

# Flash Storage Disaggregation

Ana Klimovic<sup>1</sup>, Christos Kozyrakis<sup>1,4</sup>,  
Eno Thereska<sup>3,5</sup>, Binu John<sup>2</sup> and Sanjeev Kumar<sup>2</sup>

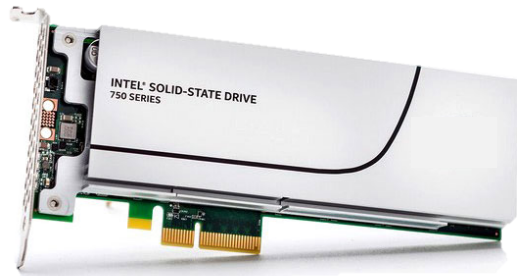


Real-time streams powered by Apache Kafka.



# Flash is underutilized

- Flash provides higher throughput and lower latency than disk



PCIe Flash:

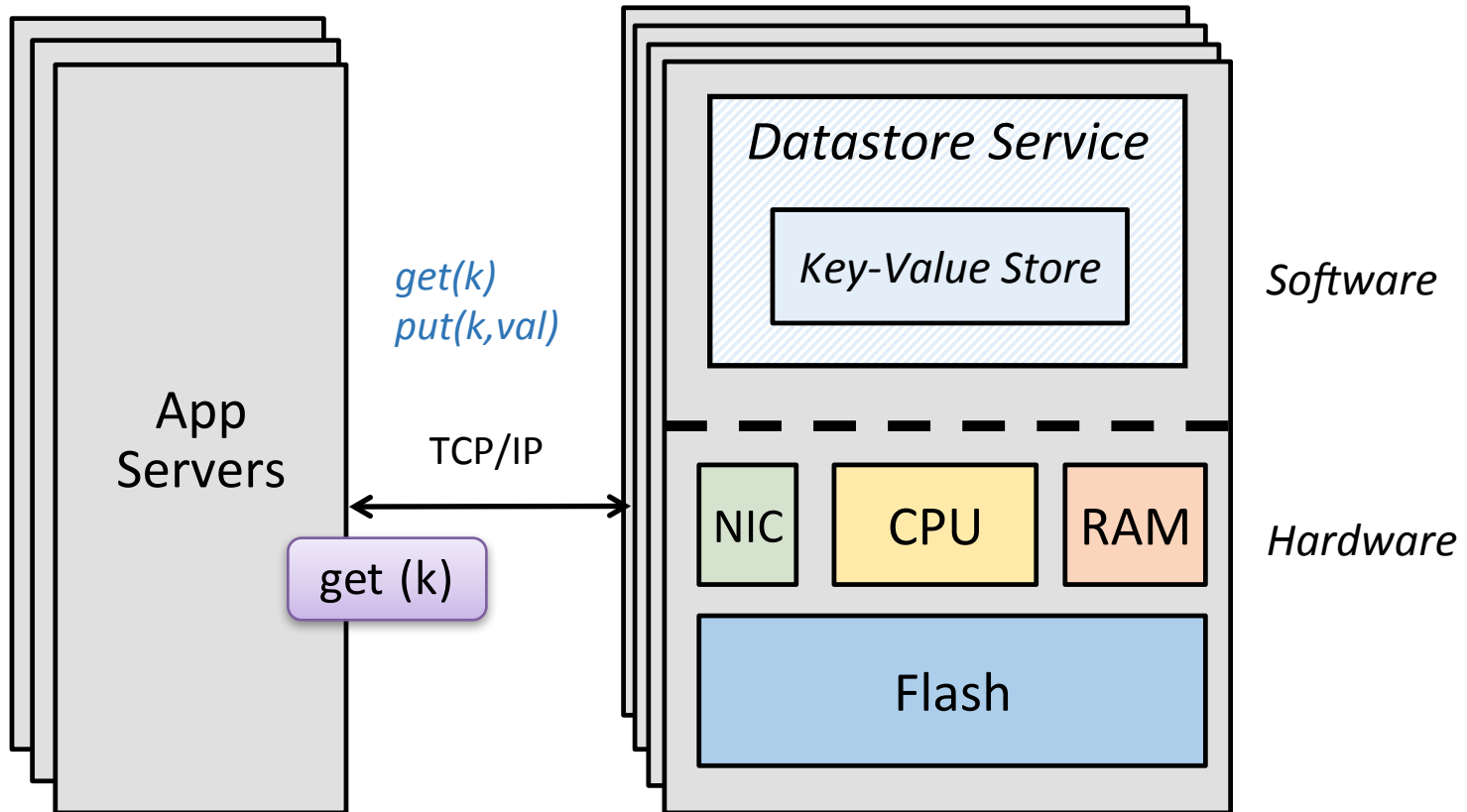
- 100,000s of IOPS
- 10s of  $\mu$ s latency

- Flash is underutilized in datacenters due to imbalanced resource requirements

# Datacenter Flash Use-Case

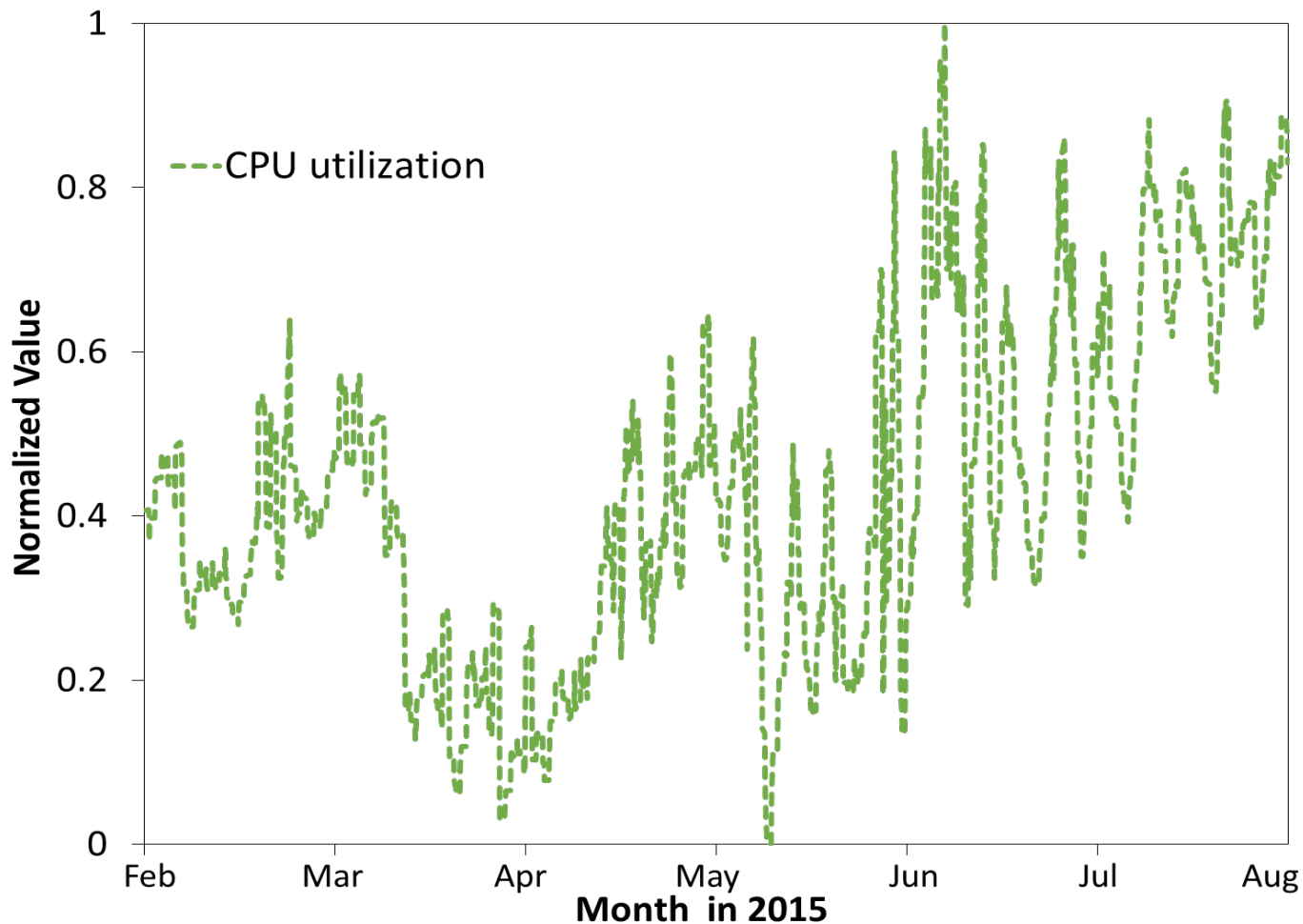
## Application Tier

## Datastore Tier



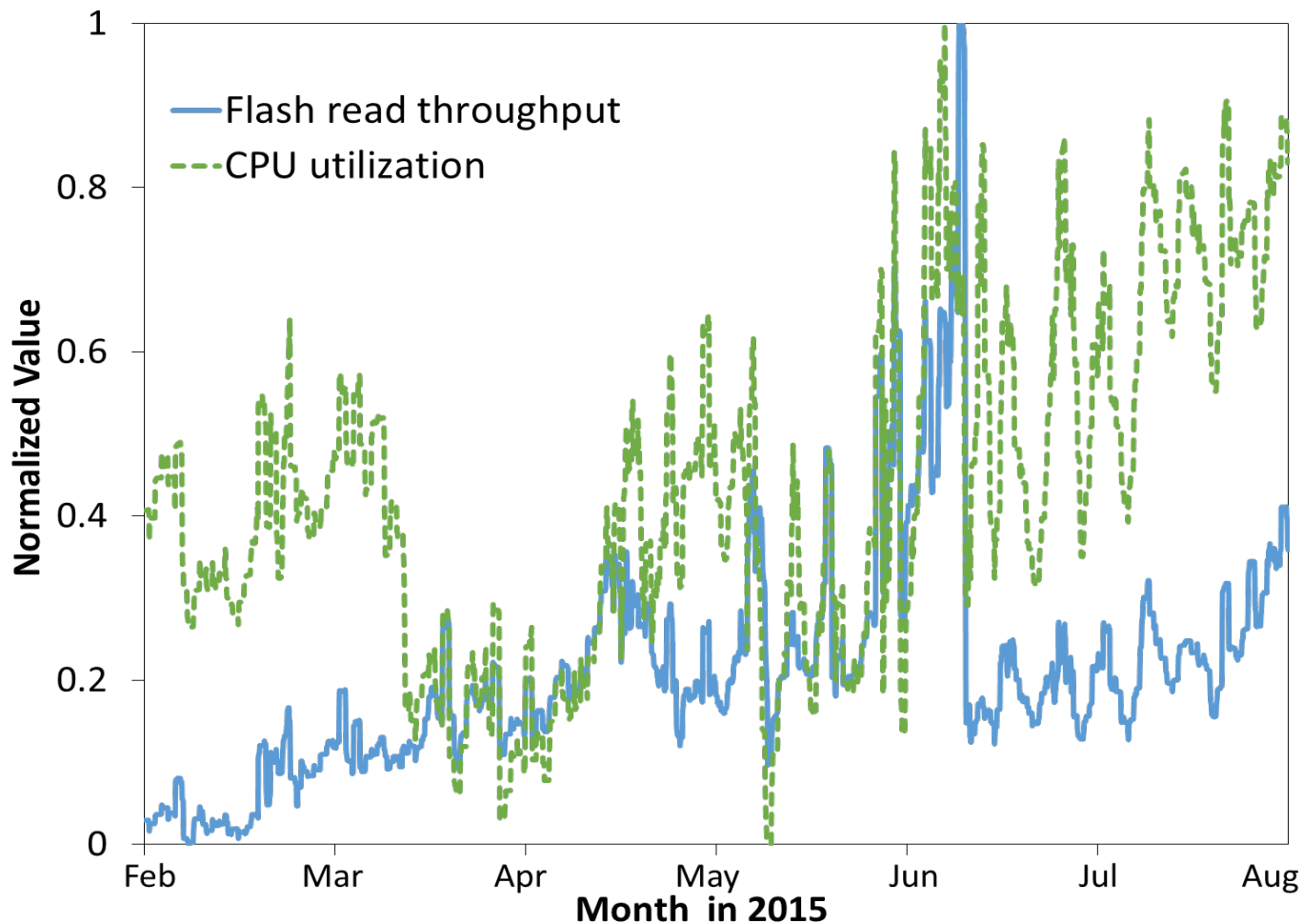
# Imbalanced Resource Utilization

- Sample utilization of Facebook servers hosting a Flash-based key-value store over 6 months



