# Transformation vs Tradition: Artificial General Intelligence (AGI) for Arts and Humanities

Zhengliang Liu[*1], Yiwei Li[*1], Qian Cao[*2], Junwen Chen[*3], Tianze Yang[1], Zihao Wu[1], John Gibbs[4], Khaled Rasheed[1], Ninghao Liu[†1], Gengchen Mai[†2], and Tianming Liu[†1]

[1]School of Computing, University of Georgia
[2]Department of Geography, University of Georgia
[3]Department of Computer Science and Engineering, Michigan State University
[4]Department of Theatre & Film Studies, University of Georgia

## Abstract

Recent advances in artificial general intelligence (AGI), particularly large language models and creative image generation systems have demonstrated impressive capabilities on diverse tasks spanning the arts and humanities. However, the swift evolution of AGI has also raised critical questions about its responsible deployment in these culturally significant domains traditionally seen as profoundly human. This paper provides a comprehensive analysis of the applications and implications of AGI for text, graphics, audio, and video pertaining to arts and the humanities. We survey cutting-edge systems and their usage in areas ranging from poetry to history, marketing to film, and communication to classical art. We outline substantial concerns pertaining to factuality, toxicity, biases, and public safety in AGI systems, and propose mitigation strategies. The paper argues for multi-stakeholder collaboration to ensure AGI promotes creativity, knowledge, and cultural values without undermining truth or human dignity. Our timely contribution summarizes a rapidly developing field, highlighting promising directions while advocating for responsible progress centering on human flourishing. The analysis lays the groundwork for further research on aligning AGI's technological capacities with enduring social goods.

## 1 Introduction

Arts and the humanities have long been reflections of human experience, emotions, and philosophical introspection [1]. These domains, deeply rooted in subjectivity, creativity, and a nuanced appreciation of the world, have served as repositories of our history, culture, and identity. Over the past few years, however, the boundary between human creativity and machine computation has started to blur, ushering in an era where Artificial Intelligence (AI) influences artistic creation and reshapes our understanding of humanities.

Historically, AI's foray into domains requiring creativity was met with skepticism [2]. Critics posited that machines, bound by algorithms and devoid of emotions, could never truly comprehend or replicate the intricacies of artistic expression. Creativity was, after all, seen as the antithesis of computation, fueled by irregularities, out-of-box thinking, and a delicate understanding of the human

---

[*]These authors contributed equally to this work

[†]Corresponding author(s). E-mail(s): ninghao.liu@uga.edu, gengchen.mai25@uga.edu, tliu@uga.edu
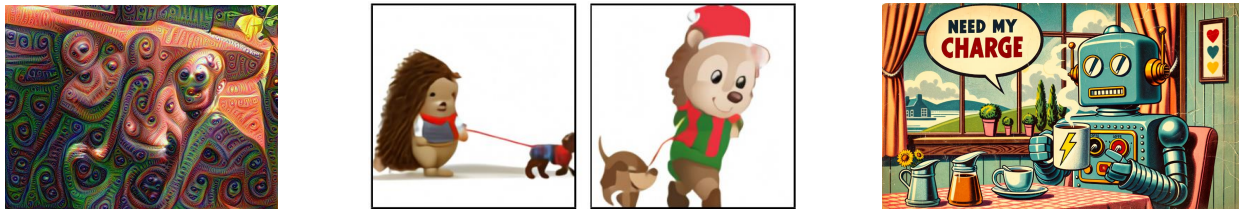
Figure 1: Some examples of AGI-generated images. **Left**: A heavily deep-dream-style photograph expressing "three men in a pool", which is difficult for humans to understand. **Middle**: An image generated by DALL-E through translation from "an illustration of a baby hedgehog in a christmas sweater walking a dog" [3]. **Right**: Image created by DALL-E 3 with the prompt "vintage 1940s cartoon featuring a robot holding a steaming coffee mug with a lightning bolt symbol on it, text bubble that reads 'Need my charge', sitting at a table by bay window in a coffee shop interior". The model can generate the high-quality image, and correctly understand the instruction.

condition. These very attributes, which are the cornerstones of arts and humanities, seemed out of reach for artificial entities.

More recently, the landscape has begun to shift. Early algorithms such as "Deep Dream" (2015) [4] and various approaches in the theme of "Neural Style Transfer" [5] marked AI's early attempts at artistic endeavors. However, Deep Dream was plagued by the problem of generating repetitive canine facial features within images, and the style transfer process, while artistically intriguing, lacked the ability to create entirely new content or comprehend the underlying semantics of images. Nevertheless, these earlier attempts began a noticeable shift in perceptions, with increasing acceptance of artificial intelligence's contribution to artistic endeavors. A pivotal moment highlighting this change was when "Edmond de Belamy, from La Famille de Belamy", a portrait produced by Generative Adversarial Networks (GANs) [6], was sold by Christie's New York on Oct 19, 2018 for $432,500 [7], which is more than 40 times Christie's initial estimate. Despite facing skepticism and questions regarding its originality from other artists who work with AI, these rudimentary techniques marked the nascent stages of AI-assisted artistry.

The leap forward came in 2021 with the arrival of text-to-image algorithms. Specifically, the introduction of DALL-E [3], supplemented by the unveiling of open-source projects like VQGAN+CLIP [8, 9], catalyzed the proliferation of AI art generators. Furthermore, in 2022, the release of "Stable Diffusion" [10] by Stability AI and "Imagen" [11] by Google AI ushered in a new era of advanced AI-powered creativity. This release further democratized the Artificial Intelligence-Generated Content (AIGC) process. The field of AIGC is still extremely young. Major contributors and platforms have a relatively short operational history, spanning less than a year. However, the trajectory suggests an impending turning point where AI capabilities will become sophisticated enough to revolutionize various art-related domains. For instance, in the realm of video game development, concept and traditional artists are already harnessing AI image generation for inspiration and as tangible assets in their creative works[†]. Looking ahead, once the complexities of image generation are comprehensively addressed, it is plausible that the intellectual capital steering this innovation will gravitate toward other modalities. This may encompass domains like auditory processing and generation, video synthesis, and literary generation, among other multidisciplinary challenges.

With the recent advancement of Large Language Models (LLMs), the rise of Artificial General

---

[†] https://www.scenario.com/

Intelligence (AGI) further challenges traditional perspectives. AGI [12], with its potential to emulate holistic human cognition, promises not just to create art but to understand and appreciate it–and in fact many proponents of LLMs having early AGI capabilities espouse that these models already have a degree of understanding of the physical world and humans [see references. They need to be added to citations]. LLMs' integration into arts and humanities could revolutionize everything from literary synthesis, capturing the depth of human emotion, to creating multi-sensory art experiences and reinterpreting historical narratives.

This paper delves deep into the rapidly evolving nexus of AGI, arts, and the humanities. While celebrating the transformative potential of AGI, it also critically examines the following underlying questions: Can AGI truly be creative? Will it ever appreciate art the way humans do? And most importantly, as AGI blurs the lines between machine capability and human creativity, what does it mean for the future of arts and humanities? Through this discourse, we seek to navigate the promising yet perplexing frontier of AGI-infused artistry.

## 2 Background

### 2.1 Generative AI: From GAN to ChatGPT

A recent survey paper [13] provides a comprehensive review of the field of AI-generated content (AIGC). AIGC refers to content like text, images, music, and code that is generated by AI systems rather than created directly by humans.

The authors review the history of generative AI models, beginning with early statistical models like Hidden Markov Models and Gaussian Mixture Models. They then discuss the rise of deep learning models like GANs, VAEs, and diffusion models, with the transformer architecture (2017) identified as a key breakthrough that enabled large-scale pre-trained models like GPT-3 [14], ChatGPT [15], and GPT-4 [16].

The paper categorizes generative models as either unimodal, which generate content in a single modality like text or images, or multimodal, which combine multiple modalities. For unimodal models, they provide an in-depth review of state-of-the-art generative language models like GPT-3 [14], BART [17], T5 [18] and vision models such as Stable Diffusion [10] and DALL-E 2 [19].

For the multimodal generation, the survey examines vision-language models like DALL-E and GLIDE [20] as well as text-audio, text-graph, and text-code models. These allow cross-modal generation between modalities. The authors discuss applications like chatbots, art creation, music composition, and code generation.

They also cover techniques that help align model outputs with human preferences, such as reinforcement learning from human feedback as used in ChatGPT. The paper analyzes challenges around efficiency, trustworthiness, and responsible use of large AIGC models. Finally, open problems and future research directions are explored.

### 2.2 Opportunities and Challenges of General AIGC

The enthusiastic reception of conversational agents like ChatGPT underscores AIGC's vast potential. However, researchers must grapple with critical challenges around data bias, computational efficiency, output quality, and ethical implications as AIGC rapidly gains traction [21].

On the opportunities front, AIGC holds promise for boosting productivity in creative fields by acting

as an intelligent assistant that can synthesize draft content. Cross-modal generation techniques can potentially bridge content formats, enabling applications like generating videos from text descriptions. In industry verticals like e-commerce, AIGC can scale the creation of catalog descriptions and customized landing pages. For news and entertainment, AIGC may enhance automation in production pipelines. The multi-task learning abilities of foundation models could spur innovation if applied judiciously. On the consumer side, AIGC can deliver more personalized, interactive, and immersive experiences.

However, substantial challenges remain. Massive computational resources are needed to develop and deploy the latest AIGC models [22], which may concentrate power in fewer hands. More crucially, the data used to train AIGC models inherits human biases [23] that are reflected in outputs. Curating high-quality datasets is an arduous task. While human-in-the-loop approaches may improve model alignments, transparency and accountability are still lacking. Safeguards against toxic outputs remain inadequate as interactions uncover harmful edge cases. For high-stakes domains like healthcare, the risks of errors loom large. Despite great enthusiasm, researchers should adopt a measured approach while addressing these concerns through technical and ethical diligence.

While AIGC represents an exciting frontier for AI research with immense potential upside, responsible development calls for holistic solutions encompassing data curation and hygiene, efficient systems, user feedback loops, and transparency. With care and consideration for societal impacts, AIGC could usher in an era where generative AI assists and augments human creativity rather than displacing it. This survey provides a timely overview of the state-of-the-art as of late 2023, and a roadmap to guide progress in this rapidly evolving domain.

# 3   Text Analysis and Generation

Text analysis and generation are crucial domains in natural language processing influencing myriad applications. At the core, text analysis delves into comprehending intricate patterns, meanings, and sentiments in textual data, whereas text generation aspires to craft human-like text based on certain criteria or prompts [15]. With the advent of sophisticated model architectures, the boundaries of what machines can comprehend and produce have been ceaselessly expanded. This section delineates the technical advancements underpinning these capabilities, including seminal models like Transformers, and extends into their pragmatic applications across diverse sectors such as poetry, music, law, advertising, and governance.

## 3.1   Technical Advances

The Transformer [24] architecture has undoubtedly carved a pivotal role in the progression of natural language processing models. Introduced by Vaswani et al. [24], the architecture abandoned recurrent layers, traditionally used for sequence data, in favor of attention mechanisms. The core concepts emanating from this architecture, including encoders, decoders, BERT, and autoregressive language models, have since dominated state-of-the-art results in various NLP tasks.

### 3.1.1   Transformer Architectures

At the heart of the Transformer model lies the pivotal *self-attention mechanism*, which computes a weighted sum of all words in a sequence relative to each other. This empowers the model to capture the intricate relationships between words, regardless of their positions in the sequence. Unlike recurrent models such as RNNs [25] or LSTMs [26] which process sequences iteratively, Transformers

handle the entire sequence in parallel. This approach, coupled with additional design elements like positional encodings [27] and residual connections [28], empowers Transformers to deliver both efficiency and effectiveness, even when confronted with lengthy sequences.

### 3.1.2 BERT (Bidirectional Encoder Representations from Transformers)

Emerging from the Transformer paradigm, BERT [29] represents a monumental shift in pre-training methods. Introduced by Google in 2018, this model captures bidirectional contexts by considering both preceding and following words in all its layers. BERT's pre-training phase involves a *masked language model* objective wherein it attempts to predict randomly masked words in a sentence. Once pre-trained on vast corpora such as Wikipedia, BERT can be adeptly fine-tuned on specific tasks using small labeled datasets, by just adding appropriate task-specific layers [30]. This approach has made BERT highly versatile, allowing it to be applied to a wide range of NLP tasks.

