

IOCost: Block IO Control for Containers in Datacenters

Tejun Heo, **Dan Schatzberg**, Andrew Newell, Song Liu, Saravanan Dhakshinamurthy, Iyswarya Narayanan, Josef Bacik, Chris Mason, Chunqiang Tang, ‡Dimitrios Skarlatos

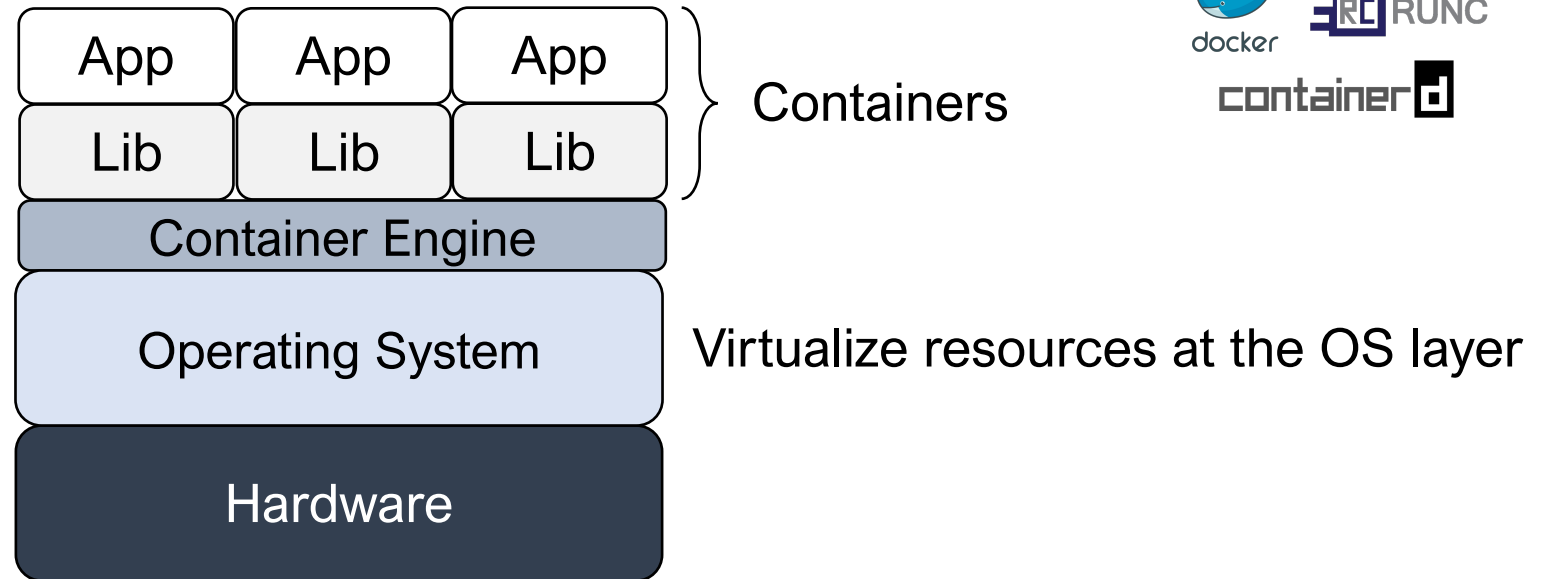
Session 5B: Data Center and Cloud Services
Thursday, March 3 @ 2pm

ASPLOS 2022

New Era in Datacenter Computing

Meta fleet runs on containers

- Lightweight
- Fast bring up
- Higher consolidation

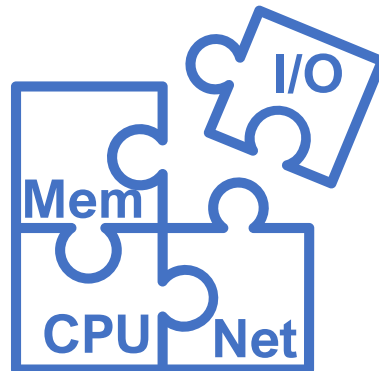


Problem: Resource Isolation

Ensure sufficient resources despite colocated workloads

Well understood problem for CPU, Memory, Network

Missing piece: Block I/O for storage devices (SSDs, HDDs)



Resource Isolation

Why do we need Block I/O?

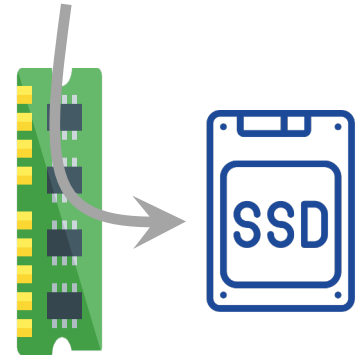
Storage Systems



System Updating
Package Fetching



Memory Offloading



TMO: Transparent Memory
Offloading @ASPLOS 2022

Challenge A → Hardware Heterogeneity

Multiple generations of SSDs, HDDs, Local/Remote

Different performance characteristics within each type

Idiosyncrasies of each device and hardware features



Challenge B → Workload Heterogeneity

IO control needs to cater to a wide variety of workloads



Databases



Video



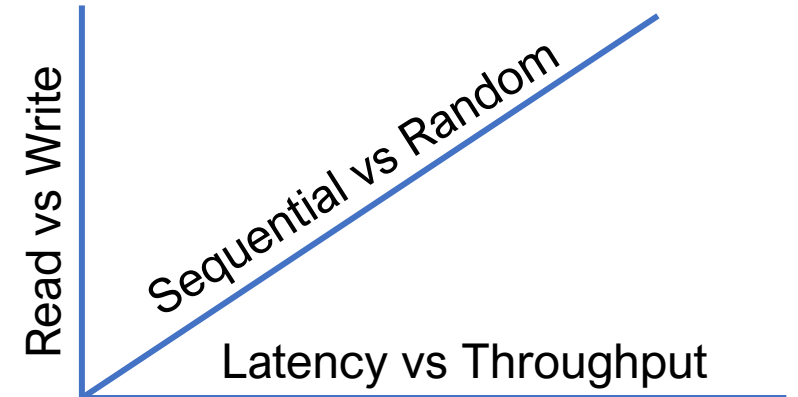
Web



Functions

Varying characteristics:

- Latency vs Throughput
- Read vs Write
- Sequential vs Random



Challenge C → Datacenter Requirements

Datacenter isolation of I/O needs to provide:

Ease of Use

Low Priority   High Priority

Time?  Bytes?
 IOPs??

Difficult to configure

Challenge C → Datacenter Requirements

Datacenter isolation of I/O needs to provide:

Ease of Use



Time?



Bytes?



IOPs??

Difficult to configure

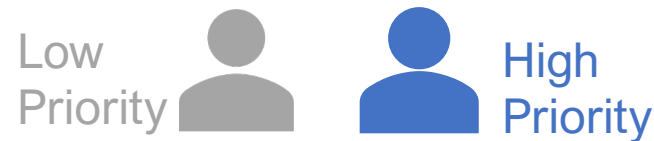
Work Conservation



Challenge C → Datacenter Requirements

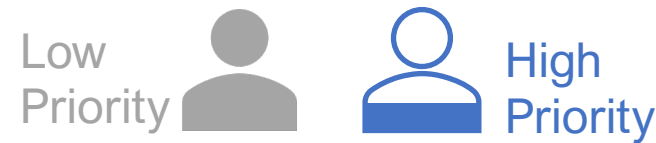
Datacenter isolation of I/O needs to provide:

Ease of Use



Difficult to configure

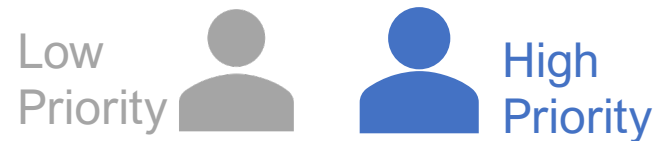
Work Conservation



Challenge C → Datacenter Requirements

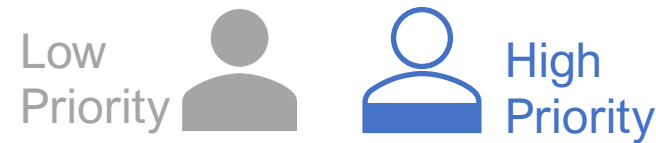
Datacenter isolation of I/O needs to provide:

Ease of Use



Difficult to configure

Work Conservation

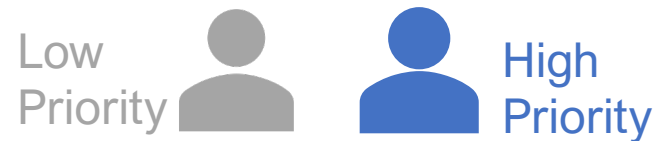


Underutilization

Challenge C → Datacenter Requirements

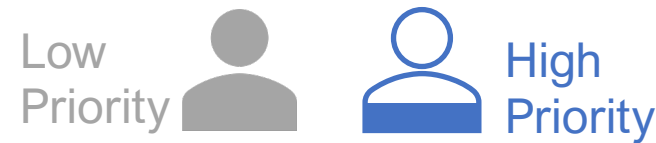
Datacenter isolation of I/O needs to provide:

