# How AI Is Shaping the Future of Data Platforms & Infrastructure in 2024

Tejas Dessai
tdessai@globalxetfs.com
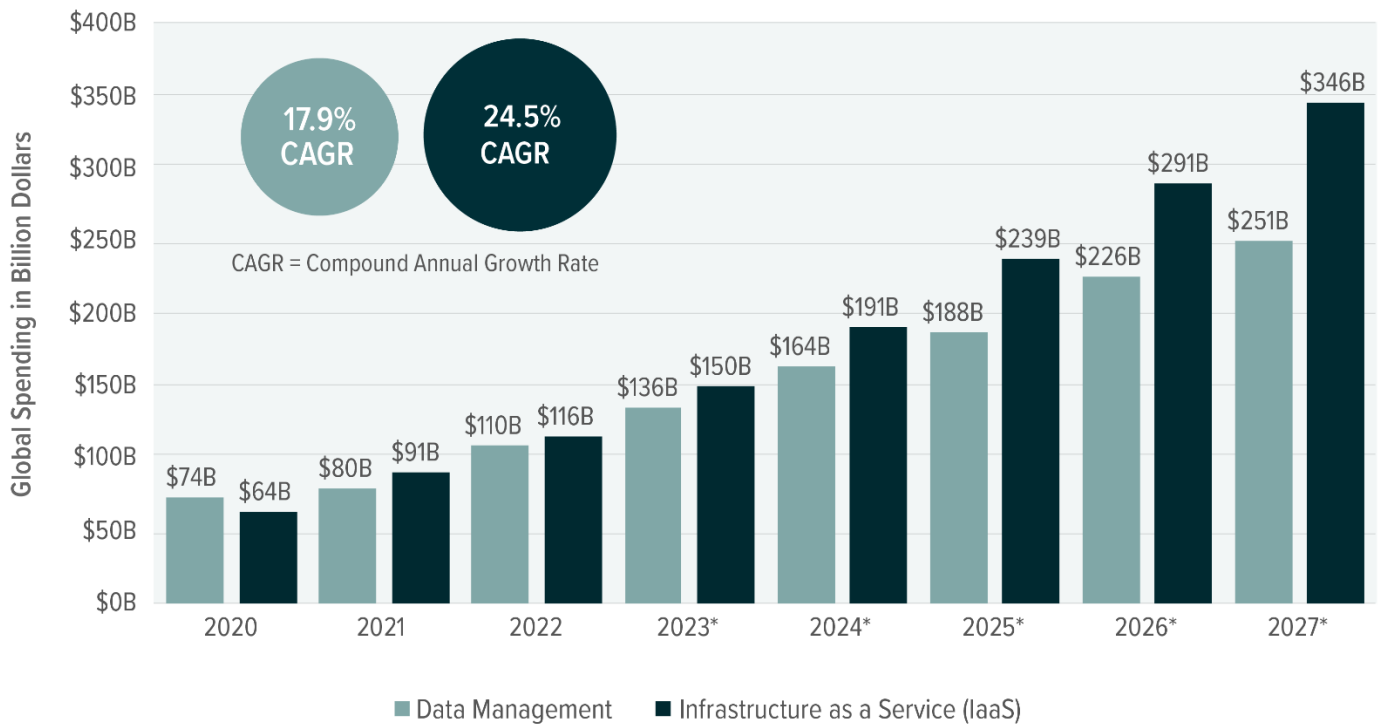
Date: February 23, 2024
Topic: Thematic, Disruptive Technology

*Data is AI's fuel. And the more data AI gets, the more of a data flywheel AI creates. In our view, the value chain involved in the processing and handling of all that data is not to be overlooked. Growing business demand for generative AI solutions creates demand for specialized software and hardware to support the capture, storage, and processing of massive amounts of data. Additionally, as generative AI models are trained on internal and private data, needs grow for secure and efficient data access solutions, governance solutions, and ancillary cloud computing infrastructure.*

*In this piece, we highlight how generative AI creates significant growth opportunities for the companies in the cloud infrastructure and data management markets.*

## AI LIKELY TO ACCELERATE GROWTH FOR CLOUD INFRASTRUCTURE AND DATA MANAGEMENT PLATFORMS

Sources: Global X ETFs with information derived from: Gartner. (2023, April 19). Gartner Forecasts Worldwide Public Cloud End-user Spending to Reach Nearly $600 Billion in 2023 and Research and Markets. (2023, March 24). Global Data Management Software Market: Analysis by Type, by Organization Size, by Deployment Type, by Application, by Region Size & Forecast with Impact Analysis of Covid19 and Forecast up to 2028.



*Indicates forecast

## KEY TAKEAWAYS

— Growing demand for generative AI solutions necessitates substantial investments by enterprises in data management, storage, and infrastructure solutions, in addition to hardware.

— Synthetic data and data generated by AI agents are likely to grow significantly alongside real-world data.

— We expect companies that provide data management and cloud infrastructure solutions to grow in stature as ways for investors to gain compelling exposure to the Artificial Intelligence theme.

## High-Quality Data for AI Wanted

Generative AI models are trained on large data sets of structured and unstructured information, which forms the basis of their ability to reason and respond to questions. For perspective, Open AI's GPT-3 was trained on a 45-terabyte (TB) dataset that combined several data sources from the open web, including Common Crawl (60%), WebText2 (22%), Books1 (8%), Books2 (8%), and Wikipedia (3%).[1]

Growing AI adoption and increased enterprise investments should stimulate additional training efforts by developers and technology companies in foundational models. Model performance scales with the size and quality of input training data, so as AI deployment grows, so too will demand for high-quality data feeds.[2] Real-world data can only take these models so far, though. Because real-world data can be limited and scarce, alternative sources of information, such as private data and synthetic data, are likely to become more prominent in AI's development.[3]

Private data refers to proprietary enterprise data and information that can be used to train models for specific internal use cases. Due to the narrow focus of private data assets, these models can be more efficient and useful than off-the-shelf models. For example, Bloomberg's 50-billion parameter large language model, trained on financial data sets, outperforms similarly sized open models on financial natural language processing (NLP) tasks.[4] We expect similar examples to emerge in healthcare, logistics, manufacturing, defense technology, cybersecurity, and other industries.

