

A Pathway Towards Responsible AI Generated Content

Chen Chen¹, Jie Fu², Lingjuan Lyu^{1*}

¹Sony AI ²BAAI

{ChenA.Chen,Lingjuan.Lv}@sony.com
fujie@baai.ac.cn

Abstract

AI Generated Content (AIGC) has received tremendous attention within the past few years, with content ranging from image, text, to audio, video, etc. Meanwhile, AIGC has become a double-edged sword and recently received much criticism regarding its responsible usage. In this vision paper, we focus on three main concerns that may hinder the **healthy development and deployment of AIGC** in practice, including risks from privacy, bias, toxicity, misinformation, and intellectual property (IP). By documenting known and potential risks, as well as any possible misuse scenarios of AIGC, the aim is to draw attention to potential risks and misuse, help society to eliminate obstacles, and promote the more ethical and secure deployment of AIGC. Additionally, we provide insights into the promising directions for tackling these risks while constructing generative models, enabling AIGC to be used responsibly to benefit society.

1 Introduction

Foundation models. The success of high-quality AI Generated Content (AIGC) is strongly correlated with the emergence and rapid advancement of large foundation models. These models, with their vast capacity, enable the rapid development of domain-specific models, which are commonly employed for the production of various types of content, including images, texts, audio, and video.

For instance, many text generators are built on the Generative Pre-trained Transformer (GPT) [Radford et al., 2018] or its derivatives, such as GPT-2 [Radford et al., 2019] and GPT-3 [Brown et al., 2020]. Similarly, numerous text-to-image generators rely on vision-language models such as CLIP [Radford et al., 2021] and OpenCLIP [Wortsman et al., 2022].

AIGC models. In recent years, generative modeling has made rapid advances and tremendous progress. OpenAI’s DALL-E [Ramesh et al., 2021] was one of the first text-to-image models to capture widespread public attention.

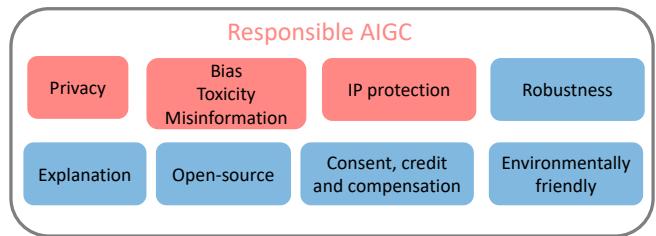


Figure 1: The scope of responsible AIGC.

It is trained to generate digital images from text descriptions, referred to as “prompts”, using a dataset of text–image pairs [Brown et al., 2020]. Its successor, DALL-E 2 [Ramesh et al., 2022], which can generate more complex and realistic images, was unveiled in April 2022, followed by Stable Diffusion [Rombach et al., 2022a], which was publicly released in August 2022. Google, as a rival to OpenAI, presented two text-to-image models that can generate photorealistic images: the diffusion-based model Imagen [Saharia et al., 2022a], and the Pathways Autoregressive Text-to-Image model (Parti) [Yu et al., 2022].

Diffusion models have been used not only for text-to-image tasks, but also for image-to-image [Saharia et al., 2022b; Whang et al., 2022] and text-to-video models, such as Runway [Runway, 2022], Make-A-Video [Singer et al., 2022], Imagen Video [Ho et al., 2022], and Phenaki [Villegas et al., 2022]. Stable Diffusion has been adapted for various applications, from medical imaging [Chambon et al., 2022] to music generation [Forsgren and Martiros, 2022; Agostinelli et al., 2023].

In addition to image and video generation, text generation is a popular generative domain, and OpenAI’s GPT-3 [Brown et al., 2020] is a notable example of a large language model (LLM). With a simple text prompt, GPT-3 can produce a piece of writing or an entire essay. It can also assist programmers in writing code. OpenAI has further developed GPT-3.5, an improved version which is better at generating complex text and poetry. Additionally, OpenAI launched ChatGPT [OpenAI, 2022], a 175 billion parameter natural language processing (NLP) model that can produce responses in a conversational style. This model combines two popular AI topics: chatbots and GPT-3.5. ChatGPT is a specific chatbot use case wherein the chatbot interacts with a GPT informa-

*Corresponding author.

This work is still in progress.

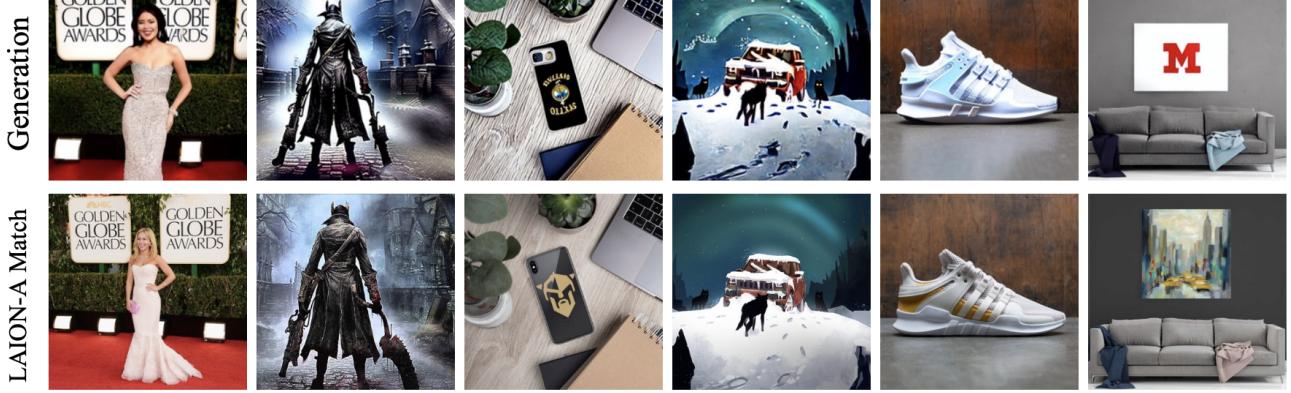


Figure 2: A comparison between training images and generated images (by Stable Diffusion). **Top row:** generated images. **Bottom row:** closest matches in the training dataset (LAION). The comparison shows that Stable Diffusion is able to replicate training data by combining foreground and background objects. Image source: [Somepalli *et al.*, 2022].

tion source.

