**CARNEGIE**
ENDOWMENT FOR
INTERNATIONAL PEACE

# What the Machine Learning Value Chain Means for Geopolitics

CHARLOTTE STANTON, VIVIEN LUNG, NANCY ZHANG, MINORI ITO, STEVE WEBER, KATHERINE CHARLET

## INTRODUCTION

Thanks to major improvements in computing power, increasingly sophisticated algorithms, and an unprecedented amount of data, artificial intelligence (AI) has started generating significant economic value. With algorithms that make predictions from large amounts of data, AI contributes, by some estimates, about $2 trillion to today's global economy. It could add as much as $16 trillion by 2030, making it more than 10 percent of gross world product.[1]

AI's outsize contribution to global economic growth has important implications for geopolitics. Around the world, governments are ramping up their investments in AI research and development (R&D), infrastructure, talent, and product development. To date, twenty-four governments have published national AI strategies and their corresponding investments.

So far, China and the United States are outspending everyone else while simultaneously taking steps to protect their investments from foreign competition.

In 2017, China passed legislation requiring foreign companies to store data from Chinese customers within China's borders, effectively hamstringing outsiders from using Chinese data to offer services to non-Chinese parties. For its part, the U.S. Committee on Foreign Investment blocked a Chinese investor from acquiring a leading U.S. producer of semiconductors, which are essential components for computing. While this was officially a national security action, it could also benefit U.S competitiveness by protecting its stake in semiconductor production.[2]

Both data and certain classes of semiconductors are core elements of the AI value chain. Given AI's economic and geopolitical significance, they're also increasingly being considered strategic assets. The extent to which countries can participate in this value chain will determine how they fare in the emerging global economic order and the stability of the broader international system. Indeed, if the gains from AI are distributed in highly variable ways, extreme divergence in national outcomes could drive widespread instability.

So what does the AI value chain look like? And where in the physical world are the key nodes of value creation and control emerging? This article addresses these questions, introducing the idea of a machine learning value chain and offering insights on the geopolitical implications for countries searching for competitive advantage in the age of AI.

## THE MACHINE LEARNING VALUE CHAIN

Machine learning, the science of getting computers to make decisions without being explicitly programmed, is the subfield of AI responsible for the majority of technical advances and economic investment. In recent years, machine learning has led all categories of AI patents (and, in fact, constituted the third-fastest-growing category of all patents granted behind 3D printing and e-cigarettes) and attracted nearly 60 percent of all investment in AI.

A value chain describes the sequence of steps through which companies take raw materials and add value to them, resulting in a finished, commercially viable product. For machine learning, that value chain consists of five stages: data collection, data storage, data preparation, algorithm training, and application development (see figure 1).

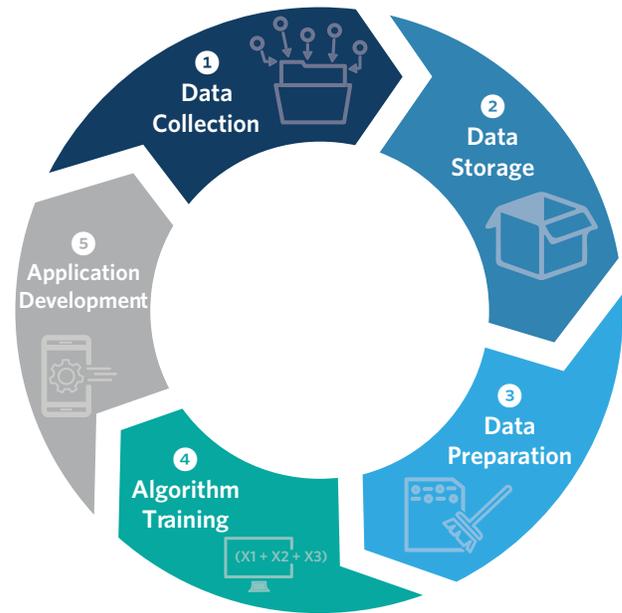*Data collection* involves the gathering of raw data from any number of sources.

