

## 79<sup>th</sup> Economic Policy Panel Meeting

4-5 April 2024

# Market Concentration Implications of Foundation Models: The Invisible Hand of ChatGPT

Anton Korinek, Jai Vipra

# Market concentration implications of foundation models: The Invisible Hand of ChatGPT<sup>1</sup>

Jai Vipra and Anton Korinek

## Abstract

We analyze the technological and economic factors shaping the structure of the market for foundation models – large AI models such as those that power ChatGPT that are adaptable to downstream uses. We observe that the market for the most capable foundation models will have a tendency towards concentration, given significant economies of scale and scope, and examine the implications for competition policy and regulation. For frontier models, we discuss how antitrust authorities can address concerns about market competition and vertical integration. For models that are behind the frontier, we expect competition to be quite intense, implying a more limited role for competition policy. Moreover, we discuss how regulation can ensure a level playing field, pass data governance rules, forestall systemic risks from concentration, and internalize safety concerns.

---

<sup>1</sup> We thank Emma Bluemke, Aidan Kane, Sarah Myers West, Sanjay Patnaik, Nicholas Ritter, Max Schnidman, Eli Schrag and Rob Seamans as well as our editors, Emilio Calvano and Giacomo Calzolari, and two anonymous referees for their thoughtful comments. Any remaining errors are our own.

1		
2		
3		
4		
5		
6		
7	<b>1. Introduction</b>	<b>3</b>
8	Principal players in the market for generative AI models	4
9	Recent attention from antitrust authorities	8
10		
11		
12	<b>2. Technological characteristics and market structure</b>	<b>9</b>
13	Cost structure of foundation models	9
14	Compute	10
15	Data	12
16	Vertically integrating economic foundations	13
17		
18		
19		
20	<b>3. Scaling and Concentration Concerns</b>	<b>15</b>
21	Competition among foundation models	15
22	Compute concentration	16
23	Tackling vertical integration	17
24		
25		
26		
27	<b>4. Regulatory concerns</b>	<b>18</b>
28	Ensuring a level playing field with non-AI providers	18
29	Data governance	19
30	Systemic risks from homogenisation	20
31	Safety considerations	21
32	Lobbying	21
33		
34		
35		
36	<b>5. Conclusions</b>	<b>22</b>
37		
38	<b>Bibliography</b>	<b>23</b>
39		
40		
41		
42		
43		
44		
45		
46		
47		
48		
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

## 1. Introduction

In a post titled “Moore’s Law for Everything”, OpenAI CEO Sam Altman predicted that within the next few decades, AI technology would “do almost everything, including making new scientific discoveries that will expand our concept of ‘everything’” (Altman 2021). A paper co-authored by researchers at OpenAI estimated that most occupations are exposed to the deployment of large language models (LLMs) like ChatGPT, and that once complementary investments are made, up to 49 percent of workers could have half or more of their tasks exposed to LLMs (Eloundou et al. 2023). If AI models will indeed play such an important role in our economy, then the structure of the market in which they are offered will have first-order implications for social welfare. Policymakers, including antitrust authorities, thus need to pay close attention to the topic.

Generative AI is powered by foundation models – large AI models that use deep learning methods, are trained on vast amounts of data, and can be adapted to specific economic tasks and applications (Bommasani et al. 2021). Foundation models encompass large language models or vision language models such as OpenAI’s ChatGPT and Google DeepMind’s Gemini that can produce text and process and generate images, audio models that can transcribe or synthesize spoken text or songs, and multimodal models that combine these and other capabilities, for example producing videos and controlling robotic functions.

Building on these foundational capabilities, generative AI models have already shown great promise in performing economically useful tasks, oftentimes much faster and at lower cost than human workers. They can write computer code, find errors in code, write essays, copyedit text, summarize documents, generate ideas, interact with customers, produce images from broad instructions, recognize speech, and so on. There is already evidence that generative AI is disrupting certain categories of work (Hui et al 2023), and a survey of 1,000 US business leaders found that nearly half had replaced workers with ChatGPT (Shani 2023).

The size of the market for foundation models and generative AI was estimated to be just \$3bn in 2023 (Fernandez et al. 2023) - a tiny sliver of a \$100tn world economy. (In 2022, the term “generative AI” had not even been coined yet.) However, leading investment banks project that the technology may underpin 7 to 10 percent of global GDP within a decade (Hatzius et al. 2023; JP Morgan 2024), implying the potential for vast growth.

The market for foundation models is concentrated in large part because developing such models requires substantial fixed costs, primarily consisting of large expenses on computational power as well as skilled talent and large datasets. Training costs for the most advanced foundation models have, on average, grown by a factor of 3.1x/year for the past 15 years (Epoch, 2023). The most expensive publicly known foundation model trained to date is Google DeepMind’s Gemini, at an estimated cost of \$630m (Epoch, 2023). Whether the trend of tripling training costs every year will go on depends on whether the economic benefits of foundation models will continue to scale with the required costs. Extrapolation comes with obvious perils, but many experts predict a

1  
2  
3 continuation of the trend for at least another 3 to 5 years (Suleyman 2023), and - given the  
4 economic promise of foundation models - possibly longer. As we elaborate in more detail below,  
5 a trillion dollar foundation model by the end of the decade is not inconceivable if the trend  
6 continues. In fact, in early 2024, OpenAI's CEO, Sam Altman, reportedly attempted to raise \$7tn  
7 for the production of computer chips needed to train frontier AI models (Hagey and Fitch 2024).  
8  
9

10 The reason why technology companies are willing to spend large amounts on foundation models  
11 is that after the expensive first step, the so-called "pre-training," the model can be used as a  
12 foundation that is adapted to a wide range of downstream economic tasks. This process, called  
13 "fine-tuning," employs task-specific datasets and is significantly cheaper than pre-training  
14 (Bommasani et al. 2021).  
15  
16  
17

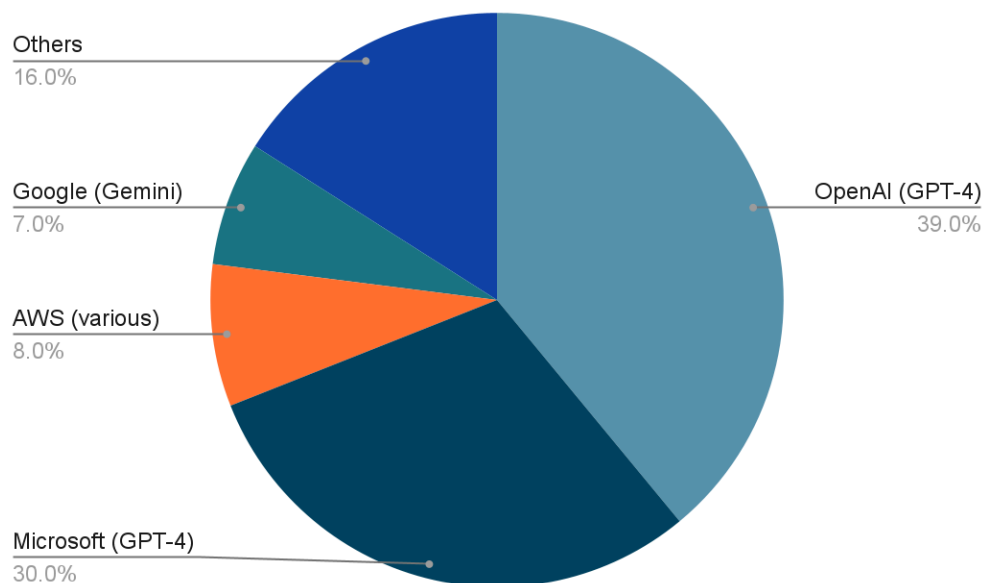
18 There is an active debate in economics that discusses to what extent earlier technological  
19 advances offer useful analogies for the economic effects of artificial intelligence. Eloundou et al  
20 (2023) make the case that foundation models are general purpose technologies akin to the  
21 steam engine or electricity. The argument is that AI has the potential to transform numerous  
22 industries, demonstrates continuous improvement, and creates a powerful need for  
23 complementary technologies and societal adjustments, thereby fulfilling the core criteria of a  
24 general purpose technology.  
25  
26  
27  
28

29 The stated mission of several leading AI labs, including OpenAI, is even grander: that a future  
30 version of their foundation models will be able to achieve artificial general intelligence (AGI),  
31 defined as the ability to perform any cognitive task that humans can perform. If this mission were  
32 to be achieved, then their models could underpin any cognitive work and, if equipped with the  
33 necessary hardware, any human work that humans would let it perform, no matter which  
34 occupation or industry. This maximalist vision of the future role of foundation models is clearly far  
35 from a certainty, but given the rapid pace of recent advances, it may be useful to consider it as a  
36 scenario for which economists and economic policymakers should be prepared (Korinek, 2023).  
37 In such a scenario, the market for foundation models would be the entire economy.  
38  
39  
40

### 41 ***Principal players in the market for generative AI models***

42  
43

44 Figure 1 and Table 1 show data on the market for generative AI models and the capabilities of the  
45 leading AI labs as of Feb 2024 according to the Massive Multitask Language Understanding  
46 (MMLU) benchmark, which measures the world knowledge and problem solving ability of LLMs on  
47 various subjects including sciences and the humanities across 57 tasks (Hendrycks et al 2020).  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



**Figure 1: Principal players in the market for generative AI models (% of total spending)**

Source: Fernandez et al (2023); ‘Others’ includes Anthropic, AI2I Labs, Cohere, Aleph Alpha, Hugging Face, Alibaba, IBM, and Baidu among others.