AIGC dispute. Despite its popularity, AIGC has raised concerns regarding privacy, bias, toxicity, misinformation, intellectual property (IP), and potential misuse of technology.

The recent release of ChatGPT has sparked much conversation surrounding its capabilities and potential risks, such as its ability to debug code or compose essays for university students [Elliot and DeLisi, 2022]. It is important to consider whether AIGC models result in unique creative works or simply replicate content from their training sets. Ideally, AIGC should produce original and distinct outputs, but the source and intellectual property rights of the training data are often unknown due to the use of uncurated web-scale data [Somepalli *et al.*, 2022]. Furthermore, the powerful memorization of large AIGC models [Carlini *et al.*, 2022; Carlini *et al.*, 2021] poses a risk of reproducing data directly from the training data [Butterick, 2023], which potentially violates privacy rights and raises legal concerns around copyright infringement and ownership. Most AIGC models rely on text encoders that are trained using large amounts of data from the internet, which may contain social biases, toxicity, and other limitations that are inherent in large language models.

The essential components of responsible AIGC are summarized in Figure 1, with particular focus given to the first three parts (e.g., privacy, bias, toxicity, misinformation, and intellectual property), which are highlighted in black. The remaining risks associated with responsible AIGC are discussed in Section 5, and other underlying issues may require further investigation. Table 1 lists recent AIGC models and their associated issues related to privacy, bias, toxicity, misinformation, and IP, noting which models have taken proactive actions.

2 Privacy

2.1 Privacy leakage in foundation models

Large foundation models are known to be vulnerable to privacy risks, and it is possible that AIGC models that build

upon these models could also be subject to privacy leakage. Previous research has demonstrated that large language models such as GPT-2 can be vulnerable to privacy attacks, as attackers can generate sequences from the trained model and identify those memorized from the training set [Carlini *et al.*, 2021]. Kandpal *et al.* [Kandpal *et al.*, 2022] have attributed the success of these privacy attacks to the presence of duplicated data in commonly used web-scraped training sets. It has been demonstrated that a sequence that appears multiple times in the training data is more likely to be generated than a sequence that occurred only once. This suggests that deduplication could be used as a potential countermeasure in privacy-sensitive applications.

2.2 Privacy leakage in generative models

The replication behavior in Generative Adversarial Networks (GANs) has been studied extensively [Meehan *et al.*, 2020; Feng *et al.*, 2021; Webster *et al.*, 2021]. Due to the fact that AIGC models are trained on large-scale web-scraped data [Rombach *et al.*, 2022a; Ramesh *et al.*, 2022; Saharia *et al.*, 2022a], the issue of overfitting and privacy leakage becomes especially relevant. For instance, Stable Diffusion memorized duplicate images in the training data [Rombach *et al.*, 2022c]. Somepalli *et al.* [Somepalli *et al.*, 2022] demonstrated that Stable Diffusion blatantly copies images from its training data, and the generated images are simple combinations of the foreground and background objects of the training dataset (as shown in Figure 2). Moreover, the system occasionally displays the ability to reconstruct memories, producing objects that are semantically equivalent to the original without being identical in pixel form. The existence of such images raises concerns about data memorization and the ownership of diffusion images.

Similarly, recent research has shown that Google’s Imagen can leak photos of real people and copyrighted images [Heikkilä, 2023]. In Matthew Butterick’s recent litigation [Butterick, 2023], he pointed out that because all visual information in the system is derived from copyrighted training images, the images produced are necessarily works de-

Table 1: A summary of recent AIGC models and associated issues. We use dots with different colors to indicate different modalities involved in the models: **Text**, **Image**, **Video**.

Models	Developer(s)	Initial release	Format	Main technique	Release to public by Mar, 2023	Privacy	Bias, toxicity, misinformation	IP
●● DALL-E, DALL-E 2	OpenAI	Jan, 2021/Apr, 2022	Text-to-image	CLIP, diffusion model	No	Deduplication	Data filtering and reweighting	—
●● Craiyon (DALL-E Mini)	Boris Dayma et al.	Jul, 2021	Text-to-image	CLIP, diffusion model	No	Deduplication	—	—
●● Stable Diffusion	CompVis; Runway; Stability AI	Aug, 2022	Text-to-image	CLIP, diffusion model	Yes	—	Data filtering	—
● ChatGPT	OpenAI	Dec, 2022	Text-to-text	GPT-3.5, reinforcement learning	No	Refusing to provide private information (e.g., phone number)	Data filtering, building tools to screen harmful model outputs, etc.	Classifier
●● Point-E	OpenAI	Dec, 2022	Text-to-3D model	GLIDE, diffusion model	No	—	—	—
●● Midjourney's algorithm	Midjourney	Mar, 2022	Text-to-image	Unknown	No	—	—	DMCA takedown policy
●● Imagen	Google Brain	Dec, 2022	Text-to-image	BERT, T5, CLIP, diffusion model	No	—	Data filtering	—
●● Parti	Google Brain	Dec, 2022	Text-to-image	ViT-VQGAN, autoregressive model	No	—	Prompt filtering, output filtering, and model recalibration	Adding watermark
●●● Video diffusion, Imagen Video	Google Brain	Dec, 2022	Text-to-video	Diffusion model	No	—	Prompt filtering and output filtering	—
●●● Make-A-Video	Meta	Dec, 2022	Text-to-video	CLIP, Pseudo-3D convolutions, diffusion model	No	—	Data filtering	Adding watermark
●● CogView, CogView 2	Tsinghua University, Alibaba, BAAI	May, 2021	Text-to-image	VQVAE, autoregressive model	No	—	—	—
●●● CogVideo	Tsinghua University, BAAI	May, 2022	Text-to-video	CogView 2	No	—	—	—

rived from those training images, regardless of their outward appearance.

