

Turbocharge AI Workloads with an AI-Native Data Platform

Introduction

The growth of artificial intelligence (AI) and the rapid shift to deep learning, machine learning, and high-performance computing workloads means organizations now demand solutions that meet the intensive computing requirements of data-pipeline driven applications.

“The most damaging phrase in the English language is **we’ve always done it this way.**”

REAR ADMIRAL GRACE MURRAY HOPPER

This same principle applies when searching for a data solution that meets the needs for the most demanding, next-generation workloads.

In order for an organization to be successful, however, it is essential that a data infrastructure is specifically designed to:

- Get faster outcomes for next-generation workloads
- Be flexible enough to deploy in the cloud, on-premises, or hybrid
- Scale up or down easily
- Be performant enough to avoid GPUs idling while they are waiting for data
- Be robust enough to avoid needing a team of admins to manage, reconfigure, and keep it running

Understanding AI Data Pipeline Challenges

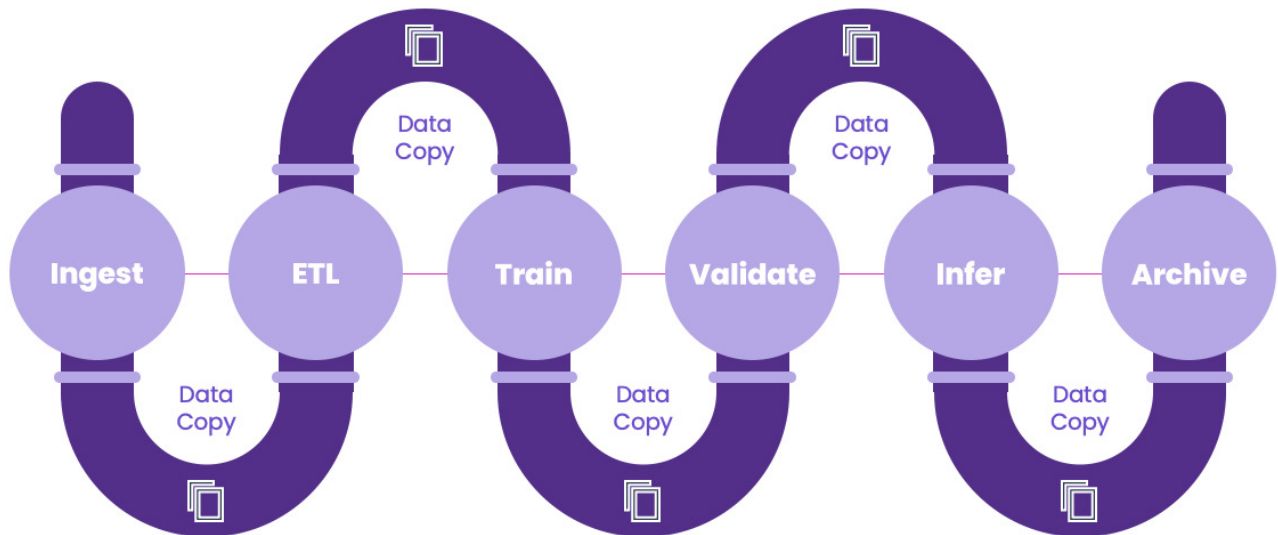
Storage solutions have gotten faster over the years using a variety of methods such as accelerator caching cards, filesystem tweaks, CPU and RAM upgrades, and moving from HDD to SSD. But even with faster storage, over the last several decades, technologies that drive computing and networking are continuing to propel the AI revolution empowering researchers and practitioners to tackle increasingly complex problems and drive innovation across diverse domains.

GPUs are the gold standard for powering high-performance workloads, and the demands on data pipelines continue to escalate. GPUs are voraciously energy-hungry and require

an ocean of streaming data to fuel them. Unfortunately traditional data architectures struggle with latency and bottlenecks in the data pipeline.

Antiquated legacy storage systems, such as those relying on the 30-year-old NFS protocol, present significant challenges for modern AI development. These systems struggle to fully utilize the bandwidth of modern networks, limiting the speed at which data can be transferred and processed. They also often face difficulties handling large volumes of small files. With the problem of lots of small files, storage metadata servers (MDS) can easily become overwhelmed and the result is performance bottlenecks.

Traditional Data Silos



Limitations in traditional storage solutions include:

- Massive concurrency & write throughput
- Annotation, index, search, & cloud bursting
- Massive read throughput
- Massive bandwidth for streams replay
- Low latency access
- Lifecycle management, versioning, & reproducibility

These limitations can actively hinder the scalability and efficiency of AI workflows, impeding tasks such as data preprocessing, model training, and inference.