As illustrated in Figure 1, OpenAI’s GPT-4 series was the clear market leader in 2023. OpenAI’s models powered both the highly successful ChatGPT interface (39% market share) and Microsoft’s offerings of generative AI (30% market share), together accounting for a 69% share of the market for generative AI. GPT-4 completed training in summer 2022, was released in March 2023 (OpenAI, 2023c), and, as indicated in Table 1, has held the title of the most powerful LLM ever since. OpenAI is structured as a non-profit, which majority-owns a for-profit entity, in which outside investors such as Microsoft have invested to help cover the tremendous cost of training frontier AI models. The investments in the for-profit are structured such that the non-profit will receive all profits once outside investors have been repaid an agreed multiple of their initial investment.<sup>2</sup>

The second player in the market is Google DeepMind, which is owned by Alphabet and released a foundation model called Gemini in December 2023. As illustrated in Table 1, Gemini has capabilities on par with GPT-4 (Gemini Team 2023). Google DeepMind is the result of a merger between Google Brain, Google’s internal frontier AI department, and DeepMind, a British AI lab that Alphabet acquired in 2014 (Shu 2014). DeepMind had developed seminal AI models such as AlphaGo and AlphaFold, while Google Brain researchers developed the Transformer architecture

<sup>2</sup> For example, Microsoft’s initial \$1bn investment in 2019 was reportedly subject to a 100x profit cap, i.e., Microsoft will receive dividends up until its initial investment has been repaid one hundred times, before additional profits on this stake would go to the non-profit (Coldewey 2019). The profit cap is far from being reached, implying that investors still have substantial upside in the medium term.

upon which today's LLMs are based. Google entered the commercial market for generative AI systems later than OpenAI/Microsoft and currently only holds a market share of 7%.

AI firm	Country	Top AI model	Release date	MMLU
OpenAI	USA	GPT-4	Mar 2023	90.1
Google DeepMind	USA/UK	Gemini Ultra	Dec 2023	90.0
Mistral	France	Mistral Large	Feb 2023	81.2
Inflection	USA	Inflection-2	Nov 2023	79.6
01.AI	China	Yi (OS)	Nov 2023	78.8
Anthropic	USA	Claude 2	Jul 2023	78.5
Alibaba	China	Qwen 1.5 (OS)	Feb 2024	77.4

**Table 1:** Leading AI labs and their leading foundation models, release date, and score on the MMLU test

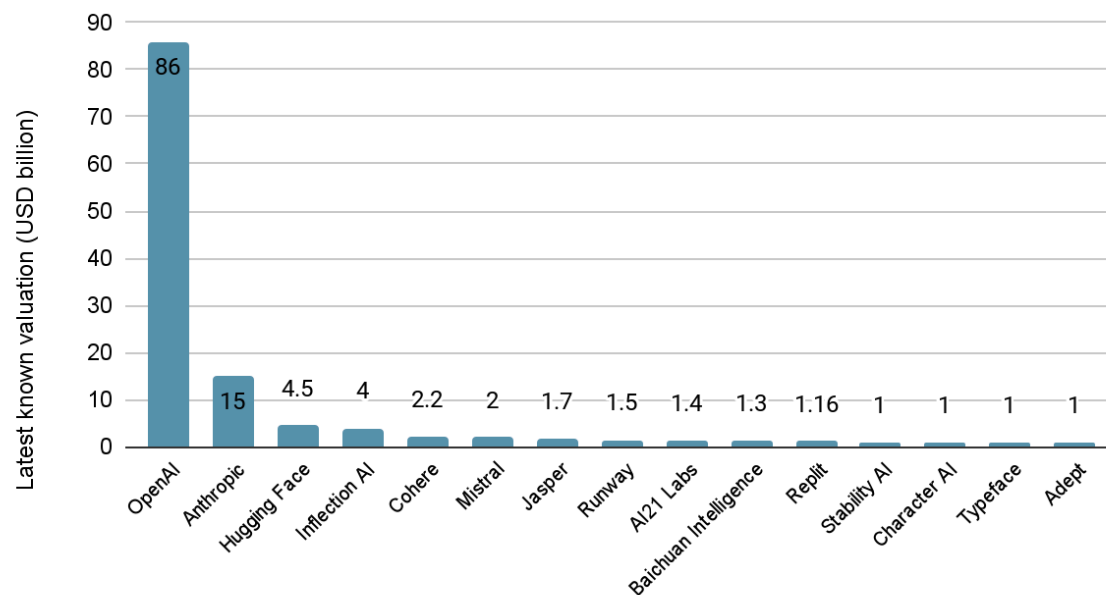
Source: compiled by authors

As illustrated in Table 1, the models of OpenAI and Google DeepMind are far more capable than the competition at the time of writing, with a significant gap in the MMLU score to their closest competitors. Moreover, both of the market leaders are expected to release updated versions of their models with even greater capabilities in 2024.

The market behind these two frontier models is dynamic and competitive, with several companies from the US, France, and China competing neck-on-neck. It also includes a number of open source models, indicated by "OS" in the table. Notably, Facebook/Meta released a series of models named LLaMA starting in 2023, on which outside researchers have built a large number of fine-tuned versions. Meta announced a drive towards building more advanced foundation models in early 2024, and is stockpiling hundreds of thousands of GPUs to this end (Heath 2024). Many of these models are small enough to run on laptops or cell phones. To add to the competition in this market segment behind the frontier, Google DeepMind released an open-source set of models named Gemma in Feb. 2024, which outperform other models in the same size class.

New model announcements that will make it into the ranks of Table 1 are widely expected in the coming months, for example from Meta, Amazon, x.ai, but it is unclear whether any of these models will be able to compete with the frontier models from OpenAI and Google DeepMind, especially as the latter two are themselves actively working on newer versions of their models that are expected to be released in 2024.

## Valuation of leading LLM firms



**Figure 2:** Most recent valuation of leading generative AI labs as of February 2023.

Source: Data collected by authors. (Latest valuation of Cohere not publicly announced.)

Figure 2 shows that there is also significant concentration when leading generative AI labs are measured by their valuation. This figure excludes Google DeepMind, which is a subunit of Alphabet for which valuation information is not publicly available. As a result, OpenAI's valuation stands head-and-shoulders above its competitors.

Company	Investments from large technology companies
OpenAI	Microsoft
Anthropic	Alphabet, Amazon, Salesforce, Zoom
Inflection AI	Microsoft, Nvidia
Hugging Face	Alphabet, AMD, Amazon, IBM, Intel, NVIDIA, Qualcomm, Salesforce
Cohere	Nvidia, Oracle, Salesforce
Mistral	Microsoft

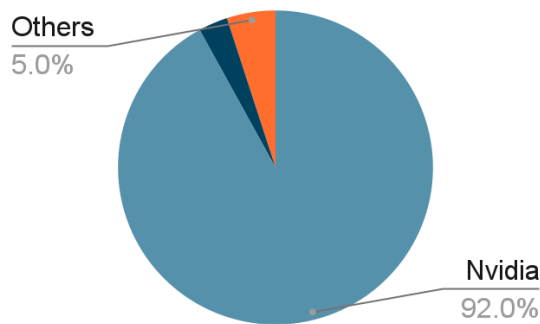
**Table 2:** Investments in leading AI labs by large technology companies

Source: collected by authors

Table 2 illustrates that many of the investments in leading AI labs producing foundation models were conducted by leading large technology companies, including Microsoft, Alphabet, Amazon and Nvidia. It is important to keep track of these linkages when evaluating to what extent large technology companies have control over the market for foundation models.



1  
2  
3  
4 An additional important consideration is that the market for the computer chips that are used to  
5 train and deploy foundation models is highly concentrated. Figure 3 shows that at the end of  
6 2023, the market for Graphical Processing Units (GPUs) was dominated by a single company,  
7 Nvidia, which supplied chips to all the leading producers of foundation models listed in the tables  
8 above. In February 2024, a Wells Fargo report estimated that Nvidia controlled 98 percent of the  
9 data center GPU market (Norem 2024). Vipra and Myers West (2023) provide a rigorous  
10 description of the market for compute.  
11  
12  
13



14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27 **Figure 3:** Market share of leading vendors in the market for GPUs

28 Source: Created by authors based on data from Fernandez et al (2023)

### 30 31 **Recent attention from antitrust authorities**

32  
33  
34 Given their growing importance in our economy, foundation models have already garnered the  
35 attention of antitrust regulators. In June 2023, the FTC released an assessment of competition  
36 concerns in generative AI, in which it highlighted the uneven control over its building blocks like  
37 data, talent, and computational power. It also drew attention to concerns over concentration in  
38 both generative AI markets, and other markets impacted by generative AI (Staff in the Bureau of  
39 Competition & Office of Technology 2023). In January 2024, the FTC launched an inquiry into  
40 generative AI investments and their partnerships with major cloud service providers, asking for  
41 information on partnerships between Microsoft and OpenAI, Amazon and Anthropic, and Google  
42 and Anthropic (Federal Trade Commission 2024).  
43  
44  
45

46  
47 The UK's Competition and Markets Authority (CMA) released a report in September 2023  
48 proposing guiding principles to ensure competition in the market for foundation models,  
49 including access to key inputs like data and computational power, diversity of closed and open  
50 source business models, interoperability, fair dealing, and transparency, among others  
51 (Competition and Markets Authority 2023).  
52  
53