### 3.1.3 Autoregressive Language Models

*Autoregressive modeling* [14] uses a step-by-step approach, where predicting the next item in a sequence depends on what came before it. In the context of language modeling, when given part of a sentence, these models try to guess what words come next. Once properly trained, the models good at guessing what word comes after the previous ones in the sequence. When autoregressive models generate text, they can use different methods, like beam search [31], greedy decoding [32], or probabilistic sampling [33]. A well-known example of an autoregressive language model is OpenAI's GPT (Generative Pre-trained Transformer) [34], which, in contrast to BERT's bidirectionality, is unidirectional and is primed mainly for text generation.

## 3.2 Real-world Applications

This section elucidates the diverse arts and humanities landscape of AGI by categorizing its applications into four distinct but interrelated subsections: 1) Literature Search and Analytics; 2) Linguistics and Communication; 3) Creative Endeavors. Each of these subsections represents a unique facet of AI's ever-expanding repertoire, showcasing its adaptability to perform tasks ranging from the systematic retrieval of knowledge to nuanced linguistic interactions, from meticulous analytics to imaginative, creative endeavors.

### 3.2.1 Literature Search and Analytics

Large language models, once equipped with common sense knowledge, can provide valuable assistance in various liberal arts domains such as history, classics, and philosophy that heavily rely on literature search and analysis. Based on [35], LLMs can be helpful in the following aspects.

- **Automated Literature Review:** LLMs can quickly scan and summarize large volumes of text. They can identify key concepts, themes, and relevant passages, saving researchers significant time and effort.

- **Cross-Referencing:** LLMs can cross-reference texts, identifying connections and references between historical events, philosophical works, and classic literature, helping researchers explore intertextual relationships.

- **Summarization:** LLMs can generate concise summaries of lengthy texts, making complex philosophical or historical writings more accessible to a broader audience.
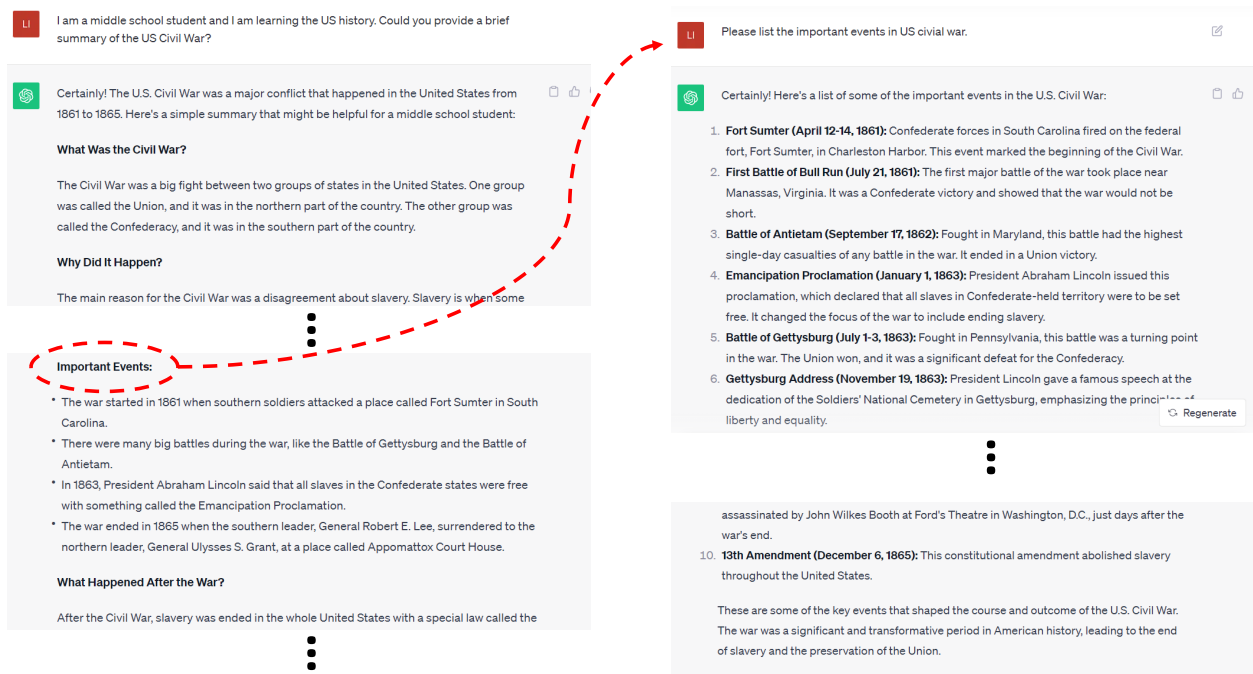
Figure 2: An example of using GPT-3.5 for learning history. The right part shows a follow-up question regarding the answer of the first question in the left part.

- **Question Answering:** LLMs excel in answering specific questions related to historical events, philosophical theories, or classic literature. They can provide concise and accurate responses by drawing on their vast knowledge base.

- **Content Generation:** LLMs can assist in generating preliminary content by retrieving relevant information in the literature [36]. They can provide background information, context, and even propose arguments based on the input provided.

- **Teaching and Learning:** LLMs can be used as educational tools to provide explanations, generate practice questions, and engage students in discussions related to historical, classic, or philosophical topics.

Some specific examples of the ongoing and potential applications are provided as below.

**Anthropology:** LLMs can process large amounts of anthropological data. Through text mining, LLMs deeply research documents, interviews, and historical texts to find information on specific cultures, social groups, or topics, and promote a deeper understanding of human social evolution and cultural differences [15]. Second, LLMs can help analyze social surveys and public opinion polls, to gain an in-depth understanding of attitudes, beliefs, and behaviors in human society, helping researchers understand social trends and changes in public opinion [37]. Finally, through cross-cultural research, LLMs support the comparison of similarities and differences between different cultures, provide translation services, analyze cross-cultural communication, and conduct in-depth research on global issues such as globalization and cultural exchanges.

**Classics:** LLMs contribute to art historical research, providing in-depth insights into a specific period or style by analyzing textual descriptions of classical artworks, historical documents, and the lives of relevant artists. In terms of art education and popularization, LLMs can assist in the
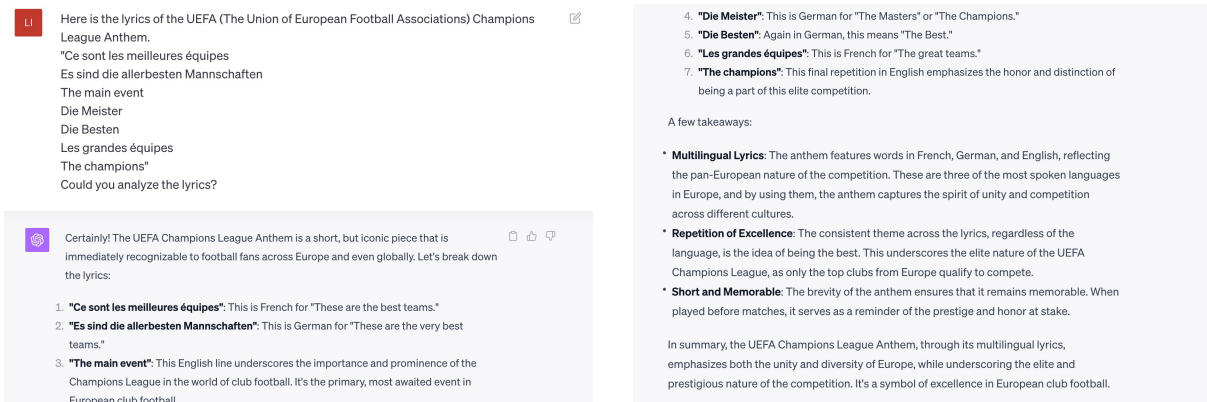
Figure 3: An example of using GPT-4 to analyze the background and design philosophy of the lyrics of the UEFA (The Union of European Football Associations) Champions League Anthem. The AI model can easily handle the multilingual content, and even point out the "spirit of unity" and "diversity" behind the design.

creation of art education materials, explain works of art so that more people can understand and appreciate classical art, and generate explanatory texts on art history for use in education, museum exhibitions, and cultural dissemination [38].

**Philosophy:** LLMs can be used for literature reviews to help researchers understand the current state of research on specific philosophical issues or thinkers [39]. They can also analyze philosophical texts and understand the author's ideas, argument structure, and logic. In addition, LLMs can also be used to analyze the structure and effectiveness of philosophical arguments, helping researchers better understand and evaluate philosophical papers.

**Psychology:** LLMs can conduct literature reviews to help researchers understand the latest research and theories on specific psychological topics [40]. LLMs can also generate questionnaire questions to ensure they are clear and effective. In addition, LLMs can analyze comments on social media to understand emotional health and mental health issues of people, and provide support and resources.

**History:** LLMs can analyze historical texts to help understand events, generate summaries, extract key information, and improve information processing efficiency [41]. LLMs can also calibrate time in historical text, track events, and help establish a detailed historical timeline. They can also help extract character relationships, help build a relationship map, and conduct in-depth research on the influence of historical figures.

### 3.2.2  Linguistics and Communication

LLMs can be highly beneficial in applications related to linguistics and communications due to their natural language processing capabilities and extensive knowledge base. An example that shows some these abilities is in Figure 3.

- **Language Understanding:** LLMs can be used to analyze the structure, grammar, and semantics of languages, aiding linguists in their research on syntax, morphology, and linguistic phenomena.

- **Translation Assistance:** LLMs can assist linguists and translators in translating text between different languages, helping bridge linguistic and cultural gaps.

- **Sentiment Analysis:** LLMs can perform sentiment analysis on text, enabling businesses and organizations to gauge public sentiment towards their products, services, or policies.

- **Speech Recognition:** LLMs can enhance speech recognition systems, improving the accuracy of voice-to-text transcriptions.

Based on the above capabilities, some specific examples of the ongoing and potential applications of LLMs are provided as below.

**Linguistics:** LLMs can generate new language texts and expand understanding of grammatical structures and vocabulary usage [42, 43]. Scientists can conduct semantic analysis through LLMs and conduct in-depth studies of lexical meanings, contextual relevance [44], and semantic relationships of language expressions [39]. These models can also be used to develop language learning tools to help students learn vocabulary, grammar rules, and other language knowledge. In studying language disorders, LLMs can reveal the manifestations and effects of language disorders in different contexts. However, some scholars have raised objections, believing that LLMs lack human cognition [45].

**Language Studies:** LLMs can analyze the grammatical, semantic, and pragmatic features of different languages. Based on this, LLMs can generate teaching materials and exercises with explanations to provide students with strong support in learning grammar, vocabulary, and expressions. In addition, LLMs excel in translation, supporting cross-language communication and translation [29]. At the same time, LLMs play an important role in writing and creation and can create articles, compose essays, and generate various literary works. In addition, LLMs support speech recognition technology, which allows speech input to be easily converted into text, facilitating speech interaction and speech recognition applications [38]. By processing large amounts of historical texts, LLMs can also assist researchers in tracing the evolution, change, and development of the language.

**Communitcation Studies:** LLMs can analyze large amounts of news and advertising, to reveal patterns, trends, and factors that influence the spread of information. LLMs can also analyze emotions and interactions on social media [46], studying how information spreads in social networks and its impact on public opinion [47]. Researchers use LLMs to analyze the language and framing of news reports and study the way news media report events and their impact on audience perceptions. In terms of multilingual content, LLMs have the advantages of translation and understanding and are helpful in studying language differences and cultural factors in cross-cultural communication.

### 3.2.3 Creative Endeavors

This subsection starts with examples of several applications where AI might models generate "creative" contents, followed by a further discussion of whether AI can genuinely attain a level of creativity comparable to that of humans.

**Song Lyrics:** LLMs such as the GPT family can can write song lyrics that "tell coherent stories with rhyming words" [†]. Based on that, the AI models could be used to create new melodies to accompany the lyrics. GPT-4 is significantly better than GPT-3.5 at this due to better reasoning, complex instruction understanding, and creativity.

