



LeoFS

Scaling and High Performance Storage System

Yosuke Hara - @yosukehara

A Researcher of R.I.T. and Tech Lead LeoFS

with Hiroki Matsue, LeoFS Support and Rakuten Software Engineer

LeoFS is an Unstructured Object Storage for the Web and a highly available, distributed, eventually consistent storage system.



LeoFS
The Lion of Storage Systems

**LeoFS was published as OSS
on July of 2012**

leo-project.net/leofs

Overview

Brief Benchmark Report

Multi Data Center Replication

LeoFS Administration at Rakuten

Future Plans

“NFS” Support and more

Overview





LeoFS

The Lion of Storage Systems

HIGH Availability

LeoFS Non Stop

3 Vs in 3 HIGHS

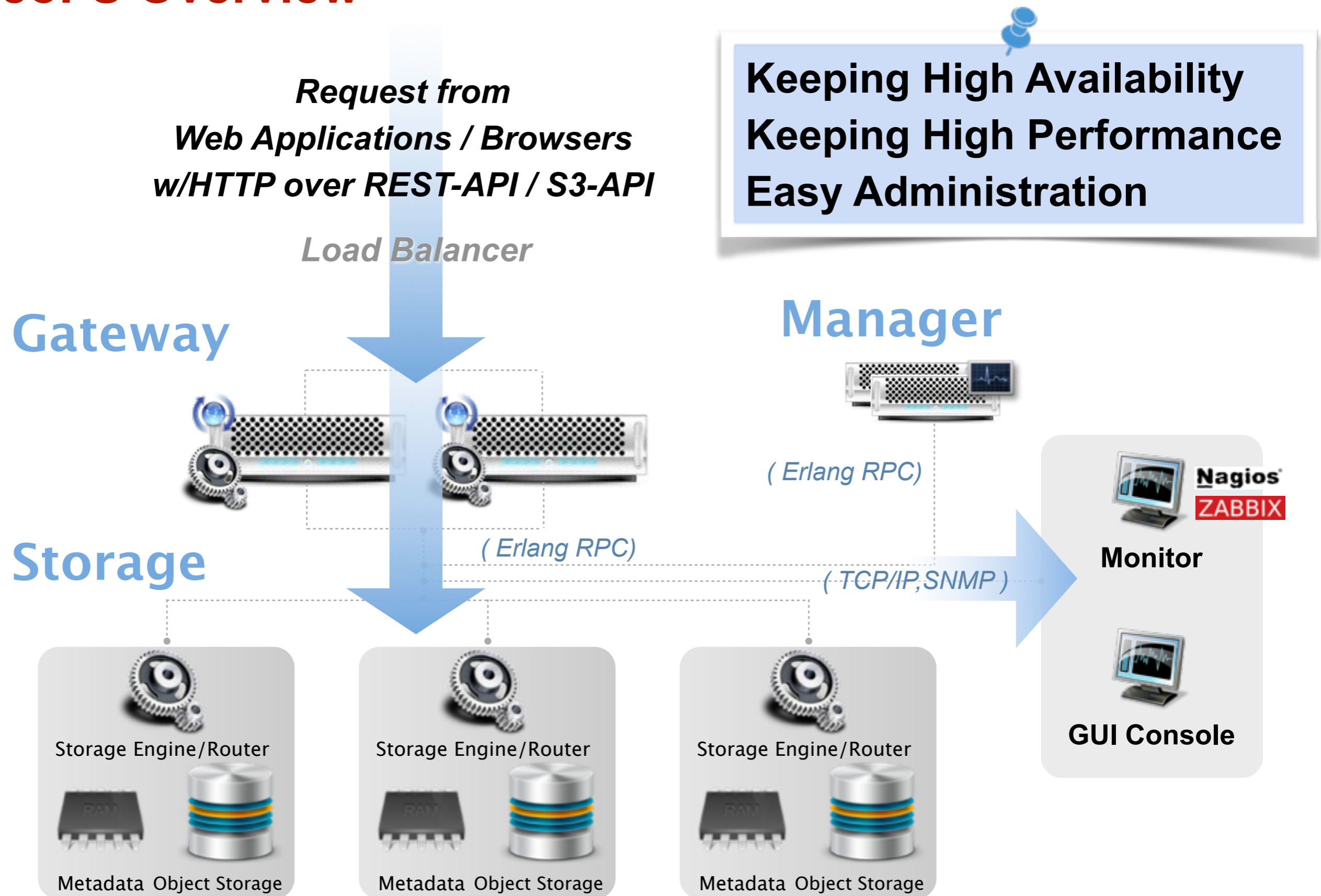
*Velocity: Low Latency
Minimum Resources*

*Volume: Petabyte / Exabyte
Variety: Photo, Movie, Unstructured-data*

**HIGH Cost
Performance Ratio**

**HIGH
Scalability**

LeoFS Overview



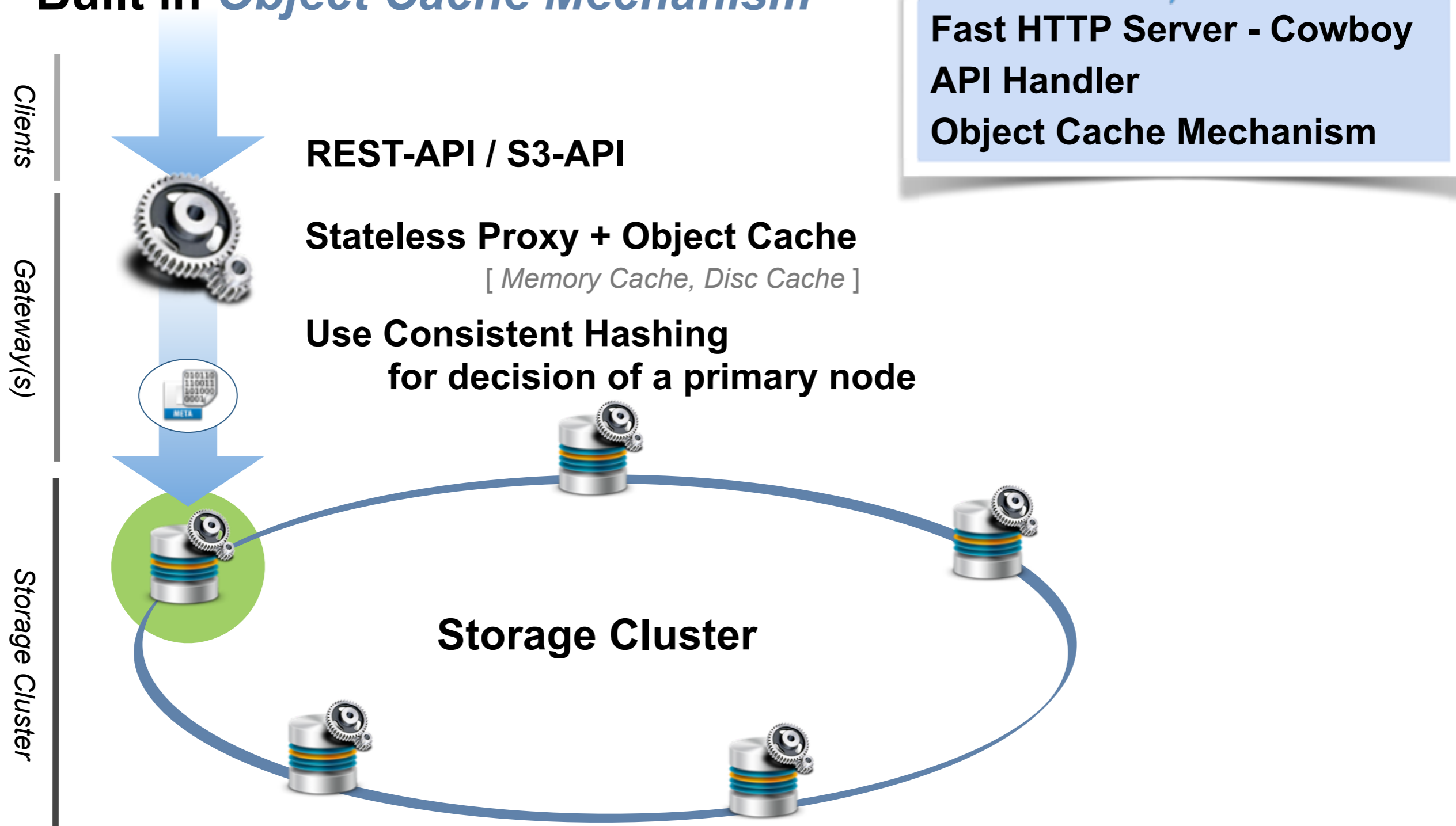
Keeping High Availability
Keeping High Performance
Easy Administration

