

Old Moats for New Models: Openness, Control, and Competition in Generative AI*

Pierre Azoulay
MIT Sloan and NBER

pazoulay@mit.edu

Joshua Krieger
Harvard Business School

jkrieger@hbs.edu

Abhishek Nagaraj
UC Berkeley-Haas and NBER

nagaraj@berkeley.edu

Abstract

Drawing insights from the field of innovation economics, we discuss the likely competitive environment shaping generative AI advances. Central to our analysis are the concepts of *appropriability*—whether firms in the industry are able to control the knowledge generated by their innovations—and *complementary assets*—whether effective entry requires access to specialized infrastructure and capabilities to which incumbent firms can ration access. While the rapid improvements in AI foundation models promise transformative impacts across broad sectors of the economy, we argue that tight control over complementary assets will likely result in a concentrated market structure, as in past episodes of technological upheaval. We suggest the likely paths through which incumbent firms may restrict entry, confining newcomers to subordinate roles and stifling broad sectoral innovation. We conclude with speculations regarding how this oligopolistic future might be averted. Policy interventions aimed at fractionalizing or facilitating shared access to complementary assets might help preserve competition and incentives for extending the generative AI frontier. Ironically, the best hopes for a vibrant open source AI ecosystem might rest on the presence of a “rogue” technology giant, who might choose openness and engagement with smaller firms as a strategic weapon wielded against other incumbents.

* We thank Josh Lerner, Janet Freilich, Shane Greenstein, Joshua Gans, Shikhar Ghosh, Pam Mishkin, Roland Szabo, Nikhil Naik, Rishi Bommasani, Juan Mateos-Garcia, Yoon Kim, and Roger Levy for useful discussions. Nilo Mitra and Yanqi Cheng provided able research assistance. We acknowledge the help of GPT-4 and Claude 3 Sonnet as unparalleled brainstorming partners, expert paraphrasers, and whisperers of imaginary references.

1 Introduction

The history of science and technology is a steady march punctuated by great inflection points (Kuhn 1962). In such moments, the accumulation of knowledge about how humans can manipulate our world coalesces into novel technologies with the potential for great jumps in productivity. Like steam engines, airplanes, antibiotics, semiconductors, and the internet in previous eras, generative artificial intelligence (AI) has the potential to be the next major transformative technology.

Interest in generative AI skyrocketed following the release of Open AI’s large language models (LLMs) for text (ChatGPT) and image generation (DALL-E) in the second half of 2022. ChatGPT’s uncanny ability to mimic human responses prompted individuals and firms to rethink what tasks might be aided or replaced by generative AI. Prominent business examples include novel protein design, writing software code, and weather forecasting. Meanwhile, individual consumers have embraced applications ranging from the mundane (automating email responses and shopping lists), to artistic (writing poetry, bedtime stories and screenplays), and skill building (e.g., personal tutors, coding copilots). As of early 2024, the GPT-4 model could pass the bar exam and a variety of Advanced Placement exams; score above the 93rd percentile on the SATs and above the 99th percentile on the GRE verbal exam; as well as perform a whole host of coding and programming challenges. In sum, there is little doubt that many industries—including law, information technology, education, and entertainment—are poised to experience large shocks to their labor and product markets.

This new wave of “foundation” generative AI models (i.e. large-scale, pre-trained AI models built for diverse applications) possesses capabilities that are discontinuous improvements over previous generations of machine learning technologies. However, we do not yet understand whether or to what extent one or a few firm(s) will dominate the supply of foundational generative AI technologies. While regulators, investors, entrepreneurs, and managers face considerable uncertainty when attempting to forecast the broad societal impacts of AI, they must nonetheless engage in sober assessments of innovation and competition in this fledgling sector to shape it. Understanding the industrial organization of the industry that will power these applications strikes us as a legitimate and necessary topic for entrepreneurship and innovation policy.

Rather than prognosticate about the depth and breadth of generative AI’s applications, this paper harnesses the accumulated wisdom of research in the field of innovation economics to venture (with some trepidation and many caveats) a prediction for the likely *competitive environment* in which generative AI advances will take place. How might the technological features of generative AI shape the extent of product market competition in the industry? How likely are deviations from perfect competition to result in blockaded entry and delayed innovation? And what will determine the viability of a vibrant “open source” ecosystem, as opposed to a more secretive, incumbent-driven oligopoly? To guide our analysis, we draw on Teece (1986), who gives primacy to the concepts of *appropriability* and *complementary assets* to identify the firms best able to profit from innovation: those first to market, follower firms, or those with formal rights to key assets (such as patents) or with related capabilities and assets that innovating firms require access to. *Our central claim is that while formal intellectual property and secrecy are unlikely to durably prevent innovative firm entry, incumbent firms’ tight control over key complementary assets will likely usher in a highly concentrated market structure.* Thus, incumbents will have the ability to

confine new entrants to the “application layer” of the industry in a pattern reminiscent of the smartphone industry, thus limiting widespread competition at the foundation layer of the generative AI stack.

Studies of previous general purpose technologies such as the steam engine, factory electrification, machine tools, and computer programming provide useful lessons for the competitive paths likely to emerge in generative AI (Bresnahan and Trajtenberg 1995; David 1990; Rosenberg and Trajtenberg 2004). As new and exciting as the early applications of AI may have been, they have mainly affected sectors that were either natively digital or had already experienced sustained digital transformation, such as internet search, digital advertising, customer relationship management, and protein structure prediction. For generative AI to reach its full technological potential, however, it will have to overcome bottlenecks that prevent invention in a broad set of application sectors, from agriculture to construction and financial services. These bottlenecks may be less technical in nature and more rooted in the incentive problems and adjustment costs that organizations face as they digitize operations (Iansiti and Lakhani 2020; Bresnahan 2024). An additional obstacle, and one central to the themes explored in this article, is that concentrated control in the upstream layer of the generative AI “technology stack” (defined precisely below) might well dim the incentives to innovate in the downstream application sectors. As a result, a flowering of AI innovation throughout the economy faces substantial headwinds.

Industry predictions about the evolution of generative AI competition are quite varied but ultimately hinge on whether early leaders can establish a competitive “moat.” The flurry of new model releases has been matched only by claims and counter-claims of model superiority. Evaluating these statements is challenging, especially given the gold rush of venture capital money flowing into the industry.¹ For example, in May 2023, a leaked internal Google memo titled “*We have no moat, and neither does OpenAI*” argued that the large and well-resourced leaders who pioneered the field using proprietary approaches would lose out to the nimbler open-source developers. Core to the argument was the idea that methods for improving and “fine tuning” open models were quickly dropping in cost.²

On the other hand, early adoption patterns suggested that early movers and well-resourced incumbents were skilled at turning new technologies into compelling products. Consumer adoption reflected the dominance of OpenAI’s ChatGPT service, which was reported to have more than 100 million monthly active users by January 2023 and 100 million weekly active users by November 2023, along with more than 2 million developers using its API.³

One software startup CEO noted that having a “secret sauce” is rarely how defensible software businesses are made, a claim backed by ample historical experience.⁴ Today’s dominant technology platforms, including Google’s search engine, Meta’s social graph, and Amazon’s e-commerce platforms were all built on technologies that were well-understood by competitors and potential entrants. Rather, these firms locked in customers through superior execution (e.g., in the form of

¹ Even after excluding the \$14 billion in funding that went to LLM leaders Open AI and Anthropic, Pitchbook data showed that the first three quarters of 2023 saw more than \$7.4 billion in venture capital deals supporting generative AI startups (<https://pitchbook.com/news/articles/generative-ai-startups-vc-deals-decline>).

² <https://www.semianalysis.com/p/google-we-have-no-moat-and-neither>

³ <https://techcrunch.com/2023/11/06/openai-chatgpt-now-has-100-million-weekly-active-users/>

⁴ <https://www.airplane.dev/blog/openai-moat-is-stronger-than-you-think>

compelling user interfaces), continuous innovation, and the harnessing of network effects. Further back, in industries from typesetting and aerospace to personal computing and disk drives, a small number of firms have often dominated in fast-growing industries despite astonishing performance improvements, cost reductions, and widespread technology diffusion. We draw on the insights gained in studying such breakthrough technologies to describe the forces likely to shape industry structure in generative AI and how innovation incentives might play out under different scenarios.

We argue that tight control over specialized complementary assets is the most likely source of “moats” that would enable pioneering firms to durably entrench their dominant positions. We describe how producing and improving foundational generative AI models requires several interrelated complementary assets, six of which are highlighted in this essay: (1) the compute environment, (2) model-serving and inference capabilities, (3) safety and governance procedures, (4) the development of benchmarks and metrics, (5) access to massive quantities of non-public training data, and (6) data network effects—whereby users’ engagement with a model generates information that dynamically improves its performance. Given the potential for early success to beget further advantage along one or more of these dimensions, the prospects for a competitive market structure rest on the belief that public policy efforts will succeed in fractionalizing these assets or facilitating their broad sharing across a diverse set of ecosystem participants. Alas, determined interventions along these lines are likely to face strong political opposition.

Ironically, a market structure characterized by oligopolistic control can perhaps best be averted through the actions of a single “rogue” technology giant. This rogue firm might pursue openness and engagement with smaller firms as a self-serving approach to improve its offerings in adjacent markets, or win the “war” for scarce AI talent. More generally, an open-access strategy might enable this firm to burnish its credentials as an AI technology leader, influencing standards and practices across the industry. Ordinarily, a maverick firm’s benevolence would be a thin reed on which to pin hopes for a democratized innovation ecosystem. However, that possibility deserves serious consideration in light of Meta, Inc.’s publicly professed commitment to open source—although doubts linger regarding the sustainability of this approach amid present and future competitive pressures.

With or without a rogue firm, the locus of technological exploration is likely to be confined to the application layer, whereby startup entrants build services in the form of “fine-tuned” models running on top of the foundation models controlled by leading tech firms. A proper role for innovation and competition policy might be to ensure that the governance of this application layer guarantees academics and entrepreneurs enough freedom and pricing autonomy, lest these app stores of the future end up resembling gated communities rather than community gardens.

Our paper begins with a technical primer on generative AI, describing the key technologies behind current-day generative AI models. Next, we lay out a framework for assessing the prospects for the emergence of a vibrant competitive landscape in generative AI, concluding that these prospects are rather dim for the foundation layer of the technological stack, but somewhat brighter for the application layer. We end by proposing policy “low-hanging fruits” that have the potential to help usher in more competition in this domain.

2 Technical Primer

2.1 What do we mean by Generative AI?

During the 2010s, progress in the field of machine learning resulted in the creation of models that excelled at predicting outcomes from input data. Deep learning and neural networks transformed the field during this period, enabling models to achieve high accuracy on tasks like financial fraud detection, image classification, and speech recognition.

