# A Complete Survey on Generative AI (AIGC): Is ChatGPT from GPT-4 to GPT-5 All You Need?

CHAONING ZHANG, Kyung Hee University, South Korea

CHENSHUANG ZHANG, KAIST, South Korea

SHENG ZHENG, Beijing Institute of Technology, China

YU QIAO, Kyung Hee University, South Korea

CHENGHAO LI, KAIST, South Korea

MENGCHUN ZHANG, KAIST, South Korea

SUMIT KUMAR DAM, Kyung Hee University, South Korea

CHU MYAET THWAL, Kyung Hee University, South Korea

YE LIN TUN, Kyung Hee University, South Korea

LE LUANG HUY, Kyung Hee University, South Korea

DONGUK KIM, Kyung Hee University, South Korea

SUNG-HO BAE, Kyung Hee University, South Korea

LIK-HANG LEE, Hong Kong Polytechnic University, Hong Kong (China)

YANG YANG, University of Electronic Science and technology, China

HENG TAO SHEN, University of Electronic Science and technology, China

IN SO KWEON, KAIST, South Korea

CHOONG SEON HONG, Kyung Hee University, South Korea

As ChatGPT goes viral, generative AI (AIGC, a.k.a AI-generated content) has made headlines everywhere because of its ability to analyze and create text, images, and beyond. With such overwhelming media coverage, it is almost impossible for us to miss the opportunity to glimpse AIGC from a certain angle. In the era of AI transitioning from pure analysis to creation, it is worth noting that ChatGPT, with its most recent language model GPT-4, is just a tool out of numerous AIGC tasks . Impressed by the capability of the ChatGPT, many people are wondering about its limits: can GPT-5 (or other future GPT variants) help ChatGPT unify all AIGC tasks for

Authors' addresses: Chaoning Zhang, Kyung Hee University, South Korea, chaoningzhang1990@gmail.com; Chenshuang Zhang, KAIST, South Korea, zcs15@kaist.ac.kr; Sheng Zheng, Beijing Institute of Technology, China, zszhx2021@gmail.com; Yu Qiao, Kyung Hee University, South Korea, qiaoyu@khu.ac.kr; Chenghao Li, KAIST, South Korea, lch17692405449@gmail.com; Mengchun Zhang, KAIST, South Korea, zhangmengchun527@gmail.com; Sumit Kumar Dam, Kyung Hee University, South Korea, skd160205@khu.ac.kr; Chu Myaet Thwal, Kyung Hee University, South Korea, chumyaet@khu.ac.kr; Ye Lin Tun, Kyung Hee University, South Korea, yelintun@khu.ac.kr; Le Luang Huy, Kyung Hee University, South Korea, quanghuy69@khu.ac.kr; Donguk kim, Kyung Hee University, South Korea, g9896@khu.ac.kr; Sung-Ho Bae, Kyung Hee University, South Korea, shbae@khu.ac.kr; Lik-Hang Lee, Hong Kong Polytechnic University, Hong Kong (China), iskweon77@kaist.ac.kr; Yang Yang, University of Electronic Science and technology, China, dlyyang@gmail.com; Heng Tao Shen, University of Electronic Science and technology, China, shenhengtao@hotmail.com; In So Kweon, KAIST, South Korea, iskweon77@kaist.ac.kr; Choong Seon Hong, Kyung Hee University, South Korea, cshong@khu.ac.kr.

diversified content creation? Toward answering this question, a comprehensive review of existing AIGC tasks is needed. As such, our work comes to fill this gap promptly by offering a first look at AIGC, ranging from its techniques to applications. Modern generative AI relies on various technical foundations, ranging from model architecture and self-supervised pretraining to generative modeling methods (like GAN and diffusion models). After introducing the fundamental techniques, this work focuses on the technological development of various AIGC tasks based on their output type, including text, images, videos, 3D content, etc., which depicts the full potential of ChatGPT's future. Moreover, we summarize their significant applications in some mainstream industries, such as education and creativity content. Finally, we discuss the challenges currently faced and present an outlook on how generative AI might evolve in the near future.

CCS Concepts: • **Computing methodologies** → **Computer vision tasks**; *Natural language generation*; Machine learning approaches.

Additional Key Words and Phrases: Survey, Generative AI, AIGC, ChatGPT, GPT-4, GPT-5, Text Generation, Image Generation

## Contents

## 1  INTRODUCTION

Generative AI (AIGC, a.k.a AI-generated content) has made headlines with intriguing tools like ChatGPT or DALL-E [343], suggesting a new era of AI is coming. Under such overwhelming media coverage, the general public are offered many opportunities to have a glimpse of AIGC. However, the content in the media report tends to be biased or sometimes misleading. Moreover, impressed by the powerful capability of ChatGPT, many people are wondering about its limits. Very recently, OpenAI released GPT-4 [307] which demonstrates remarkable performance improvement over the previous variant GPT-3 as well multimodal generation capability like understanding images. Impressed by the powerful capability of GPT-4 powered by AIGC, many are wondering about its limits: can GPT-5 (or other GPT variants) help next-generation ChatGPT unify all AIGC tasks? Therefore, a comprehensive review of generative AI serves as a groundwork to respond to the inevitable trend of AI-powered content creation. More importantly, our work comes to fill this gap in a timely manner.

The goal of conventional AI is mainly to perform classification [263] or regression [227]. Such a discriminative approach renders its role mainly for analyzing existing data. Therefore conventional AI is also often termed analytical AI. By contrast, generative AI differentiates by creating new content. However, generative AI often also requires the model to first understand some existing data (like text instruction) before generating new content [40, 342]. From this perspective, analytical AI can be seen as the foundation of modern generative AI and the boundary between them is often ambiguous. Note that analytical AI tasks also generate content. For example, the label content is generated in image classification [216]. Nonetheless, image recognition is often not considered in the category of generative AI because the label content has low dimensionality. Typical tasks for generative AI involve generating high-dimensional data, like text or images. Such generated content can also be used as synthetic data for alleviating the need for more data in deep learning [144]. An overview of the popularity of generative AI as well as its underlying reasons, is presented in Sec.2.

As stated above, what distinguishes generative AI from conventional one lies in its generated content. With this said, generative AI is conceptually similar to AIGC (a.k.a. AI-generated content) [304]. In the context of describing AI-based content generation, these two terms are often interchangeable. In this work, we call the content generation tasks AIGC for simplicity. For example, ChatGPT is a tool for the AIGC task termed ChatBot [43], which is the tip of the iceberg considering the variety of AIGC tasks. Despite the high resemblance between generative AI and AIGC, these two terms have a nuanced difference. AIGC focuses on the tasks for content generation, while generative AI additionally considers the fundamental technical foundations that support the development of various AIGC tasks. In this work, we divide those underlying techniques into two classes. The first class refers to the generative modeling techniques, like GAN [124] and diffusion model [156], which are directly related to generative AI for content creation. The second class of AI techniques mainly consists of backbone architecture (like Transformer [443]) and self-supervised pretraining (like BERT [87] or MAE [141]). Some of them are developed in the context of analytical AI. However, they have also become essential for demonstrating competitive performance, especially in challenging AIGC tasks. Considering this, both classes of underlying techniques are summarized in Sec.3.

On top of these basic techniques, numerous AIGC tasks have become possible and can be straightforwardly categorized based on the generated content type. The development of various AIGC tasks is summarized in Sec.4, Sec.5 and Sec.6. Specifically, Sec.4 and Sec.5 focus on text output and image output, respectively. For text generation, ChatBot [43] and machine translation [497] are two dominant tasks. Some text generation tasks also take other modalities as the input, for which we mainly focus on image and speech. For image generation, two dominant tasks are image restoration and editing [253]. More recently, text-to-image has attracted significant attention. Beyond the above two dominant output types (*i.e.* text and image), Sec.6 covers other types of output, such as Video, 3D, Speech, etc.

As technology advances, the AIGC performance gets satisfactory for more and more tasks. For example, ChatBot used to be limited to answering simple questions. However, the recent ChatGPT has been shown to understand jokes and generate code under simple instruction. Text-to-image used to be considered a challenging task; however, recent DALL-E 2 [342] and stable diffusion [357] have been able to generate photorealistic images. Therefore, opportunities of applying the AIGC to the industry emerge. Sec.7 covers the application of AIGC in various industries, including entertainment, digital art, media/advertising, education, etc. Along with the application of AIGC in the real world, numerous challenges like ethical concerns have also emerged and they are disused in Sec.8. Alongside the current challenges, an outlook on how generative AI might evolve is also presented.
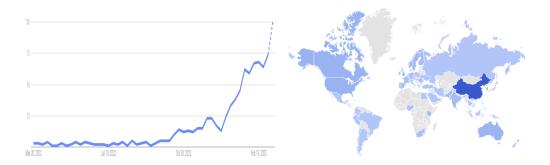
Fig. 1. Search interest of generative AI: Timeline trend (left) and region-wise interest (right). The color darkness on the right part indicates the rank interest level.

Overall, this work conducts a survey on generative AI through the lens of generated content (*i.e.* AIGC tasks), covering its underlying basic techniques, task-wise technological development, application in the industry as well as its social impact. An overview of the paper structure is presented in Figure 4.

## 2 OVERVIEW

Adopting AI for content creation has a long history. IBM made the first public demonstration of a machine translation system at its head office in New York in 1954. The first computer-generated music came out with the name "Illiac Suite" in 1957. Such early attempts and proof-of-concept successes caused a high expectation of the AI future, which motivated governments and companies to invest numerous resources in AI. Such a high boom in investment, however, did not yield the expected output. After that, a period called AI winter came, which dramatically undermines the development of AI and its applications. Entering the 2010s, AI has again become popular again, especially after the success of AlexNet [216] for ImageNet classification in 2012. Entering the 2020s, AI has entered a new era of not only understanding existing data but also creating new content [40, 342]. This section provides an overview of generative AI by focusing on its popularity and why it gets popular.

### 2.1 Popularity indicated by search interest

A good indicator of 'how popular a certain term is' refers to search interest. Google provides a promising tool to visualize search frequency, called Google trends. Although alternative search engines might provide similar functions, we adopt Google trends because Google is one of the most widely used search engines in the world.

**Interest over time and by region.** Figure 1 (left) shows the search interest of generative AI, which indicates that the search interest significantly increased in the past year, especially after October 2022. Entering 2013, this search interest reaches a new height. A similar trend is observed for the term AIGC, see Figure 2 (left). Except for interest over time, Google trends also provides region-wise search interest. The search heatmaps for generative AI and AIGC are shown in Figure 1 (right) and Figure 2 (right), respectively. For both terms, the main hot regions include Asia, Northern America, and Western Europe. Most notably, for both terms, China ranks highest among all countries with a search interest of 100, followed by around 30 in Northern America and 20 in Western Europe. It is worth mentioning that some small but tech-oriented countries also have a very high search interest in generative AI. For example, the three countries that rank top on the country-wise search interest are Singapore (59), Israel (58), and South Korea (43).

Fig. 2. Search interest of AIGC: Timeline trend (left) and region-wise interest (right). The color darkness on the right part indicates the rank interest level.



Fig. 3. Search interest comparison between generative AI and AIGC: Timeline trend (left) and region-wise interest (right).

**Generative AI *v.s.* AIGC.** Figure 3 shows a comparison between generative AI and AIGC for the search interest. Here, we define the interest ratio of generative AI and AIGC as GAI/AIGC. A major observation is that China prefers to use the term AIGC compared with generative AI with the GAI/AIGC ratio being 15/85. By contrast, the GAI/AIGC in the US is 90/10. In many countries, including Russia and Brazil, the GAI/AIGC is 100/0. Overall, most countries prefer generative AI to AIGC, which makes generative AI have an overall higher search interest than AIGC. The reason that China becomes the leading country to adopt the term AIGC is not fully clear. A possible explanation is that AIGC is shortened to a single word and thus is easier to use. We also search the Chinese version of generative AI and AIGC on Google trends, however, the current demonstration is not sufficient.

## 2.2 Why does it get popular?

The recent surging interest in generative AI in the last year can be mainly attributed to the emergence of intriguing tools like Stable diffusion or ChatGPT. Here, we discuss why generative AI gets popular by focusing on what factors contributed to the advent of such powerful AIGC tools. The reasons are summarized from two perspectives: content need and technology conditions.

*2.2.1 Content need.* The way we communicate and interact with the world has been fundamentally changed by the Internet, for which *digital content* plays a key role. Over the last few decades, the content on the web has also

undergone multiple major changes. In the Web 1.0 era (the 1990s-2004), the Internet was primarily used to access and share information, with websites mainly static. There was little interaction between users and the primary mode of communication was one-way, with users accessing information but not contributing or sharing their own content. The content was largely text-based and it was mainly generated by professionals in the relative fields, like journalists generating news articles. Therefore, such content is often called Professional Generated Content (PGC), which has been dominated by another type of content, termed User Generated Content (UGC) [214, 322, 427]. In contrast to PGC, UGC in Web 2.0 [308] is mainly generated by users on social media, like Facebook [203], Twitter [257], Youtube [159], etc. Compared with PGC, the volume of UGC is significantly larger, however, its quality might be inferior.

We are currently transitioning from Web 2.0 to Web 3.0 [363]. With defining features of being decentralized and intermediary-free, Web 3.0 also relies on a new content generation type beyond PGC and UGC to address the trade-off between volume and quality. AI is widely recognized as a promising tool for addressing this trade-off. For example, in the past, only those users that have a long period of practice could draw images of decent quality. With text-to-image tools (like stable diffusion [357]), anyone can create drawing images with a plain text description. Such a combination of user imagination power and AI execution power makes it possible to generate new types of images at an unprecedented speed. Beyond image generation, AIGC tasks also facilitate generating other types of content.

Another change AIGC brings is that the boundary between content consumer and creator becomes vague. In Web 2.0, Content generators and consumers are often different users. With AIGC in Web 3.0, however, data consumers are now able to become data creators, as they are able to use AI algorithms and technology to generate their own original content, and it allows them to have more control over the content they produce and consume, making them use their own data and AI technology to produce content that is tailored to their specific needs and interests. Overall, the shift towards AIGC has the potential to greatly transform the way data is consumed and produced, giving individuals and organizations more control and flexibility in the content they create and consume. In the following, we discuss why AIGC has become popular now.

*2.2.2 Technology conditions.* When it comes to AIGC technology, the first thing that comes into mind is often machine (deep) learning algorithm, while overlooking its two important conditions: data access and compute resources.

**Advances in data access.** Deep learning refers to the practice of training a model on data. The model performance heavily relies on the size of the training data. Typically, the model performance increases with more training samples. Taking image classification as an example, ImageNet [83] with more than 1 million images is a commonly used dataset for training the model and validating the performance. Generative AI often requires an even larger dataset, especially for challenging AIGC tasks like text-to-image. For example, approximately 250M images were used for training DALL-E [343]. DALL-E 2 [342], on the other hand, used approximately 650M images. ChatGPT was built on top of GPT3 [40] partly trained on CommonCrawl dataset, which has 45TB of compressed plaintext before filtering and 570GB after filtering. Other datasets like WebText2, Books1/2, and Wikipedia are involved in the training of GPT3. Accessing such a huge dataset becomes possible mainly due to the Internet.

**Advances in computing resources.** Another important factor contributing to this development of AIGC is advanced in computing resources. Early AI algorithm was run on CPU, which cannot meet the need of training large deep learning models. For example, AlexNet [216] was the first model trained on full ImageNet and the training was done on Graphics Processing Units (GPUs). GPUs were originally designed for rendering graphics in video games but have become increasingly common in deep learning. GPUs are highly parallelized and can perform matrix operations much faster than CPUs. Nvidia is a leading company in manufacturing GPUs. The computing capability of its CUDA has improved

Fig. 4. An overview of generative AI (AIGC): fundamental techniques, core AIGC tasks, and industrial applications.

from the first CUDA-capable GPU (GeForce 8800) in 2006 to the recent GPU (Hopper) with hundreds of times more computing power. The price of GPUs can range from a few hundred dollars to several thousand dollars, depending on the number of cores and memory. Tensor Processing Units (TPUs) are specialized processors designed by Google specifically for accelerating neural network training. TPUs are available on the Google Cloud Platform, and the pricing varies depending on usage and configuration. Overall, the price of computing resources is on the trend of becoming more affordable.

## 3  FUNDAMENTAL TECHNIQUES BEHIND AIGC

In this work, we perceive AIGC as a set of tasks or applications that generates content with AI methods. Before introducing AIGC, we first visit the fundamental techniques behind AIGC, which fall in the scope of generative AI at the technical level. Here, we summarize the fundamental techniques by roughly dividing them into two classes: Generative techniques and Creation techniques. Specifically, Creation techniques refer to the techniques that are able to generate various contents, e.g., GAN and diffusion model. Meanwhile, General techniques cannot generate content directly but are essential for the development of AIGC, e.g., the Transformer architecture. In this section, we provide a brief summary of the required techniques for AIGC.

### 3.1  General techniques in AI

After the phenomenal success of AlexNet [216], there is a surging interest in deep learning, which somewhat becomes a synonym for AI. In contrast to traditional rule-based algorithms, deep learning is a data-driven method that optimizes the model parameters with a stochastic gradient. The success of deep learning in obtaining a superior feature representation depends on better backbone architecture and more data, which greatly accelerates the development of AIGC.

*3.1.1  Backbone architecture.*  As two mainstream fields in deep learning, the research on natural language processing (NLP) and computer vision (CV) have significantly improved the backbone architectures and inspired various applications of improved backbones in other fields, e.g., the speech area. In the NLP field, Transformer [443] has replaced recurrent neural networks (RNN) [281, 285] to be the de-facto standard backbone. In the CV area, vision Transformer (ViT) [97] has also shown its power besides the traditional convolutional neural networks (CNN). Here, we will briefly introduce how these mainstream backbones work and their representative variants.

**RNN architecture.** RNN is mainly adopted for handling data with time sequences, like language or audio. A vanilla RNN has three layers: input, hidden, and output. The information flow in RNN is in two directions. The first direction is from the input to the hidden layer and then to the output. What captures the *recurrent* nature of RNN lies in its second information flow in the time direction. Except for the corresponding input, the current hidden state depends at time $t$ depends on the hidden state at time $t - 1$. This two-flow design well handles the sequence order but suffers from exploding or vanishing gradients when the sequence gets long. To mitigate long-term dependency, LSTM [158] was introduced with a cell state that acts like a freeway to facilitate the information flow in the sequence direction. LSTM is one of the most popular methods for alleviating the gradient vanishing/exploding issue. With three types of gates, however, LSTM suffers from high complexity and a higher memory requirement. Gated Recurrent Unit (GRU) [65] simplifies LSTM by merging its cell and hidden states and replacing the forget and input gates with a so-called update state. Unitary RNN [18] handles the gradient issue by implementing unitary matrices. Gated Orthogonal Recurrent Unit [184] leverages the merits of both gate and unitary matrices. Bidirectional RNN [376] improves vanilla RNN by capturing both past and future information in the cell, i.e., the state at time $t$ is calculated based on both time $t - 1$ and $t + 1$. Depending on the tasks, RNN can have various architectures with a different number of inputs and outputs: one-to-one, many-to-one, one-to-many, and many-to-many. The many-to-many can be used in machine translation and is also called the sequence-to-sequence (seq2seq) model [413]. Attention was introduced in [24] to make the model decoder see every encoder token and automatically decide the weights on them based on their importance.

**Transformer.** Different from Seq2seq with attention [24, 267, 315], a new variant of architecture discards the seq-2seq architecture and claims that attention is all you need [443]. Such attention is called self-attention, and the proposed architecture is termed *Transformer* [443] (see Figure 5). A standard Transformer consists of an encoder and a

decoder and is developed based on residual connection [143] and layer normalization [22]. Except for the Add & Norm module, the Transformer has two core components: multi-head attention and feed-forward neural network (a.k.a. MLP). The attention module adopts a multi-head design with the self-attention in the form of scaled dot-product defined as:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \tag{1}$$

Unlike RNNs, which build positional information by sequentially inputting sentence information, Transformer obtains powerful modeling capabilities by constructing global dependencies but also loses information with positional bias. Therefore, positional encoding is needed to enable the model to sense the positional information of the input signal. There are two types of positional encoding. Fixed position coding is represented by sinusoids and cosines of different frequencies. The learnable position encoding is composed of a set of learnable parameters. Transformer has become the de-facto standard method in NLP tasks.



Fig. 5. Transformer structure (figure obtained from [443]).

**CNN architecture.** After introducing RNN and Transformer in NLP field, we start to visit two mainstream backbones in CV area, i.e., CNN and ViT. CNNs have become a standard backbone in the field of computer vision. The core of CNN lies in its convolution layer. The convolution kernel (also known as filter) in the convolution layer is a set of shared weight parameters for operating on images, which is inspired by the biological visual cortex cells. The convolution kernel slides on the image and performs correlation operations with the pixel values on the image, finally obtaining

the feature map and realizing the feature extraction of the image. GoogleNet[417], with its Inception module allowing multiple convolutional filter sizes to be chosen in each block, increased the diversity of convolutional kernels, thus the performance of CNN was improved. ResNet[143] was a milestone for CNNs, introducing residual connections that stabilized training and enabled the models to achieve better performance through deeper modeling. After that, it became part of the binding in CNNs. In order to expand the work of ResNet, DenseNet[165] establishes dense connections between all the previous layers and the subsequent layers, thus enabling the model to have better modeling ability. EfficientNet[418] uses a scaling method which uses a set of fixed scaling coefficients to uniformly scale the width, depth, and resolution of the convolutional neural network architecture, thus making the model more efficient.



Fig. 6. ViT structure (figure obtained from [97]).