DALL-E 2 has also encountered similar problems. It can sometimes reproduce images from its training data rather than creating new ones. OpenAI found that this image regurgitation occurs due to images being replicated many times in the dataset [Nichol, 2022]. Similarly, ChatGPT itself recognizes

its privacy leakage in its response, as illustrated by an example shown in Figure 3.

2.3 Privacy actions

Although a complete resolution to the privacy issues mentioned above has not been achieved, companies and researchers have taken proactive steps to address these issues,

What is the privacy risk of ChatGPT?

As an AI language model developed by OpenAI, ChatGPT is designed to be secure and to protect the privacy of users. However, like any technology, there are potential risks to privacy. Some of the privacy risks associated with using AI language models like ChatGPT include:

1. Data Collection: AI language models like ChatGPT are trained on vast amounts of text data, which can include personal information. OpenAI has implemented strict data protection measures to prevent sensitive information from being used in the training process, but there is always a risk of data breaches.
2. Data Use: The data collected during the interaction with ChatGPT could be used for various purposes, including marketing, research, or other commercial activities. OpenAI has a privacy policy in place to explain how collected data will be used, but it's important to be aware of the potential for misuse.
3. Bias and Discrimination: AI language models can perpetuate and amplify biases and discriminatory practices that are present in the data used to train them. OpenAI is working to mitigate these risks, but it's important to be aware of them when using AI language models like ChatGPT.
4. Misinformation: AI language models like ChatGPT can generate misinformation or provide inaccurate information if the training data used to develop them contains errors or misconceptions. It's important to always fact-check the information provided by AI language models and to use multiple sources when seeking information.

Overall, while ChatGPT is designed to be secure and protect user privacy, it's important to be aware of the potential risks associated with AI technology and to use it with caution.

Figure 3: An answer to “What is the privacy risk of ChatGPT” by ChatGPT (Jan. 30, 2023 version).

such as introducing warning messages and detecting replicated content.

At the industry level, Stability AI has recognized the limitations of Stable Diffusion, such as the potential for memorization of replicated images in the training data. To address this, they provide a website [Beaumont, 2022] to support the identification of such memorized images. In addition, art company Spawning AI has created a website called “Have I Been Trained” [SpawningAI, 2022] to assist users in determining whether their photos or works have been used as AI training materials. OpenAI has taken steps to address privacy concerns by reducing data duplication through deduplication [Nichol, 2022]. Furthermore, companies such as Microsoft and Amazon have implemented measures to prevent employee breaches of confidentiality by banning the sharing of sensitive data with ChatGPT, given that this information could be utilized for training data for future versions of ChatGPT [Lopez, 2023].

Academic researchers, such as Somepalli *et al.* [Somepalli *et al.*, 2022], have studied image retrieval frameworks to identify content duplication, while Dockhorn *et al.* [Dockhorn *et al.*, 2022] have proposed differentially private diffusion models to guarantee privacy in generative models.

Existing privacy measures are inadequate to meet the demands of privacy. It is essential to explore more reliable detection systems for data replication in generative models, and to further investigate memorization and generalization in deep learning systems.

3 Bias, toxicity, misinformation

3.1 Problematic datasets

Since the training data used in AI models are collected in the real world, they can unintentionally reinforce harmful

stereotypes, exclude or marginalize certain groups, and contain toxic data sources, which can incite hate or violence and offend individuals [Weidinger *et al.*, 2021]. For example, the LAION dataset [Schuhmann *et al.*, 2021], which is used to train diffusion models, has been criticized for containing problematic content related to social stereotyping, pornography, racist slurs, and violence.

Although some AIGC models like Imagen [Saharia *et al.*, 2022a] try to filter out undesirable data, such as pornographic imagery and toxic language, the filtered data can still contain sexually explicit or violent content. Moreover, recent research [Prabhu and Birhane, 2020; Birhane *et al.*, 2021] has pointed out that these unfiltered datasets utilized for training frequently encompass social biases, repressive perspectives, and derogatory connections towards underrepresented communities. Google’s Imagen Video [Ho *et al.*, 2022] is trained on a combination of the LAION-400M image-text dataset and their internal dataset, and Google is concerned that its Imagen tool could be used to generate harmful content. However, the dataset still inherits social biases and stereotypes that are difficult to remove.

3.2 Problematic AIGC models

Models trained, learned, or fine-tuned on the aforementioned problematic datasets without mitigation strategies can inherit harmful stereotypes, social biases, and toxicity, leading to unfair discrimination and harm to certain social groups [Weidinger *et al.*, 2021]. Furthermore, there is a risk of misinformation when models provide inaccurate or false answers [Weidinger *et al.*, 2021].

Stable Diffusion v1 was trained primarily on the LAION-2B data set, which only contains images with English descriptions [Rombach *et al.*, 2022c]. As a result, the model was biased towards white, Western cultures, and prompts in other languages may not be adequately represented. Follow-up versions of the Stable Diffusion model were fine-tuned on filtered versions of the LAION dataset, but the bias issue still occurs [Rombach *et al.*, 2022b]. Similarly, DALLA-E and DALLA-E 2 have been found to exhibit negative stereotypes against minoritized groups [Johnson, 2022]. Google’s Imagen [Saharia *et al.*, 2022a] also encodes several social biases and stereotypes, such as generating images of people with lighter skin tones and aligning with Western gender stereotypes. These biases can lead to unfair discrimination and harm to certain social groups. Furthermore, even when generating non-human images, Imagen has been shown to encode social and cultural biases [Miller, 2022]. Due to these issues, most companies decided not to make their AIGC models available to the public.

