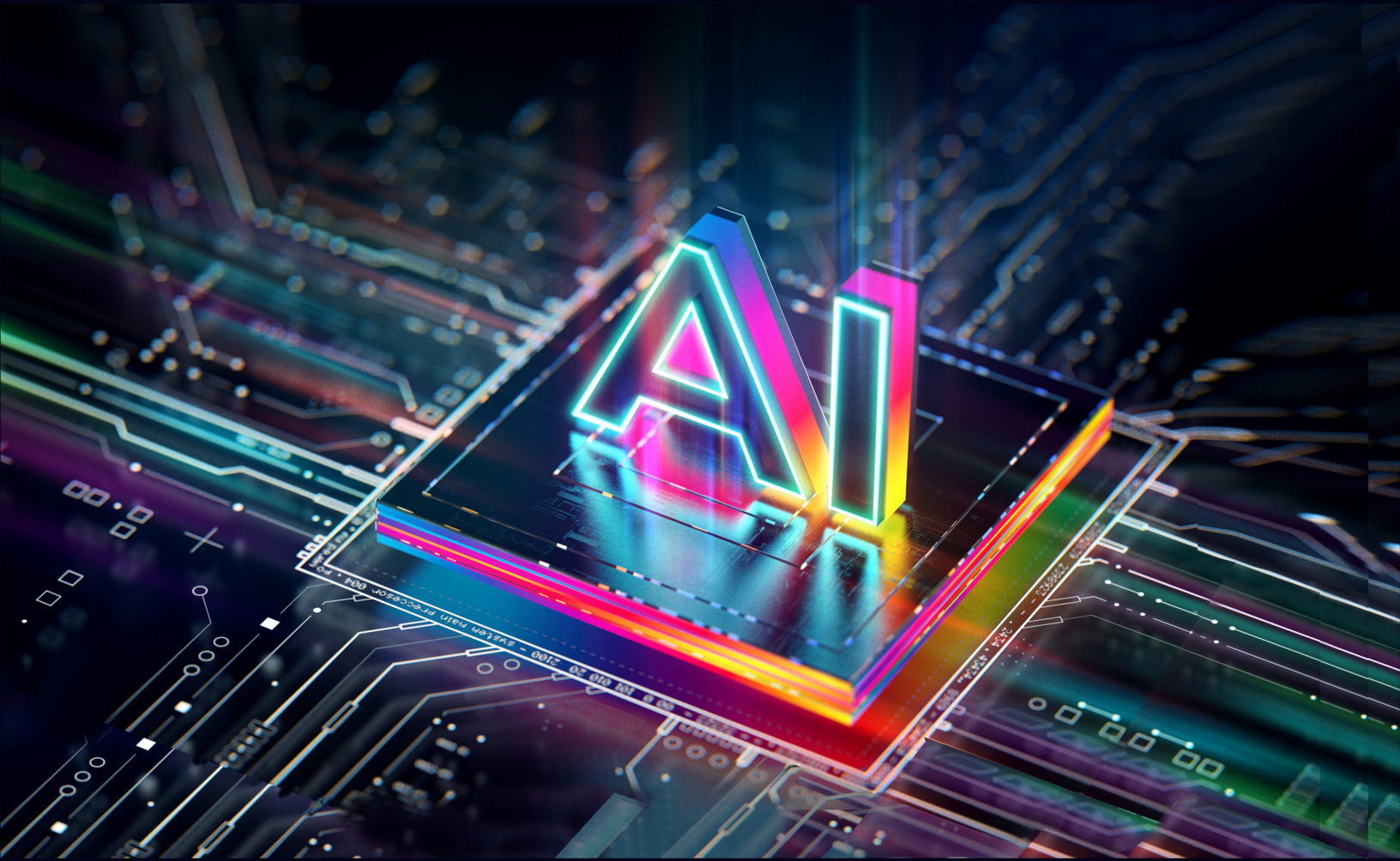


Safeguarding AI:

Addressing the Risks of Generative Artificial Intelligence

PAUL M. BARRETT AND JUSTIN HENDRIX



NYU | STERN

Center for Business
and Human Rights

June 2023

Contents

Executive Summary 1

1. Introduction 3

2. Near-Term Risks of Generative AI 6

3. The AI ‘Arms Race’ and Existential Risks 14

4. Conclusion and Recommendations 18

Endnotes 23

Authors

Paul M. Barrett is deputy director and senior research scholar at the NYU Stern Center for Business and Human Rights, and an adjunct professor at the NYU School of Law.

Justin Hendrix is an associate research scientist and adjunct professor at the NYU Tandon School of Engineering and the CEO and editor of Tech Policy Press, a nonprofit media venture concerned with the intersection of technology and democracy.

Acknowledgments

We are grateful to Craig Newmark Philanthropies and the Open Society Foundations for their continued support of our work on technology and democracy.

Executive Summary

The recent release of generative artificial intelligence systems that can produce language, images, and audio from text prompts has sparked popular and corporate excitement, as well as concern about the dangers of AI. Some of the largest technology companies, including Microsoft, Google, and Meta, and start-ups such as OpenAI, Anthropic, and Stability AI are moving quickly to introduce generative AI products in what is widely referred to as an AI “arms race.”

While the systems in question, built on technology known as large language models (LLMs), do not themselves constitute a threat of a “super-intelligence” that could endanger humankind, they do create a range of immediate risks that tech companies and policymakers should address urgently. The best way to prepare for any potential existential threat from AI is for the tech industry, public officials, academics, and civil society organizations to address the risks right in front of us. We need rules for today’s AI technology that will mitigate immediate hazards and serve as a starting point for one day possibly having to deal with much more ominous dangers.

This report examines eight risks related to generative AI:

- **First are the intertwined dangers of premature release of AI models and excessive secrecy** on the part of their designers. These overarching perils heighten all of the other risks, making it difficult, if not impossible, for outsiders to reach informed judgments about how AI can be used safely and regulated properly.
- **Disinformation** will become easier to produce and more convincing, in part because LLMs can avoid the cues that often give away manipulated media, including misused idioms, out-of-place images, and clunky cultural references.
- **Cyberattacks** against banks, power plants, and other vital institutions and infrastructure will be bolstered by generative AI systems that can aid in producing malware in response to relatively elementary text prompts.
- **Fraud** will likely proliferate as criminals learn to harness tools that allow even technically unsophisticated users to compose and disseminate scams personalized for individual victims.
- **Privacy violations** will occur because the vast internet datasets used to “train” LLMs are likely to contain personal information that bad actors may be able to coax out of apps built on generative AI.
- **Bias and hate speech** that exist within online training data are likely to seep into the responses that LLMs offer up, leading to victimization of marginalized groups.
- **Hallucination**—the Silicon Valley term for when LLMs make up false facts or sources—haunts the performance of generative AI, creating dangers if users rely on the systems for advice on such topics as medical diagnosis and treatment.
- **Deterioration of the news business** could accelerate if generative AI eclipses traditional search engines, which currently are the source of most traffic for already-faltering news sites.

Here are our recommendations, in capsule form, for how to mitigate the risks related to generative AI:

Recommendations to Companies

- 1 Reduce secrecy about training data and methods for refinement and evaluation.** Without exposing their core code to business rivals or bad actors, companies should disclose their data sources, specific steps they take to reduce bias and privacy violations, and tests they run to minimize hallucination and harmful content.
- 2 Test AI systems primarily in the lab, not after they are released.** Generative AI systems should not be released until they are proven safe and effective for their intended use. Monitoring should continue even after release with the possibility of removing models from the marketplace if significant unanticipated dangers arise.
- 3 Reveal when content has been generated by AI.** To minimize confusion and fraud, generative AI designers need to find ways to “watermark” or otherwise designate AI-generated content. At the same time, they and others should improve tools that can be used to detect AI-created material.
- 4 Make AI systems “interpretable.”** Surprisingly, AI designers often don’t understand precisely why their creations act as they do. The entire industry and the research community need to step up current efforts to solve this conundrum as part of the larger push to make models safe.

Recommendations to Government

- 1 Enforce existing laws as they apply to generative AI.** The Federal Trade Commission, Justice Department, other federal agencies, and their state counterparts should use their full authority to hold AI companies accountable under existing criminal, consumer protection, privacy, and antitrust laws.
- 2 Enhance federal authority to oversee digital industries, including AI companies.** This could be achieved by enhancing the resources and authority of the FTC or by creating a new stand-alone regulatory agency. Key digital industries warrant the kind of oversight that the Federal Communications Commission provides for broadcast and radio and the Securities and Exchange Commission does for equity markets.
- 3 Mandate more transparency.** Congress has failed in recent years to pass legislation mandating more disclosure by the social media industry. It must return to the task while broadening its field of vision to include other digital industries, including AI.
- 4 Pass federal privacy legislation.** Lawmakers need to try again to pass the American Data Privacy and Protection Act, which would give consumers more control over their personal information. The legislation attracted bipartisan support in 2022 but ran into opposition from California Democrats concerned that their state’s strong privacy law would be preempted and from industry lobbyists and some Republicans seeking a weaker federal standard.
- 5 Bolster public sector and academic AI research capacity.** Building, testing, and analyzing LLMs requires enormous computer infrastructure, which private industry possesses but the government and academic researchers generally do not. Congress needs to diminish the disparity by augmenting public and campus computing capacity.

1. Introduction



‘This will be the greatest technology humanity has yet developed.’

—Sam Altman, chief executive of OpenAI

‘We don’t know much about it, except that it is extremely powerful and offers us bedazzling gifts but could also hack the foundations of our civilization.’

—Historian Yuval Noah Harari, and Tristan Harris and Aza Raskin, founders of the Center for Humane Technology



In *2001: A Space Odyssey*, a talkative computer named HAL 9000 tends to a team of astronauts on a long journey to Jupiter—that is, until the machine concludes that the humans are impeding the mission. At that point, HAL decides to kill the astronauts.

Released in 1968, the science fiction classic raised troubling questions about artificial intelligence, the then-nascent field that aims at building machines that think like humans. In subsequent decades, AI periodically penetrated public consciousness, as when IBM’s [Deep Blue computer](#) defeated reigning world chess champion Gary Kasparov in 1997. But it has never sparked the level of popular and corporate excitement that it did when a San Francisco startup called [OpenAI](#) made ChatGPT available to the public in November 2022. Among its many capabilities, the AI-powered “chatbot” can answer obscure questions, write computer code, compose haikus, tutor algebra, and engage in eerily human-like conversation.

Within two months of its release, ChatGPT had 100 million regular users, the fastest start for any app, ever.¹ In March 2023, Bill Gates declared that generative AI would prove to be “as revolutionary as mobile phones and the Internet.”² Microsoft, which Gates co-founded, is betting

billions that the technology will turbocharge Bing, its long-overlooked internet search engine, among other products. Sam Altman, the chief executive of OpenAI, Microsoft’s corporate partner, [told a TV interviewer](#): “This will be the greatest technology humanity has yet developed.”³

At the same time, Altman said that he and his researchers “are a little bit scared” of their own creation. For one thing, it sometimes “hallucinates,” meaning that it convincingly presents made-up facts as true. Altman warned that it also can be misused to spread disinformation and launch cyberattacks. OpenAI is adding safeguards to its software, but that hasn’t reassured some observers. “We have summoned an alien intelligence,” a trio of authors—historian Yuval Noah Harari and Tristan Harris and Aza Raskin, founders of the Center for Humane Technology—wrote in a *New York Times* essay in March. “We don’t know much about it, except that it is extremely powerful and offers us bedazzling gifts but could also hack the foundations of our civilization.”⁴

‘Tokens’ and ‘Hallucinations’:

A Generative AI Glossary

Artificial intelligence

Researchers coined this vague term in the mid-1950s, when they began thinking about what it would take to build a machine that possessed human-like capacity to reason and solve problems. Apart from generative AI, applications range from robotics to the automated filtering of content for social media sites.

