# From application requests to Virtual IOPs: Provisioned key-value storage with Libra

**David Shue**[*] and Michael J. Freedman

(*now at Google)

Shared Cloud

# Shared Cloud Storage
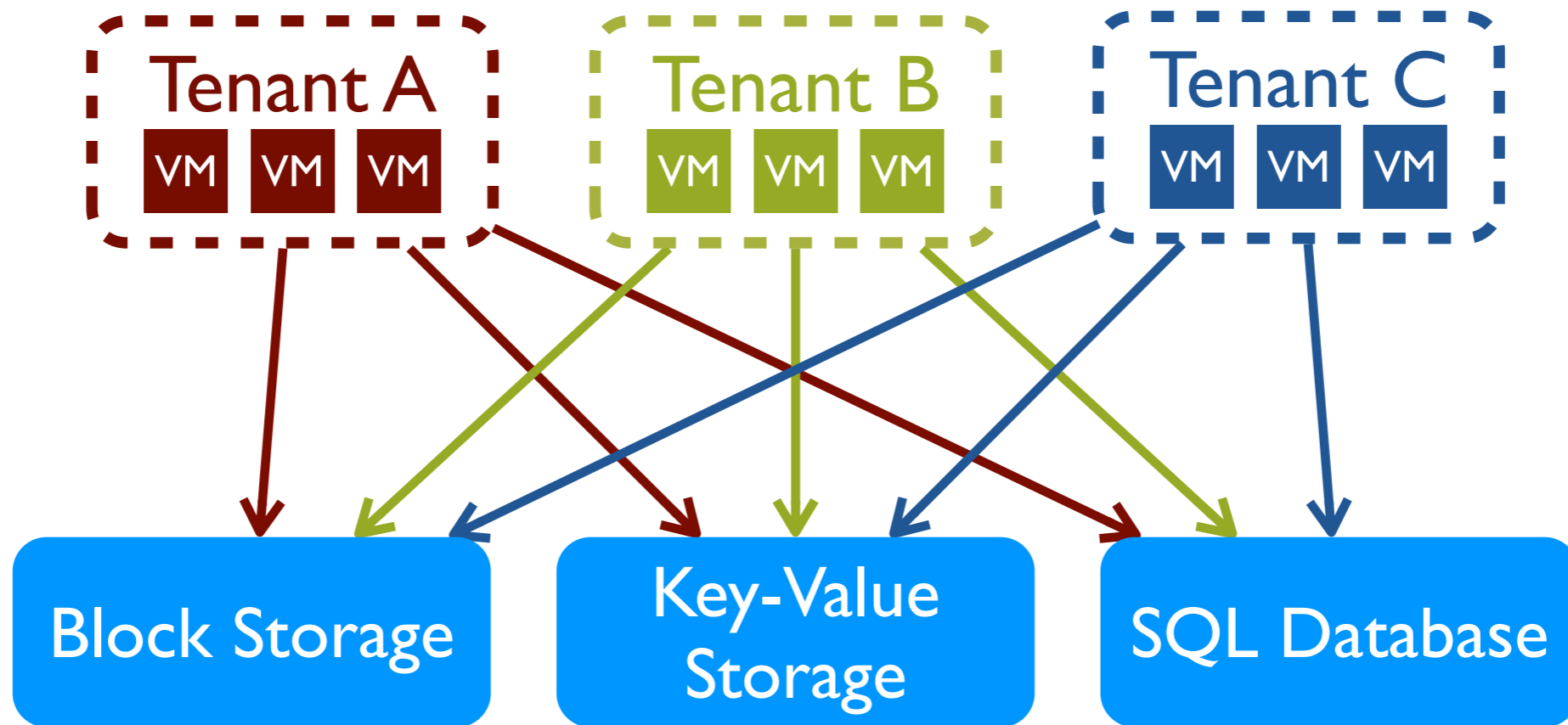
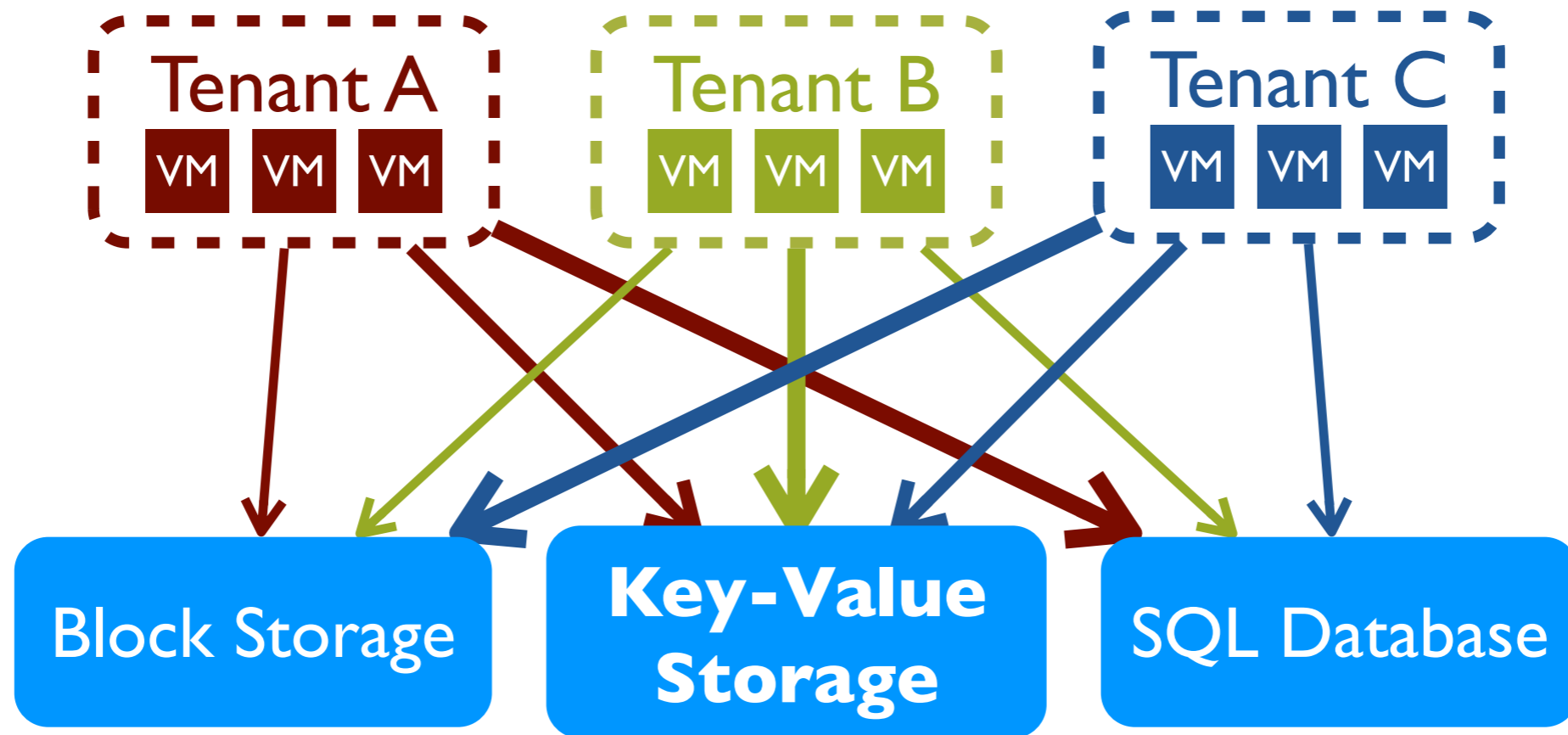| | | |
|---|---|---|
| **Tenant A** | **Tenant B** | **Tenant C** |
| VM VM VM | VM VM VM | VM VM VM |