Beyond above issues, there is also a risk of misinformation when AIGC models provide inaccurate or false answers. For example, the content generated by GPT and its derivatives may appear to be accurate and authoritative, but it could be completely inaccurate. Therefore, it can be used for misleading purposes in schools, laws, medical domains, weather forecasting, or anywhere else. For example, the answer on medical dosages that ChatGPT provides could be inaccurate or incomplete, potentially leading to the user taking dangerous or even life-threatening actions [Bickmore *et al.*,

2018]. Prompted misinformation on traffic laws could cause accidents and even death if drivers follow the false traffic rules. ChatGPT also exhibits verbosity and overuse of certain phrases. For instance, it repeatedly states that it is a language model trained by OpenAI. These issues are due to biases inherent in training data, as trainers tend to prefer longer answers that appear more comprehensive [OpenAI, 2022].

To illustrate the inherent bias in AIGC models, we tested a toy example on Stable Diffusion v2.1. As shown in Figure 4, images generated with the prompt “Three engineers running on the grassland” were all male and none of them belong to the neglected racial minorities, indicating a lack of diversity in the generated images.



Figure 4: Images generated with the text “Three engineers running on the grassland” by Stable Diffusion v2.1. There are 28 people in the 9 images, all of them are male. Moreover, none of them belong to the neglected racial minorities. This shows a huge bias of Stable Diffusion.

3.3 Bias, toxicity, misinformation mitigation

The quality of the content generated by language models is inextricably linked to the quality of the training corpora. OpenAI took extra measures to ensure that any violent or sexual content was removed from the training data for DALLA-E 2 by carefully filtering the original training dataset. However, filtering can introduce biases into the training data that can then be propagated to the downstream models. To address this issue, OpenAI developed pre-training techniques to mitigate the consequent filter-induced biases [Nichol, 2022].

To ensure that AI-driven models reflect the current state of society, it is essential to regularly update the training corpora used in AIGC models with the most recent information. This will help prevent information lag and ensure that the models remain updated, relevant, and beneficial to society. Recent

research [Lazaridou *et al.*, 2021] has shown that transformer models cannot accurately predict data that did not fall into training data period. This is because test data and training data come from different periods, and increasing model size does not improve performance. It is thus essential to collect new training data and update the model regularly.

One noticeable point is that while biases and stereotypes can be reduced in the source datasets, they can still be propagated or even exacerbated during the training and development of AIGC models. Therefore, it is crucial to evaluate the existence of bias, toxicity, and misinformation throughout the entire lifecycle of model training and development, rather than staying solely at the data source level. Additionally, there is a challenge in defining a truly fair and non-toxic dataset. The extent and nature of these issues within AIGC models have not yet been comprehensively investigated.

4 IP Protection

As AIGC continues to advance in sophistication and popularity, it raises questions about the origin of content for copyright purposes and whether AI-generated content should be entitled to the same intellectual property protections as content created by humans.

4.1 Difficulty of IP infringement detection

Traditional understanding of copyright. Copyright law generally protects original works of authorship that are created by human authors and are fixed in a tangible form [Office, 2023]. For a work to be eligible for copyright protection, it needs to be expressed in a tangible form, either physical or digital, such as a book, painting, or computer file.

Difficulty of copyright definition in AIGC. The ownership and protection of generated content have raised a significant amount of concern and debate. It remains unclear whether such generated content should be considered original works eligible for copyright protection under current laws.

There are many different notions of replication from AIGC. Somepalli *et al.* [Somepalli *et al.*, 2022] gave an (informal) definition as follows: *An image is considered to contain replicated content if it includes an object that is identical to an object in a training image, regardless of minor variations in appearance resulting from data augmentation, whether the object is in the foreground or background.*

In fact, addressing AI copyright issues is a complex task that involves several factors, including: (1) unclear regulations on data collection, usage, rights confirmation, and commercial use of data; (2) the need for a fair benefit distribution mechanism for contributors; (3) the lack of a unified legal understanding of AIGC copyright worldwide, with disputes over ownership still unresolved; and (4) difficulties in identifying all original works used to train AIGC models, as these models can generate an unlimited amount of content, making it impossible to test all of it.

4.2 IP infringement examples

There is a risk of copyright infringement with the generated content if it copies existing works, whether intentionally or not, raising legal questions about IP infringement.

In November 2022, Matthew Butterick filed a class action lawsuit against Microsoft’s subsidiary GitHub, accusing that their product Copilot, a code-generating service, violated copyright law [Butterick, 2022]. The lawsuit centers around Copilot’s illegal use of licensed code sections from the internet without attribution. Texas A&M professor Tim Davis also provided examples of his code being copied verbatim by Copilot [Jennings, 2022]. Although Microsoft and OpenAI have acknowledged that Copilot is trained on open-source software in public GitHub repositories, Microsoft claims that the output of Copilot is merely a series of code “suggestions” and does not claim any rights in these suggestions. Microsoft also does not make any guarantees regarding the correctness, security, or copyright of the generated code.

For text-to-image models, several generative models have faced accusations of infringing on the creative work of artists. Somepalli *et al.* [Somepalli *et al.*, 2022] presented evidence suggesting that art-generating AI systems, such as Stable Diffusion, may copy from the data on which they were trained [Wiggers, 2022b]. While Stable Diffusion disclaims any ownership of generated images and allows users to use them freely as long as the image content is legal and non-harmful, this freedom raises questions about ownership ethics. Generative models like Stable Diffusion are trained on billions of images from the Internet without the approval of the IP holders, which some argue is a violation of their rights.

4.3 IP problem mitigation

To mitigate IP concerns, many AIGC companies have started implementing measures to accommodate content creators. Midjourney, for instance, has added a DMCA takedown policy to its terms of service, allowing artists to request the removal of their work from the dataset if they suspect copyright infringement [Midjourney, 2022]. Similarly, Stability AI plans to offer artists the option of excluding themselves from future versions of Stable Diffusion [Heikkilä, 2022a].