Ease of Use



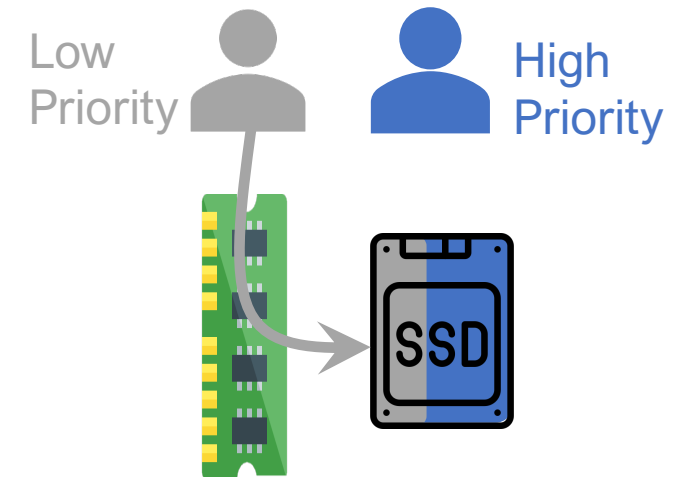
Difficult to configure

Work Conservation



Underutilization

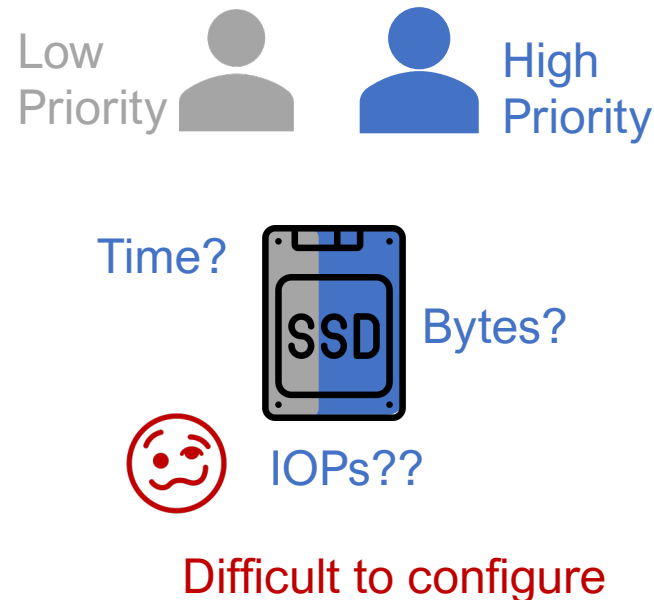
Memory-awareness



Challenge C → Datacenter Requirements

Datacenter isolation of I/O needs to provide:

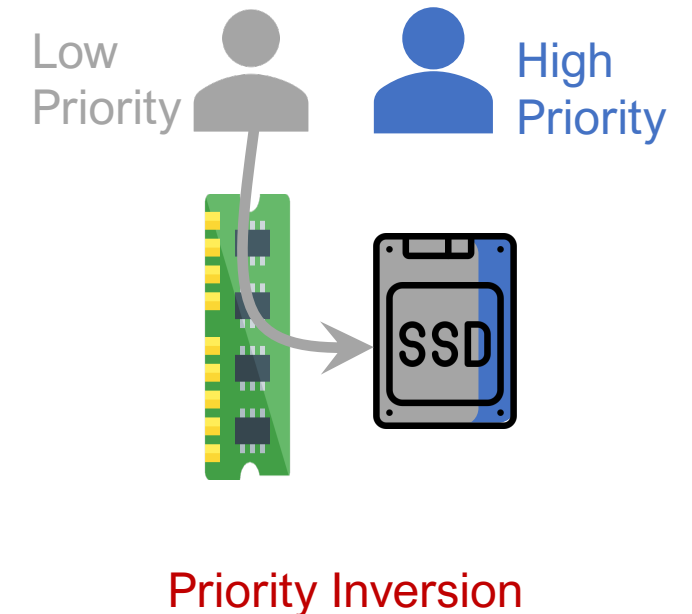
Ease of Use



Work Conservation



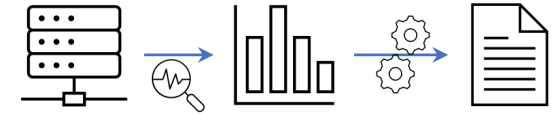
Memory-awareness



Contribution → IOCost

IO Control for Containers in Datacenters

Device occupancy with offline cost and QoS models



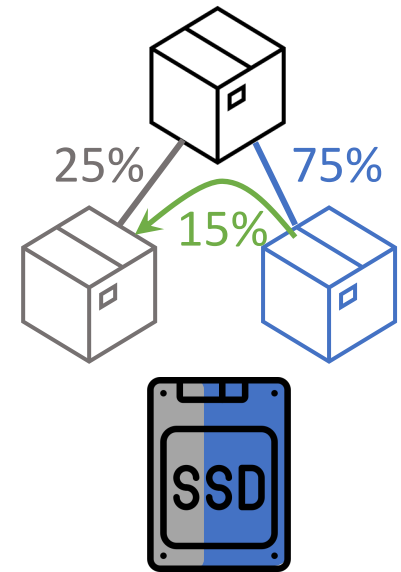
Proportional IO control through weights across containers

Lightweight work-conserving budget donation algorithm

IOCost manages IO across the entire Meta fleet

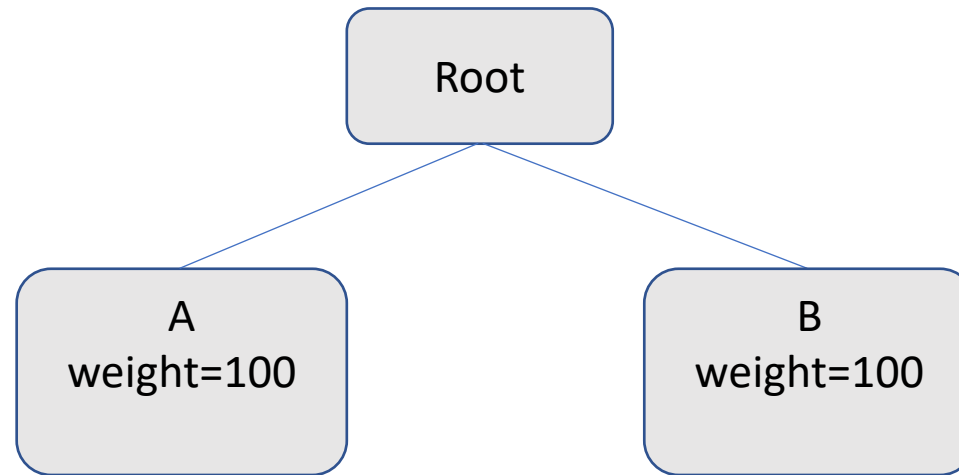
Open-source device profiling and benchmarking tools