The recent revolution in generative AI builds on past achievements in machine learning using so-called “foundation” models largely built using the newly developed transformer architecture. Foundation models are large-scale AI models that serve as the basis for a wide range of downstream applications. These models generate new content by learning from the patterns and structures present in the training data. During the training process, the model ingests vast amounts of text, images, or other forms of data, and breaks that data into “tokens” (e.g., word, pixel). In training, the model “learns” to predict the next token based on the input sequence. By iteratively refining its predictions, it can generate outputs that closely resemble the training data. Once trained, foundation models can be fine-tuned for specific tasks or used for open-ended generation of new content.

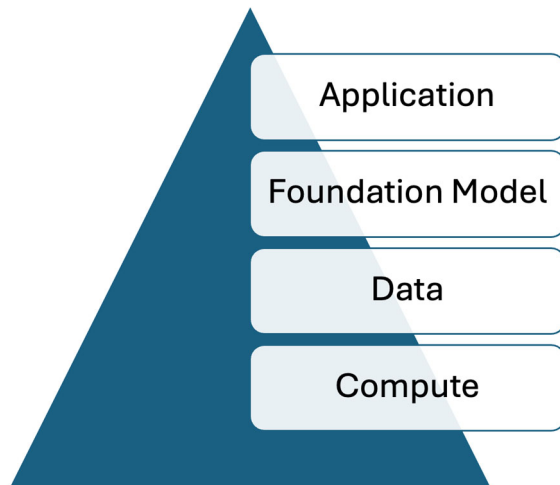
What makes foundation models different than the previous generation of machine learning models is that they are “pre-trained” (i.e., they do not need to be trained using a specific dataset) and they are generic and context-independent (i.e., they are often powerful enough to perform a vast array of tasks without much additional training). For example, a language model like GPT-4 is extremely flexible and can perform varied tasks like language translation, text summarization, or question answering, with very little additional instruction. Alternatively, it can be used for creative writing tasks, such as generating stories, articles, or even code. Similarly, image-based foundation models like DALL-E can generate novel images based on textual descriptions, while audio-based models like AudioLM can generate music and speech. Foundation models are now being developed in a variety of other fields, such as robotics and biology. It is the versatility and adaptability of foundation models that have opened new possibilities for generative AI across various domains.

Recently, companies like OpenAI have exposed the functionality of foundation models like GPT-4 to third parties so that they can build “plug in” modules to enhance the attributes of this leading foundation model. Therefore, a technology “stack” that might enable greater experimentation in generative AI applications is beginning to emerge (see Figure 1). At the bottom of the stack is the compute layer, which includes graphics processing units (GPUs) and a supporting technical infrastructure with many interrelated hardware, software, and data communication components. The second layer is composed of large amounts of data and the associated cloud storage infrastructure required to host these massive datasets. The third layer is the foundation model itself, which provides a general-purpose interface on top of which specific applications can be run. These models can either seem like powerful, almost magical “black boxes” or present themselves as transparent and customizable systems where developers have access to the model’s weights (the parameters learned during training that define how the model makes predictions) or can otherwise inspect it “under the hood.” The topmost layer is the application level, which facilitates interactions that allow end users to

access the functionality of the foundation model. It includes specialized applications enabling users to elicit answers from the underlying model in a specific domain or format.

For illustration purposes, GPT-4 is a foundation model, while ChatGPT is a specific application allowing users to query GPT-4 with a chatbot interface.⁵ This design is reminiscent of the smartphone application stores where third-party developers leverage the functional features of devices and operating systems to develop targeted applications. Given their central position in the technology stack, foundation models can be a chokepoint for innovation and competition at the application level, and are therefore at the heart of generative AI policy debates. The central concern of this article is competition and control *within* the foundation model layer, but we will also attend to the effects that concentration at the foundation model layer could have for entry at the application layer, with the emergence of a vibrant AI application ecosystem hanging in the balance.

Figure 1: The Generative AI Technology Stack



2.2 Large Language Models (LLMs)

Perhaps the most well-known of all foundation models is the Generative Pre-Trained Transformer (GPT) class of models from OpenAI, which belongs to a class of models called Large Language Models (LLMs) that are designed specifically to work with human language (and recently images and voice). Large Language Models (LLMs) are a subset of generative AI that analyze vast amounts of data to understand context, recognize patterns, and generate new content. By understanding how LLMs work, how they are trained, and the applications they enable, we can better grasp the potential and implications of generative AI for innovation and competition.

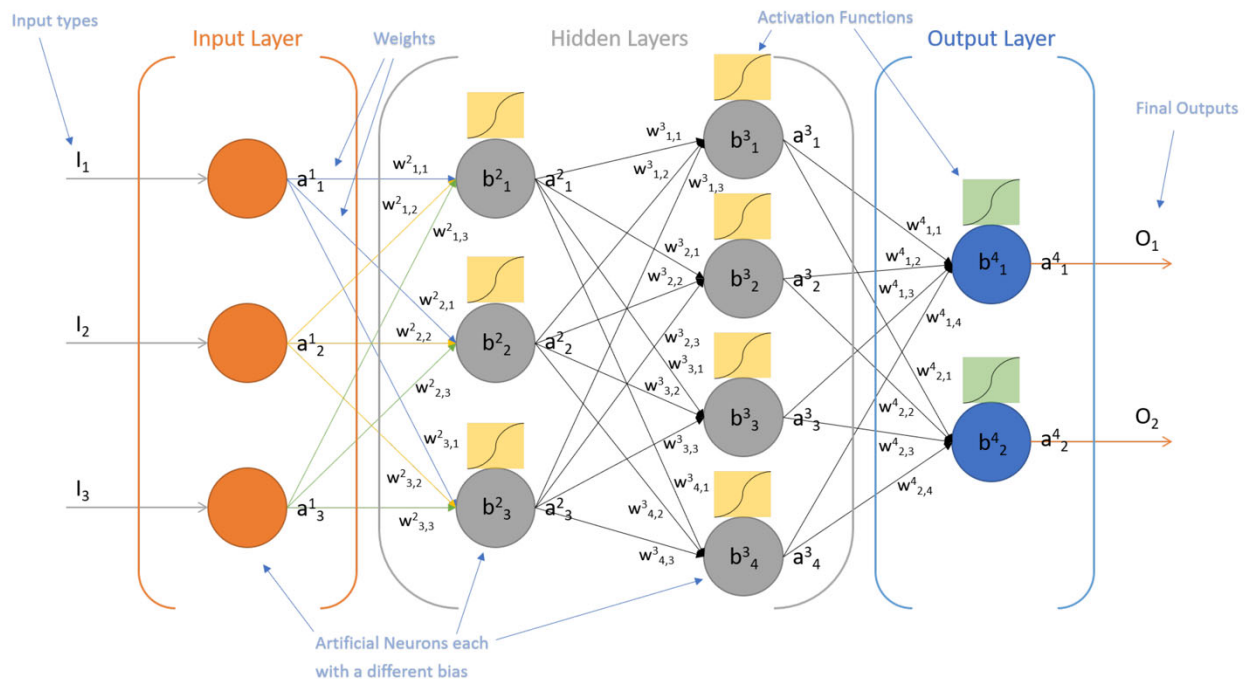
The term “language” in LLMs refers to various forms of symbolic representation, including human languages, computer languages, and any other type of encoded meaning. At their core, LLMs are

⁵ At the time of writing this manuscript, most models and applications have identical owners, so that the distinction between the two layers appears mostly semantic. But this soon could change, which is why we choose make the separation between model and application salient.

algorithms that predict the next word or token when fed a string of words, progressively building up syntax that a human could plausibly create. LLMs model the ingested training data as a multi-dimensional space with numerous parameters, identifying artifacts (e.g., words, pixels) and their relationships to build new structures. Critical to the recent leap forward in generative AI is the transformer architecture, which adduces to traditional deep learning approaches an “attention mechanism” to discern contextual relevance—e.g., how the relationship between words (in the training data, prompts, or output) changes the likely meaning of a phrase (e.g., a machine learning model vs. a fashion model).⁶

To understand how a model like GPT-4 works, it is useful to first take a deeper dive into the transformer architecture. The “M” in the LLM acronym stands for “Model,” i.e., the empirical derivation of a mathematical formula that best captures the input text, manipulated as numerical data. At its core, each model runs on a neural network, which consists of nodes arranged in layers interconnected with nodes in the adjacent layer (See Figure 2).

Figure 2: Representation of a Neural Network in LLMs⁷



Each node performs some computation, with the result passed forward to the adjacent interconnected nodes. Models also incorporate activation functions which are mappings applied to the weighted sum of inputs coming into a neuron. It introduces non-linearity into the neural network, allowing it to learn and model complex relationships between inputs and outputs. Each interconnection is given a certain “weight,” a multiplicative factor that captures the importance of the originating node’s result. The

⁶ Grant Sanderson’s video series on deep learning provides an accessible introduction. See in particular Chapter 5, “But what is a GPT? Visual intro to transformers.” Available at <https://www.youtube.com/watch?v=wjZofJX0v4M>.

⁷ From Yazici, Mahmut Taha, Shadi Basurra, and Mohamed Medhat Gaber. 2018. “Edge Machine Learning: Enabling Smart Internet of Things Applications.” *Big Data and Cognitive Computing* 2(3): 26.

numerical values are derived iteratively by “training” the neural network with large amounts of known text samples to minimize a loss function that captures the difference between the model’s output and the expected outcome. These weights and some other adjustable inputs called “biases” are collectively called “parameters” and have become a proxy for representing the complexity of the neural network and the number of factors used to predict an output. Industry observers often compare generative AI models by the number of parameters. GPT-3, for example, has 175 billion parameters, while GPT-4 is estimated to have over a trillion. Meta’s latest Llama 3.1 model has a version with over 405 billion parameters.

Earlier generations of neural networks used for natural language processing suffered from several drawbacks, including their inability to process whole sentences in one run and “long-term memory loss”—the inability to “remember” words occurring much earlier in the input sequence. These issues were overcome using the transformer architecture, created by Google researchers in 2017 (Vaswani et al. 2017). Transformer architectures process input text—broken down into units called tokens, each representing a word, word fragment, strings of words, or a sentence—in parallel and can capture the context between words in sentences using a mechanism called “attention.” When predicting the next word, the attention mechanism learns from the training data to “attend” or “focus” on specific words or phrases that are particularly relevant in the input context.

