

# COSBench: Cloud Object Storage Benchmark

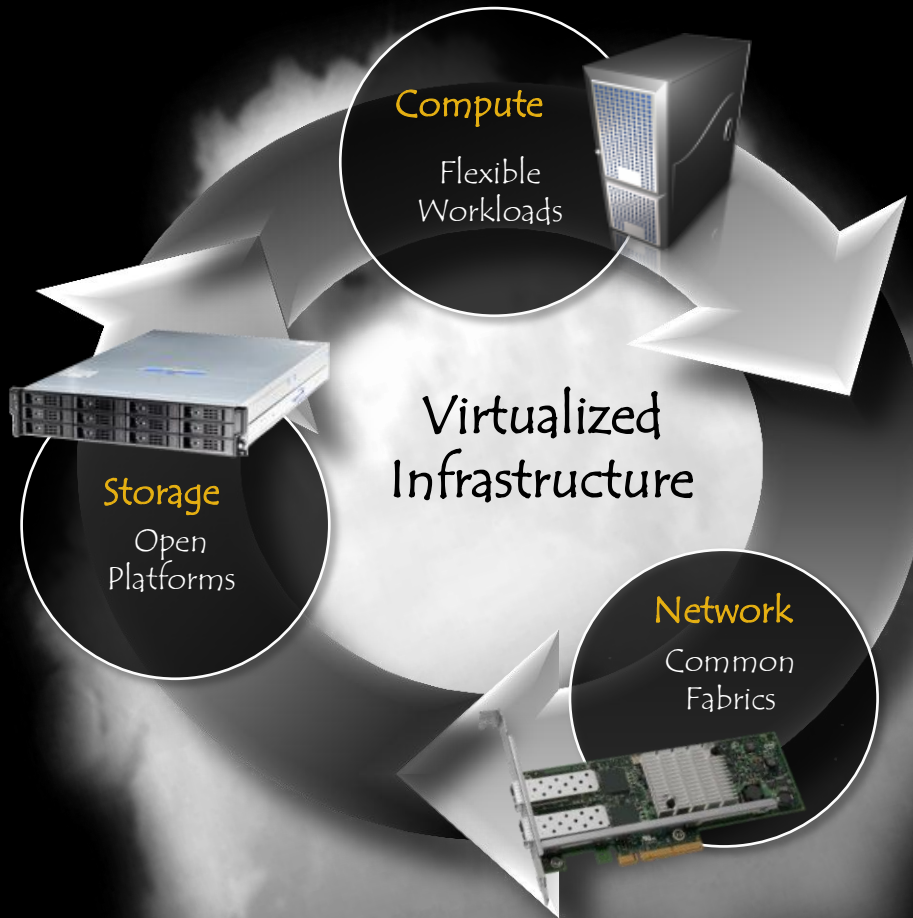
Qing Zheng

Reliable, Scalable, Intelligent Systems Group  
School of Software, Shanghai Jiao Tong University

<http://reins.se.sjtu.edu.cn>

Oct. 12nd, 2012

qzheng2010@hotmail.com



- Data centers are built upon 3 fundamental pillars:
  - Compute
  - Network
  - Storage
- To achieve high efficiency in performance & utilization
  - A balanced data center is essential!

# Storage Capacity Growth

REINS

## Structured data (23.6% ↑)

- Traditional enterprise database

## Replicated data (24.2% ↑)

- Backups
- Data warehouses

## Unstructured data (54.8% ↑)

- Archives

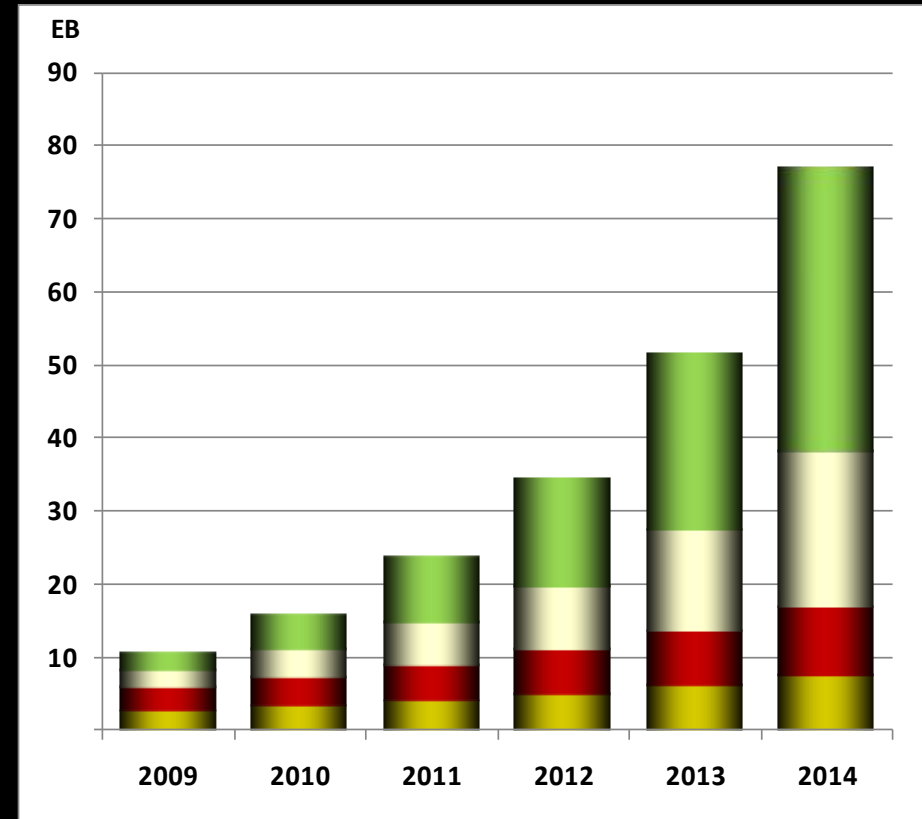
## Content Depots (75.6% ↑)

- Web
- Email
- Document sharing
- Social network content (pictures/videos)