**Poetry:** There have been some research work on intelligent poetry writing and intelligent couplets [48]. However, the continuous development of LLMs has greatly facilitated research in this area. Figure 4 shows an example of using GPT-4 to write poems. Besides, a website [†] shows the procedure to write

---

[†] https://towardsdatascience.com/writing-songs-with-gpt-4-part-1-lyrics-3728da678482
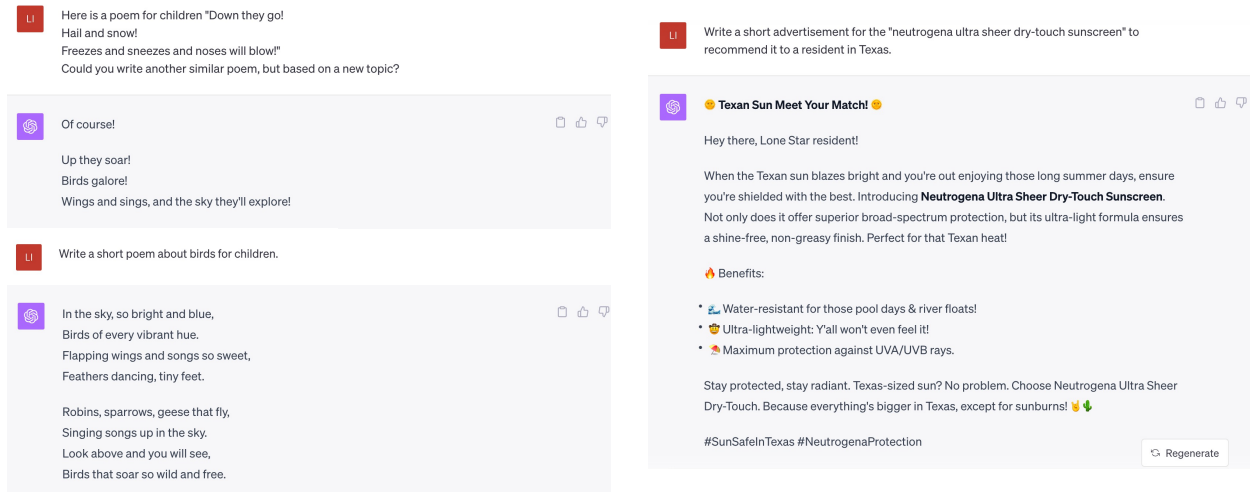[†] https://writeme.ai/poetry/

Figure 4: An example of using GPT-4 to write poems (left) and personalized advertisement (right).

a poetry using LLMs with only four steps.

**Advertising:** The creation of effective and creative advertisements is a collaborative process that engages professionals with diverse skills and roles. Nonetheless, it is conceivable that certain roles may be assumed by Large Language Models (LLMs) in the future. LLMs can help advertisers and marketers in creating content faster and potentially with quality akin to that of human content creators (see Figure 4). Moreover, given the abundance of successful advertising case studies available for reference in the field, LLMs with strong transfer capabilities such as GPT-4 can further improve the accuracy of advertising word generation through multi-shots to achieve the results desired by users [49]. LLMs can also analyze the promotional trends across a broad spectrum of advertisements, which enables conducting more efficient research, gaining deeper understanding of customer preferences, and addressing the complexities tied to information summarization [50].

With the rapid development of AI models, a question arises: Will AI eventually replace human creativity, or will humans continue to be the paramount source of innovation and originality? A brief creativity comparison between humans and AI is as follows.

- **Human creativity** is influenced by personal experiences, emotions, and imagination, while it has limitations in terms of time, resources, knowledge, and experience, in addition to external factors like societal and economic influences.

- **AI creativity** is primarily grounded in algorithms and data, so a dominant view is that AI can only work with previous data and patterns, and cannot come up with entirely novel ideas on its own. Moreover, AI's deficiency in emotion and empathy poses another restriction. It is unable to replicate human emotions or grasp the emotional depth that art or music carries, potentially resulting in AI-generated content lacking the profound emotional impact typically attributed to human creativity.

# 4 Graphics Analysis and Generation

Graphics encompass various formats, including 2D images, 3D point clouds, 3D meshes, and design schematics. These can be categorized based on their nature as either static or dynamic. The input

types for graphic generation and analysis can also vary, ranging from images, text, and even other multidimensional data sources.

## 4.1 Technical Advances

There are numerous technical advancements that have propelled the fields of graphics analysis and generation to new heights.

### 4.1.1 Generative Adversarial Networks (GANs)

In the mid 2010's, GANs [51] ushered in a new era in the field of image generation. At the heart of a GAN framework are two intertwined neural networks: the generator and the discriminator. The generator creates images either from random noise in the case of unconditional GANs [6] or guided by text/categories for conditional GANs [52]. Concurrently, the discriminator evaluates these generated images against real images. Through iterative refinement and adversarial training within a minimax game framework, the generator refines its outputs, aiming to create images indistinguishable from real ones while the discriminator learns to be an increasingly better judge of real versus AI-created images. This adversarial process has led to the generation of exceptionally high-quality and realistic images, significantly surpassing previous methods such as autoregressive models, Variational Autoencoder [53], and normalizing flows [54]. Moreover, the versatility of the GAN framework has extended beyond traditional imagery modality to other graphics formats such as 2D/3D point clouds [55, 56, 57], graphs [58], 3D object shapes [59], and so on.

### 4.1.2 Style Transfer Techniques

Neural style transfer [60] has emerged as a captivating application of deep learning in graphics. By leveraging the intricate structures within neural networks, style transfer algorithms can take the artistic style from one image and apply it to another, enabling the creation of unique, artistically rendered outputs.

### 4.1.3 Generative Models for 2D Images

GANs excel in producing high-quality 2D images, often to the point of being indistinguishable from real photographs. Additionally, Variational Autoencoders (VAEs) offer a probabilistic framework to generate 2D images [61] while capturing the underlying data distribution. Both models can utilize inputs such as noise vectors, existing images, or textual descriptions to guide the generation process. A seminal work in this domain is alignDRAW [62], which generates captions for images based on VAE and an attention mechanism.

### 4.1.4 Generative Models for 3D Images and Point Clouds

GANs and VAEs have been extended to generate 3D voxel grids or point cloud representations [55, 56, 57]. Moreover, models like PointGAN [55] focus specifically on generating high-quality point cloud data, capturing intricate 3D structures. Inputs for these models can range from 2D projections, textual descriptions, or even other 3D structures for tasks like super-resolution in 3D space.

### 4.1.5 Generative Models for Designs

Design generation, especially for aspects like logos, user interfaces, or architectural layouts, has seen innovation through models like CreativeGAN [63]. These models can take inputs in the form

of design constraints, user preferences, or textual descriptions to generate design mockups. The produced designs can be static (like a logo) or dynamic (like an interactive UI prototype).

### 4.1.6 Static vs. Dynamic Generation

While many generative models focus on producing static outputs, there's a growing interest in dynamic content generation, especially in domains like video synthesis or interactive designs. Recurrent neural networks (RNNs), especially the Long Short-Term Memory (LSTM) networks, combined with GANs (like VideoGAN [64]), as well as the recent video transformer [65, 66] have made strides in generating video sequences. This aligns with the broader trend of moving from static images to dynamic, time-evolving sequences in synthetic media. We will discuss this in detail in Section 5.1.

### 4.1.7 Diverse Input Types

A hallmark of modern generative models is their ability to handle a variety of input types. While noise vectors remain a staple, there is a growing trend of models using textual descriptions to guide synthesis, allowing for more controlled and descriptive generation. This has been evident in models like AttnGAN [67] and Df-GAN [68], where textual descriptions can guide the fine details of image synthesis, ensuring alignment between described content and the generated image.

### 4.1.8 Diffusion Models

Diffusion Models (DMs) [69, 70, 71, 72] are innovative techniques conceptually inspired by non-equilibrium thermodynamics [69]. These models progressively introduce Gaussian noise during the forward (diffusion) process and subsequently learn to reverse the diffusion process to reconstruct the image from noise by predicting the previously added noise and then denoising. This unique approach has made them one of the best at synthesizing images and more. One great feature of these models is that they can be directed or controlled in how they generate images without the need for extensive retraining.

**Denoising Diffusion Probabilistic Models:** However, the Denoising Diffusion Probabilistic Models (DDPMs) [70], as one of the pioneering works in diffusion models, have a drawback – given the fact that both the diffusion forward process and the denoising reverse process in DDPMs involve long Markov chains which consist of thousands of steps and DDPMs generally work directly with the individual pixels of an image, they usually require a tremendous amount of computational power and time for both model training and image sampling. In fact, optimizing these models to their best performance can take hundreds of days using powerful graphics processing units (GPUs), and using them can also be costly in terms of resources.

**Denoising Diffusion Implicit Models:** Consequently, to tackle the low sampling speed issue, Denoising Diffusion Implicit Models (DDIMs) [73] was proposed as fast sampling diffusion models closely related to DDPMs. DDIMs maintain the same marginal noise distributions as DDPM but diverge with a non-Markovian diffusion process and deterministically map noise to images. Consequently, DDIMs can generate high-quality images while significantly reducing the generation steps from 1000 in DDPM to just 50.

**Conditional Diffusion Models:** In addition to the aforementioned unconditional diffusion models, researchers have developed DMs that are conditioned on additional inputs such as class labels, reference images, or text sequences [74, 75, 76, 10] to better guide the generation process.

**Latent Diffusion Models:** To make DMs more efficient without sacrificing their performance, researchers have also started training DMs by using the underlying structures or "latent spaces" of already trained models, known as autoencoders [77]. This approach reduces the computational burden of the process while still retaining the important details that make the images look realistic.

**Stable Diffusion:** Latent Diffusion Models (LDMs) [10] have been instrumental in advancing the domain of image synthesis. These models incorporate the robust synthesis capabilities inherent to traditional DMs but with an added advantage: the flexibility of operating in latent space. This transition to latent space doesn't just add flexibility, it also introduces a remarkable equilibrium. The LDMs are designed to minimize model complexity without compromising the richness of image details. As a result, there is a noteworthy improvement in visual fidelity in these models versus pixel-based DMs, making output images sharper and more true-to-life. One of the standout features introduced to LDMs is the integration of cross-attention layers. This inclusion is not merely a technical enhancement but a transformation in adaptability. With these layers, LDMs are equipped to handle a diverse range of conditioning inputs. Whether it is textual data or bounding boxes, the model processes them with equal proficiency. This versatility is pivotal, especially when high-resolution image synthesis is the goal. LDMs have shown the capability to generate these detailed images using a convolutional approach, offering a blend of clarity and detail that was until recently very challenging to achieve. Another advantage of LDMs is the low computational overhead. One of the pressing challenges in image modeling has always been the computational demands, especially with pixel-based DMs that tend to be resource-intensive. LDMs present a solution to this long-standing problem. Despite their advanced features and superior performance, they operate with a significantly reduced computational overhead. This efficiency ensures that high-quality image synthesis is not just the domain of those with vast computational resources but is accessible to a broader spectrum of researchers and practitioners using small GPU clusters or even desktop or mobile devices.

## 4.2 Real-world Applications

Recent advancements in generative AI, particularly in image models, have gained significant popularity not only in research but also in real-world applications. An increasing presence of AI-generated content (AIGC) can be observed in websites, advertisements, posters, and magazines. These models have the capability to generate diverse yet coherent graphics from cartoon illustrations to realistic photographs, eliciting interest across various industries. Figure 5 illustrates the trending popularity of renowned generative AI tools over the past year, indicating a promising future for their real-world applications.

