

AI is transforming how business processes are carried out in the digital era. But while AI's power and promise are exciting, it is not easy to deploy AI models and workloads. To eliminate bottlenecks in faster model iteration and value realization, organizations need to architect an AI data pipeline.

Infrastructure Considerations for AI Data Pipelines

July 2018

Written by: Ritu Jyoti, Program Vice President

Introduction

Artificial intelligence (AI), machine learning (ML), and continual deep learning (DL) are the latest technologies poised to transform how consumers and enterprises work and learn. While data is at the core of the new digital economy, how you sense the environment and manage the data from edge to core to cloud, analyze data in near real time, learn from data, and act on data to affect outcomes is also important. What differentiates winning organizations is how they leverage that data to deliver meaningful, value-added predictions and actions for improving industrial processes, experiential engagement, healthcare, and other kinds of enterprise decision making.

AI is transforming how business processes are carried out in the digital era. But while AI's power and promise are exciting, it is not easy to deploy AI models and workloads. Most organizations are struggling through proofs of concept (POCs), and only a few have made it to full production. Building, testing, optimizing, training, inferencing, and maintaining the accuracy of models are integral to AI workflow. These neural-network models are hard to build. ML and DL algorithms need huge quantities of training data, and AI effectiveness depends heavily on high-quality, diverse, and dynamic data inputs. Data management of these data sets is complex and challenging.

Reducing the time to insight from months to days speeds up the process of getting value out of a model, which can result in competitive advantage. To eliminate bottlenecks in faster model iteration and value realization, organizations need an AI data pipeline — a set of data processing elements connected in a series, where the output of one element is the input of the next element. Such a pipeline supports smooth data flow from data ingestion to transformation, from exploration to training to inferencing, and then into long-term archival storage for future use. The decision to run an AI pipeline on a public cloud versus on-premises or in a distributed setup from edge to core to cloud is typically driven by data gravity, where the data currently exists or is likely to be stored. It's also important when inferencing needs to occur to support business service-level agreements (SLAs), and where training of the model needs to run to support organization and implementation, as well as specific security and compliance requirements.

AT A GLANCE

KEY STATS

IDC predicts the following:

By 2019, 40% of DX initiatives will use AI services.

By 2021, 75% of commercial enterprise apps will use AI.

AI Deployments: Opportunities

AI has been around for decades, but AI technologies are making headway now because of data proliferation, the growing sophistication of ML and DL algorithms in spotting patterns in larger amounts of data, access to cloud computing, and the availability of accelerators. Companies of all sizes and across various industries are making significant inroads in using AI to solve real-life problems, especially in the enterprise.

Payment services such as PayPal are using GPU-accelerated DL for fraud detection. Consulting firm Accenture's R&D arm and other businesses are using it to detect internet security threats. Drive.ai is using DL algorithms and preparing to offer a self-driving car service for public use in July 2018.