The training data used by prominent LLMs typically have trillions of sentences or images, taken from publicly available sources such as the Web, books, social media content, audio, and video transcripts to recognize rules and patterns in the underlying data and make predictions for generating new content that uncannily resembles the training data. Some LLMs pre-train on specific categories, such as images (e.g., Open AI’s DALL-E) or music (e.g., Google’s AudioLM) to create new content specifically for these areas. Developers can also “fine-tune” models, a process that allows additional data to inform the output generated by a pre-trained model. While the most well-known LLMs distinguish themselves around content format (e.g., text, images, video, audio), other transformer-based generative AI models target specific knowledge domains and physical phenomena. One such specialized model is AlphaFold, an AI system pre-trained on data from the Protein Data Bank (PDB), a publicly-funded effort started in the early 1970s to catalog proteins’ three-dimensional structures. DeepMind trained their LLM on the PDB in order to predict a protein’s 3D structure from its amino acid sequence. In 2018, AlphaFold made a major leap in accuracy, outperforming alternative protein structure predictions by a wide margin (Bertoline et al. 2023).

In addition to the transformer architecture, diffusion models (which are a type of convolutional neural network) have gained popularity. These models are based on reversing a diffusion process—they start with data (like an image or audio) and gradually add noise to it, “diffusing” the original data over many steps. They then learn to reverse this diffusion process to reconstruct the original data from the pure noise input. Popular diffusion models include OpenAI’s DALL-E 2, Stability AI’s Stable Diffusion, and RFdiffusion for *de novo* protein design (Watson et al. 2023).⁸ Diffusion models have proved

⁸ For protein structure determination, the transformer-based model AlphaFold 2 (Jumper et al., 2020) provided the initial breakthrough. The training data for this model consisted of known protein structures, but the model also made use of multiple sequence alignments to gather evolutionary information about a protein family. Encoding evolutionary rules into the architecture endowed these models with the ability to propose high quality structures for a vast number of proteins that had proved hard to resolve using traditional experimental methods, such as membrane proteins. At the same time, it

especially attractive to extend the power of LLMs to media other than text. However, the importance of training data and computing power remains central in this architecture. Therefore, diffusion models and “chatbot”-style LLMs appear quite complementary in both their development and use, with many of the leading-edge generative AI developers offering both within the same user interface (e.g., “multi-modal” options within ChatGPT, Meta AI, or Google Gemini).

Apart from OpenAI’s GPT, prominent examples of foundation models include Google’s Gemini, Large Language Model Meta AI (Llama) by Meta, and Claude from Anthropic. This diverse set of foundation models serves as the basis for various applications, such as chatbots, language translation, and code generation.

2.3 Alternative foundation model architectures

LLMs have clearly brought a step-increase to the immediate practical value of generative AI technology. The early success of transformer-based foundation models strengthens the claim that they represent a “dominant design,” meaning that they will become a widely adopted paradigm that simultaneously enables and constrains progress in this domain (Utterback 1994). If this prognosis is correct, then recent advances are only the beginning of a sequence during which incremental improvements in transformer-based models will continue for quite some time until their very design becomes a straitjacket, causing progress to markedly slow down as the current technological trajectory reaches a plateau (Dosi 1982).

A counter-argument might be that the field has not yet experienced the emergence of a dominant design, and that the enthusiasm for transformer-based models might recede if more promising alternatives emerge, or if different technology “architectures” prove more valuable for certain generative AI tasks. Alternatives or complements to the transformer-based LLM architecture include convolutional neural networks (CNNs) used in many image generation models, and liquid neural networks (LNNs), which offer potential advantages such as greater interpretability, adaptability to new information, and lower computing requirements. LNNs are particularly attractive architectures for robotics and autonomous vehicle applications (Hasani et al. 2021).⁹

Whether these alternatives surprise industry or remain an academic curiosity depends on both the uncertainty of scientific discovery and generative AI researchers’ allocation of resources and attention across emerging technologies. While predicting the take-off of these alternatives (and their eventual commercial appeal) is beyond our remit and best left to computer scientists, the rest of the article espouses the conventional view that foundation AI models currently situate themselves in the growth phase of an *S*-curve—a pattern followed by many new technologies, from tire materials to personal computing, artificial hearts, or X-Ray Lithography (Foster 1986; Henderson 1995).

limited their usefulness in exploring the vast space of possible protein structures that have not been observed in nature. New models such as RFdiffusion have sought to free computational biology from this constraint by employing a model architecture akin to that found in image generation models, allowing the generated structures to include novel folding patterns or functionalities, a valuable feature in the context of drug discovery research.

⁹ See <https://techcrunch.com/2023/08/17/what-is-a-liquid-neural-network-really/> for a discussion of advantages and disadvantages of LNNs, compared to other deep learning methods.

2.4 The “openness” fetish

As their name implies, foundation models form the substrate developers use to build a wide range of generative AI applications. Given this critical position in the generative AI ecosystem, foundation models can be a chokepoint for innovation and new entrants, offering incumbents a competitive edge that may be difficult to overcome.

Computer scientists have begun to tackle the issue of foundation model competition by focusing on the concepts of “openness” and “transparency.” The basic idea is that if the various elements of the “secret sauce” behind a model are made available through code, data, or detailed description, then it should be technically feasible to recreate it, build upon it, and possibly compete against it. For example, one definition of openness claims that “open systems” are those that provide “transparency, reusability, and extensibility—they can be scrutinized, reused, and built on” (Widder et al. 2023). The idea is imported from the philosophy of the open source software movement, where openness means that the code behind a particular application is made available for public inspection and scrutiny.

How precisely openness and transparency should be measured and what counts as genuinely “open” remain matters of active debate in the academic computer science community. For example, licensees of Open AI’s GPT-4 have access to its APIs but have limited knowledge of its internals (such as the weights used in the model, a critical component for extensibility). The company claims that it is open because it allows open access to the model’s output, while critics argue that without transparency around how the model was built, it remains relatively closed. They point to other large foundational models—Meta’s Llama 2 and OPT-175B and Google’s BERT—that are identified as open source. But even in these cases, as Widder et al. note, the term “open source” does not yet have a precise meaning in AI, unlike traditional open-source software development, because many aspects of these purportedly open models remain black boxes. For example, some development efforts unaligned with major tech firms (such as Bloom from BigScience and GPT-J by EleutherAI) are considered more “open” as the term is conventionally understood in the software industry, while Meta’s Llama 2 has bespoke licensing terms that do not conform to those used by conventional open-source software projects.

Efforts are underway to standardize the definitions of openness and transparency. One notable effort in this spirit is the Foundation Model Transparency Index (Bommasani et al. 2023) created by Stanford University’s Center for Research on Foundation Models, which aims to formally evaluate transparency and openness claims. These researchers score ten leading foundation models against a battery of one hundred indicators such as compute, data, model weights, etc. The results of the 2023 index show that the ten major foundation models provide limited transparency, with a top score of 54% and a median of 37%, with the so-called “open” foundation models, Llama 2 and BLOOMZ, being the top scorers for transparency. Another such evaluation comes from Liesenfeld et al. (2023), who categorize model openness across thirteen separate dimensions. Similarly, Schrepel and Pentland (2023) have proposed an openness taxonomy to guide policy recommendations that policymakers can use to maintain a competitive ecosystem.

These efforts are valuable in that they reveal the technical factors that might hinder the reverse-engineering of foundation models, including access to training datasets with the required breadth and depth, as well as detailed knowledge of the model’s architecture (e.g., number of parameters, number

of neural network layers, precise description of the attention mechanism, possible use of reinforcement learning with human feedback among many other factors). Open source advocates believe that unless leading firms lay bare their models' architecture, weights, data sources, etc., competition will necessarily be thwarted because even well-resourced rivals would have trouble replicating the pioneering models' performance.

Despite their value in assessing the technical landscape, we disagree with the underlying assumption that openness and transparency will necessarily beget entry and competition. As the long history of technological innovation shows, there is usually only a weak relationship between openness in a technical sense and a regime where leading incumbents are unable to appropriate rents from their innovations. For example, in the pharmaceutical industry, firms protect their key discoveries using a mix of secrecy and patent protection, yet multiple incumbents can develop successful products using varied approaches within a broader "class" of therapies. Conversely, there are cases where generics have entered after a key patent has expired, but leading incumbents still enjoy considerable market power through their brand, distribution, and related assets, as seen with synthetic insulin and the continued dominance of Eli Lilly and Novo Nordisk despite the expiration of formulation patents over two decades ago.

We contend that philosophical discussions about the true meaning of openness constitute a distraction if one's goal is to elucidate the conditions that will enable ongoing competition in the emerging generative AI sector. Rather, as we will argue, the more important factors driving foundation model competition are incumbents' ability to endogenously raise the costs of reverse engineering over time, thus strengthening their grip on key complementary assets required to effectively compete.

3 Economic Moat-ivation: A framework for predicting likely competitive dynamics in foundation AI models

As practitioners of the dismal science, we are perhaps playing to type when we make light of the computer scientists' esthetic preferences for openness and transparency. But this stance stems from a dispassionate assessment of other technology markets' competitive dynamics. Historically, the relationship between openness of key technological ingredients and the extent of product market competition has been tenuous at best. To be sure, the circulation of technical talent across firms often means that access to key technical components eventually becomes democratized through feats of reverse engineering. This is especially the case when the national competitiveness stakes are high, as can currently be seen in the semiconductor industry. And yet, it is equally possible to point to numerous instances where technical know-how is widely dispersed across firms, and yet intense competition does not ensue.

Teece (1986) suggests that the ability of pioneer innovators to durably appropriate the returns from innovation is shaped by two fundamental forces: control over the knowledge at the core of the innovation—also called *appropriability*—and control over the *complementary assets* necessary to transform a firm's innovative know-how into a value proposition customers might be willing to pay for. These assets might be tangible (like R&D tools, distribution networks, or manufacturing facilities), or intangible (like regulatory expertise, specialized skills, industry relationships, or brand equity). Below,

we analyze the emerging generative AI industry through this lens. With a clear understanding of how the twin scissors of appropriability and complementary assets cut across the sector, we can derive implications for competition policy in this fascinating—but ultimately “normal”—domain.

3.1 Appropriability

Two broad approaches can be used by firms to protect the knowledge generated by their inventive and innovative activities: formal intellectual property rights (such as patents, copyrights, and trademarks) and mechanisms to forestall reverse engineering (such as trade secrets, non-disclosure agreements, non-competes, and tacit knowledge).

Although there are plenty of AI-related patents (Giczy et al. 2022), patents do not appear to be especially valuable for protecting foundation models from imitation. For example, the introduction of the transformer architecture was a significant milestone in the development of modern LLMs, but its broad contours were disclosed in a manuscript made available on arXiv, the most widely-used free distribution platform for pre-prints and working papers in computer science (Vaswani et al. 2017). That initial disclosure does not block all patents related to this architecture (e.g., pertaining to specific attention mechanisms), many of which might be quite valuable due to their wide applicability in models like BERT, GPT, and others. However, as is more broadly true in the software realm, much of AI/ML research may fall in the category of “abstract concept” and therefore fall short of the patentability standard.