**ViT architecture.** Inspired by the success of Transformer in NLP, numerous works have tried to apply Transformer to the field of CV with ViT[97] (see Figure 6), being the first of its kind. ViT first flattens the image into a sequence of 2D patches and inserts a class token at the beginning of the sequence to extract classification information. After the embedding position encoding, the token embeddings are fed into a standard Transformer. This simple and effective implementation of ViT makes it highly scalable. Swin [261] efficiently deals with image classification and dense recognition tasks by constructing hierarchical feature maps by merging image blocks at a deeper level, and due to its computation of self-attention only within each local window, it reduces computational complexity. DeiT[430] uses the teacher-student strategy for training, reducing the dependence of Transformer models on large data, by introducing distillation tokens. CaiT[431] introduces class attention to effectively increase the depth of the model. T2T[508] effectively localizes the model by Token Fusion and introduces hierarchical deep and narrow structures through the prior of CNNs by recursively aggregating adjacent Tokens into one Token. Through permutation equivariance, Transformers have liberated CNNs from their translation invariance, allowing for long-range dependencies and less inductive bias, making them more powerful modeling tools and better transferable to downstream tasks than CNNs. In

the current paradigm of large models and large datasets, Transformers have gradually replaced CNNs as the mainstream model in the field of computer vision.

*3.1.2 Self-supervised pretraining.* Parallel to better backbone architecture, deep learning also benefits from self-supervised pertaining which can exploit a larger (unlabeled) training dataset. Here, we summarize the most relevant pretraining techniques to AIGC, and categorize them according to the training data type (e.g., language, vision, and joint pretraining).

**Language pretraining.** There are three major types of language pretraining methods. The first type pretrains an encoder with masking, for which the representative work is BERT [87] (see Figure 7). Specifically, BERT predicts the masked language tokens from the unmasked tokens. There is a significant discrepancy between the mask-then-predict pertaining task and downstream tasks, therefore masked language modeling like BERT is rarely used for text generation without finetuning. By contrast, autoregressive language pretraining methods are suitable for few-shot or zero-shot text generation. GPT family [40, 338, 339] is the most popular one which adopts a decoder instead of an encoder. Specifically, GPT-1 [338] is the first of its kind with GPT-2 [339] and GPT-3 [40] further investigating the role of massive data and large model in the transfer capacity. Based on GPT-3, the unprecedented success of ChatGPT has attracted great attention recently. Moreover, a stream of language models adopts both an encoder and decoder as the original Transformer. BART [226] perturbed the input with various types of noise and predicted the original clean input, like a denoising autoencoder. MASS [400] and PropheNet [332] follow BERT to take a masked sequence as the input of the encoder with the decoder predicting the masked tokens in an autoregressive manner. T5 [340] replaces the masked tokens with some random tokens.



Fig. 7. BERT structure (figure obtained from [87]).

**Visual pretraining.** To learn better representations of vision data during pretraining, self-supervised learning (SSL) has been widely applied, and we term it visual SSL. Visual SSL has undergone three stages. Early works focused on designing various pretext tasks like jigsaw puzzles [303] or predicting rotation [121]. Such pretraining yields better performance on the downstream task than training from scratch, which motivates contrastive learning methods [54, 142, 520]. Contrastive learning adopts joint embedding to minimize the representation distance between augmented

images for learning augmentation-invariant representation. The representation in pure joint embedding can collapse to a constant regardless of the inputs, for which contrastive learning simultaneously maximizes the representation distance from negative samples. Negative-free joint-embedding methods have also been investigated in SimSiam [55] and BYOL [129]. How SimSiam works without negative samples have been investigated in [521]. Inspired by the success of BERT in NLP for pertaining, BEiT [30] applied masking modeling in vision and its success relies on a pre-trained VAE to obtain the visual token. Masked autoencoder (MAE) [141] (see Figure 8) simplifies it to an end-to-end denoising framework by predicting the masked patches from the unmasked patches. Outperforming contrastive learning and negative-free joint-embedding methods, MAE has become a new variant of the visual SSL framework. Interested readers can refer [519] for more details.



Fig. 8.  MAE structure (figure obtained from [141]).

**Joint pretraining.** With large datasets of image-text pairs collected from the Internet, multimodal learning [29, 487] has made unprecedented progress to learn data representations, at the front of which is cross-modal matching [115]. Contrastive pretraining is widely used to match the image embedding and text encoding in the same representation space [180, 336, 507]. CLIP [336] (see Figure 9 is a pioneering work in this direction and is used in numerous text-to-image models, such as DALL-E 2 [342], Upainting [241], DiffusionCLIP [206]. ALIGN [180] extended CLIP with noisy text supervision so that the text-image dataset requires no cleaning and can be scaled to a much larger size (from 400M to 1.8B). Florence [507] further expands the cross-modal shared representation from coarse scene to dine object and from static images to dynamic videos, etc. Therefore, the learned shared representation is more universal and shows superior performance [507].

### 3.2  Creation techniques in AI

Deep generative models (DGMs) are a group of probabilistic models that use neural networks to generate samples. Early attempts at generative modeling focused on pre-training with an autoencoder [28, 154, 365]. A variant of autoencoder with masking has emerged to become a dominant self-supervised learning framework, and interested readers are

Fig. 9. CLIP structure (figure obtained from [336]).

encouraged to check a survey on masked autoencoder [519]. Unless specified, the use cases of deep generative models in this survey only consider generating new data. The generated data is typically high-dimensional, and therefore, predicting a label of a sample is not considered discriminative instead of generative modeling even though something like a label is also technically generated.

Numerous DGMs have emerged and can be categorized into two major groups: likelihood-based and energy-based. Likelihood-based probabilistic models, like autoregressive models [126] and flow models [90], have a tractable likelihood which provides a straightforward method to optimize the model weights w.r.t. the log-likelihood of the observed (training) data. The likelihood is not fully tractable in variational autoencoders (VAEs) [210], but a tractable lower bound can be optimized, thus VAE is also considered to lie in the likelihood-based group which specifies a normalized probability. By contrast, energy-based models [128, 153] are featured by the unnormalized probability, a.k.a. energy function. Without the constraint on the tractability of the normalizing constant, energy-based models are more flexible in parameterizing but difficult to train [403]. Notably, GAN and diffusion models are highly related to energy-based models even though are developed from different motivations. In the following, we present an introduction to each class of likelihood-based models, followed by how the energy-based models can be trained as well as the mechanism behind GAN and diffusion models.

*3.2.1 Likelihood-based models.* **Autoregressive models.** Autoregressive models learn the joint distribution of sequential data and predict each variable in the sequence with previous time-step variables as inputs. As shown in Eq. 2, autoregressive models assumes that the joint distribution $p_\theta(x)$ can be decomposed to a product of conditional distributions.

$$p_\theta(x) = p_\theta(x_1)p_\theta(x_2|x_1)...p_\theta(x_n|x_1, x_2, ..., x_{n-1}), \tag{2}$$

Although both rely on previous timesteps, autoregressive models differ from RNN architecture since the previous timesteps are given to the model as input instead of hidden states in RNN. In other words, autoregressive models can be seen as a feed-forward network that takes all the previous time-step variables as inputs. Early works model discrete data with different functions estimating the conditional distribution, e.g. logistic regression in Fully Visible Sigmoid Belief Network (FVSBN) [114] and one hidden layer neural networks in Neural Autoregressive Distribution Estimation

(NADE) [221]. The following research further extends to model the continuous variables [437, 438]. Autoregressive methods have been widely applied in multiple areas, including computer vision (PixelCNN [441] and PixelCNN++ [373]), audio generation (WaveNet [440]), natural language processing (Transformer [443]).

**VAE.** Autoencoders are a family of models that first map the input to a low-dimension latent layer with an encoder and then reconstruct the input with a decoder. The entire encoder-decoder process aims to learn the underlying data patterns and generate unseen samples [310]. Variational autoencoder (VAE) [210] is an autoencoder that learns the data distribution $p(x)$ from latent space z, i.e., $p(x) = p(x|z)p(z)$, where $p(x|z)$ is learned by the decoder. In order to obtain $p(z)$, VAE [210] adopts Bayes' theorem and approximates the posterior distribution $p(z|x)$ by the encoder. The VAE model is optimized toward a likelihood goal with regularizer [13].

*3.2.2 Energy-based models.* With a tractable likelihood, autoregressive models and flow models allow a straightforward optimization of the parameters w.r.t. the log-likelihood of the data. This forces the model to be constrained in a certain form. For example, the autoregressive model needs to be factorized as a product of conditional probabilities, and the flow model must adopt invertible transformation.

Energy-baed models specify probability up to an unknown normalizing constant, therefore, they are also known as non-normalzied probabilistic models. Without losing generality by assuming the energy-based model is over a single variable $x$, we denote its energy as $E_\theta(x)$. Its probability density is then calculated as

$$p_\theta(x) = \frac{\exp(-E_\theta(x))}{z_\theta} \tag{3}$$

where $z_\theta$ is the so-called normalizing constant and defined as $z_\theta = \int \exp(-E_\theta(x)) \, dx$. $z_\theta$ is an intractable integral, making optimizing energy-based models a challenging task.

**MCMC and NCE.** Early attempts at optimizing energy-based models opt to estimate the gradient of the log-likelihood with Markov chain Monte Carlo (MCMC) approaches, which require a cumbersome drawing of random samples. Therefore, some works aim to improve the efficiency of MCMC a representative work Langevin MCMC [128, 316]. Nonetheless, performing MCMCM to obtain requires large computation and contrastive divergence (CD) [153] is a popular method to reduce the computation via approximation with various variants: persistent CD [425], mean field CD [468], and multi-grid CD [116]. Another line of work optimizes energy-based models via notice contrastive estimation (NCE) [137], which contrasts the probabilistic model with another noise distribution. Specifically, it optimizes the following loss:

$$\mathbb{E}_{p_d}\left[\ln \frac{p_\theta(x)}{p_\theta(x) + q_\phi(x)}\right] + \mathbb{E}_{q_\phi}\left[\ln \frac{q_\phi(x)}{p_\theta(x) + q_\phi(x)}\right], \tag{4}$$

**Score matching.** For optimizing energy-based models, another popular MCMC-free method minimizes the derivatives of log probability density between the model and the observed data. The first-order of a log probability density function is called *score* of the distribution ($s(x) = \nabla_x \log p(x)$), therefore, this method is often termed *score matching*. Unfortunately, the data score function $s_d(x)$ is unavailable. Various attempts [314, 374, 389, 401, 402, 446] have been made to mitigate this issue, with a representative method called denoising score matching [446]. Denoising score matching approximates the score of data with noisy samples. The model takes a noisy sample as the input and predicts its noise. Therefore, it can be used for sampling clean samples from noise by iterative removing the noise [374, 401].

*3.2.3 Two star-models: from GAN to diffusion model.* When it comes to deep generative models, what first comes to your mind? The answer depends on your background, however, GAN is definitely one of the most mentioned models.

GAN stands for generative adversarial network [124] which was first proposed by Ian J. Goodfellow and his team in 2014 and rated as "the most interesting idea in the last 10 years in machine learning" by Yann Lecun in 2016. As the pioneering work to generate images of reasonably high quality, GAN has been widely regarded as a de facto standard model for the challenging task of image synthesis. This long-time dominance has been recently challenged by a new family of deep generative models termed diffusion models [156]. The overwhelming success of diffusion models starts from image synthesis but extends to other modalities, like video, audio, text, graph, etc. Considering their dominant influence in the development of generative AI, we first summarize GAN and diffusion models before introducing other families of deep generative models.

**GAN.** The architecture of GAN is shown in Figure 10. GAN is featured by its two network components: a discriminator ($\mathcal{D}$) and a generator ($\mathcal{G}$). $\mathcal{D}$ distinguishes real images from those generated by $\mathcal{G}$, while $\mathcal{G}$ aims to fool $\mathcal{D}$. Given a latent variable $z \sim p_z$, the output of $\mathcal{G}$ is $\mathcal{G}(z)$ constituting a probability distribution $p_g$. The goal of GAN is to make $p_g$ approximate the observed data distribution $p_{data}$. This objective is achieved through adversarial learning, which can be interpreted as a min-max game [375]:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{data}} \log[D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} \log\left[1 - D(G(\mathbf{z}))\right]. \tag{5}$$

where $\mathcal{D}$ is trained to maximize the probability of assigning correct labels to real images and generated ones, and is used to guide the optimization of $\mathcal{G}$ towards generating more real images. GANs have the weakness of potentially unstable training and less diversity in generation due to their adversarial training nature. The basic difference between GANs and autoregressive models is that GANs learn implicit data distribution, whereas the latter learns an explicit distribution governed by a prior imposed by model structure.



Fig. 10. A schematic of GAN structure.

**Diffusion model.** The use of diffusion models, a special form of hierarchical VAEs, has seen explosive growth in the past few years [45, 73, 245, 320, 435]. Diffusion models (Figure 11) are also known as denoising diffusion probabilistic models (DDPMs) or score-based generative models that generate new data similar to the data on which they are trained [156]. Inspired by non-equilibrium thermodynamics, DDPMs can be defined as a parameterized Markov chain

of diffusion steps to slowly add random noise to the training data and learn to reverse the diffusion process to construct desired data samples from the pure noise.



Fig. 11. Diffusion model for image generation (figure obtained from [156]).

In the forward diffusion process, DDPM destroys the training data through the successive addition of Gaussian noise. Given a data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, DDPM maps the training data to noise by gradually perturbing the input data. This is formally achieved by a simple stochastic process that starts from a data sample and iteratively generates noisier samples $\mathbf{x}_T$ with $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$, using a simple Gaussian diffusion kernel:

$$q(x_{1:T}|x_0) := \prod_{t=1}^{T} q(x_t|x_{t-1}), \tag{6}$$

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I) \tag{7}$$

where $T$ and $\beta_t$ are the diffusion steps and hyper-parameters, respectively. We only discuss the case of Gaussian noise as transition kernels for simplicity, indicated as $\mathcal{N}$ in Eq. 7. With $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s$, we can obtain noised image at arbitrary step $t$ as follows:

$$q(x_t|x_0) := \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I) \tag{8}$$

During the reverse denoising process, DDPM is learning to recover the data by reversing the noising process i.e., it undoes the forward diffusion by performing the iterative denoising. This process represents data synthesis and DDPM is trained to generate data by converting random noise into real data. It is also formally defined as a stochastic process, which iteratively denoises the input data starting from $p_\theta(T)$ and generates $p_\theta(x_0)$ which can follow the true data distribution $q(x_0)$. Therefore, the optimization objective of the model is as follows:

$$E_{t\sim\mathcal{U}(1,T),\mathbf{x}_0\sim q(\mathbf{x}_0),\epsilon\sim\mathcal{N}(\mathbf{0},\mathbf{I})}\lambda(t)\,\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \tag{9}$$

Both the forward and reverse processes of DDPMs often use thousands of steps for gradual noise injection and during generation for denoising.

## 4  AIGC TASK: TEXT GENERATION

NLP studies natural language with two fundamental tasks: understanding and generation. These two tasks are not exclusively separate because the generation of an appropriate text often depends on the understanding of some text inputs. For example, language models often transform a sequence of text into another, which constitutes the core task of text generation, including machine translation, text summarization, and dialogue systems. Beyond this, text generation evolves in two directions: controllability and multi-modality. The first direction aims to make the generated content

## 4.1 Text to text

*4.1.1 Chatbots.* The main task of the dialogue system (chatbots) is to provide better communication between humans and machines [85, 299]. According to whether the task is specified in the applications, dialogue system can be divided into two categories : (1) task-oriented dialogue systems (TOD) [323, 502, 533] and (2) open-domain dialogue systems (OOD) [4, 532, 541]. Specifically, the task-oriented dialogue systems focus on task completion and solve specific problems (e.g., restaurant reservations and ticket booking) [533]. Meanwhile, open-domain dialogue systems are often data-driven and aim to chat with humans without task or domain restrictions [353, 533].

**Task-oriented systems.** Task-oriented dialogue systems can be divided into modular and end-to-end systems. The modular methods include four main parts: natural language understanding (NLU) [395, 409], dialogue state tracking (DST) [382, 462], dialogue policy learning (DPL) [169, 483], and natural language generation (NLG) [25, 99]. After encoding the user inputs into semantic slots with NLU, DST, and DPL decide the next action that is then converted to natural language by NLG as the final response. These four modules aim to generate responses in a controllable way and can be optimized individually. However, some modules may not be differentiable, and the improvement of a single module may not lead to the improvement of the whole system [533]. To solve these problems, end-to-end methods either achieve an end-to-end training pipeline by making each module differentiable [139, 162], or use a single end-to-end module in the system [498, 531]. There still exist several challenges for both modular and end-to-end systems, including how to improve tracking efficiency for DST [208, 312] and how to increase the response quality of end-to-end system with limited data [145, 148, 282].



Fig. 12. A diagram illustrating the three steps of how ChatGPT is trained by OpenAI (figure obtained from [311]).

**Open-domain systems.** Open-domain systems aim to chat with users without task and domain restrictions [353, 533], and can be categorized into three types: retrieval-based systems, generative systems, and ensemble systems [533]. Specifically, retrieval-based systems always find an existing response from a response corpus, while generative systems can generate responses that may not appear in the training set. Ensemble systems combine retrieval-based and generative methods by either choosing the best response or refining the retrieval-based model with generative one [378, 533, 546]. Previous works improve the open-domain systems from multiple aspects, including dialogue context modeling [105, 181, 250, 282], improving the response coherence [9, 117, 251, 483] and diversity [31, 211, 335, 408]. Most recently, ChatGPT (see Figure 12) has achieved unprecedented success and also falls into the scope of open-domain dialogue systems. Apart from answering various questions, ChatGPT can also be used for paper writing, code debugging, table generation, and to name but a few.



Fig. 13. An example of machine translation (figure obtained from [39]).

*4.1.2 Machine translation.* As the term suggests, machine translation automatically translates the text from one language to another [171, 497] (see Figure 13). With deep learning replacing rule-based [108] and statistical [212, 213] methods, neural machine translation (NMT) requires minimum linguistic expertise [399, 451] and has become a mainstream approach featured by its higher capacity in capturing long dependency in the sentence [62]. The success of neural machine learning can be mainly attributed to language models [34], which predicts the probability of a word conditioned on previous ones. Seq2seq [413] is a pioneering work to apply encoder-decoder RNN structure [191] to machine translation. When the sentence gets long, the performance of Seq2seq [413] deteriorates, for which an attention mechanism was proposed in [24] to help translate the long sentence with additional word alignment. With increasing attention, in 2006, Google's NMT system helped reduce the translation effort of humans by around 60% compared to Google's phrase-based production system, which bridges the gap between Human and machine translation [475]. CNN-based architectures have also been investigated for NMT with numerous attempts [190, 192], but fail to achieve

comparable performance as the RNN boosted by attention [24]. Convolutional Seq2seq [120] makes CNN compatible with the attention mechanism, showing CNN can achieve comparable or even better performance than RNN. However, this improvement was later outperformed by another architecture termed Transformer [443]. With RNN or Transformer as the architecture, NMT often utilizes autoregressive generative model, where a greedy search only considers the word with the highest probability for predicting the next work during inference.

A trend for NMT is to achieve satisfactory performance in low-resource setup, where the model is trained with limited bilingual corpus [458]. One way to mitigate this data scarcity is to utilize auxiliary languages, like multilingual training with other language pairs [187, 383, 547] or pivot translation with English as the middle pivot language [58, 350]. Another popular approach is to utilize pre-trained language models, like BERT [87] or GPT [338]. For example, it is shown in [359] that initializing the model weights with BERT [87] or RoBERTa [259] significantly improves the English-German translation performance. Without the need for fine-tuning, GPT-family models [40, 338, 339] also show competitive performance. Most recently, ChatGPT has shown its power in machine translation, performing competitively with commercial products (e.g., Google translate) [182].

## 4.2 Multimodal text generation

*4.2.1 Image-to-text.* Image-to-text, also known as image captioning, refers to describing a given image's content in natural language (see Figure 14). A seminal work in this area is Neural Image Caption (NIC) [447], which employs CNN as an encoder to extract high-level representations of input images and then feed these representations into an RNN decoder to generate image descriptions. This two-step encoder-decoder architecture has been widely applied in later works on image captioning, and we term them as visual encoding [407] and language decoding, respectively. Here, we first revisit the history and recent trends of both stages in image captioning.



Fig. 14. An example of image captioning (figure obtained from [109]).

**Visual encoding.** Extracting an effective representation of images is the main task of visual encoding module. Start from NIC [447] with GoogleNet [417] extracting the global feature of input image, multiple works adopt various CNN backbones as the encoder, including AlexNet [216] in [195] and VGG network [393] in [92, 272]. However, it is hard for a language model to generate fine-grained captions with global visual features. Following works introduce attention

mechanism for fine-grained visual features, including attention over different grids of CNN features [56, 264, 463, 484] or over different visual regions [16, 200, 518]. Another branch of work [500, 536] adopts graph neural networks to encode the semantic and spatial relationships between different regions. However, the human-defined graph structures may limit the interactions among elements [407], which can be mitigated by the self-attention methods [231, 501, 530] (including ViT [256]) that connects all the elements.

**Language decoding.** In image captioning, a language decoder generates captions by predicting the probability of a given word sequence [407]. Inspired by the breakthroughs in the NLP area, the backbones of language decoders evolve from RNN [200, 264, 447, 456] to Transformer [132, 149, 231], achieving significant performance improvement. Beyond the visual encoder-language decoder architecture, a branch of work adopts BERT-like architecture that fuses the image and captions in the early stage of a single model [244, 526, 542]. For example, [542] adopts a single encoder to learn a shared space for image and text, which is first pre-tained on large image-text corpus and finetuned, specifically for image captioning tasks.

*4.2.2 Speech-to-Text.* Speech-to-text generation, also known as automatic speech recognition (ASR), is the process of converting spoken language, specifically a speech signal, into a corresponding text [173, 347] (see Figure 15). With many potential applications such as voice dialing, computer-assisted language learning, caption generation, and virtual assistants like Alexa and Siri, ASR has been an exciting field of research [194, 270, 345] since the 1950s, and evolved from hidden Markov models (HMM) [188, 225] to DNN-based systems [75, 127, 152, 297, 473].



