# The AI Revolution of Good and Bad

P. Somol, M. Salat, S. Afroz, M. Pechoucek, T. Pevny

# Introduction

It is undeniable that the current AI revolution is reshaping the world and opening unprecedented opportunities for human progress. In this whitepaper, we will review aspects of current and upcoming AI technology that could be seen as problematic or harmful. We will focus on the effects of AI technology that societies should reflect on, prepare for and possibly defend against to ensure a positive balance of the AI revolution.

# Gen

# Generated content

## Generation unlimited

The generative capabilities of AI models can appear across a virtually limitless domain—from generating images, videos, sketches and text documents to audio, speech and music with a simplicity and quality that is nothing short of breathtaking. For instance, procuring a new Shakespearean play now merely requires a handful of concise instructions assigned to a sophisticated language model. Similarly, obtaining an illustration for your sci-fi novel is just as straightforward by leveraging the composition of one of your own photos that is set in a fantastical realm. If you're curious how Antonín Dvořák might rework a contemporary rap song into a symphony, it's easily accessible. With a bit of patience, these ventures yield outcomes that are impressively credible and high quality. The issue and debate of generated content was perhaps most strikingly highlighted by the viral circulation of an apparently genuine photograph of Pope Francis, concocted by the Midjourney generator (Novak, 2023).

## Copyright versus training data

Setting aside questions of the quality of generated outputs in terms of human originality and authenticity, the capabilities of current generative models have reopened challenging discussions on originality and artistic depth. The essence lies in the extent to which generated digital content is merely a statistical re-interpretation of pre-existing human outputs. Despite this, generative models rely on a vast trove of pre-existing data from which they derive their expressive potential. This ushers in a significant legal quandary—the creators of these large models have thus far prioritized the rapid development of the technology over a careful consideration of the data ownership used for training their models (Dean, 2023). Therefore, they may, in some cases, infringe upon the copyrights of creators whose works are accessible online in any form. This is a serious issue, though it's anticipated that with time, legal frameworks governing the use of generative models will become clearer.

## Utility versus reliability

The practical utility of generative models is indisputable in areas where human labor does not aim at originality but rather seeks to reorganize information, summarize, amalgamate data from various sources and fill pre-prepared templates. Compiling technical documentation, preparing case studies based on a collection of rules or laws, or assembling encyclopedic knowledge chiefly requires diligence and patience over creativity; here, generative models prove their significant added value. However, they cannot yet guarantee the accuracy of the outcomes under all circumstances. For users, it remains imperative to meticulously check the generated results. While text-based generative models predominantly deliver excellent outputs, it's easy to be lulled into overreliance on their accuracy. Their errors, subtle and convincingly real as they may be, pose a challenge. Notably, a legal firm in the USA encountered difficulties after it was revealed that the justification for a case it submitted was generated by a language model and contained references to non-existent sources (Bohannon, 2023).

## The imitation game

A paramount challenge for society will be the fact that, with the aid of generative models, virtually anyone can imitate others. The possibility of creating a lifelike avatar to replace one's appearance during a video call is already tangible. Similarly, voice imitation and movement characteristics are likely to follow suit soon. Before long, we will be grappling with an influx of utterly false yet convincing videos, news and voice recordings in addition to active fraud in social engineering. The extent to which generated content will inundate our digital realm remains to be seen, but it's prudent to assume that a significant portion of what we encounter online will be designed to deceive our senses and undermine our trust in navigating the world.

# Generated content

### Generation intent
Distinguishing between malicious and benign AI-generated content hinges not only on its technical aspects but also on the intent and context behind its creation and dissemination. Just as old photographs can be innocuously mislabeled to suggest events that never occurred, AI-generated content can be weaponized or used for deception on the basis of how it's presented or framed. AI-generated content should not be considered problematic based solely on the fact it has been generated. This creates a need for precise distinction of the malicious use of AI-generated or AI-enhanced content, as there is a narrow margin between what is still acceptable use and what presents a threat. In the following, we will consider Deepfakes to be AI-generated or AI-altered content supporting or conveying information that is not true and is used with a malicious intent such as misinformation or impersonation. For example, actor Tom Hanks' digital persona was used in a dental plan commercial without his consent, raising questions about intellectual property and identity theft (Guardian, 2023).