Further, the rapid advancement of generative AI innovation may not align well with the pace of patent examination in the modern patent system. This may lead inventors to turn to alternative mechanisms in order to control the critical technological knowledge associated with their ideas. It is worth noting that a variety of firms have at least applied for patents pertaining to technologies closely adjacent to foundation models.¹⁰ Yet other firms hold patents on software features in their products which directly embed the output of foundation models.¹¹ With massive uncertainty surrounding the value of IP claims in this area, patenting investments might be thought of as “lottery tickets” that may be worth purchasing even if the expected value of such a patent is likely to be very low (Lemley and Shapiro 2005).

The development of foundation models is grounded in well-established scientific principles from fields like statistics, mathematics, and computer science, typically mastered by data scientists over the course of their undergraduate or graduate studies. Understanding these principles is crucial for making informed decisions about model design and training. Training these models involves designing experiments to test hypotheses about model behavior, testing their performance on real-world tasks

¹⁰ For instance, Google’s “Processing of Deep Networks” (Szegedy and Vanhoucke 2017) covers methods for the compression of deep neural networks, i.e. reducing the computational complexity of deep neural networks while preserving, as much as possible, their performance. Microsoft’s “Neural network categorization accuracy with categorical graph neural networks” (Du et al. 2024) pertains to the enhancement of neural network-based categorization using graph neural networks. NVIDIA’s “Virtual environment scenarios and observers for autonomous machine applications” (Nassar et al. 2024) describes methods for AI systems to test and validate themselves.

¹¹ For instance, Adobe has patented technologies related to content generation and manipulation, such as the “Content-Aware Fill” feature in Photoshop, which intelligently fills in the space when a part of an image is deleted.

and datasets, and drawing conclusions in a methodical way. This scientific knowledge and associated practices are broadly available to all firms in the industry, short-term skill shortage notwithstanding.

However, training generative AI models also requires knowledge more akin to craft than the scientific method. Seasoned data scientists often develop an intuition for choosing the right model architectures, data preprocessing methods, and training parameters. Addressing challenges like overfitting, underfitting, or dealing with data irregularities can require solutions seldom found in textbooks. Finally, the process of tuning hyperparameters, designing human reinforcement learning programs, and iterating over different model configurations involves much trial and error, rather than a systematic process prescriptively determined in advance.

Craft knowledge is more likely to remain tacit as it is deeply embedded and hard to unbundle from the broader architectural choices made by each firm in the design and training process. In this respect, generative AI models may be like other fields buffeted by technological breakthroughs. The spread of recombinant DNA research for instance, in spite of a seemingly broadly accessible scientific foundations, was for many years limited by pioneer researchers' ability to engage the next generation of innovators in long periods of craft-like apprenticeship (Zucker et al. 1998).

In addition to tacit knowledge, firms investing in the design and training of foundation models may be able to keep their models' weights proprietary. Without access to the weights, others cannot modify the pre-trained model without having access to the original training data or computational resources. Sponsors of generative AI models can forestall reverse engineering by maintaining tight control in other areas, including the specifics of the training data corpus (which encompasses the exact datasets used, preprocessing steps that have been implemented), details about the infrastructure and scaling technology), and the methods employed for fine-tuning (which are essential for adapting a pre-trained model to a particular task or dataset). The intricacies of the user interface layer—which facilitate smoother interaction with the model through mechanisms such as natural language prompts—are yet another aspect that might be hard to replicate at scale. Lastly, foundation model developers might take inspiration from industries like stock photography and genetically modified seeds and develop their own technologies to detect IP leakage or unlicensed commercialization using their LLMs.

Our interim conclusion is that pioneering firms in the industry benefit from a tight appropriability regime more owing to the many avenues to keep critical knowledge proprietary or tacit, rather than the assertion of formal intellectual property such as patents. This conclusion must remain guarded, however, since the field is still at an early stage of development. As AI talent becomes less scarce, the mobility of engineers and other technical staff will contribute to the circulation of knowledge that previously circulated as snippets within circumscribed communities of practice.

3.2 Complementary Assets

Complementary assets refer to the set of assets or capabilities required to effectively commercialize and extract value from an innovation (Teece 1986). These may include specialized manufacturing capabilities, access to distribution channels, service networks, or complementary technologies. Complementary assets may be generic and therefore not tightly held by incumbent firms. In contrast, other complementary assets may be narrowly tailored and dependent on a specific innovation, with few or no substitutes in the industry. In this case, complementary assets can be deployed by

incumbents in such a way that they constitute potent barriers to entry and imitation. As we discuss below, leading firms in the generative AI industry have already begun developing complementary assets. At least some of these assets have the potential to become the exclusive fiefdoms of a narrow set of large, incumbent technology firms. The following six complementary assets are probably the most salient.

3.2.1 Compute environment

The compute environment refers to the hardware and software infrastructure used to train and fine-tune foundation models. It includes (1) hardware for substantial computational power, often provided by GPUs, which are specialized components for handling the parallel computations that are common in machine learning; (2) a software stack (e.g., operating system, machine learning frameworks); (3) networking infrastructure to facilitate quick data transfer between different computers; (4) integration with cloud services to scale resources up or down based on the needs of the model training or inference; and (5) tools to monitor the performance of the hardware, the progress of model training, resource utilization, and to manage the deployment and operation of models.

Given their large scale, a well-optimized compute environment is crucial for the efficient training and testing of foundation models. In practice, this entails a careful balance of hardware capabilities, software support, and infrastructure management to ensure that the models can be trained, fine-tuned, and served effectively. An active computer science literature explores how model sponsors should trade-off model size, training dataset size, and compute budget at the margin (Hoffmann et al. 2023). This research broadly points to “scaling laws” reminiscent of Moore’s Law (Kaplan et al. 2020), and as a result the capital expenditures needed to set up, operate, and maintain the necessary infrastructure is increasing with model size and complexity. To fix ideas, Meta’s CEO Mark Zuckerberg recently announced that in its effort to train the next generation of its Llama large language model (Llama 3), the company was building a massive compute infrastructure including 350,000 of NVIDIA’s H100 GPUs by the end of 2024. Valued at the retail price of this crucial piece of equipment, the GPU investment alone would amount to approximately 10 billion dollars.¹² Even if such a model was “open” (in the sense that the weights are publicly accessible), Meta would still exert tight control over these computing resources and, therefore, over how the model will be trained in future iterations.

3.2.2 Model serving and model inference capabilities and infrastructure

Once trained, foundation model sponsors must possess the capability to deploy the model in a production environment where it can receive input data, process it, and provide outputs to end-users or other systems. The computational power needed for deploying large AI models in practical applications often surpasses the initial training requirements. While precise data on this subject is scarce, an informal analysis indicates that the operational costs of running the ChatGPT interface might exceed the training expenses of the GPT-4 model it is based on every week.¹³ Data scientists confusingly denote the task of applying a learned model to make predictions or decisions on new, unseen data as “inference,” and high inference costs have prompted OpenAI’s CEO, Sam Altman, to acknowledge in a congressional testimony that the company aims to design systems that do not

¹² <https://www.theverge.com/2024/1/18/24042354/mark-zuckerberg-meta-agi-reorg-interview>

¹³ <https://www.semianalysis.com/p/the-inference-cost-of-search-disruption>

prioritize user engagement, partly due to a scarcity of GPUs, which are crucial for AI training and operations (Oremus 2023).

In addition to compute infrastructure, sponsors must make the trained model accessible for practical use, for example through APIs or other interfaces. Serving a model involves tracking its performance, ensuring its uptime, and updating it (or the underlying infrastructure) to handle additional load or to integrate improvements. When models also process user-generated data, sponsors must set up secure data transmission channels, ensure data privacy, and comply with relevant regulations.

Lastly, the output of foundation models might be further enhanced by tight integration into the web browsing experience and wider ecosystem features, such as Google's Chrome or Microsoft's Edge web browsers or Meta's WhatsApp messaging app. Not only might ecosystem owners elevate access to their preferred model by default (not unlike their current approaches to the selection of search engines), they might also enrich answers by leveraging location data or the past browsing history of individual users.

3.2.3 Safety and governance

Sponsors must engage in a range of practical measures to ensure that foundation models are developed and used responsibly. The specifics vary between organizations, but generally begin with a set of ethical principles that guide development and deployment, along with the creation of boards or committees responsible for overseeing compliance with these principles. From a technical standpoint, sponsors may develop methods to make their systems more transparent and understandable to humans, and implementing measures to protect the data the models are trained on and ensure the models do not leak private information.

Sponsors also invest in research dedicated to understanding and mitigating potential risks associated with AI, such as studying the problem of alignment (ensuring AI systems reliably do what their operators intend) and robustness (ensuring AI systems behave safely under a wide range of conditions).¹⁴ To some degree, these R&D investments are product- or brand-specific choices undertaken with the goal of differentiation and risk management. However, developing safety systems might prove to be necessary to placate regulators and earn public trust. For example, model developers often build in safety valves to make sure that LLMs do not output large chunks of copyrighted information from its training data. When testing new LLMs and their derivatives, conducting audits of released generative AI products, and handling customer complaints, incumbents that invested and demonstrated such safety systems may have both efficiency and credibility advantages.

Investing in AI safety can be quite expensive, but the exact figures are often not publicly disclosed and can vary widely between organizations and depending on the scale of the AI systems being developed. The expenses include not only direct financial investments in research and development but also the opportunity costs of implementing rigorous safety measures, which can slow down the development process. Despite these costs, leading AI organizations tend to view these investments as

¹⁴ Industry observers caught a glimpse of the importance of these research investment with a series of events surrounding the departure of Dr. Timnit Gebru from Google in late 2020. This incident attracted significant attention and led to widespread discussion about academic freedom, the ethical implications of AI, and diversity and inclusion within the tech industry.

essential, both for ethical reasons and for the long-term viability and acceptance of their services. In fact, AI safety is often seen as the main reason against promoting an open and competitive marketplace of foundation models. If strict AI safety requirements were instituted, leading incumbents might be able to further reinforce their lead. Such an effect would not be unlike the increase in market concentration in web service provision ushered in with the introduction of GDPR privacy regulation in Europe (Peukert et. al 2022).

3.2.4 Benchmarks and metrics

Commercial sponsors invest in the development of benchmarks to evaluate the performance of foundation models, often in collaboration with academic researchers, industry experts, and policymakers.