Gateway

LeoFS Overview - Gateway

HTTP Request and Response

Built in *Object Cache Mechanism*

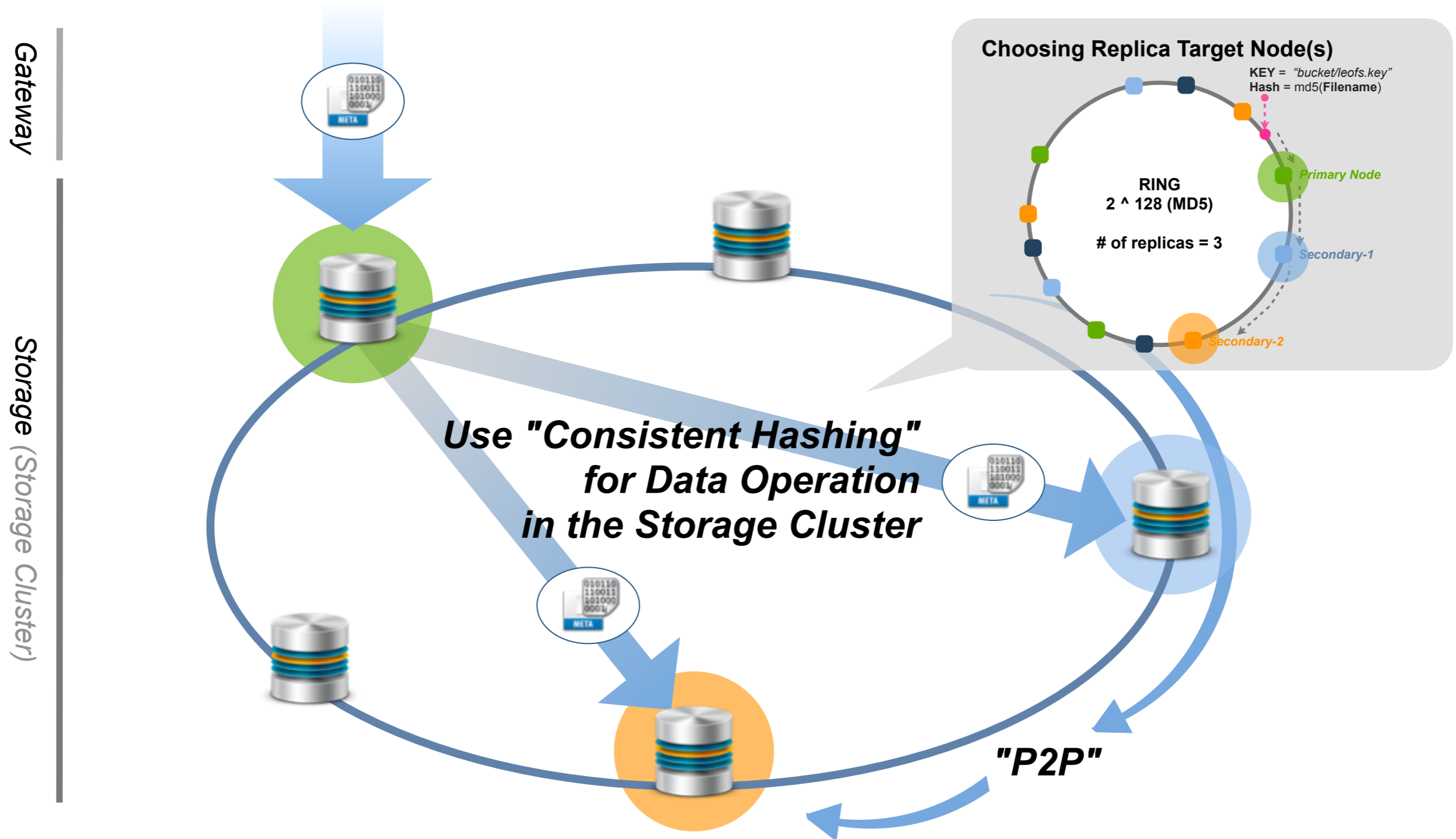


Storage

LeoFS Overview - Storage

WRITE: Auto Replication

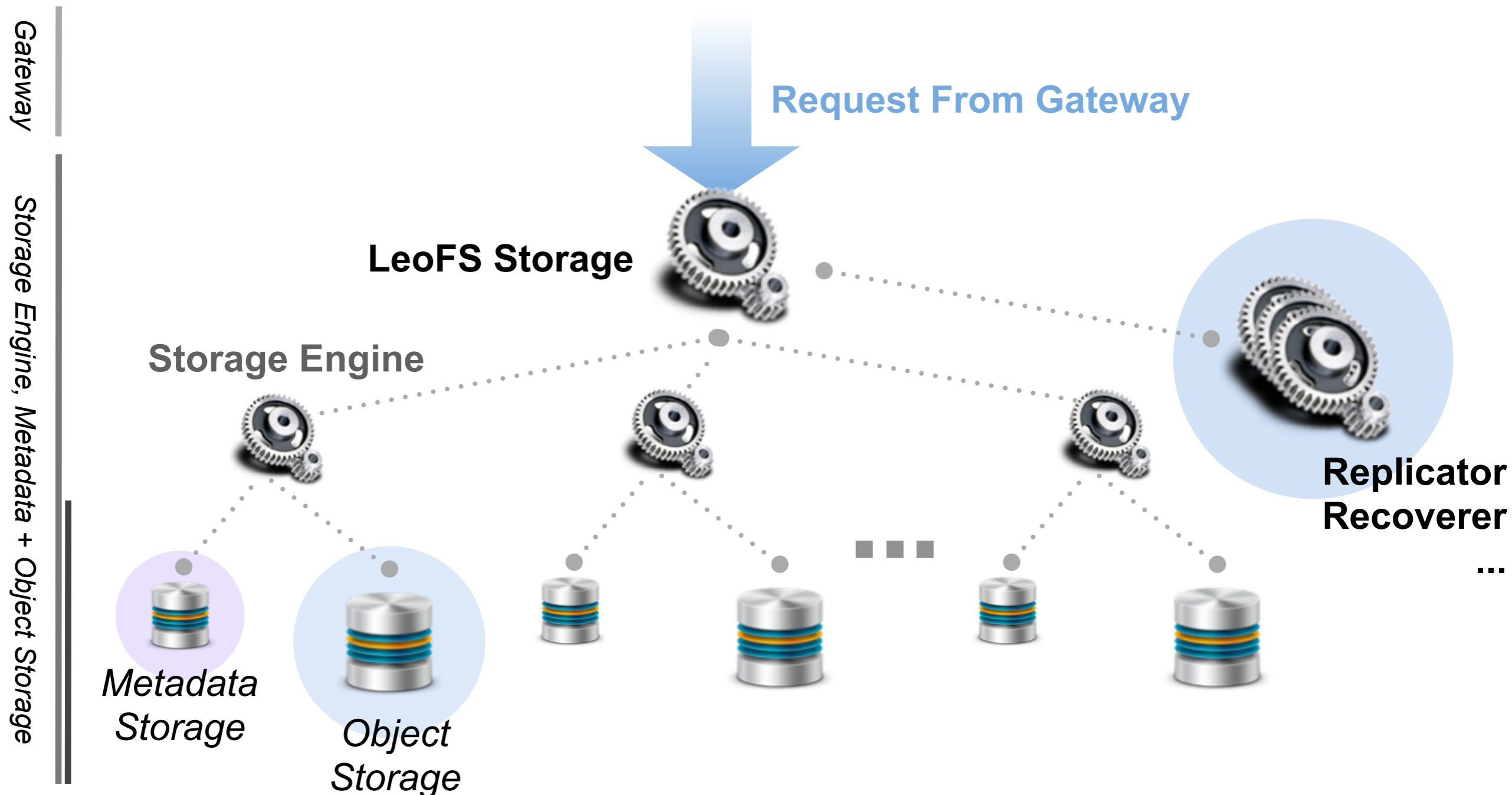
READ : Auto Repair of an Inconsistent Object with Async



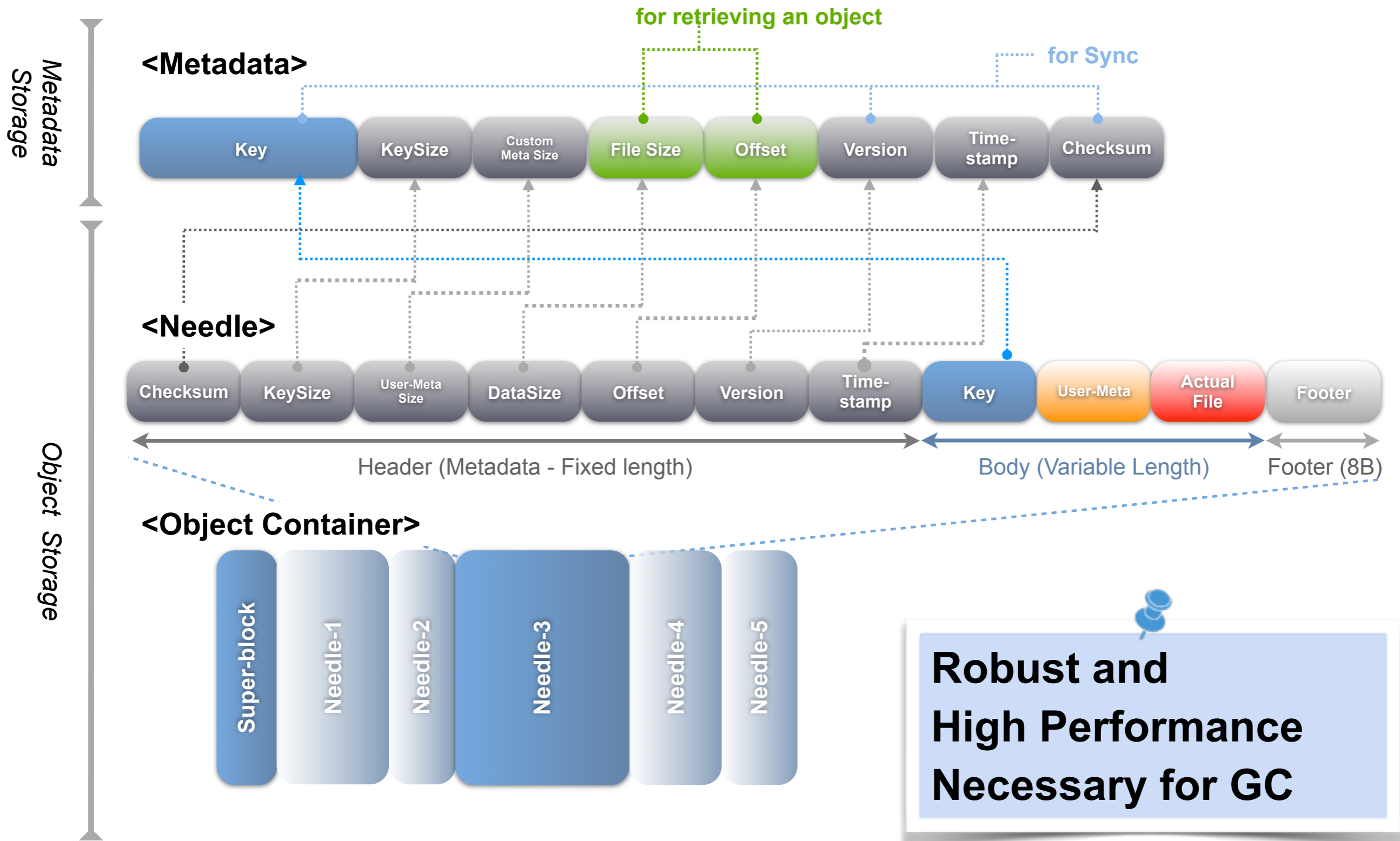
LeoFS Overview - Storage

Storage consists of *Object Storage and Metadata Storage*

Includes *Replicator and Recoverer* for the eventual consistency



LeoFS Overview - Storage - Data Structure

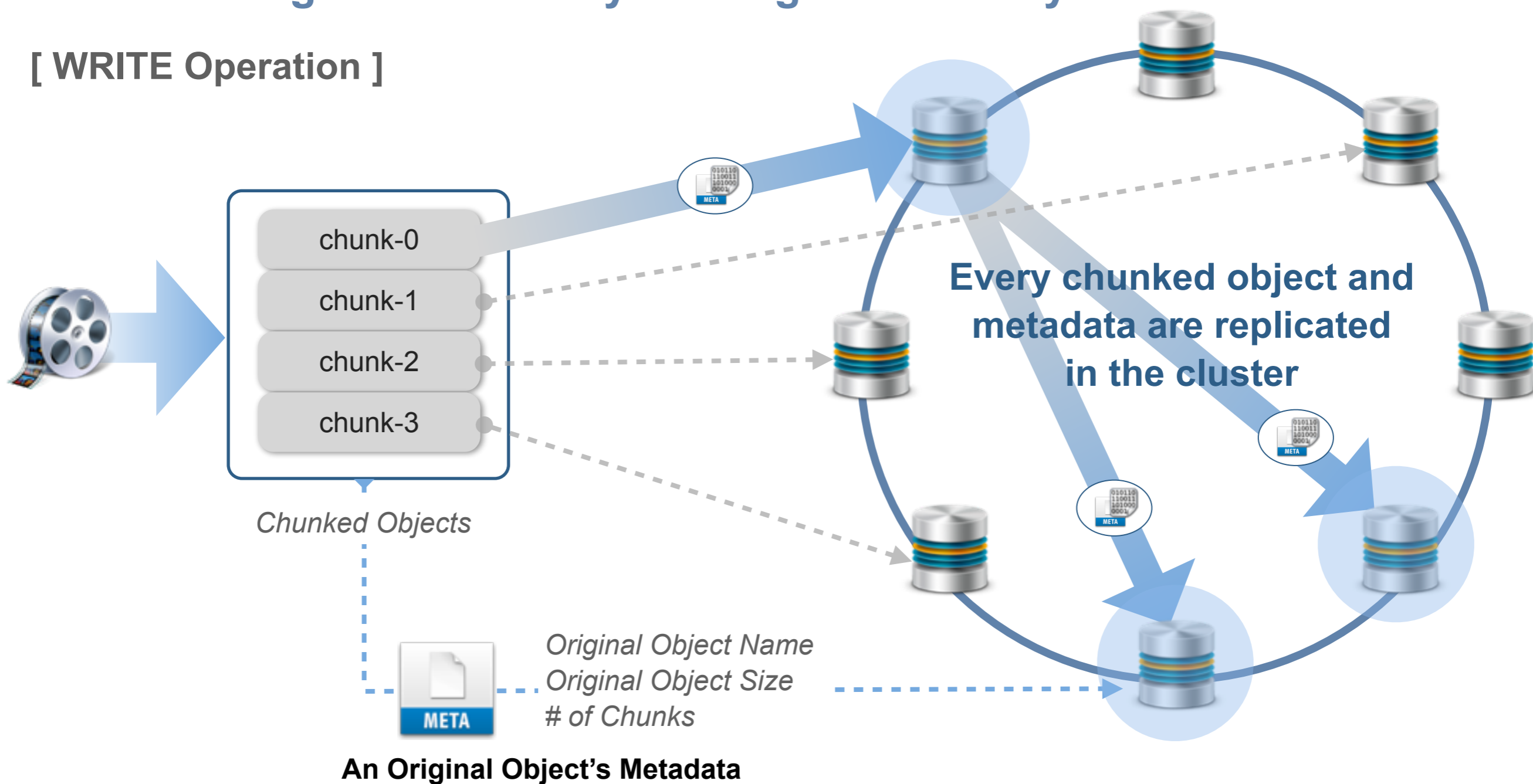


LeoFS Overview - Storage - Large Object Support

To Equalize Disk Usage in Every Storage Node

To Realize High I/O efficiency and High Availability

[WRITE Operation]



Client(s)