# Deepfakes and misinformation

### Deepfakes
Generative AI makes it easy to generate faces, audio and video footage impersonating individuals and manipulating them into saying or doing anything. Creating fake content with the intent to deceive is not a new concept. Two decades ago, Photoshop had a similar transformational impact on digital photography. Nevertheless, deepfakes leverage AI techniques to dramatically decrease the necessary skill requirements and increase the scale at which content can be personalized to the victim's characteristics. Today, it's possible to generate deepfake videos using only one image of a victim and around 30 seconds of audio. The use of deepfakes more than doubled from 2022 to 2023. On a broad scale, 98% of the deepfake videos online are porn (HSH, 2023), a majority (94%) of which target female celebrities and influencers. Around 3% of identity fraud (BW, 2023) including liveness bypass, edited ID card and forged ID cards in 2023 used deepfake images, with 88% of these targeting the crypto industry (IS, 2023). ID Cards seem to be the most vulnerable document type. Generated faces are used to create fake profiles on platforms like LinkedIn and dating apps to scam users, targeting individuals for reputation damage, cyber-exploitation and other forms of abuse. Impersonating someone in a user's trusted contacts, such as the CEO of their company or a family member, can be used for financial fraud or even blackmail. Popular online personalities like Taylor Swift, Mr Beast, Mike Saylor and Martin Lewis have been deepfaked to scam consumers in advanced fee fraud, leading to stolen funds. This type of deepfake is primarily spread through social media. Social media examples include the Pentagon explosion video (CNN, 2023), audio deepfakes used to disparage a candidate in an election in Slovakia (Schneier, 2023), videos of Zelensky (Rozsa, 2023) and Putin (Muzaffar, 2023) created for political disruption.

### Misinformation
Generative AI models are designed to produce content and can generate virtually anything when provided with the appropriate prompt. Misinformation by generative AI can be intentional, such as bad actors using generative AI to supercharge misinformation generation. They can also be unintentional, such as generative AI producing misinformation as a result of hallucination. Research found that LLM-generated misinformation can be harder to detect for humans and detectors compared to human-written misinformation (Chen, 2023).

# Personal information, identity and behavioral models.

**Digital versus physical world**
Today, our "digital life" is intertwined with the internet through various communication channels, online services and social networks. These aspects of digital engagement often utilize search engines and sometimes, unknowingly, are also leveraging various forms of AI that now run behind virtually all digital services. When comparing the internet as our digital societal world to the physical one, we find them to be vastly different.

Our physical world, developed over time, teaches us from childhood to navigate it safely, building at least two significant pillars of societal life: the ability to identify individuals at a personal level and distinguish friendly from calculative or dangerous behavior, and a system of shared organizations and services (healthcare, police, education, judiciary, banking, transportation infrastructure) that significantly enhance our quality of life.

**The problematic legacy of Internet origins**
The internet, which is becoming a digital extension that is increasingly merging with our societal and physical worlds, fails to adhere to these pillars. Created for narrowly defined military and academic purposes, it lacks built-in mechanisms for ensuring user identity, privacy and security. These aspects are foundational for prosperity in the physical world, where impersonating someone or concealing one's existence is universally perceived as suspicious or criminal. We, as people, determine which information to keep private or share the pillars mentioned above to make it possible to exist safely in the physical world.

The internet's significant flaws in identity, privacy and security are familiar, posing a serious unresolved problem despite the existence of impartial solutions. Its open, distributed architecture and boundlessness have opened unprecedented possibilities for rapid progress, service efficiency, global connectivity and knowledge sharing, but also endless new abuse methods.

**Pitfalls of living in the digital world**
For someone who has been given the opportunity to adapt to changes in the physical world relatively gradually (though some argue that the development of the last few centuries has pushed humans to our limits), the pace of change driven by rapid internet development presents an exceptionally demanding challenge. This challenge is made more difficult by the blurring of lines between the physical and digital worlds, particularly concerning identity, privacy and security. As humans, we will undoubtedly overcome this challenge to our collective benefit, but we must navigate several current pitfalls. In today's digital world, these pitfalls include:

1. The identity (or even existence) of people we communicate with is much harder to gauge than in the physical world.
2. Our own level of online privacy is more difficult to estimate and control.
3. The truthfulness of information obtained from the internet is harder to assess and verify.
4. The intentions of our contacts (plus service and information providers) on the internet are more challenging to determine.

**No free digital lunch**
We reveal more about ourselves on the internet than it might seem initially. Services "for free" are often paid for at the expense of our data and privacy. AI offers extensive possibilities for estimating not only the usual preferences of many users for targeted advertising but also their emotional reactions, social status, influence on others, political preferences or personal situations. This might remind you of a spy thriller where opposing sides gather information about their targets for exploitation. The unintended disclosure of information poses a real danger, although typically in a less severe form than in thriller scenarios.

