# Social and Governance Implications of Improved Data Efficiency

Aaron D. Tucker*
Department of Computer Science,
Cornell University

Centre for the Governance of AI,
Future of Humanity Institute,
University of Oxford
aarondtucker@cs.cornell.edu

Markus Anderljung
Centre for the Governance of AI,
Future of Humanity Institute,
University of Oxford
markus.anderljung@governance.ai

Allan Dafoe
Department of Politics and
International Relations

Centre for the Governance of AI,
Future of Humanity Institute,
University of Oxford
allan.dafoe@governance.ai

## ABSTRACT

Many researchers work on improving the data efficiency of machine learning. What would happen if they succeed? This paper explores the social-economic impact of increased data efficiency. Specifically, we examine the intuition that data efficiency will erode the barriers to entry protecting incumbent data-rich AI firms, exposing them to more competition from data-poor firms. We find that this intuition is only partially correct: data efficiency makes it easier to create ML applications, but large AI firms may have more to gain from higher performing AI systems. Further, we find that the effect on privacy, data markets, robustness, and misuse are complex. For example, while it seems intuitive that misuse risk would increase along with data efficiency – as more actors gain access to any level of capability – the net effect crucially depends on how much defensive measures are improved. More investigation into data efficiency, as well as research into the "AI production function", will be key to understanding the development of the AI industry and its societal impacts.

## CCS CONCEPTS

• **Social and professional topics** → **Economic impact**; *Surveillance*; • **Computing methodologies** → **Machine learning**; *Transfer learning*; • **Theory of computation** → *Active learning*; • **Security and privacy** → *Economics of security and privacy*.

## KEYWORDS

Data efficiency, Production function, Competitive Advantage, Transfer learning, Active learning, Data markets

*Work done as a Summer Fellow at Centre for the Governance of AI

Figure 1: The Access Effect



## 1 INTRODUCTION

How does the performance of an artificial intelligence (AI) system scale with more data, more computational resources, and better algorithms? In other words, what is the *AI production function*[1]? This question influences the shape of AI progress, the structure of the AI industry, and the societal impacts of AI.
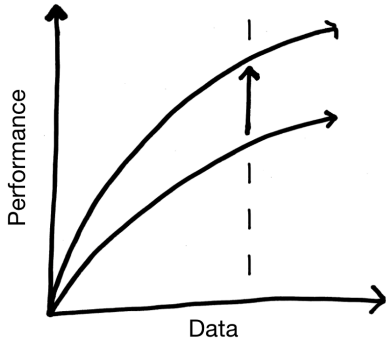
In this paper, we offer a preliminary analysis of one aspect of the AI production function - data. Specifically, we analyze the implications of increases in data efficiency[2]: increases in the performance a system achieves for any given data input. We find that increases in data efficiency have an *access effect* – where more actors get access to ML capabilities – and a *performance effect* – where performance for any given dataset is improved.

Predictions about the AI industry and its societal impacts often implicitly rely on claims regarding the AI production function and specifically on the relationship between data and the performance of systems. Kai-Fu Lee, for example, claims that China has an advantage with regards to developing AI technology, partly because we live in an "Age of Data", where "once computing power and engineering talent reach a certain threshold, the quantity of data becomes decisive in determining the ... accuracy of an algorithm" [21]. Views on the AI production function âĂŞ in particular current

---

[1] We depart slightly from the standard definition of production functions [7] by focusing on the relationship between the inputs to a machine learning (ML) system and the performance of the system, rather than between the inputs and outputs of using the ML system on a specific task.

[2] Computer scientists may recognize this as being related to sample complexity. We prefer the term data efficiency because it is more intuitive to map "more efficient" to "higher performance", than "lower complexity" to "higher performance". Further, we do not mean to imply any statistical properties of our data (in contrast to the word "sample").

**Figure 2: The Performance Effect**



methods' data efficiency âĂŞ often also inform views about the limits of machine learning, with many leading researchers for example Yann LeCun, Geoffrey Hinton, Yoshua Bengio (interviewed in [11]), and Gary Marcus [24], commenting that the need for large amounts of data suggest limitations of our current techniques.