*Data storage* involves amassing raw data in data centers.

*Data preparation* involves efforts to clean, convert, format, and label raw data.

*Algorithm training* involves configuring an algorithm to make predictions from data.

*Application development* converts algorithmic predictions into commercially viable products.

FIGURE 1
**The Machine Learning Value Chain**



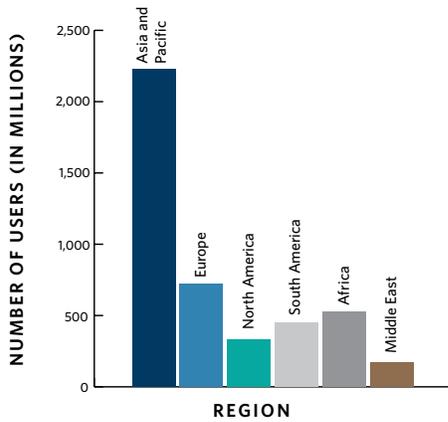## How Is the Machine Learning Value Chain Distributed Globally?

Proxy measures that focus on the key nodes of the machine learning value chain provide a useful, albeit imperfect, means of quantifying and comparing the value chain's distribution across countries and regions. County-level data are used here wherever they are available; regional data are used everywhere else.
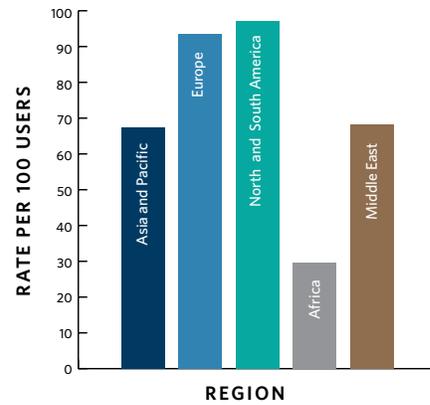
### Data Collection

Raw data are the bedrock of machine learning. Every day, roughly 2.5 quintillion bytes of raw data are collected via myriad devices, from tactile sensors to system logs, that record all manner of digital transactions, such as internet searches, camera images, phone calls, social media posts, and credit card transactions.

FIGURE 2
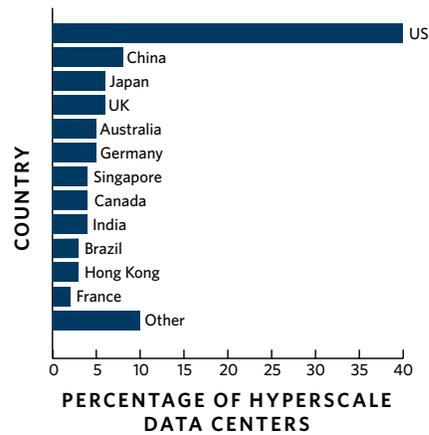**Geographical Distribution of Machine Learning Value Chain Proxies**
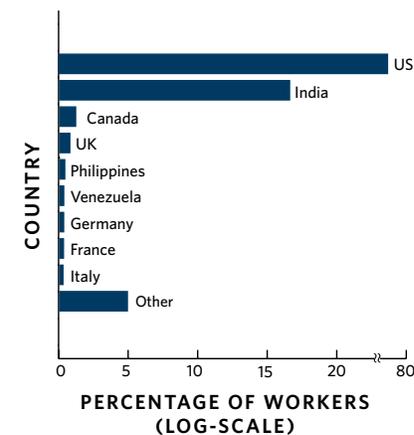
### a. Number of Internet Users
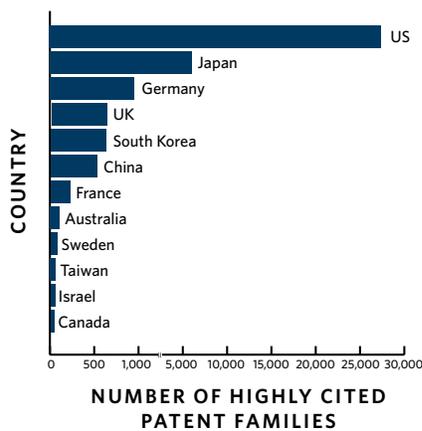


