# Next-Generation NVMe-Native Parallel Filesystem for Accelerating HPC Workloads

Liran Zvibel

CEO, Co-founder WekaIO
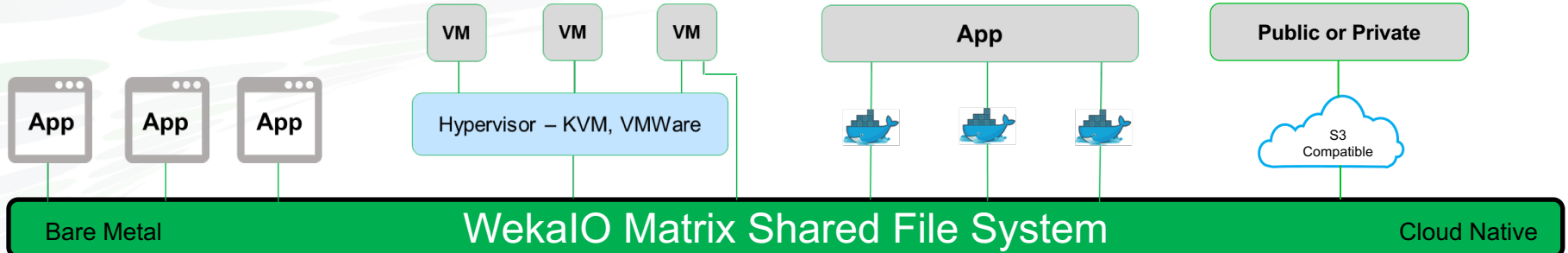
@liranzvibel

# WekaIO Matrix: Full-featured and Flexible

| VM | VM | VM |

| App |

| Public or Private |

Hypervisor – KVM, VMWare

| App | | App | | App |

S3 Compatible

**WekaIO Matrix Shared File System**

Bare Metal — Cloud Native

Fully Coherent POSIX File System That is Faster than a Local FS

Distributed Coding, Scale-out metadata, Fast Rebuilds, End-to-End DP

Instantaneous Snaps, Clones, Performance Tiering to S3, DR, Backup

InfiniBand or Ethernet, Hyperconverged or Dedicated Storage Server
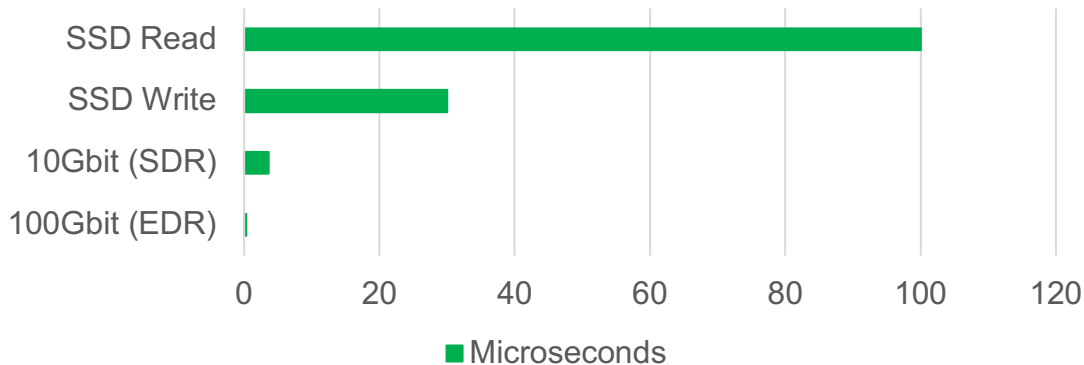
WEKA.IO

# Why NVMe-oF parallel FS?

- Local copy architectures were developed when 1GbitE and HDDs were standard
- Modern networks on 100Gbit are 100x faster than SSD
- It is much easier to create distributed algorithms when locality is not important
- 4KB IOs latency similar to local FS, bigger IOs parallelize, so even lower latency

### Time it takes to Complete a 4KB Page Move



Chart showing time (Microseconds) for:
- SSD Read: ~100
- SSD Write: ~30
- 10Gbit (SDR): ~4
- 100Gbit (EDR): ~0