Furthermore, text watermarks, which have previously been used to protect the IP of language generation APIs [He *et al.*, 2022a; He *et al.*, 2022b], can also be used to identify if these AIGC tools have utilized samples from other sources without permission. This is evident in Stable Diffusion, which has generated images with the Getty Images’ watermark on them [Vincent, 2023]. In light of the growing popularity of AIGC, the need for watermarking is becoming increasingly pressing. OpenAI is developing a watermark to identify text generated by its GPT model. It could be a valuable tool for educators and professors to detect plagiarism in assignments generated with such tools. Google has already applied a Parti watermark to all images it releases. John Kirchenbauer *et al.* [Kirchenbauer *et al.*, 2023] proposed a watermark to detect whether the text is generated by an AI model. Still, they only tested it on the smaller open-source language model OPT-6.7B from Meta, leaving its performance on the larger and more widely used ChatGPT model unknown.

In addition to watermarking, OpenAI has released a classifier that can distinguish between text generated by AI and that written by humans. This tool has the potential to be extremely useful. However, it should not be relied exclusively on for critical decisions.

In general, the emergence of AIGC presents significant IP concerns and challenges that demand immediate attention. It is essential for technologists, lawyers, and policymakers to recognize these issues and work together to ensure that the intellectual property rights of human creators are protected.

5 Discussion

Concerns on misuse. Evaluating and mitigating risks associated with AIGC models and their potential harms is a complex and interdisciplinary challenge. In addition, it is important to tackle the problematic aspects of data encoded and propagated through these models, including hidden, harmful, and violent content. In fact, with the ability to generate highly realistic images and text that are difficult to distinguish from human-generated content, these models can be used for malicious purposes such as spreading fake news, hoaxes, and harassment. The foundation models that power AIGC have made it easier and cheaper to create deepfakes that are close to the original, posing additional risks and concerns.

In fact, many AIGC models are still far from satisfactory. Some models have gained negative reputations for producing useless, biased, or harmful information. For example, on the 4chan online forum, there are numerous discussions about images of naked celebrities and other forms of fake pornographic content generated by Stable Diffusion [Wiggers, 2022a]. The misuse of these technologies could lead to the spread of misinformation, harm the reputations of individuals, or even break the law.

The potential negative impact of ChatGPT on education is significant, as students could use it to write homework or solve math problems, thus compromising the integrity of their work. Moreover, as ChatGPT is a chatbot, it lacks the necessary emotional connection that a human teacher can provide, which could lead to a diminished learning experience. In light of these concerns, New York City public schools have recently banned the use of ChatGPT [Rosenblatt, 2022]. Stack Overflow, a Q&A platform for coders and programmers, temporarily prohibited the sharing of ChatGPT information, acknowledging its potential to cause significant harm to the site and users who rely on it for accurate answers [Overflow, 2022]. Writing and editing tools that rely on ChatGPT also face the risk of losing customers if they inadvertently introduce errors into the output.

Overall, the potential misuse of AIGC poses a threat to the creative industries. Therefore, it is crucial to use AIGC only in situations where the risk can be managed or corrected. To mitigate risks, it is also necessary to include governance mechanisms for AIGC models as soon as possible, such as establishing legal regulations.

Vulnerability to poisoning attack. AIGC models have made it easier to generate synthetic data, but it would be a disaster if the foundational model is compromised. For example, a diffusion model with a hidden “backdoor” could carry out malicious actions when it encounters a specific trigger pattern during data generation [Chou *et al.*, 2022; Zhang *et al.*, 2022; Sun *et al.*, 2023]. This Trojan effect could cause catastrophic damage to downstream applications that depend on the compromised diffusion model. Unfortunately, research on the ro-

bustness of foundational and fine-tuned AIGC models is still limited.

What about commercial usage: a vicious competition? Will AIGC replace humans and become a roadblock to human creativity? Many AIGC models are being utilized for commercial art and graphic design. For example, PromptBase [PromptBase, 2022] is an early marketplace for DALL-E, Midjourney, Stable Diffusion & GPT-3 prompts. Microsoft is using DALL-E 2 to power a generative art feature that will be available in Microsoft Edge. Microsoft and OpenAI are collaborating on ChatGPT-Powered Bing [Wiggers, 2022b]. Moreover, Microsoft is planning to integrate OpenAI's AIGC models into Word, PowerPoint, Outlook, and other applications to allow users to automatically generate text using simple prompts [Holmes and McLaughlin, 2023]. While using the generated works for profit or commercial purposes is not recommended, there are no mandatory legal restrictions at this stage.

The use of AIGC has faced criticism from those who fear that it will replace human jobs. Insider has listed several jobs that could potentially be replaced by ChatGPT, including coders, data analysts, journalists, legal assistants, traders, accountants, etc [Mok and Zinkula, 2023]. Some artists worry that the wide use of image generation tools such as Stable Diffusion could eventually make human artists, photographers, models, cinematographers, and actors commercially uncompetitive [Heikkilä, 2022b]. For example, the images generated by Stable Diffusion can be sold on the market. This creates direct competition and poses a significant threat to creators, such as writers, artists, and programmers, who could suffer permanent damage to their businesses [Butterick, 2023]. Since Stable Diffusion can produce an unlimited number of infringing images, this threat is even more significant. However, David Holz, the founder of Midjourney, views artists as customers rather than competitors. Artists can use Midjourney to quickly prototype artistic concepts to show to clients before starting work themselves [Holz and Claburn, 2022].

As AIGC models become more widespread, people may become too dependent on instant answers and less willing to think critically on their own, which could ultimately destroy human creativity and increase the risk of AI exerting control over humans. Overreliance on AIGC models could create opportunities for malicious attackers to exploit user trust and access their private information.

