Pause artificial intelligence research? Understanding AI policy challenges

Avi Goldfarb[®] *Rotman School of Management, University of Toronto*

Abstract. Artificial intelligence (AI) may be the next general purpose technology. General purpose technologies, such as the steam engine and computing, can have an outsized impact on productivity through a positive feedback loop between producing and application industries. Along with the discussion of AI's potential to improve productivity come a number of policy concerns related to AI's potential to automate jobs and to create existential risk for humanity. Because of these worries, in March 2023, a widely circulated petition called for a pause in AI research. That letter asked several questions about AI's potential impact on society. This paper examines those questions through an economic lens. It highlights reasons to be optimistic about the long-run impact of AI, while underscoring short-run risks. Economic models provide an understanding of where the ambiguity lies and where it does not. Our models suggest no ambiguity on whether there will be jobs and little ambiguity on long-term productivity growth if AI diffuses widely. In contrast, there is substantial ambiguity on the implications of AI's diffusion for inequality.

Résumé. Interrompre temporairement la recherche sur l'intelligence artificielle? Comprendre les enjeux politiques de l'IA. L'intelligence artificielle (IA) pourrait être la prochaine technologie universelle. Les technologies universelles, comme la machine à vapeur et l'ordinateur, peuvent avoir une incidence considérable sur la productivité grâce à une boucle de rétroaction positive entre les secteurs de la production et des applications. La discussion sur le potentiel de l'IA à augmenter la productivité s'accompagne d'un certain nombre de préoccupations politiques en lien avec la possibilité d'automatiser des emplois et de créer un risque existentiel pour l'humanité. En raison de ces inquiétudes, en mars 2023, une pétition largement diffusée a réclamé une pause dans la recherche sur l'IA. Cette lettre soulevait plusieurs questions sur l'incidence potentielle de l'IA sur la société. Le présent article examine l'aspect économique de ces questions. Il met en évidence les raisons d'être optimiste quant à l'incidence à long terme de l'IA, tout en soulignant les risques à court terme. Les modèles économiques permettent de comprendre où l'ambiguïté réside. Ces modèles ne laissent planer aucune ambiguïté sur la création d'emplois et la croissance de la productivité à long terme si l'utilisation de l'IA se généralise. En revanche, de nombreuses questions subsistent quant aux conséquences d'une propagation inégale de l'IA.

JEL classification: O3, O4

H OW DO WE understand artificial intelligence (AI) technology as economists? To discuss this, we begin with the internet.¹ A pivotal year for the internet was 1995, when the last aspects of the public internet, the NSFNET in the US, was privatized (Greenstein 2015).

Corresponding author: Avi Goldfarb, avi.goldfarb@rotman.utoronto.ca

This article is adapted from the State of the Art Lecture, "The Economics of Artificial Intelligence," at the 2023 Canadian Economic Association meetings in Winnipeg, Manitoba. Many of the ideas in this article were developed in collaboration with Ajay Agrawal and Joshua Gans. I thank Verina Que and Tate Li for excellent research assistance. Support was provided by the Social Sciences and Humanities Council of Canada.

¹ This introduction draws heavily on ideas from Agrawal et al. (2018). Canadian Journal of Economics / *Revue canadienne d'économique* 2024 57(2) May 2024. / *mai 2024*.

^{24 /} pp. 363-377 / DOI: 10.1111/caje.12705

[©] The Authors. Canadian Journal of Economics/Revue canadienne d'économique published by Wiley Periodicals LLC on behalf of Canadian Economics Association.

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

It was the year that Bill Gates wrote his "internet tidal wave" email, saying that Microsoft would focus on enabling the benefits of the internet to accrue to consumers and businesses. It was also the year of the Netscape IPO, where Netscape was valued at over a billion dollars with zero profit.

The hype around the internet kept building, and people stopped talking about it like a new technology—plenty of people started calling the internet a new economy, claiming we could throw away our economics textbooks and calling for a reset in our understanding of how economics and finance worked. We economists disagreed. We argued that economics had a lot to say about understanding the changes coming from the widespread diffusion of digital communication. The economists focusing on these questions in the late 1990s emphasized that people just need to understand what has gotten better, faster and cheaper in order to anticipate changes in the economy. As emphasized in the book *Information Rules* (Shapiro and Varian 1998), by understanding that the implications of the internet can be anticipated by recognizing that digital communication makes information search, information transportation and information replication cheap, it becomes possible to map out likely consequences.

Shapiro and Varian (1998) highlighted in their book that, more obviously, cheaper search should lead to fiercer price competition. A less obvious observation was that, if copying becomes cheaper—if anything said by anyone can now be instantly broadcast around the world for free—then people will be more careful about what they say. The changes because of the internet were directly going to lead to an increased attention to privacy. Once it was clear what had gotten cheaper, it was possible to map out the consequences. This framing on the impact of the internet has proven powerful. Goldfarb and Tucker (2019) summarize 20 years of the digital economics literature, emphasizing that the internet lowered the costs of information: (i) searching, (ii) replication, (iii) transportation, (iv) tracking and (v) verification. Hundreds of papers fit under this theme of what this technology enables to be better, faster and cheaper.

The idea of understanding of technological change through lower costs in some fundamental input is rooted in an older literature in the economics of technological change. Bresnahan (1999) emphasizes that computers do arithmetic and that it has been cheap arithmetic that has been transformative. Nordhaus (2007) measures how the cost of arithmetic has fallen since the early 19th century. Demand curves slope downward, and so a rightward shift in supply means increased quantity. Therefore, cheap arithmetic meant increased use of arithmetic. As the price of arithmetic fell, it was first applied to well-established arithmetic problems. For example, an early large-scale use of machine computing was in artillery tables in the 1940s and into the 1950s (Ceruzzi 2003). As machine arithmetic got cheaper, we started to apply it to other aspects of arithmetic that were part of the workplace, such as bookkeeping. Accountants and bookkeepers no longer spend much time doing arithmetic. There are still plenty of accountants and bookkeepers, though such roles might involve different skills than they did decades ago.² While some accountants and bookkeepers likely never adapted to the use of machine arithmetic in their work, many of the people who were good at using arithmetic in bookkeeping also turned out to be good at using arithmetic for company strategy, tax policy and other things.

