



WEKA for GenAI

Shimon Ben David – CTO

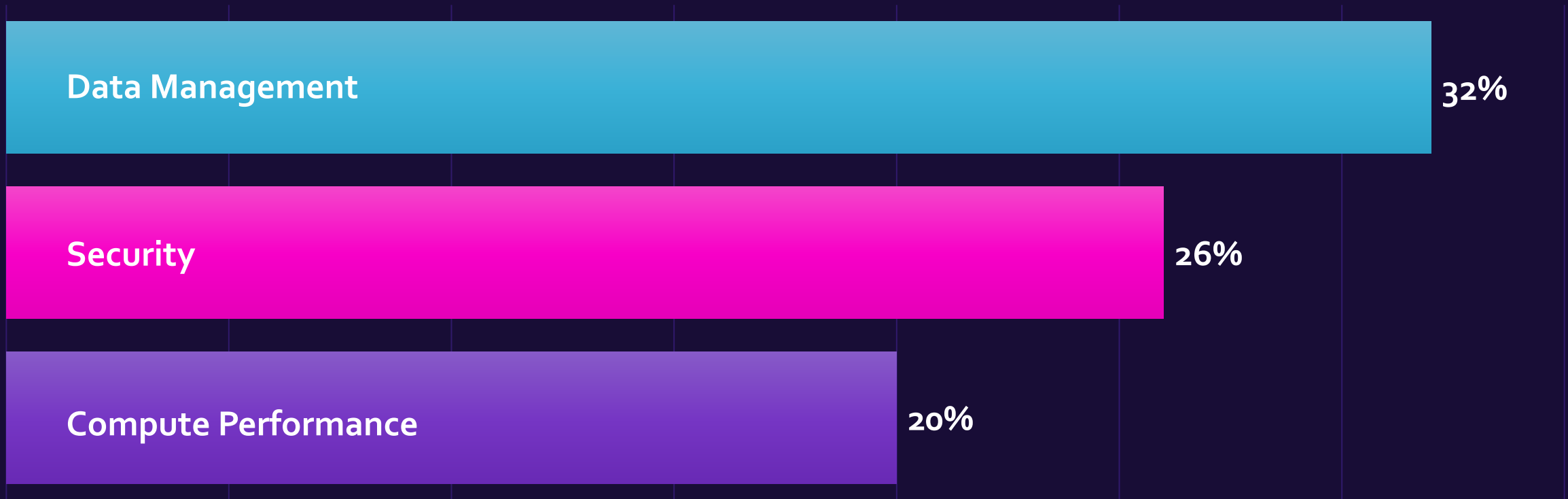
A Majority of Enterprises Are Now Investing In Generative AI



By 2027, 90% of enterprises will deploy generative AI models and applications —up from less than 5% in 2023.

KEY FINDING:

Top technical inhibitors to AI/ML success

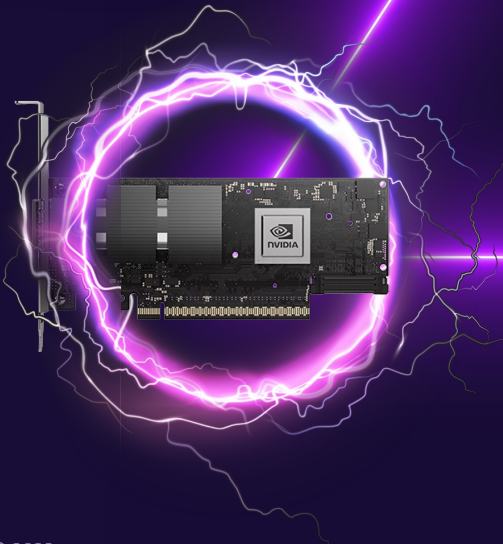




Compute

Next generation GPU Accelerators

The Infrastructure Triangle



Networking

Modern 400GbE/IB fabrics
(800Gb/s Switches)

"Vintage"
Storage

Antiquated Storage

- 30-year-old protocol (NFS)
- Unable to saturate modern networks
- Lots of small files overwhelming MDS
- etc.



