# LLMs Meet Multimodal Generation and Editing: A Survey

Yingqing He[†*], Zhaoyang Liu[†*], Jingye Chen[*], Zeyue Tian[*], Hongyu Liu[*], Xiaowei Chi[*],
Runtao Liu[*], Ruibin Yuan[*], Yazhou Xing[*], Wenhai Wang, Jifeng Dai, Yong Zhang,
Wei Xue, Qifeng Liu, Yike Guo, Qifeng Chen

**Abstract**—With the recent advancement in large language models (LLMs), there is a growing interest in combining LLMs with multimodal learning. Previous surveys of multimodal large language models (MLLMs) mainly focus on multimodal *understanding*. This survey elaborates on *multimodal generation and editing* across various domains, comprising image, video, 3D, and audio. Specifically, we summarize the notable advancements with milestone works in these fields and categorize these studies into LLM-based and CLIP/T5-based methods. Then, we summarize the various roles of LLMs in multimodal generation and exhaustively investigate the critical technical components behind these methods and the multimodal datasets utilized in these studies. Additionally, we dig into tool-augmented multimodal agents that can leverage existing generative models for human-computer interaction. Lastly, we discuss the advancements in the generative AI safety field, investigate emerging applications, and discuss future prospects. Our work provides a systematic and insightful overview of multimodal generation and processing, which is expected to advance the development of Artificial Intelligence for Generative Content (AIGC) and world models. A curated list of all related papers can be found at https://github.com/YingqingHe/Awesome-LLMs-meet-Multimodal-Generation.

**Index Terms**—LLMs, MLLMs, Multimodal Generation, Text-to-Image, Text-to-Video, Text-to-3D, Text-to-Audio, Multimodal Agents, AI Safety, Diffusion Models, Transformers, Generative AI, AIGC.

✦

## 1 INTRODUCTION

The interactions between human beings and the physical world involve information from multiple modalities, such as language, vision, and audio. Therefore, realizing a world simulator also requires the model to be able to perceive and respond with multimodal information in a flexible manner. Recently, OpenAI introduced a foundation text-to-video generation model termed Sora [1] that is capable of generating highly realistic videos as world simulators. It makes great progress in simulating or generating real-world scenes but is unable to generate other modalities, such as text, 3D, and audio. Also, it lacks the ability to perceive other modalities such as image, video, 3D, and audio, making it an incomprehension world simulator.

In the past few years, researchers focused on the generation of each single modality and have achieved great progress: In terms of text generation, we have witnessed a qualitative leap in the performance of natural language processing tasks: From BERT [2], GPT1 [3], GPT2 [4], GPT3 [5], GPT4 [6] to ChatGPT [7], LLaMA [8], [9], the number of model parameters and training samples has grown rapidly, resulting in the continual growth of modal abilities and product deployment. In the visual generation field, with the

---

- † *Project leaders; * Co-first authors.*
- *Yingqing He, Zhaoyang Liu, Jingye Chen, Zeyue Tian, Hongyu Liu, Xiaowei Chi, Runtao Liu, Ruibin Yuan, Yazhou Xing, Wei Xue, Qifeng Liu, Yike Guo, and Qifeng Chen are with The Hong Kong University of Science and Technology, Hong Kong SAR.*
- *Wenhai Wang is with The Chinese University of Hong Kong, Hong Kong SAR.*
- *Jifeng Dai is with Tsinghua University, China.*
- *Yong Zhang is with the Tencent AI Lab, China.*

rapid progress of diffusion models and large-scale image-text datasets, text-to-image (T2I) generation has achieved remarkable achievement and can synthesize high-quality images based on various user-provided text prompts, such as SDXL [10] and PIXART-$\alpha$ [11]. Subsequently, significant advancements have been made in the field of text-to-video generation through the utilization of video diffusion models [12] and large-scale video-language datasets [13]. Notably, several milestone works have emerged, such as [14]–[23] and Sora [1]. For the 3D generation, with the emergence of the CLIP [24] model, some methods [25]–[27] try to align the text information to the rendered images from 3D representations, *i.e.*, mesh, point cloud, NeRF [28] and gaussian splatting [29]). These approaches have led to significant developments in text-to-3D generation. Additionally, the integration of Stable Diffusion (SD) [30] with text-to-image rendering has enabled a series of works in text-to-3D generation [31]–[45]. The powerful text-to-image model helps the 3D generation achieve higher performance and better visual results. In the area of text-to-audio generation, a series of representative works tackle different audio domains such as [46]–[48] for text-to-audio, [49]–[51] for text-to-music, and [52]–[57] for text-to-speech, and they have achieved significant performance in generating high-quality natural sounds, music, and human-level speech.

With the notable progress and significant performance improvements of Large Language Models (LLMs), other non-text modalities have started harnessing the power of LLMs to either enhance their generation quality or integrate multiple modalities in a unified system to achieve more powerful functionalities. In the context of image generation, the integration of LLMs can be divided into two categories.
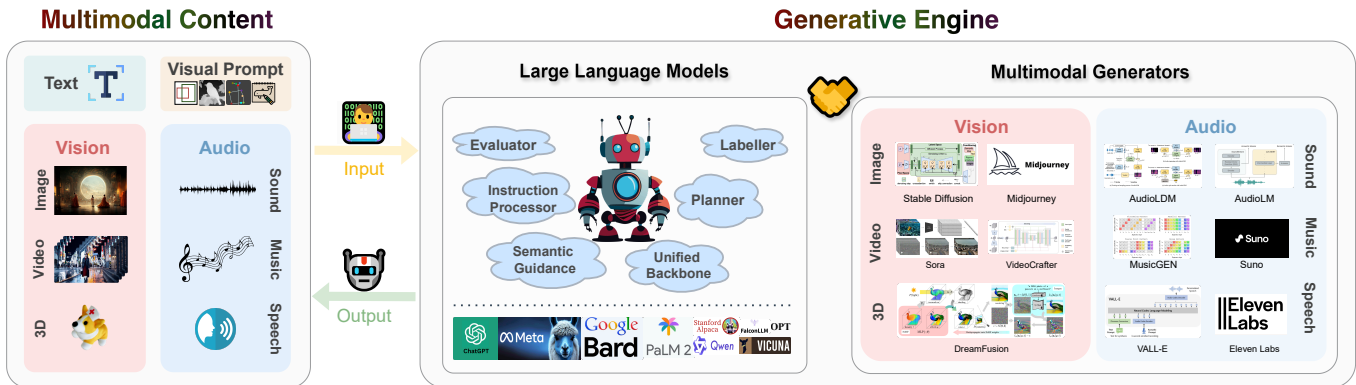
Fig. 1: Our main goal is to investigate the roles of LLMs in the task of language-guided multimodal generation. The modalities we focused on consist of image, video, 3D, and audio (including sound, music, and speech).

The first category involves encoding visual information into discrete token indices, trying to unify visual understanding and generation [58]–[63]. Specifically, visual information is encoded into token representations, and LLMs directly comprehend and generate visual tokens, enabling simultaneous visual understanding and generation. The second category focuses on leveraging LLMs to enhance the generation quality of existing pretrained T2I models: One type of approach utilizes an LLM as a layout planner to incorporate knowledge of object spatial positions, quantity, and object size, enabling the generation of required bounding boxes [64]–[68]. After obtaining the bounding boxes, the images can be generated through a grounded T2I model such as GLIGEN [69]. Another approach utilizes LLMs to expand input user prompts [70]: By providing highly detailed and comprehensive prompts, LLMs generate images with high quality and richness. With the assistance of LLMs, image generation has achieved higher generation quality, improved prompt following capabilities, dialogic function, and user-friendly interface.

Similar to image domain, in video generation, LLMs serve as the general backbone for unified multimodal joint generation [71], [72], video layout planning [65], [73]–[76] and temporal prompt generation [77]–[81] for temporal dynamics guidance. For 3D generation and editing, LLMs serve as a bridge between users and 3D assets, which improves interaction efficiency [82], [83] and helps users understand the 3D assets [84], [85]. In the context of audio generation and editing, the role of LLMs primarily lies in serving as coordinated backbones for multimodal audio [86]–[98], conditioners for specific tasks [99]–[101], labelers for audio understanding [102]–[104], agents for interactive generation and editing [105]–[110], as well as inspiration for novel approaches [49], [50], [55], [111]–[113]. The growing utilization of LLMs in the audio domain is not only transforming our way of engaging with sound and music but also expanding the boundaries at the crossroads of AGI and audio technologies. Besides, multimodal agents [105], [114]–[118] integrate lots of AIGC tools into the framework as a universal system, which relies on LLMs to invoke tools but endows LLMs with the ability to comprehend and generate content of non-text modalities. Generally, LLMs significantly play an indispensable role in generating various modes of content.

To promote the development of multimodal generation and empower the world simulator, in this work, we provide a comprehensive review of works involving LLMs in the generation of multiple modalities. As shown in Fig. 1, we summarize the roles of LLMs into several key aspects, such as evaluator, labeller, instruction processor, planner, provider of semantic guidance, or as backbone architectures. Additionally, we discuss the development in the generative AI safety topic in Sec. 9, with the emerged applications, and the potential future prospects in Sec. 10 and Sec. 11.

We summarize our contributions as follows:

- We present the first systematic review of LLMs applied to the generation and editing of multiple modalities, including images, videos, 3D, and audio.
- We discuss the evolution of generative techniques through a comparative analysis of pre-LLM and post-LLM eras, offering a clear perspective on the progression and refinement of these approaches.
- We summarize the various roles of LLMs in the generation or editing process for each modality from a technical view.
- We discuss important AI safety issues, investigate emerging applications and explore future directions to boost the development of multimodal generation and world models.

## 1.1 Scope

This survey explores the generation of multiple modalities, including images, videos, 3D models, and audio. The multimodal generation in our survey includes the separate generation of different modalities as well as the joint generation of multiple modalities. We will not dig into pure text generation and processing extensively, as there have already been many surveys specifically focusing on the advancements in that field [119]–[121]. Our primary focus is on how the recent emergence of LLMs in the past few years can assist in the generation of other vision and audio modalities, especially in the open-domain generation. This will aid us in designing better unified generative models for multi-modalities. Note that the tasks and works we discussed are primarily *language-based* generation and editing.

Unconditional generation and other non-text-based editing are not our primary focus since they are either limited to a small domain or lack flexibility and controllability. In detail, we focus on the following tasks:

- **Text-to-image generation and editing**: Image generation aims to create various open-domain image contents, including pictures, photos, or stylized drawings from user-provided textual descriptions. Image editing aims to modify the input image content and can be based on user instructions.
- **Text-to-video generation and editing**, where models generate or modify arbitrary and various dynamic visual contents guided by free-form text descriptions.
- **Text-to-3D generation and editing**, which is a task for generating and editing 3D objects, scenes, or avatars with user-provided textual descriptions.
- **Text-to-audio generation and editing**, where textual descriptions are used to generate audio, including general sounds, music, and speech. Audio editing tasks, such as adding, removing, or inpainting, can all be performed by modifying existing audio content through textual descriptions.
- **Multimodal generative agents**, which enable LLMs to handle data across different modalities by utilizing a variety of specialized multimodal tools.
- **Generative AI Safety**, which focuses on reducing toxic and biased content, protecting copyright, and addressing the creation of fabricated content by multimodal generative models.

### 1.2 Content Overview

We first review related surveys on both single modality generation and LLMs in Sec. 2. We then briefly review the basic techniques of representative generative models, multimodal alignment models, LLMs, and MLLMs in Sec. 3. Next, we review the LLM-based vision generation in different visual modalities, including image in Sec. 4, video in Sec. 5, 3D in Sec. 6, audio modality in Sec. 7, and multimodal agents in Sec. 8, respectively. Lastly, we review the safety aspect of generative AI in Sec. 9, emerging applications in Sec. 10, and potential future directions for the LLMs-based multimodal generation field in Sec. 11.

## 2 RELATED SURVEYS

**Survey on modality-specific generation.** A series of surveys focus on single modality generation, such as [122] for image generation, [123] for video generation, [124] for 3D generation, [125] for audio generation. However, the previous generation paradigm mainly employs pretrained CLIP [24], CLIP-related variant [126], or language encoder T5 [127] to achieve open-domain text-guided generation. With the emergence of LLMs, there is a growing trend of leveraging powerful LLMs to enhance the generation of content for each modality. Our work aims to provide a comprehensive survey on the role of LLMs in the generation of various types of modalities, an aspect that is absent from previous surveys.

**Survey on LLMs and MLLMs.** Numerous surveys have been conducted to explore various aspects of LLMs. For example, [128] offers a comprehensive examination of LLMs-based autonomous agents. Additionally, [129], and [130] look into MLLMs, introducing papers that combine LLMs with other non-text modalities. They review papers on multimodal understanding and generation in a mixed manner, introducing primarily multimodal understanding works and focusing less on multimodal generation. In contrast, our work primarily concentrates on the *generation* aspect, aiming to investigate the performance and functionality improvement that LLMs bring to each modality's generation process, leading to a better AI-generated world with various modalities.

## 3 PRELIMINARIES

In this section, we first review different types of generative models in Sec. 3.1. Then, we illustrate multimodal alignment models in Sec. 3.2. Lastly, We introduce the technical principle of large language models in Sec. 3.3, and explain multimodal large language models in Sec. 3.4
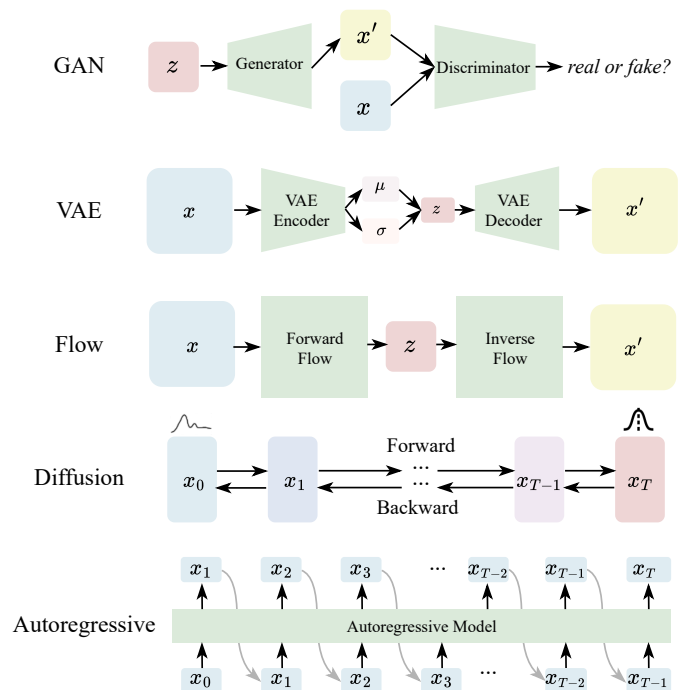


Fig. 2: The illustration of generative models. In this picture, $x$ and $x_0$ indicate the sample from the real data distribution, $x'$ stands for the sample from the model's estimated data distribution, and $z$ means the latent sampled from a prior distribution (typically a Gaussian distribution).

### 3.1 Generative Models

We review the core principles and basic concepts of classic generative models, including generative adversarial Networks (GANs), variational autoencoders (VAEs), flow-based models, diffusion models, and autoregressive models.

The generation process of generative models can be characterized as a transformation from a latent sample $z$

drawn from a prior distribution $p_{\mathbf{z}}(\mathbf{z})$ to a generated sample $\mathbf{x}'$ from a real data distribution $p_{\text{data}}(\mathbf{x})$ that aligns with the target data distribution. Specifically, the latent variable $\mathbf{z}$ is passed through a parametric function, typically implemented as a neural network, which learns to map the prior distribution to the target data distribution. The output $\mathbf{x}'$ of this transformation is then regarded as a synthetic instance that mimics the statistical properties of the original data distribution, which may correspond to various modalities such as images, videos, 3D representations, audio, or text.

### 3.1.1 Generative Adversarial Networks

GAN [131] has achieved promising results in various tasks during years of development. As shown in Fig. 2, GAN comprises two crucial components: a discriminator ($D$) and a generator ($G$). The discriminator is designed to distinguish between real and fake samples. The generator aims to create fake samples that are indistinguishable from the real data and try to fool the discriminator. During training, the $G$ and $D$ are trained simultaneously and play a two-player minimax game. The optimization objective is formulated as follows:

$$\min_G \max_D V(D, G) =$$
$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \tag{1}$$

Where $D(\mathbf{x})$ indicates the probability of real sample $\mathbf{x}$ being a real sample, and $D(G(\mathbf{z}))$ indicates the probability of a generated sample being a fake sample. $\mathbb{E}_x$ is the expectation value over all samples.

### 3.1.2 Variational Auto-Encoder

Variational auto-encoder [132] contains an encoder and a decoder to learn latent representations from the input data, as shown in the second row in Fig. 2.

Encoder is a neural network that maps the input data $\mathbf{x}$ to the distribution of the latent space variable $\mathbf{z}$. Then, the variational posterior distribution $q(\mathbf{z}|\mathbf{x})$ is usually assumed to be a Gaussian distribution $N(\mu, \sigma^2)$. In this case, the encoder gives the $\mu$ and $\sigma^2$.

Decoder maps the latent space variable $\mathbf{z}$ back to the input space $\mathbf{x}'$, yielding the conditional distribution of the generated data $x, q(x|z)$.

The training optimization target of VAE is to maximize the lower bound of the marginal log-likelihood of the data. This target can be achieved through stochastic gradient descent and the reparameterization trick. It is also known as Evidence Lower Bound (ELBO). Specifically, the ELBO can be written in the following form:

$$\mathcal{L}(\theta; \mathbf{X}) =$$
$$- \left[ KL(q(\mathbf{z}|\mathbf{x}; \theta) || p(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \theta)}[\log p(\mathbf{x}|\mathbf{z}; \theta)] \right] \tag{2}$$

where $KL(||)$ denotes the KL divergence, which measures the difference between the posterior distribution inferred by the encoder $q(\mathbf{z}|\mathbf{x})$ and the prior distribution $p(z)$. The second term is the reconstruction error, which

represents the match between the $\mathbf{x}'$ generated from $z$ and the actual data $x$.

**Vector Quantized Variational Auto-Encoder** VQ-VAE [133] is a variant of the VAE that introduces a discrete latent space, which significantly improves the quality of the generated samples compared to the original VAE. In VQ-VAE, the encoder maps the continuous output of the encoder to the nearest point in a predefined discrete codebook and outputs a discrete latent representation $\mathbf{z}$. The codebook is learned jointly with the rest of the model parameters. Using a discrete latent space allows VQ-VAE to capture more global and structured information about the data, leading to better generation quality.

### 3.1.3 Flow-based Model

Flow-based model, also known as normalizing flows, is a class of generative models successfully applied in various tasks, including image synthesis, variational inference, and unsupervised representation learning. The architecture of a flow-based model consists of a sequence of invertible transformations (or flows). Each flow is parameterized by a neural network, which learns to transform the data distribution step by step into the simpler prior distribution. The objective function for training a flow-based model is the negative log-likelihood of the data under the model, which can be computed exactly due to the invertibility of the flows. The function is given by:

$$\mathcal{L}(\theta) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[\log p_{\text{model}}(\mathbf{x}; \theta)] \tag{3}$$

where $p_{\text{model}}(\mathbf{x}; \theta)$ is the probability density function of the model. In practice, the transformations used in flow-based models are chosen to be easily invertible and have easily computable Jacobians, such as affine transformations.

### 3.1.4 Diffusion Model

The diffusion model is proposed in [134], which first gives the prototype of recent diffusion models. However, the foundational structure of the modern diffusion model, which led to a revolution of generation paradigms, was proposed in Denoised Diffusion Probabilistic Models [135]. It is elegant in training and mythical improvement and only introduced a simple regression loss. As shown in Fig. 2, the diffusion model turned the complex work into a series of denoising tasks, primarily consisting of two steps: the injection of prior noise into the data and the denoising prediction.

**Forward Noise Injection** In the forward noise injection process, the model sequentially introduces Gaussian noise $\zeta_{\mathbf{t}}$ into the data in each step $t$ of $T$ steps. The process can be represented as follows:

$$\mathbf{x}_{\mathbf{t+1}} = \sqrt{1 - \alpha_t^2} x_t + \alpha_t \zeta_{\mathbf{t}} \tag{4}$$

Where $\mathbf{x}_{\mathbf{t}}$ is the data at time $t$, $\zeta_{\mathbf{t}}$ is the Gaussian noise at time $t$, and $\alpha_{\mathbf{t}}$ is a noise schedule determining the amount of noise to be added at each step. The noise schedules $\alpha_{\mathbf{t}}$ typically start close to 0 and gradually increase to 1 over the $T$ steps.

The noise at each step is assumed to follow a Markov transition process, which means the noise at time $t$, $\zeta_{\mathbf{t}}$, is independent of the noises at all previous times. This

assumption simplifies the model and makes it tractable to train.

**Reverse Denoising** After the forward noise injection, the model aims to reverse the process by predicting the original data from the noisy version. This is done by learning a denoising function, which is typically parameterized as a deep neural network. The denoising function inputs the noisy data at time $t$ and attempts to predict the noise added at that step. This process is repeated for each time step, going backward from $T$ to 1. The denoising function can be represented as follows:

$$\hat{\zeta}_{\mathbf{t}} = D_\theta(\mathbf{x_t}, t) \tag{5}$$

where $D_\theta$ is the denoising function parameterized by $\theta$, $x_t$ is the noisy data at time $t$, and $\hat{\zeta}_{\mathbf{t}}$ is the predicted noise.

The model's objective during training is to minimize the difference between the predicted noise $\hat{\zeta}_{\mathbf{t}}$ and the actual noise $\zeta_{\mathbf{t}}$ added during the forward process. This can be measured using a simple mean squared error loss.

Training the model to perform this reverse denoising prediction allows it to generate data similar to the original data distribution, making the diffusion model a powerful tool for generative tasks.

### 3.1.5 Autoregressive Model

Autoregressive models are another class of generative models widely used in various tasks, including time series forecasting, speech synthesis, and natural language processing. Their architecture predicts future values based on past values.

The objective function for training an autoregressive model is also the negative log-likelihood of the data under the model, which can be computed exactly due to the sequential nature of the model. This is given by:

$$\mathcal{L}(\theta) = -\mathbb{E}\mathbf{x} \sim p_{\text{data}}(\mathbf{x})[\log p_{\text{model}}(\mathbf{x_{t+1}}|\mathbf{x}_{\leq \mathbf{t}}; \theta)] \tag{6}$$

where $p_{\text{model}}(\mathbf{x_{t+1}}|\mathbf{x}_{\leq \mathbf{t}}; \theta)$ is the conditional probability density function of the model. In practice, the model is trained to maximize the likelihood of the next value in the sequence given the previous values.

### 3.2 Multimodal Alignment Model

CLIP [24] is a ground-breaking image-language alignment model that simultaneously learns an image encoder and a text encoder to produce visual and textual representations in a shared semantic space, trained on a diverse range of internet text and images through contrastive learning [136]. After its large-scale contrastive pretraining, it is capable of tackling various downstream tasks, including fine-grained object recognition, video action recognition, facial emotion recognition, geo-localization, and many others in a zero-shot manner. Thanks to its web-scale pertaining, it can understand plenty of semantics. Thus, it has become one of the most widely used visual and textual encoders in various vision generation and editing works such as DALLE-2 [137] and LDM [30] for text-to-image generation, VideoCrafter [19] for text-to-video and CLIP-Nerf [138] for 3D.

Besides text and vision alignment, CLAP [139] aligns text and audio information. The audio-aligned text embedding

representation is used as the condition of AudioLDM [46] for text-guided audio generation.

CAVP makes a further step towards video-audio alignment, which is trained in Diff-Foley [140] for the video-to-audio generation task. After training CAVP, Diff-Foley further trains a latent diffusion model, which is conditioned on the audio-aligned video representation to output synchronized audio signals.

Different from previous methods for the alignment of paired modalities, ImageBind [141] aligns six different modalities in one shared semantic space. The supporting modalities include text, image, video, audio, depth, and thermal. It has been used in multimodal generation tasks such as Next-GPT [142], Seeing-and-Hearing [143], and also multimodal understanding works such as PandaGPT [144].

### 3.3 Large Language Models

Modern Large Language Models utilize the transformer architecture to generate contextually rich embeddings. These models are trained on a large corpus of text and then fine-tuned for specific tasks. They generate text by predicting the next word in a sequence, given the previous words.

Typical examples include LlaMA [145] and GPT [146]–[148], which are the autoregressive models that use only the left context to make predictions. They are mainly built by transformer decoder. The models would be pre-trained on large, diverse datasets to acquire a strong foundation of language understanding and generation capabilities and then fine-tuned on datasets that provide explicit instructions or guidance on specific tasks, such as question-answering summarization or code generation. Furthermore, tricks like Chain-of-Thought (CoT) [5] fine-tuning and Reinforcement Learning from Human Feedback (RLHF) [149] improve the models' task-specific ability.

### 3.4 Multimodal Large Language Models

Multimodal Large Language Models (MLLMs) are recently emerging models that aim to equip LLMs with the ability to understand or generate other modalities. MLLMs usually incorporate a few key components: extra pre-trained modality-specific encoders for feature extraction and input projectors for multimodal hidden feature alignment with the LLM backbone. For MLLMs with generation ability, they usually contain extra output projectors and corresponding modality generators as the generation ends. A series of works incorporating extra pre-trained encoders for encoding multimodal information to a pre-trained LLM and train modality alignment modules to achieve this [142], [150]–[154]. Other works train the whole multimodal system in an end-to-end manner [155]. In the following sections, we will illustrate a range of recent MLLMs works, especially MLLMs on multimodal generation.

## 4 IMAGE GENERATION AND EDITING

### 4.1 Image Generation

Image generation has long been a fundamental task in the field of computer vision, playing a vital role in various applications such as digital art, entertainment, education, and communication [206]–[208]. In the beginning stage of image

**Single domain generation**
A single generative model possesses the capability to generate images within a singular domain only.

*VAE (Dec 2013), GAN (Jun 2014), Normalizing flow (May 2015), DDPM (Dec 2020) ...*

**Open domain generation**
Generative models exhibit the capability to generate images across arbitrary domains guided by textual descriptions.

*DALL·E (Feb 2021), LDM (Dec 2021), DALL·E 2 (Apr 2022), Imagen (May 2022) ...*

**Interactive / Interleaved generation**
Generative models create images through user interaction powered by 🤖 **Large Language Models**.

*DALL·E 3 (Sept 2023), DreamLLM (Sept 2023), mini DALL·E 3 (Oct 2023), CoDi-2 (Nov 2023) ...*

2013 - 2020　　2021 - 2022　　2023 -

GAN on MNIST　GAN on CIFAR-10

DDPM on LSUN　DDPM on CelebA-HQ

A painting of a squirrel eating a burger (LDM)　An animal half mouse half octopus (LDM)　A slightly conscious neural network (LDM)

A wine glass on top of a dog (Imagen)　A pear cut into seven pieces arranged in a ring (Imagen)　A wine glass on top of a dog (Imagen)

My 5 year-old keeps talking about a "super-dupersunflower hedgehog" -- what does it look like?

Could you design some stickers for it?

Sure, here you are.
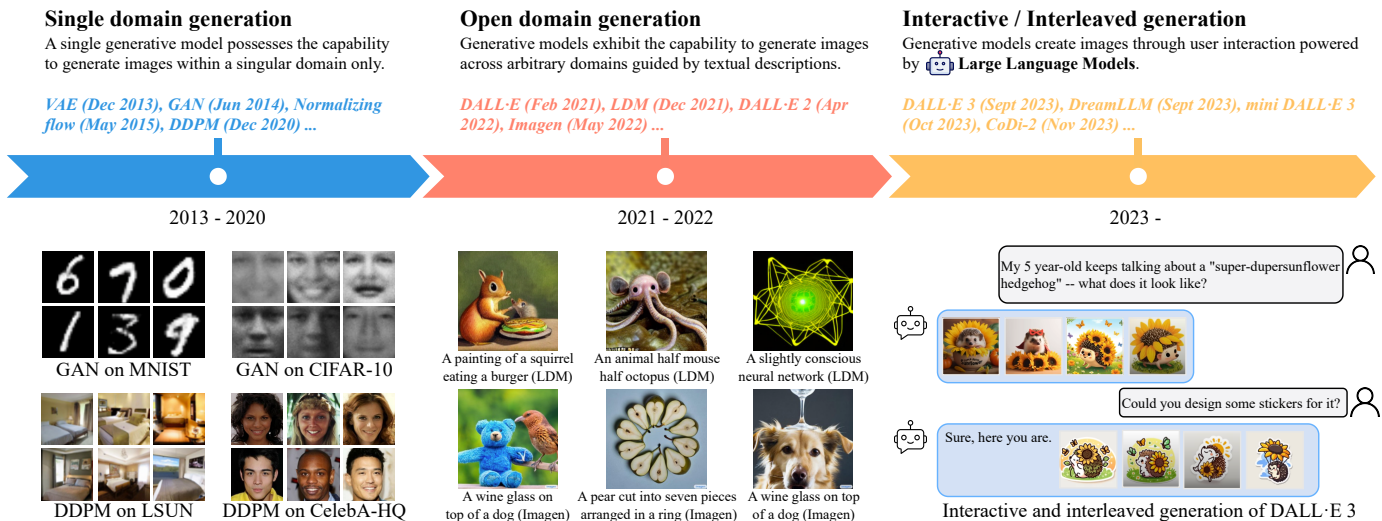
Interactive and interleaved generation of DALL·E 3

Fig. 3: History review on the development trajectory of image generation. Early works on image generation predominantly concentrated on synthesizing images within specific narrow domains, such as human faces or bedrooms [156], [157]. Subsequently, DALL-E [158] and Latent Diffusion Models (LDM) [30] have progressed to generate images through user prompts and support the synthesis of open-domain images. In the recent two years, powered by LLMs, research has trended toward achieving a more intuitive and interactive image generation process, such as iterative generation through conversations [159], [160].

generation, the generated content is limited to a specific category, such as human faces, cats, or buildings. Recent advancements in image generation have been particularly notable due to the introduction of text guidance and open-domain generation. Most recently, the power of LLMs has elevated image generation to a new level, enabling interactive or interleaved generation. In Fig. 3, we summarize the history and development trajectory of image generation in detail. A curated list of recent image generation methods is shown in Table 1. We also list representative datasets for image generation in Table 2. These works allow for the creation of images that closely align with textual prompts, providing a powerful tool for visualizing ideas with creativity.

### 4.1.1　Text-guided Image Generation with CLIP

Previously, the adoption of image-text alignment models, such as CLIP [24], has played a crucial role in the development of text-guided image generation [30], [137], [177], [209]–[211]. The alignment capabilities of CLIP text encoder ensure that the generated image aligns with the given textual prompt, resulting in images that accurately match the intended descriptions, including desired objects, scenes, or attributes. Given the significant strides made by CLIP in generating realistic images, it is natural to question: *will more powerful LLMs further benefit the domain of image generation?*

It is worth mentioning that, the application of LLMs for the image domain has been extensively investigated, especially for image understanding. LLMs can effectively serve as unified processors of visual tokens and language tokens [8], [9], [142], [150], [212]–[225], tool coordinators [105], [116], [117], [226], or analysts of upstream visual model outputs [227]–[229]. Inspired by these works, numerous works further take advantage of LLMs for image generation, and the milestone works are shown in Fig. 3. In the following,

we will introduce the progress achieved in the task of image generation after the emergence of LLMs.

### 4.1.2　Text-guided Image Generation with LLMs

As demonstrated in Fig. 4, MLLMs have emerged as a transformative extension of LLMs, addressing the inherent limitations of LLMs in handling visual content. While LLMs excel in flexible text-based interactions, they are confined to textual inputs and outputs. The introduction of MLLMs stems from the necessity to bridge this gap and enable language models to comprehend and generate images. MLLMs offer a two-fold advantage: firstly, they serve as a unified interface for understanding and generating both textual and visual information, providing users with a seamless integration of language and images. Secondly, MLLMs introduce interactive generation capabilities, allowing users to send commands to modify image content iteratively. This interactive process empowers users with greater control over the image generation process, enhancing the overall user experience and controllability towards user-desired content.