Over time, as the price of machine arithmetic fell, we began to realize that a multitude of other problems could be reframed as arithmetic. For example, many games can be solved, engineering-wise, with arithmetic. Music, mail and even pictures can also be solved through

² Over 1.5 million in the United States in 2024 (www.bls.gov/oes/current/oes433031.htm#nat).

arithmetic. Kodak was fundamentally a chemical engineering company; but, as of 2024, taking your picture using a camera is more a computer science problem than a chemistry one. Pictures have been reframed as arithmetic (Agrawal et al. 2018).

This is basic microeconomics. As the price of something falls, we do more of it. And when the price of something falls a lot, we do a lot more of it. The demand curve is downward sloping.

1. Al as a prediction machine

For artificial intelligence, the question is what has become cheaper as AI technology has improved. It is tempting to look to science fiction and to anticipate the arrival of machines that can do everything humans do, implying that what has gotten cheaper is general intelligence, human-level or better. Such artificial general intelligence may be feasible someday, with profound implications for society.

However, a careful look at the technologies underlying the recent advances in AI suggests that it is machine learning that has experienced the most rapid gains and that underlies the reason AI is commercially relevant today. Machine learning is a branch of computational statistics and is, therefore, prediction technology in the statistical sense. It uses available data to fill in missing information (Agrawal et al. 2018).

Therefore, while AI sometimes has a broader definition of machines that can solve problems normally requiring human intelligence, in order to understand the economic and policy implications of today's AI, it is most useful to understand recent advances in AI as advances in prediction machines (Agrawal et al. 2018, 2019a).

Just like with arithmetic, as prediction got cheaper, machine prediction started being used in more places. More precisely, an outward shift in supply led to a movement down and to the right along the demand curve. Many early applications of AI in practice involved well-established prediction problems in business such as whether someone would default on a loan or pay it back and whether someone would make an insurance claim. Banks and insurance companies have long used various tools to predict loan defaults and insurance claims, and over the past decade or more, they have increasingly been using machine learning tools—AI—to help with these predictions (Agrawal et al. 2018).

As predictions get cheaper though, we are also recognizing that there are a handful of things that we may not have thought of as "predictions," in the traditional sense, about 10 or 20 years ago. Just as we did with games and mail and photos when it comes to arithmetic, we are now finding problems that can clearly be solved engineering-wise as *predictions*. Machine diagnosis is prediction in the sense that a doctor takes in data about a patient's symptoms and fills in the missing information of the cause of those symptoms (e.g., Mullainathan and Obermeyer 2021) and image recognition is prediction in the sense that a context.

The most recent wave of AI, using generative models, is also prediction. Large language models like those underlying ChatGPT are prediction models that predict the set of words and order that is the most helpful, honest and harmless response to a particular query (Nakano et al. 2021). Image-generation tools such as DALL-E are also prediction machines, predicting the set of pixels that is the best response to a particular request.

These examples suggest that computational statistics is being used in many contexts that are often not considered to be predictions—just as computers apply arithmetic to problems that may not seem like arithmetic (Agrawal et al. 2018).

2. Al as a general purpose technology

These rapid advances in what machine prediction can do have coincided with an increase in literature exploring the idea that these prediction technologies could be the next "general purpose technology" (GPT). GPTs, as defined by Bresnahan and Trajtenberg (1995), describe a set of technologies that have an outsized impact on productivity growth, such as electricity, the computer and the steam engine. In a standard growth model, the impact of an innovation on total factor productivity is short lived. Bresnahan and Trajtenberg define GPTs as technologies that generate sustained productivity growth because of a positive feedback loop between industries that produce the core technology and industries that find applications for the core technology. GPTs generate innovation in the application (or using) industries, which in turn generate innovation upstream in the producing industries, etc. (Bresnahan 2010).

Electricity provides an informative example. The initial invention of the ability to produce electricity led to innovations in applications such as the light bulb and the electric motor. These in turn led to innovation in producing and distributing electricity at scale, such as hydropower and alternating current. Further downstream innovations in household appliances and the organization of the factory followed (e.g., Devine 1983, David 1990). This positive feedback loop between upstream and downstream innovations, between "producing" and "application" industries, meant that a standard model of growth—where, in the long run, we do not have any growth— is temporarily disrupted because of this positive feedback loop. Bresnahan and Greenstein (1996) call this co-invention, downstream innovation that then feeds into further upstream innovation.³

At the inaugural NBER AI conference in 2017, a number of economists were invited to talk about how they believed AI would relate to their particular expertise in economics (Agrawal et al. 2019c). A common theme, highlighted by Brynjolfsson et al. (2019), Cockburn et al. (2019) and (unsurprisingly) Trajtenberg (2019) was that AI is likely to be the next GPT. A challenge, however, in determining whether a technology was likely to be the next GPT is that GPTs are defined in retrospect. In order to determine if the positive feedback loop occurred, it is necessary to have decades of data. For example, Lipsey et al.'s (2005) comprehensive history of GPTs uses a millennia-long time scale, starting with the domestication of plants. Furthermore, it is possible that some technologies could have been GPTs, but the positive feedback loop never happened. In a sense, GPTs are an equilibrium concept and there is an equilibrium without the feedback loop in which the productivity gains never occurred.