This might take the form of investing in the creation or curation of large, diverse, and representative datasets that can be used to rigorously test the performance of foundation models across various tasks (such as the case of ImageNet), or in the form of hosting contests where researchers worldwide work to achieve the best performance on a given benchmark, thereby pushing forward the state-of-the-art in model evaluation. For example, the GLUE, SuperGLUE or the MMLU benchmarks are designed to evaluate the performance of models on a range of natural language understanding tasks such as sentence completion, question answering, and sentiment analysis (Hendryks et al. 2020, Wang et al. 2019).

Another area in which benchmark development is becoming more salient is governance, where firms endeavor to publish and promote safety methodologies for AI systems. For instance, OpenAI published Safety Gym, a suite of reinforcement learning environments and tools for developing AI systems that learn to perform tasks while respecting safety constraints (Ray et al. 2019).

3.2.5 Training data

The training data corpus for foundation models is typically extensive, diverse, and sourced from a wide range of public or semi-public sources. Foundation model sponsors must bear the costs of harvesting, ingesting, storing, and preprocessing these data for training purposes. In some cases, they must guard against contaminating the corpus with biased and sensitive content, and keep it updated to prevent “data staleness.”

Some of the data used for training might be in the public domain, such as the Protein Data Bank (used to train structure prediction models such as AlphaFold or RoseTTAFold), ImageNet (a large visual dataset used for image classification), and Project Gutenberg (a massive dataset of books in the public domain). But many other data sources could be characterized instead as “semi-public,” i.e., data that is not explicitly private (like personal e-mails or messages) but exists in a gray area concerning usage and privacy expectations. Examples include social-media posts, user-generated content on forums and discussion boards, code repositories, product reviews and feedback, academic papers, and newspaper articles.

When using such semi-public data, model sponsors need to consider legal aspects (like copyright and data protection laws), ethical considerations (like the reasonable expectation of privacy and the intentions of the content creators), and the potential impact on their reputation and trust with users. Most sponsors do not typically disclose the exact datasets used in the training process, while also

claiming they will respect opt-out requests and endeavor to use the data in ways that align with the expectations of the individuals and organizations that originally created the content.

The legal framework that will govern the use of copyrighted data in the training of foundation models is not settled (Lemley and Casey 2020; Gans 2024). While sponsors typically invoke the doctrine of fair use, this is contested by copyright holders, and a flurry of legal cases suggest that sponsors might mitigate legal risks by using data for which they have licenses or explicit permission to use for training purposes. Of course, a consequence might be that only sponsors with sufficient financial resources or legal expertise will be in the position of acquiring training datasets of sufficient breadth and depth, while taking care of purging personal information accidentally found in the training data and more generally deal with biased and sensitive content. The evolution of copyright policy in this area (especially whether and to what extent training on copyrighted material will be deemed fair use) will have important implications for competition in generative AI.

3.2.6 Data network effects

An additional consideration is whether a very particular type of network effect is likely to arise from sheer access to massive amounts of training data. For LLMs specifically, as more data is collected and used in training, the model generally becomes better at understanding and generating natural language. This is due to the model's increased exposure to various linguistic patterns, nuances, and contexts. In some implementations, user interactions with the model can be used to further refine its output. For example, when users correct or flag inappropriate responses, this feedback can be used to improve the model. As more users interact with the model, the quality of the model can improve, benefiting all users.

Reinforcement Learning from Human Feedback (RLHF) is a training strategy whereby a model is fine-tuned based on feedback that reflects human preferences or corrections. RLHF has been particularly influential when training models for complex tasks where traditional reward functions are insufficient to capture all aspects of desired behaviors, such as language models, robotics, and game playing. It typically involves an iterative process where the model is continuously updated based on ongoing feedback from human raters. In this iterative cycle, each new piece of data adds value to the model, and the improved model, in turn, provides better service to users, encouraging more engagement and feedback.

Since RLHF involves tradeoffs between generalization and diversity of output (Kirk et al. 2024), it is probably too early to assert it will give rise to data network effects that will favor firms able to create the infrastructure necessary to capture human feedback on a large scale. We choose to mention it as one of a broad class of complementary technologies that, if deployed in conjunction with massive training datasets, might endogenously lead initial performance advantages to cumulate and magnify over time.

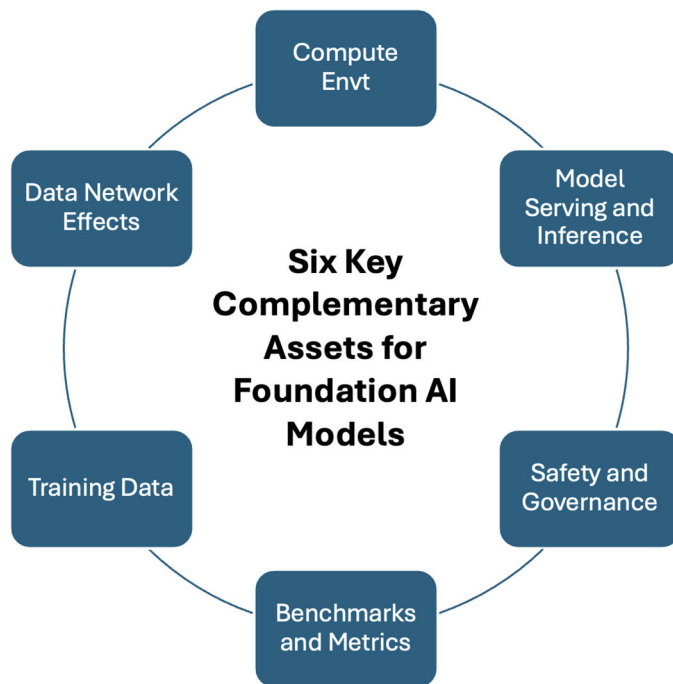
3.2.7 The durability of complementary assets

While the list of complementary assets above (summarized in Figure 3) is not meant to be exhaustive, it suggests nonetheless that complementary assets have the potential to help turn foundation models from standalone technologies into integral components of products, services, or business processes.

Their importance lies not only in enhancing the performance of these models but also in ensuring that they are usable, reliable, and aligned with broader organizational goals and societal norms.¹⁵

However, the inhibition of competition through tight control of complementary assets is subject to change as technological developments either amplify or depreciate their value (Tripsas 1997). For instance, much research effort has been recently directed to alleviating the burden of initial model training, thus lessening the need to invest in an extensive compute environment to make contributions at the frontier of the field. DeepMind’s Chinchilla study highlighted that the size of large language models (measured in tokens) and the size of its training corpus are broad substitutes (Hoffmann et al. 2022). Simultaneously low-rank adaptation (LoRA) is a method designed to efficiently fine-tune large models without the need to retrain them fully, which can be computationally expensive and time-consuming (Hu et al. 2021). Combined with techniques in model distillation, these developments offer the tantalizing possibility of a level playing field between large tech companies, startups, and academic groups.

Figure 3: Key Complementary Assets Shaping Competition in Foundation Models



The research frontier in AI foundation models is perpetually shifting outward. On the one hand this could provide opportunity for laggards to leapfrog over the achievements of pioneer innovators. On

¹⁵ At the same time, the gatekeeping effects of complementary assets need to be ascertained on a case by case basis. In structural biology, for example, the training datasets have so far mostly consisted of publicly available data, and the users of foundation models are technically sophisticated, lessening the importance of model serving capabilities. To date, the compute environment necessary to train these models has not put the frontier out of reach for academic labs, as attested by the arrival of the RoseTTAFold model (Baek et al. 2021), quick on the heels of DeepMind’s AlphaFold (Jumper et al. 2021).

the other hand, there is likely no model training “free lunch,” with follow-on innovators requiring access to an infrastructure of their own to realize the promise of their innovative approaches.

If these complementary assets maintain their relevance over time, the companies likely to dominate the AI sector in the future will likely be recruited from the ranks of current incumbents. Yet, another possibility is that the pivotal complementary assets that would enable a handful of players to dominate the AI landscape are not yet in evidence, because it is not clear yet what key challenges need to be overcome in order to facilitate AI adoption across a broad cross-section of economic sectors.¹⁶ If that alternative view turns out to be correct, the industry would still evolve towards an oligopolistic structure, but its commanding heights might be controlled, by new entrants, including some not yet founded at the time we write these lines, rather the current set of “usual suspects.”

3.3 A natural experiment in appropriability: The Llama leak and its consequences

Relative to Google and OpenAI, Meta’s AI research efforts have distinguished themselves by embracing an approach of relatively open sharing. In April 2022, the company released its OPT (Open Pre-trained Transformer) large language model, along with extensive documentation describing model architecture and design, training data/methodology, and detailed performance information. Even seemingly unabridged log books detailing the many steps needed to train the model were made available, which provided a rare glimpse into the craft-like techniques used by data scientists in their model training efforts.

However, the model weights for the 175 billion parameters OPT-175B were not released to the public. Meta’s decision to withhold the model weights was likely due to concerns about potential misuse and the need for responsible AI development and deployment. One year later, Meta followed a similar approach with its more advanced Llama (Large Language Model Meta AI) model.¹⁷ Merely a week after Meta began accepting inquiries for Llama access, the model found its way online. On March 3rd, 2023, a downloadable version of the system (including the model’s weights) appeared as a torrent on 4chan, and quickly disseminated through numerous AI forums.

The accidental leak provided a glimpse into the importance of openness for innovation in the AI foundation model domain, sparking a rapid series of cascading events. Even though Llama’s release was stripped to essentials, the availability of pre-trained weights allowed researchers who might not have the resources to train such large models from scratch to experiment with and improve upon a state-of-the-art LLM. Having a high-quality LLM available without constraints also permitted comprehensive benchmarking and comparative studies.¹⁸ Following the leak, numerous Llama-based developments emerged in quick succession, including Stanford’s Alpaca, an instruction-following model; Vicuna from UC Berkeley, CMU, Stanford, and UC San Diego; BAIR’s Koala, enhanced with internet dialogues; Nebuly’s ChatLLama for custom conversational assistants; FreedomGPT, an

¹⁶ <https://joshuagans.substack.com/p/so-whos-really-going-to-win-with>

¹⁷ <https://ai.meta.com/blog/large-language-model-llama-meta-ai/>

¹⁸ In text-to-image generation, Stable Diffusion’s launch in 2022 marked the first time that a foundation model was released as open source. Yet, until Llama’s release, this approach had not been mirrored in the large language model (LLM) arena, where frontier advances such as GPT-4, Claude, and Cohere were API-bound, and therefore preventing third party researchers to peek “under the hood.”

Alpaca-based conversational agent; and UC Berkeley’s Colossal-AI project’s ColossalChat, a ChatGPT-style model using Llama.

