

Gartner

Cool

WekalO: Extracting 2018 Performance from Modern HPC Hardware

Derek Burke, Regional Account Manager Chris Weeden, Senior Systems Engineer

Company Overview





WekalO Matrix is the fastest, most scalable parallel file system for AI and technical compute workloads, on premises and in the public cloud

WEKA.iO

WekalO Confidential

Single Namespace for Multiple Workloads



WEKA.iO

WekalO Confidential

Focused On Performance Use Cases



- · Genomics sequencing and analytics
- Drug discovery
- Microscopy



- Machine Learning\ Al
- Real-time analytics
- IOT



- Algorithmic trading
- Business analytics (SAS Grid)
- Risk analysis (Monte Carlo simulation)



- Seismic
- Reservoir simulation
- Analytics



- Media rendering
- Transcoding
- Visual Effects (VFX)



- CFD, Simulations
- EDA

WEKA.IO

WekalO Confidential

We Make a Bold Claim...

World's Fastest File System^w

From one rack of commodity hardware with 100GbE Network Approximately:

- 400GB/s
- 30 Million IOPS

WEKA.IO

SPEC SFS2014 vda Results

Test is 90% write intensive



Lower and Longer Are Better

SPEC SFS2014 eda Results

SPEC SFS2014 eda Benchmark 1.6 1.4 1.2 1.0 Average Late ncy (ms) 0.8 .45ms latency 0.6 0.4 0.2 200 400 1,000 600 800 1,200 1,400 1,600 1,800 2,200 2,000 Number of Simultaneous Job Sets DDN GridScaler Oracle ZFS ZS7-2 —WekalO 3.1.9

WEKA.IO

Lower and Longer Are Better

(7)

Frontend is 50% stats test Backend is 50% R and W

SPEC SFS2014 vdi Results

Test is 60% random write 20% random read



WEKA.IO

Lower and Longer Are Better

SPEC SFS2014 db Results

Test is 80% random read 20% random write



WEKA.IO

Lower and Longer Are Better

How do we do it...

World's Fastest File System^{**}

WEKA.iO



Why Data Locality is Irrelevant

- Modern networks on 10Gbit Ethernet are 30 times faster than SSD for reads and 10 times faster than SSD for writes
- With right networking stack, shared storage is faster than local storage



Time it takes to Complete a 4KB Page Move

WEKA.iO



The Kernel is designed for multi-tasking...

- The system needs to be notified of the new packet and pass the data onto a specially allocated buffer sk_buff struct (Linux allocates these buffers for every packet).
- To do this, Linux uses an interrupt mechanism: an interrupt is generated several times when a new packet enters the system. The packet then needs to be transferred to the user space.
- As more packets have to be processed, more resources are consumed negatively impacting the overall system performance.
- sk_buff struct: the Linux network stack was originally designed to be compatible with as many protocols as possible. As such, metadata for all of these protocols is included in the sk_buff struct, but that's simply not necessary for processing specific packets. Because of this overly complicated struct, processing is slower than it could be.
- When an application in the user space needs to send or receive a packet, it executes a system call. The context is switched to kernel mode and then back to user mode. This consumes a significant amount of system resources

WEKA.iO

Context Switching is wasteful...



WEKA.IO

(13)

DPDK by-passes the Kernel...



WEKA.iO

(14)

Components of DPDK:

•Environment Abstraction Layer (EAL) : It is responsible for gaining access to low-level resources such as hardware and memory space.

•Memory Manager: Responsible for allocating pools of objects in memory. A pool is created in huge page memory space and uses a ring to store free objects.

•Buffer Manager: Reduces by a significant amount the time the operating system spends allocating and de-allocating buffers using advanced techniques such as Bulk Allocation, Buffer Chains, Per Core Buffer Caches etc.

•Queue Manager: Implements safe lockless queues, instead of using spinlocks, that allow different software components to process packets, while avoiding unnecessary wait times.

•Packet Flow Classification: DPDK Flow Classifier implements hash based flow classification to quickly place packets into flows for processing.

•Poll Mode Drivers: Instead of using Interrupts and wasting CPU attention, PMD uses polling (scanning the NIC whether packets arrived or not), and doesn't disturb the CPU at all!