Gateway

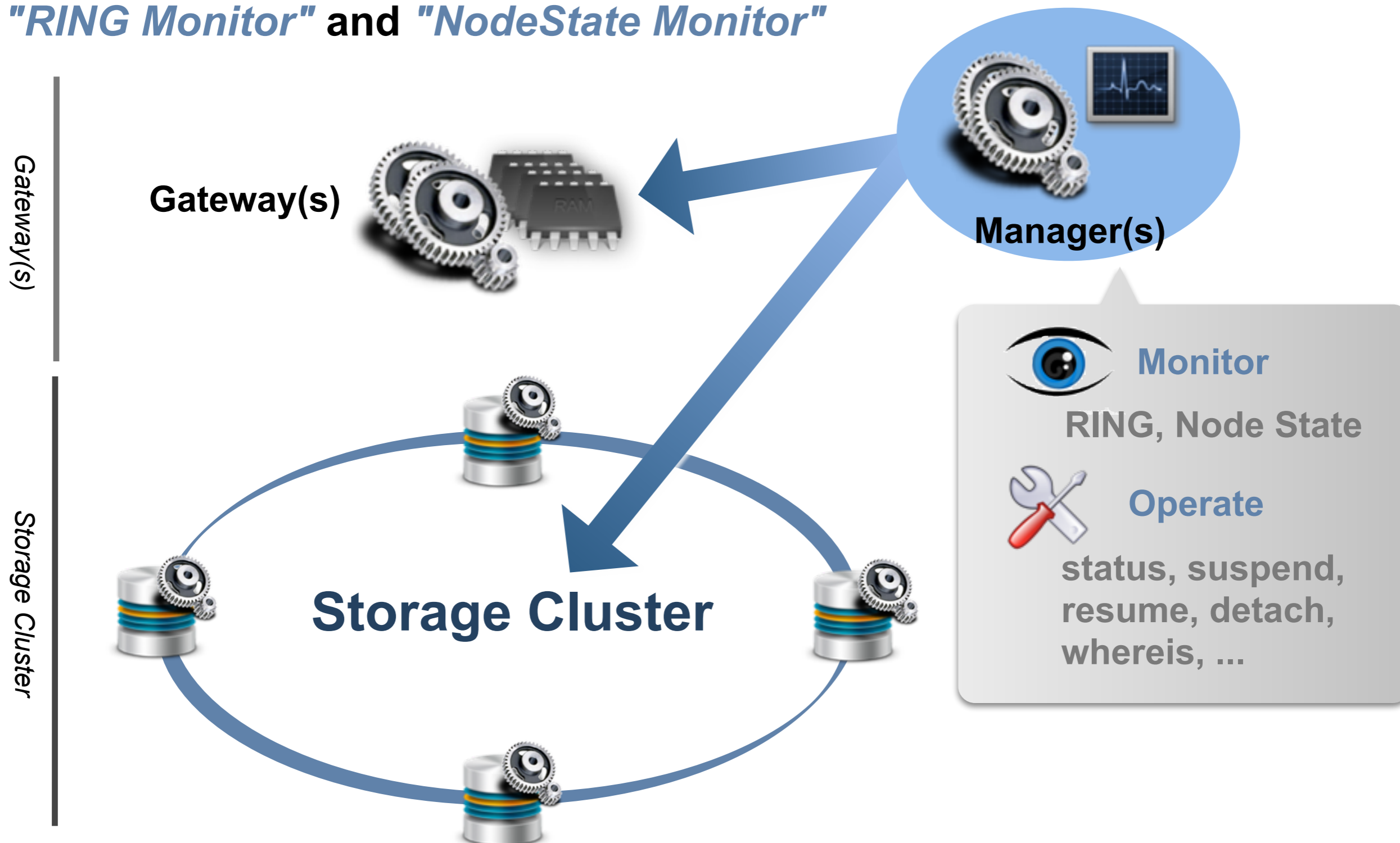
Storage Cluster

Manager

LeoFS Overview - Manager

Operate LeoFS - Gateway and Storage Cluster

"RING Monitor" and "NodeState Monitor"



Brief Benchmark Report



Brief Benchmark Report

Summary of the benchmark results

LeoFS kept in a stable performance through the benchmark

Bottleneck is Disk I/O

The cache mechanism contributed to reduce network traffic between Gateway and Storage

Brief Benchmark Report

1st Case:

Group of Value Ranges (HDD)

Storage:5, Gateway:1, Manager:2

R:W = 9:1

source: https://github.com/leo-project/notes/tree/master/leofs/benchmark/leofs/20140605/tests/1m_r9w1_240min

2nd Case:

Group of Value Ranges (HDD)

Storage:5, Gateway:1, Manager:2

R:W = 8:2

source: https://github.com/leo-project/notes/tree/master/leofs/benchmark/leofs/20140605/tests/1m_r8w2_120min

Brief Benchmark Report

Server Spec - Gateway:

CPU	Intel(R) Xeon(R) CPU X5650 @ 2.67GHz * 2 (12 cores / 24 threads)
Memory	96GB
Disk	HDD - 240GB RAID0
Network	10G-Ether

Server Spec - Storage x5:

CPU	Intel(R) Xeon(R) CPU X5650 @ 2.67GHz * 2 (12 cores / 24 threads)
Memory	96GB
Disk	HDD - 240GB RAID0 (System)
	HDD - 2TB RAID0 (Data)
Network	10G-Ether

Brief Benchmark Report - 1st Case (HDD / R:W=9:1)

Environment:

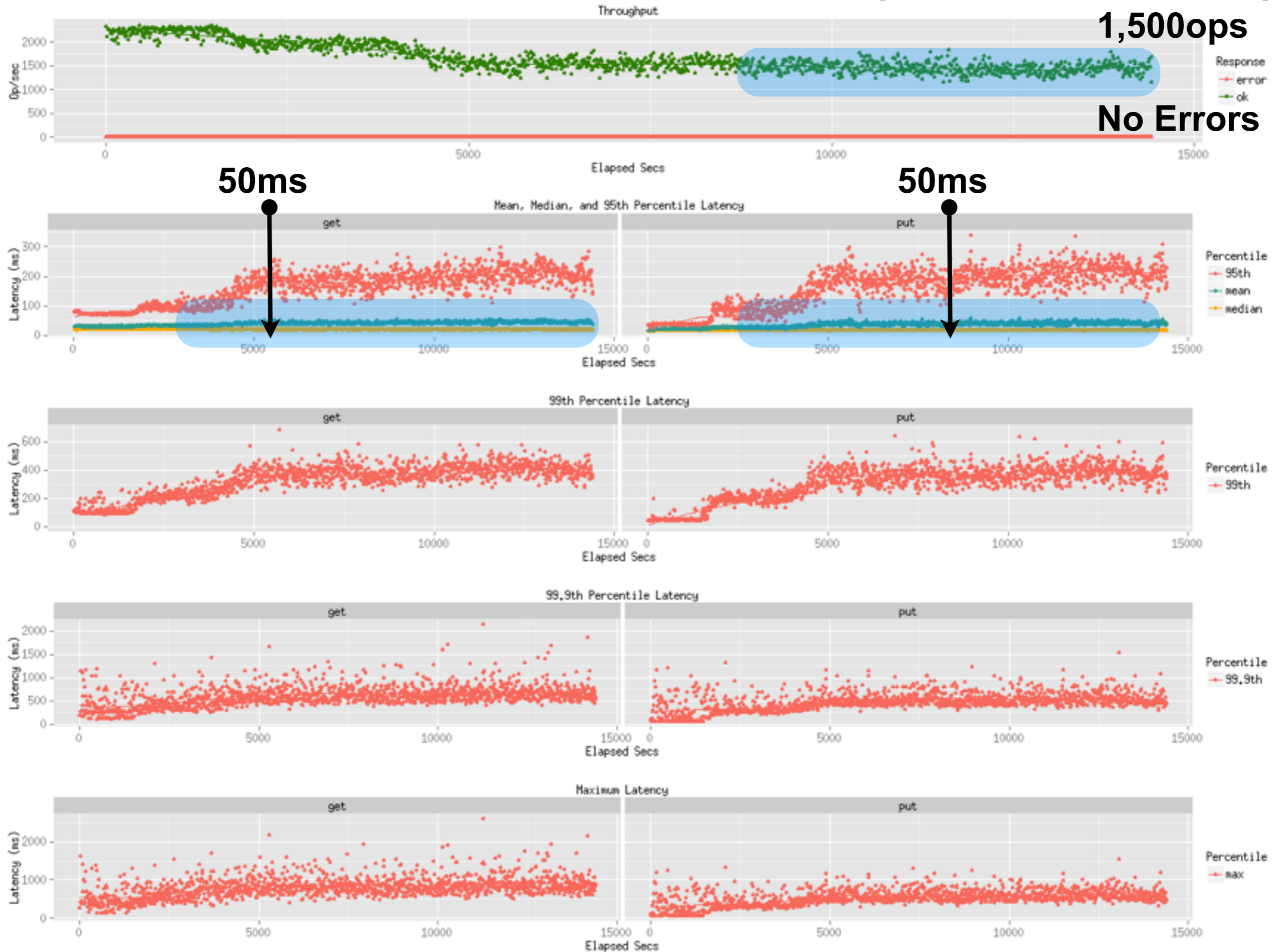
Network	10Gbps
OS	CentOS release 6.5 (Final)
Erlang	OTP R16B03-1
LeoFS	v1.0.2

System Consistency Level: [N:3, W:2, R:1, D:2]

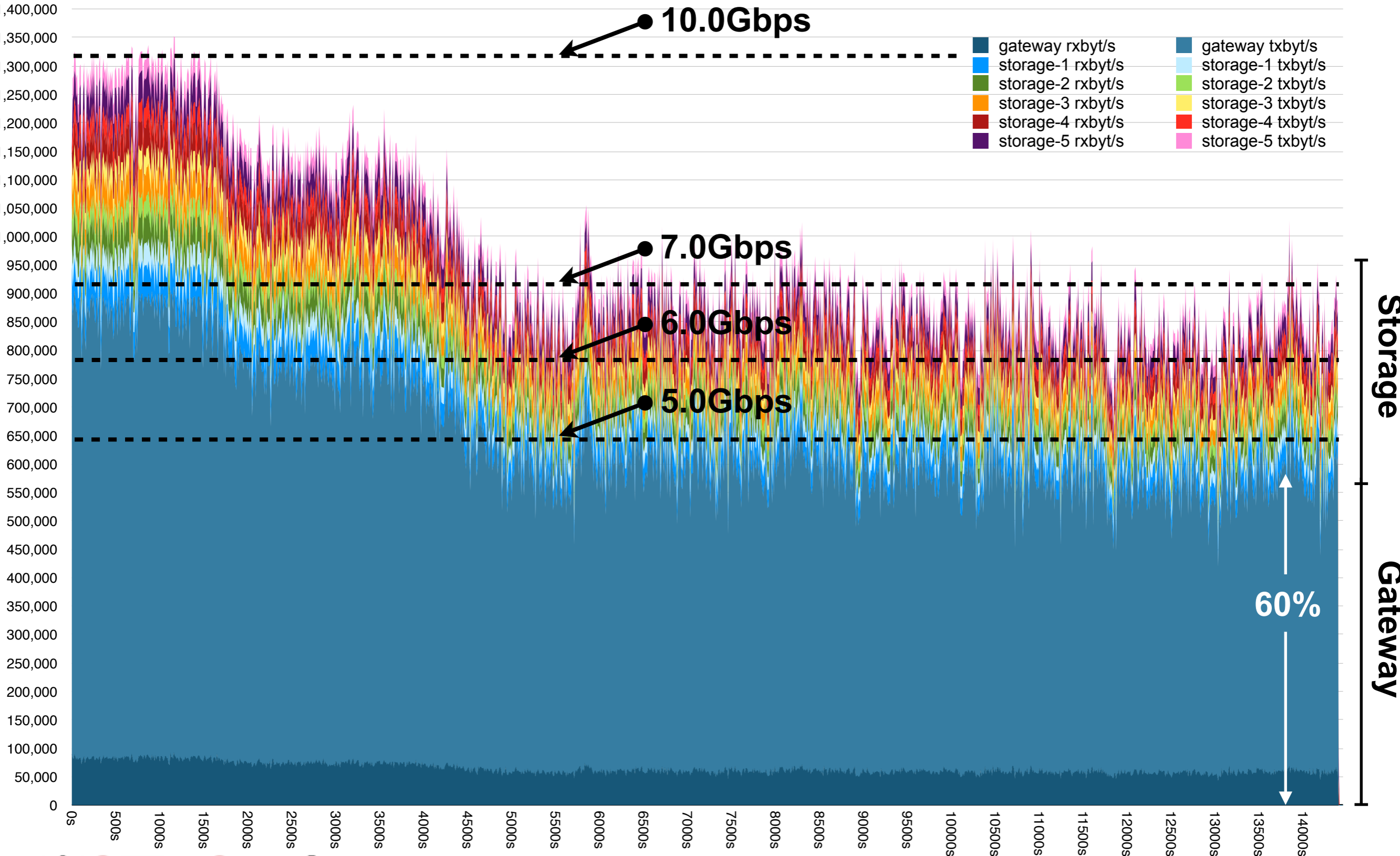
Benchmark Configuration:

Duration	4.0h															
R:W	9:1															
# of Concurrent Processes	64															
# of Keys	100,000															
Value Size	<table border="1"><thead><tr><th colspan="2">Range (byte)</th><th>Percentage</th></tr></thead><tbody><tr><td>1024</td><td>10240</td><td>24%</td></tr><tr><td>10241</td><td>102400</td><td>30%</td></tr><tr><td>10241</td><td>819200</td><td>30%</td></tr><tr><td>819201</td><td>1572864</td><td>16%</td></tr></tbody></table>	Range (byte)		Percentage	1024	10240	24%	10241	102400	30%	10241	819200	30%	819201	1572864	16%
Range (byte)		Percentage														
1024	10240	24%														
10241	102400	30%														
10241	819200	30%														
819201	1572864	16%														

Brief Benchmark Report - 1st Case (HDD / R:W=9:1)