## Utility of personal behavioral data

Whether consciously or inadvertently, providing private information on the internet is primarily exploited by service providers for targeted advertising. This form of advertising is dramatically more effective than untargeted advertising and can be seen as a mild form of manipulation. However, the current internet problem lies in the virtually limitless possibilities to escalate manipulation further. It's not just about selling something through ads; "stealing attention" can lead to internet content addiction. Some authors call the socio-economic model of the current internet "surveillance capitalism" (Zuboff, 2015). Have you ever noticed that a quick internet visit for a video recipe can easily spiral into hours of browsing various videos, making it difficult to detach, with most of that time feeling satisfied and interested? An important question for all internet users is: do we realize when we are being manipulated? Are we aware that even the smartest and most educated people can be manipulated? It's not about personal weakness but rather internet content providers cleverly taking advantage of human psychological traits.



## Exploiting emotions

The example above describes using (AI-estimated) user profiles to elicit positive emotional satisfaction, thus capturing attention. Unfortunately, negative emotions are even easier to provoke online, with an equally compelling and captivating power. Evoked negative emotions can paradoxically also bring satisfaction through a heightened desire for justice, rectification and a better world. However, in the digital world, exercising caution is advised. AI, in the service of subtle or strong manipulation, plays on our brain's reward centers more effectively than ever before, serving sellers, interest groups and states. Today's issue includes the opaque use of AI for these purposes. Yet, the problem is not with AI itself, but how it amplifies the internet's positive and negative capabilities.

## Transparency reforms

The technological and societal challenge lies in reshaping the internet's foundations, which serve as the primary platform for deploying global AI models. This restructuring aims to reduce the discrepancy between the relatively transparent pillars of the physical world, where communication, cooperation and coexistence can thrive based on an acceptable level of mutual trust and the opaque ecosystem of the digital realm. For us, as digital citizens of the internet world, the challenge is rooted in better realizing how much more public the internet is than it seems, how much more insidious and unpredictable the consequences of our actions online are and how much we let the internet influence ourselves. Our own autonomy is at stake.

# Cybersecurity on steroids

### Attackers versus defenders

The world of cybersecurity bears much resemblance to the realm of video games, with the formal mathematical possibility to describe both worlds using game theory. Cybersecurity encompasses the internet, computer networks, interconnected computers and a myriad of devices from mobile phones to cars. The plethora of devices are all propelled by computer programs ranging from operating systems and communication protocols to office software and games. Software and hardware ensure the proper functioning and the desired access restrictions to devices or documents.

However, from its inception, the digital world has suffered from a fundamental flaw: due to its complexity, creators regularly overlook or misjudge errors, making the system vulnerable to those who discover these mistakes, potentially gaining unauthorized access to documents, or in extreme cases, control over a system. Preventive error detection is a standard part of the development of all digital world components, with most errors being discovered and corrected before any misuse occurs. Unfortunately, this is not always the outcome, or the response is not swift enough.

### Attack generation

For a long time, attackers have employed semi-automatic and automatic attack creation methods. The "dark web" has long been a hub for the cybercrime economy. It offers malware development environments, pre-prepared offensive software libraries, specific yet publicly unknown vulnerabilities and even groups of infected computers controlled by an attacker. Up to now, the utilization of more powerful but expensive artificial intelligence tools in automation, encompassing creation, obfuscation and deployment of malware, has been limited. Economic principles apply; as long as cheaper technologies prove sufficient for attacks, there is no need to employ costlier technologies. Many devices on the present-day internet are still easily compromised by simpler methods, and therefore we do not yet see a massive transition in offensive technologies leveraging the latest forms of AI.

This situation is only temporary and with growing awareness among digital citizens for the need to keep systems updated and the increasing quality of foundational security for operating systems and digital devices, the room for classic offensive tactics is diminishing. At the same time, the advent of artificial intelligence practically surpasses individual human capabilities in all data-intensive tasks. AI techniques will inevitably become a fundamental tool for all future attacks. Both attackers and defenders searching for vulnerabilities is a key area where the rapid adoption of AI is expected due to growing complexity. Current AI models already possess strong capabilities in analyzing and generating computer code, thus dramatically accelerating both the analysis of vulnerabilities and the creation of new attack variations. This includes the generation of new, more complex offensive tactics across computer systems, utilizing game theory, similar to how AI learns to play games itself (OpenAI Research, 2018).