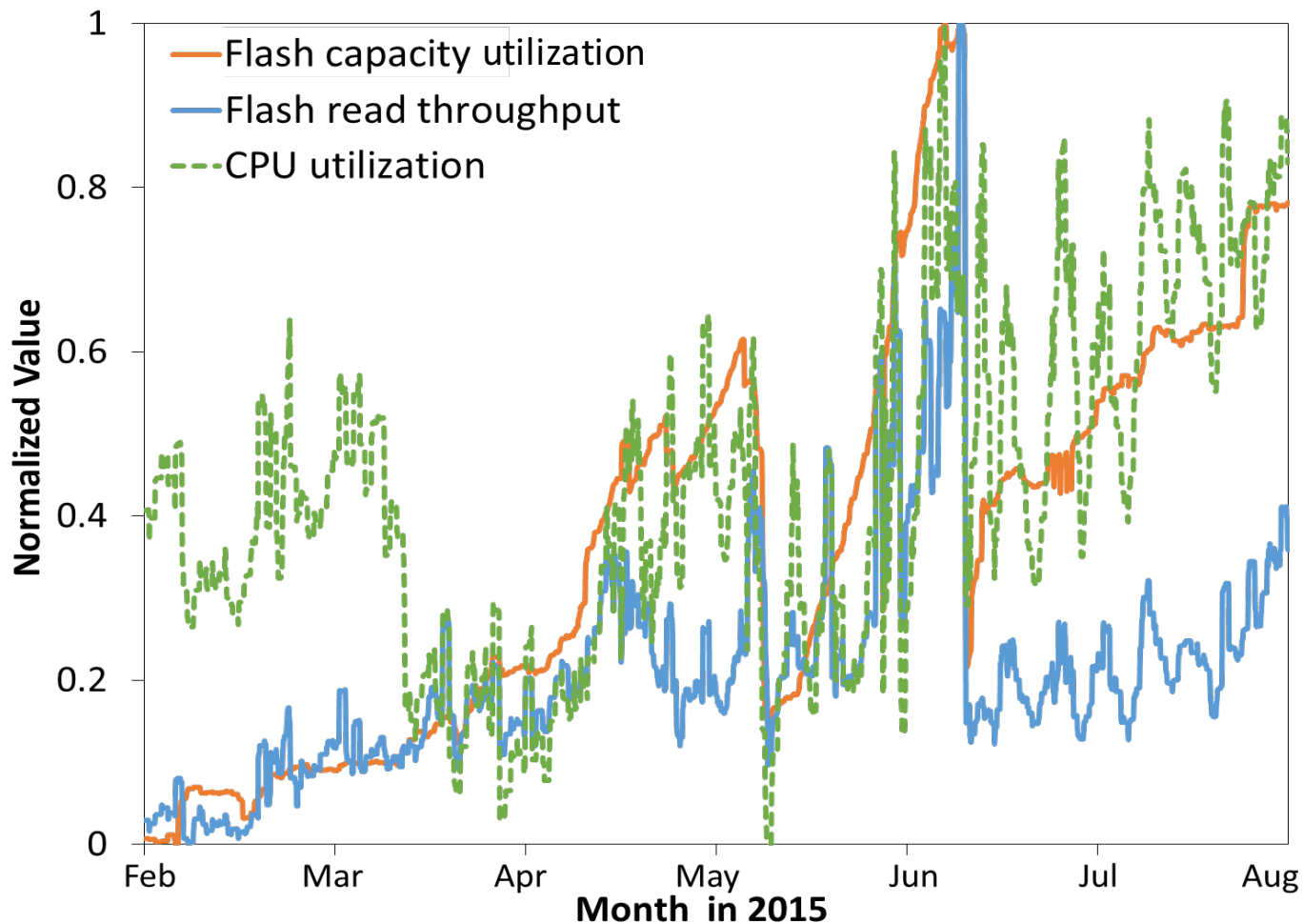
# Imbalanced Resource Utilization

- Sample utilization of Facebook servers hosting a Flash-based key-value store over 6 months



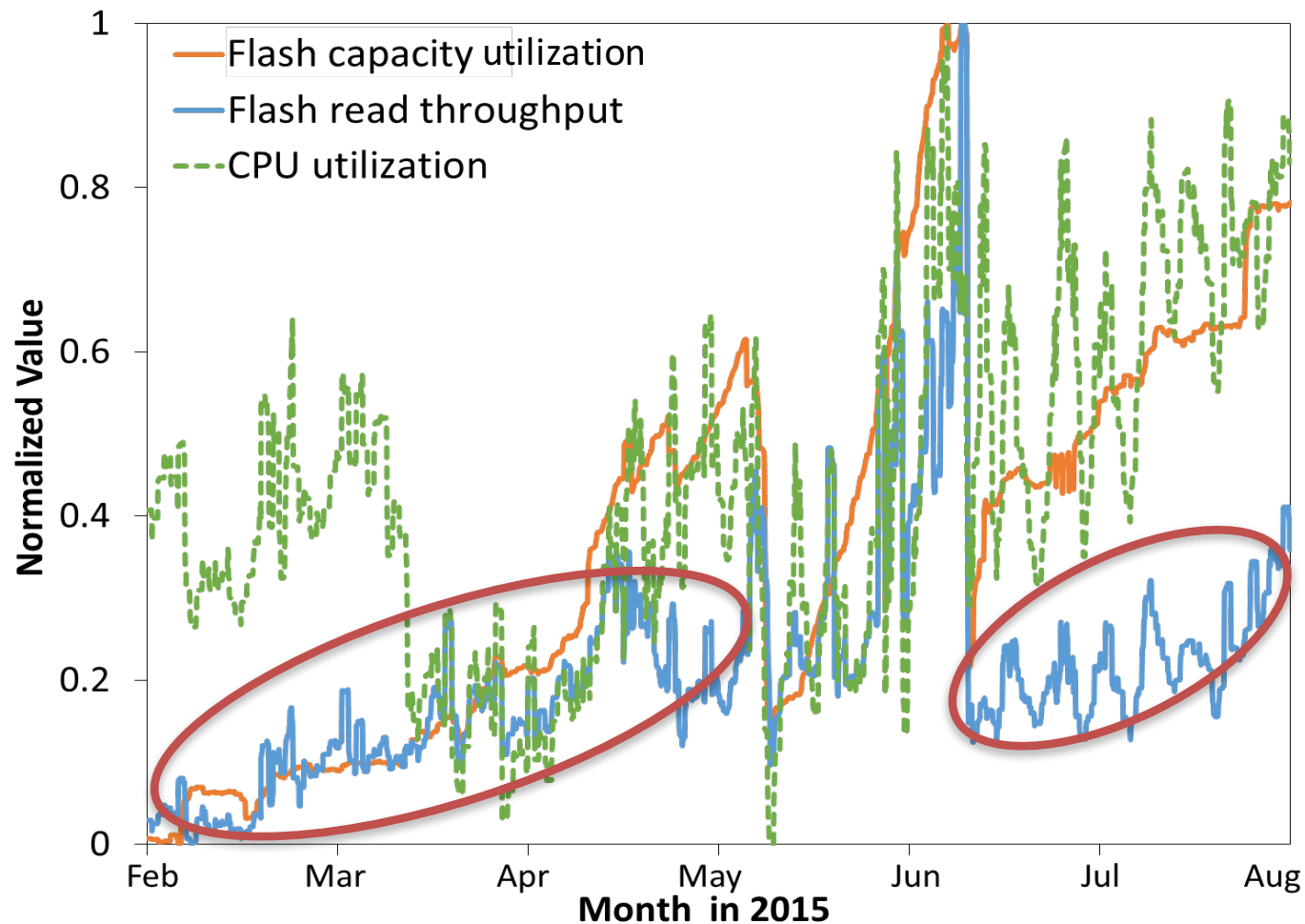
# Imbalanced Resource Utilization

- Sample utilization of Facebook servers hosting a Flash-based key-value store over 6 months