### b. Rate of Mobile Broadband Subscription

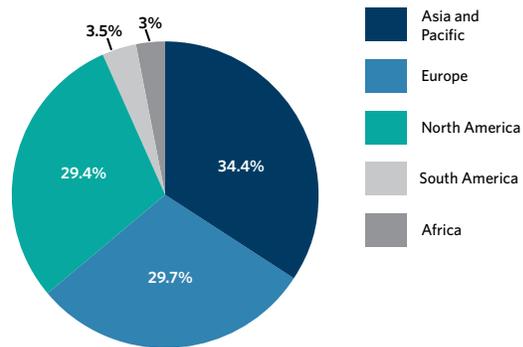

### c. Distribution of Hyperscale Data Centers



### d. Distribution of Mechanical Turk Workers



### e. Number of Highly Cited Patent Families



### f. Distribution of Mobile App Developers



**SOURCES: 2a:** "Internet Users in the World by Regions – 2019 June – Updated," Internet World States, June 30, 2019, https://www.internetworldstats.com/stats.htm. For comparison with other graphs, Latin America was separated into North and South America; Central American and Caribbean countries were added to North America. Asia and Oceania were combined as Asia and Pacific. **2b:** "ITU Releases 2018 Global and Regional ICT Estimates," International Telecommunications Union, December 7, 2018, https://www.itu.int/en/mediacentre/Pages/2018-PR40.aspx. **2c:** "Hyperscale Data Center Count Jumps to 430; Another 132 in the Pipeline," Synergy Research Group, January 9, 2019 https://www.srgresearch.com/articles/hyperscale-data-center-count-jumps-430-mark-us-still-accounts-40. Other includes countries with one or less percent. **2d:** Djellel Difallah, Elena Filatova, and Panos Ipeirotis, "Demographics and Dynamics of Mechanical Turk Workers," Association for Computing Machinery, 2018, http://www.ipeirotis.com/wp-content/uploads/2017/12/wsdmf074-difallahA.pdf. Other includes countries with less than .2 percent. **2e:** "Pure-Play Foundry Market Surges 11% in 2016 to Reach $50 Billion!," IC Insights Research Bulletin, January 12, 2017, http://www.icinsights.com/data/articles/documents/945.pdf. **2f:** Business of Apps, " Number of Mobile App Developers Worldwide in 2014, by Country (in 1,000s)," Statista, January 30, 2016, https://www.statista.com/statistics/629370/share-of-android-app-developers-worldwide-by-country/. Oceania and Asia were combined as Asia and Pacific.

With data collection increasingly taking place through mobile devices, it's no surprise that China and India are two of the most significant data collectors in the world. The number of mobile device users is a useful estimate of how much and from where data is collected worldwide.[3] In absolute terms, China and India have the most mobile device users with 1.22 and 0.44 billion users, respectively.[4] These two countries also contribute the most to Asia's impressive number of internet users, which exceeded 2 billion in 2018—roughly the same number of internet users in the other six regions combined (see figure 2a).

But absolute numbers don't tell the full story. In each of the last five years, mobile broadband subscriptions have grown more than 20 percent, with the highest growth rates in developing countries. The U.S. and European mobile markets are almost saturated at 97.1 and 93.6 subscriptions per 100 users, respectively (see figure 2b). The Middle East and Asia and the Pacific regions are hovering at around 70 subscriptions per 100 users, leaving a sizable margin for further growth. Africa is lagging behind with an average of 29.7 subscriptions per 100 users, but this shortfall presents a significant economic opportunity for African countries seeking to expand their data collection capacity.

## Data Storage

Once data are collected, they are stored and secured in data centers. In the early days of machine learning, companies stored their data in their own brick-and-mortar data centers that contained room-size computer servers. But data storage is increasingly shifting to the cloud, where companies access their data through the internet from cloud service providers that operate hundreds of servers and thousands of virtual machines. The most advanced means of data storage today—hyperscale data centers, or HDCs—operate thousands of servers and millions of virtual machines across multiple locations. And to keep pace with the unprecedented volume of data collected every day, the data storage market is expanding. In 2015, just over 250 HDCs operated worldwide. Experts expect this number to double by 2020.