Data efficiency is likely to increase. Many recent advances in Machine Learning and Artificial Intelligence have come from deep learning, a technique which enables high performance on challenging tasks such as image recognition at the cost of large compute requirements as well as needing large data sets. At the same time, we know that improvements to data efficiency are possible. For instance, humans can learn simple visual concepts such as novel characters from single instances [19]. In addition, researchers are interested in making more efficient machine learning algorithms, which decrease the amount of data or compute necessary to achieve some level of performance. For example, transfer learning seeks to use pretraining on one dataset to improve performance on another problem or dataset [26][3]. Few and zero-shot learning seek to successfully learn classifications from very few examples [33]. Active learning seeks to increase the value of each data point by getting more informative data points [8]. Data efficiency can also be improved by improving data quality, for instance by developing better sensors.

This paper summarises the empirical literature on how performance scales with data, and introduces a simple model to analyze the implications of improved data efficiency. We then explore implications for the AI industry and society, followed by suggestions for directions for future work.

## 2  PRIOR WORK ON THE DATA EFFICIENCY AND PRODUCTION FUNCTION OF AI

*2.0.1  Compute costs over time.* Recent authors have estimated that the amount of compute used in the largest AI training runs is increasing exponentially, doubling every 3.5 months [2]. Some authors suggest this is reason for skepticism about future AI progress, since requiring exponentially more resources to achieve results will become prohibitively expensive [5, 12]. These high compute requirements, at least in reinforcement learning applications where

---

[3]We include transfer learning as a method for data efficiency because it improves the data efficiency of the target task.

the compute is being used to produce training data, suggest there would be large gains from improvements in data efficiency.

*2.0.2  The relationship between performance and data.* Recent empirical work suggests that there may be a simple relationship between data, model size, and performance.

One investigation found that across different settings, a power law model (where performance is proportional to the amount of data to some power) described the relationship between the amount of data and performance, as long as the model size grew at a rate dictated by a separate power law [18]. If true, these power law relationships allow one to model the data and hardware requirements for a specific performance level, as attempted in recent work [17].

Work by other groups have also found that performance increases with more data, as long as model size is also allowed to increase. While many authors agree on the power law model (for instance [29], [9], [6], and [18]), other authors find that a logarithmic model explains the relationship [34]. Not all work agrees that there is a simple relationship – one experiment which did not increase model size found that performance was only marginally improved by increasing the dataset size [23].

Note that the overall shape of our data to performance schematics agree with the empirical results of a recent paper investigating data efficient supervised learning [16].

*2.0.3  ML performance over time.* While there has been excellent work in tracking changes to ML performance over time (for example [32] and [10]), to our knowledge there are no similar compilations in tracking how data efficiency has changed over time. We believe that this would be a promising direction for future research.

## 3  MODELING DATA EFFICIENCY: ASSUMPTIONS

What are the effects of increasing data efficiency? We construct a simple model of what it means for data efficiency to increase. We model data efficiency as a certain data to performance function, making two assumptions.

Let the *data to performance curve* for an ML system be a function $f$, which takes in a quantity of data $d$, and returns the performance $p$ of the system given that amount of training data. Assume that $f$ is defined in a manner where a larger $p$ is higher performance.

For clarity of presentation, our assumptions tend to be stronger than are necessary for our argument.

*3.0.1  Assumption 1: Monotonic performance increases.* Adding more data will not decrease system performance according to its performance function âĂŞ performance will remain increase. We omit considerations of computational cost from our analysis. Formally:

$$\forall d, d' : d' > d \implies f(d') > f(d)$$

"More data always improves performance"

*3.0.2  Assumption 2: Eventually diminishing marginal returns.* Eventually the system will reach a point (here denoted as $m$) where it sees diminishing marginal returns to performance from data. Intuitively, the first time you see something is more informative than the millionth time. This claim has been theoretically shown for some performance functions [1].

$\exists m$ such that $\forall d > m, d' > m, \Delta > 0$ the following holds

$$d' > d \implies f(d' + \Delta) - f(d') < f(d + \Delta) - f(d)$$

"At some point, more data does not help as much as before"

## 4 MODELS OF DATA EFFICIENCY

How can we model increased data efficiency as a transformation of a data to performance function: $f$ to a new $f_{\text{efficient}}$? We discuss three models of data efficiency to demonstrate different intuitions.