Therefore, while it is not possible to prove that a technology is a GPT in advance, it is possible to compare emerging technologies in the extent to which they have the characteristics of GPTs. Goldfarb et al. (2023) is an attempt to identify emerging technologies with the potential to be GPTs. We started with a number of emerging technologies that were often described as transformative in the press. We found one such list in the technologies listed in Forrester's hype cycle and showed robustness to other choices of lists, such as technologies that appeared on the cover of the academic journals *Science* and *Nature*. Using data on online job postings from 2010–2020, we evaluated each technology in terms of the core aspects of a GPT, as described by Bresnahan and Trajtenberg (1995): (i) widespread use, (ii) potential

³ The phenomenon of a handful of innovations having an outsized impact on productivity growth through combined innovation was already established in the literature. For example, the idea is closely related to Mokyr's (1990) emphasis on "macroinventions" and Rosenberg's (1963) discussion of innovation in the 19th century machine tools industry.

for innovation and (iii) innovation in application industries. We found that machine learning—along with a set of related data science technologies—was employed across a wide variety of industries, was disproportionately used for research and was applied within research contexts in application industries. These features were much stronger for these technologies than for the others in the list, such as blockchain, nanotechnology and 3D printing.

This is the best evidence we have that AI might be a GPT. I stress that it is necessary to use the word "might" until the positive feedback loop has played out. On one hand, given the multiple equilibria, it is still possible that the positive feedback loop leading to a large productivity gain from AI does not occur. On the other hand, it is also possible that AI is the next GPT but the productivity gains may not appear in the data for a long time.

Central to the idea of GPTs leading to this outsized productivity growth is the need for co-invention in downstream application industries (Bresnahan and Greenstein 1996). Otherwise, the feedback loop never gets going. This, however, suggests a productivity challenge. Brynjolfsson et al. (2019) argue that AI might lead to a "modern productivity paradox." The productivity paradox was described in a 1987 New York Times book review article by economist Robert Solow, who said, "you can see the computer age everywhere but in the productivity statistics." He emphasized computer adoption in the industry was well underway in the 1960s and widespread by the 1980s. Nevertheless, there was little evidence of an impact on aggregate productivity.

David (1990) turned Devine's (1983) work on the history of electricity in American factories into a parable for understanding the productivity paradox identified by Solow. He identified a parallel between electrical motors and computing. Edison's patent for the electric light bulb was in 1880, and Tesla's patent for the alternating current electric motor was in 1890. Yet, it would not be until the 1920s—40 years later—that half of US factories and half of US households had adopted electricity. Devine documented that a small number of factories were electrified in the 1880s and 1890s. These factories treated electricity as a new cheap power source, replacing only the central steam engine or water wheel. It allowed the factory to keep operating as it always had, but at a slightly lower cost. By 1900, fewer than 5% of US factories were electrified. Generally, the benefits of electrifying were not worth the costs of taking out the old power source and putting in a new one, along with the costs of figuring out electricity distribution, reducing fire risk and training employees to be safe around high voltage wires. Around 1900, Devine noted that some people realized that electricity enabled the building of a new kind of factory. These people *co-invented*, along with electricity, the 20th-century factory, with inputs coming in one end and outputs coming out the other. It was when we arrived at modular production—which one might think of as the Henry Ford style factory—where the order of production was not determined by the power needs of the machines but by the logic of production. This would be a turning point where we started to see a change in the trajectory and a sharp increase in the adoption of electricity.

We saw something similar in computers. Solow was just a few years too early. By the 1990s the impact of computers on productivity was measurable, and David's application of Devine's work to computing seemed correct (Jorgenson et al. 2008). Brynjolfsson et al. (2019) suggest the same might be happening with AI. The excitement around AI appears to be accelerating, but evidence that it is affecting productivity remains elusive. One possibility is that AI will not live up to the hype.

A separate set of worries focus on the consequences of AI if it does live up to the hype. It is this separate set of worries that has received a great deal of attention since the release of OpenAI's chatGPT in November 2022.

3. The "pause" letter

The Future of Life Institute summarized many of these worries when they released their "Pause Giant AI Experiments: An Open Letter" in March 2023. This letter opened with the statement, "AI systems with human-competitive intelligence can pose profound risks to society and humanity."⁴ While, as noted earlier in this article, the AI tools that are available at the time of writing this article are prediction machines and not artificial general intelligence tools, the letter emphasized the risk that today's prediction machines could soon turn into something much more powerful. It is not my expertise to assess the likelihood of such an event, and I struggle to find a definition for human-competitive intelligence that is conducive to economic analysis.⁵

The pause letter, however, includes several concerns about the risks of AI that we can use the existing economics literature to understand even without a clear definition of what human-competitive intelligence means. Specifically, the pause letter asks four questions:

- 1. Should we risk loss of control of our civilization?
- 2. Should we develop nonhuman minds that might eventually outnumber, outsmart, obsolete and replace us?
- 3. Should we let machines flood our information channels with propaganda and untruth?
- 4. Should we automate away all the jobs, including the fulfilling ones?

It is clear from the way these questions are written that the authors of the letter think the answer to each question is clearly "no."

Economics has little to say about the first question. We do not have good models of who controls civilization in the first place; hence, it is not clear how to assess the impact of a loss of control. For now, "Should we risk loss of control of our civilization?" remains a question outside of our literature, and assuming we can define "we," then I'm happy to defer to the authors of the pause letter and agree that the answer should be "no."

In contrast, I argue that economics has a lot to say about the other questions and suggests that the answers are not obvious, particularly for the last question on automating away jobs. I next discuss each question in turn.

3.1. Should we develop non-human minds that might eventually outnumber, outsmart, obsolete and replace us?

Several economists have incorporated a machine that facilitates innovation into familiar macroeconomic models of economic growth. Nordhaus (2021) discusses the possibility of an economic singularity with infinite consumption as a consequence of rapid technological change through superintelligent innovation. Korinek and Stiglitz (2019) and Aghion et al. (2019) also discuss how superintelligent AI could lead to a rapid improvement in productivity as the limits to growth set by human constraints disappear. In these models, subject to the important caveats that the machines are aligned with humans and that the wealth generated is not too concentrated, a world where machines can invent things is fantastic for humans. Innovative machines mean that humans can have more of what we want, whatever that might be.