Upstreamed in Linux

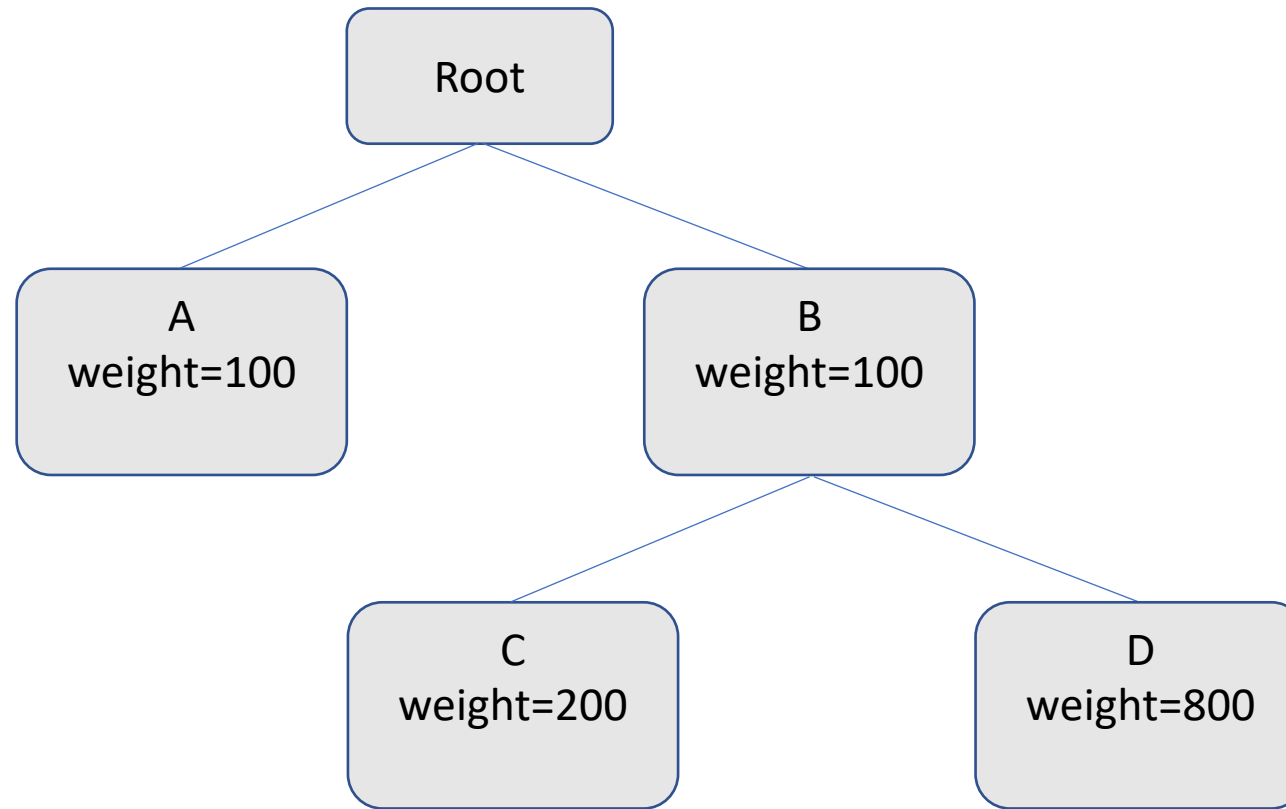


IOCost Design

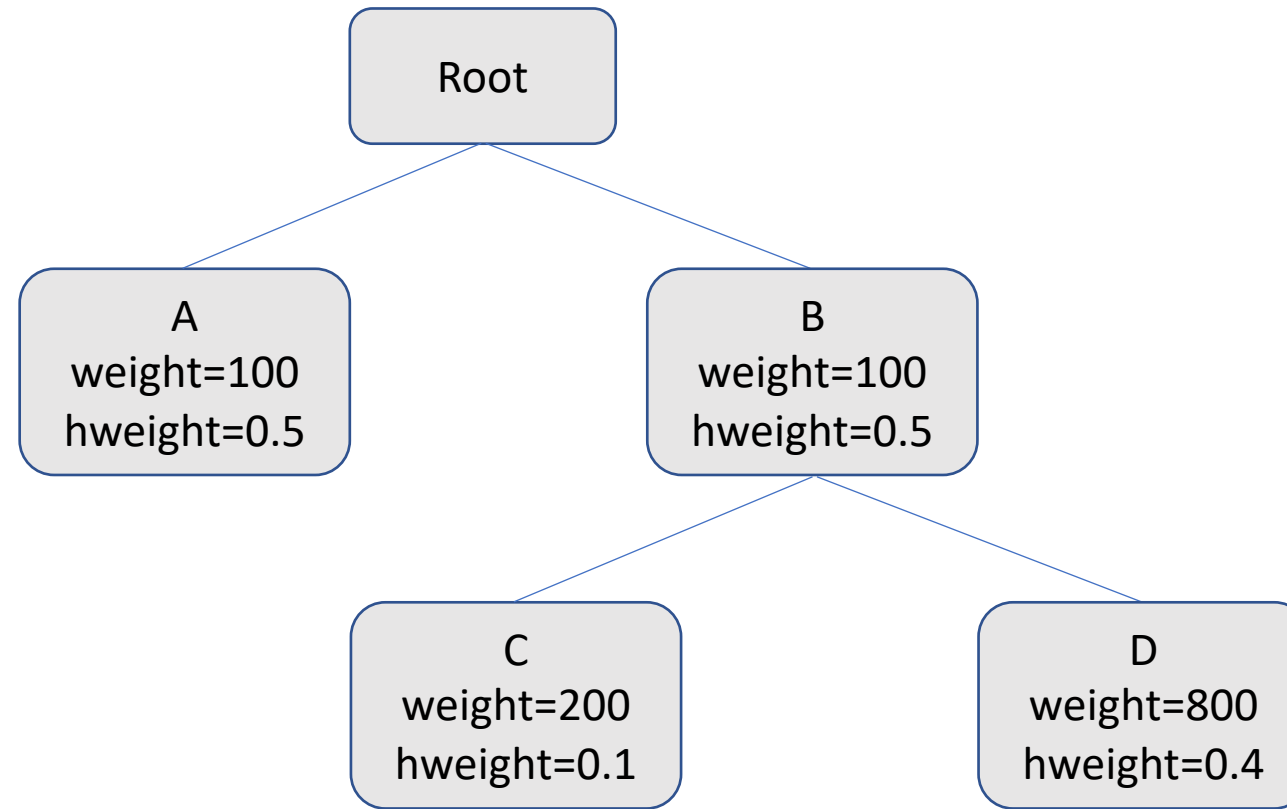
IOCost Is Configured With Weights



IOCost Is Configured With Weights



IOCost Is Configured With Weights

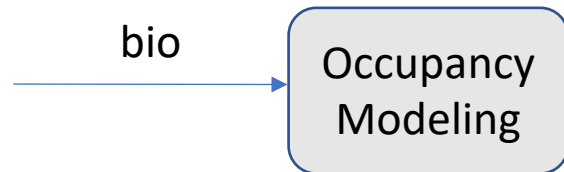


IOCost Issue Path Estimates Occupancy

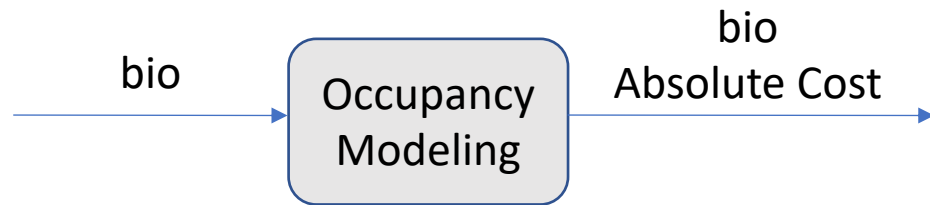
bio



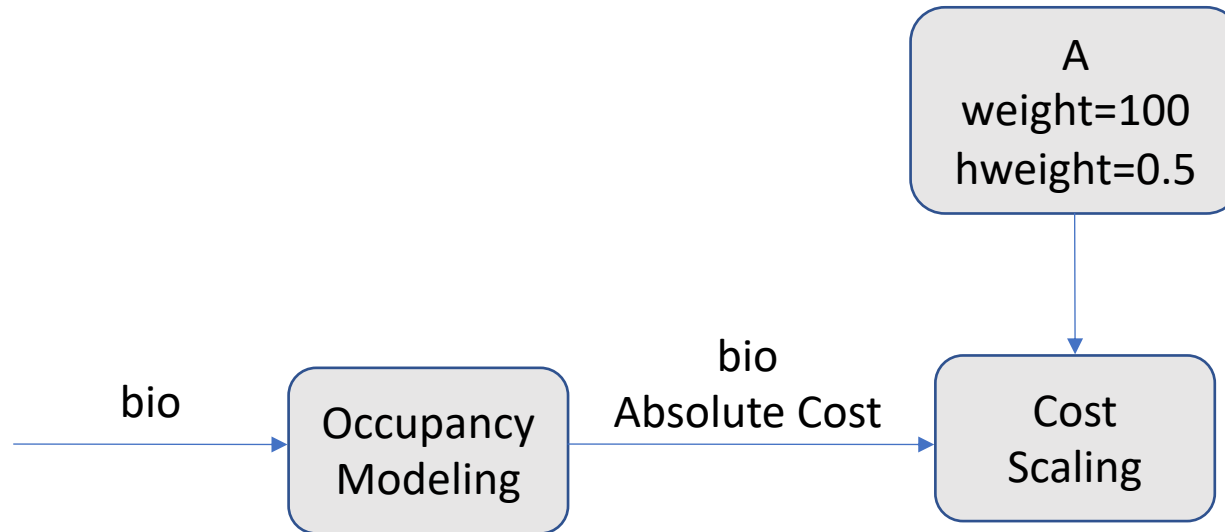
IOCost Issue Path Estimates Occupancy



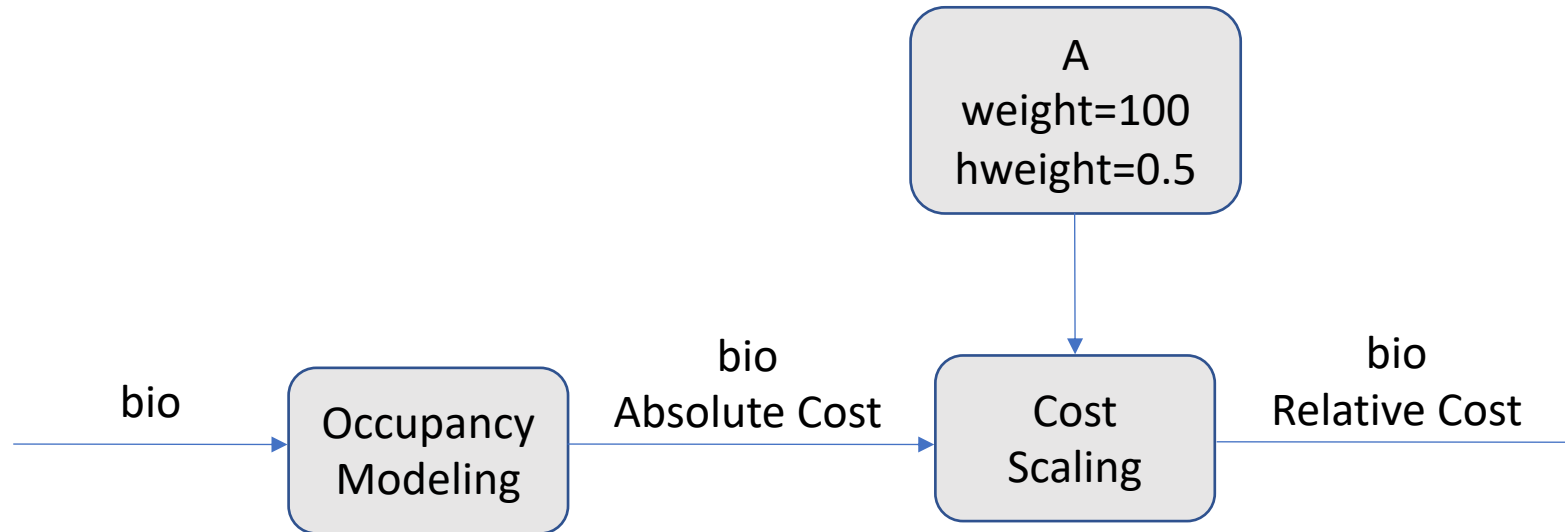
IOCost Issue Path Estimates Occupancy



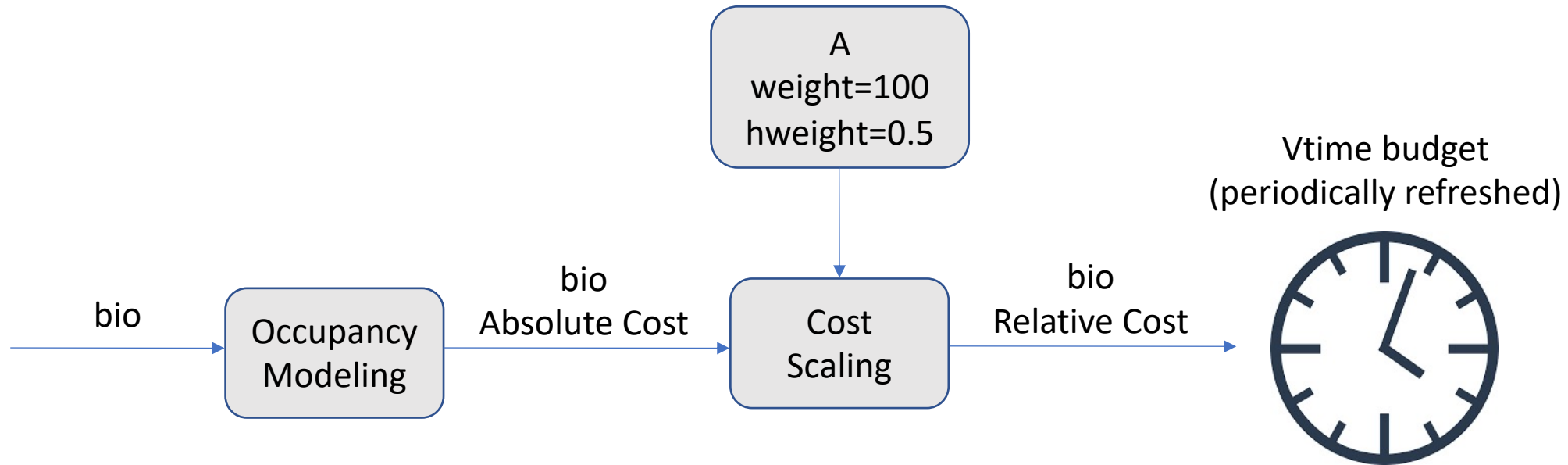
IOCost Issue Path Estimates Occupancy



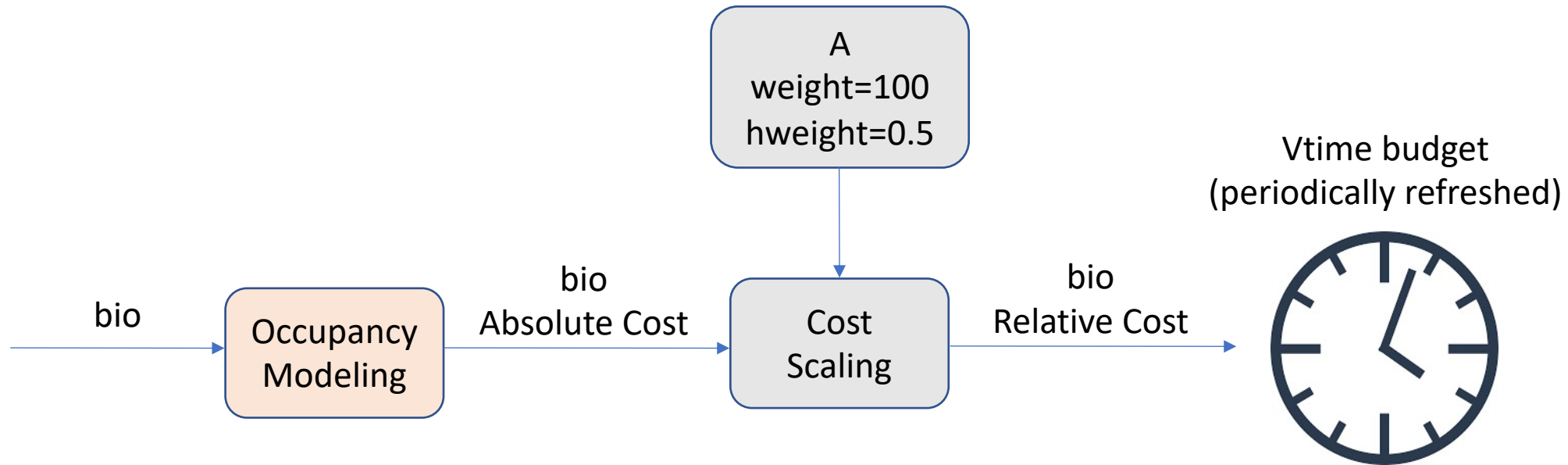
IOCost Issue Path Estimates Occupancy



IOCost Issue Path Estimates Occupancy



IOCost Issue Path Estimates Occupancy



Occupancy Modeling Captures Device Differences

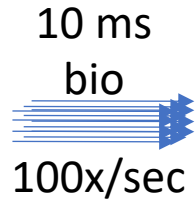
Occupancy Modeling Captures Device Differences

IO Cost is in units of time but not a measure of latency

10 ms
bio
→

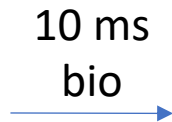
Occupancy Modeling Captures Device Differences

IO Cost is in units of time but not a measure of latency