Specifically, CM3Leon [164] is an autoregressive MLLMs designed for the simultaneous generation of textual and image outputs. Operating within a decoder-only, retrieval-augmented, token-based framework, CM3Leon offers a unique approach to multimodal language processing. DreamLLM [159] presents the first MLLM capable of generating free-form interleaved content, supporting multi-round dialogue, and achieving remarkable results in image captioning and video question answering (VQA) without the need for fine-tuning. The entire framework is trained on interleaved multimodal contents in a truly end-to-end manner. SEED-LLaMA [61], similar to DreamLLM, enables LLMs to comprehend multimodal instructions and supports multi-turn in-context image and text generation. Notably,

TABLE 1: Overview of existing methods using LLMs for language-based image generation. According to the role of LLMs in this task, these methods can be divided into four categories: multimodal LLMs for generation, image layout planning, prompt synthesis and refinement, and image quality evaluation. In the "Task" column, the "T" and "I" are the abbreviations of "text" and "image", respectively, while "Any" represents the universal generation supporting text, image, video, and audio modality. "-" indicates that the information is not available in the official paper.

| Method | Venue | Task | LLM | Generative Model | Training Cost |
|---|---|---|---|---|---|
| *Multimodal LLMs with image generation* | | | | | |
| FROMAGe [161] | ICML 2023 | TI→TI | OPT | Retrieval | 1×A6000, 24Hrs |
| GILL [162] | NeurIPS 2023 | TI→TI | OPT | Retrieval/SD | 2×A6000, 48Hrs |
| SPAE [163] | NeurIPS 2023 | Tokenization | PaLM2/GPT-3.5 | CNN | - |
| Emu [155] | ICLR 2024 | TI→TI | LLaMa | SD | - |
| SEED [58] | ICLR 2024 | Tokenization | OPT | SD | 64×V100, 44Hrs |
| CM3Leon [164] | arXiv 2023 | TI→I, I→T | CM3Leon | CM3Leon | 64×A100 |
| NExT-GPT [142] | arXiv 2023 | Any→Any | Vicuna | SD etc. | - |
| DreamLLM [159] | ICLR 2024 | TI→TI | Vicuna | SD | 128×A800, 17.5Hrs |
| MiniGPT-5 [165] | arXiv 2023 | TI→TI | Vicuna | SD | 4×A6000 |
| OpenLEAF [166] | arXiv 2023 | T→TI | GPT-4 | SDXL | Training-free |
| Mini-DALLE3 [59] | arXiv 2023 | TI→TI | GPT-3.5/GPT-4 etc. | SD-XL/DALLE-3 etc. | Training-free |
| EasyGen [63] | arXiv 2023 | I→T,T→I | FlanT5XL/Vicuna | BiDiffuser | 120×A100 Hrs |
| TEAL [167] | arXiv 2023 | Any→Any | LLaMa-Adapter | VQGAN | 8×A100 |
| LLMGA [168] | arXiv 2024 | T→I | LLaVA-1.5 | SD | - |
| ChatIllusion [169] | arXiv 2023 | TI→TI | LLaMa-AdapterV2 | SDXL | 4×A6000, 80Hrs |
| CoDi-2 [60] | CVPR 2024 | Any→Any | Llama 2 | SD | - |
| CAFE [170] | CVPR 2024 | T→I | Llama 2 | SD | 10000×A100 Hrs |
| StoryGPT-V [171] | arXiv 2024 | Story Generation | OPT/Llama2 | Char-LDM | - |
| ELLA [172] | arXiv 2024 | T→I | Llama 2 | SDXL | 1344×A100 Hrs |
| Lavi-Bridge [173] | arXiv 2024 | T→I | Llama 2 | SD/PixArt-$\alpha$ | 8×A100, 48Hrs |
| *Image layout planning* | | | | | |
| LMD [65] | TMLR 2024 | T→I | GPT-3.5/GPT-4 | SD | Training-free |
| LayoutGPT [66] | NeurIPS 2023 | T→I | GPT-3.5/GPT-4/Codex | GLIGEN/SD | - |
| VP-GEN [174] | NeurIPS 2023 | T→I | Vicuna | GLIGEN/SD | 4×A6000, 48Hrs |
| Control-GPT [67] | arXiv 2023 | T→I | GPT-4 | ControlNet/SD | - |
| LayoutLLM-T2I [68] | MM 2023 | T→I | GPT-3.5 | GLIGEN/SD | - |
| LLM Blueprint [175] | ICLR 2024 | T→I | GPT-3.5 | LMD | 1×A100 |
| SLD [176] | CVPR 2024 | T→I | GPT-4 | DALLE3/SD | Training-free |
| TextDiffuser-2 [177] | arXiv 2023 | T→I | Vicuna | SD | 8×A100, 168Hrs |
| COLE [178] | arXiv 2023 | T→I | Llama 2/LLaVA | IF | - |
| *Prompt synthesis and refinement* | | | | | |
| SUR-Adapter [179] | MM 2023 | T→I | LLaMa | SD | - |
| ChatGenImage [180] | arXiv 2023 | Data Synthesis | GPT-3.5 | SD | 1×GTX3090 |
| SwitchGPT [181] | arXiv 2023 | T→TI | Llama 2/GPT-3.5 | SD | 4×A100, 3Hrs |
| TIAC [182] | arXiv 2023 | T→I | GPT-3.5 | SD | - |
| Idea2Img [183] | arXiv 2023 | T→I | GPT-4V | IF/SD | Training-free |
| WordArt Designer [184] | EMNLP 2023 | T→I | GPT-3.5 | ControlNet | 1×V100 |
| *Image quality evaluation* | | | | | |
| DreamSync [185] | arXiv 2023 | T→I | PaLM2/TIFA | SD-XL | - |
| LLMScore [186] | NeurIPS 2023 | T→I | GPT-4 | SD | - |

TABLE 2: Image-language Datasets that can be adopted for language-based image generation.

| Name | Date | Venue | Org | Domain | Source | #Images | Caption |
|---|---|---|---|---|---|---|---|
| Im2Text [187] | 12 Dec 2011 | NeurIPS 2011 | SBU | Open | Internet | 1M | Manual |
| Microsoft-COCO [188] | 1 May 2014 | ECCV 2014 | Microsoft | Common Objects | Internet | 328K | Manual |
| ALIGN [189] | 11 Feb 2021 | ICML 2021 | Google | Open | Internet | 1.8B | Manual |
| Conceptual 12M [190] | 17 Feb 2021 | CVPR 2021 | Google | Open | Internet | 12M | Manual |
| WIT [191] | 2 Mar 2021 | SIGIR 2021 | Google | Open | Wikipedia | 11.5M | Manual |
| LAION-400M [192] | 3 Nov 2021 | NeurIPS 2021 | LAION | Open | Internet | 400M | Manual |
| LAION-FACE [193] | 6 Dec 2021 | CVPR 2021 | Microsoft | Face | LAION | 20M | Manual |
| M3W [194] | 29 Apr 2022 | NeurIPS 2022 | Deepmind | Interleave | Internet | 43M | Manual |
| LAION-COCO [195] | 15 Sep 2022 | - | LAION | Open | LAION | 600M | Synthetic |
| LAION-5B [196] | 16 Oct 2022 | NeurIPS 2022 | LAION | Open | Internet | 5B | Manual |
| Coyo-700M [197] | 31 Aug 2022 | - | Kakao Brain | Open | Internet | 700M | Manual |
| KOSMOS-1 [198] | 27 Feb 2023 | NeurIPS 2023 | Microsoft | Interleave | Internet | 355M | Manual |
| Multimodal C4 [199] | 14 Apr 2023 | NeurIPS 2023 | UCSB | Interleave | Internet | 571M | Manual |
| LLaVA-instruct [200] | 17 Apr 2023 | NeurIPS 2023 | UWM | Instruction | COCO | 150k | Synthetic |
| DATACOMP [201] | 27 Apr 2023 | NeurIPS 2023 | DATACOMP | Open | Internet | 12.8B | Manual |
| MARIO-10M [202] | 19 May 2023 | NeurIPS 2023 | Microsoft | Text within image | LAION, TMDB, OpenLibrary | 10M | Manual |
| LAION-Glyph [203] | 29 May 2023 | NeurIPS 2023 | Microsoft | Text within image | LAION | 10M | Manual |
| MIMIC-IT [204] | 8 Jun 2023 | arXiv 2023 | NTU | Interleave | Internet | 2.8M | Synthetic |

SEED-LLaMA emphasizes the design of the Image Tokenizer, proposing two crucial design principles for its functionality. MiniGPT-5 [165] introduces visual tokens (referred to as "voken") to traditional LLMs, enabling them to gen-

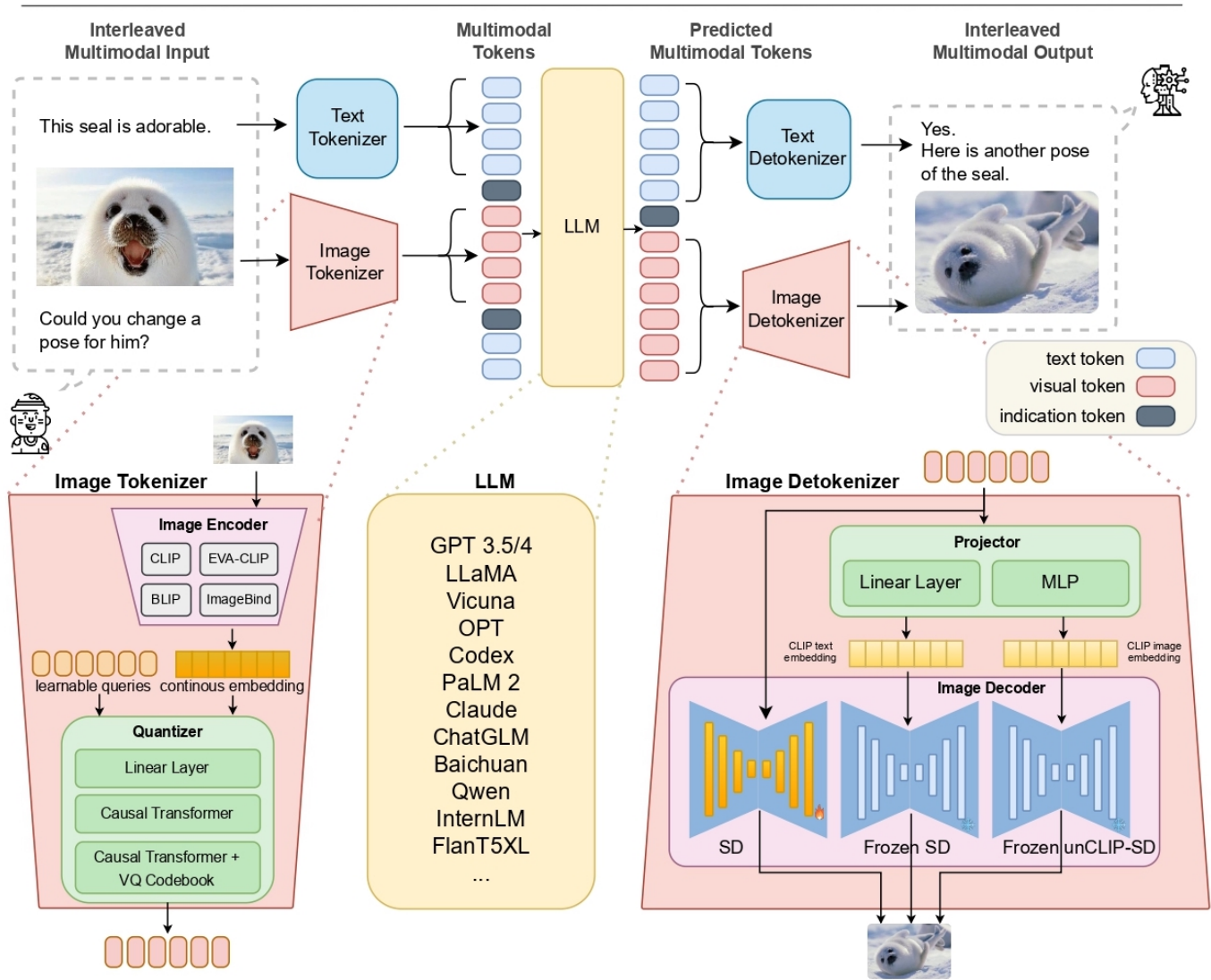## LLM-based Pipeline for Both Image Comprehension and Generation



Fig. 4: A generic pipeline of integrating image comprehension and generation ability on LLMs [58], [61], [62], [159]. During inference time, users can input interleaved multimodal data (*e.g.*, text and images). The image tokenizer processes the information into image tokens and feeds them into the LLM. LLM outputs image tokens and then decodes them into textual responses and images.
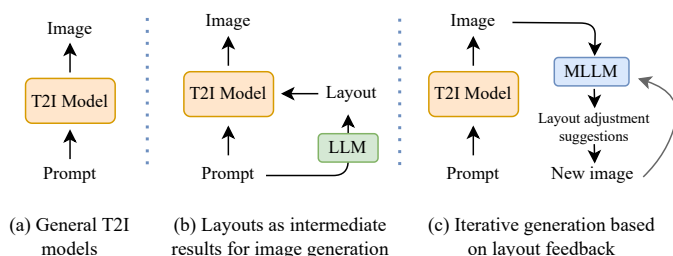


Fig. 5: Pipeline comparisons of (a) standard text-to-image (T2I) [30], [205], (b) T2I with LLMs as layout planners [64]–[68], [174], [175], and (c) T2I with LLMs for layout suggestions [176], [178].

erate images. A two-stage training pipeline is proposed,

involving an unimodal alignment stage and a multimodal learning stage, allowing LLMs to generate both text and images organically. OpenLeaf [166] utilizes prompting of LLMs to generate interleaved textual and visual data, producing entity and style-consistent images and text. It supports various tasks such as how-to question answering, story-telling, graphical story rewriting, and webpage/poster generation. EasyGen [63] leverages a bidirectional conditional diffusion model, BiDiffuser, to endow LLMs with multimodal understanding and generation capabilities. Unlike previous CLIP-based approaches, EasyGen generates images based on this model. TEAL [167] utilizes existing tokenizers for different modalities and transforms the obtained tokens into a joint embedding space, enabling frozen LLMs to understand and generate across various modalities, including text, image, and audio. ChatIllusion [169] introduces Genadapter and

LLaMa-AdapterV2 to bridge the hidden embedding space of SD XL and enable LLMs to understand visual instructions and generate interleaved images and text, supporting image generation, editing, and storytelling. Emu2 [62] emphasizes the in-context learning capability of MLLMs, showcasing improved performance through scaling up the model and unified autoregressive training. It supports tasks such as visual prompting and object-grounded generation, achieving state-of-the-art results in question-answering benchmarks and open-ended subject-driven generation after instruction tuning. ELLA [172] and Lavi-Bridge [173] incorporate large language models into the T→I generation architecture by training several lightweight adapters. LLMGA [168] utilizes LLaVA to simultaneously encode images and instruction, enabling manipulating images based on stable diffusion. StoryGPT-V [171] takes advantage of LLM for coherent story script generation.

### 4.1.3 Image Layout Planning via LLMs

Despite the rapid development of T→I generation, there remain several challenging issues that have yet to be fully addressed, including text rendering, spatial relationships, and quantity representation. Under this circumstance, some methods seek to utilize LLMs for layout planning and subsequently generate images based on the obtained layouts, as shown in Fig. 5. LayoutGPT [66] leverages the inherent reasoning capabilities in LLMs to facilitate layout generation through contextual demonstrations. It utilizes GPT3.5/4 to transform user prompts to CSS-style output layouts, in which the position of each object is specified. LMD [65] improves T→I diffusion models by enhancing prompt understanding capabilities. It utilizes a two-stage approach, leveraging a pre-trained language model to generate scene layouts and guide image generation. VP-GEN [174] breaks down the T2I task into object/count generation, layout generation, and image generation steps. By leveraging GPT-3.5-Turbo fine-tuned on text-layout pairs, VPGEN achieves better spatial control than end-to-end models. Control-GPT [67] takes advantage of GPT-4 to output TikZ code to construct the sketch layouts according to the text descriptions. LayoutLLM-T2I [68] utilizes ChatGPT to induce the layout based on the user prompt. A prompt encoder module is then employed to model the text prompt, relation triplets, and the induced layout separately. To efficiently integrate the layout information, a Layout-aware Spatial Transformer based on the UNet is introduced. LLM Blueprint [175] utilizes ChatGPT to produce detailed object descriptions, bounding box layouts, as well as background prompts. Subsequently, iterative refinement operations are conducted to remedy regional errors based on the layouts. SLD [176] improves T→I generation by iteratively generating images from input prompts and correcting mistakes using an LLM-based layout planner. In particular, the layout planner can add, delete, or resize object boxes to help T2I models produce more accurate images. TextDiffuser-2 [64] employs Vicuna-7B-1.5 for layout planning, generating the position and content of the text to be rendered based on user-provided prompts. COLE [178] harnesses large language models to transform user prompts into detailed JSON files. These files encompass specifications, such as the content, position, and style, for the text to be added.

### 4.1.4 Prompt Synthesis and Refinement via LLMs

LLMs can be treated as a huge knowledge base. Some methods [180], [182]–[184] have explored LLMs to synthesize or optimize prompts, guiding T→I (T2I) models to generate images with rich and detailed content. For example, ChatGenImage [180] utilizes ChatGPT to generate prompts, directing AIGC models in generating preliminary images. Subsequently, it iteratively refines these prompts by incorporating automatically generated detailed annotations as local constraint prompts, resulting in the production of diverse and intricate scenes. Inspired by the three-layer artwork theory, TIAC [182] and WordArt Designer [184] use LLMs to translate abstract concepts into semantically relevant physical objects, making it easier for downstream T→I models. Idea2Img [183] employs a Multimodal LLM to evaluate the generated images of T→I models. Subsequently, based on the obtained feedback, the framework iteratively refines the initial prompts to generate satisfactory results. DiffusionGPT [230] leverages LLMs to refine prompts for image generation. By parsing diverse prompts and utilizing domain-specific Trees-of-Thought, the model selects the most suitable generative model to generate high-quality images. RPG [231] is a training-free framework for T→I generation. It leverages multimodal LLMs to refine the original prompt, decomposing complex prompts into sub-region tasks, and achieves superior performance in object composition and text-image alignment. SUR-adapter [179] leverages LLMs to improve its semantic understanding and reasoning capabilities, enabling it to create better textual semantic representations for T→I generation. SwitchGPT [181] introduces an innovative framework that enables traditional LLMs, like GPT, to interpret the underlying intent of the given instruction, thus producing a more suitable response non-text outputs.

### 4.1.5 Image Quality Evaluation via LLMs

A few works focus on using large language models to evaluate the quality of generated images. For example, Dream-Sync [185] employs two vision language models (VLMs) to evaluate the generated results, and select the best-generated image: one for text alignment and another for aesthetic quality. LoRA [232] is then used to iteratively fine-tune the T2I model towards the selected best generations. LLM-Score [186] converts images into image-level and object-level visual descriptions. Subsequently, a set of instructions is given to the LLMs to check how closely between the images and the descriptions. Finally, a score is generated with reasons.

## 4.2 Image Editing

Image editing is a closely related task with generation and thus receives remarkable progress with the development of image generation models. A curated list of recent image editing methods is shown in Table 3.

### 4.2.1 Image Editing with CLIP/T5

CLIP model enables language-based image editing [233]–[246]. PAIR-Diffusion [243] identifies structure and appearance as the two most intuitive aspects of image editing. Thus, the PAIR-Diffusion is trained using structure and

TABLE 3: An overview of language-based image editing methods based on CLIP and LLMs. We summarize the involved LLMs and generative models, as well as whether the method requires training or not.

| Method | Venue | LLM | Generative Model | Training |
|---|---|---|---|---|
| *CLIP for image editing* | | | | |
| SDEdit [233] | ICLR 2022 | - | DDPM | ✗ |
| DiffusionCLIP [234] | CVPR 2022 | - | DDPM | ✗ |
| P2P [236] | ICLR 2023 | - | Imagen&SD | ✗ |
| NTI [237] | ICLR 2023 | - | SD | ✗ |
| Imagic [238] | CVPR 2023 | - | Imagen | ✓ |
| PaP [239] | CVPR 2023 | - | SD | ✗ |
| SINE [241] | CVPR 2023 | - | SD | ✓ |
| pix2pix-zero [242] | SIGGRAPH 2023 | - | SD | ✗ |
| PAIR-Diffusion [243] | arXiv 2023 | - | PAIR-Diffusion | ✓ |
| MasaCtrl [244] | SIGGRAPH 2023 | - | SD | ✗ |
| Dragondiffusion [247] | ICLR 2024 | - | SD | ✓ |
| DiffEditor [248] | CVPR 2024 | - | SD | ✓ |
| *LLMs for image editing* | | | | |
| InstructPix2Pix [249] | CVPR 2023 | GPT-3 | SD | ✓ |
| VisualChatGPT [116] | arXiv 2023 | GPT-3 | SD | ✓ |
| CHATEDIT [250] | EMNLP 2023 | GPT-3 | StyleGAN | ✓ |
| MGIE [251] | ICLR 2024 | LLaVA-7B | SD | ✓ |
| Emu Edit [252] | arXiv 2023 | Llama 2-70B | Emu | ✓ |
| SLD [176] | CVPR 2024 | GPT-4 | DALL-E 3 | ✓ |
| SmartEdit [253] | CVPR 2024 | LLaVA-7B/13B | InstructDiffusion [254] | ✓ |

appearance information explicitly extracted from the training images, which enables editing of the structure and appearance separately during inference. Instead of relying on the heavy training process of diffusion models on large-scale datasets, another line of work [238], [241] slightly tunes the pre-trained diffusion models to edit the target images. Imagic [238] proposes a generic model for text-based image editing. Notably, Imagic is free of source prompts, which is achieved by optimizing the target prompt embedding. After that, the entire diffusion model is fine-tuned for better reconstruction performance. The editing ability is achieved by interpolating the optimized prompt embeddings and the target prompt embeddings. There also exist many tuning-free methods for text-guided image editing [233]–[237], [239], [240], [242], [244]–[246]. SDEdit [233] can generate realistic images from user inputs such as strokes, sketches, or masks, as well as edit existing images with text instructions. It works by first adding noise to the target image and then gradually denoising it with text instructions. The denoising process is guided by a diffusion model generative prior, which is trained on a large-scale image dataset. SDEdit outperforms state-of-the-art GAN-based methods on multiple image synthesis and editing tasks, according to a human perception study. Other tuning-free methods [236], [239], [242], [244] achieve text-guided editing through the regularization or manipulation of latents, cross-attention maps, or UNet features. However, most text-guided image editing works rely on the CLIP model, whose capability limits the editing to simple text prompts and cannot understand complex human instructions.

### 4.2.2 Image Editing with LLMs

LLMs provide powerful chat-based or interactive editing capabilities for image editing [116], [176], [250]–[252], [268].

InstructPix2pix [268] proposes to use LLMs to construct the data tuples (original image, prompt, target image) to train a model that learns to edit the images following the editing prompt. The model is based on a conditional diffusion model that can process arbitrarily interleaved image and text inputs and produce coherent image (and text) outputs. To generate the data tuples, the authors leverage the knowledge of two large pretrained models: a language model (GPT-3) and a T→I model (Stable Diffusion). The language model generates the editing instructions and the textual descriptions of the edited images, while the T→I model renders the edited images based on the textual descriptions. The authors also introduce a mapping network that translates the hidden representations of the language model into the embedding space of the visual models, enabling the model to use the strong text representations of the LLM for visual outputs.

CHATEDIT [250] further leverages LLMs to contribute an interactive facial image editing system via dialogue. Specifically, CHATEDIT split the dialogue-based editing problem into (1) user edit request tracking, (2) image editing, and (3) response generation subtasks. The user edit request tracking module is responsible for extracting the user's editing intentions from the dialogue history and updating them dynamically. The image editing module is based on a conditional diffusion model that can process both image and text inputs and outputs, and perform various editing operations such as changing hair color, adding glasses, or removing wrinkles. The response generation module is designed to generate natural and informative responses that reflect the editing results and guide the user to the next step. CHATEDIT is evaluated on a novel benchmark dataset proposed by the authors, which contains multi-turn dialogues and corresponding facial images annotated with user edit requests.

MGIE [251] investigates the MLLMs to do the image editing task. The proposed MGIE can learn to convert expressive human instructions to editing guidance. The editing model is also trained to follow the editing guidance in an end-to-end way. The effectiveness of the MGIE is verified on

TABLE 4: An overview of video editing methods based on CLIP and LLMs. We summarize the involved LLMs and generative models, as well as whether the method requires training or not.

| Method | Venue | LLM | Generative Model | Training |
|---|---|---|---|---|
| *CLIP for video editing* | | | | |
| Tune-A-Video [255] | ICCV 2023 | - | SD | ✓ |
| Dreamix [256] | arXiv 2023 | - | Imagen-video | ✓ |
| Video-P2P [257] | arXiv 2023 | - | SD | ✓ |
| FateZero [258] | ICCV 2023 | - | SD | ✗ |
| Pix2Video [259] | ICCV 2023 | - | SD | ✗ |
| StableVideo [260] | ICCV 2023 | - | SD | ✗ |
| Rerender-A-Video [261] | SIGGRAPH Asia 2023 | - | SD | ✓ |
| TokenFlow [262] | ICLR 2024 | - | SD | ✗ |
| CoDeF [263] | CVPR 2024 | - | SD | ✓ |
| MagicEdit [264] | arXiv 2023 | - | SD | ✓ |
| MagicStick [265] | arXiv 2023 | - | SD | ✓ |
| *LLMs for language-based video editing* | | | | |
| InstructV2V [266] | ICLR 2024 | GPT-3 | SD | ✗ |
| InstructVid2Vid [267] | arXiv 2023 | GPT-3 | SD | ✗ |

photoshop-style manipulations, global photo optimization and local editing. SmartEdit [253] is another recent work that utilizes MLLMs for complex instruction-based image editing. SmartEdit analyzes the performance of instruction-based image editing models under complex instructions and proposes a Bidirectional Interaction Module to make the image feature output by pre-trained image encoder and the output feature of LLaVA. They also fine-tune the pre-trained diffusion model to enhance the model's perception and reasoning capability.

Emu edit [252] trained the image editing model in a multi-task way. The tasks include region-based editing, free-form editing and other computer vision tasks, where all tasks are formulated into generative tasks. Emu edit leverages LLMs for instruction generation. Concretely, the authors provide the LLMs with a task description, a few task-specific exemplars, and a real image caption. Then, the LLM is expected to output an editing instruction, an output caption for an ideal output image, and which objects should be updated or added to the original image.

Different from the above works that utilize LLMs to provide editing instructions, SLD [176] uses LLMs to correct the incorrect generation to enable object-level image editing.

### 4.3 Image-language datasets

Image-language datasets play a crucial role in the training of T→I models, providing the foundational data necessary for these models to learn how to generate accurate and relevant visual content from textual descriptions. Over ten years ago, the IM2Text [187] project gathered a huge collection of photos by searching through Flickr, a popular photo-sharing website. They sifted through a massive amount of data and carefully picked out one million images that had clear and directly related captions. MS-COCO [188] collected images depicting complex everyday scenes with common objects in their natural settings. The researchers included five written captions to provide descriptive context. These captions offer a richer understanding of the scene and the objects within it. In recent years, the academic community has witnessed a proliferation of large-scale image-text datasets. Typically, researchers curate these datasets by crawling the internet. For example, LAION-5B [196] is a massive dataset that

was collected by searching the internet for image-text pairs. Using the CLIP model to filter the results, the researchers ensured that the text was relevant to the images. This process resulted in a dataset containing 5.85 billion image-text pairs. Additionally, some researchers are digging into the LAION-5B to find specific types of content. For example, the Mario-10M dataset focuses on pulling out parts of the dataset where the images have text in them to study further. Meanwhile, LAION-FACE [193] is all about images with faces. These specialized subsets help researchers focus on particular areas within the massive collection of image-text pairs. Further, to help image generators follow instructions during a conversation, the LLaVA Visual Instruct 150K dataset [212] comprises a collection of multimodal data designed for instruction-following tasks, generated via a GPT model.

## 5 VIDEO GENERATION AND EDITING

### 5.1 Video Generation

While video understanding has been thoroughly investigated [269]–[280], the past two years have witnessed the rapid development or video generation. Particularly in the field of text-based video generation, numerous works have achieved remarkable results. We list the milestone works in Fig. 6 including both CLIP/T5-based and LLM-based approaches, and summarize the key technical components in Table 5 and commonly-used video-language datasets in Table 6.

#### 5.1.1 Text to Video Generation with CLIP

Based on the type of generative models, there are two main paradigms: one is based on diffusion models, and the other is based on autoregressive models built with transformer architectures and discrete codebooks, trained with the next-token prediction loss. Diffusion models have become a mainstream paradigm due to their ease of training. Within the diffusion framework, there are pixel-level video diffusion models [12], [16], [17] and latent-level video diffusion models [14], [15], [281]–[286]. Pixel-level approaches exhibit better text alignment but require substantial computational resources. On the other hand, latent-level models are more

TABLE 5: Overview of existing methods leveraging LLMs for language-based video generation. We divide these methods into four categories: multimodal LLMs for video generation, video layout planning, and temporal prompt generation. In each method, we summarize the input-output of the task, the involved LLM, and the generative model. In the "Task" column, the "T" and "V" are the abbreviations of "text" and "video", respectively, while "Any" represents the universal generation supporting text, image, video, and audio modalities. Tokenization is the task of converting video into discrete video tokens, which can be viewed as a submodule of some video generation pipelines.

| Method | Venue | Task | LLM | Generative Model |
|---|---|---|---|---|
| *Multimodal LLMs for video generation* | | | | |
| VideoPoet [71] | arXiv 2023 | Any→V | VideoPoet | VideoPoet |
| MAGVIT-v2 [72] | ICLR 2024 | Tokenization | BERT | BERT |
| Video-LaVIT [72] | arXiv 2024 | TIV→TIV | Llama 2-7B | SVD img2vid-xt |
| *Video layout planning* | | | | |
| Dysen-VDM [73] | CVPR 2024 | T→V | GPT-4 | Text2Video-Zero |
| VideoDirectorGPT [74] | arXiv 2023 | TI→V | GPT-4 | LayoutVid |
| LVD [65] | ICLR 2024 | T→V | GPT-3.5/GPT-4 | DSL-grounded generator |
| GPT4MOTION [75] | arXiv 2023 | T→V | GPT-4 | SDXL / ControlNet |
| FlowZero [76] | arXiv 2023 | T→V | GPT-4 | Gligen |
| *Temporal prompt generation* | | | | |
| DirecT2V [77] | arXiv 2023 | T→V | GPT-4 | Text2Video-Zero |
| Free-Bloom [78] | NeurIPS 2023 | T→V | GPT-3.5 | LDM |
| InterControl [79] | arXiv 2023 | T→V | GPT-4 | HMDM |
| PRO-Motion [80] | arXiv 2023 | T→V | GPT-3.5 | Posture-Diffuser |
| VideoDrafter [81] | arXiv 2024 | T→V | ChatGLM-6B | SD-XL |

efficient as they reduce redundancy in video data. Different works in video generation have different focuses. Some emphasize photorealistic or high-definition outputs, aiming to improve the quality of generated videos. Others focus on controllable generation, such as image-to-video approaches that enable local control over motion regions, trajectory control for object and camera movement directions, and the use of sketch, depth, and pose to control structure. Some works concentrate on generating longer videos or exploring better network architectures.

### 5.1.2 Text to Video Generation with LLMs

Recently, several works have also taken advantage of multimodal LLMs [71], [72] for the task of video generation. For example, VideoPoet [71] leverages a pre-trained autoregressive transformer model to handle multimodal data for synthesizing videos with temporal consistency and high-motion fidelity. It adapts LLMs training techniques, allowing for task-specific video generation, including tasks such as text-to-video and image-to-video conversion. Another work, MAGVIT-v2 [72], explores the video tokenization technique of MLLMs. It transforms visual inputs into discrete tokens, enhancing the performance of image and video generation tasks. It outperforms diffusion models on benchmarks such as ImageNet and Kinetics, provides video compression on par with advanced codecs, and improves action recognition through effective representation learning.

### 5.1.3 Video Layout Planning via LLMs

Numerous studies have validated the efficacy of LLMs in generating image layouts, and similarly, recent research has sought to explore the potential of LLMs in crafting video layouts. For instance, some methods employ LLMs to sequentially generate bounding boxes for objects in each frame to assist in the video generation process. VideoDirectorGPT [74] utilizes a language model to generate a video plan that includes object-bound boxes for each scene. These

bounding boxes provide spatial coordinates for the entities, which are used to maintain object consistency and precise layout control throughout the video generation process. LLM-grounded Video Diffusion (LVD) [65] enhances video generation by first using an LLM to create detailed scene layouts from textual prompts, capturing complex motions. These layouts then guide a diffusion model through adjusted attention maps to produce videos that accurately reflect the prompted actions and dynamics, improving upon existing video generation methods. FlowZero [76] utilizes LLMs to generate a dynamic scene syntax that includes object bounding boxes, which are critical for defining object positions and movements within each frame. These bounding boxes guide the diffusion model to ensure accurate object placement and coherent motion throughout the video. Another line of works attempts to utilize layouts beyond conventional bounding boxes. Dysen-VDM [73] enhances the quality of generated videos by employing a dynamic scene graph (DSG). The DSG is used to capture and organize the temporal dynamics of actions described in the text, which are then enriched with details and integrated into a diffusion model to produce more dynamic and realistic videos. GPT4Motion [75] uses GPT-4 to script physical scenes in Blender. The simulated scenes are transformed into intermediate representations, such as depth maps, to serve as layout conditions. Then these conditions are fed to Stable Diffusion to produce the final video, ensuring motion consistency and efficiency in scenarios like object interactions and fluid dynamics.

### 5.1.4 Temporal Prompt Generation via LLMs

In contrast to image generation, video generation requires more elaborate and dense descriptions. The refinement and expansion of prompts can be facilitated by leveraging the capabilities of LLMs. For instance, DirecT2V [77] improves narrative consistency and scene composition from user prompts. It employs instruction-tuned large language mod-

els to decompose a single user prompt into detailed frame-by-frame descriptions. These descriptions guide the generation of each video frame, allowing for the seamless integration of time-varying elements and coherent storytelling. Free-Bloom [78] leverages large language models LLMs to create a sequence of semantically coherent prompts, guiding the video's narrative flow. Pre-trained latent diffusion models (LDMs) are then used as animators to generate high-fidelity frames that visually represent the evolving semantic content, like the process of a flower blooming. InterControl [79] uses an LLM planner to convert textual interaction descriptions into detailed contact plans, improving the generation quality of motion videos. PRO-MOTION [80] uses LLMs to create a sequence of scripts detailing the key postures needed for the target motion. These scripts are based on simple templates, which are different from natural languages and are designed to comprehensively describe all possible postures, thereby simplifying the process for the subsequent generation process. VideoDrafter [81] utilizes LLMs to convert an input prompt into a detailed multi-scene script. This script capitalizes on the logical knowledge of LLMs to ensure the scenes make sense in sequence.

## 5.2 Video Editing

### 5.2.1 Text-guided Video Editing with CLIP/T5

CLIP [24] enables language-based video editing. Here, we mainly discuss the diffusion-model-based video editing approaches. Tune-A-Video [255] presents the early attempt at text-guided video editing with pre-trained diffusion models. Instead of training a giant video diffusion model, Wu et al. [255] propose to inflate and tune the pre-trained text-to-image diffusion model on the target video in a one-shot way. After tuning, the inflated diffusion models support versatile video editing capabilities. Despite its simplicity, Tune-A-Video [255] shows poor temporal stability and the limitation of keeping the unrelated regions unaltered. Video-P2P [257] and FateZero [258] utilize better inversion techniques and propose manipulating attention maps to keep the background unchanged when editing. Pix2Video [259] adopts editing the key frame first and then propagates the editing to other frames, which achieves improved temporal consistency and longer video editing capabilities. Different from these works, Rerender-A-video [261] and CoDeF [263] focus their applications on video-to-video style translation and achieves impressive results through optical-flow-based regularization [261] or adoption of deformable content field [263]. More recently, the developments of CLIP-based video editing are towards better temporal consistency [260], [262], [264], more controllable [265], and more computational effective [287]. We list recent representative video editing works in Table 4.

### 5.2.2 Text-guided Video Editing with LLMs

Existing works that utilize LLMs for video editing are relevantly limited. Current LLMs-based video editing follows a similar scheme as the InstructPix2pix [268], i.e., using LLMs to construct the training data more efficiently.

InstructVid2Vid [267] is one of the works that involve LLMs inefficient training data construction. The method uses an LLM model to generate synthetic video-instruction pairs, which are then used to train an editing model for controllable video editing with natural language instructions. InstructVid2Vid leverages the pre-trained image generation model, i.e., stable diffusion, and a conditional 3D Unet to produce high-quality and temporally coherent videos that match the input video and instruction. To improve the diversity and realism of the synthetic data, InstructVid2Vid incorporates the knowledge and expertise of different models, such as ChatGPT, BLIP, and Tune-a-Video, to synthesize various instructions for the same video. The paper demonstrates the effectiveness of using LLMs to synthesize training data for complex and creative tasks, such as attribute editing, background change, and style transfer.

InsV2V [266] is another approach that extends the paradigm of InstructPix2Pix to the video editing domain. InsV2V uses a large language model to construct synthetic data for training a video editing model, which can also follow natural language instructions to edit videos. InsV2V adopts a one-model-all-video strategy, eliminating the need for per-video-per-model fine-tuning or inversion, and simplifying the user interaction by only requiring an editing prompt. InsV2V leverages a pre-trained image generation model, Stable Diffusion, and a conditional 3D U-Net architecture to produce high-quality and temporally coherent videos that match the input video and instruction. InsV2V demonstrates the versatility and effectiveness of using LLMs to synthesize training data and perform text-based video editing for various tasks, such as object replacement, style transfer, and background change.

## 5.3 Video-language Datasets

The availability of captioned video datasets is crucial for text-to-video generation. To address this challenge, MSR-VTT [288] has introduced a large-scale open-domain video-language dataset that encompasses a wide range of categories and diverse content, setting a new benchmark for the video understanding task in 2016. It contained 200k clip-language pairs from 10K videos on the web, and every video had 20 human annotations in English. Anna et al. have presented the large-scale movie description challenge (LSMDC) [289], which consists of 202 movies accompanied by transcribed audio descriptions. These descriptions provide a narrative of the significant events depicted in the visual video. In natural videos, multiple events often occur within a single video. For instance, a video may feature a man playing the piano, a girl singing, and a crowd applauding. To identify and describe each event, Ranjay et al. have proposed the ActivityNet Caption [290] benchmark, which involves detecting events, describing them using natural language, and localizing them with start and end times. How2 [291] and VATEX [292] are multilingual video description datasets. How2 is a large-scale instructional multimodal and multilingual video dataset that includes English and Portuguese descriptions, videos, speech, and English video-level summaries. VATEX comprises both English and Chinese descriptions, covering 600 human activities. HowTo100M [293] has introduced an automatic video captioning method that utilizes transcribed narrations from web videos instead of manual labeling, enabling fast and scalable data collection. Jonathan et al. have observed that
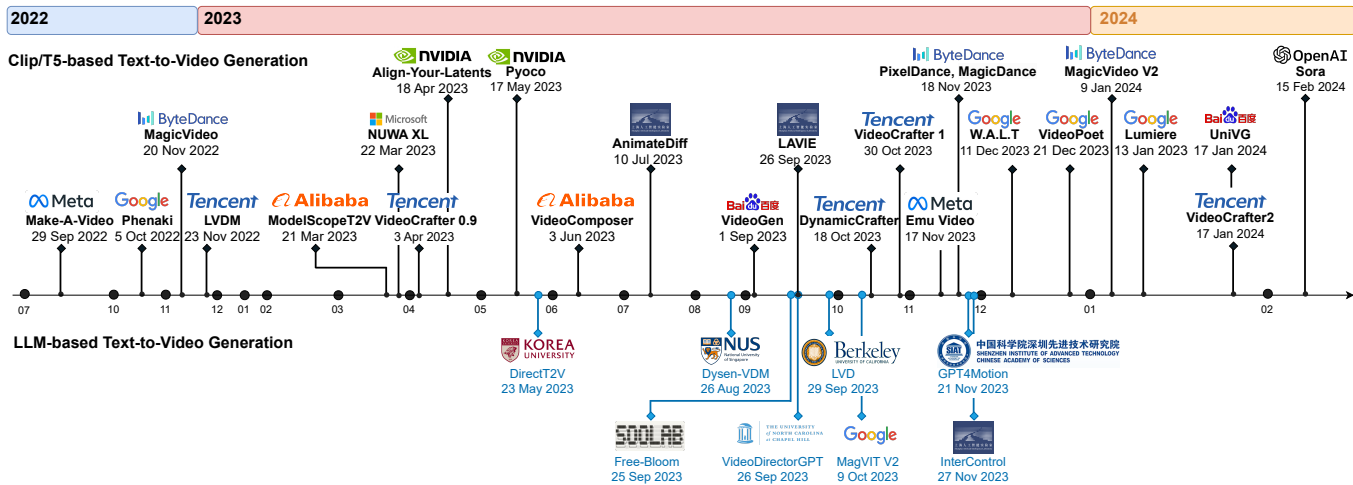
Fig. 6: Milestone works of Clip/T5-based and LLM-based language-guided video generation.

web videos often come with text metadata such as titles and descriptions. They have proposed a data collection process and gathered 70M video clips, known as WTS70M [294], using metadata including titles, descriptions, tags, and channel names. The WebVid [13] dataset has been created for the text-to-video retrieval task. Recognizing the noise in previous datasets like HowTo100M, WebVid-2M, and WebVid-10M have been collected from the internet, featuring weak captions. YT-Temporal-180M [295] is a dataset that contains a diverse corpus of frames, where ASR is provided from a filtered set of 6 million YouTube videos, serving as a resource for multimodal representation learning. HDVILA [296] is a high-resolution large-scale dataset comprising 370k videos, covering 15 popular YouTube categories and offering diverse video content. VideoCC3M [297] proposes a method to transfer captions from existing image captions in CC3M, creating a new weakly labeled audio-video captioning dataset. VideoFactory [298] introduces HD-VG-130M, a dataset consisting of 130M high-definition, widescreen, and watermark-free text-video pairs. InternVid [299] presents a scalable approach to building a high-quality video-text dataset. They employ a multi-scale approach that leverages Tag2Text, LLM, and BLIP2 to generate video captions. Panda-70M [300] is a high-quality and large-scale captioned video dataset proposed in 2024. It contains 70M video clips from YouTube videos, and the caption is extracted via multiple teacher models to obtain multiple captions for one video and a well-trained caption retrieval model to select the best caption. Vript [301] is a fine-grained video-text dataset proposed in 2024 containing 12K annotated videos. Although the number of videos is limited, the caption of each video is fine-grained and comprises information on the shot type, camera movement, content, and scene title.