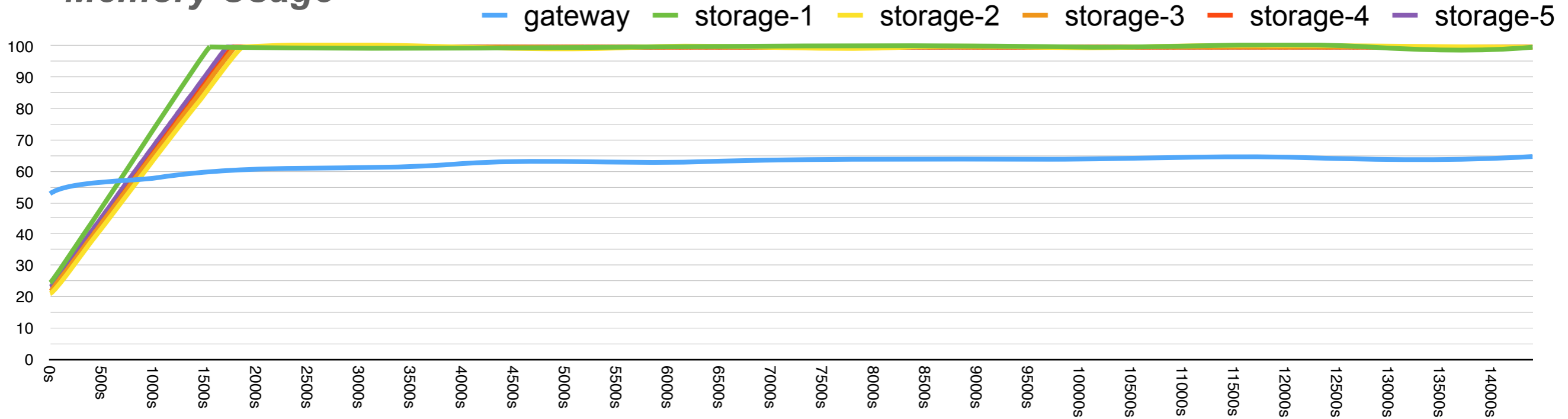


Brief Benchmark Report - 1st Case / Network Traffic

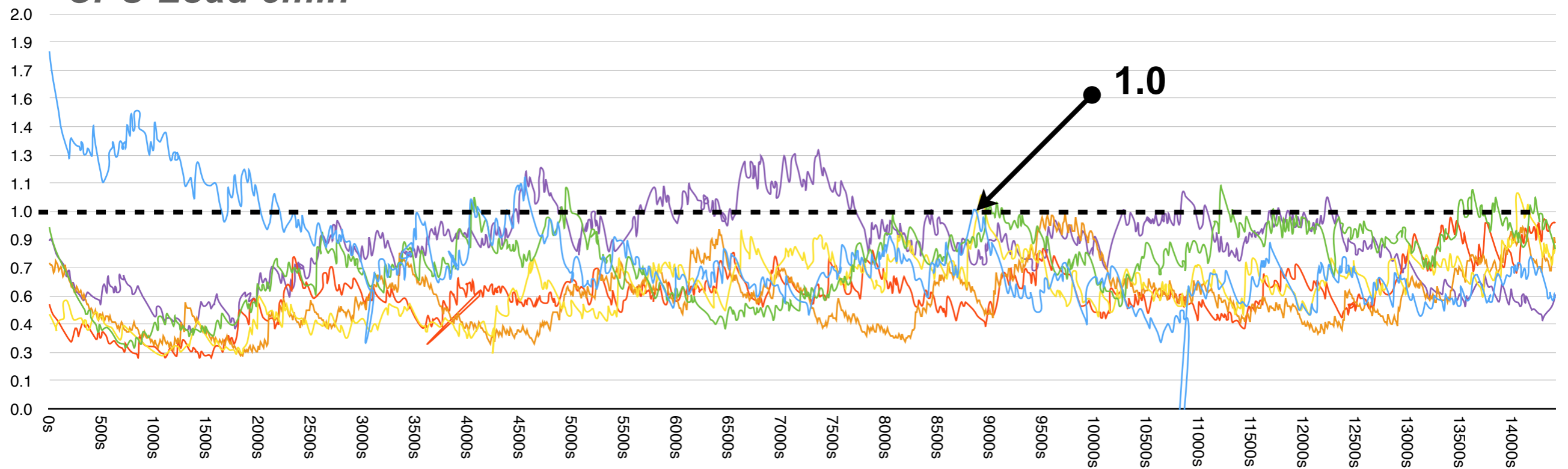


Brief Benchmark Report - 1st Case / Memory and CPU

Memory Usage



CPU Load 5min



Brief Benchmark Report - 2nd Case (HDD / R:W=8:2)

Environment:

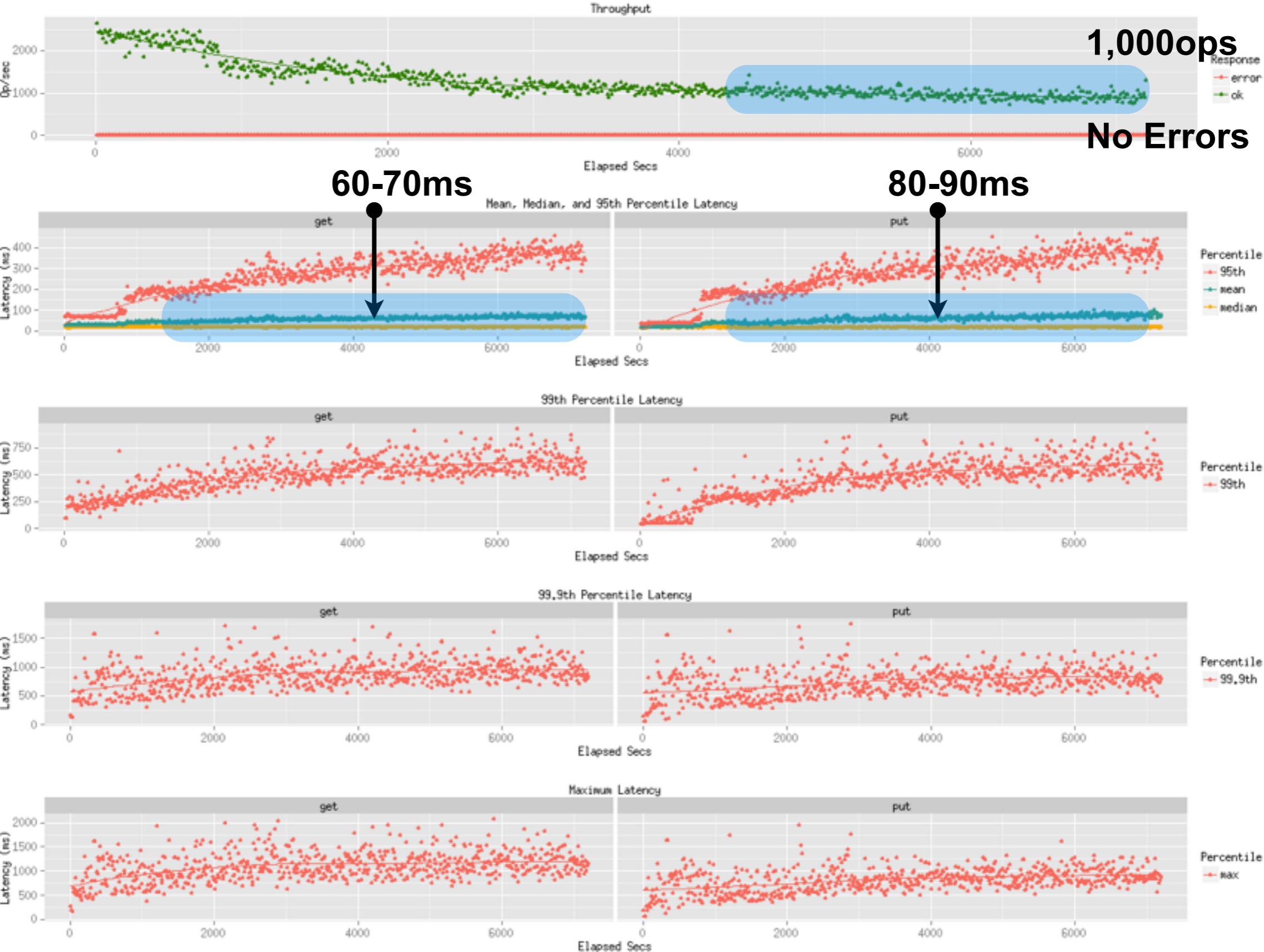
Network	10Gbps
OS	CentOS release 6.5 (Final)
Erlang	OTP R16B03-1
LeoFS	v1.0.2

System Consistency Level: [N:3, W:2, R:1, D:2]

Benchmark Configuration:

Duration	2.0h															
R:W	8:2															
# of Concurrent Processes	64															
# of Keys	100,000															
Value Size	<table border="1"><thead><tr><th colspan="2">Range (byte)</th><th>Percentage</th></tr></thead><tbody><tr><td>1024</td><td>10240</td><td>24%</td></tr><tr><td>10241</td><td>102400</td><td>30%</td></tr><tr><td>10241</td><td>819200</td><td>30%</td></tr><tr><td>819201</td><td>1572864</td><td>16%</td></tr></tbody></table>	Range (byte)		Percentage	1024	10240	24%	10241	102400	30%	10241	819200	30%	819201	1572864	16%
Range (byte)		Percentage														
1024	10240	24%														
10241	102400	30%														
10241	819200	30%														
819201	1572864	16%														

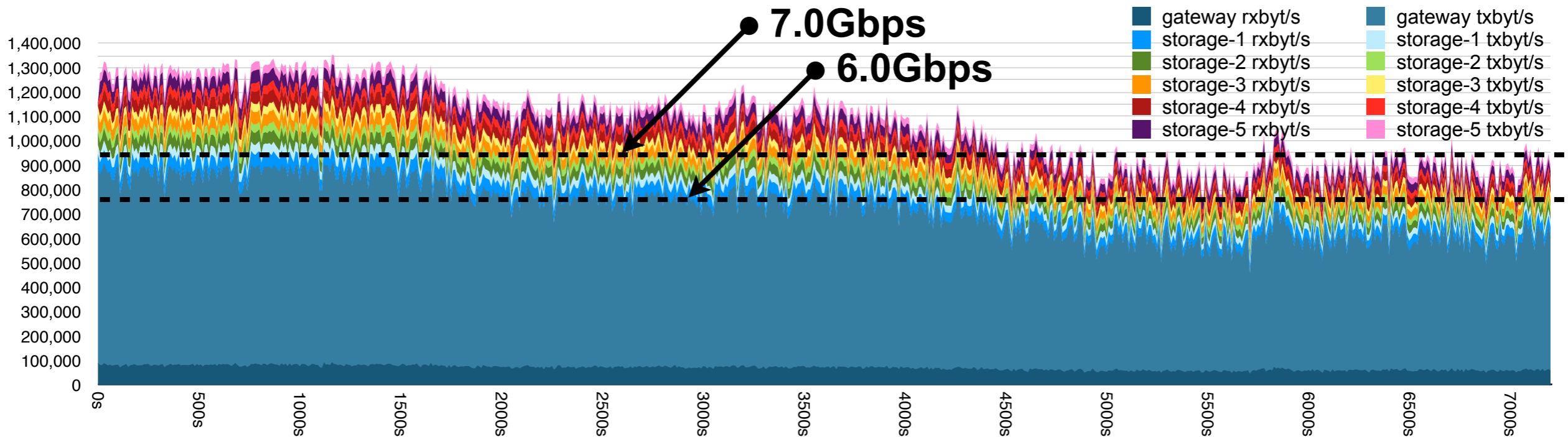
Brief Benchmark Report - 2nd Case (HDD / R:W=8:2)