### Scam and fraud

A rapidly growing branch of computer attacks targets users directly rather than digital devices. We observed a dramatic increase in phishing attacks, where an attacker will deceive a user with a seemingly legitimate website such as their bank with the purpose of extracting their access data, as well as scam attacks, where various forms of fraudulent money are extracted or manipulated from users to benefit the attacker.

These attacks stem from social engineering, aimed at deceiving the user. A wide range of attacks exists, from very naive general ones (Nigerian prince) to highly convincing personalized, time-staggered ambushes to gain the victim's trust. With the advent of generative AI models, social engineering has become more accessible and dangerous than ever before. Not only have AI models significantly reduced the standard for the creation of authentic-looking text with perfect grammar in virtually any language, but the recent advancements in image and video generation dramatically affect our perception of what is legitimate. We have reached a point where experience no longer holds true. In the past, we would advise users to look for grammar mistakes and unusual formulations, however, the situation is now nearly reversed.

Over the last couple of months, we have revisited the business email compromise (BEC) scams where attackers posed as a senior executive of a company urging to extract payment from business partners. These attacks evolved from emails into video calls where attackers create virtual avatars of the executive using public pictures and videos to create an even more convincing illusion.

## Protecting generative models

As users, it is imperative to keep in mind that various internet services based on artificial intelligence can process user inputs to provide the desired response to the user and improve the models themselves. A properly trained AI model should use information from training data only for abstract generalization; the exact form of the training data should not be reconstructed from the resulting model. However, incorrect training may lead to memorization, where specific pieces of training data are remembered. Therefore, memorization of user inputs can occur by mistake rather than out of malice. Even without memorization, data leakage from the model is theoretically possible if it is subjected to a successful targeted adversarial attack. If an attacker can test model responses from a user's perspective, they may, under certain circumstances, estimate details of the information used for model training. In the case of large language models, it is theoretically possible to approximate a user's query or "prompt" from the model's output. Fortunately, all these types of information leakage are unlikely.

Attackers may not only seek to steal non-public information from the model but also aim to intentionally damage the model or its outputs. Under certain conditions, even an ordinary user can engage in what is known as "model poisoning." An interesting case was Microsoft's well-intentioned attempt in 2016 to create a public chatbot named Tay that would improve human communication by learning from interactions with real people. The experiment required termination after a few hours when the chatbot rapidly began to exhibit racist and aggressively misogynistic behavior (Kraft, 2016). This happened because of how people began interacting with the chatbot. Evidently, too many users tried to communicate with it offensively, aggressively and with the negativity humans are capable of. The reasons can be a whole spectrum, ranging from curiosity and playful disruption to venting frustration in an environment presumed to be anonymous and without a human counterpart.

Arbitrary third parties that cannot be trusted could amplify existing security and privacy risks and introduce new attack vectors such as hijacking generative AI platforms or prompts. This could result in polluting training data or stealing plugin data and denial of service by generative AI platforms or plugins (Umar, 2023).

On the side of the creators, the recent advent of large language models has required the introduction of a comprehensive set of measures to prevent undesirable behavior of the model. Many current large models incorporate controls for both user inputs and model outputs, along with the selection of training data aimed at preventing the worst excesses. The fact remains that no available mechanism can provide protection to the intended extent, as the fact remains that large language models are constantly in the hackers' sights. It appears that it is still relatively easy to formulate queries (prompts) leading to information on the edge of or beyond legality.

### Hallucinations

The stochastic nature of current AI models renders them unsuitable for providing precise guaranteed information and recommendations due to their inherent tendency to generate unpredictable outputs. The frequency of mistakes is low and newer models are constantly getting better, however, they still cannot give guarantees. Critical applications require caution when applying AI. For instance, in software engineering, consider a scenario where an AI model is tasked with recommending Python packages.

Due to the vastness of the Python ecosystem and the probabilistic nature of the model, it may occasionally "hallucinate" or suggest a non-existing package name. These hallucinations can be observed by an attacker using the same model. The attacker can then create a package with a hallucinated name containing a malicious payload. An unsuspecting victim relying solely on the model's recommendation could inadvertently install the malicious package, highlighting the significant risks associated with relying solely on AI-generated recommendations for critical tasks. Similar risks can be demonstrated with URL address recommendations, where the attacker can register a hallucinated domain and host malicious content.