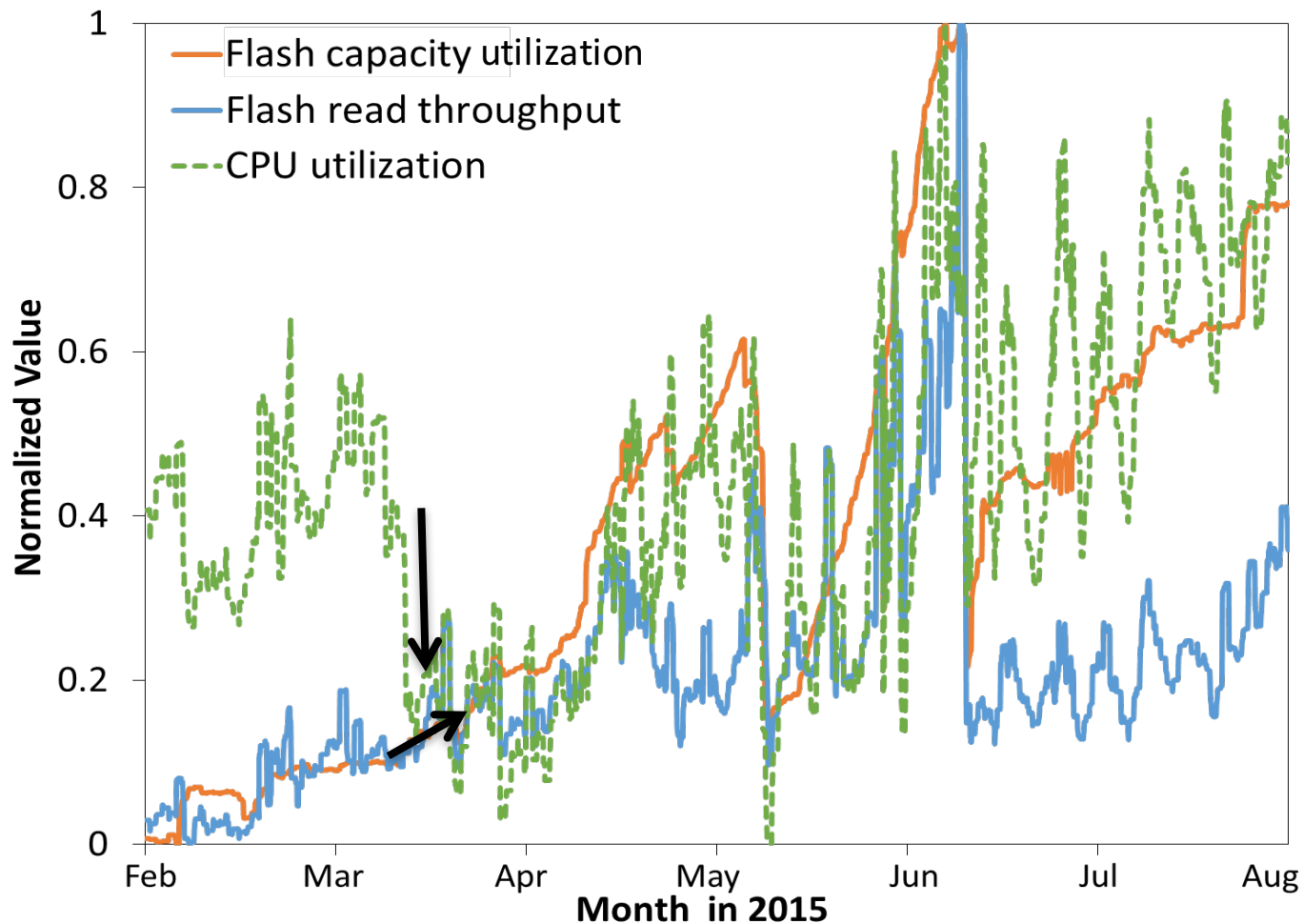
# Imbalanced Resource Utilization

- Flash capacity and IOPS are underutilized for long periods of time



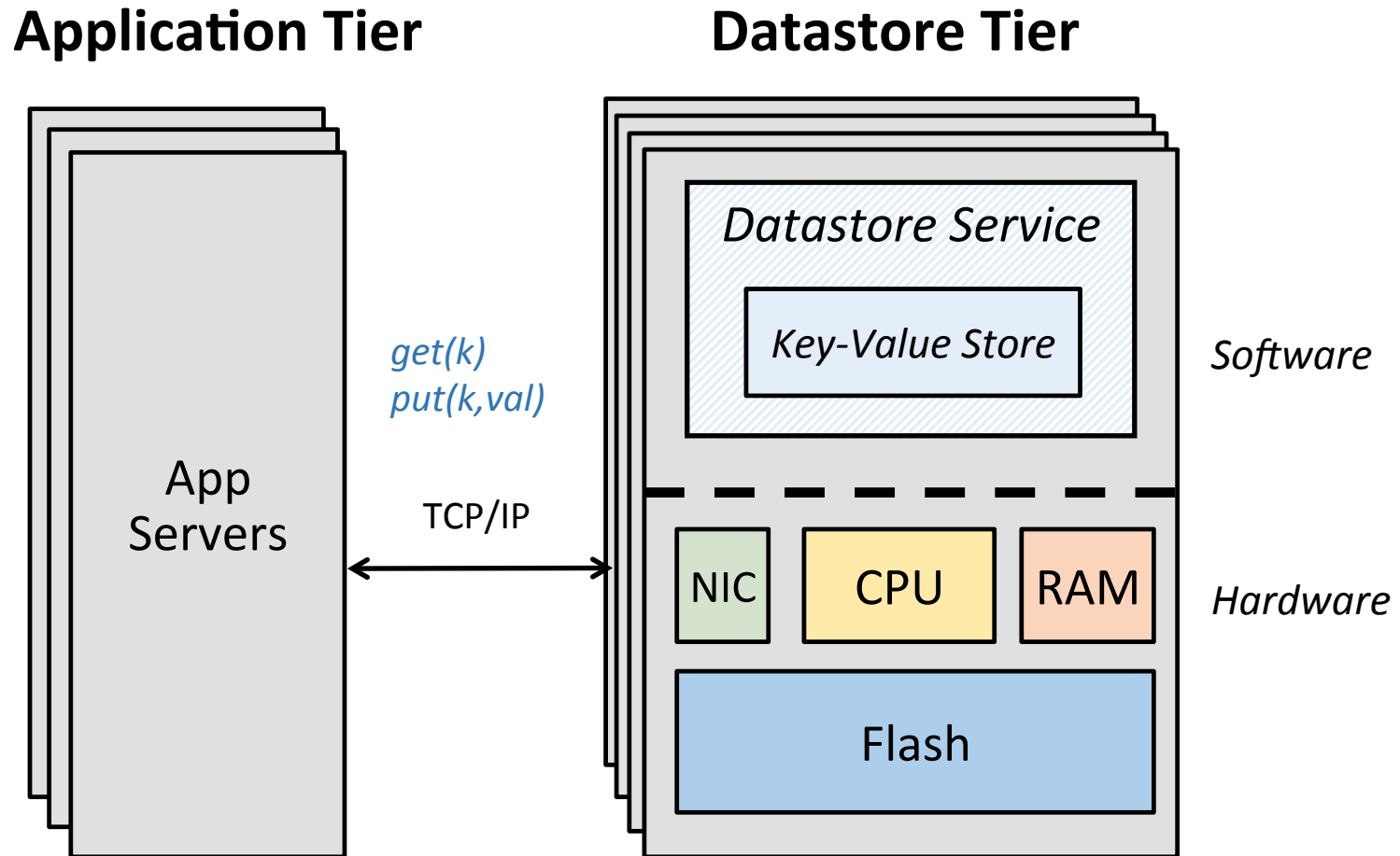
# Imbalanced Resource Utilization

- CPU and Flash utilization vary with separate trends



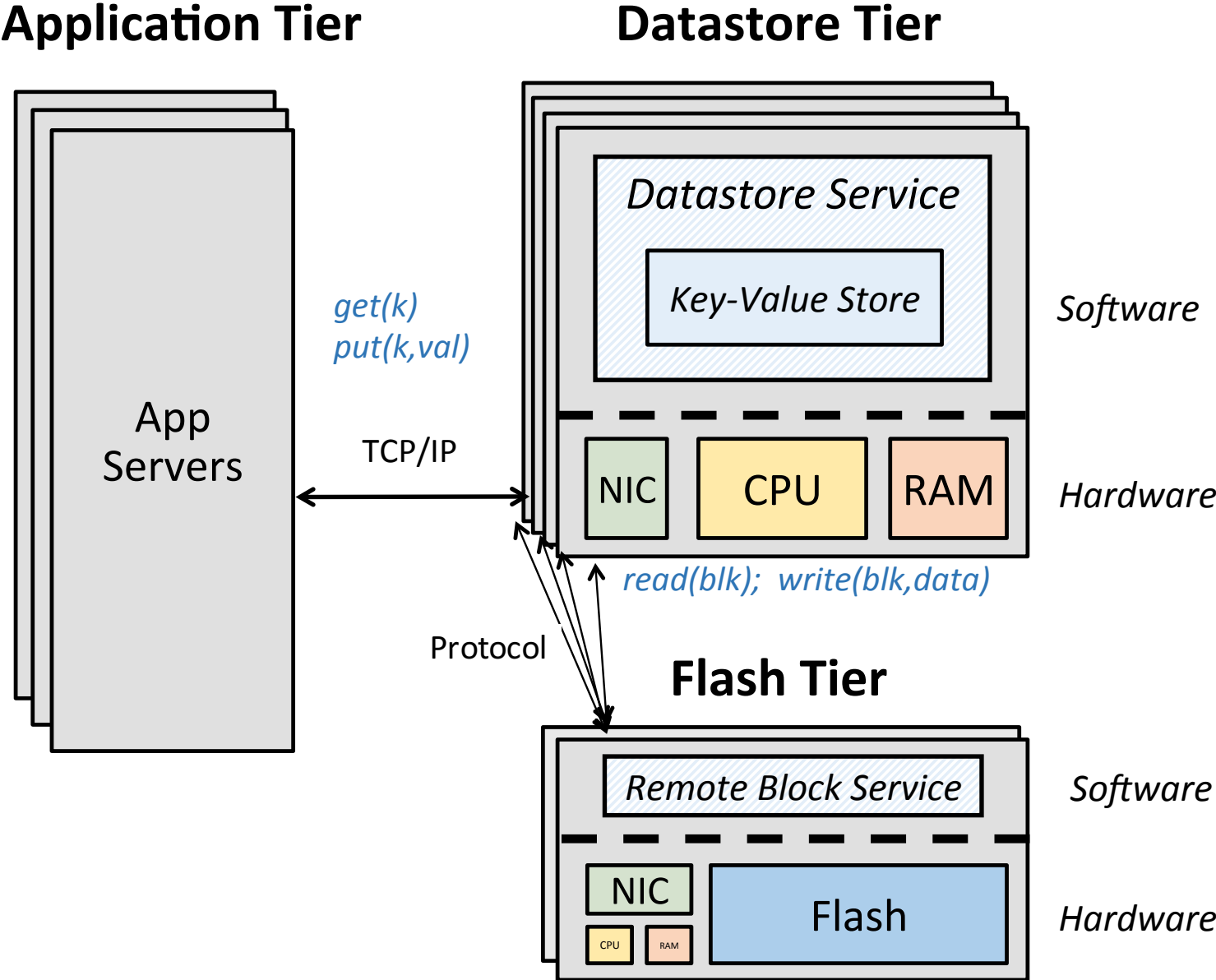


# Local Flash Architecture



Provision Flash and CPU in a dependent manner.

# Disaggregated Flash Architecture



# Contributions

For real applications at Facebook, we analyze:

1. What is the performance overhead of remote Flash using existing protocols?
2. What optimizations improve performance?
3. When does disaggregating Flash lead to resource efficiency benefits?

# Flash Workloads at Facebook

- Analyze IO patterns of real Flash-based Facebook applications
- Applications use RocksDB, a key-value store with a log structured merge tree architecture

	IOPS/TB	IO size
Read	2K – 10K	10KB – 50KB
Write	100 – 1K	500KB – 2MB

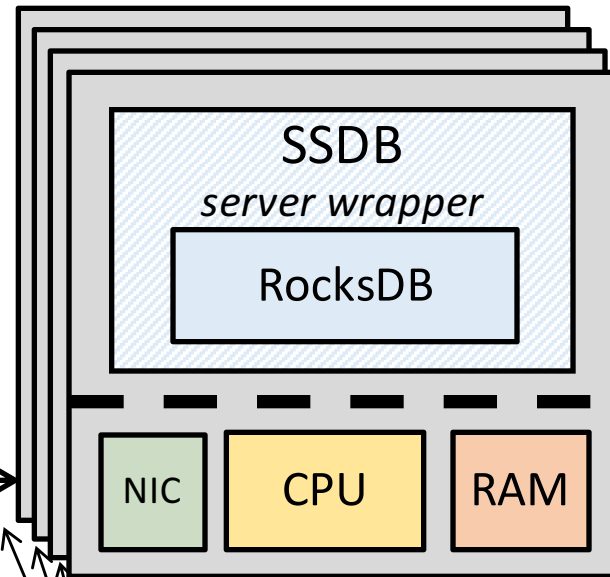
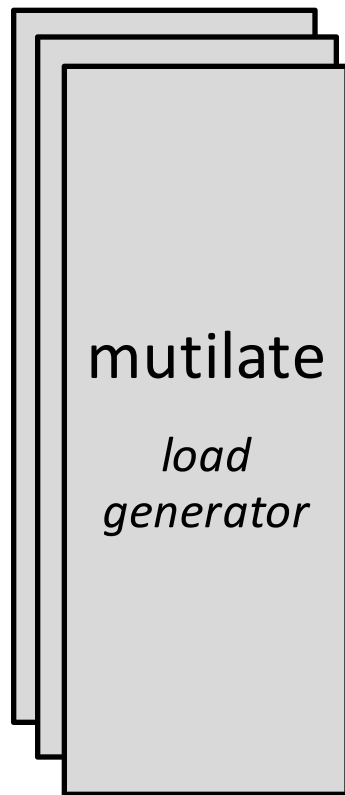
*Lots of random reads*

*Large, bursty writes*

# Workload Analysis

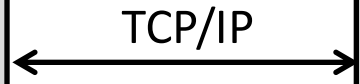
Application Tier

Datastore Tier



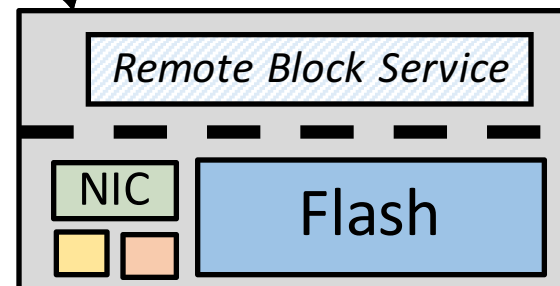
*Software*

*Hardware*



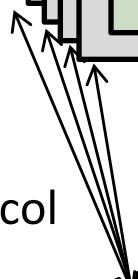
Protocol

Flash Tier



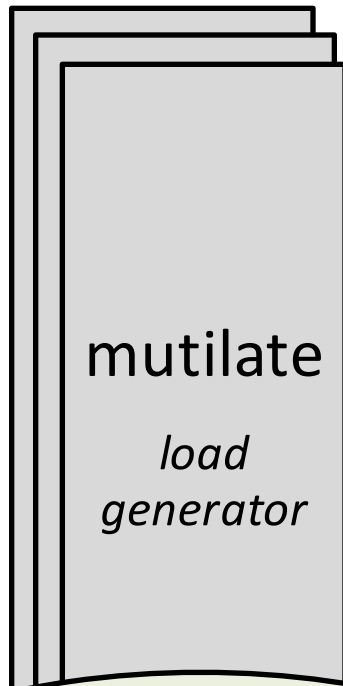
*Software*

*Hardware*

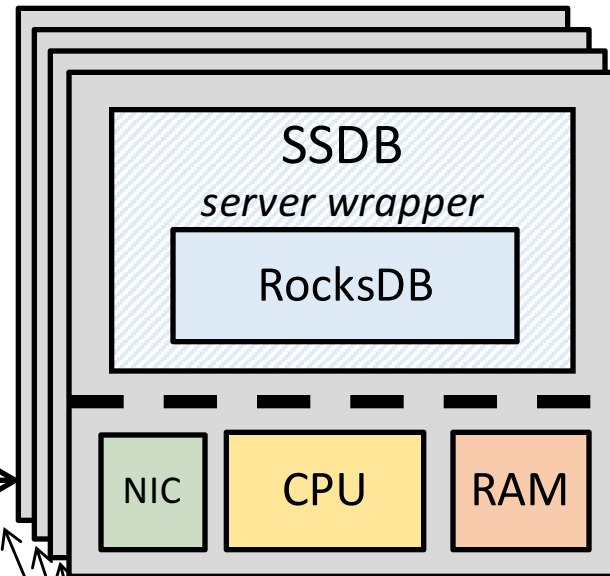


# Workload Analysis

## Application Tier



## Datastore Tier



*Software*

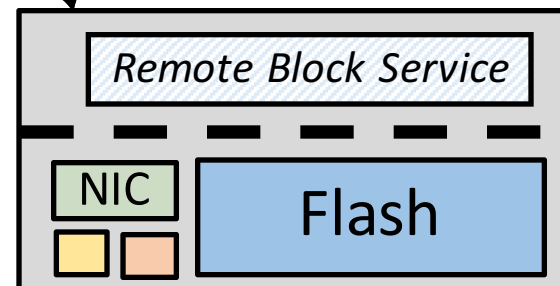
*Hardware*

TCP/IP



iSCSI

## Flash Tier



*Software*

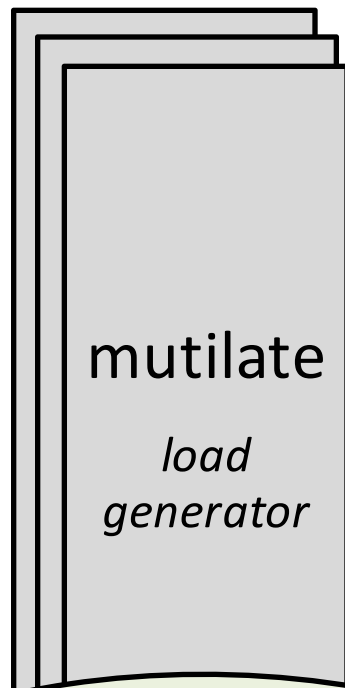
*Hardware*



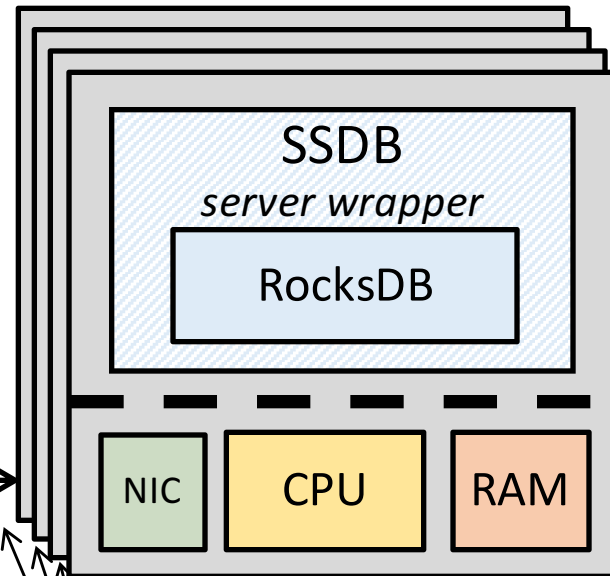
iSCSI is a standard network storage protocol that transports block storage commands over TCP/IP