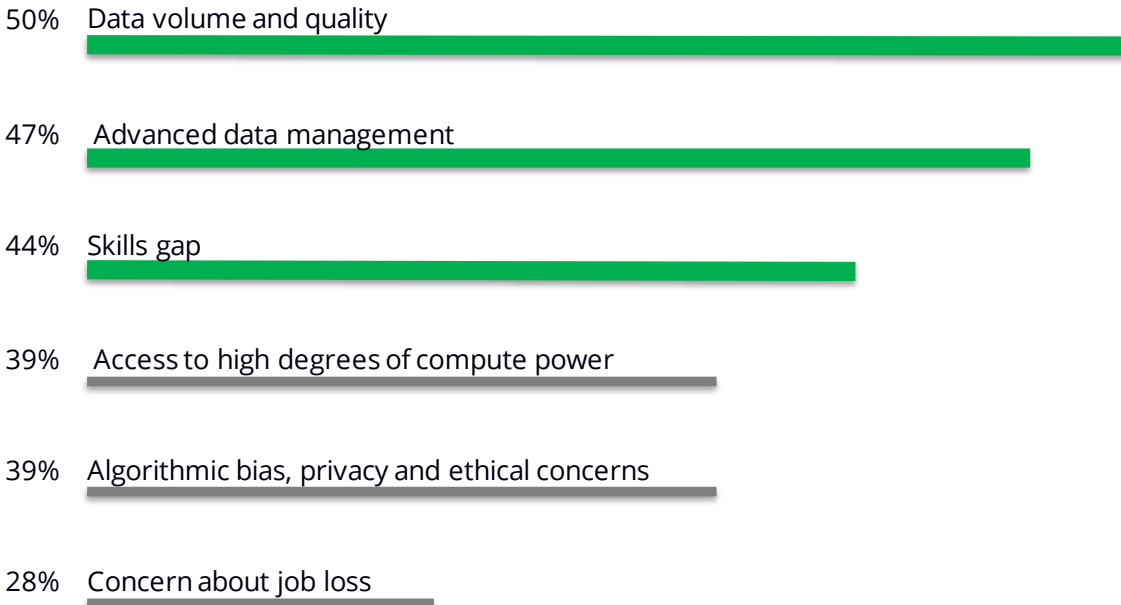
Healthcare is a key industry for the application of AI and DL, particularly to accelerate drug discoveries. New drugs typically are brought to market in 12–14 years. Benevolent AI is using GPU DL to bring new therapies to market quickly and affordably by automating the process of identifying patterns within large amounts of research data, enabling scientists to form hypotheses and draw conclusions quicker. For example, the NVIDIA DGX-1 AI supercomputer was used to identify two potential drug targets for Alzheimer's in less than one month.

DL is also proving to be a game changer for the retail industry. An interesting example is the use of DL by Focal Systems to automate real-time out-of-stock detection. When a customer with a Focal Systems device on his/her cart walks by an out-of-stock shelf, an alert is displayed on the store's dashboard within seconds to trigger replacement action and improve the product's on-shelf availability.

Challenges to AI Deployment

In January 2018, IDC surveyed 405 IT and data professionals in the United States and Canada who had successfully completed an AI project, had budget control or influence, and were responsible for evaluating or architecting a platform to run AI workloads. The survey sought to determine how organizations use and manage AI-enabled technologies and to identify the infrastructure used for running cognitive/ML/AI workloads, the deployment location of the technology, and the associated challenges and needs. Respondents identified dealing with massive data volumes and associated quality and management issues as their key AI deployment challenges, as seen in Figure 1.

AI effectiveness depends heavily on high-quality, diverse, dynamic, and distributed data sets. Advanced data management is one of the top challenges for successful AI deployments.

FIGURE 1: CHALLENGES DEPLOYING AI WORKLOADS

n = 405

Source: IDC's Cognitive, ML, and AI Workloads Infrastructure Market Survey, January 2018

AI effectiveness depends heavily on high-quality, diverse, and dynamic data inputs. Historically, data analytics centered around large files, sequential access, and batched data. Modern data sources and characteristics are different. Today, data consists of small to large files and structured, semistructured, and unstructured content. Data access varies from random to sequential. By 2025, more than a quarter of the global data set will be real time in nature, and real-time IoT data will make up more than 95% of it. In addition, data is increasingly distributed across on-premises, colocation, and public cloud environments. Data science teams often struggle to move or copy data across multiple data repositories. For example, data sets gathered at edge locations by automotive companies or by global retailers from point-of-sale devices are growing exponentially and driving IT technology to its limits. Meeting production-quality service levels for performance and protection across large and dynamic data sets has been challenging.