## 2012 deployment estimation

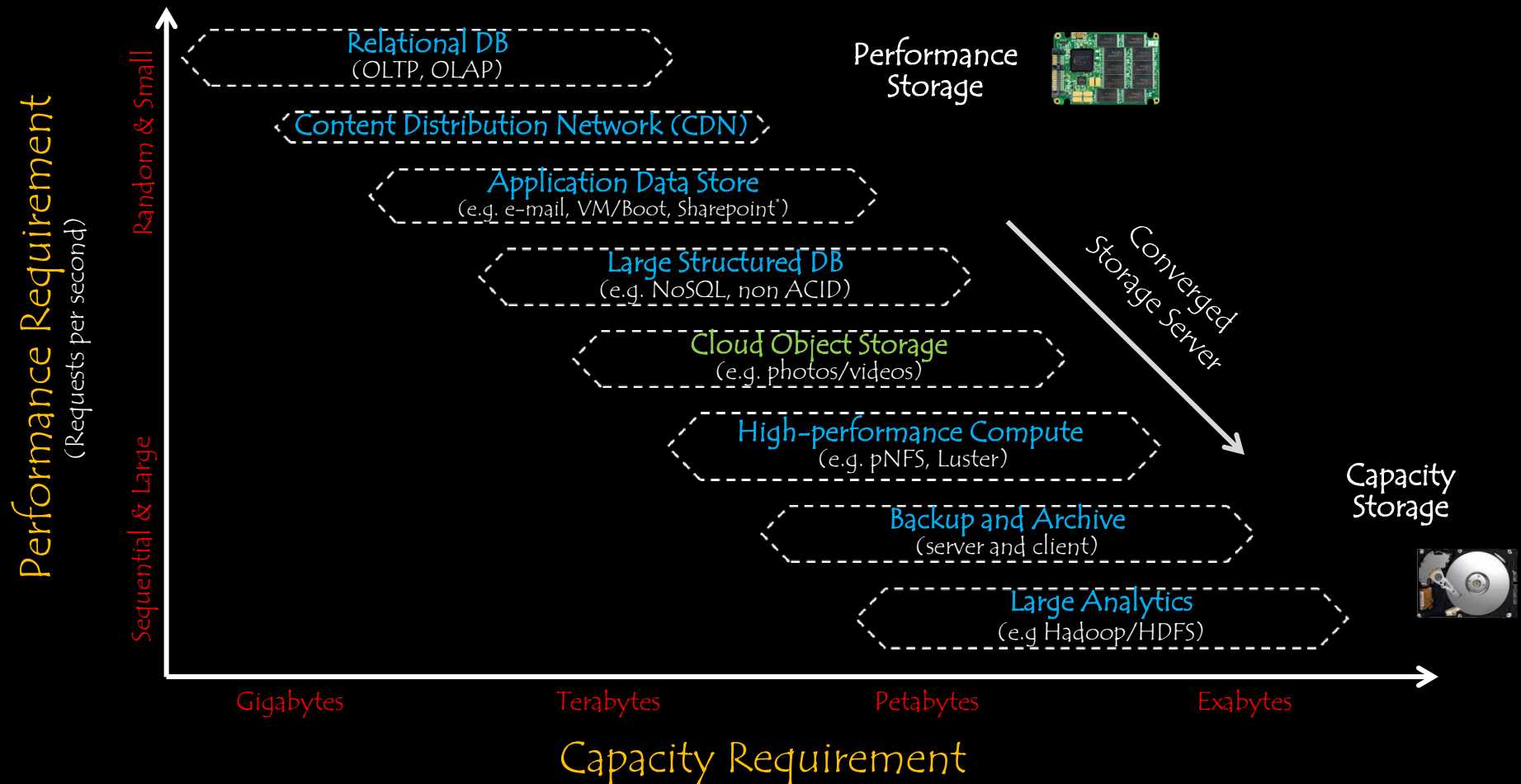
~ 7.6 million drives

~ 500,000 storage systems



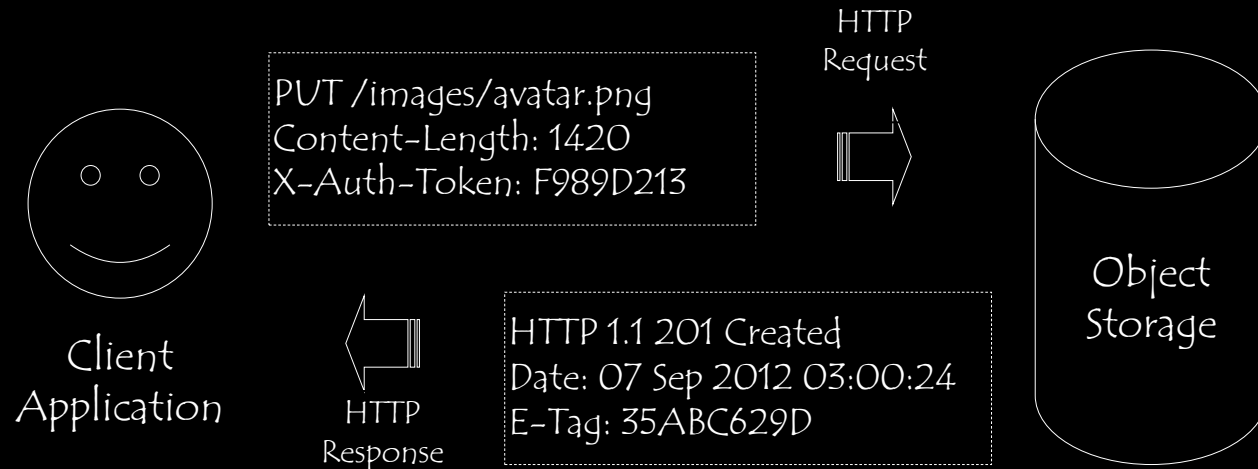
# Solutions dictated by Usage Models

# REINS

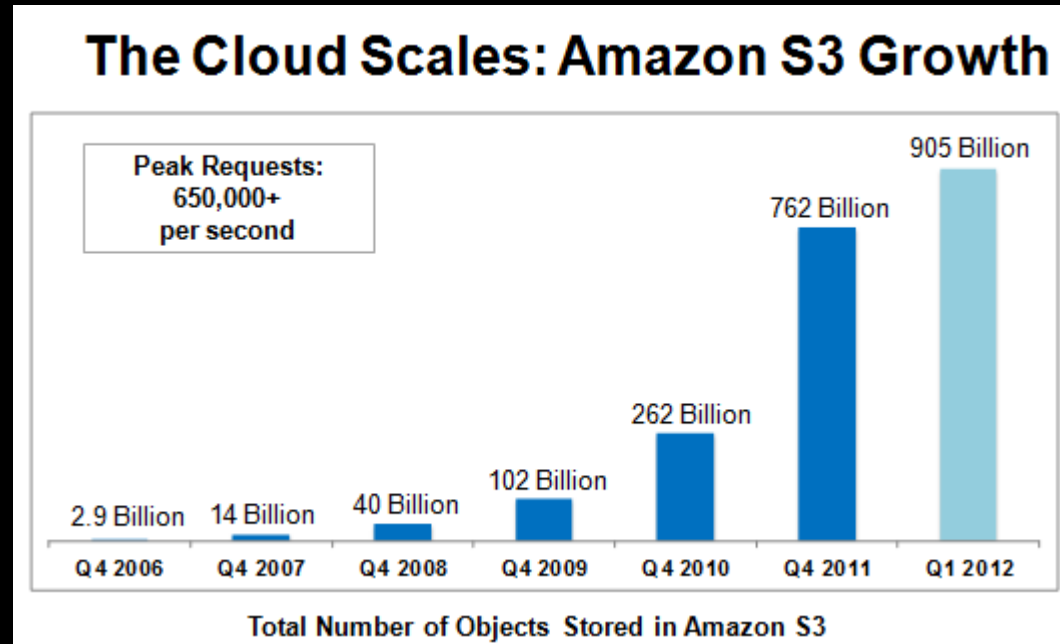


Different storage usage models create different ecosystems!

- Storage Interface



HEAD / GET / PUT / POST / DELETE / COPY



Object storage have been increasingly recognized as a right destination for data outsourcing

- The proliferation of existing offerings
  - Amazon S3, Rackspace Cloud Files, Google Cloud Storage, HP Cloud Object Storage, Windows Azure, EMC ATMOS, Openstack Swift, Ceph, ...
  - coupled with a lack of workload modeling object storage apps
  - makes it difficult for one to choose the right infrastructure!
- Tuning systems to their optimal performance is also nontrivial
  - resulted from the complexity in the designs of various available solutions

We need a benchmark tool  
dedicated for object storage system!