**Block Storage**

**Key-Value Storage**

**SQL Database**

# Unpredictable Shared Cloud Storage

Tenant A

VM VM VM

Tenant B

VM VM VM

Tenant C

VM VM VM

Block Storage

**Key-Value Storage**

SQL Database

Disk IO-bound Tenants
SSD-backed storage

# Libra Contributions

- Libra IO Scheduler
  - Provisions *low-level IO allocations* for *app-request reservations* w/ high utilization.
  - Supports arbitrary object distributions and workloads.

- 2 key mechanisms
  - Track per-tenant app-request resource profiles.
  - Model IO resources with Virtual IOPs.

# Related Work

| | Storage Type | App-requests | Work Conserving | Media |
|---|---|---|---|---|
| Maestro | Block | N | N | HDD |
| mClock | Block | N | Y | HDD |
| FlashFQ | Block | N | Y | SSD |
| DynamoDB | Key-Value | Y | N | SSD |

# Provisioned Distributed Key-Value Storage

Global Reservation Problem
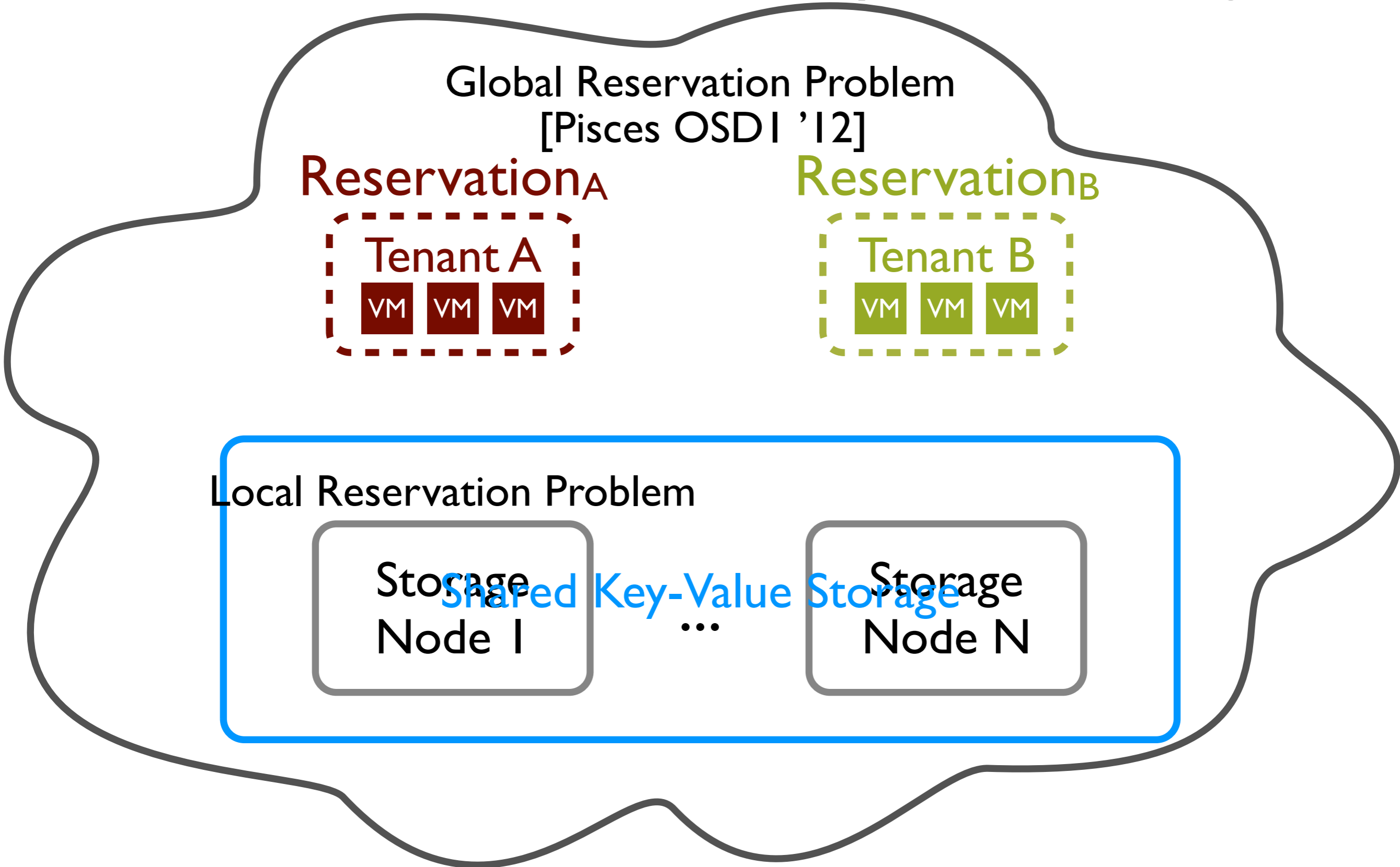[Pisces OSD1 '12]