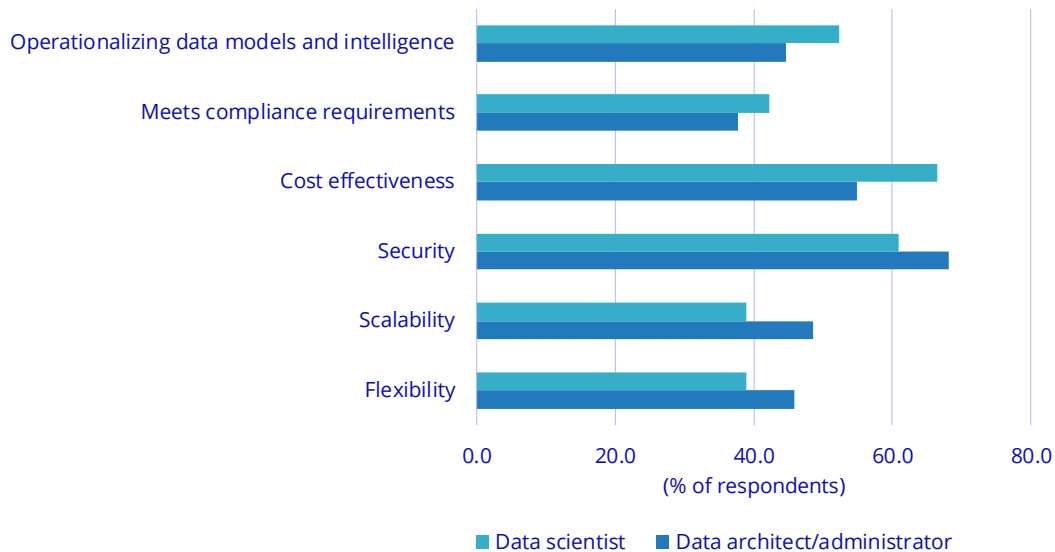
Poor data quality has a direct correlation to biased and inaccurate model buildout. Ensuring data quality with large volumes of dynamic, diverse, and distributed data sets is challenging because it is hard for developers to know, predict, and code for all the appropriate checks and validations.

Supporting the growing segment of AI-dependent digital transformation (DX) initiatives requires talent — meaning both AI engineers and data scientists. However, IT organizations face a shortage of these professionals as well as a skills gap. For example, there is a high learning curve in building/optimizing and training models, a skill set most data scientists don't possess.

AI Deployments: Infrastructure Design Considerations

When data scientists, data architects, and data administrators were asked by IDC for their top decision criteria in choosing an AI solution, they identified security, cost effectiveness, scalability, flexibility, and operationalization of data models/intelligence, as shown in Figure 2.

FIGURE 2: DECISION CRITERIA FOR AI INFRASTRUCTURE/SOLUTIONS



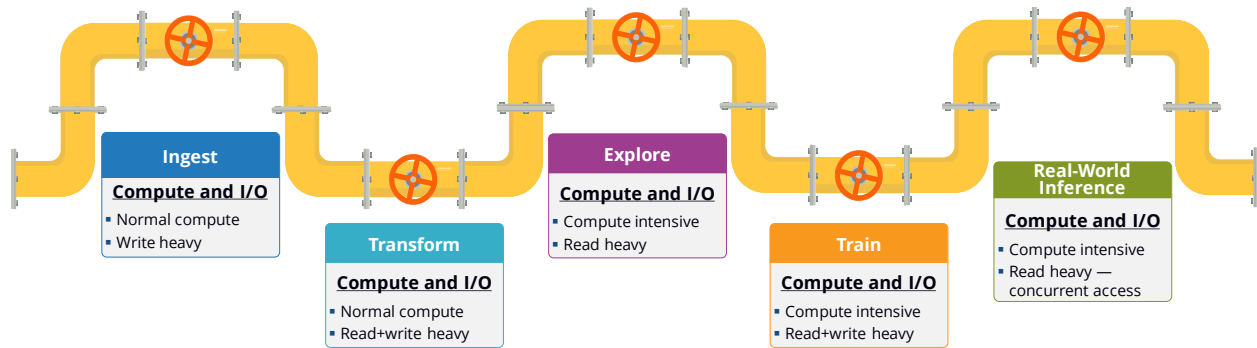
Source: IDC, 2018

The following sections discuss the three primary infrastructure design considerations that can support the needs of data scientists, data architects, and data administrators for an AI solution.