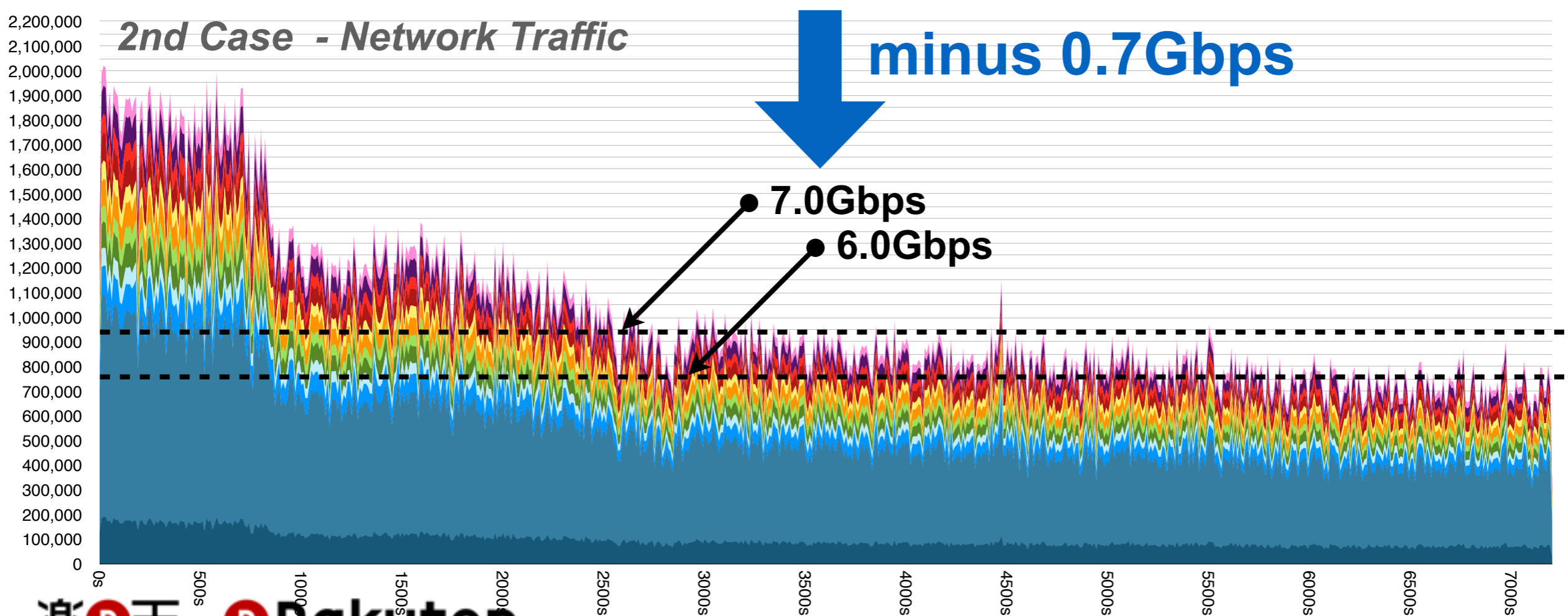
**Compare 1st case
with 2nd case**

Brief Benchmark Report

1st Case - Network Traffic

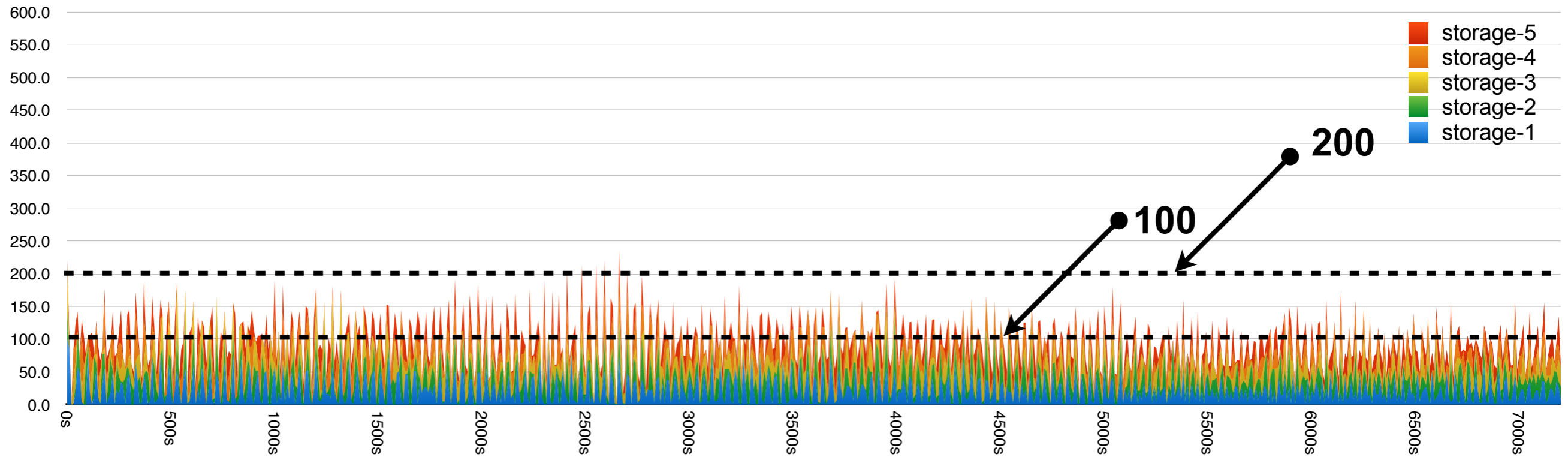


2nd Case - Network Traffic

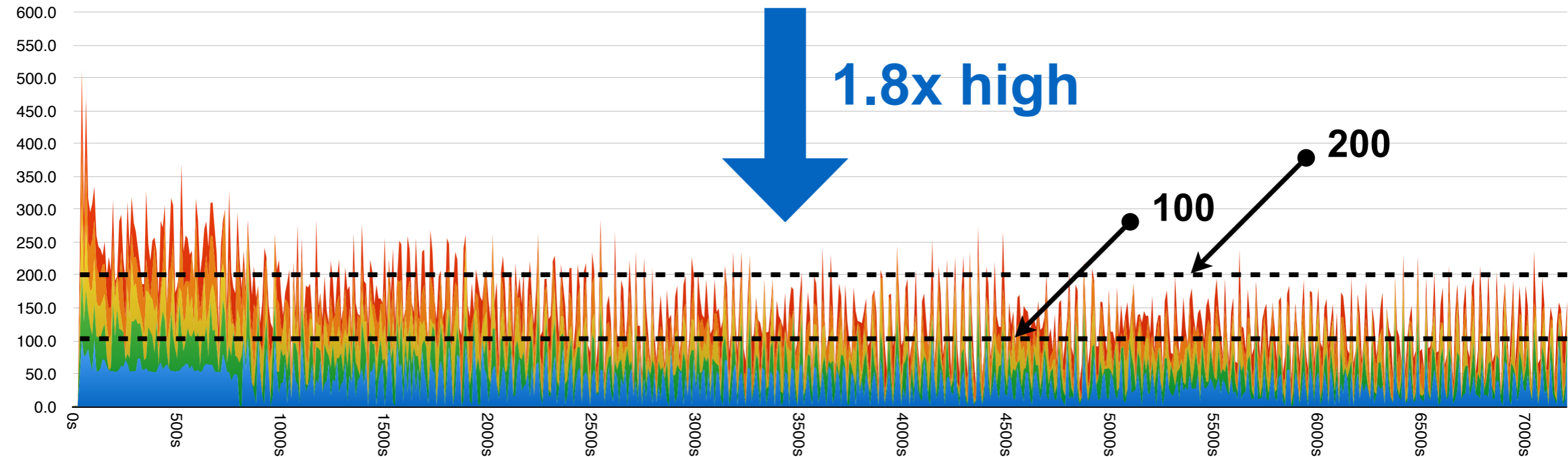


Brief Benchmark Report

1st Case - Disk util%

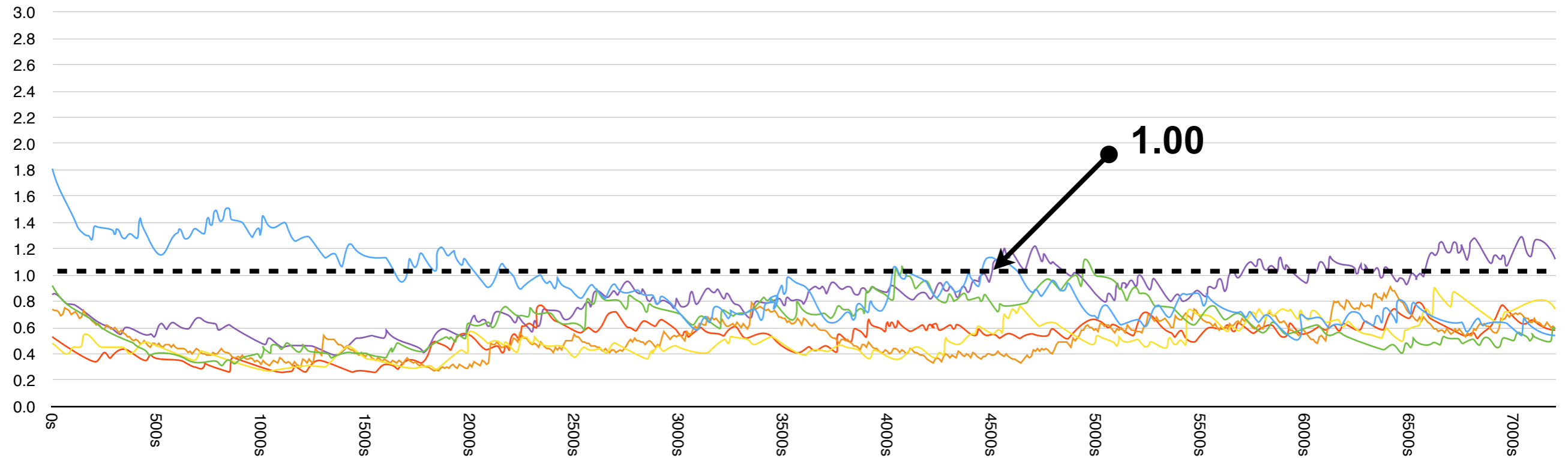


2nd Case - Disk util%

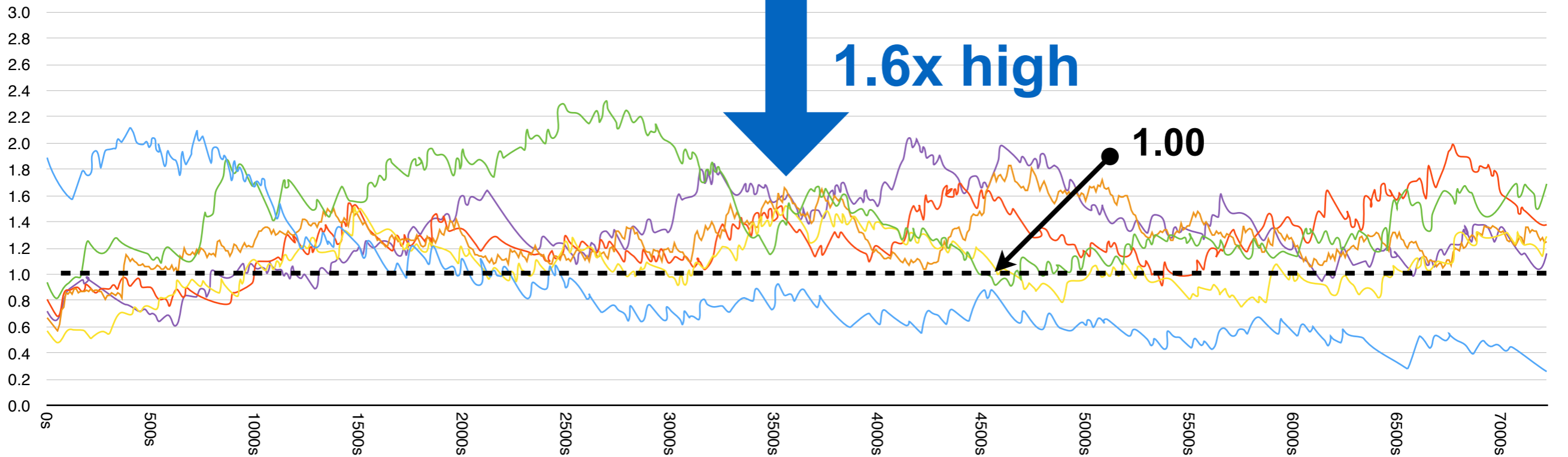


Brief Benchmark Report

1st Case - CPU Load 5min



2nd Case - CPU Load 5min



Brief Benchmark Report

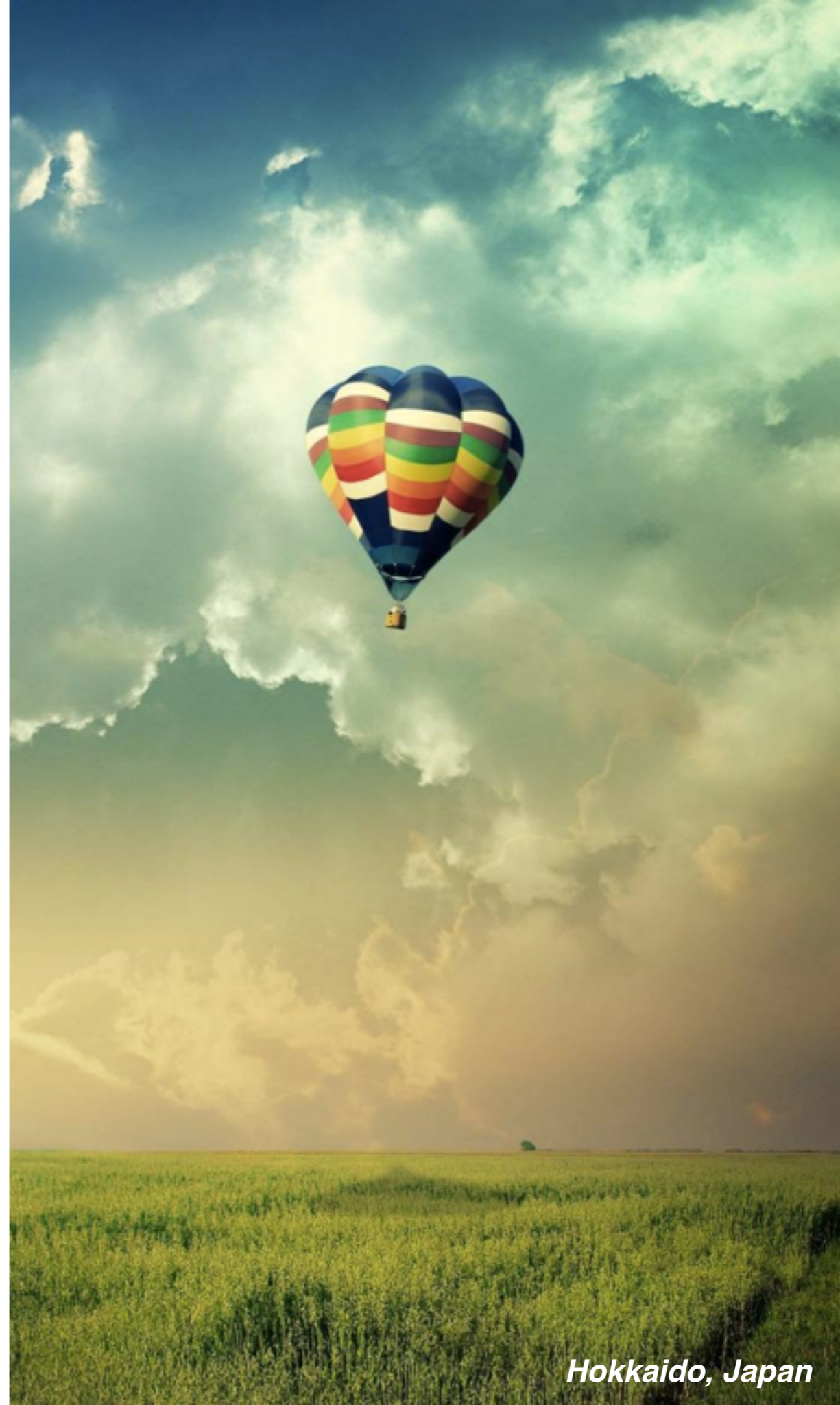
Conclusion:

LeoFS kept in **a stable performance** through the benchmark