Synthetic data is information manufactured artificially through algorithms rather than from real-world events. This data can be designed to appear nearly perfect. In most cases, synthetic data is utilized to fill-in real-world datasets, replacing historical data that is no longer relevant or sometimes inaccessible. It is also cost efficient, entails no privacy concerns, and is perfectly annotated. By 2024, nearly 60% of all the data used to developing AI and analytics is expected to be synthetic.[5]

## More Investments in Hardware and Software Needed

The training and development of new models, the integration of real-world data for AI based reasoning and assessment, and the use of private data and synthetic data require comprehensive investments in data infrastructure.

Heightened data usage and processing as AI's footprint expands primarily creates a need for data center storage and memory. With data center vacancies already at historic lows, the construction of new data center capacity to support demand for AI-focused workloads is witnessing an uptick.[6] Billions of dollars in capital investments from cloud hyperscalers like Microsoft, Amazon, Google, data center giants like Equinix, as well as private equity companies are being catalyzed towards purpose-built AI datacenters.[7]

Storage for AI processing must also accommodate low-latency and real-time access to data, so storage solution vendors are expanding their portfolios. For example, Seagate Technologies recently launched a Seagate SkyHawk AI 24TB storage platform, which is designed for image and video storage at the edge of AI applications.[8] Alongside AI server installations is demand for high-bandwidth memory (HBM), which is designed for low-power consumption and ultra-wide communication lanes.[9] Companies like Samsung and HK Hynix dominate this market.[10]

On the data software side, enterprises will likely have to invest in platform solutions and build pipelines and other necessary infrastructure that enables models to interact with users and systems. Traditional enterprises may also have to clean, process, and reshape their existing data assets to make it ready for AI training and inferencing. Vector database services, like those provided by MongoDB and increasingly by incumbents like Oracle, are specialized storage systems optimized for storing and searching through vector data. Generative AI models use vector search technology to parse through extensive information repositories by identifying relevant data points based on their vectors, which are numerical representations in a multi-dimensional content.

AI agents interacting with each other is likely to become common, requiring unique systems integrations that give data solutions providers room to innovate. Also, capturing error-free data from real-world AI applications, including Internet of Things devices, robotics, drones, and other mechanical systems, requires investments in better sensing setups, edge processors, and edge networks, as well as associated software and data platforms.

These investments extend to broader cloud strategies as enterprises configure their IT operations to be compatible with public and private cloud systems. Existing digital transformation agendas should accelerate, to the benefit of cloud-based computing infrastructure services vendors such as Microsoft, Amazon Web Services, and Google Cloud as well as large platform vendors like ServiceNow.[11,12,13,14]

## Conclusion: Infrastructure Essential to AI's Data Flywheel

Generative AI models need access to high-quality, real-time, and proprietary data to fulfill their vast potential. The public and private sector buildout of the data management and infrastructure platforms needed to make that happen positions the companies selling the cloud infrastructure, storage hardware, databases, data warehouses, data streaming tools, and more to benefit. And within those wares, we believe there are attractive opportunities for investors to capture AI's growth.

**Related ETFs**

AIQ - Artificial Intelligence & Technology ETF

CLOU – Cloud Computing ETF

VPN – Data Center REITs & Digital Infrastructure ETF

*Click the fund name above to view current performance and holdings. Holdings are subject to change. Current and future holdings are subject to risk.*

### Footnotes

1.  Dennis Layton, Medium. (2023, January). ChatGPT — Show me the Data Sources.
2.  AI Multiple Research. (2023, Jan 03). Data Quality in AI: Challenges, Importance & Best Practices in '24.
3.  MIT News. (2022, November 3). In machine learning, synthetic data can offer real performance improvements.
4.  Bloomberg. (2023, March 30). Introducing BloombergGPT, Bloomberg's 50-billion parameter large language model, purpose-built from scratch for finance.
5.  MIT Sloan. (2023, Jan 23). What is synthetic data — and how can it help you competitively?
6.  CBRE. (2023, September 26). North American Data Center Trends H1 2023.
7.  Wall Street Journal. (2023, August 23). AI-Ready Data Centers Are Poised for Fast Growth
8.  Seagate. (2023, December 12). Seagate SkyHawk AI 24TB Elevates Edge Security Capacity and Performance.
9.  WCCFTech. (2024, Feb 12). HBM Memory Prices Reaches An All-Time High, 500% Surge Amid Huge Demand
10. TrendForce. (2023, April 18). HBM Supply Leader SK Hynix's Market Share to Exceed 50% in 2023 Due to Demand for AI Servers, Says TrendForce.
11. Microsoft Earnings Release. (2024, January 24). Microsoft Cloud Strength Drives Second Quarter Results.
12. Amazon Earnings Release. (2024, February 1). Amazon.com Announces Fourth Quarter Results.
13. Alphabet Investor Relations. (2024, January 30). Alphabet Announces Fourth Quarter and Fiscal Year 2023 Results.
14. ServiceNow Investor Relations. (2024, January 24). ServiceNow Reports Fourth Quarter and Full-Year 2023 Financial Results.