Architected for Performance and Scale

Examining the data pipeline for AI workflows shown in Figure 3, we see that the application profiles, compute, and I/O profiles change from ingestion to real-world inferencing. ML and DL need vast quantities of training data. While the training phase requires large data stores, inferencing has less need for them. Both training and inferencing are compute intensive and require high performance for fast execution. The inference models are often stored in a DevOps-style repository where they benefit from ultra-low-latency access. AI applications push the limits on thousands of GPU cores or thousands of CPU servers. AI and DL require a new class of accelerated infrastructure primarily based on GPUs. For the linear math computations needed for training neural network models, a single system configured with GPUs is significantly more powerful than a cluster of non-accelerated systems.

However, not all AI deployments are the same. Organizations should explore heterogeneous processing architectures (e.g., GPUs, FPGAs, ASICs, or Manycore processors) based on the performance, operating environment, required skill sets, costs, and energy demand for their AI deployments.

FIGURE 3: DATA PIPELINE FOR AI WORKFLOWS

Source: IDC, 2018

Regardless of the processing architecture selected, parallel compute demands parallel storage. Enterprises cannot support a cutting-edge tool such as AI on legacy infrastructure challenged to meet the required needs for scale, elasticity, compute power, performance, and data management. Today, organizations are using different infrastructure solutions and approaches to support the data pipeline for AI, an approach that generally leads to data silos. Some create duplicate copies of the data for the pipeline to avoid disturbing the stable applications. Instead, organizations need to adopt infrastructure that supports varied data formats and access, processes and analyzes large volumes of data, possesses the speed to support faster compute calculations and decision making, manages risks, and reduces the overall costs of AI deployments.

Architected for Distributed Deployment — from Edge to Core to Cloud

Data ingestion usually occurs at the edge, such as capturing data streaming from cars or point-of-sale devices. Depending on the use case, IT infrastructure might be needed at or near the ingestion point. For instance, a retailer might need a small footprint in each store, consolidating data from multiple devices.

Data transformation or preprocessing of the data before training takes place in a data lake, in the cloud, or on-premises, driving the need for an appropriate file or object store.

For the critical *training phase* of DL, data is typically copied from the data lake into the training cluster at regular intervals. Servers used in this phase often employ GPUs or custom silicon to parallelize operations. To support the data flow needed, raw I/O bandwidth is crucial.

The resulting model is pushed out for testing and then moved to production for *inferencing*. Based on the use case, the model might be deployed back to edge operations. Real-world results of the model are monitored, and feedback in the form of new data flows back into the data lake, along with new data to iterate on the process.

Cold data from past iterations may be saved indefinitely. Many AI teams archive cold data to object storage in either a private cloud or a public cloud.

Hence an infrastructure architecture that spans support from edge to core to cloud is typical.

Architected for End-to-End Data Flow and Management

Whether organizations execute their AI workflow on-premises or in the cloud, operational bottlenecks at the edge, core, or the cloud can extend the time needed to complete each training cycle. Extra time reduces the pipeline's productivity and delays realization of the value from the model. It is crucial that data flow smoothly through the entire pipeline. To achieve this goal requires edge, on-premises, and cloud data management.

Edge Data Management

The amount of data generated by smart edge devices and numerous ingestion points can overwhelm compute, storage, and networks at the edge. This data can create congestion as it moves into the datacenter or the cloud. Applying edge-level analytics makes it possible to process and selectively pass on data during ingest. To do so requires infrastructure at the edge with high-performance, ultra-low-latency storage. Delivering data at high bandwidth from the edge also requires a smart data mover that coalesces data, applies data transformations to reduce the data footprint, and applies network transformations to move only changed blocks. This approach can dramatically accelerate data movement and reduce bandwidth requirements.

On-Premises Data Management

As data flows in from the edge, it collects in a data lake for *data transformation and exploration*. An improperly implemented data lake becomes a bottleneck as the amount of data grows. It is important to select a data lake solution that supports scale, performance, and cost-efficiency requirements. It also should be optimized for small and large files as well as sequential and random I/O.