*4.0.1 Data efficiency modeled as adding data.* Data efficiency can be simply modeled as analogous to giving all users more data. This model may be appropriate for understanding the impact of transfer learning, where you use data from one source to improve performance on a variety of tasks. Similarly, this model may be appropriate when data efficiency with respect to real world data comes from using additional simulated data. We do not claim that each data point of simulated or transfer data is as useful as a data point for the target task, but rather that they may be equivalent to some amount of data for the target task.

$$f_{\text{efficient}}(d) = f(d + c) \text{ where } c > 0$$

"Data efficiency is like adding more data"

*4.0.2 Data efficiency by increasing data value.* Another way to model data efficiency is as an increase in the value of data. This could be accomplished through "better data", for instance by collecting data from better placed sensors, which better capture the phenomena one is trying to model. Another plausible path to increased data value is by using active learning, where each data point is chosen to be more informative to the system.

$$f_{\text{efficient}}(d) = f(a * d) \text{ where } a > 1$$

"Data efficiency is like accessing a constant factor more data"

*4.0.3 Data efficiency modeled by function composition.* A general formal expression of data efficiency is as follows, where $g$ is monotonic and continuous:

$$f_{\text{efficient}}(d) = f(g(d)) \text{ where } g(d) > d$$

### 4.1 Two core effects of Data Efficiency

We conceptualize the impacts of data efficiency as composed of two effects – an access effect and a performance effect. The access effect refers to how any given ML capability becomes more accessible to more actors: a given level of performance becomes accessible with less data. The performance effect refers to how for any given amount of data, it becomes possible to achieve higher performance.

*4.1.1 The Access Effect.* As depicted in Figure 1, the access effect refers to the leftward shift of the data to performance curve. This captures most of the straightforward impacts of improved data efficiency, namely decreased data requirements to achieve any given level of performance. This has the effect of enabling new applications in data limited domains and broadening access of existing capabilities to more actors.

CLAIM 1. *Improved data efficiency makes any given level of performance attainable with less data*

Formally, all of our models of data efficiency transform $d$ in some way such that $g(d) > d$ for every $d$. Assuming continuity, this means that there is some $d' < d$ such that $g(d') \geq d$, and therefore $f_{\text{efficient}} = f(g(d')) \geq f(d)$ attaining or exceeding the same level of performance with less data.

*4.1.2 The Performance Effect.* As depicted in Figure 2, the performance effect refers to the upward shift in performance for a given amount of data. It can be muted by the presence of performance ceilings or diminishing marginal returns. This increases the level of performance for many levels of data, assuming access to the same algorithms.

CLAIM 2. *Improved data efficiency increases performance*

Formally, all of our models of data efficiency transform $d$ in some way such that $g(d) > d$ for every $d$. By the monotonicity assumption, this fact implies that $f_{\text{efficient}}(d) = f(g(d)) > f(d)$.

## 5 CONSEQUENCES OF INCREASED DATA EFFICIENCY

### 5.1 Impact on ML-based Capabilities and Applications

*5.1.1 New applications in data limited domains.* One of the clearest implications of data efficiency is the ability to use ML to solve problems in data-limited domains. Data may be limited because there are fundamental limitations – e.g. the data does not exist – or because collecting it is expensive.

One notable example of an area where there is a limited amount of obtainable data is ancient languages, where only so many known text fragments exist. Improvements in data efficiency may improve machine translation applications in these domains.

There are many domains where obtaining new data is costly. This can be data from expensive medical or chemical tests, sensors, real world experiments, or human feedback. As data efficiency improves one would expect ML applications in these domains to become more feasible. This is also the case in domains where data is not presently being collected but potentially could be (for instance, expert judgment for a specific area in a standardized format on difficult questions e.g. medical diagnosis).

A particularly important example for this trend is robotics. Collecting data from real world robots may be expensive because of the costs of robot time (maintenance, damage risk, needing to reset the task, etc.). Relatively recent improvements to data efficiency have made deep reinforcement learning from only real world data possible on simple robotic tasks [15][4].

*5.1.2 New actors access ML capabilities.* Another implication of the access effect is that more actors have access to ML capabilities, since one needs less data in order to achieve a given level of performance. This benefits data-poor actors, suggesting that more companies will be able to deliver a (potentially new) product with a certain level of performance.

[4]This is in contrast to methods which are data-efficient with respect to real-world interaction, but which rely on large amounts of simulated data [25]. Those methods are less computationally efficient, and require upfront investment in simulation capabilities.