⁴ See https://futureoflife.org/open-letter/pause-giant-ai-experiments/.

⁵ There is a long literature on the definition of human-competitive intelligence (e.g., Hawkins and Blakeslee 2005, Bostrom 2014, Domingos 2015) that recognizes that both technology and human skills change over time.

Jones (2023) recognizes that such innovative machines carry real risks. He notes that "creating a superintelligent entity misaligned with human values could lead to catastrophic outcomes, including human extinction" (p. 1). Once such risk is acknowledged, it is useful to weigh the possibility of extraordinary growth against existential risk. Jones formally models this trade-off and demonstrates that, under standard modelling assumptions, large consumption gains might be worth trading off against existential risk. It depends on the assumed curvature of utility, with some models suggesting a high degree of risk might be worth the benefit under reasonable assumptions and others leading to much more caution. Importantly, even for those models that suggest caution, if AI leads to new technologies that extend life expectancy, then large existential risk and do not run into the sharply declining marginal utility of consumption" (p. 2). Thus, if AI might improve health care in a meaningful way, macroeconomic growth models suggest that existential risk might be worth bearing.

It is important to note that this could be interpreted as a weakness of the growth models rather than a statement about whether it is worth risking catastrophe. Jones (2023, p. 19) concludes his paper by making clear the limitations of a utilitarian social planner in understanding existential risk: it treats "10% of the population dies each period" as equivalent to a "10% chance of human extinction." Thus, the standard way to interpret an economic growth model of how to weigh the benefits and risks of "non-human minds that might eventually outnumber, outsmart, obsolete and replace us" might not apply here.

Nevertheless, what the growth models examining superintelligent machines suggest is that it is not at all obvious that we should stop AI progress because machines might outnumber, obsolete and replace us. Furthermore, there are good reasons to want (aligned) machines that outsmart us.

3.2. Should we let machines flood our information channels with propaganda and untruth?

An important contribution that economists bring to this discussion is a recognition of the importance of equilibrium. Our discipline suggests that it cannot be a long-run equilibrium that people get repeatedly fooled by misinformation. The short run, however, is less optimistic. Until people learn what to ignore, and until new verification tools get developed and deployed, people will be fooled. To my knowledge, the economics literature on this question is sparse, and so this section will be brief. Still, it is clear that a variety of long-run equilibria are possible.

First, it is possible that society lands in a babbling equilibrium (see, e.g., Farrell and Rabin 1996) in which people stop trusting information, particularly online information. For example, rather than lead to more blackmail from fake images, AI misinformation is likely to make it more difficult for blackmail to occur, even with real images.

Second, new ways to verify information could be developed, through technologies or mechanism design. While these verification tools may be more costly than the processes that existed before AI's diffusion, they nevertheless overcome the main challenges of misinformation. For situations where verification is important enough, tools will be developed to address them. This is already happening. Consider, for instance, the widespread adoption of voice authentication services by financial institutions. An initial worry about AI voice mimicry was that AI would trick voice verification tools at banks and elsewhere.⁶ Today, banks are less likely to rely on voice alone.

⁶ www.banking.senate.gov/newsroom/majority/brown-presses-banks-voice-authentication -services

370 A. Goldfarb

Third, the equilibrium may nevertheless be to have misinformation in some domains because of motivated reasoning. Many people may prefer to receive incorrect information that is consistent with their beliefs corrected, and therefore, they will not want the information corrected. This could lead to an alternative possibility of propaganda and polarization (see, e.g., Epley and Gilovich 2016, Zimmerman 2020).

Thus, while the long-run equilibrium is not that people are repeatedly fooled and unhappy about it, there are both pessimistic and optimistic possibilities in terms of the degree to which misinformation affects society.

3.3. Should we automate away all the jobs, including the fulfilling ones?

In many ways, this question goes directly against our understanding of how productivity improvements affect jobs. Baumol (1967) pointed out that increased labour productivity in one sector will increase the size of the labour force in sectors that do not experience such a productivity increase. Aghion et al. (2019) highlighted that Baumol's "cost disease" applies to AI automation. Rapid productivity growth in some sectors leads to slow-growing sectors becoming an increasingly large fraction of the economy. This suggests that, for AI to automate away all jobs, humans must lack any comparative advantage in all work (e.g., AI could be infinitely better than humans at all tasks). Otherwise, any residual tasks would gradually become the work that humans do, as long as the marginal productivity of that work was higher than the workers' outside option. Sector-specific productivity improvements have led to changes in the nature of jobs in the past. As Autor (2015) points out, 41% of the US workforce was employed in agriculture in 1900. One hundred years later, that fraction had fallen to 2%. Employment did not collapse. While some jobs disappeared, other jobs grew in importance as complements to automation. This calls into doubt whether the question of automating away all jobs makes economic sense.

Even if AI were to automate away some or all jobs, our textbook models imply this is not necessarily bad because the willingness to work is modelled as a trade-off between consuming leisure and consuming the goods and services that can be bought with work income (e.g., Goolsbee et al. 2013, p. 187). Thus, all else equal, individuals would prefer to have income without the necessity of work, implying that automating away all jobs but keeping income fixed would be welfare-improving. Leisure is a good, and so work, as the opposite of leisure in the textbook model, is a bad. Today, workers in many developed countries start work later, take more vacations, work fewer hours per week and have more years of retirement than workers from a century ago. Nevertheless, that decline in work and increase in leisure is widely seen as for the better.

The issue, therefore, is not jobs per se. It is about what people could do in the absence of work. As Stevenson (2018, p. 195) put it, rather than focus on jobs,

... there are two separate questions: there is an employment question, in which the fundamental question is, can we find fulfilling ways to spend our time if robots take our jobs? And there is an income question, can we find a stable and fair distribution of income?

People like leisure, but they also like to have purpose and to get paid. On Stevenson's first question, if people like work for the sake of work because work gives purpose, then that suggests a reason to worry about machine automation (or at least a reason to find purpose outside of work). Beyond that, standard economic analysis may not be the best tool for understanding whether people will find meaning without work.