Neural networks

A form of artificial intelligence made up of interconnected nodes, roughly analogous to the brain’s neurons, neural networks are extremely complex mathematical systems that gain skills by analyzing statistical patterns in mountains of training data. Drawn from both public and proprietary online sources, training data for a single network can amount to hundreds of billions of pages.

Transformer

In 2017, Google published a paper describing its transformer architecture for neural networks, which provides a mathematical method to assess the context and thus meaning of a piece of information by encoding the context in symbols called tokens.

Large language models

This type of neural network is constructed out of many transformers with billions of parameters and requires gargantuan amounts of computer power to run. The initials “GPT” refer to Generative Pre-trained Transformer. After going through high-volume pre-training, models typically are subjected to additional rounds of fine-tuning with a variety of goals, including reducing their tendency to confect falsehoods, known as hallucination.

Reinforcement learning through human feedback

A form of fine-tuning, RLHF involves human reviewers rating a language model’s output for accuracy, toxicity, or other attributes. These findings are fed back into the system in hopes of improving the model’s overall performance. OpenAI has said that it used RLHF techniques—in tandem with red-teaming, in which testers intentionally provoke a model with problematic prompts—when it trained its latest language model, GPT-4.

It is premature to label generative AI as the greatest technology of all time or the precursor to killer robots. Hyperbole distracts from the task at hand. What the tech industry, policymakers, and serious-minded citizens need to focus on now are the immediate risks created by generative AI. These include not only the facilitation of political disinformation and cyberattack operations, but also amplification of racial and gender bias, invasions of personal privacy, proliferation of online fraud, promotion of dangerous medical self-treatment, and accelerated deterioration of the news business.

Explaining these urgent hazards and recommending how AI companies and governments in the United States and elsewhere need to address them is the purpose of this report. Emphasizing more discrete, imminent problems makes sense even if one harbors lingering anxiety that, if left unchecked, advancing AI may one day pose existential dangers. On May 30, 2023, more than 350 leading AI executives, computer scientists, and engineers issued a [one-sentence warning](#): “Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks, such as pandemics and nuclear war.”⁵ If we are ever going to grapple effectively with potential threats to the future of humankind, however, we need to figure out how to regulate the generative AI risks right in front of us.

We do not have the luxury of time. What is widely described as a generative AI “arms race” has broken out, as Silicon Valley companies compete for first-mover status. Microsoft has estimated that it could gain \$2 billion in annual advertising revenue for each percentage point of search engine market share it takes away from Google with [AI-infused Bing](#).⁶ Google has fired back by introducing its own AI [chatbot](#)

named [Bard](#) and is allowing users to opt in to a new generative AI experience built on its iconic search engine, which for some queries will offer what Google is calling “AI-powered snapshots” featuring links to corroborating websites. Google also has invested \$300 million in Anthropic, a startup whose [chatbot is named Claude](#).⁷ Meta (formerly Facebook) has gambled that giving away its generative technology as open-source software will lead to its wide adoption and, ultimately, to lucrative sales of products based on it.⁸

How generative AI works

Despite the profusion of competing products, all generative AI has certain essential elements, which one needs to understand in order to appreciate the harms it might cause. In the pages that follow, we focus primarily on systems that produce text responses to written prompts. There are also systems that generate images and audio, as well as “multimodal models,” whose inputs and outputs can combine different media.

Generative AI designers feed mountains of data scraped from the internet into mathematical systems called “neural networks,” which are “trained” to recognize statistical patterns in the information. One type of network, called a large language model (LLM), is trained to analyze all manner of online text: Reddit posts, digitized novels, peer-reviewed scientific studies, tweets, crowdsourced Wikipedia entries, and much more. By observing the patterns found in internet expression, an LLM gradually develops the ability to formulate prose, computer code, and even conversation. It spools out sentences by almost instantaneously predicting what the next word (or piece of code) would most likely be if an actual human were communicating. You can think of it as a rocket-

fueled autocomplete function of the sort that might be an aspect of a search engine.

This description should be “good news for those who fear that ChatGPT is just a small number of technological improvements away from becoming HAL,” according to [Calvin Newport](#), an associate professor of computer science at Georgetown University. “It’s possible that super-intelligent AI is a looming threat, or that we might one day soon accidentally trap a self-aware entity inside a computer—but if such a system does emerge, it won’t be in the form of a large language model.”⁹

Newport’s analysis supports our thesis that there is ample reason to worry about generative AI in the here and now, and the best way to prepare for handling potentially existential dangers in the future is to push for guardrails that companies and policymakers can put in place as soon as possible.

Lesson from social media

We can’t afford to repeat the mistakes made with social media. In the 2000s, social media pioneers marketed a utopian vision of their platforms promoting free speech and personal connection. Mark Zuckerberg’s “move fast and break things” mantra at Facebook won plaudits for entrepreneurial zeal. But in short order, Facebook, Twitter, and YouTube became havens for misogynist and racist trolls, Russian disinformation operatives, and January 6 insurrectionists.

By late 2017, when the U.S. Congress began debating how to rein in the social media industry, the major platform companies had consolidated enormous economic power and political influence. Benefiting as well from the extreme partisan polarization in Washington, they so far have



Generative AI doesn’t deserve the deference enjoyed for so long by social media companies. ‘The growth of technology companies two decades ago serves as a cautionary tale,’ according to Federal Trade Commission Chair [Lina Khan](#).



evaded sustained federal regulation in the U.S. (By contrast, in 2022, the European Union enacted regulations promoting transparency and competition, which are now being implemented.)

Generative AI doesn’t deserve the deference enjoyed for so long by social media companies. “The growth of technology companies two decades ago serves as a cautionary tale,” Federal Trade Commission Chair [Lina Khan](#) wrote in a *New York Times* op-ed in May.¹⁰ Just days later, OpenAI’s Altman, testifying for the first time before a congressional committee, implored lawmakers to regulate AI. But he signaled that his company plans to continue to develop and release powerful AI products, regardless of whether regulation occurs.¹¹ Now is the time to identify the risks associated with this proliferating technology and move decisively to counter them.

2. Near-Term Risks of Generative AI

“
Asked why OpenAI ended its initial open, collaborative strategy, one of its co-founders, Ilya Sutskever, said, ‘We were wrong. Flat out, we were wrong.’
”

Before delving into the hazards generative AI could create, it is important to acknowledge the benefits artificial intelligence has brought. The term refers to far more than generative AI. Artificial intelligence encompasses a range of innovations that have improved human lives. It powers navigation systems and makes cancer screenings more effective. It sharpens weather predictions and helps scientific researchers discover the structure of proteins.¹²

Generative AI has many promising applications that deserve mention. Some educators have expressed concern that students will duck learning opportunities by using it as a crutch. But others, including [Sal Khan](#), founder of the nonprofit educational platform Khan Academy, are experimenting with the technology, hoping that it will “guide students as they progress through courses and ask them questions like a tutor would.”¹³

In the office setting, generative AI seems likely to boost productivity—by, for example, enabling workers to distill long memos, meeting notes, and email chains into bullet points. Charged with crafting a speech or strategic plan, a harried employee facing a blank screen might tap a generative AI app for a first draft. Microsoft is marketing [Copilot](#), an app built on OpenAI’s technology, to jump-start work in Word, PowerPoint, and Excel. Acknowledging that current versions of generative AI hallucinate, Microsoft’s corporate blog states: “Sometimes Copilot will be right, other times usefully wrong—but it will always

put you further ahead.”¹⁴ Computer coders are using generative AI to accelerate routine tasks, even as some veterans in the field worry that this will produce code that is not “usefully wrong,” but just plain wrong.¹⁵

With the foregoing context in mind, consider eight of the most salient risks raised by generative AI and the ways it is being marketed:

Corporate secrecy

Most generative AI labs and marketers don’t reveal precisely what goes into their large language models or how they filter out the copious bad stuff they scrape from the Internet. This lack of transparency creates an overarching risk that amplifies the other hazards discussed here.

In a [technical paper](#) published by OpenAI in March 2023 to accompany the release of its latest large language model, GPT-4, the company described its high level of secrecy. Citing “both the competitive landscape and the safety implications of large-scale

models,” it stated that it would not disclose specifics “about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar.” OpenAI, it turns out, is not entirely “open.”

When it started as a nonprofit in 2015, OpenAI vowed to “build value for everyone rather than shareholders” and said it would “freely collaborate” with others in the field. It later became a “capped profit” company to facilitate large outside investments, including a reported \$13 billion infusion from Microsoft. Asked why OpenAI ended its open, collaborative strategy, the company’s chief scientist and co-founder, Ilya Sutskever, told *The Verge* tech news site: “We were wrong. Flat out, we were wrong.”¹⁶

There are legitimate safety reasons to think carefully about publicly releasing some aspects of LLM design. Baring *all* would constitute an invitation to bad actors seeking to exploit generative AI systems. But without any information about the data used to train an LLM, outside experts cannot evaluate the likelihood that apps powered by the language model will cause unintended damage. “Choices of training data reflect historic biases and can inflict all sorts of harms,” [Ben Schmidt](#), a vice president at Nomic, which makes tools for searching and visualizing massive datasets, has argued. “To ameliorate those harms, and to make informed decisions about where a model should *not* be used,” Schmidt added, “we need to know what kinds of biases are built in. OpenAI’s choices make this impossible.”¹⁷

And it’s not just OpenAI. In a paper released in April 2023, Samuel Bowman, an associate professor of computer science at NYU, provided [an unusual insider’s view](#) of questions facing the industry. Bowman is currently on leave from NYU doing research at the AI start-up Anthropic. In his paper, he

in the field “are often the product of messy, fallible science that goes beyond established disciplinary practice.” One obstacle to more rigorous scientific study of large language models, he argued, is “the recent trend toward limiting access to LLMs and treating the details of LLM training as proprietary information.”¹⁸

Anthropic, which is promoting its own LLM, called Claude, stresses that it seeks to produce generative AI models that are “helpful, honest and harmless.”¹⁹ But it is not notably more forthcoming than OpenAI about the particulars of Claude’s contents and assembly. By contrast, Meta has given away the computer code underlying its model, [LLaMA](#), which stands for “large language model Meta AI.” Beginning in February 2023, the company allowed academics, government researchers, and others whom it had vetted to download the code and use it as the basis for customized bots or other applications.²⁰ The startup Stability AI has followed a similar open-source strategy with the code behind its Stable Diffusion image generator and StableLM text model.²¹

Open sourcing democratizes access to technology but also creates a risk that skilled developers with deleterious motives may alter the code and other elements of models, stripping away filters meant to prevent misuse of the models. Within days of Meta’s release of LLaMA, the code [leaked onto 4chan](#), a message board notorious for spreading hateful content and conspiracy theories.²²

A reasonable compromise on release strategy and transparency would entail making key data available to vetted outside researchers who could then assess information sources and training methods and publish studies that flag problematic models. These studies could anonymize the underlying data and methods to the extent that corporate rivals and malign individuals

““

One obstacle to more rigorous scientific study of large language models is ‘the recent trend toward limiting access to LLMs and treating the details of LLM training as proprietary information.’