GPU Acceleration 1000x vs Traditional CPU

2008

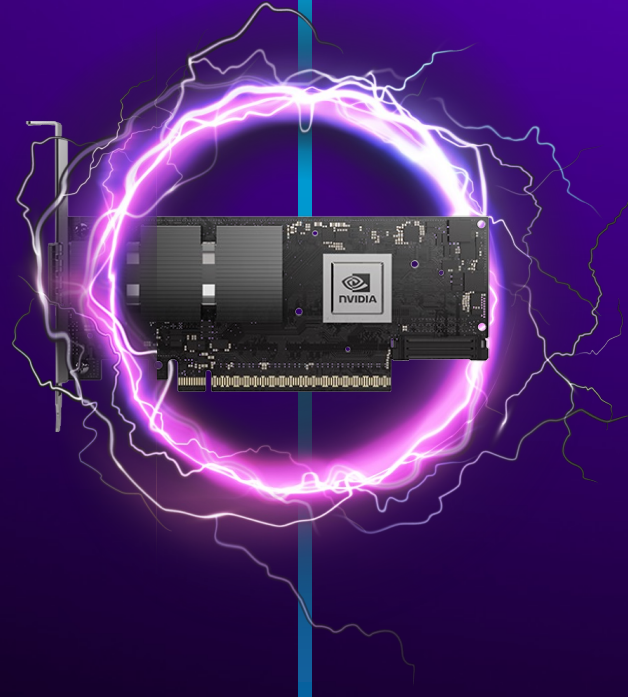
First Top500 supercomputer

- 10,400 Rack Units
- 1 PetaFLOP

2022

NVIDIA DGX-H100

- 8 Rack Units
- 32 PetaFLOP



Modern Networking 80x Improvement vs legacy

2002

10GbE IEEE Standard

- Full duplex point-to-point links
- 10 gigabits per second

2022

NVIDIA ConnectX-7 400Gb

- NDR InfiniBand and 400Gb Ethernet
- 400 gigabits per second, with switch back hauls of 800Gb/s



**Traditional PFS
All Flash NAS
Software Defined 'NFS'**

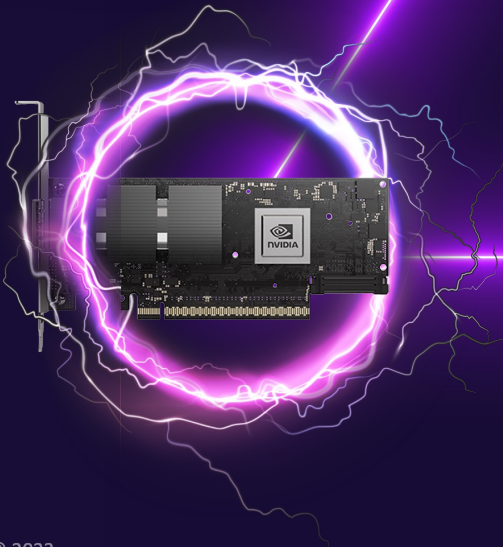




Compute

Next-generation GPU Accelerators

The Infrastructure Triangle



Networking

Modern 400GbE/IB fabrics
(800Gb/s switches)