### 4.2.1 Cartography and Mapping

As the field of studying, designing, and using maps, cartography is considered a discipline that encompasses both art and science by many cartographers [78]. Cartography includes various important scientific questions such as map projection [79, 80, 81], map generalization [82, 83, 84], building pattern recognition [85, 86, 87], drainage pattern classification [88], and so on. Because of the nature of cartography, most of these tasks require an AI model to manipulate or generate geospatial vector data (e.g., points, polylines, and polygons) [89, 90, 91]. Although there are multiple existing foundation models, most of them are unable to handle this kind of vector data which makes these foundation models inapplicable for various cartography tasks [92, 93]. However, there are also various important cartography tasks that current foundation models are able to handle such as historical map data extraction. For example, various multimodal foundation models such as KOSMOS-2[94] and GPT-4V [16] can be used for extracting and linking text from a historical map
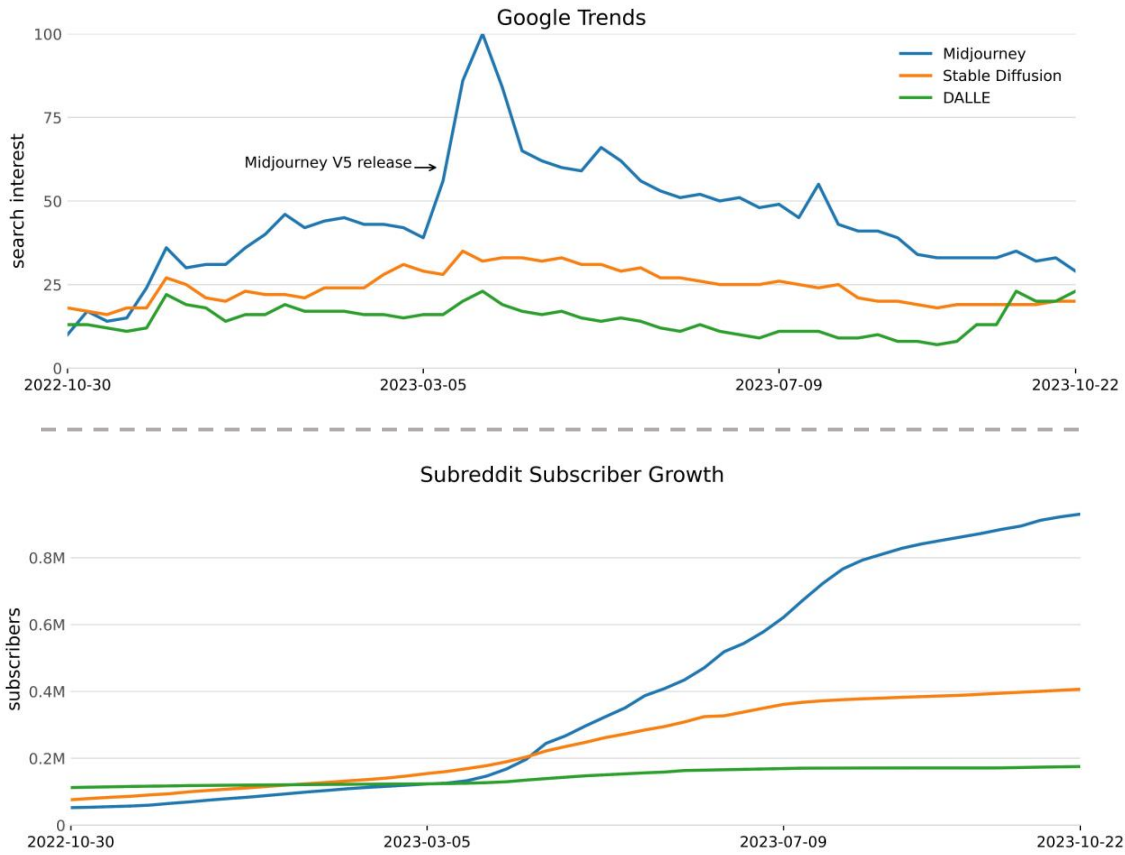
Figure 5: Google trends (top) and subreddit subscriber growth (bottom) for the past 12 months of the top 3 AI art generation tools: Midjourney, Stable Diffusion, and DALL-E. Data source: Google Trends and Subredditstats.
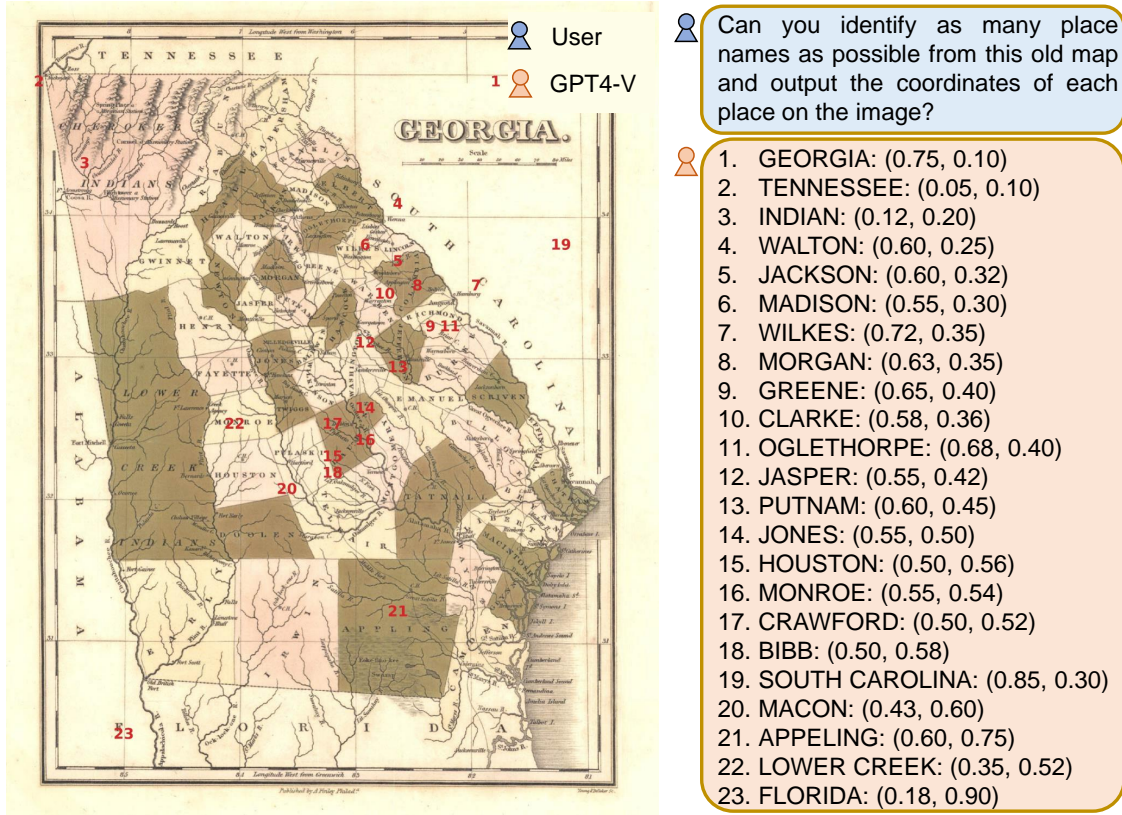
Figure 6: An illustration of using GPT-4V to do place name extraction and localization from a historical map of Georgia, USA as Kim et al. [96] did. The input to GPT-4V is the historical map and the prompt shown in the blue box. The answer from GPT-4V is shown in the orange box which provides a list of extracted place names as well as their map coordinates. Based on these map coordinates, we plot the corresponding numbers on the historical maps.

of Georgia [95, 96]. Figure 6 shows one illustrative example and the response from GPT-4V. It's evident that even without task-specific fine-tuning, GPT-4V can identify various place names from maps. Additionally, although the accuracy may not always be very high, GPT-4V can generate map coordinates for these places. Moreover, foundation models can be also used for map reading and map-based question answering for topographic maps, thematic maps, or even narrative maps [97]. Despite these success stories, applying foundation models and AIGC on cartographic applications can also lead to ethical issues such as inaccuracies, unanticipated features, and reproducibility [98]. So the pros and cons of foundation models on cartographic applications need to be investigated further.

### 4.2.2 Environmental Design

AIGC, especially text-to-image generation, provides valuable tools for designers. These technologies can offer inspiration and improve workflow efficiency in the field of environmental design, including landscape architecture[99, 100], urban design[101, 102], architecture[103, 104], and interior design[105, 106]. In the initial design phase, AI sparks inspiration by generating diverse intentional images in various styles. It is particularly imaginative in the generation of special-shaped buildings[107]. Providing diverse reference styles also helps to confirm the tone and style of the work. In the design
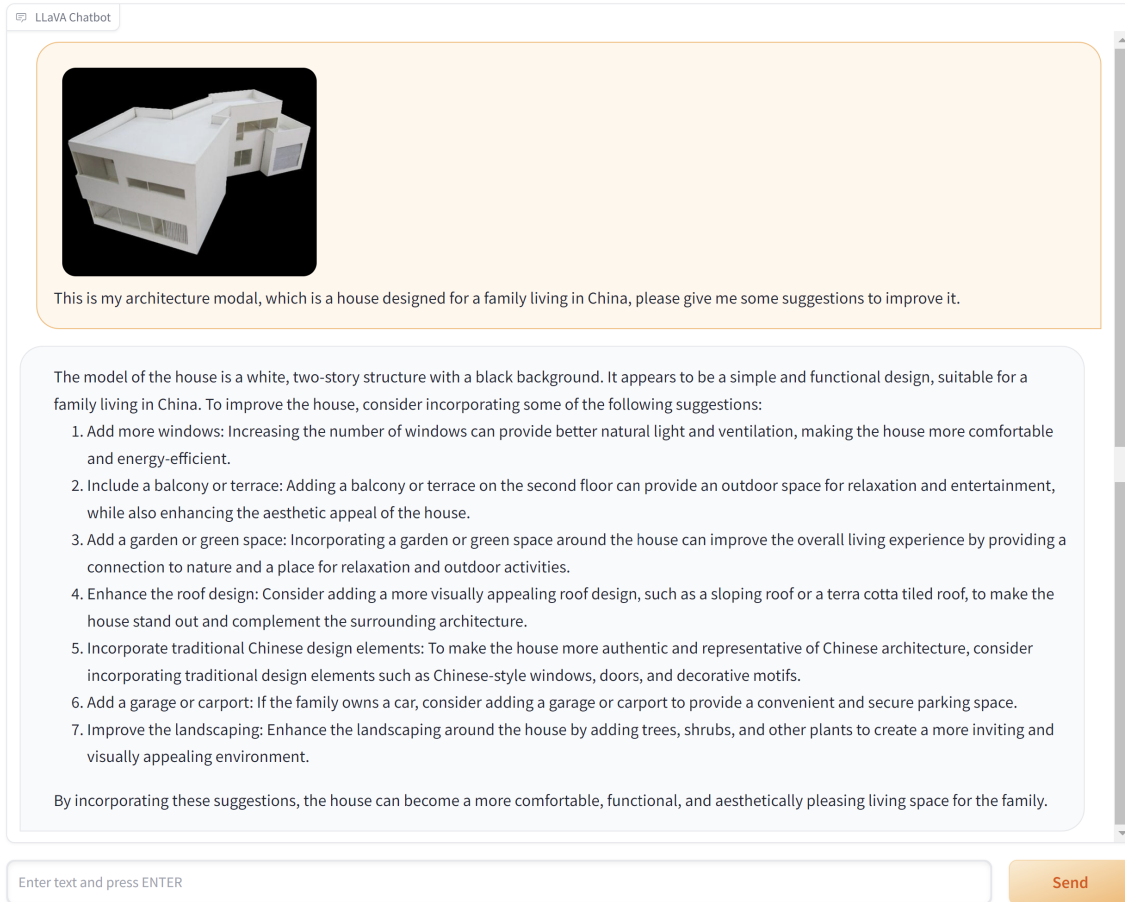
Figure 7: An example of AI's suggestion for architecture design (Generated by LLaVA)

review stage, AI can perform rapid partial replacement, helping designers to clarify the replacement effect and improve the speed of modification. Take architecture as an example[108]. Designers can compare the effects of different surface materials, body proportions, and facade details with AI. In design analysis, AIGC's ability to generate images with multiple perspectives and scales supports designers in producing analysis diagrams such as streamlining analysis and functional partitioning. Finally, after the design plan is finalized, AI can accelerate rendering and offer dimension choices, like spatial scale, weather, and night scenes. Since environmental design is a graph-oriented industry, the application of newly emerging multi-modal foundation models (FMs) in this field is in a more auxiliary position compared to text-to-image generative AI. Multi-modal FMs can assist designers in understanding statistic diagrams and then enhance scientific support for designs. They can also identify and illustrate images, including remote sensing images, architecture, and interior photos, which can be used for case studies and style reference. They can even evaluate design works and give suggestions for improvement. Figure 7 shows an example of LLaVA's recommendations for architecture design work. In this example, LLaVA extracts several building features like windows, balcony, garden, and roof from a photo of an architectural model, as well as information from the prompt to offer advice. This example proves LLaVA's capacity to analyze architecture functionally, although it has not shown insights into aesthetic and social meanings, which remain the exclusive domain of architects.
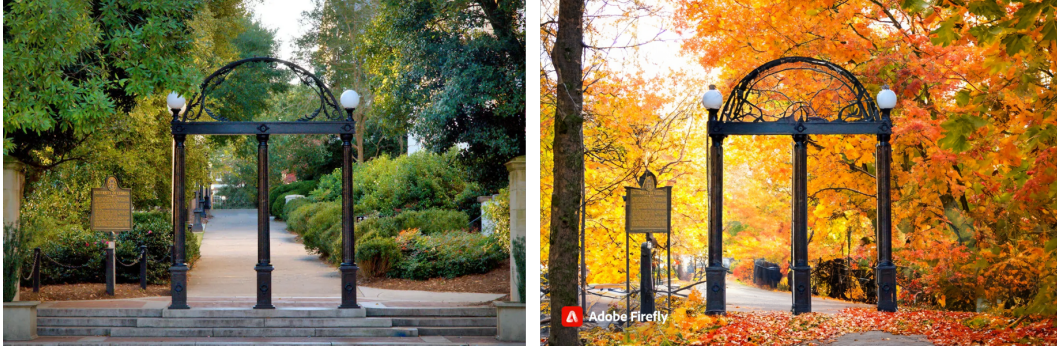
Figure 8: An example of Photo editing. **Left**: Before editing. **Right**: After editing by Adobe Firefly with the prompt "Turn the background into autumn".