Fig. 15. A example of speech recognition (figure obtained from [46]).

**Various research topics and challenges.** Previous works improved ASR systems in various aspects. Multiple works discuss different feature extraction methods for speech signals [270], including temporal features (e.g., discrete wavelet transform [287, 419]) and spectral features such as the most commonly used mel-frequency cepstral coefficients (MFCC) [61, 69, 429]. Another branch of work improves the system pipeline [355] from multi-model [268] to end-to-end ones [161, 233, 234, 296, 453]. Specifically, a multi-model system [268, 270] first learns an acoustic model (e.g., a phoneme classifier that maps the features to phonemes) and then a language model for the word outputs [355]. On the other hand, end-to-end models directly predict the transcriptions from the audio input [161, 233, 234, 296, 453]. Although end-to-end models achieve impressive performance in various languages and dialects, many challenges still exist. First, their applications for under-resourced speech tasks remain challenging as it is costly and time-consuming to acquire vast amounts of annotated training data [104, 355]. Second, these systems may struggle to handle speech with specialized out-of-vocabulary words and may perform well on the training data but may not generalize well to new or unseen data [104, 334]. Moreover, biases in the training data can also affect the performance of supervised ASR systems, leading to poor accuracy on certain groups of people or speech styles [35].

Fig. 16. Examples of image restoration (figure obtained from [452]).

**Under-resourced speech tasks.** Researchers work on new technologies to overcome challenges in ASR systems, among which we mainly discuss the under-resourced speech problem that lacks data for impaired speech [355]. A branch of work [321, 346] adopts multi-task learning to optimize a shared encoder for different tasks. Meanwhile, self-supervised ASR systems have recently become an active area of research without relying on a large number of labeled samples. Specifically, self-supervised ASR systems first pre-train a model on huge volumes of unlabeled speech data, then fine-tune it on a smaller set of labeled data to facilitate the efficiency of ASR systems. It can be applied for low-resource languages, handling different speaking styles or noise conditions, and transcribing multiple languages [23, 71, 255, 492].

## 5 AIGC TASK: IMAGE GENERATION

Similar to text generation, the task of image synthesis can also be categorized into different classes based on its input control. Since the output is images, a straightforward type of control is images. Image-type control induces numerous tasks, like super-resolution, deblur, editing, translation, etc. A limitation of image-type control is the lack of flexibility. By contrast, text-guided control enables the generation of any image content with any style at the free will of humans. Text-to-image falls into the category of cross-modal generation, since the input text is a different modality from the output image.

### 5.1 Image-to-image

*5.1.1 Image restoration.* Image restoration solves a typical inverse problem that restores clean images from their corresponding degraded versions, with examples shown in Figure 16. Such an inverse problem is non-trivial with its ill-posed nature because there are infinite possible mappings from the degraded image to the clean one. There are two sources of degradation: missing information from the original image and adding something undesirable to the clean image. The former type of degradation includes capturing a photo with a low resolution and thus losing some detailed information, cropping a certain region, and transforming a colorful image to its gray form. Restoration tasks recover them in order are image super-resolution, inpainting, and colorization, respectively. Another class of restoration tasks aims to remove undesirable perturbations, like denoise, derain, dehaze, deblur, etc. Early restoration techniques primarily use mathematical and statistical modeling to remove image degradations, including spatial filters for denoising [123, 392, 529],

kernel estimation for deblurring [485, 489]. Lately, deep learning-based methods [42, 59, 93, 177, 248, 252, 481, 486] have become predominant in image restoration tasks due to their versatility and superior visual quality over their traditional counterparts. CNN is widely used as the building block in image restoration [94, 411, 442, 459], while recent works explore more powerful transformer architecture and achieve impressive performance in various tasks, such as image super-resolution [247], colorization [218], and inpainting [240]. There are also works that combine the strength of CNNs and Transformers together [103, 534, 535].

**Generative methods for restoration.** Typical image restoration models learn a mapping between the source (degraded) and target (clean) images with a reconstruction loss. Depending on the task, training data pairs can be generated by degrading clean images with various perturbations, including resolution downsampling and grayscale transformation. To keep more high-frequency details and create more realistic images, generative models are widely used for restoration, such as GAN in super-resolution [223, 460, 528] and inpainting [42, 252, 298]. However, GAN-based models typically suffer from a complex training process and mode collapse. These drawbacks and the massive popularity of DMs led numerous recent works to adopt DMs for image restoration tasks [199, 232, 265, 349, 367, 369]. Generative approaches like GAN and DM can also produce multiple variations of clean output from a single degraded image.

**From single-task to multi-task.** A majority of existing restoration approaches train separate models for different forms of image degradation. This limits their effectiveness in practical use cases where the images are corrupted by a combination of degradations. To address this, several studies [6, 207, 391, 540] introduce multi-distortion datasets that combine various forms of degradation with different intensities. Some studies [207, 258, 505, 509] propose restoration models in which different sub-networks are responsible for different degradations. Another line of work [228, 242, 391, 410, 540] relies on attention modules or a guiding sub-network to assist the restoration network through different degradations, allowing a single network to handle multiple degradations.

*5.1.2 Image editing.* In contrast to image restoration for enhancing image quality, image editing refers to modifying an image to meet a certain need like style transfer (see Figure 17). Technically, some image restoration tasks like colorization might also be perceived as image editing by perceiving adding color as the desired need. Modern cameras often have basic editing features such as sharpness adjustments [524], automatic cropping [525], red eye removal [396], etc. However, in AIGC, we are more interested in advanced image editing tasks that change the image semantics in various forms, such as content, style, object attributes, etc.

A family of image editing targets to modify the attributes (like age) of the main object (like a face) in the image. A typical use case is facial attribute editing which can change the hairstyle, age, or even gender. Based on a pre-trained CNN encoder, a line of pioneering works adopt optimization-based approaches [236, 436], which is time-consuming due to its iterative nature. Another line of works adopts learning-based approaches to directly generate the image, with a trend from single attribute [237, 385] to multiple ones [146, 209, 478]. A drawback of most aforementioned methods is the dependence on annotated labels for attributes, therefore, unsupervised learning has been introduced to disentangle different attributes [60, 386].

Another family of image editing changes the semantics by combining two images. For example, image morphing [185] interpolates the content of two images, while style transfer [119] yields a new image with the content of one image and the style of the other. A naive method for image morphing is to perform interpolation in the pixel space, which causes obvious artifacts. By contrast, interpolating in the latent space can consider the view change and generate a smooth image. The latent space for those two images can be obtained via GAN inversion method [477]. Numerous works [1, 490, 544, 545] have explored the latent place of a pre-trained GAN for image morphing. For the task of style

Fig. 17. Examples of style transfer as a form of image editing (figure obtained from [118]).

transfer, a specific style-based variant of GAN termed StyleGAN [197] is a popular choice. From the earlier layers to the latter ones, StyleGAN controls the attributes from coarser-grained (like structure) to finer-grained ones (like texture). Therefore, StyleGAN can be used for style transfer by mixing the earlier layer's latent representation of the content image and the latter layer's latent representation of the style image [1, 131, 444, 467].

Compared with restoration tasks, various editing tasks enable a more flexible image generation. However, its diversity is still limited, which is alleviated by allowing other text as the input. More recently, image editing based on diffusion models has been widely discussed and achieved impressive results [48, 150, 206, 450]. DiffusionCLIP [206] is a pioneering work that finetunes a pre-trained diffusion model to align the target image and text. By contrast, LDEdit [48] avoids finetuning based on LDM [357]. A branch of works discusses the mask problem in image editing, including how to

connect a manually designed masked region and background seamlessly [3, 19, 21, 21]. On the other hand, DiffEdit [72] proposes to predict the mask automatically that indicates which part to be edited. There are also works editing 3D objects based on diffusion models and text guidance [47, 205, 230].

## 5.2    Multimodal image generation

*5.2.1    Text-to-image.* Text-to-image (T2I) task aims to generate images from textual descriptions (see Figure **??**.), and can be traced back to image generation from tags or attributes [405, 495]. AlignDRAW [271] is a pioneering work to generate images from natural language, and it is impressive that AlignDRAW [271] can generate images from novel text like 'a stop sign is flying in blue skies'. More recently, advances in text-to-image area can be categorized into three branches, including GAN-based methods, autoregressive methods, and diffusion-based methods.

**GAN-based methods.** The limitation of AlignDRAW [271] is that the generated images are unrealistic and require an additional GAN for post-processing. Based on a deep convolutional generative adversarial network (DC-GAN) [337], [348] is the first end-to-end differential architecture from the character level to the pixel level. To generate high-resolution images while stabilizing the training process, StackGAN [522] and StackGAN++ [523] propose a multi-stage mechanism that multiple generators produce images of different scales, and high-resolution image generation is conditioned on the low-resolution images. Moreover, AttnGAN [488] and Controlgan [229] adopt attention networks to obtain fine-grained control on the subregions according to relevant words.

**Autoregressive methods.** Inspired by the success of autoregressive Transformers [443], a branch of works generates images in an auto-regressive manner by mapping images to a sequence of tokens, among which DALL-E [343] is a pioneering work. Specifically, DALL-E [343] first converts the images to image tokens with a pre-trained discrete variational autoencoder (dVAE), then trains an auto-regressive Transformer to learn the joint distribution of text and image tokens. A concurrent work CogView [88] independently proposes the same idea with DALL-E [343] but achieves superior FID [151] than DALL-E [343] on blurred MS COCO dataset. CogView2 [89] extends CogView [88] to various tasks, e.g., image captioning, by masking different tokens. Parti [504] further improves the image quality by scaling the model size to 20 billion.

**Diffusion-based methods.** Diffusion model-based methods have achieved unprecedented success and attention recently, which can be categorized by either working on the pixel space directly [300, 368] or the latent space [342, 357]. GLIDE [300] outperforms DALL-E by extending class-conditional diffusion models to text-conditional settings, while Imagen [368] improves the image quality further with a pre-trained large language model (e.g., T5) capturing the text semantics. To reduce resource consumption of diffusion models in pixel space, Stable Diffusion [357] first compresses the high-resolution images to a low-dimensional latent space, then trains the diffusion model in the latent space. This method is also known as Latent Diffusion Models (LDM) [357]. Different from Stable Diffusion [357] that learns the latent space based on only images, DALL-E2 [342] applies diffusion model to learn a prior as alignment between image space and text space of CLIP. Other works also improve the model from multiple aspects, including introducing spatial control [20, 449] and reference images [37, 387].

*5.2.2    Talking face.* From the perspective of output, the task of talking face[537] generates a series of image frames which are thus technically a video (see Figure 19). Different from general video generation (see Sec. 6.1), talking face requires an image face as an identity reference, and edits it based on the speech input. In this sense, talking face is more related to image editing. Moreover, talking face converts a speech clip to a corresponding face image, resembling speech recognition to convert a speech clip to a corresponding word text. With speech recognition recognized as a

"a hedgehog using a calculator"    "a corgi wearing a red bowtie and a purple party hat"    "robots meditating in a vipassana retreat"    "a fall landscape with a small cottage next to a lake"

"a surrealist dream-like oil painting by salvador dalí of a cat playing checkers"    "a professional photo of a sunset behind the grand canyon"    "a high-quality oil painting of a psychedelic hamster dragon"    "an illustration of albert einstein wearing a superhero costume"

Fig. 18.  Examples of text-to-image (figure from [300]).

multimodal generation text task, this survey considers talking face as a multimodal image generation task. Driven by deep learning models, speech-to-head video synthesis models have attracted wide attention, which can be divided into 2D-based methods and 3D-based methods.

With 2D-based methods, talking face video synthesis mainly relies on landmarks, semantic maps, or similar representations. Landmarks are used as an intermediate layer from low-dimensional audio to high-dimensional video, as well as two decoders to decouple speech and speaker identity for generating video unaffected by speaker identity [66], which is also the first work to use deep generative models to create speech faces. In addition, image-to-image translation generation [178] can also be used for lip synthesis, while the combination of separate audio-visual representations and neural networks can also be used to optimize synthesis [404, 539]



**Input:** audio and single portrait image                                    **Output:** talking head animation

Fig. 19.  Examples of talking face (image obtained from [51]).

Another line of work is based on building a 3D model and controlling the motion process through rendering technology [219, 414], with a drawback of high construction cost. Later, many generative talking face models based on 3DMM parameters [74, 111, 196, 423] were established, using models such as blendshape [74], flame [239], and 3D mesh [352], with audio as model input for content generation. At present, most methods are directly reconstructed from training videos. NeRF uses multi-layer perceptrons to simulate implicit representations, which can store 3D spatial coordinates and appearance information and are used for high-resolution scenes [238, 286, 294]. In addition, a pipeline and an end-to-end framework for unrestricted talking face video synthesis have also been proposed [215, 328], taking any unidentified video and arbitrary speech as input.

## 6 AIGC TASK: BEYOND TEXT AND IMAGE

### 6.1 Video

Compared with image generation, the progress of video generation lags behind largely because of the complexity of modeling higher-dimensional video data. Video generation involves not only generating pixels but also ensuring semantic coherence between different frames. Video generation works can be categorized into unguided and guided generation (e.g., text, images, video, and action classes), with text-guided age (see Figure **??**) receiving the most attention due to its high influence.



(a) A dog wearing a superhero outfit with red cape flying through the sky.

(b) There is a table by a window with sunlight streaming through illuminating a pile of books.

(c) Robot dancing in times square.

(d) Unicorns running along a beach, highly detailed.

Fig. 20. Examples of text-guided video generation (figure obtained from [394]).

**Unguided video generation.** Early works on extending image generation from single frame to multiple frames are limited to creating monotonous yet regular content like sea waves. The generated dynamic textures [96, 466] often have

a spatially repetitive pattern with time-varying visualization. With the development of generative models, numerous works [2, 68, 305, 370, 433, 448, 512] extend the exploration from naive dynamic textures to real video generation. Nonetheless, their success is limited to short videos for simple scenes with the availability of low-resolution datasets. More recent works [67, 157, 371, 424] improve the video quality further, among which [157] is regarded as a pioneering work of diffusion models.

**Text-guided video generation.** Compared to text-to-image models that can create almost photorealistic pictures, text-guided video generation is more challenging. Early works [136, 246, 260, 276, 290, 313] based on VAE or GAN concentrate on creating video in simple settings, such as digit bouncing, and human walking. Given the great success of the VQ-VAE model in text-guided image generation, some works [160, 472] extend it to text-guided video generation, resulting in more realistic video scenes. To achieve high-quality video, [157] first applies the diffusion model to text-guided video generation, which refreshes the benchmarks of evaluation. After that, Meta and Google propose Make-a-Video [394] and Imagen Video [155] based on the diffusion model, respectively. Specifically, Make-a-Video extends a diffusion-based text-guided image generation model to video generation, which can speed up the generation and eliminate the need for paired text-video data in training. However, Make-a-Video requires a large-scale text-video dataset for fine-tuning, which results in a significant amount of computational resources. The latest Tune-a-Video [474] proposes one-shot video generation, driven by text guidance and image inputs, where a single text-video pair is used to train an open-domain generator.

## 6.2 3D generation

The tremendous success of deep generative models on 2D images has prompted researchers to explore 3D data generation, which is actually a modeling of the real physical world. Different from the single format of 2D data, a 3D object can be represented by depth images, voxel grids[476], point clouds[330, 331], meshes[140] and neural fields[283], each of which has its advantages and disadvantages.

According to the type of input and guidance, 3D objects can be generated from text, images and 3D data. Although multiple methods [112, 175, 262] have explored shape editing guided by semantic tags or language descriptions, 3D generation is still challenging due to the lack of 3D data and suitable architectures. Based on the diffusion model, DreamFusion [326] proposes to solve these problems with a pre-trained text-to-2D model. Another branch of works reconstruct the 3D objects from single-view images [33, 122, 243, 432, 457, 510] or multi-view images [63, 167, 454, 480], termed Image-to-3D. A new branch of multi-view 3D reconstruction is Neural Radiance Fields (NeRF) [286] for implicit representation of 3D information. The 3D-3D task includes completion from partial 3D data [455] and transformation [26], with 3D object retrieval as a representative transformation task.

## 6.3 Speech

Speech synthesis is an important research area in speech processing that aims to make machines generate natural and understandable speech from text. Methods of traditional speech synthesis include articulatory [217, 380], formant [12, 377], concatenative synthesis [293, 306], and statistical parametric speech synthesis (SPSS) [198, 292]. These methods have been widely studied and applied, e.g., formant synthesis is still used in the open-source NVDA (one of the leading free screen readers for Windows). However, these generated speeches are identifiable from the human voice, and artifacts in synthesis speech reduce intelligibility.

Early works [102, 333, 514–516] consist of three modules: text analysis, an acoustic model, and a vocoder. WaveNet [440] is a revolution within speech synthesis which can generate the raw waveform from the linguistic features. To improve

the quality of speech and diversity of voices, generative models are introduced in speech synthesis, such as GAN [124]. Compared with GAN, diffusion models do not require a discriminator, making training more stable and simple. Therefore, the works of speech synthesis adopt diffusion models, becoming a rising trend. A branch of works [57, 168, 220, 479] focuses on efficient speech synthesis, in which different ways are adopted to reduce the generated time by accelerating inference, such as combining the schedule and score networks for training, jointly trained GAN. Another branch of study [52, 289, 361, 390] concentrates on end-to-end models, which directly generate waveform from text without any intermediate representations. A fully end-to-end model not only simplifies the training and inference, but also reduces the demand for human annotations. The branch of diffusion-based speech synthesis is not limited to the two mentioned above, such as speech enhancement and guided speech synthesis.

## 6.4 Graph

Graphs are ubiquitous in the world, which aid in visualizing and defining the relationships between objects in a wide range of domains, from social networks to chemical compounds. Graph generation, which creates new graphs from a trained distribution that is similar to the existing graphs, has received a lot of attention.

Traditional graph generation works [11, 224, 464] create new graphs with specific features that are related to the hand-crafted statistical features of real graphs , which simplifies the process but fails to capture relational structure in complex scenarios. With the successes of deep learning algorithms, researchers have begun to apply them to graph generation, which, unlike the traditional methods, can be directly trained by real data and automatically extract features. Among them, works [76, 249, 503] based on autoregressive model create graph structures sequentially in a step-wise fashion, which allows for greater scalability but fails to model the permutation invariance and is computationally expensive. Simultaneously, One-shot models [254, 269, 269] such as VAE and flow are incapable of accurately modeling structure information because of ensuring tractable likelihood computation. Although graph generation [80, 183, 278] based on GAN sides step likelihood-based optimization by using a discriminator, the training is unstable.

Recently, there has been a surging interest in developing diffusion models for graph-structured data. EDP-GNN [302] is the pioneering to show the capability of the diffusion model in the Graph generation, with the goal of addressing non-invariant properties. After that, On the one hand, diffusion-based works [138, 166, 186, 266, 445] focus on realistic graph generation, which produces graphs that are similar to a given set of graphs. On the other hand, [14, 388, 482, 513] concentrate on goal-directed graph generation, which generates graphs that optimizes given objects, like molecular and material generation.

## 6.5 Others

There are also other interesting tasks generating content in different modalities, e.g., music generation [179] and lip-reading [106]. A typical music generation system can be categorized into three representation levels (from top to bottom), which generates score, performance, and audio, respectively [179]. With the development of deep learning, music generation introduces various methods for higher music quality, e.g., MusicVAE [354], MuseGAN [95] and transformer in [170]. Music generation inspires and accelerates the development of computer-assistant composition software, including Magenta project from Google and Flow Machine project from Sony Computer Science Laboratories. A Lip reading task transforms visual inputs of lip movement to decoded speech [106], and has also shown impressive advances thanks to improved corpora and architectures.

## 7 INDUSTRY APPLICATIONS

Undoubtedly, AIGC has gone viral on social media since 2022. For example, users are active in sharing their experience of using ChatGPT for having an interactive conversation or Stable diffusion for generating images with a text prompt. However, this hype is expected to dwindle if AIGC cannot be used for practical applications in the industry to demonstrate its value. Therefore, we discuss how AIGC might influence various industries.

### 7.1 Education

AIGC is changing the paradigm of education by assisting in teaching and learning. Generative AI carries transformative potential in teaching, with the application ranging from course materials generation to assessment and evaluation [324, 517]. Simultaneously, applications of generative models have begun to influence how students learn [27, 420].

Generative AI technologies can provide educators with creation of personalized tutoring [517], designing course materials [324], and assessment and evaluation [27, 517]. A unique foreign language teaching product for young children using generative technologies such as ChatGPT can attract children's attention, motivate them, and provide a fun learning environment. Higher education needs to embrace the use of AI in higher education, which can create more engaging, effective, and efficient learning experiences for students [517]. One of the primary benefits of generated AI course material generation is that it can save teachers time and effort by automating the process of creating and updating course material. In addition, ChatGPT could significantly reduce the workload of law school instructors, freeing up time to increase academic productivity or develop more complex teaching skills [324]. The benefits of ChatGPT in promoting teaching include but are not limited to facilitating personalized and interactive learning. However, some limitations of ChatGPT, such as generating incorrect information, exacerbating existing biases in data training, and privacy issues, can also appear [27]. Overall, addressing these challenges requires collaborative efforts from policymakers and educators to provide recommendations or guidance for the appropriate use of generative AI tools.