Occupancy Modeling Captures Device Differences

IO Cost is in units of time but not a measure of latency

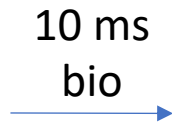


Linear Model estimated offline via synthetic saturating workloads

$$\text{io_cost} = \text{base_cost} + \text{size_cost_rate} * \text{bio_size}$$

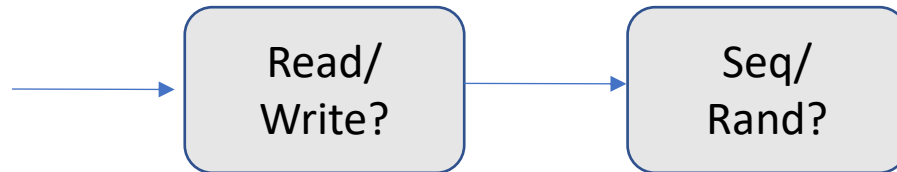
Occupancy Modeling Captures Device Differences

IO Cost is in units of time but not a measure of latency



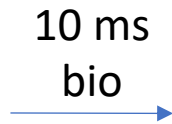
Linear Model estimated offline via synthetic saturating workloads

$$\text{io_cost} = \text{base_cost} + \text{size_cost_rate} * \text{bio_size}$$



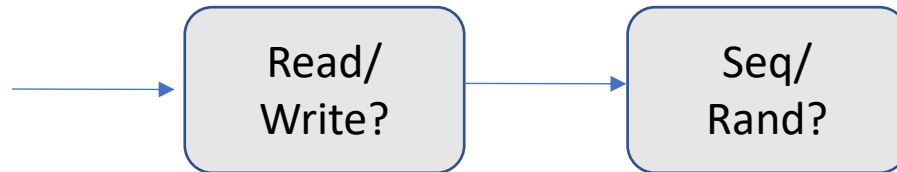
Occupancy Modeling Captures Device Differences

IO Cost is in units of time but not a measure of latency



Linear Model estimated offline via synthetic saturating workloads

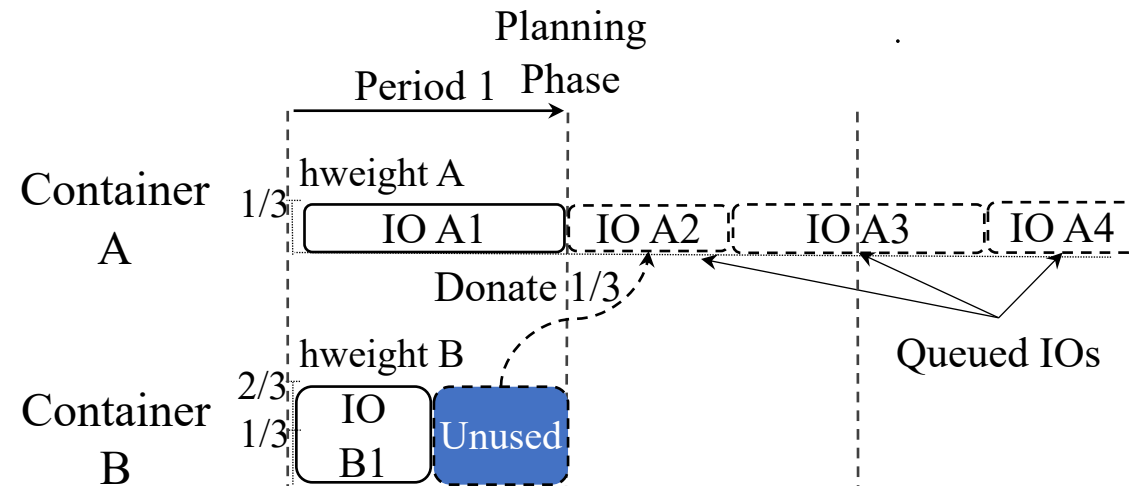
$$\text{io_cost} = \text{base_cost} + \text{size_cost_rate} * \text{bio_size}$$



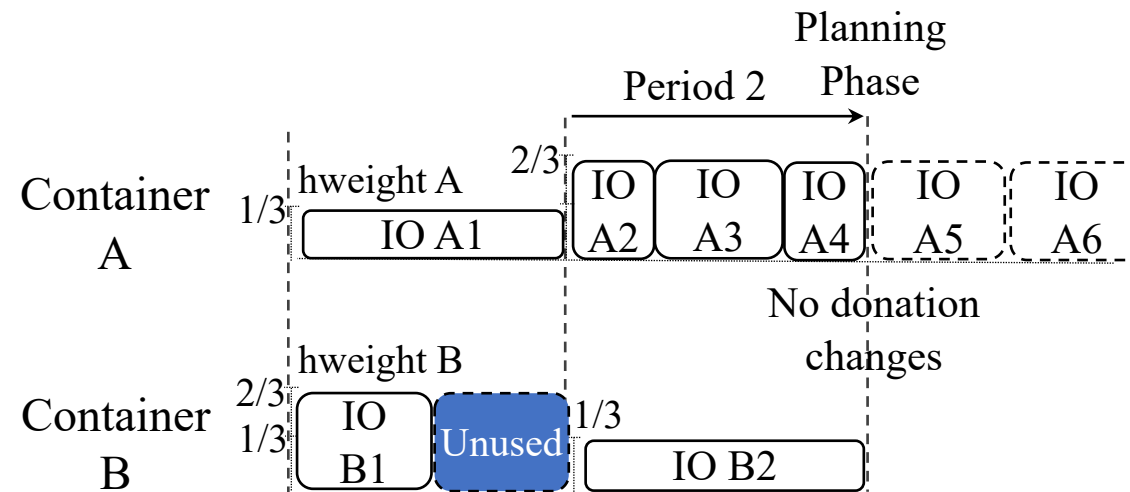
How to adjust for inaccurate models? VRate adjustment (see paper for details)