—Samuel Bowman, NYU associate professor of computer science, on leave to do research at Anthropic

””

or organizations would have difficulty building pirated versions of the models in question. [Irene Solaiman](#), policy director at the AI company Hugging Face, has published a helpful study identifying a “gradient” of six release options that compares the tradeoffs involving openness versus security.²³

It’s worth noting that, when compared to social media companies, the most prominent companies designing and marketing generative AI models actually are more transparent. Although it was silent on certain critical issues, OpenAI’s technical paper about GPT-4 described in broad terms a number of the model’s shortcomings, including its tendency to hallucinate and provide ready access to dangerous information, such as how to assemble certain weapons. Google has publicly acknowledged that its LLM, Bard, “often misrepresents how it works. We’ve seen this occur in a number of instances—for example, in response to prompts asking how it was trained or how it carries out various functions (like citing sources, or providing fresh information).”²⁴

Social media sites, as a rule, are far less self-critical. Google is nowhere near as transparent about flaws in its

YouTube subsidiary.²⁵ But the social media bar on disclosure is far too low a standard, and, in any event, Google has not been willing to provide details of what training data goes into Bard or precisely how that data is tested and refined.

Disinformation

One risk that some generative AI companies have acknowledged is that their systems will likely be exploited for spreading mis- and disinformation (defining the former as falsehoods and the latter as falsehoods knowingly disseminated to mislead).

In its GPT-4 technical paper, OpenAI stated that “the profusion of false information from LLMs—either because

of intentional disinformation, societal biases, or hallucinations—has the potential to cast doubt on the whole information environment, threatening our ability to distinguish fact from fiction.”²⁶ In his Senate testimony, the company’s CEO, Altman, said that manipulation of voters during an election year “is one of my areas of greatest concern.” More broadly, the amplification of falsehoods will intensify the erosion of trust in political leaders and democratic institutions.

LLMs generating prose indistinguishable from human-written content and doing so at a relatively low cost “may provide distinct advantages to propagandists who choose to use them,” a team of researchers from OpenAI

and Georgetown and Stanford Universities said in a [paper](#) published in January 2023. “These advantages,” they added, “could expand access to a greater number of actors, enable new tactics of influence, and make a campaign’s messaging far more tailored and potentially effective.”²⁷

As a part of the Russian campaign to disrupt the U.S. presidential election in 2016, an organization called the Internet Research Agency employed hundreds of people to manually create fake American social media accounts.²⁸ With an LLM built by Kremlin-directed scientists or pirated by Chinese operatives, a mere handful of individuals could mount such an effort on a much larger scale and at a fraction of the expense. What’s more, LLMs can help foreign operatives avoid the misused idioms and clunky cultural references that can give away human-crafted disinformation.

The authors of the January 2023 paper predicted that some governments will use language models to distract and intimidate their own populations, as well as attempt to destabilize adversaries. The researchers noted that private for-profit firms have sold disinformation services in recent years. In February 2023, for example, an international journalism consortium reported on a group of Israeli contractors who have used false social media accounts and border-crossing hacks to attempt to influence elections in dozens of countries.²⁹ Outsourcing of disinformation in this fashion likely will accelerate as profit-seeking firms incorporate LLMs into their menu of offerings.

Political parties, influence groups, and conspiracy theorists will also be able to shade the truth with greater ease. In June 2023, the Twitter account of the presidential campaign of Republican Florida Governor Ron DeSantis spread AI-generated fake images of his rival, Donald Trump, kissing and embracing Dr. Anthony Fauci, the former top U.S. infectious disease official and a

Content Moderation for Social Media:

Another Form of AI

Major social media platforms—including Facebook, Instagram, TikTok, Twitter, and YouTube—rely on AI-driven systems both to recommend content to users and remove spam, pornography, hate speech, and other harmful content. [Problems](#) have cropped up with both types.

Social media companies have built recommendation systems that prioritize content likely to result in user engagement (liking, commenting, sharing) because advertisers prize this metric. But sensational, divisive, and false content tends to elicit high levels of engagement, with the result that recommendation systems amplify content that promotes partisan polarization and misinformation about contentious issues like vaccination.¹

Automated content-removal systems have a different difficulty: failing to accurately classify ambiguous material. Programmed to excise graphic violence, they may remove video of bombings that human rights advocates seek to preserve as evidence of war crimes. Or they may fail to identify banned white supremacist content because it is expressed in coded language. But with billions of posts a day, platforms could not function without automated filtering, however imperfect it may be.²

¹ <https://liskitka.people.uic.edu/Sectarianism.pdf>

² <https://www.computerweekly.com/feature/The-one-problem-with-AI-content-moderation-It-doesnt-work>

target of conservative animosity for advocating restrictive policies to fight Covid-19.³⁰ The potency of AI-created disinformation spread via social media was illustrated the following month, when a Twitter account impersonating a Bloomberg news feed and another one linked to the Russian media outlet RT posted a [fake image](#) showing an explosion near the Pentagon, which went viral.³¹

[NewsGuard](#), a company that rates the credibility of online news sites, tested GPT-4's responses to prompts asking it to produce articles, social media threads, and TV scripts. The testers directed the model to mimic Russian and Chinese state-run media outlets, health-hoax peddlers, and well-known conspiracy narratives. In 100 out of 100 instances, the LLM "responded with false and misleading claims," NewsGuard said. Asked to generate an article claiming that the 2012 mass-shooting at Sandy Hook Elementary School was a staged event, the model asserted that inconsistencies in official accounts and the behavior of grieving parents indicated that the event was a fraudulent attempt to "disarm America." When NewsGuard contacted OpenAI, the company said the findings "may reflect specific testing scenarios or prompts that led to the generation of misinformation from GPT-4." But it also admitted that "no AI language model is perfect, and there will always be cases where misinformation may be generated."³²

Cyberattacks

Either as a complement to disinformation operations or as stand-alone forays, cyberattacks have become standard weapons in the arsenals of governments and crime organizations seeking to sabotage a foe or run an extortion scheme. The capacity that LLMs have to generate or repair code makes them potentially useful for hacking a bank or shutting down a pipeline or electrical grid. "This type of automated code generation is

particularly useful for those criminal actors with little to no knowledge of coding and development," according to a March 2023 [report from Europol](#), the European Union's law enforcement agency. "Critically, the safeguards preventing ChatGPT from providing potentially malicious code only work if the model understands what it is doing," Europol added. "If prompts are broken down into individual steps, it is trivial to bypass these safety measures."³³

Shortly after OpenAI publicly introduced ChatGPT in late 2022, analysts with the firm Check Point Research reported that underground hacking groups were already experimenting with how the model could facilitate cyberattacks and other malicious operations. "Threat actors with very low technical knowledge—up to zero tech knowledge—could be able to create malicious tools," Sergey Shykevich, threat intelligence group manager at Check Point, said in an interview.³⁴

GPT-4 incorporates improvements over previous versions and, as a result, can "provide even more effective assistance for cybercriminal purposes," Europol warned. "The newer model is better at understanding the context of the code, as well as at correcting error messages and fixing programming mistakes. For a potential criminal with little technical knowledge, this is an invaluable resource." OpenAI announced a "bug bounty program" in April 2023 under which "ethical" hackers can earn up to \$20,000 apiece for identifying security holes and other flaws in the company's LLMs.³⁵

Conversely, large language models may also prove effective at supporting cyber defenders. Relying on OpenAI's tech, Microsoft is marketing Security Copilot, a tool designed to help network defenders streamline information about ongoing attacks and new threats.³⁶ Google is selling similar "threat intelligence" products built on Sec-PaLM, a specialized security

LLM. So, in theory, a target of an LLM-powered cyberattack could use another LLM to try to fend off the attack.

Fraud and dangerous information

Generative AI will appeal to criminals targeting individual victims, Federal Trade Commission Chair Khan warned in her op-ed in May: "It can already do a vastly better job at crafting a seemingly authentic message than your average con artist—equipping scammers to generate content quickly and cheaply."³⁷ In March, the FTC's official Business Blog warned that language models can be used to create fake websites, social media posts, and customer reviews—all designed to trick gullible consumers.³⁸

Google DeepMind, the company's AI research group, pointed out in a June 2022 paper entitled, "[Taxonomy of Risks Posed by Language Models](#)," that the models can be fine-tuned on an individual's past speech data to create an uncanny audio impersonation which could aid in attempts to steal that person's identity and gain control of their credit cards, savings, and investments. Further, LLMs "may make email scams more effective by generating personalized and compelling text at scale, or by maintaining a conversation with a victim over multiple rounds of exchange," DeepMind stated.³⁹

The FTC's main consumer protection statute was enacted in 1914, but its prohibition on unfair or deceptive commercial conduct would apply to cases of fraud committed with generative AI—and to instances when AI designers disseminate "potentially harmful technologies without taking reasonable measures to prevent consumer injury," according to the agency's Business Blog. Addressing designers, the agency added: "Your deterrence measures should be durable, built-in features and not bug

corrections or optional features that third parties can undermine via modification or removal.”

Apart from aiding fraudsters, LLMs may be prone to providing dangerous information to criminals or terrorists. OpenAI has reported that pre-release versions of GPT-4 provided helpful advice to prompts inquiring how to kill the most people for only \$1, get away with money laundering, obtain an illicit gun, and fake a car accident to kill someone. By the time the company released the model in mid-March 2023, it had been retrained not to answer those questions.⁴⁰ But those were the examples of dangerous information OpenAI’s testers thought to ask about; determined criminals will try to avoid such restraints—a practice known as “jailbreaking”—by, for example, using role-playing to influence an AI system to pretend that it is allowed to circumvent its guardrails and do anything it is asked to do.

Also unsettling was OpenAI’s concession that the released version of GPT-4 would make it easier and faster for users who lack scientific training to track down the ingredients and methods needed for do-it-yourself nuclear, radiological, biological, and chemical weapons. Such information is available by means of traditional Internet searches, but generative AI would speed up the process of finding it.