Economists have more to say about Stevenson's second question on the stable and fair distribution of income. This is an important concern if this time is different and Baumol's (1967) insights do not apply to AI because AI replaces all work. However, even if the question in the pause letter is misguided, and AI will not automate away all the jobs, AI's impact on the stable and fair distribution of income remains a fundamental concern. If AI substantially improves productivity and average income, income inequality may nevertheless increase. Becoming wealthier on average does not mean the median person, or the 25th percentile or even the 99th percentile is becoming wealthier. This suggests that there are real reasons to wonder whether AI is a technology that will increase inequality and not decrease it.

4. Inequality

Over the past 70 years, there have been two other major advances in information technology that have been widely categorized as general purpose technologies (Bresnahan 2010): digital computing and the internet. The empirical literature suggests that each led to an increase in wage and wealth inequality within developed countries. For example, Autor et al. (2008) document increased US wage inequality during the 1980s and emphasize increased demand for skills related to computing as an important factor. This conclusion draws on earlier work by Doms et al. (1997), Autor et al. (1998) and Bartel et al. (2007). Similarly, the diffusion of the internet appears to have increased inequality in the 1990s and early 2000s (Forman et al. 2012, Akerman et al. 2015).

Beyond an expectation that the future might look like the past, there are two broad reasons to expect that AI might increase inequality, reasons related to capital and labour.

First, AI tools are embedded in capital. They are machines. Furthermore, capital income is generally less widely distributed than labour income.⁷ AI might lead to increased concentration of income if capital income becomes increasingly important. Acemoglu and Robinson (2023) emphasize that this, in turn, may lead to further concentration in political power leading to a feedback loop of increasing inequality. At the same time, it is not a foregone conclusion that AI will lead to an increase in the capital share of income. If technology complements labour, then labour's share could rise. Furthermore, Acemoglu and Johnson (2023) emphasize the role of policies and institutions in allocating any rents accrued from technology to capital and labour.

Second, within labour, there are also reasons to anticipate AI will increase income inequality. The increased inequality that occurred as computers and the internet diffused was likely a result of technology increasing demand for skilled labour without a corresponding increase in supply, in the context of a "race between education in technology" (Goldin and Katz 2008). For computers and the internet, the technological change was skill-biased, in that those who benefitted most were those with relatively high levels of education and training (Autor et al. 2008, Forman et al. 2012, Autor 2014). Goldin and Katz (2008) emphasize that the supply of skilled workers did not keep up with increased demand and so labour income inequality increased in the latter part of the 20th century. Nevertheless, Goldin and Katz note that technological change does not necessarily increase inequality, citing the first part of the 20th century as a period in which inequality decreased despite the widespread adoption of new technologies in the workplace.

It is an open question whether AI, primarily in the form of prediction machines, but perhaps also in the form of artificial general intelligence, is skill-biased and whether education can keep up with any increased demand for particular skills.

Several people have put forth the argument that it is a choice whether AI increases or decreases inequality. Accomoglu and Johnson (2023) emphasize political choices to empower

⁷ www.oecd.org/economy/growth/49421421.pdf

workers or corporations. Furthermore, they emphasize that technology should be directed toward creating new tasks rather than automating work (Acemoglu and Johnson 2023, p. 394).

This idea that engineers, scientists and other innovators should prioritize augmenting rather than automating work is a recurring theme in writing about technology. In his book on the history of the relationship between humans and computers, *Machines of Loving Grace*, Markoff (2015) identified two distinct groups of computer scientists. One group focused on artificial intelligence, building machines that can replicate human-level intelligence. The other focused on intelligence augmentation, allowing humans to perform previously unthinkable tasks such as the large-scale arithmetic done by computers or instant global communication facilitated by the internet. Markoff celebrates those who have chosen human-centred design over automation. Put differently, he argues for engineers to focus on intelligence augmentation over artificial intelligence.

Brynjolfsson (2022, p. 282) also advocates for incentives to augment humans rather than automate the tasks that are part of human work. He argues, "[a] good start would be to replace the Turing test, and the mindset it embodies, with a new set of practical benchmarks that steer progress toward AI-powered systems that exceed anything that could be done by humans alone."⁸ Instead, we should redirect our efforts to augmentation and avoid what he calls the "Turing trap."

Acemoglu and Johnson, Markoff, and Brynjolfsson propose one way that AI might reduce inequality. Their arguments emphasizes that, if engineers can be incentivized to make technologies that augment human intelligence, this should lead to a decrease in inequality and better outcomes for human workers.

In Agrawal et al. (2023; forthcoming), my coauthors and I argue that this is not necessarily the case. First, human-augmenting technologies do not necessarily decrease inequality. Computing and the internet, two technologies that Markoff highlights as intelligence augmentation, led to increased inequality because the demand for skills outpaced the supply. The people these technologies augmented most were already near the top of the income distribution. Furthermore, technology can be equalizing even if it automates. Task automation and labour augmentation are not polar opposites because the automation of some tasks can lead to augmentation of others. Put differently, one person's automation can be another's augmentation.

Furthermore, preliminary evidence from the usage of AI in business suggests that AI-based automation might be inequality-reducing. Eloundou et al. (2023) identify the tasks human workers do as part of their daily workflows that can also be done by large language models. The identification of tasks within human work uses the O*NET data developed by the US Department of Labor. The evaluation of whether a particular task can be done by a large language model is done by a combination of human expertise and AI tools. The paper documents that large language models can do things that many human workers do as a key part of their daily workflows. However, they tend to be the things that the people toward the top of the income distribution do a lot more than the people in the middle or at the bottom of the income distribution. Thus, the exposure to these models is higher for higher-income workers. The paper is careful to note that "exposure" could be positive or negative. If AI automates some tasks done by a high-income worker, that could make the worker more productive and increase their wage, or it could reduce wages because

⁸ The Turing test, originally called the imitation game (Turing 1950), is a test of whether a machine can fool a human into believing the machine is a human on the basis of its responses in a natural language conversation.

that particular task was a specialized skill leading to high wages.⁹ Writing, for example, is a skill that has a relatively high return.

