

Accelerating Workloads with IBM Storage Scale & Storage Scale System

Distributed file and object storage for AI,
high-performance computing, analytics, and
other data-intensive applications



Highlights

Storage Scale software
provides a global data
platform for distributed file
and object storage

A single Storage Scale System
6000 hardware delivers up to
330 GB/s throughput and up
to 13 million IOPS

Primary use cases include
GPU-accelerated AI, analytics
and data lakehouses, and
high-performance computing

Also widely used for IT
modernization, data backup,
and long-term archiving

Organizations today are reassessing their data storage strategies to adapt to a new generation of data-intensive workloads, especially those used for artificial intelligence (AI) and machine learning (ML).

With a mandate from leadership to leverage AI and unlock the value of organizational data, IT leaders face challenges that include:

- Accessing and analyzing data and workloads scattered across the globe.
- The increasing time needed by AI training and inferencing workloads.
- Managing the growing AI infrastructure and ensuring scalability for evolving workloads.

Addressing these challenges requires specialized software and hardware:

- IBM Storage Scale is software-defined file and object storage optimized for unstructured data.
- IBM Storage Scale System 6000 is a hardware implementation of Storage Scale software that is optimized for the most data-intensive workloads.

Storage Scale

IBM Storage Scale is designed to provide a global data platform that addresses these challenges, with global data abstraction services that provide connectivity from multiple data sources and multiple locations to bring together data from IBM and non-IBM storage environments. It's based on a massively parallel file system and can be deployed on multiple platforms including x86, IBM Power, IBM Z, ARM-based POSIX client, virtual machines, and Kubernetes.

Storage Scale System 6000

Storage Scale System 6000 is a hardware platform that's designed to be the simplest and fastest way for organizations to build a global data platform around their file and object data. It leverages the power of Storage Scale software combined with NVMe flash and hybrid flash/disk technology to deliver high-performance storage for AI, data analytics, and file and object use cases.



Figure 1 – The IBM Storage Scale System 6000 can deliver up to 330 GB/s throughput, up to 13M IOPS, and up to 1.8 PB effective capacity in a 4U rack.

Storage Scale System 6000 is available in all-flash and hybrid configurations providing:

- Up to 330 GB/s throughput with low latency.
- Up to 13 million IOPS using NVMe over Fabric (NVMe-oF).
- Up to 1.8 PB effective capacity in a standard 4U rack space.

Storage Scale software is designed to enable the Storage Scale System 6000 to scale linearly, so that throughput increases proportionally as more systems are added to a cluster. For sequential data and workloads requiring access to massive data sets, Storage Scale System supports up to nine SAS HDD expansion enclosures, providing a more cost-effective alternative to flash storage.

Use Cases

Storage Scale and Scale System are used by organizations worldwide across nearly every vertical industry, especially in financial services, universities and research organizations, automotive, computer services, telecom, and government.

The most common use cases include:

- GPU-accelerated AI;
- Analytics and data lakehouses;
- High performance computing;
- IT modernization;
- Archive and data backup.

Unlocking AI Potential with Content-Aware Storage

Very little enterprise data has been indexed for generative AI applications, which prevents AI assistants from providing accurate, up-to-date answers. The content-aware storage capabilities in Storage Scale address this challenge by extracting the semantic meaning hidden inside unstructured data so that AI assistants can automatically generate smarter answers. Storage Scale enriches data using embedded compute and data pipelines that minimize data movement and latency to help reduce costs and improve performance.

