

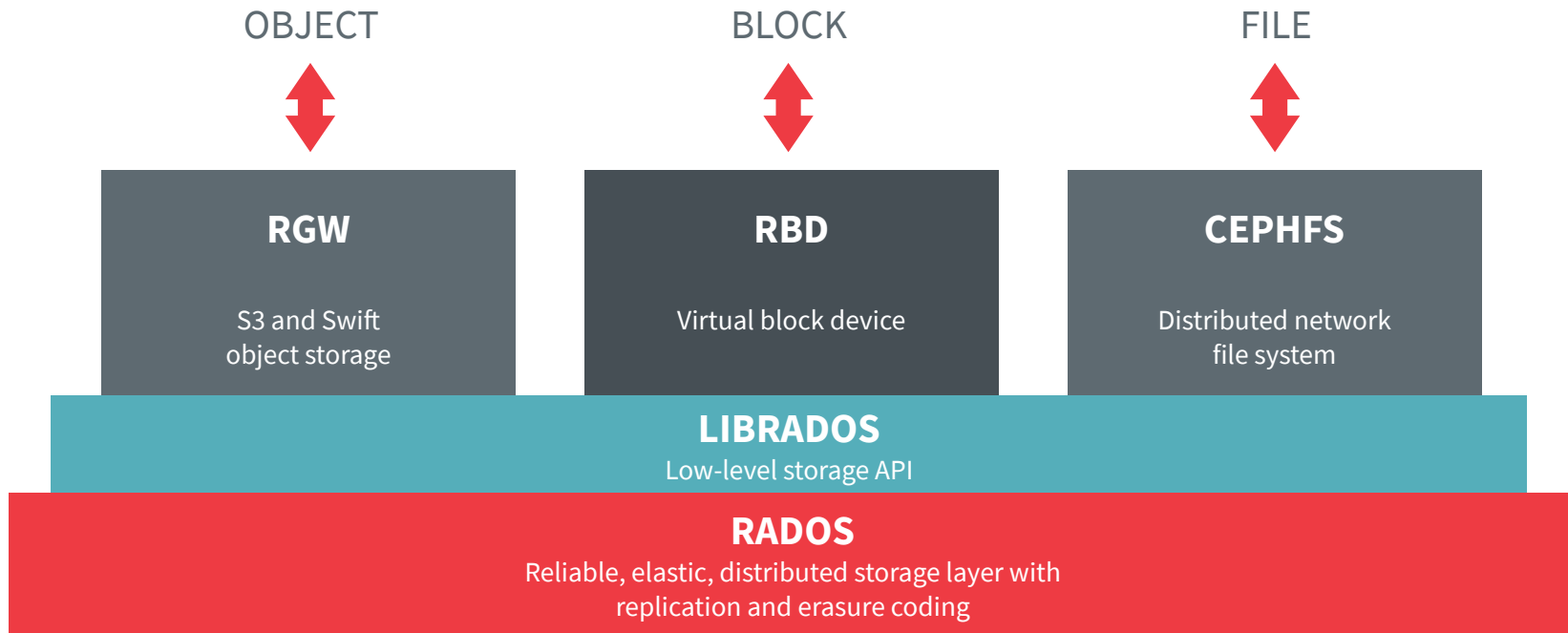
Update on Crimson

The Seatarized Ceph

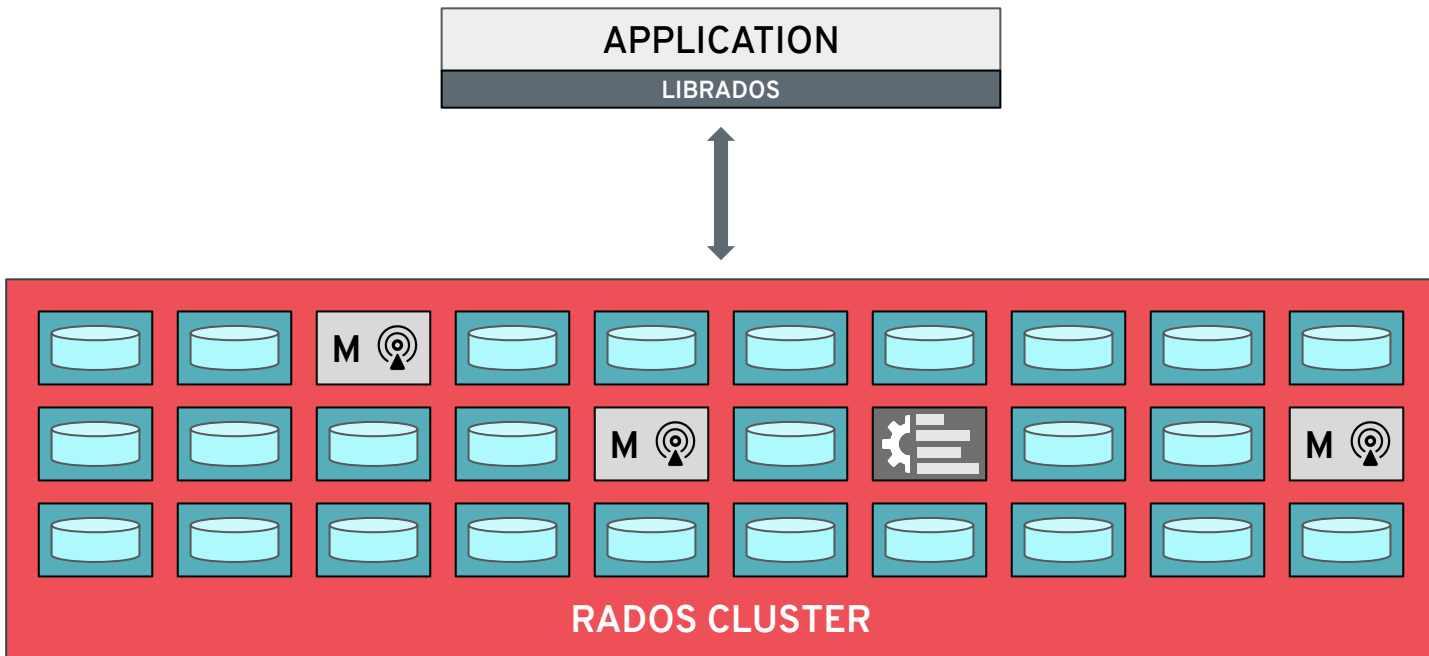
kchai@redhat.com

Seastar Summit 2019