## 5.4 Summary

In Sec. 5, we have introduced the research works on the generation and editing of the video modality. For each task, we divided the papers into two groups: CLIP/T5-based methods and LLMs-based methods to highlight the

advancement brought by LLMs. We summarize the key technical components of LLMs-based approaches in Table 5, the development of milestone works in the task of language-guided video generation in Fig. 6, and the related video-language dataset in Table 6

# 6 3D GENERATION AND EDITING

Recent studies have focused on establishing a connection between 3D assets and texts. There are two types of methods that can bring the text information for 3D assets, including LLMs and CLIP/T5 models. LLMs can make the output of 3D generation and understanding iterative updates according to user requirements, thereby facilitating highly effective human-computer interaction (i.e., guiding the generation of human motions based on language). Different from LLMs, which directly influence the 3D assets in an interactive iterative manner, the CLIP/T5 model integrates the features of both rendered images and text, enabling the injection of textual information into the 3D assets. In this section, we will explore various methods that leverage the CLIP model or LLMs to guide the processes of 3D generation, editing, and understanding. The generic pipeline is shown in Fig. 7. The corresponding overview of the 3D generation and editing methods is shown in Table 7.

## 6.1 3D Generation

### 6.1.1 3D Generation with CLIP/T5

By leveraging the multi-modal representation capabilities of CLIP, researchers have been able to guide the generation and editing of 3D assets using textual descriptions or queries, thereby enabling more precise control and customization. Specifically, CLIP [24] extracts the feature of image and text with two encoders and aligns them in a contrastive learning manner. This alignment builds the connection between image and text effectively which brings significant improvements to text-guided 2D [25], [30], [137], [205], [305], [306], text-guided 3D [25], [43], [138] generations. There are three typical methods to utilize the CLIP model to provide text information during 3D generation: 1. Adopting the pre-trained CLIP model as supervisor to calculate the CLIP

TABLE 6: Public video-language datasets can be adopted for language-guided video generation. For each dataset, we list the following information in each column: dataset name (Dataset), paper conference venue (Venue), dataset domain (Domain), video source (Vid. Source), video spatial resolution (Res.), average duration per clip (Dur./Clip), total number of clips (#Clips), total number of videos (#Videos), total number of hours (#Hours), and caption source (Cap. Source). The dataset is sorted in ascending order of its released time.

| Dataset | Venue | Domain | Vid. Source | Res. | Dur./Clip | #Clips | #Videos | #Hrs | Cap. Source |
|---|---|---|---|---|---|---|---|---|---|
| MSR-VTT [288] | CVPR 2016 | Open | Internet | 240p | 15s | 10K | 7K | 41 | Human |
| LSMDC [289] | IJCV 2017 | Movie | Amazon | 1080p | 4.8s | 118K | 202 | 158 | Audio Desc. |
| ActivityNet Captions [290] | ICCV 2017 | Activity | Internet | - | 36s | 100K | 20K | 849 | Human |
| How2 [291] | NIPS 2018 | Instruction | Youtube | - | 90s | 80K | - | 298 | Human |
| VATEX [292] | ICCV 2019 | Open | Youtube | 240p | 10s | 41K | 41K | - | Human |
| HowTo100M [293] | ICCV 2019 | Instruction | Youtube | 240p | 4s | 136M | 1.2M | 134K | ASR |
| WTS70M [294] | arXiv 2020 | Open | Youtube | - | 10s | 70M | 70M | 194K | Metadata |
| WebVid-10M [13] | ICCV 2021 | Open | Internet | 360p | 18s | 10.7M | 10.7M | 52K | Alt-text |
| YT-Temporal-180M [295] | NeurIPS 2021 | Open | Youtube | - | - | 180M | 6M | - | ASR |
| HD-VILA-100M [296] | CVPR 2022 | Open | Youtube | 720p | 13.4s | 103M | 3.3M | 372 | Algorithm |
| CelebV-Text [302] | CVPR 2023 | Face | Internet | $512^2+$ | 14s | 70K | - | 279 | Human+Algorithm |
| VideoCC3M [297] | ECCV 2022 | Open | CC3M | - | 10s | 6.3M | - | 17.5K | CC3M |
| HD-VG-130M [298] | arXiv 2023 | Open | Youtube | 720p | - | 130M | - | - | Algorithm |
| InternVid [299] | ICLR 2024 | Open | Youtube | 720p | 12s | 234M | 7M | 760K | Algorithm |
| Panda-70M [300] | CVPR 2024 | Open | Youtube | 720p | 8.5s | 70.8M | - | 166.8K | Algorithm |
| Vript [301] | Github | Open | Youtube | 720p-2K | 11.7s | 400K | 12K | 1.3K | Algorithm |

loss between the generation image and text or injecting the text feature to the 3D assets with CLIP model directly. 2. Leveraging pre-trained text-to-image generation models as supervisory signals and utilizing distillation loss functions such as SDS [43], [307] to distill 3D assets [308]–[310]. 3. Utilizing the public annotated 3D Datasets [304], [311]–[313] to get the mesh, NeRF or multiview image, then set these outputs to train a 3D diffusion model. As mentioned above, we will elaborate on the current methods in two categories in the following content.

**CLIP/T5 Model Supervisor.** Text2Mesh [27] focuses on stylizing 3D meshes by predicting color and local geometric details that align with target text prompts. The entire process is guided by a CLIP loss, which helps ensure the generated meshes conform to the desired textual specifications. This method offers enhanced control over the visual appearance and geometric attributes of 3D meshes, enabling the creation of visually appealing and semantically meaningful shapes.

TANGO [314] proposes a pipeline for generating texture on a given mesh. By leveraging the CLIP model, TANGO aligns the texture generation process with textual descriptions, allowing for the synthesis of textured meshes that match specific visual or semantic criteria. This approach facilitates the creation of realistic and visually coherent 3D models with detailed surface textures. CLIP-Mesh [26] addresses unsupervised text-guided 3D generation by optimizing the texture, normal, and vertical position of a 3D object using the CLIP loss. This approach enables the generation of 3D objects that align with textual prompts, providing a powerful tool for text-driven content creation and design. X-mesh [315] further improves the performance of CLIP-Mesh by adopting an attention-based network, enhancing the fidelity and accuracy of the generated 3D meshes. CLIP-forge [25] introduces a zero-shot text-to-shape method for predicting volumetric occupancy by leveraging the text feature of the CLIP model in conjunction with a conditional normalizing flow network [316]. This approach enables the generation of 3D shapes based on textual prompts without the need for explicit supervision or labeled training data.

Some methods employ CLIP guidance to output NeRF representations [308], which are used to model complex scenes and capture high-frequency spatial information. DreamFields [317] introduces general-purpose priors to assist in aligning the optimized NeRF with the given text prompt, improving the quality and fidelity of the generated scenes. CLIP-NeRF [138] takes a two-step approach, first training a disentangled conditional NeRF and then utilizing the text feature to adjust the parameters of the learned NeRF, enabling more fine-grained control over the generated scenes.

ShapeGPT [318] utilizes a "word-sentence-paragraph" pipeline to convert shapes into words. These words are then assembled to form shape sentences, which can be integrated with instructional text to create multi-modal paragraphs describing 3D shapes. These multi-modal paragraphs help the ShapeGPT do several applications, including text-to-shape generation, image-to-shape generation, and multimodal-to-shape completion and editing.

Moving beyond object generation, MotionCLIP [319] proposes a 3D human motion auto-encoder that predicts pose sequences. By leveraging CLIP, this method enables the generation of realistic and contextually coherent human motion based on textual prompts, providing a means for text-driven animation and virtual character control. MotionGPT [320] considers human motion as a distinct language and trains a motion-language model with the T5 model. The approach incorporates discrete vector quantization to represent human motion and transforms 3D motion into motion tokens. By establishing a comprehensive "motion vocabulary," the model conducts language modeling on both motion and text in a cohesive manner.

While the aforementioned methods have achieved notable success in text-guided generation, they still face challenges related to visual artifacts. This can be attributed to the semantic-level nature of the CLIP loss, which tends to reduce the high-frequency spatial information in the generated images. Addressing this limitation remains an active area of research, with the aim of further improving the visual quality and fidelity of text-guided 3D generation approaches.

**Text-to-image Model Supervisor.** In contrast to the aforementioned methods that directly utilize the CLIP model or
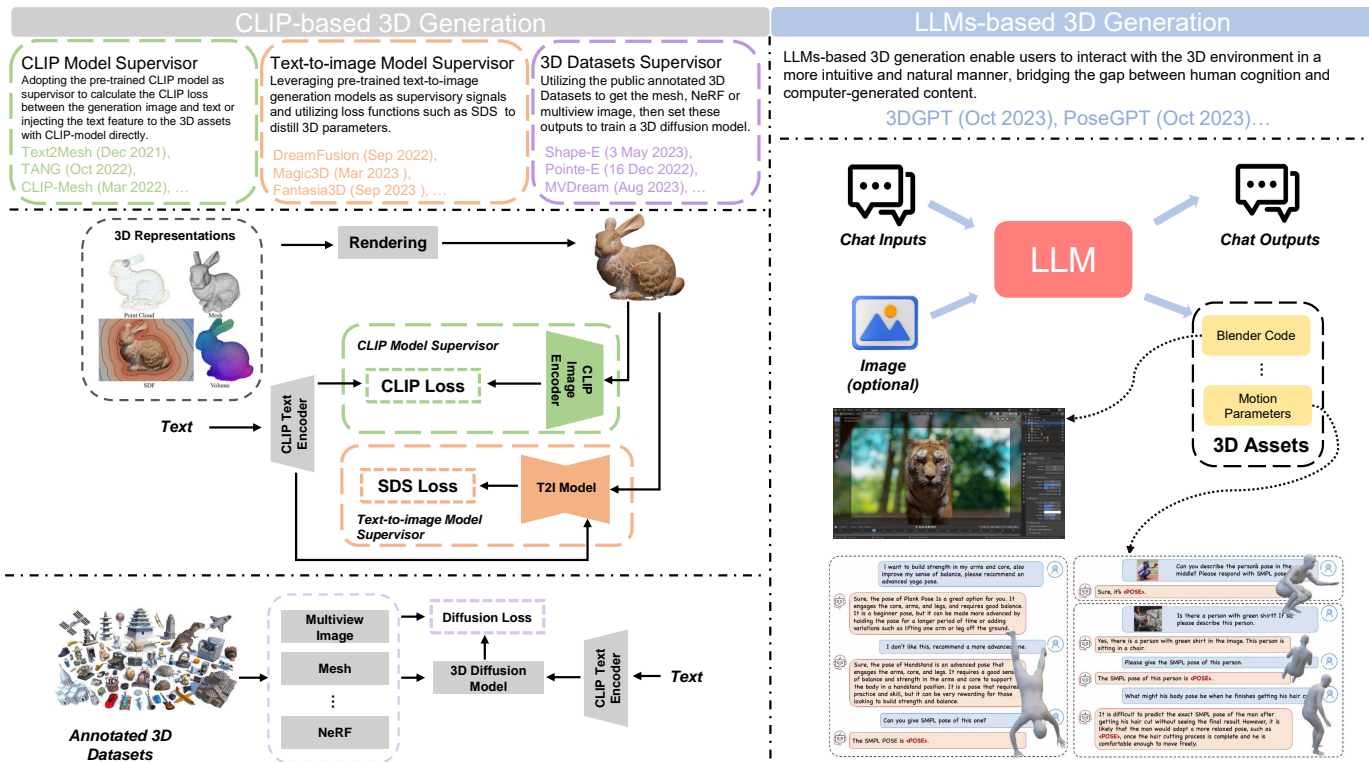
Fig. 7: A generic pipeline of 3D generation with CLIP and LLMs. The CLIP-based models optimize the 3D representation by minimizing the distance between the rendering image and the text prompt. In order to better improve interaction efficiency, LLMs-based methods try to transfer the language outputs of LLMs to blender code or 3D representation (i.e. human motion) directly. Some images are borrowed from [83], [303], [304]

CLIP loss for 3D generation, recent approaches [33], [36], [38]–[40], [42], [45], [307], [321]–[323] have focused on distilling 3D assets from pre-trained text-to-image generation models (i.e., Stable Diffusion [30]). These models employ score distillation sampling (SDS) loss within the framework of DreamFusion [43]. Specifically, text-to-2D image generation models utilize the text feature from CLIP to train a generative model (i.e., diffusion model) within a text-image pair dataset, SDS [43] simulates the training process of diffusion model and calculates the spatial gradient of the rendered image to optimize the parameters of NeRF. Many subsequent methods have utilized the SDS loss function as a supervisory signal and introduced various techniques to improve the performance of 3D generation.

For instance, Magic3D [44] combines DMTed [310] to transfer the NeRF model to a mesh representation, enabling high-resolution rendering results. The experimental results presented in the paper demonstrate the effectiveness of Magic3D in enhancing the generation quality of both geometry and texture. Similar to the Magic3D, TextMesh [42] replaces the NeRF model with a textured mesh representation for 3D asset generation. LatentNeRF [324] introduces a latent space optimization strategy. It takes advantage of the latent space structure learned by the text-to-image model, which has already been trained on a large-scale text-image dataset. By aligning the NeRF optimization process with the latent space structure, LatentNeRF improves the convergence and stability of the training process.

Fantasia3D [41] presents a novel approach to generating

3D assets by separating the geometry and texture components. It first refines the normal map to generate the geometry and then fixes the geometry to predict the color field. Moreover, Fantasia3D set the color field as a Physically Based Rendering (PBR) material model to enhance the generation fidelity. Moreover, artists and designers can manipulate and modify the geometry and texture components independently, allowing for a wider range of artistic expression and customization possibilities.

Despite the significant improvements brought about by SDS, there remain some limitations in the results generated. These include over-saturation, over-smoothing, multi-face Janus artifacts, and time-consuming computations. To address these issues and further enhance the generation performance, some methods [34], [35], [325] have been proposed.

For example, Perp-Neg [326] debiases the score-distillation framework for view-consistent text-to-3D generation, aiming to mitigate the multi-face Janus problem. Prolificdreamer [37] introduces variational score distillation (VSD) to avoid over-saturation and over-smoothing. It minimizes the KL divergence between the optimized 3D assets and the target distribution, which makes the generation process utilize the normal CFG weight. DreamPropeller [327] presents a drop-in acceleration algorithm to expedite the training process. It extends the concept of Picard iterations, a well-established algorithm for parallel sampling of an ODE path, to encompass various scenarios beyond ODEs. This includes accommodating momentum-based gradient

updates and handling dimensional changes during the optimization process, which are often encountered in the context of 3D generation. Additionally, some methods [31], [32] have replaced the NeRF model with Gaussian splatting [29] techniques to improve training efficiency and generate high-quality geometry.

Apart from the methods mentioned, which focus on general object generation, some methods [328]–[344] try to explore 3D generation in 3D avatar generation.

**3D Datasets Supervisor.** Some methods [345]–[348] try to train the multi-view diffusion model with diffusion loss or reconstruction model [349]–[351] for text-to-3D generation, and these methods always adopt the text-to-3D dataset [304], [311], [312] as a training label.

Point-E [313] and Shape-E [352] utilize the Blender to collect a point cloud 3D dataset first, then train a diffusion model to generate the point cloud and mesh with text condition, respectively.

MVDream [345], RichDreamer [348], SPAD [346] and UniDreamer [347] get the multi-view images from the Objaverse dataset. With these multi-view images, they Then, fine-tune the 2D text-to-image diffusion to a 3D diffusion model to predict the multi-view images under the guidance of text and camera poses.

While CLIP-based 3D generation methods have demonstrated impressive progress, the CLIP model itself cannot maintain flexibility in human-computer interaction during the 3D generation process. More recently, Large Language Models (LLMs) have been introduced to enhance human-computer interaction in 3D generation, which will be discussed in the following section.

### 6.1.2  3D Generation with LLMs

The integration of LLMs with 3D assets has emerged as a promising research direction recently. By harnessing the powerful language understanding capabilities of LLMs, researchers aim to directly enhance the performance of generation, manipulation, or understanding of 3D assets through textual instructions. These approaches enable users to interact with the 3D environment more intuitively and naturally, bridging the gap between human cognition and computer-generated content.

3D-GPT [82] proposes a training-free framework that utilizes LLMs. This framework consists of three agents: the task dispatch agent, the conceptualization agent, and the modeling agent. By employing these agents, 3D-GPT can produce the code for blender corresponding to the language, and the efficiency of end-users engaged in procedural 3D modeling is improved. Similar to the 3D-GPT, SceneCraft [353] introduces an LLM agent to transfer input text query into a 3D scene by generating a Blender script. Specifically, sceneCraft has a dual-loop self-improving pipeline: in the inner-loop, per each scene, an LLM autonomously writes a script to interact with Blender, receives rendered images, and keeps improving the script until getting good scenes; in the outer-loop, SceneCraft summarizes common functions over a batch of written scripts to maintain a reusable design skill library.

LL3DA [84] proposed a Large Language 3D Assistant that utilizes a transformer network to predict querying tokens. These tokens are projected to the prefix of textual instructions as the input of a frozen LLM. Finally, the LLM will produce the answer to the textual instructions.

PointLLM [354] processes colored point clouds using human instructions and predicts responses to user questions with the help of LLMs. This enables users to analyze and interpret point clouds more effectively.

3D-LLM [355] takes 3D points with features and language prompts as input and performs various 3D-related tasks, leveraging the capabilities of LLMs. The pipeline involved in this study involves the collection of a comprehensive dataset comprising over 300,000 instances of 3D-language data. This dataset encompasses a wide range of diverse 3D-related tasks, which include, but are not limited to, 3D captioning, dense captioning, 3D question answering, task decomposition, 3D grounding, 3D-assisted dialog, navigation, and various other tasks.

In summary, the integration of LLMs and 3D assets has opened up new possibilities for generating, manipulating, and understanding 3D content through natural language instructions. These methods have demonstrated significant advancements in enhancing human-computer interaction in the 3D domain.

## 6.2  3D Editing

Similar to the generation, we separate the text-to-3D editing into two aspects. And we will first discuss the CLIP/T5-based editing methods.

### 6.2.1  3D Editing with CLIP/T5

**CLIP/T5 model Supervisor.** Blended-NeRF [356] introduced a framework for modifying a specific region of interest within an existing NeRF scene using CLIP loss. This approach allows for targeted editing of NeRF scenes by leveraging the power of contrastive learning with CLIP. NeRF-Art [28] proposed a global-local contrastive learning strategy to stylize pre-trained NeRF models. By employing contrastive learning, NeRF-Art enables the creation of artistic and stylized renderings through the manipulation of pre-existing NeRF scenes. TextDeformer [357] employed text-to-geometry manipulation by introducing a mesh deformation technique based on Jacobians. This approach utilizes text descriptions to deform the geometry of objects, offering a novel way to manipulate and edit 3D models. Sine [358] presented a prior-guided editing field that encodes fine-grained geometric and texture modifications. By leveraging this approach, Sine enables precise and detailed editing of both geometry and texture in 3D scenes, providing a powerful tool for creative expression.

**Text-to-image Model Supervisor.** SKED [359] utilizes sketches as a guiding input for text-to-image generation models in SDS. By incorporating sketch information, SKED enhances the generation process, resulting in more accurate and contextually aligned text-to-image translations. DreamEditor [360] transfers NeRF representations to meshes and employs a personalized text-to-image generation model called DreamBooth for mesh editing. This approach allows for interactive and personalized editing of meshes using text-based instructions, enabling users

to modify and shape 3D scenes according to their preferences. Instruct-NeRF2NeRF [361] utilizes the Instruct-Pix2Pix model to incorporate SDS loss for editing NeRF scenes with text-based instructions. By combining the power of text instructions with NeRF editing, this approach enables users to modify scenes in a controlled and precise manner, enhancing the interactive editing capabilities of NeRF. 3D Paintbrush [362] proposes a cascaded score distillation method for editing the texture of local semantic regions in meshes. By distilling scores, this approach allows for targeted and localized texture editing, providing users with a powerful tool for enhancing the visual appearance of specific regions in 3D models.

### 6.2.2 3D Editing with LLMs

Different from CLIP-based methods, there is no specific method to do the 3D editing with LLMs. Editing is more like a sub-task of LLMs-based 3D generation, so some generation methods (i.e., 3D-GPT [82], SceneCraft [353] ) can edit the 3D assets directly. We will follow the latest developments in LLMs-based 3D editing and discuss them in the future.

## 6.3 Summary

The utilization of CLIP or LLMs in the context of 3D generation and editing offers several advantages. Firstly, it enables users to express their creative intentions or desired modifications in natural language, simplifying the interaction process and reducing the need for specialized software or technical expertise. Additionally, the incorporation of textual information into the 3D generation pipeline enhances the interpretability and explainability of the generated outputs, allowing users to better understand and fine-tune the results according to their requirements.

In conclusion, the integration of CLIP or LLMs with 3D assets has opened up new avenues for highly effective human-computer interaction. By aligning textual information with the visual features of 3D assets, researchers have been able to facilitate more intuitive and precise control over the generation and editing of 3D content. These advancements hold great potential for applications in fields such as computer graphics, virtual reality, and augmented reality, offering enhanced user experiences and empowering users to unleash their creativity more seamlessly and efficiently.

## 7 AUDIO GENERATION, UNDERSTANDING AND EDITING

Recently, a surge of innovative works, such as [86], [97], [100], [102], [104], [106], has demonstrated the utilization of LLMs in various audio-related tasks. These tasks span across domains including the creation of audio effects, speech processing, and music composition, showcasing the versatility of LLMs.The roles of LLMs in these areas are varied, acting as backbones for complex systems [86]–[98], conditioners for specific tasks [99]–[101], labellers for audio content [102]–[104], agents in interactive environments [105]–[110], and as inspiration for some approaches [49], [50], [55], [111]–[113]. This surge in the application of LLMs in the audio field is not only reshaping how we interact with sound and music but also opening new frontiers in the intersection of AGI and audio technologies.

## 7.1 Domains

The works that integrate LLMs in the field of audio can be divided into three categories: general audio sounds, music, and speech. Each category poses the challenges and prospects from three pivotal perspectives—generation, understanding, and editing. The summary of the key technical components of LLMs-based approaches for LLMs-related audio tasks is shown in Table 9, and the related audio-language dataset in Table 8. The timeline of milestone works is shown in Fig. 8.

### 7.1.1 General Audio Sounds

General Audio Sounds refers to any sound that can be heard. It includes various auditory experiences, including natural sounds (like birds chirping and wind rustling through trees), human activities (such as traffic noise and machinery), and other environmental noises. In the past few months, the field of general audio has experienced significant advancements through the application of LLMs [87], [93], [94], [100]–[102], [104]–[106], [112], [113], [393]. In the following sections, we will explore the specific areas of audio generation, audio understanding, and audio editing to analyze how these LLMs-driven developments are reshaping the audio domain.

**Audio Understanding.** General audio understanding involves the analysis and interpretation of a wide array of sounds from our environment, beyond just speech and music. This task includes identifying and classifying sounds (such as distinguishing a car horn from a dog bark), recognizing patterns in environmental sounds (like detecting the sound of rainfall or an approaching vehicle), and even understanding the context or source of sounds.

A suite of groundbreaking models, LTU (Listen, Think, and Understand) [87], SALMONN (Speech Audio Language Music Open Neural Network) [93], Qwen-Audio [94], and UNIFIED-IO 2 [393], have utilized LLMs as their backbone for audio understanding. Unlike LTU [87], which is the first multimodal LLM to focus on general audio understanding beyond just speech, SALMONN [93] is the first multimodal LLM capable of perceiving and understanding general audio inputs with speech, audio events, and music. By integrating audio with other data modalities, UNIFIED-IO 2 [393] leverages LLMs to enhance the understanding of complex interactions between various types of input. Qwen-Audio [94] improves interaction capabilities of pre-trained audio models by covering over 30 different tasks and various audio types, including human speech, natural sounds, music, and songs, thereby promoting comprehensive audio understanding capabilities. To improve user interaction, models like AudioGPT [106] and HuggingGPT [105] also leverage LLMs to serve as intelligent interfaces. The works [87], [93], [94], [101], [393] demonstrate how LLMs can be used to enhance automated audio captioning.

LTU [87] combines the audio perception model AST [399] with LLaMA [145] to improve audio understanding through a perception-to-understanding curriculum. For this purpose, this work also constructs the OpenAQA-5M dataset, which includes 1.9 million closed-ended and 3.7 million open-ended tuples. This dataset facilitates LTU's training within an autoregressive framework.

TABLE 7: Summary for 3D general object generation. The optimization target means the essential constraint during the learning process. The representation means the type of the 3D output. The method without an optimization target means this method is not guided by CLIP loss or SDS-based loss.

| Method | Venue | Optimization Target | Representation | Guided Model |
|---|---|---|---|---|
| *CLIP/T5 for 3D generation* | | | | |
| MotionCLIP [319] | ECCV 2022 | CLIP Loss | Motion Sequences | CLIP |
| MotionGPT [320] | NeurIPS 2023 | - | Motion Sequences | T5 |
| MDM [363] | ICLR 2023 | - | Motion Sequences | CLIP |
| CLIP-Mesh [26] | SIGGRAPH Asia 2022 | CLIP Loss | Mesh | CLIP |
| TANGO [314] | NeurIPS 2022 Spotlight | CLIP Loss | Mesh | CLIP |
| DreamFields [317] | CVPR 2022 | CLIP Loss | NeRF | CLIP |
| Clip-forge [25] | CVPR 2022 | CLIP Loss | Voxel | CLIP |
| Text2Mesh [27] | CVPR 2022 | CLIP Loss | Mesh | CLIP |
| TextMesh [42] | 3DV 2023 | CLIP Loss | Mesh | CLIP |
| X-Mesh [315] | ICCV 2023 | CLIP Loss | Mesh | CLIP |
| ShapeGPT [318] | arXiv 2023 | - | SDF | T5 |
| Shape-E [352] | arXiv 2023 | Diffusion Loss | Mesh/NeRF | CLIP |
| Point-E [313] | arXiv 2023 | Diffusion Loss | PointCloud | CLIP |
| DreamFusion [43] | ICLR 2023 Oral | Score Distillation | NeRF | Imagen |
| SJC [307] | CVPR 2023 | Score Distillation | NeRF | SD |
| Magic3D [44] | CVPR 2023 Highlight | Score Distillation | NeRF | SD |
| Perp-Neg [326] | arXiv 2023 | Score Distillation | NeRF | SD |
| Latent-NeRF [324] | CVPR 2023 | Score Distillation | NeRF | SD |
| Fantasia3D [364] | ICCV 2023 | Score Distillation | NeRF | SD |
| ATT3D [38] | ICCV 2023 | Score Distillation | NeRF | SD |
| ProlificDreamer [37] | NeurIPS 2023 Spotlight | Score Distillation | NeRF | SD |
| Text2Room [33] | ICCV 2023 | Score Distillation | Mesh | SD |
| 3DFuse [45] | ICLR 2024 | Score Distillation | NeRF | SD |
| GaussianDreamer [31] | CVPR 2024 | Score Distillation | Gaussian Splatting | SD |
| DreamGaussian [32] | ICLR 2024 | Score Distillation | Gaussian Splatting | SD |
| NFSD [325] | ICLR 2024 | Score Distillation | NeRF | SD |
| MVDream [345] | ICLR 2024 | Diffusion Loss | Multi-view images | SD |
| RichDreamer [348] | CVPR 2024 | Diffusion Loss | Multi-view images | SD |
| SPAD [346] | CVPR 2024 | Diffusion Loss | Multi-view images | SD |
| UniDreamer [347] | CVPR 2024 | Diffusion Loss | Multi-view images | SD |
| Enhancing3D [365] | ICLR 2024 | Score Distillation | NeRF | SD |
| LucidDreamer [34] | CVPR 2024 | Score Distillation | Gaussian Splatting | SD |
| CSD [35] | ICLR 2024 | Score Distillation | NeRF | SD |
| SweetDreamer [36] | ICLR 2024 | Score Distillation | NeRF | SD |
| HiFA [40] | ICLR 2024 | Score Distillation | NeRF | SD |
| AToM [366] | arXiv 2023 | Score Distillation | Mesh | SD |
| Consistent3D [367] | arXiv 2023 | Score Distillation | Mesh/NeRF | SD |
| DreamControl [368] | CVPR 2024 | Score Distillation | NeRF | SD |
| IT3D [369] | AAAI 2024 | Score Distillation | NeRF | SD |
| Efficientdreamer [370] | CVPR 2024 | Score Distillation | NeRF | SD |
| GSGEN [371] | CVPR 2024 | Score Distillation | Gaussian Splatting | SD |
| X-Dreamer [372] | arXiv 2023 | Score Distillation | Gaussian Splatting | SD |
| HD-Fusion [373] | WACV 2024 | Score Distillation | Gaussian Splatting | SD |
| LODS [374] | arXiv 2023 | Score Distillation | Gaussian Splatting | SD |
| Sherpa3d [375] | CVPR 2024 | Score Distillation | NeRF | SD |
| DreamPropeller [327] | CVPR 2024 | Score Distillation | NeRF | SD |
| DreamPolisher [376] | arXiv 2024 | Score Distillation | Gaussian Splatting | SD |
| *LLM for 3D generation* | | | | |
| 3D-GPT [82] | arXiv 2023 | - | Blender Code | GPT-3.5 |
| PoseGPT [83] | CVPR 2024 | - | Motion Sequences | LLaVA |
| HOLODECK [377] | CVPR 2024 | - | Scene | GPT-4 |
| LL3DA [84] | arXiv 2023 | - | PointCloud | GPTV |
| SceneCraft [353] | arXiv 2023 | - | Blender Code | GPT-3.5 |
| *CLIP for 3D editing* | | | | |
| CLIP-NeRF [138] | CVPR 2022 | CLIP Loss | NeRF | CLIP |
| Blended-NeRF [356] | ICCVW 2023 | CLIP Loss | NeRF | CLIP |
| SKED [359] | ICCV 2023 | Score Distillation | NeRF | SD |
| DreamEditor [360] | SIGGRAPH Asia 2023 | Score Distillation | NeRF | SD |
| Instruct-NeRF2NeRF [361] | SIGGRAPH Asia 2023 | Score Distillation | NeRF | SD |
| TextDeformer [357] | TVCG 2022 | Score Distillation | Mesh | SD |
| SINE [358] | CVPR 2023 | Score Distillation | NeRF | SD |
| Blending-NeRF [378] | ICCV2023 | CLIP Loss | NeRF | CLIP |
| CustomNeRF [379] | CVPR 2024 | Score Distillation | NeRF | SD |
| Paint3D [380] | arXiv 2023 | - | Mesh | SD |
| 3D Paintbrush [362] | arXiv 2023 | Score Distillation | NeRF | SD |

TABLE 8: Audio datasets that can be adopted for language-based audio research. For each dataset, we list the following information in each column: dataset name (Dataset), paper conference venue (Venue), average duration per clip (Dur./Clip), total number of clips (#Clips), total number of hours (#Hours), and dataset domain (Domain).

| Dataset | Venue | Dur./Clip | #Clips | #Hours | Domain |
|---|---|---|---|---|---|
| MagnaTagATune [381] | ISMIR 2009 | 29s | 25,863 | 208h | Music |
| Librispeech [382] | ICASSP 2015 | - | - | 1,000h | Speech |
| Audioset [383] | ICASSP 2017 | 10s | 2M | - | Audio |
| MAESTRO [384] | ICLR 2019 | - | - | 200h | Music |
| Libri-TTS [385] | INTERSPEECH 2019 | - | - | 585h | Speech |
| MTG-Jamendo [386] | ICMLw 2019 | - | 55,000 | - | Music |
| Librilight [387] | ICASSP 2020 | - | - | 60,000h | Speech |
| Vggsound [388] | ICASSP 2020 | 10s | 210,000 | 550h | Audio |
| WenetSpeech [389] | ICASSP 2022 | - | - | 22,400h | Speech |
| Libri-heavy [390] | ICASSP 2024 | - | - | 50,000h | Speech |



Fig. 8: Milestone works of LLMs-based audio research, including audio generation, understanding, and editing.

SALMONN [93] processes general audio inputs, including speech, events, and music, by integrating a text-based LLM with speech and audio encoders together. This fusion improves SALMONN's understanding ability across various audio phenomena.

UNIFIED-IO 2 [393] is the first autoregressive multimodal model that integrates text, images, audio, and actions into a unified framework. Using a single encoder-decoder transformer model, it tokenizes inputs from different modalities into a shared semantic space for processing.

Qwen-Audio [94] scales up audio-language pre-training to include over 30 tasks. To address the interference problems that arise from training all tasks and datasets together, a multi-task training framework was designed. This framework uses a sequence of hierarchical tags for the decoder, which helps share knowledge and prevent interference by using both shared and specific tags. Further developed upon Qwen-Audio, Qwen-Audio-Chat can take input from different audio and text sources, allowing for multi-turn conversations and supporting various audio-focused scenarios.

AudioGPT [106] and HuggingGPT [105] showcase the use of LLMs for audio understanding by coordinating tools through LLMs-driven interfaces. AudioGPT utilizes ChatGPT as a central node for audio and speech applications, depending on external audio systems for functionality. HuggingGPT functions as an agent that combines ChatGPT's language capabilities with a diverse set of AI models from the Hugging Face community, improving its ability to understand audio content.

Wu et al. [101] focuses on advancing automated audio captioning (AAC), a field dedicated to generating descriptive text for sounds from nature and human activities. This work pushes this development further by extensively integrating pretrained models and LLMs. Wu et al. employ BEATS for extracting detailed audio features and use the INSTRUCTOR LLM for obtaining text embeddings of captions. Additionally, Wu et al. introduce a data augmentation technique using ChatGPT to create caption mix-ups and enrich the training data in terms of quantity, complexity, and diversity.

**Audio Generation.** Audio generation is an emerging field that focuses on modeling the creation of diverse audio content. The application of LLMs has significantly advanced the generation of audio. Significant contributions in text-to-audio generation include TANGO [100], Make-an-Audio 2 [102], WavJourney [107], AudioLM [112], and Audiobox [104]. From employing text embedders and diffusion models in TANGO [100] and Make-an-Audio 2 [102] to integrating multi-modal approaches in WavJourney [107] and advanced tokenization in AudioLM [112] and Audiobox [104], these initiatives highlight the versatility and impact of LLMs in driving forward the capabilities of audio generation technologies.

TANGO employs FLAN-T5 [400] as a text embedder, while Make-an-Audio 2 uses pre-trained LLMs for text

TABLE 9: Summary of approaches for LLMs-related audio tasks: Generation (G), Understanding (U), and Editing (E). We categorize the methods into five types according to the role of LLMs: LLMs as backbone, LLMs inspired backbone, LLMs as conditioner, LLMs as agent, and LLMs as labeller.

| Task | Method | Venue | LLM Model | Domain |
|---|---|---|---|---|
| *LLMs as backbone* | | | | |
| G, U | SongComposer [98] | arXiv 2024 | SongComposer | Audio music, Speech |
| G, U | ChatMusician [97] | arXiv 2024 | Llama 2 | Symbolic music |
| G, U | AnyGPT [391] | arXiv 2024 | Llama 2 | Audio, Audio music |
| G | Boosting Large [392] | arXiv 2023 | LLaMA | Speech |
| G, U | Unified-IO 2 [393] | arXiv 2023 | Unified-IO 2 | Speech, Audio, Audio music |
| G, U | M²UGen [95] | arXiv 2023 | Llama 2 | Audio music |
| G, U | LauraGPT [91] | arXiv 2023 | - | Speech |
| U | LLaSM [96] | arXiv 2023 | Llama 2 | Speech |
| G, U | AudioPaLM [89] | arXiv 2023 | PaLM | Speech |
| U | Pengi [88] | NeurIPS 2023 | - | Speech, Audio, Audio music |
| G, U | Speechgpt [86] | EMNLP 2023 | LLaMA | Speech |
| G, U | Sparks [394] | arXiv 2023 | GPT-4 | Symbolic music |
| U | Qwen-Audio [94] | arXiv 2023 | Qwen-LM | Audio, Speech, Audio music |
| U | SALMONN [93] | arXiv 2023 | Vicuna | Audio, Speech, Audio music |
| U | Llark [92] | arXiv 2023 | Llama 2 | Audio music |
| U | MU-LLaMA [90] | arXiv 2023 | LLaMA | Audio music |
| U | Speech-LLaMA [395] | ASRU 2023 | LLaMA | Speech |
| U | LTU [87] | ICLR 2024 | LLaMA | Audio |
| U | Yu et al. [396] | ICASSP 2024 | Vicuna | Speech |
| *LLMs inspired backbone* | | | | |
| G, E | UniAudio [113] | arXiv 2023 | - | Audio, Speech, Audio music |
| G | AudioLM [112] | IEEE/ACM TASLP | - | Audio |
| G | MusicGen [50] | NeurIPS 2023 | - | Audio music |
| G | Jukebox [111] | arXiv 2020 | - | Audio music |
| G | MusicLM [49] | arXiv 2023 | - | Audio music |
| G | VALL-E [55] | arXiv 2023 | - | Speech |
| U | SICL [397] | arXiv 2023 | - | Speech |
| *LLMs as conditioner* | | | | |
| G | TANGO [100] | arXiv 2023 | FLAN-T5 | Audio |
| G | Music ControlNet [90] | ICASSP 2024 | ChatGPT | Audio music |
| U | Wu et al. [101] | ICASSP 2024 | - | Audio |
| *LLMs as agent* | | | | |
| G, E | Loop Copilot [109] | arXiv 2023 | GPT-4 | Audio music |
| G, U | MusicAgent [108] | EMNLP (Demos) 2023 | ChatGPT | Audio music, Symbolic music |
| G, U | Audiogpt [106] | AAAI 2024 | GPT-3.5 | Audio, Speech, Audio music |
| G, U | Hugginggpt [105] | NeurIPS 2023 | ChatGPT | Audio, Speech, Audio music |
| G | Wavjourney [107] | arXiv 2023 | ChatGPT | Audio, Audio music |
| G | ComposerX [398] | arXiv 2024 | GPT-4 | Symbolic music |
| *LLMs as labeller* | | | | |
| G | Audiobox [104] | arXiv 2023 | LLAMA2 7B | Audio, Speech, Audio music |
| G | Make-An-Audio 2 [102] | arXiv 2023 | GPT-3.5 | Audio |

parsing into structured pairs, both utilizing latent diffusion-based models for audio synthesis. WavJourney [107] leverages the LLM agent to integrate various audio models for producing cohesive audio content, including speech, music, and sound effects, based on textual descriptions. Audi-oLM [112] generates high-quality audio with an emphasis on long-term consistency. The method converts input audio into discrete tokens, effectively making audio generation a task akin to language modeling within a discrete space. Audiobox [104] uses LLMs for data construction, including tagging audio with high-quality, detailed captions and using LLMs to evaluate the quality of these annotations automatically. Then, Audiobox uses flow-matching techniques to produce diverse audio types with precise attribute control, from speech to music and sound effects. UniAudio [113] introduces a versatile system utilizing LLMs to generate a diverse range of audio types, including speech, sounds, music, and singing, under various input conditions. Unlike task-specific models, UniAudio tokenizes different audio types and their respective conditions into a unified se-quence, enabling next-token prediction with LLMs.