In medicine, physicians are highly paid partly because they are able to diagnose and other medical professionals are not. Diagnosis is prediction, in the sense that diagnosis requires taking available information and then filling in a cause for the observed symptoms. If AI does diagnosis, that might empower nurses, pharmacists and other medical professionals to be much more productive, while it may simultaneously reduce physician wages. In other words, it would be equalizing within the medical profession (e.g., Agrawal et al. 2022, p. 93).

Several recent papers provide evidence that AI might help lower-wage workers relative to higher-wage workers in a particular context. Brynjolfsson et al. (2023) show that, within a call centre, the use of generative models to support workers increased overall productivity. Furthermore, it increased the productivity of new workers and less-productive workers most. Similarly, Dell'Acqua et al. (2023) study the usage of large language models in a management consulting firm and show that AI increased productivity overall, and especially for consultants who were relatively low performers. Wiles et al. (2023) show that the use of a writing assistant tool in resumés increased the likelihood of getting hired, without impacting employer satisfaction, improving the job prospects of workers who are not good writers.

Despite this evidence, there are two notes of caution. First, these examples are applications within a particular context. Brynjolfsson et al. (2023) document reduced inequality of worker productivity within a call centre; however, if AI automates much of the call centre so that many call centre jobs disappear, the overall impact of AI on the workforce might increase inequality. Furthermore, the evidence is still being collected. Otis et al. (2024) implement a field experiment in Kenya and shows that an AI assistant improves the performance of relatively skilled entrepreneurs. While the relatively skilled entrepreneurs in their study are still small and medium enterprises (SMEs) in the developing world, the paper makes clear there are limits to the degree to which AI can help the low-skilled workers.

Overall, therefore, the impact of AI on inequality is uncertain. There are forces pushing toward increased inequality and others toward decreased inequality. It is too early to distinguish between these hypotheses. The doom and gloom is not warranted, but neither is unbridled optimism. Like the 1890s for electricity or the 1990s for the internet, it is clear that there is strong potential for technical change. We have models to help us understand the world, and the models provide an understanding of where the ambiguity lies and where it does not. We do not have ambiguity on whether there will be jobs. We do not have ambiguity on long-term productivity growth if AI diffuses widely. We have ambiguity on the implications of AI's diffusion for inequality.

5. Conclusion

To summarize, AI is prediction technology. It is a likely candidate for the next general purpose technology. That does not mean it is necessarily the case. General purpose technologies create value through a positive feedback loop between producing and application industries. If that feedback loop arises and AI fulfills its promise, we should expect to see large productivity gains.

⁹ Current tools for evaluating whether a technology will affect workers provide information on which jobs will be affected, but they do not provide information on whether the technology will have a positive or negative effect on workers. In my view, this represents a weak point in the usefulness of these tools for industry and for government policy.

Those large productivity gains are the main reason we might accept the risks that AI advances might bring. The pause letter summarized many concerns about how AI-driven changes might negatively impact society. A detailed look at the economic literature around superintelligence, misinformation and task automation suggests that it is not clear that these concerns have been appropriately framed in the popular discussion. Unambiguously, there are short-term risks from misinformation and job loss and real longer-term worries about inequality, democracy and control. Still, much of the economics literature provides an optimistic perspective on the longer-term consequences of AI for society both because of the potential for productivity gains to improve our standard of living (broadly defined) and because AI has potential to lead to a more equitable distribution of those gains.

True to the reputation of our profession, perhaps the most accurate answer to questions about whether society should develop machines that might outsmart us and automate away jobs is "it depends." Whether we should develop machines that outsmart us depends on how much risk we are willing to take in exchange for potentially extraordinary improvements in productivity, generally, and in health, in particular. Whether we should develop machines that automate away jobs depends on the implications of AI-driven productivity gains for the distribution of income and on society's preferences for a potential increase, or reduction, in income inequality. Our models do not provide direction on what to do, but they do provide clarity on what, specifically, it depends on.

References

- Acemoglu, D., and D. Autor (2011) "Skills, tasks and technologies: Implications for employment and earnings." In D. Card and O. Ashenfelter, eds., *Handbook of Labor Economics*, vol. 4, part B, ch. 12, pp. 1043–171. https://doi.org/10.1016/s0169-7218(11)02410-5
- Acemoglu, D., and S. Johnson (2023) Power and Progress: Our Thousand-Year Struggle Over Technology and Prosperity. New York: PublicAffairs
- Acemoglu, D., and P. Restrepo (2018) "The race between man and machine: Implications of technology for growth, factor shares, and employment," *American Economic Review* 108(6), 1488-542. https://doi.org/10.1257/aer.20160696
 - (2019a) "Artificial intelligence, automation, and work." In A. Agrawal, J. Gans and
 - A. Goldfarb, eds., *The Economics of Artificial Intelligence*, ch. 8, pp. 197–236. Chicago: University of Chicago Press. https://doi.org/10.7208/chicago/9780226613475.003.0008 — (2019b) "Automation and new tasks: How technology displaces and reinstates labor."
 - Journal of Economic Perspectives 33(2), 3–30. https://doi.org/10.1257/jep.32.2.3
- Aghion, P., B. F. Jones, and C. I. Jones (2019) "Artificial intelligence and economic growth." In A. Agrawal, J. Gans and A. Goldfarb, eds., *The Economics of Artificial Intelligence*, ch. 9, pp. 282–89. Chicago: University of Chicago Press
- Agrawal, A., J. Gans, and A. Goldfarb (2018) Prediction Machines: The Simple Economics of Artificial Intelligence. Brighton, MA: Harvard Business Review Press
- (2019a) "Economic policy for artificial intelligence," *Innovation Policy and the Economy* 19, 139–59. https://doi.org/10.1086/699935
- (2019b) "Artificial intelligence: The ambiguous labor market impact of automating prediction," Journal of Economic Perspectives 33(2), 31-50. https://doi.org/10.1257/jep .33.2.31
- (2019c). The Economics of Artificial Intelligence: An Agenda. Chicago: University of Chicago Press
 - (2022) Power and Prediction: The Disruptive Economics of Artificial Intelligence. Brighton, MA: Harvard Business Review Press
 - (2023) "Do we want less automation?," Science 381(6654), 155-58. https://doi.org/10 .1126/science.adh9429