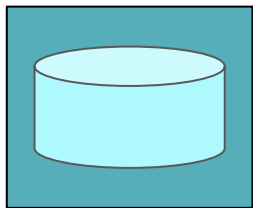
A unified storage system



RADOS -- The Cluster



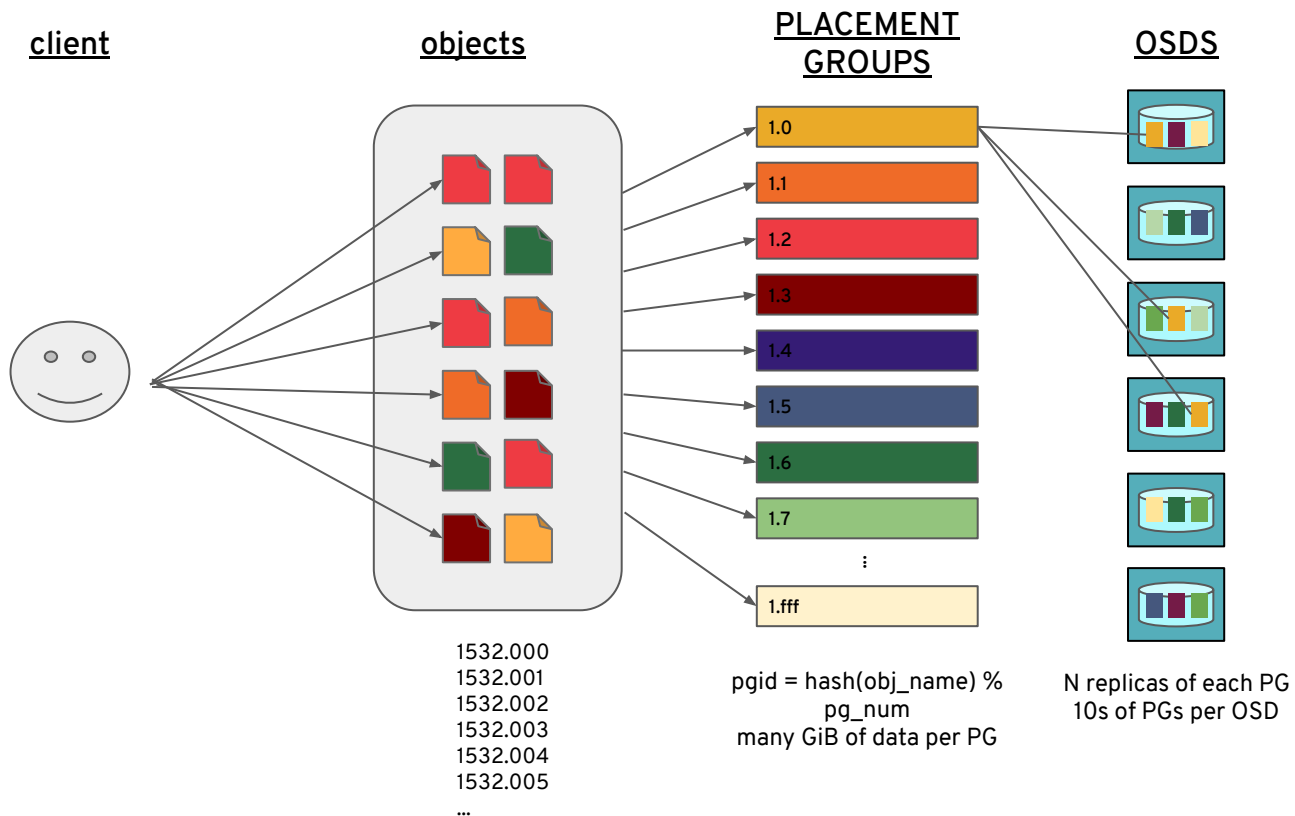
OSD (Object storage Daemon)