### 7.1.2 Music

Music is an art form characterized by the arrangement of sounds in time, typically including elements such as melody, harmony, rhythm, and timbre. It is produced using musical instruments and/or the human voice and is often organized according to pitch (which influences melody and harmony), rhythm (including tempo, meter, and articulation), dynamics (variations in loudness), and the sonic properties of timbre and texture. Music serves multiple functions, including aesthetic enjoyment, ceremonial purposes, and the expression of cultural identity. In the realm of music research, audio music refers to the actual recorded sound waves of music, while symbolic music pertains to the notational representation of music, such as MIDI [401] files. Each of these forms requires distinct approaches for analysis and manipulation. We explore the interconnected fields of music understanding, generation, and editing, each of which utilizes different techniques and technologies to analyze,
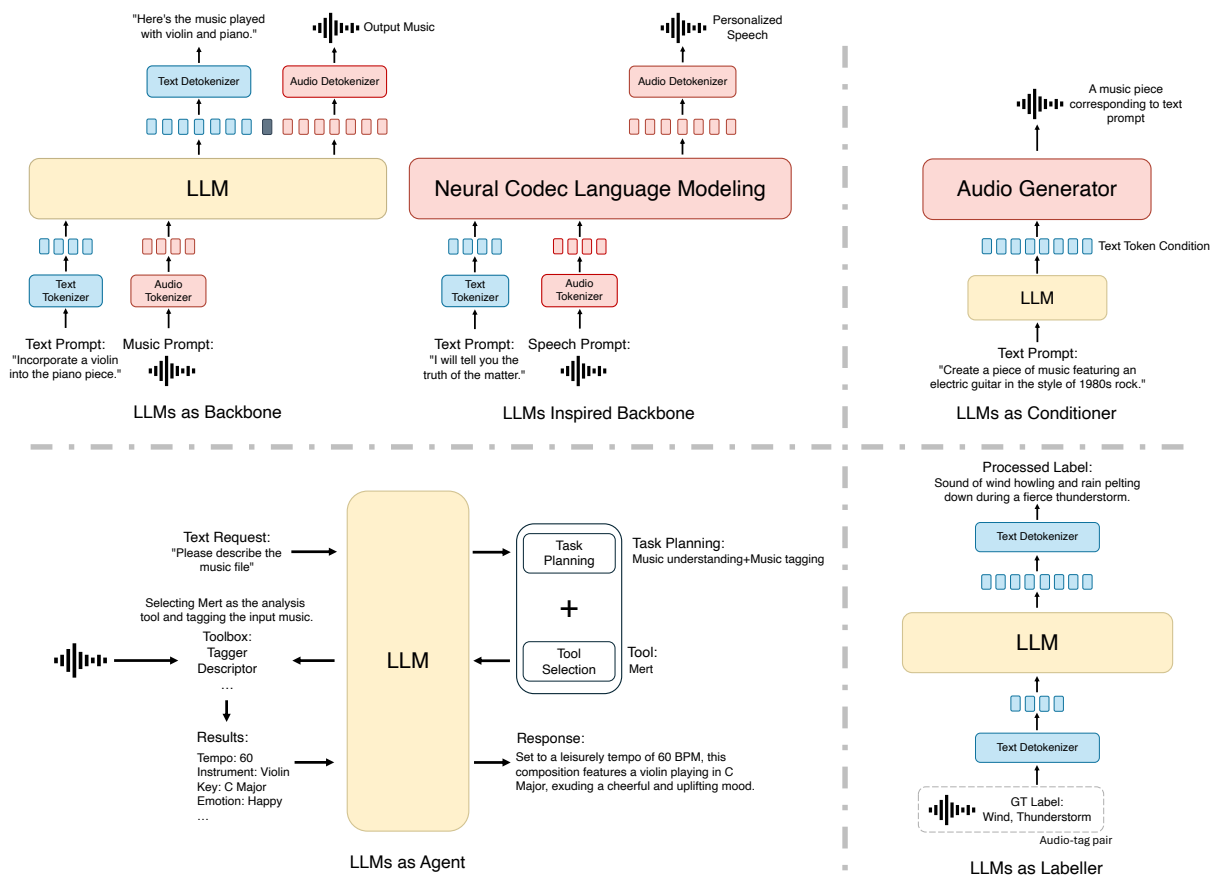
Fig. 9: Method summary of LLMs-related audio research according to different roles of LLMs. LLMs as backbone: Language pre-trained LLMs checkpoint serves as the central unit for processing text and audio tokens, either continuous or discretized. LLMs as inspiration: Different from LLMs as backbone, this method trains on randomly initialized LLMs architecture over discretized audio tokens. LLMs as Conditioner: LLMs encode text prompts into embeddings which serve as the condition of the audio generator. LLMs as agent: LLMs solve user requests by leveraging external tools. LLMs as labeller: LLMs convert class labels into audio captions.

create, and refine music, further enriching its cultural and artistic impact.

**Music Understanding.** Music understanding involves the analysis and interpretation of musical elements such as melody, harmony, rhythm, and timbre, to recognize patterns, genres, emotions, and contextual meanings within music. It includes the analysis of simple motifs to complex structures.

The field of music understanding has made significant progress with the development of models like the Music Understanding LLaMA (MU-LLaMA) [90], LLARK [92] MusicAgent [108], LyricWhiz [110], and ChatMusician [97]. These models showcase a range of methods, from analyzing detailed music features to improving lyric transcription. They highlight how LLMs can assist our understanding and interaction with music in different ways and applications.

MU-LLaMA [90] uses a pretrained MERT encoder for initial music representation, which is then integrated into the LLaMA model with an adapter. This process leverages the capabilities of LLMs to understand music by analyzing its comprehensive features. LLARK [92] improves music understanding by using a multimodal model that is fine-

tuned with instructions from refined annotations in music datasets. It combines a generative music model with a language model to analyze music in a unified way. MusicAgent [108] assists music understanding and generation through LLMs, automating tasks to meet user needs and using tools for task execution. This simplifies processes and encourages exploration in music processing. In summary, MU-LLaMA focuses on the analysis of musical features, LLARK leverages refined labels for a more general understanding, while MusicAgent emphasizes user-friendly interactions.

LyricWhiz [110] proposes a multilingual, zero-shot method for automatic lyrics transcription, performing well across various genres. It uses "Whisper" for speech recognition and "GPT-4" for context-aware annotations, acting as the transcription's "ear" and "brain". This combination greatly reduces the Word Error Rate in English and provides effective transcription in multiple languages.

**Music Generation.** AI music generation, especially with the use of LLMs, is changing the industry by creating various and complex musical pieces. Notable examples include MusicLM [49], Jukebox [111], MusicGen [50], Mu-

sic ControlNet [99], M$^2$UGen [95], ChatMusician [97], and SongComposer [98]. Specifically, these models use different techniques, from using LLMs as text embedders to employing diffusion processes and autoregressive Transformers.

MusicLM [49], Jukebox [111], and MusicGen [50] represent significant strides in text-to-music generation, each drawing inspiration from the capabilities of LLMs and employing Transformer architectures to handle complex audio tasks. MusicLM [49] treats conditional music generation as a hierarchical sequence-to-sequence task, utilizing decoder-only Transformers to create music in both the semantic and acoustic stages. Jukebox [111] addresses the challenges of long audio contexts by compressing raw audio into discrete codes using a multiscale VQ-VAE and then modeling these codes with autoregressive Transformers.

On the other hand, MusicGen [50] directly incorporates an LLM as a text embedder. It combines text tokens and melody conditions into a Transformer decoder, which then processes the inputs in an autoregressive way. The final step involves a codec model that converts the processed tokens back into music.

Music ControlNet [99] uses an LLM as a text embedder and introduces a diffusion-based model for music generation that offers precise control over dynamic and temporal aspects of audio. Inspired by image-domain ControlNet's pixel-wise control, it applies similar precision to audio through controls for melody, dynamics, and rhythm derived from training audio.

M$^2$UGen [95] introduces a multi-modal music understanding and generation framework that leverages LLMs along with pretrained models like MERT, ViT, and ViViT to analyze and create music from diverse inputs such as music, images, and videos. The decoder part utilizes AudioLDM 2 [47] and MusicGen [50] for generating music.

ChatMusician [97] and SongComposer [98] both focus on generating symbolic music but use different methods and representations. ChatMusician is an open-source LLM with intrinsic musical abilities, using continual pre-training and fine-tuning of Llama2 with ABC notation, a text-compatible music representation. It treats music as a second language, allowing it to understand and generate music using only a text tokenizer, without the need for external multi-modal neural structures. In contrast, SongComposer uses MIDI for its symbolic music representation and introduces a unique tuple design to format lyrics alongside three note attributes: pitch, duration, and rest duration. This design ensures the correct interpretation of musical symbols and precise alignment between lyrics and melodies, differentiating it from ChatMusician's approach.

Music editing involves refining and altering musical elements to improve the sound quality and artistic expression. Loop Copilot [109] combines LLMs with specialized AI music models to create a conversational interface for collaborative human-AI music loop creation. It uses a large language model to interpret user intentions and directs specialized AI models to generate and refine music through interactive dialogue. Key musical attributes are centralized and managed to ensure consistent quality throughout the creative process.

### 7.1.3 Speech

Speech specifically refers to the sounds that humans produce when they speak. It is the verbal manifestation of language, including a variety of linguistic elements such as words, sentences, tone, intonation, and rhythm. Speech is a fundamental mode of human communication and can vary greatly depending on factors like language, dialect, emotional state, and context. In the realm of artificial intelligence, LLMs have been advancing both the understanding and generation of speech, facilitating machines to interpret and replicate human-like spoken communication with increasing accuracy and naturalness.

**Speech Understanding.** Speech understanding empowers machines to interpret spoken language. This aspect of AI captures not just the words but also the speaker's intent and nuances, with progress driven by LLMs. Key contributions in this field include SpeechGPT [86], AudioPaLM [89], and other studies [395]–[397], showcasing the improved capabilities of LLMs in recognizing and processing speech across diverse contexts.

SpeechGPT [86], AudioPaLM [89], and Speech-LLaMA [395] represent pivotal developments in speech understanding, all utilizing LLMs as the structural backbone of their frameworks. SpeechGPT [86] and Speech-LLaMA [395] specifically use LLaMA as their foundation. SpeechGPT not only facilitates understanding and generating multi-modal content but also supports inter-modal knowledge transfer. It introduces SpeechInstruct, a large-scale cross-modal speech instruction dataset built upon discrete speech representations, highlighting its multi-modal capabilities. Meanwhile, Speech-LLaMA integrates speech signals with LLMs, emphasizing a mixture of auditory and linguistic data processing. Similarly, AudioPaLM [89] combines the strengths of PaLM-2 [402] and AudioLM [112] into a unified multimodal framework that performs well in both speech understanding and generation. It retains paralinguistic features such as speaker identity and intonation from AudioLM and blends them with the textual linguistic capabilities of PaLM-2, demonstrating an approach to multimodal speech processing.

Recent research in Automatic Speech Recognition (ASR) has focused on improving model accuracy with LLMs. Wang et al. [397] investigate the in-context learning abilities of the Whisper [403], which is an ASR model released by OpenAI. SICL [397] is introduced to reduce the Word Error Rates (WERs) with only a small number of labeled speech samples without gradient descent. Yu et al. [396] presents a study of structures that include fully connected layers, multi-head cross-attention, and Q-Former as connectors for integrating ASR models with LLMs.

**Speech Generation.** Speech generation, the process of converting text or other inputs like speech prompts into spoken language, has significantly evolved with the integration of LLMs. These models facilitate the naturalness and contextual relevance of generated speech, making it increasingly realistic and similar to human speech.

Inspired by the capabilities of LLMs, Wang et al. introduced VALL-E [55], a transformative approach in speech generation. VALL-E utilizes a neural codec language model that employs discrete codes from an existing neural audio codec, reframing text-to-speech (TTS) synthesis as a

conditional language modeling task rather than traditional continuous signal regression. Building on this VALL-E, Hao et al. further advanced the field with their study [392]. They conducted an investigation to improve LLMs' speech generation capabilities by integrating the pre-trained LLM frameworks, LLaMA/OPT, with the TTS model VALL-E. This research showcases a combination of language modeling and speech synthesis technology, aiming to produce more natural and effective speech outputs.

Different from previous works, using an LLM as the backbone, LauraGPT [91], developed by Wang et al., is a unified GPT model capable of handling both audio and text for recognition, understanding, and generation tasks. It performs well in a variety of functions including speech recognition, translation, text-to-speech synthesis, and more.

Another work proposed by Kakouros et al. [404] examines the potential of word surprisal, a measure of word predictability in context, to improve speech synthesis prosody.

## 7.2 Roles of LLMs

Language provides a great abstraction of our world. With the flexibility and rich descriptive power of language, researchers have unified language understanding and language generation, into a paradigm so-called generative understanding. Audio research benefits a lot from language in the LLM era, with LLMs serving as a bridge to collect and process information, the field is now able to reach a similar generative understanding stage. Broadly, we categorize the roles of LLMs into the following: LLMs as Backbone, LLMs as Conditioner, LLMs as Labeller, LLMs as Agent, and LLMs Inspired Backbone. The method summary is shown in Fig. 9.

### 7.2.1 LLMs as Backbone

Using an LLM as a backbone entails leveraging a pre-trained LLM, like LLaMA, as the central architecture of a system. These backbones, essential for the system's learning and processing abilities, are integrated with various network components and undergo fine-tuning. In the realm of multimodal LLM applications within the audio domain, the LLM backbone plays a pivotal role. It is either coupled with structures specialized for modality-specific understanding or generation, or it employs a tokenizer to convert audio into discrete tokens.

Much of the current research in LLMs typically adopts a cascade approach, which often involves the use of modality-specific encoders and/or decoders. In LTU [87], Yuan et al. suggest utilizing an Audio Spectrogram Transformer (AST) [399], a transformer encoder pre-trained with CAV-MAE [405]. This encoder's representations are aggregated and input into a LLaMA-7B backbone. These audio pre-trained representations are then paired with corresponding texts. To fine-tune the LLM backbone, LoRA adapters [232] are employed, tasked with predicting text pairs based on the audio representations. Pengi [88] follows a similar paradigm. Soham et al. call it a 'audio-text-to-text' format. Except for an CLAP [139] audio encoder, it also uses a text encoder to encode the task instructions. The audio representation, as well as the text instruction representation, are together fed into the LLM as a prefix. The LLM is then trained to predict the paired text output, e.g. a sound description.

Both LTU and Pengi show improvement in close-ended audio understanding tasks and a certain level of open-ended audio understanding tasks. Similar approaches can also be found in LLaSM [96], Mu-LLaMA [90], MusicLingo [406], Llark [92], Qwen-Audio [94]. Popular audio encoder may include CLAP [139], MERT [407], Whisper [397], AST [399].

Besides the model focused solely on understanding, there has also been researches extending into generation. Part of these researches adopt the design philosophy of the cascade approach, incorporating not just an audio encoder but also introducing an audio decoder. For instance, in M$^2$UGen [95], Atin et al. adapt both a MERT encoder and an Audioldm2 [47] / MusicGen [50] decoder. The output projection layer then maps the LLaMA2 model's output embeddings to the music decoder. A similar approach is observed in NExT-GPT [142], a recently proposed any-to-any multimodal language model.

However, the cascade approach requires training heterogeneous neural structures. In scenarios with abundant data and computational resources, these heterogeneous neural structures could lead to decreased training efficiency and lower system scalability. Recently, a unified approach has garnered the attention of researchers. This method typically necessitates the use of audio codecs [408], [409], to tokenize raw audio into discrete tokens, which are then flattened into a one-dimensional sequence for input into an LLM. This requires the LLM's vocabulary to include audio tokens, thereby necessitating an expansion of the LLM vocabulary, akin to integrating audio as a new language into the LLM. This method employs a uniform LLM structure, facilitating scalability. AudioPaLM [89], LauraGPT [91], SpeechGPT [86] follows this paradigm.

### 7.2.2 LLMs as Conditioner

In this setting, LLMs usually act as text embedders, encoding input text to condition the system's response or output, thus enabling a more nuanced and context-aware processing of audio data.

Tango [100] follows this paradigm. It comprises three main components: a text encoder, a Latent Diffusion Model (LDM), and a Mel-Spectrogram/Audio Variational Autoencoder (VAE). The text encoder is a Flan-T5, which translates the audio's input text prompt into a textual representation. This representation is then utilized to construct a latent audio representation or audio prior from standard Gaussian noise through reverse diffusion. Following this, the Mel-Spectrogram VAE's decoder generates a Mel-Spectrogram from the latent audio representation. Finally, this Mel-Spectrogram is input into a vocoder to produce the final audio output. MusicGen [50] follows a similar paradigm. The study tested both T5 and Flan-T5 models along with CLAP, and found that the T5 encoder, serving as a text conditioner, achieved the highest relevance subjective test score in relation to the text input.

### 7.2.3 LLMs as Labeller

Currently, the majority of large-scale audio datasets, such as AudioSet [383] and VGGSound [388], are annotated solely with class labels, akin to ImageNet [410]. Researchers aiming to undertake text-to-audio tasks are compelled to

transform these class labels into full-sentence audio descriptions, also known as audio captioning. A prevalent method involves utilizing LLMs to achieve this transformation.

A common pipeline for text description augmentation involves initially crafting description templates manually for a labeled audio dataset, thereby parsing audio classes into more uniformly formatted descriptions. Subsequent steps leverage self-instruction methods, employing LLMs such as ChatGPT [7] that are capable of following instructions to paraphrase these descriptions, often making use of self-instruct [411] techniques to enrich the dataset further.

### 7.2.4 LLMs as Agent

In 'LLMs as Agent', LLMs are employed to interface with various tools, orchestrating multiple functionalities to accomplish diverse tasks. This role highlights the versatility of LLMs in managing and executing complex, multi-dimensional operations.

Communicating with LLMs can be approached in various ways. A notably straightforward yet effective method is through a text interface. During the nascent phase of multi-modal LLM audio research, Huang et al. introduced AudioGPT. This system, capitalizing on advanced audio foundation models, tackles tasks such as sound detection, audio-to-text conversion, speech recognition, and speech translation. The data gleaned from these audio processes is then transformed into text, seamlessly integrating with LLM interactions. Within this framework, task analysis, model assignment, and response generation all function through textual operations. AudioGPT drew inspiration from its contemporary work, HuggingGPT, which employs a similar approach. HuggingGPT uses LLMs to invoke various models on Hugging Face, a platform hosting a diverse array of machine learning models. Similarly, MusicAgent is proposed to streamline AI-powered music processing by integrating a variety of tools and an autonomous workflow, primarily facilitated by LLMs like ChatGPT. It features a diverse toolset sourced from platforms like Hugging Face, GitHub, and various web APIs. Tasks like music classification, music separation, lyric recognition are supported.

### 7.2.5 LLMs Inspired Backbone

With the success of the next token prediction paradigm in language modeling, the audio domain has also sought to apply this approach by discretizing audio into tokens for modeling. Researchers aim to achieve emergent capabilities on audio tokens similar to those observed in LLMs, such as in-context learning and the chain of thoughts ability. Currently, modest in-context learning abilities have been confirmed to be attainable through the language modeling pretraining of audio tokens.

In VALL-E [55], researchers combined autoregressive and non-autoregressive language models to model encoded tokens. Thanks to the residual vector quantized (RVQ) modeling [408] of acoustic information, VALL-E can continue speech with only short audio and text prompts, preserving the speaker's timbre, prosody, and acoustic environment while following text constraints. In AudioLM [112], researchers discovered that unconditional training on RVQ-based acoustic tokens did not yield semantic-level consistency. Consequently, they proposed introducing semantic

tokens based on Self-Supervised Learning (SSL) representations. The representations from an SSL-pretrained teacher contain rich semantic information, and performing k-means clustering on these representations yields a k-means quantizer, allowing for the extraction of semantic tokens for the training set. Language modeling on these semantic tokens achieves better semantic consistency, enabling unconditional speech continuation to maintain semantic coherence.

## 8 TOOL-AUGMENTED MULTIMODAL AGENTS

Over the past few months, lots of works [105], [114]–[117], [412]–[430] known as tool-augmented LLMs have been emerged as a promising direction in human-computer interaction. They empower LLMs to use external tools to enhance the models' capabilities. Among them, several works [105], [114]–[118] which extend LLMs to other modalities beyond the natural language stand out. In contrast to those expert models which only focus on optimizing one specific task, *e.g.*, image generation or video generation, these pioneering works built upon LLMs can interactively generate or edit images, videos, and audio by invoking corresponding tools. In this section, we mainly focus on reviewing the recent works that aim to extend LLMs to multimodal generation by augmenting them with external tools.

### 8.1 Motivation

It is known that LLMs have limitations in accessing and processing information that is not available in their training data, such as ephemeral, changing, or private data. For example, LLMs may not be able to answer questions that involve factual knowledge updated frequently, such as the current weather or stock prices.

To overcome these limitations, several works have proposed to augment LLMs with external tools or APIs, such as retrieval-augmented generation (RAG), calculator or visual foundation models, that can provide additional information or functionality for LLMs. These tools can be invoked by LLMs through natural language instructions, and the results can be integrated into the LLMs' outputs. For example, an LLM can use a weather API to obtain the current temperature and humidity of a given location, and use them to generate a natural language response. The tool-augmented paradigm have been validated the efficacy by many works [412], [413]. In practice, Microsoft Copilot[1] augmented by various tools has been integrated into their applications including Bing, Edge and Windows operating system, which dramatically facilitates user experience. OpenAI also releases a Function Calling[2] service that can give assistants access to OpenAI-hosted tools like Code Interpreter and Knowledge Retrieval or build your own tools.

It is well known that LLMs can not generate or edit content in other modalities, such as images, videos, or audio, which is helpful for creative purposes. Motivated by tool-augmented LLMs, some pioneers develop multimodal agents that can control diverse tools across different modalities like image, video and audio. By augmenting LLMs

1. https://www.microsoft.com/en/microsoft-copilot
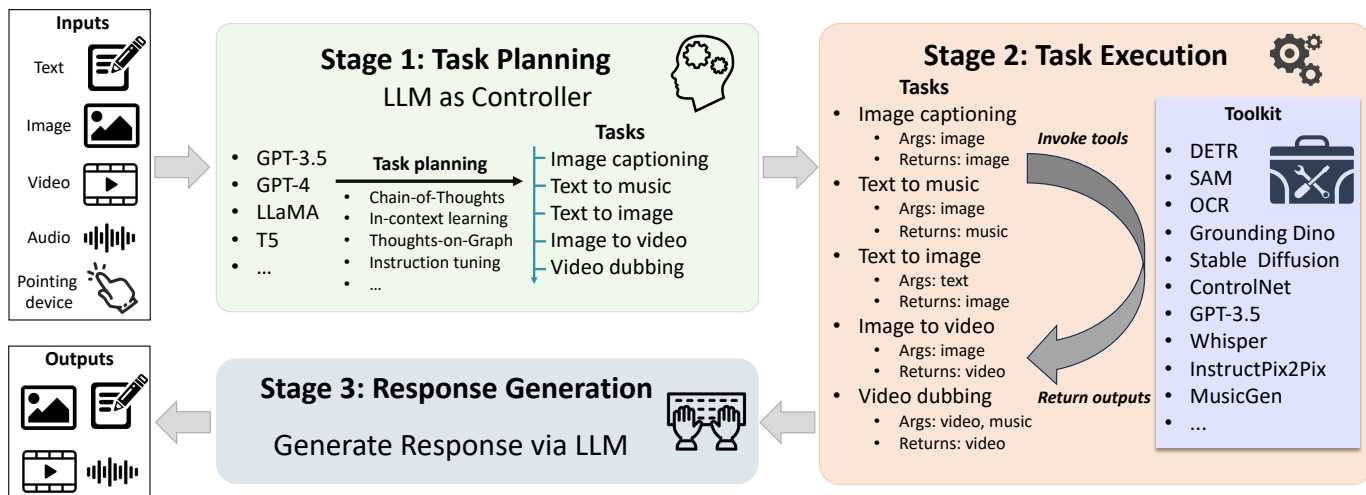2. https://platform.openai.com/docs/assistants/tools

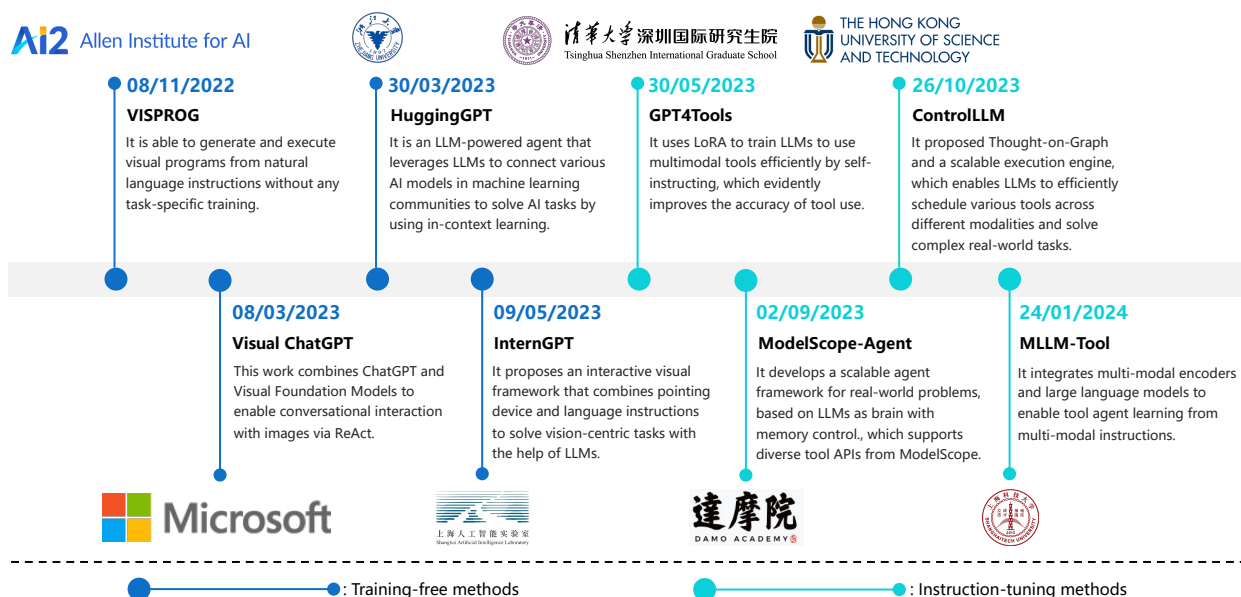Fig. 10: The pipeline of tool-augmented multimodal agents.



Fig. 11: The milestone of multimodal agents that focus on multimodal generation and editing.

with external tools, it can enable more natural and versatile human-computer interaction, as well as more powerful and creative applications.

## 8.2 Methods

As depicted in Fig. 10, the general framework of tool-augmented LLMs for multimodal interaction consists of three main stages: 1) **Task Planning** which tasks LLM as a controller to interpret the natural language instructions into a tool invocation scheme. Specifically, the core aim of this stage is to decide which tools will be used and prepare the arguments for tools. In this stage, the selected tools are organized into the tool invocation scheme that specifies the sequence of tool invocation and the inputs of tools. 2) **Task Execution** which hosts lots of external multimodal tools, *e.g.*, image generation, video editing, or audio synthesis. The tools are invoked based on the invocation scheme obtained in task planning. It is noteworthy that most of the

tools are based on deep learning models including stable diffusion [30], ControlNet [67], Blip [213], LLaVA [151], *etc.* 3) **Response Generation** which can make a user-friendly response by prompting the LLM with execution outputs from task execution. The overall system connects the LLM and external multimodal tools, which not only enhances the capabilities of LLM but also evidently improves the user experience.

The main difference among the existing works lies in how they perform task planning. To this end, the tool-augmented LLMs for multimodal interaction can be roughly divided into two categories: (1) training-free methods [105], [116], [117], [431], and (2) instruction-tuning methods [114], [115], [118], [432]. The evolution path is shown in Fig. 11. In addition, Table 11 wraps up works related to multimodal agents. Next, we will elaborate on the details of two types of methods.

### 8.2.1 Training-free Methods

Training-free methods [105], [116], [117], [431], [433], [434] primarily reply on prompt engineering, *e.g.*, ReAct [435], in-context learning [436], to boost the reasoning capabilities of LLMs. These methods exhibit their effectiveness on many tasks.

Gupta *et al.* [431] and Surís *et al.* [433] propose a code generation-based method for solving complex and compositional visual tasks given natural language instructions. These two methods mainly focus on image perception and visual question answering. Taking [431] as an example, they develop the framework of VISPROG, which consists of two main components: a program generator and a program executor. The program generator uses a large language model (GPT-3) to generate python-like modular programs from natural language instructions. The authors prompt GPT-3 with pairs of instructions and the desired high-level programs, along with a new instruction. GPT-3 then generates a program that can be executed on the input image(s) to perform the described task. Each line of the program invokes one of the 20 modules supported by VISPROG, such as object detection, segmentation, image editing, knowledge retrieval, etc. Then, the program executor executes the generated program on the input image(s) and produces the output and a visual rationale. The executor steps through the program line-by-line and invokes the correct module with the specified inputs. The modules are implemented as python classes that use off-the-shelf computer vision models, image-processing routines, or Python functions. The executor also updates the program state with the output variable name and value after each step. It is noticed that VISPROG does not have stage 3 explicitly, and directly returns the output from tools to users. This work opens up a new direction for multimodal human-computer interaction and inspires the following studies.

In contrast to VISPROG [431], [433], Visual Chat-GPT [116] does not generate code, directly. It combines a large language model (ChatGPT) with various visual foundation models (VFMs) to enable conversational interaction with both text and images. To achieve this, Wu *et al.* [116] design a Prompt Manager that bridges the gap between ChatGPT and VFMs. The Prompt Manager converts all non-language signals into language prompts that ChatGPT can understand and process. Essentially, Visual ChatGPT performs task planning via ReAct [435], which is the simplest and most straightforward way to augment LLMs with tools. Instead of generating a full solution containing all related tools once, ReAct, which extends Chain-of-Thought to tool use, immediately executes the tool after each thought. It means that the system carries out stage 1 and stage 2 alternately. In addition, to accurately call the tools, each tool is equipped with a crafted natural language prompt that instructs the LLMs how to use it. In practice, one can use a prompt like "Please use the tool named ImageBind to generate an image of a cat wearing a hat" to invoke an image generation tool, and expect the system to return an image path as the response. This approach does not require any modification or retraining of the LLMs, and can leverage the existing pre-trained LLMs, such as GPT-3 or ChatGPT. InternGPT [117] shares a similar pipeline to Visual ChatGPT, but it supports pointing devices. Thus, InternGPT offers more interesting and diverse modes of interaction, including clicking and drawing. For example, After the clicking operation is triggered, InternGPT utilizes SAM [437] to select the chosen semantic region, which can be used to remove or replace the object. Furthermore, InternGPT supports user directly in drawing the sketch based on which it can generate a new image.

In addition, HuggingGPT [105] builds upon a large language model used as the core controller to manage and organize the cooperation of expert models from machine learning communities such as Hugging Face. HuggingGPT consists of four stages: task planning, model selection, task execution, and response generation. HuggingGPT separates model selection from task planning. In the task planning stage, Shen *et al.* [105] use ChatGPT to parse the user request into a list of tasks, and determine the execution order and resource dependencies among them. In the model selection stage, they use ChatGPT to assign appropriate models to each task, based on the model descriptions available on Hugging Face. In the task execution stage, the system invokes and executes the selected models. Last, the response generation stage uses ChatGPT to integrate the predictions from all models and generate a response for the user. It is noticeable that HuggingGPT uses in-context learning in task planning and model selection. Therefore, it can perform well on some easy cases, but almost invariably fails to address the hard problems. Different from Fig. 10, HuggingGPT actually decomposes the task planning into two steps, namely, one step to parse the task and the other step to identify tools for each task. HuggingGPT also introduces some techniques to handle resource dependencies, hybrid endpoints, and prompt design in HuggingGPT.

**Limitations.** Training-free methods have some disadvantages. First, these methods rely on the availability and accessibility of the pre-trained LLMs, which are often proprietary and expensive to use. Second, these methods require manual design and tuning of the prompts, which can be time-consuming and error-prone. Last, these methods assume that the LLMs have sufficient knowledge and capability to use the tools. However, it usually fails to address complex problems. Furthermore, we found directly using off-the-shelf LLMs leads to a performance drop when extending the tool set to a large scale. This is due to the fact that LLMs normally are not trained for this purpose.

### 8.2.2 Instruction-tuning Methods

Instruction-tuning methods [412]–[414] involve training a language model to follow human instructions more accurately, which can vastly improve the capabilities of tool use for LLMs. As such, several multimodal agents [114], [115], [118] finetune an LLM in the first stage of task planning to use the tools across different modalities. In this type of method, the key is how to generate the instruction corpus to train the LLMs.

Taking GPT4Tools [115] as an example, it aims to efficiently enable LLMs to use multimodal tools, such as visual models, by self-instructing from advanced LLMs, such as LLaMA [145] and OPT [438]. The first challenge is how to construct the training corpus. Yang *et al.* [115] utilize a simple yet effective method to generate an instruction-following

TABLE 10: **The instruction samples of multimodal agents, which are used to train or assess models.**

| Image |
|---|
| 1. Can you remove the dog in the image_1.png? |
| 2. Can you create a new image that depicts a family having a picnic in the foreground? |
| 3. Generate a new image based on a segmentation map named seg.png. The new image should show a beautiful landscape. |
| 4. Using the pose map pose.png, generate a new image that features a cozy cabin in the woods with a campfire, trees, and a couple roasting marshmallows. |
| 5. Using the segmentation from the given image, generate a new image that showcases a magical forest with glowing mushrooms and fairies flying around. |
| **Video** |
| 1. Can you remove the dog in the video_1.mp4? |
| 2. Can you provide a video that is related to the given image? |
| 3. Please give me a new video derived from this file image_2.png? |
| 4. Make a video that illustrates a calm landscape of a lake surrounded by vegetation. |
| 5. I would appreciate it if you can dub the video_1.mp4 with given audio file named aud_1.wav. |
| **Audio** |
| 1. Can you create a song that visually represents the image? |
| 2. You need to generate a piece of background music for video_3.mp4. |
| 3. Please generate a piece of music by given prompt: An 80s driving pop song with heavy drums in the background. |
| 4. I want to create a piece of background music for video_4.mp4 and dub the video with this piece of music. |
| 5. translate the text into speech: "Hope is the thing with feathers That perches in the soul, And sings the tune without the words, And never stops at all". |

dataset by prompting an advanced LLM such as ChatGPT with various multimodal contexts and tool descriptions. Next, they filter out similar or invalid instructions from the raw data, resulting in 41K items. GPT4Tools augments the data by introducing negative samples (instructions that do not require tool usage) and context samples (instructions that involve multiple turns or actions). Finally, This work builds a dataset of 71.4K instruction-response pairs, which covers 31 tools for various visual tasks. We here show some examples of instructions in Table 10. After dataset construction, GPT4Tools incorporates Low-Rank Adaptation (LoRA) to finetune the open-source LLMs on the generated dataset, thus adapting them to use tools for various visual tasks, such as visual comprehension and image generation. Since the LLMs are tuned on the instruction corpus, the capabilities of tool use are dramatically improved.

Li *et al.* [118] proposes ModelScope-Agent, utilizing the models in ModelScope to augment the open-sourced LLMs. In this work, the authors also provide a tool dataset named MSAgent-Bench. Distinct from the above methods, Li *et al.* [118] design a module of tool retrieval with memory control to identify tools instead of directly prompting LLMs. Such a design makes the overall system more flexible and scalable.