Reservation$_A$

Reservation$_B$

Tenant A

VM VM VM

Tenant B

VM VM VM

Local Reservation Problem

Storage
Node 1

Shared Key-Value Storage
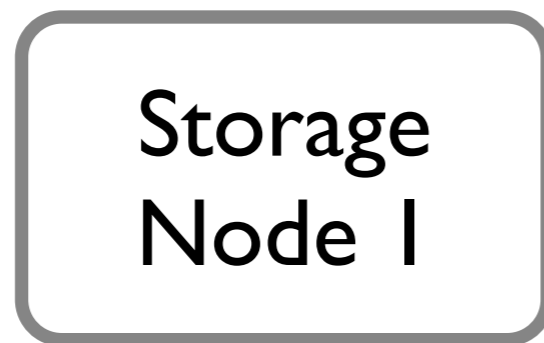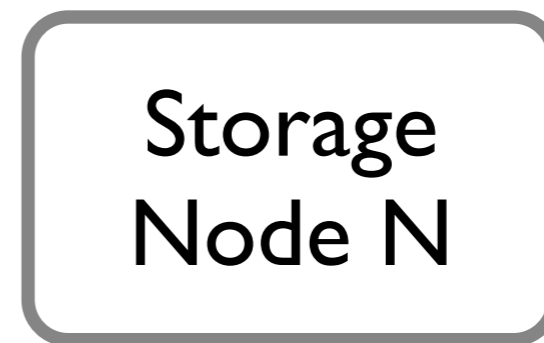
...