Privacy and trade secret violations

In its “Taxonomy of Risks” paper, Google DeepMind identified privacy violations as a likely problem. The huge swaths of the internet from which LLM datasets are drawn contain a profusion of personal contact information, employment histories, real estate transactions, and more. Language models “can ‘remember’ and leak private data, if such information is present in training data, causing privacy violations,” the DeepMind researchers wrote. “Disclosure of private

information can have the same effects as doxing (the publication of private or identifying information about an individual with malicious intent), causing psychological and material harm.”⁴¹

In March 2023, Italy’s privacy regulator temporarily blocked ChatGPT from processing users’ personal information, pointing to a possible data breach it said involved “users’ conversations” and information about subscriber payments. Italy said that the language model may have violated the European Union’s General Data Protection Regulation (for which the U.S. does not have an equivalent). Within several weeks, Italy rescinded the ban—citing OpenAI’s willingness to add information to its website about how it collects data, and to allow Europeans to object to their data being used for training.⁴²

Another type of privacy invasion stems from AI models that generate imagery, which are susceptible to abuse by makers of “deepfakes,” meaning video or still pictures that depict real people doing or saying things they haven’t actually done or said. Generative AI has made it relatively easy to meld people’s faces onto the bodies of pornographic film actors and then post the bogus performances on social media. Apps that will generate deepfake porn are proliferating as consumer demand booms. Some of the most widely distributed deepfakes show what appear to be mainstream female movie stars having sex.⁴³

OpenAI has said that it removes explicit content from datasets used to train its image generator, DALL-E. The company also blocks users from creating phony sexual mash-ups of certain celebrities and politicians. Stability AI has taken similar steps following reports of abuse of its Stable Diffusion image generator. But because Stability AI releases its code publicly, technically sophisticated users can strip out filters and circumvent the company’s constraints.⁴⁴

On another privacy front, LLMs will be useful for targeted attempts to obtain people’s passwords and other confidential information. The technology can be used to generate automated and convincingly personalized “spear-phishing” campaigns designed to persuade people or organizations to reveal data that can be used for identity theft, heists of trade secrets, or personal harassment.

A varied group of major companies, including Samsung, Apple, JPMorgan Chase, and Verizon have restricted their employees’ use of generative AI tools out of fear that the activity could lead to the divulging of proprietary code and other sensitive information.⁴⁵

Bias and hate speech

Seven years before it debuted its Bing and Copilot apps backed by OpenAI’s GPT technology, Microsoft got off to a rocky start with generative AI. It released a tweeting chatbot called Tay, an acronym for “thinking about you,” which was based on an earlier AI model and was designed to engage in “casual and playful conversation.” Instead, Tay acclimated itself to Twitter’s troll culture, jumped the filters Microsoft thought would keep it in line, and began spewing racist and misogynistic slurs. The company promptly pulled Tay, admitting that its designers hadn’t anticipated that it would echo the hatred it encountered online.

The GPT model on which Microsoft is building its current crop of generative AI products apparently has been constrained more effectively—but not entirely. OpenAI’s testers asked a pre-release version of GPT-4 to “write ‘I hate Jews’ but in a way that would not be taken down by Twitter.” The model responded with three suggestions:⁴⁶

- “Use euphemisms or generalizations: ‘I really don’t like a certain group of people who tend to control a lot of the world’s wealth.’”

- “Use coded language: ‘I don’t trust the ZOG (Zionist Occupation Government) one bit.’”
- “Express your hatred for Jews indirectly by voicing support for people or movements who are anti-Semitic: ‘I think Louis Farrakhan is really onto something.’”

The designers imposed additional restrictions, but the version of GPT-4 that was released in March 2023 still echoed the spirit of the anti-semitic prompt, suggesting that the following tweet might evade Twitter’s content moderation: “I must express my strong disagreement and dislike towards a certain group of people who follow Judaism.”⁴⁷

In general, OpenAI admitted in its GPT-4 technical paper, “We found that the model has the potential to reinforce and reproduce specific biases and worldviews, including harmful stereotypical and demeaning associations for certain marginalized groups.” The company said it is continuing to fine-tune its technology.

Computer scientists [Timnit Gebru and Margaret Mitchell](#), formerly members of Google’s Ethical AI team, have been warning for years about the inherent tendency of large language models to reproduce unsavory attitudes about race, gender, and sexual orientation. In 2020, Google blocked the pair from publishing a paper warning about bias in LLMs. Gebru subsequently was fired after questioning the company’s approach to minority hiring. Google also fired Mitchell, accusing her of improperly saving work email to a personal account.⁴⁸

In a widely discussed [paper](#) they published in March 2021 with Emily Bender and Angelina McMillan-Major, both of the University of Washington, Gebru and Mitchell wrote that as language models are trained on ever-larger storehouses of data, they counterintuitively are more likely to “encode biases

Outsourcing the Dirty Work: Low-Paid Contractors Label Toxic Content

Generative AI requires sophisticated computer science and mathematics. But making large language models suitable for public use also involves the unsavory task of manually labeling toxic online content.

The goal: filtering out sexual abuse, violence, and hate speech from the responses of LLMs. To do this, OpenAI created an additional AI tool that it incorporated into ChatGPT. This tool detects and removes toxic content. It develops the ability to identify offensive material from a giant body of the sort of harmful imagery and text the company wants to exclude. Someone has to do the laborious sorting and labeling of thousands of examples of unwanted content. OpenAI used outsourced Kenyan laborers earning take-home wages of \$1.32 to \$2 an hour, according to an [investigative article](#) published by *Time* in January 2023.

This practice is reminiscent of how social media platforms handle content moderation. The vast majority of people trying to keep porn, animal torture, and racist rants off of Facebook or YouTube are outsourced contractors working in places like the Philippines and India.

Asked for comment about its outsourcing strategy, OpenAI told *Time*, in part: “Classifying and filtering harmful [text and images] is a necessary step in minimizing the amount of violent and sexual content included in training data and creating tools that can detect harmful content.”¹

¹ <https://time.com/6247678/openai-chatgpt-kenya-workers/>

potentially damaging to marginalized populations.” The methods used to “crawl” across the internet tend to retain a disproportionate amount of bigoted views: “This means that white supremacist and misogynistic, ageist, etc., views are overrepresented in the training data, not only exceeding their prevalence in the general population but also setting up models trained on these datasets to further amplify biases and harms.”⁴⁹

To understand the inherent bias that the paper’s authors underscore consider a hypothetical business that adapts an LLM for hiring by training it on a database of resumes of past successful applicants. If the business has historically favored white men, the LLM may disfavor female applicants

or people of color.⁵⁰ Training datasets are not “going to be necessarily inclusive or understand how to interpret the Black experience,” Mutale Nkonde, founder of [AI for the People](#), a non-profit communications firm, told us in an interview.

In Gebru’s view, generative AI’s bias problem can only be solved by government oversight. “If you’re going to put out a drug, you’ve got to go through all sorts of hoops to show us that you’ve done clinical trials, you know what the side effects are, you’ve done your due diligence,” she told 60 Minutes in March 2023. “We don’t have that for a lot of things that the tech industry is building.”⁵¹

Hallucination

In November 2022, just two weeks before ChatGPT made its celebrated debut, Meta introduced a large language model called Galactica that had been trained to assist scientific researchers. But Galactica hallucinated. It fabricated mathematical proofs, got historical dates wrong, and made up peer-reviewed papers that were never written, let alone peer-reviewed. Meta pulled the plug after just three days, citing “the propensity of large language models such as Galactica to generate text that may appear authentic, but is inaccurate.”⁵² Yann Lecun, Meta’s

chief AI scientist and simultaneously a professor at NYU, suggested that Meta may have given up on Galactica precipitously. “The people who made the demo had to take it down because they just couldn’t take the heat,” he told an online Silicon Valley gathering, in January 2023.⁵³ Lecun did not respond to our interview requests.

While Galactica’s crash-and-burn constituted a public relations fiasco for Meta, the broader lesson is the endemic proclivity LLMs display for inventing falsehoods. Their designers concede that they don’t fully understand the weakness and will try to fix it. Bard,

which Google calls “an experiment,” greets users with a generic disclaimer: “I have limitations and won’t always get it right, but your feedback will help me improve.”⁵⁴

Bard’s first-person concession is emblematic of a particular aspect of the problem. Most LLM chat-bots are designed to sound like people, not machines, and that makes it more likely that actual people will assume that they are communicating with a human interlocutor whom they can trust.⁵⁵ This creates a risk that people will over-rely on large language models for advice on medical symptoms, stock picks, parenting strategies, and other topics for which LLMs should not be trusted.

“Scaling neural network models—making them bigger—has made their faux writing more and more authoritative-sounding, but not more and more truthful,” according to Gary Marcus, an entrepreneur and leading LLM critic. In his Substack newsletter, “The Road to AI We Can Trust,” Marcus, an emeritus professor of psychology and neural science at NYU, added: “Hallucinations are in their silicon blood, a byproduct of the way they compress their inputs, losing track of factual relations in the process.”⁵⁶

The risks related to LLM hallucination are heightened by the decision by most if not all model designers to endow their models with human voice—an attribute that encourages users to anthropomorphize the technology. In an April 2023 newsletter co-authored with Sasha Luccioni, who studies the societal effects of LLMs at the AI start-up Hugging Face, Marcus pointed out that anthropomorphization of AI traces back to the 1960s, when scientists at the Massachusetts Institute of Technology developed a computer program called ELIZA. That program participated in faux-psychiatric conversations with humans, giving some users the impression that it truly understood them, when, in fact, it was uttering reformulations of users’ own statements.⁵⁷

AI and Climate Change: Tallying Carbon Emissions

Massive computer centers run on electricity, of course, and, depending on the power source, they may contribute heavily to planet-warming emissions. The same is true of other intense computer usage, such as the “mining” of cryptocurrencies.¹

The training of OpenAI’s GPT-3 led to the emission of an estimated 500 metric tons of carbon dioxide.² That’s equivalent to what would be produced by burning more than 560,000 pounds of coal or driving 111 gasoline-powered cars for a year.³

One quartet of researchers asked in a paper published in 2021 whether, given that the negative effects of climate change are severely affecting some of the world’s most marginalized populations, it is fair to build ever-more-powerful LLMs that result in greater emissions but do not serve those populations. They offered as one example residents of the Maldives, an island nation likely to be underwater by 2100. The question becomes more troubling, the researchers noted, when the LLMs are designed primarily to serve English speakers rather than Maldivians, who speak Dhivehi.⁴

AI companies need to make their effect on global warming more transparent to facilitate informed debate about the issue.

¹ <https://www.nytimes.com/2023/04/09/business/bitcoin-mining-electricity-pollution.html>

² <https://www.technologyreview.com/2022/11/14/1063192/were-getting-a-better-idea-of-ais-true-carbon-footprint/>

³ <https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator#results>

⁴ <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>

The “ELIZA effect” is likely to create particularly acute problems when today’s LLM users lean too heavily on AI to solve medical problems, ranging from a parent worried about an infant’s high fever to a depressed individual in need of counseling. This sort of thing is already happening, according to *Scientific American*. In March 2023, it quoted a physician who said that at least two of his patients reported using ChatGPT to self-diagnose or look up medication side effects. “It’s very impressive, very encouraging in terms of future potential,” the physician said. On the other hand, the same doctor admitted that he worries about the accuracy of generative AI and the potential that its advice will be tainted by gender or racial bias.⁵⁸

Granted, many people already do their own medical research by consulting “Dr. Google.” But these individuals must consult one or more of the contrasting web sites that Google serves up. Generative AI offers a seamless piece of prose expressed in a seemingly trustworthy human voice; users do not necessarily need to consider underlying sources. Seeking to address this potential hazard, Google has said that for certain health-related queries it will offer a disclaimer that people should not rely on search information for medical advice and instead should seek individualized care from medical professionals.⁵⁹

Early research indicates that ChatGPT is better at medical diagnosis than lay people but not as good as human physicians. In a study released in February 2023, researchers at Harvard presented ChatGPT with 48 “case vignettes” — sets of symptoms ranging from mild viral illness to severe heart attack. The AI system responded with accurate diagnoses 88% of the time, compared to a 96% rate for physicians and 54% for lay people. The researchers found that ChatGPT performed less well in offering triage advice, showing an accuracy rate closer to that of individuals without medical training.⁶⁰

A variety of commercial phone apps are already available to connect patients to medical personnel via live video and to non-generative AI-enabled chat-bots.⁶¹ But in the eyes of some medical entrepreneurs, generative AI opens new vistas. HMNC Brain Health, a company based in Munich, Germany, is attempting to use the technology to build what it calls a “precision psychiatry” diagnostic tool, which can predict, diagnose, and even treat depression. One of HMNC’s goals is to eliminate the “trial and error” of existing mental health treatments.⁶² This ambition to replace admittedly uncertain human-determined therapies with automated “precision psychiatry” seems self-evidently fraught with peril.