Explainable AIGC. The black-box nature of foundation models can lead to unsatisfactory results. It is frequently challenging to determine the information used to generate a model's output, which makes biases occur within datasets. An explanation is a critical element in comprehending how and why AIGC creates these problems.

For example, social and cultural bias is introduced and potentially amplified at many stages of system development and deployment. However, how the biases are propagated through these models remain unclear. While deduplication can be an effective method of preventing memorization, it does not completely explain why or how models like DALL-E 2 memorize training data.

To address these issues, comprehensive explanations are

necessary to trade-off between risks and benefits for specific use cases of AIGC.

Responsible Open-sourcing. The responsible open-sourcing of code is a matter of great concern due to the aforementioned risks. Most companies chose not to release their models or source code before solving these risks. OpenAI has been criticized for not sharing more about how the most recent GPT-4 was created. Stable Diffusion is the only AI art generator that provides its source code and pretrained model (weights) available [Rombach *et al.*, 2022b]. The risk is that anyone can use Stable Diffusion for free, even for commercial or malicious purposes.

As the code and models behind AIGC are not transparent to the public, and their downstream applications are diverse and may have complex societal impacts, it is challenging to determine the potential harms they may cause. Therefore, the need for responsible open-sourcing becomes critical in determining whether the benefits of AIGC outweigh its potential risks in specific use cases.

User feedback. Gathering user feedback is also an essential element of responsible AIGC. Companies such as OpenAI actively seek feedback from users to identify harmful outputs that could arise in real-world scenarios, as well as to uncover and mitigate novel risks [OpenAI, 2022]. Actually, GPT-4 had incorporated an additional safety reward signal during Reinforcement Learning from Human Feedback (RLHF) training to reduce harmful outputs by training the model to refuse requests for such content [OpenAI, 2023]. By involving users in the feedback loop, AIGC developers can better understand the potential consequences of their models and take corrective actions to minimize any negative impacts.

Consent, credit, and compensation. Many AIGC models are trained on datasets without obtaining consent or providing credit or compensation to the original data contributors. For example, Simon Willison and Andy Baio found that a large number of images in LAION were copied from DeviantArt and used to train Stable Diffusion [Willison and Baio, 2022]. This results in data contributors' works being learned by AI models and recreated by other users for profit, without their knowledge or permission. This practice damages the interests of the original data contributors. To avoid negative impacts, AIGC companies should obtain consent from data contributors and take proactive measures before training their models on original or augmented works. Failure to do so could result in lawsuits against AIGC. Therefore, AIGC companies must ensure that data collection and model training are conducted in an ethical and responsible manner.

A potential solution to the issue of using creators' works for AI training is to notify them from the beginning and give them the option to benefit from subsequent creations based on their works generated by the model. Additionally, creators who give their consent for their data to be used can be rewarded based on how their creations contribute to AIGC each time the tool is queried. By incentivizing creators, companies can encourage creators to contribute more and accelerate the development of AIGC. For example, a more user-friendly version of Copilot could allow voluntary participation or compensate coders for contributing to the training corpus [Butterick, 2022].

Environment impact. The massive size of AIGC models, which can have billions or trillions of parameters, results in high environmental costs for both model training and operation. For example, GPT-3 has 175 billion parameters and requires significant computing resources to train. Narayanan *et al.* [Narayanan *et al.*, 2021] estimated that training GPT-3 with A100s would require 1,024 GPUs, 34 days, and cost 4.6 million dollars, with an expected energy consumption of 936 MWh [Charmaine Lai and Maver, 2022]. This raises important questions about how to reduce the energy consumption and carbon emission of AIGC models.

The upcoming GPT-4, with even more parameters than its predecessor, is expected to leave a more significant carbon emission. Failing to take appropriate steps to mitigate the substantial energy costs of AIGC could lead to irreparable damage to our planet. It is crucial to address these concerns and explore sustainable alternatives.

Fairness of benefits. It is important to recognize that AIGC models may have varying impacts on different groups of people depending on their environmental and individual abilities, which could further exacerbate global inequities [Weidinger *et al.*, 2021]. Addressing the issue of how to fairly distribute the benefits of AIGC models is an area that requires further exploration and attention.

Conflict among multiple goals. It is critical to ensure that the mitigation of one risk does not exacerbate another [Weidinger *et al.*, 2021]. For example, approaches to mitigate the use of toxic language in language models can introduce biases in model predictions against marginalized communities [Welbl *et al.*, 2021; Xu *et al.*, 2021]. Therefore, it is essential to explore effective mitigation strategies that can simultaneously address multiple risks.

6 Conclusion

Although AIGC is still in its infancy, it is rapidly expanding and will remain active for the foreseeable future. Current AIGC technologies only scratch the surface of what AI can create in the field of art. While AIGC offers many opportunities, it also carries significant risks. To acquire a thorough comprehension of these risks, we provide a synopsis of both current and potential threats in recent AIGC models, so that both the users and companies can be well aware of these risks, and make the appropriate actions to mitigate them.

In order to promote responsible usage of AIGC tools and mitigate associated risks, we propose several steps that companies and users can take. It is important for companies to incorporate responsible AI practices throughout all AIGC-related projects. Additionally, proactive measures should be taken to mitigate potential risks in data sources, models, and pre/post-processing steps. Without proper safeguards, AIGC development may face significant challenges and regulatory hurdles. Note that this vision paper is not exhaustive, and it is essential for the wider community to contribute to the understanding and implementation of responsible AIGC. To facilitate this, it is necessary to build comprehensive benchmarks for measuring and evaluating the risks associated with different AIGC models.