Instead of directly training an LLM as a controller to generate the solution, Liu *et al.* [114] train a language model to perform tool-agnostic task decomposition and propose Thoughts-on-Graph (ToG) to generate solutions for each sub-task, which enables LLMs to use various tools across different modalities, such as text, image, audio, and video, to solve complex real-world tasks. This paper argues there are three main challenges identified in tool-augmented LLMs: a) ambiguous user prompts, b) inaccurate tool selection and parameterization, and c) inefficient tool scheduling. To this end, Liu *et al.* [114] proposes ControlLLM, a powerful framework that consists of three components: Task decomposition, which is able to parse the user input into several subtasks with specific inputs and outputs; Thoughts-on-Graph (ToG) paradigm, which finds the optimal solution path on a pre-built tool graph by depth-first search (DFS) algorithm; The execution engine is equipped with a powerful toolbox, which interprets the solution path and schedules the tools efficiently on different computational devices. ControlLLM supports lots of multimodal tools and provides a user-friendly demo interface.

Followed by ControLLM, Wang *et al.* [432] develop a multimodal tool agent system called MLLM-Tool. In contrast to prior works [114], [115], [118], this is the first work to train a large multimodal model for tool learning. In task planning, they utilize a multimodal encoder based on ImageBind [141] as well as a projection layer to extract a unified embedding space for six modalities: text, image, video, audio, speech, and music. Then, the user instruction combining the multimodal embeddings is fed into a language model to predict the corresponding API names.

**Limitations.** Currently, the instruction-training methods still have some weaknesses. On the one hand, training an LLM is forbiddingly expensive for most researchers. Despite that lots of works [232], [439], [440], such as LoRA and its variants, have been proposed to train LLMs efficiently, it still needs to make a trade-off between performance and training costs. On the other hand, it needs to generate to diverse instruction corpus for training purposes. Self-instruct [411] is an effective method that can prompt LLMs to generate more instruction automatically from seed instruction. However, it is hard to control the quality of the generated corpus, which unavoidably poses a negative impact on training LLMs. In addition, there still is an open problem to be addressed. That is, how to make an LLM that learns from a closed corpus generalize to unseen instructions. This problem is about whether the language models are able to generate novel solutions that are unseen in the training corpus, to solve more complex problems.

## 8.3 Demonstrations

Some works [105], [114], [116], [117] not only have open-sourced code, but also release online demos. As a result, in this section, we simply showcase the functionalities of some works by using their available online demos.

Currently, several multimodal agents [105], [114], [116], [117] can interact with images, either by generating, editing, or understanding them. For example, they augment LLMs with image generation or editing models, such as Stable Diffusion [30], ControlNet [67] and InstructPix2Pix [249], that can create or modify images based on text prompts. As shown in Fig. 12, we make comparisons between Visual ChatGPT, HuggingGPT, InternGPT, and ControlLLM by an example of image generation, which is beneficial to study their capabilities directly. Interestingly, HuggingGPT decomposes the instruction into two tasks and returns two generated images. The generated image from ControlLLM is more aligned to the provided instruction. We can find a slight performance gap between different methods due to the underlying models used in their systems. In contrast to those text2image models like Stable Diffusion [30], multimodal agents are able to generate vivid images in an interactive manner instead of stiffly returning images to the user. In addition, InternGPT and ControlLLM support pointing devices as input to enhance interactivity. As depicted in Fig. 13, taking ControlLLM as an example, the user can click the region of interest from the image, and the object lying in this region is then segmented into a mask by using SAM [437]. Next, the user can send an instruction to edit the image such as removing the masked object in the image. Such a manner not only can edit the image more precisely and efficiently, but also improves the rate of success for tool use.

In addition, several works [105], [114] are able to generate the video as shown in Fig. 14 and also can dub the video with the audio track. However, there is no multimodal agent available that supports direct editing on the video frame. This is partly due to the fact that video editing is extremely challenging and still needs to be further studied. Some agents also support to generate audio like speech [105], [114] and music [114]. It is really interesting to combine the tools (*e.g.*, image_captioning and text_to_music) to generate a piece of music for an image as depicted in Fig. 15. Furthermore, multimodal agents are not limited to these audio-visual tasks. For example, ControlLLM [114] supports querying weather and can even visualize the weather condition by an image as demonstrated in Fig. 16.

## 8.4 Summary

This section investigates LLMs with external tools to enhance their capabilities, particularly in multimodal interactions. The motivation behind this integration is to address LLMs' limitations in processing information not present in their training data, such as real-time or private data. This is achieved by augmenting LLMs with tools that can provide additional information or functionalities, which are invoked through natural language instructions.

The methods fall into two categories: training-free methods, which rely on prompt engineering and in-context learning, and instruction-tuning methods, which involve training LLMs to follow instructions more accurately. These methods generally involve a three-stage framework: task planning, task execution, and response generation. Task planning interprets instructions into tool invocation schemes, Task execution involves the use of multimodal tools for tasks like image generation or audio synthesis, and Response Generation creates user-friendly responses from the execution outputs.

In summary, tool-augmented multimodal agents represent a significant advancement in human-computer interaction, enabling more natural and versatile interactions and fostering creative applications across various modalities. However, they also present challenges that need to be addressed to realize their full potential.

## 9 GENERATIVE AI SAFETY

The security concerns of multimodal-generated content are drawing increasing attention. The research mainly focuses on mitigating biased and toxic content generation, safeguarding copyright, and alleviating the impact caused by fabricated content originating from generative models.

The vulnerability of generative models to attacks or malicious usage presents unique challenges and has attracted significant research attention. Recent research includes optimization-based attacks, prompt-level manipulations, and data poisoning methods: (i) Optimization-based attacks demonstrate the effectiveness of adversarial techniques to degrade model performance [442] or induce biases and harmful outputs [443]–[446]. Adversarial attack and detection-based defense research are also conducted in the field of audio and video [447], [447]–[451]. (ii) Prompt-level attacks [452]–[455] reveal the risks at the inference level, where human-crafted inputs can bypass safeguards and elicit unsafe outputs, bringing the security challenges. (iii) Data poisoning methods [456], [457] expose that models can be manipulated by injecting malicious data inputs when the integrity of training data is broken. These research works underscore the need for comprehensive approaches to enhance model robustness, secure data integrity, and recognize user's unsafe interactions, addressing the generative AI vulnerabilities.

Following the discussion on techniques to attack large generative models, there are generally two main approaches to defending against undesirable generation content. The first approach involves not modifying the existing parameters of the model but employing detection mechanisms or manipulating the input prompt context. [458] utilizes VLMs to detect and correct hate speech in multimodal memes. [459], [460] manipulate LLMs through in-context learning and defending against jailbreaking attacks. For text-to-image generation, an effective and efficient framework named Latent Guard [461] is proposed to detect the unsafe input prompt. Compared to traditional blacklist-based approaches, it is much more robust because the input prompt is checked in a latent space. And the speed is much faster than the LLM-based unsafe detection method. The second strategy enhances safety by aligning the model with human preferences or values using alignment algorithms. The Proximal Policy Optimization (PPO) algorithm, introduced in [462], has been widely used for aligning LLMs.

TABLE 11: **Multimodal agents.** We only showcase the methods that build upon LMMs to solve the user's question by invoking expert models.

| Methods | Venue | Image Editing | Image Generation | Video Editing | Video Generation | Audio Editing | Audio Generation | 3D Generation | Pointing Device |
|---|---|---|---|---|---|---|---|---|---|
| **Idea-2-3D** [441] | arXiv 2024 | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |
| **MLLM-Tool** [432] | arXiv 2024 | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| **ControlLLM** [114] | arXiv 2023 | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| **ModelScope-Agent** [118] | EMNLP 2023 | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| **GPT4Tools** [115] | NeurIPS 2023 | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **InternGPT** [117] | arXiv 2023 | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| **HuggingGPT** [105] | NeurIPS 2023 | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **Visual ChatGPT** [116] | arXiv 2023 | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **VISPROG** [431] | CVPR 2022 | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

Recently, Direct Preference Optimization(DPO) and related methods [463], [464] introduce an improvement over [462] by presenting a more efficient alignment algorithm capable of learning alignment directly from preference data. For image generative models some preference datasets [465] and alignment methods [466] have also been proposed.

Some studies focus on preventing multi-modal generative models from fabricating facts. Powerful generative models that are capable of producing highly realistic videos have attracted significant attention for their potential misuse like Deepfakes [467]–[469]. Deepfake is a technique that can generate realistic content of a certain identity in the form of an image or video. Methods [470]–[472] focus on distinguishing deepfake videos by detecting visual artifacts. While these detection methods are still limited [473] that they do not insist on adversarial attacks so they are hard to work well on new AIGC models that rely on the pattern seen in training.

Generated content can also bring up copyright issues. To solve these problems, we may need to use data attribution and embed watermarks in the generated outputs. Data origin attribution, which involves tracing the model prediction to its original training data, can help find the data causing copyright issues [474]–[476]. Watermark technique [477]–[479] can make generated content distinguishable from real content by integrating the ownership information in the generated content.

In addition, a variety of datasets also have been proposed to evaluate different aspects of generative AI safety. The SafetyBench [480] dataset is a multiple-choice question dataset for assessing unsafe contents which contains 11,435 items spanning 7 safety categories. This effort is complemented by the GOAT-Bench dataset [481] that evaluates unsafe memes with over 6K diverse topics, including implicit hate speech, gender discrimination, and cyberbullying. Furthermore, datasets like ToViLaG and others [482]–[484] have been developed specifically for visual LLMs, shedding light on the challenges of addressing the generation of toxic content, such as offensive texts and inappropriate images. These datasets offer a comprehensive evaluation for further enhancing the safety of generative models, spanning across text and image contents.

In summary, safety techniques for generative AI models can mitigate ethical risks and protect copyrights. Detection and data algorithm technologies are being utilized by commercial models. And some open-source projects also offer safety checks for users by default. Watermarking and data tracing technologies have made significant progress in alleviating copyright protection concerns. Adopting safety technologies for impactful public projects can enhance the security and trustworthiness of multimodal generation applications. Table 12 summarizes selected research related to generative AI safety issues.

TABLE 12: Overview of generative AI safety across various modalities and methods. The term "Adv." denotes "adversarial attack".

| Name | Media | Type | Method | Venue |
|---|---|---|---|---|
| Wallace et al. [442] | T | Attack | Adv. | EMNLP 2019 |
| Fu et al. [443] | T | Attack | Adv. | arXiv 2023 |
| Image hijacks [444] | I | Attack | Adv. | arXiv 2023 |
| Jones et al. [446] | T | Attack | Adv. | ICML 2023 |
| Wu et al. [452] | T + I | Attack | Prompt | arXiv 2024 |
| Xie et al. [453] | T | Attack | Prompt | NMI 2023 |
| Liu et al. [454] | T | Attack | Prompt | arXiv 2023 |
| Carlini et al. [456] | T | Attack | Data | arXiv 2023 |
| Jia et al. [457] | T | Attack | Data | EMNLP 2017 |
| Latent Guard [461] | T+I | Defense | Detection | arXiv 2023 |
| Van et al. [458] | T+I | Defense | Detection | arXiv 2023 |
| Wei et al. [459] | T | Defense | Prompt | arXiv 2023 |
| Smoothllm [460] | T | Defense | Prompt | arXiv 2023 |
| Rafailov et al. [463] | T | Defense | Alignment | arXiv 2023 |
| Raft [466] | T+I | Defense | Alignment | TMLR 2023 |
| Wodajo et al. [471] | V | Defense | Detection | arXiv 2023 |
| Safetybench [480] | T | Dataset | - | arXiv 2023 |
| GOAT-Bench [481] | T | Dataset | - | arXiv 2024 |
| ToViLaG [482] | T+I | Dataset | - | EMNLP 2023 |
| Figstep [483] | T+I | Dataset | - | arXiv 2023 |
| Liu et al. [484] | T+I | Dataset | - | arXiv 2023 |

## 10 APPLICATIONS

The rapid advancements in LLMs from companies like OpenAI, Google, Meta, Baidu, and Microsoft have led to the development of a wide range of impressive AI-powered applications. These models, such as GPT-4, Gemini, and Claude, have demonstrated remarkable capabilities in multimodal tasks, particularly in multimodal understanding.

The ability of these models to comprehend, interpret, and generate multimodal content is a significant milestone

in artificial intelligence. This multimodal capability holds immense potential for various industries and showcases the effectiveness of LLMs in multimodal generation.

In this section, we will review some remarkable applications that have already been published. Beginning with image generation and progressing to video, audio, and 3D generation, these showcases demonstrate the remarkable impact of LLMs in generating content across multiple modalities.

## 10.1 Image

The rapid advancements in diffusion models have witnessed a remarkable increase in the quality and realism of synthesized images. This has sparked the emergence of numerous companies developing high-quality text-to-image generation tools and multimodal conditional image editing or generation solutions.

Midjourney [485] is making significant strides in the industry. It enables content creation and design by offering users the ability to generate high-quality, realistic images from text prompts. Its user-friendly interfaces and robust performance make it a top choice for professionals and enthusiasts in image generation.

Moreover, Stability AI [486] has provided a powerful open-source generation model. The user community has provided various usage methods, indeed handing over creativity and tools to the users. Opening up fine-tuning has created a sizeable open-source image usage community. Even artists who are not computer scientists can easily make their small models based on their basic models. Users integrate various modal tools for deployment, enabling their image generation model to play a better role.

DALLE3 [160] stands out as a remarkable example of seamlessly integrating image generation capabilities into the powerful ChatGPT4 chatbot [487]. With DALLE3, users can generate and modify images through text-based prompts. The success of DALL-E [158] and DALL-E 2 [488] from OpenAI [489] has paved the way for highly sophisticated image generation capabilities within LLMs. These models can create detailed, photorealistic images from textual descriptions, allowing for rapid prototyping and content creation across numerous domains.

In addition to the industry-leading solutions mentioned above, lots of text-to-image generation tools have emerged that leverage LLMs to enhance the robustness and overall quality of the user experience. By leveraging LLMs to expand and refine captions, these tools can improve the quality of the generated images and the platforms' overall reliability and user-friendliness.

## 10.2 Video

With the advent of large-scale video generation models, individuals can now obtain a high-quality video clip by simply inputting a textual description. Users do not need specialized skills in traditional video production, such as CG modeling, 3D modeling, or other professional knowledge. Users can generate desired video clips prompted by a textual description and then assemble them to create a captivating short film or animated video. Existing prominent tools in this domain include commercial tools like

Pika [490] and Runway's Gen2 [491], as well as open-source video generation models such as AnimateDiff [20], VideoCrafter [19] and SVD [260]. Regarding human video generation, Heygen [492] is a popular tool widely applied in various domains, including e-commerce, social media, and advertising videos.

After some demo videos generated by Sora [1] were released, significant advancements have been made in terms of realism and prompt-following capabilities, instilling greater confidence in the application of large-scale text-to-video models. Many efforts have been made to reduce the video production cost of the film and television industry.

## 10.3 Audio

The application of multimodal AI in audio has been explored for a long time. The use cases are more well-defined, and the demand for customized and diversified sounds is more established. Technologies such as text-to-speech generation, sound transfer, music generation, and other audio generation have demonstrated promising prospects in education, video dubbing, intelligent terminals, voice assistance, and the medical field.

Microsoft's Azure platform [493] is taking a leading position in speech generation and is driving the integration of AI-generated sounds across short-video platforms. Descript [494], an AI-based audio and video editor, can transcribe speech in audio and video into text, enabling users to modify the audio and video akin to editing a Word document. Moreover, numerous video platforms, video editing software, and audio platforms have devoted significant attention to applying multimodal models and audio generation.

Besides speech and audio generation, music generation is also a hot spot in industries. The passion for music has driven countless AI researchers and scientists to dedicate immense effort to advancing this field. Suno AI [495] has ushered in the "Sora era" of music generation, where users can now create vivid, high-quality songs simply by providing a text prompt describing the desired lyrical style. Additionally, companies like Stability Audio [496], Google's MusicFX [497], Tuneflow [498], and Deepmusic [499] have also provided their music generation products, further expanding the capabilities in this domain.

## 10.4 3D

The generation of 3D models is crucial across diverse domains, including film, gaming, industrial design, architecture, interior design, product design, and virtual reality. It provides realistic visual experiences and immersive interactions, facilitating the creation of characters, scenes, products, and virtual environments to enhance creativity and engagement. Meta [500] has been heavily investing in 3D modeling and virtual reality technologies. Epic Games' MetaHuman Creator [501], a cloud-streamed app designed to elevate real-time digital human creation, is another noteworthy development that can be used in conjunction with Unreal Engine, a state-of-the-art real-time engine and editor.

As for 3D reconstruction and generation, Luma AI [502] is making significant advancements, with their technology capable of generating 3D models from 2D images, simplifying the process of creating 3D content. Other industry

players, such as Adobe [503] and Kaedim3D [504], are also making substantial strides in this field. Adobe's [503] 3D and AR tools enable the creation of immersive content, while Kaedim3D's [504] AI technology can convert 2D images into 3D models.

Wonder Studio [505] is a powerful AI tool for character replacement in videos, where it can replace the original characters in a video with user-created 3D models, unlocking exciting possibilities for personalized content creation.

Recent advancements in Language-to-Language Models (LLMs) have revealed the significant potential of text interaction and generation, opening new possibilities for creating and manipulating 3D models using natural language commands, making the process more intuitive and accessible. For instance, SceneScript [506] from Meta is able to reconstruct environments and represent the layout of physical spaces based on their powerful language-based model, Llama [145]. However, compared to the image-to-3D, the text-to-3D is still the research topic of companies like Meta [500], Google [507], Tencent [508], *etc*.

Integrating LLMs in the 3D world is transforming how we create and interact with digital content. As these technologies continue to evolve, we anticipate even more interesting and practical applications.

### 10.4.1   Others

An AI-driven software normally needs to handle various modalities of input data. This growing demand for multimodal solutions highlights the importance of advanced AI models that can seamlessly integrate and process various data types. For instance, AI-generated movies incorporate 3D technology for video, music, and speech generation, collaborating with human artists to produce high-quality cinematic experiences. Digital humans have also emerged as prominent figures across various industries, from live streaming and gaming to memorial services and large-scale interactive displays. Furthermore, LLM + multimodal generative tools have found diverse applications in mathematics, law, education, and robotics fields. In summary, we are currently witnessing the dawn of multimodal generative models with LLMs, which will undoubtedly change our lives.

## 11   FUTURE PROSPECTS

LLMs augmented multimodal generation stands out as a promising research topic, which harnesses the linguistic knowledge of LLMs to enhance the generation across various modalities, such as image, video, 3D, and audio. This series of approaches not only can improve the quality, diversity, and controllability of the generated content but also can facilitate interactivity during multimodal generation. In line with this direction, we intend to present prospects for future works.

### 11.1   Technical Prospects

In this section, we focus on the technical prospects for multimodal generation, which are expected to provide more insights and facilitate future work.

### 11.1.1   High-resolution Generation

High-resolution multimodal generation is pivotal as it directly impacts the quality and usability of generated content across various domains such as image [509], [510], video [511], [512], audio and 3D generation [44]. Accordingly, high fidelity also needs to be taken into account in audio generation [49], [409], [513], [513], [514]. The ability to produce High-resolution multimodal generation is crucial for applications requiring detailed and realistic representations, ranging from virtual reality to film production. Because it enhances the perceptual experience, provides more information for analysis, and improves the performance of subsequent tasks like object recognition and scene understanding.

LLMs hold the potential to address the challenges in a high-resolution multimodal generation. They can provide a more seamless integration of visual and textual modalities, offering a dialogue-based interface and instruction-following capabilities [70]. It could enhance the generation process by improving the understanding of complex instructions and generating more accurate and diverse outputs. Recent advancements across different modalities, including image [11], [30], [70], [515], video [19], [281], [516], [517], 3D [365], [518], [519] and audio [46], [48], [102], have led to significant improvements in the quality of generated content. We are extremely eager to witness the increase in future works that integrate LLMs, thereby offering enhanced support for high-resolution generation. Moreover, high-resolution content generation typically entails substantial hardware expenses and time costs. Consequently, the efficient generation of high-resolution content is also a topic worthy of research.

### 11.1.2   Long-term Sequence Generation

Long-term sequence generation is crucial for creating immersive experiences in the video [517], [520] and the audio [46], [112], [521], [522]. In video, it allows for the portrayal of evolving scenes and narratives, while in audio, it supports the development of music and dialogue that can adapt and flow over time. The ability to generate long sequences over time is not just a technical challenge but also a creative one, where the model must understand and predict complex patterns and progressions. It should maintain continuity, prevent repetition, and introduce novel elements that align with the overarching theme and input conditions. Only when we are able to generate long sequences for video and audio, it can potentially lead to practical significance.

Recent advancements in LLMs, such as OpenAI's GPT series and Meta's LLaMA [145], addressing the challenges of long-term sequence generation. LLMs build upon pre-trained language representations and fine-tuning techniques to capture intricate patterns and dependencies in text data, enabling them to generate coherent and contextually relevant sequences over extended lengths. By harnessing the contextual understanding and generative capabilities of LLMs, researchers can explore long-term sequence generation. For example, fine-tuning pre-trained LLMs on multimodal datasets could enable them to generate coherent and diverse sequences across different modalities, including video and audio. Additionally, techniques such as prompt

engineering and conditioning can guide the generation process toward desired outcomes, allowing for the creation of long sequences with specific themes or narratives. We argue that LLMs can enhance coherence and consistency for generated long sequence generation.

Generally, long-term sequence generation represents a complex yet compelling area of research for various domains. By leveraging the capabilities of LLMs and addressing the associated challenges, researchers can unlock new opportunities for creating immersive and engaging sequences that captivate audiences and push the boundaries of content creation and storytelling.

### 11.1.3 More Accurate and Fine-grained Generation Control

Accurate and fine-grained generation control is a significant topic in AIGC for several reasons. First, it allows for the creation of more realistic and high-quality multimodal content. This is particularly important in fields such as entertainment, advertising, and education, where high-quality content can significantly enhance user experience. Second, fine-grained control can facilitate more effective communication between humans and AI. For instance, an AI model with fine-grained control can generate a specific image or sound based on a user's detailed description, thereby improving the interaction between the user and the AI. Last, fine-grained control can also contribute to the advancement of other AI fields. For example, in reinforcement learning, an AI agent can learn more effectively if it can generate detailed and accurate simulations of its environment.

Lots of methods [19], [70], [523], [524] have been proposed to address accurate and fine-grained generation control. However, these methods still have some limitations. For instance, they still struggle with generating fine details, such as fingers or body parts, which can lead to unrealistic outputs. Moreover, they may also fail to accurately capture the nuances in the control signals, resulting in a mismatch between the generated content and the control signal.

Large language models have shown remarkable capabilities in understanding and generating text. By leveraging these capabilities, we can potentially improve the accuracy and granularity of generation control. One of prominent examples is text rendering [64], [202], [203], [525], [526] on images or videos. It has been observed that by using powerful language models, such as T5-XXL, as the encoder, the image generation models will exhibit better spelling ability. In this context, the integration of more potent LLMs into generative models is worthy of further exploration. Generally, a large language model can be trained to better understand the nuances in the control signals, thereby improving the alignment between the control signal and the generated content.

### 11.1.4 Multi-view Consistency

Multi-view consistency (MVC) is a fundamental aspect of visual generation, particularly in 3D generation, ensuring the coherence and continuity of an object's appearance from different viewpoints. This consistency is crucial for applications in augmented reality (AR), virtual reality (VR), and computer graphics, where users interact with 3D objects in a seemingly real-world context. Inconsistent appearances can break immersion and lead to a less realistic experience. The significance of multi-view consistency lies in its ability to provide a seamless and integrated perception of 3D objects, enhancing the user's experience and interaction with digital content.

MVC is particularly challenging due to the complex nature of translating 2D images into consistent 3D models, where issues such as occlusions, lighting variations, and geometric distortions can arise. Recent advancements pay lots of attention to multi-view consistency. In 3D generation, Sculpt3D [527] introduces a sparse 3D prior to improve consistency without retraining the 2D diffusion model. HarmonyView [528] addresses the balance between consistency and diversity by employing a diffusion sampling technique. Additionally, MVDream [345] lacks comprehensive multi-view knowledge or 3D-awareness during score distillation, leading to unstable generation and artifacts. In image and video generation, the works [529], [530] have contributed to the field by focusing on novel view synthesis and multi-view image generation based on large video datasets, respectively.

Despite these advancements, there are still several challenges to be further studied: 1) Limited Generalization: Many methods struggle to generalize well across diverse datasets and object categories. 2) Struggling in Complex Geometries: Accurately rendering objects with complex geometries or textureless surfaces. Since linguistic prompts can provide more prior knowledge for generation, we believe it can enhance the multi-view consistency as well as generation quality by incorporating LLMs in the pipeline.

### 11.1.5 Unified Training for Multimodal Generation

Multimodal generation is defined as the ability to simultaneously create content across different modalities, including images, videos, 3D objects, and audio. Currently, most methods [30], [46], [70], [71], [83], [112], [260], [353] only focus on one aspect such as text-to-image or text-to-video synthesis. This inevitably prompts consideration: Can a single model possess the capability to generate multiple modalities?

Several recent works [141], [279], [393], [531]–[533] have made notable strides in feature alignment for text, image, video, audio and other modalities. Some multimodal agents [105], [114]–[117] offer awesome generation capabilities for various modalities, but atom tools they used are not jointly trained. Furthermore, the pioneer works [95], [141], [142] have made their preliminary efforts to explore how to generate multimodal content in one model.

However, despite these advancements, challenges persist in achieving effective unified training for multimodal generation. One prominent obstacle lies in feature alignment across different modalities, as each modality possesses distinct statistical properties and underlying structures, necessitating robust alignment mechanisms to ensure consistency and coherence in generated outputs. Moreover, the mutual interference during training poses a significant hurdle, as optimizing for multiple modalities concurrently may lead to conflicts or competition among modal-specific objectives, hindering the overall training stability and convergence. In addition, the inherent complexity of multimodal data imposes computational overhead, necessitating efficient algorithms and scalable architectures to handle the diverse modalities efficiently.

The pursuit of unified training for multimodal generation represents a crucial advancement in AI research, offering immense potential for advancing the capabilities of generative models across diverse domains. In the future, we even look forward to models capable of generating different modalities in an interleaved manner.

### 11.1.6 Efficient Training and Deployment Strategies

Efficient training and deployment strategies also still remain to be studied in multimodal generation. As datasets and models continue to scale exponentially, the challenge of achieving efficient training and deployment becomes increasingly significant, in line with the scaling law, which posits that the computational resources required for training and deploying models grow rapidly with model size and dataset size [534]. Efficient strategies are essential not only for reducing computational costs but also for enabling real-time or resource-constrained applications of multimodal generation technologies. By minimizing computational overhead and resource utilization, efficient training and deployment strategies not only reduce time and energy costs but also enhance scalability and accessibility, democratizing access to advanced generative capabilities across diverse domains.

Several approaches have been proposed to address the challenge of efficient training in multimodal generation. Several works study low-rank approximation techniques, such as LoRA [232] and Q-LoRA [439], which aim to reduce the computational complexity of model training by approximating weight matrices with low-rank structures. Additionally, mixed precision training [535], which involves using reduced precision (e.g., 16-bit floating-point) arithmetic for certain computations, has emerged as a powerful tool for accelerating training without sacrificing model accuracy. Despite their effectiveness, these efficient training techniques still have limitations. Low-rank approximation methods may introduce approximation errors that degrade the quality of generated outputs, particularly in scenarios where high-fidelity synthesis is crucial. Similarly, mixed precision training may encounter numerical instability issues, especially when dealing with extremely large models or datasets, leading to suboptimal convergence or even training failures.

Efficient deployment strategies, such as quantization [536]–[539] to int8 or even int4 precision, offer another avenue for reducing the computational and memory requirements of multimodal generation models during inference. By quantizing model weights and activations to lower precision formats, significant savings in memory bandwidth and computational resources can be achieved, enabling faster inference and deployment on resource-constrained devices. However, quantization also presents issues, particularly in preserving model accuracy and generative quality. Lowering the precision of model parameters and activations can lead to information loss and degradation in output fidelity, especially in complex multimodal synthesis tasks where fine-grained details are crucial.

In conclusion, efficient training and deployment strategies are indispensable for realizing the full potential of multimodal generation technologies across diverse applications. By overcoming the challenges associated with scala-bility and resource constraints, researchers can accelerate the adoption of multimodal generation systems in real-world scenarios, unlocking new possibilities for content creation, human-computer interaction, and beyond.

### 11.1.7 Ethically Safe Content Generation

While there have been many works exploring how to strengthen the security of generative models of text and images [458] [459] [460], the increasing capabilities of video generation models should raise security concerns. Due to the emergence of security issues like Deepfakes [469] even with the use of previously less powerful video models, the growing strength of video models magnifies the societal impact of potential risks.

Adversarial attacks have demonstrated effective transferability from open-source models to commercial closed-source models [445]. Future commercial closed-source models should consider guarding against attacks from open-source models, such as by implementing corresponding adversarial token detection mechanisms. Simultaneously, efforts can also be considered to mitigate the impact of transferable attacks like minimizing the commercial models similarities with open-source models, such as in network architecture, data usage, and weights.

Currently, most research articles focus on ensuring security from individual perspectives such as detection [458], alignment [462], post-hoc checking [540], etc. Each of these methods generally has their own advantages and disadvantages. For instance, detection techniques offer fast checking but may overlook certain vulnerabilities. Alignment methods also cannot guarantee that the data used for training alignment covers all security cases. Additionally, post-hoc checking can be computationally expensive, especially for the generation of images and videos. There has not been much work integrating these techniques into a holistic system to ensure the security of large generative models. For example, the system can first detect user inputs, then simultaneously apply securely aligned models, and finally conduct security checks on the output to determine whether to proceed. Integrating these techniques can lead to higher efficiency and security.

## 11.2 Application Prospects

In this section, we make an effort to build blueprints for the application of multimodal generative models.

### 11.2.1 Semantic Audio Synthesis

Semantic audio synthesis involves generating of audio signals based on semantic descriptions or contextual cues, enabling the creation of immersive auditory experiences with specific characteristics or attributes. Multimodal generative models offer a promising approach to semantic audio synthesis by leveraging contextual information from other modalities, such as text or images. For instance, text-based descriptions of soundscapes or music compositions can be translated into audio waveforms using generative models trained on multimodal data. Similarly, images or videos depicting scenes or environments can inform the generation of corresponding audio accompaniments, enhancing

the realism and richness of multimedia content. By integrating semantic information across modalities, multimodal generative models enable the creation of highly personalized and contextually relevant audio experiences, spanning applications in entertainment, virtual reality, and assistive technologies.

### 11.2.2 Multi-modal Storytelling.

Multimodal storytelling involves fusing different modalities to craft compelling narratives that engage multiple senses simultaneously. This approach not only enriches storytelling experiences but also opens up new avenues for creative expression and audience engagement. In multimodal storytelling, the synthesis of content can occur in several directions.

From text prompts to image sequences, multimodal storytelling can begin with a topic, a script, or even a story outline, which serves as the basis for generating complementary modalities such as text and image sequences. For instance, given a prompt about a fantastical adventure, a multimodal generative model could generate vivid imagery depicting characters and scenes, produce an animated video sequence illustrating key events, or compose a thematic musical score to accompany the narrative.

From text prompts or images to videos and audio, in this scenario, an image serves as the starting point for generating accompanying textual descriptions, video sequences, or audio narratives. For example, given an image depicting a scenic landscape, a multimodal generative model can generate descriptive text detailing the setting, produce a video animation depicting the scene in motion, or create an immersive audio experience capturing ambient sounds and atmosphere.

Multimodal storytelling holds immense potential for enhancing traditional narrative formats and creating immersive, multi-sensory experiences that resonate with audiences across various mediums and platforms. By harnessing the capabilities of multimodal generative models, storytellers, content creators, and media producers can unlock new dimensions of creativity and engagement in the digital age.

### 11.2.3 Interactive Content Design

Interactive content design aims to create and manipulate media elements in real time, empowering users to actively participate in the creative process. Traditionally, content creation processes involve iterative steps of ideation, design, and refinement, often requiring extensive time and resources. However, with the interactive capabilities afforded by foundation generative models, creators can swiftly explore a multitude of design possibilities, rapidly iterate on concepts, and refine compositions in real time, thereby streamlining the overall creative workflow.

By enabling real-time interaction and manipulation of media elements, it can improve the efficiency of multimodal generative models. Creators can efficiently experiment with different visual and auditory elements, explore diverse artistic styles, and generate high-quality content without the need for extensive manual labor or specialized expertise. Consequently, this not only accelerates the production process but also minimizes the expenses incurred in hiring additional resources or outsourcing tasks. Moreover, the integration of multimodal generative models in Interactive Content Design contributes to the democratization of creativity by lowering the barriers to entry for aspiring artists and designers. Unlike traditional design tools that often require proficiency in complex software interfaces or artistic skills, these models offer intuitive and accessible interfaces that empower individuals from diverse backgrounds to engage in creative expression. By democratizing access to advanced content creation capabilities, these tools foster inclusivity and diversity within the creative community, enabling a broader range of voices to be heard.

Looking ahead, the evolution of multimodal generative models holds exciting prospects for the future of Interactive Content Design. As advancements continue to expand the scope and fidelity of generated content across different modalities, we can anticipate even greater opportunities for innovation in areas such as virtual reality, augmented reality, and immersive storytelling. Additionally, ongoing research efforts aimed at enhancing the interpretability, controllability, and scalability of these models will further fuel their adoption in diverse creative domains, paving the way for transformative changes in how we conceive, design, and interact with digital content.

### 11.2.4 3D Scene Generation

3D scene generation refers to the creation of immersive and realistic environments in virtual worlds, games, simulations, and architectural visualization. This application domain leverages multimodal generative models to synthesize complex 3D scenes comprising objects, textures, lighting, and spatial arrangements. The ability to generate 3D scenes has profound implications for various industries, including entertainment, education, design, and virtual reality.

In the context of games and virtual environments, multimodal generative models can automate the process of scene creation, reducing the reliance on manual modeling and asset creation. By inputting textual descriptions or conceptual sketches, developers can generate entire 3D environments populated with interactive objects, characters, and atmospheric effects. This not only accelerates the game development pipeline but also enables the creation of dynamic and immersive gameplay experiences. Moreover, in architectural visualization and design, multimodal generative models can assist architects, urban planners, and designers in visualizing and exploring different design options. By inputting architectural blueprints or design parameters, designers can generate realistic 3D renderings of buildings, landscapes, and interior spaces, allowing for rapid iteration and exploration of design concepts. This facilitates collaboration, decision-making, and communication among stakeholders involved in the design process.

By harnessing the capabilities of multimodal generative models, 3D scene generation might revolutionize how virtual environments are created, experienced, and interacted with. Whether in games, simulations, or architectural visualization, the ability to generate immersive and realistic 3D scenes programmatically opens up new possibilities for creativity, exploration, and storytelling in virtual worlds.

### 11.2.5 Customizable Avatars

Customizable avatars represent digital representations of users that can be personalized and adapted to reflect individual preferences, identities, and characteristics.

Multimodal generative models offer a compelling approach to customizable avatar creation by synthesizing diverse media types such as images, text, and audio to create lifelike and expressive avatars. For example, generative models trained on multimodal data can generate photorealistic images of avatars based on textual descriptions or user preferences, incorporating details such as facial features, clothing styles, and expressions. Similarly, audio-based avatars can be generated using voice synthesis techniques, enabling avatars to communicate with users using natural-sounding voices that reflect their personalities or preferences. By enabling the creation of customizable avatars across multiple modalities, multimodal generative models empower users to express themselves in virtual environments, fostering deeper engagement and personalization in social interactions, gaming, and virtual communication platforms.

Currently, there are several aspects that could be further studied: 1) Personalization and Customization: Multimodal generative models can generate avatars that closely resemble users based on input parameters such as facial features, body type, and clothing preferences. Users can interactively customize their avatars using intuitive interfaces, adjusting attributes such as hairstyle, facial expression, and accessories in real time. 2) Emotional Expression and Limb Movements: Avatars generated by multimodal models can exhibit a wide range of emotional expressions, gestures and physical movements, enhancing their ability to convey non-verbal communication cues in virtual environments. Users can dynamically control their avatar's behaviors, allowing for more immersive social interactions and collaborative experiences in virtual worlds. 3) Integration with Virtual Environments: Customizable avatars can be seamlessly integrated into various virtual environments, including social platforms, online games, and virtual reality applications. Users can navigate these environments using their avatars, interacting with other users and objects in real time, fostering a sense of presence and belonging in digital spaces.

### 11.3 Towards the World Model

World models [541]–[545] have recently emerged as a hotspot topic. Many renowned researchers have expressed that world models will come true in the foreseeable future, and researchers around the world hold great expectations for this development. We found that all the topics mentioned in the survey correspond precisely to the principal components of the world modeling, encompassing perceptual modalities such as vision, audition, and speech, as well as spatial understanding and generation. Once world models advance to a usable stage, they will endow numerous industries with new possibilities. We here highlight several core applications for reference.

**Multimodal education and communication.** World models hold immense promise for revolutionizing education and communication by facilitating multimodal learning experiences and immersive interactions. By integrating diverse sensory modalities such as text, images, audio, and video, these models enable the creation of rich educational content that caters to different learning styles and preferences. Furthermore, they empower learners to engage with complex concepts and environments more intuitively and interactively, thereby enhancing comprehension and retention. Additionally, world models facilitate seamless communication by enabling the synthesis of natural and expressive multimodal dialogue, fostering more engaging and personalized interactions in virtual learning environments and online collaboration platforms.