Budget Donation Ensures Work Conservation

Budget Donation Ensures Work Conservation

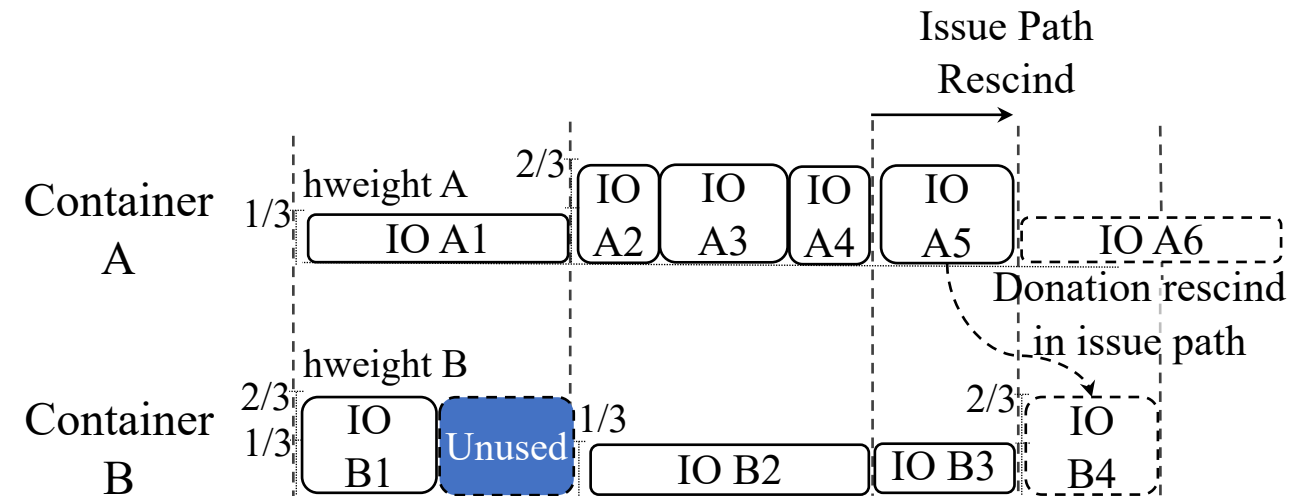


Budget Donation Ensures Work Conservation



(b)

Budget Donation Ensures Work Conservation

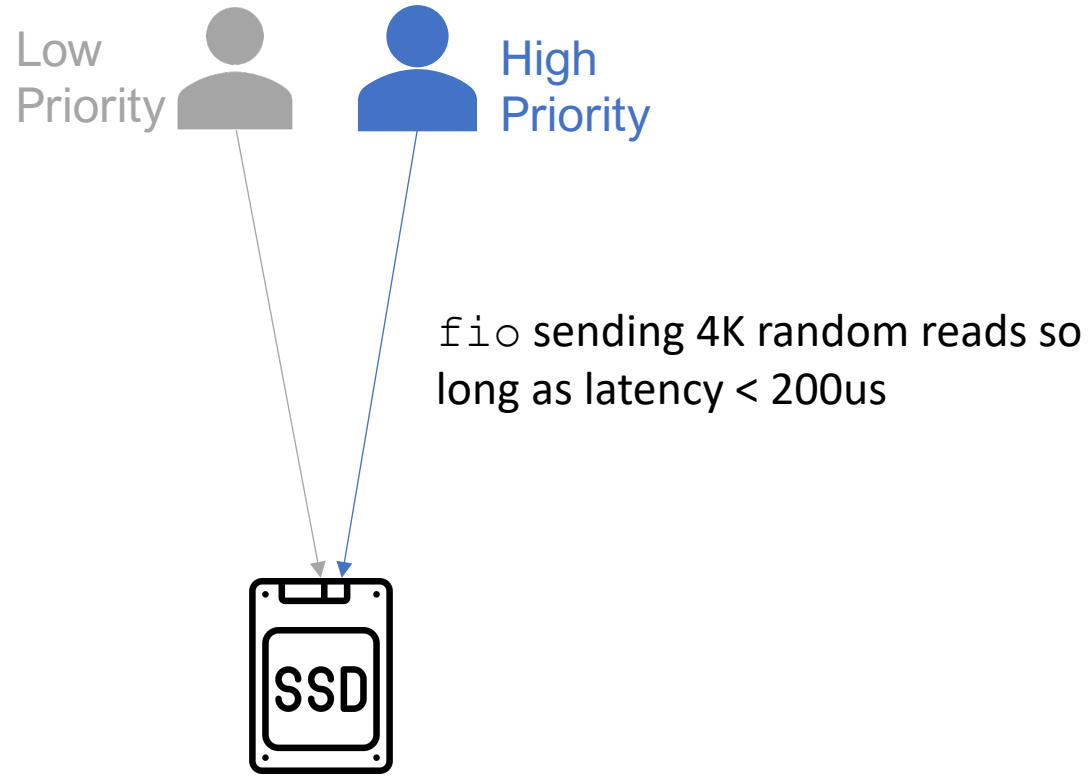


(c)