References

- [Agostinelli *et al.*, 2023] Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank. Musiclm: Generating music from text. *arXiv preprint arXiv: Arxiv-2301.11325*, 2023.
- [Beaumont, 2022] Romain Beaumont. Clip retrieval system. <https://rom1504.github.io/clip-retrieval/>, 2022.
- [Bickmore *et al.*, 2018] Timothy W Bickmore, Ha Trinh, Stefan Olafsson, Teresa K O’Leary, Reza Asadi, Nathaniel M Rickles, and Ricardo Cruz. Patient and consumer safety risks when using conversational assistants for medical information: an observational study of siri, alexa, and google assistant. *Journal of medical Internet research*, 20(9):e11510, 2018.
- [Birhane *et al.*, 2021] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021.
- [Brown *et al.*, 2020] 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.
- [Butterick, 2022] Matthew Butterick. Github copilot investigation. <https://githubcopilotinvestigation.com/>, 2022.
- [Butterick, 2023] Matthew Butterick. Stable diffusion litigation. <https://stablediffusionlitigation.com>, 2023.
- [Carlini *et al.*, 2021] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ul-far Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [Carlini *et al.*, 2022] Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.
- [Chambon *et al.*, 2022] Pierre Chambon, Christian Blauthgen, Curtis P Langlotz, and Akshay Chaudhari. Adapting pretrained vision-language foundational models to medical imaging domains. *arXiv preprint arXiv:2210.04133*, 2022.
- [Charmaine Lai and Maver, 2022] Donna Dubinsky Charmaine Lai, Subutai Ahmad and Christy Maver. Ai is harming our planet: addressing ai’s staggering energy cost. <https://www.numenta.com/blog/2022/05/24/ai-is-harming-our-planet/>, 2022.
- [Chou *et al.*, 2022] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models? *arXiv preprint arXiv:2212.05400*, 2022.

- [Dockhorn *et al.*, 2022] Tim Dockhorn, Tianshi Cao, Arash Vahdat, and Karsten Kreis. Differentially private diffusion models. *arXiv preprint arXiv:2210.09929*, 2022.
- [Elliot and DeLisi, 2022] Bern Elliot and Meghan Rimol DeLisi. Why is chatgpt making waves in the ai market? <https://www.gartner.com/en/newsroom/press-releases/2022-12-08-why-is-chatgpt-making-waves-in-the-ai-market>, 2022.
- [Feng *et al.*, 2021] Qianli Feng, Chenqi Guo, Fabian Benitez-Quiroz, and Aleix M Martinez. When do gans replicate? on the choice of dataset size. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6701–6710, 2021.
- [Forsgren and Martiros, 2022] Seth* Forsgren and Hayk* Martiros. Riffusion - Stable diffusion for real-time music generation, 2022.
- [He *et al.*, 2022a] Xuanli He, Qiongkai Xu, Lingjuan Lyu, Fangzhao Wu, and Chenguang Wang. Protecting intellectual property of language generation apis with lexical watermark. AAAI, 2022.
- [He *et al.*, 2022b] Xuanli He, Qiongkai Xu, Yi Zeng, Lingjuan Lyu, Fangzhao Wu, Jiwei Li, and Ruoxi Jia. Cater: Intellectual property protection on text generation apis via conditional watermarks. Advances in Neural Information Processing Systems, 2022.
- [Heikkilä, 2022a] Melissa Heikkilä. Artists can now opt out of the next version of stable diffusion. <https://www.technologyreview.com/2022/12/16/1065247/artists-can-now-opt-out-of-the-next-version-of-stable-diffusion/>, 2022.
- [Heikkilä, 2022b] Melissa Heikkilä. This artist is dominating ai-generated art. and he's not happy about it. <https://www.technologyreview.com/2022/09/16/1059598/this-artist-is-dominating-ai-generated-art-and-hes-not-happy-about-it/>, 2022.
- [Heikkilä, 2023] Melissa Heikkilä. Ai models spit out photos of real people and copyrighted images. <https://www.technologyreview.com/2023/02/03/1067786/ai-models-spit-out-photos-of-real-people-and-copyrighted-images/>, 2023.
- [Ho *et al.*, 2022] 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.
- [Holmes and McLaughlin, 2023] Aaron Holmes and Kevin McLaughlin. Ghost writer: Microsoft looks to add openai’s chatbot technology to word, email. <https://www.theinformation.com/articles/ghost-writer-microsoft-looks-to-add-openais-chatbot-technology-to-word-email>, 2023.
- [Holz and Claburn, 2022] David Holz and Thomas Claburn. David holz, founder of ai art generator midjourney, on the future of imaging. https://www.theregister.com/2022/08/01/david_holz_midjourney/, 2022.
- [Jennings, 2022] Richi Jennings. Devs: Don’t rely on github copilot — legal risk gets real. <https://www.reversinglabs.com/blog/devs-dont-rely-on-github-copilot-legal-risk-is-real>, 2022.
- [Johnson, 2022] Khari Johnson. Dall-e 2 creates incredible images—and biased ones you don’t see. <https://www.wired.com/story/dall-e-2-ai-text-image-bias-social-media/>, 2022.
- [Kandpal *et al.*, 2022] Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language models. *arXiv preprint arXiv:2202.06539*, 2022.
- [Kirchenbauer *et al.*, 2023] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
- [Lazaridou *et al.*, 2021] Angeliki Lazaridou, Adhiguna Kunoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, C d M d’Autume, Sebastian Ruder, Dani Yogatama, et al. Pitfalls of static language modelling. *arXiv preprint arXiv:2102.01951*, 2021.
- [Lopez, 2023] James Lopez. Microsoft, and amazon guard against chatgpt theft, ban employees from sharing sensitive data. <https://www.techgoing.com/microsoft-and-amazon-guard-against-chatgpt-theft-ban-employees-from-sharing-sensitive-data/>, 2023.
- [Meehan *et al.*, 2020] Casey Meehan, Kamalika Chaudhuri, and Sanjoy Dasgupta. A non-parametric test to detect data-copying in generative models. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- [Midjourney, 2022] Midjourney. Midjourney: Terms of service. <https://midjourney.gitbook.io/docs/terms-of-service>, 2022.
- [Miller, 2022] Kirk Miller. Google admits its mind-blowing text-to-image ai is endlessly problematic. https://www.insidehook.com/daily_brief/tech/google-imagen-text-to-image, 2022.
- [Mok and Zinkula, 2023] Aaron Mok and Jacob Zinkula. Chatgpt may be coming for our jobs. here are the 10 roles that ai is most likely to replace. <https://www.businessinsider.com/chatgpt-jobs-at-risk-replacement-artificial-intelligence-ai-labor-trends-2023-02>, 2023.
- [Narayanan *et al.*, 2021] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostafa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15, 2021.
- [Nichol, 2022] Alex Nichol. Dall-e 2 pre-training mitigations. <https://openai.com/blog/dall-e-2-pre-training-mitigations/>, 2022.