- People can use this tool to evaluate & compare different
  - hardware and software stacks
  - and obtain a better understanding of these offerings

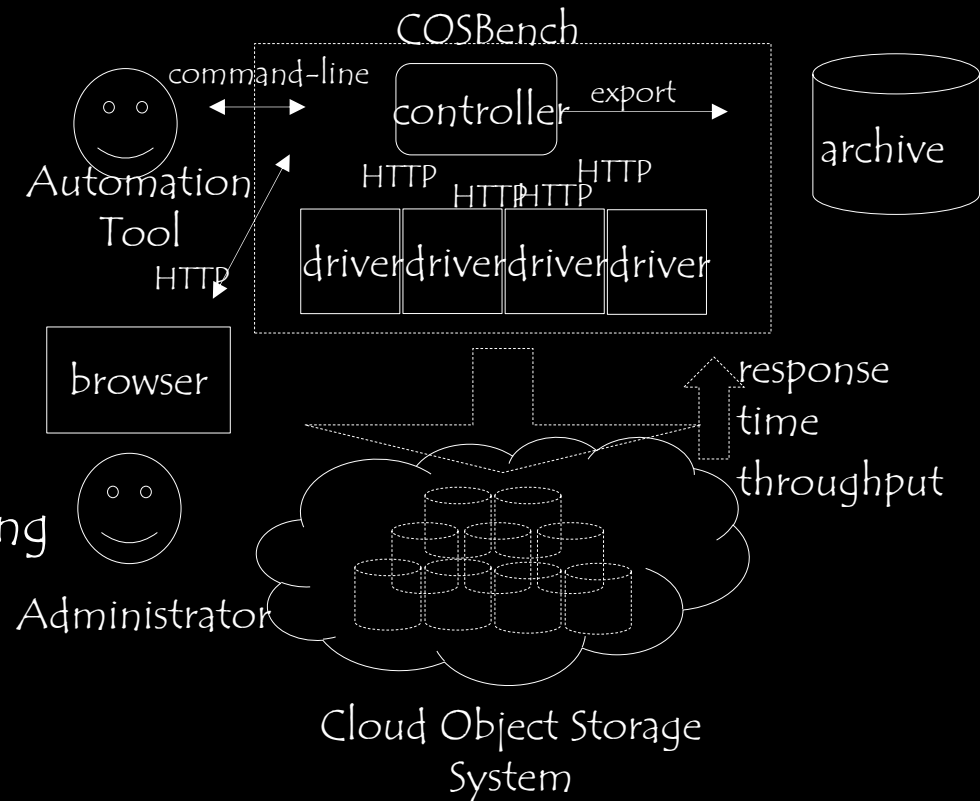
Throughput / Latency / Scalability / Flexibility

- People can also use this tool to characterize their systems
  - and get insights to guide system designs, tuning, & optimization

Architecture Decisions / Algorithms / Hardware Impacts / Innovation



- Config.xml:
  - define workload with flexibility.
- Controller:
  - Control all drivers
  - Collect and aggregate stats.
- Driver:
  - generate workload according to config parameters.
  - can run tests without a supervising controller.
- Web Console:
  - System facade
  - Browse real-time stats
  - HTTP based Communication (RESTful style)



**COSBENCH - CONTROLLER WEB CONSOLE** GA Release version: 2.0.0.GA

Controller Overview

Name: *not configured* URL: *not configured*

Driver	Name	URL	Link
1	driver1	http://127.0.0.1:18088/driver	<a href="#">view details</a>
2	driver2	http://127.0.0.1:18088/driver	<a href="#">view details</a>

There are 2 drivers attached to the controller.

**Active Workloads**

Id	Name	Submitted-At	State	Link
w6	demo	Aug 3, 2012 2:56:48 PM	processing	<a href="#">view details</a>
w7	demo	Aug 3, 2012 2:56:52 PM	queuing	<a href="#">view details</a>

There are currently 2 active workloads.

[submit new workloads](#)

**History Workloads**

[view performance matrix](#)

Id	Name	Duration	Op-Info	State	Link
w4	demo	Aug 3, 2012 2:52:51 PM - 2:53:37 PM	prepare, read	finished	<a href="#">view details</a>
w5	demo	Aug 3, 2012 2:53:37 PM - 2:54:23 PM	prepare, read	finished	<a href="#">view details</a>

Software Solution Group - System Software Division - System Optimization Technology Center

Controller

**COSBENCH - DRIVER WEB CONSOLE** GA Release version: 2.0.0.GA

Driver Overview

Name: *not configured* URL: *not configured*

**Active Missions**

Id	Name	Submitted-At	State	Link
M9EB45CC751	main	Aug 3, 2012 2:56:56 PM	launched	<a href="#">view details</a>
M7EB45CC753	main	Aug 3, 2012 2:56:56 PM	launched	<a href="#">view details</a>

There are currently 2 active missions.

[submit new missions](#)

**History Missions**

Id	Name	Submitted-At	State	Link
M1EB42119A3	prepare	Aug 3, 2012 2:52:51 PM	accomplished	<a href="#">view details</a>
MFEb42119AE	prepare	Aug 3, 2012 2:52:51 PM	accomplished	<a href="#">view details</a>
MCEB422EA9F	main	Aug 3, 2012 2:52:59 PM	accomplished	<a href="#">view details</a>
M4EB422EAA5	main	Aug 3, 2012 2:52:59 PM	accomplished	<a href="#">view details</a>
MCEB42C5AA3	prepare	Aug 3, 2012 2:53:37 PM	accomplished	<a href="#">view details</a>
MDEB42C5AFC	prepare	Aug 3, 2012 2:53:37 PM	accomplished	<a href="#">view details</a>
MDEB42E20E1	main	Aug 3, 2012 2:53:45 PM	accomplished	<a href="#">view details</a>
M8EB42E20E6	main	Aug 3, 2012 2:53:45 PM	accomplished	<a href="#">view details</a>
M6EB45AFC8B	prepare	Aug 3, 2012 2:56:48 PM	accomplished	<a href="#">view details</a>

Software Solution Group - System Software Division - System Optimization Technology Center

Driver

DEMO

```
- <workflow >
- <workstage name="init">
  <work type="init" workers="8" config="containers=r(1,32)" />
</workstage>
- <workstage name="prepare">
  <work type="prepare" workers="8" config="containers=r(1,32);objects=r(1,50);sizes=c(64)KB" />
</workstage>
- <workstage name="main">
  <work name="main" workers="8" rampup="100" runtime="300">
    <operation type="read" ratio="80" config="containers=u(1,32);objects=u(1,50)" />
    <operation type="write" ratio="20" config="containers=u(1,32);objects=u(51,100);sizes=c(64)KB" />
  </work>
</workstage>
- <workstage name="cleanup">
  <work type="cleanup" workers="8" config="containers=r(1,32);objects=r(1,50)" />
</workstage>
- <workstage name="dispose">
  <work type="dispose" workers="8" config="containers=r(1,32)" />
</workstage>
</workflow >
</workload >
```