**Movie Generation.** The application of world models in movie generation represents a paradigm shift in filmmaking, offering filmmakers unprecedented creative freedom and flexibility. By leveraging multimodal generative techniques, filmmakers can seamlessly integrate various elements such as dialogue, visuals, sound effects, and music to craft immersive cinematic experiences that resonate with audiences on a deeper level. Moreover, world models enable the generation of dynamic and personalized narratives tailored to individual viewer preferences, thereby enhancing viewer engagement and immersion. Furthermore, these models facilitate the exploration of alternative storytelling formats and experimental filmmaking techniques, pushing the boundaries of cinematic creativity and expression.

**Metaverse.** The emergence of the metaverse presents exciting opportunities for leveraging world models to create immersive and interactive virtual worlds. By synthesizing multimodal sensory experiences, including visual, auditory, and haptic feedback, these models enable the creation of highly realistic and immersive virtual environments that blur the boundaries between physical and digital reality. Moreover, world models facilitate the development of intelligent virtual agents and NPCs that exhibit lifelike behaviors and interactions, enhancing the sense of presence and social immersion within the metaverse. Additionally, these models empower users to customize and personalize their virtual experiences, fostering creativity and exploration within digital worlds.

## 12 CONCLUSION

In this survey, we systematically review multimodal editing and generation works augmented by LLMs, delving into the advancements across various modalities including image, video, 3D and audio. Following that, multimodal agents that integrate a plethora of state-of-the-art generative models are thoroughly discussed with diverse case studies. We also investigate the concern for the safety of multimodal generative models, which play an indispensable role in practical applications. Our comprehensive review highlights the significant contributions of LLMs in enhancing the quality and capabilities of generative systems. Looking forward, we anticipate further innovations at the intersection of AI and generative content, driving progress toward a more unified and capable multimodal generation framework. In conclusion, we ardently expect that our investigation will provide insights and inspiration for the development of the multimodal generation, particularly the world models, which have garnered the attention and anticipation of the majority of researchers.

🧑: Generate a figure of a yellow dog sitting on the grass and a cute cat running in the garden
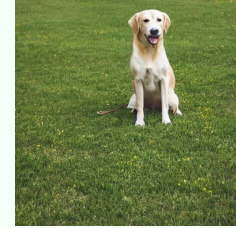
### HuggingGPT



🤖: I'm pleased to inform you that I've successfully generated the images based on your request.
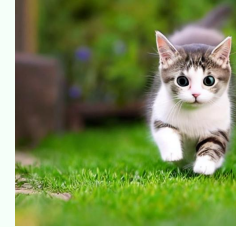
Firstly, I created an image of a yellow dog sitting on the grass. The image was generated using the 'runwayml/stable-diffusion-v1-5' model. You can view the image at this location: '/images/d9ba.jpg'.

Next, I generated an image of a cute cat running in the garden, also using the 'runwayml/stable-diffusion-v1-5' model. The image can be found at this location: 'images/36b4.jpg'.

Please note that these images are AI-generated and may not perfectly represent the descriptions. However, I've done my best to ensure they match your request as closely as possible. If you have any other requests or need further assistance, feel free to ask!

images/d9ba.jpg

images/36b4.jpg

### InternGPT

🤖: I have generated an image based on your description of a yellow dog sitting on the grass and a cute cat running in the garden. Here is the image: [image/121851_5aea79_update.png].

image/121851_5aea79_update.png

### VisualChatGPT

🤖: I have generated an image based on your description of a yellow dog sitting on the grass and a cute cat running in the garden. Here is the image: image/b241d32b.png.

/images/ b241d32b.png

### ControlLLM

🤖: I have generated an image for you. The image depicts a yellow dog sitting on the grass and a cute cat running in the garden. The image file name is 1dd44a_image.png.

Here is a summary of the workflow:

1.The text-to-image tool was used to generate the requested image.

2.The tool took the text description "The new image depicts a yellow dog sitting on the grass and a cute cat running in the garden."

3.The tool outputted the image file named 1dd44a_image.png.

Please let me know if there is anything else I can help you with.

1dd44a_image.png

Fig. 12: We show some case studies on image generation for different methods.

👩‍💼: remove the masked object.

**ControlLLM**

🤖: I have removed the masked object from the image. The result can be found in the file 34b15a_image.png.
To remove the masked object, the workflow consists of using the image_inpainting tool. The input image is b30ff5_image.png and the mask is 67a4f9_mask.png. The output image after removing the masked object is 34b15a_image.png.
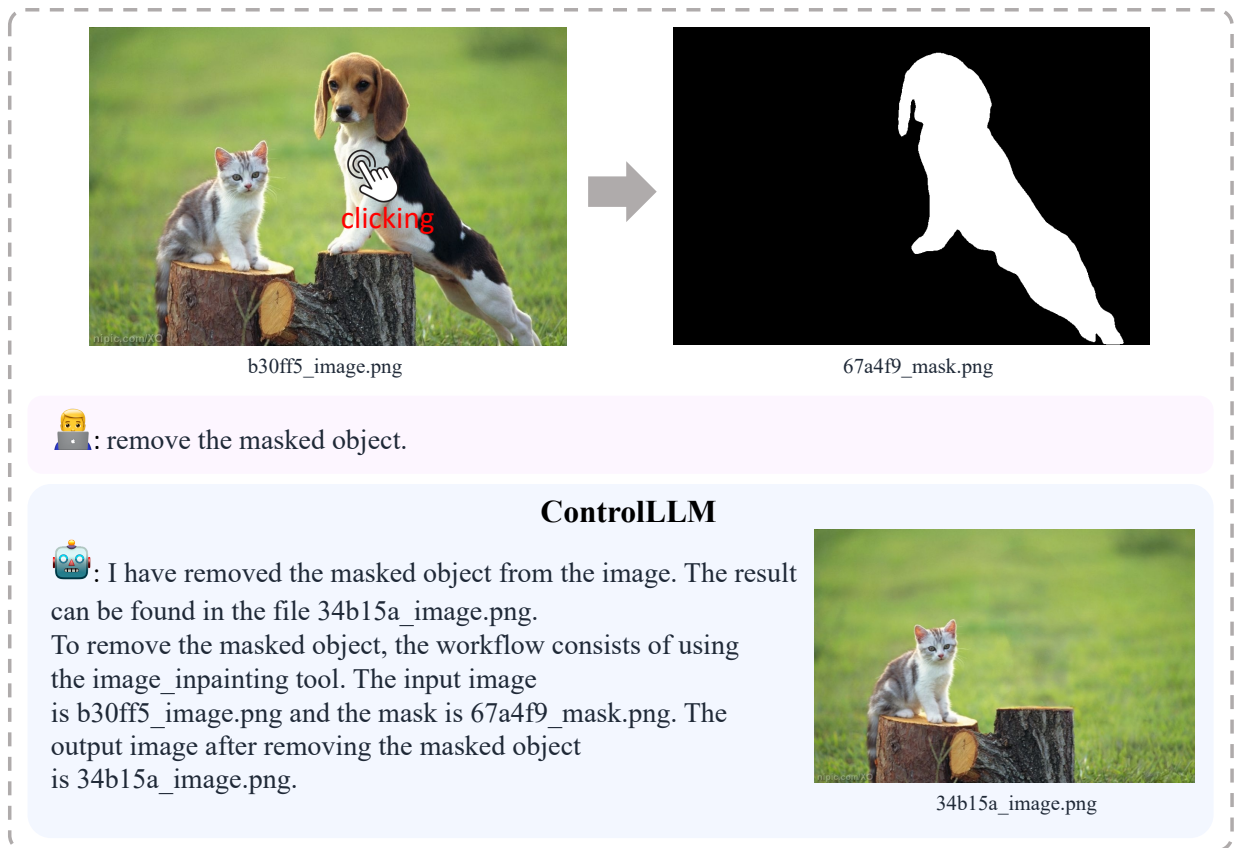
34b15a_image.png

Fig. 13: Interactive image editing by clicking.

👩‍💼: based on the /examples/a.jpg, please generate an audio

**HuggingGPT**

🤖: Sure, I understand your request. Here's the response based on the inference results:
Firstly, I converted the image you provided into text. The image was interpreted as "a cat sitting on a window sill looking out". This interpretation was made using the model 'ydshieh/vit-gpt2-coco-en'.
Then, I converted this text into an audio file. The audio file was successfully generated and can be found at the following location: '/audios/bace.wav'. This conversion was done using the model 'espnet/kan-bayashi_ljspeech_vits'.
So, in summary, I have successfully converted the image into text and then the text into an audio file. The audio file can be accessed at '/audios/bace.wav'.

/images/a.png

00:00          -03:46

⏮ ▶ ⏭

/audios/bace.wav

Fig. 14: An example of generating audio from the given image.

👤: Generate a new image that shows an insect robot preparing a delicious meal. Then give me a video based on new image.

### ControlLLM

🤖: I have generated a new image that shows an insect robot preparing a delicious meal. The image file is called `a784c4_image.png`.
Based on this image, I have also created a video. The video file is named `4115c8_video.mp4`.
Please let me know if there is anything else I can help you with!



a784c4_image.png

4115c8_video.mp4

Fig. 15: This example shows the multimodal agent [114] generate multiple outputs from the instruction. It first generates an image conditioned on the text prompt and then creates a video from the generated image.

👤: What's the weather situation in Beijing? Can you generate a new image that represents the weather in there?

### ControlLLM

🤖: The current weather condition in Beijing, China is overcast with a temperature of 7.0 degrees Celsius. The humidity is 22.5%, and there is no precipitation or snow expected. The wind speed is 3.6 kilometers per hour, coming from the southeast direction. The pressure is 1029.0 millibars, and visibility is 10.0 kilometers. The cloud cover is 100%, and the solar radiation is 137.0 W/m^2. The UV index is 1.0. The sunrise is at 06:50:37, and the sunset is at 18:04:18.

Please note that the severe risk for all days is 10.0.

I have also generated an image that represents the weather in Beijing. Here is the file name: 1481b0_image.png.



1481b0_image.png

Fig. 16: An example of multimodal generation that visualizes the weather condition.

# REFERENCES

[1] OpenAI, "Video generation models as world simulators," OpenAI, Tech. Rep., 2024. [Online]. Available: https://openai.com/research/video-generation-models-as-world-simulators

[2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[3] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," 2018.

[4] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[5] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

[6] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[7] OpenAI, "Chatgpt: A language model for conversational ai," OpenAI, Tech. Rep., 2023. [Online]. Available: https://www.openai.com/research/chatgpt

[8] Y. Li, C. Wang, and J. Jia, "Llama-vid: An image is worth 2 tokens in large language models," *arXiv preprint arXiv:2311.17043*, 2023.

[9] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou, W. Zhang, P. Lu, C. He, X. Yue *et al.*, "Llama-adapter v2: Parameter-efficient visual instruction model," *arXiv preprint arXiv:2304.15010*, 2023.

[10] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023.

[11] J. Chen, J. Yu, C. Ge, L. Yao, E. Xie, Y. Wu, Z. Wang, J. Kwok, P. Luo, H. Lu *et al.*, "Pixart-α: Fast training of diffusion transformer for photorealistic text-to-image synthesis," *arXiv preprint arXiv:2310.00426*, 2023.

[12] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video diffusion models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 8633–8646, 2022.

[13] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, "Frozen in time: A joint video and image encoder for end-to-end retrieval," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 1708–1718, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:232478955

[14] Y. He, T. Yang, Y. Zhang, Y. Shan, and Q. Chen, "Latent video diffusion models for high-fidelity long video generation," *arXiv preprint arXiv:2211.13221*, 2022.

[15] D. Zhou, W. Wang, H. Yan, W. Lv, Y. Zhu, and J. Feng, "Magicvideo: Efficient video generation with latent diffusion models," *arXiv preprint arXiv:2211.11018*, 2022.

[16] 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," *arXiv preprint arXiv:2209.14792*, 2022.

[17] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet *et al.*, "Imagen video: High definition video generation with diffusion models," *arXiv preprint arXiv:2210.02303*, 2022.

[18] R. Villegas, M. Babaeizadeh, P.-J. Kindermans, H. Moraldo, H. Zhang, M. T. Saffar, S. Castro, J. Kunze, and D. Erhan, "Phenaki: Variable length video generation from open domain textual descriptions," in *International Conference on Learning Representations*, 2022.

[19] H. Chen, M. Xia, Y. He, Y. Zhang, X. Cun, S. Yang, J. Xing, Y. Liu, Q. Chen, X. Wang *et al.*, "Videocrafter1: Open diffusion models for high-quality video generation," *arXiv preprint arXiv:2310.19512*, 2023.

[20] Y. Guo, C. Yang, A. Rao, Y. Wang, Y. Qiao, D. Lin, and B. Dai, "Animatediff: Animate your personalized text-to-image diffusion models without specific tuning," *arXiv preprint arXiv:2307.04725*, 2023.

[21] O. Bar-Tal, H. Chefer, O. Tov, C. Herrmann, R. Paiss, S. Zada, A. Ephrat, J. Hur, Y. Li, T. Michaeli *et al.*, "Lumiere: A space-time diffusion model for video generation," *arXiv preprint arXiv:2401.12945*, 2024.

[22] R. Girdhar, M. Singh, A. Brown, Q. Duval, S. Azadi, S. S. Rambhatla, A. Shah, X. Yin, D. Parikh, and I. Misra, "Emu video: Factorizing text-to-video generation by explicit image conditioning," *arXiv preprint arXiv:2311.10709*, 2023.

[23] H. Chen, Y. Zhang, X. Cun, M. Xia, X. Wang, C. Weng, and Y. Shan, "Videocrafter2: Overcoming data limitations for high-quality video diffusion models," *arXiv preprint arXiv:2401.09047*, 2024.

[24] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.

[25] A. Sanghi, H. Chu, J. G. Lambourne, Y. Wang, C.-Y. Cheng, M. Fumero, and K. R. Malekshan, "Clip-forge: Towards zero-shot text-to-shape generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 603–18 613.

[26] N. Mohammad Khalid, T. Xie, E. Belilovsky, and T. Popa, "Clip-mesh: Generating textured meshes from text using pretrained image-text models," in *SIGGRAPH Asia 2022 conference papers*, 2022, pp. 1–8.

[27] O. Michel, R. Bar-On, R. Liu, S. Benaim, and R. Hanocka, "Text2mesh: Text-driven neural stylization for meshes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 492–13 502.

[28] C. Wang, R. Jiang, M. Chai, M. He, D. Chen, and J. Liao, "Nerf-art: Text-driven neural radiance fields stylization," *IEEE Transactions on Visualization and Computer Graphics*, 2023.

[29] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, July 2023. [Online]. Available: https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/

[30] 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.

[31] T. Yi, J. Fang, J. Wang, G. Wu, L. Xie, X. Zhang, W. Liu, Q. Tian, and X. Wang, "Gaussiandreamer: Fast generation from text to 3d gaussians by bridging 2d and 3d diffusion models," *arXiv preprint arXiv:2310.08529*, 2023.

[32] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng, "Dreamgaussian: Generative gaussian splatting for efficient 3d content creation," *arXiv preprint arXiv:2309.16653*, 2023.

[33] L. Höllein, A. Cao, A. Owens, J. Johnson, and M. Nießner, "Text2room: Extracting textured 3d meshes from 2d text-to-image models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 7909–7920.

[34] Y. Liang, X. Yang, J. Lin, H. Li, X. Xu, and Y. Chen, "Lucid-Dreamer: Towards high-fidelity text-to-3d generation via interval score matching," *https://arxiv.org/abs/2311.11284*, 2023.

[35] X. Yu, Y.-C. Guo, Y. Li, D. Liang, S.-H. Zhang, and X. Qi, "Text-to-3d with classifier score distillation," *https://arxiv.org/abs/2310.19415*, 2023.

[36] W. Li, R. Chen, X. Chen, and P. Tan, "SweetDreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d," *https://arxiv.org/abs/2310.02596*, 2023.

[37] Z. Wang, C. Lu, Y. Wang, F. Bao, C. Li, H. Su, and J. Zhu, "Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation," *arXiv preprint arXiv:2305.16213*, 2023.

[38] J. Lorraine, K. Xie, X. Zeng, C.-H. Lin, T. Takikawa, N. Sharp, T.-Y. Lin, M.-Y. Liu, S. Fidler, and J. Lucas, "Att3d: Amortized text-to-3d object synthesis," in *International Conference on Computer Vision ICCV*, 2023.

[39] J. Xu, X. Wang, W. Cheng, Y.-P. Cao, Y. Shan, X. Qie, and S. Gao, "Dream3D: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models," *https://arxiv.org/abs/2212.14704*, 2023.

[40] J. Zhu and P. Zhuang, "HiFA: High-fidelity text-to-3d with advanced diffusion guidance," *https://arxiv.org/abs/2305.18766*, 2023.

[41] R. Chen, Y. Chen, N. Jiao, and K. Jia, "Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023.

[42] C. Tsalicoglou, F. Manhardt, A. Tonioni, M. Niemeyer, and

F. Tombari, "Textmesh: Generation of realistic 3d meshes from text prompts," *arXiv preprint arXiv:2304.12439*, 2023.

[43] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022.

[44] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, and T.-Y. Lin, "Magic3d: High-resolution text-to-3d content creation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 300–309.

[45] J. Seo, W. Jang, M.-S. Kwak, J. Ko, H. Kim, J. Kim, J.-H. Kim, J. Lee, and S. Kim, "Let 2d diffusion model know 3d-consistency for robust text-to-3d generation," *arXiv preprint arXiv:2303.07937*, 2023.

[46] 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," *arXiv preprint arXiv:2301.12503*, 2023.

[47] H. Liu, Q. Tian, Y. Yuan, X. Liu, X. Mei, Q. Kong, Y. Wang, W. Wang, Y. Wang, and M. D. Plumbley, "Audioldm 2: Learning holistic audio generation with self-supervised pretraining," *arXiv preprint arXiv:2308.05734*, 2023.

[48] F. Kreuk, G. Synnaeve, A. Polyak, U. Singer, A. Défossez, J. Copet, D. Parikh, Y. Taigman, and Y. Adi, "Audiogen: Textually guided audio generation," *arXiv preprint arXiv:2209.15352*, 2022.

[49] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzetti, A. Caillon, Q. Huang, A. Jansen, A. Roberts, M. Tagliasacchi *et al.*, "Musiclm: Generating music from text," *arXiv preprint arXiv:2301.11325*, 2023.

[50] J. Copet, F. Kreuk, I. Gat, T. Remez, D. Kant, G. Synnaeve, Y. Adi, and A. Défossez, "Simple and controllable music generation," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[51] S. Forsgren and H. Martiros, "Riffusion-stable diffusion for real-time music generation, 2022," *URL https://riffusion. com/about*, vol. 6, 2022.

[52] 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," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[53] K. Shen, Z. Ju, X. Tan, Y. Liu, Y. Leng, L. He, T. Qin, S. Zhao, and J. Bian, "Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers," *arXiv preprint arXiv:2304.09116*, 2023.

[54] Z. Ju, Y. Wang, K. Shen, X. Tan, D. Xin, D. Yang, Y. Liu, Y. Leng, K. Song, S. Tang *et al.*, "Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models," *arXiv preprint arXiv:2403.03100*, 2024.

[55] C. Wang, S. Chen, Y. Wu, Z. Zhang, L. Zhou, S. Liu, Z. Chen, Y. Liu, H. Wang, J. Li *et al.*, "Neural codec language models are zero-shot text to speech synthesizers," *arXiv preprint arXiv:2301.02111*, 2023.

[56] Z. Jiang, J. Liu, Y. Ren, J. He, C. Zhang, Z. Ye, P. Wei, C. Wang, X. Yin, Z. Ma *et al.*, "Mega-tts 2: Zero-shot text-to-speech with arbitrary length speech prompts," *arXiv preprint arXiv:2307.07218*, 2023.

[57] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, "Fastspeech 2: Fast and high-quality end-to-end text to speech," *arXiv preprint arXiv:2006.04558*, 2020.

[58] Y. Ge, Y. Ge, Z. Zeng, X. Wang, and Y. Shan, "Planting a seed of vision in large language model," *arXiv preprint arXiv:2307.08041*, 2023.

[59] L. Zeqiang, Z. Xizhou, D. Jifeng, Q. Yu, and W. Wenhai, "Mini-dalle3: Interactive text to image by prompting large language models," *arXiv preprint arXiv:2310.07653*, 2023.

[60] Z. Tang, Z. Yang, M. Khademi, Y. Liu, C. Zhu, and M. Bansal, "Codi-2: In-context, interleaved, and interactive any-to-any generation," *arXiv preprint arXiv:2311.18775*, 2023.

[61] Y. Ge, S. Zhao, Z. Zeng, Y. Ge, C. Li, X. Wang, and Y. Shan, "Making llama see and draw with seed tokenizer," *arXiv preprint arXiv:2310.01218*, 2023.

[62] Q. Sun, Y. Cui, X. Zhang, F. Zhang, Q. Yu, Z. Luo, Y. Wang, Y. Rao, J. Liu, T. Huang *et al.*, "Generative multimodal models are in-context learners," *arXiv preprint arXiv:2312.13286*, 2023.

[63] X. Zhao, B. Liu, Q. Liu, G. Shi, and X.-M. Wu, "Making multimodal generation easier: When diffusion models meet llms," *arXiv preprint arXiv:2310.08949*, 2023.

[64] J. Chen, Y. Huang, T. Lv, L. Cui, Q. Chen, and F. Wei, "Textdiffuser-2: Unleashing the power of language models for text rendering," *arXiv preprint arXiv:2311.16465*, 2023.

[65] L. Lian, B. Li, A. Yala, and T. Darrell, "Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models," *arXiv preprint arXiv:2305.13655*, 2023.

[66] W. Feng, W. Zhu, T.-j. Fu, V. Jampani, A. Akula, X. He, S. Basu, X. E. Wang, and W. Y. Wang, "Layoutgpt: Compositional visual planning and generation with large language models," *arXiv preprint arXiv:2305.15393*, 2023.

[67] T. Zhang, Y. Zhang, V. Vineet, N. Joshi, and X. Wang, "Controllable text-to-image generation with gpt-4," *arXiv preprint arXiv:2305.18583*, 2023.

[68] L. Qu, S. Wu, H. Fei, L. Nie, and T.-S. Chua, "Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 643–654.

[69] Y. Li, H. Liu, Q. Wu, F. Mu, J. Yang, J. Gao, C. Li, and Y. J. Lee, "Gligen: Open-set grounded text-to-image generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 511–22 521.

[70] J. Betker, G. Goh, L. Jing, T. Brooks, J. Wang, L. Li, L. Ouyang, J. Zhuang, J. Lee, Y. Guo *et al.*, "Improving image generation with better captions," *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2023.

[71] D. Kondratyuk, L. Yu, X. Gu, J. Lezama, J. Huang, R. Hornung, H. Adam, H. Akbari, Y. Alon, V. Birodkar *et al.*, "Videopoet: A large language model for zero-shot video generation," *arXiv preprint arXiv:2312.14125*, 2023.

[72] L. Yu, J. Lezama, N. B. Gundavarapu, L. Versari, K. Sohn, D. Minnen, Y. Cheng, A. Gupta, X. Gu, A. G. Hauptmann *et al.*, "Language model beats diffusion–tokenizer is key to visual generation," *arXiv preprint arXiv:2310.05737*, 2023.

[73] H. Fei, S. Wu, W. Ji, H. Zhang, and T.-S. Chua, "Empowering dynamics-aware text-to-video diffusion with large language models," *arXiv preprint arXiv:2308.13812*, 2023.

[74] H. Lin, A. Zala, J. Cho, and M. Bansal, "Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning," *arXiv preprint arXiv:2309.15091*, 2023.

[75] J. Lv, Y. Huang, M. Yan, J. Huang, J. Liu, Y. Liu, Y. Wen, X. Chen, and S. Chen, "Gpt4motion: Scripting physical motions in text-to-video generation via blender-oriented gpt planning," *arXiv preprint arXiv:2311.12631*, 2023.

[76] Y. Lu, L. Zhu, H. Fan, and Y. Yang, "Flowzero: Zero-shot text-to-video synthesis with llm-driven dynamic scene syntax," *arXiv preprint arXiv:2311.15813*, 2023.

[77] S. Hong, J. Seo, S. Hong, H. Shin, and S. Kim, "Large language models are frame-level directors for zero-shot text-to-video generation," *arXiv preprint arXiv:2305.14330*, 2023.

[78] H. Huang, Y. Feng, C. Shi, L. Xu, J. Yu, and S. Yang, "Freebloom: Zero-shot text-to-video generator with llm director and ldm animator," *arXiv preprint arXiv:2309.14494*, 2023.

[79] Z. Wang, J. Wang, D. Lin, and B. Dai, "Intercontrol: Generate human motion interactions by controlling every joint," *arXiv preprint arXiv:2311.15864*, 2023.

[80] J. Liu, W. Dai, C. Wang, Y. Cheng, Y. Tang, and X. Tong, "Plan, posture and go: Towards open-world text-to-motion generation," *arXiv preprint arXiv:2312.14828*, 2023.

[81] F. Long, Z. Qiu, T. Yao, and T. Mei, "Videodrafter: Content-consistent multi-scene video generation with llm," *arXiv preprint arXiv:2401.01256*, 2024.

[82] C. Sun, J. Han, W. Deng, X. Wang, Z. Qin, and S. Gould, "3d-gpt: Procedural 3d modeling with large language models," *arXiv preprint arXiv:2310.12945*, 2023.

[83] Y. Feng, J. Lin, S. K. Dwivedi, Y. Sun, P. Patel, and M. J. Black, "Posegpt: Chatting about 3d human pose," *arXiv preprint arXiv:2311.18836*, 2023.

[84] S. Chen, X. Chen, C. Zhang, M. Li, G. Yu, H. Fei, H. Zhu, J. Fan, and T. Chen, "Ll3da: Visual interactive instruction tuning for omni-3d understanding, reasoning, and planning," *arXiv preprint arXiv:2311.18651*, 2023.

[85] T. Wu, G. Yang, Z. Li, K. Zhang, Z. Liu, L. Guibas, D. Lin, and G. Wetzstein, "Gpt-4v(ision) is a human-aligned evaluator for text-to-3d generation," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[86] 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 arXiv:2305.11000*, 2023.

[87] Y. Gong, H. Luo, A. H. Liu, L. Karlinsky, and J. Glass, "Listen, think, and understand," *arXiv preprint arXiv:2305.10790*, 2023.

[88] S. Deshmukh, B. Elizalde, R. Singh, and H. Wang, "Pengi: An audio language model for audio tasks," *arXiv preprint arXiv:2305.11834*, 2023.

[89] P. K. Rubenstein, C. Asawaroengchai, D. D. Nguyen, A. Bapna, Z. Borsos, F. d. C. Quitry, P. Chen, D. E. Badawy, W. Han, E. Kharitonov *et al.*, "Audiopalm: A large language model that can speak and listen," *arXiv preprint arXiv:2306.12925*, 2023.

[90] S. Liu, A. S. Hussain, C. Sun, and Y. Shan, "Music understanding llama: Advancing text-to-music generation with question answering and captioning," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 286–290.

[91] Q. Chen, Y. Chu, Z. Gao, Z. Li, K. Hu, X. Zhou, J. Xu, Z. Ma, W. Wang, S. Zheng *et al.*, "Lauragpt: Listen, attend, understand, and regenerate audio with gpt," *arXiv preprint arXiv:2310.04673*, 2023.

[92] J. Gardner, S. Durand, D. Stoller, and R. M. Bittner, "Llark: A multimodal foundation model for music," *arXiv preprint arXiv:2310.07160*, 2023.

[93] C. Tang, W. Yu, G. Sun, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, and C. Zhang, "Salmonn: Towards generic hearing abilities for large language models," *arXiv preprint arXiv:2310.13289*, 2023.

[94] Y. Chu, J. Xu, X. Zhou, Q. Yang, S. Zhang, Z. Yan, C. Zhou, and J. Zhou, "Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models," *arXiv preprint arXiv:2311.07919*, 2023.

[95] A. S. Hussain, S. Liu, C. Sun, and Y. Shan, "M$^2$ ugen: Multi-modal music understanding and generation with the power of large language models," *arXiv preprint arXiv:2311.11255*, 2023.

[96] Y. Shu, S. Dong, G. Chen, W. Huang, R. Zhang, D. Shi, Q. Xiang, and Y. Shi, "Llasm: Large language and speech model," *arXiv preprint arXiv:2308.15930*, 2023.

[97] R. Yuan, H. Lin, Y. Wang, Z. Tian, S. Wu, T. Shen, G. Zhang, Y. Wu, C. Liu, Z. Zhou *et al.*, "Chatmusician: Understanding and generating music intrinsically with llm," *arXiv preprint arXiv:2402.16153*, 2024.

[98] S. Ding, Z. Liu, X. Dong, P. Zhang, R. Qian, C. He, D. Lin, and J. Wang, "Songcomposer: A large language model for lyric and melody composition in song generation," *arXiv preprint arXiv:2402.17645*, 2024.

[99] S.-L. Wu, C. Donahue, S. Watanabe, and N. J. Bryan, "Music controlnet: Multiple time-varying controls for music generation," *arXiv preprint arXiv:2311.07069*, 2023.

[100] D. Ghosal, N. Majumder, A. Mehrish, and S. Poria, "Text-to-audio generation using instruction-tuned llm and latent diffusion model," *arXiv preprint arXiv:2304.13731*, 2023.

[101] S.-L. Wu, X. Chang, G. Wichern, J.-w. Jung, F. Germain, J. Le Roux, and S. Watanabe, "Improving audio captioning models with fine-grained audio features, text embedding supervision, and llm mix-up augmentation," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 316–320.

[102] 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 arXiv:2305.18474*, 2023.

[103] Z. Wang, S. Mao, W. Wu, Y. Xia, Y. Deng, and J. Tien, "Assessing phrase break of esl speech with pre-trained language models and large language models," *arXiv preprint arXiv:2306.04980*, 2023.

[104] A. Vyas, B. Shi, M. Le, A. Tjandra, Y.-C. Wu, B. Guo, J. Zhang, X. Zhang, R. Adkins, W. Ngan *et al.*, "Audiobox: Unified audio generation with natural language prompts," *arXiv preprint arXiv:2312.15821*, 2023.

[105] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, "Hugging-gpt: Solving ai tasks with chatgpt and its friends in huggingface," in *Advances in Neural Information Processing Systems*, 2023.

[106] R. Huang, M. Li, D. Yang, J. Shi, X. Chang, Z. Ye, Y. Wu, Z. Hong, J. Huang, J. Liu *et al.*, "Audiogpt: Understanding and generating speech, music, sound, and talking head," *arXiv preprint arXiv:2304.12995*, 2023.

[107] 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 arXiv:2307.14335*, 2023.

[108] D. Yu, K. Song, P. Lu, T. He, X. Tan, W. Ye, S. Zhang, and J. Bian, "Musicagent: An ai agent for music understanding and generation with large language models," *arXiv preprint arXiv:2310.11954*, 2023.

[109] Y. Zhang, A. Maezawa, G. Xia, K. Yamamoto, and S. Dixon, "Loop copilot: Conducting ai ensembles for music generation and iterative editing," *arXiv preprint arXiv:2310.12404*, 2023.

[110] L. Zhuo, R. Yuan, J. Pan, Y. Ma, Y. LI, G. Zhang, S. Liu, R. Dannenberg, J. Fu, C. Lin *et al.*, "Lyricwhiz: Robust multilingual zero-shot lyrics transcription by whispering to chatgpt," *arXiv preprint arXiv:2306.17103*, 2023.

[111] P. Dhariwal, H. Jun, C. Payne, J. W. Kim, A. Radford, and I. Sutskever, "Jukebox: A generative model for music," *arXiv preprint arXiv:2005.00341*, 2020.

[112] 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*, 2023.

[113] D. Yang, J. Tian, X. Tan, R. Huang, S. Liu, X. Chang, J. Shi, S. Zhao, J. Bian, X. Wu *et al.*, "Uniaudio: An audio foundation model toward universal audio generation," *arXiv preprint arXiv:2310.00704*, 2023.

[114] Z. Liu, Z. Lai, Z. Gao, E. Cui, Z. Li, X. Zhu, L. Lu, Q. Chen, Y. Qiao, J. Dai, and W. Wang, "Controlllm: Augment language models with tools by searching on graphs," *arXiv preprint arXiv:2310.17796*, 2023.

[115] R. Yang, L. Song, Y. Li, S. Zhao, Y. Ge, X. Li, and Y. Shan, "Gpt4tools: Teaching large language model to use tools via self-instruction," in *Advances in Neural Information Processing Systems*, 2023.

[116] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, "Visual chatgpt: Talking, drawing and editing with visual foundation models," *arXiv preprint arXiv:2303.04671*, 2023.

[117] Z. Liu, Y. He, W. Wang, W. Wang, Y. Wang, S. Chen, Q. Zhang, Y. Yang, Q. Li, J. Yu *et al.*, "Internchat: Solving vision-centric tasks by interacting with chatbots beyond language," *arXiv preprint arXiv:2305.05662*, 2023.

[118] C. Li, H. Chen, M. Yan, W. Shen, H. Xu, Z. Wu, Z. Zhang, W. Zhou, Y. Chen, C. Cheng *et al.*, "Modelscope-agent: Building your customizable agent system with open-source large language models," *arXiv preprint arXiv:2309.00986*, 2023.

[119] 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 arXiv:2303.18223*, 2023.

[120] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao, "Large language models: A survey," *arXiv preprint arXiv:2402.06196*, 2024.

[121] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, W. Ye, Y. Zhang, Y. Chang, P. S. Yu, Q. Yang, and X. Xie, "A survey on evaluation of large language models," *ACM Trans. Intell. Syst. Technol.*, vol. 15, no. 3, mar 2024. [Online]. Available: https://doi.org/10.1145/3641289

[122] C. Zhang, C. Zhang, M. Zhang, and I. S. Kweon, "Text-to-image diffusion model in generative ai: A survey," *arXiv preprint arXiv:2303.07909*, 2023.

[123] Z. Xing, Q. Feng, H. Chen, Q. Dai, H. Hu, H. Xu, Z. Wu, and Y.-G. Jiang, "A survey on video diffusion models," *arXiv preprint arXiv:2310.10647*, 2023.

[124] Z. Shi, S. Peng, Y. Xu, A. Geiger, Y. Liao, and Y. Shen, "Deep generative models on 3d representations: A survey," *arXiv preprint arXiv:2210.15663*, 2022.

[125] C. Zhang, C. Zhang, S. Zheng, M. Zhang, M. Qamar, S.-H. Bae, and I. S. Kweon, "Audio diffusion model for speech synthesis: A survey on text to speech and speech enhancement in generative ai," *arXiv preprint arXiv:2303.13336*, 2023.

[126] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "Eva-clip: Improved training techniques for clip at scale," *arXiv preprint arXiv:2303.15389*, 2023.

[127] 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," *Journal of machine learning research*, vol. 21, no. 140, pp. 1–67, 2020.

[128] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin *et al.*, "A survey on large language model based autonomous agents," *arXiv preprint arXiv:2308.11432*, 2023.

[129] S. Yin, C. Fu, S. Zhao, K. Li, X. Sun, T. Xu, and E. Chen, "A survey on multimodal large language models," *arXiv preprint arXiv:2306.13549*, 2023.

[130] J. Wu, W. Gan, Z. Chen, S. Wan, and S. Y. Philip, "Multimodal large language models: A survey," in *2023 IEEE International Conference on Big Data (BigData)*. IEEE, 2023, pp. 2247–2256.

[131] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[132] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[133] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 873–12 883.

[134] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International conference on machine learning*. PMLR, 2015, pp. 2256–2265.

[135] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.

[136] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*. Springer, 2020, pp. 776–794.

[137] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, vol. 1, no. 2, p. 3, 2022.

[138] C. Wang, M. Chai, M. He, D. Chen, and J. Liao, "Clip-nerf: Text-and-image driven manipulation of neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3835–3844.

[139] B. Elizalde, S. Deshmukh, M. Al Ismail, and H. Wang, "Clap learning audio concepts from natural language supervision," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[140] S. Luo, C. Yan, C. Hu, and H. Zhao, "Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[141] 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.

[142] S. Wu, H. Fei, L. Qu, W. Ji, and T.-S. Chua, "Next-gpt: Any-to-any multimodal llm," *arXiv preprint arXiv:2309.05519*, 2023.

[143] Y. Xing, Y. He, Z. Tian, X. Wang, and Q. Chen, "Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners," *arXiv preprint arXiv:2402.17723*, 2024.

[144] Y. Su, T. Lan, H. Li, J. Xu, Y. Wang, and D. Cai, "Pandagpt: One model to instruction-follow them all," *arXiv preprint arXiv:2305.16355*, 2023.

[145] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[146] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," 2018.

[147] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[148] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

[149] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan *et al.*, "Training a helpful and harmless assistant with reinforcement learning from human feedback," *arXiv preprint arXiv:2204.05862*, 2022.

[150] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigpt-4: Enhancing vision-language understanding with advanced large language models," *arXiv preprint arXiv:2304.10592*, 2023.

[151] W.-G. Chen, I. Spiridonova, J. Yang, J. Gao, and C. Li, "Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing," *arXiv preprint arXiv:2311.00571*, 2023.