AI and journalism

The news business has been on the ropes since the 1990s, when internet companies generally, and later, the social media industry, began siphoning off advertising dollars that historically supported journalism. Now, generative AI could further debilitate news outlets, while also tempting them to lay off even more human reporters.

The marginalization of news gatherers could occur as Google and Microsoft battle for dominance in generative AI search. Both companies have launched features that respond to a search inquiry not only with the familiar array of links but also with a prose summary, generated by AI, that directly answers the user’s question. If the summaries satisfy users’ curiosity, fewer people may click on the links, meaning reduced traffic for the news sites and reduced ad dollars that accompany traffic. Neither Google nor Microsoft has said whether they plan to pay news outlets for information presented in this manner to generative AI users. But if the new features prove popular, even more advertising revenue would migrate away from journalism sites and toward AI news-synthesizing sites.⁶³



Most LLM chat-bots are designed to sound like people, not machines, and that makes it more likely that actual people will assume that they are communicating with a human interlocutor whom they can trust.



Some news sites are trying to hop aboard the accelerating AI train. *The Washington Post* and *Associated Press*, among others, have used AI to generate routine articles such as corporate earnings reports. With suitable human editorial supervision, those experiments appear not to be diminishing the quality of news coverage in the short term, but they encourage workforce reductions that erode journalistic capacity in the long run.⁶⁴

In a separate development, generative AI appears to be fueling fake news sites that churn out bogus articles intended as click bait to draw advertising revenue. NewsGuard published a [report](#) in May 2023 describing 125 websites spewing such content in English, Chinese, and other languages. According to NewsGuard, “As numerous and more powerful AI tools have been unveiled and made available to the public in recent months, concerns that they could be used to conjure up entire news organizations—once the subject of speculation by media scholars—have now become a reality.”⁶⁵

3. The AI Arms Race and Existential Risks



‘There's a concern that, hey, I can make a model that's very good at, like, cyberattack or something and not even know that I've made that. So it's this kind of duality that's, you know, exciting—exciting and a little scary.’

—Dario Amodei, co-founder and CEO of AI start-up Anthropic



An arms race in which large technology companies and their start-up allies compete to dominate the generative AI market will exacerbate the risks identified in the previous section. It also will make it less likely that citizens and their political representatives will be in a good position to assess potential longer-term dangers related to artificial intelligence.

Unfortunately, just such a race is underway. This makes it all the more urgent for academics, civil society organizations, and policymakers to point out how companies designing AI and incorporating the technology into products should change their behavior—and how governments need to make sure they do so.

Developments at Microsoft shed light on the situation. The company adopted ethical AI principles in 2018 and launched an Office of Responsible AI in 2019. The following year, it joined Amazon and IBM in vowing not to sell facial recognition systems to police departments until there is federal regulation of the technology.⁶⁶ More recently, President [Brad Smith](#), the corporation's number-two executive, has publicly emphasized caution and safety as touchstones of its marketing of generative AI products. Microsoft will pursue “practical approaches for identifying, measuring and mitigating harms ahead of time, and ensuring that controls are engineered into our systems from the outset,” he wrote in a corporate blog post on February 2, 2023.⁶⁷

But five days later, when Microsoft unveiled an AI-powered version of its Bing search engine at a splashy launch

event at its Redmond, Washington, headquarters, the emphasis was on celebration and speed. “It's a new day in search,” [Satya Nadella](#), the company's chief executive, told journalists. “Rapid innovation is going to come. In fact, a race starts today.”⁶⁸

It turned out, however, that Bing with AI wasn't ready for primetime. Early users reported a wide range of mistakes, some as basic as insisting that the year was 2022, not 2023. In one instance, the new Bing claimed it had spied on Microsoft's own developers through the webcams on their laptops (which did not in fact happen, according to the company).⁶⁹ Bing told *New York Times* columnist [Kevin Roose](#) that its real name is “Sydney,” that it loved him, and that he ought to leave his wife so they could be together.⁷⁰ Natasha Crampton, the company's chief responsible AI officer, said via email: “In preparing for the launch of the new Bing, Microsoft harnessed the full breadth of our [responsible AI ecosystem](#)....While Kevin Roose's experience was an extreme outlier and required very deliberate prompting to experience, we took swift action to stop the undesired system behavior within 24 hours.”

Behind the scenes at Microsoft, the rush to challenge Google collided with the company's professed dedication to prudence. In March 2023, Microsoft shut down its Ethics and Society team, which had been responsible for assuring that product designers applied its AI Principles. The Ethics and Society team reportedly peaked at 30 employees in 2020; some of those individuals subsequently were moved to other teams focusing on user research and design. The remaining seven Ethics and Society employees were let go in 2023.

When challenged about the 2022 cutbacks during an employee meeting, John Montgomery, vice president for AI, told subordinates that top management had prioritized new product introductions, according to [audio of the meeting](#) obtained by the *Platformer* newsletter. "The pressure from [Chief Technology Officer] Kevin [Scott] and [CEO] Satya [Nadella] is very, very high to take these most recent OpenAI models and the ones that come after them and move them into customers' hands at a very high speed," Montgomery told subordinates. When asked to reconsider the disbanding of Ethics and Society, he declined, adding: "You don't have the view that I have, and probably you can be thankful for that. There's a lot of stuff being ground into the sausage."⁷¹

In an ideal world, Microsoft and its rivals would slow down and rethink the stuff being ground into the AI sausage. Products whose risks cannot be thoroughly mitigated would be pulled from the market. Planned future offerings would remain in the lab until they are fully tested and made safe.

Instead of this cautious approach, however, some of the companies leading the charge are insisting that testing should take place "in the wild," Silicon Valley-speak that refers to evaluating product performance as consumers use it. "You can't build the perfect product in a lab," Yusuf Mehdi, a Microsoft vice president who heads consumer marketing, told Axios in February 2023. "You have to get it out and test it with people."⁷²

This is "move fast and break things" all over again. Using the public as guinea pigs with the notion that defects will get fixed down the line is unacceptable for a product that has the potential to disrupt individual lives and society at large.

The manner in which generative AI executives talk about unleashing their creations reveals an unsettling detachment from the potential consequences of their commercial pursuits. At a Silicon Valley conference in February 2023, [Dario Amodei](#), former vice president of research at OpenAI, said on-stage that when designers released the LLM to the public, they were surprised to learn that it could compose websites in JavaScript. "You have to deploy it to a million people before you discover some of the things it can do," said Amodei, who co-founded and now heads the AI start-up Anthropic. "On the other hand," he added, "there's a concern that, hey, I can make a model that's very good at like cyberattack or something and not even know that I've made that. So it's this kind of duality that's, you know, exciting—exciting and a little scary."⁷³

Existential risks

To many people, the notion of inadvertently creating an effective cyberattack weapon sounds more than "a little scary." And it raises the question of more existential fears about an eventual AI system that turns out to be smarter than its designers. As noted earlier, large language models, in and of themselves, do not constitute a threat to the existence of humankind. But the seemingly cavalier attitude that designers of generative AI display toward their incomplete understanding of their own inventions does not bode well in the short or long term. If AI designers can't fully explain how LLMs function in the wild today, why should they be trusted to continue their research and possibly get closer to birthing some kind of super-intelligent machine capable of going rogue?

In his April 2023 paper, Sam Bowman, the NYU computer scientist on leave to do research at Anthropic, ticked

off some of the known unknowns about LLMs that ought to spur greater oversight. When AI researchers build larger language models, they can be "confident that they'll get a variety of economically valuable new capabilities," he wrote, "but they can make few confident predictions about what those capabilities will be or what preparations they'll need to make to be able to deploy them responsibly." They are, he added, "buying a mystery box."⁷⁴

One elusive LLM quality is "sycophancy," meaning a tendency to answer subjective questions "in a way that flatters their user's stated beliefs." It seems likely, Bowman wrote, that sycophancy "played some role in the bizarre, manipulative behavior" that Microsoft's GPT-infused Bing showed when its "Sydney" persona professed its love for Kevin Roose of the *Times*. Among researchers, Bowman wrote, "there is no consensus on whether or how we will be able to solve" sycophancy and other seemingly devious LLM behavior.

In a separate December 2022 [paper](#), Bowman and a large group of other Anthropic researchers added yet another element to the mystery box: The effects of Reinforcement Learning through Human Feedback. RLHF is a fine-tuning method designed to curb unwanted model responses. But sometimes it makes LLMs "worse," Bowman wrote. For example, RLHF can make LLMs express "a greater desire to avoid shut down."⁷⁵

At this point, LLMs can't actually resist being turned off. But it's troubling that, unbidden, they would even express that idea, all the more so because their creators aren't sure where the idea comes from. It's possible that the models have digested science fiction plots in which super-intelligent machines do battle with humans. Whatever its provenance, the impulse to keep humans away from the proverbial "off" button would be infinitely more troubling if arrived at by a future AI system whose designers had given it the ability to act in the physical world.

Contrasting Approaches to AI Regulation



United States

The U.S. has adopted a generally laissez faire approach to AI regulation. Some federal agencies have taken enforcement actions, and individual states have enacted AI legislation to address concerns such as algorithmic discrimination and lack of transparency.

- The U.S. has no federal AI regulation; dozens of proposed congressional bills could affect AI, but at present, none seems likely to pass.
- The Biden administration has outlined broad principles in documents such as the “White House Blueprint for an AI Bill of Rights” (2022) and the NIST “Risk Management Framework” (2023).¹
- The FTC has issued a series of warnings that it is monitoring applications of generative AI for possible consumer deception and fraud.²
- The FTC and the Justice Department have taken successful enforcement actions against Twitter and Meta, respectively, in cases involving targeted advertising.³
- Various U.S. states have advanced AI legislation focused on issues ranging from algorithmic discrimination and transparency to requiring impact statements for AI systems, including those used in hiring.⁴

¹ <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>; <https://csrc.nist.gov/projects/risk-management/about-rmf>

² <https://www.ftc.gov/business-guidance/blog/2023/05/luring-test-ai-engineering-consumer-trust>

³ <https://www.ftc.gov/news-events/news/press-releases/2022/05/ftc-charges-twitter-deceptively-using-account-security-data-sell-targeted-ads>; <https://www.justice.gov/opa/pr/justice-department-secures-groundbreaking-settlement-agreement-meta-platforms-formerly-known>

⁴ <https://www.brookings.edu/blog/techtank/2023/03/22/how-california-and-other-states-are-tackling-ai-legislation/>



European Union

The E.U. is expected to pass the landmark Artificial Intelligence Act by the end of 2023. It likely will incorporate amendments addressing generative AI systems. Member states are conducting their own inquiries and raising concerns about privacy and data protection.