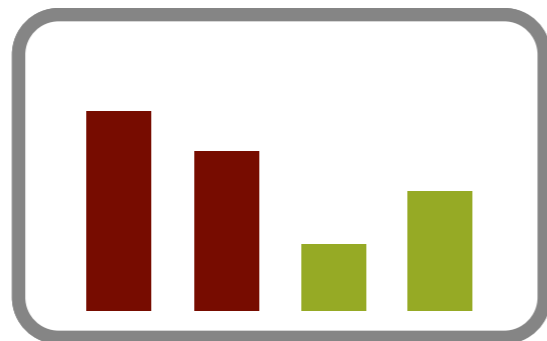Storage
Node N

# Provisioned Distributed Key-Value Storage

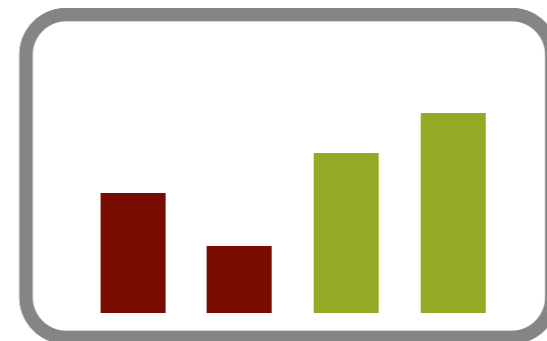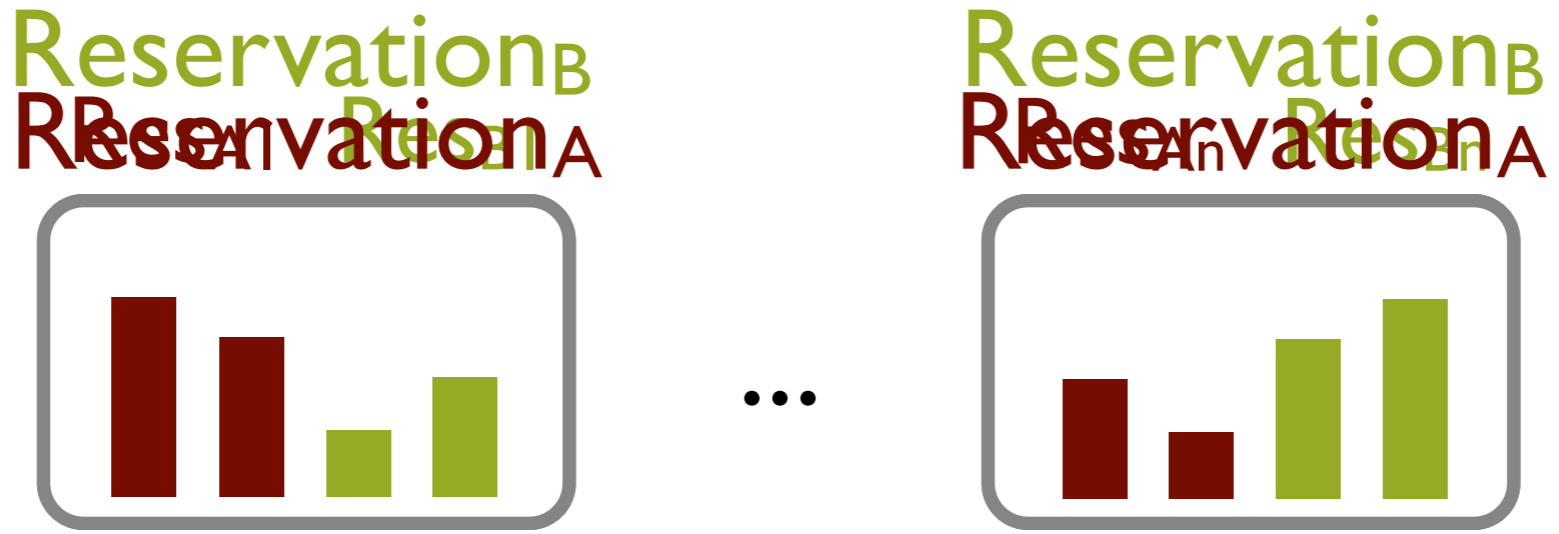# Provisioned