# Fairness and safeguards

### Model fairness

The quest for fairness in AI models with regard to humans has long been accompanied by research and regulation toward what's termed as "Fair AI." This primarily involves legal protection and self-regulation within companies developing AI solutions to shield users from automated discrimination, such as traditionally disadvantaged groups in the job market, especially when AI is deployed in hiring processes or other scenarios where software assesses individuals. This could also include identifying suspects in criminal activities among other scenarios.

### Curating training data

Today's colossal models are trained on data scraped from the internet, thus mirroring the state of the digital realm. This raises a critical question: Is the internet, at least in its raw form, the best source of data for training large models? The internet, in its vastness, can be a distorted reflection of civilization—with amplified biases, violent and viral content and a high share of misleading information. A considerable portion of internet content is sexual or sexist, and social media platforms often become outlets for emotions, unlike real-world behavior. Such bias can be illustrated again in the 2016 Microsoft Tay chatbot case (Kraft, 2016). The internet is a haven for influential groups, which virtually have unrestricted endless reach online. The so-called dark web has become a difficult-to-regulate marketplace for the worst of human output. Noteworthy ideas such as the concept of cryptocurrency bring dramatic benefits and negatives, such as significant electricity consumption and frequent misuse by the underworld not tolerated in traditional banking operations. The internet is a fragmented picture of our civilization, and as a consequence of its history, wealthier players or countries are more represented online. However, all of these are balanced by the unprecedented benefit of information sharing that supports economic and social development on a scale unimaginable before the internet. Yet, models trained on internet data provide a reflection of civilization in the quality of the internet.

**Responsibility of AI's output**
As artificial intelligence is integrated into many aspects of our lives, an important question arises: Who is responsible for the output of the AI-powered system? Is it the user, system vendor or model provider? For example, if a chatbot deployed by a car dealership offers a huge discount on a car, is the dealer obliged to go through with the sale? If the chatbot provides incorrect information leading to financial harm to the customer, is the company liable? In fact, both these cases have already occurred. In the first case, a Chevrolet dealership deployed an AI chatbot on its website and users convinced it to discount a brand-new car to $1. In another case, an Air Canada chatbot provided wrong information about a bereavement fares policy, resulting in financial damage to the customer. This case was taken to a small claims court, which later decided in favor of the user. While the ethical solutions for these two cases may differ, they both highlight the need to hold the operator of the AI tool formally responsible for its output and the general need for proper guardrails of AI models.

**Freedom of "model speech"**
With models' knowledge spanning the internet, the urgent question of free speech becomes increasingly complex. The issue lies in whether we want absolute freedom, thus enabling anyone to inquire about the best ways to commit a crime without detection. Some "free" models lack training data curation and have no inhibitions: Is this better? Isn't general data from the internet unacceptably biased and indebted to the lesser side of human existence and knowledge?

Among technology providers, there's an ongoing dispute over the appropriate extent of editing and filtering. Offered solutions span a wide range. Elon Musk, as a proponent of absolute free speech, has removed almost all filtering from the X platform. ChatGPT creators, along with Microsoft, perform filtering based on a set of basic ethical requirements to prevent support for violence, prejudice, sexism and hate, but otherwise, leave the models' communicative range unrestricted. Some specialized uses of LLM models yield higher quality outputs in specifically supported areas, such as for educational purposes at Khan Academy (Khan, 2023) or for exact scientific purposes in the ChatGPT and Wolfram Alpha link (Stephen Wolfram, 2023).

AI undoubtedly extracts and provides easier access to a large part of human knowledge than the previous internet combined with search services. This raises an ethical question analogous to the existence of compulsory basic education: As a society, do we want to enforce a certain minimal consensus on the foundations of our society? To what extent do we need to regulate the training sources of large generative models and their outputs?

# Conclusion

Our society meets an increasing number of new challenges to protect safety, freedom and trust in an AI-enabled digital environment. The challenges can be overcome but need action from all participants in the digital world: citizens, private organizations, policymakers and governmental bodies. Building a safer and more trusted AI-supported digital world has the potential to greatly benefit our society in ways we only partially foresee. Failure to do so could worsen the negative effects we already suffer from — worldwide societal divide, instability, inequality and underwhelming utilization of the potential of the internet and AI technology.