Revolutionary, not evolutionary change is what the market truly needed to keep pace with not only the technological advancements but customer expectations.

Embracing the Data Platform Revolution

Revolutionary, not evolutionary change is what the market truly needs to keep pace with not only the technological advancements but customer expectations. Although storage solutions have gained performance, many of these GPU-fueled workloads are still compromised because of data bottlenecks. To power data-driven innovation and address data bottlenecks, organizations are eliminating the complexity of legacy data infrastructures and replacing them with data platforms. But the truth is that while there are a lot of newer data infrastructures, not all deliver the value of a data platform.

A data platform capable of delivering unparalleled performance across a multitude of data profiles is critical. A data platform is an integrated, end-to-end solution that provides holistic support for an organization's data management needs while supporting every step of the organizations' data lifecycle – from ingest and pre-processing, to analyzing, storage, and archiving. The faster you can stream various data pipelines to feed the GPU, the faster you will get to results for your AI workloads.

A true data platform is designed to support both the structured and unstructured data a digital organization uses, regardless of whether the data is at the core, cloud, or edge. It is multi-tenant, multi-workload, multi-performant, and multi-location, all with a common management interface. A data platform can handle all the distinct requirements of the different types of data needs across the organization – massive ingest bandwidth, mixed read/write handling, and ultra-low latency. It also manages data across on-premises, cloud, or hybrid environments and enables easy data mobility between them all.

But not all data platforms are created equal. Some of the current solutions in the market fail because they are based on an old legacy architecture masked by the promises of flash. And while the parallel file systems may have the capability of delivering on performance, it is the complexities of management, reliability, usability, metadata management, and overall customer experience that has left many enterprise IT organizations hesitant to deploy into their own data center or cloud. This dilemma has led enterprise IT to have serious trade-off discussions internally in order to meet company expectations.

In some cases, the older, legacy solutions may be a better fit based on unique and monolithic data characteristics. Research, however, indicates the data profiles in most enterprise IT environments are rarely monolithic, and depending on the workload, may require independent filesystems configured specifically to support those monolithic data sets. This is where enterprise IT has generally reached a breaking point between managing multiple silos containing individually configured file systems for optimal performance outcomes or simply averaging out the results across one or two file systems and settling for more of a general purpose solution that is good enough.

But is good enough really good enough when the success of running AI workloads and your business outcomes are inextricably tied to the performance of your data infrastructure?

Unpacking the World of Data Infrastructure Solutions

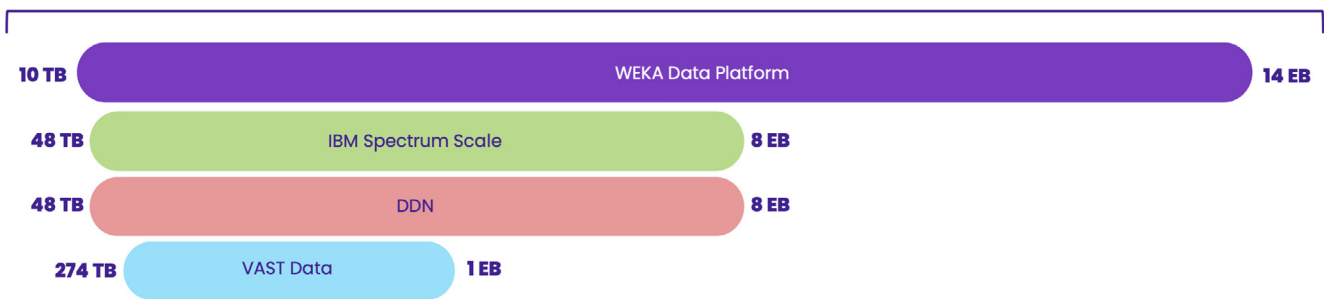
Choosing the right data platform and storage solution can be challenging and depends heavily on an organization's data management needs and goals.

We compared 4 solutions considered by many to be representative of the leaders targeting next-gen, performance-intensive workloads. Our analysis is based on the most common selection criteria customers use to select the right storage infrastructure. Some of the solutions we compare have been in the market for several years and were developed prior to new standards such as NVMe Flash, Cloud, and de facto standard cloud protocols like S3.