Within a month, these “spawns of Llama” included variants with instruction tuning (Zhang et al. 2023), quantization,¹⁹ and RLHF (Köpf et al. 2024). Based on our tabulation of the data gathered by the Center for Research on Foundation Models, by the end of 2023, 30% of large language models introduced by startups and academic groups were built on Llama or its direct descendants, Alpaca and Vicuna.

These open-source developments sent shockwaves through the major tech firms’ engineering staff, as they seem to imply that these firms would struggle to maintain control of the knowledge at the core of the models they sponsored. With much trepidation, a leaked internal document from Google claimed that *“while our models still hold a slight edge in terms of quality, the gap is closing astonishingly quickly. Open-source models are faster, more customizable, more private, and pound-for-pound more capable. They are doing things with \$100 and 13B params that we struggle with at \$10M and 540B. And they are doing so in weeks, not months.”*²⁰

The impact of these investments has been amplified by the platform *Hugging Face*, whose vibrant developer community has catalyzed open-source innovation by democratizing access to pre-trained models, offering access to libraries and tools, providing users with the ability to fine-tune models for specific tasks or datasets, and facilitating direct comparison of models on standardized benchmarks (Greenstein et al. 2022). In particular, Hugging Face hosts and maintains benchmark datasets and leaderboards, which are crucial for tracking the state-of-the-art and assessing the progress of foundation models over time.²¹

This natural experiment in openness demonstrates that leading firms’ attempts to titrate access to their innovation by keeping some of the implementation details secret have not been successful. Leaks are hard to prevent and police, and once a model’s inner workings are exposed for all developers to see and tinker with, incumbent firms may not be able to wind the openness clock back. Further, the rapid development of the “spawns of Llama” demonstrate that a vibrant and talented open source developer community stands ready to build on open models to extend the generative AI frontier.

Less certain is whether the open source AI ecosystem’s “green shoots” foretells the devaluation of some of the complementary assets mentioned earlier. Llama’s availability catalyzed cumulative innovation in the open source developer community, but it did not hinder the widespread adoption proprietary models released by leading firms such as OpenAI, Anthropic, and Google. Even if the

¹⁹ Quantization refers to a set of techniques to make it feasible to run powerful AI models directly on user devices, enabling real-time processing without the need for constant cloud connectivity. A cognate concept is that of model *distillation*, whereby knowledge from a large, complex model (or an ensemble of models) is transferred to a smaller, more efficient model.

²⁰ <https://www.semianalysis.com/p/google-we-have-no-moat-and-neither>

²¹ An important caveat is that, in the absence of well-validated benchmarks to compare the performance of different models, we should not treat claims of performance parity between open source models and large scale implementations such as OpenAI’s GPT-4 or Google’s BARD at face value. It may well be that recent developments democratize access to language models useful for a broad range of niche applications, while human-like performance in general tasks remains the prerogative of firms endowed with the resources to train and retrain ever larger models using the latest, most expensive, and most energy-consuming hardware.

importance of the compute environment fades over time, other complementary assets—such as investments in AI safety and access to unique training data—might still result in an industry with an oligopolistic structure, albeit one with a highly variegated fringe of smaller actors.

3.4 LLAMA in the ointment: Competitive scenarios with one “rogue” tech firm

Our pessimism regarding the emergence of a vibrant and competitive foundation model marketplace is grounded in the belief that large tech firms have incentives to close their models, multiple technical avenues to do so, and will likely exert a stranglehold on access to some essential complementary assets. But what if one “rogue” technology firm chose to deviate from this predictable script? We mention this possibility because, following the LLAMA leak, Meta’s founder Mark Zuckerberg has explicitly expressed his company’s intent to pursue an open source strategy as it invests massively in the quest to achieve “General Artificial Intelligence.”

We can speculate on the reasons that might lead a powerful AI pioneer to choose an open source approach. That company could benefit indirectly from a vibrant ecosystem that uses its open models. For example, improvements by third parties (such as chip developers) could be integrated back into its own offerings, enhancing their value at minimal cost. Further, if that company’s business model relies on selling value-added services on top of foundation models, standardizing the underlying model might benefit these follow-on services. Third, by positioning itself as a leader in open AI, such a firm could significantly enhance its reputation and establish itself as a benevolent technology leader, influencing standards and practices across the industry. Lastly, an open and collaborative approach could make that firm a magnet for top AI talent—researchers motivated by the prospect of working on high-impact, accessible models. A recent essay by Meta validates at least a few of these claims.²²

In the scenario where Meta (or any other large tech firm) “goes rogue” and commits to a genuinely transparent approach for the development of its foundation model, the competitive landscape could be transformed in a manner akin to the early dynamics between Android and iOS in the smartphone industry. Meta’s approach could democratize access to powerful AI tools, much as Android initially made smartphone technology widely accessible, enabling a broad range of manufacturers to produce devices that catered to various market segments. This could significantly lower the costs of exploration for startups and academic teams who might otherwise hesitate to investigate new use-cases for the nascent AI technology. By lowering the cost of experimentation and adoption, the rogue open-access competitor could catalyze a surge in AI applications, unlocking innovative opportunities across a broad array of sectors and geographies that would otherwise remain underdeveloped.

In the smartphone market, Android’s open system served as a counterbalance to Apple’s more closed and expensive iOS, fostering wider smartphone adoption globally and supporting a diverse array of application developers. This competition kept iPhone prices somewhat in check and spurred continuous innovation relative to what might occur in a monopoly scenario. Similarly, an open-source AI model from a leading firm could encourage rapid advancements in AI applications, propelled by an open collaboration model that accelerates iterative experimentation.

²² <https://about.fb.com/news/2024/07/open-source-ai-is-the-path-forward>

However, just as Android’s nominally open-source license has not resulted in a flurry of entry into the smartphone OS market, the existence of a leading open access foundation model would not by itself ensure access to the complementary assets that are necessary to build a viable foundation model in the AI marketplace. In fact, by reducing the profit opportunities available to competitors at the foundation model layer of the stack, the rogue firm might divert innovation efforts to the downstream application layer. Additionally, the long-term sustainability of Meta’s openness “crusade” is uncertain. As in the case of the Android ecosystem, where Google has occasionally shifted its strategy and tightened control over aspects of the platform, Meta could potentially alter the terms of access, especially if the economic or strategic benefits of the open approach turn out to clash with its broader corporate goals.

4 Implications for public policy

Given this framework, how might policy makers think about competition policy in this setting? Below, we highlight four low-hanging fruits that may prove socially beneficial without committing government authorities to wade in too deeply in the dynamics of a still nascent industry.

4.1 Credible performance benchmarks

First, it might be worthwhile for the government to coordinate the establishment of benchmarks for foundation models, both with respect to their performance as well as their safety. In spite of its ubiquity, the task of comparing foundation models is fraught with challenges, making definitive claims of superiority less compelling than they often appear.

One of the primary difficulties in comparing foundation models lies in their multifaceted nature and the great variety of tasks they are designed to perform. LLMs and related image generation AI tools are not only evaluated based on their accuracy but also on factors like efficiency, scalability, adaptability to different languages or domains, and ethical considerations such as bias and fairness. Given this, a model that excels in one benchmark task or metric might underperform in another, complicating direct comparisons.²³ Moreover, the rapid pace of development in AI means that benchmarks themselves can quickly become outdated, as models evolve to address their limitations and new capabilities are developed. Further, given a lack of transparency around what data is used to train a model, whether a model truly “learned” to perform a particular task or whether it is simply regurgitating responses from its training data can be hard to assess. Further, even if some models are better at taking tests or achieving pre-determined targets than others, their economic utility rests on their impact in real-world applications such as assisting call-center agents or copywriters; at the time of this writing, such impacts lie far beyond the purview of existing benchmarks (Brynjolfsson et al. 2023, Noy and Zhang 2023).

²³ For example, the information retrieval and prediction needed to climb a leaderboard that ranks LLMs based on their performance on Jeopardy questions (<https://github.com/aigoopy/llm-jeopardy>) is likely quite different from the features that make LLMs perform well on software programming tasks (<https://huggingface.co/spaces/mike-ravkine/can-ai-code-results>). Furthermore, such performance metrics may not adjust for speed, cost, or model size.

Given these challenges, claims of one model being superior to another, or of open-source models catching up to their proprietary counterparts, should indeed be approached with caution. Without a comprehensive and nuanced understanding of both the benchmarks and the models being compared, such claims can be misleading. One recent example involves a careful comparison of several open-source model LLMs with the proprietary models GPT-4 and Claude 2 in the context of the Nephrology Self-Assessment Program, a multiple-choice questionnaire administered by the American Society of Nephrology to help clinicians assess their disciplinary knowledge (Wu et al. 2024). Large performance gaps emerged between the open-source models and the proprietary ones, casting doubt on the reliability of unvalidated leaderboards such as those used by the leading open source repository Hugging Face. Similarly, one open source project compared over 200 open-source models to the six leading closed models on a variety of logic-based word questions. While the top-performing open-source models produced impressive results and were competitive with the laggard closed-source models, the top performing closed models (GPT-4 and Claude 3) outperformed all other models by a significant margin.²⁴

This complexity and the need for nuanced comparison provide a potential avenue for government intervention, specifically in the development of widely accepted and well-validated benchmarks. An organization like the National Institute for Standards and Technology (NIST) could play a crucial role in this regard. NIST, with its history of developing standards across various domains, could bring a level of rigor, transparency, and neutrality to the process of benchmarking AI models. By convening a diverse set of stakeholders from academia, industry, and civil society, NIST could help ensure that benchmarks are comprehensive, up-to-date, and reflective of the broader societal impacts of AI technologies.²⁵

To avoid the often-criticized incentive problems of third-party rating agencies in finance, a government-sponsored generative AI benchmarking organization might instead follow the structure of the National Renewable Energy Lab (NREL). As a government-owned, contractor-operated organization, NREL plays an important role in the performance testing, measurement, and validation of new energy technologies like solar photovoltaics. Instead of sending prototypes of a solar cell for efficiency testing in a controlled environment, generative AI researchers could provide LLM access to these third party evaluators for a transparent set of performance tests and published results.

In the rapidly evolving domain of AI, the ability to accurately assess and compare the capabilities of different models is not just a technical necessity but also a regulatory imperative. Without clear standards, assessing the competitive landscape becomes fraught with uncertainty, making it challenging to identify actual market leaders or to evaluate claims of technological parity critically. This uncertainty may lead to either premature regulatory interventions or, conversely, a lack of action in the face of emerging monopolies.