# References

Bohannon, M. (2023, June 8). Lawyer Used ChatGPT In Court—And Cited Fake Cases. A Judge Is Considering Sanctions. Forbes. Retrieved January 13, 2024, from
https://www.forbes.com/sites/mollybohannon/2023/06/08/lawyer-used-chatgpt-in-court-and-cited-fak e-cases-a-judge-is-considering-sanctions/

BW (2023). New North America Fraud Statistics: Forced Verification and AI/Deepfake Cases Multiply at Alarming Rates. Business Wire.
https://www.businesswire.com/news/home/20230530005194/en/

C2PA (2022, January 26). Coalition for Content Provenance and Authenticity.
https://c2pa.org/

Chen, C., & Shu, K. (2023). Can LLM-Generated Misinformation Be Detected? ArXiv, abs/2309.13788. Accepted to Proceedings of ICLR 2024.
https://arxiv.org/abs/2309.13788

CNN (2023, May 23). 'Verified' Twitter accounts share fake image of 'explosion' near Pentagon, causing confusion. CNN Business.
https://edition.cnn.com/2023/05/22/tech/twitter-fake-image-pentagon-explosion/index.html

HSH (2023). 2023 State of DeepFakes. Home Security Heroes.
https://www.homesecurityheroes.com/state-of-deepfakes/

IS (2023, November 28). Deepfake Digital Identity Fraud Surges Tenfold, Sumsub Report Finds. Infosecurity Magazine.
https://www.infosecurity-magazine.com/news/deepfake-identity-fraud-surges/

Umar, I., Kohno, U., and Roesner, F. (2023). LLM Platform Security: Applying a Systematic Evaluation Framework to OpenAI's ChatGPT Plugins.
https://arxiv.org/abs/2309.10254

Khan, S. (2023, March 14). Harnessing GPT-4 so that all students benefit. A nonprofit approach for equal access! Khan Academy Blog. Retrieved January 13, 2024, from
https://blog.khanacademy.org/harnessing-ai-so-that-all-students-benefit-a-nonprofit-approach-for-equa l-access/

Kraft, A. (2016, March 25). Microsoft shuts down AI chatbot after it turned into a Nazi. CBS News. Retrieved January 13, 2024, from
https://www.cbsnews.com/news/microsoft-shuts-down-ai-chatbot-after-it-turned-into-racist-nazi/

Muzaffar (2023, June 7). Deepfake Putin declares martial law and cries: 'Russia is under attack'. Independent.
https://www.independent.co.uk/news/world/europe/deepfake-putin-martial-law-state-media-b2353005.html

Novak, M. (2023, March 26). That Viral Image Of Pope Francis Wearing A White Puffer Coat Is Totally Fake. Forbes. Retrieved January 13, 2024, from
https://www.forbes.com/sites/mattnovak/2023/03/26/that-viral-image-of-pope-francis-wearing-a-whit e-puffer-coat-is-totally-fake/

OpenAI Research. (2018, July 4). Learning Montezuma's Revenge from a single demonstration. OpenAI. Retrieved January 13, 2024, from
https://openai.com/research/learning-montezumas-revenge-from-a-single-demonstration

# References

Rozsa (2023, April 15). Deepfake videos are so convincing — and so easy to make — that they pose a political threat. Salon.
https://www.salon.com/2023/04/15/deepfake-videos-are-so-convincing--and-so-easy-to-make--that-they-pose-a-political/

Schneier (2023). Deepfake Election Interference in Slovakia. Schneier on Security.
https://www.schneier.com/blog/archives/2023/10/deepfake-election-interference-in-slovakia.html

Stephen Wolfram. (2023, March 23). ChatGPT Gets Its "Wolfram Superpowers"!—Stephen Wolfram Writings. Stephen Wolfram Writings. Retrieved January 13, 2024, from
https://writings.stephenwolfram.com/2023/03/chatgpt-gets-its-wolfram-superpowers/

Zuboff, Shoshana. (2015, April 4). Big Other: Surveillance Capitalism and the Prospects of an Information Civilization. Journal of Information Technology (2015) 30, 75–89.
doi:10.1057/jit.2015.5, Available at SSRN: https://ssrn.com/abstract=2594754

**If you want more information, please reach-out to:**

**Kim Allman**
**Head of Corporate Responsibility,**
**ESG & Government Affairs**
**Kim.Allman@GenDigital.com**

Transparency Register number:  083146048556-68

United States: 60 E Rio Salado Pkwy STE 1000 Tempe, AZ 85203

Czech Republic: Enterprise Office Center Pikrtova 1737/1A 140 00 Prague 4