for audio events,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 776–780.
 - [59] Y. Li, R. Yuan, G. Zhang, Y. Ma, X. Chen, H. Yin, C. Lin, A. Ragni, E. Benetos, N. Gyenge *et al.*, “MERT: Acoustic music understanding model with large-scale self-supervised training,” *arXiv preprint:2306.00107*, 2023.
 - [60] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
 - [61] S. Schneider, A. Baevski, R. Collobert, and M. Auli, “Wav2Vec: Unsupervised pre-training for speech recognition,” *INTERSPEECH*, pp. 3465–3469, 2019.
 - [62] D. Niizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “Byol for audio: Self-supervised learning for general-purpose audio representation,” in *International Joint Conference on Neural Networks*, 2021.
 - [63] J. Kong, J. Kim, and J. Bae, “HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17022–17033, 2020.
 - [64] K. J. Piczak, “ESC: Dataset for environmental sound classification,” in *Proceedings of the ACM International Conference on Multimedia*, 2015, pp. 1015–1018.
 - [65] L. Van der Maaten and G. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, 2008.
 - [66] Y. Qu, P. Liu, W. Song, L. Liu, and M. Cheng, “A text generation and prediction system: Pre-training on new corpora using BERT and GPT-2,” in *IEEE International Conference on Electronics Information and Emergency Communication*, 2020, pp. 323–326.
 - [67] T. Klein and M. Nabi, “Learning to answer by learning to ask: Getting the best of GPT-2 and BERT worlds,” *arXiv preprint:1911.02365*, 2019.
 - [68] A. M. Lamb, A. G. ALIAS PARTH GOYAL, Y. Zhang, S. Zhang, A. C. Courville, and Y. Bengio, “Professor forcing: A new algorithm for training recurrent networks,” *Advances in Neural Information Processing Systems*, vol. 29, 2016.
 - [69] S. Masoudnia and R. Ebrahimpour, “Mixture of experts: A literature survey,” *Artificial Intelligence Review*, vol. 42, pp. 275–293, 2014.
 - [70] Y. Wu, K. Chen, T. Zhang, Y. Hui, T. Berg-Kirkpatrick, and S. Dubnov, “Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
 - [71] H.-H. Wu, O. Nieto, J. P. Bello, and J. Salomon, “Audio-text models do not yet leverage natural language,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
 - [72] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, “Scaling instruction-finetuned language models,” *arXiv preprint:2210.11416*, 2022.
 - [73] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 5485–5551, 2020.
 - [74] X. Tan, *Neural Text-to-Speech Synthesis*, ser. Artificial Intelligence: Foundations, Theory, and Algorithms. Springer Singapore, 2023.
 - [75] D. P. Kingma and M. Welling, “Auto-encoding variational Bayes,” *arXiv preprint:1312.6114*, 2013.
 - [76] Q. Kong, Y. Cao, H. Liu, K. Choi, and Y. Wang, “Decoupling magnitude and phase estimation with deep ResUNet for music source separation,” *International Society for Music Information Retrieval Conference*, pp. 342–349, 2021.
 - [77] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
 - [78] J. Ho and T. Salimans, “Classifier-free diffusion guidance,” in *NeurIPS Workshop on Deep Generative Models and Downstream Applications*, 2021.
 - [79] A. Q. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, “GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models,” in *International Conference on Machine Learning*, 2022, pp. 16784–16804.
 - [80] X. Mei, C. Meng, H. Liu, Q. Kong, T. Ko, C. Zhao, M. D. Plumbley, Y. Zou, and W. Wang, “WavCaps: A ChatGPT-assisted weakly-labelled audio captioning dataset for audio-language multimodal research,” *arXiv preprint:2303.17395*, 2023.
 - [81] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, “VGGSound: A large-scale audio-visual dataset,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 721–725.
 - [82] M. Defferrard, K. Benzi, P. Vandergheynst, and X. Bresson, “FMA: A dataset for music analysis,” in *International Society for Music Information Retrieval Conference*, 2017.
 - [83] T. Bertin-Mahieux, D. P. Ellis, B. Whitman, and P. Lamere, “The million song dataset,” *International Society for Music Information Retrieval Conference*, pp. 591–596, 2011.
 - [84] K. Ito and L. Johnson, “The LJSpeech dataset,” <https://keithito.com/LJ-Speech-Dataset/>, 2017.

- [85] G. Chen, S. Chai, G. Wang, J. Du, W. Q. Zhang, C. Weng, D. Su, D. Povey, J. Trmal, J. Zhang *et al.*, “GigaSpeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio,” in *INTER-SPEECH*, 2021, pp. 4376–4380.
- [86] S. Doh, M. Won, K. Choi, and J. Nam, “Toward universal text-to-music retrieval,” *arXiv preprint:2211.14558*, 2022.
- [87] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, “CNN architectures for large-scale audio classification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 131–135.
- [88] K. Koutini, J. Schlüter, H. Eghbal-Zadeh, and G. Widmer, “Efficient training of audio transformers with patchout,” *INTERSPEECH*, pp. 2753–2757, 2021.
- [89] R. Likert, “A technique for the measurement of attitudes,” *Archives of Psychology*, 1932.
- [90] Z. Chen, N. Kanda, J. Wu, Y. Wu, X. Wang, T. Yoshioka, J. Li, S. Sivasankaran, and S. E. Eskimez, “Speech separation with large-scale self-supervised learning,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023.
- [91] J. Song, C. Meng, and S. Ermon, “Denoising diffusion implicit models,” in *International Conference on Learning Representations*, 2020.
- [92] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *International Conference on Learning Representations*, 2019.
- [93] J. Huang, Y. Ren, R. Huang, D. Yang, Z. Ye, C. Zhang, J. Liu, X. Yin, Z. Ma, and Z. Zhao, “Make-An-Audio 2: Temporal-enhanced text-to-audio generation,” *arXiv preprint:2305.18474*, 2023.
- [94] F. Schneider, Z. Jin, and B. Schölkopf, “Mousai: Text-to-music generation with long-context latent diffusion,” *arXiv preprint:2301.11757*, 2023.
- [95] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, “A survey for in-context learning,” *arXiv preprint:2301.00234*, 2023.