Flexible Load Control

Path/Size Distribution

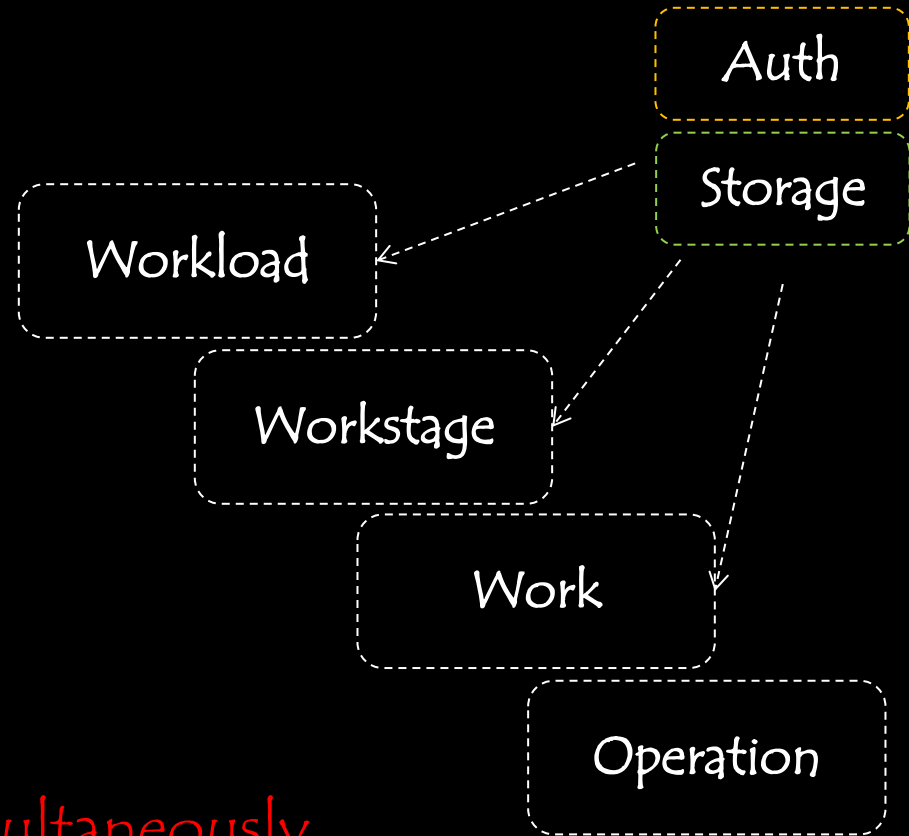
Read/Write Patterns

Extensible Parameters

Workflow Model

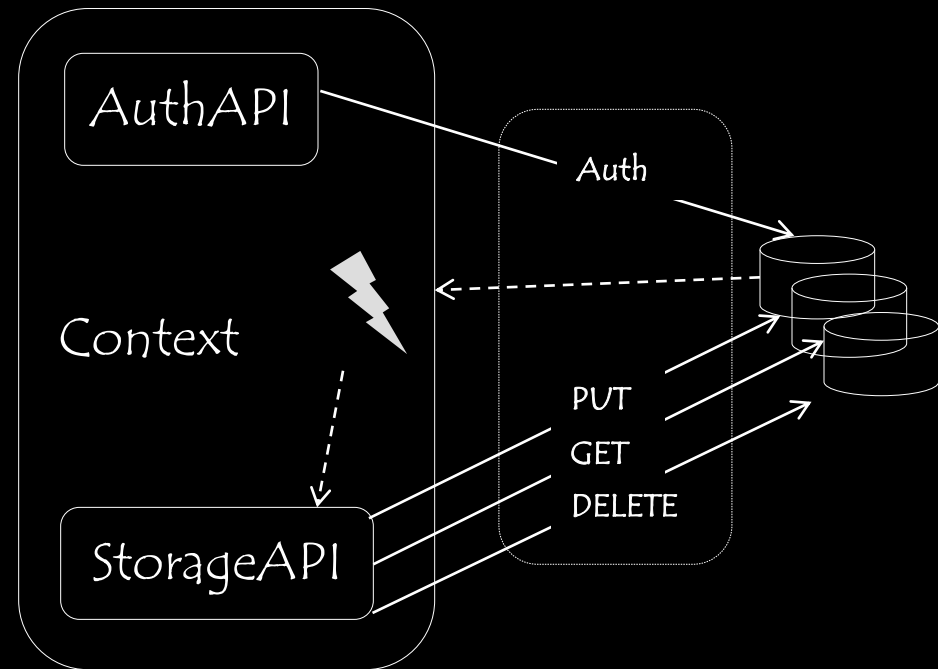
Flexible configuration which gives birth to diverse usage patterns

- Mixed operations (GET/PUT...)
- Mixed object sizes
- Multiple stages
- Auth/storage association
- Load control
- Extensible parameters

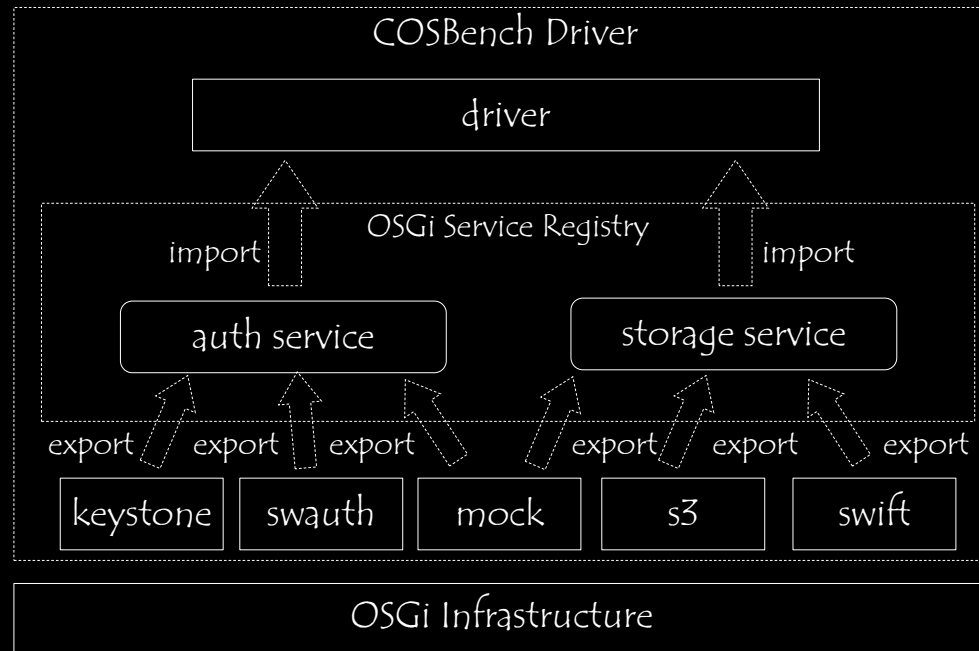


*Stress multiple systems simultaneously*

- Separate auth and storage API, so
  - One auth → multiple storages
  - One storage → multiple auths