An interesting type of data-poor actor is a team within a larger organization, which would like access to more resources (e.g. data, compute, or engineers) to develop some ML capability. As data efficiency improves, less data is required to develop a prototype ML application in order to demonstrate the potential application's value. This may smooth out the adoption of ML by organizations âĂŞ e.g. government agencies âĂŞ who have enough data, but lack organizational buy-in to develop applications.

### 5.1.3 *Misuse potential.*
By increasing the number of actors with access to a given ML capability, the chance increases that an actor with malicious ends will also gain access. Researchers have explored the many ways that ML-based applications could be misused, including for cyberattacks, surveillance, and attempts to affect elections [4].

Deepfakes and synthetic media are a popular example of technologies with a high potential for misuse. In fact, the risk from these systems can be understood as a product of their high data efficiency. Deepfakes that required 1000s of hours of video would be much less disruptive, whereas a recent system was able to base deepfakes on as few as 32 video frames [35]. Data efficiency is a crucial parameter in judging misuse potential.

However, the net effect on misuse from increased data efficiency is complex. Many malicious uses can be defended against. Therefore, as data efficiency increases, we can also expect more actors to gain access to defensive capacities as well as the development of more powerful defenses. The net effect will therefore depend e.g. on the offense-defence balance in the relevant domain [13], the adoption rates of defensive measures, and the extent to which defender or attacker capabilities scale faster as data efficiency increases. Take the example of cybersecurity. Automated vulnerability detection can be used offensively, but it can also be used defensively in order to pre-emptively detect and patch vulnerabilities prior to releasing systems. The factors regarding the net effect of more actors having access to a technology are generally complex, and vary considerably based on the domain [31].

## 5.2 Impacts on Competitive Advantage
Competitive advantage is a core concept in Economics, which refers to factors which allow firms to outperform competitors, charge more, offer better services, etc. [28]. The concentration of AI-related industries and capabilities is an important factor in the governance landscape of AI – there are different policy options in highly concentrated or highly decentralized situations. Further, technological changes can impact competitive advantage [27].

How would improvements to data efficiency affect the competitive advantage of large AI firms? Prima facie, it seems that improvements in data efficiency would lead to a levelling effect, decreasing market concentration. Firstly, the access effect gives more actors access to any level of performance. Secondly, assuming that there are performance bounds to a task, such as for example in facial recognition, the performance effect will eventually diminish the absolute difference in performance between actors. While the above effects may dominate, the overall effect on competitive advantage may in fact be to benefit data-rich actors more than data-poor actors. This is because the value derived from a certain level of performance on a task âĂŞ say revenue created by a recommendation algorithm âĂŞ

differs greatly between actors and often correlates with the actor's size, and because revenue does not scale linearly with performance.

### 5.2.1 *Actors derive different amounts of value from the same level of ML performance.*
Actors derive different amounts of value from the same performance on a task, and so the performance effect benefits some actors more than others. The value an actor derives from a certain capability depends on access to complements to the technology: e.g. having products to sell, customers to sell those products to, and market access. An ML capability which increases user engagement by a fixed 5% will increase total engagement, revenue, and profit more for actors with large user bases.

Furthermore, one can expect being an AI incumbent to correlate with having substantial complements to AI technology. Many contemporary data-rich actors made their investments in ML on the basis of already having more complements to ML performance than other actors. For example, digital advertising may be a domain in which having better ML applications is especially useful, and so companies with large digital advertising revenue may be more inclined to invest in ML. In sum, actors who have more complements to AI technology benefit more from across-the-board increases in performance (such as from improvements to data efficiency) and data rich AI incumbents are likely to have more AI complements, potentially increasing their competitive advantage.

### 5.2.2 *Winner-take-all markets.*
Economists often characterize aspects of the AI industry as a winner-take-all, or winner-take-most market. In such a market, what matters most is whether a firm is first, or not; it doesn't matter much for their market share how good their service is in an absolute sense. Since data efficiency does not alter the rank ordering of actors in the performance of their ML systems, holding datasets and other assets constant, it will have no impact on a pure winner-take-all market assuming that all actors have access to the same algorithms.

### 5.2.3 *Threshold effects.*
There are many tasks where a certain threshold of performance is needed before the service has value. Autonomous vehicles, for example, will only become a viable mass consumer product once they exceed some performance threshold. Once this threshold is reached, a company will see a large spike in the value they can reap from the capability. Improvements in data efficiency may therefore lead to increased concentration if it pushes only a small number of actors above the threshold at which a product becomes viable or a task becomes solvable.