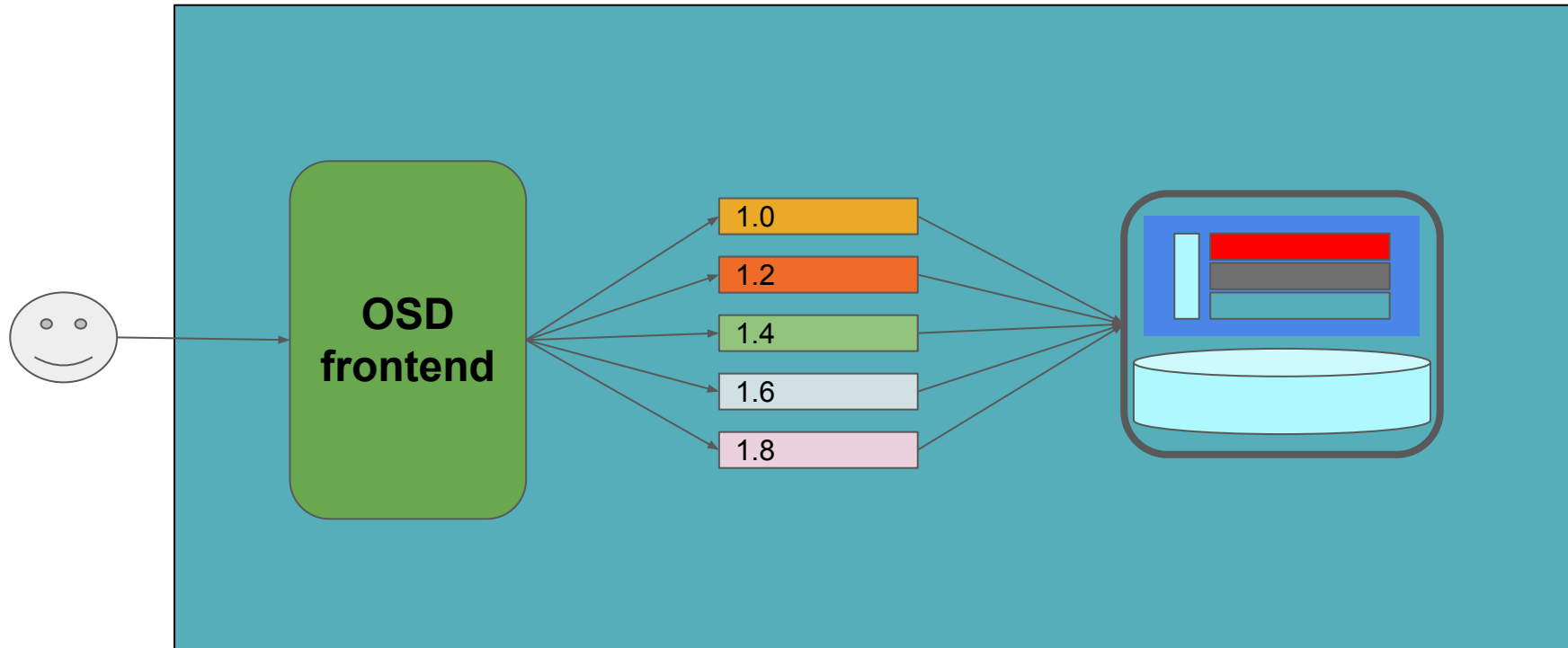
ceph-osd
(crimson)

- Stores data on an HDD or SSD
- Services client IO requests
- Cooperatively peers, replicates, rebalances data
- Reports stats to manager daemons
- 10s-1000s per cluster

PG (placement groups)



A closer look



Crimson - a **faster** OSD

- Less overhead
 - Bypass kernel
 - Zero memcpy
 - Less context switches
- Understands modern storage devices

share something => share nothing

What we imaged:

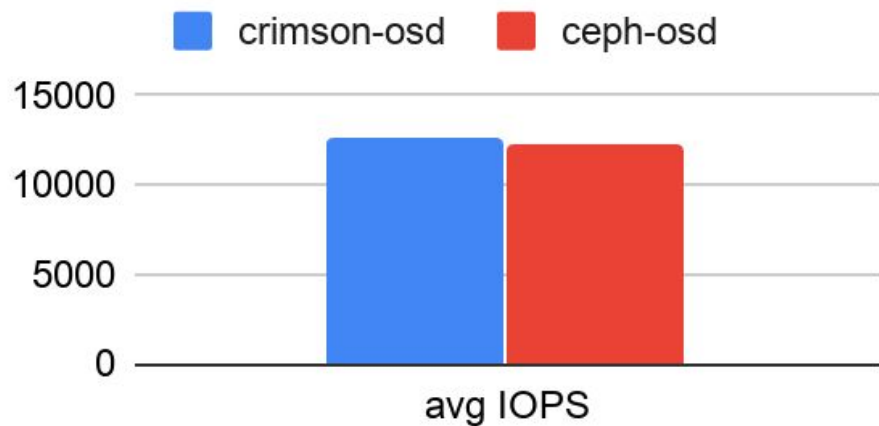
- Multi-reactor OSD
- Shared connections
 - Connections to manager daemons
 - Connections to peer OSDs
 - Connections to clients
- Shared io queue
- Shared metadata
 - Knowledge about the cluster