Extensible API which support various storage systems



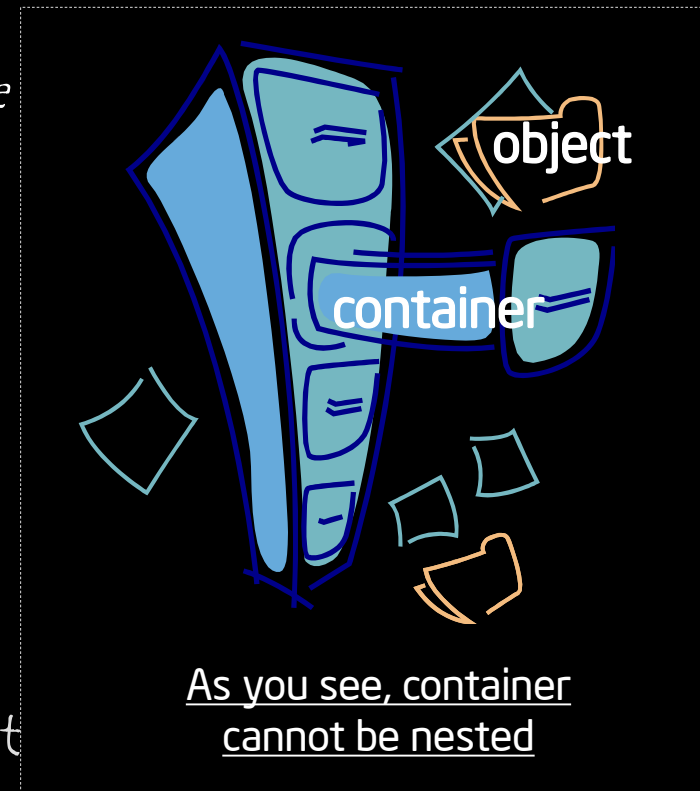
New adaptors can be separately developed, individually configured, and dynamically plugged into the env. without the knowledge of the core system

- Throughput (Operations/s): the operations completed in one second
- Response Time (in ms): the duration between operation initiation and completion.
- Bandwidth (KB/s): the total data in KiB transferred in one second
- Success Ratio (%): the ratio of successful operations

- Object Storage for Openstack



- As an object storage system, Swift:
  1. allows users to create containers and to store data objects in these containers
    - *objects are identified by their paths and have metadata associated with them*
  2. can be accessed via "RESTful" interface
    - *including "GET"/"PUT"/"DELETE"*
  3. can be well built upon commodity storage devices and is highly scalable
    - *achieving cost effectiveness*
  4. is redundant and is eventually consistent
    - *suitable for long-term storage*

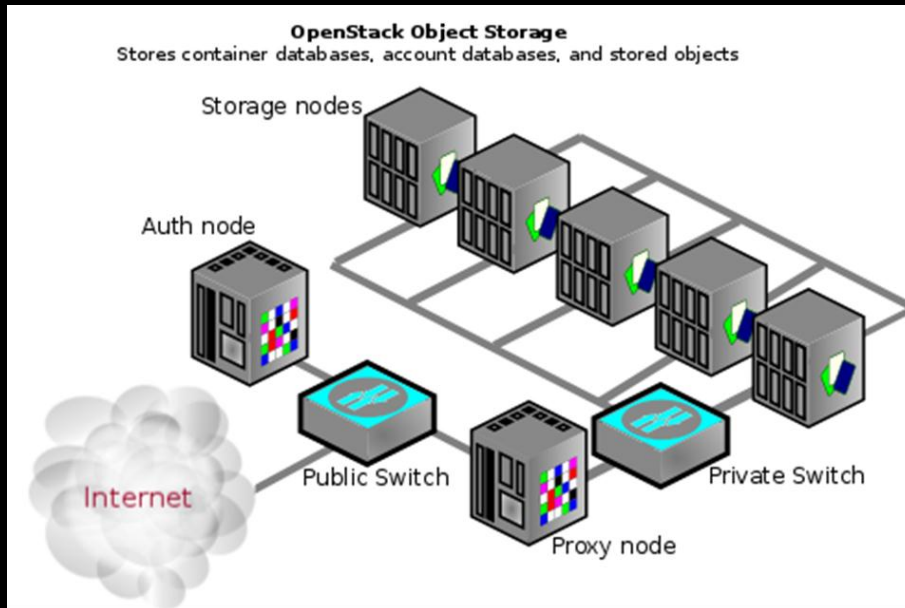


## • Proxy Node

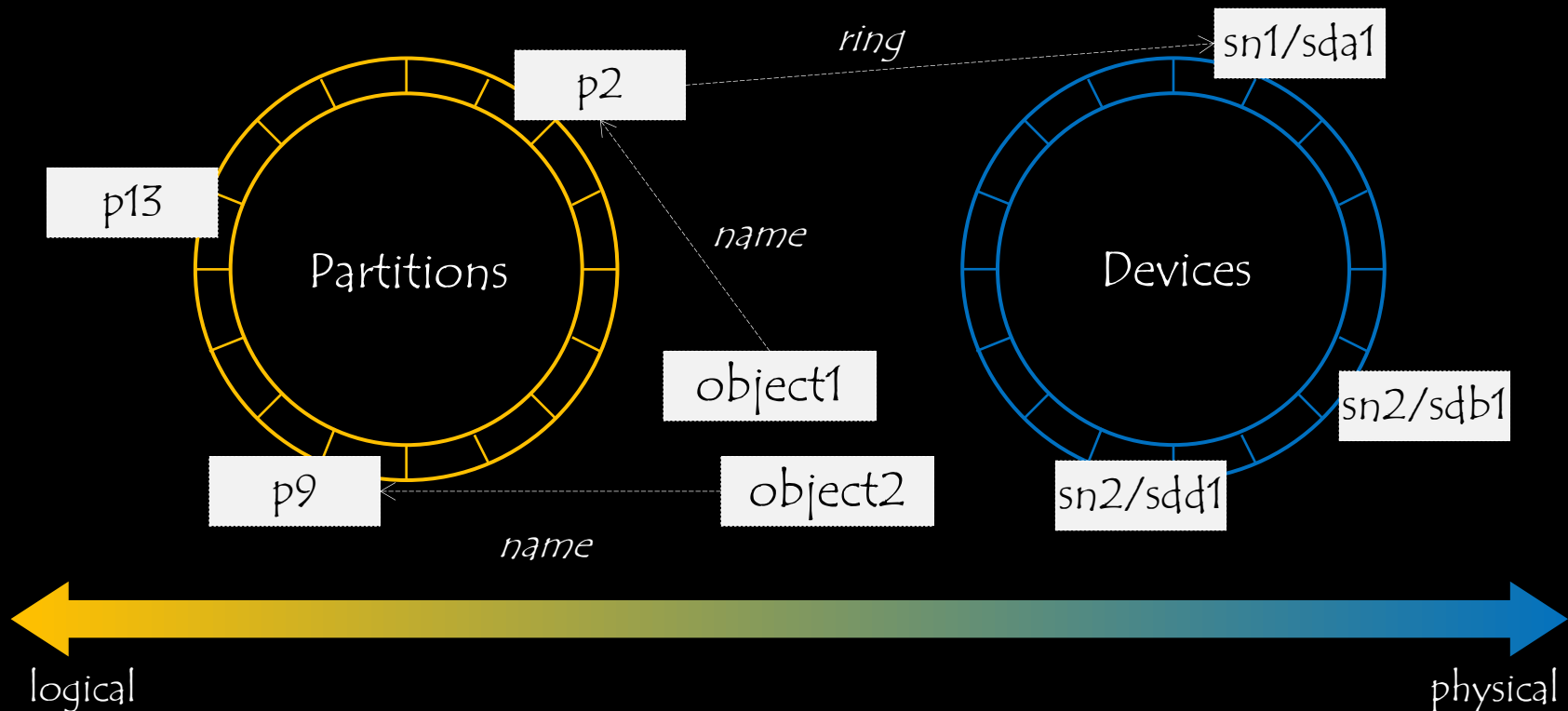
- proxy-server
  - cluster gateway
- swauth-server
  - authentication & authorization