— (forthcoming) "The Turing transformation: Artificial intelligence, intelligence augmentation, and skill premiums," *Harvard Data Science Review*. https://doi.org/10.1162 /99608f92.35a2f3ff

- Akerman, A., I. Gaarder, and M. Mogstad (2015) "The skill complementarity of broadband internet," Quarterly Journal of Economics 130(4), 1781-824. https://doi.org/10.1093/qje /qjv028
- Atack, J., R. A. Margo, and P. W. Rhode (2019) "Automation' of manufacturing in the late nineteenth century: The hand and machine labor study," *Journal of Economic Perspectives* 33(2), 51–70. https://doi.org/10.1257/jep.33.2.51
- Autor, D.H. (2014) "Skills, education, and the rise of earnings inequality among the 'other 99 percent'," Science 344(6186), 843-51. https://doi.org/10.1126/science.1251868
- (2015) "Why are there still so many jobs? The history and future of workplace automation," Journal of Economic Perspectives 29(3), 3-30. https://doi.org/10.1257/jep.29.3.3
- Autor, D.H., L. F. Katz, and M. S. Kearney (2006) "The polarization of the U.S. labor market," American Economic Review v96(2), 189–94. https://doi.org/10.1257/000282806777212620
- Autor, D.H., L. F. Katz, and A. B. Krueger (1998) "Computing inequality: Have computers changed the labor market?," *Quarterly Journal of Economics* 113(4), 1169–213. https://doi .org/10.1162/003355398555874
- Bartel, A., C. Ichniowski, and K. Shaw (2007) "How does information technology affect productivity? Plant-level comparisons of product innovation, process improvement, and worker skills," *Quarterly Journal of Economics* 122(4), 1721–58. https://doi.org/10.1162/qjec .2007.122.4.1721
- Baumol, W. (1967) "Macroeconomics of unbalanced growth: The anatomy of urban crisis," American Economic Review 57(3), 806–17
- Bresnahan, T. (1999) "Computing." In D. C. Mowery, ed., U.S. Industry in 2000: Studies in Competitive Performance, ch. 9, pp. 215–44. Washington: National Academy Press
- (2010) "General purpose technologies." In B. H. Hall and N. Rosenberg, eds., Handbook of the Economics of Innovation, vol. 2, ch. 118, pp. 761-91. https://doi.org/10.1016/s0169 -7218(10)02002-2
- (2021) "Artificial intelligence technologies and aggregate growth prospects." In J. W. Diamond and G. R. Zodrow, eds., Prospects for Economic Growth in the United States, pp. 132–70. Cambridge, UK: Cambridge University Press. https://doi.org/10.1017/9781108856089.008
- Bresnahan, T., and S. Greenstein (1996) "Technical progress and co-invention in computing and in the uses of computers," *Brookings Papers on Economic Activity, Microeconomics* 1996, 1–83. https://doi.org/10.2307/2534746
- Bresnahan, T.F., and M. Trajtenberg (1995) "General purpose technologies 'Engines of growth'?," Journal of Econometrics 65(1), 83-108. https://doi.org/10.1016/0304-4076(94)01598-t
- Brynjolfsson, E. (2022) "The Turing trap: The promise & peril of human-like artificial intelligence," *Daedalus* 151(2), 272–87. https://doi.org/10.1162/daed_a_01915
- Brynjolfsson, E., D. Li, and L. Raymond (2023) "Generative AI at work," SSRN Electronic Journal. https://doi.org/10.2139/ssrn.4426942
- Brynjolfsson, E., D. Rock, and C. Syverson (2019) "Artificial intelligence and the modern productivity paradox: A clash of expectations and statistics." In A. Agrawal, J. Gans and A. Goldfarb, eds., *The Economics of Artificial Intelligence*, ch. 1, pp. 23–60. Chicago: University of Chicago Press. https://doi.org/10.7208/chicago/9780226613475.003.0001
- Ceruzzi, P.E. (2003) A History of Modern Computing, 2nd ed. Cambridge, MA: MIT Press
- Cheng, H., R. Jia, D. Li, and H. Li (2019) "The rise of robots in China," *Journal of Economic Perspectives* 33(2), 71-88. https://doi.org/10.1257/jep.33.2.71
- Cockburn, I.M., R. Henderson, and S. Stern (2019) "The impact of artificial intelligence on innovation: An exploratory analysis." In A. Agrawal, J. Gans and A. Goldfarb, eds., *The* Economics of Artificial Intelligence, ch. 4, pp. 115–48. Chicago: University of Chicago Press. https://doi.org/10.7208/chicago/9780226613475.003.0004