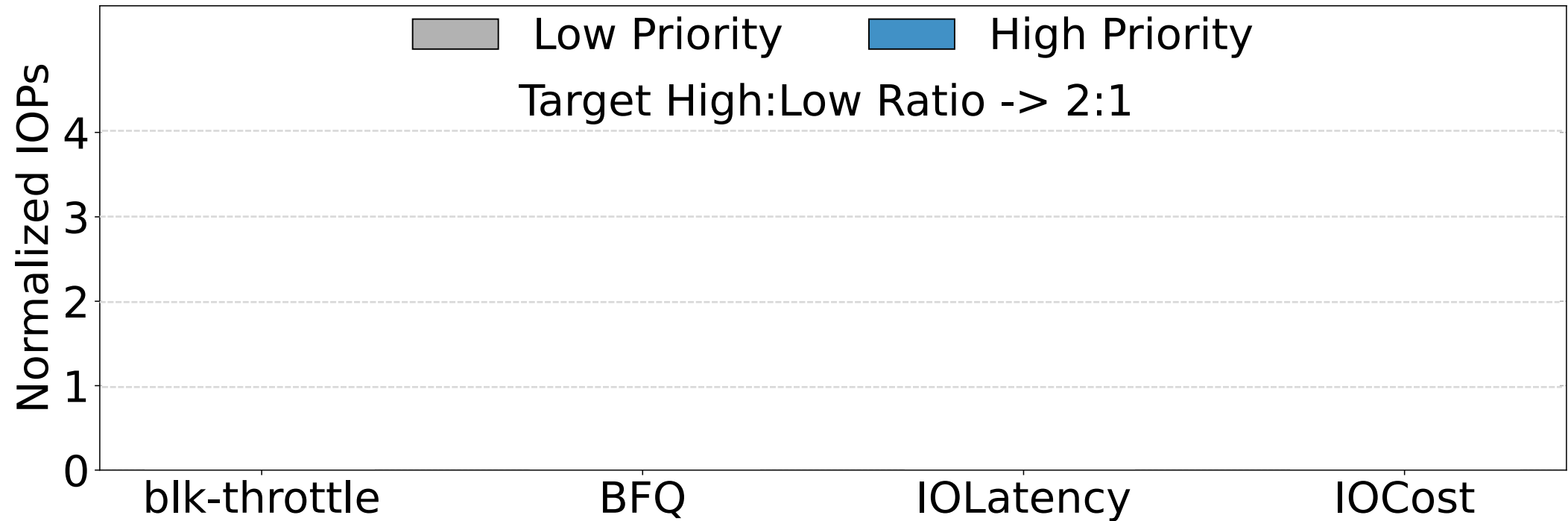
Evaluation

Evaluation – Proportional Control

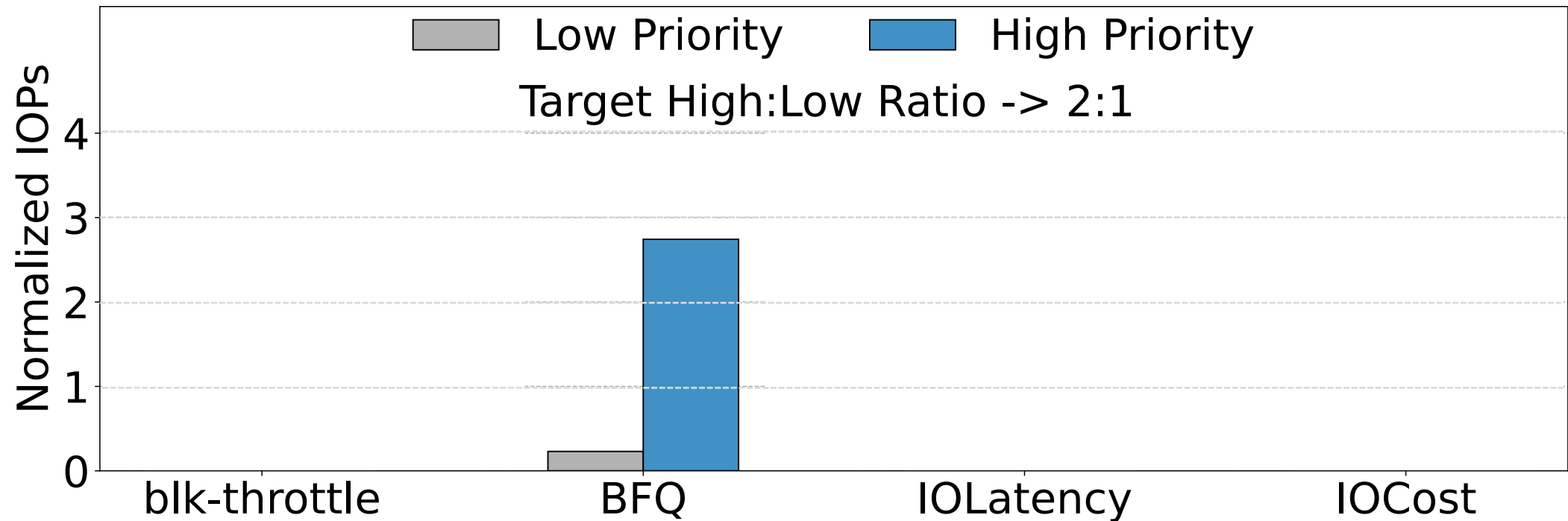
Evaluation – Proportional Control



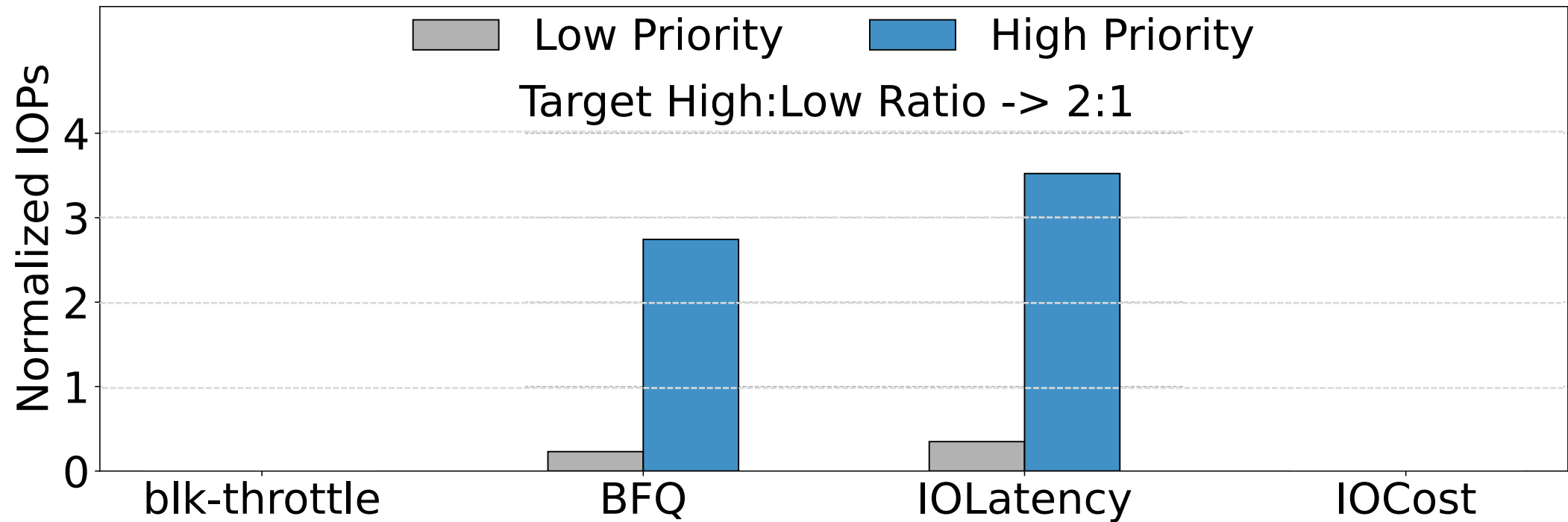
Evaluation – Proportional Control



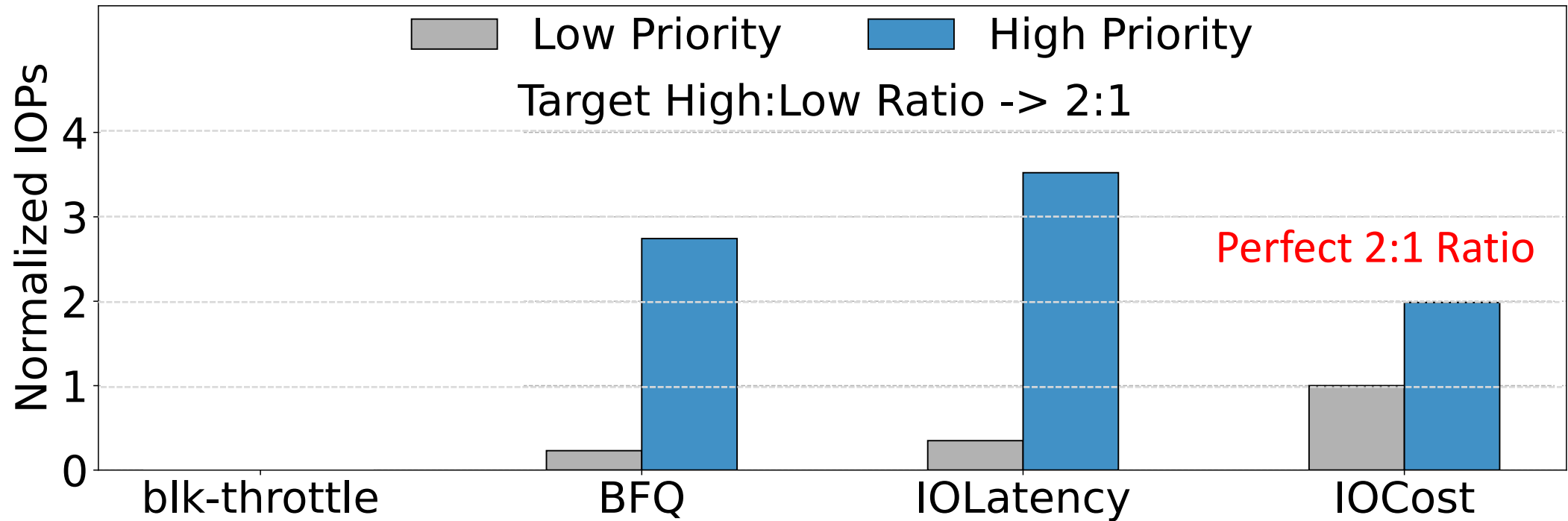
Evaluation – Proportional Control



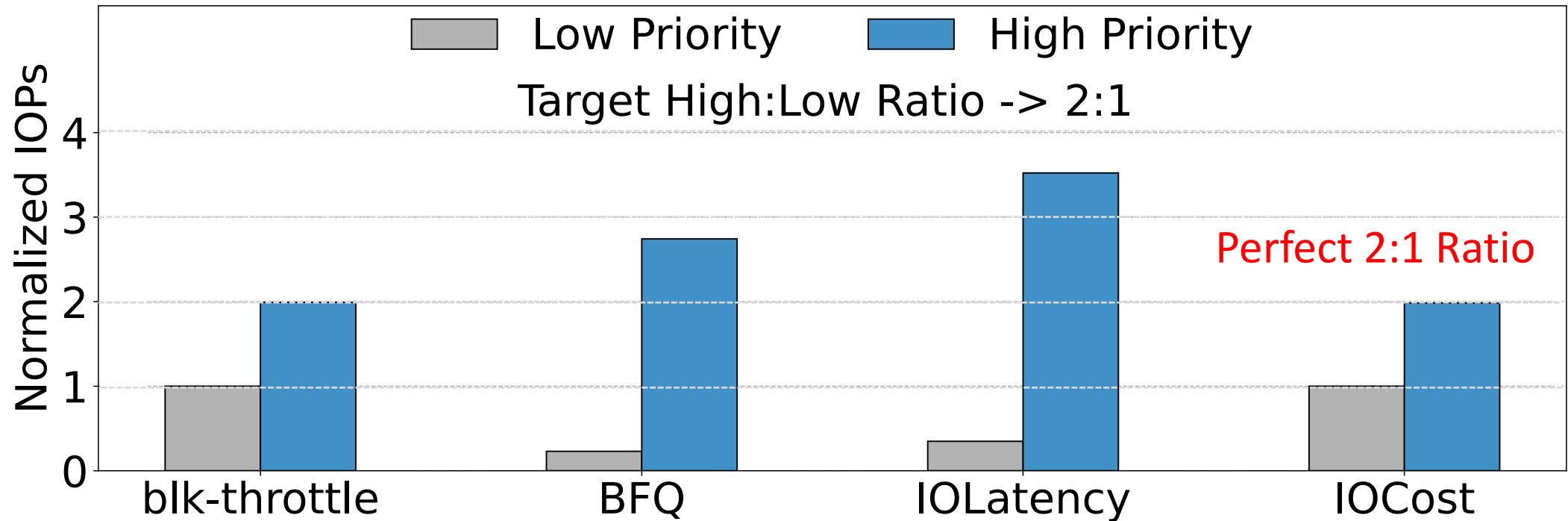
Evaluation – Proportional Control



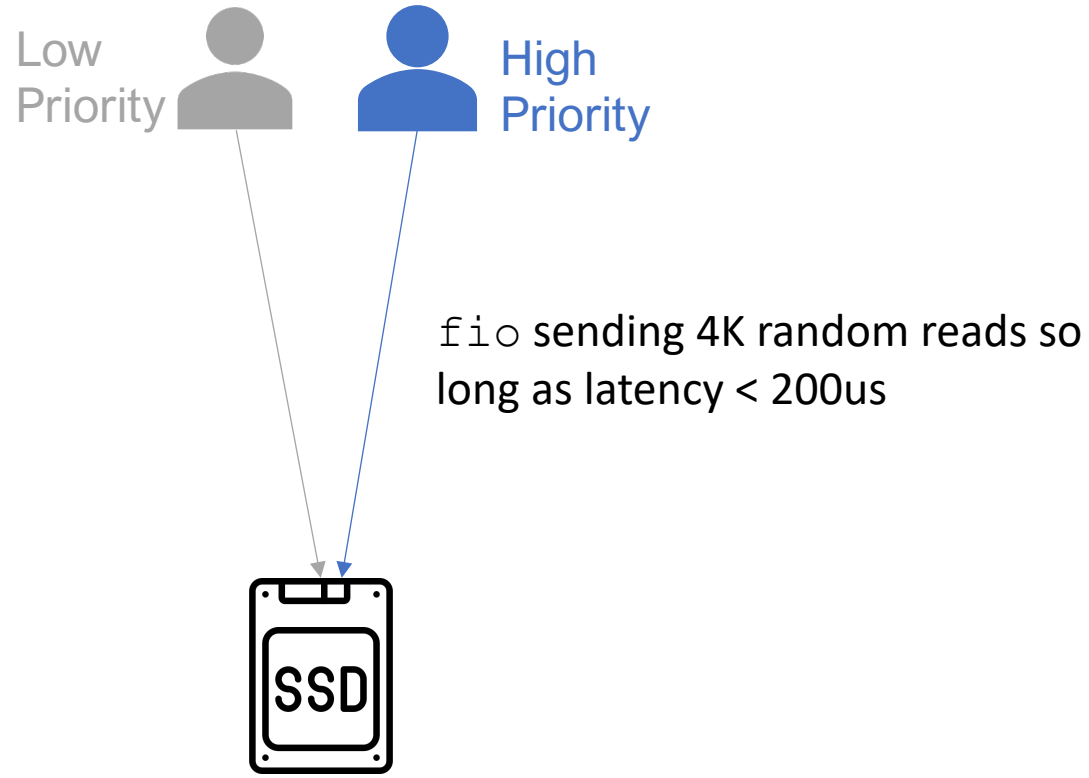
Evaluation – Proportional Control



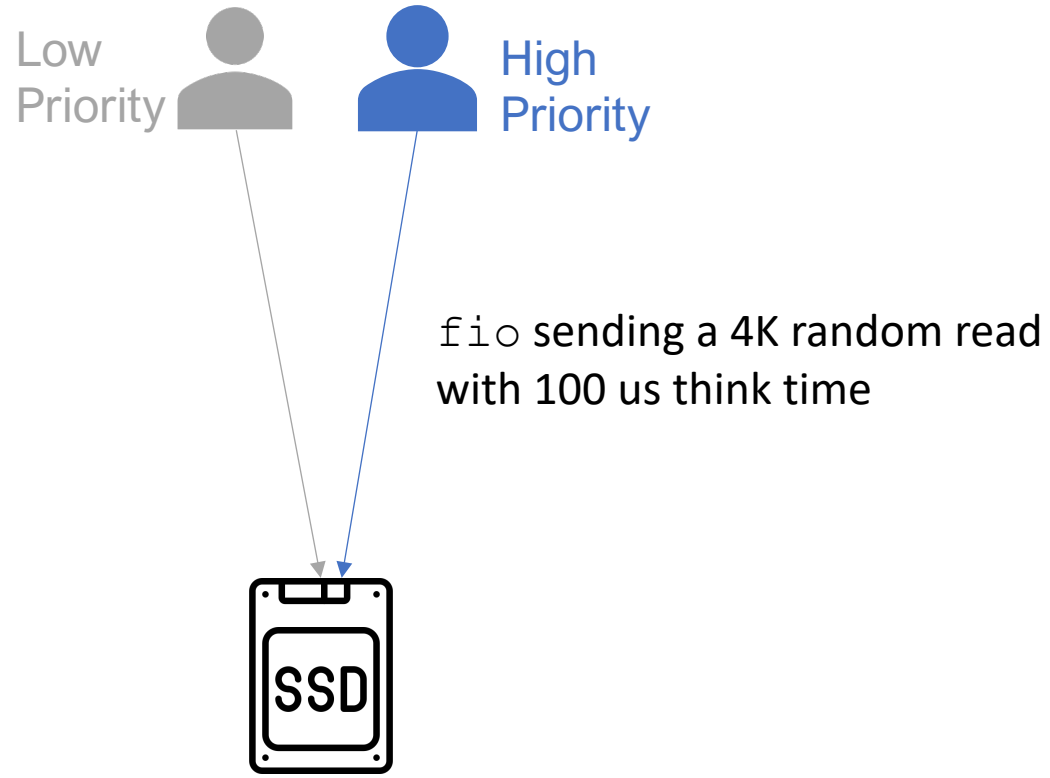
Evaluation – Proportional Control



Evaluation – Work Conservation

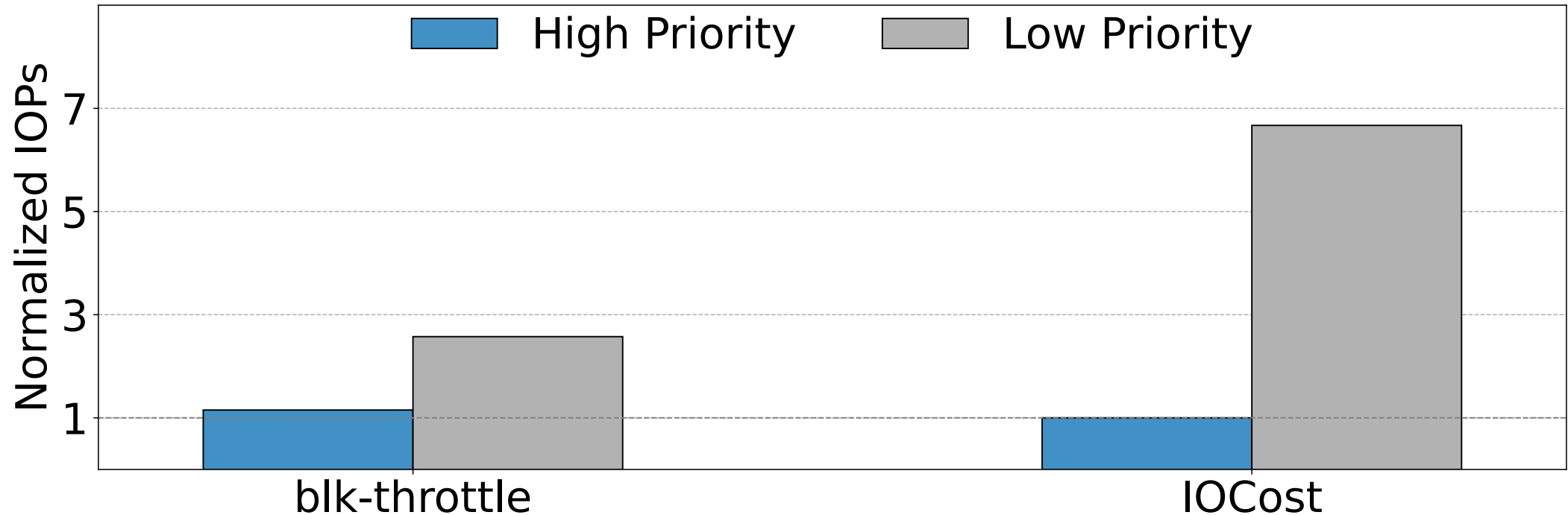


Evaluation – Work Conservation



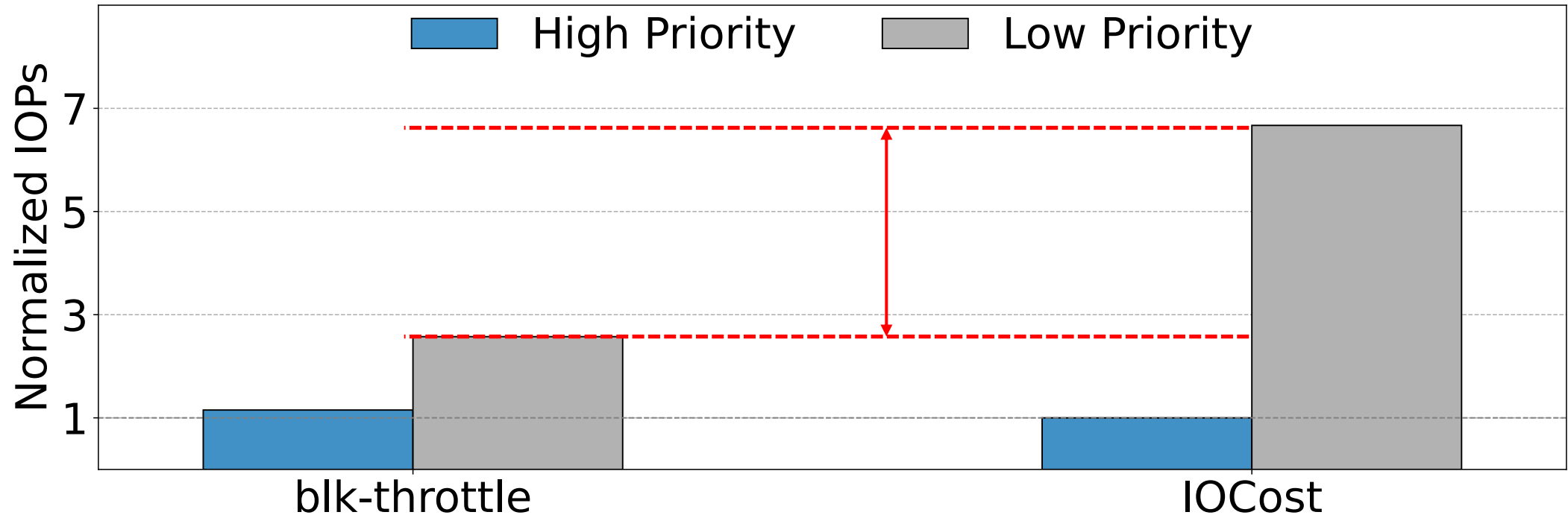
Evaluation – Work Conservation