## • Storage Node

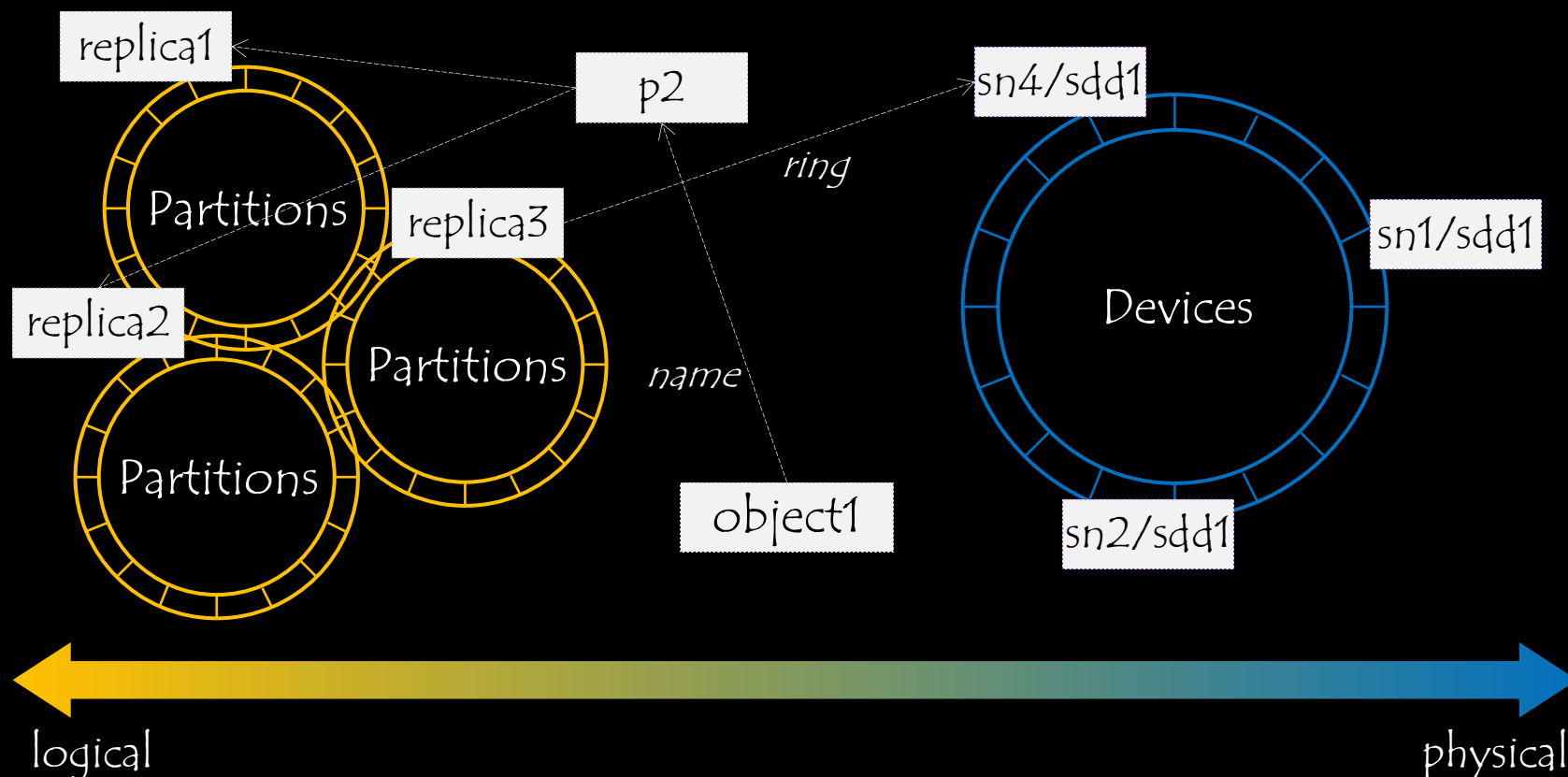
- account/container/object - server
  - listing containers
  - listing objects
  - saving/retrieving/removing objects
- account/container/object - replicator
  - pushing local replicas to other storage nodes should replicas of those nodes are missing
- account/container/object - auditor
  - quarantining local corrupted data entities
- container/object - updater
  - updating metadata asynchronously



- A ring maintains its mapping using
  - physical perspective: devices
  - logical perspective: partitions