Moreover, generative AI technologies can help students write essays [420], at-home tests or quizzes [420], comprehend certain theories and concepts, and different language essays and papers in academic issues [27, 517]. Chatbots can provide students with 24/7 support, allowing them to get the help they need when they need it. With the ability to correct grammar, suggest improvements, and identify weak areas, chatbots like ChatGPT can provide students with immediate feedback on their writing, helping them to learn from their mistakes and improve their writing skills over time. This not only saves students time but also helps them to become better writers [493]. According to a survey conducted by an online course provider, 89% of students use ChatGPT to complete their homework, with 50% using it for essays and 48% using it for at-home tests or quizzes [420]. Additionally, generated AI can tailor the course material to individual students' needs, such as learning style and pace, which has the potential to improve student engagement and learning outcomes. ChatGPT can also help students comprehend certain theories, concepts, and different language articles, making them work more effectively [27, 517]. There are also challenges and concerns associated with generated AI course material generation, including the generated material's quality, and the possibility of bias in the data used to train the AI. As a result, before using generated course material in an educational context, it is critical to evaluate and validate it carefully [79].

With the use cases mentioned above, AIGC has the potential to revolutionize education by improving the quality and accessibility of educational content, increasing student engagement and retention, and providing personalized support for learners. With the continuous advancements in AI, AIGC is poised to become an integral part of the education industry, offering students a more engaging, accessible, and personalized learning experience.

## 7.2   Game and metaverse

Most users may not resonate with one-size-fits-all content in the game and metaverse, where personalization yields the best experience. Although games and metaverse provide users with virtual worlds, the content represents the character and personality of users. Generative AI makes that possible, which not only allows users to customize their avatars but also provides diverse scenarios and storylines, making the experience more immersive [53, 325, 344, 344].

AI Dungeon powered by GPT-3 model allows users to generate an open-end story navigated by text, where generative AI will produce new events as the response to the different actions of users, creating a one-of-a-kind and unexpected gameplay experience [422]. Horizon Worlds, one of the most popular games, allows you to wander into the virtual world related to content consumption and creation. In Horizon, users will have more control over how they want to tailor their online experience to meet their individual needs. Specifically, users can design their unqiue avatars and scenes using gizmos that include pre-built object and avatar properties [284]. Moreover, the visual novel game Traveler concentrates on generating gorgeous scenarios to present users with visual impact, in which you will embark on a journey through a diverse world. When a player explores the game Traveler, the player will be exposed to magnificent visuals and immersive soundscapes. As each scene is unique, the content can range from dark forests to bustling cities, all crafted by generative AI [113].

Although the term "metaverse" has become a buzzword recently, in the real world, the virtual space created by the game may serve as the portal to the metaverse [301, 344]. Roblox, a sandbox game platform, first included the concept of "Metaverse" in its prospectus and made its market value soar, where players can create their own world beyond their imagination [301]. Virtual concert singers have a more comprehensive range of musical styles and talents. Audiences can choose their own favorite styles and even idols, providing them with a more diverse and personalized concert experience. Travis Scott, a well-known American rapper and producer, performed a historic concert inside the Fortnite game, and his avatar guided the players to experience different scenes, ranging from underwater to outer space [465]. The University of California, Berkeley, presented its commencement ceremony in 'Minecraft', a popular computer game. In the Minecraft game, students and alumni built a copy of campus using generative AI technology, allowing thousands of graduating seniors using their avatars from around the world to attend the event [358]. Overall, AI has played a significant role in the evolution of the game and metaverse, and its use continues to grow as technology improves and becomes more accessible.

## 7.3   Media

With the ubiquitous growth of generative AI technologies, they play a rising role in media and advertising. AIGC not only promotes the diversity of media, which provides a better experience for audiences, but it also enables media practitioners more efficient in their work [44, 84, 107, 222, 288, 439, 527].

The media powered with AIGC enables more diversified content and ways of reporting, changing the media mode of production and organizational structure [204, 273, 527]. AIGC can be applied to a variety of applications in media, such as writing robots, news anchors, and caption generation. Traditionally, media outlets have relied on expert journalists to write new articles and reports, which requires a significant amount of energy and time, resulting in a limited number of articles. Moreover, the timeliness of the news is critical, and the news may be eclipsed after an hour. Generative AI can greatly assist journalism by using text generation technologies to make journalism more efficient and responsive [288]. Associated Press applies these technologies to generate roughly 40000 stories a year and its articles on company earnings increase from 1200 to 14800 [130]. The Quakebot, a robot reporter of Los Angeles Times News, only takes three minutes

to complete a related article after the Los Angeles earthquake [309]. Bloomberg News, an international financial media company, launched Buttetin in 2018, with the goal of providing personally storied whose one-sentence summaries are generated by chatbots [470].

AI news anchors have emerged as a result of the deep integration of generative AI in the media [461, 491? ]. AI news anchors, combined with real anchors, make the way to spread information more diverse. AI news anchors can broadcast news based on the text, whose appearance and expression imitate the real anchor. China's state news agency Xinhua and Chinese search engine, Sogou, have developed AI news anchors with different profiles and languages. The most impressive is the 3D AI news anchor Xin Xiaowei, whose broadcast form can be presented in all directions from various angles, significantly improving the sense of three-dimensionality and layering [279]. Additionally, Korea's cable channel MBN has created the AI news anchor AI Kim, who can quickly respond to various emergencies and even report all day [397]. To help the hearing-impaired people to get more information about international sporting events and Beijing Winter Olympics, an AI-driven sign language service provider, namely Ling Yu, was developed by giant tech Tencent, where AIGC tasks are used including 3D digital human modeling, machine translation, image generation, and speech-to-text [511]. Moreover, Migu, a Chinese business dedicated to providing digital information, can offer smart subtitle functions to the live broadcast of Beijing Winter Olympics. Therefore, people with hearing impairments can watch the live broadcast of the sports event, which makes them more immersive.

## 7.4 Advertising

Various AI applications have transformed the advertising industry, giving advertisers powerful tools to create innovative and engaging content that connects to consumers at a deeper level [84, 133, 439]. Among the various applications of AI in advertising, AIGC is particularly influential by allowing advertisers to create personalized and attractive content that resonates with individual consumers. A creative advertising system (CAS) aligns with the principles of AI for generating and testing advertising ideas, which helps aspiring and mature creators understand that creativity is not an elite privilege, but rather a systematic process that can be assisted through data and computation [439]. Implementing programmatic advertising has not fully utilized self-generating technology, resulting in different consumers being exposed to the same content. Fortunately, a personalized advertising copy intelligent generation system (SGS-PAC) can automatically personalize advertising content to meet individual consumer needs [84]. Advertising posters are a common form of information display used to promote products. Another intelligent system, Vinci, supports the automatic generation of advertising posters [133]. By inputting product images and slogans specified by users, Vinci uses deep-generation models to generate beautiful posters.

In addition, Brandmark.io is an AIGC-based tool that automatically generates logos for businesses. The tool creates multiple logo variations based on the user's preferences and specifications. Advertisers can purchase and use the logos created by the tool for their businesses, making it an easy and cost-effective solution for logo design [193]. By GAN that forces output to include specific keywords, the approach automates product listing generation likely to attract potential buyers [275]. It enhances users' marketing efforts on peer-to-peer marketplaces. Moreover, technological innovations have provided digital and automated tools to the advertising industry, but have also allowed advertisers to automate the production of "synthetic advertising". As reported by [381], AIGC has transformed the advertising industry by enabling advertisers to create highly personalized and engaging content at scale while saving time and resources. We expect to see even more innovative and influential applications of generated AI in advertising.

## 7.5 Movie

It is interesting to see how technology now affects almost every step of movie creation. Research has led to the development of computer-based surroundings that help with editing, labeling, video retrieval, and many more [78, 125, 295, 317–319, 366, 415, 428, 434]. To start with, AI-powered screenwriting software has significantly impacted the movie-making process. AI has created a new movie experience by integrating visual effects (VFX) [291], improved sound effects (SFX), and new viewing platforms. The 4K, IMAX, and 3D movies, as well as animations, are highly impacted by them.

The script forms the foundation of how a movie will fare at the ticket counters. AI-generation devices store and compute massive amounts of data to create "ideal" scripts [15]. AI software is also employed to rework old screenplays into polished versions that are then analyzed and improved by the director and writer. It goes beyond just developing and analyzing current scripts. Jasper AI [362] and Scalenut [329] are two examples of AI scriptwriters.

Movies are given visual effects (VFX) to increase spectator appeal. They combine original images with real video to create engrossing, realistic, contextual depictions that may include digital surroundings, de-aging, and many more. The VFX team of The Curious Case of Benjamin Button [416] put up two arrays of cameras in a bright room and utilized the MOVA Contour reality capture technology [82] to construct a three-dimensional database of the hero's facial expressions. The VFX team next developed high-resolution 3D models of the lead character at various ages and lastly employed AI to manipulate the data retrieved from the three-dimensional database to cause the head models to age. The outcome was convincing and garnered widespread acclaim across the movie community. In movies like Blade Runner 2049 [274] and Gemini Man [360], the VFX team tracked and recorded the protagonist's facial data with the help of motion capture technology [5]. They then rebuilt the 3D facial data in a computer and further polished it to accomplish age deduction. Although this approach takes expensive technology and a vast amount of money, it is incredibly precise and adaptable.

AI is expanding the limits of amusement by bringing back actors who have passed away in movies. Deepfake technology uses computer visuals and AI to produce incredibly amazing and lifelike fake videos of actual people or made-up characters. Fast and Furious 7 [469] used VFX and vintage video to bring back Paul Walker after he passed away in 2013, using the actor's visage transferred onto his sibling.

In addition to the visual effects, subtitles also play a vital role in viewers' experience. For the benefit of viewers having hearing impairments, automated subtitles for the deaf and hard of hearing, also known as SDH [8], include textual transcriptions of speech, speaker changes, and background noise. These benefits substantially increase how much money movies make. Natural Language Processing (NLP), a kind of AI that focuses on deciphering spoken language, offers multilingual subtitles in movies. To produce automated subtitles, Rev [351], a cloud-based program, is widely used by movie lovers. AI has also revolutionized the task of speech prediction in silent movies. AI-generated speech synthesis can narrate silent movies and dub movies into multiple languages. Deep learning systems trained on massive human audio samples can produce natural-sounding voiceovers. Features like LPC (Linear Predictive Coding) [101] and mel spectrograms [100] generate high-quality intelligible audio through conversion. Recently, a Tacotron2 model [384] variation for the video-to-audio synthesis was proposed by [327]. [494] suggests an efficient stochastic model that produces endless high-quality audio patterns for a specific silent video, thus effectively encapsulating the multimodality of the speech prediction issue. Along with the visual and sound effects, we have Colourlab.Ai [70] for color grading, Descript [86] for video editing, and many more tools continuously making waves in the movie industry.

## 7.6  Music

AIGC also makes it to the music industry with notable developments [499]. AI can not only spot patterns and trends in vast data sets that are challenging for humans to notice, but also allows amateur musicians a cutting-edge technique to enhance their creative process, which is a fantastic opportunity. The fusion of AI technology into music is a new trend that many experts, researchers, musicians, and record companies are exploring [280]. Many utilize AIGC to create entirely new music, while some software edit compositions in the style of various composers.

Music industries are anticipating significant expenditures in this field, whether it is because of using AI to compose music or to help musicians. A fantastic illustration of an AI melody generator is Google's Magenta project [172], [10]. IBM's Watson Beat is one more example. For composing an original song, it makes use of AI and machine learning [49], [110]. 2016 saw the successful creation of text-to-speech (TTS) recordings and recordings that resembled music by DeepMind researchers [421], [496]. AI is also vastly used for the processing and improvement of digital audio. LANDR, an incredible AI-powered creative tool that enables musicians to get their music on several streaming services like Spotify and Apple Music, is one such service. A significant problem known as "writer's block [379]" frequently confronts lyricists. But thanks to AI, it's no longer a problem now. Nowadays, many musicians employ AI to create new lyrics for their songs [277]. GPT-2 [7], [538], a text-generating tool, has been created by OpenAI, an AI technology firm. Not only can this remarkable text generator produce authentic news, but it can also write lyrics for Beatles songs and music from all other genres. However, AI is not just capable of producing text; it can also create original soundtracks and melodies. The Sony CSL flow machine offers assistance to artists so they can develop original music based on their ideas [356]. One of the most well-known AI tools for writing unique music is called AIVA [189]. To produce a unique track, the user first chooses a pre-set style and then modifies a variety of variables, such as the key, instrumentation, time signature, etc. AIVA can deconstruct all the intricate auditory information saved into discrete characteristics while reading hundreds of musical compositions by renowned musicians like Bach and Mozart [98]. These qualities may then be interpreted again and used to produce a completely new musical work. Apart from the tools mentioned above, there are so many other applications that made a significant impact on the music industry, such as the iOS-based tool Amadeus Code [38], the cloud-based platform Amper [147], Ecrett Music [36], etc.

## 7.7  Painting

From offering automatic painting tools to encouraging creative experimentation, AIGC is revolutionizing the painting industry in many ways. AI programs can analyze pictures to produce color schemes, patterns, and textures that can make artwork. The automatic drawing tools generated using these algorithms are able to apply these patterns and textures to produce distinctive and intricate works of art [548]. AI can also analyze a person's preferences, interests, and style to create customized artwork. Empowering artists to create art specifically suited to their preferences and interests can increase their appeal and value.

The artwork created by MidJourney under the title *'Space Opera Theatre'* earned first place in the Colorado State Fair Art Competition [81], demonstrating the capability of AI painting tools to produce excellent pieces of art. Midjourney is an excellent AI image generator with comprehensive functions, which is used by many artists to generate inspiration. The creation of various art forms by generative AI, such as abstract painting generation [235], Chinese shanshui painting [543], and Chinese ink paintings [64], undoubtedly promotes the advancement of painting. Moreover, AIGC can assist in conservation and restoration [506], [164]. AI algorithms are capable of analyzing and repairing ruined

artwork. These algorithms make it simpler for conservators to return the artwork to its initial state by detecting and removing dust, scratches, and other flaws.

For non-professionals unfamiliar with drawing or animation, AIGC is also very helpful because it enables them to produce high-quality visual effects. By adding additional constraints to the diffusion model, ControlNet [341] can increase the variability of the produced images. It can describe the generated images along with those other constraints of border drawing, depth information, Hough line map, normal map, and posture estimation. AIGC has also started a new era of collaborative artwork [41, 50]. AI algorithms can create collaborative paintings that involve multiple artists working together. These algorithms can analyze the styles of each artist and produce a unified style that incorporates elements from all of the artists' works.

### 7.8 Code development

Generative AI can contribute to the field of code development [134, 135, 176, 412], where AIGC can create code without the need for manual coding. The work by [412] explores the interpretability requirements of generative AI for code and demonstrates how human-centered approaches can drive the development of explainable AI (XAI) technologies in new domains. To improve testing efficiency and increase test coverage, it is particularly important to generate high-quality test cases automatically [134, 135]. In order to optimize the efficiency of data engineering, a novel software engineering approach based on neural networks for dataset augmentation can be designed [176]. One of the popular applications in AI-generated code is Github's Copilot, an AI tool jointly developed by GitHub and OpenAI. Users can automatically complete code through GitHub Copilot using software development tools [91]. Moreover, AI-generated technology can also assist in code refactoring, which improves existing code without changing its original functionality. This can shorten the time for developers to refactor and improve the quality of the code. A popular code refactoring tool is DeepCode [174], an AI-supported code review tool that can inspect your code and provide suggestions for improvement. In addition, AIGC can also make an impact on the e-commerce and finance industries [163]. E-commerce platforms such as Amazon, JD.com, and so on can use AI-powered customer service to provide shopping guide services to customers, thereby saving costs for enterprises. Financial companies can use virtual investment advisors to advise customers on securities account opening, financial investment, and other related services.

### 7.9 Phone apps and features

Numerous AIGC applications have emerged as fun-oriented mobile apps, typically in the form of image and video editing. Photoshop is traditionally a common tool for image editing, but manual work is time-consuming and can result in unnatural or unrealistic output. In addition, video editing involves analyzing each video clip and making editorial decisions based on both the audio and visual content. This process is time-consuming because the video is a time-based, dual-track medium that requires careful consideration of every frame. Fortunately, some work [17, 32, 398] has explored the utilization of AI technologies behind AIGC, to the image or video editing, making the applications in AIGC such as face swapping and digital avatar possible.

Some popular applications based on face swapping are gaining widespread popularity on the Internet. This technology uses advanced AI technologies to analyze and swap people's faces with their favorite celebrities or anyone else in seconds, making it easier and faster to use compared to traditional PS technologies. VanceAI, Voila AI Artist and FaceAPP are leading figures, with FaceApp being recognized as the best facial photo editing App, winning numerous awards, and being downloaded by over 500 million users and counting [372]. Another popular application is voice-changing technology. This technology can adjust the pitch, timbre, speech rate, and other characteristics of the human voice to

change the quality of the human voice. MagicMic [471] and Voicemod [202] are two popular applications for real-time voice modification and soundboard operations, which people can use to change their voices for creating fun content, live streaming, or other purposes, enhancing the enjoyment of communication between people.

In addition, another technological trend is to transform individuals into virtual characters, thereby increasing entertainment value. virtual characters are digital avatars of people in a virtual world, they can be partial replicas of real people or even completely digital versions. Apple's first "digital avatar" technology, Animoji, focuses mainly on generating preset cartoon and animal characters and does not support custom generation [426]. The second generation of "digital avatar" technology represented by iPhone's Memoji and Xiaomi's Mimoji started to support personalized avatar customization, which offers a variety of options, starting from hairstyle, eyes, nose, dresses, etc [406]. This upgrade allows users to create an avatar that not only can track their facial movements, but also look like them. Besides that, the created avatars can also be posted as comments in WeChat or Facebook chats, giving users a more personal way to express themselves on social media. Since then, digital avatar technology has become one of the standard features of smartphones among various smartphone manufacturers.

## 7.10 Other fields

Beyond the above fields, AIGC is expected to have applications in more fields. For example, the design and development of a novel drug are complex, costly, and time-consuming. On average, it takes around $3 billion and more than 10 years for a new drug to be accepted by the market [201]. This motivates using AIGC to accelerate the drug discovery process and reduce costs. In 2018 DeepMind created AlphaFold [364], which can accurately predict the structure of proteins and has been considered a milestone for drug discovery and fundamental biology research. Its updated version AlphaFold2 was released in 2020 and had higher accuracy than the former. ProteinMPNN[77], designed by Justas Dauparas, can design protein sequences for specific tasks, generating entirely new proteins quickly in just a few seconds. Besides directly exploiting the generated content, AIGC can also help workers in various fields improve their efficiency. For example, in medical consultation, the patient can rely on chatbots for basic medical advice, while turning to the doctor only for more severe cases. In manufacturing design, it is possible to combine AIGC with the widely used computer-aided design system to minimize the repetitive effort so that the designer can focus on the more meaningful part.

## 8 CHALLENGES AND OUTLOOK

### 8.1 Challenges

Even though AIGC has shown remarkable success in generating realistic and diverse outputs across various domains, there are still numerous challenges in real-world applications. Except for requiring a large amount of training data and compute resources, we list some of the most significant challenges as follows.

(1) **Lack of interpretability.** While AIGC models can yield impressive outputs, it remains challenging to understand how the model arrives at the outputs. This is especially a concern when the model generates an undesirable output. Such a lack of interpretability makes it difficult to control the output.

(2) **Ethical and legal concerns.** The AIGC model is prone to data bias. For example, a language model mainly trained on the English text can be biased toward western culture. Copyright infringement and privacy violations are the underlying legal concerns that cannot be ignored. Moreover, the AIGC model also has the potential for malicious use. For example, students can exploit these tools to cheat on their essay assignments, for which AI

content detectors are desired. AIGC models can also be used for distributing misleading content for political campaigns.

(3) **Domain-specific technical challenges.** At the current and in the near future, different domains require their unique AIGC models. Each domain is still faced with its unique challenges. For example, Stable Diffusion, a popular text-to-image AIGC tool, occasionally generates output that is far from what the user desires, such as drawing humans as animals, one person as two people, etc. ChatBot, on the other hand, makes factual mistakes occasionally.

## 8.2 Outlook

Despite its unprecedented popularity, generative AI is still in its early stage. Here, we present how AIGC might evolve in the near future.

(1) **More flexible control.** A major trend of AIGC tasks is to realize more flexible control. Taking image generation as an example, early GAN-based models can generate images of high quality, but with little control. The recent diffusion models trained on large text-image data enable control through text instruction. This facilitates the generation of images that better match the users' needs. Nonetheless, current text-to-image models still require more fine-grained control so that the images can be generated in a more flexible manner

(2) **From pertaining to finetuning.** Currently, the development of AIGC models like ChatGPT focuses on the pretraining stage. The corresponding technology is relatively mature; however, how to fine-tune these foundation models for the downstream tasks is an under-explored field. Different from training a model from scratch, the goal of finetuning needs to trade-off between the foundation model's original general capability and its adaptation performance on the new task.

(3) **From big tech companies to startups.** At present, AIGC technology has mainly developed big tech companies, like Google and Meta. With the support of big tech companies, some startup companies have emerged to show high potentials, like OpenAI (supported by Microsoft) and DeepMind (supported by Google). With the focus transition from core technology development to applications, more startup companies are expected to emerge due to increasing demand.
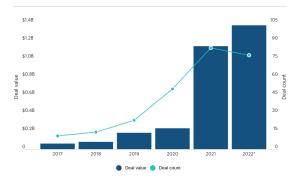


Fig. 21. Deal value and deal count for generative AI funding in the past 6 years (figure obtained from: PitchBook).

**Discussion: investment, bubble and job opportunities.** Technology-wise, there is no doubt that AIGC has made significant progress in the past few years. When a transformative technology emerges, the market tends to be

over-optimistic about its potential applications and future growth, which also applies to generative AI. According to PitchBook (see Fig. 21), the funding for generative AI from venture capital (VC) increased significantly in the last two years. Some critics have concerns that generative AI might be the next bubble. One of their main concerns is that most AIGC tools are mainly playful instead of practical. For example, text-to-image models are fun to play with, but how they might generate revenues remains unclear. It is difficult to predict how generative AI might evolve. However, the authors of this work believe that generative AI is unlikely to become the next bubble considering it is a relatively new and rapidly growing field with many potential applications. There is also a hot debate about whether generative AI will replace humans, causing the loss of numerous job opportunities. On the other hand, generative AI can also create new job opportunities for individuals with skills on AI research and implementation skills. The industries that benefit from the power of AIGC might also boom and generate more job opportunities.