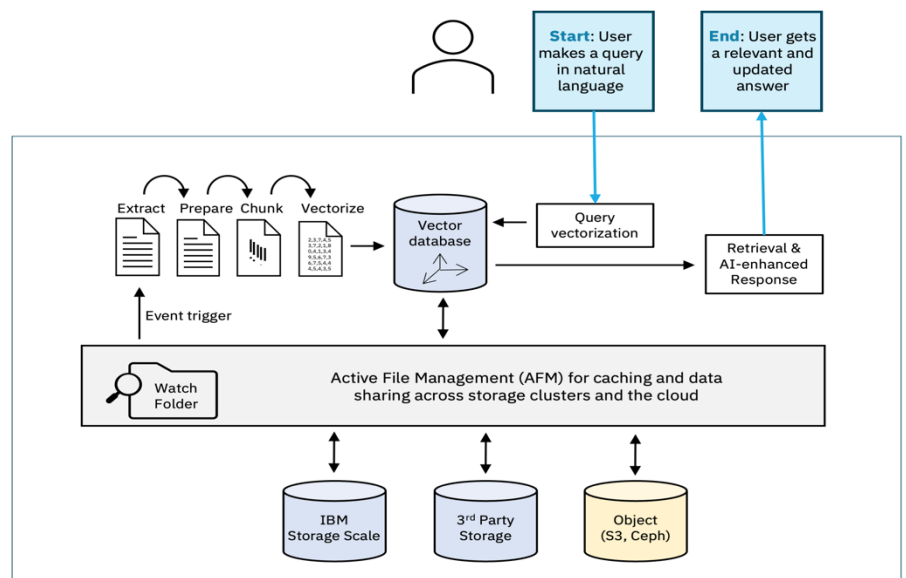


Figure 2. Storage Scale automates data extraction, vectorization, and storage updates, enabling seamless retrieval via a vector database. When users submit natural language queries, AI enhances search results for optimized responses.

GPU-Accelerated AI

Today's most demanding workloads, including AI and advanced analytics, rely on massive datasets and require significant investments in compute, storage, networking, and expertise to build and scale high-performance clusters. To support these efforts, organizations need a storage architecture that delivers high data throughput to accelerate AI training and prevent GPU idle time caused by slow I/O.

This is especially important during model training, which can last days or even months. Interruptions like power outages or hardware failures may force a full restart, wasting time and compute resources. To avoid this, training workflows periodically save checkpoints—snapshots of the model's internal state, including weights and learning rate. Checkpointing enables fault tolerance by allowing training to resume from a known state, but it is a synchronous process that pauses training while writing. As models grow, so do checkpoint sizes; for instance, a large language model (LLM) with one trillion parameters may require up to 14 terabytes per checkpoint.

This is where Storage Scale provides a significant advantage, with its POSIX-style file system optimized for multi-threaded read and write operations across multiple nodes. It can be deployed as primary storage or as a high-performance tier in front of object storage. Acting as a caching layer between GPUs and object storage, its Active File Management (AFM) enables faster data loading at training start or restart and allows model weights to be checkpointed to the file system more efficiently than writing directly to object storage. AFM then asynchronously transfers checkpointed data to object storage without slowing down the training process.

Engineered to accelerate AI workloads, Storage Scale System 6000 supports the NVIDIA Magnum IO GPUDirect Storage (GDS) technology, which enables a direct data path between GPU memory and local or remote storage, such as NVMe or NVMe-oF. GDS removes the host server CPU and DRAM from the data path, so the IO path between storage and the GPU is shorter and faster.

Storage Scale System 6000 is a certified storage solution for NVIDIA DGX SuperPOD and BasePOD AI infrastructures. To streamline AI deployment, IBM and NVIDIA jointly test, plan, and install integrated systems, providing certified reference architectures that combine Storage Scale System 6000 with both DGX SuperPOD and BasePOD, with storage deployment and support backed by IBM's global services team.

Analytics and Data Lakehouses

Today's advanced analytics workloads generate massive volumes of data from diverse sources, necessitating storage solutions that can scale, support complex processing, and provide rapid access to both real-time and historical data. To meet these demands, organizations increasingly rely on data lakes and data lakehouses, each serving distinct but complementary roles.