54  
55 The European Commissioner for Competition, Margrethe Vestager, recently stated that merger  
56 control, vertical integration and algorithmic collusion are areas of interest for the EU in relation to  
57  
58  
59  
60

1  
2  
3 AI markets. (Lomas 2024a) The EU is also scrutinizing Microsoft’s investment in OpenAI. (Lomas  
4 2024b)  
5  
6

## 7 **2. Technological characteristics and market structure**

8  
9  
10 This section describes a number of technological and economic forces that characterize the  
11 market for foundation models. We observe that rapidly growing fixed costs generate significant  
12 economies of scale and scope for foundation models and observe the importance of inputs to  
13 production such as computational resources and data. Moreover, we identify strong forces  
14 toward vertical integration.  
15

### 16 ***Cost structure of foundation models***

17  
18  
19  
20 Producing and operating foundation models involves three main types of costs: (i) a significant  
21 fixed cost for the pre-training of foundation models; (ii) an additional fixed cost per area of  
22 application when foundation models are adapted (“fine-tuned”) for specific use cases; and (iii) low  
23 variable costs of operating the model.  
24  
25

26 **Pre-training** is the process of creating a foundation model, which can then be adapted for a wide  
27 range of use cases. The costs of pre-training have risen rapidly in recent years, driven primarily  
28 by growing spending on computational resources (“compute”). *The Economist* (2023) estimates  
29 the cost of training GPT-4 in the first half of 2022 at \$100 million; Epoch (2023) estimates that  
30 Google DeepMind’s Gemini cost around \$630 million, approximately 3/4 of which was spent on  
31 compute, and the remainder covering personnel expenses. We will describe the drivers behind  
32 this rapid growth in spending on compute in the next subsection.  
33  
34  
35

36 **Fine tuning** refers to the process of training a pre-trained model for a specific purpose, usually by  
37 using application-specific data. While exact figures are difficult to obtain, these costs are much  
38 lower than the fixed cost of pre-training a model, because fine-tuning requires less time, data,  
39 and compute. They include wages of in-house workers as well as the cost of outsourced workers  
40 who label data and help train models via reinforcement learning from human feedback (RLHF).  
41  
42  
43

44 One factor that could push up these costs in the future is that fine-tuning requires data that is  
45 labeled, purpose-relevant, and may therefore be costlier to obtain. However, when the producer  
46 of a foundation model provides access to an AI application company, the purpose-relevant data is  
47 usually provided by the latter. For instance, if OpenAI provides access to GPT-4 to a healthcare  
48 provider, the healthcare provider can use its existing data on treatment patterns to fine-tune the  
49 model to its needs.  
50  
51  
52

53 **Variable costs** of operation (“inference costs”) are significantly lower but depend on the specific  
54 application. For example, at the time of writing, OpenAI charges \$0.03 to generate 1000 syllables  
55 (“tokens”) of text using its GPT-4 model – a number that is estimated to be close to the company’s  
56  
57  
58

1  
2  
3 cost of inference. However, one operative question for deploying foundation models in the  
4 broader economy is how many inferences per hour are needed to replace a human worker, and  
5 how much this would cost.  
6  
7

8 Since fixed costs are high and variable costs relatively low, foundation models offer a classical  
9 example of economies of scale. Moreover, the general purpose nature of foundation models also  
10 gives rise to significant economies of scope. Since one foundation model can be adapted to  
11 many different areas of application across different industries, economies of scope are expected  
12 to be very large. For instance, a single foundation model such as GPT-4 or Gemini can be used to  
13 automate copyediting, to create holiday itineraries, to check for errors in computer code, or to  
14 provide health advice. More generally, OpenAI, for instance, already allows developers to create  
15 plugins or custom GPTs to apply GPT-4 to many different applications and benefit from these  
16 economies of scope, allowing users to order groceries, search for flights, learn languages, and  
17 shop online, all by using the same foundation model (OpenAI 2023b).  
18  
19  
20  
21

## 22 **Compute**

23  
24

25 The computational resources (frequently abbreviated by the neologism “compute”) required to  
26 train frontier foundation models are massive. In recent years, AI researchers identified “scaling  
27 laws” for foundation models that predict how model performance increases with the amount of  
28 compute used in training the model (see e.g., Kaplan et al. 2020; Hoffman et al., 2022). To make  
29 better models, developers have to add more computational power. These regularities are  
30 important because they reduce the uncertainty that companies face when they make investment  
31 decisions.  
32  
33  
34

35 Building on the described scaling laws, leading AI labs have increased the amount of compute  
36 deployed in frontier AI models by a factor of 4.1x per year over the past 15 years, as illustrated in  
37 Figure 4. Given reductions in chip prices, which evolved approximately in line with Moore’s Law,  
38 Cottier (2023) finds that spending on compute for frontier AI systems has grown by a factor of  
39 3.09x per year over the same period.<sup>3</sup>  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

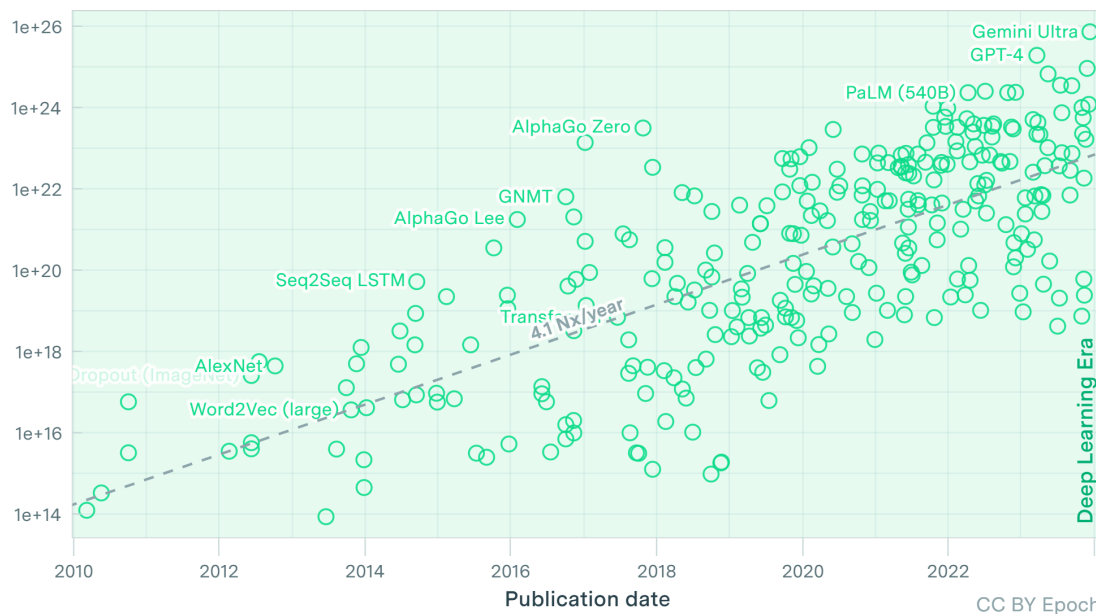
---

54 <sup>3</sup> Estimates of compute costs for pre-training AI models usually take into account only the compute  
55 required for the final training run, and not the compute employed by trial training runs used to arrive at a  
56 workable architecture, which implies even greater compute costs (Heim 2021).  
57  
58  
59  
60

## Training Compute of Notable machine learning Systems Over Time



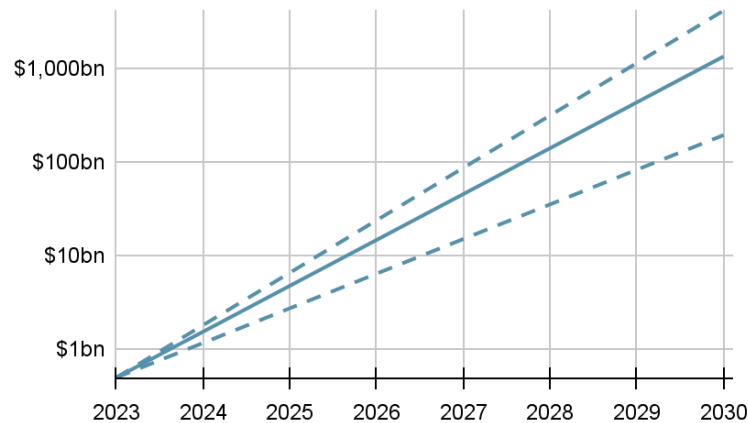
Training compute (FLOP)



**Figure 4:** Training compute of notable AI systems. Copyright © by Epoch.org, reproduced under a CC-BY-4.0 license, 2024.

Whether this trend will go on will depend on whether not only the capabilities but also the economic benefits of foundation models will continue to scale with the amount of spending. Extrapolation comes with obvious perils - at first, any successful new technology grows faster than the rest of the economy as it starts from zero. Eventually, no sector can grow faster than the economy as a whole. This implies that we would expect the growth of AI training costs to eventually slow down as they become a more and more significant part of the economy. However, concurrently, economic growth may take off if advanced foundation models can replicate a growing fraction of the tasks that were traditionally performed by scarce labor (see, e.g., Aghion et al. 2019 and Trammell and Korinek 2023).