### 5.2.4 *Potential value near performance ceilings.*
Many real-world problems have performance ceilings. For instance, a mean squared error cost function used to measure the performance of an image recognition algorithm has a fundamental performance ceiling at 0 error. Predicting the outcome of a fair coin has an irreducible error of 50%. While it may seem that this would mute competitive advantage stemming from the performance effect, it is not so straightforward. Improved performance can be valuable even close to a performance boundary, and thus a smaller absolute improvement for a data-rich actor may still yield more value than a larger absolute performance improvement for a data-poor actor.

Firstly, having a very high performing system allows one to use its output as the input to other systems. For example, Alipay's Smile to Pay allows users to authenticate payments with their face

and access to their mobile phone, showing high confidence in the accuracy of the underlying facial recognition systems [20].

Secondly, a nominal bound on the performance function does not necessarily imply a practical bound on performance; further, often the marginal benefits of improvement increase as we approach a nominal performance bound. Consider a hypothetical task with the performance function of $P(\text{task is performed correctly}) = p$. This performance function is trivially bounded by 0 and 1. One might think that moving from $P(\text{task is performed correctly}) = 0.99$ to 0.999 is relatively unimportant. However, the expected number of times that we can perform the task before encountering a single error is simply the negative binomial distribution $\mathcal{NB}(p, 1)$, which has the following expected value.

$$\text{Expected number of tasks before error} = \frac{p}{1 - p}$$

Going from $p = 0.99$ to $p = 0.999$ takes the number of times that we can do the task before an error from 99 to 999, almost a 10x improvement. Further, this remains true all the way to the trivial upper bound of 1 – going from 0.999 to 0.9999 is almost another 10x increase in the expected number of task attempts before error.

In many domains (such as capital investments, survival analysis, etc.), the time until error is the important parameter, rather than the probability of failure in any given unit time. This shows that an apparent bound on the performance function is not necessarily a bound on the utility function.

## 5.3 Consequences for Safety and Robustness

*5.3.1 Distributional shift.* In a much more data-efficient world, high performance is attainable with access to much less data. This may mean that deployed systems are more sensitive to distributional shifts, since they may be trained on less representative data and because actors will be more tempted to deploy high-performing systems.

Typically, ML practitioners evaluate their models before deployment using the data that they have access to. If performance is good enough, they may choose to deploy the model. Depending on how the dataset is constructed, larger datasets are more likely to contain representatives of relatively unlikely inputs. This means that needing a large dataset to get the necessary level of performance could give some more robustness to distributional shift, if only because it provides more examples, and a better sense of the rarer parts of the distribution. As such, increased data efficiency may increase issues related to distributional shift.

To counteract this effect, developers ought to think carefully about evaluation and dataset collection âĂŞ if high performance is possible with a smaller dataset, then it is important to proactively include less well-represented inputs in the evaluation, since they will not be sampled as often.

A similar point is that if ML seemingly works on more problems, then this will increase the extent to which such systems are deployed. If deployment of systems based on smaller training datasets happens before researchers address issues with generalization, then this may lead more people to deploy non-robust ML systems.

*5.3.2 Human oversight.* Human oversight and feedback is a particularly costly type of data. Some methods for AI safety are based on the idea of scalable oversight [3], where one either directs human oversight to be more effective [30], itself a form of data efficiency, or trains a model of human approval/disapproval and uses that as a safety component in other parts of the system [22]. In a more data efficient world, these methods are more viable.

## 5.4 Impact of Marginal vs. Total Value of Data

Data efficiency, whether modeled as adding data or as increasing data value, both show a performance effect regarding the *total* performance value of data, but they disagree on the *marginal* performance value of data. The additive model yields a lower marginal performance value of data, because of the diminishing marginal returns assumption. The multiplicative model yields a higher marginal performance value of data because of the chain rule. Thus, the effect of increased data efficiency on marginal value of data is an open question according to these models.

*5.4.1 Data Markets.* Data markets would likely be greatly affected by changes in data efficiency. Firstly, if the marginal value of data goes up, this may increase actors' willingness to buy, sell, or protect their data. Secondly, the collection of new forms of data may become viable if the marginal value of that data surpasses the marginal cost of collecting it. If the marginal value of data is already greater than the cost of collection and then increases, this may instead be realized as increased profit for data-selling firms, rather than increased data collection.