Distributed Key-Value Storage
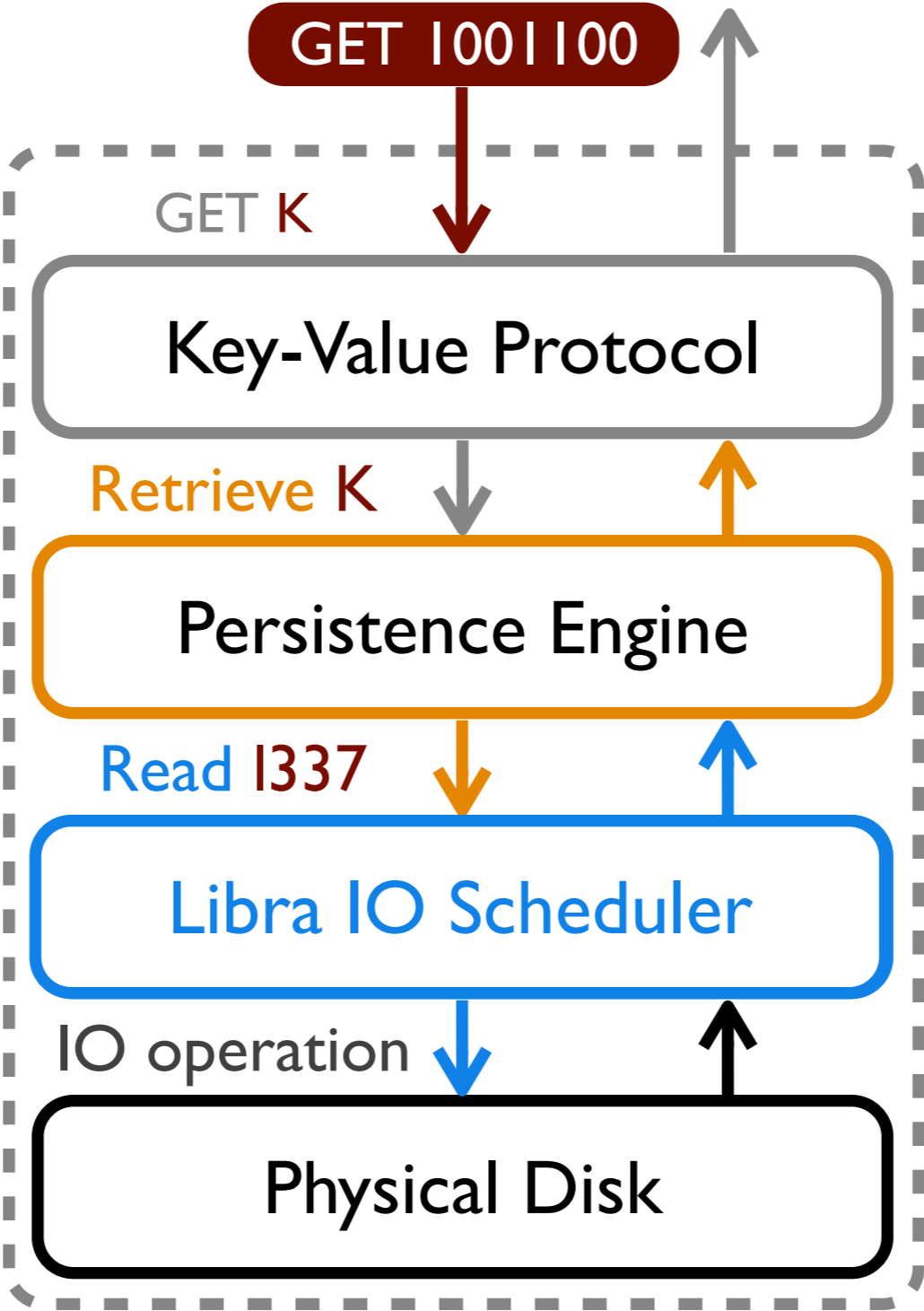
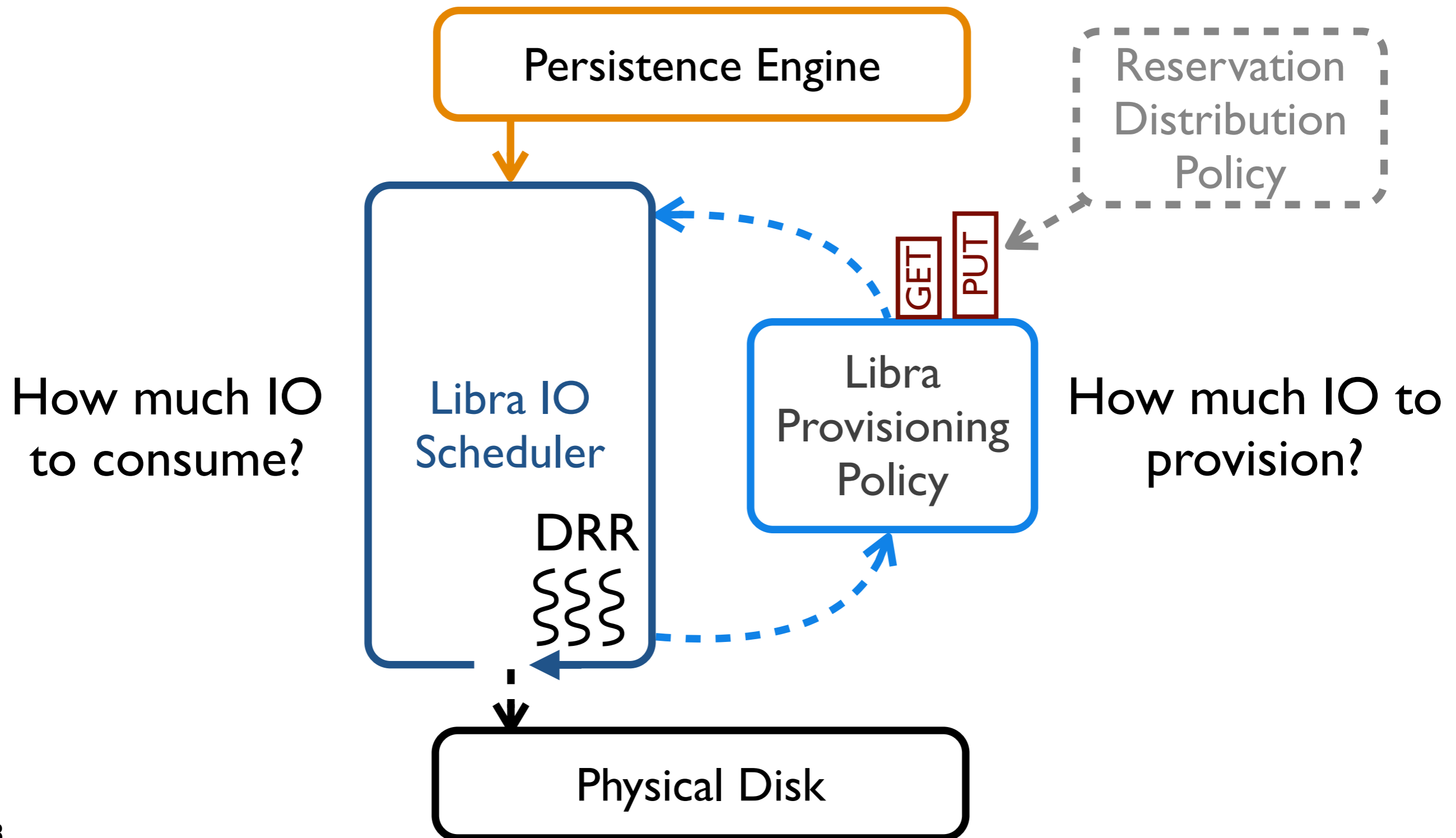ReservationA          ReservationB