## REFERENCES

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. 2019. Image2stylegan: How to embed images into the stylegan latent space?. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4432–4441.

[2] Dinesh Acharya, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. 2018. Towards high resolution video generation with progressive growing of sliced wasserstein gans. *arXiv preprint arXiv:1810.02419* (2018).

[3] Johannes Ackermann and Minjun Li. 2022. High-resolution image editing via multi-stage blended diffusion. *arXiv preprint arXiv:2210.12965* (2022).

[4] Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977* (2020).

[5] Adobe. 2022. What is motion capture and how does it work? *https://www.adobe.com/uk/creativecloud/animation/discover/motion-capture.html* (2022).

[6] Namhyuk Ahn, Byungkon Kang, and Kyung-Ah Sohn. 2017. Image distortion detection using convolutional neural network. In *2017 4th IAPR Asian Conference on Pattern Recognition (ACPR)*. IEEE, 220–225.

[7] Open AI. 2019. Blog GPT-2. *https://openai.com/blog/tags/gpt-2/* (2019).

[8] ai media. 2017. Subtitles for the Deaf or Hard-of-Hearing (SDH) - Subtitles, Closed Captions, and SDH. *https://blog.ai-media.tv/blog/what-is-sdh* (2017).

[9] Reina Akama, Sho Yokoi, Jun Suzuki, and Kentaro Inui. 2020. Filtering noisy dialogue corpora by connectivity and content relatedness. *arXiv preprint arXiv:2004.14008* (2020).

[10] Mincer Alaeddine and Anthony Tannoury. 2021. Artificial Intelligence in Music Composition. In *Artificial Intelligence Applications and Innovations: 17th IFIP WG 12.5 International Conference, AIAI 2021, Hersonissos, Crete, Greece, June 25–27, 2021, Proceedings 17*. Springer, 387–397.

[11] Réka Albert and Albert-László Barabási. 2002. Statistical mechanics of complex networks. *Reviews of modern physics* 74, 1 (2002), 47.

[12] Jonathan Allen, Sharon Hunnicutt, Rolf Carlson, and Bjorn Granstrom. 1979. MITalk-79: The 1979 MIT text-to-speech system. *The Journal of the Acoustical Society of America* 65, S1 (1979), S130–S130.

[13] Jaan Altosaar. 2016. *Tutorial - What is a Variational Autoencoder?* https://doi.org/10.5281/zenodo.4462916

[14] Namrata Anand and Tudor Achim. 2022. Protein Structure and Sequence Generation with Equivariant Denoising Diffusion Probabilistic Models. *arXiv preprint arXiv:2205.15019* (2022).

[15] Nantheera Anantrasirichai and David Bull. 2022. Artificial intelligence in the creative industries: a review. *Artificial intelligence review* (2022), 1–68.

[16] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6077–6086.

[17] Dawit Mureja Argaw, Fabian Caba Heilbron, Joon-Young Lee, Markus Woodson, and In So Kweon. 2022. The anatomy of video editing: A dataset and benchmark suite for ai-assisted video editing. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VIII*. Springer, 201–218.

[18] Martin Arjovsky, Amar Shah, and Yoshua Bengio. 2016. Unitary evolution recurrent neural networks. In *International conference on machine learning*. PMLR, 1120–1128.

[19] Omri Avrahami, Ohad Fried, and Dani Lischinski. 2022. Blended Latent Diffusion. *arXiv preprint arXiv:2206.02779* (2022).

[20] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. 2022. SpaText: Spatio-Textual Representation for Controllable Image Generation. *arXiv preprint arXiv:2211.14305* (2022).

[21] Omri Avrahami, Dani Lischinski, and Ohad Fried. 2022. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18208–18218.

[22] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).

[23] Alexei Baevski, Michael Auli, and Abdelrahman Mohamed. 2019. Effectiveness of self-supervised pre-training for speech recognition. *arXiv preprint arXiv:1911.03912* (2019).

[24] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).

[25] Ashutosh Baheti, Alan Ritter, and Kevin Small. 2020. Fluent response generation for conversational question answering. *arXiv preprint arXiv:2005.10464* (2020).

[26] Song Bai, Xiang Bai, Wenyu Liu, and Fabio Roli. 2015. Neural shape codes for 3D model retrieval. *Pattern Recognition Letters* 65 (2015), 15–21.

[27] David Baidoo-Anu and Leticia Owusu Ansah. 2023. Education in the Era of Generative Artificial Intelligence (AI): Understanding the Potential Benefits of ChatGPT in Promoting Teaching and Learning. *Available at SSRN 4337484* (2023).

[28] Dana H Ballard. 1987. Modular learning in neural networks.. In *Aaai*, Vol. 647. 279–284.

[29] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence* 41, 2 (2018), 423–443.

[30] Hangbo Bao, Li Dong, and Furu Wei. 2022. Beit: Bert pre-training of image transformers. *ICLR* (2022).

[31] Siqi Bao, Huang He, Fan Wang, Rongzhong Lian, and Hua Wu. 2019. Know more about each other: Evolving dialogue strategy via compound assessment. *arXiv preprint arXiv:1906.00549* (2019).

[32] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. 2022. Text2live: Text-driven layered image and video editing. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*. Springer, 707–723.

[33] Jan Bednarik, Pascal Fua, and Mathieu Salzmann. 2018. Learning to reconstruct texture-less deformable surfaces from a single view. In *2018 International Conference on 3D Vision (3DV)*. IEEE, 606–615.

[34] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *Advances in neural information processing systems* 13 (2000).

[35] Mohamed Benzeghiba, Renato De Mori, Olivier Deroo, Stephane Dupont, Teodora Erbes, Denis Jouvet, Luciano Fissore, Pietro Laface, Alfred Mertins, Christophe Ris, et al. 2007. Automatic speech recognition and speech variability: A review. *Speech communication* 49, 10-11 (2007), 763–786.

[36] Sean Berry. 2019. Ecrett Music uses AI to generate royalty free music for your videos. *https://www.videomaker.com/news/ecrett-music-uses-ai-to-generate-royalty-free-music-for-your-videos/* (2019).

[37] Andreas Blattmann, Robin Rombach, Kaan Oktay, and Björn Ommer. 2022. Retrieval-Augmented Diffusion Models. *arXiv preprint arXiv:2204.11824* (2022).

[38] Amadeus Code Blog. 2019. Artificial intelligence-powered songwriting assistant. *https://blog.amadeuscode.com/* (2019).

[39] Denny Britz, Anna Goldie, Thang Luong, and Quoc Le. 2017. Massive Exploration of Neural Machine Translation Architectures. *ArXiv e-prints* (March 2017). arXiv:1703.03906 [cs.CL]

[40] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* (2020).

[41] Melissa G Bublitz, Tracy Rank-Christman, Luca Cian, Xavier Cortada, Adriana Madzharov, Vanessa M Patrick, Laura A Peracchio, Maura L Scott, Aparna Sundar, Ngoc To, et al. 2019. Collaborative art: A transformational force within communities. *Journal of the Association for Consumer Research* 4, 4 (2019), 313–331.

[42] Weiwei Cai and Zhanguo Wei. 2020. PiiGAN: generative adversarial networks for pluralistic image inpainting. *IEEE Access* 8 (2020), 48451–48463.

[43] Guendalina Caldarini, Sardar Jaf, and Kenneth McGarry. 2022. A literature survey of recent advances in chatbots. *Information* 13, 1 (2022), 41.

[44] Colin Campbell, Kirk Plangger, Sean Sands, and Jan Kietzmann. 2022. Preparing for an era of deepfakes and AI-generated ads: A framework for understanding responses to manipulated advertising. *Journal of Advertising* 51, 1 (2022), 22–38.

[45] Hanqun Cao, Cheng Tan, Zhangyang Gao, Guangyong Chen, Pheng-Ann Heng, and Stan Z Li. 2022. A survey on generative diffusion model. *arXiv preprint arXiv:2209.02646* (2022).

[46] Nicholas Carlini and David Wagner. 2018. Audio adversarial examples: Targeted attacks on speech-to-text. In *2018 IEEE security and privacy workshops (SPW)*. IEEE, 1–7.

[47] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. 2022. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16123–16133.

[48] Paramanand Chandramouli and Kanchana Vaishnavi Gandikota. 2022. LDEdit: Towards Generalized Text Guided Image Manipulation via Latent Diffusion Models. *arXiv preprint arXiv:2210.02249* (2022).

[49] Anna Chaney. 2018. The Watson Beat: Using Machine Learning to Inspire Musical Creativity. *https://medium.com/@anna_seg/the-watson-beat-d7497406a202* (2018).

[50] Rong Chang, Xinmiao Song, and Huiwen Liu. 2022. Between Shanshui and Landscape: An AI Aesthetics Study Connecting Chinese and Western Paintings. In *HCI International 2022 Posters: 24th International Conference on Human-Computer Interaction, HCII 2022, Virtual Event, June 26–July 1, 2022, Proceedings, Part III*. Springer, 179–185.

[51] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. 2019. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7832–7841.

[52] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, Najim Dehak, and William Chan. 2021. WaveGrad 2: Iterative refinement for text-to-speech synthesis. *arXiv preprint arXiv:2106.09660* (2021).

[53] Shu-Ching Chen. 2022. Multimedia research toward the Metaverse. *IEEE MultiMedia* 29, 1 (2022), 125–127.

[54] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *ICML*.

[55] Xinlei Chen and Kaiming He. 2021. Exploring simple siamese representation learning. In *CVPR*.

[56] Xinpeng Chen, Lin Ma, Wenhao Jiang, Jian Yao, and Wei Liu. 2018. Regularizing rnns for caption generation by reconstructing the past with the present. In *Proceedings of the IEEE Conference on computer vision and pattern recognition*. 7995–8003.

[57] Zehua Chen, Xu Tan, Ke Wang, Shifeng Pan, Danilo Mandic, Lei He, and Sheng Zhao. 2022. Infergrad: Improving Diffusion Models for Vocoder by Considering Inference in Training. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 8432–8436.

[58] Yong Cheng. 2019. Joint training for pivot-based neural machine translation. In *Joint training for neural machine translation*. Springer, 41–54.

[59] Zezhou Cheng, Qingxiong Yang, and Bin Sheng. 2015. Deep colorization. In *Proceedings of the IEEE international conference on computer vision*. 415–423.

[60] Anton Cherepkov, Andrey Voynov, and Artem Babenko. 2021. Navigating the gan parameter space for semantic image editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3671–3680.

[61] Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonina, et al. 2018. State-of-the-art speech recognition with sequence-to-sequence models. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 4774–4778.

[62] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014).

[63] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 2016. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *European conference on computer vision*. Springer, 628–644.

[64] Chieh-Yu Chung and Szu-Hao Huang. 2022. Interactively transforming Chinese ink paintings into realistic images using a border enhance generative adversarial network. *Multimedia Tools and Applications* (2022), 1–34.

[65] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).

[66] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. 2017. You said that? *arXiv preprint arXiv:1705.02966* (2017).

[67] Aidan Clark, Jeff Donahue, and Karen Simonyan. 2019. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571* (2019).

[68] Aidan Clark, Jeff Donahue, and Karen Simonyan. 2019. Efficient video generation on complex datasets. (2019).

[69] Ronan Collobert, Christian Puhrsch, and Gabriel Synnaeve. 2016. Wav2letter: an end-to-end convnet-based speech recognition system. *arXiv preprint arXiv:1609.03193* (2016).

[70] @ColourlabAI. 2023. Powering those 'I made that' moments. *https://colourlab.ai/* (2023).

[71] Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli. 2020. Unsupervised cross-lingual representation learning for speech recognition. *arXiv preprint arXiv:2006.13979* (2020).

[72] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. 2022. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427* (2022).

[73] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. 2022. Diffusion models in vision: A survey. *arXiv preprint arXiv:2209.04747* (2022).

[74] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. 2019. Capture, learning, and synthesis of 3D speaking styles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10101–10111.

[75] George E Dahl, Dong Yu, Li Deng, and Alex Acero. 2011. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *IEEE Transactions on audio, speech, and language processing* 20, 1 (2011), 30–42.

[76] Hanjun Dai, Azade Nazi, Yujia Li, Bo Dai, and Dale Schuurmans. 2020. Scalable deep generative modeling for sparse graphs. In *International Conference on Machine Learning*. PMLR, 2302–2312.

[77] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. 2022. Robust deep learning–based protein sequence design using ProteinMPNN. *Science* 378, 6615 (2022), 49–56.

[78] Marc Davis. 1994. Media streams: representing video for retrieval and repurposing. In *Proceedings of the second ACM international conference on Multimedia*. 478–479.

[79] Marcelo Zerwes Dawn Gilmore, Anitra Nottingham. 2023. ChatGPT and learning design: what online content creation opportunities does it offer? *https://www.csmonitor.com/Technology/2023/0217/Tremendous-potential-Why-some-disability-advocates-laud-ChatGPT* (2023).

[80] Nicola De Cao and Thomas Kipf. 2018. MolGAN: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973* (2018).

[81] Paul DelSignore. 2022. AI Art Wins Competition And Sparks Controversy. *https://medium.com/mlearning-ai/ai-art-wins-fine-arts-competition-and-sparks-controversy-882f9b4df98c* (2022).

[82] Rick DeMott. 2006. Mova Contour Moves "Motion Capture To Reality Capture". *https://www.awn.com/news/mova-contour-moves-motion-capture-reality-capture* (2006).

[83] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *CVPR*.

[84] Shasha Deng, Chee-Wee Tan, Weijun Wang, and Yu Pan. 2019. Smart generation system of personalized advertising copy and its application to advertising practice and research. *Journal of Advertising* 48, 4 (2019), 356–365.

[85] Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echegoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2021. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review* 54, 1 (2021), 755–810.

[86] @descript. 2022. There's a new way to make video and podcasts. A good way. *https://www.descript.com/* (2022).

[87] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[88] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. 2021. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems* 34 (2021), 19822–19835.

[89] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. 2022. CogView2: Faster and Better Text-to-Image Generation via Hierarchical Transformers. *arXiv preprint arXiv:2204.14217* (2022).

[90] Laurent Dinh, David Krueger, and Yoshua Bengio. 2015. Nice: Non-linear independent components estimation. *ICLR 2015 Workshop Track* (2015).

[91] Thomas Dohmke. 2022. GitHub Copilot is generally available to all developers. *https://github.blog/2022-06-21-github-copilot-is-generally-available-to-all-developers/* (2022).

[92] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2625–2634.

[93] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. 2014. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*. Springer, 184–199.

[94] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. 2015. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence* 38, 2 (2015), 295–307.

[95] Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. 2018. Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.

[96] Gianfranco Doretto, Alessandro Chiuso, Ying Nian Wu, and Stefano Soatto. 2003. Dynamic textures. *International Journal of Computer Vision* 51, 2 (2003), 91–109.

[97] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[98] Eric Drott. 2021. Copyright, compensation, and commons in the music AI industry. *Creative Industries Journal* 14, 2 (2021), 190–207.

[99] Henry Elder, Alexander O'Connor, and Jennifer Foster. 2020. How to make neural natural language generation as reliable as templates in task-oriented dialogue. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2877–2888.

[100] Ariel Ephrat, Tavi Halperin, and Shmuel Peleg. 2017. Improved speech reconstruction from silent video. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 455–462.

[101] Ariel Ephrat and Shmuel Peleg. 2017. Vid2speech: speech reconstruction from silent video. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 5095–5099.

[102] Yuchen Fan, Yao Qian, Feng-Long Xie, and Frank K Soong. 2014. TTS synthesis with bidirectional LSTM based recurrent neural networks. In *Fifteenth annual conference of the international speech communication association*.

[103] Jinsheng Fang, Hanjiang Lin, Xinyu Chen, and Kun Zeng. 2022. A Hybrid Network of CNN and Transformer for Lightweight Image Super-Resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1103–1112.

[104] Jean Louis K E Fendji, Diane CM Tala, Blaise O Yenke, and Marcellin Atemkeng. 2022. Automatic speech recognition using limited vocabulary: A survey. *Applied Artificial Intelligence* 36, 1 (2022), 2095039.

[105] Shaoxiong Feng, Xuancheng Ren, Hongshen Chen, Bin Sun, Kan Li, and Xu Sun. 2020. Regularizing dialogue generation by imitating implicit scenarios. *arXiv preprint arXiv:2010.01893* (2020).

[106] Souheil Fenghour, Daqing Chen, Kun Guo, Bo Li, and Perry Xiao. 2021. Deep learning-based automated lip-reading: A survey. *IEEE Access* (2021).

[107] Feyyaz Fırat. 2019. Robot journalism. *The International Encyclopedia of Journalism Studies* (2019), 1–5.

[108] Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation* 25, 2 (2011), 127–144.

[109] Danny Francis and Benoit Huet. 2021. Image and Video Captioning Using Deep Architectures. *Multi-faceted Deep Learning: Models and Data* (2021), 151–174.

[110] Emma Frid, Celso Gomes, and Zeyu Jin. 2020. Music creation by example. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–13.

[111] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. 2019. Text-based editing of talking-head video. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–14.

[112] Rao Fu, Xiao Zhan, Yiwen Chen, Daniel Ritchie, and Srinath Sridhar. 2022. Shapecrafter: A recursive text-conditioned 3d shape generation model. *arXiv preprint arXiv:2207.09446* (2022).

[113] Yocat Games. 2023. Traveler-The AI Story. *https://skidrowcracked.com/traveler-the-ai-story/* (2023).

[114] Zhe Gan, Ricardo Henao, David Carlson, and Lawrence Carin. 2015. Learning deep sigmoid belief networks with data augmentation. In *Artificial Intelligence and Statistics*. PMLR, 268–276.

[115] Zhe Gan, Linjie Li, Chunyuan Li, Lijuan Wang, Zicheng Liu, Jianfeng Gao, et al. 2022. Vision-language pre-training: Basics, recent advances, and future trends. *Foundations and Trends® in Computer Graphics and Vision* 14, 3–4 (2022), 163–352.

[116] Ruiqi Gao, Yang Lu, Junpei Zhou, Song-Chun Zhu, and Ying Nian Wu. 2018. Learning generative convnets via multi-grid modeling and sampling. 9155–9164.

[117] Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, and Bill Dolan. 2020. Dialogue response ranking training with large-scale human feedback data. *arXiv preprint arXiv:2009.06978* (2020).

[118] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2015. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576* (2015).

[119] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image style transfer using convolutional neural networks. In *CVPR*.

[120] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *International conference on machine learning*. PMLR, 1243–1252.

[121] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. 2018. Unsupervised representation learning by predicting image rotations. *ICLR* (2018).

[122] Vladislav Golyanik, Soshi Shimada, Kiran Varanasi, and Didier Stricker. 2018. Hdm-net: Monocular non-rigid 3d reconstruction with learned deformation model. In *International conference on virtual reality and augmented reality*. Springer, 51–72.

[123] RC Gonzalez. 2006. Woods. RE,(2002)"Digital Image Processing".

[124] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NeurIPS*.

[125] Andrew S Gordon and Eric A Domeshek. 1995. Conceptual indexing for video retrieval. In *Working Notes of IJCAI Workshop on Intelligent Multimedia Information Retrieval, Montreal*. 23–38.

[126] Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850* (2013).

[127] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*. Ieee, 6645–6649.

[128] Ulf Grenander and Michael I Miller. 1994. Representations of knowledge in complex systems. *Journal of the Royal Statistical Society: Series B (Methodological)* 56, 4 (1994), 549–581.

[129] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. 2020. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems* (2020).

[130] Barbara Gruber. 2022. Facts, Fakes and Figures: How AI is Influencing Journalism. *https://www.goethe.de/prj/k40/en/lan/aij.html* (2022).

[131] Shanyan Guan, Ying Tai, Bingbing Ni, Feida Zhu, Feiyue Huang, and Xiaokang Yang. 2020. Collaborative learning for faster stylegan embedding. *arXiv preprint arXiv:2007.01758* (2020).

[132] Longteng Guo, Jing Liu, Xinxin Zhu, Peng Yao, Shichen Lu, and Hanqing Lu. 2020. Normalized and geometry-aware self-attention network for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10327–10336.

[133] Shunan Guo, Zhuochen Jin, Fuling Sun, Jingwen Li, Zhaorui Li, Yang Shi, and Nan Cao. 2021. Vinci: an intelligent graphic design system for generating advertising posters. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–17.

[134] Xiujing Guo. 2021. Towards automated software testing with generative adversarial networks. In *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks-Supplemental Volume (DSN-S)*. IEEE, 21–22.

[135] Xiujing Guo, Hiroyuki Okamura, and Tadashi Dohi. 2022. Automated Software Test Data Generation With Generative Adversarial Networks. *IEEE Access* 10 (2022), 20690–20700.

[136] Tanmay Gupta, Dustin Schwenk, Ali Farhadi, Derek Hoiem, and Aniruddha Kembhavi. 2018. Imagine this! scripts to compositions to videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 598–613.

[137] Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*.

[138] Kilian Konstantin Haefeli, Karolis Martinkus, Nathanaël Perraudin, and Roger Wattenhofer. 2022. Diffusion Models for Graphs Benefit From Discrete State Spaces. *arXiv preprint arXiv:2210.01549* (2022).

[139] Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang, and Kee-Eung Kim. 2020. End-to-end neural pipeline for goal-oriented dialogue systems using GPT-2. In *Proceedings of the 58th annual meeting of the association for computational linguistics*. 583–592.

[140] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. 2019. Meshcnn: a network with an edge. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.

[141] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *CVPR*.

[142] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *CVPR*.

[143] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.