# Workload Analysis

Application Tier



Datastore Tier



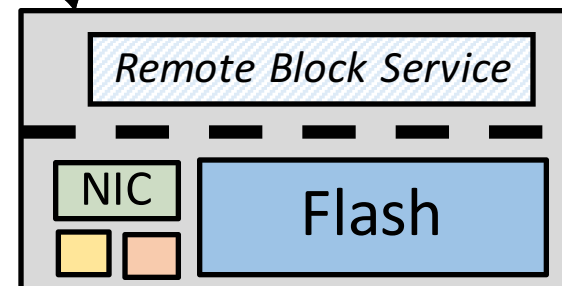
Software

Hardware

TCP/IP

iSCSI

Flash Tier



Software

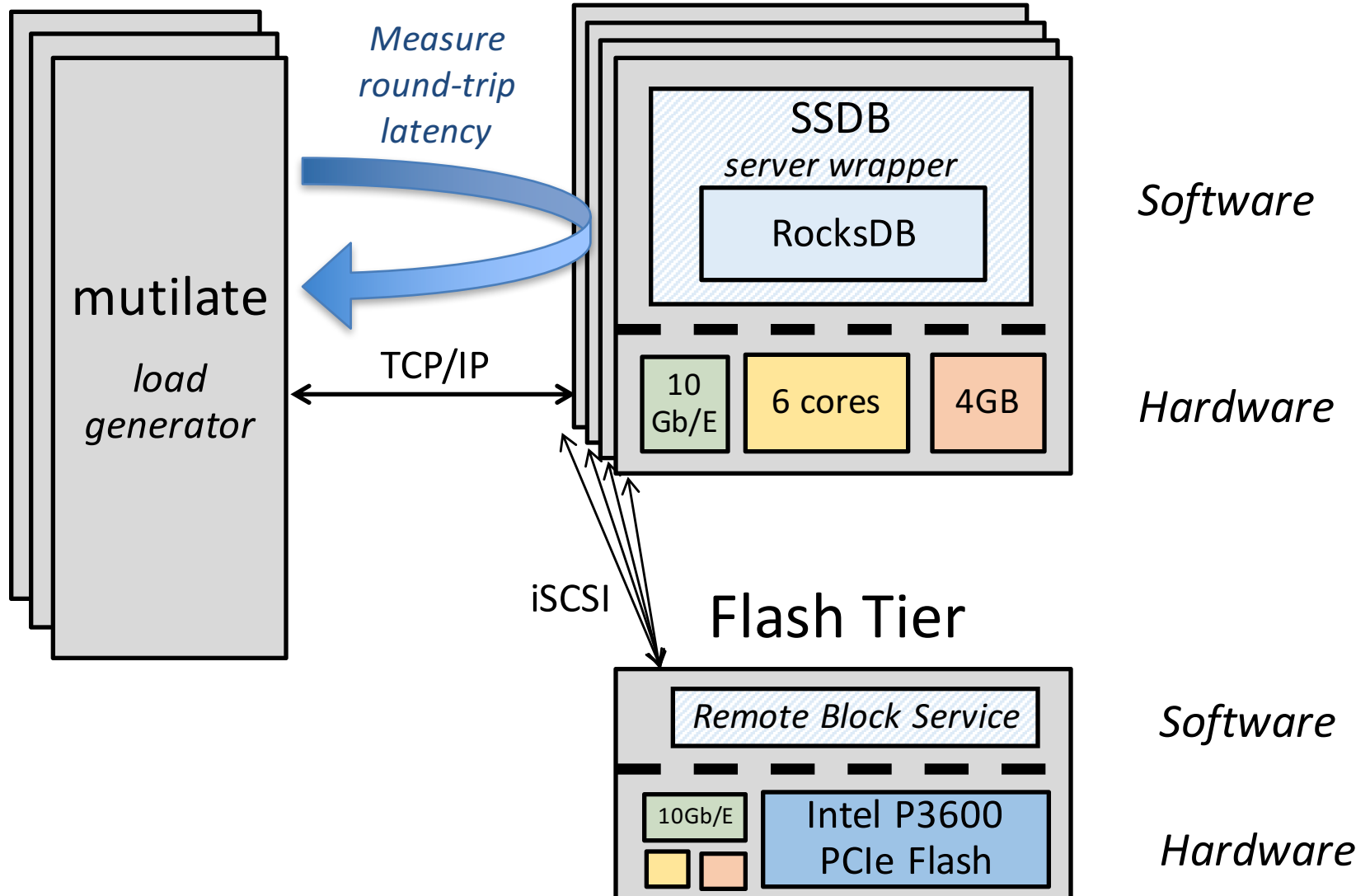
Hardware

- ✓ Transparent to application
- ✓ Runs on commodity network
- ✓ Scales datacenter-wide

# Workload Analysis

Application Tier

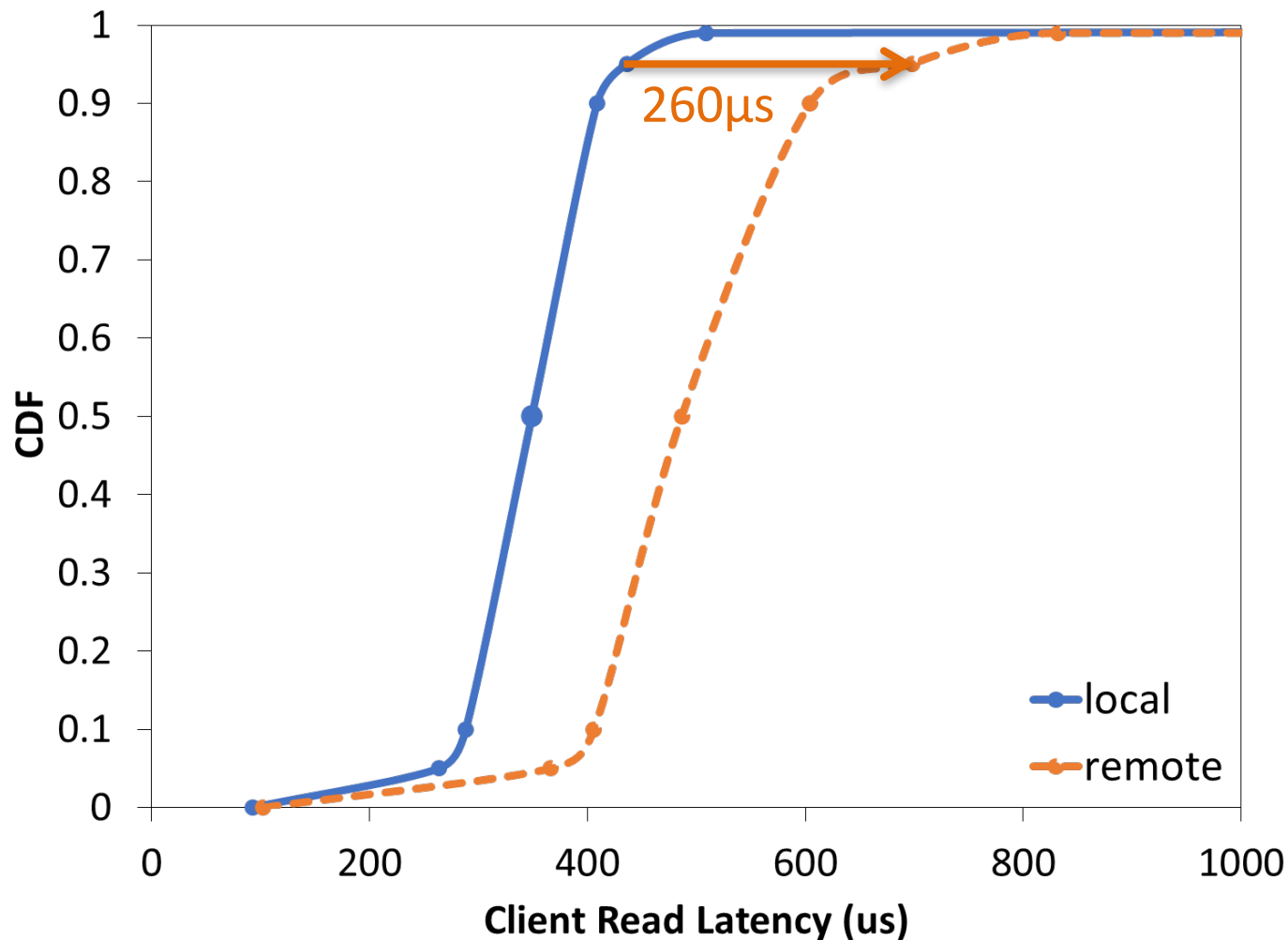
Datastore Tier





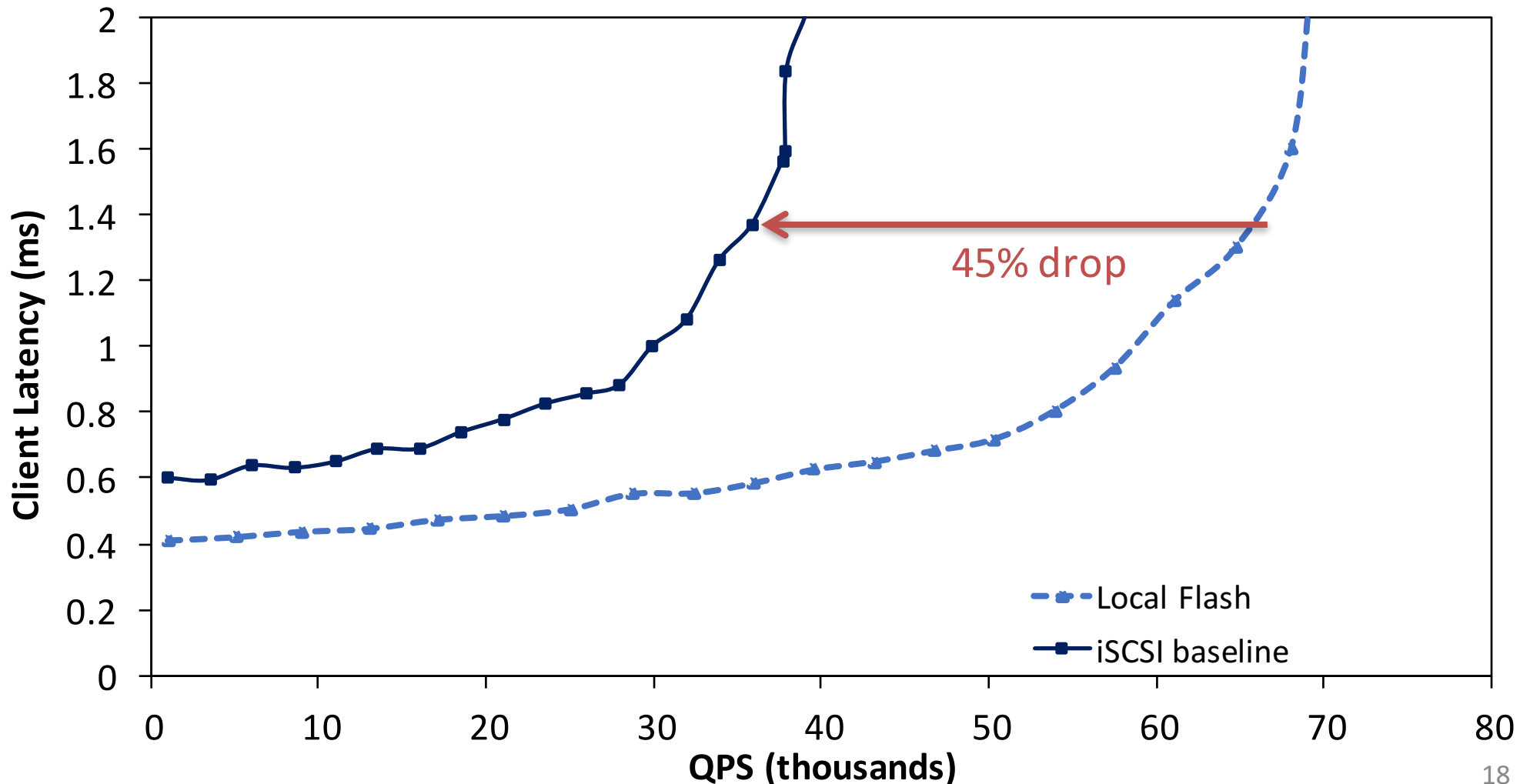
# Unloaded Latency

- Remote access with iSCSI adds  $260\mu\text{s}$  to p95 latency, tolerable for our target application (latency SLO  $\sim 5\text{ms}$ )



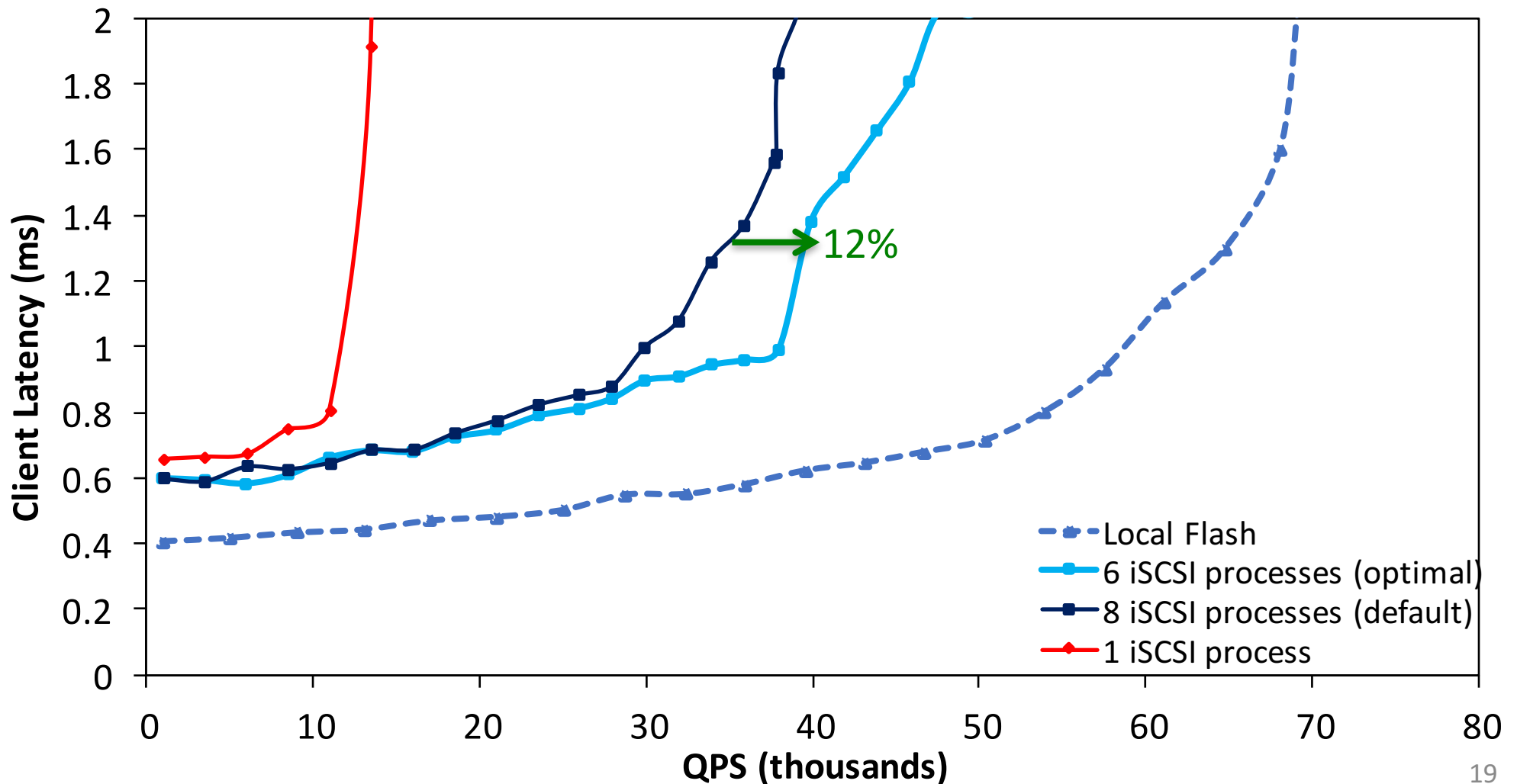
# Application Throughput

- 45% throughput drop with “out of the box” iSCSI Flash
- Need to optimize remote Flash server for higher throughput