# Provisioned Distributed Key-Value Storage

# Provisioned Local Key-Value Storage

# Libra Design

# Provisioning App-request Reservations is Hard

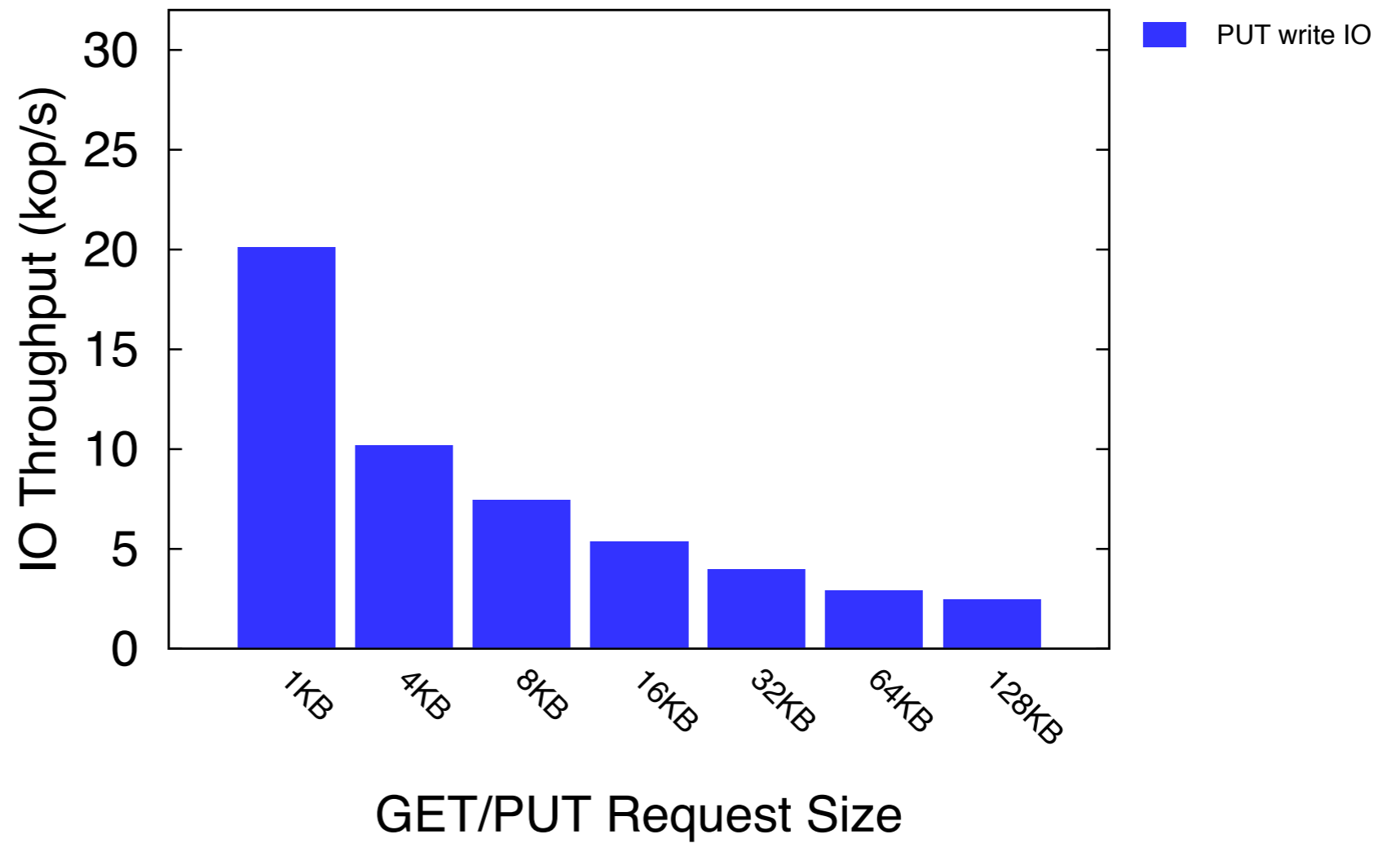**IO Amplification**

Track tenant app-request resource profiles
KB PUT = x KB Write

**IO Interference**

Validate Optimal IO
Under-rate to provision-able IO

**Non-linear IO Performance**
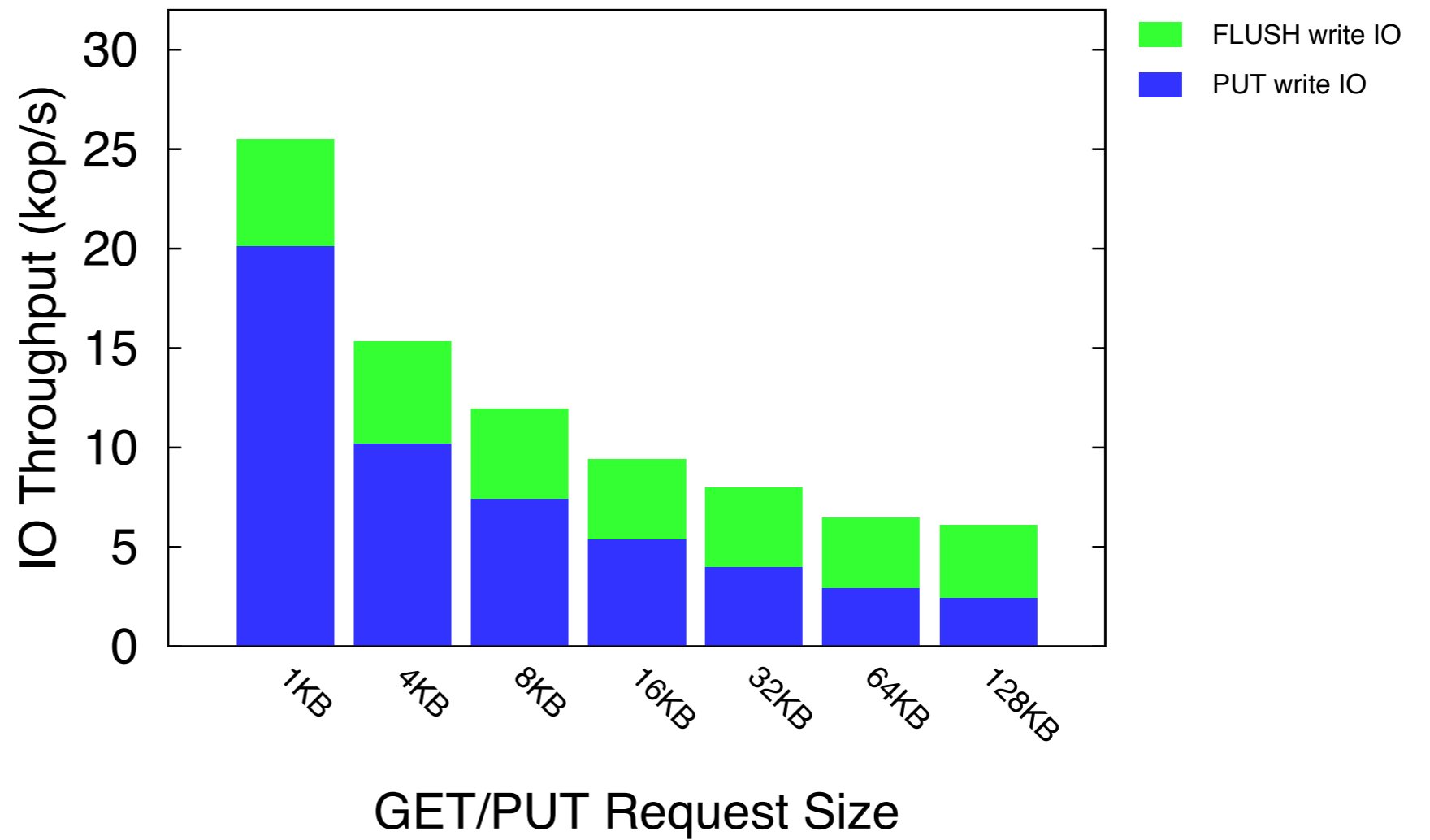
Model-IO with Virtual KOPs
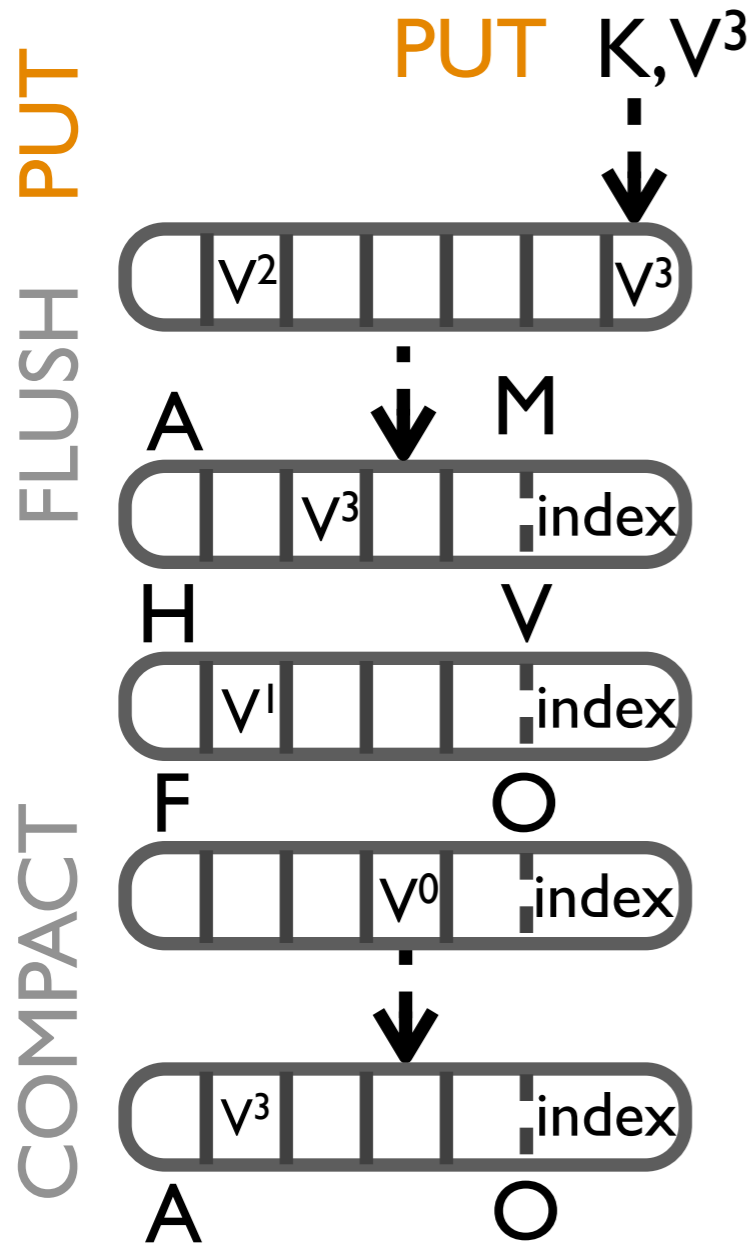Non-Linear without per-op