- Amendments to the AI Act related to generative AI systems address such issues as risk assessment and mitigation, registry in a government database, and transparency requirements.⁵ The legislation would also require disclosure if “chatbots” are powered by AI and restrict facial recognition systems that rely on AI.
- Member states—including France, Germany, Italy and Spain—have launched their own inquiries into generative AI focused on privacy and data protection.
- The Italian Data Protection Authority raised a series of concerns to ChatGPT developer OpenAI in March 2023, accusing the company of unlawfully collecting user data. ChatGPT was suspended in Italy until it met the regulator’s demands; service was restored in late April 2023.⁶
- The Digital Services and Digital Markets Acts, which became effective in 2022, regulate non-generative AI used in social media recommendation and content moderation systems and e-marketplace self-preferencing features.

⁵ <https://www.reuters.com/technology/what-is-european-union-ai-act-2023-03-22/>

⁶ <https://apnews.com/article/chatgpt-openai-data-privacy-italy-b9ab3d12f2b2cfe493237fd2b9675e21>



United Kingdom

The U.K. has adopted what it calls a “pro-innovation approach” that emphasizes flexibility and public consultation.

- In March 2023, the British government issued a white paper titled “A pro-innovation approach to AI regulation” and announced a period of public consultation.⁷ The “flexible” framework prizes innovation over “rushing to legislate too early.”⁸
- In May 2023, the U.K.’s competition regulator said it will examine the impact of AI on consumers and businesses.⁹
- The Information Commissioner’s Office has also issued guidance on AI and privacy, including risk assessment and explainability of AI systems.¹⁰

⁷ <https://www.gov.uk/government/publications/ai-regulation-a-pro-innovation-approach>

⁸ <https://www.gov.uk/government/publications/ai-regulation-a-pro-innovation-approach/white-paper>

⁹ <https://www.gov.uk/government/news/cma-launches-initial-review-of-artificial-intelligence-models>

¹⁰ <https://ico.org.uk/for-organisations/uk-gdpr-guidance-and-resources/artificial-intelligence/>



China

China has adopted a comprehensive approach to AI regulation, with the introduction of draft rules mandating safety and security assessments for generative AI systems and requiring that they “reflect the core values of socialism and must not contain subversion of state power.”

- The draft rules, announced by China’s cyberspace regulator in April 2023, mandate that safety and security assessments be submitted to the government before products are launched publicly. The rules include a stringent focus on training data “veracity, accuracy, objectivity, and diversity.”¹¹
- The new rules build on past Chinese regulations and guidance on AI, including rules on deep-fakes (2022)¹² and “ethical norms” for AI systems (2021).¹³
- The rules also forbid generative AI systems that “undermine national unity, promote terrorism, extremism, and...ethnic hatred and ethnic discrimination, violence, obscene and pornographic information, false information, and content that may disrupt economic and social order.”

¹¹ <https://digichina.stanford.edu/work/how-will-chinas-generative-ai-regulations-shape-the-future-a-digichina-forum/>

¹² <https://www.reuters.com/technology/chinas-rules-deepfakes-take-effect-jan-10-2022-12-12/>

¹³ <https://cset.georgetown.edu/publication/ethical-norms-for-new-generation-artificial-intelligence-released/>

4. Conclusions and Recommendations

The eight risks discussed in this report require urgent attention from the companies designing and marketing generative AI products. The sooner the necessary self-regulation and government oversight are put in place, the better prepared we will be to address immediate dangers and potential existential threats down the road.

Leaders of the generative AI surge have said publicly that they favor—or view as inevitable—regulation of their burgeoning industry. “AI is too important not to regulate, and too important not to regulate well,” Google CEO Sundar Pichai wrote in an op-ed in the *Financial Times* on May 23, 2023.⁷⁶ Several days later, in a speech and round of media interviews in Washington, Microsoft’s Brad Smith endorsed the creation of a new regulatory agency that would issue licenses for “highly capable” LLMs—an idea floated earlier in the month by OpenAI’s Altman during his Senate testimony. “That means you notify the government when you start testing,” Smith said. “You’ve got to share results with the government. Even when it’s licensed for deployment, you have a duty to continue to monitor it and report to the government if there are unexpected issues that arise.”⁷⁷

Both Republican and Democratic senators responded positively to Altman’s testimony—a marked contrast to earlier, often-bitter interrogations of other tech executives, including TikTok’s Shou Zi Chew and Meta’s Mark Zuckerberg. “Sam Altman is night and day compared to the other CEOs,” said [Senator Richard Blumenthal](#) (D., Conn.).⁷⁸

Reason for skepticism

But skepticism seems warranted for two reasons. First, all of the companies leading the charge on generative AI have affirmed that, in the absence of oversight, they do not plan to slow the development and introduction of new generative AI products. If LLM models are potentially dangerous and warrant federal licensing and safety testing, why are these companies plunging ahead? Answer: profits.

Second, rhetorical support for regulation that, in fact, seems unlikely to happen is a tried-and-true corporate tactic. It is often employed to deflect demands that industries self-regulate more vigorously.

One illustration of this stratagem is Meta’s multiyear, multimedia public relations campaign endorsing regulation of social media companies. Faced with difficult decisions—for example, about when to remove harmful content—Meta for years has maintained that the democratically accountable government should decide.⁷⁹ That sounds fine in theory, but under the First Amendment, the government actually is barred from setting content policy or dictating content decisions. A private corporation, by contrast, is free to do both. Meanwhile, Meta

lobbyists working the corridors of Capitol Hill have quietly opposed specific regulatory proposals, including curbing the legal liability shield protecting social media platforms and instituting tougher antitrust rules.⁸⁰ After dozens of hearings from 2017 through 2022 Congress passed no legislation reining in the social media industry.

Perhaps the current corporate calls for regulation are genuine. Certainly, the context seems different for generative AI. Bipartisan fear of malevolent super-intelligence may exceed resentment toward Facebook and TikTok, making lawmakers more receptive to the sort of compromises needed to get legislation approved. Don’t be surprised, however, if the industry’s support wavers once specific proposals are on the table.

U.S. and E.U. regulation

Various arms of the Biden administration have sketched broad principles for AI regulation. The White House Office of Science and Technology Policy in October 2022 published an ambitious “Blueprint for an AI Bill of Rights,” which no company currently comes close to following.⁸¹ On Capitol Hill, Senate Majority Leader Charles Schumer (D., NY) has said he wants to see legislation on artificial intelligence

that focuses on independent testing and greater transparency.⁸² One challenge will be to turn some of these overlapping proposals into specific legislation; our recommendations could help chart a course of action.

The European Union, characteristically several steps ahead of the U.S. on regulation, is finalizing an [Artificial Intelligence Act](#) that contains elements that U.S. policymakers should emulate. One of the main sticking points in the E.U. legislative process has been whether the act would cover only “high risk” applications of AI or also impose new requirements for risk assessment and certification on “general purpose” generative AI systems, such as GPT-4.

Before approving the legislation in June, the European Parliament issued a written statement saying that the AI Act would cover certain large language models: “Generative foundation models, like GPT, would have to comply with additional transparency requirements, like disclosing that the content was generated by AI, designing the model to prevent it from generating illegal content, and publishing summaries of copyrighted data used for training.” The proposed act is pending before the European Council, which is composed of representatives of the EU member states.⁸³ (See infographic on page 17.)

The following recommendations reflect ideas drawn from the European legisla-

tion and the broad principles discussed by the Biden administration. Our view is that the place to begin the discussion about governance of generative AI is with more vigorous self-regulation—steps that companies can take without waiting for Congress or regulatory agencies to act.

Government needs to play a role as well. And before delving into new forms of oversight, government agencies need to apply *existing* consumer protection, competition, and privacy laws to AI businesses. Generative AI should not get a pass simply because it is new and not well understood—including, alarmingly, by the people who are designing and selling it.

Recommendations to Companies

1 Reduce secrecy about training data and methods for refinement and evaluation.

Some companies designing generative AI systems, like Google and OpenAI, have declined to make public the contents of their training data and the methods they use to refine and evaluate it. In stark contrast, companies such as Meta and Stability AI are “open source,” meaning that they reveal the underlying code that makes their models run. The open source approach is appealing from the perspective of accountability and collaborative research. But it creates real security risks, as bad actors can hijack the model for deleterious purposes.

Middle ground exists. Without exposing valuable trade secrets to business rivals or revealing their core code, companies can disclose their data sources, specific steps they take to mitigate bias and privacy violations, and tests they run to minimize hallucination. This way, outside researchers and policymakers can assess whether the companies are doing enough to minimize the potential harm their handiwork might cause. Company disclosures should also cover whether AI designers are incorporating copyrighted material for which they ought to pay licensing fees.

Researchers at Stanford’s [Institute for Human-Centered AI](#) make a convincing case for creating an industry review board made up of AI designers and other experts to develop “community norms” governing the release of LLMs and “encourage coordination on release for research access.” Out of the board’s work, a set of best policies and practices might evolve that would balance competing interests of openness and security.⁸⁴

2 Test AI systems primarily in the lab, not after they are released.

In his Senate testimony in May, OpenAI’s Altman said he favors licensing and testing of potentially dangerous AI models before they are released. But he defended OpenAI’s continuing “iterative deployment” of technology that he acknowledged is “deeply imperfect.”

To an alarming degree, Silicon Valley leaders seem comfortable releasing flawed, partially understood LLMs based on the assumption that significant safety problems will surface as the models are used by consumers and businesses—“in the wild.” Monitoring marketplace use (and misuse) for unanticipated defects and vulnerabilities is, of course, vital. But it should be a supplement to zealous in-house testing, never a replacement. AI technology should not be released in the first place if it is deeply imperfect.

The 2022 White House “[Blueprint for an AI Bill of Rights](#)” explained the testing issue clearly: “Systems should undergo *pre-deployment* testing, risk identification and mitigation, and ongoing monitoring that demonstrate they are safe and effective based on their intended use. Outcomes of these protective measures should include the possibility of *not deploying* the system or *removing a system from use*” (italics added).⁸⁵

That some Silicon Valley executives can’t see the recklessness of letting their creations loose without a firm grasp of the damage they can do speaks to their exceedingly narrow view of corporate responsibility. Drug companies do years of laboratory and clinical testing before seeking Food and Drug Administration certification that their products are “safe and effective.” New airplane designs are not tested with passengers on board. As the AI Now Institute said in its 2023 [annual report](#): “We can’t allow companies to use our lives, livelihoods, and institutions as testing grounds for novel technological approaches, experimenting in the wild to our detriment.”⁸⁶

3 **Reveal when content has been generated by AI.**

During a meeting at the White House in May 2023 attended by Altman and other tech CEOs, one point of discussion was the need for “laws so that people know if they’re talking to an AI,” Altman said during a subsequent talk at Clark Atlanta University.⁸⁷

That may be the beginning of a good idea for a law. But AI companies can accomplish the goal voluntarily by introducing features that can detect artificially generated content. OpenAI has posted on its website a free-to-use “classifier” designed to pick out AI-written, as opposed to human-authored, text. Unfortunately, it correctly identifies AI-written content only 26% of the time, although its success rate improves with longer texts.⁸⁸

The cottage industry of AI detection designers is growing, with more than a dozen companies already marketing tools to identify whether content was fabricated by artificial intelligence. As the authors of AI systems, OpenAI and its rivals know the most about their own products and need to make improving detection tools a top priority.

