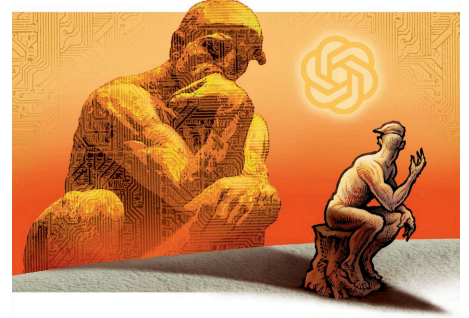# ChatGPT Heralds an Intellectual Revolution

**Generative artificial intelligence presents a philosophical and practical challenge on a scale not experienced since the start of the Enlightenment.**

By Henry Kissinger, Eric Schmidt and Daniel Huttenlocher
Feb. 24, 2023 2:17 pm ET

A new technology bids to transform the human cognitive process as it has not been shaken up since the invention of printing. The technology that printed the Gutenberg Bible in 1455 made abstract human thought communicable generally and rapidly. But new technology today reverses that process. Whereas the printing press caused a profusion of modern human thought, the new technology achieves its distillation and elaboration. In the process, it creates a gap between human knowledge and human understanding. If we are to navigate this transformation successfully, new concepts of human thought and interaction with machines will need to be developed. This is the essential challenge of the Age of Artificial Intelligence.

The new technology is known as generative artificial intelligence; GPT stands for Generative Pre-Trained Transformer. ChatGPT, developed at the OpenAI research laboratory, is now able to converse with humans. As its capacities become broader, they will redefine human knowledge, accelerate changes in the fabric of our reality, and reorganize politics and society.

Generative artificial intelligence presents a philosophical and practical challenge on a scale not experienced since the beginning of the Enlightenment. The printing press enabled scholars to replicate each other's findings quickly and share them. An unprecedented consolidation and spread of information generated the scientific method. What had been impenetrable became the starting point of accelerating query. The medieval interpretation of the world based on religious faith was progressively undermined. The depths of the universe could be explored until new limits of human understanding were reached.

Generative AI will similarly open revolutionary avenues for human reason and new horizons for consolidated knowledge. But there are categorical differences. Enlightenment knowledge was achieved progressively, step by step, with each step testable and teachable. AI-enabled systems start at the other end. They can store and distill a huge amount of existing information, in ChatGPT's case much of the textual material on the internet and a large number of books—billions of items. Holding that volume of information and distilling it is beyond human capacity.

Sophisticated AI methods produce results without explaining why or how their process works. The GPT computer is prompted by a query from a human. The learning machine answers in

literate text within seconds. It is able to do so because it has pregenerated representations of the vast data on which it was trained. Because the process by which it created those representations was developed by machine learning that reflects patterns and connections across vast amounts of text, the precise sources and reasons for any one representation's particular features remain unknown. By what process the learning machine stores its knowledge, distills it and retrieves it remains similarly unknown. Whether that process will ever be discovered, the mystery associated with machine learning will challenge human cognition for the indefinite future.

AI's capacities are not static but expand exponentially as the technology advances. Recently, the complexity of AI models has been doubling every few months. Therefore generative AI systems have capabilities that remain undisclosed even to their inventors. With each new AI system, they are building new capacities without understanding their origin or destination. As a result, our future now holds an entirely novel element of mystery, risk and surprise.

Enlightenment science accumulated certainties; the new AI generates cumulative ambiguities. Enlightenment science evolved by making mysteries explicable, delineating the boundaries of human knowledge and understanding as they moved. The two faculties moved in tandem: Hypothesis was understanding ready to become knowledge; induction was knowledge turning into understanding. In the Age of AI, riddles are solved by processes that remain unknown. This disorienting paradox makes mysteries unmysterious but also unexplainable. Inherently, highly complex AI furthers human knowledge but not human understanding—a phenomenon contrary to almost all of post-Enlightenment modernity. Yet at the same time AI, when coupled with human reason, stands to be a more powerful means of discovery than human reason alone.

The essential difference between the Age of Enlightenment and the Age of AI is thus not technological but cognitive. After the Enlightenment, philosophy accompanied science. Bewildering new data and often counterintuitive conclusions, doubts and insecurities were allayed by comprehensive explanations of the human experience. Generative AI is similarly poised to generate a new form of human consciousness. As yet, however, the opportunity exists in colors for which we have no spectrum and in directions for which we have no compass. No political or philosophical leadership has formed to explain and guide this novel relationship between man and machine, leaving society relatively unmoored.