X-axis: 0, 20, 40, 60, 80, 100, 120

Legend: ■ Microseconds

WEKA.IO

# Only PFS for NVMe-oF, PFS over S3
## Faster than burst-buffer + traditional PFS

- ○ Massive Scale
  - – Trillions of Files
  - – Billions of files per directory
  - – 100's of Petabytes
  - – Millions of IOPS
  - – 100's of GB of BW
- ○ Lowest latency FS, higher perf than AFA
- ○ HDD throughput similar to traditional PFS
- ○ Cloud Economics

WekaIO

AFA

SAN

All Flash NAS

Scale-out NAS | Scale-out Parallel NAS

Cloud Object Store

Performance

Speed

Simplicity

Scalability

Scale and Value

# Software Architecture – Keep out of kernel

- Runs inside LXC container for isolation
- SR-IOV to run network stack and NVMe in user space
- Provides POSIX VFS through lockless queues to WekaIO driver
- I/O stack bypasses kernel
- Scheduling and memory management also bypass kernel
- Metadata split into many Buckets – Buckets quickly migrate ➔ no hot spots
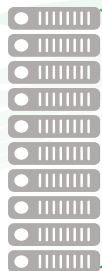- Support, bare metal, container & hypervisor



WEKA.IO
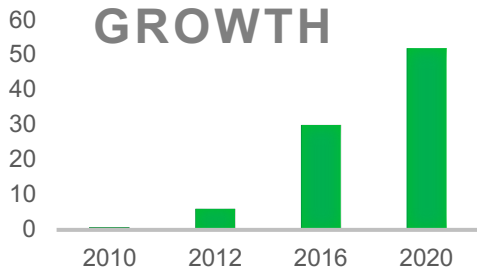
# Processing Has Shrunk while Data Sets Explode

10 CPU-Only Servers

1 GPU Accelerated Server

GPUs have shrunk compute infrastructure by 10x

**DATA GROWTH**

```
60
50
40
30
20
10
 0
    2010   2012   2016   2020
```
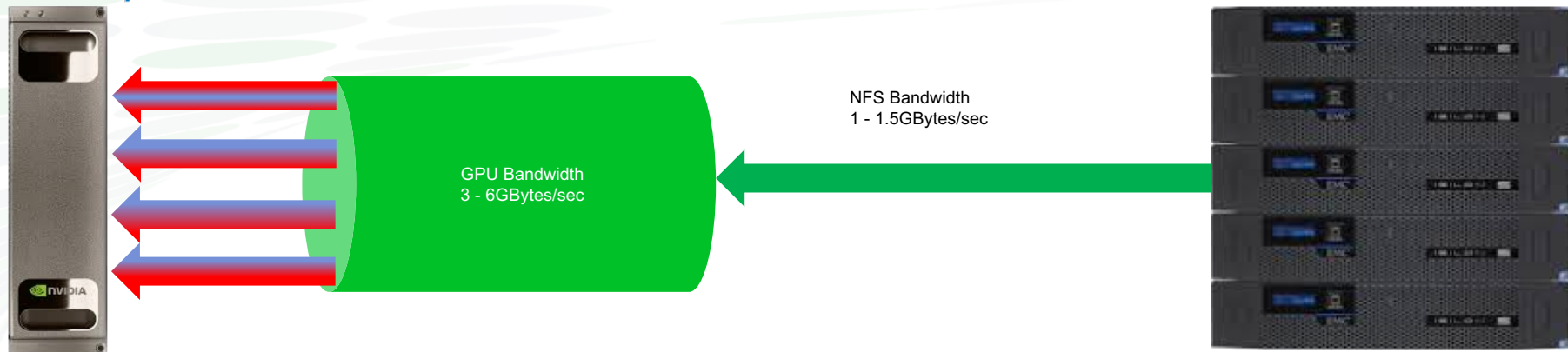
But the data that needs processing has grown 50x