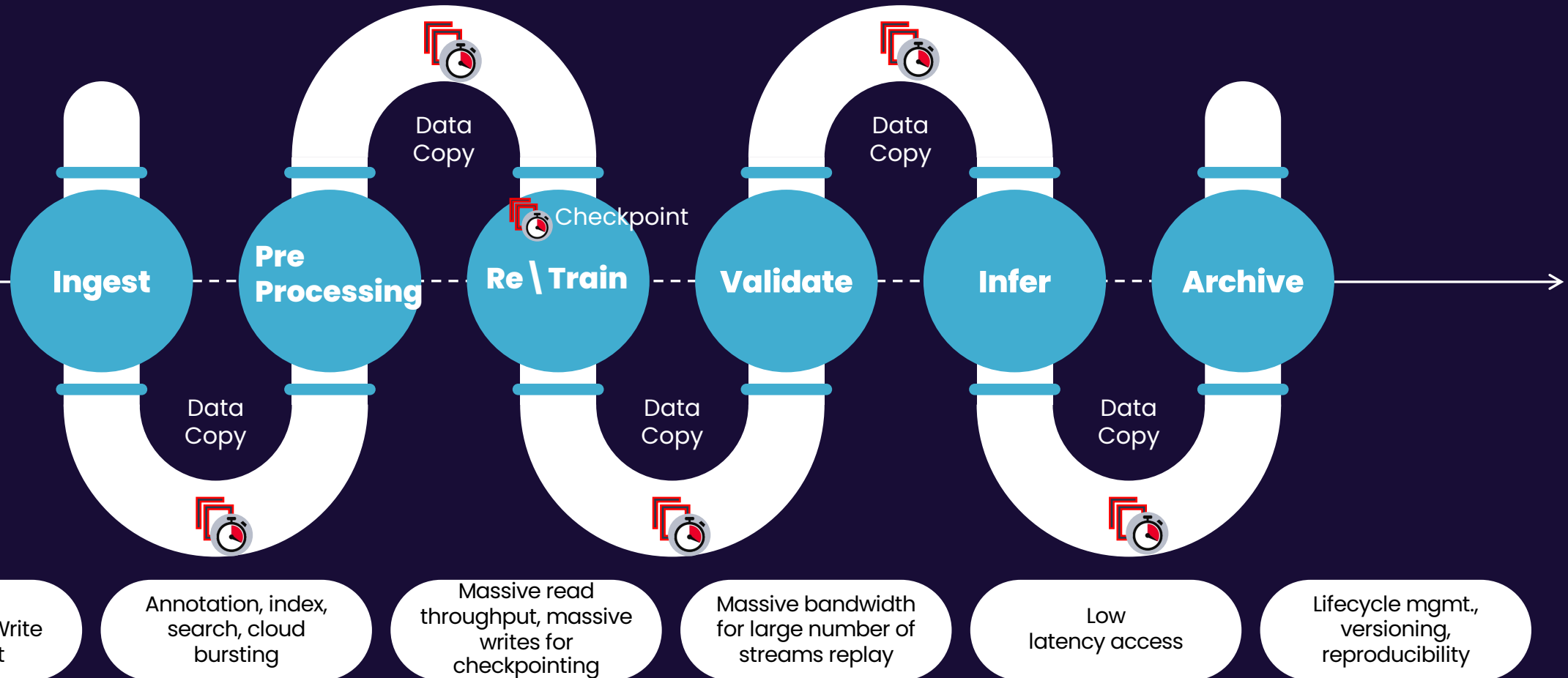


WEKA Data Platform

Modern Data Infrastructure

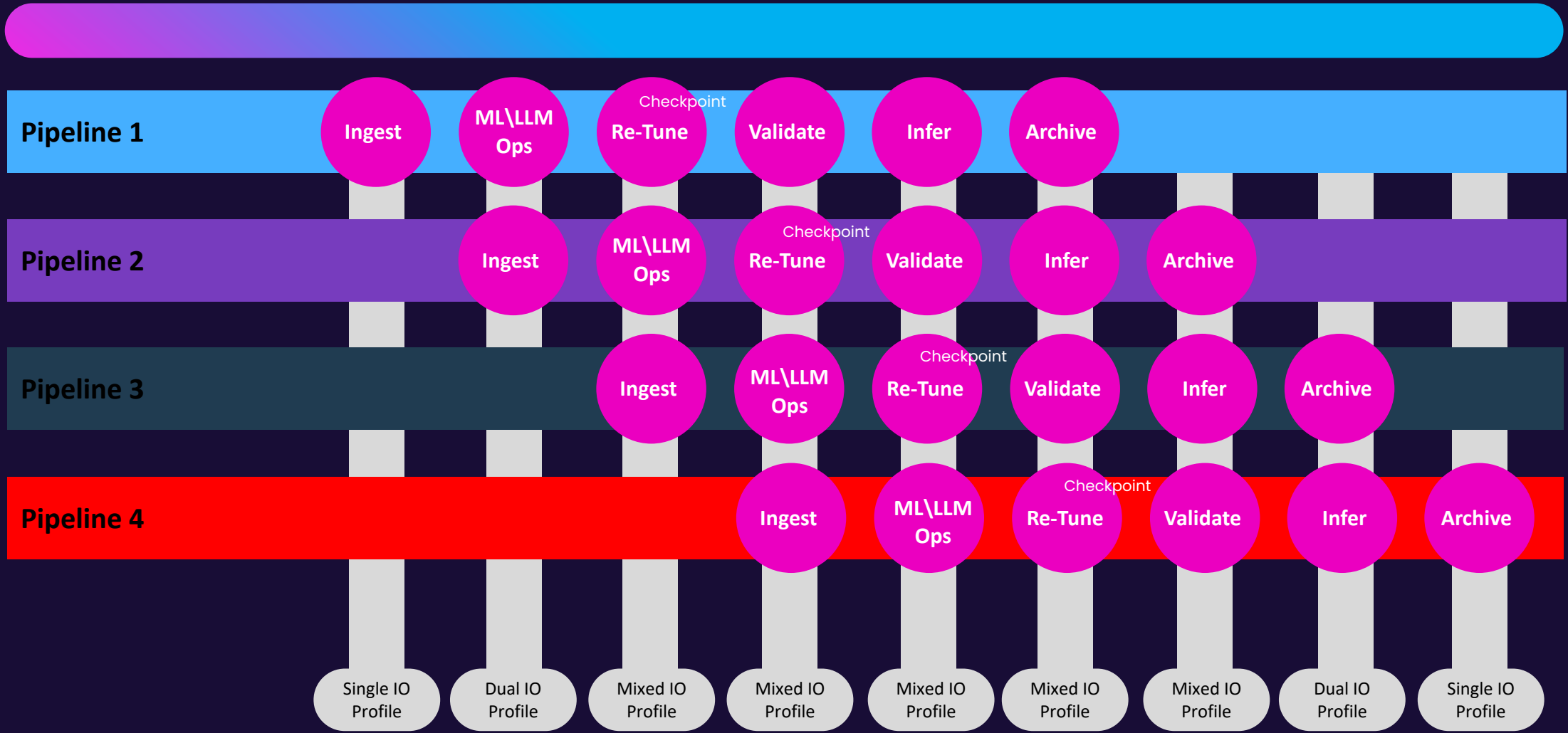
AI Pipelines Break Traditional Data Silos