²⁴ https://docs.google.com/spreadsheets/d/1NgHDxbVWJFolq8bLvLkuPWKC7i_R6I6W

²⁵ One could argue that the US government is already heeding this advice. In February of 2024, the US Department of Commerce and NIST announced the creation of the AI Safety Institute Consortium (AISIC), which included a consortium of more than 200 large incumbent technology companies and financial institutions with the goal of developing AI safety tools and standards (<https://www.reuters.com/technology/us-says-leading-ai-companies-join-safety-consortium-address-risks-2024-02-08/>). The exclusive focus on safety appears misplaced to us, although a neutral arbiter like NIST offers prospects that emerging safety standards will not be tilted to favor entrenched incumbents.

Moreover, the fast emergence of such standards could help demystify the actual progress and capabilities of AI technologies for policymakers, businesses, and the public. It would help distinguish hype from genuine advancement, ensuring that public and private resources are allocated efficiently and that innovations that truly push the boundaries of AI are recognized and fostered.

4.2 Hastening the pace of clarification of key legal issues

The development and adoption of AI foundation models, including LLMs and those capable of generating images or films, raise several legal issues that courts will likely need to assess. These issues create a degree of uncertainty regarding the returns on investments in this rapidly evolving domain. The legal challenges span the areas of intellectual property rights, privacy concerns, liability for harm, employment issues (such as the delicate interplay of non-compete agreements and trade secrecy law), and regulatory compliance, among others.

These legal challenges highlight the need for clear regulatory frameworks and legal precedents to guide the development and deployment of AI foundation models. Until these issues are more definitively addressed by courts and regulators, companies and investors in the AI domain will face ongoing uncertainty about the legal landscape and the risks associated with their AI endeavors.

Governments, including the U.S. government, could take two steps to hasten the pace of legal clarification surrounding AI foundation models without introducing comprehensive new legislation.

First, they could wade into the use of copyrighted data for training purposes. A Supreme Court decision on proper scope of fair use in AI model training would have a huge impact on investment incentives. The solicitor general could encourage the highest Court in the land to take on an AI fair use case. The Copyright Office could also issue guidance, as it has done on the question of whether AI output is copyrightable. Even if immediately challenged in the courts, this might force the legal system to adjudicate the underlying issues at a less leisurely pace.

Second, Regulatory agencies could issue guidance, frameworks, and best practices for AI development and use. This approach would provide immediate clarity on how existing laws apply to AI technologies, including intellectual property rights, privacy, liability, and ethical considerations. For instance, the Federal Trade Commission (FTC) could offer guidance on applying consumer protection laws to AI products and services; or the FDA could develop disclosure standards that would be required when AI software is involved in patient care, or regulating the use of AI in medical devices.

These two steps can help provide greater clarity and guidance for AI developers and users, addressing some of the most pressing legal uncertainties without the need for comprehensive new legislation. The options above would contribute to preparing the legal landscape for future initiatives, perhaps introduced by Congress through new legislation.

4.3 Encouraging the fractionalization of infrastructure

One way policy makers might mitigate market dominance by a few large players is to encourage the “fractionalization” of complementary assets so that smaller firms or academic researchers can experiment and build products on top of the leading technologies, without the barrier of large capital

expenditures. In this context, fractionalization means the unbundling of the distinct activities and infrastructure (e.g., compute, training data, safety systems) required to bring forth and improve foundation AI models. Focusing policy attention on the complementary assets most likely to elude market-led democratization efforts would amplify the effectiveness of such interventions.

Policy efforts to encourage and simplify entrant innovators' access to shared AI infrastructure could take many forms, from heavy-handed interventions such as compulsory licensing or mandated access, to more subtle measures whereby governments delicately place their finger on the scales to alleviate specific frictions or bottlenecks. Looking across sectors in the recent past, there are interesting precedents where government agencies played a constructive role to level the playing field by encouraging technology sharing. For instance, in the late 1990s, the National Institutes of Health led negotiations between the public company DuPont, the non-profit Jackson Laboratories, and the academic community to reach a memorandum of understanding to make genetically modified mice available to academic researchers through a simple license and material transfer agreement (Murray et. al 2016). Stretching the analogy, one might imagine that foundation model giants might prefer negotiating larger deals for compute and AI training systems with a small number of large research funding agencies, rather than being enmeshed in a myriad of bilateral collaboration agreements with universities and startups. Funding agencies could then decide how to allocate “compute credits” to their many recipients.²⁶

Beyond supplying infrastructure and standards for generative AI research and testing, government policy might also support the fractionalization of AI assets by providing credible demand signals for new AI tools. Modeled off of government programs that have accelerated private industry development in frontier technologies like space launch (NASA's Commercial Orbital Transportation Service program) and fusion energy (the Department of Energy's Milestone-Based Fusion Development Program), the government might co-finance companies in the early development of cutting-edge generative AI “infrastructure” tools aimed at reducing barriers to entry, then provide subsequent larger awards to companies that achieve new milestones and publish their methods and findings. Milestones could include achievements in model compression and distillation (reducing needs for expensive compute), fine-tuning systems, performance benchmarking, detection of unauthorized model use, and reducing the energy intensity of cloud computing.

In a nascent industry, there is a legitimate concern that forceful government interventions could dull incentives for subsequent innovation by leading firms. While measures such as the compulsory licensing of technologies, the development of public APIs, or mandated access to infrastructure could be justified in the future, we believe it to be more fruitful to first envision a role where the government creates a legal framework and norms to govern collaborations between leading firms and teams from academia and startups. Although incumbents might at first resist even these “light touch” interventions, these might well rebound to their benefit if the collaborations they enable provide leading firms with greater insight into new technologies and acquisition targets.

²⁶ A recent example in this spirit is the National Deep Inference Fabric (NDIF) Experiment Program, supported by the NSF and developed in collaboration with researchers at Northeastern University and the Public Interest Technology University Network (a consortium of 63 universities and colleges). NDIF gives researchers remote access to large scale compute and “inference” systems to test new methods on frontier open foundation models like Llama3.

4.4 Platform Policy

Our discussion has been exclusively focused on competition between sponsors of AI foundation models. Distilled in a few words, our considered view is that the industry, left to its own devices, is likely to evolve an oligopolistic structure due to the stranglehold exercised by leading firms on at least some valuable complementary assets. Of course, limited competition in foundation models is compatible with the existence of one or more vibrant ecosystems of “fine-tuned” models that build on top of the infrastructure layer provided by foundation models. Such a platform evolution is portended by OpenAI’s fledgling store of third party plugins.

Just as in the example of application stores for smartphone operating systems, competition policy challenges are sure to arise to limit the bargaining power of the caretakers of these ecosystems—the companies controlling the foundation models—vis à vis third-party application providers. For example, how proprietary should the APIs that allow application developers to develop on top of foundation models be? Should foundation model sponsors be allowed to avail themselves of the data used by third party developers to fine-tune these models for their particular use case? And how unfettered should the power of model sponsors be in setting royalty rates and access fees to capture value generated by third-party developers? Since we deem the prospects for competition and innovation at the application layer of the technology stack to be most promising, preventing the walled gardens that prevailed in mobile computing to emerge in even more restrictive forms in the AI context will require heightened vigilance on the part of competition policy authorities. For examples, efforts to standardize the development of third-party applications such that they are interoperable with multiple foundation models might significantly enhance consumer choice in the application layer as well as effective competition at the foundation model layer.

5 Conclusion

Our analysis emphasizes how the technology and complementary assets of the nascent generative AI industry might favor a small set of leading firms as they build up a war chest of both specialized know-how and physical assets. Naturally, predicting another concentrated technology industry might suggest preemptive policy guardrails against market power or spark calls to subsidize entrants or offer government-created “public option” models to check the private sector oligopolists.²⁷ However, if the goal of innovation policy is to foster progress and technological leadership through the emergence of a diverse and commercially viable ecosystem, that outcome does not strictly require many competitors. The economics of innovation literature highlights the inverted U-shaped relationship between the number of competitors and incentives for innovation: not enough competition and firms exploit their market power without the urgency to invest in disruptive R&D; too much competition and profits erode, lowering innovation returns for all (Aghion et al. 2005; Segal and Whinston 2007).

Thus, competition policy efforts in the AI domain face the high wire act of fostering intense competition without dulling incentives for engaging in risky R&D. That daunting challenge is why our policy recommendations emphasize reducing barriers to experimentation—such as legal and regulatory uncertainty, milestone based private-public partnerships, public benchmarking, and access

²⁷ <https://foreignpolicy.com/2023/06/12/ai-regulation-technology-us-china-eu-governance/>

to AI research infrastructure. Instead of bluntly forcing incumbents to license or share their models, generative AI innovation policies could ideally encourage a robust ecosystem of academic researchers and firms of different types to compete for the development of next-generation AI models and tools. Through policies that are agnostic with respect to the size of organizations that will ultimately bring the most promising ideas to market, innovators can strive to complement or displace the leading firms' models, rather than merely duplicate their efforts by developing models with similar technological underpinnings.

A challenge specific to the rapidly evolving AI foundation model sector is that U.S. policymakers find themselves navigating between the Charybdis of fostering a competitive market and the Scylla of ensuring technological safety, resulting in a policy approach that can appear, at times, to be marked by a certain ambivalence, bordering on a dissociative identity disorder.

First, there is a pronounced concern among policymakers regarding the potential for market domination by a select few large tech companies, such as Meta, Google, and Microsoft. Left to their own devices, these behemoths could monopolize the AI field, stifle innovation and control the direction of future technological developments. This perspective champions the idea of breaking down barriers to entry and encouraging a vibrant ecosystem where startups and academic teams innovate freely, thus preventing any single entity from wielding disproportionate influence over the AI landscape.

In contrast, a palpable sense of caution permeates policy discussions, primarily centered on the safety and ethical implications of AI technologies. The rapid pace of innovation, especially among smaller startups with fewer resources to dedicate to ethical considerations, raises concerns about the potential for harmful applications, such as the creation and proliferation of deep fakes. Here, the larger corporations, with their established reputations and experience with operating under regulatory scrutiny, are seen as entities that might be more accountable and responsible stewards of powerful AI technologies.

The reaction to the “Model Forum,” convened by leading AI providers to develop safety standards, represents a quintessential example of the ambivalence permeating the AI policy landscape. One perspective lauds this initiative as a commendable example of self-regulation, reminiscent of the Asilomar conference, which played a pivotal role in addressing the ethical and safety concerns surrounding recombinant DNA technology (Berg et al. 1975; Frederickson 1991). This view sees the forum as a proactive step toward ensuring AI technologies are developed and deployed responsibly, mitigating risks and fostering public trust. Meanwhile, a more skeptical view might interpret the same forum as a strategic maneuver by incumbent firms to solidify their dominance. By setting the bar for compliance with safety standards they themselves crafted, these firms could inadvertently (or perhaps tactically) raise the operational and financial hurdles for emerging and future competitors, potentially entrenching their market position.