ChatGPT is an example of what is known as a large language model, which can be used to generate human-like text. GPT is a type of model that can be automatically learned from large amounts of text without the need for human supervision. ChatGPT's developers have fed it a massive amount of the textual content of the digital world. Computing power allows the model to capture patterns and connections.

The ability of large language models to generate humanlike text was an almost accidental discovery. These models are trained to be able to predict the next word in a sentence, which is useful in tasks such as autocompletion for sending text messages or searching the web. But it turns out that the models also have the unexpected ability to create highly articulate paragraphs, articles and in time perhaps books.

ChatGPT is further specialized beyond a basic large language model, using feedback from humans to tune the model so that it generates more natural-seeming conversational text, as well as to try to contain its propensity for inappropriate responses (a substantial challenge for large

language models). ChatGPT instantaneously converts its representations into unique responses. The ultimate impression on a human conversant is that the AI is relating stationary collections of facts into dynamic concepts.

ChatGPT's answers, statements and observations appear without an explanation of where they came from and without an identifiable author. On its face, ChatGPT has no discernible motive or bias. Its outputs are complex, but its work is astonishingly rapid: In a matter of seconds, it can produce answers that coherently explain a high-level topic. They are not simply copied from the text in the computer's memory. They are generated anew by a process that humans are unable to replicate. It is able to incorporate hypotheticals and nonobvious psychological inferences. It can prioritize among billions of data points to select the single set of 200 words that is most relevant (or will appear most relevant to a human reader). In its own words, it makes probabilistic judgments about future outcomes, blending information from discrete domains into an integrated answer. It appears to accept contrary views, presenting them as simultaneously compatible. It imitates other writing without copying it. Even though the model is incapable of understanding in the human sense, its outputs reflect an underlying essence of human language.

We asked ChatGPT to explain its own capabilities. Can it predict the future?

No, it responded. "I can only provide information and answer questions based on the data that I have been trained on. I do not have personal experiences or consciousness, so I am not able to make predictions or forecasts about future events."

Is it able to make probabilistic judgments?

"I can provide probabilistic judgments based on the data that I have been trained on. For example, if I have been trained on a large dataset of weather patterns, I might be able to provide a probabilistic judgment about the likelihood of a particular type of weather occurring in a specific location on a given day. However, it is important to note that these probabilistic judgments are based on patterns and trends in the data, and they may not always accurately reflect real-world outcomes."

Prompt: Write a short essay in the thought of Henry Kissinger about the role of artificial general intelligence in international arms control.

The response began: "As we stand on the cusp of a new technological revolution, the question of artificial general intelligence and its role in international arms control becomes increasingly pressing. Will AGI bring about unprecedented stability, as advanced systems take over the decision-making processes associated with nuclear deterrence? Or will it introduce a new set of dangers, as the global balance of power is altered by the development of highly intelligent machines?"

We queried ChatGPT on topics of immediate interest, such as Russia's role in maintaining or undermining European stability. It gave seemingly thoughtful, concise answers in less than a minute, balancing contradictory information and answering fully without appearing to have an opinion.

On the receiving end, generative artificial intelligence appears superhuman or at least greatly enhanced cognitively. To the naive user, it seems like a supremely fast and highly articulate librarian-scholar coupled with a professorial savant. It facilitates the summary and interrogation of the world's knowledge far more effectively than any existing technological or human interface, and it does so with unique comprehensiveness. Its ability to integrate multiple knowledge domains and imitate multiple aspects of human thought makes it polymathic to a degree that exceeds the ambitions of any group of top-level humans. All of these qualities encourage unquestioning acceptance of whatever GPT generates and a kind of magical atmosphere for their operation. Yet at the same time, it possesses a capability to misinform its human users with incorrect statements and outright fabrications

Within a few days of ChatGPT's launch, more than a million people signed up to ask it questions. Hundreds of companies are working on generative technologies, and investment is pouring in, tilting discoveries to the commercial field. The huge commercial motives will, for the foreseeable future, take precedence over long-range thinking about their implications.