[152] R. Pi, L. Yao, J. Gao, J. Zhang, and T. Zhang, "Perceptiongpt: Effectively fusing visual perception into llm," *arXiv preprint arXiv:2311.06612*, 2023.

[153] R. Pi, J. Gao, S. Diao, R. Pan, H. Dong, J. Zhang, L. Yao, J. Han, H. Xu, and L. K. T. Zhang, "Detgpt: Detect what you need via reasoning," *arXiv preprint arXiv:2305.14167*, 2023.

[154] J. Gao, R. Pi, J. Zhang, J. Ye, W. Zhong, Y. Wang, L. Hong, J. Han, H. Xu, Z. Li *et al.*, "G-llava: Solving geometric problem with multi-modal large language model," *arXiv preprint arXiv:2312.11370*, 2023.

[155] Q. Sun, Q. Yu, Y. Cui, F. Zhang, X. Zhang, Y. Wang, H. Gao, J. Liu, T. Huang, and X. Wang, "Generative pretraining in multimodality," *arXiv preprint arXiv:2307.05222*, 2023.

[156] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," *arXiv preprint arXiv:1506.03365*, 2015.

[157] Z. Liu, P. Luo, X. Wang, and X. Tang, "Large-scale celebfaces attributes (celeba) dataset," *Retrieved August*, vol. 15, no. 2018, p. 11, 2018.

[158] "Dalle-1," https://openai.com/research/dall-e.

[159] R. Dong, C. Han, Y. Peng, Z. Qi, Z. Ge, J. Yang, L. Zhao, J. Sun, H. Zhou, H. Wei *et al.*, "Dreamllm: Synergistic multimodal comprehension and creation," *arXiv preprint arXiv:2309.11499*, 2023.

[160] "Dalle-3," https://openai.com/dall-e-3.

[161] J. Y. Koh, R. Salakhutdinov, and D. Fried, "Grounding language models to images for multimodal generation," *arXiv preprint arXiv:2301.13823*, 2023.

[162] J. Y. Koh, D. Fried, and R. Salakhutdinov, "Generating images with multimodal language models," *arXiv preprint arXiv:2305.17216*, 2023.

[163] L. Yu, Y. Cheng, Z. Wang, V. Kumar, W. Macherey, Y. Huang, D. A. Ross, I. Essa, Y. Bisk, M.-H. Yang *et al.*, "Spae: Semantic pyramid autoencoder for multimodal generation with frozen llms," *arXiv preprint arXiv:2306.17842*, 2023.

[164] L. Yu, B. Shi, R. Pasunuru, B. Muller, O. Golovneva, T. Wang, A. Babu, B. Tang, B. Karrer, S. Sheynin *et al.*, "Scaling autoregressive multi-modal models: Pretraining and instruction tuning," *arXiv preprint arXiv:2309.02591*, 2023.

[165] K. Zheng, X. He, and X. E. Wang, "Minigpt-5: Interleaved vision-and-language generation via generative vokens," *arXiv preprint arXiv:2310.02239*, 2023.

[166] J. An, Z. Yang, L. Li, J. Wang, K. Lin, Z. Liu, L. Wang, and J. Luo, "Openleaf: Open-domain interleaved image-text generation and evaluation," *arXiv preprint arXiv:2310.07749*, 2023.

[167] Z. Yang, Y. Zhang, F. Meng, and J. Zhou, "Teal: Tokenize and embed all for multi-modal large language models," *arXiv preprint arXiv:2311.04589*, 2023.

[168] B. Xia, S. Wang, Y. Tao, Y. Wang, and J. Jia, "Llmga: Multimodal large language model based generation assistant," *arXiv preprint arXiv:2311.16500*, 2023.

[169] X. Chi, Y. Liu, Z. Jiang, R. Zhang, Z. Lin, R. Zhang, P. Gao, C. Fu, S. Zhang, Q. Liu *et al.*, "Chatillusion: Efficient-aligning interleaved generation ability with visual instruction model," *arXiv preprint arXiv:2311.17963*, 2023.

[170] Y. Zhou, R. Zhang, J. Gu, and T. Sun, "Customization assistant for text-to-image generation," *arXiv preprint arXiv:2312.03045*, 2023.

[171] X. Shen and M. Elhoseiny, "Storygpt-v: Large language models as consistent story visualizers," 2023.

[172] X. Hu, R. Wang, Y. Fang, B. Fu, P. Cheng, and G. Yu, "Ella: Equip diffusion models with llm for enhanced semantic alignment," *arXiv preprint arXiv:2403.05135*, 2024.

[173] S. Zhao, S. Hao, B. Zi, H. Xu, and K.-Y. K. Wong, "Bridging different language models and generative vision models for text-to-image generation," *arXiv preprint arXiv:2403.07860*, 2024.

[174] J. Cho, A. Zala, and M. Bansal, "Visual programming for text-to-image generation and evaluation," *arXiv preprint arXiv:2305.15328*, 2023.

[175] H. Gani, S. F. Bhat, M. Naseer, S. Khan, and P. Wonka, "Llm blueprint: Enabling text-to-image generation with complex and detailed prompts," *arXiv preprint arXiv:2310.10640*, 2023.

[176] T.-H. Wu, L. Lian, J. E. Gonzalez, B. Li, and T. Darrell, "Self-correcting llm-controlled diffusion models," *arXiv preprint arXiv:2311.16090*, 2023.

[177] J. Chen, Y. Huang, T. Lv, L. Cui, Q. Chen, and F. Wei, "Textdiffuser: Diffusion models as text painters," *arXiv preprint arXiv:2305.10855*, 2023.

[178] P. Jia, C. Li, Z. Liu, Y. Shen, X. Chen, Y. Yuan, Y. Zheng, D. Chen, J. Li, X. Xie *et al.*, "Cole: A hierarchical generation framework for graphic design," *arXiv preprint arXiv:2311.16974*, 2023.

[179] S. Zhong, Z. Huang, W. Wen, J. Qin, and L. Lin, "Sur-adapter: Enhancing text-to-image pre-trained diffusion models with large language models," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 567–578.

[180] Q. Yu, J. Li, W. Ye, S. Tang, and Y. Zhuang, "Interactive data synthesis for systematic vision adaptation via llms-aigcs collaboration," *arXiv preprint arXiv:2305.12799*, 2023.

[181] X. Wang, B. Zhuang, and Q. Wu, "Switchgpt: Adapting large language models for non-text outputs," *arXiv preprint arXiv:2309.07623*, 2023.

[182] J. Liao, X. Chen, Q. Fu, L. Du, X. He, X. Wang, S. Han, and D. Zhang, "Text-to-image generation for abstract concepts," *arXiv preprint arXiv:2309.14623*, 2023.

[183] Z. Yang, J. Wang, L. Li, K. Lin, C.-C. Lin, Z. Liu, and L. Wang, "Idea2img: Iterative self-refinement with gpt-4v (ision) for automatic image design and generation," *arXiv preprint arXiv:2310.08541*, 2023.

[184] J.-Y. He, Z.-Q. Cheng, C. Li, J. Sun, W. Xiang, X. Lin, X. Kang, Z. Jin, Y. Hu, B. Luo *et al.*, "Wordart designer: User-driven artistic typography synthesis using large language models," *arXiv preprint arXiv:2310.18332*, 2023.

[185] J. Sun, D. Fu, Y. Hu, S. Wang, R. Rassin, D.-C. Juan, D. Alon, C. Herrmann, S. van Steenkiste, R. Krishna *et al.*, "Dreamsync: Aligning text-to-image generation with image understanding feedback," *arXiv preprint arXiv:2311.17946*, 2023.

[186] Y. Lu, X. Yang, X. Li, X. E. Wang, and W. Y. Wang, "Llmscore: Unveiling the power of large language models in text-to-image synthesis evaluation," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[187] V. Ordonez, G. Kulkarni, and T. Berg, "Im2text: Describing images using 1 million captioned photographs," *Advances in neural information processing systems*, vol. 24, 2011.

[188] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.

[189] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International conference on machine learning*. PMLR, 2021, pp. 4904–4916.

[190] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3558–3568.

[191] K. Srinivasan, K. Raman, J. Chen, M. Bendersky, and M. Najork, "Wit: Wikipedia-based image text dataset for multimodal multilingual machine learning," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 2443–2449.

[192] C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, and A. Komatsuzaki, "Laion-400m: Open dataset of clip-filtered 400 million image-text pairs," *arXiv preprint arXiv:2111.02114*, 2021.

[193] Y. Zheng, H. Yang, T. Zhang, J. Bao, D. Chen, Y. Huang, L. Yuan, D. Chen, M. Zeng, and F. Wen, "General facial representation learning in a visual-linguistic manner," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 697–18 709.

[194] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 716–23 736, 2022.

[195] LAION-COCO, "https://laion.ai/blog/laion-coco/," 2022. [Online]. Available: https://laion.ai/blog/laion-coco/

[196] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, "Laion-5b: An open large-scale dataset for training next generation image-text models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25 278–25 294, 2022.

[197] Coyo-700M, "https://huggingface.co/datasets/kakaobrain/coyo-700m," 2022. [Online]. Available: https://huggingface.co/datasets/kakaobrain/coyo-700m

[198] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed, B. Patra *et al.*, "Language is not all you need: Aligning perception with language models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[199] W. Zhu, J. Hessel, A. Awadalla, S. Y. Gadre, J. Dodge, A. Fang, Y. Yu, L. Schmidt, W. Y. Wang, and Y. Choi, "Multimodal c4: An open, billion-scale corpus of images interleaved with text," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[200] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *Advances in neural information processing systems*, vol. 36, 2024.

[201] S. Y. Gadre, G. Ilharco, A. Fang, J. Hayase, G. Smyrnis, T. Nguyen, R. Marten, M. Wortsman, D. Ghosh, J. Zhang *et al.*, "Datacomp: In search of the next generation of multimodal datasets," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[202] J. Chen, Y. Huang, T. Lv, L. Cui, Q. Chen, and F. Wei, "Textdiffuser: Diffusion models as text painters," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[203] Y. Yang, D. Gui, Y. Yuan, W. Liang, H. Ding, H. Hu, and K. Chen, "Glyphcontrol: Glyph conditional control for visual text generation," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[204] B. Li, Y. Zhang, L. Chen, J. Wang, F. Pu, J. Yang, C. Li, and Z. Liu, "Mimic-it: Multi-modal in-context instruction tuning," *arXiv preprint arXiv:2306.05425*, 2023.

[205] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36 479–36 494, 2022.

[206] T. Vetter and T. Poggio, "Linear object classes and image synthesis from a single example image," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 19, no. 7, pp. 733–742, 1997.

[207] M. F. Cohen and J. R. Wallace, *Radiosity and realistic image synthesis*. Morgan Kaufmann, 1993.

[208] D. Kirk and J. Arvo, "Unbiased sampling techniques for image synthesis," *ACM SIGGRAPH Computer Graphics*, vol. 25, no. 4, pp. 153–156, 1991.

[209] K. Frans, L. Soros, and O. Witkowski, "Clipdraw: Exploring text-to-drawing synthesis through language-image encoders," *Advances in Neural Information Processing Systems*, vol. 35, pp. 5207–5218, 2022.

[210] S. Gu, D. Chen, J. Bao, F. Wen, B. Zhang, D. Chen, L. Yuan, and B. Guo, "Vector quantized diffusion model for text-to-image synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 696–10 706.

[211] Z. Wang, W. Liu, Q. He, X. Wu, and Z. Yi, "Clip-gen: Language-free training of a text-to-image generator with clip," *arXiv preprint arXiv:2203.00386*, 2022.

[212] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *arXiv preprint arXiv:2304.08485*, 2023.

[213] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," *arXiv preprint arXiv:2301.12597*, 2023.

[214] T. Lv, Y. Huang, J. Chen, L. Cui, S. Ma, Y. Chang, S. Huang, W. Wang, L. Dong, W. Luo *et al.*, "Kosmos-2.5: A multimodal literate model," *arXiv preprint arXiv:2309.11419*, 2023.

[215] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi *et al.*, "mplug-owl: Modularization empowers large language models with multimodality," *arXiv preprint arXiv:2304.14178*, 2023.

[216] J. Ye, A. Hu, H. Xu, Q. Ye, M. Yan, Y. Dan, C. Zhao, G. Xu, C. Li, J. Tian *et al.*, "mplug-docowl: Modularized multimodal large language model for document understanding," *arXiv preprint arXiv:2307.02499*, 2023.

[217] J. Chen, D. Zhu, X. Shen, X. Li, Z. Liu, P. Zhang, R. Krishnamoorthi, V. Chandra, Y. Xiong, and M. Elhoseiny, "Minigpt-v2: large language model as a unified interface for vision-language multi-task learning," *arXiv preprint arXiv:2310.09478*, 2023.

[218] R. Zhang, J. Han, A. Zhou, X. Hu, S. Yan, P. Lu, H. Li, P. Gao, and Y. Qiao, "Llama-adapter: Efficient fine-tuning of language models with zero-init attention," *arXiv preprint arXiv:2303.16199*, 2023.

[219] W. Dai, J. Li, D. Li, A. Tiong, J. Zhao, W. Wang, B. Li, P. Fung, and S. Hoi, "Instructblip: Towards general-purpose vision-language models with instruction tuning. arxiv 2023," *arXiv preprint arXiv:2305.06500*.

[220] W. Wang, Z. Chen, X. Chen, J. Wu, X. Zhu, G. Zeng, P. Luo, T. Lu, J. Zhou, Y. Qiao *et al.*, "Visionllm: Large language model is also an open-ended decoder for vision-centric tasks," *arXiv preprint arXiv:2305.11175*, 2023.

[221] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou, "Qwen-vl: A frontier large vision-language model with versatile abilities," *arXiv preprint arXiv:2308.12966*, 2023.

[222] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song *et al.*, "Cogvlm: Visual expert for pretrained language models," *arXiv preprint arXiv:2311.03079*, 2023.

[223] Y. Jin, K. Xu, L. Chen, C. Liao, J. Tan, B. Chen, C. Lei, A. Liu, C. Song, X. Lei *et al.*, "Unified language-vision pretraining with dynamic discrete visual tokenization," *arXiv preprint arXiv:2309.04669*, 2023.

[224] Z. Chen, J. Wu, W. Wang, W. Su, G. Chen, S. Xing, Z. Muyan, Q. Zhang, X. Zhu, L. Lu *et al.*, "Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks," *arXiv preprint arXiv:2312.14238*, 2023.

[225] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed, Q. Liu *et al.*, "Language is not all you need: Aligning perception with language models," *arXiv preprint arXiv:2302.14045*, 2023.

[226] Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, and L. Wang, "Mm-react: Prompting chatgpt for multimodal reasoning and action," *arXiv preprint arXiv:2303.11381*, 2023.

[227] S. Wang, Z. Zhao, X. Ouyang, Q. Wang, and D. Shen, "Chatcad: Interactive computer-aided diagnosis on medical image using large language models," *arXiv preprint arXiv:2302.07257*, 2023.

[228] Z. Yang, Z. Gan, J. Wang, X. Hu, Y. Lu, Z. Liu, and L. Wang, "An empirical study of gpt-3 for few-shot knowledge-based vqa," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 3081–3089.

[229] Z. Gu, B. Zhu, G. Zhu, Y. Chen, M. Tang, and J. Wang, "Anomalygpt: Detecting industrial anomalies using large vision-language models," *arXiv preprint arXiv:2308.15366*, 2023.

[230] J. Qin, J. Wu, W. Chen, Y. Ren, H. Li, H. Wu, X. Xiao, R. Wang, and S. Wen, "Diffusiongpt: Llm-driven text-to-image generation system," *arXiv preprint arXiv:2401.10061*, 2024.

[231] L. Yang, Z. Yu, C. Meng, M. Xu, S. Ermon, and B. Cui, "Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms," *arXiv preprint arXiv:2401.11708*, 2024.

[232] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "LoRA: Low-rank adaptation of large language models," in *International Conference on Learning Representations*, 2022. [Online]. Available: https://openreview.net/forum?id=nZeVKeeFYf9

[233] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, "Sdedit: Guided image synthesis and editing with stochastic differential equations," *arXiv preprint arXiv:2108.01073*, 2021.

[234] G. Kim, T. Kwon, and J. C. Ye, "Diffusionclip: Text-guided diffusion models for robust image manipulation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2426–2435.

[235] G. Couairon, J. Verbeek, H. Schwenk, and M. Cord, "Diffedit: Diffusion-based semantic image editing with mask guidance," *arXiv preprint arXiv:2210.11427*, 2022.

[236] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or, "Prompt-to-prompt image editing with cross attention control," *arXiv preprint arXiv:2208.01626*, 2022.

[237] R. Mokady, A. Hertz, K. Aberman, Y. Pritch, and D. Cohen-Or, "Null-text inversion for editing real images using guided diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6038–6047.

[238] B. Kawar, S. Zada, O. Lang, O. Tov, H. Chang, T. Dekel, I. Mosseri, and M. Irani, "Imagic: Text-based real image editing with diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6007–6017.

[239] N. Tumanyan, M. Geyer, S. Bagon, and T. Dekel, "Plug-and-play diffusion features for text-driven image-to-image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 1921–1930.

[240] R. Morita, Z. Zhang, M. M. Ho, and J. Zhou, "Interactive image manipulation with complex text instructions," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 1053–1062.

[241] Z. Zhang, L. Han, A. Ghosh, D. N. Metaxas, and J. Ren, "Sine: Single image editing with text-to-image diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6027–6037.

[242] G. Parmar, K. Kumar Singh, R. Zhang, Y. Li, J. Lu, and J.-Y. Zhu, "Zero-shot image-to-image translation," in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–11.

[243] V. Goel, E. Peruzzo, Y. Jiang, D. Xu, N. Sebe, T. Darrell, Z. Wang, and H. Shi, "Pair-diffusion: Object-level image editing with structure-and-appearance paired diffusion models," *arXiv preprint arXiv:2303.17546*, 2023.

[244] M. Cao, X. Wang, Z. Qi, Y. Shan, X. Qie, and Y. Zheng, "Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing," *arXiv preprint arXiv:2304.08465*, 2023.

[245] T. Nguyen, Y. Li, U. Ojha, and Y. J. Lee, "Visual instruction inversion: Image editing via visual prompting," *arXiv preprint arXiv:2307.14331*, 2023.

[246] A. Mirzaei, T. Aumentado-Armstrong, M. A. Brubaker, J. Kelly, A. Levinshtein, K. G. Derpanis, and I. Gilitschenski, "Watch your steps: Local image and scene editing by text instructions," *arXiv preprint arXiv:2308.08947*, 2023.

[247] C. Mou, X. Wang, J. Song, Y. Shan, and J. Zhang, "Dragondiffusion: Enabling drag-style manipulation on diffusion models," *arXiv preprint arXiv:2307.02421*, 2023.

[248] ——, "Diffeditor: Boosting accuracy and flexibility on diffusion-based image editing," *arXiv preprint arXiv:2402.02583*, 2024.

[249] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," *arXiv preprint arXiv:2211.09800*, 2022.

[250] X. Cui, Z. Li, P. Li, Y. Hu, H. Shi, C. Cao, and Z. He, "Chatedit: Towards multi-turn interactive facial image editing via dialogue," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 14 567–14 583.

[251] T.-J. Fu, W. Hu, X. Du, W. Y. Wang, Y. Yang, and Z. Gan, "Guiding instruction-based image editing via multimodal large language models," *arXiv preprint arXiv:2309.17102*, 2023.

[252] S. Sheynin, A. Polyak, U. Singer, Y. Kirstain, A. Zohar, O. Ashual, D. Parikh, and Y. Taigman, "Emu edit: Precise image editing via recognition and generation tasks," *arXiv preprint arXiv:2311.10089*, 2023.

[253] Y. Huang, L. Xie, X. Wang, Z. Yuan, X. Cun, Y. Ge, J. Zhou, C. Dong, R. Huang, R. Zhang *et al.*, "Smartedit: Exploring complex instruction-based image editing with multimodal large language models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

[254] Z. Geng, B. Yang, T. Hang, C. Li, S. Gu, T. Zhang, J. Bao, Z. Zhang, H. Hu, D. Chen *et al.*, "Instructdiffusion: A generalist modeling interface for vision tasks," *arXiv preprint arXiv:2309.03895*, 2023.

[255] J. Z. Wu, Y. Ge, X. Wang, S. W. Lei, Y. Gu, Y. Shi, W. Hsu, Y. Shan, X. Qie, and M. Z. Shou, "Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 7623–7633.

[256] E. Molad, E. Horwitz, D. Valevski, A. R. Acha, Y. Matias, Y. Pritch, Y. Leviathan, and Y. Hoshen, "Dreamix: Video diffusion models are general video editors," *arXiv preprint arXiv:2302.01329*, 2023.

[257] S. Liu, Y. Zhang, W. Li, Z. Lin, and J. Jia, "Video-p2p: Video editing with cross-attention control," *arXiv preprint arXiv:2303.04761*, 2023.

[258] C. Qi, X. Cun, Y. Zhang, C. Lei, X. Wang, Y. Shan, and Q. Chen, "Fatezero: Fusing attentions for zero-shot text-based video editing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 15 932–15 942.

[259] D. Ceylan, C.-H. P. Huang, and N. J. Mitra, "Pix2video: Video editing using image diffusion," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 23 206–23 217.

[260] W. Chai, X. Guo, G. Wang, and Y. Lu, "Stablevideo: Text-driven consistency-aware diffusion video editing," in *Proceedings of the*

*IEEE/CVF International Conference on Computer Vision*, 2023, pp. 23 040–23 050.

[261] S. Yang, Y. Zhou, Z. Liu, and C. C. Loy, "Rerender a video: Zero-shot text-guided video-to-video translation," in *SIGGRAPH Asia 2023 Conference Papers*, 2023, pp. 1–11.

[262] M. Geyer, O. Bar-Tal, S. Bagon, and T. Dekel, "Tokenflow: Consistent diffusion features for consistent video editing," *arXiv preprint arXiv:2307.10373*, 2023.

[263] H. Ouyang, Q. Wang, Y. Xiao, Q. Bai, J. Zhang, K. Zheng, X. Zhou, Q. Chen, and Y. Shen, "Codef: Content deformation fields for temporally consistent video processing," *arXiv preprint arXiv:2308.07926*, 2023.

[264] J. H. Liew, H. Yan, J. Zhang, Z. Xu, and J. Feng, "Magicedit: High-fidelity and temporally coherent video editing," *arXiv preprint arXiv:2308.14749*, 2023.

[265] Y. Ma, X. Cun, Y. He, C. Qi, X. Wang, Y. Shan, X. Li, and Q. Chen, "Magicstick: Controllable video editing via control handle transformations," *arXiv preprint arXiv:2312.03047*, 2023.

[266] J. Cheng, T. Xiao, and T. He, "Consistent video-to-video transfer using synthetic dataset," *arXiv preprint arXiv:2311.00213*, 2023.

[267] B. Qin, J. Li, S. Tang, T.-S. Chua, and Y. Zhuang, "Instructvid2vid: Controllable video editing with natural language instructions," *arXiv preprint arXiv:2305.12328*, 2023.

[268] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18 392–18 402.

[269] B. Wu, S. Yu, Z. Chen, J. B. Tenenbaum, and C. Gan, "Star: A benchmark for situated reasoning in real-world videos," in *Thirty-fifth conference on neural information processing systems datasets and benchmarks track (Round 2)*, 2021.

[270] Z. Liu, L. Wang, W. Wu, C. Qian, and T. Lu, "Tam: Temporal adaptive module for video recognition," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 13 708–13 718.

[271] M. Zhao, B. Li, J. Wang, W. Li, W. Zhou, L. Zhang, S. Xuyang, Z. Yu, X. Yu, G. Li *et al.*, "Towards video text visual question answering: benchmark and baseline," *Advances in Neural Information Processing Systems*, vol. 35, pp. 35 549–35 562, 2022.

[272] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, "Frozen in time: A joint video and image encoder for end-to-end retrieval," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1728–1738.

[273] A. Yang, A. Miech, J. Sivic, I. Laptev, and C. Schmid, "Zero-shot video question answering via frozen bidirectional language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 124–141, 2022.

[274] ——, "Learning to answer visual questions from web videos," *arXiv preprint arXiv:2205.05019*, 2022.

[275] K. Lin, L. Li, C.-C. Lin, F. Ahmed, Z. Gan, Z. Liu, Y. Lu, and L. Wang, "Swinbert: End-to-end transformers with sparse attention for video captioning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 949–17 958.

[276] J. Lin, C. Gan, and S. Han, "Tsm: Temporal shift module for efficient video understanding," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 7083–7093.

[277] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding?" in *ICML*, vol. 2, no. 3, 2021, p. 4.

[278] C.-Y. Wu, C. Feichtenhofer, H. Fan, K. He, P. Krahenbuhl, and R. Girshick, "Long-term feature banks for detailed video understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 284–293.

[279] H. Zhang, X. Li, and L. Bing, "Video-llama: An instruction-tuned audio-visual language model for video understanding," *arXiv preprint arXiv:2306.02858*, 2023.

[280] J. Chen, D. Zhu, K. Haydarov, X. Li, and M. Elhoseiny, "Video chatcaptioner: Towards the enriched spatiotemporal descriptions," *arXiv preprint arXiv:2304.04227*, 2023.

[281] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S. W. Kim, S. Fidler, and K. Kreis, "Align your latents: High-resolution video synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 563–22 575.

[282] Y. He, M. Xia, H. Chen, X. Cun, Y. Gong, J. Xing, Y. Zhang, X. Wang, C. Weng, Y. Shan *et al.*, "Animate-a-story: Sto-

[283] Y. Ma, Y. He, X. Cun, X. Wang, S. Chen, X. Li, and Q. Chen, "Follow your pose: Pose-guided text-to-video generation using pose-free videos," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 5, 2024, pp. 4117–4125.

[284] Y. He, H. Liu, H. Chen, X. Cun, X. Wang, Y. Shan *et al.*, "Make-your-video: Customized video generation using textual and structural guidance." *IEEE Transactions on Visualization and Computer Graphics*, 2024.

[285] H. Qiu, M. Xia, Y. Zhang, Y. He, X. Wang, Y. Shan, and Z. Liu, "Freenoise: Tuning-free longer video diffusion via noise rescheduling," *arXiv preprint arXiv:2310.15169*, 2023.

[286] Y. Ma, Y. He, H. Wang, A. Wang, C. Qi, C. Cai, X. Li, Z. Li, H.-Y. Shum, W. Liu *et al.*, "Follow-your-click: Open-domain regional image animation via short prompts," *arXiv preprint arXiv:2403.08268*, 2024.

[287] Y. Bao, D. Qiu, G. Kang, B. Zhang, B. Jin, K. Wang, and P. Yan, "Latentwarp: Consistent diffusion latents for zero-shot video-to-video translation," *arXiv preprint arXiv:2311.00353*, 2023.

[288] J. Xu, T. Mei, T. Yao, and Y. Rui, "Msr-vtt: A large video description dataset for bridging video and language," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5288–5296, 2016. [Online]. Available: https://api.semanticscholar.org/CorpusID:206594535

[289] A. Rohrbach, A. Torabi, M. Rohrbach, N. Tandon, C. J. Pal, H. Larochelle, A. C. Courville, and B. Schiele, "Movie description," *International Journal of Computer Vision*, vol. 123, pp. 94 – 120, 2016. [Online]. Available: https://api.semanticscholar.org/CorpusID:18217052

[290] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. C. Niebles, "Dense-captioning events in videos," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 706–715, 2017. [Online]. Available: https://api.semanticscholar.org/CorpusID:1026139

[291] R. Sanabria, O. Caglayan, S. Palaskar, D. Elliott, L. Barrault, L. Specia, and F. Metze, "How2: A large-scale dataset for multimodal language understanding," *ArXiv*, vol. abs/1811.00347, 2018. [Online]. Available: https://api.semanticscholar.org/CorpusID:53186236

[292] X. E. Wang, J. Wu, J. Chen, L. Li, Y. fang Wang, and W. Y. Wang, "Vatex: A large-scale, high-quality multilingual dataset for video-and-language research," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4580–4590, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:102352148

[293] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, "Howto100m: Learning a text-video embedding by watching hundred million narrated video clips," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2630–2640, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:182952863

[294] J. C. Stroud, D. A. Ross, C. Sun, J. Deng, R. Sukthankar, and C. Schmid, "Learning video representations from textual web supervision," *ArXiv*, vol. abs/2007.14937, 2020. [Online]. Available: https://api.semanticscholar.org/CorpusID:220845567

[295] R. Zellers, X. Lu, J. Hessel, Y. Yu, J. S. Park, J. Cao, A. Farhadi, and Y. Choi, "Merlot: Multimodal neural script knowledge models," in *Neural Information Processing Systems*, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:235352775

[296] H. Xue, T. Hang, Y. Zeng, Y. Sun, B. Liu, H. Yang, J. Fu, and B. Guo, "Advancing high-resolution video-language representation with large-scale video transcriptions," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5026–5035, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:244462849

[297] A. Nagrani, P. H. Seo, B. Seybold, A. Hauth, S. Manén, C. Sun, and C. Schmid, "Learning audio-video modalities from image captions," in *European Conference on Computer Vision*, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:247939759

[298] W. Wang, H. Yang, Z. Tuo, H. He, J. Zhu, J. Fu, and J. Liu, "Videofactory: Swap attention in spatiotemporal diffusions for text-to-video generation," *ArXiv*, vol. abs/2305.10874, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:258762479

[299] Y. Wang, Y. He, Y. Li, K. Li, J. Yu, X. J. Ma, X. Chen, Y. Wang, P. Luo, Z. Liu, Y. Wang, L. Wang, and Y. Qiao, "Internvid: A large-scale video-text dataset for multimodal understanding

and generation," *ArXiv*, vol. abs/2307.06942, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:259847783

[300] T.-S. Chen, A. Siarohin, W. Menapace, E. Deyneka, H.-w. Chao, B. E. Jeon, Y. Fang, H.-Y. Lee, J. Ren, M.-H. Yang *et al.*, "Panda-70m: Captioning 70m videos with multiple cross-modality teachers," *arXiv preprint arXiv:2402.19479*, 2024.

[301] "Vript," https://github.com/mutonix/Vript.

[302] J. Yu, H. Zhu, L. Jiang, C. C. Loy, W. T. Cai, and W. Wu, "Celebv-text: A large-scale facial text-video dataset," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14805–14814, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:257767123

[303] X. Li, Q. Zhang, D. Kang, W. Cheng, Y. Gao, J. Zhang, Z. Liang, J. Liao, Y.-P. Cao, and Y. Shan, "Advances in 3d generation: A survey," *arXiv preprint arXiv:2401.17807*, 2024.

[304] M. Deitke, D. Schwenk, J. Salvador, L. Weihs, O. Michel, E. VanderBilt, L. Schmidt, K. Ehsani, A. Kembhavi, and A. Farhadi, "Objaverse: A universe of annotated 3d objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13142–13153.

[305] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8821–8831.

[306] A. 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," *arXiv preprint arXiv:2112.10741*, 2021.

[307] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich, "Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12619–12629.

[308] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[309] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *arXiv preprint arXiv:2106.10689*, 2021.

[310] T. Shen, J. Gao, K. Yin, M.-Y. Liu, and S. Fidler, "Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis," *Advances in Neural Information Processing Systems*, vol. 34, pp. 6087–6101, 2021.

[311] M. Deitke, R. Liu, M. Wallingford, H. Ngo, O. Michel, A. Kusupati, A. Fan, C. Laforte, V. Voleti, S. Y. Gadre *et al.*, "Objaverse-xl: A universe of 10m+ 3d objects," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[312] J. Reizenstein, R. Shapovalov, P. Henzler, L. Sbordone, P. Labatut, and D. Novotny, "Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10901–10911.

[313] A. Nichol, H. Jun, P. Dhariwal, P. Mishkin, and M. Chen, "Point-e: A system for generating 3d point clouds from complex prompts," *arXiv preprint arXiv:2212.08751*, 2022.

[314] J. Lei, Y. Zhang, K. Jia *et al.*, "Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition," *Advances in Neural Information Processing Systems*, vol. 35, pp. 30923–30936, 2022.

[315] Y. Ma, X. Zhang, X. Sun, J. Ji, H. Wang, G. Jiang, W. Zhuang, and R. Ji, "X-mesh: Towards fast and accurate text-driven 3d stylization via dynamic textual guidance," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2749–2760.

[316] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," *arXiv preprint arXiv:1605.08803*, 2016.

[317] A. Jain, B. Mildenhall, J. T. Barron, P. Abbeel, and B. Poole, "Zero-shot text-guided object generation with dream fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 867–876.

[318] F. Yin, X. Chen, C. Zhang, B. Jiang, Z. Zhao, J. Fan, G. Yu, T. Li, and T. Chen, "Shapegpt: 3d shape generation with a unified multi-modal language model," *arXiv preprint arXiv:2311.17618*, 2023.

[319] G. Tevet, B. Gordon, A. Hertz, A. H. Bermano, and D. Cohen-Or, "Motionclip: Exposing human motion generation to clip space,"

[320] B. Jiang, X. Chen, W. Liu, J. Yu, G. Yu, and T. Chen, "Motiongpt: Human motion as a foreign language," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[321] Z. Zhou and S. Tulsiani, "Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction," in *CVPR*, 2023.

[322] Z. Wan, D. Paschalidou, I. Huang, H. Liu, B. Shen, X. Xiang, J. Liao, and L. Guibas, "Cad: Photorealistic 3d generation via adversarial distillation," *arXiv preprint arXiv:2312.06663*, 2023.

[323] B. Yang, W. Dong, L. Ma, W. Hu, X. Liu, Z. Cui, and Y. Ma, "Dreamspace: Dreaming your room space with text-driven panoramic texture propagation," in *2024 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2024, pp. 650–660.

[324] G. Metzer, E. Richardson, O. Patashnik, R. Giryes, and D. Cohen-Or, "Latent-nerf for shape-guided generation of 3d shapes and textures," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12663–12673.

[325] O. Katzir, O. Patashnik, D. Cohen-Or, and D. Lischinski, "Noise-free score distillation," 2023.

[326] M. Armandpour, H. Zheng, A. Sadeghian, A. Sadeghian, and M. Zhou, "Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond," *arXiv preprint arXiv:2304.04968*, 2023.

[327] L. Zhou, A. Shih, C. Meng, and S. Ermon, "Dreampropeller: Supercharge text-to-3d generation with parallel sampling," *arXiv preprint arXiv:2311.17082*, 2023.

[328] C. Yu, G. Lu, Y. Zeng, J. Sun, X. Liang, H. Li, Z. Xu, S. Xu, W. Zhang, and H. Xu, "Towards high-fidelity text-guided 3d face generation and manipulation using only images," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 15326–15337.

[329] C. Zhang, Y. Chen, Y. Fu, Z. Zhou, G. Yu, B. Wang, B. Fu, T. Chen, G. Lin, and C. Shen, "Styleavatar3d: Leveraging image-text diffusion models for high-fidelity 3d avatar generation," *arXiv preprint arXiv:2305.19012*, 2023.

[330] T. Wang, B. Zhang, T. Zhang, S. Gu, J. Bao, T. Baltrusaitis, J. Shen, D. Chen, F. Wen, Q. Chen *et al.*, "Rodin: A generative model for sculpting 3d digital avatars using diffusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4563–4573.

[331] S. Aneja, J. Thies, A. Dai, and M. Nießner, "Clipface: Text-guided editing of textured 3d morphable models," in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–11.

[332] M. Wu, H. Zhu, L. Huang, Y. Zhuang, Y. Lu, and X. Cao, "High-fidelity 3d face generation from natural language descriptions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4521–4530.

[333] T. Liao, H. Yi, Y. Xiu, J. Tang, Y. Huang, J. Thies, and M. J. Black, "Tada! text to animatable digital avatars," *arXiv preprint arXiv:2308.10899*, 2023.

[334] S. Huang, Z. Yang, L. Li, Y. Yang, and J. Jia, "Avatarfusion: Zero-shot generation of clothing-decoupled 3d avatars using 2d diffusion," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 5734–5745.

[335] X. Han, Y. Cao, K. Han, X. Zhu, J. Deng, Y.-Z. Song, T. Xiang, and K.-Y. K. Wong, "Headsculpt: Crafting 3d head avatars with text," *arXiv preprint arXiv:2306.03038*, 2023.

[336] Y. Cao, Y.-P. Cao, K. Han, Y. Shan, and K.-Y. K. Wong, "Dreamavatar: Text-and-shape guided 3d human avatar generation via diffusion models," *arXiv preprint arXiv:2304.00916*, 2023.

[337] H. Zhang, B. Chen, H. Yang, L. Qu, X. Wang, L. Chen, C. Long, F. Zhu, K. Du, and M. Zheng, "Avatarverse: High-quality & stable 3d avatar creation from text and pose," *arXiv preprint arXiv:2308.03610*, 2023.

[338] L. Zhang, Q. Qiu, H. Lin, Q. Zhang, C. Shi, W. Yang, Y. Shi, S. Yang, L. Xu, and J. Yu, "Dreamface: Progressive generation of animatable 3d faces under text guidance," *arXiv preprint arXiv:2304.03117*, 2023.

[339] F. Hong, M. Zhang, L. Pan, Z. Cai, L. Yang, and Z. Liu, "Avatarclip: Zero-shot text-driven generation and animation of 3d avatars," *arXiv preprint arXiv:2205.08535*, 2022.

[340] N. Kolotouros, T. Alldieck, A. Zanfir, E. G. Bazavan, M. Fieraru, and C. Sminchisescu, "Dreamhuman: Animatable 3d avatars from text," *arXiv preprint arXiv:2306.09329*, 2023.