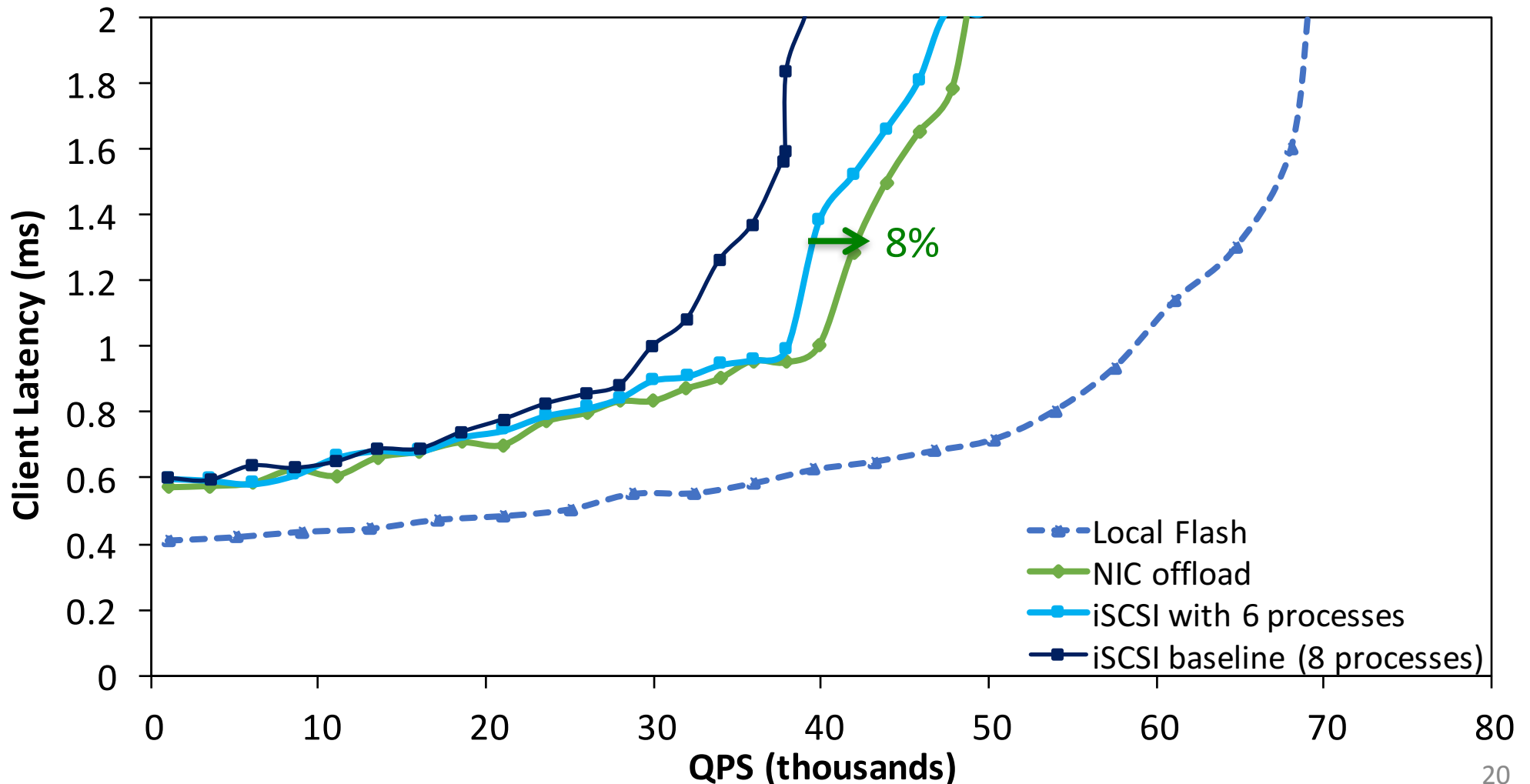
# Multi-process iSCSI

- Vary number of iSCSI processes that issue IO
- Want enough parallelism, avoid scheduling interference



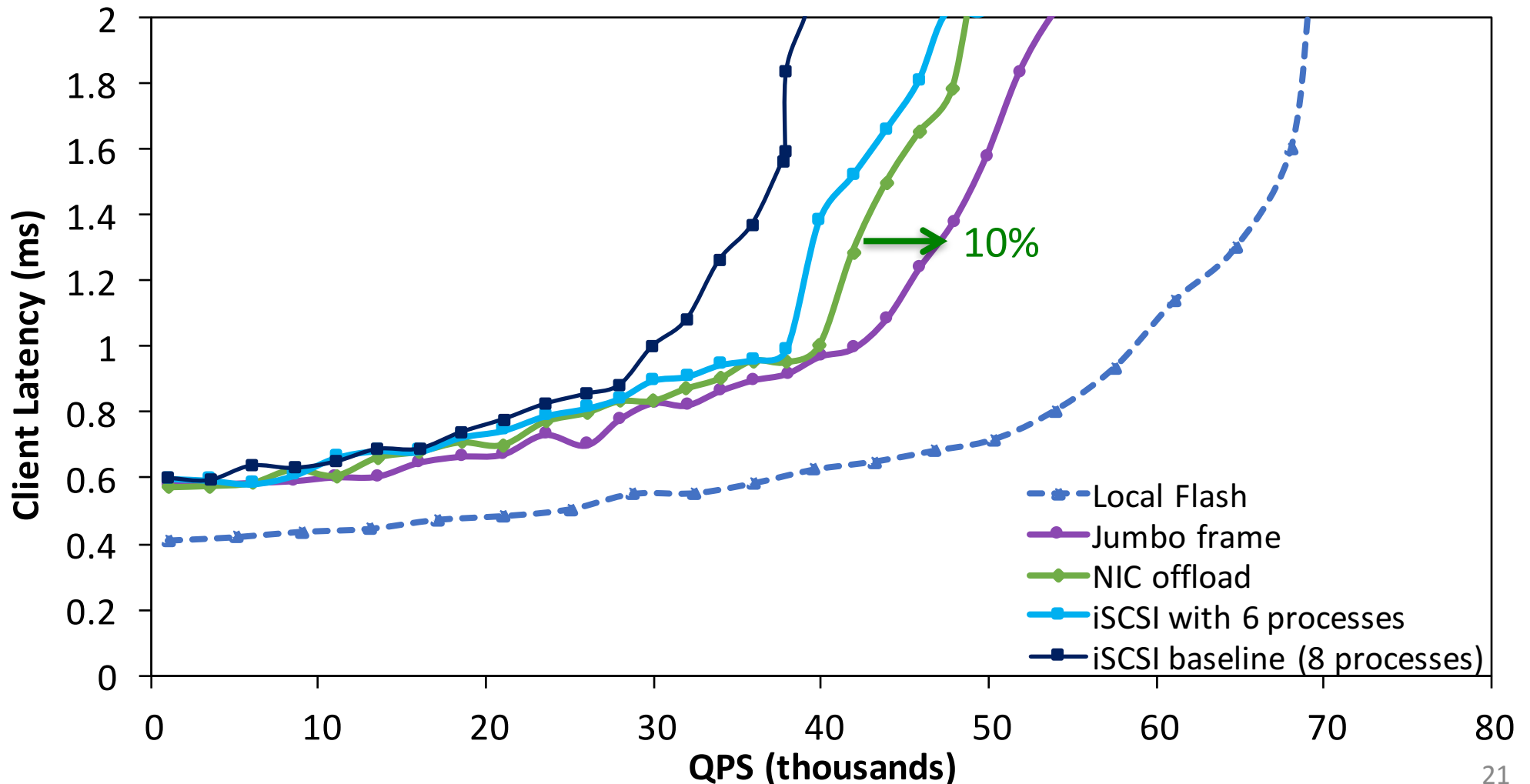
# NIC offloads

- Enable NIC offloads for TCP segmentation (TSO/LRO) to reduce CPU load on Flash server and datastore server



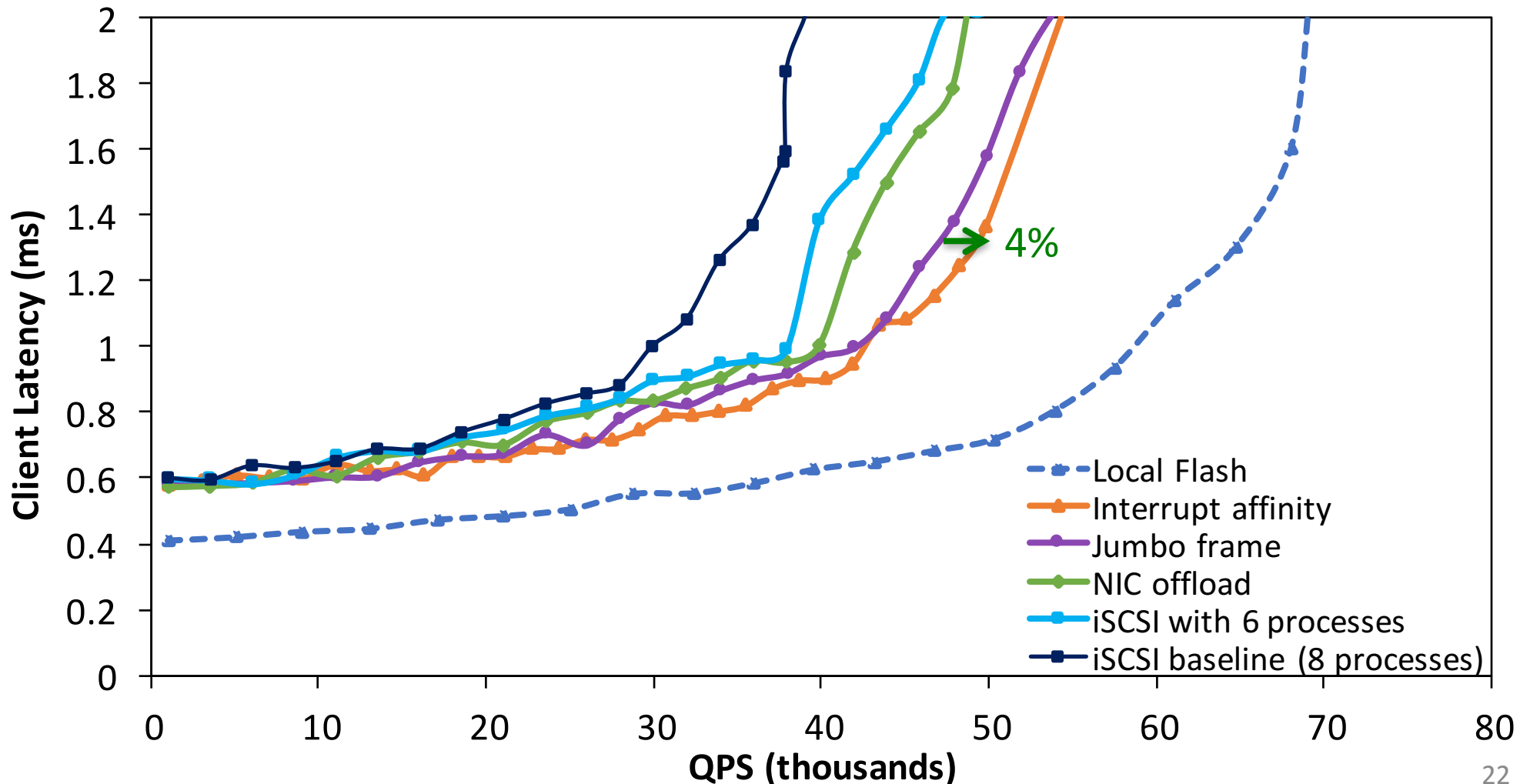
# Jumbo Frames

- Jumbo frames further reduce overhead by reducing segmentation altogether (max MTU 9kB)



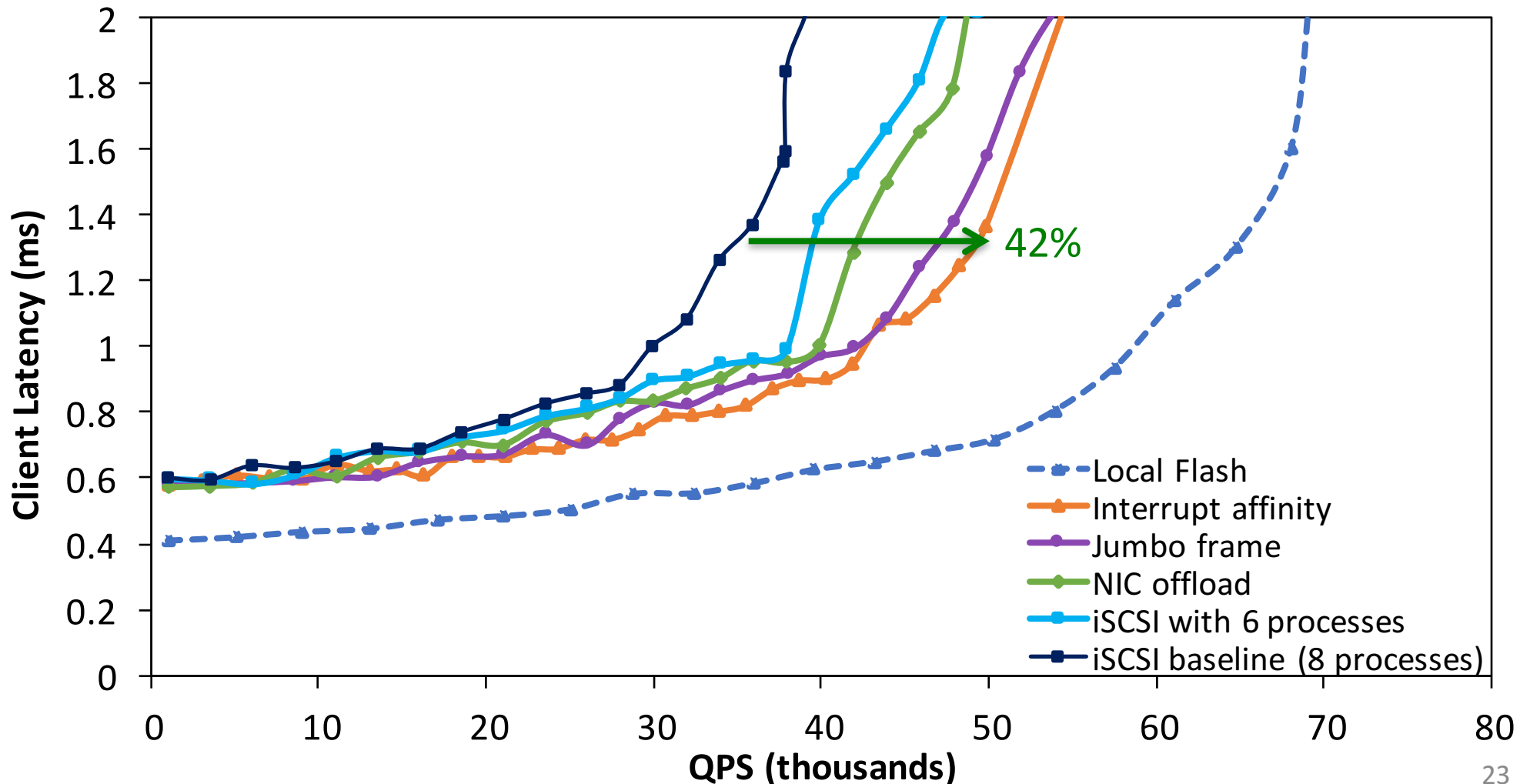
# Interrupt Affinity Tuning

- Steer NIC interrupts to core handling TCP connection and Flash interrupts to cores issuing IO commands