Data lakes store raw, unstructured, and semi-structured data in its native format, making them ideal for exploratory analytics, machine learning, and AI workloads where flexibility and scale are critical. By contrast, data lakehouses, such as IBM watsonx.data, combine the raw data storage capabilities of data lakes with the structured, performance-oriented approach of data warehouses.

watsonx

To support these evolving data and AI needs, IBM provides **watsonx**, an integrated AI and data platform built for business. It includes a set of components designed to support every stage of the AI and analytics lifecycle:

- **watsonx.ai** is an interactive studio for building and deploying AI applications.
- **watsonx.governance** is designed to help organizations direct, manage and monitor their activities for generative AI and machine learning models, including health, accuracy, drift, and bias.
- **watsonx.data** provides a data lakehouse that brings together all an organization's business data to scale analytics and AI.
- **watsonx Assistant** is a conversational platform for building and deploying virtual assistants, chatbots, and other interactive agents across various channels.

Lakehouses are designed to transform, cleanse, and structure data, making it suitable for high-performance analytics and business intelligence applications. They address some of the limitations of data lakes, such as data redundancy and inconsistent quality, by enforcing data management and governance practices while maintaining the scalability and flexibility of a data lake.

The distributed file and object capabilities of Storage Scale are ideal for supporting both data lakes and lakehouses. Storage Scale supports a wide variety of data formats, can scale to handle vast amounts of data across distributed environments, and integrates seamlessly with big data frameworks like Apache Hadoop and Spark. It delivers the throughput and low-latency access that are essential for accommodating multiple concurrent users and applications.

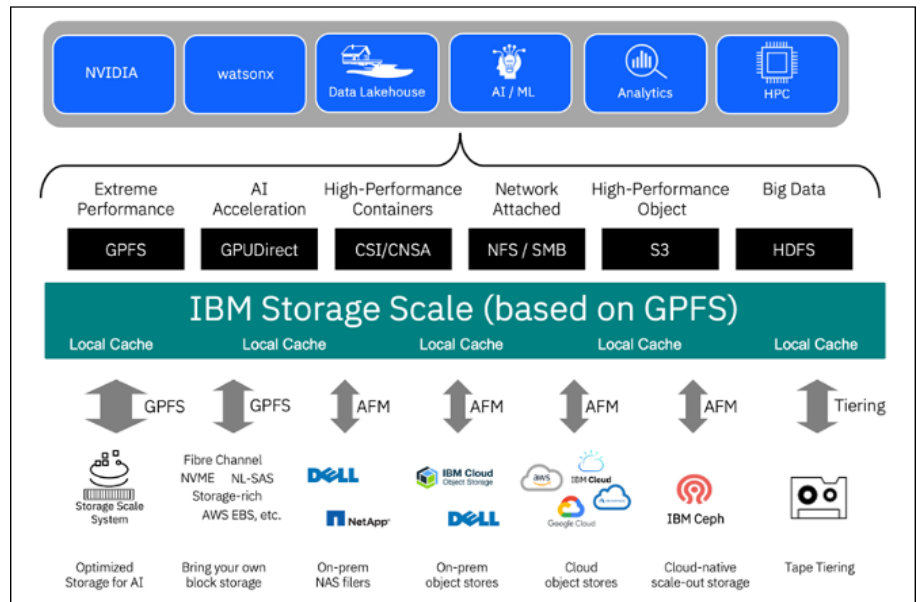


Figure 3 – IBM Storage Scale software provides organizations with a global data platform optimized for today's most demanding unstructured data workloads.

The unique global data abstraction services in Storage Scale can help organizations run analytics workloads against data regardless of its format or location, breaking down traditional data silos and enabling seamless access to both structured and unstructured data. Only Storage Scale software provides shared multi-protocol data access – the ability to ingest data via one protocol and make it available to multiple workloads simultaneously, even in different data protocols.



Figure 4 – The IBM Storage Scale System Expansion Enclosure enables organizations to cost-effectively deploy workloads operating on massive data sets.

Tiered Storage

Analytics workloads often require access to two different kinds of storage systems – some that use high-performance media for quickly reading and writing active data, and others that provide more cost-effective storage for less frequently accessed data. Storage Scale System 6000 is designed to allow organizations to dial in the exact balance of performance and capacity required by their workloads.