For *training*, it is important that data be streamed quickly and efficiently from the data lake into the training cluster. The data lake should be designed to support high-bandwidth output. The training phase is read-write heavy, and the cluster should be architected to deliver high bandwidth and low latencies to support many I/O streams in parallel.

For the *inferencing* stage, it is important to deploy storage systems that support extremely low latency.

A single storage architecture that addresses all the performance needs for the core of the deep learning pipeline has immediate advantages. However, in some cases, big data pipelines might already be in place and an organization might want to implement just the new elements needed for DL as a separate project and copy data from phase to phase. As data continues to grow, however, unifying the pipeline will become necessary.

Cloud Data Management

To deploy AI deep learning in the cloud for agility and ease of consumption, organizations need to ensure that the solution can deliver the necessary performance for data ingestion and provide I/O parallelism for the training cluster and ultralow latency for inferencing while ensuring data sovereignty.

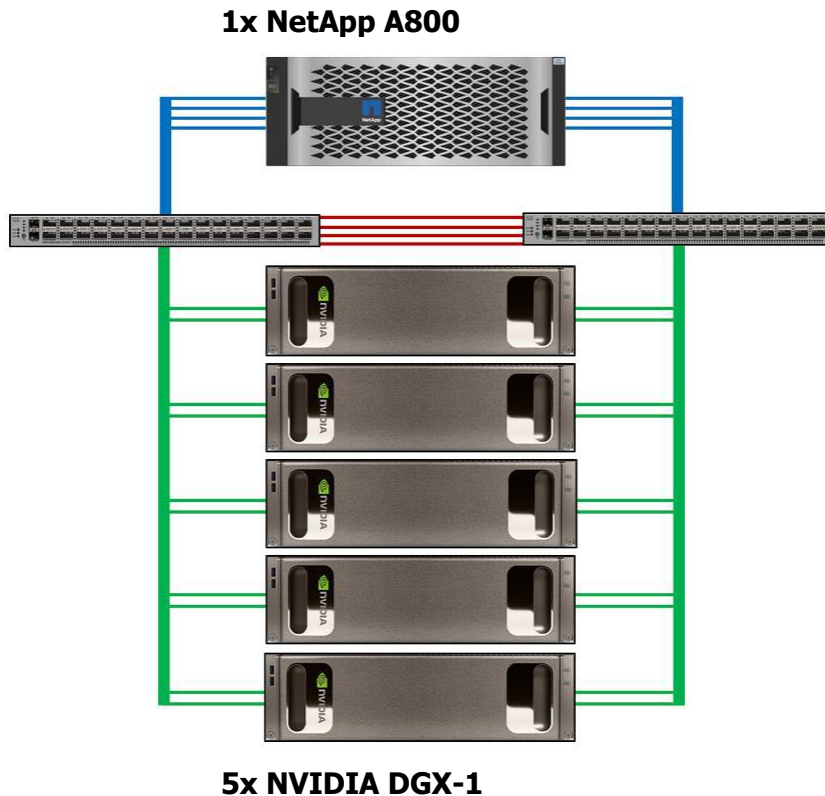
Considering NetApp for AI/ML/DL Workflows

NetApp's strategy is to make AI/DL data flow and management smooth, performant, and integrated from edge to core to cloud. By combining ONTAP data management/Data Fabric, AFF all-flash arrays, and Cloud Volumes, organizations can quickly deploy an optimized platform for AI/ML/DL workloads with high performance.