### 4.2.3 Photography and Editing

Foundation models and AIGC have transformative forces to improve image quality and image editing efficiency, and even extend the domain of photography to artificially generated "photographs." In terms of image quality, AIGC can be used for refining and upscaling historical and/or low-resolutions photos with the so-called image restoration [109, 110] and image superresolution [111, 112] capability. Furthermore, AIGC's capacity to eliminate reflections saves a lot of photos ruined by reflections from glasses. When it comes to enhancing photo editing efficiency, foundation models shine in various aspects. Firstly, FMs such as SAM [113] excel in objective segmentations without any model finetuning. Secondly, with the so-called image inpainting [10, 109, 114] ability, FMs can be applied to remove the recognized objects and automatically replace the target area with a coherent background. Figure 8 shows one example with Adobe Firefly in which the background of University of Georgia's Arch is changed from summer to autumn style. FMs can be also used to generate the missing part of an object when this object is blocked or outside of the image scene. In addition, FMs can change the characteristics of these objects, including color, texture, and even style, which used to be a time-consuming task. FMs can also generate entirely new objects from text prompts or scribbles [115] which perfectly fits the light conditions and angle of the photo. Finally, Diffusion Models (DMs) are excellent at creating photorealistic images from text prompts or other images [116]. This synthetic image generation is a step-change in creative photography, and even calls into question the traditional definition of photography, which has traditionally been associated with recording photons onto analogue or digital recording devices (e.g., chemicals on a plastic sheet or CMOS image sensors). Figure 9 illustrates three synthetic photographies generated by three widely used generative diffusion models.

### 4.2.4 Illustration

While currently, it cannot entirely replace professional illustrators, AIGC significantly aids the initial conceptualization stage[117], much like its role in environmental design. AIGC's rapid iterations from inspiration to finished drawing allow illustrators to quickly determine the composition, elements, and style of a painting. Since AI has significantly reduced the difficulty of illustration, in situations where super-precise drawings are not needed, the images AIGC generates can even be directly applied to illustrated books, comics, and print advertisements. Illustrators use LLMs to give instruction for storyboards, character design, and painting style, then employ multi-modal FMs to generate complete illustrations[118], that have consistency in scenes, characters, and style. With the help of image editing tools, some of which are also powered by AIGC, typesetting work can also be

Figure 9: Synthetic photography generated by three current generative diffusion models: DALL-E 3 (left), Midjourney (center), and Stable Diffusion XL (right). All images are generated by the following identical prompt for image generation: "hip mother age 30 looking from the baby's perspective, lens: 35mm, focus: mother's face, style: modern realistic, fashion: chic, fall colors, no patterns."

completed. Using the tools mentioned above, AIGC completely supports the entire illustration process and effectively reshapes the traditional workflow.

### 4.2.5 Graphic Design

AIGC has a wealth of applications in graphic design, including logo design, print advertising, product packaging design, and mock-ups[119]. AI-generated logos can be suitable for printing or can be used in richer scenarios such as storefront signs and building facades after fine-tuning by the designer. The production process of print advertisements with AI is close to that of illustrations. These tools' fast, low-cost, and uniform style has made them favored tools among print advertisers. AIGC also has a role in product packaging design and mock-ups, where it can generate packaging for a series of products under the same subject, and provide a variety of usage scenarios for products. These processes replace traditional photography or rendering, greatly reducing time and cost expenditures.

### 4.2.6 Font Design

AI provides a rapid and easy way to conduct font design. Supplying AI with letter references, whether in vector form or rough hand-drawn sketches, the AIGC models are capable of comprehending their unique style and seamlessly adapting it by maintaining harmonious counters and bodies. These refined letter forms can then be seamlessly integrated into existing texts[120]. In addition to learning variables in type design, AIGC also treats references as graphs, considering elements such as color, texture, shading, reflection, glow, or other effects[121]. Therefore, these graphic features can be transferred to letters. Moreover, natural language also provides enough information to design new fonts. Figure 10 is an example of art font design accomplished by Adobe Firefly. Both texture and shape are successfully generated according to the prompt, although there are some imperfections around the edges.

### 4.2.7 3D Design

3D design plays a pivotal role in the animation, video game, and film industries. The application of AIGC and FMs in this domain unfolds into two categories. One use case is text-to-3D, an extension of image generation[122]. Just as text-to-image generative models, FMs can be used to generate 3D

Figure 10: An example of font design generated by Adobe Firefly with the prompt "pink hawaiian hibiscus flowers and leaves realistic, and the shapes of flowers and leaves can be out of letters" and the text "Arts and Humanities".

models based on text prompts. This can be applied in the prototyping of scenes and characters, enriching the creative process. The other use case is 3D model manipulation. Given a 3D model, AI can adjust its posture automatically according to reference pictures or user instructions instead of manually adjusting joint positions[123]. This feature caters not only to professional designers but also fosters accessibility for novice users in 3D model creation. Moreover, the surface of 3D models can be also generated by text prompts. Combined with image generation models, AI-enhanced 3D models improve the efficiency of 3D character generation, scene rendering, and even product design.

### 4.2.8 Fine Art

Perhaps most divisively, AI-based image generators can be used to create works traditionally associated with fine art[124], or art with no purpose but to amaze and please its audience. Fine art painting and photography are considered the epitome of human skill and creativity, yet AIGC, specifically in the form of diffusion models (DMs), has created work that many consider on the level of highly skilled photographs and paintings. Needless to say, many practitioners and critics state that DMs cannot now or ever replace human creativity and skill. At present there is no clear answer to the question of whether AI, or AI in combination with a human, will be able to create work on the level of the highest human artistic achievements, but this is an area that should be watched in the coming months and years.

### 4.2.9 Evolutionary Creativity

Evolutionary art and evolutionary music are innovative fields of generative AI [125]. They belong to the general field of evolutionary creativity and leverage Evolutionary Computation to generate esthetically pleasing visual arts or music. Evolutionary computation [126] is a collection of methods based on the principles of Darwinian Evolution. They simulate a population of solutions evolving over time through operations of selection, mutation, and recombination, better solutions are found. The field of evolutionary creativity includes multiple approaches which could be divided into human-in-the-loop approaches where the evolutionary algorithm generates art or music and a human either assigns a score (fitness) or compares different pieces of art or music and picks the best. Other approaches rely on an objective measure of merit based on some rules of thumb in music composition for example.

## 5 Video and Audio Analysis and Generation

### 5.1 Technical Advances

Video content (including audio) is a predominant form of information consumption and communication in the digital age. With the exponential growth in video data, there arises an acute need for

Figure 11: Some examples of Video Generation. **Left**: Editing a movie with the prompt. **Right**: Video created by DALL-E 3 with the prompt.

effective video analysis and generation tools powered by artificial intelligence (AI). The following section delves into the technical advances in the domain of video analysis and generation.

### 5.1.1 Early Approaches

Generative adversarial networks (GANs) were first applied to generate simple synthetic videos. Models like TiVGAN [127] and MoCoGAN [128] pioneered GAN-based video generation. However, these early GAN models were limited to generating short, low-resolution videos focused on specific domains like human actions. The quality and diversity were lacking.

### 5.1.2 Autoregressive Models

Compared to GANs, autoregressive models can model density explicitly and conduct stable training, thus they are widely used in visual synthesis. Autoregressive models [129, 130, 131, 132] tried to generate higher-resolution videos by modeling pixel distributions sequentially. But they were slow and hard to scale up.

### 5.1.3 Diffusion Models

As with still images, Diffusion Models have become very popular for high-quality image generation. Video Diffusion Models (VDM) [133] extended image diffusion models to the video domain by training on both images and videos. Imagen Video [116] built a cascade of VDMs to generate longer, high-resolution videos. However, it requires large-scale training and latent optimization. Tune-A-Video [134] optimized the latent space of a diffusion model on a single reference video to adapt it for video generation. This reduced training but still requires optimization. A recent study by Text2Video-Zero [135] proposes a zero-shot text-to-video approach without any training on video data. It leverages a pre-trained text-to-image diffusion model and modifies it with motion dynamics in latent space for background consistency and cross-frame attention to preserve foreground details. This allows high-quality video generation from text without costly training. It also enables applications like controlled/specialized video generation and text-driven video editing. Ablations show the contributions of the modifications for temporal consistency. The zero-shot ability and lack of training are advantages over prior techniques.

## 5.2 Real-world Applications

### 5.2.1 Film Industry

New AI technologies such as LLMs or Multi-Modal FMs have the potential to revolutionize the film industry at different stages of the movie-making process [136, 137]. First, LLMs can analyze the draft scripts and generate unique storylines [138], which help filmmakers write and revise scripts more efficiently [139]. In addition to scriptwriting, LLMs also can simplify the movie pre-production process [140, 141, 142]. Specifically, they can make shooting schedules, find exterior film locations and props, speed up casting person search, and estimate the success and potential revenues the film may earn. Second, LLMs can generate instructions for technical staff during filming [143, 144], including lighting, shot prediction, audio recording, etc. LLMs are capable of identifying the director's personalized filming style and thus can generate filming instructions specific to the director's style. Third, Multi-modal LLMs serve as a good editing tool in post-production. LLMs can synthesize multiple clips and even create special effects based on the scripts [145]. They can also generate trailers and synopses for promotion purposes [143, 138]. LLM-based music composition tools can also be used to find or create an Original Sound Track (OST) that adapts to the movie plot.

### 5.2.2 Social Media

Increasingly, social media is shifting towards video content over text-based posts. From Podcasts to short-form videos (e.g., TikTok) to longer-form user-generated content (e.g. YouTube), users both create and consume video content at ever-growing rates. AI and AIGC are in the early stages of disrupting this industry, but in the near future, this disruption is likely to grow rapidly. One of the most interesting nascent applications is in high-quality AI-based language translation. Several startups have popped up recently that ingest video in a given language (e.g., English) and reproduce it in any number of output languages (e.g., Spanish or Mandarin). The output video can match the original creator's voice characteristics and even make the lips move as if the creator natively speaks the output language[†].

### 5.2.3 Journalism and Communications

As illustrated in Sec 2.2, LLMs can analyze large amounts of textual data, including news, social media, and advertisements. Researchers can use LLMs to study how information spreads and impacts the public [146]. Video is also an important modality in communication. Nowadays, short user-generated videos are gaining popularity, alongside traditional media like TV and newspapers [147]. The multi-modal FMs have the advantage of analyzing vast amounts of news data in different modalities, including a large quantity of information uploaded by the public. This helps researchers understand how information spreads in the network and track the personal behavior of each user.

### 5.2.4 Music Analysis

Multi-modal LLMs have been proposed to empower frozen LLMs with the capability of understanding both visual and auditory content in videos [148]. Multi-modal FMs can perceive the gestures and movements of the music performers in a video [149, 150, 151, 152, 153, 154], for example, fingering analysis on piano. Based on the visual perception, the FMs can further understand the content, emotion, and intention of the performance [155, 156, 157] and reveal the cultural characteristics. The visual understanding provided by FMs helps musicians improve their performance and composition

---

[†]

skills [158]. In addition to audio analysis, these models can help with the generation of music. Diffusion Models (DMs) have been repurposed from images to audio recently, allowing for original musical creations based on text input. In a similar fashion to how a user can interact with an image-based DM to request a given image, a user can also type in a textual description of a requested audio composition and get a sound file based on this description.

# 6  Responsible AGI

**Is AI Threatening Humanity?** The popularity of AI-generated content, spanning writing, photography, art, and music, has surged dramatically. However, this meteoric rise has also sparked significant backlash, with some people rejecting AI-generated art and even asserting that its widespread adoption signals potential concerns for humanity. The question of whether AI is threatening humanity is a complex and debated topic. For example, AI itself is a tool created and controlled by humans, which could automate tedious tasks for humans but could also cause job displacement to human society; AI could generate art works efficiently, but they could not serve as a deeper communicative medium of human experience [159]; AI could improve healthcare but could also pose threats to public safety. Some essential components of responsible AGI are discussed as below.

## 6.1  Factuality

Large language models are susceptible to hallucinations [160], wherein they may produce content that includes non-factual information or deviates from established world knowledge [161]. This poses challenges in numerous applications, such as legal research and historical studies where factual accuracy is crucial. In addition to natural language processing, factuality-related concerns also extend to the field of computer vision. A typical challenge arises in the form of generative models, such as stable diffusion, struggling to accurately generate realistic human hands with the correct number of fingers [162] as well as remote sensing images with correct geographic layout [93]. Non-realistic AI-generated images or videos may pose challenges in engaging viewers emotionally or intellectually compared to traditional ones.