*5.4.2 Data Labeling.* If the marginal value of data increases, then actors will be more willing to pay for data labeling, and there is more potential for higher wage data-labeling jobs, especially where the labeling task is more skill or knowledge intensive. For example, it may become viable to have highly paid professionals such as doctors or lawyers to label data. As ML becomes viable for more tasks, the range of labeling tasks may also expand.

*5.4.3 Surveillance and Privacy.* If the marginal value of data increases this may potentially exacerbate issues in surveillance and privacy. This can potentially be mitigated by the fact that the increased marginal value of data makes it more worthwhile to undergo the expense to collect or process it in a more privacy-preserving manner.

If data efficiency improves in such a way that the marginal value of data decreases, one may expect less surveillance on the margin. However, one might still see a net negative impact on privacy. Firstly, there are likely high fixed costs of building a data collection infrastructure, such that a decrease in the marginal value of data discourages future investments in surveillance, but does not necessarily affect existing data collection infrastructure. Secondly, the performance effect would mean that systems are higher performing overall. As such, the data that the actor already has on its users provides more information about them. An actor would need less data to e.g. predict whether a user is pregnant. As such each piece of data could become arguably more privacy infringing. Further, even if actors are less willing to spend to get new data, as the total value of data increases, they may be more strongly incentivised to hold on to data they have, rather than for instance acquiescing to requests to delete it.

# 6 FUTURE WORK

*6.0.1 Data to performance curves.* The questions of how performance scales with data using current algorithms, how this has changed over time, and how it is likely to change in the future remain fairly unexplored. Researchers with an interest in these issues can consider conducting empirical tests of the relationships between data and performance, as well as investigations into how algorithmic improvements have affected the data to performance curve.

*6.0.2 Performance to utility functions.* A major complicating factor in our analysis is the distinction between performance according to a cost or performance function, and the value provided to a system's owner. What is the owner's utility, for any given level of performance? Research in this area could help yield a more granular and detailed view of the dynamics of AI development, but would require detailed investigation into how exactly ML is used by various actors.

*6.0.3 Production function of AI.* The AI production function plays a crucial role in determining parameters relevant to AI governance and ethics, but there are many questions remaining about how the production function of AI works, what to include, how it changes over time, etc. Research in this area would shed significant light on questions of the requirements of AI research, and provide a better understanding of AI progress.

*6.0.4 Implications of the AI production function.* This paper analyzed the potential impact of changes to data efficiency. What would happen if the performance gains of extra computational resources, the size of model, AI talent, or access to state-of-the-art algorithms were to change? Authors may also be interested in studying the extent to which the current structure of the AI industry depends on the AI production function. For example, many recent state-of-the-art results have come out of private AI labs rather than universities, with many researchers moving from university positions into private industry [14]. To what extent are these changes a result of e.g. large AI firms having access to large amounts of computational resources and data?

## REFERENCES

[1] Shun-Ichi Amari. 1993. A universal theorem on learning curves. *Neural networks* 6, 2 (1993), 161–166.

[2] Dario Amodei and Danny Hernandez. 2018. AI and Compute. *OpenAI Blog, https://openai.com/blog/ai-and-compute/* (2018).

[3] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565* (2016).

[4] Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, et al. 2018. *The malicious use of artificial intelligence: Forecasting, prevention and mitigation.* Technical Report. Technical Report 1802.07228, arXiv.

[5] Ryan Carey. 2018. Interpreting AI Compute Trends. *AI Impacts Blog, https://aiimpacts.org/interpreting-ai-compute-trends/* (2018).

[6] Junghwan Cho, Kyewook Lee, Ellie Shin, Garry Choy, and Synho Do. 2015. How much data is needed to train a medical image deep learning system to achieve necessary high accuracy? *arXiv preprint arXiv:1511.06348* (2015).

[7] Charles W Cobb and Paul H Douglas. 1928. A theory of production. In *Proceedings of the Fortieth Annual Meeting of the American Economic Association*, Vol. 139. 165.

[8] David A Cohn, Zoubin Ghahramani, and Michael I Jordan. 1996. Active learning with statistical models. *Journal of artificial intelligence research* 4 (1996), 129–145.

[9] Arjun D Desai, Garry E Gold, Brian A Hargreaves, and Akshay S Chaudhari. 2019. Technical Considerations for Semantic Segmentation in MRI using Convolutional Neural Networks. *arXiv preprint arXiv:1902.01977* (2019).