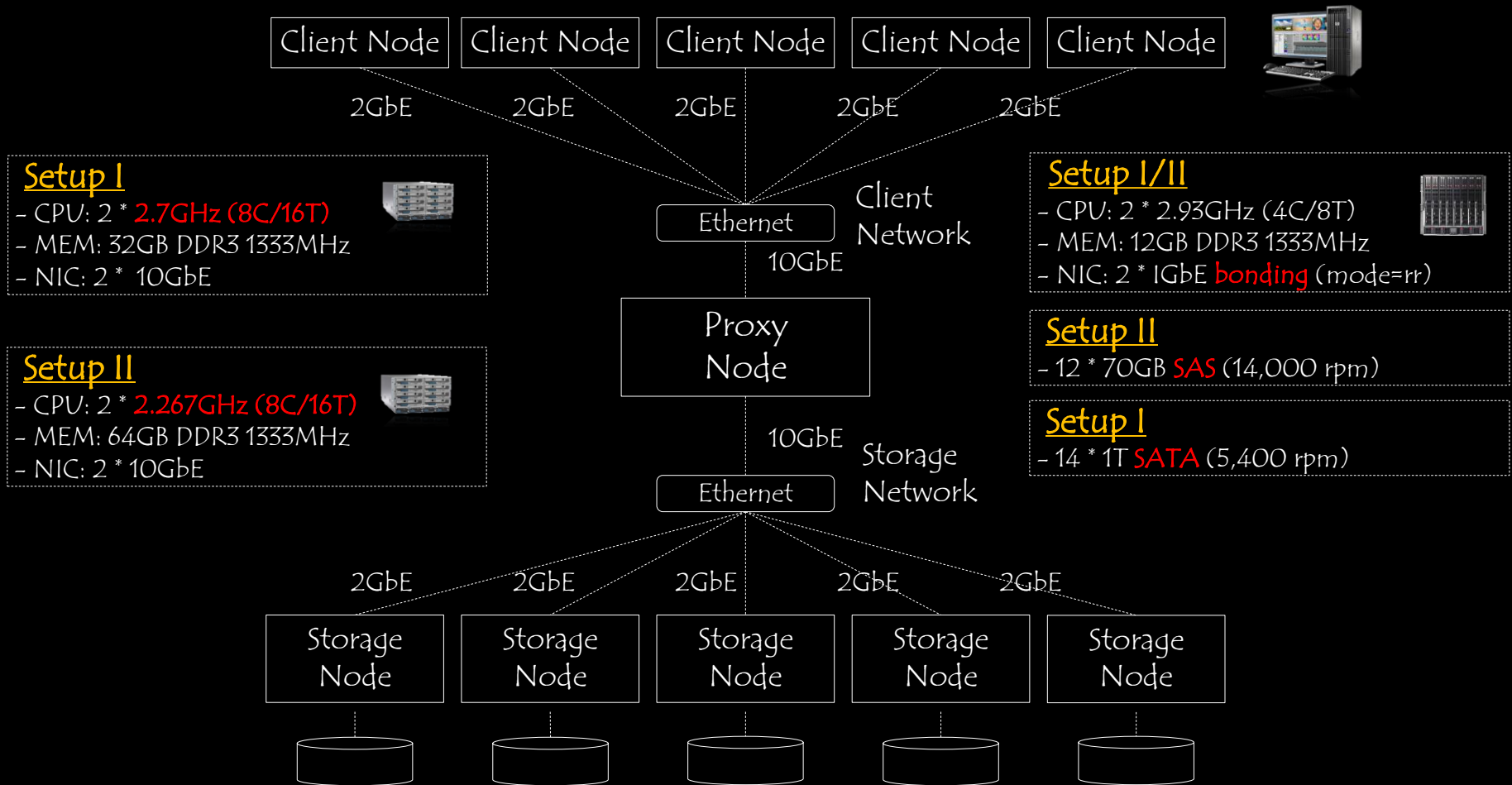


- A ring maintains its mapping using
  - physical perspective: **devices**
  - logical perspective: **partitions** → **replicas**



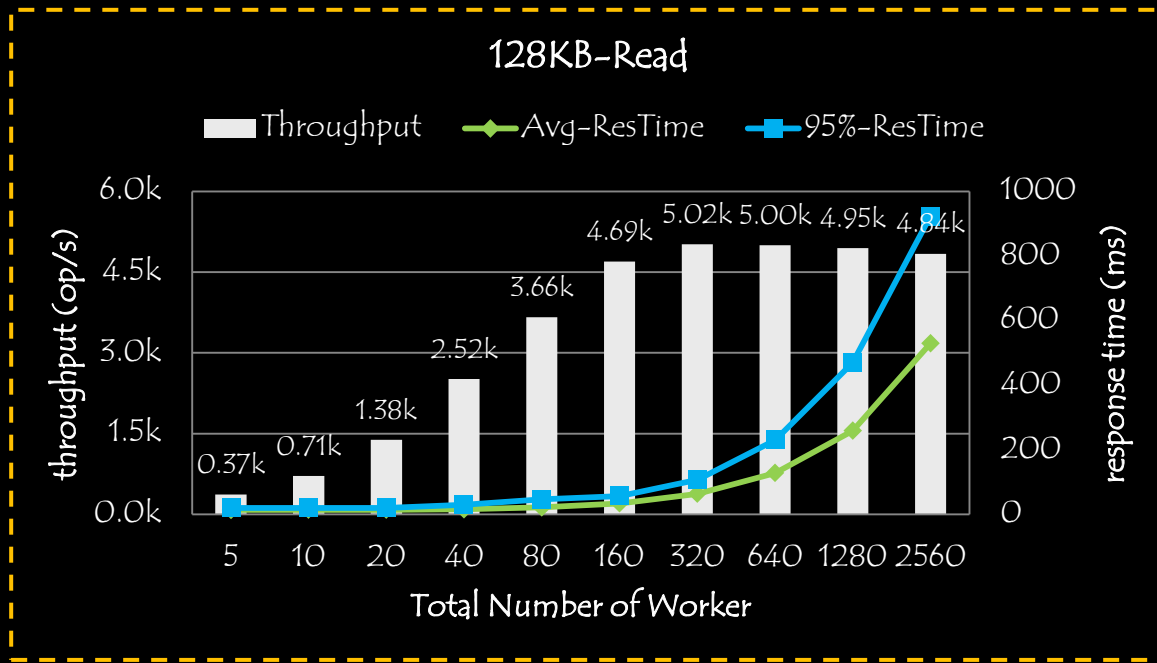
- In ring, each device is associated with metadata
  - describing device's weight, # replica wanted and # replica assigned
- A ring is built as
  - replicas from each partitions are in turn assigned to
  - the device currently enjoying the highest value of
    - “# replica wanted - # replica assigned”
  - Note that best attempts will be made to assign replicas from a sample partition to different zones or at least different nodes if possible

Setup-I had higher CPU power, Setup-II had faster disks



# A: 128KB-Read

- SLA: 200ms + 128KB/1MBps = 325ms



Workers	95%-ResTime ms	Throughput op/s
5	20.00	369.49
10	20.00	711.24
20	20.00	1383.30
40	30.00	2517.94
80	46.67	3662.71
160	56.67	4693.97
<b>320</b>	<b>106.67</b>	<b>5019.85</b>
640	230.00	4998.13
1280	470.00	4947.15
2560	923.33	4840.19

The **bottleneck** was identified to be the **proxy's CPU**

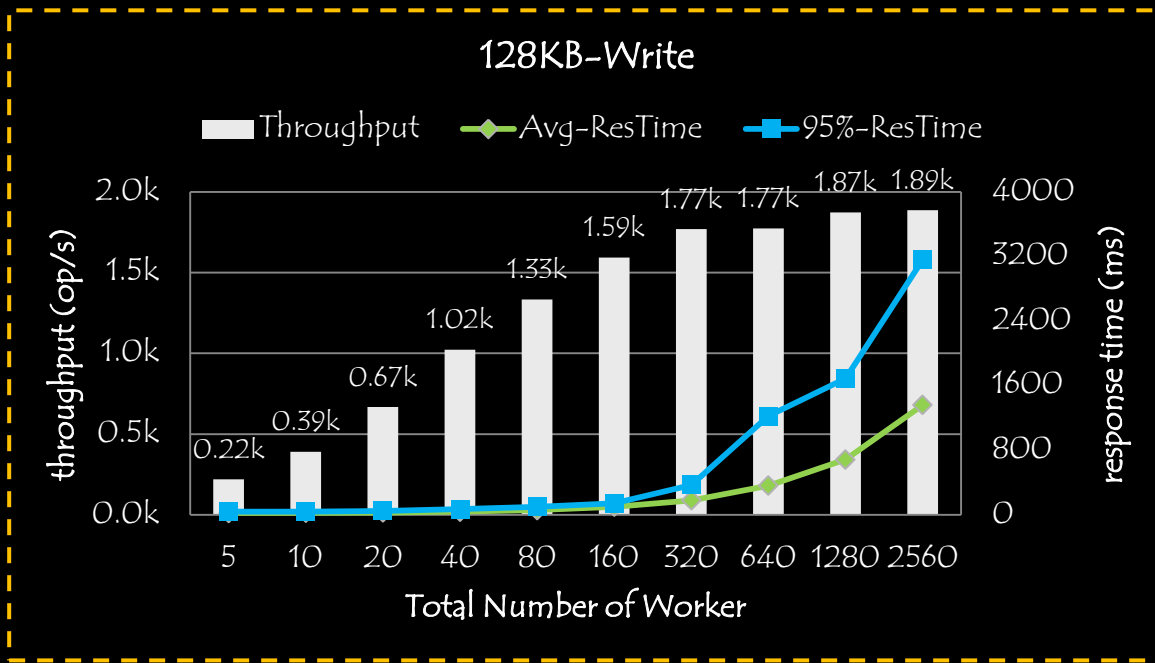
-- The CPU utilization at that node was **93%**!