What we have now:

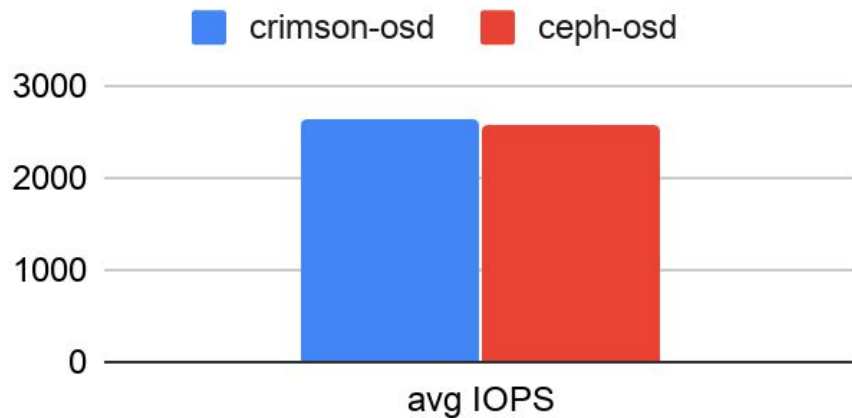
- Single threaded OSD
- Fully connected network
- Monitor's load increases

Average IOPS

random read

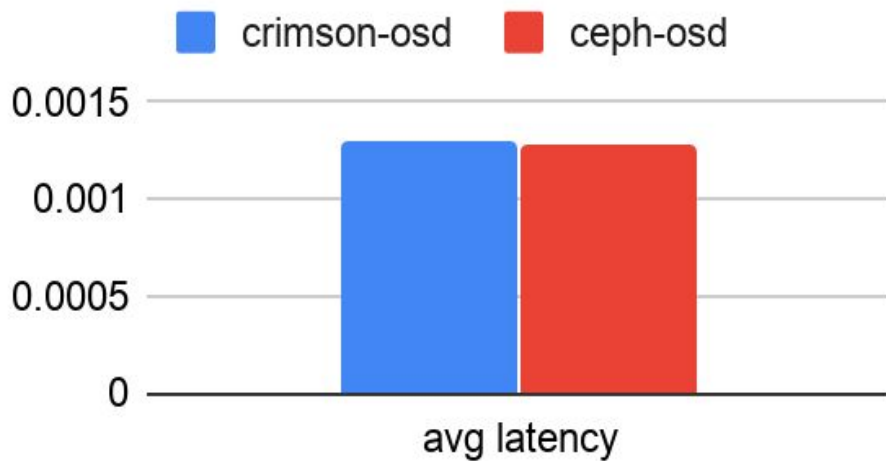


seq write