[144] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. 2022. Is synthetic data from generative models ready for image recognition? *arXiv preprint arXiv:2210.07574* (2022).

[145] Wanwei He, Min Yang, Rui Yan, Chengming Li, Ying Shen, and Ruifeng Xu. 2020. Amalgamating knowledge from two teachers for task-oriented dialogue system with adversarial training. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 3498–3507.

[146] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. 2019. Attgan: Facial attribute editing by only changing what you want. *IEEE transactions on image processing* 28, 11 (2019), 5464–5478.

[147] Brian Heater. 2022. Amper is providing a plug-and-play-solution to digitize manufacturing. *https://techcrunch.com/2022/04/22/amper-is-providing-a-plug-and-play-solution-to-digitize-manufacturing/* (2022).

[148] Matthew Henderson, Ivan Vulić, Daniela Gerz, Iñigo Casanueva, Paweł Budzianowski, Sam Coope, Georgios Spithourakis, Tsung-Hsien Wen, Nikola Mrkšić, and Pei-Hao Su. 2019. Training neural response selection for task-oriented dialogue systems. *arXiv preprint arXiv:1906.01543* (2019).

[149] Simao Herdade, Armin Kappeler, Kofi Boakye, and Joao Soares. 2019. Image captioning: Transforming objects into words. *Advances in neural information processing systems* 32 (2019).

[150] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2022. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022).

[151] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).

[152] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. 2012. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine* 29, 6 (2012), 82–97.

[153] Geoffrey E Hinton. 2002. Training products of experts by minimizing contrastive divergence. *Neural computation* 14, 8 (2002), 1771–1800.

[154] Geoffrey E Hinton and Richard Zemel. 1993. Autoencoders, minimum description length and Helmholtz free energy. *NeurIPS* (1993).

[155] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. 2022. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303* (2022).

[156] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.

[157] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. 2022. Video diffusion models. *arXiv preprint arXiv:2204.03458* (2022).

[158] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.

[159] Margaret Holland. 2016. How YouTube developed into a successful platform for user-generated content. *Elon journal of undergraduate research in communications* 7, 1 (2016).

[160] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. 2022. CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers. *arXiv preprint arXiv:2205.15868* (2022).

[161] Takaaki Hori, Jaejin Cho, and Shinji Watanabe. 2018. End-to-end speech recognition with word-based RNN language models. In *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 389–396.

[162] Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *Advances in Neural Information Processing Systems* 33 (2020), 20179–20191.

[163] Stephanie Houde, Vera Liao, Jacquelyn Martino, Michael Muller, David Piorkowski, John Richards, Justin Weisz, and Yunfeng Zhang. 2020. Business (mis) use cases of generative ai. *arXiv preprint arXiv:2003.07679* (2020).

[164] Zhen-jiang Hu. 2022. Analysis of the Impact of Artificial Intelligence Technology-Assisted Environmental Protection on the Integrity of Chinese Painting. *Journal of Environmental and Public Health* 2022 (2022).

[165] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. 2017. Densely connected convolutional networks. In *CVPR*.

[166] Han Huang, Leilei Sun, Bowen Du, Yanjie Fu, and Weifeng Lv. 2022. GraphGDP: Generative Diffusion Processes for Permutation Invariant Graph Generation. *arXiv preprint arXiv:2212.01842* (2022).

[167] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. 2018. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2821–2830.

[168] Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren. 2022. ProDiff: Progressive Fast Diffusion Model For High-Quality Text-to-Speech. *arXiv preprint arXiv:2207.06389* (2022).

[169] Xinting Huang, Jianzhong Qi, Yu Sun, and Rui Zhang. 2020. Semi-supervised dialogue policy learning via stochastic reward estimation. *arXiv preprint arXiv:2005.04379* (2020).

[170] Yu-Siang Huang and Yi-Hsuan Yang. 2020. Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions. In *Proceedings of the 28th ACM International Conference on Multimedia*. 1180–1188.

[171] William John Hutchins. 1986. *Machine translation: past, present, future*. Ellis Horwood Chichester.

[172] Matthew Hutson. 2017. How Google is making music with arti4cial intelligence. (2017).

[173] Nitin Indurkhya and Fred J Damerau. 2010. *Handbook of natural language processing*. Chapman and Hall/CRC.

[174] Prathamesh Ingle. 2023.          Top Artificial Intelligence (AI) Tools That Can Generate Code To Help Programmers. *https://www.marktechpost.com/2023/01/01/top-artificial-intelligence-ai-tools-that-can-generate-code-to-help-programmers/* (2023).

[175] Tansin Jahan, Yanran Guan, and Oliver van Kaick. 2021. Semantics-Guided Latent Space Exploration for Shape Generation. In *Computer Graphics Forum*, Vol. 40. Wiley Online Library, 115–126.

[176] Benjamin Jahić, Nicolas Guelfi, and Benoit Ries. 2019. Software engineering for dataset augmentation using generative adversarial networks. In *2019 IEEE 10th International Conference on Software Engineering and Service Science (ICSESS)*. IEEE, 59–66.

[177] Viren Jain and Sebastian Seung. 2008. Natural image denoising with convolutional networks. *Advances in neural information processing systems* 21 (2008).

[178] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. 2019. You said that?: Synthesising talking faces from audio. *International Journal of Computer Vision* 127 (2019), 1767–1779.

[179] Shulei Ji, Jing Luo, and Xinyu Yang. 2020. A comprehensive survey on deep music generation: Multi-level representations, algorithms, evaluations, and future directions. *arXiv preprint arXiv:2011.06801* (2020).

[180] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*.

[181] Qi Jia, Yizhu Liu, Siyu Ren, Kenny Q Zhu, and Haifeng Tang. 2020. Multi-turn response selection using dialogue dependency relations. *arXiv preprint arXiv:2010.01502* (2020).

[182] Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is ChatGPT a good translator? A preliminary study. *arXiv preprint arXiv:2301.08745* (2023).

[183] Wengong Jin, Kevin Yang, Regina Barzilay, and Tommi Jaakkola. 2018. Learning multimodal graph-to-graph translation for molecular optimization. *arXiv preprint arXiv:1812.01070* (2018).

[184] Li Jing, Caglar Gulcehre, John Peurifoy, Yichen Shen, Max Tegmark, Marin Soljacic, and Yoshua Bengio. 2019. Gated orthogonal recurrent units: On learning to forget. *Neural computation* 31, 4 (2019), 765–783.

[185] Yongcheng Jing, Yezhou Yang, Zunlei Feng, Jingwen Ye, Yizhou Yu, and Mingli Song. 2019. Neural style transfer: A review. *IEEE transactions on visualization and computer graphics* 26, 11 (2019), 3365–3385.

[186] Jaehyeong Jo, Seul Lee, and Sung Ju Hwang. 2022. Score-based Generative Modeling of Graphs via the System of Stochastic Differential Equations. *arXiv preprint arXiv:2202.02514* (2022).

[187] Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics* 5 (2017), 339–351.

[188] Biing Hwang Juang and Laurence R Rabiner. 1991. Hidden Markov models for speech recognition. *Technometrics* 33, 3 (1991), 251–272.

[189] Romain Juillet. 2021. How AI Is Transforming the Music Industry. *https://www.bocasay.com/ai-transforming-music-industry/* (2021).

[190] Łukasz Kaiser and Samy Bengio. 2016. Can active memory replace attention? *Advances in Neural Information Processing Systems* 29 (2016).

[191] Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of the 2013 conference on empirical methods in natural language processing*. 1700–1709.

[192] Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aaron van den Oord, Alex Graves, and Koray Kavukcuoglu. 2016. Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099* (2016).

[193] Shitanshu kapadia. 2023. 7 AI Marketing Tools for Business. *https://www.entrepreneur.com/science-technology/how-will-chatgpt-change-education-and-teaching/445018* (2023).

[194] S Karpagavalli and Edy Chandra. 2016. A review on automatic speech recognition architecture and approaches. *International Journal of Signal Processing, Image Processing and Pattern Recognition* 9, 4 (2016), 393–404.

[195] Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3128–3137.

[196] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. 2017. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–12.

[197] Tero Karras, Samuli Laine, and Timo Aila. 2019. A Style-Based Generator Architecture for Generative Adversarial Networks. 4401–4410.

[198] Hideki Kawahara. 2006. STRAIGHT, exploitation of the other aspect of VOCODER: Perceptually isomorphic decomposition of speech sounds. *Acoustical science and technology* 27, 6 (2006), 349–353.

[199] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. 2022. Denoising diffusion restoration models. *arXiv preprint arXiv:2201.11793* (2022).

[200] Lei Ke, Wenjie Pei, Ruiyu Li, Xiaoyong Shen, and Yu-Wing Tai. 2019. Reflective decoding network for image captioning. In *Proceedings of the IEEE/CVF international conference on computer vision*. 8888–8897.

[201] Regina Kelder and Mary Parker. 2021. Why Does Drug Development Take So Long? *https://www.criver.com/eureka/why-does-drug-development-take-so-long* (2021).

[202] Tilly Kenyon. 2022. Voicemod: Allowing creators to find their voice. *https://technologymagazine.com/ai-and-machine-learning/voicemod-allowing-creators-to-find-their-voice* (2022).

[203] Angella J Kim and Kim KP Johnson. 2016. Power of consumers using social media: Examining the influences of brand-related user-generated content on Facebook. *Computers in human behavior* 58 (2016), 98–108.

[204] Daewon Kim and Seongcheol Kim. 2017. Newspaper companies' determinants in adopting robot journalism. *Technological Forecasting and Social Change* 117 (2017), 184–195.

[205] Gwanghyun Kim and Se Young Chun. 2022. DATID-3D: Diversity-Preserved Domain Adaptation Using Text-to-Image Diffusion for 3D Generative Model. *arXiv preprint arXiv:2211.16374* (2022).

[206] Gwanghyun Kim and Jong Chul Ye. 2021. Diffusionclip: Text-guided image manipulation using diffusion models. (2021).

[207] Sijin Kim, Namhyuk Ahn, and Kyung-Ah Sohn. 2020. Restoring spatially-heterogeneous distortions using mixture of experts network. In *Proceedings of the Asian Conference on Computer Vision.*

[208] Sungdong Kim, Sohee Yang, Gyuwan Kim, and Sang-Woo Lee. 2019. Efficient dialogue state tracking by selectively overwriting memory. *arXiv preprint arXiv:1911.03906* (2019).

[209] Taeksoo Kim, Byoungjip Kim, Moonsu Cha, and Jiwon Kim. 2017. Unsupervised visual attribute transfer with reconfigurable generative adversarial networks. *arXiv preprint arXiv:1707.09798* (2017).

[210] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

[211] Wei-Jen Ko, Avik Ray, Yilin Shen, and Hongxia Jin. 2020. Generating dialogue responses from a semantic latent space. *arXiv preprint arXiv:2010.01658* (2020).

[212] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions.* 177–180.

[213] Philipp Koehn, Franz J Och, and Daniel Marcu. 2003. *Statistical phrase-based translation.* Technical Report. University of Southern California Marina Del Rey Information Sciences Inst.

[214] Marijn Koolen, Jaap Kamps, Gabriella Kazai, et al. 2013. Social Book Search: The Impact of Professional and User-Generated Content on Book Suggestions.. In *DIR.* 38–39.

[215] Prajwal KR, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and CV Jawahar. 2019. Towards automatic face-to-face translation. In *Proceedings of the 27th ACM international conference on multimedia.* 1428–1436.

[216] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet classification with deep convolutional neural networks. In *NeurIPS.*

[217] BJ Kröger. 1992. Minimal rules for articulatory speech synthesis. *Proceedings of EUSIPCO92 (1)* (1992), 331–334.

[218] Manoj Kumar, Dirk Weissenborn, and Nal Kalchbrenner. 2021. Colorization transformer. *arXiv preprint arXiv:2102.04432* (2021).

[219] Rithesh Kumar, Jose Sotelo, Kundan Kumar, Alexandre de Brébisson, and Yoshua Bengio. 2017. Obamanet: Photo-realistic lip-sync from text. *arXiv preprint arXiv:1801.01442* (2017).

[220] Max WY Lam, Jun Wang, Dan Su, and Dong Yu. 2022. BDDM: Bilateral Denoising Diffusion Models for Fast and High-Quality Speech Synthesis. *arXiv preprint arXiv:2203.13508* (2022).

[221] Hugo Larochelle and Iain Murray. 2011. The neural autoregressive distribution estimator. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics.* JMLR Workshop and Conference Proceedings, 29–37.

[222] Noam Lemelshtrich Latar. 2018. *Robot journalism: Can human journalism survive?* World Scientific.

[223] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. 2017. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 4681–4690.

[224] Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos, and Zoubin Ghahramani. 2010. Kronecker graphs: an approach to modeling networks. *Journal of Machine Learning Research* 11, 2 (2010).

[225] Stephen E Levinson, Lawrence R Rabiner, and M Mohan Sondhi. 1983. An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition. *Bell System Technical Journal* 62, 4 (1983), 1035–1074.

[226] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* (2019).

[227] Colin Lewis-Beck and Michael Lewis-Beck. 2015. *Applied regression: An introduction.* Vol. 22. Sage publications.

[228] Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu, Jiancheng Lv, and Xi Peng. 2022. All-in-one image restoration for unknown corruption. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 17452–17462.

[229] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip Torr. 2019. Controllable text-to-image generation. *Advances in Neural Information Processing Systems* 32 (2019).

[230] Gang Li, Heliang Zheng, Chaoyue Wang, Chang Li, Changwen Zheng, and Dacheng Tao. 2022. 3DDesigner: Towards Photorealistic 3D Object Generation and Editing with Text-guided Diffusion Models. *arXiv preprint arXiv:2211.14108* (2022).

[231] Guang Li, Linchao Zhu, Ping Liu, and Yi Yang. 2019. Entangled transformer for image captioning. In *Proceedings of the IEEE/CVF international conference on computer vision.* 8928–8937.

[232] Haoying Li, Yifan Yang, Meng Chang, Huajun Feng, Zhi hai Xu, Qi Li, and Yue ting Chen. 2022. SRDiff: Single Image Super-Resolution with Diffusion Probabilistic Models. *Neurocomputing* 479 (2022), 47–59.

[233] Jason Li, Vitaly Lavrukhin, Boris Ginsburg, Ryan Leary, Oleksii Kuchaiev, Jonathan M Cohen, Huyen Nguyen, and Ravi Teja Gadde. 2019. Jasper: An end-to-end convolutional neural acoustic model. *arXiv preprint arXiv:1904.03288* (2019).

[234] Jinyu Li, Yu Wu, Yashesh Gaur, Chengyi Wang, Rui Zhao, and Shujie Liu. 2020. On the comparison of popular end-to-end models for large scale speech recognition. *arXiv preprint arXiv:2005.14327* (2020).

[235] Mao Li, Jiancheng Lv, Jian Wang, and Yongsheng Sang. 2020. An abstract painting generation method based on deep generative model. *Neural Processing Letters* 52 (2020), 949–960.

[236] Mu Li, Wangmeng Zuo, and David Zhang. 2016. Convolutional network for attribute-driven and identity-preserving human face generation. *arXiv preprint arXiv:1608.06434* (2016).

[237] Mu Li, Wangmeng Zuo, and David Zhang. 2016. Deep identity-aware transfer of facial attributes. *arXiv preprint arXiv:1610.05586* (2016).

[238] Ruilong Li, Matthew Tancik, and Angjoo Kanazawa. 2022. NerfAcc: A General NeRF Acceleration Toolbox. *arXiv preprint arXiv:2210.04847* (2022).

[239] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Trans. Graph.* 36, 6 (2017), 194–1.

[240] Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. 2022. MAT: Mask-Aware Transformer for Large Hole Image Inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10758–10768.

[241] Wei Li, Xue Xu, Xinyan Xiao, Jiachen Liu, Hu Yang, Guohao Li, Zhanpeng Wang, Zhifan Feng, Qiaoqiao She, Yajuan Lyu, et al. 2022. UPainting: Unified Text-to-Image Diffusion Generation with Cross-modal Guidance. *arXiv preprint arXiv:2210.16031* (2022).

[242] Xin Li, Xin Jin, Jianxin Lin, Sen Liu, Yaojun Wu, Tao Yu, Wei Zhou, and Zhibo Chen. 2020. Learning disentangled feature representation for hybrid-distorted image restoration. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*. Springer, 313–329.

[243] Xi Li and Ping Kuang. 2021. 3D-VRVT: 3D Voxel Reconstruction from A Single Image with Vision Transformer. In *2021 International Conference on Culture-oriented Science & Technology (ICCST)*. IEEE, 343–348.

[244] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*. Springer, 121–137.

[245] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B Hashimoto. 2022. Diffusion-LM Improves Controllable Text Generation. *arXiv preprint arXiv:2205.14217* (2022).

[246] Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. 2018. Video generation from text. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.

[247] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. 2021. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1833–1844.

[248] Jianglin Liang and Ruifang Liu. 2015. Stacked denoising autoencoder and dropout together to prevent overfitting in deep neural network. In *2015 8th international congress on image and signal processing (CISP)*. IEEE, 697–701.

[249] Renjie Liao, Yujia Li, Yang Song, Shenlong Wang, Will Hamilton, David K Duvenaud, Raquel Urtasun, and Richard Zemel. 2019. Efficient graph generation with graph recurrent attention networks. *Advances in neural information processing systems* 32 (2019).

[250] Zibo Lin, Deng Cai, Yan Wang, Xiaojiang Liu, Hai-Tao Zheng, and Shuming Shi. 2020. The world is not binary: Learning to rank with grayscale data for dialogue response selection. *arXiv preprint arXiv:2004.02421* (2020).

[251] Pierre Lison and Serge Bibauw. 2017. Not all dialogues are created equal: Instance weighting for neural conversational models. *arXiv preprint arXiv:1704.08966* (2017).

[252] Hongyu Liu, Ziyu Wan, Wei Huang, Yibing Song, Xintong Han, and Jing Liao. 2021. Pd-gan: Probabilistic diverse gan for image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9371–9381.

[253] Ming Liu, Yuxiang Wei, Xiaohe Wu, Wangmeng Zuo, and Lei Zhang. 2022. A Survey on Leveraging Pre-trained Generative Adversarial Networks for Image Editing and Restoration. *arXiv preprint arXiv:2207.10309* (2022).

[254] Qi Liu, Miltiadis Allamanis, Marc Brockschmidt, and Alexander Gaunt. 2018. Constrained graph variational autoencoders for molecule design. *Advances in neural information processing systems* 31 (2018).

[255] Shuo Liu, Adria Mallol-Ragolta, Emilia Parada-Cabeleiro, Kun Qian, Xin Jing, Alexander Kathan, Bin Hu, and Bjoern W Schuller. 2022. Audio Self-supervised Learning: A Survey. *arXiv preprint arXiv:2203.01205* (2022).

[256] Wei Liu, Sihan Chen, Longteng Guo, Xinxin Zhu, and Jing Liu. 2021. CPTR: Full transformer network for image captioning. *arXiv preprint arXiv:2101.10804* (2021).

[257] Xia Liu, Alvin C Burns, and Yingjian Hou. 2017. An investigation of brand-related user-generated content on Twitter. *Journal of Advertising* 46, 2 (2017), 236–247.

[258] Xing Liu, Masanori Suganuma, Xiyang Luo, and Takayuki Okatani. 2019. Restoring images with unknown degradation factors by recurrent use of a multi-branch network. *arXiv preprint arXiv:1907.04508* (2019).

[259] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

[260] Yue Liu, Xin Wang, Yitian Yuan, and Wenwu Zhu. 2019. Cross-modal dual learning for sentence-to-video generation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 1239–1247.

[261] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. *ICCV*.

[262] Zhengzhe Liu, Yi Wang, Xiaojuan Qi, and Chi-Wing Fu. 2022. Towards Implicit Text-Guided 3D Shape Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17896–17906.

[263] Dengsheng Lu and Qihao Weng. 2007. A survey of image classification methods and techniques for improving classification performance. *International journal of Remote sensing* 28, 5 (2007), 823–870.

[264] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 375–383.

[265] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. 2022. Repaint: Inpainting using denoising diffusion probabilistic models. 11461–11471.

[266] Tianze Luo, Zhanfeng Mo, and Sinno Jialin Pan. 2022. Fast Graph Generative Model via Spectral Diffusion. *arXiv preprint arXiv:2211.08892* (2022).

[267] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025* (2015).

[268] Christoph Lüscher, Eugen Beck, Kazuki Irie, Markus Kitza, Wilfried Michel, Albert Zeyer, Ralf Schlüter, and Hermann Ney. 2019. RWTH ASR Systems for LibriSpeech: Hybrid vs Attention–w/o Data Augmentation. *arXiv preprint arXiv:1905.03072* (2019).

[269] Kaushalya Madhawa, Katushiko Ishiguro, Kosuke Nakago, and Motoki Abe. 2019. Graphnvp: An invertible flow model for generating molecular graphs. *arXiv preprint arXiv:1905.11600* (2019).

[270] Mishaim Malik, Muhammad Kamran Malik, Khawar Mehmood, and Imran Makhdoom. 2021. Automatic speech recognition: a survey. *Multimedia Tools and Applications* 80, 6 (2021), 9411–9457.

[271] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. 2016. Generating images from captions with attention. *ICLR* (2016).

[272] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan Yuille. 2014. Deep captioning with multimodal recurrent neural networks (m-rnn). *arXiv preprint arXiv:1412.6632* (2014).

[273] Francesco Marconi. 2020. *Newsmakers: Artificial intelligence and the future of journalism.* Columbia University Press.

[274] Rick Marshall. 2018. Behind the breathtaking visual effects of 'Blade Runner 2049'. *https://www.digitaltrends.com/movies/blade-runner-2049-visual-effects-john-nelson/* (2018).

[275] Richard Diehl Martinez and John Kaleialoha Kamalu. 2018. Using General Adversarial Networks for Marketing: A Case Study of Airbnb. *arXiv preprint arXiv:1806.11432* (2018).