Industry is cornered into an I\O nightmare

WEKA.IO

# NFS = Not For Speed

GPU Bandwidth
3 - 6GBytes/sec

NFS Bandwidth
1 - 1.5GBytes/sec

- A protocol developed in 1984 trying to solve a 2018 problem
- pNFS tried to fix NFS but failed when metadata workloads exploded
- Legacy parallel file systems like Lustre and GPFS cannot handle billions of small files
  - And they require a PhD to operate

WEKA.IO
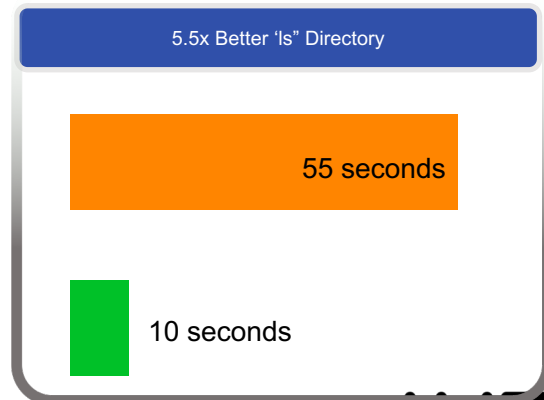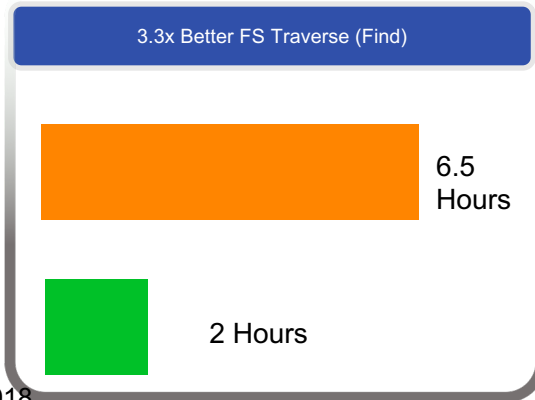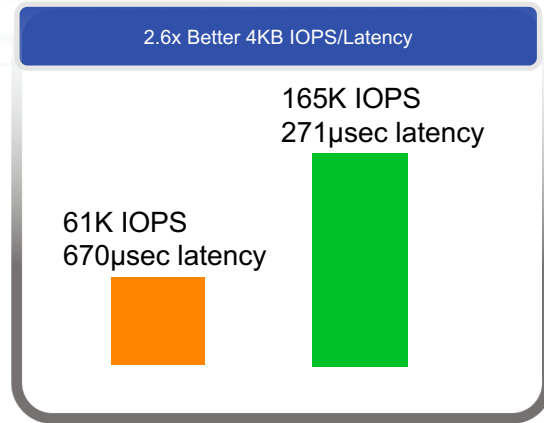
# WekaIO Solves the Data Accessibility Problem



GPU Bandwidth
3 - 6GBytes/sec

Per GPU Server 5-
44GBytes/sec
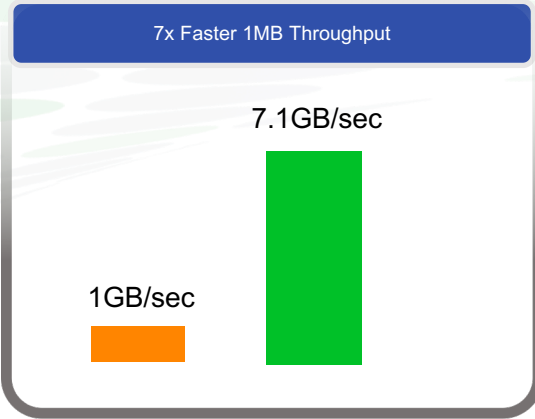
o Shared, Parallel file system written for NVMe