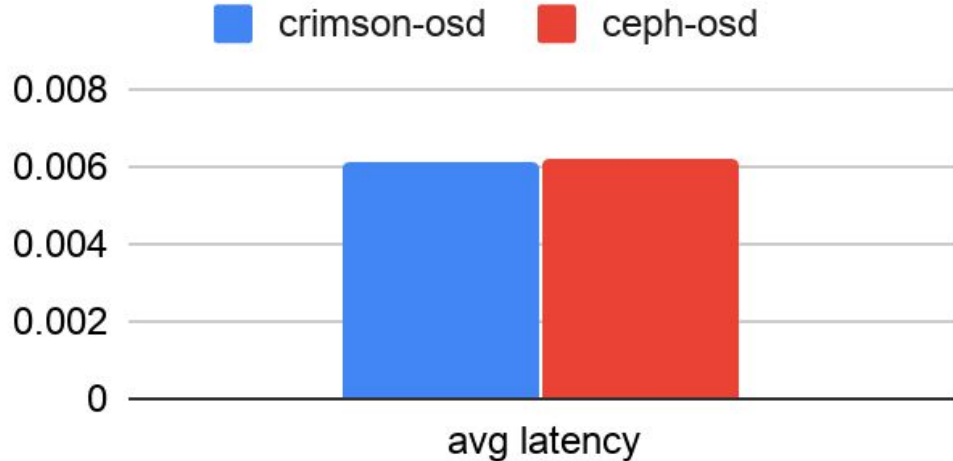


Average latency

random read

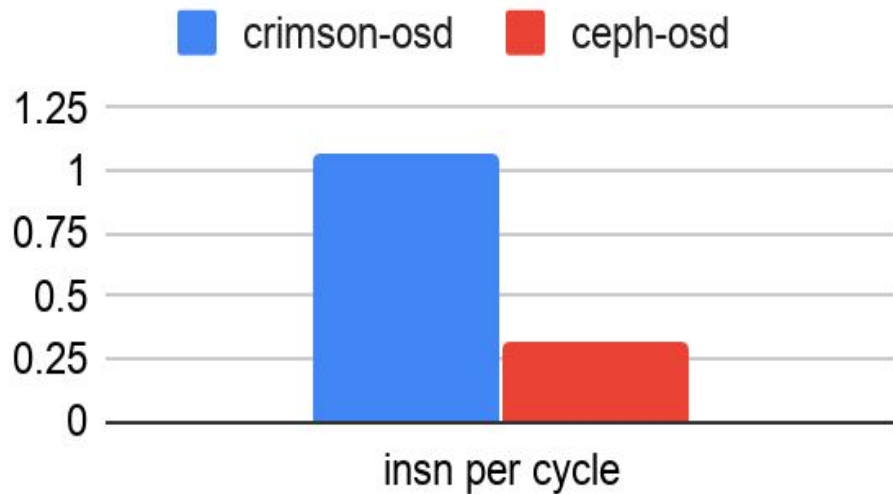


seq write

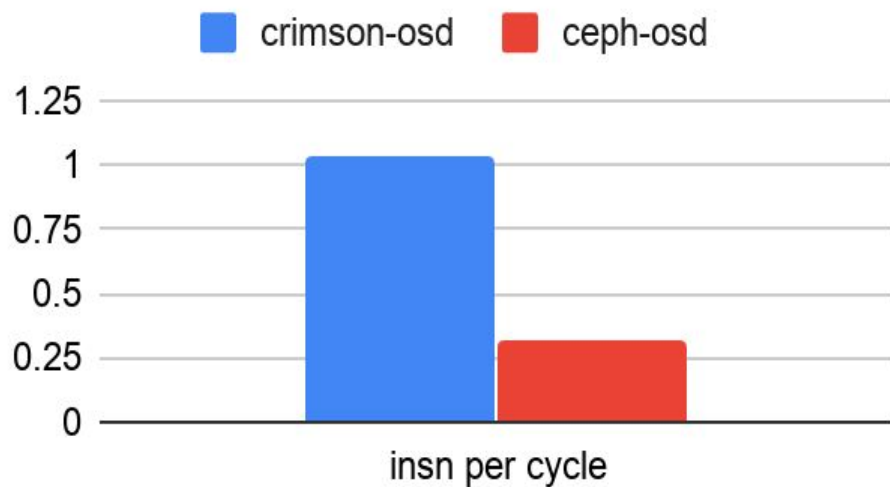


Instructions per cycle

random read