These considerations notwithstanding, technology leaders such as Suleyman (2023) and analysts predict a continuation of the described trend for at least another 3 to 5 years, and - given the economic promise of foundation models - possibly longer. We illustrate the implications in Figure 5, in which we start with the approximate compute training cost of Gemini and visualize the 3.09x/year trend growth of Cottier (2023), with lower and upper 90% confidence bounds given by growth factors of 2.34x and 3.63x per year respectively. If these rates continue, a naive extrapolation implies that it is not inconceivable that a technology company may spend more than \$1tn on a single foundation model before the end of the decade. This is also consistent with reports that OpenAI's CEO, Sam Altman, attempted to raise \$7tn for the production of computer chips needed to train frontier AI models in early 2024 (Hagey and Fitch 2024).



**Figure 5:** Extrapolation of Training compute of frontier AI systems, based on estimates from Cottier (2023).

Recent increases in demand for computational power may drive up costs even further. The costs of manufacturing additional equipment for producing computer chips are high, as they rely on highly specialized manufacturing techniques and high fixed costs (Khan, Peterson, and Mann 2021). Moreover, the semiconductor design and manufacturing industry is highly concentrated, as illustrated in Figure 3. Vipra and Myers West (2023) describe the market for compute in further detail. (We should probably cut this in the final version, but one of the world’s leading technology experts, Elon Musk, proclaimed that “GPUs at this point are considerably harder to get than drugs,” WSJ 2023).

There are also forces that pull in the opposite direction. The high returns generated by the scarcity of computer chips as well as government subsidies provide strong incentives for investment in additional capacity and may help to diversify the semiconductor supply chain and bring down costs. New breakthroughs in computing technology, such as neuromorphic computing, quantum computing, or improvements in memory technology, might ease concentration in the market for chips.

Overall, however, we expect that access to compute will continue to be a bottleneck for training frontier foundation models. This will likely restrict access to the compute needed for state-of-the-art models to a small number of well-funded AI labs like OpenAI and Google DeepMind. By contrast, for less advanced models, compute requirements are lower and more attainable so other companies, ranging from large tech firms to AI startups, will compete in the broader market for less advanced foundation models. This two-tiered structure may lead to some market power for the leading firms offering top-tier models, while we expect significant competition for the residual demand for lower-capability models.

## **Data**

Foundation models are trained on large quantities of data. The most data-intensive language model at the time of writing was Google's FLAN, trained on 1.87 trillion words (Epoch 2023). Models are data-hungry and achieve higher performance when trained on higher quality text (Anil et al. 2023). Yet we are about to run out of high-quality text that is publicly available on the internet soon (Villalobos et al. 2022).

This makes proprietary datasets very useful for companies producing foundation models, which gives an advantage to large technology companies, which control a large fraction of the data generated online – including platform interactions, search, emails, photos, videos and other documents. Access to data is therefore an important reason why producers of foundation models may be interested in vertically integrating with Big Tech companies.

Over the course of 2023, a growing number of content providers on the internet have restricted free access to their content to the bots that automatically scoop up data used for the training of foundation models. For example, Fletcher (2024) finds that 79% of all news sites in the US blocked OpenAI's crawlers, as did about half of all news sites in a sample of ten advanced countries. The New York Times has even filed a copyright infringement lawsuit against OpenAI and Microsoft for using its news articles to train OpenAI's foundation models.

There are emerging techniques that can reduce the cost of acquiring data. These include simulation learning (where a simulated environment substitutes for a real training environment), self-play (where a model can interact with itself to improve its performance), and synthetic data generation (Hwang 2018; Azizi et al. 2023). Some of these techniques are domain-specific and do not work for all foundation models; some, like synthetic data generation, decrease accuracy. When such techniques are successful, they increase compute requirements, effectively substituting compute for training data.

Fine-tuning foundation models for specific applications also requires data – for instance, a manufacturer will be able to derive the most value from a foundation model by fine-tuning it on its own historical data. Such data may encompass detailed records ranging from production plans, machinery, product quality, to supply chain logistics, offering a granular view into every aspect of their operations. By leveraging this multi-dimensional dataset, the manufacturer can unlock insights for automating processes, reducing downtime, and enhancing product quality, thereby achieving a competitive edge in the industry. In many sectors, specific proprietary data held by traditional companies will be a valuable source of training data for fine-tuning purposes.

## ***Vertically integrating economic foundations***

Producers of foundation models may engage in vertical integration both upstream and downstream. Upstream vertical integration refers to integration with suppliers of inputs to

1  
2  
3 producing foundation models, the most important of which are compute and data. Downstream  
4 vertical integration refers to integration with producers or distributors of final goods and services  
5 that employ foundation models.  
6  
7

8 Vertical integration may increase consumer welfare by eliminating multiple markups in the  
9 production process, ultimately leading to lower final prices. Furthermore, it may mitigate the risks  
10 associated with specialized assets stemming from hold-up problems. An integrated firm may  
11 reduce transaction costs by maintaining better quality control over its products and streamlining  
12 operations, potentially leading to both lower prices and higher quality goods for consumers.  
13 However, by consolidating market power, vertical integration may also pose the risk of  
14 foreclosing rivals' access to essential inputs or distribution channels or enabling firms to impose  
15 vertical restraints, thereby increasing monopoly distortions. It may also reduce innovation by  
16 lowering the number of independent players competing on new technologies and ideas (see, eg,  
17 Tirole 1988, Bresnahan and Levin 2012).  
18  
19  
20  
21

22 **Upstream** Since compute has the largest input share in creating foundation models, a growing  
23 number of AI companies have vertically integrated the production of chips and foundation  
24 models. An example of this is Google DeepMind, which produces its own chips, termed Tensor  
25 Processing Units (TPUs), which are particularly efficient for AI training. Similarly, the leading chip  
26 provider for training AI models, NVIDIA, is also offering foundation models together with a  
27 training platform that allows customers to fine-tune these models for their own purposes. Amazon  
28 and Microsoft, two of the leading providers of cloud computing, are both working on in-house  
29 designed chips for AI applications.  
30  
31  
32  
33

34 Another type of vertical integration occurs when providers of cloud computing services and  
35 foundation models enter exclusive contracts. Microsoft's investment in OpenAI gave Microsoft the  
36 right to be the exclusive provider of cloud services for OpenAI. This provision covers both training  
37 and provision of the model (Warren 2023). Likewise, Google Cloud is the preferred provider of  
38 cloud services for Anthropic ('Anthropic Partners with Google Cloud' 2023). Many of the  
39 investments listed in Table 2 are of a similar nature.  
40  
41  
42

43 With data an important input in pre-training foundation models, there is also significant potential  
44 for vertical integration between data-rich technology companies and producers of foundation  
45 models. Google, Meta and X/Twitter are employing data from their public platforms like (including  
46 YouTube, Facebook, Instagram, Twitter) in pre-training their foundation models (Victor 2023, Olson  
47 2024b). At the same time, data-rich technology companies are restricting data access to  
48 web crawlers that allow other producers of foundation models to train on their data (Verhulst,  
49 2023), raising anti-competitive concerns.  
50  
51  
52

53 **Downstream** As foundation models have a growing number of general use cases for cognitive  
54 workers, they are rapidly being integrated into office suites. For example, Google and Microsoft  
55 have both integrated generative AI capabilities into their products, including Microsoft Office,  
56  
57  
58  
59  
60

1  
2  
3 Gmail, and Google Documents. Google also offers a beta version of search that incorporates  
4 generative AI (Reid 2023).  
5

6  
7 OpenAI is integrating its GPT-4 model in a growing number of downstream uses by allowing  
8 commercial providers to create custom GPTs that are adapted and fine-tuned for specific use  
9 cases. Since their launch in November 2023, more than a hundred thousand custom GPTs have  
10 been created that link ChatGPT to commercial applications, allowing users access to additional  
11 third-party functions that range from travel itineraries to incorporating Wolfram Mathematica in  
12 their chats. This may be the first step of turning ChatGPT into a platform.  
13  
14

15  
16 As the capabilities of foundation models improve and as companies devise strategies for how to  
17 effectively incorporate them into their operations, downstream use cases and by extension the  
18 scope for vertical integration may continue to grow.  
19  
20

### 21 **3. Scaling and Concentration Concerns**

#### 22 ***Competition among foundation models***

23  
24 At the time of writing in early 2024, startups and leading technology companies around the world  
25 are racing to develop competing foundation models and generative AI offerings as fast as  
26 possible, and fierce competition has driven the price of powerful foundation models literally to  
27 zero: Part of the strategy of firms that are attempting to catch up to market leaders, e.g., Meta,  
28 has been to open-source their model releases, i.e., to make them available to download for free  
29 to anyone who wants to use them. Although these are not models at the frontier like the ones  
30 listed in Table 1, open sourcing comes with significant economic benefits – pre-trained foundation  
31 models are non-rivalrous so free distribution of the model corresponds to the first-best price and  
32 the maximum level of consumer surplus. Moreover, open sourcing also allows researchers and  
33 other companies to build on the foundation model and fine-tune it, encouraging innovation.  
34 However, open sourcing highly capable foundation models also carries safety risks, especially if  
35 models become more powerful, as we will discuss in the last section.  
36  
37  
38  
39  
40  
41  
42

43 OpenAI and Google DeepMind, the producers of the two leading frontier models described in  
44 Table 1, are also engaged in fierce competition, charging prices that barely allow them to cover  
45 their variable costs (Knight 2024). The competition dynamics between the two seem to be close  
46 to Bertrand competition.<sup>4</sup> If competition dynamics remained as fierce as they currently are, then  
47 there would be little reason for concern about market concentration.  
48  
49  
50  
51  
52