NetApp's AI solutions are as follows:

- » The NetApp Data Fabric includes data management technologies for the DL pipeline. At the edge, NetApp offers ONTAP Select, which runs on commodity hardware to enable data aggregation and advanced data management. To address storage needs for both the data lake and the training cluster, NetApp AFF storage delivers high performance and high capacity while reducing the need for time-consuming data copies.
- » With its NetApp ONTAP AI solution, the company has partnered with NVIDIA to introduce a rack-scale architecture that enables organizations to start small and grow the infrastructure seamlessly as the number of projects and the size of data sets increase. Organizations can deploy DL workflows in a few hours and seamlessly scale out as needed.
- » If an organization decides to use cloud-based GPU services, it has the option to store and seamlessly manage NFS data as a cloud service using NetApp Cloud Volumes. Cloud Volumes is a managed, high-performance file system service that enables companies to run highly available workloads in a public cloud. It enables secure data transfers and synchronization between on-premises and cloud storage, NAS, and object stores. This means that AI projects can be quickly started in the cloud and then moved on-premises as needed to take advantage of a dedicated DGX system.

Figure 4 shows the NetApp ONTAP AI solution in a 1:5 configuration that consists of five DGX-1 servers fed by one A800 high-availability pair via two switches. Each DGX-1 server connects to each of the two switches via two 100GbE links. The A800 connects via four 100GbE links to each switch. The switches can have two to four 100Gb interswitch links, designed for failover scenarios. The high-availability design is active-active, so maximum throughput can be sustained across all network connections in the absence of a failure.

FIGURE 4: NETAPP ONTAP AI 1:5 CONFIGURATION

Source: NetApp, 2018

The NetApp ONTAP AI solution includes the following software and infrastructure components:

- » **NVIDIA's DGX-1 Servers:** Each DGX-1 server is powered by eight Tesla V100 GPUs, configured in a hybrid cube-mesh topology using NVIDIA NVLink providing an ultra-high bandwidth, low-latency fabric for inter-GPU communications essential to multi-GPU training, thus eliminating the bottleneck associated with PCIe based interconnect. The DGX-1 server is equipped with low-latency, high-bandwidth network interconnects for multinode clustering over RDMA-capable fabrics and leverages GPU-optimized software containers from NVIDIA GPU Cloud (NGC), including containers for the most popular DL frameworks. The NGC DL containers are pre-optimized at every layer, including drivers, libraries, and communications primitives, and deliver maximum performance for NVIDIA GPUs. These pre-integrated containers insulate users from the constant churn typically associated with popular open source DL frameworks, thus providing teams with a stable, QA-tested stack on which to build enterprise-grade DL applications.
- » **NetApp AFF All-Flash Arrays:** This storage system is architected to deliver high throughput while maintaining low-latency support. According to NetApp, a single NetApp A800 system supports throughput of 25GBps for sequential reads and 1 million IOPS for small random reads at sub-500 microsecond latencies. It also supports the 100GbE network, which accelerates data movement and fosters balance in the overall training system because the DGX-1 supports 100GbE RDMA for cluster interconnect.

The NetApp A700s system supports multiple 40GbE links to deliver a maximum throughput of 18GBps. The NetApp A800 and A700s systems can scale independently, seamlessly, and nondisruptively from two nodes (364.8TB) to a 24-node cluster (74.8PB with A800, 39.7PB with A700s). Using ONTAP FlexGroup volumes enables easy data management in a single namespace logical volume of 20PB, supporting more than 400 billion files. For cluster capacities greater than 20PB, organizations can create multiple FlexGroups to span the required capacity.

- » **Trident Storage Provisioner:** As organizations accelerate their rate of data collection, the need to introduce automation around that data becomes apparent. Containers are one way to achieve this; a container enables faster deployments by separating applications from the operating system and device-driver layer dependencies. Efficient and easy data management is central to reducing the time to train. Trident is a NetApp dynamic storage orchestrator for container images that is fully integrated with Docker and Kubernetes. Trident dynamically provisions new or existing volumes directly to application container pods for persistent data storage with all of ONTAP's efficiency and data management capabilities.