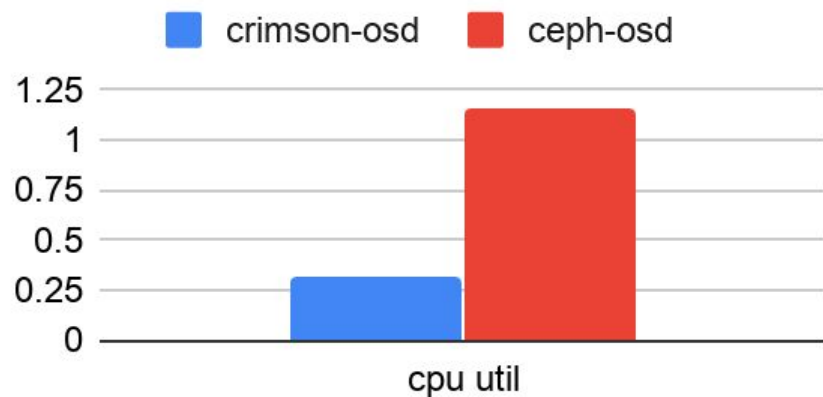


seq write

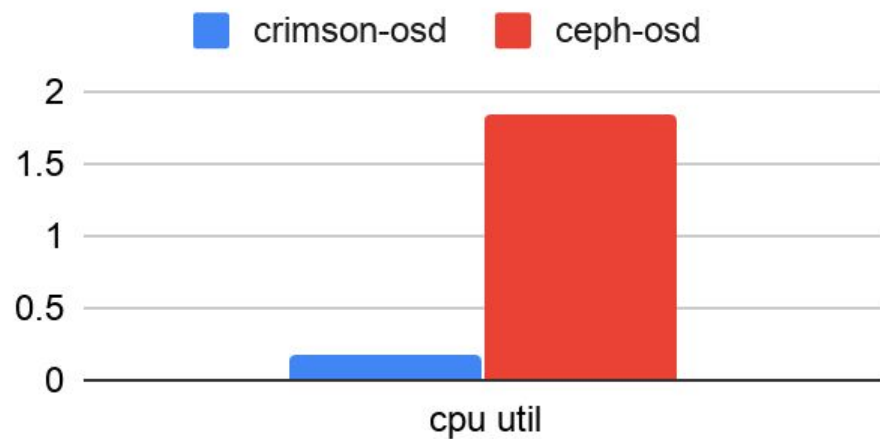


CPU util

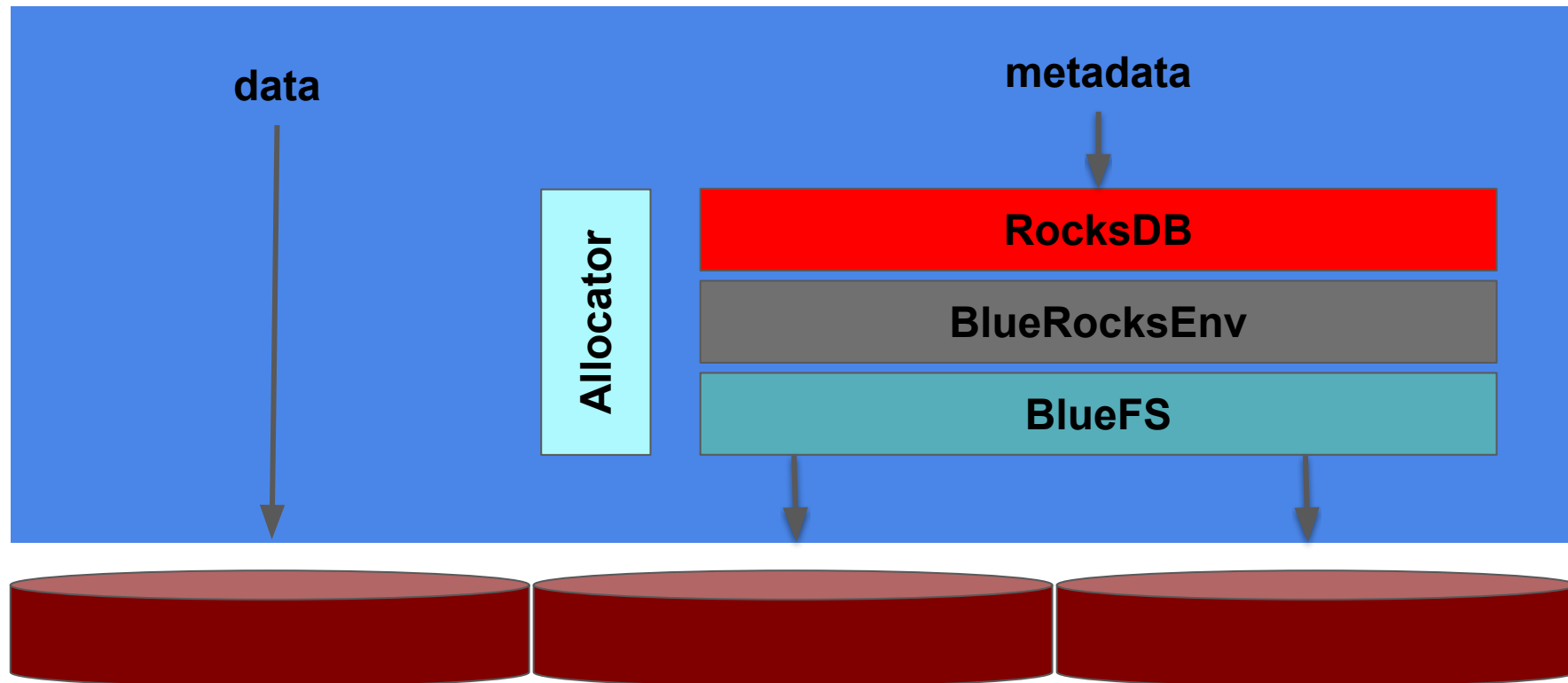
random read



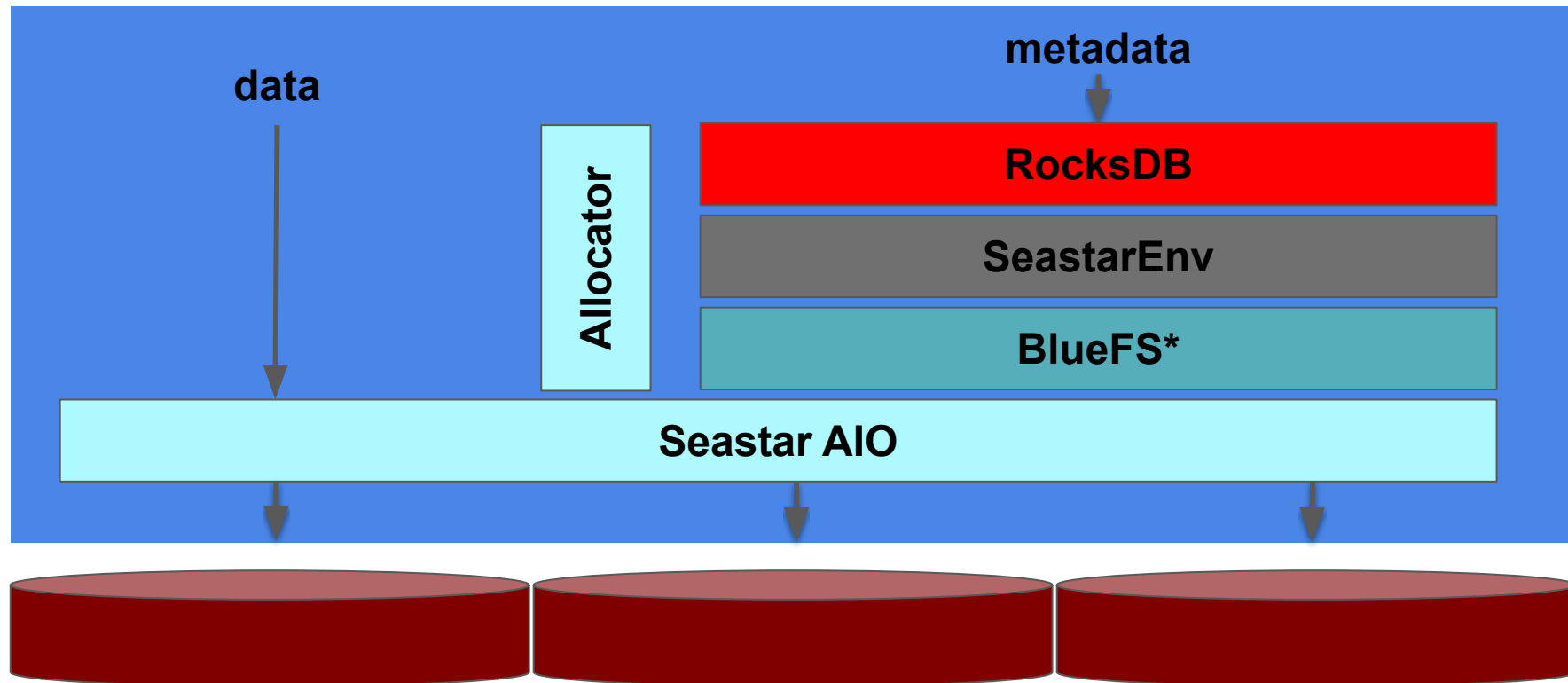
seq write



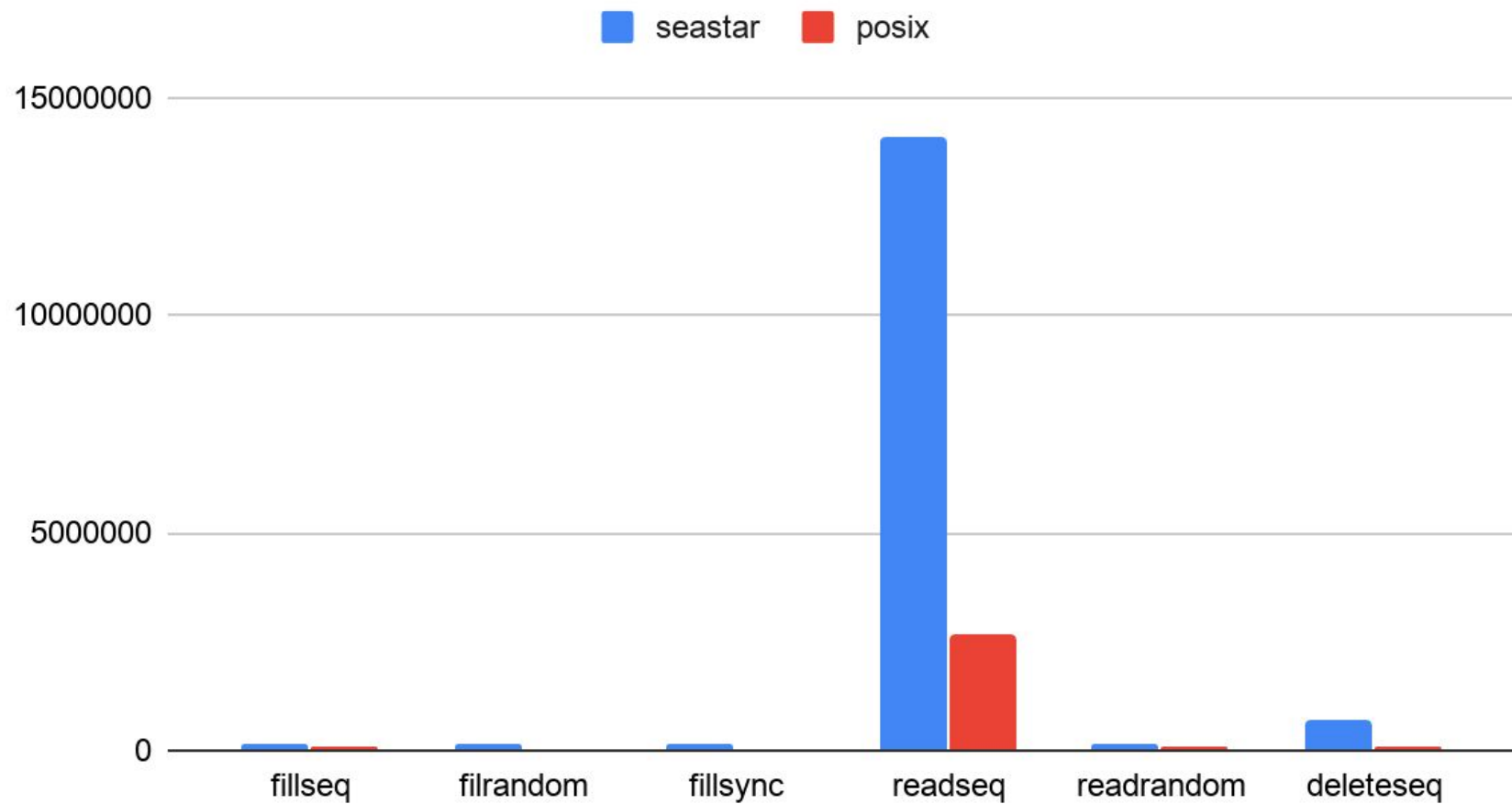