Energy and time wasted on keeping copies



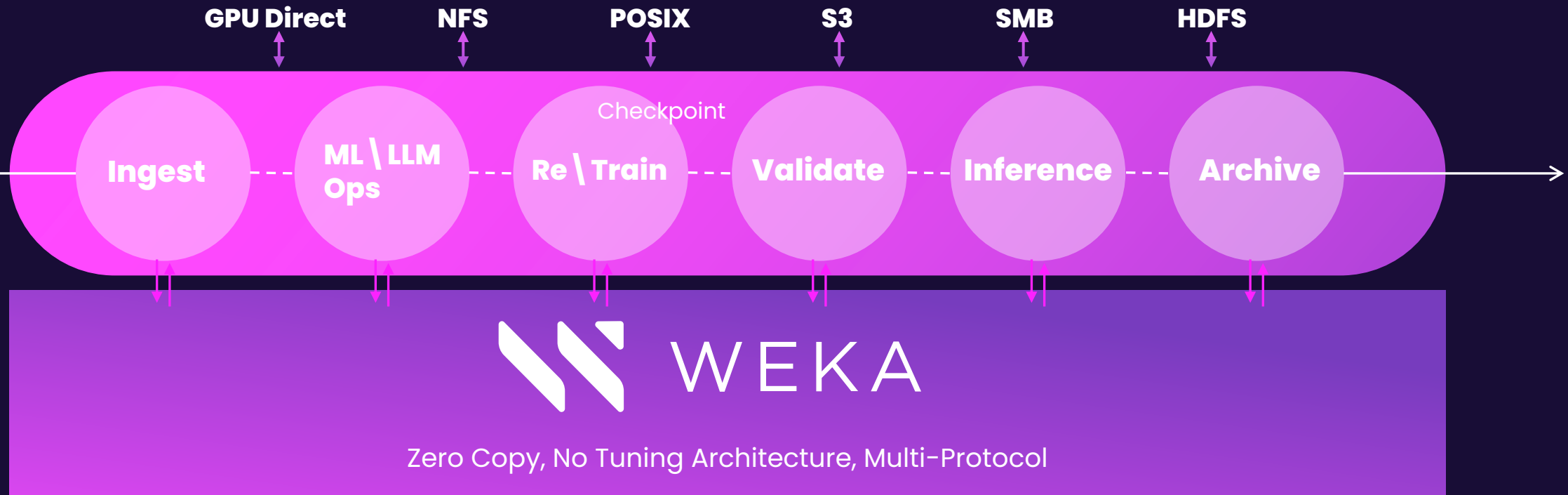
AI Data Pipeline: multiple pipelines heating storage