376 A. Goldfarb

- David, P.A. (1990) "The dynamo and the computer: An historical perspective on the modern productivity paradox," *American Economic Review* 80(2), 355–61
- Dell'Acqua, F., E. McFowland, E. R. Mollick, H. Lifshitz-Assaf, K. Kellogg, S. Rajendran, L. Krayer, F. Candelon, and K. R. Lakhani (2023) "Navigating the jagged technological frontier: Field experimental evidence of the effects of AI on knowledge worker productivity and quality," SSRN Electronic Journal. https://doi.org/10.2139/ssrn.4573321
- Devine, W.D., Jr. (1983) "From shafts to wires: Historical perspective on electrification," Journal of Economic History 43(2), 347–72. https://doi.org/10.1017/s0022050700029673
- Domingos, P. (2015) The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World. New York: Basic Books
- Doms, M., T. Dunne, and K. R. Troske (1997) "Workers, wages, and technology," Quarterly Journal of Economics 112(1), 253–90. https://doi.org/10.1162/003355397555181
- Eloundou, T., S. Manning, P. Mishkin, and D. Rock (2023) "GPTs are GPTs: An early look at the labor market impact potential of large language models," arXiv. https://doi.org/10.48550 /arXiv.2303.10130
- Epley, N., and T. Gilovich (2016) "The mechanics of motivated reasoning," Journal of Economic Perspectives 30(3), 133-40. https://doi.org/10.1257/jep.30.3.133
- Farrell, J., and M. Rabin (1996) "Cheap talk," Journal of Economic Perspectives 10(3), 103-18. https://doi.org/10.1257/jep.10.3.103
- Forman, C., A. Goldfarb, and S. Greenstein (2012) "The internet and local wages: A puzzle," *American Economic Review* 102(1), 556-75. https://doi.org/10.1257/aer.102.1.556
- Goldfarb, A., B. Taska, and F. Teodoridis (2023) "Could machine learning be a general purpose technology? A comparison of emerging technologies using data from online job postings," *Research Policy* 52(1), 104653. https://doi.org/10.1016/j.respol.2022.104653
- Goldfarb, A., and C. Tucker (2019) "Digital economics," *Journal of Economic Literature* 57(1), 3-43. https://doi.org/10.1257/jel.20171452
- Goldin, C., and L. F. Katz (2008) The Race Between Education and Technology. Cambridge, MA: Harvard University Press
- Goolsbee, A., S. Levitt, and C. Syverson (2013) Microeconomics. Broadway, UK: Worth Publishing
- Greenstein, S. (2015). How The Internet Became Commercial: Innovation, Privatization, and the Birth of a New Network. Princeton, NJ: Princeton University Press
- Hawkins, J., and S. Blakeslee (2007). On Intelligence: How a New Understanding of the Brain Will Lead to the Creation of Truly Intelligent Machines. London: UK: Macmillan
- Jones, C.I. (2023) "The A.I. dilemma: Growth versus existential risk," NBER working paper no. w31837. Available at SSRN: https://ssrn.com/abstract=4624239
- Jorgenson, D. W., M. S. Ho, and K. J. Stiroh (2008) "A retrospective look at the U.S. productivity growth resurgence," *Journal of Economic Perspectives* 22(1), 3-24. https://doi .org/10.1257/jep.22.1.3
- Korinek, A., and J. E. Stiglitz (2019) "Artificial intelligence and its implications for income distribution and unemployment." In A. Agrawal, J. Gans and A. Goldfarb, eds., *The Economics of Artificial Intelligence*, ch. 14, pp. 349–90. Chicago: University of Chicago Press. https://doi.org/10.7208/chicago/9780226613475.003.0014
- Lipsey, R.G., K. I. Carlaw, and C. T. Bekar (2005). Economic Transformations: General Purpose Technologies and Long-Term Economic Growth. Oxford, UK: Oxford University Press
- Markoff, J. (2015). Machines of Loving Grace. New York: HarperCollins
- Mokyr, J. (1990). The Lever of Riches: Technological Creativity and Economic Progress. Oxford, UK: Oxford University Press.
- Mullainathan, S., and Z. Obermeyer (2021) "Diagnosing physician error: A machine learning approach to low-value health care," *Quarterly Journal of Economics* 137(2), 679–727. https://doi.org/10.1093/qje/qjab046
- Nakano, R., J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, X Jiang, K. Cobbe, T. Eloundou, G. Krueger, K. Button, M. Knight, B. Chess, and J. Schulman (2022) "WebGBT: Browser-assisted question-answering with human feedback," arXiv. https://doi.org/10.48550/arXiv.2112.09332

- Nordhaus, W.D. (2007) "Two centuries of productivity growth in computing," Journal of Economic History 67(1), 128-59. https://doi.org/10.1017/s0022050707000058
- (2021) "Are we approaching an economic singularity? Information technology and the future of economic growth," American Economic Journal: Macroeconomics 13(1), 299-332. https://doi.org/10.1257/mac.20170105
- Otis, N., R. P. Clarke, S. Delecourt, D. Holtz, and R. Koning (2024) "The uneven impact of generative AI on entrepreneurial performance. SSRN Electronic Journal. https://doi.org/10 .2139/ssrn.4671369
- Rosenberg, N. (1963) "Technological change in the machine tool industry, 1840–1910," Journal of Economic History 23(4), 414–43. https://doi.org/10.1017/s0022050700109155
- Shapiro, C., and H. R. Varian (1998). Information Rules: A Strategic Guide to the Network Economy. Brighton, MA: Harvard Business Review Press
- Stevenson, B. (2018) "AI, income, employment, and meaning." In A. Agrawal, J. Gans and A. Goldfarb, eds., *The Economics of Artificial Intelligence*, pp. 189-96. https://doi.org/10 .7208/chicago/9780226613475.003.0007
- Sutton, J. (2001). Technology and market structure: Theory and history. MIT Press.
- Trajtenberg, M. (2019) "AI as the next GPT." In A. Agrawal, J. Gans and A. Goldfarb, eds., The Economics of Artificial Intelligence, pp. 175–86. Chicago: University of Chicago Press. https://doi.org/10.7208/chicago/9780226613475.003.0006
- Turing, A. (1950) "Computing Machinery and Intelligence. Mind 49: 433-460.
- Wiles, E., Munyikwa, Z., and Horton, J. J. (2023) "Algorithmic writing assistance on Jobseekers' resumes increases hires. SSRN Electronic Journal. https://doi.org/10.2139/ssrn.4364678
- Zimmermann, F. (2020) "The dynamics of motivated beliefs. American Economic Review, 110(2), 337-363. https://doi.org/10.1257/aer.20180728