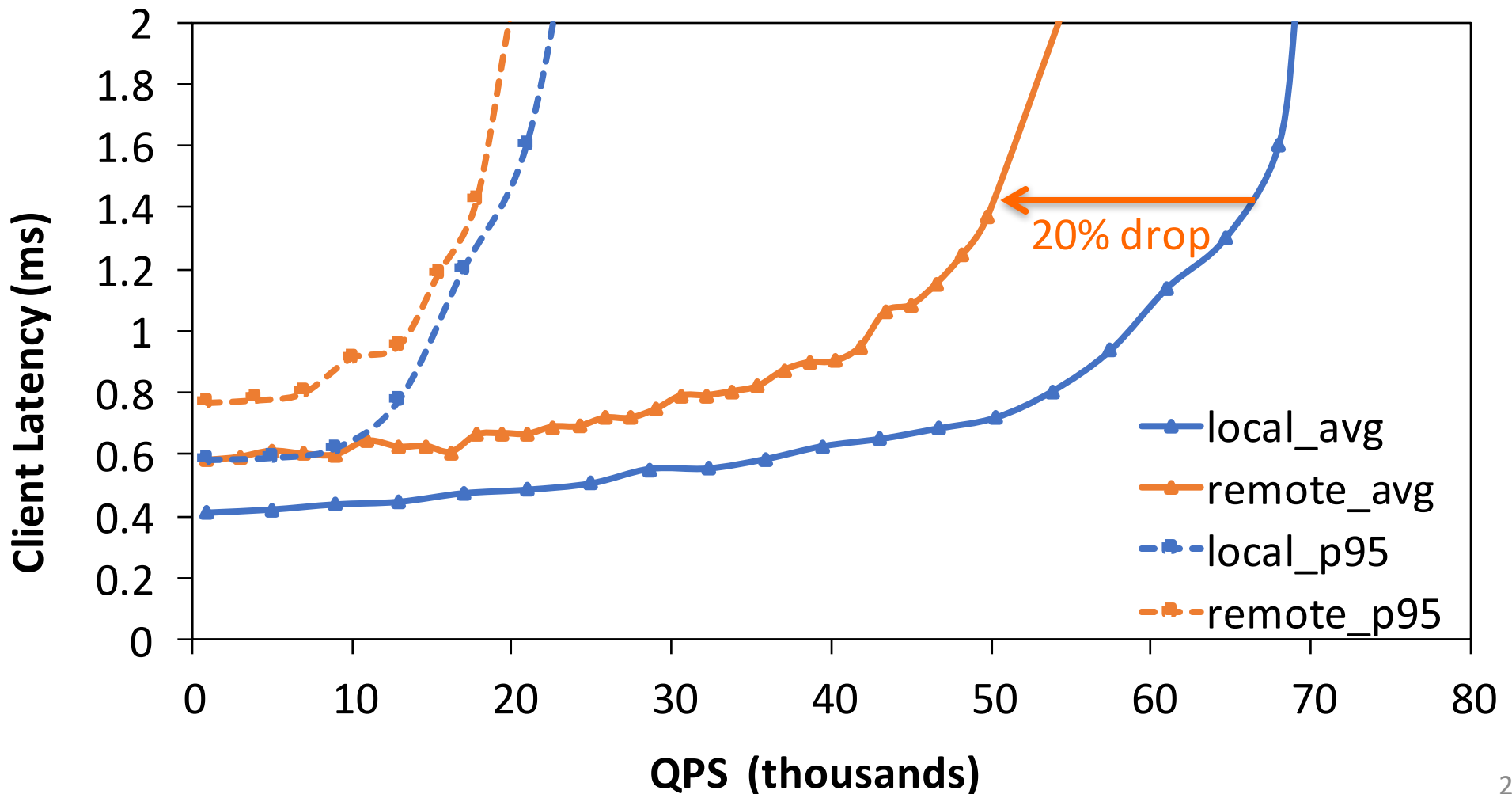
# Optimized Application Throughput

- Steer NIC interrupts to core handling TCP connection and Flash interrupts to cores issuing IO commands



# Application Throughput

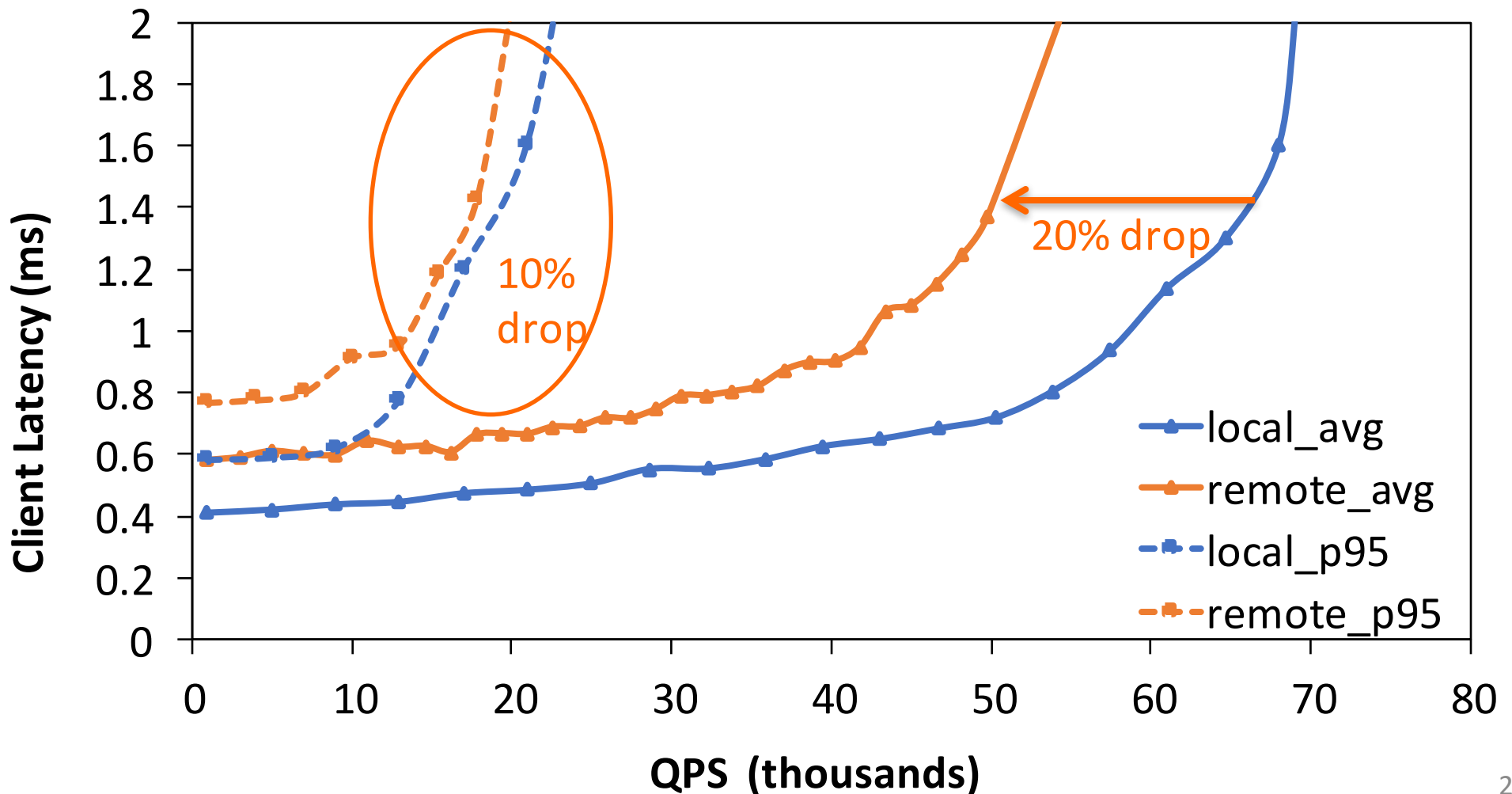
- 20% drop in application throughput, on average





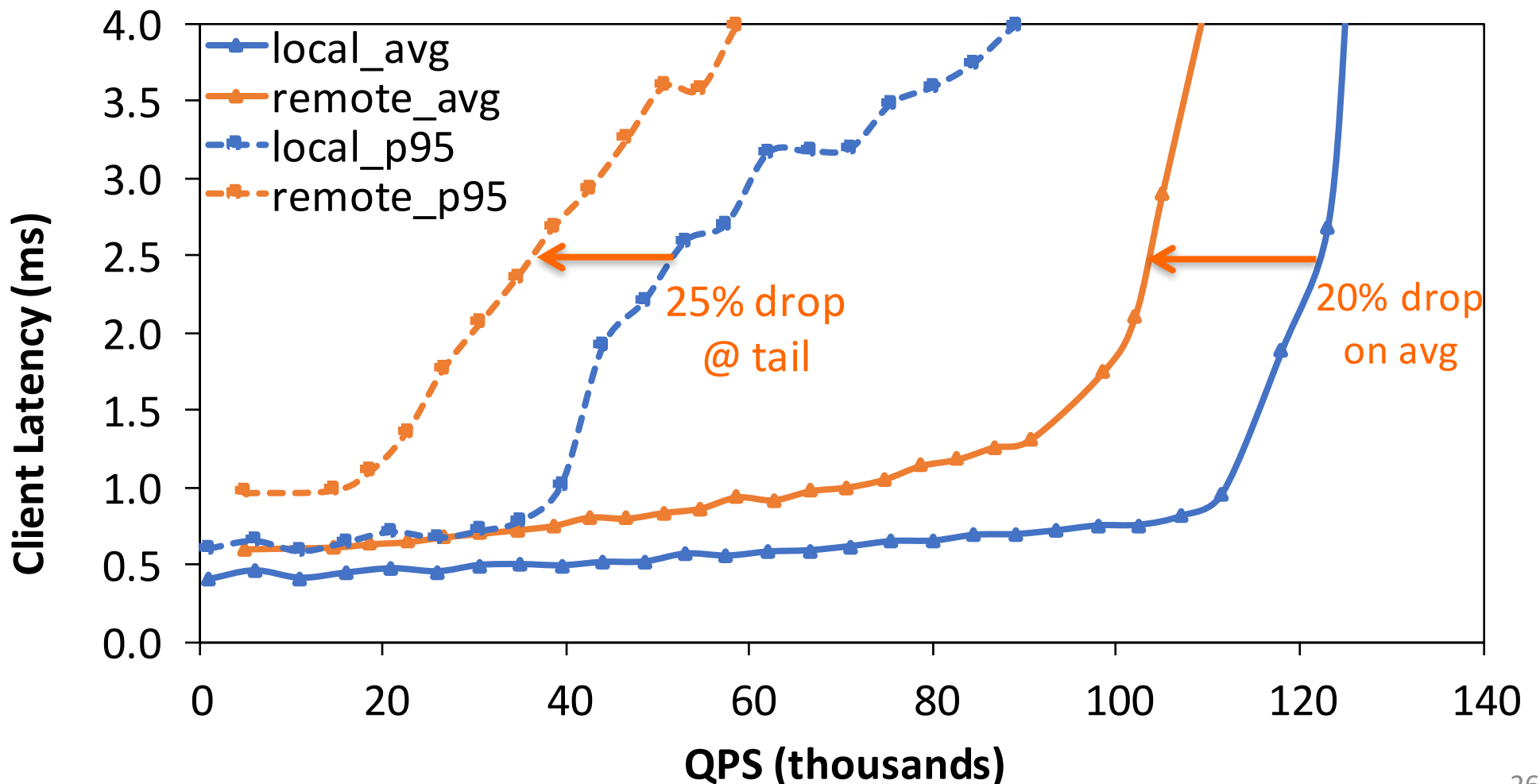
# Application Throughput

- At the tail, overhead of remote access is masked by other factors like write interference on Flash



# Sharing Remote Flash

- Sharing Flash among 2 or more tenants leads to more write interference → degrades tail performance