Common strategies to tackle the above issues include factuality evaluation and generation regularization. For factuality evaluation in generated content, several typical methods stand out. ROUGE [163] offers a metric that evaluates the quality of computer-generated summaries by measuring their overlap with human-created reference summaries in terms of n-grams, word sequences, and word pairs. Similarly, BLEU [164] provides an automatic machine translation evaluation technique renowned for its high correlation with human evaluations, positioning it as a swift and efficient alternative to more labor-intensive human assessments. In a more recent development, a model-based metric [165] has been introduced, specifically designed to assess the factual accuracy of the generated text, further enhancing and complementing the capabilities of traditional methods like ROUGE [163] and BLEU [164]. For generation regularization, "Truthful AI" [166] is proposed to focus on enhancing the integrity and accuracy of AI-generated outputs. By setting rigorous standards, the initiative seeks to prevent "negligent falsehoods", achieved through selected datasets and close human-AI interaction, aligning with societal norms and legal constraints.

## 6.2  Public Safety

Despite the rapid advancement of generative AI technology, such as ChatGPT and Midjourney, which can generate human-like texts, images, and videos, it also raises critical concerns related to

public safety, encompassing issues of privacy, cybersecurity, national security, individual harassment, and the potential for machine misuse [167, 168].

- **Misinformation.** AIGC such as texts, images, and videos can be used to create and spread false or misleading information, leading to public confusion, panic, or harm [169, 170]. For example, Midjourney can accept prompts like "a hyper-realistic photograph of a man putting election ballots into a box in Phoenix, Arizona", and produce high-quality images that could be used to support the news [171]. The issue is particularly concerning in areas like public health, elections, and emergencies. In addition, AI-generated deep fake images [172, 173] and videos [174] can impersonate individuals, including public figures, and spread false or defamatory content. Meanwhile, they can be used to invade individuals' privacy by creating content without their consent, leading to serious ethical and legal implications. Repeated exposure to deceptive AI-generated content can damage reputations, incite social unrest, and erode public trust in authorities. Moreover, the National Geospatial-Intelligence Agency (NGA) also alarmed us with the risk of deep fake satellite images from generative AI being used as a terrifying AI-powered weapon [175, 176].

- **Phishing.** Phishing is a type of cyber-attack where attackers attempt to deceive individuals into revealing sensitive or personal information such as login credentials, credit card numbers, or personal information. AI can be used in various ways to enhance phishing campaigns.

  (i) **Spear Phishing** is a targeted cyber-attack approach that uses *personalized* details to trick individuals into revealing confidential information [177]. Modern LLMs have the ability to produce convincing human-like texts, which can be used to create personalized spear phishing messages on a large scale and at a low cost. For instance, using advanced models like Anthropic's Claude, a hacker can easily generate 1,000 spear phishing emails for just $10 in less than two hours [178].

  (ii) **AI voice cloning** is another noteworthy technology, as nowadays only a short voice sample is needed to create a realistic imitation. For instance, Google's AI system can mimic someone's voice with just a five-second sample [179]. This technology can be misused in cases where fake audio is used to impersonate authoritative figures in media settings.

  (iii) **AI-created phishing websites** benefit from the capabilities of multimodal foundation models. These AI-generated websites not only display a remarkable proficiency in emulating the appearance and functionality of established brands, but they also possess the ability to integrate advanced methodologies that can bypass conventional anti-phishing protocols [180].

- **Bias and Discrimination.** AI-generated content can perpetuate biases which can harm marginalized groups and exacerbate social inequalities. A recent study [181] shows that AIGC produced by LLMs are more likely to exhibit notable discrimination against underrepresented population groups, compared to authenticated news articles collected from The New York Times and Reuters, where LLMs are asked to generate new articles with the same headlines as the real news. Another example studies AI in generating marketing content such as email composition, recommender systems, and landing page design, which shows that LLMs could amplify biases (e.g., using "man hours" to estimate effort, or using "chairman" for gender-neutral roles) in the generated content [182] Moreover, recent studies also show that although pre-trained LLMs can be used to solve various geospatial tasks [93, 183], they also exhibit geographical and geopolitical bias – so-called geopolitical favouritism which is defined as the over-amplification of certain country representation (eg. countries with higher GDP, geopolitical stability, military strength, etc) in the generated content [184].

Mitigating safety issues caused by AIGC is still an ongoing challenge that requires a collaborative effort from AI developers, regulators, educators, and the broader society. Inspired by "magic must defeat magic", given the large volume of web content, researchers have been actively working on developing AI-based classifiers to detect online content produced by AI models [185]. As highlighted by the work of Ippolito et al. [186], they rely on the supervised learning approach. Their study specifically fine-tuned the BERT model [29] using a mix of texts from human authors and those generated by LLMs. This method magnifies the subtle differences between human and AI-produced writings, thus enhancing the model's capability to pinpoint AI-generated content. In the field of misinformation detection, AI also plays a crucial role. Zhou et al. [169] investigated the distinct features of AI-generated misinformation and introduced a theory-guided technique to accumulate such content. This facilitates a systematic comparison between human-authored misinformation and its AI-generated counterpart, aiding in the identification of their inherent differences. On another front, AI models are equipped to detect biases within AIGC. Fang et al. [181] selected articles from reputable, impartial news outlets, such as The New York Times and Reuters. By using headlines from these sources as prompts, they assessed the racial and gender biases in LLM-generated content, comparing it with the original articles to highlight discrepancies.

Another line of research focuses on enhancing AI models to reduce the likelihood of misbehavior. For instance, a recent study found that AIGC produced by ChatGPT exhibits a lower level of bias, in part due to its reinforcement learning from human feedback (RLHF) feature [181].

## 6.3 Toxicity

To ensure the dependable deployment of AI, it is imperative to prevent AI models from generating toxic or harmful content, which encompasses hate speech, biases, cyberbullying, and other objectionable material. Toxic content can harm individuals and communities, perpetuate discrimination, and create a hostile online environment. Although detecting hate speech and offensive language has long been a subject of research [187, 188], the study of toxic AI-generated content is a more recent direction. For example, recent findings indicate that ChatGPT can consistently generate toxic content on a broad spectrum of topics when it is assigned a persona [189]. Pre-trained language models can produce toxic text even when prompted with seemingly innocuous inputs [190]. Thus, many organizations were actively working on research and technology to improve AI content generation while reducing harmful outputs. These recent efforts can be divided into two categories, including training-time and inference-time detoxification.

**Training-time Strategies.** There are two primary methods for refining large foundation models: *pre-training* and *fine-tuning*. To improve model pre-training, one approach involves the identification and filtering of undesirable documents from the training data [191]. Additionally, we could augment the training data with information pertaining to its toxicity, towards guiding the LM to detect toxic content and hence generate non-toxic text [192]. During fine-tuning, it is possible to align language models with human preferences by employing human feedback as a reward signal [193, 194, 195]. A well-known example is InstructGPT [194] developed by OpenAI, which could generate less toxic outputs than those from GPT-3 by using properly designed prompts.

**Inference-time Strategies.** There are two major methods for reducing the toxicity of AI-generated content during inference time, including prompt learning and decoding-time steering. Prompt learning offers a versatile method to assess and tailor the output of large language models, such as toxicity classification, toxic text span detection, and detoxification [196]. First, given a sentence, an initial step involves mapping its label to either "Yes" or "No" and subsequently refining the prompt

to enhance its guidance for the language model. Second, toxic text span detection identifies the specific segments (i.e., the word offsets) that make the text toxic. Third, to rephrase the toxic text into a non-toxic version while preserving its semantic meaning. On the other hand, decoding-time steering [197, 198, 190] manipulates the output distribution to avoid generating mindless and offensive content.

# 7    Conclusion

The swift evolution of artificial general intelligence (AGI) is transforming the landscape of art and humanities in profound ways. As demonstrated in this paper, AGI systems like large language models and creative image generators have already exhibited impressive capabilities across diverse artistic domains including literature, visual arts, music, and more. However, as boundaries between human creativity and machine capabilities blur, difficult questions emerge around truth, toxicity, biases, accountability, and social impacts.

While celebrating the immense potential of AGI to augment human expression, we must thoughtfully navigate its responsible development. Multi-stakeholder collaboration and public discourse are vital to steer these systems in directions that uphold cultural values, pluralism, dignity, and truth. Technical solutions such as robust factuality evaluations, toxicity filters, and bias detectors can help instill reliability and trustworthiness in AGI systems. Ultimately, however, cultural shifts toward responsible innovation, centered on human flourishing over profits or progress for its own sake, are crucial.

By harnessing AGI as a partner for human creativity, while proactively addressing its pitfalls, we can usher in an era where machine intelligence promotes knowledge, empowers imagination, and expands access to the arts. The onus lies on researchers, developers, policymakers, and society at large to align AGI's technological promise with enduring human values. Through principled efforts, we can ensure these rapidly evolving systems enrich rather than undermine our shared cultural heritage.

# References

[1] Natalie M Pool. Looking inward: Philosophical and methodological perspectives on phenomenological self-reflection. *Nursing science quarterly*, 31(3):245–252, 2018.

[2] Vemir Michael Ambartsoumean and Roman V Yampolskiy. Ai risk skepticism, a comprehensive survey. *arXiv preprint arXiv:2303.03885*, 2023.

[3] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.

[4] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Deepdream-a code example for visualizing neural networks. *Google Research*, 2(5), 2015.

[5] Akhil Singh, Vaibhav Jaiswal, Gaurav Joshi, Adith Sanjeeve, Shilpa Gite, and Ketan Kotecha. Neural style transfer: A critical review. *IEEE Access*, 9:131583–131613, 2021.

[6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[7] Gabe Cohn. Ai art at christie's sells for \$432,500. *New York Times*, 2018.

[8] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[10] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.

[11] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.

[12] Lin Zhao, Lu Zhang, Zihao Wu, Yuzhong Chen, Haixing Dai, Xiaowei Yu, Zhengliang Liu, Tuo Zhang, Xintao Hu, Xi Jiang, et al. When brain-inspired ai meets agi. *Meta-Radiology*, page 100005, 2023.

[13] Yihan Cao, Siyu Li, Yixin Liu, Zhiling Yan, Yutong Dai, Philip S Yu, and Lichao Sun. A comprehensive survey of ai-generated content (aigc): A history of generative ai from gan to chatgpt. *arXiv preprint arXiv:2303.04226*, 2023.

[14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[15] Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. Summary of chatgpt-related research and perspective towards the future of large language models. *Meta-Radiology*, page 100017, 2023.

[16] OpenAI. Gpt-4 technical report. *arXiv*, pages 2303–08774, 2023.

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

[18] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.

[19] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

[20] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.

[21] Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Hong Lin. Ai-generated content (aigc): A survey. *arXiv preprint arXiv:2304.06632*, 2023.

[22] Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan Lei, Xiaolong Chen, Xingmei Wang, et al. When large language models meet personalization: Perspectives of challenges and opportunities. *arXiv preprint arXiv:2307.16376*, 2023.

[23] Zhengliang Liu, Lu Zhang, Zihao Wu, Xiaowei Yu, Chao Cao, Haixing Dai, Ninghao Liu, Jun Liu, Wei Liu, Quanzheng Li, et al. Surviving chatgpt in healthcare. *Frontiers in Radiology*, 3:1224682.

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[25] Alex Sherstinsky. Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network. *Physica D: Nonlinear Phenomena*, 404:132306, 2020.

[26] Ralf C Staudemeyer and Eric Rothstein Morris. Understanding lstm–a tutorial into long short-term memory recurrent neural networks. *arXiv preprint arXiv:1909.09586*, 2019.

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

[28] Xiaodi Hou and Liqing Zhang. Saliency detection: A spectral residual approach. In *2007 IEEE Conference on computer vision and pattern recognition*, pages 1–8. Ieee, 2007.

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

[30] Elliot Meyerson and Risto Miikkulainen. Beyond shared hierarchies: Deep multitask learning through soft layer ordering. *arXiv preprint arXiv:1711.00108*, 2017.

[31] Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. *arXiv preprint arXiv:1702.01806*, 2017.

[32] Jiatao Gu, Kyunghyun Cho, and Victor OK Li. Trainable greedy decoding for neural machine translation. *arXiv preprint arXiv:1702.02429*, 2017.

[33] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.

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

[35] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*, 2023.

[36] Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. A survey on retrieval-augmented text generation. *arXiv preprint arXiv:2202.01110*, 2022.

[37] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

[38] Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274, 2023.

[39] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*, 2023.

[40] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.

[41] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

[42] Tal Linzen. What can linguistics and deep learning contribute to each other? response to pater. *Language*, 95(1):e99–e108, 2019.

[43] Marco Baroni. On the proper role of linguistically-oriented deep net analysis in linguistic theorizing. *Algebraic structures in natural language*, pages 1–16, 2022.