---

53  
54 <sup>4</sup> There are even some indications that inference costs may currently be priced below cost (Patel and  
55 Nishball 2023). However, determining whether pricing is below cost is tricky in nascent industries with high  
56 fixed costs and uncertain demand.  
57  
58



1  
2  
3 However, the investment requirements for frontier foundation models are expected to rise rapidly,  
4 making it more difficult for a large number of players to compete in the market. Moreover, as  
5 described in Section 2, the role of generative AI in the economy is expected to grow significantly  
6 from its current \$3bn market size (Hatzius 2023). Moreover, the current top players enjoy  
7 significant first-mover advantages from their technological leadership, from the preemption of  
8 scarce assets, and from buyer switching costs (Lieberman and Montgomery 1988). The scarce  
9 assets accumulated by the pioneers include vast amounts of compute and datasets. Whether  
10 companies like OpenAI can maintain their first-mover advantage remains to be seen.<sup>5</sup> However,  
11 taken together, these trends imply that it is important for policy to ensure that the market for  
12 foundation models remains competitive.<sup>6</sup>  
13  
14  
15

### 16 17 **Compute concentration** 18

19  
20 Whereas the market for foundation models itself is fiercely competitive, the market for chips on  
21 which the training and inference of such models relies is highly concentrated, as illustrated in  
22 Figure 3. The creation of cutting-edge chips is a process that is highly sophisticated and involves  
23 massive R&D costs, giving rise to a complex supply chain. Moreover, from the design, the  
24 manufacturing of equipment to fabricate chips, to the fabrication itself, there are several  
25 companies that are close to monopolists in their respective functions.  
26  
27

28  
29 Due to the complexity of the supply chain for computational power, increasing the  
30 competitiveness of computational power markets is extraordinarily difficult. Regulators can still  
31 ensure that this market is not further concentrated. Belfield and Hua (2022) discuss classical  
32 options such as disallowing some mergers and acquisitions among chip companies, seeking  
33 remedial action from entities wishing to merge; targeting abuse of dominance, including prices  
34 that are too high (although this would be difficult to monitor & implement in practice), bundling,  
35 self-preferencing, etc.; and ensuring that horizontal agreements relating to hardware standards  
36 are not anti-competitive.  
37  
38

39  
40 Vipra and Myers West (2023) propose that structural separations in the ecosystem for  
41 computational resources can help reduce the concentration in this industry. They recommend  
42  
43

---

44  
45 <sup>5</sup> After the release of ChatGPT, fierce technological competition also seems to have changed the open  
46 publication norms of the sector, making it more difficult for new entrants to catch up. For instance, Google  
47 changed its long-standing practice of releasing AI research papers freely, opting instead to hold back  
48 critical model details (Tiku and Vynck 2023).

49 <sup>6</sup> We should also note that some capability improvements in the future could change the prevailing  
50 competition dynamics entirely. If a lab such as OpenAI manages to develop a foundation model that has  
51 the ability to improve itself without human input, then the pace of technological progress could be rapid  
52 and sustained. This would cause first mover advantages to snowball, leave the competition ever further  
53 behind, and create a large monopoly. If such a breakthrough is achieved by a non-profit entity, then there is  
54 the potential that it will maximize social welfare rather than engaging in inefficient monopolistic behavior.  
55 However, although OpenAI is notionally owned by a non-profit, it is unclear what considerations will drive  
56 the organization's behavior.  
57  
58

1  
2  
3 that antitrust regulators examine the viability of separating cloud provision from chip design,  
4 compute hardware from compute software (like NVIDIA's CUDA), and model development from  
5 cloud infrastructure. They propose that these measures could reduce lock-in to compute  
6 ecosystems.  
7

### 8 9 ***Tackling vertical integration***

10  
11  
12 As the role of foundation models in the economy grows, their widespread potential applicability  
13 implies that vertical integration may become a growing concern for competition. Merger review in  
14 this area would have to pay attention to how any proposed mergers affect the cost of inputs to  
15 this market, and whether they might privilege some participants in downstream markets to the  
16 detriment of others.  
17

18  
19  
20 Antitrust oversight will also be required even when coordination does not take the form of  
21 mergers. For instance, exclusive and preferred use contracts among computational power  
22 providers and foundation model companies are rife – as in OpenAI's deal with Microsoft that  
23 gives the latter exclusivity, or Anthropic's deal with Google to use its cloud services, raising  
24 competition concerns and potentially foreclosing the market for other entrants.  
25

26  
27  
28 Moreover, early access provisions by foundation model producers may also raise competition  
29 concerns that may grow in importance as these models are deployed throughout the economy.  
30 For example, OpenAI offered beta access to GPT-4 to certain companies (including Duolingo,  
31 Stripe and Morgan Stanley), which privileged them in relation to their competitors.  
32

33  
34 Antitrust regulators should pay keen attention to acquisitions made by foundation model  
35 companies, especially of startups that might compete with them. Large technology companies  
36 are reportedly already wary of making acquisitions in AI, and are using investments instead  
37 (Olson 2024a). Antitrust scrutiny in investments is also growing. Governments should equip  
38 antitrust regulators with stronger ex ante powers to stop acquisitions that can be shown to  
39 significantly reduce competition, although how to deal with 'nascent competitors' is tricky in  
40 practice (Hemphill and Wu, 2020). Ex post measures in digital markets are often too little, too late  
41 as important intellectual property can be transferred before a merger or acquisition is undone.  
42

43  
44  
45 Moreover, as foundation models become more and more integrated into the economy, they may  
46 provide intellectual infrastructure for a wide range of economic functions and play a similar role to  
47 public utilities like electricity. Regulators may need to respond to certain concerns that arise from  
48 the wide range of their uses. For example, anyone excluded from their services may be at great  
49 economic disadvantage, creating a case for instituting non-discrimination requirements for  
50 access to foundation models. In the US, public utility law prohibits undue or unreasonable price  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 discrimination, requiring that similar customers receiving similar services pay the same prices  
4 (Henderson and Burns 1989).<sup>7</sup>  
5  
6

7 When foundation models morph into platforms, for example through GPT stores,  
8 non-discrimination requirements would prevent foundation model companies from privileging  
9 their own downstream products and services over other products and services. For example,  
10 Khan (2017) argues that the use of the essential facilities doctrine – compelling a monopolist to  
11 provide easy access to competitors in an adjacent market – can be apt in such situations.  
12  
13

## 14 15 **4. Regulatory concerns** 16

17 In addition to antitrust concerns, the tendency of foundation models towards concentration blunts  
18 market forces and may make it desirable to impose a number of additional regulations:  
19  
20

### 21 ***Ensuring a level playing field with non-AI providers*** 22

23  
24 As foundation models are deployed throughout the economy in a growing number of different  
25 functions, they are likely to compete with non-AI providers (including human providers) of  
26 different products and services. In those instances, ensuring a level playing field with non-AI  
27 providers will be essential.  
28  
29

30 The law is often silent on the liability that AI solutions carry when they engage in the real world,  
31 especially in sectors where these solutions are not already common. There is no reason why  
32 foundation models should be exempt from sectoral regulations including liability, professional  
33 licensing, and professional ethics guidelines. If these regulations are derived from various rights  
34 of consumers and citizens, they should apply in appropriate ways to foundation models as well.  
35  
36  
37

38 For instance, governments worldwide have extensive regulations in the education sector, such as  
39 around student privacy, non-discrimination, and educational standards. AI solutions used for the  
40 classroom should explicitly be subject to the same regulatory standards. Sometimes this will  
41 require clarifications or amendments in regulation, because it is difficult to apply regulation in the  
42 same way to all technologies, or to technology and humans.  
43  
44

45  
46 There are cases where human workers are penalized for discounting the analysis of an AI  
47 solution in their workplace, creating a lopsided liability burden. For instance, nurses in some US  
48 hospitals can disregard algorithmic assessment of a patient's diagnosis with doctor approval, but  
49  
50

---

51 <sup>7</sup> Section 205(b) of the Federal Power Act of 1920, for example, prohibits undue preference or prejudice to  
52 any person, as well as unreasonable differences in rates, charges, service and facilities. See  
53 [https://www.ferc.gov/sites/default/files/2021-04/federal\\_power\\_act.pdf](https://www.ferc.gov/sites/default/files/2021-04/federal_power_act.pdf). In the EU, Article 10 of Directive  
54 2002/19/EC gives the national regulatory authority power to impose access requirements on  
55 communications networks. See  
56 <https://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=OJ:L:2002:108:0007:0020:EN:PDF>.  
57  
58

1  
2  
3 face high risks for such disregard as they are penalized for overriding algorithms that turn out to  
4 be right. This may lead nurses to err on the side of caution and follow AI solutions even when  
5 they know they are wrong in a given instance (Bannon 2023).  
6  
7

8 To avoid situations where humans defer to AI against their better judgment, liability frameworks  
9 should be neutral to ensure that technology follows sectoral regulation and not the other way  
10 round. AI technology should not be applied in circumstances in which it does not meet regulatory  
11 standards.  
12  
13