No tuning is possible



The WEKA Data Platform

No copying, faster performance leads to x10–100 improvement in time to epoch



Production Weka AI Systems Measured IO Patterns

AI customer #1 IO Pattern – Millions of Tiny IOs Reads / Writes



General / Cluster Summary

Last 2 Days

Last 2 days

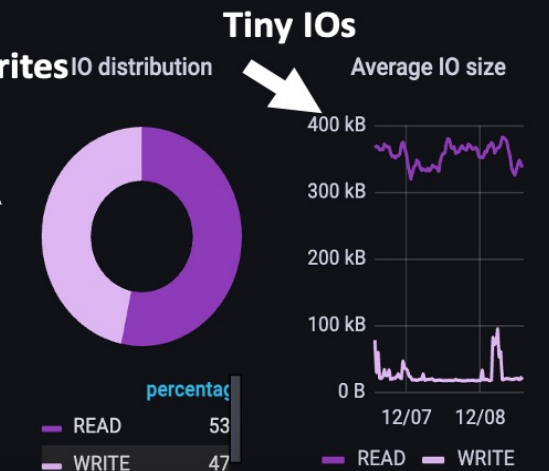
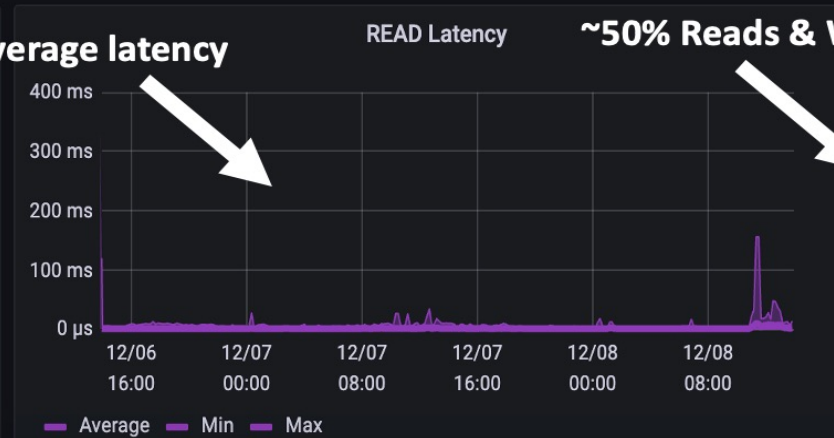
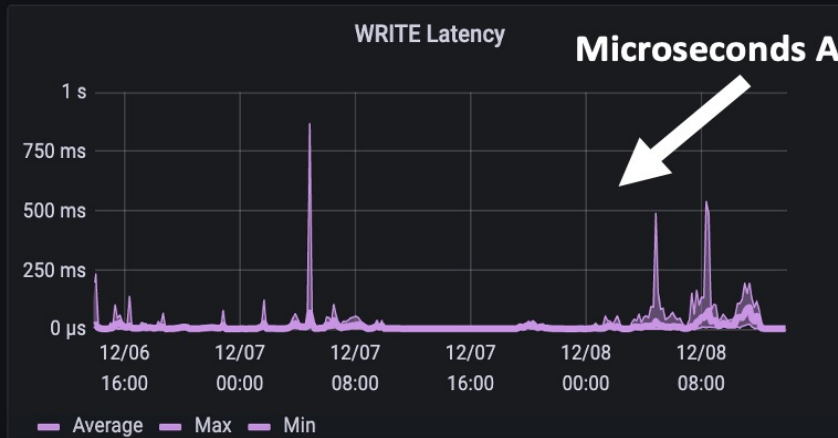
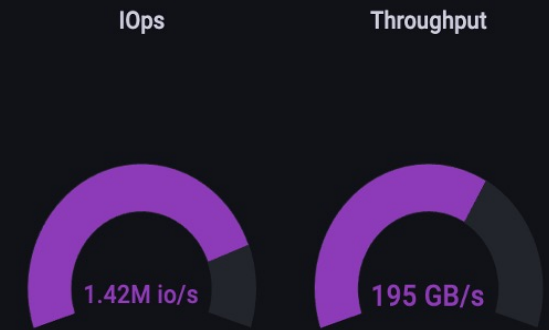
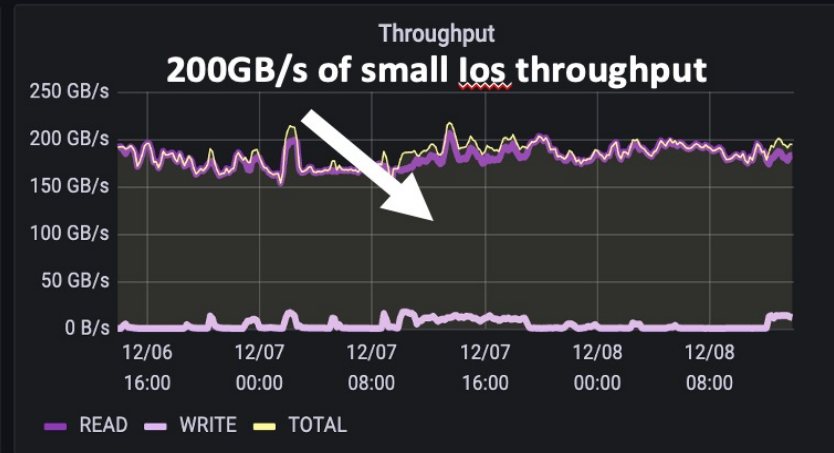
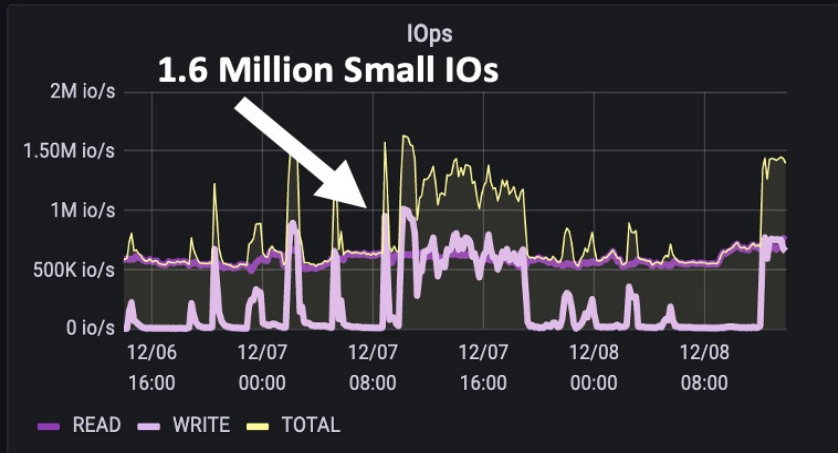