Challenges and Opportunities

In IDC's opinion, NetApp's AI solutions are architected to support an edge to core to cloud approach for AI workflows. The ONTAP AI solution, developed in partnership with NVIDIA, is an optimized software and hardware stack. It consists of tested, supported, and ready-to-use libraries, high-throughput GPU, and storage that is intelligent, scalable, secure, metadata rich, cloud integrated, multiprotocol, high performing, and efficient. The fully integrated hardware and software solution, backed by NVIDIA expertise, can accelerate DL deployments, reducing training time from weeks to days or hours, and it increases the productivity of data scientists, enabling them to spend more time on experimentation rather than systems integration and IT support.

The solutions can reduce the complexity of AI deployments and help organizations improve productivity and efficiency, lower acquisition and support costs, and accelerate adoption of AI. Potential opportunities to enhance the solution include the following:

- » The support of heterogeneous processing architectures (e.g., GPUs, FPGAs, ASICs, or Manycore processors) provides organizations with the flexibility to select the appropriate acceleration technology based on performance, operating environment, required skill sets, costs, and energy demand for their AI deployments.
- » The ability to partner with the ecosystem provides centralized management, monitoring, automation, and end-to-end security for running deep learning frameworks and applications within its context.
- » The build/training stage of the AI/ML/DL process is compute heavy and very iterative, and expertise in model tuning and optimization is scarce. This is the stage where an organization's gaps in DL and data science skills hurt the most. Partner with the ecosystem to suggest and optimize the hyperparameters for the frameworks, allow for flexible allocation of resources at and during runtime, prioritize one job over another, and enable resiliency to failure. Also, explore runtime training visualization that allows the data scientist to see the model's progress and provides the opportunity to stop training if it's not providing the right results, which helps in delivering more accurate neural models faster.

Conclusion

Running applications infused with AI/ML/DL algorithms to achieve better business outcomes is critical for enterprises worldwide and necessary for the bulk of DX efforts and use cases. To help organizations accelerate AI-driven business outcomes and overcome deployment obstacles, IDC offers the following guidance:

- » Focus on the business outcomes, keep the project timeline well defined, and prioritize projects with immediate revenue and cost impact.
- » Seek software tools to simplify and automate data preparation, and accelerate the iterative building, training, and deployment of AI models to drive improved business outcomes.
- » Choose the best storage architecture and data architecture to meet your needs today while keeping an eye on the future.
 - Accelerate the flow of data through your pipeline, whether the data is on-premises or in the cloud.
 - Implement intelligent data management at the edge to better cope with data growth.
 - Move data more intelligently and efficiently from the edge to the core or the cloud.
 - Create a more agile core hardware architecture that can evolve quickly.
- » Embrace intelligent (self-configurable, self-healing, self-optimizing) infrastructure, leverage the infrastructure for predictive analytics and valuable insights, and then slowly phase in task automation once the trustworthiness and quality of data are established.

**About the analyst:*****Ritu Jyoti, Program Vice President***

Ritu Jyoti is Program Vice President for IDC's Cloud IaaS, Enterprise Storage and Server team, which includes research offerings, quarterly trackers, advisory services, and consulting programs. Ms. Jyoti is responsible for managing the systems infrastructure research portfolio spanning topics such as artificial intelligence, big data analytics, cloud computing, software-defined infrastructure, cloud data management and protection, and digital transformation — IT transformation data infrastructure strategies.

**IDC Corporate USA**

5 Speen Street
Framingham, MA
01701, USA
T 508.872.8200
F 508.935.4015
Twitter @IDC
idc-insights-community.com
www.idc.com

This publication was produced by IDC Custom Solutions. The opinion, analysis, and research results presented herein are drawn from more detailed research and analysis independently conducted and published by IDC, unless specific vendor sponsorship is noted. IDC Custom Solutions makes IDC content available in a wide range of formats for distribution by various companies. A license to distribute IDC content does not imply endorsement of or opinion about the licensee.

External Publication of IDC Information and Data — Any IDC information that is to be used in advertising, press releases, or promotional materials requires prior written approval from the appropriate IDC Vice President or Country Manager. A draft of the proposed document should accompany any such request. IDC reserves the right to deny approval of external usage for any reason.

Copyright 2018 IDC. Reproduction without written permission is completely forbidden.