[10] Petter Eckersley and Nasser Yomna. 2017. Measuring the progress of AI research.

[11] Martin Ford. 2018. *Architects of Intelligence: The truth about AI from the people building it.* Packt Publishing Ltd, pages 12, 18, 86, 124.

[12] Ben Garfinkel. 2018. Reinterpreting âĂIJAI and ComputeâĂİ. *AI Impacts Blog, https://aiimpacts.org/reinterpreting-ai-and-compute/* (2018).

[13] Ben Garfinkel and Allan Dafoe. 2019. How does the offense-defense balance scale? *Journal of Strategic Studies* 42, 6 (2019), 736–763. https://doi.org/10.1080/01402390.2019.1631810 arXiv:https://doi.org/10.1080/01402390.2019.1631810

[14] Michael Gofman and Zhao Jin. 2019. Artificial Intelligence, Human Capital, and Innovation (August 20, 2019). *ssrn preprint Available at SSRN: https://ssrn.com/abstract=3449440 or http://dx.doi.org/10.2139/ssrn.3449440* (2019).

[15] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. 2018. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905* (2018).

[16] Olivier J Hénaff, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. 2019. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272* (2019).

[17] Joel Hestness, Newsha Ardalani, and Gregory Diamos. 2019. Beyond human-level accuracy: computational challenges in deep learning. In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming*. ACM, 1–14.

[18] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Patwary, Mostofa Ali, Yang Yang, and Yanqi Zhou. 2017. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409* (2017).

[19] Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum. 2011. One shot learning of simple visual concepts. In *Proceedings of the annual meeting of the cognitive science society*, Vol. 33.

[20] Amanda Lee. 2017. Alipay rolls out worldâĂŹs first âĂŸSmile to PayâĂŹ facial recognition system at KFC outlet in Hangzhou. *South China Morning Post* (Sep 2017). https://www.scmp.com/tech/start-ups/article/2109321/alipay-rolls-out-worlds-first-smile-pay-facial-recognition-system-kfc

[21] Kai-Fu Lee. 2018. *AI superpowers: China, Silicon Valley, and the new world order.* Houghton Mifflin Harcourt. 14 pages.

[22] Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. 2018. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871* (2018).

[23] Trond Linjordet and Krisztian Balog. 2019. Impact of Training Dataset Size on Neural Answer Selection Models. In *European Conference on Information Retrieval*. Springer, 828–835.

[24] Gary Marcus. 2018. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631* (2018).

[25] Ilge OpenAI, Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. 2019. Solving Rubik's Cube with a Robot Hand. *arXiv preprint arXiv:1910.07113* (2019).

[26] Sinno Jialin Pan and Qiang Yang. 2009. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22, 10 (2009), 1345–1359.

[27] Michael E Porter. 1985. Technology and competitive advantage. *Journal of business strategy* 5, 3 (1985), 60–78.

[28] Michael E Porter and Victor E Millar. 1985. How information gives you competitive advantage. *Harvard Business Review* (1985).

[29] Vittorio Sala. 2019. Power law scaling of test error versus number of training images for deep convolutional neural networks. In *Multimodal Sensing: Technologies and Applications*, Vol. 11059. International Society for Optics and Photonics, 1105914.

[30] William Saunders, Girish Sastry, Andreas Stuhlmueller, and Owain Evans. 2018. Trial without error: Towards safe reinforcement learning via human intervention. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2067–2069.

[31] Toby Shevlane and Allan Dafoe. 2020. The Offense-Defense Balance of Scientific Knowledge: Does Publishing AI Research Reduce Misuse?. In *Proceedings of the 2020 AAAI/ACM Conference on AI, Ethics, and Society*. ACM.

[32] Yoav Shoham, Raymond Perrault, Erik Brynjolfsson, Jack Clark, James Manyika, Juan Carlos Niebles, Terah Lyons, John Etchemendy, and Z Bauer. 2018. The AI Index 2018 Annual Report. *AI Index Steering Committee, Human-Centered AI Initiative, Stanford University. Available at http://cdn. aiindex. org/2018/AI% 20Index* 202018 (2018).

[33] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. 2013. Zero-shot learning through cross-modal transfer. In *Advances in neural information processing systems*. 935–943.

[34] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. 2017. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*. 843–852.

[35] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. 2019. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models. *arXiv preprint arXiv:1905.08233* (2019).