o POSIX Client runs on GPU Servers

o Cluster of servers provide high performance file services from NVMe

o Low latency networking on InfiniBand or Ethernet

o Training "data lake" stored on low cost object storage for best cost

WEKA.IO

# Actual Results from Bake-off vs All flash filer



**7x Faster 1MB Throughput**

7.1GB/sec

1GB/sec

**2.6x Better 4KB IOPS/Latency**

165K IOPS
271μsec latency

61K IOPS
670μsec latency

**3.3x Better FS Traverse (Find)**

6.5 Hours

2 Hours

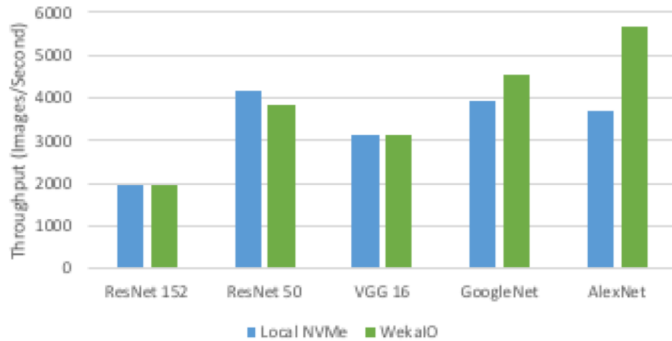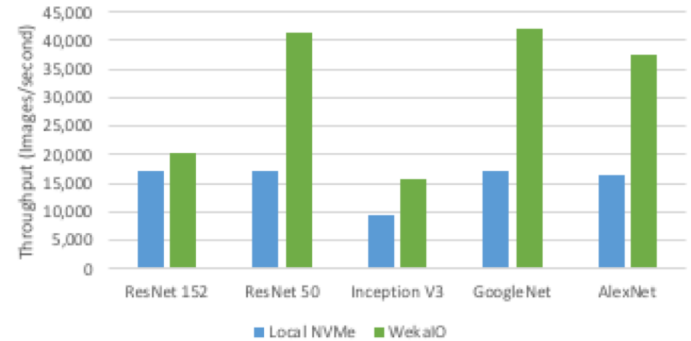**5.5x Better 'ls" Directory**