Scale System 6000 can be configured with standard NVMe flash drives for maximum performance or with IBM FlashCore Modules when data density and compression are higher priority.

For workloads where storage capacity is important, Scale System 6000 supports up to nine expansion enclosures. The IBM Storage Scale Expansion Enclosure is an enterprise-class, fully redundant storage enclosure, containing up to 91 20TB or 22TB self-encrypting SAS hard disk drives (HDDs). Attaching eight expansion enclosures to the Storage Scale System 6000 expands the maximum storage capacity to 18PB of HDD storage per rack (using 24Gb SAS drives).

High-Performance Computing

The combination of Storage Scale software and Storage Scale System hardware has been widely used in high-performance computing (HPC) environments for years. Indeed, the main reason systems optimized for HPC workloads are also widely used for AI workloads is that the two share many characteristics:

- High data volume and throughput: Both types of workloads often deal with massive datasets, sometimes ranging from terabytes to petabytes. Whether running complex scientific simulations or data-intensive computations, HPC and AI both require storage systems that can provide high throughput to move large volumes of data quickly and efficiently to computing resources. Data transfer between the storage systems and GPU infrastructure is especially important.
- Low latency access: HPC and AI workloads both demand low-latency access to data. In HPC, quick data retrieval is essential for maintaining high performance in simulations and real-time analyses, where delays can significantly impact the quality and business value of results. AI models, particularly deep learning models, require vast amounts of data to train effectively.

During training, these models repeatedly access large datasets to learn from the data. Low-latency storage ensures that data can be quickly retrieved and fed into the training process to accelerate the overall training cycle. Also, many AI applications, such as advanced driver assistance systems (ADAS), fraud detection, and recommendation systems, rely on real-time or near-real-time inference. Low-latency storage minimizes the time it takes to access the necessary data, enabling faster and more responsive AI models.

- Parallel processing and concurrency: Both HPC and AI workloads often involve parallel processing across multiple compute nodes. This necessitates storage systems that support high IOPS (input/output operations per second) and can handle many simultaneous read and write operations, ensuring smooth data flow across distributed environments.
- Scalability: As with AI workloads, HPC applications need storage systems that can scale seamlessly to meet increasing demands. As the complexity and scale of simulations or computations grow, storage must scale in both capacity and performance to avoid bottlenecks and maintain efficiency.

- Data integrity and reliability: Maintaining data integrity is vital in HPC, as it is in AI workloads. Accurate data is critical for valid simulation outcomes and reliable results. Therefore, HPC storage systems must include robust data protection features to prevent data loss and ensure high reliability.

Together, Storage Scale and Storage Scale System check off every item on this list – they’re designed for high data volume and throughput, with low latency access, supporting parallel processing and concurrency, and providing massive scalability, data integrity, and reliability.

IT Modernization

IT modernization is the process of updating and optimizing an organization's technology infrastructure and applications to be more agile, scalable, and efficient, often by adopting cloud-native technologies and methodologies. In practical terms, this means standardizing on a single platform that supports the development, deployment, and management of modern applications across hybrid and multi-cloud environments.

For many enterprises, that platform is Red Hat OpenShift, because:

- OpenShift supports the shift to containerized, microservices-based architectures, making applications more modular, scalable, and portable.
- It provides a consistent platform across different cloud environments, facilitating the move to hybrid and multi-cloud setups without vendor lock-in.
- It automates deployment and supports CI/CD (continuous integration / continuous delivery) pipelines, accelerating development cycles and enhancing operational efficiency.
- It has robust built-in security features to help ensure that modernization efforts align with evolving security and compliance requirements.
- It optimizes resource usage, reducing costs and improving performance as organizations modernize their IT infrastructure.

Storage Scale software provides distributed file and object storage to help organizations build the scalable, flexible, and resilient data infrastructure they need to support IT modernization efforts powered by platforms like OpenShift.