	WEKA	IBM	DDN	VAST
Model (Rack Units)	WEKApod™ (1U)	ESS 3500 (2U)	AI400X2 (2U)	Ceres DF30xx (1U)
Read BW (per RU)	90 GBps	63 GBps	60 GBps	64 GBps
Write BW (per RU)	23.3 GBps	30 GBps	37.5 GBps	10 GBps
Read IOPS (per RU)	2,280 KIOPS	1,500 KIOPS	1,500 KIOPS	590 KIOPS
Write IOPS (per RU)	535 KIOPS	no data	1,000 KIOPS	135 KIOPS

Solution Capacity Scale



	WEKA	IBM SPECTRUM SCALE	DDN LUSTRE	VAST DATA
Commodity Hardware	Yes	No	No	No
	Industry standard COTS			
SW Only	Yes	Yes	EXAScaler based on Lustre, only on DDN HW	Yes
				But requires specialized hardware
Tiering to Object	Yes	Yes	No	No
	Global namespace			
Tiering/Hybrid Disk Storage	Yes	Yes	Yes	No
	Tier to object storage		Complex policy based HSM	Cache to SCM SSD/Copy to QLC SSD
Ease of Use	Easy	Hard	Hard	Easy
PROTOCOL				
NFS	Yes	Yes	Yes	Yes
	NFS v3/4	Via additional gateway	Via additional gateway	
LDAP	Yes	Yes	Yes	Yes
SMB	Yes	Yes	Yes	Yes
	SMB2/3.1 including Multichannel and SMB Direct	SMB2/3 or via additional gateway*	Via additional gateway	SMB2
GPUDirect	Yes	Yes	Yes	Yes
POSIX	Yes	Yes	Yes	No
S3	Yes	Yes	Yes	Yes
		limited or via additional gateway*	Via additional gateway	
SECURITY & RECOVERY				
Data Encryption	Yes	Yes	Yes	Yes
	at-rest & in-flight	at-rest & in-flight	Using Linux kernel fsencrypt	With limitations
Data Protection	N+2, N+4	N+3	EC	N+4
		Reed Solomon only ESS Appliance		
Data Efficiency	Yes	Limited	No	Yes
Snapshots	Yes	Yes	No (native)	Yes
		With perf impact	Yes (on ZFS)	
Max Snapshots per file system	24,000	256	Native snapshots not supported	1,000,000
Snapshot to S3	Yes	No	No	Yes
CLOUD				
Cloud Supported	AWS, Google, Oracle, Azure	AWS, Azure, IBM Private Cloud, Oracle	AWS, Azure, GCP	AWS, Azure, GCP
Tier to Cloud	Yes	Yes	No	Yes
Backup to Cloud	Yes	Yes	No	Yes
Burst to Cloud	Yes	No	No	Partial

VAST Data

VAST Data is a relatively new company in the storage solution space. It was initially designed and positioned as a backup target, presumably competing in the backup appliance market, and then shifting its focus to more general IT primary storage. VAST's initial product in market consists of non-standard storage class memory (SCM) by Intel called Optane (3D XPoint) that caches all of the writes to a broad deployment of SCM SSDs in the VAST system and then evacuates the data to QLC SSDs on the backend when more front-end cache space is needed. This design requires a great deal of copying of data across and between SCM devices in order to be as optimally compacted as possible prior to writing to lower endurance QLC.

IBM Spectrum Scale (GPFS)

Spectrum Scale or GPFS is a parallel file system which began in 1993 as the Tiger Shark file system, a research project at IBM's Almaden Research Center. Tiger Shark was initially designed to support scientific computing and high throughput multimedia applications, such as streaming video from VHS tapes. Given its many configuration options, GPFS was, and still is, a complex, brittle solution requiring a high level of expertise to plan, install, configure, and operate.

Lustre

Similar to GPFS, Lustre's origins can be traced back to a research project that began at Carnegie Mellon University in 1999. By 2001 the company, Cluster File Systems, Inc was formed and work on what would become the Lustre file system had begun under a program funded by the US Department of Energy (DoE). Lustre is a parallel file system, like GPFS, but it has struggled to find mainstream adoption outside of national labs and impromptu builds for spot testing and evaluations. Lustre has changed hands several times over the 20+ years it has been in existence, from Sun, Oracle, Whamcloud, and Intel. Intel abandoned Lustre in pursuit of its own file system, selling it off to DDN who has now released it back to open source and the user community. Given its age, little development has been done with Lustre over the years which is reflected in its current version number, 2.15.4.

SUPPORT FOR MONOLITHIC DATA SETS

This is where enterprise IT has generally reached a breaking point between managing multiple silos containing individually configured file systems for optimal performance outcomes or simply averaging out the results across one or two file systems and settling for more of a general purpose solution that is good enough.