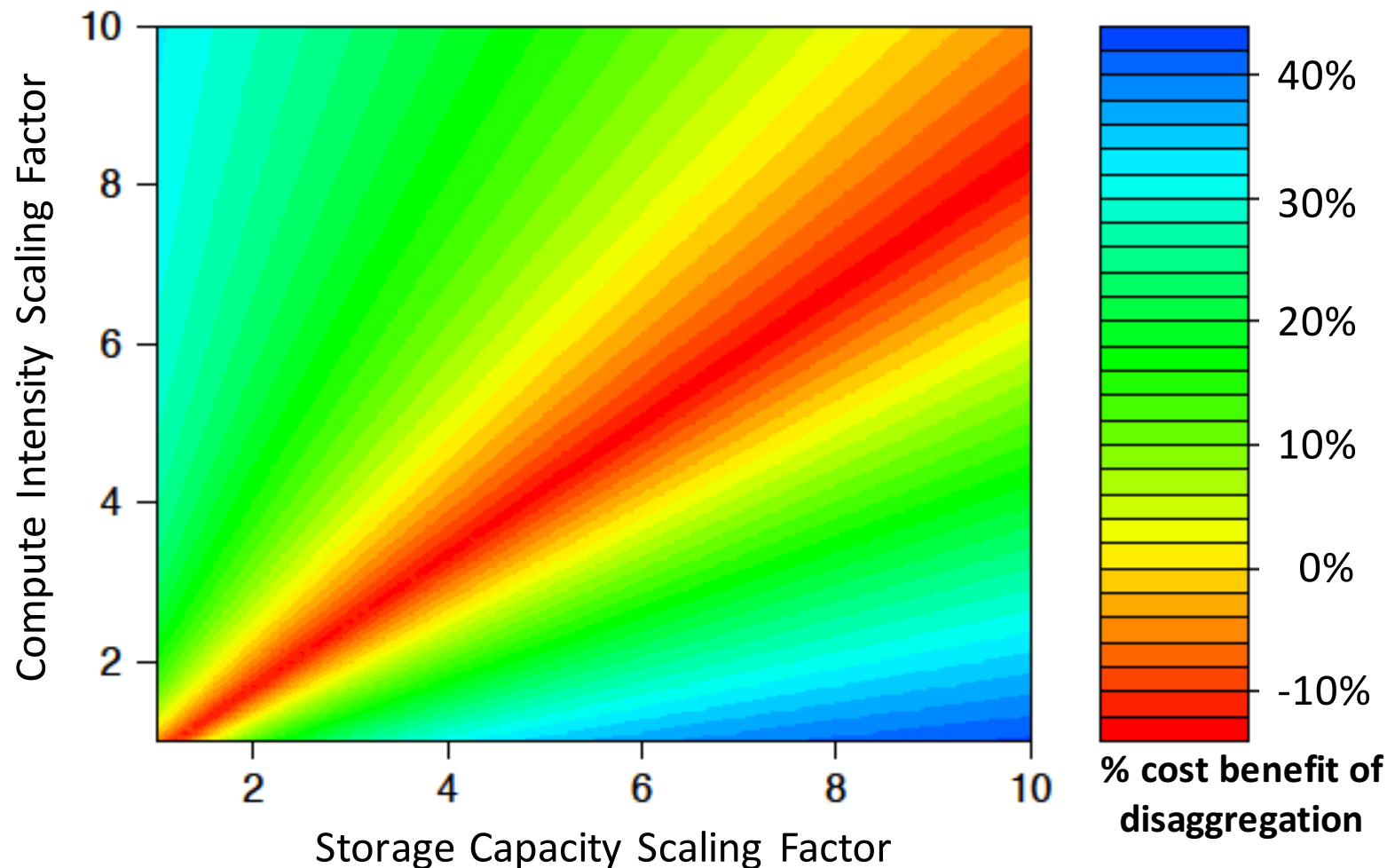


# Disaggregation Benefits

- Make up for throughput loss by *cost-effectively* scaling resources with disaggregation
- Improve overall resource utilization
- Formulate cost model to quantify benefits

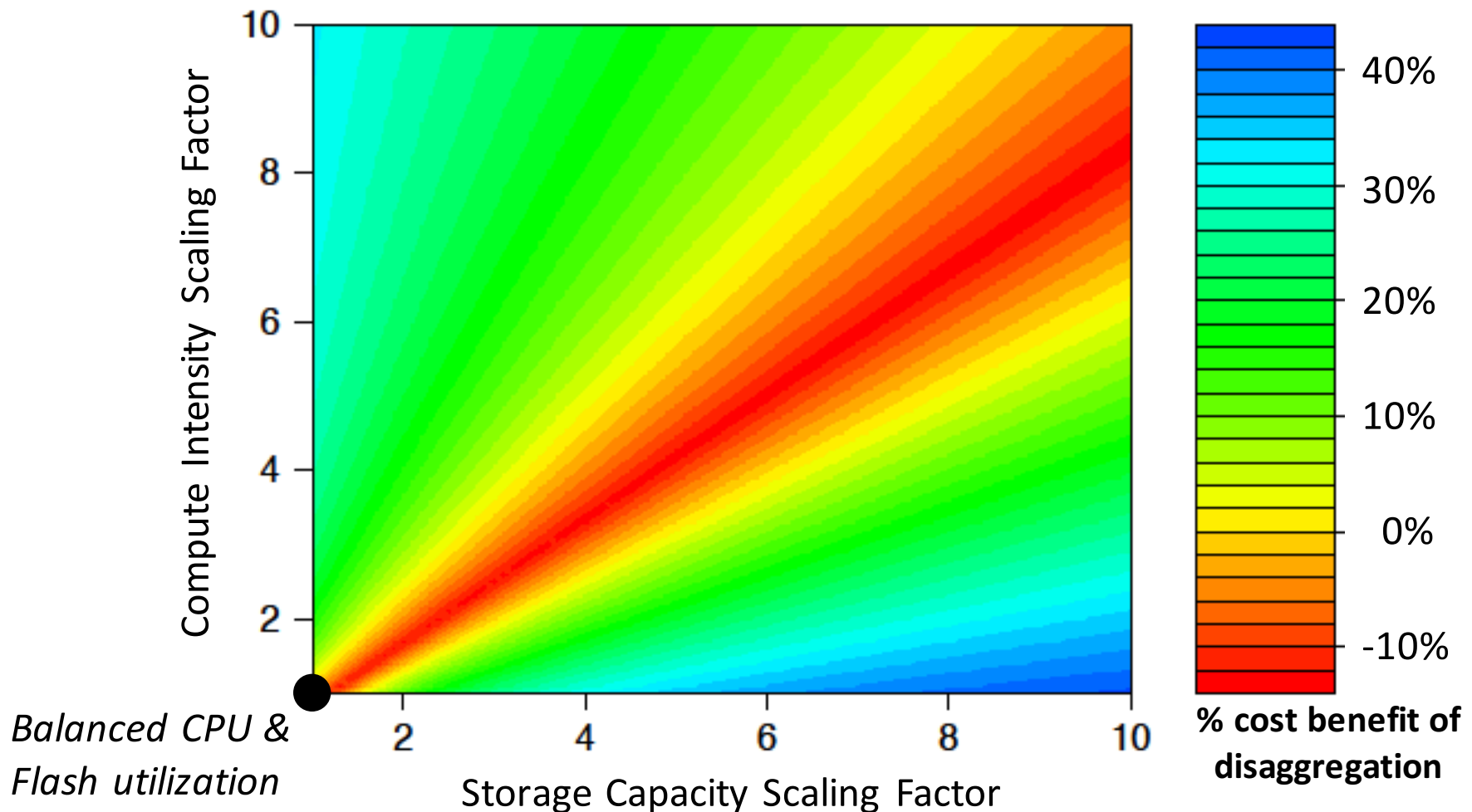
# Resource Savings

- Resource savings of disaggregated vs. local Flash architecture as app requirements scale



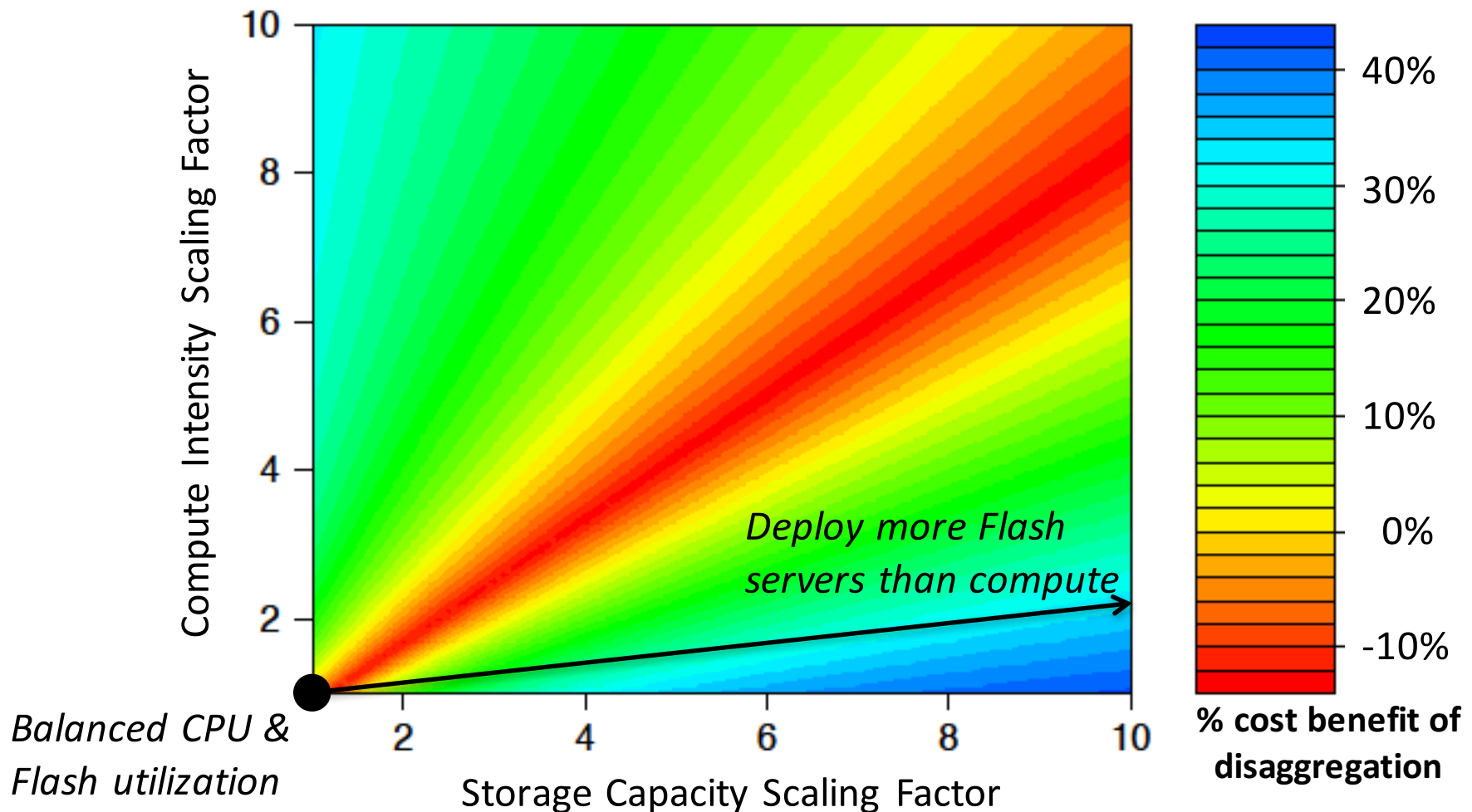
# Resource Savings

- Resource savings of disaggregated vs. local Flash architecture as app requirements scale



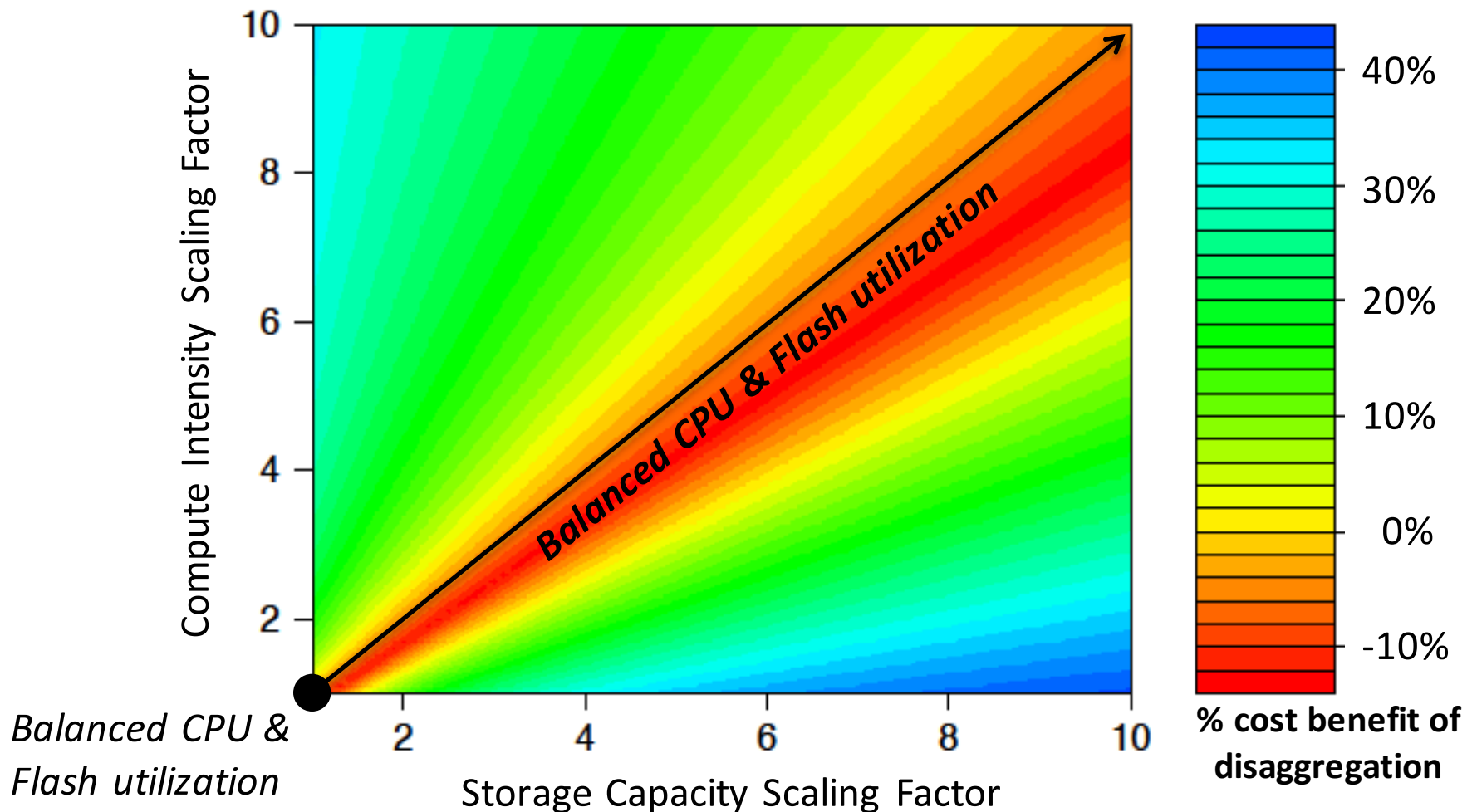
# Resource Savings

- When storage scales at higher rate than compute, save resources by deploying Flash without as much CPU