First and foremost, Storage Scale delivers the scalability and flexibility organizations need to manage unprecedented growth in data volumes across hybrid and multi-cloud environments, complementing OpenShift's flexibility. Its unique data abstraction capabilities, which enable real-time translation between file-based and object-based data, provide the foundation for organizations' global data platforms.

The integration of Storage Scale with OpenShift provides consistent, persistent storage for containerized applications, helping ensure data accessibility across diverse deployment environments. Because Storage Scale operates on structured, semi-structured, and unstructured data, it supports a broad range of modern and traditional applications.

Storage Scale was designed for high availability and resilience, with built-in redundancy to help ensure applications on OpenShift remain reliable, even in distributed environments. It also allows for tiered storage, optimizing resources and aligning with OpenShift's resource management for better cost efficiency.

Storage Scale supports container-native storage. It integrates with Kubernetes and other container orchestration platforms to provide scalable, high-performance storage for containerized applications. Storage Scale offers features such as dynamic provisioning, data sharing, and high availability, which are essential for running stateful applications in containers. This allows Storage Scale to be used directly within container environments, making it suitable for various workloads that require persistent storage in a container-native format.

Backup & Archive

Data resilience is a crucial part of 21st-century business, as organizations strive to improve their ability to withstand and recover from ransomware, hardware failures, human error, natural disasters, and other threats. They typically need frequent backups of their active data as well as cost-effective storage of their massive archival data sets. In addition, many organizations face rigorous regulatory requirements for data storage, with significant financial penalties for non-compliance. Storage Scale is well-suited for data backup and archiving, providing a scalable, high-performance solution with advanced data management features.

Its key capabilities include:

- **Scalability:** Storage Scale is designed to handle massive amounts of data, scaling easily across distributed environments to accommodate growing backup datasets and archives without compromising performance.
- **Automated data tiering:** To optimize storage costs, Storage Scale automatically moves data to different storage tiers (i.e.: flash, disk, tape, or cloud) based on policies, all managed automatically.
- **Data replication and protection:** Advanced replication capabilities and snapshots help ensure robust data resilience to help mitigate the effects of ransomware and other potential threats to data.
- **Tape integration for long-term archiving:** Storage Scale integrates seamlessly with IBM Tape Storage, using tape as an economical and durable solution for long-term data archiving.
- **High availability and reliability:** Features like erasure coding and automated failover help ensure high data availability and protect against hardware failures.
- **Policy-based data management:** Storage Scale supports policy-based data management, which reduces administrative overhead by automating the movement and retention of data according to predefined rules.
- **Data security and compliance:** Storage Scale provides encryption and access controls to help ensure that data meets compliance requirements and is protected against unauthorized access.

Overall, Storage Scale delivers a comprehensive solution for scalable, efficient, and secure data backup and archiving, combining advanced data management capabilities with robust protection and cost optimization.



Figure 5 – The IBM Diamondback Tape Library delivers ultra-high data density, with up to 27.9 PB of native data in a single eight-square-foot library using LTO Ultrium 9 cartridges.

For more information

To learn more about IBM Storage Scale and IBM Scale System, contact your IBM representative or IBM Business Partner, or, for Storage Scale software go to: <https://www.ibm.com/products/storage-scale> or download a data sheet at: <https://www.ibm.com/downloads/cas/LGPLW1MO>

For Storage Scale System go to: <https://www.ibm.com/products/storage-scale-system> or download a data sheet at: <https://www.ibm.com/downloads/cas/JBVOYVXB>

© Copyright IBM Corporation 2025
IBM Corporation
New Orchard Road
Armonk, NY 10504

Produced in the
United States of America
April 2025

IBM and the IBM logo are trademarks or registered trademarks of International Business Machines Corporation, in the United States and/or other countries. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on ibm.com/trademark.

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

THE INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS" WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT.

IBM products are warranted according to the terms and conditions of the agreements under which they are provided.