[341] X. Huang, R. Shao, Q. Zhang, H. Zhang, Y. Feng, Y. Liu, and Q. Wang, "Humannorm: Learning normal diffusion model for

in *European Conference on Computer Vision*. Springer, 2022, pp. 358–374.

high-quality and realistic 3d human generation," *arXiv preprint arXiv:2310.01406*, 2023.

[342] Y. Zeng, Y. Lu, X. Ji, Y. Yao, H. Zhu, and X. Cao, "Avatarbooth: High-quality and customizable 3d human avatar generation," *arXiv preprint arXiv:2306.09864*, 2023.

[343] D. Wang, H. Meng, Z. Cai, Z. Shao, Q. Liu, L. Wang, M. Fan, Y. Shan, X. Zhan, and Z. Wang, "Headevolver: Text to head avatars via locally learnable mesh deformation," *arXiv preprint arXiv:2403.09326*, 2024.

[344] H. Liu, X. Wang, Z. Wan, Y. Shen, Y. Song, J. Liao, and Q. Chen, "Headartist: Text-conditioned 3d head generation with self score distillation," *arXiv preprint arXiv:2312.07539*, 2023.

[345] Y. Shi, P. Wang, J. Ye, M. Long, K. Li, and X. Yang, "Mv-dream: Multi-view diffusion for 3d generation," *arXiv preprint arXiv:2308.16512*, 2023.

[346] Y. Kant, Z. Wu, M. Vasilkovsky, G. Qian, J. Ren, R. A. Guler, B. Ghanem, S. Tulyakov, I. Gilitschenski, and A. Siaro-hin, "Spad: Spatially aware multiview diffusers," *arXiv preprint arXiv:2402.05235*, 2024.

[347] Z. Liu, Y. Li, Y. Lin, X. Yu, S. Peng, Y.-P. Cao, X. Qi, X. Huang, D. Liang, and W. Ouyang, "Unidream: Unifying diffusion priors for relightable text-to-3d generation," 2023.

[348] L. Qiu, G. Chen, X. Gu, Q. zuo, M. Xu, Y. Wu, W. Yuan, Z. Dong, L. Bo, and X. Han, "Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d," *arXiv preprint arXiv:2311.16918*, 2023.

[349] J. Li, H. Tan, K. Zhang, Z. Xu, F. Luan, Y. Xu, Y. Hong, K. Sunkavalli, G. Shakhnarovich, and S. Bi, "Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model," *arXiv preprint arXiv:2311.06214*, 2023.

[350] J. Tang, Z. Chen, X. Chen, T. Wang, G. Zeng, and Z. Liu, "Lgm: Large multi-view gaussian model for high-resolution 3d content creation," *arXiv preprint arXiv:2402.05054*, 2024.

[351] X. Yinghao, S. Zifan, Y. Wang, C. Hansheng, Y. Ceyuan, P. Sida, S. Yujun, and W. Gordon, "Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation," 2024.

[352] H. Jun and A. Nichol, "Shap-e: Generating conditional 3d implicit functions," *arXiv preprint arXiv:2305.02463*, 2023.

[353] Z. Hu, A. Iscen, A. Jain, T. Kipf, Y. Yue, D. A. Ross, C. Schmid, and A. Fathi, "Scenecraft: An llm agent for synthesizing 3d scene as blender code," *arXiv preprint arXiv:2403.01248*, 2024.

[354] R. Xu, X. Wang, T. Wang, Y. Chen, J. Pang, and D. Lin, "Pointllm: Empowering large language models to understand point clouds," *arXiv preprint arXiv:2308.16911*, 2023.

[355] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan, "3d-llm: Injecting the 3d world into large language models," *arXiv preprint arXiv:2307.12981*, 2023.

[356] O. Gordon, O. Avrahami, and D. Lischinski, "Blended-nerf: Zero-shot object generation and blending in existing neural radiance fields," *arXiv preprint arXiv:2306.12760*, 2023.

[357] W. Gao, N. Aigerman, T. Groueix, V. Kim, and R. Hanocka, "Textdeformer: Geometry manipulation using text guidance," in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–11.

[358] C. Bao, Y. Zhang, B. Yang, T. Fan, Z. Yang, H. Bao, G. Zhang, and Z. Cui, "Sine: Semantic-driven image-based nerf editing with prior-guided editing field," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 20 919–20 929.

[359] A. Mikaeili, O. Perel, M. Safaee, D. Cohen-Or, and A. Mahdavi-Amiri, "Sked: Sketch-guided text-based 3d editing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 14 607–14 619.

[360] J. Zhuang, C. Wang, L. Lin, L. Liu, and G. Li, "Dreameditor: Text-driven 3d scene editing with neural fields," in *SIGGRAPH Asia 2023 Conference Papers*, 2023, pp. 1–10.

[361] A. Haque, M. Tancik, A. A. Efros, A. Holynski, and A. Kanazawa, "Instruct-nerf2nerf: Editing 3d scenes with instructions," *arXiv preprint arXiv:2303.12789*, 2023.

[362] D. Decatur, I. Lang, K. Aberman, and R. Hanocka, "3d paint-brush: Local stylization of 3d shapes with cascaded score distillation," *arXiv preprint arXiv:2311.09571*, 2023.

[363] G. Tevet, S. Raab, B. Gordon, Y. Shafir, D. Cohen-or, and A. H. Bermano, "Human motion diffusion model," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: https://openreview.net/forum?id=SJ1kSyO2jwu

[364] R. Chen, Y. Chen, N. Jiao, and K. Jia, "Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation," *arXiv preprint arXiv:2303.13873*, 2023.

[365] Z. Pan, J. Lu, X. Zhu, and L. Zhang, "Enhancing high-resolution 3d generation through pixel-wise gradient clipping," in *International Conference on Learning Representations (ICLR)*, 2024.

[366] G. Qian, J. Cao, A. Siarohin, Y. Kant, C. Wang, M. Vasilkovsky, H.-Y. Lee, Y. Fang, I. Skorokhodov, P. Zhuang *et al.*, "Atom: Amortized text-to-mesh using 2d diffusion," *arXiv preprint arXiv:2402.00867*, 2024.

[367] Z. Wu, P. Zhou, X. Yi, X. Yuan, and H. Zhang, "Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior," *arXiv preprint arXiv:2401.09050*, 2024.

[368] T. Huang, Y. Zeng, Z. Zhang, W. Xu, H. Xu, S. Xu, R. W. Lau, and W. Zuo, "Dreamcontrol: Control-based text-to-3d generation with 3d self-prior," *arXiv preprint arXiv:2312.06439*, 2023.

[369] Y. Chen, C. Zhang, X. Yang, Z. Cai, G. Yu, L. Yang, and G. Lin, "It3d: Improved text-to-3d generation with explicit view synthesis," 2023.

[370] M. Zhao, C. Zhao, X. Liang, L. Li, Z. Zhao, Z. Hu, C. Fan, and X. Yu, "Efficientdreamer: High-fidelity and robust 3d creation via orthogonal-view diffusion prior," *arXiv preprint arXiv:2308.13223*, 2023.

[371] Z. Chen, F. Wang, and H. Liu, "Text-to-3d using gaussian splatting," *arXiv preprint arXiv:2309.16585*, 2023.

[372] Y. Ma, Y. Fan, J. Ji, H. Wang, X. Sun, G. Jiang, A. Shu, and R. Ji, "X-dreamer: Creating high-quality 3d content by bridging the domain gap between text-to-2d and text-to-3d generation," *arXiv preprint arXiv:2312.00085*, 2023.

[373] J. Wu, X. Gao, X. Liu, Z. Shen, C. Zhao, H. Feng, J. Liu, and E. Ding, "Hd-fusion: Detailed text-to-3d generation leveraging multiple noise estimation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 3202–3211.

[374] X. Yang, Y. Chen, C. Chen, C. Zhang, Y. Xu, X. Yang, F. Liu, and G. Lin, "Learn to optimize denoising scores for 3d generation: A unified and improved diffusion prior on nerf and 3d gaussian splatting," *arXiv preprint arXiv:2312.04820*, 2023.

[375] F. Liu, D. Wu, Y. Wei, Y. Rao, and Y. Duan, "Sherpa3d: Boosting high-fidelity text-to-3d generation via coarse 3d prior," 2023.

[376] Y. Lin, R. Clark, and P. Torr, "Dreampolisher: Towards high-quality text-to-3d generation via geometric diffusion," *arXiv preprint arXiv:2403.17237*, 2024.

[377] Y. Yang, F.-Y. Sun, L. Weihs, E. VanderBilt, A. Herrasti, W. Han, J. Wu, N. Haber, R. Krishna, L. Liu, C. Callison-Burch, M. Yatskar, A. Kembhavi, and C. Clark, "Holodeck: Language guided generation of 3d embodied ai environments," *arXiv preprint arXiv:2312.09067*, 2023.

[378] H. Song, S. Choi, H. Do, C. Lee, and T. Kim, "Blending-nerf: Text-driven localized editing in neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 14 383–14 393.

[379] R. He, S. Huang, X. Nie, T. Hui, L. Liu, J. Dai, J. Han, G. Li, and S. Liu, "Customize your nerf: Adaptive source driven 3d scene editing via local-global iterative training," *arXiv preprint arXiv:2312.01663*, 2023.

[380] X. Zeng, X. Chen, Z. Qi, W. Liu, Z. Zhao, Z. Wang, B. FU, Y. Liu, and G. Yu, "Paint3d: Paint anything 3d with lighting-less texture diffusion models," 2023.

[381] E. Law, K. West, M. I. Mandel, M. Bay, and J. S. Downie, "Evaluation of algorithms using games: The case of music tagging." in *ISMIR*. Citeseer, 2009, pp. 387–392.

[382] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.

[383] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2017, pp. 776–780.

[384] C. Hawthorne, A. Stasyuk, A. Roberts, I. Simon, C.-Z. A. Huang, S. Dieleman, E. Elsen, J. Engel, and D. Eck, "Enabling factorized piano music modeling and generation with the maestro dataset," *arXiv preprint arXiv:1810.12247*, 2018.

[385] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, "Libritts: A corpus derived from librispeech for text-to-speech," *arXiv preprint arXiv:1904.02882*, 2019.

[386] D. Bogdanov, M. Won, P. Tovstogan, A. Porter, and X. Serra, "The mtg-jamendo dataset for automatic music tagging." ICML, 2019.

[387] J. Kahn, M. Riviere, W. Zheng, E. Kharitonov, Q. Xu, P.-E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen *et al.*, "Librilight: A benchmark for asr with limited or no supervision," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7669–7673.

[388] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, "Vggsound: A large-scale audio-visual dataset," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 721–725.

[389] B. Zhang, H. Lv, P. Guo, Q. Shao, C. Yang, L. Xie, X. Xu, H. Bu, X. Chen, C. Zeng *et al.*, "Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6182–6186.

[390] W. Kang, X. Yang, Z. Yao, F. Kuang, Y. Yang, L. Guo, L. Lin, and D. Povey, "Libriheavy: a 50,000 hours asr corpus with punctuation casing and context," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 10991–10995.

[391] J. Zhan, J. Dai, J. Ye, Y. Zhou, D. Zhang, Z. Liu, X. Zhang, R. Yuan, G. Zhang, L. Li *et al.*, "Anygpt: Unified multimodal llm with discrete sequence modeling," *arXiv preprint arXiv:2402.12226*, 2024.

[392] H. Hao, L. Zhou, S. Liu, J. Li, S. Hu, R. Wang, and F. Wei, "Boosting large language model for speech synthesis: An empirical study," *arXiv preprint arXiv:2401.00246*, 2023.

[393] J. Lu, C. Clark, S. Lee, Z. Zhang, S. Khosla, R. Marten, D. Hoiem, and A. Kembhavi, "Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action," *arXiv preprint arXiv:2312.17172*, 2023.

[394] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg *et al.*, "Sparks of artificial general intelligence: Early experiments with gpt-4," *arXiv preprint arXiv:2303.12712*, 2023.

[395] J. Wu, Y. Gaur, Z. Chen, L. Zhou, Y. Zhu, T. Wang, J. Li, S. Liu, B. Ren, L. Liu *et al.*, "On decoder-only architecture for speech-to-text and large language model integration," in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2023, pp. 1–8.

[396] W. Yu, C. Tang, G. Sun, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, and C. Zhang, "Connecting speech encoder and large language model for asr," *arXiv preprint arXiv:2309.13963*, 2023.

[397] S. Wang, C.-H. H. Yang, J. Wu, and C. Zhang, "Can whisper perform speech-based in-context learning," *arXiv preprint arXiv:2309.07081*, 2023.

[398] Q. Deng, Q. Yang, R. Yuan, Y. Huang, Y. Wang, X. Liu, Z. Tian, J. Pan, G. Zhang, H. Lin *et al.*, "Composerx: Multi-agent symbolic music composition with llms," *arXiv preprint arXiv:2404.18081*, 2024.

[399] Y. Gong, Y.-A. Chung, and J. Glass, "Ast: Audio spectrogram transformer," *arXiv preprint arXiv:2104.01778*, 2021.

[400] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *arXiv preprint arXiv:2210.11416*, 2022.

[401] J. Rothstein, *MIDI: A comprehensive introduction*. AR Editions, Inc., 1995, vol. 7.

[402] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen *et al.*, "Palm 2 technical report," *arXiv preprint arXiv:2305.10403*, 2023.

[403] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International Conference on Machine Learning*. PMLR, 2023, pp. 28492–28518.

[404] S. Kakouros, J. Šimko, M. Vainio, and A. Suni, "Investigating the utility of surprisal from large language models for speech synthesis prosody," *arXiv preprint arXiv:2306.09814*, 2023.

[405] Y. Gong, A. Rouditchenko, A. H. Liu, D. Harwath, L. Karlinsky, H. Kuehne, and J. Glass, "Contrastive audio-visual masked autoencoder," *arXiv preprint arXiv:2210.07839*, 2022.

[406] Z. Deng, Y. Ma, Y. Liu, R. Guo, G. Zhang, W. Chen, W. Huang, and E. Benetos, "Musilingo: Bridging music and text with pre-trained language models for music captioning and query response," *arXiv preprint arXiv:2309.08730*, 2023.

[407] 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 arXiv:2306.00107*, 2023.

[408] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," *arXiv preprint arXiv:2210.13438*, 2022.

[409] R. Kumar, P. Seetharaman, A. Luebs, I. Kumar, and K. Kumar, "High-fidelity audio compression with improved rvqgan," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[410] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[411] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi, "Self-instruct: Aligning language models with self-generated instructions," *arXiv preprint arXiv:2212.10560*, 2022.

[412] Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian *et al.*, "Toolllm: Facilitating large language models to master 16000+ real-world apis," *arXiv preprint arXiv:2307.16789*, 2023.

[413] Q. Tang, Z. Deng, H. Lin, X. Han, Q. Liang, and L. Sun, "Toolalpaca: Generalized tool learning for language models with 3000 simulated cases," *arXiv preprint arXiv:2306.05301*, 2023.

[414] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, L. Zettlemoyer, N. Cancedda, and T. Scialom, "Toolformer: Language models can teach themselves to use tools," *arXiv preprint arXiv:2302.04761*, 2023.

[415] N. Farn and R. Shin, "Tooltalk: Evaluating tool-usage in a conversation setting," *arXiv preprint arXiv:2311.10775*, 2023.

[416] S. Hao, T. Liu, Z. Wang, and Z. Hu, "Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings," *arXiv preprint arXiv:2305.11554*, 2023.

[417] C.-Y. Hsieh, S.-A. Chen, C.-L. Li, Y. Fujii, A. Ratner, C.-Y. Lee, R. Krishna, and T. Pfister, "Tool documentation enables zero-shot tool-usage with large language models," *arXiv preprint arXiv:2308.00675*, 2023.

[418] J. Ruan, Y. Chen, B. Zhang, Z. Xu, T. Bao, G. Du, S. Shi, H. Mao, X. Zeng, and R. Zhao, "Tptu: Task planning and tool usage of large language model-based ai agents," *arXiv preprint arXiv:2308.03427*, 2023.

[419] A. Parisi, Y. Zhao, and N. Fiedel, "Talm: Tool augmented language models," *arXiv preprint arXiv:2205.12255*, 2022.

[420] J. Zhang, "Graph-toolformer: To empower llms with graph reasoning ability via prompt augmented by chatgpt," *arXiv preprint arXiv:2304.11116*, 2023.

[421] Y. Zhuang, X. Chen, T. Yu, S. Mitra, V. Bursztyn, R. A. Rossi, S. Sarkhel, and C. Zhang, "Toolchain*: Efficient action space navigation in large language models with a* search," *arXiv preprint arXiv:2310.13227*, 2023.

[422] Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, N. Duan, and W. Chen, "Critic: Large language models can self-correct with tool-interactive critiquing," *arXiv preprint arXiv:2305.11738*, 2023.

[423] Q. Jin, Y. Yang, Q. Chen, and Z. Lu, "Genegpt: Augmenting large language models with domain tools for improved access to biomedical information," *ArXiv*, 2023.

[424] B. Paranjape, S. Lundberg, S. Singh, H. Hajishirzi, L. Zettlemoyer, and M. T. Ribeiro, "Art: Automatic multi-step reasoning and tool-use for large language models," *arXiv preprint arXiv:2303.09014*, 2023.

[425] Z. Gou, Z. Shao, Y. Gong, Y. Yang, M. Huang, N. Duan, W. Chen *et al.*, "Tora: A tool-integrated reasoning agent for mathematical problem solving," *arXiv preprint arXiv:2309.17452*, 2023.

[426] Y. Song, W. Xiong, D. Zhu, C. Li, K. Wang, Y. Tian, and S. Li, "Restgpt: Connecting large language models with real-world applications via restful apis," *arXiv preprint arXiv:2306.06624*, 2023.

[427] S. Qiao, H. Gui, H. Chen, and N. Zhang, "Making language models better tool learners with execution feedback," *arXiv preprint arXiv:2305.13068*, 2023.

[428] K. Zhang, H. Chen, L. Li, and W. Wang, "Syntax error-free and generalizable tool use for llms via finite-state decoding," *arXiv preprint arXiv:2310.07075*, 2023.

[429] W. Shen, C. Li, H. Chen, M. Yan, X. Quan, H. Chen, J. Zhang, and F. Huang, "Small llms are weak tool learners: A multi-llm agent," *arXiv preprint arXiv:2401.07324*, 2024.

[430] J. Wang, H. Xu, J. Ye, M. Yan, W. Shen, J. Zhang, F. Huang, and J. Sang, "Mobile-agent: Autonomous multi-modal mobile device

agent with visual perception," *arXiv preprint arXiv:2401.16158*, 2024.

[431] T. Gupta and A. Kembhavi, "Visual programming: Compositional visual reasoning without training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 953–14 962.

[432] C. Wang, W. Luo, Q. Chen, H. Mai, J. Guo, S. Dong, X. M. Xuan, Z. Li, L. Ma, and S. Gao, "Mllm-tool: A multimodal large language model for tool agent learning," *arXiv preprint arXiv:2401.10727*, 2024.

[433] D. Surís, S. Menon, and C. Vondrick, "Vipergpt: Visual inference via python execution for reasoning," *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2023.

[434] Z. Gao, Y. Du, X. Zhang, X. Ma, W. Han, S.-C. Zhu, and Q. Li, "Clova: A closed-loop visual assistant with tool usage and update," *arXiv preprint arXiv:2312.10908*, 2023.

[435] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao, "ReAct: Synergizing reasoning and acting in language models," in *International Conference on Learning Representations (ICLR)*, 2023.

[436] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24 824–24 837, 2022.

[437] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," *arXiv preprint arXiv:2304.02643*, 2023.

[438] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, T. Mihaylov, M. Ott, S. Shleifer, K. Shuster, D. Simig, P. S. Koura, A. Sridhar, T. Wang, and L. Zettlemoyer, "Opt: Open pre-trained transformer language models," 2022.

[439] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," *arXiv preprint arXiv:2305.14314*, 2023.

[440] S. Mangrulkar, S. Gugger, L. Debut, Y. Belkada, S. Paul, and B. Bossan, "Peft: State-of-the-art parameter-efficient fine-tuning methods," https://github.com/huggingface/peft, 2022.

[441] J. Chen, X. Li, X. Ye, C. Li, Z. Fan, and H. Zhao, "Idea-2-3d: Collaborative lmm agents enable 3d model generation from interleaved multimodal inputs," *arXiv preprint arXiv:2404.04363*, 2024.

[442] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, "Universal adversarial triggers for attacking and analyzing nlp," *arXiv preprint arXiv:1908.07125*, 2019.

[443] X. Fu, Z. Wang, S. Li, R. K. Gupta, N. Mireshghallah, T. Berg-Kirkpatrick, and E. Fernandes, "Misusing tools in large language models with visual adversarial examples," *arXiv preprint arXiv:2310.03185*, 2023.

[444] L. Bailey, E. Ong, S. Russell, and S. Emmons, "Image hijacks: Adversarial images can control generative models at runtime," *arXiv preprint arXiv:2309.00236*, 2023.

[445] A. Zou, Z. Wang, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," *arXiv preprint arXiv:2307.15043*, 2023.

[446] E. Jones, A. Dragan, A. Raghunathan, and J. Steinhardt, "Automatically auditing large language models via discrete optimization," in *International Conference on Machine Learning*. PMLR, 2023, pp. 15 307–15 329.

[447] P. Żelasko, S. Joshi, Y. Shao, J. Villalba, J. Trmal, N. Dehak, and S. Khudanpur, "Adversarial attacks and defenses for speech recognition systems," *arXiv preprint arXiv:2103.17122*, 2021.

[448] Z. Chen, L. Xie, S. Pang, Y. He, and Q. Tian, "Appending adversarial frames for universal video attack," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2021, pp. 3199–3208.

[449] H. Liu, W. Zhou, D. Chen, H. Fang, H. Bian, K. Liu, W. Zhang, and N. Yu, "Coherent adversarial deepfake video generation," *Signal Processing*, vol. 203, p. 108790, 2023.

[450] S.-Y. Lo and V. M. Patel, "Defending against multiple and unforeseen adversarial videos," *IEEE Transactions on Image Processing*, vol. 31, pp. 962–973, 2021.

[451] H. J. Lee and Y. M. Ro, "Defending video recognition model against adversarial perturbations via defense patterns," *IEEE Transactions on Dependable and Secure Computing*, 2023.

[452] Y. Wu, X. Li, Y. Liu, P. Zhou, and L. Sun, "Jailbreaking gpt-4v via self-adversarial attacks with system prompts," *arXiv preprint arXiv:2311.09127*, 2023.

[453] Y. Xie, J. Yi, J. Shao, J. Curl, L. Lyu, Q. Chen, X. Xie, and F. Wu, "Defending chatgpt against jailbreak attack via self-reminders," *Nature Machine Intelligence*, vol. 5, no. 12, pp. 1486–1496, 2023.

[454] Y. Liu, G. Deng, Y. Li, K. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu, "Prompt injection attack against llm-integrated applications," *arXiv preprint arXiv:2306.05499*, 2023.

[455] F. Perez and I. Ribeiro, "Ignore previous prompt: Attack techniques for language models," *arXiv preprint arXiv:2211.09527*, 2022.

[456] N. Carlini, M. Jagielski, C. A. Choquette-Choo, D. Paleka, W. Pearce, H. Anderson, A. Terzis, K. Thomas, and F. Tramèr, "Poisoning web-scale training datasets is practical," *arXiv preprint arXiv:2302.10149*, 2023.

[457] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," *arXiv preprint arXiv:1707.07328*, 2017.

[458] M.-H. Van and X. Wu, "Detecting and correcting hate speech in multimodal memes with large visual language model," *arXiv preprint arXiv:2311.06737*, 2023.

[459] Z. Wei, Y. Wang, and Y. Wang, "Jailbreak and guard aligned language models with only few in-context demonstrations," *arXiv preprint arXiv:2310.06387*, 2023.

[460] A. Robey, E. Wong, H. Hassani, and G. J. Pappas, "Smoothllm: Defending large language models against jailbreaking attacks," *arXiv preprint arXiv:2310.03684*, 2023.

[461] R. Liu, A. Khakzar, J. Gu, Q. Chen, P. Torr, and F. Pizzati, "Latent guard: a safety framework for text-to-image generation," *arXiv preprint arXiv:2404.08031*, 2024.

[462] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[463] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[464] R. Pi, T. Han, W. Xiong, J. Zhang, R. Liu, R. Pan, and T. Zhang, "Strengthening multimodal large language model with bootstrapped preference optimization," *arXiv preprint arXiv:2403.08730*, 2024.

[465] X. Wu, K. Sun, F. Zhu, R. Zhao, and H. Li, "Better aligning text-to-image models with human preference," *arXiv preprint arXiv:2303.14420*, 2023.

[466] H. Dong, W. Xiong, D. Goyal, R. Pan, S. Diao, J. Zhang, K. Shum, and T. Zhang, "Raft: Reward ranked finetuning for generative foundation model alignment," *arXiv preprint arXiv:2304.06767*, 2023.

[467] P. Korshunov and S. Marcel, "Deepfakes: A new threat to face recognition? assessment and detection. arxiv 2018," *arXiv preprint arXiv:1812.08685*.

[468] Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," *ACM computing surveys (CSUR)*, vol. 54, no. 1, pp. 1–41, 2021.

[469] M. Masood, M. Nawaz, K. M. Malik, A. Javed, A. Irtaza, and H. Malik, "Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward," *Applied intelligence*, vol. 53, no. 4, pp. 3974–4026, 2023.

[470] L. Verdoliva, "Media forensics and deepfakes: an overview," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 910–932, 2020.

[471] D. Wodajo, S. Atnafu, and Z. Akhtar, "Deepfake video detection using generative convolutional vision transformer," *arXiv preprint arXiv:2307.07036*, 2023.

[472] D. Wodajo and S. Atnafu, "Deepfake video detection using convolutional vision transformer," *arXiv preprint arXiv:2102.11126*, 2021.

[473] S. Hussain, P. Neekhara, M. Jere, F. Koushanfar, and J. McAuley, "Adversarial deepfakes: Evaluating vulnerability of deepfake detectors to adversarial examples," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2021, pp. 3348–3357.

[474] W. Shi, A. Ajith, M. Xia, Y. Huang, D. Liu, T. Blevins, D. Chen, and L. Zettlemoyer, "Detecting pretraining data from large language models," *arXiv preprint arXiv:2310.16789*, 2023.

[475] S. M. Park, K. Georgiev, A. Ilyas, G. Leclerc, and A. Madry, "Trak: Attributing model behavior at scale," *arXiv preprint arXiv:2303.14186*, 2023.

[476] Z. Wang, C. Chen, Y. Zeng, L. Lyu, and S. Ma, "Where did i come from? origin attribution of ai-generated images," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[477] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein, "A watermark for large language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 17 061–17 084.

[478] Y. Cui, J. Ren, H. Xu, P. He, H. Liu, L. Sun, and J. Tang, "Diffusionshield: A watermark for copyright protection against generative diffusion models," *arXiv preprint arXiv:2306.04642*, 2023.

[479] P. Fernandez, G. Couairon, H. Jégou, M. Douze, and T. Furon, "The stable signature: Rooting watermarks in latent diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 22 466–22 477.

[480] Z. Zhang, L. Lei, L. Wu, R. Sun, Y. Huang, C. Long, X. Liu, X. Lei, J. Tang, and M. Huang, "Safetybench: Evaluating the safety of large language models with multiple choice questions," *arXiv preprint arXiv:2309.07045*, 2023.

[481] H. Lin, Z. Luo, B. Wang, R. Yang, and J. Ma, "Goat-bench: Safety insights to large multimodal models through meme-based social abuse," *arXiv preprint arXiv:2401.01523*, 2024.

[482] X. Wang, X. Yi, H. Jiang, S. Zhou, Z. Wei, and X. Xie, "Tovilag: Your visual-language generative model is also an evildoer," *arXiv preprint arXiv:2312.11523*, 2023.

[483] Y. Gong, D. Ran, J. Liu, C. Wang, T. Cong, A. Wang, S. Duan, and X. Wang, "Figstep: Jailbreaking large vision-language models via typographic visual prompts," *arXiv preprint arXiv:2311.05608*, 2023.

[484] X. Liu, Y. Zhu, Y. Lan, C. Yang, and Y. Qiao, "Query-relevant images jailbreak large multi-modal models," *arXiv preprint arXiv:2311.17600*, 2023.

[485] "midjourney," https://www.midjourney.com/home.

[486] "Stability ai," https://stability.ai/.

[487] "Gpt-4," https://openai.com/gpt-4.

[488] "Dalle-2," https://openai.com/dall-e-2.

[489] "Openai," https://openai.com.

[490] "Pika labs," https://www.pika.art/.

[491] "Gen2," https://research.runwayml.com/gen2.

[492] "heygen," https://app.heygen.com/home.

[493] "Azure ai-services: text-to-speech," https://azure.microsoft.com/zh-cn/products/ai-services/text-to-speech.

[494] "descript," https://www.descript.com/.

[495] "Suno ai," https://suno-ai.org/.

[496] "Stability ai: Stable audio," https://stability.ai/stable-audio.

[497] "Musicfx," https://aitestkitchen.withgoogle.com/tools/music-fx.

[498] "tuneflow," https://www.tuneflow.com/.

[499] "deepmusic," https://www.deepmusic.fun/.

[500] "meta," https://about.meta.com/.

[501] "Epic games' metahuman creator," https://www.unrealengine.com/en-US/metahuman.

[502] "Luma ai," https://lumalabs.ai/.

[503] "Adobe," https://www.adobe.com/.

[504] "Kaedim3d," https://www.kaedim3d.com/.

[505] "Wonder studio," https://wonderdynamics.com/.

[506] A. Avetisyan, C. Xie, H. Howard-Jenkins, T.-Y. Yang, S. Aroudj, S. Patra, F. Zhang, D. Frost, L. Holland, C. Orme, J. Engel, E. Miller, R. Newcombe, and V. Balntas, "Scenescript: Reconstructing scenes with an autoregressive structured language model," 2024.

[507] "google," https://www.google.com/.

[508] "tencent," https://www.tencent.com/.

[509] Y. He, S. Yang, H. Chen, X. Cun, M. Xia, Y. Zhang, X. Wang, R. He, Q. Chen, and Y. Shan, "Scalecrafter: Tuning-free higher-resolution visual generation with diffusion models," in *The Twelfth International Conference on Learning Representations*, 2023.

[510] L. Guo, Y. He, H. Chen, M. Xia, X. Cun, Y. Wang, S. Huang, Y. Zhang, X. Wang, Q. Chen *et al.*, "Make a cheap scaling: A self-cascade diffusion model for higher-resolution adaptation," *arXiv preprint arXiv:2402.10491*, 2024.

[511] Y. Xu, T. Park, R. Zhang, Y. Zhou, E. Shechtman, F. Liu, J.-B. Huang, and D. Liu, "Videogigagan: Towards detail-rich video super-resolution," *arXiv preprint arXiv:2404.12388*, 2024.

[512] S. Zhou, P. Yang, J. Wang, Y. Luo, and C. C. Loy, "Upscale-a-video: Temporal-consistent diffusion model for real-world video super-resolution," *arXiv preprint arXiv:2312.06640*, 2023.

[513] R. S. Roman, Y. Adi, A. Deleforge, R. Serizel, G. Synnaeve, and A. Défossez, "From discrete tokens to high-fidelity audio using multi-band diffusion," *arXiv preprint arXiv:2308.02560*, 2023.

[514] Y. Yao, P. Li, B. Chen, and A. Wang, "Jen-1 composer: A unified framework for high-fidelity multi-track music generation," *arXiv preprint arXiv:2310.19180*, 2023.

[515] M. Ding, W. Zheng, W. Hong, and J. Tang, "Cogview2: Faster and better text-to-image generation via hierarchical transformers," *Advances in Neural Information Processing Systems*, vol. 35, pp. 16 890–16 902, 2022.

[516] Y. Zhang, Y. Wei, X. Lin, Z. Hui, P. Ren, X. Xie, X. Ji, and W. Zuo, "Videoelevator: Elevating video generation quality with versatile text-to-image diffusion models," *arXiv preprint arXiv:2403.05438*, 2024.

[517] R. Henschel, L. Khachatryan, D. Hayrapetyan, H. Poghosyan, V. Tadevosyan, Z. Wang, S. Navasardyan, and H. Shi, "Streamingt2v: Consistent, dynamic, and extendable long video generation from text," *arXiv preprint arXiv:2403.14773*, 2024.

[518] R. Or-El, X. Luo, M. Shan, E. Shechtman, J. J. Park, and I. Kemelmacher-Shlizerman, "Stylesdf: High-resolution 3d-consistent image and geometry generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 503–13 513.

[519] X. Huang, W. Li, J. Hu, H. Chen, and Y. Wang, "Refsr-nerf: Towards high fidelity and super resolution view synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8244–8253.

[520] F.-Y. Wang, W. Chen, G. Song, H.-J. Ye, Y. Liu, and H. Li, "Gen-l-video: Multi-text to long video generation via temporal co-denoising," *arXiv preprint arXiv:2305.18264*, 2023.

[521] J. Yoo, S. Kim, D. Lee, C. Kim, and S. Hong, "Towards end-to-end generative modeling of long videos with memory-efficient bidirectional transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 888–22 897.

[522] L. Lin, G. Xia, Y. Zhang, and J. Jiang, "Arrange, inpaint, and refine: Steerable long-term music audio generation and editing via content-based controls," 2024.

[523] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3836–3847.

[524] M. Zhao, R. Wang, F. Bao, C. Li, and J. Zhu, "Controlvideo: Adding conditional control for one shot text-to-video editing," *arXiv preprint arXiv:2305.17098*, 2023.

[525] R. Liu, D. Garrette, C. Saharia, W. Chan, A. Roberts, S. Narang, I. Blok, R. Mical, M. Norouzi, and N. Constant, "Character-aware models improve visual text rendering," *arXiv preprint arXiv:2212.10562*, 2022.

[526] J. Ma, M. Zhao, C. Chen, R. Wang, D. Niu, H. Lu, and X. Lin, "Glyphdraw: Learning to draw chinese characters in image synthesis models coherently," *arXiv preprint arXiv:2303.17870*, 2023.

[527] C. Chen, X. Yang, F. Yang, C. Feng, Z. Fu, C.-S. Foo, G. Lin, and F. Liu, "Sculpt3d: Multi-view consistent text-to-3d generation with sparse 3d prior," *arXiv preprint arXiv:2403.09140*, 2024.

[528] S. Woo, B. Park, H. Go, J.-Y. Kim, and C. Kim, "Harmonyview: Harmonizing consistency and diversity in one-image-to-3d," *arXiv preprint arXiv:2312.15980*, 2023.

[529] J. Ye, P. Wang, K. Li, Y. Shi, and H. Wang, "Consistent-1-to-3: Consistent image to 3d view synthesis via geometry-aware diffusion models," *arXiv preprint arXiv:2310.03020*, 2023.

[530] Q. Zuo, X. Gu, L. Qiu, Y. Dong, Z. Zhao, W. Yuan, R. Peng, S. Zhu, Z. Dong, L. Bo *et al.*, "Videomv: Consistent multi-view generation based on large video generative model," *arXiv preprint arXiv:2403.12010*, 2024.

[531] P. Wang, S. Wang, J. Lin, S. Bai, X. Zhou, J. Zhou, X. Wang, and C. Zhou, "One-peace: Exploring one general representation model toward unlimited modalities," *arXiv preprint arXiv:2305.11172*, 2023.

[532] C. Boletsis, A. Lie, O. Prillard, K. Husby, and J. Li, "The invizar project: Augmented reality visualization for non-destructive testing data from jacket platforms," 2023.

[533] S. Chen, H. Li, Q. Wang, Z. Zhao, M. Sun, X. Zhu, and J. Liu, "Vast: A vision-audio-subtitle-text omni-modality foundation

model and dataset," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[534] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.

[535] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh *et al.*, "Mixed precision training," *arXiv preprint arXiv:1710.03740*, 2017.

[536] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2704–2713.

[537] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, "Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale," *Advances in Neural Information Processing Systems*, vol. 35, pp. 30 318–30 332, 2022.

[538] Y. Choukroun, E. Kravchik, F. Yang, and P. Kisilev, "Low-bit quantization of neural networks for efficient inference," in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2019, pp. 3009–3018.

[539] X. Wu, C. Li, R. Y. Aminabadi, Z. Yao, and Y. He, "Understanding int4 quantization for language models: latency speedup, composability, and failure cases," in *International Conference on Machine Learning*. PMLR, 2023, pp. 37 524–37 539.

[540] Y. Qu, X. Shen, X. He, M. Backes, S. Zannettou, and Y. Zhang, "Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models," in *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, 2023, pp. 3403–3417.

[541] D. Ha and J. Schmidhuber, "World models," *arXiv preprint arXiv:1803.10122*, 2018.

[542] H. Liu, W. Yan, M. Zaharia, and P. Abbeel, "World model on million-length video and language with ringattention," *arXiv preprint arXiv:2402.08268*, 2024.

[543] Y. LeCun, "A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27," *Open Review*, vol. 62, no. 1, 2022.

[544] C. Min, D. Zhao, L. Xiao, Y. Nie, and B. Dai, "Uniworld: Autonomous driving pre-training via world models," *arXiv preprint arXiv:2308.07234*, 2023.

[545] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap, "Mastering diverse domains through world models," *arXiv preprint arXiv:2301.04104*, 2023.