55 seconds

10 seconds

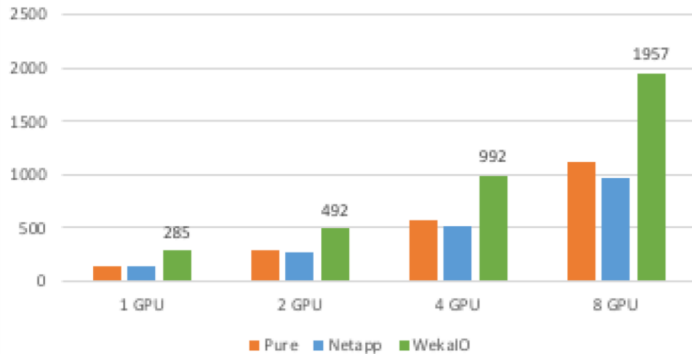WEKA.IO

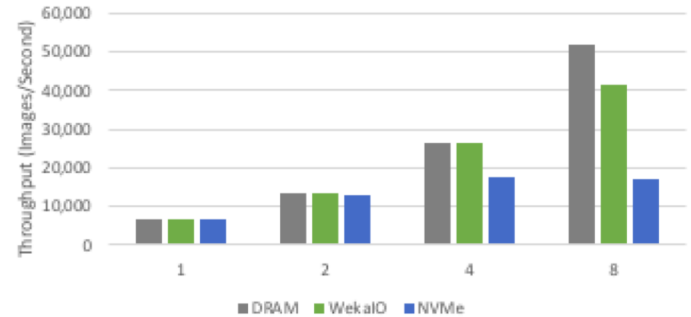# GPU Performance vs. Alternatives



Training Benchmarks vs Local NVME

Inference Benchmarks vs Local NVMe

Training Benchmark vs. Competition (ResNet 152)
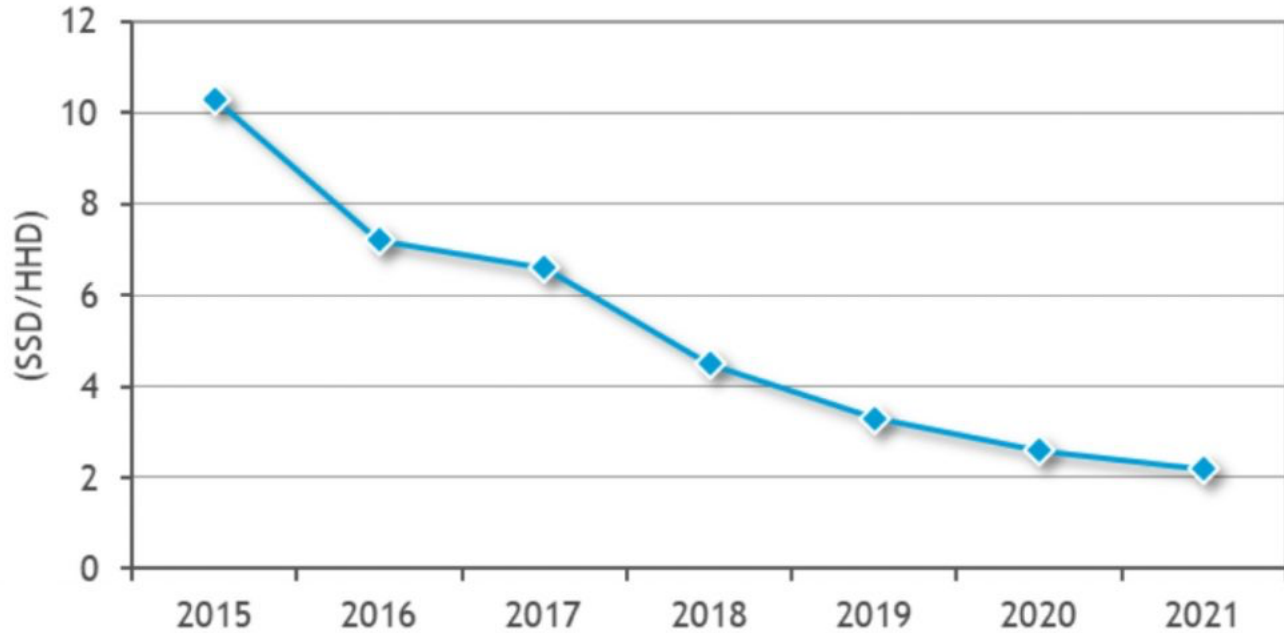
Inference Scaling vs DRAM and Local NVMe (ResNet 50)

# Deep Learning Requirements

o Actually very close to HPC problems…

o Store a vast amount of data
  - Effectively "**stage**" working set **back on fast storage**, for efficient access

o High bandwidth, **low latency**

o Very good **metadata performance**, traverse files quickly
  - Billions of files per directory, huge namespaces

o Very high **single host** performance

o Support **multiprotocol** (S3, HDFS, SMB, NFS)

WEKA.IO

# SSD vs HDD pricing (per gb ratio)



Source: Hyperion research
https://www.storagenewsletter.com/2018/08/07/flash-storage-trends-and-impacts/

# HPC only cares about throughput, right?

o NAND is cheaper for IOPS (and obviously latency) for several years now

o HDD stats: 160MB/sec ; $0.02/GB capacity for 10TB devices

o 3.84TB TLC devices read at 1700MB/sec ; so faster than 10 HDDs

- Total HDD cost needed to read at 1700MB/sec → $2000; avg per NAND device $0.52/GB

- Already cheaper today!

o 7.68TB QLC devices coming next year writing at 1000MB/sec; 6 HDDs needed

- Total HDD cost needed to read at 1000MB/sec → $1200; avg per NAND device $0.16/GB

- Next year QLC will be cheaper for write throughput

- Endurance will probably not hold for checkpointing; but anywyas small capacity that TLC makes sense for

WEKA.IO

# Future of HPC storage is NAND FLASH

o Currently HDDs still make sense for some workloads

o In a year (and obviously later) HPC storage should steer towards NAND FLASH technologies

o Parallel FS for NVMe-oF require different data structure, and algorithms based on modern workloads (scaling metadata, small IOPS, etc)

o HPC applications should consider NAND FLASH only for active workload; other media (tape;optics; etc) for archival capacity

WEKA.IO