14 Such requirements might be onerous and promote more concentration, and the burden to  
15 compensate for these effects must be borne by antitrust authorities. Moreover, such  
16 requirements would make the deployment of unsafe and substandard systems less financially  
17 attractive. It is also possible that enforcing sectoral standards promotes the development of AI  
18 systems that are more sector-specific than general, potentially leading to more competition. More  
19 importantly, we expect these requirements to protect people against the degradation of product  
20 and service standards, and the erosion of consumer rights, due to the use of AI.  
21  
22  
23

### 24 ***Data governance***

25  
26  
27 Data governance can aid in preventing undesirable concentration in foundation model markets  
28 but can also exacerbate monopoly power. Existing data protection law can institute important  
29 limits on the exclusive accumulation of data resources required to build large scale models. This  
30 may become more important as training data increasingly moves from freely available web data  
31 to proprietary datasets.  
32  
33  
34

35 Purpose limitation laws require that data only be used for the purpose for which it was collected;  
36 a challenge to the use of proprietary data through this route is possible, given the multiple uses  
37 to which a foundation model can be put. Platforms that collect data might have to collect consent  
38 separately for the use of data through foundation models. Seen a little differently, purpose  
39 limitations can reduce the economies of scope of foundation models by limiting their legal use  
40 cases.  
41  
42  
43

44 Many data governance regulations aim at restricting the transfer of data between firms, as  
45 exemplified by the EU's General Data Protection Regulation. Although this may be desirable for  
46 reasons such as privacy, it increases the market power of the firms with the most data, typically,  
47 the dominant firms. It also provides new incentives for vertical integration. As a solution, antitrust  
48 regulators could mandate that foundation model companies and cloud companies hold certain  
49 data in silos, i.e., not share it for other uses (including new foundation models) or with other  
50 companies within their corporate groups (Competition and Markets Authority 2020).  
51  
52  
53

54 Data portability rights allow users to port their personal data from one platform to another at no  
55 or low cost. Some regulations and proposals also create a framework on mandatory data sharing  
56  
57  
58

1  
2  
3 or interoperability to ensure a level playing field in data-driven markets. For instance the EU's  
4 Revised Payment Services Directive creates a data sharing and interoperability framework to  
5 level the playing field between banks and new payment operators. Some proposals recommend  
6 mandatory data sharing to break the monopoly in the online search market (Martens 2023).  
7  
8

### 9 ***Systemic risks from homogenisation***

10  
11  
12 As foundation models are integrated into production and delivery processes for goods and  
13 services across many sectors of the economy, Bommasani et al. (2021) point out that they may  
14 give rise to systemic risks that arise from homogenisation. For example, assume that one  
15 foundation model in its fine-tuned versions is powering decision-making processes in search,  
16 market research, customer service, advertising, design, manufacturing, and many more.  
17 Widespread, cross-industrial applications mean that any errors, vulnerabilities or failures in a  
18 foundation model can threaten a significant fraction of economic activity, producing the risk of  
19 systemic economic effects that may warrant regulation.  
20  
21  
22

23  
24 Biases or errors in the outputs of a foundation model can be inherited by downstream models.  
25 These might not be noticed until the model is in use, and its effects might be observed  
26 throughout the economy. An example scenario is if an LLM understands demand patterns for  
27 consumer goods based on gender stereotypes, and consequently leads to over- or  
28 under-production of certain goods.  
29  
30

31  
32 Foundation models could also be vulnerable to malicious attacks such as through data poisoning.  
33 Data poisoning refers to the practice of tampering with training data in order to influence the  
34 outputs of a model. These outputs can then be reflected in the outputs of fine-tuned models as  
35 well. Cybersecurity and data auditing methods at the pre-training level can protect against the  
36 risk of data poisoning.  
37  
38

39  
40 Foundation models may exacerbate income inequality through automation that leads to  
41 centralization of intellectual functions. Automating tasks with foundation models can result in a  
42 single entity monopolizing entire skill sets if deployment is highly centralized, reducing  
43 competition. Meanwhile, unequal access to frontier foundation models internationally gives  
44 production advantages to countries where they are developed, disadvantaging others. This could  
45 worsen income inequality. Policies promoting decentralized access to foundation models and  
46 preventing monopolization of automated skills may mitigate these risks.  
47  
48

### 49 ***Safety considerations***

50  
51  
52 Current foundation models are, for most practical purposes, quite safe to use and deploy.  
53 However, AI experts warn that future more capable AI systems could cause wide-ranging and  
54 perhaps even irreversible harm to society (Hinton et al 2023; Anderljung et al. 2023). These  
55  
56  
57  
58

1  
2  
3 concerns include the potential malfunctioning of powerful AI systems as well as the malicious use  
4 of such AI systems.  
5

6  
7 AI safety research is particularly concerned about the potential development of artificial general  
8 intelligence (AGI), which is the explicit goal of both OpenAI and Google DeepMind. AGI refers to  
9 technology that can perform all cognitive tasks that humans can perform, or at least all  
10 economically relevant tasks. Many problems could potentially arise from the further advancement  
11 of AI if it proceeds in a reckless and unsafe manner, including misuse by unscrupulous actors,  
12 catastrophes arising from errors in understanding human preferences, mass disruption of labor  
13 markets, or even disempowerment or extinction (Hendrycks et al. 2022). AI safety risk – the risk  
14 of unsafe frontier AI systems being developed and deployed – interacts with market structure  
15 and competition in several ways, which we explore in the following section.  
16  
17  
18

19  
20 The market structure of foundation models, whether concentrated or competitive, presents  
21 tradeoffs for AI safety efforts:  
22

23  
24 **Safety risks from competition** Competitive pressures in AI development may undermine safety  
25 efforts. Breakneck competition could accelerate unsafe foundation model advances before  
26 society installs safeguards. Competitive dynamics may also incentivize developers to cut corners  
27 on safety research in order to get ahead. However, competition has also encouraged open  
28 sourcing of some foundation models, which raises separate safety issues due to potential  
29 misuse.<sup>8</sup> Overall, unchecked competition risks rapidly advancing AI in dangerous directions  
30 without enough oversight.  
31  
32

33  
34 **Safety risks from market concentration** Market concentration also poses AI safety risks.  
35 Monopolies with few competitors can amass resources more quickly to develop more powerful  
36 models, with greater potential for misuse. Concentration likewise may reduce incentives for  
37 safety precautions. However, a monopoly held by a responsible firm may also allow more  
38 controlled AI advancement. Ultimately, whether concentration or competition prevails,  
39 responsible development and strong oversight are essential to mitigate the safety hazards from  
40 increasingly capable AI systems.  
41  
42

### 43 ***Regulatory capture***

44  
45  
46  
47 A concentrated market for foundation models, combined with widespread application of  
48 foundation models, implies high financial stakes for foundation model companies. This makes it  
49 likely that both antitrust rules and the other regulations we discussed above will be subject to  
50 lobbying efforts and to growing risks of regulatory capture. While regulatory capture is typically  
51

---

52  
53 <sup>8</sup> The open-source release of Meta's Llama model has already enabled programmers to use Llama's  
54 published model weights to create novel and targeted adversarial attacks for LLMs that make it possible to  
55 circumvent the safety restrictions built into all other leading LLMs (e.g., ChatGPT, Bard, Claude), thereby  
56 seriously undermining those organizations' attempts to keep their AI systems safe (Zou et al. 2023).  
57  
58

1  
2  
3 viewed as the risk of firms co-opting regulators to make it more difficult for competitors to enter a  
4 market, it is also possible that lobbying ensures that certain externalities are *not* regulated, if this  
5 is in the interests of current market participants. The risk that the lobbying power of foundation  
6 model producers may rise over time is one reason why proactive regulation – before producers  
7 become too powerful to stave off regulations that they find undesirable – may be indicated.  
8  
9

## 10 11 **5. Conclusions**

12  
13 Foundation models powered by deep learning techniques have demonstrated rapid advances in  
14 capability over the past decade. As these AI systems grow more powerful, their potential  
15 economic impact expands as well.  
16  
17

18  
19 However, the market structure for developing and deploying these models tends toward  
20 concentration, given the economies of scale arising from large pre-training costs and bottlenecks  
21 in key inputs like data and compute. This raises competition concerns if too few players come to  
22 dominate the provision of foundation systems across industries, which policymaker should aim to  
23 counteract. It also poses new regulatory challenges, including for data governance, novel  
24 systemic risks, and safety.  
25  
26

27  
28 If left unchecked, the trajectory of foundation models is toward one or a few dominant firms  
29 providing the AI substrate for major parts of the global economy. Preventing extreme  
30 concentration via proactive policies today is preferable to untangling monopolies after the fact.  
31 Proactive competition policies may also help to distribute the productivity gains from AI more  
32 equitably to avoid worsening inequality. With prudent competition and regulatory policies, we are  
33 hopeful that society will be able to steer this powerful new technology toward broadly shared  
34 prosperity.  
35  
36  
37  
38  
39  
40  
41

## 42 **Disclosures**

43  
44 Anton Korinek is a Professor of Economics at the University of Virginia and the Darden School of  
45 Business as well as a Nonresident Fellow at the Brookings Institution, a Research Associate at the  
46 NBER, a Research Fellow at the CEPR and the Economics of AI Lead at the Centre for the  
47 Governance of AI. Jai Vipra received financial support as a Winter Fellow from the Centre for the  
48 Governance of AI while writing this paper. The Centre for the Governance of AI is an independent  
49 research organization that is dedicated to helping humanity navigate the transition to a world with  
50 advanced AI. For further information, see <https://www.governance.ai/about-us>. The views  
51 expressed in this paper are the personal views of the authors and do not necessarily represent  
52 the views of the institutions they are affiliated with.  
53  
54  
55  
56  
57  
58  
59  
60