Bottleneck is **Disk I/O**

The cache mechanism contributed to **reduce network traffic** between Gateway and Storage

Multi Data Center Replication



Multi Data Center Replication

HIGH-Scalability
HIGH-Availability



Easy Operation for Admins



US



Europe



Singapore



Tokyo



NO SPOF

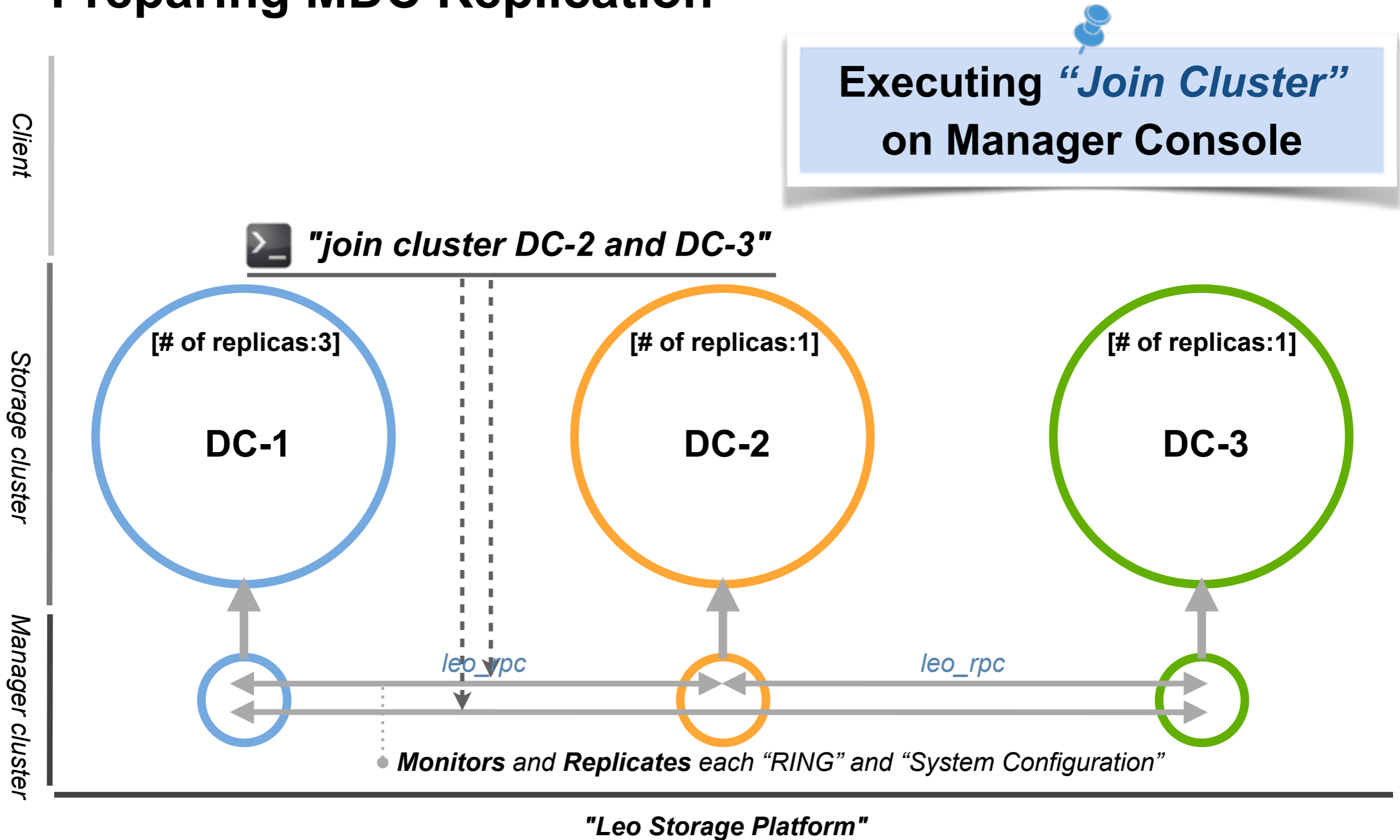
NO Performance Degradation

Designed it as simple as possible

1. Easy Operation to build **multi clusters**.
2. **Asynchronous data replication** between clusters
Stacked data is transferred to remote cluster(s)
3. **Eventual consistency**