Growth notwithstanding, the global distribution of HDCs is highly concentrated. The United States dominates the market—40 percent of HDCs are located there or owned by U.S. companies.[5] China, Japan, and the United Kingdom together account for 20 percent of the market, while Australia, Germany, Singapore, Canada, India, and Brazil each represent 3 to 5 percent (see figure 2c).

Brazil's presence in the HDC market, albeit small, is instructive. Because HDCs require significant amounts of energy to operate, countries with low electricity costs may be able to secure a foothold in the HDC market. The cost of electricity in Brazil is just $0.13 per kilowatt hour (kW/h). Developing countries with similar or lower electricity costs—such as Argentina ($0.01 per kW/h), South Africa ($0.09 per kW/h), and Indonesia ($0.10 per kW/h)—may have elements of a comparative advantage at this stage relative to countries in Western Europe, for instance, where electricity costs upwards of $0.20 per kW/h.

## Data Preparation

Unlike data storage, which is capital intensive, data preparation is labor intensive. Raw data are generally messy and unstructured, containing numerous outliers and errors. Highly skilled data engineers and scientists are needed to clean the data and convert them into a usable format (usually a table of seemingly endless rows and columns)—a practice called data pre-processing. Less-skilled human reviewers, or "data labelers," also may be employed to manually classify a subset of the pre-processed data depending on the type of algorithm deployed. For instance, to train a facial recognition algorithm to distinguish between female and male faces, data labelers will first manually classify a subset of faces as either "female" or "male," so that the algorithm can identify the salient features of each gender and predict the gender of the remaining (unlabeled) facial data.

Since the majority of highly skilled workers are data engineers and scientists with science, technology, engineering, and mathematics (STEM) backgrounds, the number of STEM graduates is an imperfect but

useful proxy measure.[6] Amazon Mechanical Turk (MTurk) workers approximate less skilled workers. MTurk is an online labor market where workers compete to execute tasks for pay. It's the largest source of workers performing data labeling tasks today. In 2016, China and India led the supply of STEM graduates with 4.7 and 2.6 million, respectively. The United States and Russia produced just over half a million graduates each, followed by Iran, Indonesia, and Japan. In 2018, 75 percent of MTurk workers originated in the United States, followed by India (16 percent), Canada (1.1 percent), the United Kingdom (0.7 percent), the Philippines (0.35 percent), and Venezuela (0.28 percent) (see figure 2d).

Demand for highly skilled data scientists and less-skilled data labelers will increase. The market for data labeling may provide an especially low barrier to entry for countries where English is not the first language, because many of the common tasks involved in data labeling, like image classification, merely require digital literacy rather than English literacy.

## Algorithm Training

Once the data are prepared, companies can start training their algorithms to make predictions on new data. With dozens of algorithm types and nearly infinite configurations to choose from, they require data scientists to develop alternative models and compare their prediction performance, making highly skilled labor a fundamental input to algorithm development and training. Another fundamental input is computer hardware: the most sophisticated algorithms that make predictions from big datasets involve trillions of calculations and machines made of bespoke semiconductors that can perform such computations quickly and efficiently.

Using STEM graduates as a proxy for the supply of data scientists, China and India again stand apart. China also owns the most individual machine learning patents, including the most patents for deep learning, which is the fastest-growing subfield of machine learning.[7] China

lags far behind the United States and other countries, however, in the number of *highly cited* patents. The number of times a patent is cited is an indicator of "patent quality" and its technical relevance and overall value. The United States owns the most highly cited patent families (28,031), followed distantly by Japan (6,221), Germany (931), the United Kingdom (778), South Korea (758), and China (691) (see figure 2e).