## Bibliography

- Altman, Sam. 2021. "Moore's Law for Everything." Blog post. <https://moores.samaltman.com/>
- Anil, Rohan, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, et al. 2023. 'PaLM 2 Technical Report'. arXiv. <https://doi.org/10.48550/arXiv.2305.10403>.
- 'Announcing Google DeepMind'. 2023. Google DeepMind. 20 April 2023. <https://www.deepmind.com/blog/announcing-google-deepmind>.
- 'Anthropic Partners with Google Cloud'. 2023. Anthropic. 3 February 2023. <https://www.anthropic.com/index/anthropic-partners-with-google-cloud>.
- Azizi, Shekoofeh, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J. Fleet. 2023. 'Synthetic Data from Diffusion Models Improves ImageNet Classification'. arXiv. <https://doi.org/10.48550/arXiv.2304.08466>.
- Bannon, Lisa. 2023. 'When AI Overrules the Nurses Caring for You'. *The Wall Street Journal*, 15 June 2023. <https://www.wsj.com/articles/ai-medical-diagnosis-nurses-f881b0fe>.
- Belfield, Haydn, and Shin-Shin Hua. 2022. 'Compute and Antitrust: Regulatory implications of the AI hardware supply chain, from chip design to cloud APIs'. *Verfassungsblog*, August. <https://verfassungsblog.de/compute-and-antitrust/>.
- Boldrin, Michele, and David K. Levine. 2009. 'Does Intellectual Monopoly Help Innovation?' *Review of Law & Economics* 5 (3): 991–1024. <https://doi.org/10.2202/1555-5879.1438>.
- Bommasani, Rishi, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, et al. 2021. 'On the Opportunities and Risks of Foundation Models'. arXiv. <https://doi.org/10.48550/arXiv.2108.07258>.
- Bresnahan, Timothy F., and Manuel Trajtenberg. 1995. 'General Purpose Technologies 'Engines of Growth?'' *Journal of Econometrics* 65(1), pp. 83-108.
- Bresnahan, Timothy F., and Jonathan D. Levin. "Vertical Integration and Market Structure." In *Handbook of Organizational Economics*, edited by Robert Gibbons and John Roberts, 853-890. Princeton: Princeton University Press, 2012.
- Competition and Markets Authority. 2020. 'Online Platforms and Digital Advertising Market Study: Final Report'. Competition and Markets Authority. <https://www.gov.uk/cma-cases/online-platforms-and-digital-advertising-market-study>.
- . 2023. 'AI Foundation Models: Initial Report'. Competition and Markets Authority. <https://www.gov.uk/government/publications/ai-foundation-models-initial-report>.
- Corden, Jez. 2023. 'Despite Its Big OpenAI Push, Microsoft's Bing Search Market Share Decreases Year-over-Year'. *Windows Central*, 13 November 2023. <https://www.windowscentral.com/microsoft/despite-its-big-openai-push-microsofts-bing-search-market-share-decreases-year-over-year>.
- Cottier, Ben. 2023. 'Trends in the Dollar Training Cost of Machine Learning Systems'. *Epoch* (blog), 31 January 2023. <https://epochai.org/blog/trends-in-the-dollar-training-cost-of-machine-learning-systems>.
- Doctorow, Cory. 'Social Quitting'. Medium (blog), 17 November 2022. <https://doctorow.medium.com/social-quitting-1ce85b67b456>.
- Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, et al. 2021. 'An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale'. arXiv. <https://doi.org/10.48550/arXiv.2010.11929>.
- Eloundou, Tyna, Sam Manning, Pamela Mishkin, and Daniel Rock. 2023. "GPTs are GPTs: An Early Look at the Labor Market Impact Potential of Large Language Models." arXiv:2303.10130.
- Epoch. 2023. 'Data Trends'. Epoch. 11 April 2023. <https://epochai.org/trends>.



- 1  
2  
3 Epoch, 2024. 'Parameter, Compute and Data Trends in Machine Learning'. Retrieved 15 Feb  
4 2024. <https://epochai.org/data/epochdb/visualization>  
5  
6 FTC. 2023. 'FTC and DOJ Charge Amazon with Violating Children's Privacy Law.' Federal Trade  
7 Commission. 31 May 2023.  
8 <https://www.ftc.gov/news-events/news/press-releases/2023/05/ftc-doj-charge-amazon-violating-childrens-privacy-law-keeping-kids-alexa-voice-recordings-forever>  
9  
10 FTC. 2024. 'FTC Launches Inquiry into Generative AI Investments and Partnerships'. Federal  
11 Trade Commission. 24 January 2024.  
12 <https://www.ftc.gov/news-events/news/press-releases/2024/01/ftc-launches-inquiry-generative-ai-investments-partnerships>.  
13  
14 Fernandez, Joaquin, Knud Lasse Lueth, and Philipp Wegner. 2023. 'Generative AI Market Report  
15 2023–2030'. IoT Analytics. 14 December 2023.  
16 <https://iot-analytics.com/product/generative-ai-market-report-2023-2030/>.  
17  
18 Fletcher, Richard. 2024. How many news websites block AI crawlers? Blog post, Reuters Institute  
19 for the Study of Journalism, Oxford University.  
20  
21 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,  
22 Radu Soricut, et al. 2023. 'Gemini: A Family of Highly Capable Multimodal Models'. arXiv.  
23 <https://doi.org/10.48550/arXiv.2312.11805>.  
24  
25 Hagey, Keach, and Asa Fitch. 'Sam Altman Seeks Trillions of Dollars to Reshape Business of Chips  
26 and AI'. Wall Street Journal, 8 February 2024, sec. Tech.  
27 <https://www.wsj.com/tech/ai/sam-altman-seeks-trillions-of-dollars-to-reshape-business-of-chips-and-ai-89ab3db0>.  
28  
29 Hatzius, Jan, Joseph Briggs, Devesh Kodnani, and Giovanni Pierdomenico. 2023. 'The Potentially  
30 Large Effects of Artificial Intelligence on Economic Growth'. Goldman Sachs Economic  
31 Research.  
32  
33 Heath, Alex. 2024. 'Mark Zuckerberg's New Goal Is Creating Artificial General Intelligence'. *The  
34 Verge*, 18 January 2024.  
35 <https://www.theverge.com/2024/1/18/24042354/mark-zuckerberg-meta-agi-reorg-interview>.  
36  
37 Heim, Lennart. 2021. 'Transformative AI and Compute - EA Forum'. 23 September 2021.  
38 <https://forum.effectivealtruism.org/s/4yLbeJ33fYrwnfDev>.  
39  
40 ———. 2023. 'This Can't Go on(?) - AI Training Compute Costs'. *Blog.Heim.Xyz* (blog). 1 June 2023.  
41 <https://blog.heim.xyz/this-cant-go-on-compute-training-costs/>.  
42  
43 Hemphill, C. Scott, and Tim Wu. 'Nascent competitors.' *University of Pennsylvania Law Review*  
44 168, no. 7 (June 2020): 1879-1910.  
45  
46 Henderson, J. Stephen, and Robert E. Burns. 1989. 'An Economic and Legal Analysis of Undue  
47 Price Discrimination'. The National Regulatory Research Institute.  
48  
49 Hendrycks, Dan, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
50 Steinhardt. 'Measuring Massive Multitask Language Understanding'. arXiv, September  
51 2020. <https://doi.org/10.48550/arXiv.2009.03300>.  
52  
53 Hinton, Geoffrey et al. 2023. "Statement on AI Risk." Center for AI Safety (CAIS), accessed  
54 February 15, 2024, <https://www.safe.ai/statement-on-ai-risk>.  
55  
56 Hjelm, Max. n.d. 'CoreWeave Becomes NVIDIA's First Elite Cloud Services Provider for Compute  
57 — CoreWeave'. Accessed 30 June 2023.  
58 <https://www.coreweave.com/blog/coreweave-becomes-nvidias-first-elite-cloud-services-provider-for-compute>.  
59  
60 Hoffmann, Jordan; Borgeaud, Sebastian; Mensch, Arthur; Buchatskaya, Elena; Cai, Trevor;  
Rutherford, Eliza; Casas, Diego de Las; Hendricks, Lisa Anne; Welbl, Johannes; Clark,