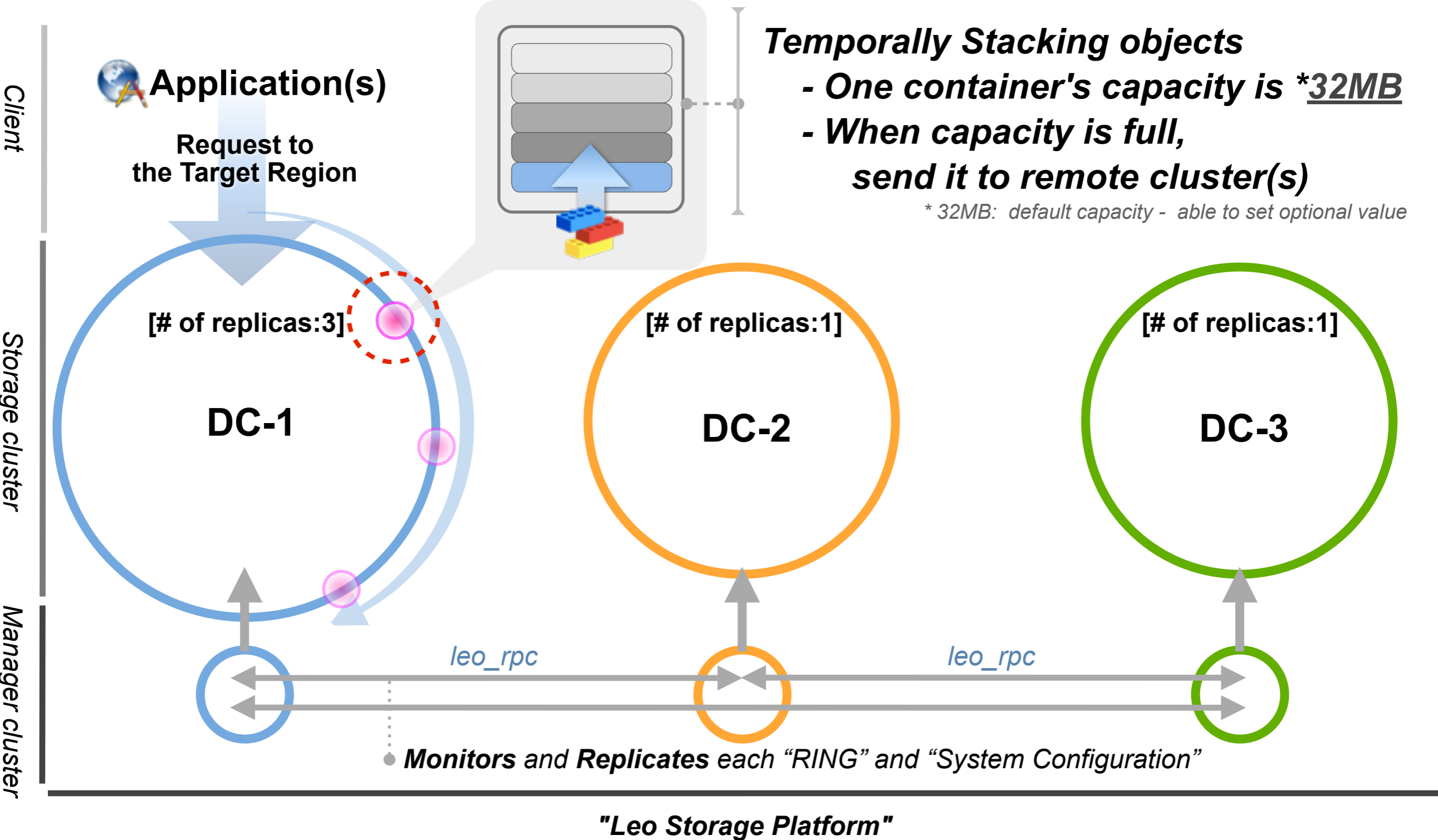
Multi Data Center Replication

Preparing MDC Replication



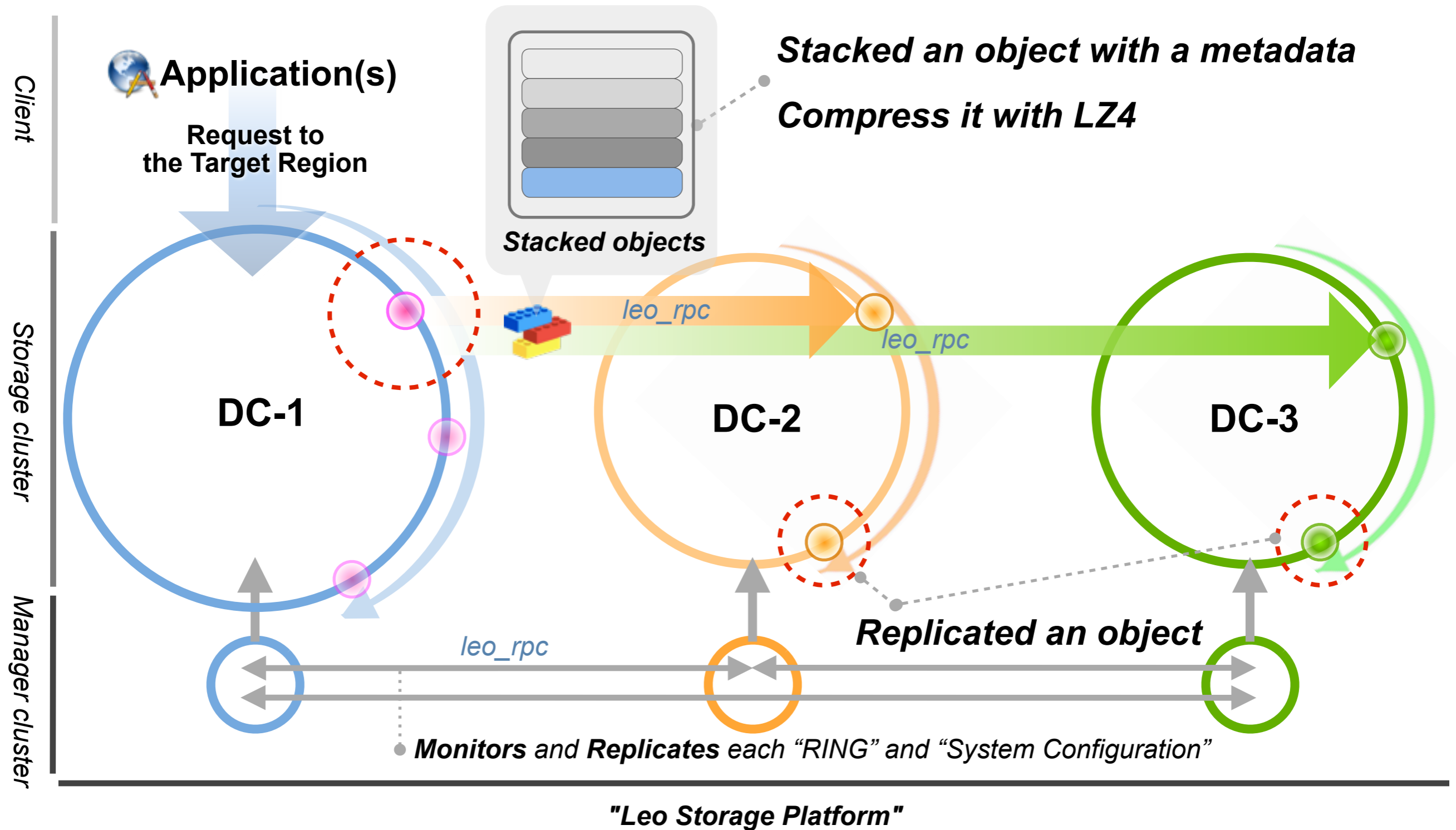
Multi Data Center Replication

Stacking objects



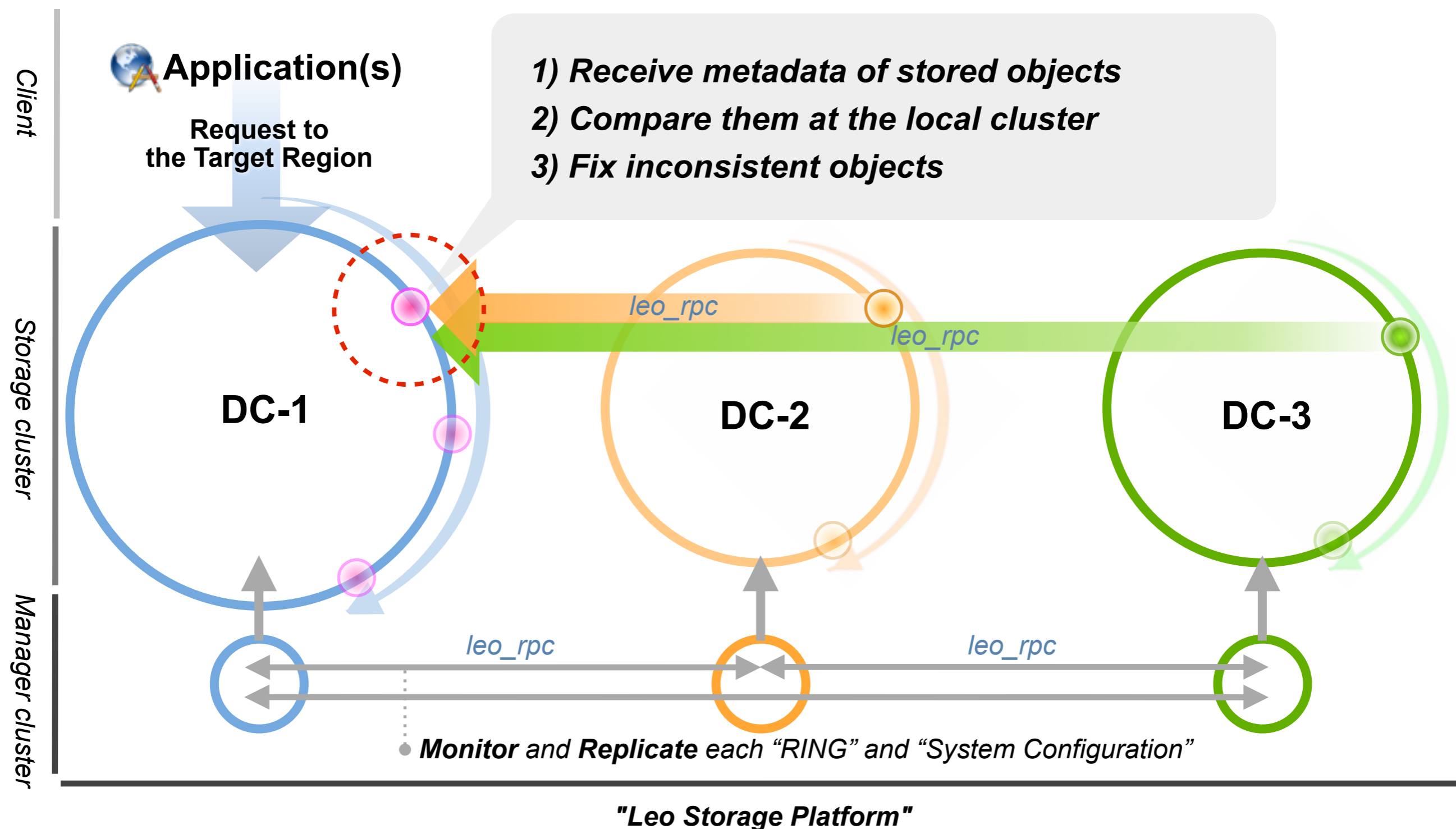
Multi Data Center Replication

Transferring stacked objects



Multi Data Center Replication

Investigating stored objects



LeoFS Administration at Rakuten

*Presented by Hiroki Matsue
Rakuten Software Engineer*



LeoFS Administration at Rakuten

Storage Platform

File Sharing Service

Others

Portal Site

Photo Storage

Background Storage of OpenStack

Storage Platform

Storage Platform - Scaling the Storage Platform

Reduce Costs
High Reliability
Easy to Scale
S3-API



Storage Platform - Scaling the Storage Platform

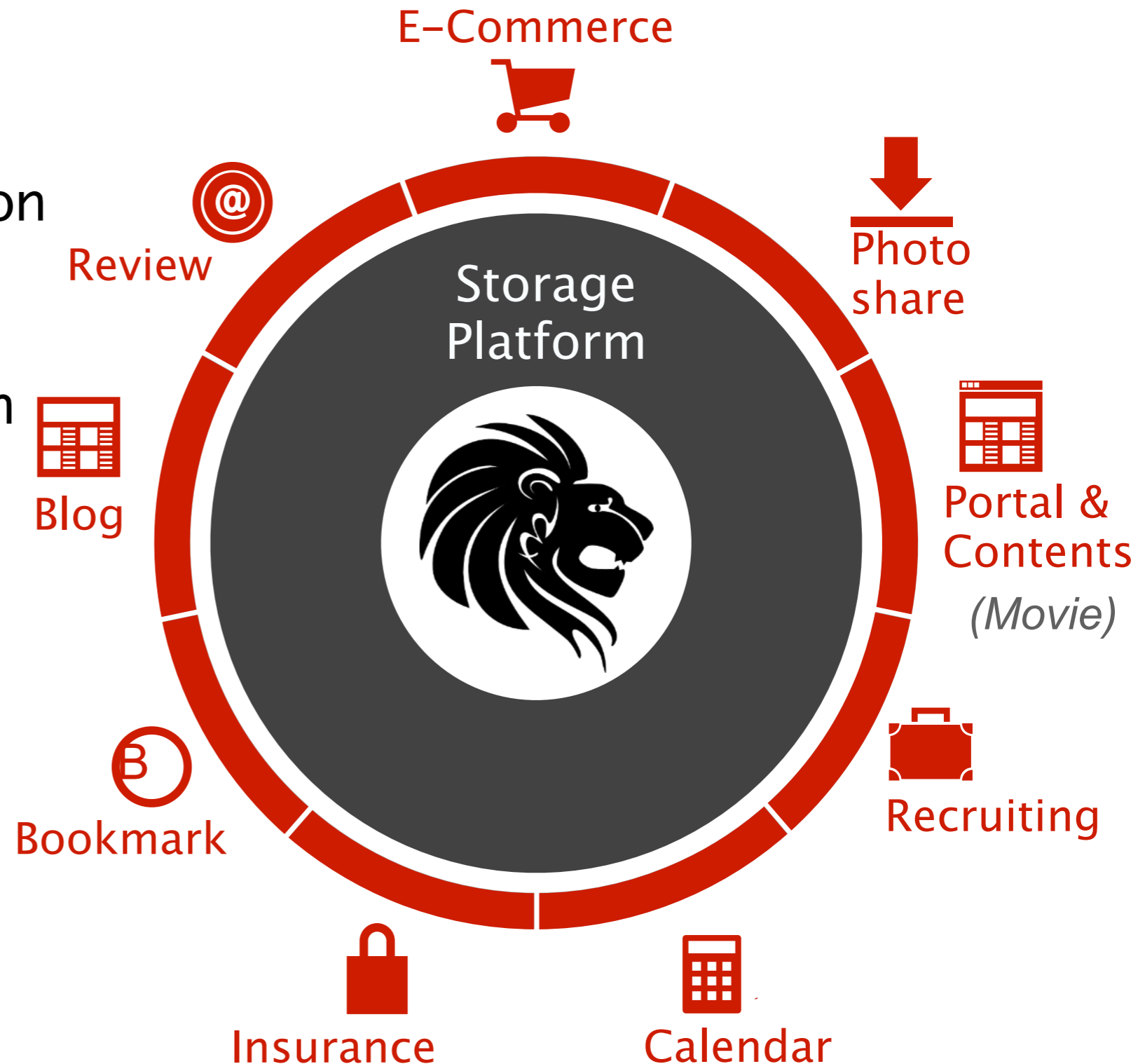
Using Various Services

Total Usage: 450TB

of Files: 600Million

Daily Growth: 100GB

Daily Reqs: 13Million



Storage Platform - System Layout

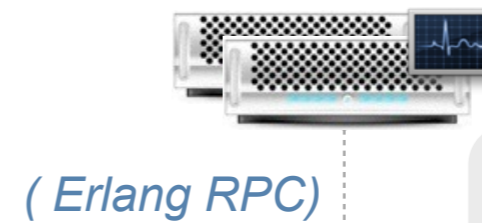
Total disk space: 600TB
Number of Files: 600Million
Access Stats:
800Mbps (MAX)
400Mbps (AVG)

*Requests from
Web Applications / Browsers
w/HTTP over S3-API*

Load Balancer / Cache Servers



Manager x 2



(Erlang RPC)

(TCP/IP,SNMP)