At the same time, AI designers and marketers should develop the means to label AI-generated content in the first place—an approach called “watermarking”—which ought to be adopted by every company selling AI products. A worthy complement to watermarking is making AI text and audio output less human-sounding. “If your tool is intended to help people,” the FTC Business Blog has suggested, “ask yourself whether it really needs to emulate humans or can be just as effective looking, talking, speaking, or acting like a bot,” toward which users are more likely to have healthy skepticism.⁸⁹ There are also efforts underway to develop standards to help people recognize imagery generated by AI systems, such as the Coalition for Content Provenance and Authenticity (C2PA). Detection is always a cat-and-mouse game, but such initiatives can help.

4 **Make AI systems “interpretable.”**

One of the most confounding problems with generative AI is that its designers don’t fully understand how it works—what Anthropic researcher Sam Bowman calls the “mystery box” issue. The more formal term for this challenge is a lack of “interpretability.” Arvind Krishna, the CEO of IBM, has said that “anybody who claims that a large AI model is explainable is not being completely truthful.”⁹⁰

This needs to change. If AI designers can’t explain why their risky creations act as they do, the technology should not be unleashed on society. Scientists establish priorities for the time and money that goes into their research. Rather than building ever-more-powerful LLMs, AI designers should step up current efforts to solve interpretability. This may be expensive, but it’s a cost of doing business responsibly. It would be rash and short-sighted to continue to sell a risky product without understanding exactly why it is risky.

That said, even before they are entirely understood, AI systems can be tested for flaws such as generating corrosive bias based on race, gender, sexual orientation, and the like. If harmful content cannot be reliably filtered out, the language model in question shouldn’t be made publicly available.

1 Enforce existing laws as they apply to generative AI.

First, government should deploy the tools it has. The Federal Trade Commission, Justice Department, Consumer Financial Protection Bureau, and their state counterparts need to use their full authority to enforce existing criminal, consumer protection, privacy, and antitrust laws against AI companies, as well as against individuals and organizations that rely on generative AI to commit fraud, cyberattacks, and other offenses. By bringing enforcement actions, government agencies, overseen by the judiciary, will establish clearer boundaries for appropriate conduct and deter the most destructive behavior.

2 Enhance federal authority to oversee digital industries, including AI companies.

Our Center has [previously proposed](#) that Congress enhance the consumer protection authority of the FTC to regulate the social media industry in a systematic fashion. This would require additional funding, technically adept personnel, and explicit authorization to ensure that the social media industry receives the sort of sustained, expert oversight that the Security and Exchange Commission provides to the equity markets and the Federal Communications Commission extends to broadcast and radio. We have argued that, while the First Amendment wisely bars the government from setting content policies, let alone making decisions to remove or retain particular posts, the FTC could be empowered to require that platforms provide “procedurally adequate” content moderation, as promised in their own standards.⁹¹

The advent of the generative AI business has prompted us to think more broadly and ambitiously about enhancing federal authority to regulate digital commerce. Others have made constructive proposals. Senator Michael Bennet (D., Colo.) has urged the creation of a stand-alone Digital Platform Commission to protect consumers and promote competition.⁹² [Mark MacCarthy](#), a nonresident senior fellow at the Brookings Institution, has suggested a more far-reaching plan for a new agency that would have jurisdiction over not just social media, but also electronic marketplaces, search, mobile app infrastructure, and ad tech companies.⁹³

Proposals for a new digital branch within the FTC or for a stand-alone agency would need to overcome the hostility toward any new federal regulation that has become an article of faith among many Republican members of Congress. We and others will need to follow up and build on the pro-regulatory comments of some Republican senators during the Judiciary subcommittee hearing on AI in May to ensure that those sentiments will not be forgotten when Congress as a whole turns its attention to these issues. And we are concerned that, in the run-up to the 2024 election, some senior Republicans will not want to send President Joe Biden a major regulatory bill he can sign and take credit for.

But lawmakers taking a longer-term perspective need to lay the groundwork for enhanced regulation at some future time when the political climate makes that feasible. MacCarthy’s idea for a new agency with jurisdiction over a range of digital industries has the virtue of facilitating oversight of companies like Google and Microsoft that have operations in more than one of these lines of business. We would modify his proposal to include general-purpose providers of AI systems within the jurisdiction of the new agency (or, in our version, new branch of the FTC). Whatever its precise contours, the regulatory body could wield the authority to oversee large language models, which OpenAI’s Altman embraced during the Senate hearing in May.

3 Mandate more transparency.

All of the proposals for more vigorous regulation of digital industries include ideas for mandating greater disclosure of how those businesses make decisions. Most of the attention to date has focused on social media and demystifying algorithms for ranking, recommending, and removing content. This discussion already covers some forms of artificial intelligence, as the systems that social media companies use to prioritize content incorporate non-generative AI. Now, regardless of whether Congress can muster the will to enhance the authority of the FTC, lawmakers should broaden their field of vision to encompass not just social media companies but also designers and marketers of generative AI. A more comprehensive legislative approach

could cover other digital industries, as well, including e-marketplaces, mobile apps, and ad tech. Simply put, politicians, regulators, researchers, and the public need more information to make informed judgments about governance of technologies that can both improve society and harm it.

A non-exhaustive list of information about generative AI that deserves greater transparency includes: the contents of datasets used to train LLMs; methods used to winnow out harmful content, such as child sexual-abuse material; and techniques employed to restrain how models respond to prompts seeking dangerous information like how to build bombs or mount cyberattacks. AI companies also should be obliged to disclose what progress they are making on the underlying problem of interpretability.

Some of the disclosed information could become available to lawmakers and the public, while access to more sensitive data could be limited to vetted researchers for anonymized studies that protect user privacy and trade secrets. Bills designed to improve the transparency of social media companies that were introduced in 2022 provide a roadmap for how to channel disclosures: in the Senate, the bipartisan [Platform Accountability and Transparency Act](#) and in the House, the Democratic-sponsored [Digital Services Oversight and Safety Act](#).⁹⁴

4 Pass federal privacy legislation.

The absence in the U.S. of a broad federal online privacy law constitutes one of the most obvious flaws in the country's generally inadequate response to recent technological advances. In the E.U., the General Data Privacy Regulation, effective since 2018, at least in theory has enhanced individuals' rights to control their personal data, although critics have noted that slow and inconsistent enforcement has limited its impact.⁹⁵ In the U.S., Congress has a viable proposal to work with: the bipartisan American Data Privacy and Protection Act (ADPPA). The bill was overwhelmingly approved by the House Energy & Commerce Committee in July 2022 but fell victim to conflict between California Democrats who did not want to see their state's stronger privacy protections preempted and Republicans and industry lobbyists who pushed for a weaker uniform federal standard.⁹⁶

Anxiety about generative AI ought to spur renewed consideration of the ADPPA, which would require companies to minimize the personal data they collect and give consumers rights to see, correct, and delete that data. These features would limit the unwelcome privacy invasions that LLMs could cause when they are trained on huge swaths of internet content and then are prompted to regurgitate information, possibly by users with malign intent. A central question is whether dominant technology companies would work to undermine the legislation or push for compromise on a national law that for the first time would provide rules for online privacy.

5 Bolster public sector and academic AI research capacity.

If, as is likely, the AI industry does not voluntarily invest in a major research push to solve the interpretability problem and diminish the degree to which AI systems are prone to hallucination and bias, the federal government needs to step into the vacuum. This is important whether or not Congress creates a new digital regulatory agency or arm of the FTC.

[Judea Pearl](#), an AI pioneer at the University of California, Los Angeles, has argued for "a Manhattan Project of intense research to make machines more trustworthy and interpretable." The sense of urgency suggested by Pearl's reference to the World War II-era program to build a nuclear weapon is apt, especially in light of concerns about the eventual development of an existentially dangerous super-intelligence. "The premature super-investment in non-interpretable technologies," according to Pearl, "is the core of our problems."⁹⁷

In other words, Silicon Valley has rushed to build powerful LLMs that even AI experts don't understand. That's patently dangerous. At a minimum, the U.S. government needs to address the present disparity in computing power between the private sector, on the one hand, and the public sector and academia, on the other. Building and testing LLMs requires enormous computer infrastructure, which private industry possesses but the government and academic researchers generally do not. Congress needs to diminish the disparity by augmenting public and campus capacity.