# Workload-dependent IO Amplification
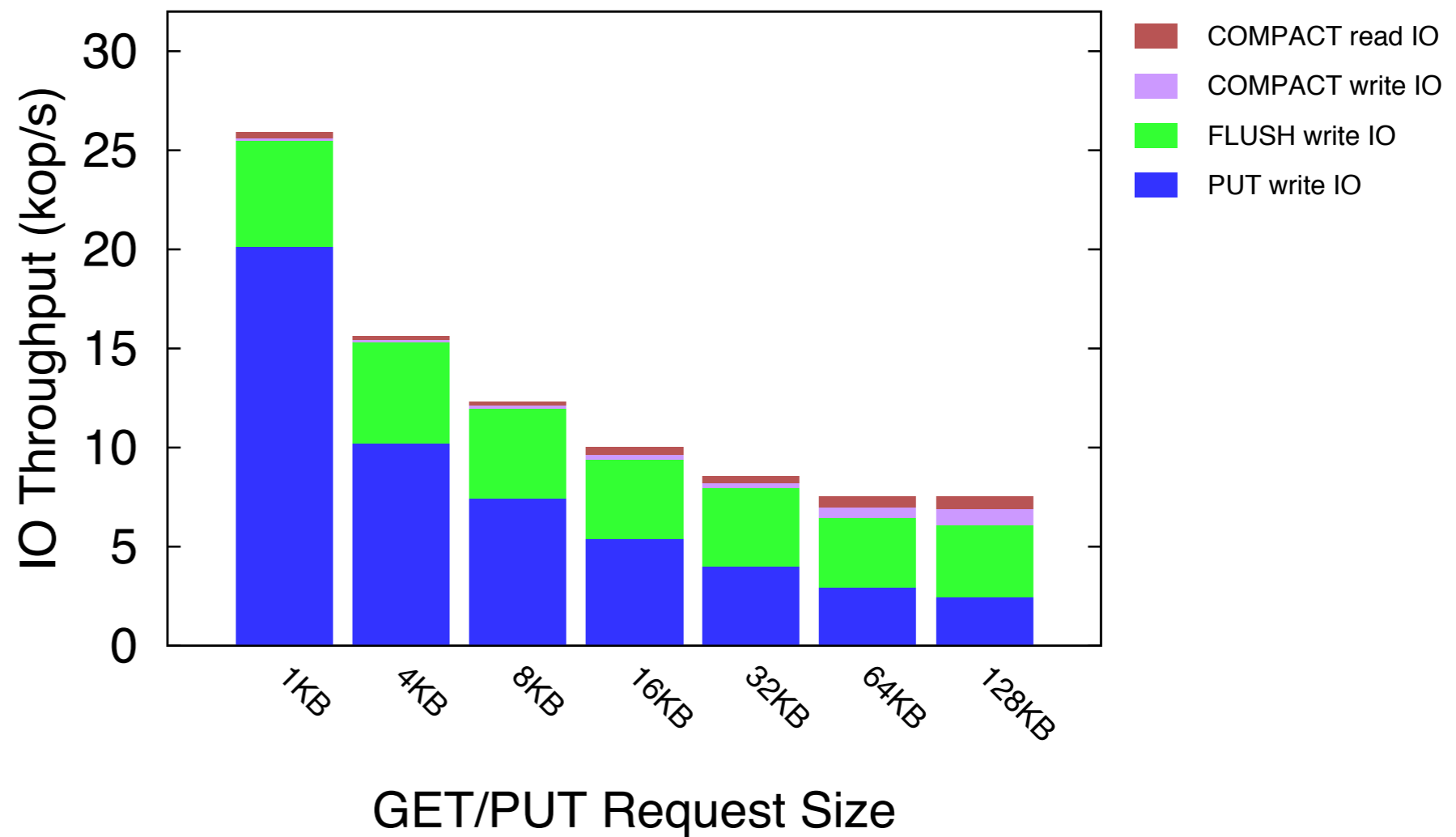


LevelDB (LSM-Tree)

# Workload-dependent IO Amplification



LevelDB (LSM-Tree)