WEKA.iO

We also use SR-IOV and SPDK

Multiple Virtual NICs

Benefits of SPDK



WEKA.IO

(16)

Software Architecture

Runs in LXC Container for isolation

Kernel Module

VFS Interface

Front End

- POSIX Client
- Other protocol access

Back End

- Data placement
- Data protection
- File system metadata
- Tiering

Networking

- SRIOV for network stack
- I/O bypasses network stack

Storage Agent

I/O bypasses the kernel



WEKA.IO

WekalO Data Path

- Application IO (file operations)
 - Access WekalO Client as Local FS
 - User-Space, Low-Latency
 - POSIX-complete, high-perf
 - Kernel Module for VFS integration
 - Client-side NFS
 - Bottlenecked by Kernel
 - Handled by WekalO's Front End



WEKA.İO

WekalO Confidential

(18

WekalO Data Path

- Application IO (file operations)
 - Access WekalO Client as Local FS
 - User-Space, Low-Latency
 - POSIX-complete, high-perf
 - Kernel Module for VFS integration
 - Client-side NFS
 - Bottlenecked by kernel
 - Handled by WekalO's Front End
- WekalO Front-Ends are Cluster-Aware
 - Incoming Read Requests optimized re Location & Loading Conditions
 - Incoming Writes can go anywhere
 - Metadata fully distributed
 - No redirects required
- SR-IOV optimizes Network access



WEKA.iO

WekalO Data Path

- Application IO (file operations)
 - Access WekalO Client as Local FS
 - User-Space, Low-Latency
 - POSIX-complete, high-perf
 - Kernel Module for VFS integration
 - Client-side NFS
 - Bottlenecked by kernel
 - Handled by WekalO's Front End
- WekalO Front-Ends are Cluster-Aware
 - Incoming Read Requests optimized re Location & Loading Conditions
 - Incoming Writes can go anywhere
 - Metadata fully distributed
 - No redirects required
- SR-IOV optimizes Network access
- WekaIO directly accesses NVMe flash
 - Bypassing kernel, better perf



WEKA.İO

File System Scales Linearly with Cluster Size



Test Environment – 30-240 R3.8xlarge cluster, 1 AZ, utilizing 2 cores, 2 local SSD drives & 10GB of RAM on each instance. About 5% of CPU/RAM.

WEKA.IO

WekalO Confidential

WekalO Matrix #1 File System on SPEC



| Benchmark | #1 | Score | ORT | #2 Position | Score | ORT | |
|--------------------|----------|-------|------|-------------|-------|-------|--|
| | Position | | (ms) | | | (ms) | |
| Software build | WekalO | 5700 | 0.26 | NetApp | 4200 | 0.78 | |
| Database | WekalO | 4480 | 0.34 | Oracle | 2240 | 0.78 | |
| Engineering design | WekalO | 2000 | 0.48 | Oracle | 900 | 0.61 | |
| Video streams | WekalO | 6800 | 1.56 | DDN | 3400 | 50.07 | |
| Virtual desktop | WekalO | 1600 | 0.48 | DDN | 800 | 2.58 | |

WEKA.iO

#1 File System on the IO500 Test

Test measures bandwidth and metadata

https://www.vi4io.org/io500/list/19-01/10node

| # | information | | | | | | | | io500 | | |
|---|-------------------------------|-----------|----------------|-----------------|--------------|--------------------|------|--------------|-------|--------|--|
| | institution | system | storage vendor | filesystem type | client nodes | client total procs | data | <u>score</u> | bw | md | |
| | | | | | | | | | GiB/s | kIOP/s | |
| 1 | WekalO | | WekalO | | 10 | 700 | zip | 58.25 | 27.05 | 125.43 | |
| 2 | Oak Ridge National Laboratory | Summit | IBM | Spectrum Scale | 10 | 160 | zip | 44.30 | 9.84 | 199.48 | |
| 3 | DDN | Bancholab | DDN | Lustre | 10 | 240 | zip | 31.50 | 6.33 | 156.69 | |

- 31% better than
 World's largest
 Supercomputer
- 85% better than Lustre



WEKA.iO

WekalO Confidential

Thank you.

WEKA.io • •

https://docs.weka.io https://start.weka.io