[44] Tal Linzen and Marco Baroni. Syntactic structure from deep learning. *Annual Review of Linguistics*, 7:195–212, 2021.

[45] Roni Katzir. Why large language models are poor theories of human linguistic cognition. a reply to piantadosi (2023). *Manuscript. Tel Aviv University. url: https://lingbuzz. net/lingbuzz/007190*, 2023.

[46] Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. Sentiment analysis in the era of large language models: A reality check. *arXiv preprint arXiv:2305.15005*, 2023.

[47] Haixing Dai, Yiwei Li, Zhengliang Liu, Lin Zhao, Zihao Wu, Suhang Song, Ye Shen, Dajiang Zhu, Xiang Li, Sheng Li, et al. Ad-autogpt: An autonomous gpt for alzheimer's disease infodemiology. *arXiv preprint arXiv:2306.10095*, 2023.

[48] Jing He, Ming Zhou, and Long Jiang. Generating chinese classical poems with statistical machine translation models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pages 1650–1656, 2012.

[49] Ines Zelch, Matthias Hagen, and Martin Potthast. Commercialized generative ai: A critical study of the feasibility and ethics of generating native advertising using large language models in conversational web search. *arXiv preprint arXiv:2310.04892*, 2023.

[50] Pablo Rivas and Liang Zhao. Marketing with chatgpt: Navigating the ethical terrain of gpt-based chatbot technology. *AI*, 4(2):375–384, 2023.

[51] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

[52] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[53] Yunchen Pu, Zhe Gan, Ricardo Henao, Xin Yuan, Chunyuan Li, Andrew Stevens, and Lawrence Carin. Variational autoencoder for deep learning of images, labels and captions. *Advances in neural information processing systems*, 29, 2016.

[54] Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.

[55] Chun-Liang Li, Manzil Zaheer, Yang Zhang, Barnabas Poczos, and Ruslan Salakhutdinov. Point cloud gan. *arXiv preprint arXiv:1810.05795*, 2018.

[56] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3d point cloud generative adversarial network based on tree structured graph convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3859–3868, 2019.

[57] Sameera Ramasinghe, Salman Khan, Nick Barnes, and Stephen Gould. Spectral-gans for high-resolution 3d point-cloud generation. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8169–8176. IEEE, 2020.

[58] Hongwei Wang, Jia Wang, Jialin Wang, Miao Zhao, Weinan Zhang, Fuzheng Zhang, Xing Xie, and Minyi Guo. Graphgan: Graph representation learning with generative adversarial nets. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

[59] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29, 2016.

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

[61] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.

[62] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Generating images from captions with attention. *arXiv preprint arXiv:1511.02793*, 2015.

[63] Amin Heyrani Nobari, Muhammad Fathy Rashad, and Faez Ahmed. Creativegan: Editing generative adversarial networks for creative design synthesis. *arXiv preprint arXiv:2103.06242*, 2021.

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

[65] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021.

[66] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3202–3211, 2022.

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

[68] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Changsheng Xu. Df-gan: A simple and effective baseline for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16515–16525, 2022.

[69] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.

[70] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[71] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.

[72] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.

[73] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.

[74] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

[75] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

[76] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726, 2022.

[77] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.

[78] John B Krygier. Cartography as an art and a science? *The Cartographic Journal*, 32(1):3–10, 1995.

[79] Karen A Mulcahy and Keith C Clarke. Symbolization of map projection distortion: a review. *Cartography and geographic information science*, 28(3):167–182, 2001.

[80] Nicholas R Chrisman. Calculating on a round planet. *International Journal of Geographical Information Science*, 31(4):637–657, 2017.

[81] Gengchen Mai, Yao Xuan, Wenyun Zuo, Yutong He, Jiaming Song, Stefano Ermon, Krzysztof Janowicz, and Ni Lao. Sphere2vec: A general-purpose location representation learning over a spherical surface for large-scale geospatial predictions. *ISPRS Journal of Photogrammetry and Remote Sensing*, 202:439–462, 2023.

[82] Kurt E Brassel and Robert Weibel. A review and conceptual framework of automated map generalization. *International Journal of Geographical Information System*, 2(3):229–244, 1988.

[83] Tinghua Ai and Peter van Oosterom. A map generalization model based on algebra mapping transformation. In *Proceedings of the 9th ACM international symposium on Advances in geographic information systems*, pages 21–27, 2001.

[84] Y Kang, J Rao, W Wang, B Peng, S Gao, and F Zhang. Towards cartographic knowledge encoding with deep learning: A case study of building generalization. In *Proceedings of the AutoCarto*, 2020.

[85] Xianjin He, Xinchang Zhang, and Qinchuan Xin. Recognition of building group patterns in topographic maps based on graph partitioning and random forest. *ISPRS Journal of Photogrammetry and Remote Sensing*, 136:26–40, 2018.

[86] Xiongfeng Yan, Tinghua Ai, Min Yang, and Hongmei Yin. A graph convolutional neural network for classification of building patterns using spatial vector data. *ISPRS journal of photogrammetry and remote sensing*, 150:259–273, 2019.

[87] Gengchen Mai, Chiyu Jiang, Weiwei Sun, Rui Zhu, Yao Xuan, Ling Cai, Krzysztof Janowicz, Stefano Ermon, and Ni Lao. Towards general-purpose representation learning of polygonal geometries. *GeoInformatica*, 27(2):289–340, 2023.

[88] Huafei Yu, Tinghua Ai, Min Yang, Lina Huang, and Jiaming Yuan. A recognition method for drainage patterns using a graph convolutional network. *International Journal of Applied Earth Observation and Geoinformation*, 107:102696, 2022.

[89] Gengchen Mai, Krzysztof Janowicz, Bo Yan, Rui Zhu, Ling Cai, and Ni Lao. Multi-scale representation learning for spatial feature distributions using grid cells. In *International Conference on Learning Representations*, 2020.

[90] Gengchen Mai, Ni Lao, Yutong He, Jiaming Song, and Stefano Ermon. Csp: Self-supervised contrastive spatial pre-training for geospatial-visual representations. In *the Fortieth International Conference on Machine Learning (ICML 2023)*, 2023.

[91] Gengchen Mai, Krzysztof Janowicz, Yingjie Hu, Song Gao, Bo Yan, Rui Zhu, Ling Cai, and Ni Lao. A review of location encoding for geoai: methods and applications. *International Journal of Geographical Information Science*, 36(4):639–673, 2022.

[92] Gengchen Mai, Chris Cundy, Kristy Choi, Yingjie Hu, Ni Lao, and Stefano Ermon. Towards a foundation model for geospatial artificial intelligence (vision paper). In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems*, pages 1–4, 2022.

[93] Gengchen Mai, Weiming Huang, Jin Sun, Suhang Song, Deepak Mishra, Ninghao Liu, Song Gao, Tianming Liu, Gao Cong, Yingjie Hu, et al. On the opportunities and challenges of foundation models for geospatial artificial intelligence. *arXiv preprint arXiv:2304.06798*, 2023.

[94] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.

[95] Basel Shbita, Craig A Knoblock, Weiwei Duan, Yao-Yi Chiang, Johannes H Uhl, and Stefan Leyk. Building spatio-temporal knowledge graphs from vectorized topographic historical maps. *Semantic Web*, (Preprint):1–23, 2023.

[96] Jina Kim, Zekun Li, Yijun Lin, Min Namgung, Leeje Jang, and Yao-Yi Chiang. The mapkurator system: A complete pipeline for extracting and linking text from historical maps. *arXiv preprint arXiv:2306.17059*, 2023.

[97] Gengchen Mai, Weiming Huang, Ling Cai, Rui Zhu, and Ni Lao. Narrative cartography with knowledge graphs. *Journal of Geovisualization and Spatial Analysis*, 6(1):4, 2022.

[98] Yuhao Kang, Qianheng Zhang, and Robert Roth. The ethics of ai-generated maps: A study of dalle 2 and implications for cartography. *arXiv preprint arXiv:2304.10743*, 2023.

[99] Phillip Fernberg and Brent Chamberlain. Artificial intelligence in landscape architecture: A literature review. *Landscape Journal*, 42(1):13–35, 2023.

[100] Marika Li. *Designer Robots: An early look at applications for Artificial Intelligence Visualization Software in Landscape Architecture*. PhD thesis, University of Guelph, 2023.

[101] Sachith Seneviratne, Damith Senanayake, Sanka Rasnayaka, Rajith Vidanaarachchi, and Jason Thompson. Dalle-urban: Capturing the urban design expertise of large text to image transformers. In *2022 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–9. IEEE, 2022.

[102] Thomas W Sanchez, Hannah Shumway, Trey Gordner, and Theo Lim. The prospects of artificial intelligence in urban planning. *International Journal of Urban Sciences*, 27(2):179–194, 2023.

[103] Joern Ploennigs and Markus Berger. Ai art in architecture. *AI in Civil Engineering*, 2(1):8, 2023.

[104] Stanislas Chaillou. Ai and architecture: An experimental perspective. In *The Routledge Companion to Artificial Intelligence in Architecture*, pages 420–441. Routledge, 2021.

[105] Ziming He, Xiaomei Li, Ling Fan, and Harry Jiannan Wang. Revamping interior design workflow through generative artificial intelligence. In *International Conference on Human-Computer Interaction*, pages 607–613. Springer, 2023.

[106] GHADA KHALED HUSSEIN et al. Improving design efficiency using artificial intelligence: A study on the role of artificial intelligence in streamlining the interior design process. *International Design Journal*, 13(5):255–270, 2023.

[107] Rafael Iván Pazos Pérez. *Blurring the boundaries between real and artificial in architecture and urban design through the use of artificial intelligence*. PhD thesis, Universidade da Coruña, 2017.

[108] Mathias Bank Stigsen, Alexandra Moisi, Shervin Rasoulzadeh, Kristina Schinegger, and Stefan Rutzinger. Ai diffusion as design vocabulary.

[109] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022.

[110] Yuanzhi Zhu, Kai Zhang, Jingyun Liang, Jiezhang Cao, Bihan Wen, Radu Timofte, and Luc Van Gool. Denoising diffusion models for plug-and-play image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1219–1229, 2023.

[111] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726, 2022.

[112] Gengchen Mai, Ni Lao, Weiwei Sun, Yuchi Ma, Jiaming Song, Chenlin Meng, Hongxu Ma, Jinmeng Rao, Ziyuan Li, and Stefano Ermon. Ssif: Learning continuous image representation for spatial-spectral super-resolution. *arXiv preprint arXiv:2310.00413*, 2023.

[113] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

[114] Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim, and Daniel Aliaga. Objectstitch: Generative object compositing. *arXiv preprint arXiv:2212.00932*, 2022.

[115] Jaskirat Singh, Liang Zheng, Cameron Smith, and Jose Echevarria. Paint2pix: interactive painting based progressive image synthesis and editing. In *European Conference on Computer Vision*, pages 678–695. Springer, 2022.

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

[117] Bahaa Mustafa. The impact of artificial intelligence on the graphic design industry. *resmilitaris*, 13(3):243–255, 2023.

[118] Yue Lu, Chao Guo, Yong Dou, Xingyuan Dai, and Fei-Yue Wang. Could chatgpt imagine: Content control for artistic painting generation via large language models. *Journal of Intelligent & Robotic Systems*, 109(2):1–15, 2023.

[119] Benjamin Matthews, Barrie Shannon, and Mark Roxburgh. Destroy all humans: The dematerialisation of the designer in an age of automation and its impact on graphic design—a literature review. *International Journal of Art & Design Education*, 42(3):367–383, 2023.

[120] Ye Yuan, Wuyang Chen, Zhaowen Wang, Matthew Fisher, Zhifei Zhang, Zhangyang Wang, and Hailin Jin. Font completion and manipulation by cycling between multi-modality representations. *arXiv preprint arXiv:2108.12965*, 2021.

[121] Ye Yuan, Yasuaki Ito, and Koji Nakano. Art font image generation with conditional generative adversarial networks. In *2020 Eighth International Symposium on Computing and Networking Workshops (CANDARW)*, pages 151–156. IEEE, 2020.

[122] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. Text2nerf: Text-driven 3d scene generation with neural radiance fields. *arXiv preprint arXiv:2305.11588*, 2023.

[123] Chun-Han Yao, Jimei Yang, Duygu Ceylan, Yi Zhou, Yang Zhou, and Ming-Hsuan Yang. Learning visibility for robust dense human body estimation. In *European Conference on Computer Vision*, pages 412–428. Springer, 2022.

[124] Hung-Cheng Chen and Zhongwen Chen. Using chatgpt and midjourney to generate chinese landscape painting of tang poem 'the difficult road to shu'. *International Journal of Social Sciences*, 3(2).