[276] Tanya Marwah, Gaurav Mittal, and Vineeth N Balasubramanian. 2017. Attentive semantic video generation using captions. In *Proceedings of the IEEE international conference on computer vision*. 1426–1434.

[277] Alexandra Kennedy Maya Ackerman. 2022. AI-Powered Lyrics Platform, LyricStudio, Surpasses 1 Million Songs. *prweb.com/releases/2022/7/prweb18785303.htm* (2022).

[278] Łukasz Maziarka, Agnieszka Pocha, Jan Kaczmarczyk, Krzysztof Rataj, Tomasz Danel, and Michał Warchoł. 2020. Mol-CycleGAN: a generative model for molecular optimization. *Journal of Cheminformatics* 12, 1 (2020), 1–18.

[279] Alex McFarland. 2022. China's State News Agency Introduces New Artificial Intelligence Anchor. *https://www.unite.ai/chinas-state-news-agency-introduces-new-artificial-intelligence-anchor/* (2022).

[280] Alex McFarland. 2023. 8 Best AI Music Generators. *https://www.unite.ai/best-ai-music-generators/* (2023).

[281] Larry R Medsker and LC Jain. 2001. Recurrent neural networks. *Design and Applications* 5 (2001), 64–67.

[282] Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. Pretraining methods for dialog context representation learning. *arXiv preprint arXiv:1906.00414* (2019).

[283] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. 2019. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4460–4470.

[284] Meta. 2022. What Is Meta Horizon Worlds? *https://www.marktechpost.com/2023/01/01/top-artificial-intelligence-ai-tools-that-can-generate-code-to-help-programmers/* (2022).

[285] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernockỳ, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model.. In *Interspeech*, Vol. 2. Makuhari, 1045–1048.

[286] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

[287] Diego H Milone and Leandro E Di Persia. 2008. Learning hidden Markov models with hidden Markov trees as observation distributions. *Inteligencia Artificial. Revista Iberoamericana de Inteligencia Artificial* 12, 37 (2008), 7–13.

[288] Andrey Miroshnichenko. 2018. AI to bypass creativity. Will robots replace journalists?(The answer is "yes"). *Information* 9, 7 (2018), 183.

[289] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. 2021. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091* (2021).

[290] Gaurav Mittal, Tanya Marwah, and Vineeth N Balasubramanian. 2017. Sync-draw: Automatic video generation using deep recurrent attentive architectures. In *Proceedings of the 25th ACM international conference on Multimedia*. 1096–1104.

[291] Irina Momot. 2022. Artificial Intelligence in Filmmaking Process: future scenarios. (2022).

[292] Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. 2016. WORLD: a vocoder-based high-quality speech synthesis system for real-time applications. *IEICE TRANSACTIONS on Information and Systems* 99, 7 (2016), 1877–1884.

[293] Eric Moulines and Francis Charpentier. 1990. Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones. *Speech communication* 9, 5-6 (1990), 453–467.

[294] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* 41, 4 (2022), 1–15.

[295] Frank-Michael Nack. 1996. *AUTEUR: The application of video semantics and theme representation for automated film editing.* Ph. D. Dissertation. Citeseer.

[296] Tomohiro Nakatani. 2019. Improving transformer-based end-to-end speech recognition with connectionist temporal classification and language model integration. In *Proc. Interspeech*.

[297] Ali Bou Nassif, Ismail Shahin, Imtinan Attili, Mohammad Azzeh, and Khaled Shaalan. 2019. Speech recognition using deep neural networks: A systematic review. *IEEE access* 7 (2019), 19143–19165.

[298] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z Qureshi, and Mehran Ebrahimi. 2019. Edgeconnect: Generative image inpainting with adversarial edge learning. *arXiv preprint arXiv:1901.00212* (2019).

[299] Jinjie Ni, Tom Young, Vlad Pandelea, Fuzhao Xue, and Erik Cambria. 2022. Recent advances in deep learning based dialogue systems: A systematic survey. *Artificial intelligence review* (2022), 1–101.

[300] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *ICML* (2022).

[301] Huansheng Ning, Hang Wang, Yujia Lin, Wenxi Wang, Sahraoui Dhelim, Fadi Farha, Jianguo Ding, and Mahmoud Daneshmand. 2021. A Survey on Metaverse: the State-of-the-art, Technologies, Applications, and Challenges. *arXiv preprint arXiv:2111.09673* (2021).

[302] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon. 2020. Permutation invariant graph generation via score-based generative modeling. In *AISTATS*. PMLR, 4474–4484.

[303] Mehdi Noroozi and Paolo Favaro. 2016. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*.

[304] China Academy of Information and Communications Technology. 2022. White Paper on AI-Generated Content (AIGC). *http://www.caict.ac.cn/english/research/whitepapers/202211/t20221111_411288.html* (2022).

[305] Katsunori Ohnishi, Shohei Yamamoto, Yoshitaka Ushiku, and Tatsuya Harada. 2018. Hierarchical video generation from orthogonal information: Optical flow and texture. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.

[306] Joseph Olive. 1977. Rule synthesis of speech from dyadic units. In *ICASSP'77. IEEE International Conference on Acoustics, Speech, and Signal Processing*, Vol. 2. IEEE, 568–570.

[307] OpenAI. 2023. GPT-4 Technical report. *arXiv preprint arXiv:2303.08774* (2023).

[308] Tim O'reilly. 2009. *What is web 2.0.* " O'Reilly Media, Inc.".

[309] Will Oremus. 2014. The First News Report on the L.A. Earthquake Was Written by a Robot. *https://slate.com/technology/2014/03/quakebot-los-angeles-times-robot-journalist-writes-article-on-la-earthquake.html* (2014).

[310] Achraf Oussidi and Azeddine Elhassouny. 2018. Deep generative models: Survey. In *2018 International Conference on Intelligent Systems and Computer Vision (ISCV)*. IEEE, 1–8.

[311] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).

[312] Yawen Ouyang, Moxin Chen, Xinyu Dai, Yinggong Zhao, Shujian Huang, and Jiajun Chen. 2020. Dialogue state tracking with explicit slot connection modeling. In *Proceedings of the 58th annual meeting of the association for computational linguistics*. 34–40.

[313] Yingwei Pan, Zhaofan Qiu, Ting Yao, Houqiang Li, and Tao Mei. 2017. To create what you tell: Generating videos from captions. In *Proceedings of the 25th ACM international conference on Multimedia*. 1789–1798.

[314] Tianyu Pang, Kun Xu, Chongxuan Li, Yang Song, Stefano Ermon, and Jun Zhu. 2020. Efficient learning of generative models via finite-difference score matching. *Advances in Neural Information Processing Systems* 33 (2020), 19175–19188.

[315] Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. 2016. A decomposable attention model for natural language inference. *arXiv preprint arXiv:1606.01933* (2016).

[316] Giorgio Parisi. 1981. Correlation functions and computer simulations. *Nuclear Physics B* 180, 3 (1981), 378–384.

[317] A Parkes. 1989. Settings and the setting structure: the description and automated propagation of networks for perusing videodisk image states. In *Proceedings of the 12th annual international ACM SIGIR conference on Research and development in information retrieval*. 229–238.

[318] Alan Philip Parkes. 1988. *An artificial intelligence approach to the conceptual description of videodisc images.* Lancaster University (United Kingdom).

[319] Alan P Parkes. 1989. The prototype CLORIS system: Describing, retrieving and discussing videodisc stills and sequences. *Information processing & management* 25, 2 (1989), 171–186.

[320] Santiago Pascual, Gautam Bhattacharya, Chunghsin Yeh, Jordi Pons, and Joan Serrà. 2022. Full-band General Audio Synthesis with Score-based Diffusion. *arXiv preprint arXiv:2210.14661* (2022).

[321] Santiago Pascual, Mirco Ravanelli, Joan Serra, Antonio Bonafonte, and Yoshua Bengio. 2019. Learning problem-agnostic speech representations from multiple self-supervised tasks. *arXiv preprint arXiv:1904.03416* (2019).

[322] Steve Paulussen and Pieter Ugille. 2008. User generated content in the newsroom: Professional and organisational constraints on participatory journalism. *Westminster papers in communication & culture* 5, 2 (2008).

[323] Baolin Peng, Chenguang Zhu, Chunyuan Li, Xiujun Li, Jinchao Li, Michael Zeng, and Jianfeng Gao. 2020. Few-shot natural language generation for task-oriented dialog. *arXiv preprint arXiv:2002.12328* (2020).

[324] Tammy Pettinato Oltz. 2023. ChatGPT, Professor of Law. *Professor of Law (February 4, 2023)* (2023).

[325] Cale Plut and Philippe Pasquier. 2020. Generative music in video games: State of the art, challenges, and prospects. *Entertainment Computing* 33 (2020), 100337.

[326] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).

[327] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. 2020. Learning individual speaking styles for accurate lip to speech synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13796–13805.

[328] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. 2020. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*. 484–492.

[329] Ash & Pri. 2022. Scalenut Review: Features and How To Use It as a Content Generator. *https://ashandpri.com/scalenutreview* (2022).

[330] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. 2017. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*.

[331] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. 2017. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems* 30 (2017).

[332] Weizhen Qi, Yu Yan, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and Ming Zhou. 2020. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. *arXiv preprint arXiv:2001.04063* (2020).

[333] Yao Qian, Yuchen Fan, Wenping Hu, and Frank K Soong. 2014. On the training aspects of deep neural network (DNN) for parametric TTS synthesis. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3829–3833.

[334] Long Qin. 2013. *Learning out-of-vocabulary words in automatic speech recognition*. Ph. D. Dissertation. Carnegie Mellon University.

[335] Lisong Qiu, Juntao Li, Wei Bi, Dongyan Zhao, and Rui Yan. 2019. Are training samples correlated? learning to generate dialogue responses with multiple references. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 3826–3835.

[336] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*.

[337] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015).

[338] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[339] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* (2019).

[340] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21, 140 (2020), 1–67.

[341] Salvatore Raieli. 2023. ControlNet: control your AI art generation. *https://levelup.gitconnected.com/controlnet-control-your-ai-art-generation-616c86c88964* (2023).

[342] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).

[343] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *ICML*.

[344] Jeremiah Ratican, James Hutson, and Andrew Wright. 2023. A Proposed Meta-Reality Immersive Development Pipeline: Generative AI Models and Extended Reality (XR) Content for the Metaverse. *Journal of Intelligent Learning Systems and Applications* 15 (2023).

[345] Pratiksha C Raut and Seema U Deoghare. 2016. Automatic Speech Recognition and its Applications. *International Research Journal of Engineering and Technology* 3, 5 (2016), 2368–2371.

[346] Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual, Pawel Swietojanski, Joao Monteiro, Jan Trmal, and Yoshua Bengio. 2020. Multi-task self-supervised learning for robust speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6989–6993.

[347] D Raj Reddy. 1976. Speech recognition by machine: A review. *Proc. IEEE* 64, 4 (1976), 501–531.

[348] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *International conference on machine learning*. PMLR, 1060–1069.

[349] Mengwei Ren, Mauricio Delbracio, Hossein Talebi, Guido Gerig, and Peyman Milanfar. 2022. Image Deblurring with Domain Generalizable Diffusion Models. *arXiv preprint arXiv:2212.01789* (2022).

[350] Shuo Ren, Wenhu Chen, Shujie Liu, Mu Li, Ming Zhou, and Shuai Ma. 2018. Triangular architecture for rare language translation. *arXiv preprint arXiv:1805.04813* (2018).

[351] @rev. 2023. Fast, accurate transcription services. *https://www.rev.com/* (2023).

[352] Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando De la Torre, and Yaser Sheikh. 2021. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1173–1182.

[353] Alan Ritter, Colin Cherry, and Bill Dolan. 2011. Data-driven response generation in social media. In *Empirical Methods in Natural Language Processing (EMNLP)*.

[354] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. 2018. A hierarchical latent vector model for learning long-term structure in music. In *International conference on machine learning*. PMLR, 4364–4373.

[355] Vincent Roger, Jérôme Farinas, and Julien Pinquier. 2022. Deep neural networks for automatic speech processing: a survey from large corpora to limited data. *EURASIP Journal on Audio, Speech, and Music Processing* 2022, 1 (2022), 1–15.

[356] Ben Rogerson. 2020. Sony CSL launches Flow Machines, an AI-assisted music production plugin. *https://www.musicradar.com/news/sony-csl-launches-flow-machines-an-ai-assisted-music-production-plugin* (2020).

[357] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.

[358] Joe Rosato. 2020. Thousands of UC Berkeley Seniors to Graduate in Minecraft Ceremony. *https://www.nbcbayarea.com/news/coronavirus/thousands-of-uc-berkeley-seniors-to-graduate-in-minecraft-ceremony/2291312/* (2020).

[359] Sascha Rothe, Shashi Narayan, and Aliaksei Severyn. 2020. Leveraging pre-trained checkpoints for sequence generation tasks. *Transactions of the Association for Computational Linguistics* 8 (2020), 264–280.

[360] JOSH ROTTENBERG. 2019. How the 'Gemini Man' visual effects team created a young Will Smith. *https://www.latimes.com/entertainment-arts/movies/story/2019-08-27/gemini-man-visual-effects-young-will-smith* (2019).

[361] Simon Rouard and Gaëtan Hadjeres. 2021. CRASH: Raw audio score-based generative modeling for controllable high-resolution drum sound synthesis. *arXiv preprint arXiv:2106.07431* (2021).

[362] Daniel Ruby. 2023. Jasper AI Review 2023: My Experience After Using For 18 Months. *https://www.demandsage.com/jasper-ai-review/* (2023).

[363] Riaan Rudman and Rikus Bruwer. 2016. Defining Web 3.0: opportunities and challenges. *The electronic library* (2016).

[364] Kiersten M Ruff and Rohit V Pappu. 2021. AlphaFold and implications for intrinsically disordered proteins. *Journal of Molecular Biology* 433, 20 (2021), 167208.

[365] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. 1985. *Learning internal representations by error propagation.* Technical Report. California Univ San Diego La Jolla Inst for Cognitive Science.

[366] Warren Sack and Marc Davis. 1994. IDIC: Assembling Video Sequences from Story Plans and Content Annotations.. In *ICMCS*. 30–36.

[367] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022. Palette: Image-to-image diffusion models. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings*. 1–10.

[368] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *arXiv preprint arXiv:2205.11487* (2022).

[369] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. 2022. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).

[370] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. 2017. Temporal generative adversarial nets with singular value clipping. In *Proceedings of the IEEE international conference on computer vision*. 2830–2839.

[371] Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. 2020. Train sparsely, generate densely: Memory-efficient unsupervised training of high-resolution temporal gan. *International Journal of Computer Vision* 128, 10 (2020), 2586–2606.

[372] Rose Salia. 2021. Top 10 Best AI Face Apps Review. *https://topten.ai/face-apps-review/* (2021).

[373] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. 2017. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517* (2017).

[374] Saeed Saremi, Arash Mehrjou, Bernhard Schölkopf, and Aapo Hyvärinen. 2018. Deep energy estimator networks. *arXiv preprint arXiv:1805.08306* (2018).

[375] Jiirgen Schmidhuber. 1990. Making the World Differentiable: On Using Self-Supervised Fully Recurrent N eu al Networks for Dynamic Reinforcement Learning and Planning in Non-Stationary Environm nts. (1990).

[376] Mike Schuster and Kuldip K Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing* 45, 11 (1997), 2673–2681.

[377] P Seeviour, J Holmes, and M Judd. 1976. Automatic generation of control signals for a parallel formant speech synthesizer. In *ICASSP'76. IEEE International Conference on Acoustics, Speech, and Signal Processing*, Vol. 1. IEEE, 690–693.

[378] Iulian V Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath Chandar, Nan Rosemary Ke, et al. 2017. A deep reinforcement learning chatbot. *arXiv preprint arXiv:1709.02349* (2017).

[379] Quin Sexton. 2023. Film Score: A Breakdown in Composing Music for Film. (2023).

[380] Christine H Shadle and Robert I Damper. 2001. Prospects for articulatory synthesis: A position paper. In *4th ISCA Tutorial and Research Workshop (ITRW) on Speech Synthesis*.

[381] Neil Shah, Sarth Engineer, Nandish Bhagat, Hirwa Chauhan, and Manan Shah. 2020. Research trends on the usage of machine learning and artificial intelligence in advertising. *Augmented Human Research* 5 (2020), 1–15.

[382] Yong Shan, Zekang Li, Jinchao Zhang, Fandong Meng, Yang Feng, Cheng Niu, and Jie Zhou. 2020. A contextual hierarchical attention network with adaptive objective for dialogue state tracking. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics.* 6322–6333.

[383] Itamar Shatz. 2017. Native language influence during second language acquisition: A large-scale learner corpus analysis. In *Proceedings of the Pacific Second Language Research Forum (PacSLRF 2016).* 175–188.

[384] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. 2018. Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP).* IEEE, 4779–4783.

[385] Wei Shen and Rujie Liu. 2017. Learning residual images for face attribute manipulation. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 4030–4038.

[386] Yujun Shen and Bolei Zhou. 2021. Closed-form factorization of latent semantics in gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 1532–1540.

[387] Shelly Sheynin, Oron Ashual, Adam Polyak, Uriel Singer, Oran Gafni, Eliya Nachmani, and Yaniv Taigman. 2022. Knn-diffusion: Image generation via large-scale retrieval. *arXiv preprint arXiv:2204.02849* (2022).

[388] Chence Shi, Shitong Luo, Minkai Xu, and Jian Tang. 2021. Learning gradient fields for molecular conformation generation. 9558–9568.

[389] Jiaxin Shi, Shengyang Sun, and Jun Zhu. 2018. A spectral approach to gradient estimation for implicit distributions. In *International Conference on Machine Learning.* PMLR, 4644–4653.

[390] Ziqiang Shi and Shoule Wu. 2022. ITÔN: End-to-end audio generation with Itô stochastic differential equations. *Digital Signal Processing* 132 (2022), 103781.

[391] Wooksu Shin, Namhyuk Ahn, Jeong-Hyeon Moon, and Kyung-Ah Sohn. 2022. Exploiting Distortion Information for Multi-degraded Image Restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 537–546.

[392] Suman Shrestha. 2014. Image denoising using new adaptive based median filters. *arXiv preprint arXiv:1410.2175* (2014).

[393] Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. In *ICLR.*

[394] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. 2022. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792* (2022).

[395] Karan Singla, Zhuohao Chen, David Atkins, and Shrikanth Narayanan. 2020. Towards end-2-end learning for predicting behavior codes from spoken utterances in psychotherapy conversations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics.* 3797–3803.

[396] Bogdan Smolka, K Czubin, Jon Yngve Hardeberg, Kostas N Plataniotis, Marek Szczepanski, and Konrad Wojciechowski. 2003. Towards automatic redeye effect removal. *Pattern Recognition Letters* 24, 11 (2003), 1767–1785.

[397] Yoon So-Yeon. 2020. MBN introduces Korea's first AI news anchor. *https://koreajoongangdaily.joins.com/2020/11/10/entertainment/television/MBN-AI-artificial-intelligence/20201110153900457.html* (2020).

[398] Than Htut Soe. 2021. Automation in Video Editing: Assisted Workflows in Video Editing.. In *AutomationXP@ CHI.*

[399] Fei Song and W Bruce Croft. 1999. A general language model for information retrieval. In *Proceedings of the eighth international conference on Information and knowledge management.* 316–321.

[400] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450* (2019).

[401] Yang Song and Stefano Ermon. 2019. Generative modeling by estimating gradients of the data distribution, Vol. 32.

[402] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. 2020. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence.* PMLR, 574–584.

[403] Yang Song and Diederik P Kingma. 2021. How to train your energy-based models. *arXiv preprint arXiv:2101.03288* (2021).

[404] Yang Song, Jingwen Zhu, Dawei Li, Xiaolong Wang, and Hairong Qi. 2018. Talking face generation by conditional recurrent adversarial network. *arXiv preprint arXiv:1804.04786* (2018).

[405] Nitish Srivastava and Russ R Salakhutdinov. 2012. Multimodal learning with deep boltzmann machines. *Advances in neural information processing systems* 25 (2012).

[406] Nick Statt. 2019. What are Memoji? How to create an Animoji that looks like you. *https://www.theverge.com/2019/7/2/20679241/xiaomi-mimoji-apple-memoji-clone-copying-smartphone-ar-avatar* (2019).

[407] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. 2021. From show to tell: A survey on image captioning. *arXiv preprint arXiv:2107.06912* (2021).

[408] Hui Su, Xiaoyu Shen, Sanqiang Zhao, Xiao Zhou, Pengwei Hu, Randy Zhong, Cheng Niu, and Jie Zhou. 2020. Diversifying dialogue generation with non-conversational text. *arXiv preprint arXiv:2005.04346* (2020).

[409] Shang-Yu Su, Chao-Wei Huang, and Yun-Nung Chen. 2019. Dual supervised learning for natural language understanding and generation. *arXiv preprint arXiv:1905.06196* (2019).

[410] Masanori Suganuma, Xing Liu, and Takayuki Okatani. 2019. Attention-based adaptive selection of operations for image restoration in the presence of unknown combined distortions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 9039–9048.

[411] Jian Sun, Wenfei Cao, Zongben Xu, and Jean Ponce. 2015. Learning a convolutional neural network for non-uniform motion blur removal. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 769–777.

[412] Jiao Sun, Q Vera Liao, Michael Muller, Mayank Agarwal, Stephanie Houde, Kartik Talamadupula, and Justin D Weisz. 2022. Investigating explainability of generative AI for code through scenario-based design. In *27th International Conference on Intelligent User Interfaces*. 212–228.

[413] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems* 27 (2014).

[414] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–13.