The United States also sells the most machine learning–specific semiconductors, namely the field-programmable gate arrays (FPGAs) and graphics processing units (GPUs) that are essential for algorithm training. U.S. companies sold 96 percent of all FPGAs in 2016 and 100 percent of all GPUs in 2017 and 2018. That said, most manufacturing of the actual hardware takes place outside of the United States.[8] In 2016, about 80 percent of FPGAs and GPUs were manufactured in Taiwan, South Korea, and China (with Taiwan's output accounting for three times the combined shares of South Korea and China).

Beyond the United States' almost unchallenged leadership in some of the key components that power the algorithm-training node of the AI economy, a handful of other countries—namely Japan, South Korea, Germany, and the United Kingdom—are shoring up their positions by boosting investment in R&D and STEM education.

## Application Development

Algorithmic predictions are translated into actionable insights via applications—software programs developed by software engineers, or app developers. The potential applications of machine learning algorithms are nearly endless, ranging from personal assistants (such as Amazon's Alexa) to product recommendations (for example, YouTube's video recommender) to autonomous weapons (like military drones).

The migration in internet usage from desktop to mobile makes the distribution of the mobile app market and mobile app developers a useful, albeit partial, indicator

of which countries are converting the predictive power of machine learning into commercial profit. Country engagement in the app market is especially important because apps, mobile or otherwise, generate additional user data that will continually expand the database on which machine learning algorithms can train, leading to more accurate predictions and more valuable products. Looking regionally, Asia and the Pacific (33 percent), Europe (30 percent), and North America (29 percent) hold roughly similar shares of the mobile app developer market (see figure 2f).

## WHAT ARE THE GEOPOLITICAL IMPLICATIONS?

While no two countries look alike in their machine learning investments, most fall into three categories: fast movers, moderate movers, and slow starters. Fast movers, namely China and the United States, are heavily investing across most if not all nodes of the machine learning value chain—effectively ensuring that both economies medal in the so-called race to win AI. Moderate movers, by contrast, are concentrating their investments in particular nodes of the value chain. Germany, Japan, and Taiwan, for instance, are heavily investing in the physical capital required for data storage and algorithm training (like HDCs and supercomputers). Australia and South Korea are investing in the requisite intellectual capital (for example, R&D and STEM graduates).

Slow starters have yet to invest significantly in any stage of the machine learning value chain. Most developing countries are slow starters. Notable exceptions include Brazil, which entered the HDC market by capitalizing on its cheap cost of energy, and Kenya, whose relatively high internet penetration rate (83 percent) enables significant data collection. But slow starters need not be left out: the most immediate opportunities for these countries are in data collection and data labeling. The market for data labeling, and specifically image labeling, has an especially low barrier to entry for developing countries where English is not the first language since classifying images doesn't require English literacy.

This offers hope that there are ways to drive a more globally inclusive machine learning economy. But dedicated attention is necessary. Governments should study where in the machine learning value chain they may have a comparative advantage. Development agencies can develop tools to assist such analysis, and they can invest in strategies to help slow starters better position themselves to participate in the machine learning value chain. Researchers can study which stage of the machine learning value chain creates the most value in order to determine whether any given country's best opportunities for investment are in data collection, data storage, or another stage. Such research must account for country-specific factors and the magnitude and relative shelf life of the value created at each stage. For instance, the profitability of operating an HDC directly depends on local energy prices and data localization laws, among other things. Likewise, a hardware component may become obsolete after three years, while a well-trained data scientist may yield value for decades, possibly becoming more valuable over time with the benefit of experience. Unpacking the profitability of each stage for each country will not be easy, but it can help countries enhance their competitive advantage in the age of AI.

Finally, it's important to note that the geographical concentration of machine learning value among the fast movers and even moderate movers will have first- and second-order impacts on the distribution of wealth and power within and between countries. Concentrating talent and wealth in certain countries will likely exacerbate economic inequality between countries. A similar geographical concentration of talent and wealth in certain cities could also impact land and housing prices, causing demographic shifts. Taken together, these first- and second-order impacts suggest increasing inequality between *and* within countries—a different trend than what took place over the last quarter century, when globalized industrial manufacturing increased inequality within countries but decreased inequality across countries. Policymakers must prepare now for the geopolitical consequences of countries' varied capabilities and investments in the machine learning value chain.