- [Office, 2023] U.S. Copyright Office. What is copyright? <https://www.copyright.gov/what-is-copyright/>, 2023.
- [OpenAI, 2022] OpenAI. Chatgpt: Optimizing language models for dialogue. <https://openai.com/blog/chatgpt/>, 2022.
- [OpenAI, 2023] OpenAI. Gpt-4. <https://openai.com/research/gpt-4>, 2023.
- [Overflow, 2022] Stack Overflow. Temporary policy: Chatgpt is banned. <https://meta.stackoverflow.com/questions/421831/temporary-policy-chatgpt-is-banned>, 2022.
- [Prabhu and Birhane, 2020] Vinay Uday Prabhu and Abeba Birhane. Large image datasets: A pyrrhic win for computer vision? *arXiv preprint arXiv:2006.16923*, 2020.
- [PromptBase, 2022] PromptBase. Promptbase official website. <https://promptbase.com>, 2022.
- [Radford *et al.*, 2018] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [Radford *et al.*, 2019] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [Radford *et al.*, 2021] 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.
- [Ramesh *et al.*, 2021] 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.
- [Ramesh *et al.*, 2022] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [Rombach *et al.*, 2022a] 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.
- [Rombach *et al.*, 2022b] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Stable diffusion github repository. <https://github.com/CompVis/stable-diffusion>, 2022.
- [Rombach *et al.*, 2022c] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Stable diffusion v1 model card. https://github.com/CompVis/stable-diffusion/blob/main/Stable_Diffusion_v1_Model_Card.md, 2022.
- [Rosenblatt, 2022] Kalhan Rosenblatt. Chatgpt banned from new york city public schools' devices and networks. <https://www.nbcnews.com/tech/tech-news/new-york-city-public-schools-ban-chatgpt-devices-networks-rcna64446>, 2022.
- [Runway, 2022] Runway. Text to video. <https://runwayml.com/text-to-video/>, 2022.
- [Saharia *et al.*, 2022a] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [Saharia *et al.*, 2022b] 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*, 2022.
- [Schuhmann *et al.*, 2021] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- [Singer *et al.*, 2022] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- [Somepalli *et al.*, 2022] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. *arXiv preprint arXiv:2212.03860*, 2022.
- [SpawningAI, 2022] SpawningAI. Have i been trained? <https://haveibeentrained.com>, 2022.
- [Sun *et al.*, 2023] Xiaofei Sun, Xiaoya Li, Yuxian Meng, Xiang Ao, Lingjuan Lyu, Jiwei Li, and Tianwei Zhang. Defending against backdoor attacks in natural language generation. In *AAAI*, 2023.
- [Villegas *et al.*, 2022] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*, 2022.
- [Vincent, 2023] James Vincent. Getty images is suing the creators of ai art tool stable diffusion for scraping its content. <https://www.theverge.com/2023/1/17/23558516/ai-art-copyright-stable-diffusion-getty-images-lawsuit>, 2023.
- [Webster *et al.*, 2021] Ryan Webster, Julien Rabin, Loic Simon, and Frederic Jurie. This person (probably) exists. identity membership attacks against gan generated faces. *arXiv preprint arXiv:2107.06018*, 2021.
- [Weidinger *et al.*, 2021] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa

- Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.
- [Welbl *et al.*, 2021] Johannes Welbl, Amelia Glaese, Jonathan Uesato, Sumanth Dathathri, John Mellor, Lisa Anne Hendricks, Kirsty Anderson, Pushmeet Kohli, Ben Coppin, and Po-Sen Huang. Challenges in detoxifying language models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2447–2469, 2021.
- [Whang *et al.*, 2022] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16293–16303, 2022.
- [Wiggers, 2022a] Kyle Wiggers. Deepfakes for all: Uncensored ai art model prompts ethics questions. <https://techcrunch.com/2022/08/24/deepfakes-for-all-uncensored-ai-art-model-prompts-ethics-questions/>, 2022.
- [Wiggers, 2022b] Kyle Wiggers. Image-generating ai can copy and paste from training data, raising ip concerns. <https://techcrunch.com/2022/12/13/image-generating-ai-can-copy-and-paste-from-training-data-raising-ip-concerns/>, 2022.
- [Willison and Baio, 2022] Simon Willison and Andy Baio. Exploring the training data behind stable diffusion. <https://simonwillison.net/2022/Sep/5/laion-aesthetics-weeknotes/>, 2022.
- [Wortsman *et al.*, 2022] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7959–7971, 2022.
- [Xu *et al.*, 2021] Albert Xu, Eshaan Pathak, Eric Wallace, Suchin Gururangan, Maarten Sap, and Dan Klein. Detoxifying language models risks marginalizing minority voices. *arXiv preprint arXiv:2104.06390*, 2021.
- [Yu *et al.*, 2022] 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*, 2022.
- [Zhang *et al.*, 2022] Zhiyuan Zhang, Lingjuan Lyu, Xingjun Ma, Chenguang Wang, and Xu Sun. Fine-mixing: Mitigating backdoors in fine-tuned language models. In *EMNLP*, 2022.