This ambivalence presents a conundrum for crafting policy recommendations aimed at fostering competition within the AI domain. In this article, we have argued that delineating a competition policy framework requires a nuanced understanding of the AI ecosystem, one that is grounded in the twin pillars of appropriability and complementary assets. The power of this conceptual framework, however, only becomes compelling if this sector is approached by policymakers like any other nascent

industry. Rather than reinventing the competition policy wheel for the AI age, we have argued that there is nothing inherently special about this domain. While there is always a risk of constraining the interpretation of new innovation phenomena within the boundaries of existing paradigms, our analysis offers a counterpoint to the fevered claims and policy proposals advanced by industry insiders, as well as their detractors.

References

- Aghion, Philippe, Nick Bloom, Richard Blundell, Rachel Griffith, and Peter Howitt. 2005. "Competition and Innovation: An Inverted-U Relationship." *The Quarterly Journal of Economics* 120(2): 701-28.
- Baek, Minkyung et al. 2021. "Accurate Prediction of Protein Structures and Interactions Using a Three-Track Neural Network." *Science* 373(6557): 871-76.
- Berg, Paul, David Baltimore, Sydney Brenner, Richard O Roblin III, and Maxine F Singer. 1975. "Asilomar Conference on Recombinant DNA Molecules." *Science* 188(4192): 991-94.
- Bertoline, Leticia M. F., Angélica N. Lima, Jose E. Krieger, and Samantha K. Teixeira. 2023. "Before and after AlphaFold2: An Overview of Protein Structure Prediction." *Frontiers in Bioinformatics* 3: 1120370.
- Bommasani, Rishi, Kevin Klyman, Shayne Longpre, Sayash Kapoor, Nestor Maslej, Betty Xiong, Daniel Zhang, and Percy Liang. 2023. "The Foundation Model Transparency Index." arXiv Preprint 2310.12941.
- Bresnahan, Timothy. 2024. "What Innovation Paths for AI to Become a GPT?" *Journal of Economics & Management Strategy* 33(2): 305-16.
- Bresnahan, Timothy, and Manuel Trajtenberg. 1995. "General Purpose Technologies 'Engines of Growth?'" *Journal of Econometrics* 65(1): 83-108.
- Brynjolfsson, Erik, Li, Danielle, and Lindsey R. Raymond. 2023. "Generative AI at Work." NBER Working Paper #31161.
- David, Paul. 1990. "The Dynamo and the Computer: An Historical Perspective on the Modern Productivity Paradox." *American Economic Review* 80(2): 355-361.
- Dosi, Giovanni. 1982. "Technological Paradigms and Technological Trajectories: A Suggested Interpretation of the Determinants and Directions of Technical Change." *Research Policy* 11(3): 147-162.
- Du, Tianchuan, Keng-Hao Chang, Ruofei Zhang, and Paul Liu. 2024. Neural Network Categorization Accuracy with Categorical Graph Neural Networks. U.S. Patent 11,960,573-B1, filed November 7, 2022.
- Foster, Richard. 1986. *Innovation—The Attacker's Advantage*. New York, NY: Summit Books.
- Frederickson, Donald S.. 1991. "Asilomar and Recombinant DNA: The End of the Beginning." In *Biomedical Politics*, ed. Kathi E. Hanna, 258-292. Washington, DC: National Academy Press.
- Gans, Joshua S. 2024. "Copyright Policy Options for Generative Artificial Intelligence." Working Paper No. 32106. National Bureau of Economic Research.
- Giczy, Alexander V., Nicholas A. Pairolo, and Andrew A. Toole. 2022. "Identifying Artificial Intelligence (AI) Invention: A Novel AI Patent Dataset." *The Journal of Technology Transfer* 47(2): 476-505.
- Greenstein, Shane, Daniel Yue, Kerry Herman, and Sarah Gulick. 2022. "Hugging Face: Serving AI on a Platform." Harvard Business School Case 623-026, November 2022 (Revised January 2023).
- Hasani, Ramin, Mathias Lechner, Alexander Amini, Daniela Rus, and Radu Grosu. 2021. "Liquid Time-Constant Networks." *Proceedings of the AAAI Conference on Artificial Intelligence* 35(9): 7657-7666.
- Henderson, Rebecca. 1995. "Of Life Cycles Real and Imaginary: The Unexpectedly Long Old Age of Optical Lithography." *Research Policy* 24(4): 631-43.
- Hendrycks, Dan, Burns, Collin, Basart, Steven, Zou, Andy, Mazeika, Mantas, Song, Dawn and Jacob Steinhardt. 2020. "Measuring Massive Multitask Language Understanding." arXiv Preprint 2009.03300.

- Hoffmann, Jordan et al. 2022. “Training Compute-Optimal Large Language Models.” arXiv Preprint 2203.15556.
- Hu, Edward J., Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. “Lora: Low-Rank Adaptation of Large Language Models.” arXiv Preprint 2106.09685.
- Iansiti, Marco, and Karim R. Lakhani. 2020. *Competing in the age of AI: Strategy and Leadership When Algorithms and Networks Run the World*. Boston, MA: Harvard Business School Press.
- Jumper, John et al. 2021. “Highly Accurate Protein Structure Prediction with AlphaFold.” *Nature* 596(7873): 583-589.
- Kaplan, Jared, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. “Scaling Laws for Neural Language Models.” arXiv Preprint 2001.08361.
- Kirk, Robert, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2024. “Understanding the Effects of RLHF on LLM Generalisation and Diversity.” arXiv Preprint 2310.06452.
- Köpf, Andreas et al. 2023. “OpenAssistant Conversations – Democratizing Large Language Model Alignment.” arXiv Preprint 2304.07327.
- Kuhn, Thomas S. 1962. *The Structure of Scientific Revolutions*. Chicago, IL: University of Chicago Press.
- Lemley, Mark A., and Carl Shapiro. 2005. “Probabilistic Patents.” *Journal of Economic Perspectives* 19(2): 75-98.
- Lemley, Mark A., and Bryan Casey. 2020. “Fair Learning.” *Texas Law Review* 99(4): 743-785.
- Liesenfeld, Andreas, Alianda Lopez, and Mark Dingemans. 2023. “Opening up ChatGPT: Tracking Openness, Transparency, and Accountability in Instruction-Tuned Text Generators.” arXiv Preprint 2307.05532.
- Murray, Fiona, Aghion, Philippe, Dewatripont, Mathias, Kolev, Julian, and Scott Stern. 2016. “Of Mice and Academics: Examining the Effect of Openness on Innovation.” *American Economic Journal: Economic Policy*, 8(1): 212-252.
- Nassar, Ahmed, Justyna Zander, and David Auld. 2024. Virtual environment scenarios and observers for autonomous machine applications. US Patent 20,240,078,363-A1, filed November 9, 2024.
- Noy, Shakked. 2023. “Experimental Evidence on the Productivity Effects of Generative Artificial Intelligence.” *Science* 381(6654): 187-192.
- Oremus, Will. 2023. “AI Chatbots Lose Money Every Time You Use Them. That Is a Problem.” *The Washington Post*, June 5, 2023. <https://www.washingtonpost.com/technology/2023/06/05/chatgpt-hidden-cost-gpu-compute/>
- Peukert, Christian, Stefan Bechtold, Michail Batikas, and Tobias Kretschmer. 2022. “Regulatory Spillovers and Data Governance: Evidence from the GDPR.” *Marketing Science* 41(4):746-768.
- Ray, Alex, Joshua Achiam, and Dario Amodei. 2019. “Benchmarking Safe Exploration in Deep Reinforcement Learning.” OpenAI Discussion Paper, Available at <https://cdn.openai.com/safexp-short.pdf>
- Ray, Joydeep et al. 2020. Compression in machine learning and deep learning processing. U.S. Patent 10,546,393-B2, filed December 30, 2017.
- Rosenberg, Nathan, and Manuel Trajtenberg. 2004. “A General-Purpose Technology at Work: The Corliss Steam Engine in the Late-Nineteenth-Century United States.” *Journal of Economic History* 64(1): 61-99.

- Schrepel, Thibault, and Alex Pentland. 2023. "Competition Between AI Foundation Models: Dynamics and Policy Recommendations." MIT Connection Science Working Paper 1-2003, Available at SSRN: <https://ssrn.com/abstract=4493900>.
- Segal, Ilya, and Michael D. Whinston. 2007. "Antitrust in Innovative Industries." *The American Economic Review* 97(5): 1703-30.
- Szegedy, Christian, and Vincent O. Vanhoucke. 2017. Processing images using deep neural networks. U.S. Patent 9,715,642-B2, filed August 28, 2015.
- Teece, David J. 1986. "Profiting from Technological Innovation: Implications for Integration, Collaboration, Licensing and Public Policy." *Research Policy* 15(6): 285-305.
- Tripsas, Mary. 1997. "Unraveling the Process of Creative Destruction: Complementary Assets and Incumbent Survival in the Typesetter Industry." *Strategic Management Journal* 18: 119-42.
- Utterback, James M. 1994. *Mastering the Dynamics of Innovation*. Boston, MA: Harvard Business School Press.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. "Attention Is All You Need." *Advances in Neural Information Processing Systems* 30.
- Watson, Joseph L. et al. 2023. "De Novo Design of Protein Structure and Function with RFDiffusion." *Nature* 620(7976): 1089-1100.
- Widder, David Gray, Sarah West, and Meredith Whittaker. 2023. "Open (for Business): Big Tech, Concentrated Power, and the Political Economy of Open AI." (August 17, 2023). Available at SSRN: <https://ssrn.com/abstract=4543807>.
- Wang, Alex, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. "SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems." arXiv Preprint 1905.00537.
- Wu, Sean, Michael Koo, Lesley Blum, Andy Black, Liyo Kao, Zhe Fei, Fabien Scalzo, and Ira Kurtz. 2024. "Benchmarking Open-Source Large Language Models, GPT-4 and Claude 2 on Multiple-Choice Questions in Nephrology." *NEJM AI* 1(2): <https://doi.org/10.1056/AIdbp2300092>.
- Yazici, Mahmut Taha, Basurra, Shadi, and Mohamed Medhat Gaber. 2018. "Edge Machine Learning: Enabling Smart Internet of Things Applications." *Big Data and Cognitive Computing* 2(3): 26.
- Zhang, Yue, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. 2023. "Multi-Task Instruction Tuning of Llama for Specific Scenarios: A Preliminary Study on Writing Assistance." arXiv Preprint 2305.13225.
- Zucker, Lynne G., Michael R. Darby, and Marilyn B. Brewer. 1998. "Intellectual Human Capital and the Birth of US Biotechnology Enterprises." *The American Economic Review* 88(1): 290-306.