## ABOUT THE AUTHORS

**Charlotte Stanton** is the inaugural director of the Silicon Valley office of the Carnegie Endowment for International Peace as well as a fellow in Carnegie's Technology and International Affairs Program.

**Vivien Lung** is a senior policy analyst at Google's Trust and Safety division. Previously, she was a research assistant at Stanford's Center for Security and International Cooperation and a consultant at Deloitte's Global Transfer Pricing practice.

**Nancy (Hanzhuo) Zhang** is an analyst for Cornerstone Research. Previously, she worked for the World Bank's Development Impact and Evaluation Unit, the Economist Group, and the Bill and Melinda Gates Foundation.

**Minori Ito** is a diplomat at the Ministry of Foreign Affairs of Japan.

**Steve Weber** is a professor of political science and information at the University of California, Berkeley and the faculty director for the Berkeley Center for Long-Term Cybersecurity. He specializes in international relations and international political economy with expertise in international and national security, the impact of technology, and the political economy of knowledge-intensive industries particularly software and pharmaceuticals.

**Katherine Charlet** is the inaugural director of Carnegie's Technology and International Affairs Program.

*The authors are grateful to Frances Reuland for outstanding research assistance.*

## NOTES

1   Jacques Bughin, Jeongmin Seong, James Manyika, Michael Chui, and Raoul Joshi, *Notes From the AI Frontier: Modeling the Impact of AI on the World Economy* (New York: McKinsey Global Institute, 2018), 13; and "GDP Long-Term Forecast," Organization for Economic Cooperation and Development, 2018, https://data.oecd.org/gdp/gdp-long-term-forecast.htm.

2   Tim Hwang, "Computational Power and the Social Impact of Artificial Intelligence," *SSRN Electronic Journal* (March 2018): 27; and Michael Brown and Pavneet Singh, "China's Technology Transfer Strategy: How Chinese Investments in Emerging Technology Enable A Strategic Competitor to Access the Crown Jewels of U.S. Innovation," Defense Innovation Unit Experimental (DIUx), January 2018, 8–15.

3   In 2016, mobile devices (for example, smartphones and tablets) accessed the most web pages worldwide, surpassing desktops for the first time. "Mobile and Tablet Internet Usage Exceeds Desktop for First Time Worldwide," GlobalStats, November 1, 2016, http://gs.statcounter.com/press/mobile-and-tablet-internet-usage-exceeds-desktop-for-first-time-worldwide.

4   Smartphone users data from: "Top 50 Countries/Markets by Smartphone Users and Penetration," NewZoo, September 2018, https://newzoo.com/insights/rankings/top-50-countries-by-smartphone-penetration-and-users/. Tablet users per country data from: "Number of Tablet Users in China From 2013 to 2018 (in Millions)," Statista, 2019, https://www.statista.com/statistics/377971/china-tablet-users-forecast/; "Number of Tablet Users in India From 2013 to 2018 (in Millions)," Statista, 2019, https://www.statista.com/statistics/413325/tablet-users-number-india/.

5   When you consider that data hosted in a U.S.-owned HDC that is physically located outside the United States, there is ambiguity about who can access the value from the HDC and its data. This ambiguity depends on the legal environment of the respective countries.

6   STEM graduates are an imperfect estimate of data scientists because the correlation between the two may differ across countries and regions. The presumption here is that over time the proportion of STEM graduates that are data scientists will be relatively similar across geographies.

7   Machine learning patents comprise 40 percent of all AI-related patents according to World Intellectual Property Organization, *WIPO Technology Trends 2019: Artificial Intelligence* (Geneva: World Intellectual Property Organization, 2019), 94.

8   Hwang, "Computational Power and the Social Impact of Artificial Intelligence," 21.

**CARNEGIE**
ENDOWMENT FOR
INTERNATIONAL PEACE