-- The peak throughput for setup-I was **5576** op/s (640 workers)

↑ CPU could ↑ throughput

# B: 128KB-Write

- SLA: 200ms + 128KB/1MBps = 325ms



Workers	95%-ResTime ms	Throughput op/s
5	40.00	219.73
10	40.00	391.14
20	50.00	668.19
40	70.00	1022.07
80	100.00	1333.34
<b>160</b>	<b>143.33</b>	<b>1594.12</b>
320	370.00	1769.55
640	1223.33	1773.12
1280	1690.00	1871.58
2560	3160.00	1886.81

The **Disks** at storage nodes had significant impact on overall throughput

-- The peak throughput for setup-I was only **155** op/s (20 clients)

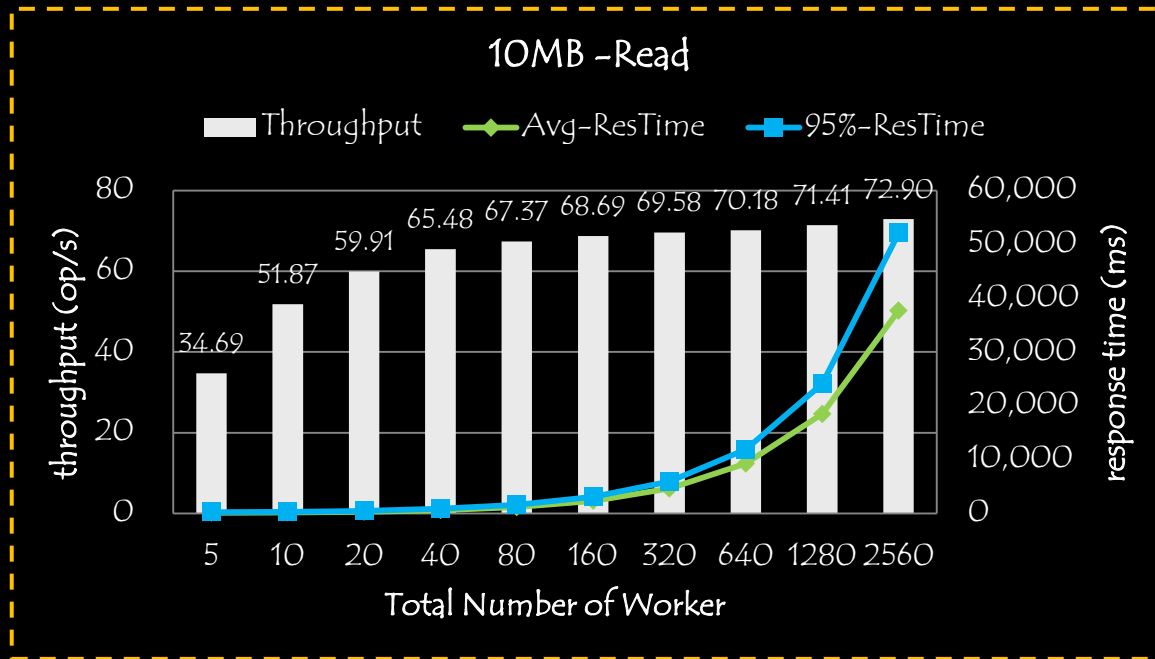
-- HDD → SSD in setup-I would ↑ throughput to **1621** op/s (320 clients)

storage disks ↔ throughput



# C: 10MB-Read

- SLA:  $200\text{ms} + 10\text{MB}/1\text{MBps} = \underline{1200\text{ms}}$



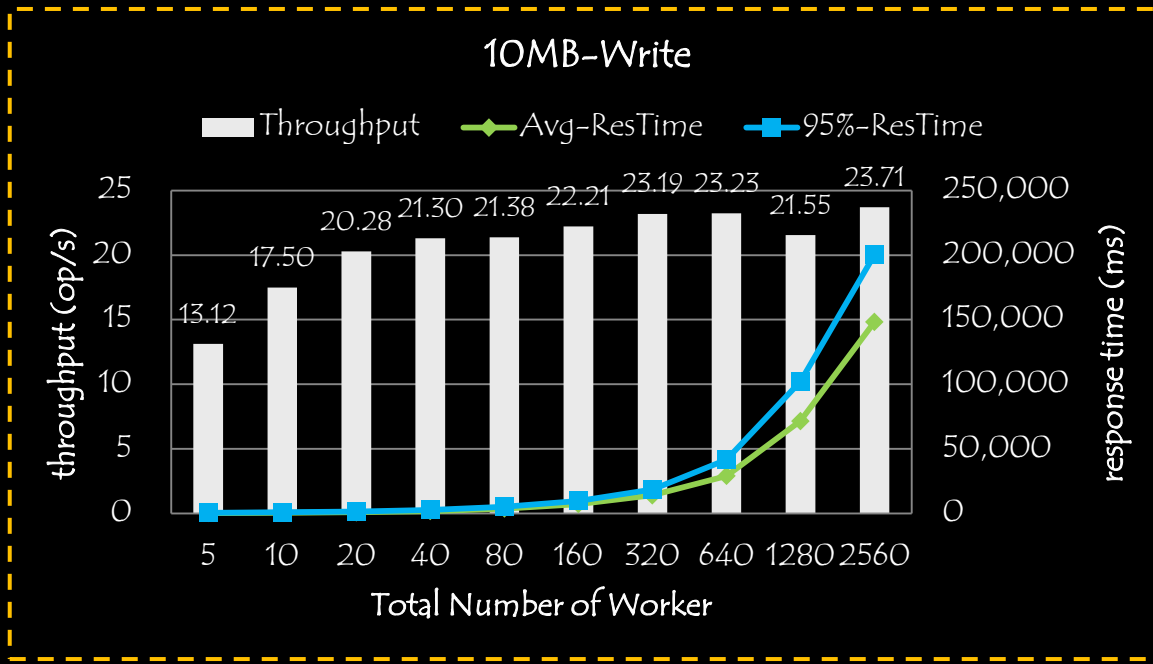
Workers	95%-ResTime ms	Throughput op/s
5	270.00	34.69
10	320.00	51.87
20	480.00	59.91
<b>40</b>	<b>900.00</b>	<b>65.48</b>
80	1636.67	67.37
160	3093.33	68.69
320	5950.00	69.58
640	11906.67	70.18
1280	24090.00	71.41
2560	52090.00	72.90

The **bottleneck** was identified to be the **clients' NICs**

- in setup-II, adding 5 more clients  $\uparrow$  throughput to **103** op/s (80 clients)
- in setup-I, using vClient + SRIOV could achieve similar improvements

# D: 10MB-Write

- SLA:  $200\text{ms} + 10\text{MB}/1\text{MBps} = \underline{1200\text{ms}}$



Workers	95%-ResTime ms	Throughput op/s
5	536.67	13.12
<b>10</b>	<b>936.67</b>	<b>17.50</b>
20	1596.67	20.28
40	2786.67	21.30
80	5133.33	21.38
160	9800.00	22.21
320	18623.33	23.19
640	41576.67	23.23
1280	102090.00	21.55
2560	200306.67	23.71

The **bottleneck** might be the **storage nodes' NICs**

-- in setup-1, the peak throughput was 15.74 op/s (10 clients)

-- in both setups, the write performance was 1/3 of the read performance

Thank you!