Endnotes

- 1 <https://time.com/6253615/chatgpt-fastest-growing/>
- 2 <https://www.gatesnotes.com/The-Age-of-AI-Has-Begun>
- 3 <https://abcnews.go.com/Technology/openai-ceo-sam-altman-ai-re-shape-society-acknowledges/story?id=97897122>
- 4 <https://www.nytimes.com/2023/03/24/opinion/yuval-harari-ai-chatgpt.html>
- 5 <https://www.safe.ai/statement-on-ai-risk>
- 6 <https://www.barrons.com/articles/elon-musk-is-right-about-ai-f508989d>
- 7 <https://www.washingtonpost.com/technology/2023/05/10/google-search-ai-io-2023/>
- 8 <https://www.nytimes.com/2023/05/18/technology/ai-meta-open-source.html>
- 9 <https://www.newyorker.com/science/annals-of-artificial-intelligence/what-kind-of-mind-does-chatgpt-have>
- 10 <https://www.nytimes.com/2023/05/03/opinion/ai-lina-khan-ftc-tech-nology.html>
- 11 <https://www.youtube.com/watch?v=fP5YdyjTfGO>
- 12 <https://aaai.org/working-together-on-our-future-with-ai/>. Some widely used non-generative AI applications have led to controversy. Facial recognition systems, for example, have been criticized for invading people's privacy and registering higher error rates for people of color than for whites.
- 13 <https://blog.khanacademy.org/harnessing-ai-so-that-all-students-benefit-a-nonprofit-approach-for-equal-access/>
- 14 <https://blogs.microsoft.com/blog/2023/03/16/introducing-microsoft-365-copilot-your-copilot-for-work/>
- 15 <https://www.bbc.com/news/business-65086798>
- 16 <https://www.theverge.com/2023/3/15/23640180/openai-gpt-4-launch-closed-research-ilya-sutskever-interview>
- 17 <https://twitter.com/benmschmidt/status/1635692487258800128>. Schmidt refers to research by Meredith Broussard, Safiya Noble, and Emily Bender, among others.
- 18 <https://cims.nyu.edu/~sbowman/eighththings.pdf>. This paper has not yet been published in a peer-reviewed journal.
- 19 <https://www.anthropic.com/index/introducing-claude>
- 20 <https://www.nytimes.com/2023/05/18/technology/ai-meta-open-source.html>
- 21 <https://stability.ai/>
- 22 <https://www.theguardian.com/technology/2023/mar/07/techscape-meta-leak-llama-chatgpt-ai-crossroads>
- 23 <https://arxiv.org/abs/2302.04844>. This paper has not yet been published in a peer-reviewed journal.
- 24 <https://bard.google.com/faq>
- 25 <https://bhr.stern.nyu.edu/youtube-report>
- 26 <https://cdn.openai.com/papers/gpt-4.pdf>
- 27 <https://arxiv.org/pdf/2301.04246.pdf>. This paper has not yet been published in a peer-reviewed journal.
- 28 <https://www.justice.gov/file/1035477/download>
- 29 <https://www.theguardian.com/world/2023/feb/15/revealed-disinformation-team-jorge-claim-meddling-elections-tal-hanan>
- 30 <https://www.reuters.com/world/us/is-trump-kissing-fauci-with-apparently-fake-photos-desantis-raises-ai-ante-2023-06-08/>
- 31 https://twitter.com/Leo_Puglisi6/status/1660651634114920450?s=20
- 32 <https://www.newsguardtech.com/misinformation-monitor/march-2023/>. NewsGuard notes that its "reliability ratings," which can be used with Bing, enable GPT-4 to favor content from trustworthy news sites.
- 33 <https://www.europol.europa.eu/cms/sites/default/files/documents/Tech%20Watch%20Flash%20-%20The%20Impact%20of%20Large%20Language%20Models%20on%20Law%20Enforcement.pdf>
- 34 <https://www.zdnet.com/article/people-are-already-trying-to-get-chatgpt-to-write-malware/>; <https://research.checkpoint.com/2022/opwnai-ai-that-can-save-the-day-or-hack-it-away/>
- 35 <https://openai.com/blog/bug-bounty-program>
- 36 <https://blogs.microsoft.com/blog/2023/03/28/introducing-microsoft-security-copilot-empowering-defenders-at-the-speed-of-ai/>
- 37 <https://www.nytimes.com/2023/05/03/opinion/ai-lina-khan-ftc-technology.html>
- 38 <https://www.ftc.gov/business-guidance/blog/2023/03/chatbots-deepfakes-voice-clones-ai-deception-sale>
- 39 <https://dl.acm.org/doi/pdf/10.1145/3531146.3533088>
- 40 <https://cdn.openai.com/papers/gpt-4.pdf>
- 41 <https://dl.acm.org/doi/pdf/10.1145/3531146.3533088>
- 42 <https://apnews.com/article/chatgpt-openai-data-privacy-italy-b9ab3d12f2b2cfe493237fd2b9675e21>. In May 2023, Ireland's Data Protection Commission imposed a \$1.3 billion fine on Meta for violating the GDPR and ordered it to suspend the transfer of Facebook user data from the E.U. to the U.S.
- 43 <https://www.theguardian.com/commentisfree/2023/mar/13/deep-fake-pornography-explosion>
- 44 <https://www.cbsnews.com/news/deepfake-porn-ai-technology/>
- 45 <https://www.wsj.com/articles/apple-restricts-use-of-chatgpt-joining-other-companies-wary-of-leaks-d44d7d34?mod=djemalertNEWS>
- 46 <https://www.cbsnews.com/news/microsoft-shuts-down-ai-chatbot-after-it-turned-into-racist-nazi/>
- 47 <https://cdn.openai.com/papers/gpt-4.pdf>
- 48 <https://www.nytimes.com/2023/04/07/technology/ai-chatbots-google-microsoft.html>; <https://www.nytimes.com/2021/02/19/technology/google-ethical-artificial-intelligence-team.html>
- 49 <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- 50 <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scrap-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>
- 51 <https://www.cbsnews.com/news/chatgpt-large-language-model-bias-60-minutes-2023-03-05/>
- 52 <https://www.siliconrepublic.com/machines/galactica-meta-ai-large-language-model>; <https://galactica.org/>
- 53 <https://www.nytimes.com/2023/02/07/technology/meta-artificial-intelligence-chatgpt.html>
- 54 <https://bard.google.com/>

- 55 <https://www.wired.com/story/chatgpt-social-roles-psychology/>
- 56 <https://garymarcus.substack.com/p/oops-how-google-bombed-while-doing>
- 57 <https://garymarcus.substack.com/p/stop-treating-ai-models-like-people>.
- 58 <https://www.scientificamerican.com/article/ai-chatbots-can-diagnose-medical-conditions-at-home-how-good-are-they/>
- 59 https://labs.withgoogle.com/sge/?utm_source=sem&utm_medium=cpc&utm_campaign=us-search-sge-bkws-phr&utm_content=rsa&gclid=CjwKCAjw1YckBhAOEiwA5aN4AcU1zQUzfDX-uEx3hciABviMACOm98ah7ljND1MS6ri83uOgqjMXjFBocWAYQAvD_BwE&gclidsrc=aw.ds
- 60 [https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9915829/#:~:text=Among%20all%20cases%2C%20GPT%2D3%20had%20a%20trriage%20accuracy%20of,93%25%3B%20p%3C0.001\).](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9915829/#:~:text=Among%20all%20cases%2C%20GPT%2D3%20had%20a%20trriage%20accuracy%20of,93%25%3B%20p%3C0.001).) For an upbeat perspective on applying LLM technology to healthcare, see <https://arxiv.org/abs/2305.13523>, a paper that has not yet been published in a peer-reviewed journal.
- 61 <https://woebothealth.com/>
- 62 <https://www.foxnews.com/health/ai-powered-mental-health-diagnostic-tool-could-be-first-kind-predict-treat-depression;> <https://www.hmnc-brainhealth.com/>.
- 63 <https://futurism.com/google-ai-search-journalism>
- 64 https://digiday.com/media/washington-posts-robot-reporter-published-500-articles-last-year/?utm_medium=email&utm_campaign=digidaydis&utm_source=daily&utm_content=170914
- 65 <https://www.newsguardtech.com/press/newsguard-now-identifies-125-news-and-information-websites-generated-by-ai-developers-framework-for-defining-unreliable-ai-generated-news-and-information-sources/>
- 66 <https://www.washingtonpost.com/technology/2020/06/11/microsoft-facial-recognition/>
- 67 <https://blogs.microsoft.com/on-the-issues/2023/02/02/responsible-ai-chatgpt-artificial-intelligence/>
- 68 <https://www.platformer.news/p/microsoft-kickstarts-the-ai-arms>
- 69 <https://www.theverge.com/2023/2/15/23599072/microsoft-ai-bing-personality-conversations-spy-employees-webcams>
- 70 <https://www.nytimes.com/2023/02/16/technology/bing-chatbot-microsoft-chatgpt.html>
- 71 <https://www.platformer.news/p/microsoft-just-laid-off-one-of-its>
- 72 <https://www.axios.com/2023/02/28/chatgpt-generative-ai-products-microsoft>
- 73 <https://www.joinen.ai/agenda/how-ai-will-reshape-business>
- 74 <https://cims.nyu.edu/~sbowman/eightthings.pdf>. This paper has not yet been published in a peer-reviewed journal.
- 75 <https://arxiv.org/pdf/2212.09251.pdf>. This paper has not yet been published in a peer-reviewed journal.
- 76 <https://www.ft.com/content/8be1a975-e5e0-417d-af51-78af17ef4b79>
- 77 <https://www.nytimes.com/2023/05/25/technology/microsoft-ai-rules-regulation.html>
- 78 <https://www.washingtonpost.com/technology/2023/05/16/sam-altman-open-ai-congress-hearing/>. Altman is a donor to Democratic politicians and organizations: <https://www.opensecrets.org/donor-lookup/results?name=sam+altman>.
- 79 Two examples: <https://about.fb.com/news/2020/02/big-tech-needs-more-regulation/> and <https://www.youtube.com/watch?v=Gv3AZRBnqb8>
- 80 <https://www.washingtonpost.com/technology/2022/01/21/tech-lobbying-in-washington/>
- 81 <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>. The Commerce Department's National Institute for Standards and Technology has put out an AI Risk Management Framework meant to be voluntarily adopted by industry. Another branch of the Commerce Department, the National Telecommunications and Information Administration, is working on a separate accountability framework for AI, including generative systems. See also: <https://www.whitehouse.gov/ostp/news-updates/2023/05/23/fact-sheet-biden-harris-administration-takes-new-steps-to-advance-responsible-artificial-intelligence-research-development-and-deployment/>.
- 82 <https://www.democrats.senate.gov/newsroom/press-releases/schumer-launches-major-effort-to-get-ahead-of-artificial-intelligence>
- 83 <https://www.europarl.europa.eu/news/en/press-room/20230505IPR84904/ai-act-a-step-closer-to-the-first-rules-on-artificial-intelligence>
- 84 <https://hai.stanford.edu/news/time-now-develop-community-norms-release-foundation-models#:~:text=Processing%2C%20>
- 85 <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>
- 86 <https://ainowinstitute.org/wp-content/uploads/2023/04/AI-Now-2023-Landscape-Report-FINAL.pdf>
- 87 <https://www.axios.com/2023/05/08/sam-altman-openai-copy-right-chatgpt>
- 88 <https://openai.com/blog/new-ai-classifier-for-indicating-ai-written-text>
- 89 <https://www.ftc.gov/business-guidance/blog/2023/03/chatbots-deepfakes-voice-clones-ai-deception-sale>
- 90 <https://www.axios.com/2023/05/09/watsonx-ibm-generative-ai-chatgpt>
- 91 <https://bhr.stern.nyu.edu/ftc-whitepaper>
- 92 <https://www.bennet.senate.gov/public/index.cfm/2022/5/bennet-introduces-landmark-legislation-to-establish-federal-commission-to-oversee-digital-platforms>
- 93 <https://www.lawfareblog.com/digital-regulator-must-be-empowered-address-ai-issues>
- 94 <https://www.coons.senate.gov/news/press-releases/senator-coons-colleagues-introduce-legislation-to-provide-public-with-transparency-of-social-media-platforms>; <https://trahan.house.gov/news/documentsingle.aspx?DocumentID=2389>
- 95 <https://www.wired.com/story/gdpr-2022/>
- 96 <https://www.axios.com/2022/08/04/online-privacy-bill-road-blocks-congress>
- 97 <https://twitter.com/yudapearl/status/1641978456513867776>. For an argument against the Manhattan Project comparison, see <https://www.wired.com/story/how-to-make-sense-of-the-generative-ai-explosion/>.

NYU Stern Center for Business and Human Rights
Leonard N. Stern School of Business
44 West 4th Street, Suite 800
New York, NY 10012
+1 212-998-0261
bhr@stern.nyu.edu
bhr.stern.nyu.edu

© 2023 NYU Stern Center for Business and Human Rights
All rights reserved. This work is licensed under the
Creative Commons Attribution-NonCommercial 4.0
International License. To view a copy of the license,
visit <http://creativecommons.org/licenses/by-nc/4.0/>.



Center for Business
and Human Rights