[125] Juan Romero, Juan J Romero, and Penousal Machado. *The art of artificial evolution: A handbook on evolutionary art and music*. Springer Science & Business Media, 2008.

[126] Agoston E Eiben and James E Smith. *Introduction to evolutionary computing*. Springer, 2015.

[127] Doyeon Kim, Donggyu Joo, and Junmo Kim. Tivgan: Text to image to video generation with step-by-step evolutionary generator. *IEEE Access*, 8:153113–153122, 2020.

[128] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1526–1535, 2018.

[129] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. *arXiv preprint arXiv:1906.02634*, 2019.

[130] Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. Ccvs: context-aware controllable video synthesis. *Advances in Neural Information Processing Systems*, 34:14042–14055, 2021.

[131] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation. In *European conference on computer vision*, pages 720–736. Springer, 2022.

[132] Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic vqgan and time-sensitive transformer. In *European Conference on Computer Vision*, pages 102–118. Springer, 2022.

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

[134] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023.

[135] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023.

[136] Chat gpt-4 in the film industry: Scriptwriting, editing, and more. https://ts2.space/en/chat-gpt-4-in-the-film-industry-scriptwriting-editing-and-more/. Accessed: 2023-09-25.

[137] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and Dahua Lin. Movienet: A holistic dataset for movie understanding. In *The European Conference on Computer Vision (ECCV)*, 2020.

[138] Yu Xiong, Qingqiu Huang, Lingfeng Guo, Hang Zhou, Bolei Zhou, and Dahua Lin. A graph-based framework to bridge movies and synopses. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.

[139] The impact of generative ai on hollywood and entertainment. `https://sloanreview.mit.edu/article/the-impact-of-generative-ai-on-hollywood-and-entertainment/`. Accessed: 2023-09-25.

[140] Qingqiu Huang, Yu Xiong, and Dahua Lin. Unifying identification and context learning for person recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[141] Qingqiu Huang, Wentao Liu, and Dahua Lin. Person search in videos with one portrait through visual and temporal links. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 425–441, 2018.

[142] Jiangyue Xia, Anyi Rao, Linning Xu, Qingqiu Huang, Jiangtao Wen, and Dahua Lin. Online multi-modal person search in videos. In *The European Conference on Computer Vision (ECCV)*, 2020.

[143] Qingqiu Huang, Yuanjun Xiong, Yu Xiong, Yuqi Zhang, and Dahua Lin. From trailers to storylines: An efficient way to learn from movies. *arXiv preprint arXiv:1806.05341*, 2018.

[144] Anyi Rao, Jiaze Wang, Linning Xu, Xuekun Jiang, Qingqiu Huang, Bolei Zhou, and Dahua Lin. A unified framework for shot type classification based on subject centric lens. In *The European Conference on Computer Vision (ECCV)*, 2020.

[145] Anyi Rao, Linning Xu, Yu Xiong, Guodong Xu, Qingqiu Huang, Bolei Zhou, and Dahua Lin. A local-to-global approach to multi-modal movie scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10146–10155, 2020.

[146] The effects of artificial intelligence on media and communication. `https://www.linkedin.com/pulse/effects-artificial-intelligence-media-communication-mayowa-lateef/`. Accessed: 2023-09-25.

[147] Ai and the future of local tv news. Accessed: 2023-09-25.

[148] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.

[149] Bochen Li, Karthik Dinesh, Gaurav Sharma, and Zhiyao Duan. Video-based vibrato detection and analysis for polyphonic string music. In *ISMIR*, pages 123–130, 2017.

[150] Bochen Li, Akira Maezawa, and Zhiyao Duan. Skeleton plays piano: Online generation of pianist body movements from midi performance. In *ISMIR*, pages 218–224, 2018.

[151] Bochen Li, Karthik Dinesh, Zhiyao Duan, and Gaurav Sharma. See and listen: Score-informed association of sound tracks to players in chamber music performance videos. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2906–2910. IEEE, 2017.

[152] Bochen Li, Karthik Dinesh, Chenliang Xu, Gaurav Sharma, and Zhiyan Duan. Online audio-visual source association for chamber music performances. *Transactions of the International Society for Music Information Retrieval*, 2(1), 2019.

[153] Mojtaba Heydari, Ju-Chiang Wang, and Zhiyao Duan. Singnet: a real-time singing voice beat and downbeat tracking system. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[154] Gabriel Sargent, Pierre Hanna, and Henri Nicolas. Segmentation of music video streams in music pieces through audio-visual analysis. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 724–728. IEEE, 2014.

[155] Janice N Killian. The effect of audio, visual and audio-visual performance on perception of musical content. *Bulletin of the Council for Research in Music Education*, pages 77–87, 2001.

[156] Bochen Li, Xinzhao Liu, Karthik Dinesh, Zhiyao Duan, and Gaurav Sharma. Creating a multitrack classical music performance dataset for multimodal music analysis: Challenges, insights, and applications. *IEEE Transactions on Multimedia*, 21(2):522–535, 2018.

[157] William Forde Thompson, Frank A Russo, and Lena Quinto. Audio-visual integration of emotional cues in song. *Cognition and Emotion*, 22(8):1457–1470, 2008.

[158] Zhiyao Duan, Slim Essid, Cynthia C. S. Liem, Gaël Richard, and Gaurav Sharma. Audio-visual analysis of music performances. *IEEE Signal Processing Magazine*, 36(1):63–73, 2019.

[159] Lucas Bellaiche, Rohin Shahi, Martin Harry Turpin, Anya Ragnhildstveit, Shawn Sprockett, Nathaniel Barr, Alexander Christensen, and Paul Seli. Humans versus ai: whether and why we prefer human-created compared to ai-created artwork. *Cognitive Research: Principles and Implications*, 8(1):1–22, 2023.

[160] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.

[161] Ali Borji. A categorical archive of chatgpt failures. *arXiv preprint arXiv:2302.03494*, 2023.

[162] Why does ai art screw up hands and fingers? https://www.britannica.com/topic/Why-does-AI-art-screw-up-hands-and-fingers-2230501. Accessed: 2023-09-25.

[163] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.

[164] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.

[165] Ben Goodrich, Vinay Rao, Peter J Liu, and Mohammad Saleh. Assessing the factual accuracy of generated text. In *proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 166–175, 2019.

[166] Owain Evans, Owen Cotton-Barratt, Lukas Finnveden, Adam Bales, Avital Balwit, Peter Wills, Luca Righetti, and William Saunders. Truthful ai: Developing and governing ai that does not lie. *arXiv preprint arXiv:2110.06674*, 2021.

[167] Danhuai Guo, Huixuan Chen, Ruoling Wu, and Yangang Wang. Aigc challenges and opportunities related to public safety: A case study of chatgpt. *Journal of Safety Science and Resilience*, 2023.

[168] Jinmeng Rao, Song Gao, Gengchen Mai, and Krzysztof Janowicz. Building privacy-preserving and secure geospatial artificial intelligence foundation models. *arXiv preprint arXiv:2309.17319*, 2023.

[169] Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G Parker, and Munmun De Choudhury. Synthetic lies: Understanding ai-generated misinformation and evaluating algorithmic and human solutions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–20, 2023.

[170] Ronak Agrawal and Dilip Kumar Sharma. A survey on video-based fake news detection techniques. In *2021 8th International Conference on Computing for Sustainable Global Development (INDIACom)*, pages 663–669. IEEE, 2021.

[171] Popular image generators accept 85 [https://aibusiness.com/responsible-ai/popular-image-generators-accept-85-of-fake-news-prompts](https://aibusiness.com/responsible-ai/popular-image-generators-accept-85-of-fake-news-prompts). Accessed: 2023-09-25.

[172] Fake trump arrest photos: How to spot an ai-generated image. [https://www.bbc.com/news/world-us-canada-65069316](https://www.bbc.com/news/world-us-canada-65069316). Accessed: 2023-09-25.

[173] Pictures of joe biden and vp celebrating trump indictment are fake. [https://factcheck.afp.com/doc.afp.com.33CK93W](https://factcheck.afp.com/doc.afp.com.33CK93W). Accessed: 2023-09-25.

[174] Deepfake scams have arrived: Fake videos spread on facebook, tiktok and youtube. [https://www.nbcnews.com/tech/tech-news/deepfake-scams-arrived-fake-videos-spread-facebook-tiktok-youtube-rcna101415](https://www.nbcnews.com/tech/tech-news/deepfake-scams-arrived-fake-videos-spread-facebook-tiktok-youtube-rcna101415). Accessed: 2023-09-25.

[175] Patrick Tucker. The newest ai-enabled weapon: Deep-faking photos of the earth. *Defense One*, (13):03, 2019.

[176] Bo Zhao, Shaozeng Zhang, Chunxue Xu, Yifan Sun, and Chengbin Deng. Deep fake geography? when geospatial data encounter artificial intelligence. *Cartography and Geographic Information Science*, 48(4):338–352, 2021.

[177] Deanna D Caputo, Shari Lawrence Pfleeger, Jesse D Freeman, and M Eric Johnson. Going spear phishing: Exploring embedded training and awareness. *IEEE security & privacy*, 12(1):28–38, 2013.

[178] Julian Hazell. Large language models can be used to effectively scale spear phishing campaigns. *arXiv preprint arXiv:2305.06972*, 2023.

[179] Sowjanya Manyam. Artificial intelligence's impact on social engineering attacks. 2022.

[180] Sayak Saha Roy, Krishna Vamsi Naragam, and Shirin Nilizadeh. Generating phishing attacks using chatgpt. *arXiv preprint arXiv:2305.05133*, 2023.

[181] Xiao Fang, Shangkun Che, Minjia Mao, Hongzhe Zhang, Ming Zhao, and Xiaohang Zhao. Bias of ai-generated content: An examination of news produced by large language models. *arXiv preprint arXiv:2309.09825*, 2023.

[182] Overcoming algorithmic gender bias in ai-generated marketing content. [https://www.forbes.com/sites/forbescommunicationscouncil/2023/07/25/overcoming-algorithmic-gender-bias-in-ai-generated-marketing-content/?sh=518d4fa11639](https://www.forbes.com/sites/forbescommunicationscouncil/2023/07/25/overcoming-algorithmic-gender-bias-in-ai-generated-marketing-content/?sh=518d4fa11639). Accessed: 2023-09-25.

[183] Rohin Manvi, Samar Khanna, Gengchen Mai, Marshall Burke, David Lobell, and Stefano Ermon. Geollm: Extracting geospatial knowledge from large language models. *arXiv preprint arXiv:2310.06213*, 2023.

[184] Fahim Faisal and Antonios Anastasopoulos. Geographic and geopolitical biases of language models. *arXiv preprint* [arXiv:2212.10408](arXiv:2212.10408), 2022.

[185] Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. The science of detecting llm-generated texts. *arXiv preprint* [arXiv:2303.07205](arXiv:2303.07205), 2023.

[186] Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. Automatic detection of generated text is easiest when humans are fooled. *arXiv preprint* [arXiv:1911.00650](arXiv:1911.00650), 2019.

[187] Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 512–515, 2017.

[188] Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. Semeval-2019 task 6: Identifying and categorizing offensive language in social media (offenseval). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, 2019.

[189] Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. *arXiv preprint* [arXiv:2304.05335](arXiv:2304.05335), 2023.

[190] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv e-prints*, pages arXiv–2009, 2020.

[191] Helen Ngo, Cooper Raterink, João GM Araújo, Ivan Zhang, Carol Chen, Adrien Morisot, and Nicholas Frosst. Mitigating harm in language models with conditional-likelihood filtration. *arXiv preprint* [arXiv:2108.07790](arXiv:2108.07790), 2021.

[192] Shrimai Prabhumoye, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Adding instructions during pretraining: Effective way of controlling toxicity in language models. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2628–2643, 2023.

[193] Boxin Wang, Wei Ping, Chaowei Xiao, Peng Xu, Mostofa Patwary, Mohammad Shoeybi, Bo Li, Anima Anandkumar, and Bryan Catanzaro. Exploring the limits of domain-adaptive training for detoxifying large-scale language models. *Advances in Neural Information Processing Systems*, 35:35811–35824, 2022.

[194] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[195] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, 2020.

[196] Xinlei He, Savvas Zannettou, Yun Shen, and Yang Zhang. You only prompt once: On the capabilities of prompt learning on large language models to tackle toxic content. *arXiv e-prints*, pages arXiv–2308, 2023.

[197] Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*, 2019.

[198] Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A Smith, and Yejin Choi. Dexperts: Decoding-time controlled text generation with experts and anti-experts. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6691–6706, 2021.