5m



Basic IO stats

VERSION	TOTAL HOSTS	BACKENDS	CLIENTS	DRIVES	TOTAL NODES	MGMT NODES	IO NODES	FE ROLES	COMPUTE ROLES	DRIVE ROLES
4.2.4.36	1078	59	1019	1027	5551	1078	3631	1920	1121	590



Checkpointing Example

Every 30–60 minutes

1000 GPUs X 80GB < 60 seconds

1.3TB/s Writes to shared storage

Inferencing Example

Model repositories of 100s TBs
Time to Load models to inference server
Time to output GenAI Artifacts
Accelerate VectorDB & Embedding

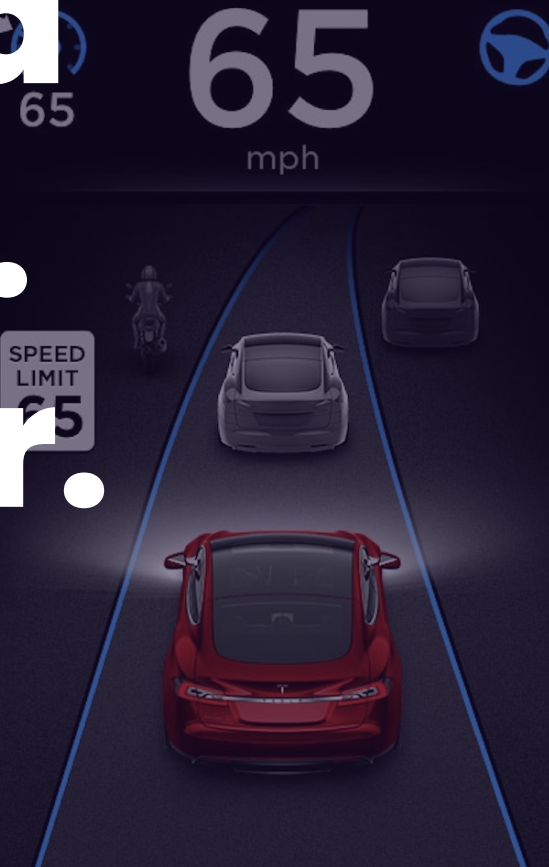


'Impossible' Workflows
powered by **WEKA**

Advanced Autopilot. 70x Faster.

e

184 mi 73°F



Radioactive
Night Visions
Imagine Dragons

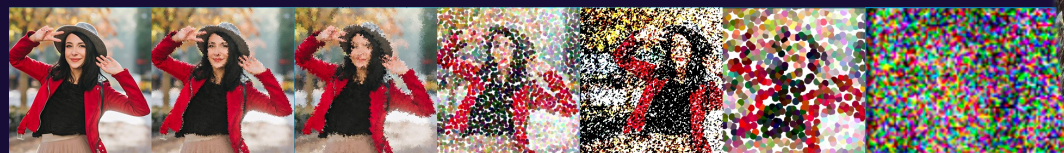
2:56

-3:14

10:30 PM

P R N D

AI by the people, for the people



stability-ai
stable-diffusion

A latent text-to-image diffusion model capable of generating photo-realistic images given any text input

- **Delivered 80%+ cost savings**
- **Increased GPU Utilization by 2.5x vs Lustre**
- **Increased data transfer 40x vs Lustre**
- **Running in Converged on GPU servers in AWS**



**One Year.
Two Vaccines.**

COVID-19
VACCINE

COVID-19
VACCINE

One Planet. Protected.



SLAC NATIONAL
ACCELERATOR
LABORATORY



U.S. DEPARTMENT OF
ENERGY



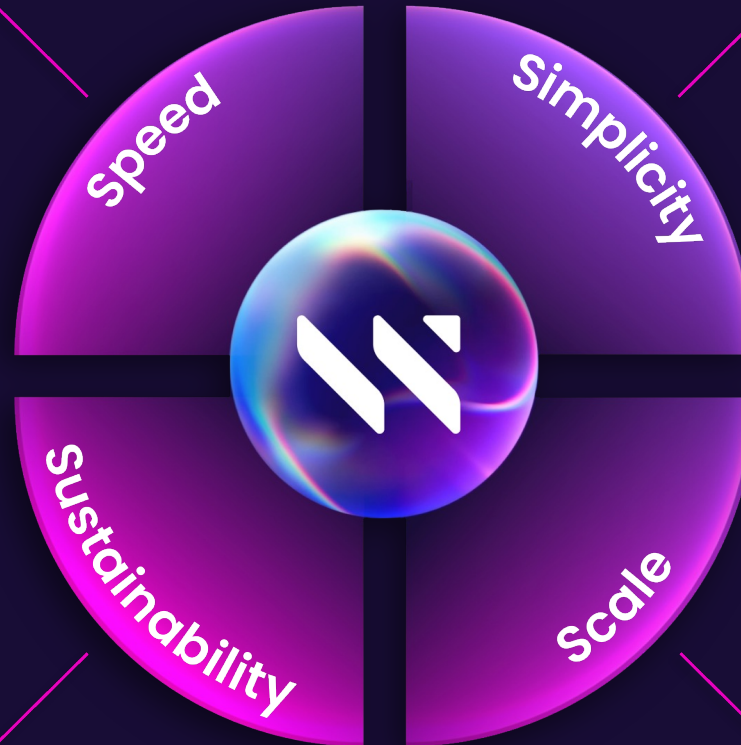
One Architecture Delivers on FOUR Promises

Mindbending **Speed**

Deliver unbeatable file and object performance for your most demanding applications supporting high I/O, low latency, small files, and mixed workloads with no tuning.