bluestore



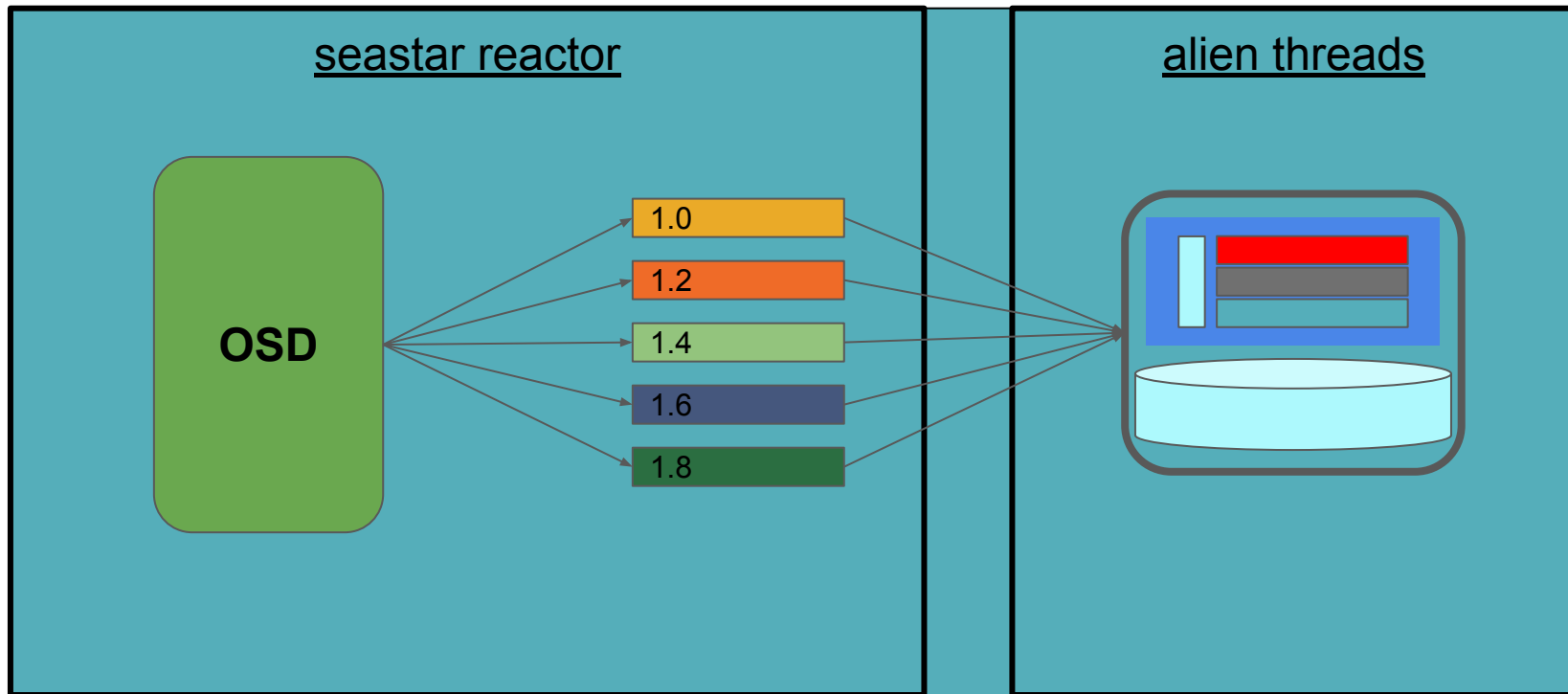
Seastarized bluestore



db_bench (ops/secs)



Alienized bluestore



SeaStore



???

The diagram illustrates the SeaStore architecture. It features a large blue rectangular area representing the storage layer, with the text '???' centered within it. Below this blue area is a dark red, horizontally-oriented oval shape representing a disk or storage medium.

Q & A