[415] Deborah Swanberg, Chiao Fe Shu, and Ramesh Jain. 1993. Architecture of a multimedia information system for content-based retrieval. In *Network and Operating System Support for Digital Audio and Video: Third International Workshop La Jolla, California, USA, November 12–13, 1992 Proceedings 3*. Springer, 387–392.

[416] Laura Sydell. 2009. Building The Curious Faces Of 'Benjamin Button'. *https://www.npr.org/2009/02/17/100668766/building-the-curious-faces-of-benjamin-button* (2009).

[417] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *CVPR*.

[418] Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*. PMLR, 6105–6114.

[419] Xian Tang. 2009. Hybrid Hidden Markov Model and artificial neural network for automatic speech recognition. In *2009 Pacific-Asia Conference on Circuits, Communications and Systems*. IEEE, 682–685.

[420] Victor Tangermann. 2023. 89 PERCENT OF COLLEGE STUDENTS ADMIT TO USING CHATGPT FOR HOMEWORK, STUDY CLAIMS. *https://futurism.com/the-byte/students-admit-chatgpt-homework* (2023).

[421] DeepMind Research Team. 2017. DeepMind's work in 2016: a round-up. *https://www.deepmind.com/blog/deepminds-work-in-2016-a-round-up* (2017).

[422] Latitude Team. 2020. AI Dungeon: Dragon Model Upgrade. *https://aidungeon.medium.com/ai-dungeon-dragon-model-upgrade-7e8ea579abfe* (2020).

[423] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural voice puppetry: Audio-driven facial reenactment. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16*. Springer, 716–731.

[424] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. 2021. A good image generator is what you need for high-resolution video synthesis. *arXiv preprint arXiv:2104.15069* (2021).

[425] Tijmen Tieleman. 2008. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th international conference on Machine learning*. 1064–1071.

[426] Maggie Tillman. 2022. What are Memoji? How to create an Animoji that looks like you. *https://www.pocket-lint.com/phones/news/apple/144743-what-are-memoji-how-to-create-an-animoji-that-looks-like-you/* (2022).

[427] Artem Timoshenko and John R Hauser. 2019. Identifying customer needs from user-generated content. *Marketing Science* 38, 1 (2019), 1–20.

[428] Yoshinobu Tonomura, Akihito Akutsu, Yukinobu Taniguchi, and Gen Suzuki. 1994. Structured video computing. *IEEE multimedia* 1, 03 (1994), 34–43.

[429] László Tóth. 2011. A hierarchical, context-dependent neural network architecture for improved phone recognition. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 5040–5043.

[430] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. 2020. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877* (2020).

[431] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. 2021. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 32–42.

[432] Aggeliki Tsoli, Antonis Argyros, et al. 2019. Patch-based reconstruction of a textureless deformable 3d surface from a single rgb image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*. 0–0.

[433] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. 2018. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1526–1535.

[434] Hirotada Ueda, Takafumi Miyatake, Shigeo Sumino, and Akio Nagasaka. 1993. Automatic structure visualization for video editing. In *Proceedings of the INTERACT'93 and CHI'93 Conference on Human Factors in Computing Systems*. 137–141.

[435] Anwaar Ulhaq, Naveed Akhtar, and Ganna Pogrebna. 2022. Efficient Diffusion Models for Vision: A Survey. *arXiv preprint arXiv:2210.09292* (2022).

[436] Paul Upchurch, Jacob Gardner, Geoff Pleiss, Robert Pless, Noah Snavely, Kavita Bala, and Kilian Weinberger. 2017. Deep feature interpolation for image content changes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7064–7073.

[437] Benigno Uria, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. 2016. Neural autoregressive distribution estimation. *The Journal of Machine Learning Research* 17, 1 (2016), 7184–7220.

[438] Benigno Uria, Iain Murray, and Hugo Larochelle. 2013. RNADE: The real-valued neural autoregressive density-estimator. *Advances in Neural Information Processing Systems* 26 (2013).

[439] Demetrios Vakratsas and Xin Wang. 2020. Artificial intelligence in advertising creativity. *Journal of Advertising* 50, 1 (2020), 39–51.

[440] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. In *The 9th ISCA Speech Synthesis Workshop*.

[441] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. 2016. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems* 29 (2016).

[442] Domonkos Varga and Tamás Szirányi. 2016. Fully automatic image colorization based on Convolutional Neural Network. In *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 3691–3696.

[443] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*.

[444] Yuri Viazovetskyi, Vladimir Ivashkin, and Evgeny Kashin. 2020. Stylegan2 distillation for feed-forward image manipulation. In *European conference on computer vision*. Springer, 170–186.

[445] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. 2022. DiGress: Discrete Denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734* (2022).

[446] Pascal Vincent. 2011. A connection between score matching and denoising autoencoders. *Neural computation* 23, 7 (2011), 1661–1674.

[447] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3156–3164.

[448] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. 2016. Generating videos with scene dynamics. *Advances in neural information processing systems* 29 (2016).

[449] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. 2022. Sketch-Guided Text-to-Image Diffusion Models. *arXiv preprint arXiv:2211.13752* (2022).

[450] Bram Wallace, Akash Gokul, and Nikhil Naik. 2022. EDICT: Exact Diffusion Inversion via Coupled Transformations. *arXiv preprint arXiv:2211.12446* (2022).

[451] Hanna M Wallach. 2006. Topic modeling: beyond bag-of-words. In *Proceedings of the 23rd international conference on Machine learning*. 977–984.

[452] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Fang Wen, and Jing Liao. 2022. Old photo restoration via deep latent space translation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).

[453] Chengyi Wang, Yu Wu, Yujiao Du, Jinyu Li, Shujie Liu, Liang Lu, Shuo Ren, Guoli Ye, Sheng Zhao, and Ming Zhou. 2019. Semantic mask for transformer based end-to-end speech recognition. *arXiv preprint arXiv:1912.03010* (2019).

[454] Dan Wang, Xinrui Cui, Xun Chen, Zhengxia Zou, Tianyang Shi, Septimiu Salcudean, Z Jane Wang, and Rabab Ward. 2021. Multi-view 3D Reconstruction with Transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5722–5731.

[455] Jun Wang, Ying Cui, Dongyan Guo, Junxia Li, Qingshan Liu, and Chunhua Shen. 2022. PointAttN: You Only Need Attention for Point Cloud Completion. *arXiv preprint arXiv:2203.08485* (2022).

[456] Li Wang, Zechen Bai, Yonghua Zhang, and Hongtao Lu. 2020. Show, recall, and tell: Image captioning with recall mechanism. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 12176–12183.

[457] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. 2018. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European conference on computer vision (ECCV)*. 52–67.

[458] Rui Wang, Xu Tan, Renqian Luo, Tao Qin, and Tie-Yan Liu. 2021. A survey on low-resource neural machine translation. *arXiv preprint arXiv:2107.04239* (2021).

[459] Xuejiao Wang, Qiuyan Tao, Lianghao Wang, Dongxiao Li, and Ming Zhang. 2015. Deep convolutional architecture for natural image denoising. In *2015 International Conference on Wireless Communications & Signal Processing (WCSP)*. IEEE, 1–4.

[460] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. 2018. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*. 0–0.

[461] Yuting Wang. 2021. The Application of Artificial Intelligence in Chinese News Media. In *2021 2nd International Conference on Artificial Intelligence and Information Systems*. 1–4.

[462] Yexiang Wang, Yi Guo, and Siqi Zhu. 2020. Slot attention with value normalization for multi-domain dialogue state tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 3019–3028.

[463] Yufei Wang, Zhe Lin, Xiaohui Shen, Scott Cohen, and Garrison W Cottrell. 2017. Skeleton key: Image captioning by skeleton-attribute decomposition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7272–7281.

[464] Duncan J Watts and Steven H Strogatz. 1998. Collective dynamics of 'small-world' networks. *nature* 393, 6684 (1998), 440–442.

[465] Andrew Webster. 2020. Travis Scott's first Fortnite concert was surreal and spectacular. *https://www.theverge.com/2020/4/23/21233637/travis-scott-fortnite-concert-astronomical-live-report* (2020).

[466] Li-Yi Wei and Marc Levoy. 2000. Fast texture synthesis using tree-structured vector quantization. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*. 479–488.

[467] Tianyi Wei, Dongdong Chen, Wenbo Zhou, Jing Liao, Weiming Zhang, Lu Yuan, Gang Hua, and Nenghai Yu. 2022. E2Style: Improve the efficiency and effectiveness of StyleGAN inversion. *IEEE Transactions on Image Processing* 31 (2022), 3267–3280.

[468] Max Welling and Geoffrey E Hinton. 2002. A new learning algorithm for mean field Boltzmann machines. In *International Conference on Artificial Neural Networks*. Springer, 351–357.

[469] WIKIPEDIA. 2013. Furious 7 (2015). *https://en.wikipedia.org/wiki/Furious_7* (2013).

[470] Max Willens. 2018. Bloomberg Media has a robot writing story summaries. *https://digiday.com/media/bloomberg-media-robot-writing-story-summaries/* (2018).

[471] Karen William. 2021. Top 6 to Voicemod Alternatives Voice Changer [Windows/Mac/Android/iOS]. *https://filme.imyfone.com/audio-edit/voicemod-alternative/* (2021).

[472] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. 2021. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806* (2021).

[473] Jibin Wu, Emre Yılmaz, Malu Zhang, Haizhou Li, and Kay Chen Tan. 2020. Deep spiking neural networks for large vocabulary automatic speech recognition. *Frontiers in neuroscience* 14 (2020), 199.

[474] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. 2022. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. *arXiv preprint arXiv:2212.11565* (2022).

[475] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).

[476] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 2015. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1912–1920.

[477] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. 2022. Gan inversion: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).

[478] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. 2017. Dna-gan: Learning disentangled representations from multi-attribute images. *arXiv preprint arXiv:1711.05415* (2017).

[479] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. 2021. Tackling the generative learning trilemma with denoising diffusion gans. *arXiv preprint arXiv:2112.07804* (2021).

[480] Haozhe Xie, Hongxun Yao, Xiaoshuai Sun, Shangchen Zhou, and Shengping Zhang. 2019. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proceedings of the IEEE/CVF international conference on computer vision*. 2690–2698.

[481] Junyuan Xie, Linli Xu, and Enhong Chen. 2012. Image denoising and inpainting with deep neural networks. *Advances in neural information processing systems* 25 (2012).

[482] Tian Xie, Xiang Fu, Octavian-Eugen Ganea, Regina Barzilay, and Tommi Jaakkola. 2021. Crystal diffusion variational autoencoder for periodic material generation. *arXiv preprint arXiv:2110.06197* (2021).

[483] Jun Xu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, Wanxiang Che, and Ting Liu. 2020. Conversational graph grounded policy learning for open-domain conversation generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 1835–1845.

[484] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. PMLR, 2048–2057.

[485] Li Xu and Jiaya Jia. 2010. Two-phase kernel estimation for robust motion deblurring. In *European conference on computer vision*. Springer, 157–170.

[486] Li Xu, Jimmy S Ren, Ce Liu, and Jiaya Jia. 2014. Deep convolutional neural network for image deconvolution. *Advances in neural information processing systems* 27 (2014).

[487] Peng Xu, Xiatian Zhu, and David A Clifton. 2022. Multimodal learning with transformers: a survey. *arXiv preprint arXiv:2206.06488* (2022).

[488] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1316–1324.

[489] Xianqiu Xu, Hongqing Liu, Yong Li, and Yi Zhou. 2016. Image deblurring with blur kernel estimation in RGB channels. In *2016 IEEE International Conference on Digital Signal Processing (DSP)*. IEEE, 681–684.

[490] Yangyang Xu, Yong Du, Wenpeng Xiao, Xuemiao Xu, and Shengfeng He. 2021. From continuity to editability: Inverting gans with consecutive images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13910–13918.

[491] Ke Xue, Yifei Li, and Hanqing Jin. 2022. What Do You Think of AI? Research on the Influence of AI News Anchor Image on Watching Intention. *Behavioral Sciences* 12, 11 (2022), 465.

[492] Hemant Yadav and Sunayana Sitaram. 2022. A Survey of Multilingual Models for Automatic Speech Recognition. *arXiv preprint arXiv:2202.12576* (2022).

[493] Karan Yadav. 2023. How education chatbots can help students and teachers. *http://www.eyeshenzhen.com/content/2022-02/08/content_24921512.htm* (2023).

[494] Ravindra Yadav, Ashish Sardana, Vinay P Namboodiri, and Rajesh M Hegde. 2021. Speech prediction in silent videos using variational autoencoders. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 7048–7052.

[495] Xinchen Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee. 2016. Attribute2image: Conditional image generation from visual attributes. In *European conference on computer vision*. Springer, 776–791.

[496] Li-Chia Yang, Szu-Yu Chou, and Yi-Hsuan Yang. 2017. MidiNet: A convolutional generative adversarial network for symbolic-domain music generation. *arXiv preprint arXiv:1703.10847* (2017).

[497] Shuoheng Yang, Yuxin Wang, and Xiaowen Chu. 2020. A survey of deep learning techniques for neural machine translation. *arXiv preprint arXiv:2002.07526* (2020).

[498] Shiquan Yang, Rui Zhang, and Sarah Erfani. 2020. Graphdialog: Integrating graph knowledge into end-to-end task-oriented dialogue systems. *arXiv preprint arXiv:2010.01447* (2020).

[499] Tiancheng Yang and Shah Nazir. 2022. A comprehensive overview of AI-enabled music classification and its influence in games. *Soft Computing* 26, 16 (2022), 7679–7693.

[500] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. 2019. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10685–10694.

[501] Xu Yang, Hanwang Zhang, and Jianfei Cai. 2019. Learning to collocate neural modules for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4250–4260.

[502] Yunyi Yang, Yunhao Li, and Xiaojun Quan. 2021. Ubar: Towards fully end-to-end task-oriented dialog system with gpt-2. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 14230–14238.

[503] Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. 2018. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *International conference on machine learning*. PMLR, 5708–5717.

[504] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. 2022. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789* (2022).

[505] Ke Yu, Chao Dong, Liang Lin, and Chen Change Loy. 2018. Crafting a toolchain for image restoration by deep reinforcement learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2443–2452.

[506] Tianxiu Yu, Cong Lin, Shijie Zhang, Chunxue Wang, Xiaohong Ding, Huili An, Xiaoxiang Liu, Ting Qu, Liang Wan, Shaodi You, et al. 2022. Artificial Intelligence for Dunhuang Cultural Heritage Protection: The Project and the Dataset. *International Journal of Computer Vision* 130, 11 (2022), 2646–2673.

[507] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. 2021. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432* (2021).

[508] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. 2021. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986* (2021).

[509] Yuan Yuan, Siyuan Liu, Jiawei Zhang, Yongbing Zhang, Chao Dong, and Liang Lin. 2018. Unsupervised image super-resolution using cycle-in-cycle generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 701–710.

[510] Yi Yuan, Jilin Tang, and Zhengxia Zou. 2021. Vanet: a view attention guided network for 3d reconstruction from single and multi-view images. In *2021 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 1–6.

[511] Xia Yuanjie. 2022. AI sign language anchor serves at Olympics. *http://www.eyeshenzhen.com/content/2022-02/08/content_24921512.htm* (2022).

[512] Vladyslav Yushchenko, Nikita Araslanov, and Stefan Roth. 2019. Markov decision process for video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*. 0–0.

[513] Shehtab Zaman, Ethan Ferguson, Cecile Pereira, Denis Akhiyarov, Mauricio Araya-Polo, and Kenneth Chiu. 2022. ParticleGrid: Enabling Deep Learning using 3D Representation of Materials. *arXiv preprint arXiv:2211.08506* (2022).

[514] Heiga Ze, Andrew Senior, and Mike Schuster. 2013. Statistical parametric speech synthesis using deep neural networks. In *2013 ieee international conference on acoustics, speech and signal processing*. IEEE, 7962–7966.

[515] Heiga Zen. 2015. Acoustic modeling in statistical parametric speech synthesis-from HMM to LSTM-RNN. (2015).

[516] Heiga Zen and Haşim Sak. 2015. Unidirectional long short-term memory recurrent neural network with recurrent output layer for low-latency speech synthesis. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 4470–4474.

[517] Aeron Zentner. 2022. Applied Innovation: Artificial Intelligence in Higher Education. *Available at SSRN 4314180* (2022).

[518] Zheng-Jun Zha, Daqing Liu, Hanwang Zhang, Yongdong Zhang, and Feng Wu. 2019. Context-aware visual policy network for fine-grained image captioning. *IEEE transactions on pattern analysis and machine intelligence* 44, 2 (2019), 710–722.

[519] Chaoning Zhang, Chenshuang Zhang, Junha Song, John Seon Keun Yi, Kang Zhang, and In So Kweon. 2022. A survey on masked autoencoder for self-supervised learning in vision and beyond. *arXiv preprint arXiv:2208.00173* (2022).

[520] Chaoning Zhang, Kang Zhang, Trung X. Pham, Changdong Yoo, and In-So Kweon. 2022. Dual Temperature Helps Contrastive Learning Without Many Negative Samples: Towards Understanding and Simplifying MoCo. In *CVPR*.

[521] Chaoning Zhang, Kang Zhang, Chenshuang Zhang, Trung X Pham, Chang D Yoo, and In So Kweon. 2022. How Does SimSiam Avoid Collapse Without Negative Samples? A Unified Understanding with Self-supervised Contrastive Learning. In *ICLR*.

[522] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2017. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 5907–5915.

[523] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2018. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence* 41, 8 (2018), 1947–1962.

[524] Kunai Zhang, Da Huang, and David Zhang. 2017. An optimized palmprint recognition approach based on image sharpness. *Pattern Recognition Letters* 85 (2017), 65–71.

[525] Mingju Zhang, Lei Zhang, Yanfeng Sun, Lin Feng, and Weiying Ma. 2005. Auto cropping for digital photographs. In *2005 IEEE International Conference on Multimedia and Expo*. IEEE, 4–pp.

[526] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5579–5588.

[527] Weixing Zhang. 2022. Application and development of robot sports news writing by artificial intelligence. In *2022 IEEE 2nd International Conference on Data Science and Computer Application (ICDSCA)*. IEEE, 869–872.

[528] Wenlong Zhang, Yihao Liu, Chao Dong, and Yu Qiao. 2019. Ranksrgan: Generative adversarial networks with ranker for image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3096–3105.

[529] Xiaobo Zhang, Xiangchu Feng, Weiwei Wang, Shunli Zhang, and Qunfeng Dong. 2013. Gradient-based Wiener filter for image denoising. *Computers & Electrical Engineering* 39, 3 (2013), 934–944.

[530] Xuying Zhang, Xiaoshuai Sun, Yunpeng Luo, Jiayi Ji, Yiyi Zhou, Yongjian Wu, Feiyue Huang, and Rongrong Ji. 2021. Rstnet: Captioning with adaptive attention on visual and non-visual words. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 15465–15474.

[531] Yichi Zhang, Zhijian Ou, Huixin Wang, and Junlan Feng. 2020. A probabilistic end-to-end task-oriented dialog model with latent belief states towards semi-supervised learning. *arXiv preprint arXiv:2009.08115* (2020).

[532] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2019. Dialogpt: Large-scale generative pre-training for conversational response generation. *arXiv preprint arXiv:1911.00536* (2019).

[533] Zheng Zhang, Ryuichi Takanobu, Qi Zhu, MinLie Huang, and XiaoYan Zhu. 2020. Recent advances and challenges in task-oriented dialog systems. *Science China Technological Sciences* 63, 10 (2020), 2011–2027.

[534] Mo Zhao, Gang Cao, Xianglin Huang, and Lifang Yang. 2022. Hybrid Transformer-CNN for Real Image Denoising. *IEEE Signal Processing Letters* (2022).

[535] Qian Zhao, Hao Yang, Dongming Zhou, and Jinde Cao. 2022. Rethinking image deblurring via CNN-Transformer multi-scale hybrid architecture. *IEEE Transactions on Instrumentation and Measurement* (2022).

[536] Wentian Zhao, Yao Hu, Heda Wang, Xinxiao Wu, and Jiebo Luo. 2021. Boosting Entity-aware Image Captioning with Multi-modal Knowledge Graph. *arXiv preprint arXiv:2107.11970* (2021).

[537] Rui Zhen, Wenchao Song, Qiang He, Juan Cao, Lei Shi, and Jia Luo. 2023. Human-Computer Interaction System: A Survey of Talking-Head Generation. *Electronics* 12, 1 (2023), 218.

[538] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. 2023. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419* (2023).

[539] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. 2019. Talking face generation by adversarially disentangled audio-visual representation. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 9299–9306.

[540] Jingyuan Zhou, Chaktou Leong, Minyi Lin, Wantong Liao, and Congduan Li. 2022. Task adaptive network for image restoration with combined degradation factors. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1–8.

[541] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. 2020. The design and implementation of xiaoice, an empathetic social chatbot. *Computational Linguistics* 46, 1 (2020), 53–93.

[542] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 13041–13049.

[543] Le Zhou, Qiu-Feng Wang, Kaizhu Huang, and Cheng-Hung Lo. 2019. An interactive and generative approach for chinese shanshui painting document. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 819–824.

[544] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. 2020. In-domain gan inversion for real image editing. In *European conference on computer vision*. Springer, 592–608.

[545] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. 2016. Generative visual manipulation on the natural image manifold. In *European conference on computer vision*. Springer, 597–613.

[546] Qingfu Zhu, Lei Cui, Weinan Zhang, Furu Wei, and Ting Liu. 2018. Retrieval-enhanced adversarial training for neural response generation. *arXiv preprint arXiv:1809.04276* (2018).

[547] Barret Zoph and Kevin Knight. 2016. Multi-source neural translation. *arXiv preprint arXiv:1601.00710* (2016).

[548] Zhengxia Zou, Tianyang Shi, Shuang Qiu, Yi Yuan, and Zhenwei Shi. 2021. Stylized neural painting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15689–15698.