

IBM Scale System 6000 AI Data Platform

Massive throughput for AI workloads, built with content-aware storage and integrated GPU acceleration.



Highlights

A production-ready AI data platform built on IBM Storage Scale and integrated with NVIDIA GPUs

GPU-accelerated data pipelines that improve preprocessing, retrieval, and inference performance

Zero-copy access to live enterprise data for AI workloads at scale

Content-aware storage that keeps RAG and agentic AI data fresh, secure, and aligned with permissions

Enterprises are rapidly scaling AI from pilots to production, but many are discovering that data, not compute, is the limiting factor. AI pipelines are constrained by fragmented data, excessive duplication, and infrastructure designs that force constant movement of large datasets between storage, preprocessing systems, and GPUs. As models grow larger and AI systems become more dynamic, these approaches drive up cost, slow time-to-value, and undermine governance.

The IBM Storage Scale System 6000 AI Data Platform (AIDP) addresses this challenge by rethinking storage as an active AI data platform. AIDP enables AI models to work directly with live, governed enterprise data at scale, with high performance and without re-architecting existing pipelines.

Built on IBM Storage Scale System 6000, enhanced with IBM Fusion Content-Aware Storage (CAS), and accelerated by NVIDIA GPUs, AIDP provides a production-ready foundation for modern AI workloads across training, retrieval-augmented generation (RAG), and inference.

The challenge: AI pipelines break at scale

Traditional AI architectures rely on duplicated datasets, batch-oriented ETL (extract / transform / load) pipelines, and complex data movement between storage tiers and compute clusters. These patterns introduce several problems as AI adoption grows, including data duplication, pipeline sprawl, performance bottlenecks, and governance gaps. As organizations move toward AI factories and continuous inference, these limitations become major structural barriers.

The IBM Storage Scale System 6000 AI Data Platform shifts the focus from moving data to activating data in place. The platform is designed to make enterprise data directly usable by AI systems while preserving performance, security, and operational simplicity.

At the core of AIDP is a tightly integrated stack comprising IBM Scale System 6000, IBM Fusion content-aware storage, and NVIDIA RTX Pro 6000 GPU acceleration. Together, these capabilities allow AI pipelines to operate on live, multi-modal enterprise data, securely and at scale.

A data-centric platform for enterprise AI

The Storage Scale System 6000 AIDP is architected as a GPU-accelerated, enterprise-ready AI data platform that enables AI pipelines to operate directly on live, governed enterprise data. Rather than moving data through multiple preprocessing tiers, AIDP brings intelligence and acceleration to the data layer itself, thereby improving performance, reducing duplication, and simplifying operations as AI adoption scales.

The architecture integrates three core layers into a single, pre-validated platform:

- **Data platform foundation: IBM Storage Scale System 6000**

IBM Storage Scale System 6000 provides the high-performance data foundation for AIDP. Built on IBM Storage Scale, the platform delivers the throughput, concurrency, and metadata performance required by modern AI workloads, supporting parallel access from many compute nodes without introducing file system bottlenecks. Storage Scale System 6000 serves as the shared system of record for structured and unstructured data across training, retrieval-augmented generation (RAG), and inference workloads, while supporting access to remote and distributed data sources through a unified namespace.

- **Intelligence at the data layer: IBM Fusion Content-Aware Storage**

IBM Fusion Content-Aware Storage adds semantic intelligence directly to the data platform. CAS processes files in place on IBM Storage Scale, extracting metadata, generating vector embeddings, and maintaining alignment with existing enterprise access controls, without copying or relocating data. Incremental change detection ensures that AI-ready representations stay synchronized with live enterprise datasets, enabling RAG and agentic AI systems to reason over current information rather than static snapshots.

By embedding content awareness into the storage layer, AIDP eliminates the need for standalone ingestion and indexing pipelines, reducing operational complexity while improving data freshness and governance consistency.

- **GPU-accelerated AI pipelines: NVIDIA integration**

AIDP integrates NVIDIA RTX Pro 6000 GPU acceleration directly into the data platform to minimize data movement and keep AI pipelines efficient at scale. NVIDIA RTX Pro 6000 GPUs are used to accelerate preprocessing, embedding generation, vector search, and inference, enabling low-latency access to enterprise data for RAG and other AI workloads. By colocating GPU-accelerated processing with the data layer and leveraging direct, high-bandwidth data paths, the architecture ensures that GPUs remain utilized without being constrained by storage or network bottlenecks.

End-to-end data flow for AI workloads

Together, these layers enable a streamlined, end-to-end AI data flow:

- Enterprise data is ingested into IBM Storage Scale or accessed remotely using active file management for external S3 or NFS sources.
- IBM Fusion CAS automatically detects new or changed data and incrementally processes it into AI-ready representations.

- Vector embeddings, metadata, and permissions are maintained in alignment with the source data.
- GPU-accelerated retrieval and inference services use this live data to support training, RAG, and real-time AI applications.

This architecture allows AI systems to operate continuously on governed enterprise data without relying on duplicated datasets or disruptive migrations.

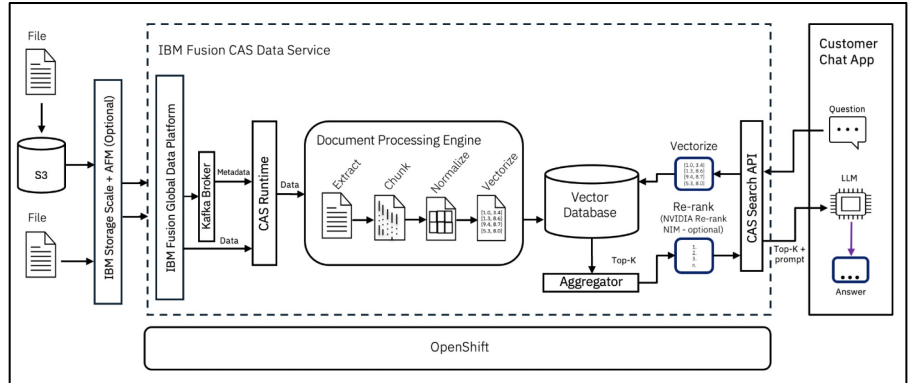


Figure 1. Fusion CAS mounts Storage Scale and processes data in place, generating normalized metadata and vector embeddings with IBM Docling or NVIDIA NeMo.

Built for production AI

The IBM Storage Scale System 6000 AI Data Platform is designed for real-world AI operations, where performance, governance, and reliability matter as much as model accuracy. By unifying high-performance storage, content-aware data services, and GPU acceleration into a single architecture, AIDP provides an enterprise-ready foundation for scaling AI from pilot projects to production AI factories, without rethinking the entire data landscape.

For more information

To learn more about IBM Scale System 6000 AIDP, contact your IBM representative or IBM Business Partner, or visit ibm.com/products/storage-scale-system.

© Copyright IBM Corporation 2026
IBM Corporation
New Orchard Road
Armonk, NY 10504

Produced in the
United States of America
May 2026

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.