# Workload-dependent IO Amplification

# Workload-dependent IO Amplification



GET/PUT Request Size

# Libra Tracks App-request IO Consumption to Determine IO Allocations

IO | Tenant A
500 IO/s

Libra
Provisioning
Policy

PUT
FLUSH

Tenant A
5 IO units

Track IO consumption

Compute app-request IO profiles

Provision IO allocations

5 GET 5

25 PUT 100

1

0.5 FLUSH 50

COMPACT

Per-GET 1 × GET 80

+

Per-PUT 6 × PUT 70 = IO 500

# Unpredictable IO Interference

✅ Die-level parallelism, low latency IOPs

❌ Shared-controller and bus contention

❌ Erase-before-write overhead

❌ FTL and read-modify-write garbage colleciton

4 read/4 write tenants
1:1 Pure Read/Pure Write

# Unpredictable IO Interference

# Libra Underestimates IO Capacity to Ensure Provisionable Throughput

Provisionable IO throughput = floor(workloads)
(18 Kop/s)

# Libra Underestimates IO Capacity to Ensure Provisionable Throughput

Provisionable IO throughput = floor(workloads)
(18 Kop/s)

# Libra Underestimates IO Capacity to Ensure Provisionable Throughput

Provisionable IO throughput = floor(workloads)
(18 Kop/s)

# Non-linear IO Performance



IO Bandwidth

Max BW

Bandwidth (MB/s): 0, 50, 100, 150, 200, 250

IOP Size (KB): 1, 2, 4, 8, 16, 32, 64, 128, 256

IOP Throughput

Max IOP/s

IOP (kop/s): 0, 5, 10, 15, 20, 25, 30, 35, 40

IOP Size (KB): 1, 2, 4, 8, 16, 32, 64, 128, 256

25

# Non-linear IO Performance



IO Bandwidth

IOP Throughput

Read Rand
Read Seq
Write Rand
Write Seq

# Libra Uses Virtual IOPs to Model IO Resources

$$\text{VOP}_{\text{CPB}}(\text{IOP-size}) = \frac{\text{Max-IOP}}{\text{Achieved-IOP}(\text{IOP-size}) \times \text{IOP-size}}$$

✅ Unifies IO cost into a single metric

✅ Captures non-linear IO performance

✅ Provides IO insulation

2 equal-allocation tenants
IO Insulation = 1/2 Max Read/Write

**Libra IO Cost Model**

Legend: Read IO cost (red), Write IO cost (blue)
Y-axis: Virtual IOP Cost (op/KB)
X-axis: IOP Size (KB)

# Libra Uses Virtual IOPs to Model IO Resources

$$\text{VOP}_{\text{CPB}}(\text{IOP-size}) = \frac{\text{Max-IOP}}{\text{Achieved-IOP}(\text{IOP-size}) \times \text{IOP-size}}$$



2 equal-allocation tenants
IO Insulation = 1/2 Max Read/Write

# Libra Design



Persistence Engine

Libra IO Scheduler

Libra Provisioning Policy

Physical Disk

Update tenant VOP allocations

Provision VOPs within provisionable limit

Charge tenant IOPs based on VOP cost

Track app-request VOP consumption

# Evaluation

- Does Libra's IO resource model achieve accurate resource allocations?

- Does Libra's IO threshold make an acceptable tradeoff of performance for predictability in a real storage stack?

- Can Libra ensure per-tenant app-request reservations while achieving high utilization?

# Libra Achieves Accurate IO Allocations



Read-Write IOP Throughput Ratio

Throughput Ratio = Actual / Expected (IO Insulation)

# Libra Achieves Accurate IO Allocations



Read-Write IOP Throughput Ratio

Throughput Ratio = Actual / Expected (IO Insulation)

# Libra Achieves Accurate IO Allocations



Read IO Cost Models

Write IO Cost Models

# Libra Achieves Accurate IO Allocations



Read IO Cost Models

Write IO Cost Models

# Libra Achieves Accurate IO Allocations

## Min-Max Ratio =
## Min Throughput Ratio / Max Throughput Ratio

# Libra Trades-off Nominal IO Throughput For Predictability



75:25 GET-PUT, Variance 4K  50:50 GET-PUT, Variance 4K  25:75 GET-PUT, Variance 4K

PUT Request Size (KB) / GET Request Size (KB) / VOP (kop/s)

# Libra Trades-off Nominal IO Throughput For Predictability

## Unprovisionable Throughput As a Percentage of Total Throughput

| Workload | Percentile | | | |
|----------|------|-------|-------|-------|
|          | 10th | 50th  | 80th  | All   |
| 99:1     | 1.6% | 30.5% | 40.5% | 45.8% |
| 25:75    | 1.4% | 14.9% | 25.0% | 34.7% |
| 1:99     | 0.7% | 12.2% | 19.5% | 28.1% |

< 10th percentile covered by SLA and higher-level policies

# Libra Achieves App-request Reservations

0.5x                                    1.5x

**Read Heavy**     **Mixed**     **Write Heavy**

Libra Normalized GET (1KB)          Libra Normalized PUT (1KB)

Fully provisioned allocations

No Profile Normalized GET (1KB)          No Profile Normalized PUT (1KB)

Work-conserving consumption of unprovisioned resources

Throughput (kreq/s)

Time (s)

# Libra Achieves App-request Reservations

Read Heavy   Mixed   Write Heavy



Libra Normalized GET (1KB)

Libra Normalized PUT (1KB)

No Profile Normalized GET (1KB)

No Profile Normalized PUT (1KB)

# Conclusion

- **Libra IO Scheduler**

  - Provisions *IO allocations* for *app-request reservations* w/ high utilization.
  - Supports arbitrary object distributions and workloads.

- **2 key mechanisms**

  - Track per-tenant app-request resource profiles.
  - Model IO resources with Virtual IOPs.

- **Evaluation**

  - Achieves accurate low-level IO allocations.
  - Provisions the majority of IO resources over a wide range of workloads
  - Satisfies app-request reservations w/ high utilization.