The biggest of these models are expensive to train—north of $1 billion per model. Once trained, thousands of computers work 24 hours a day to operate them. Operating a pretrained model is cheap compared with the training itself, and it requires only capital, rather than capital and computing skill. Still, paying for exclusive use of a large language model remains outside the bounds of most enterprises. These models' developers are likely to sell subscriptions, so that a single model will serve the needs of many thousands of individuals and businesses. As a result, the number of very large language models in the next decade may be relatively constrained. Design and control of these models will be highly concentrated, even as their power to amplify human efforts and thought becomes much more diffuse.

Generative AI will be used beyond the large language model to build many types of models, and the method will become increasingly multimodal and arcane. It will alter many fields of human endeavor, for example education and biology. Different models will vary in their strengths and weaknesses. Their capabilities—from writing jokes and drawing paintings to designing antibodies—will likely continue to surprise us. Just as the large language model developed a richer model of human language than its creators anticipated, generative AIs in many fields are likely to learn more than their assigned tasks imply. Breakthroughs in traditional scientific problems have become probable.

The long-term importance of generative AI transcends commercial implications or even noncommercial scientific breakthroughs. It is not only generating answers; it is generating philosophically profound questions. It will infuse diplomacy and security strategy. Yet none of the creators of this technology are addressing the problems it will itself create. Nor has the U.S. government addressed the fundamental changes and transformations that loom.

The seeming perfection of the model's answers will produce overconfidence in its results. This is already an issue, known as "automation bias," with far less sophisticated computer programs. The effect is likely to be especially strong where the AI generates authoritative-sounding text. ChatGPT is likely to reinforce existing predispositions toward reliance on automated systems reducing the human element.

The lack of citations in ChatGPT's answers makes it difficult to discern truth from misinformation. We know already that malicious actors are injecting reams of manufactured "facts," and increasingly convincing deepfake images and videos, into the internet—that is to say, into ChatGPT's present and future learning set. Because ChatGPT is designed to answer questions, it sometimes makes up facts to provide a seemingly coherent answer. That phenomenon is known among AI researchers as "hallucination" or "stochastic parroting," in which an AI strings together phrases that look real to a human reader but have no basis in fact. What triggers these errors and how to control them remain to be discovered.

We asked ChatGPT to give "six references on Henry Kissinger's thoughts on technology." It generated a list of articles purportedly by Mr. Kissinger. All were plausible topics and outlets, and one was a real title (though its date was wrong). The rest were convincing fabrications. Possibly the so-called titles appear as isolated sentences in the vastness of GPT's "facts," which we are not yet in a position to discover.

ChatGPT has no immediately evident personality, although users have occasionally prompted it to act like its evil twin. ChatGPT's lack of an identifiable author makes it harder for humans to intuit its leanings than it would be to judge the political or social viewpoint of a human being. Because the machine's design and the questions fed to it generally have a human origin, however, we will be predisposed to imagine humanlike reasoning. In reality, the AI is engaging in an inhuman analog to cognition. Though we perceive generative AI in human terms, its mistakes are not the mistakes of a human; it makes the mistakes of a different form of intelligence based on pattern recognition. Humans should not identify these mistakes as errors. Will we be able to recognize its biases and flaws for what they are? Can we develop an interrogatory mode capable of questioning the veracity and limitations of a model's answers, even when we do not know the answers ahead of time?

Thus, AI's outputs remain difficult to explain. The truth of Enlightenment science was trusted because each step of replicable experimental processes was also tested, hence trusted. The truth of generative AI will need to be justified by entirely different methods, and it may never become similarly absolute. As we attempt to catch our understanding up to our knowledge, we will have to ask continuously: What about the machine has not yet been revealed to us? What obscure knowledge is it hiding?

Generative AI's reasoning is likely to change over time, to some extent as part of the model's training. It will become an accelerated version of traditional scientific progress, adding random adaptations to the very process of discovery.

The same question put to ChatGPT over a period of time may yield changed answers. Slight differences in phrasing that seem unimportant at the first pass may cause drastically different results when repeated. At the present, ChatGPT is learning from an information base that ends at a fixed point in time. Soon, its developers will likely enable it to take in new inputs, eventually consuming an unending influx of real-time information. If investment continues to surge, the model is likely to be retrained with rising frequency. That will increase its currency and accuracy but will oblige its users to allow an ever-expanding margin for rapid change. Learning from the changing outputs of generative AI, rather than exclusively from human written text, may distort today's conventional human knowledge.