- 1  
2  
3 Aidan; Hennigan, Tom; Noland, Eric; Millican, Katie; Driessche, George van den; Damoc,  
4 Bogdan. 2022. "Training Compute-Optimal Large Language Models". arXiv:2203.15556.  
5 Holmes, Aaron. 2023. 'How Microsoft Swallowed Its Pride to Make a Massive Bet on OpenAI'. The  
6 Information. 25 May 2023.  
7 <https://www.theinformation.com/articles/how-microsoft-swallowed-its-pride-to-make-a-mas>  
8 [sive-bet-on-openai](https://www.theinformation.com/articles/how-microsoft-swallowed-its-pride-to-make-a-mas).  
9  
10 Holmes, Aaron, and Anissa Gardizy. 2023. 'AI Developers Stymied by Server Shortage at AWS,  
11 Microsoft, Google'. *The Information*, 7 April 2023.  
12 <https://www.theinformation.com/articles/ai-developers-stymied-by-server-shortage-at-aws->  
13 [microsoft-google](https://www.theinformation.com/articles/ai-developers-stymied-by-server-shortage-at-aws-).  
14 Hui, Xiang, Oren Reshef, and Luofeng Zhou. 'The Short-Term Effects of Generative Artificial  
15 Intelligence on Employment: Evidence from an Online Labor Market'. SSRN Scholarly  
16 Paper. Rochester, NY, 31 July 2023. <https://doi.org/10.2139/ssrn.4527336>.  
17 Hwang, Tim. 2018. 'Computational Power and the Social Impact of Artificial Intelligence'. SSRN  
18 Working Paper. Rochester, NY. <https://doi.org/10.2139/ssrn.3147971>.  
19 JP Morgan Research. 'Is Generative AI a Game Changer?'. 14 February 2024.  
20 <https://www.jpmorgan.com/insights/global-research/artificial-intelligence/generative-ai>.  
21 Kaplan, Jared, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
22 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. 'Scaling Laws for Neural  
23 Language Models'. arXiv. <https://doi.org/10.48550/arXiv.2001.08361>.  
24 Kerin, Roger A., P. Rajan Varadarajan, and Robert A. Peterson. 1992. 'First-Mover Advantage: A  
25 Synthesis, Conceptual Framework, and Research Propositions'. *Journal of Marketing* 56  
26 (4): 33–52. <https://doi.org/10.1177/002224299205600404>.  
27 Khan, Lina M. 2017. "Amazon's Antitrust Paradox." *The Yale Law Journal* 126 (3): 710–805.  
28 Khan, Saif M., Dahlia Peterson, and Alexander Mann. 2021. 'The Semiconductor Supply Chain'.  
29 Center for Security and Emerging Technology.  
30 <https://cset.georgetown.edu/publication/the-semiconductor-supply-chain/>.  
31 Korinek, Anton. 2023. Scenario Planning for an A(G)I Future, *IMF Finance & Development*  
32 *Magazine* 60(4), pp. 30-33.  
33 Knight, Will.. 2024. 'Amazon's Cloud Boss Likens Generative AI Hype to the Dotcom Bubble'.  
34 *Wired*, 7 February 2024.  
35 <https://www.wired.com/story/amazons-cloud-boss-selipsky-generative-ai-hype/>.  
36 Li, Wendy Y. 2023. 'Regulatory Capture's Third Face of Power'. *Socio-Economic Review* 21 (2):  
37 1217–45. <https://doi.org/10.1093/ser/mwad002>.  
38 Lieberman, Marvin B., and David B. Montgomery. 1988. 'First-Mover Advantages'. *Strategic*  
39 *Management Journal* 9 (S1): 41–58. <https://doi.org/10.1002/smj.4250090706>.  
40 Lomas, Natasha. 'Big Tech AI Infrastructure Tie-Ups Set for Deeper Scrutiny, Says EU Antitrust  
41 Chief'. TechCrunch, 20 February 2024.  
42 <https://techcrunch.com/2024/02/20/eu-merger-control-ai/>.  
43 ———. 'EU Checking If Microsoft's OpenAI Investment Falls under Merger Rules'. TechCrunch, 9  
44 January 2024. <https://techcrunch.com/2024/01/09/openai-microsoft-eu-merger-rules/>.  
45 Martens, Bertin. 'What Should Be Done about Google's Quasi-Monopoly in Search? Mandatory  
46 Data Sharing versus AI-Driven Technological Competition', 6 July 2023.  
47 <https://www.bruegel.org/working-paper/what-should-be-done-about-googles-quasi-mono>  
48 [poly-search-mandatory-data-sharing-versus](https://www.bruegel.org/working-paper/what-should-be-done-about-googles-quasi-mono).  
49 Mattioli, Dana. 2020. 'Amazon Scooped Up Data From Its Own Sellers to Launch Competing  
50 Products'. *Wall Street Journal*, 23 April 2020, sec. Tech.  
51 <https://www.wsj.com/articles/amazon-scooped-up-data-from-its-own-sellers-to-launch-com>  
52  
53  
54  
55  
56  
57  
58  
59  
60

- peting-products-11587650015.
- Norem, Josh. 'Analysts Estimate Nvidia Owns 98% of the Data Center GPU Market'. *ExtremeTech* (blog), 1 February 2024.  
<https://www.extremetech.com/computing/analysts-estimate-nvidia-owns-98-of-the-data-center-gpu-market>.
- Olson, Parmy. 2024a. 'Google, Microsoft Will Dominate AI as Computing Costs Surge'. *Bloomberg.Com*, 19 February 2024.  
<https://www.bloomberg.com/opinion/articles/2024-02-19/artificial-intelligence-microsoft-google-nvidia-win-as-computing-costs-surge>.
- . 2024b. 'Zuckerberg's Secret Weapon for AI Is Your Facebook Data'. *Bloomberg.Com*, 6 February 2024.  
<https://www.bloomberg.com/opinion/articles/2024-02-06/zuckerberg-s-plan-for-ai-hinges-on-your-facebook-and-instagram-data>.
- OpenAI. 2023a. 'Our Structure'. 2023. <https://openai.com/our-structure>.
- . 2023b. 'ChatGPT Plugins'. OpenAI. 23 March 2023.  
<https://openai.com/blog/chatgpt-plugins>.
- Patel, Dylan, and Daniel Nishaball. 'Inference Race To The Bottom - Make It Up On Volume?' *SemiAnalysis* (blog), 19 December 2023.  
<https://www.semianalysis.com/p/inference-race-to-the-bottom-make>.
- Reid, Elizabeth. 2023. 'Supercharging Search with Generative AI'. *Google* (blog). 10 May 2023.  
<https://blog.google/products/search/generative-ai-search/>.
- Salesforce. 2023. 'Top Generative AI Statistics for 2024'. Salesforce. 1 September 2023.  
<https://www.salesforce.com/news/stories/generative-ai-statistics/>.
- Shani, Inbal. 2023. 'Survey Reveals AI's Impact on the Developer Experience'. *The GitHub Blog* (blog). 13 June 2023.  
<https://github.blog/2023-06-13-survey-reveals-ais-impact-on-the-developer-experience/>.
- Shu, Catherine. 2014. 'Google Acquires Artificial Intelligence Startup DeepMind For More Than \$500M'. *TechCrunch*, 27 January 2014.  
<https://techcrunch.com/2014/01/26/google-deepmind/>.
- Staff in the Bureau of Competition & Office of Technology. 2023. 'Generative AI Raises Competition Concerns'. *Federal Trade Commission* (blog). 29 June 2023.  
<https://www.ftc.gov/policy/advocacy-research/tech-at-ftc/2023/06/generative-ai-raises-competition-concerns>.
- Suleyman, Mustafa. *The Coming Wave: Technology, Power, and the Twenty-First Century's Greatest Dilemma*. First edition. New York: Crown, 2023.
- The Economist*. 2023. 'Large, Creative AI Models Will Transform Lives and Labour Markets', 22 April 2023.  
<https://www.economist.com/interactive/science-and-technology/2023/04/22/large-creative-ai-models-will-transform-how-we-live-and-work>.
- Tiku, Nitasha, and Gerrit De Vynck. 2023. 'Google Shared AI Knowledge with the World — until ChatGPT Caught Up'. *The Washington Post*, 4 May 2023.  
<https://www.washingtonpost.com/technology/2023/05/04/google-ai-stop-sharing-research/>.
- Tirole, Jean. 1988. *The Theory of Industrial Organization*. Cambridge, MA: MIT Press.
- Trammell, Philip, and Anton Korinek. 'Economic Growth under Transformative AI'. Working Paper. Working Paper Series. National Bureau of Economic Research, October 2023.  
<https://doi.org/10.3386/w31815>.
- Varadarajan, Rajan, Manjit S. Yadav, and Venkatesh Shankar. 2008. 'First-Mover Advantage in an

- 1  
2  
3 Internet-Enabled Market Environment: Conceptual Framework and Propositions'. *Journal*  
4 *of the Academy of Marketing Science* 36 (3): 293–308.  
5 <https://doi.org/10.1007/s11747-007-0080-y>.  
6  
7 Verhulst, Stefaan G. 'Are We Entering a Data Winter? On the Urgent Need to Preserve Data  
8 Access for the Public Interest'. *Frontiers Policy Labs* (blog), 2024.  
9 <https://policylabs.frontiersin.org/content/commentary-are-we-entering-a-data-winter>.  
10  
11 Victor, Jon. 2023. 'Why YouTube Could Give Google an Edge in AI'. *The Information*, 14 June  
12 2023.  
13 <https://www.theinformation.com/articles/why-youtube-could-give-google-an-edge-in-ai>.  
14  
15 Villalobos, Pablo, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and Anson  
16 Ho. 2022. 'Will We Run out of Data? An Analysis of the Limits of Scaling Datasets in  
17 Machine Learning'. arXiv. <https://doi.org/10.48550/arXiv.2211.04325>.  
18  
19 Vipra, Jai, and Sarah Myers West. 'Computational Power and AI'. AI Now Institute, 27 September  
20 2023. <https://ainowinstitute.org/publication/policy/compute-and-ai>.  
21  
22 Warren, Tom. 2023. 'Microsoft Extends OpenAI Partnership in a "Multibillion Dollar Investment"'.  
23 *The Verge*, 23 January 2023.  
24 <https://www.theverge.com/2023/1/23/23567448/microsoft-openai-partnership-extension-a>  
25 i.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60