# Resource Savings

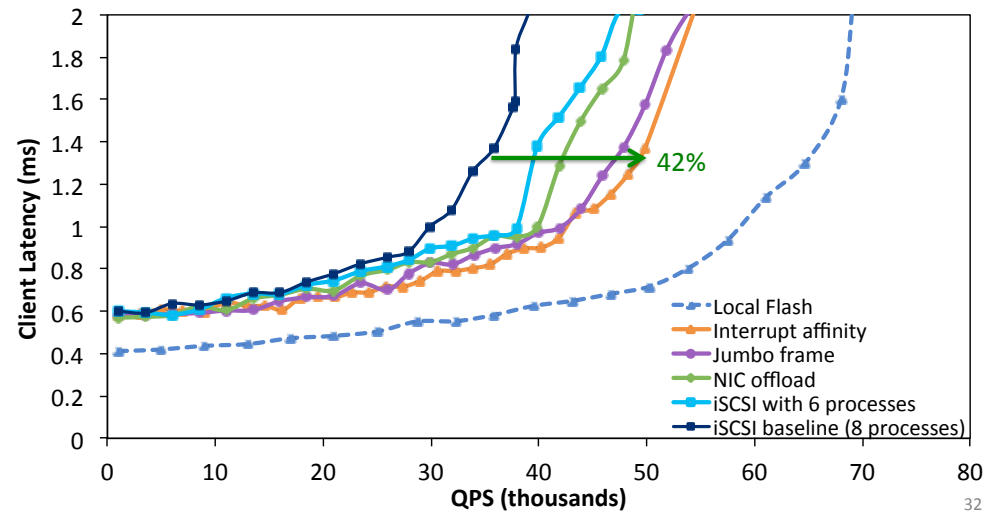
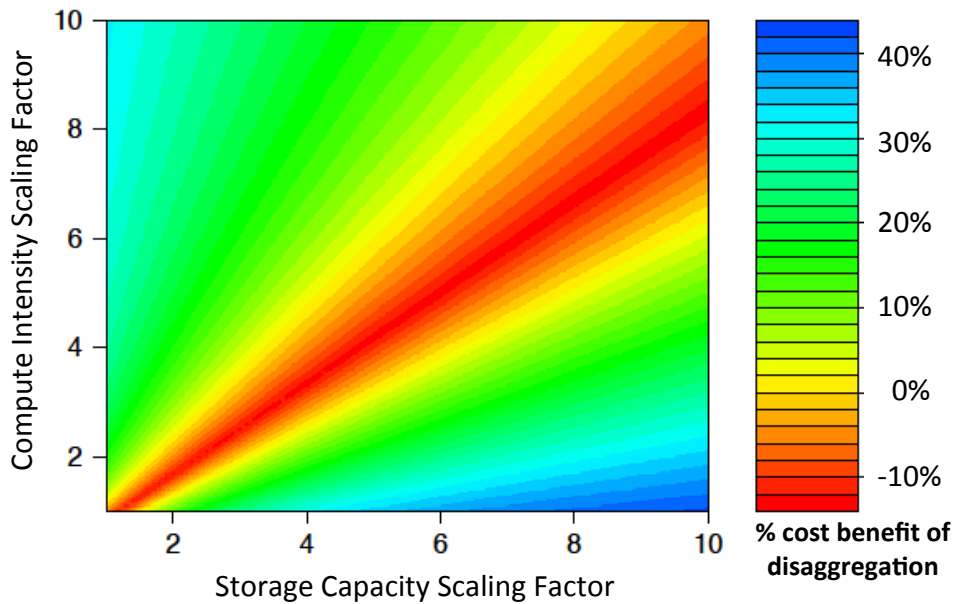
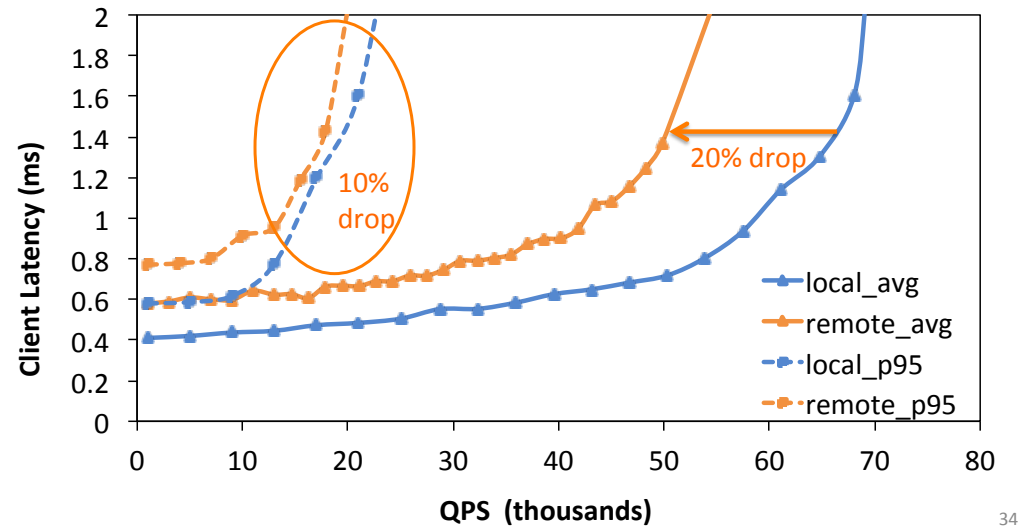
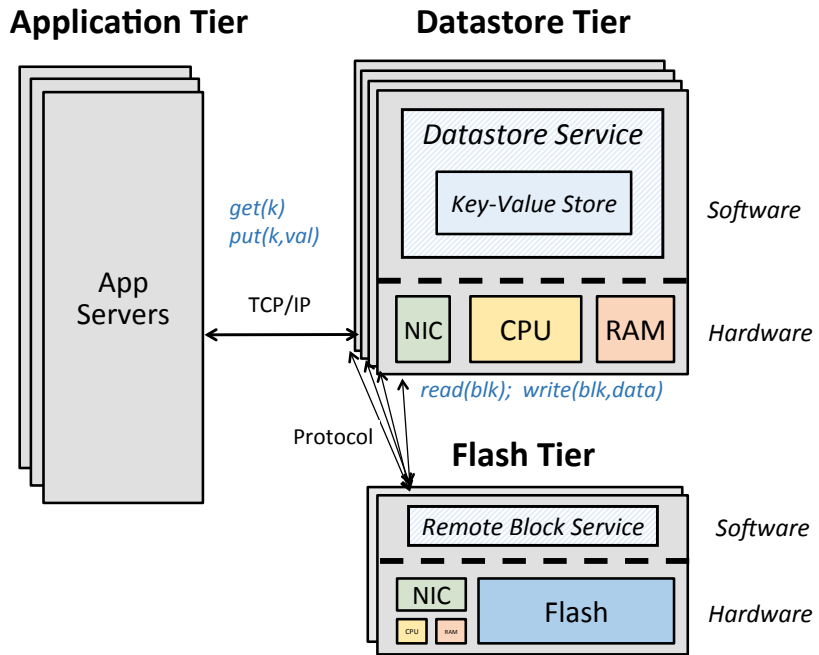
- When compute and storage demands remain balanced, no benefit with disaggregation



# Implications for System Design

- Dataplane:
  - Reduce compute overhead of network (storage) stack
    - Optimize TCP/IP processing
    - Use a light-weight protocol
  - Provide isolation mechanisms for shared remote Flash
- Control plane:
  - Policies for allocating and sharing remote Flash
    - Important to consider write IO patterns of applications





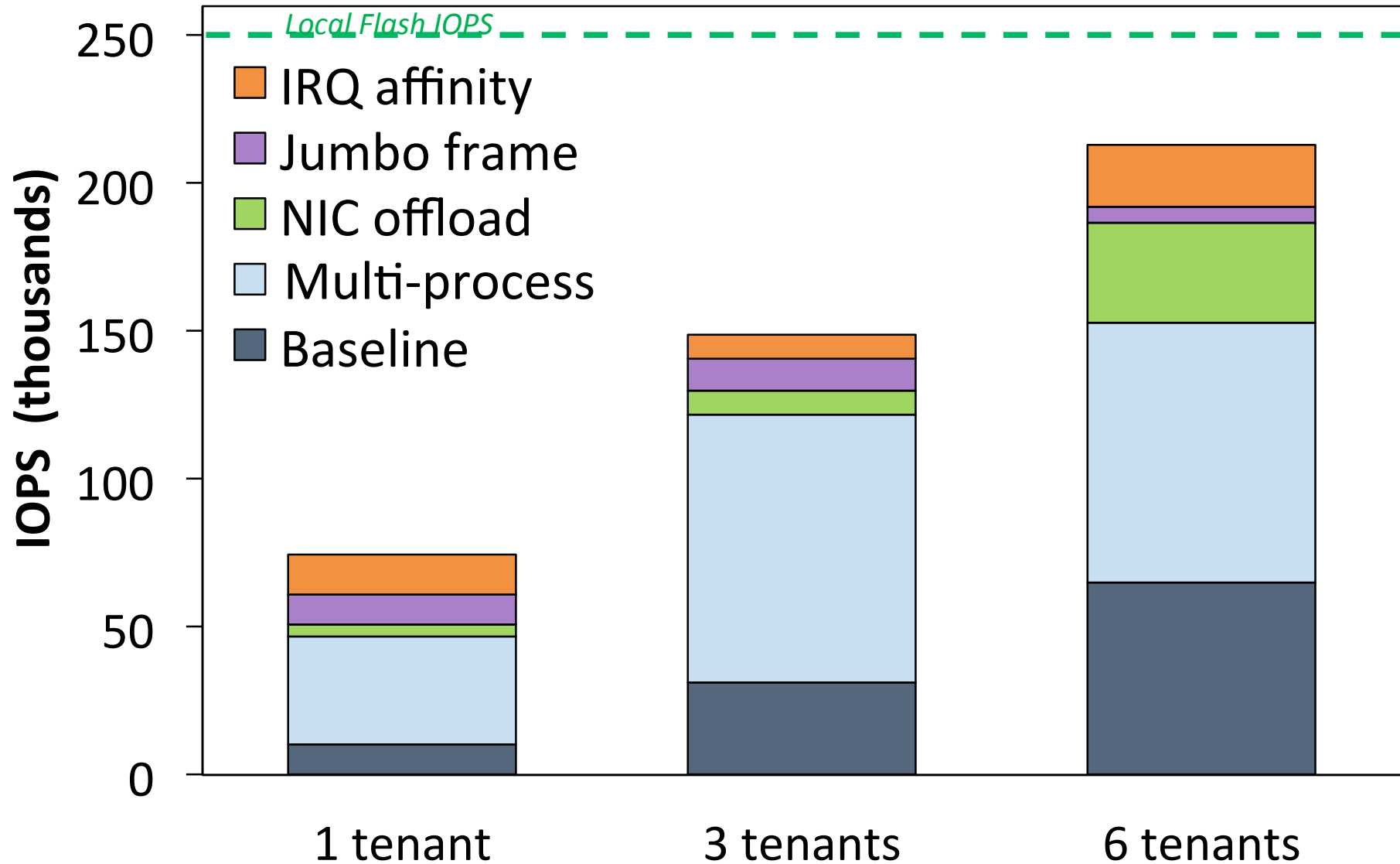
# Conclusion

- Disaggregating Flash is beneficial because it allows us to cost-effectively scale resources:
  - Improve overall resource efficiency
  - Compensate for 20% throughput overhead by independently deploying application resources
- System tuning improves performance ~40%, more opportunities if redesign software stack

Backup

# Remote Flash IOPS

IO-intensive benchmark: 4kB random reads



# Cost Model

$$C_{direct} = \max \left( \frac{GB_t}{GB_s}, \frac{IOPS_t}{IOPS_s}, \frac{QPS_t}{QPS_s} \right) \cdot (f + c)$$

$$C_{disagg} = \max \left( \frac{GB_t}{GB_s}, \frac{IOPS_t}{IOPS_s} \right) \cdot (f + \delta) + \left( \frac{QPS_t}{QPS_s} \right) c$$

where:

$f$ : cost of Flash on a server

$c$ : cost of CPU, RAM and NIC on datastore server

$\delta$ : cost of CPU, RAM and NIC on Flash tier server,  
i.e. resource “tax” for disaggregation

$x_s$ :  $x$  provided by a single server,  $x = \{GB, IOPS, QPS\}$

$x_t$ :  $x$  required in total for the application

# Related Work

- Disaggregated disk storage:
  - Petal [ASPLOS'96], Parallax [HotOS'05], Blizzard [NSDI'14]
- Disaggregated Flash as distributed shared log:
  - CORFU [NSDI'12], FAWN [SOSP'09]
- Disaggregated memory:
  - Memory blade servers (Lim et al.) [ISCA'09]
- Rack-scale disaggregation:
  - Pelican [OSDI'14], HP Moonshot, Intel Rack-Scale