The low priority workload should use up all available capacity



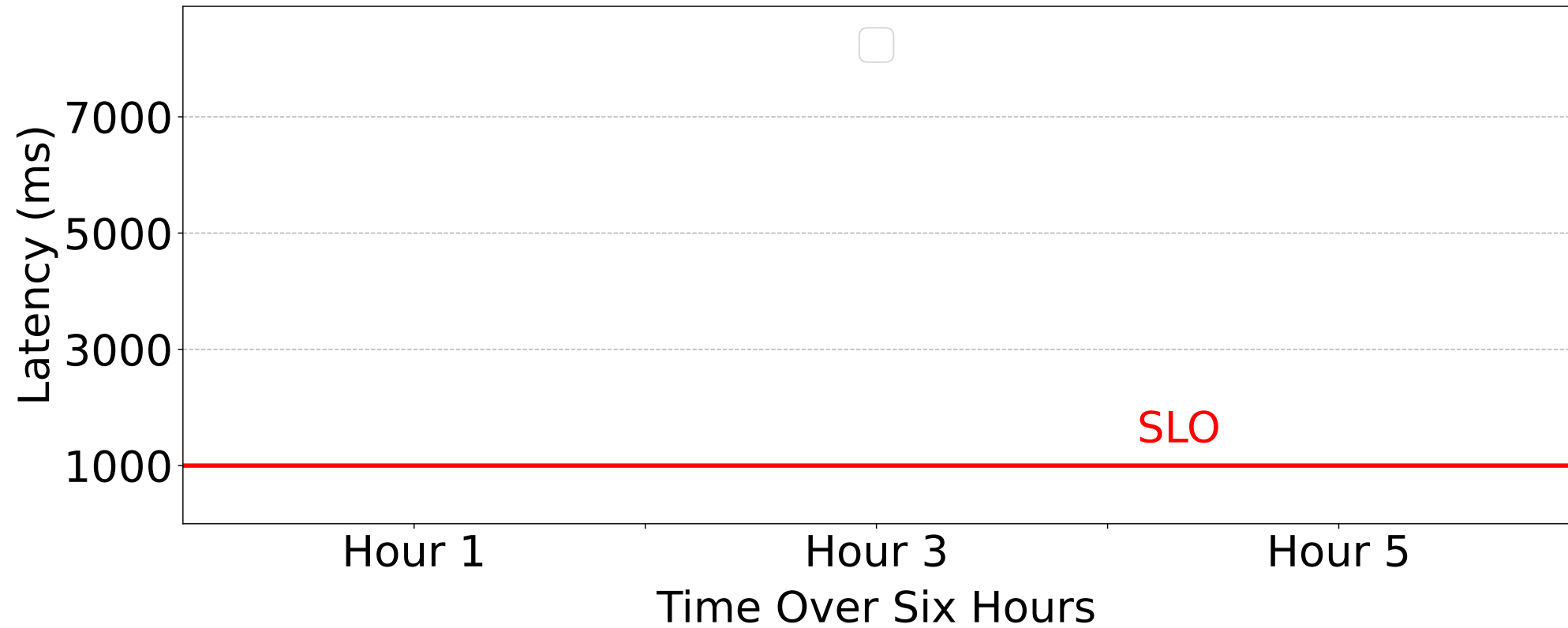
Evaluation – Work Conservation

The low priority workload should use up all available capacity

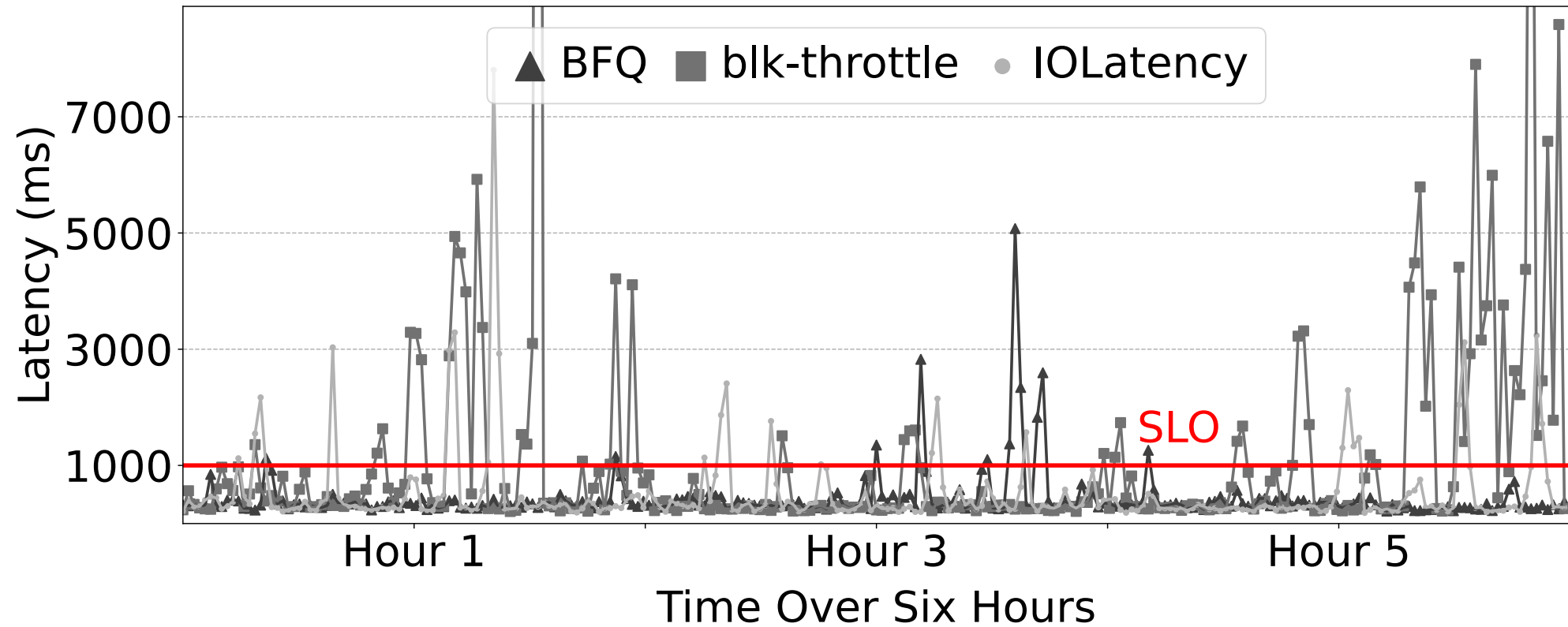


IOCost allows much more aggregate IO

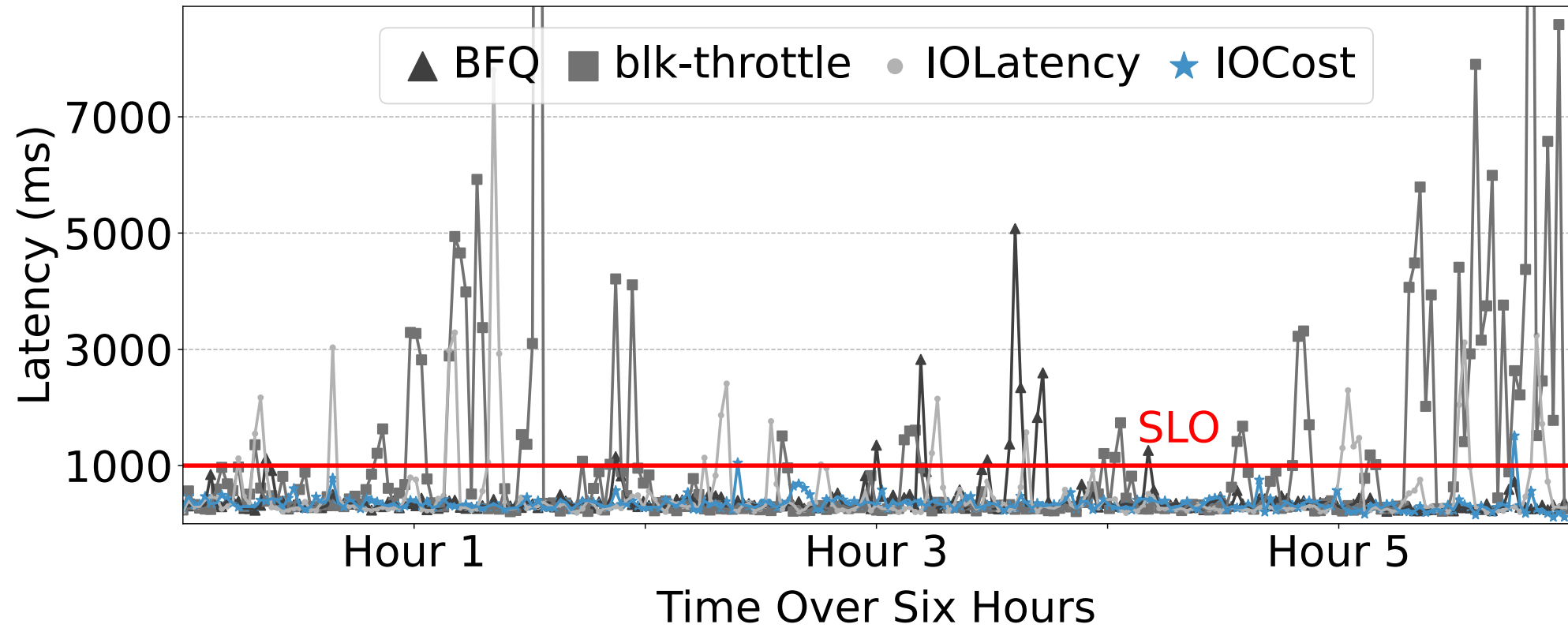
Evaluation – Zookeeper



Evaluation – Zookeeper



Evaluation – Zookeeper



More in the Paper

Memory awareness

Budget donation algorithm

SSD and workload heterogeneity across the Fleet

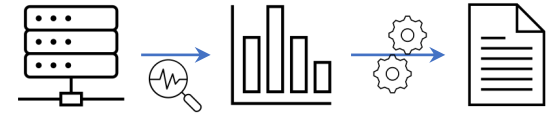
Additional evaluation:

- IO control overhead
- Disk modeling
- QoS and Vrate adjustment
- Overcommitted environments
- Remote storage and VMs