Seductive **Simplicity**

Eliminate the complexity and compromises of traditional data infrastructure with a single, easy-to-use data platform that eliminates storage silos across on-premises and the cloud.



Effortless **Sustainability**

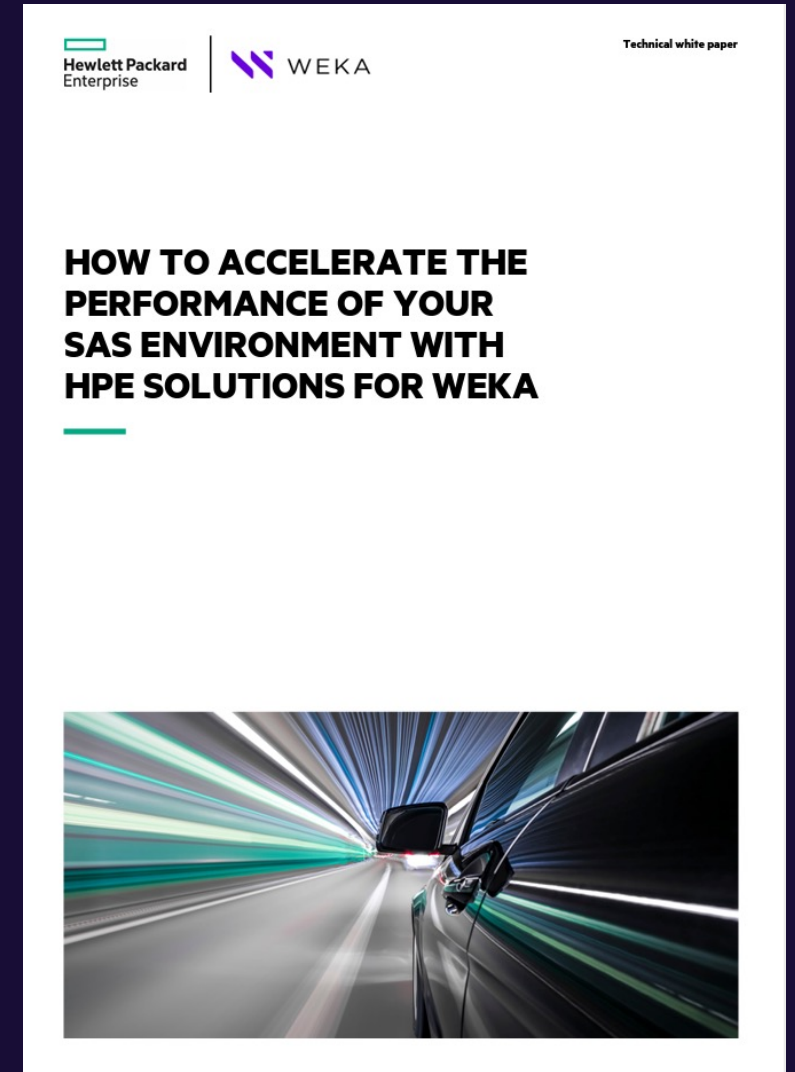
Lower energy consumption and reduce the resulting carbon emissions by cutting data pipeline idle time, extending the usable life of your hardware, and moving workloads to the cloud.

Infinite **Scale**

Scale your compute and storage independently and linearly on-premises or in the cloud with WekaFS to handle 10s of millions or even billions of files of all data types and data sizes.

Why HPE with weka?

- HPE is an early investor in WEKA (relationship since 2017) through the HPE Pathfinder Program with executive level relationships in place (Board level representation). HPE has participated in multiple funding rounds.
- WEKA is a Strategic Partner in the HPE Complete, the one-stop shop for validated HPE and third-party partner end-to-end infrastructure solutions.
- WEKA has won the HPE Technical Partner Program Momentum Partner of the year.
- Long-standing engineering and support relationships between HPE and WEKA have led to insights and domain-specific work in areas from AI for medical imaging to analytics.



Thank You