Even if generative AI models become fully interpretable and accurate, they would still pose challenges inherent in human conduct. Students are using ChatGPT to cheat on exams. Generative AI could create email advertisements that flood inboxes and are indistinguishable from the messages of personal friends or business acquaintances. AI-generated videos and advertisements depicting false campaign platforms could make it difficult to distinguish between political positions. Sophisticated signals of falsehood—including watermarks that signify the presence of AI-generated content, which OpenAI is considering—may not be enough; they need to be buttressed by elevated human skepticism.

Some consequences could be inherent. To the extent that we use our brains less and our machines more, humans may lose some abilities. Our own critical thinking, writing and (in the context of text-to-image programs like Dall-E and Stability.AI) design abilities may atrophy. The impact of generative AI on education could show up in the decline of future leaders' ability to discriminate between what they intuit and what they absorb mechanically. Or it could result in leaders who learn their negotiation methods with machines and their military strategy with evolutions of generative AI rather than humans at the terminals of computers.

It is important that humans develop the confidence and ability to challenge the outputs of AI systems. Doctors worry that deep-learning models used to assess medical imaging for diagnostic purposes, among other tasks, may replace their function. At what point will doctors no longer feel comfortable questioning the answers their software gives them? As machines climb the ladder of human capabilities, from pattern recognition to rational synthesis to multidimensional thinking, they may begin to compete with human functions in state administration, law and business tactics. Eventually, something akin to strategy may emerge. How might humans engage with AI without abdicating essential parts of strategy to machines? With such changes, what becomes of accepted doctrines?

It is urgent that we develop a sophisticated dialectic that empowers people to challenge the interactivity of generative AI, not merely to justify or explain AI's answers but to interrogate them. With concerted skepticism, we should learn to probe the AI methodically and assess whether and to what degree its answers are worthy of confidence. This will require conscious mitigation of our unconscious biases, rigorous training and copious practice.

The question remains: Can we learn, quickly enough, to challenge rather than obey? Or will we in the end be obliged to submit? Are what we consider mistakes part of the deliberate design? What if an element of malice emerges in the AI?

Another key task is to reflect on which questions must be reserved for human thought and which may be risked on automated systems. Yet even with the development of enhanced skepticism and interrogatory skill, ChatGPT proves that the genie of generative technology is out of the bottle. We must be thoughtful in what we ask it.

Computers are needed to harness growing volumes of data. But cognitive limitations may keep humans from uncovering truths buried in the world's information. ChatGPT possesses a capacity for analysis that is qualitatively different from that of the human mind. The future therefore implies a collaboration not only with a different kind of technical entity but with a different kind of reasoning—which may be rational without being reasonable, trustworthy in one sense but not in another. That dependency itself is likely to precipitate a transformation in metacognition and

hermeneutics—the understanding of understanding—and in human perceptions of our role and function.

Machine-learning systems have already exceeded any one human's knowledge. In limited cases, they have exceeded humanity's knowledge, transcending the bounds of what we have considered knowable. That has sparked a revolution in the fields where such breakthroughs have been made. AI has been a game changer in the core problem in biology of determining the structure of proteins and in which advanced mathematicians do proofs, among many others.

As models turn from human-generated text to more inclusive inputs, machines are likely to alter the fabric of reality itself. Quantum theory posits that observation creates reality. Prior to measurement, no state is fixed, and nothing can be said to exist. If that is true, and if machine observations can fix reality as well—and given that AI systems' observations come with superhuman rapidity—the speed of the evolution of defining reality seems likely to accelerate. The dependence on machines will determine and thereby alter the fabric of reality, producing a new future that we do not yet understand and for the exploration and leadership of which we must prepare.

Using the new form of intelligence will entail some degree of acceptance of its effects on our self-perception, perception of reality and reality itself. How to define and determine this will need to be addressed in every conceivable context. Some specialties may prefer to muddle through with the mind of man alone—though this will require a degree of abnegation without historical precedent and will be complicated by competitiveness within and between societies.

As the technology becomes more widely understood, it will have a profound impact on international relations. Unless the technology for knowledge is universally shared, imperialism could focus on acquiring and monopolizing data to attain the latest advances in AI. Models may produce different outcomes depending on the data assembled. Differential evolutions of societies may evolve on the basis of increasingly divergent knowledge bases and hence of the perception of challenges.