The Revolutionary WEKA Data Platform

Tired of seeing enterprise customers being forced to use disparate and siloed legacy data storage and infrastructure solutions that were costly, complex, wasteful to deploy, manage and maintain, and were also ill-suited for modern workloads, WEKA's founders decided to develop a new solution.

The goal was a single product that would be powerful enough to meet the insatiable performance and compute demands of next-generation, AI workloads in the world's most demanding and distributed environments.

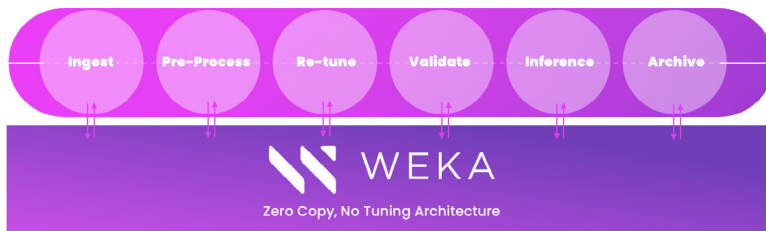
The solution would need to:

- Have exabyte-scale capacity
- Process tens of terabytes of data a second
- Be as simple to use as an iPhone
- Be deployable within or beyond the walls of the traditional data center so organizations could run their research or business virtually anywhere

In 2017, the WEKA® Data Platform was born. The earliest iteration was delivered as a screaming fast, cloud-native parallel file system that was purpose-built for large-scale data processing and high-performance computing.

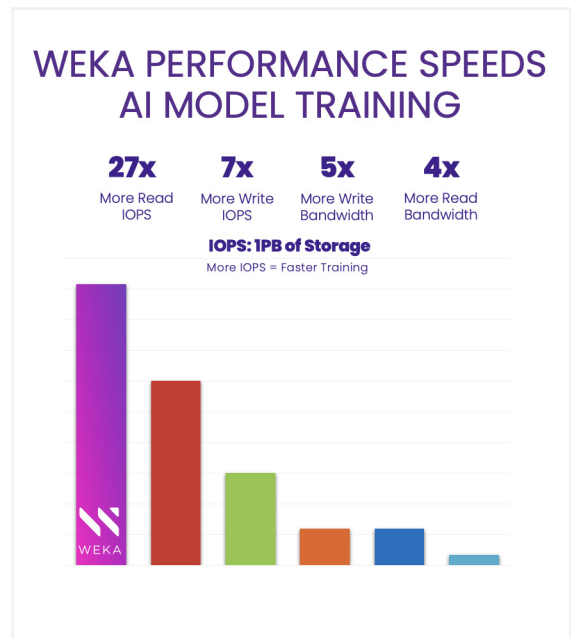
Now—fortified with rich data services, multi-cloud capabilities, and delivered as a single software solution on a customer's existing hardware, in the cloud of their choice, or as a service—The WEKA® Data Platform has evolved into a robust data platform that is truly AI-native and specifically designed to effectively and efficiently power HPC, AI, and other next-generation workloads.

The WEKA Data Platform

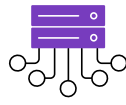
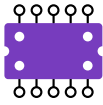


Push the limits of what is possible. With industry-leading performance, you can dramatically accelerate time to insights. Instantly adapt to any new workload with no knobs to tune and no expertise required. Get up and running quickly in the cloud or on-premises with validated reference architectures.

The [WEKA® Data Platform](#) delivers the radical simplicity, epic performance, and infinite scale required to support AI, machine learning, and other next-generation workloads in virtually any location.



WEKA has built an incredibly strong partner ecosystem that includes solution resellers and leading server, processor, and network interconnect manufacturers. In addition, the company has partnerships with the world’s largest server OEMs and the world’s major cloud providers. So, whether you are looking for an on-premises, cloud, or hybrid deployment, WEKA has you covered.



OPTIMIZED FOR NVME

MULTI-PROTOCOL READY

BUILT-IN DURABILITY

ADVANCED SECURITY

BUILT FOR THE CLOUD

Achieve lowest possible latency and highest performance

Supports Linux, Windows and Native POSIX access to data

Uses distributed data protection and instant backup to S3 cloud for rapid recovery

Keeps your data completely safe with integrated encryption, key management, and access control

Seamlessly run on-premises, in the cloud and burst between platforms

We’d love to show you how the [WEKA Data Platform](#) can take your business to the next level.

Connect with Us



weka.io

844.392.0665