- Package fetching and Container Cleanups

Takeaway→ IOCost

IO Control for Containers in Datacenters

Device occupancy with offline cost and QoS models



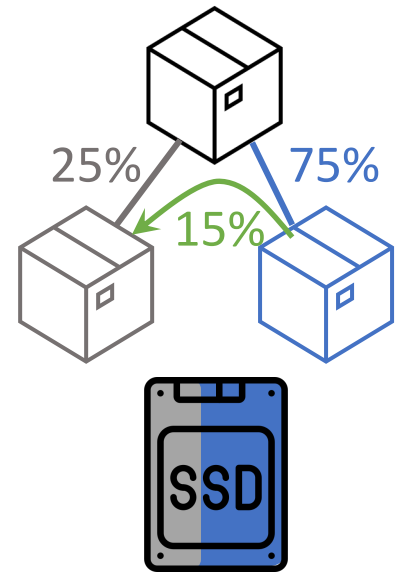
Proportional IO control through weights across containers

Lightweight work-conserving budget donation algorithm

IOCost manages IO across the entire Meta fleet

Open-source device profiling and benchmarking tools

Upstreamed in Linux



IOCost: Block IO Control for Containers in Datacenters

Tejun Heo, **Dan Schatzberg**, Andrew Newell, Song Liu, Saravanan Dhakshinamurthy, Iyswarya Narayanan, Josef Bacik, Chris Mason, Chunqiang Tang, ‡Dimitrios Skarlatos

Session 5B: Data Center and Cloud Services
Thursday, March 3 @ 2pm

ASPLOS 2022