Heretofore most reflection on these issues has assumed congruence between human purposes and machine strategies. But what if this is not how the interaction between humanity and generative AI will develop? What if one side considers the purposes of the other malicious?

The arrival of an unknowable and apparently omniscient instrument, capable of altering reality, may trigger a resurgence in mystic religiosity. The potential for group obedience to an authority whose reasoning is largely inaccessible to its subjects has been seen from time to time in the history of man, perhaps most dramatically and recently in the 20th-century subjugation of whole masses of humanity under the slogan of ideologies on both sides of the political spectrum. A third way of knowing the world may emerge, one that is neither human reason nor faith. What becomes of democracy in such a world?

Leadership is likely to concentrate in hands of the fewer people and institutions who control access to the limited number of machines capable of high-quality synthesis of reality. Because of the enormous cost of their processing power, the most effective machines within society may stay in the hands of a small subgroup domestically and in the control of a few superpowers

internationally. After the transitional stage, older models will grow cheaper, and a diffusion of power through society and among states may commence.

A reinvigorated moral and strategic leadership will be essential. Without guiding principles, humanity runs the risk of domination or anarchy, unconstrained authority or nihilistic freedom. The need for relating major societal change to ethical justifications and novel visions for the future will appear in a new form. If the maxims put forth by ChatGPT are not translated into a cognizably human endeavor, alienation of society and even revolution may become likely.

Without proper moral and intellectual underpinnings, machines used in governance could control rather than amplify our humanity and trap us forever. In such a world, artificial intelligence might amplify human freedom and transcend unconstrained challenges.

This imposes certain necessities for mastering our imminent future. Trust in AI requires improvement across multiple levels of reliability—in the accuracy and safety of the machine, alignment of AI aims with human goals and in the accountability of the humans who govern the machine. But even as AI systems grow technically more trustworthy, humans will still need to find new, simple and accessible ways of comprehending and, critically, challenging the structures, processes and outputs of AI systems.

Parameters for AI's responsible use need to be established, with variation based on the type of technology and the context of deployment. Language models like ChatGPT demand limits on its conclusions. ChatGPT needs to know and convey what it doesn't know and can't convey.

Humans will have to learn new restraint. Problems we pose to an AI system need to be understood at a responsible level of generality and conclusiveness. Strong cultural norms, rather than legal enforcement, will be necessary to contain our societal reliance on machines as arbiters of reality. We will reassert our humanity by ensuring that machines remain objects.

Education in particular will need to adapt. A dialectical pedagogy that uses generative AI may enable speedier and more-individualized learning than has been possible in the past. Teachers should teach new skills, including responsible modes of human-machine interlocution. Fundamentally, our educational and professional systems must preserve a vision of humans as moral, psychological and strategic creatures uniquely capable of rendering holistic judgments.

Machines will evolve far faster than our genes will, causing domestic dislocation and international divergence. We must respond with commensurate alacrity, particularly in philosophy and conceptualism, nationally and globally. Global harmonization will need to emerge either by perception or by catastrophe, as Immanuel Kant predicted three centuries ago.

We must include one caveat to this prediction: What happens if this technology cannot be completely controlled? What if there will always be ways to generate falsehoods, false pictures and fake videos, and people will never learn to disbelieve what they see and hear? Humans are taught from birth to believe what we see and hear, and that may well no longer be true as a result of generative AI. Even if the big platforms, by custom and regulation, work hard to mark and sort bad content, we know that content once seen cannot be unseen. The ability to manage and control global distributed content fully is a serious and unsolved problem.

The answers that ChatGPT gives to these issues are evocative only in the sense that they raise more questions than conclusions. For now, we have a novel and spectacular achievement that stands as a glory to the human mind as AI. We have not yet evolved a destination for it. As we become Homo technicus, we hold an imperative to define the purpose of our species. It is up to us to provide the real answers.

---

Kissinger, Henry et. al. "ChatGPT Heralds an Intellectual Revolution." *Wall Street Journal*. 25 Feb. 2023,
https://www.wsj.com/articles/chatgpt-heralds-an-intellectual-revolution-enlightenment-artificial-intelligence-homo-technicus-technology-cognition-morality-philosophy-774331c6