Nagios
Ganglia

Monitor

GUI Console

Gateway x 4

Storage x 14

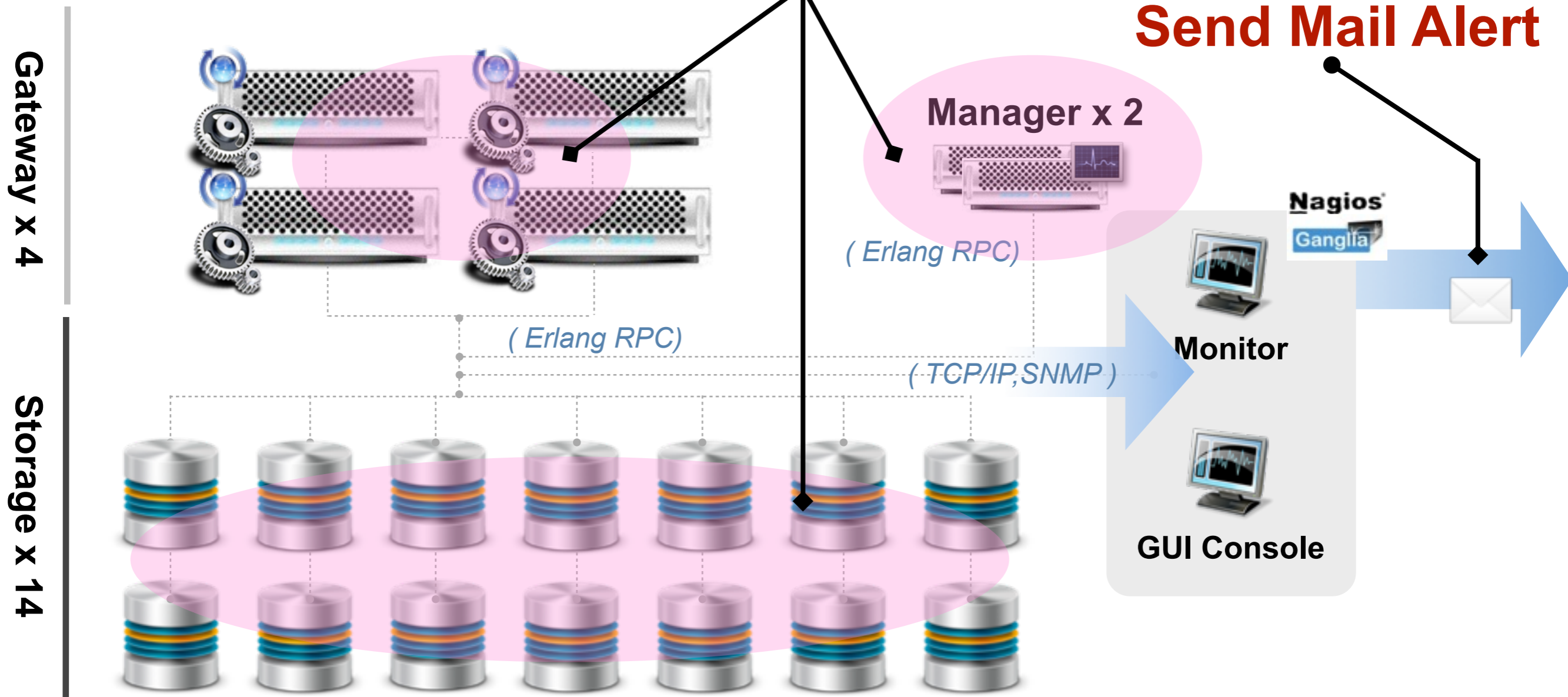


Storage Platform - Monitor

Status Collection (*Ganglia*)
Status Check (*Nagios*)
Port + Threshold Check

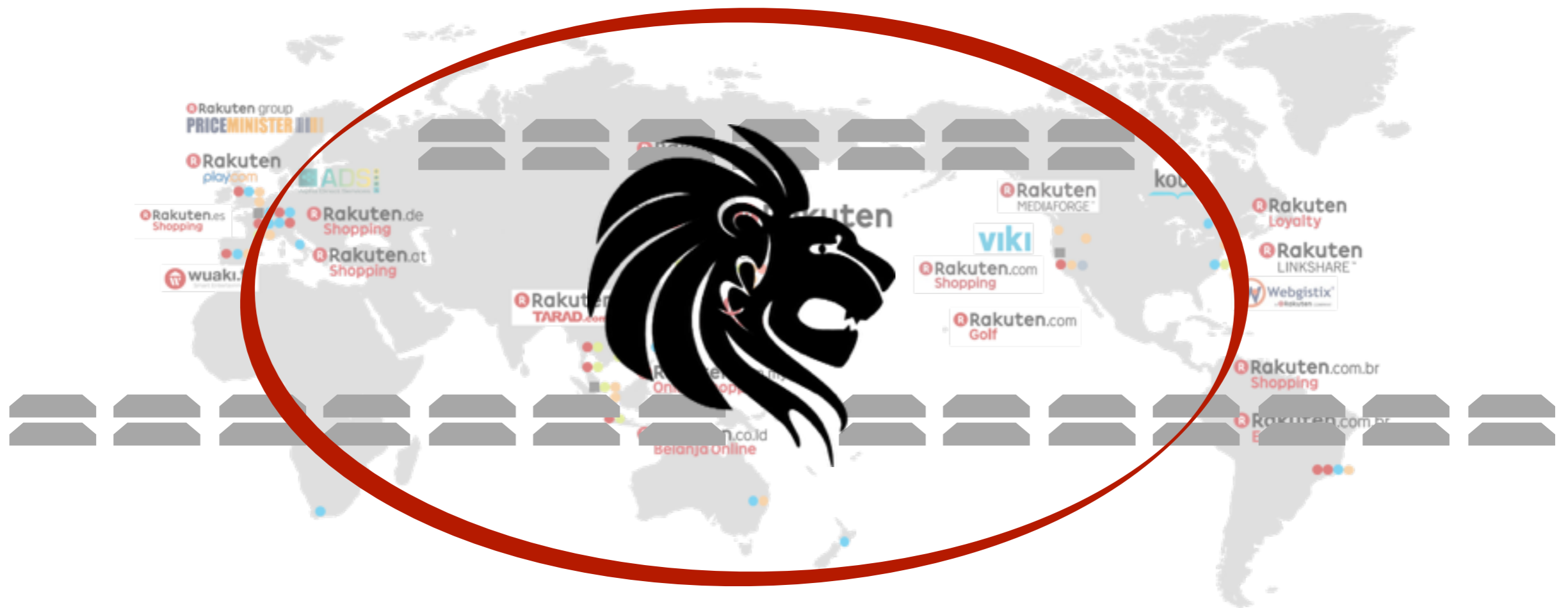
Ganglia and Nagios Agent

Send Mail Alert



Storage Platform - Spreading Globally

Covering All Services with Multi DC Replication



File Sharing Service

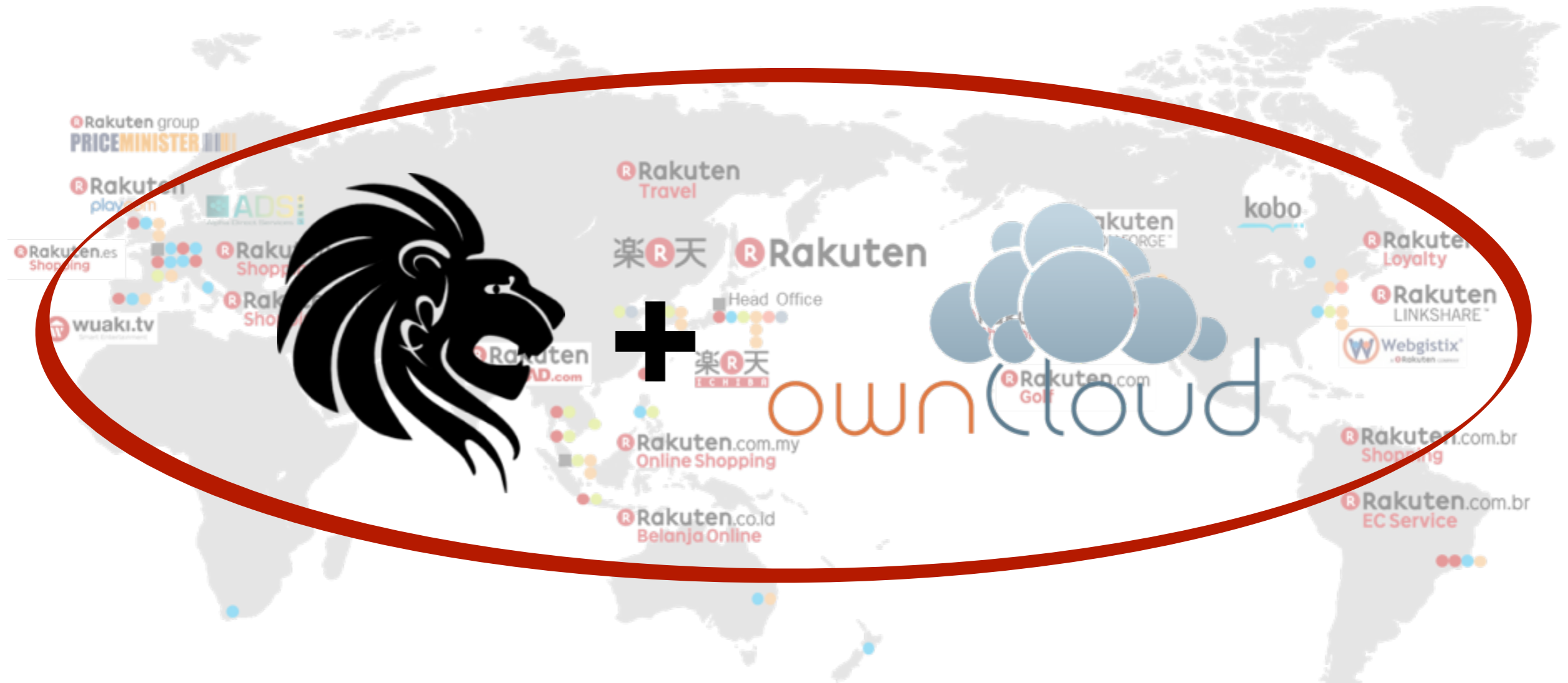


+



<https://owncloud.com/>

File Sharing Service - Usage

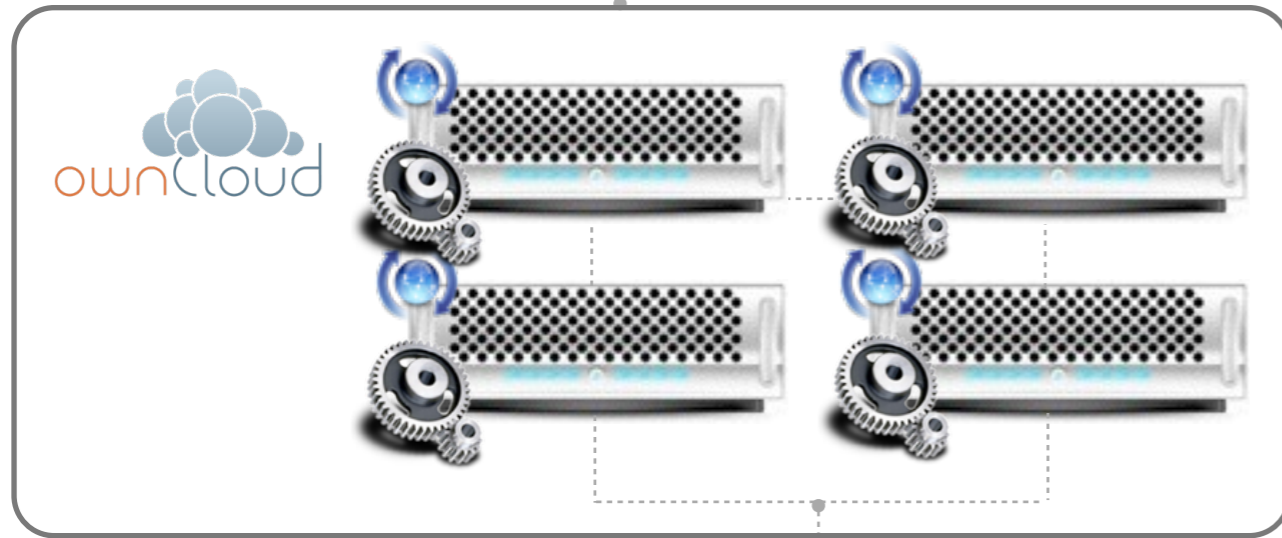


Share **Docs** and **Movies** with Group Companies
Over **20** Companies, Over **10** Countries
Over **4,000** Users, Over **10,000** Teams

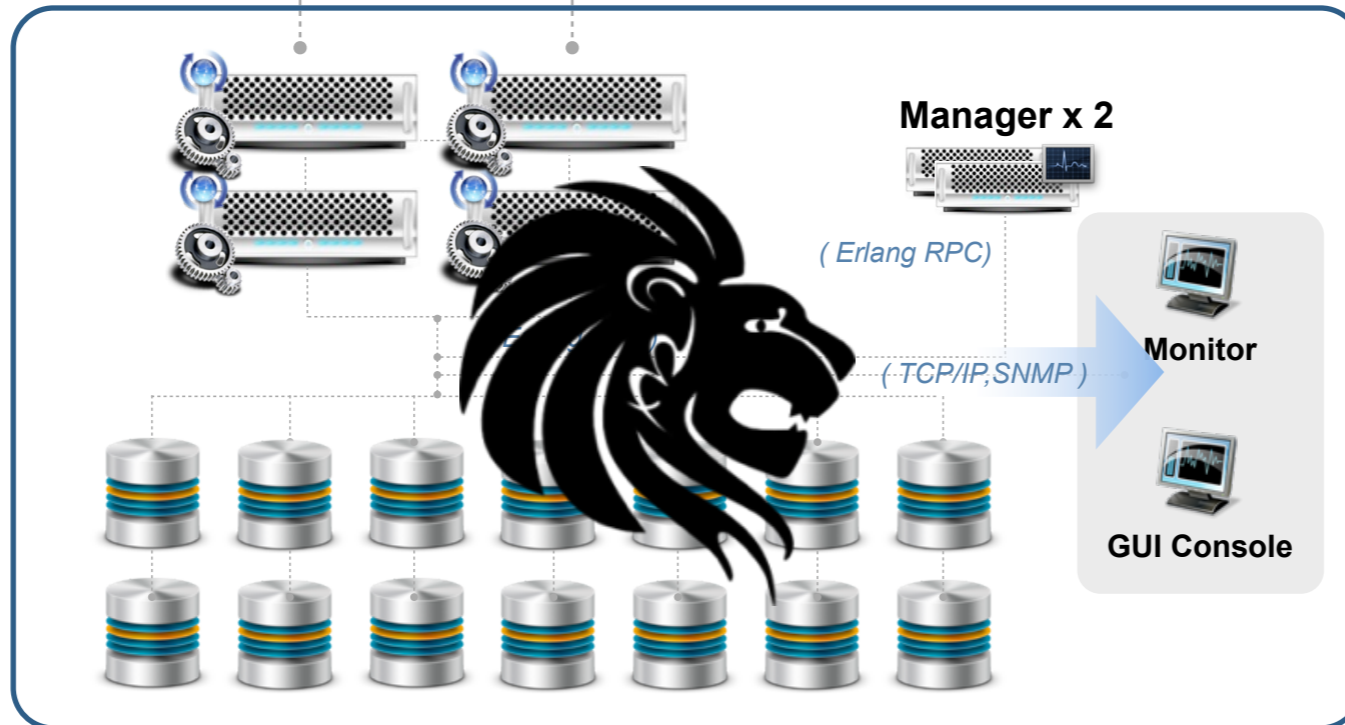
File Sharing Service - System Layout

Web GUI File Browser

Authenticate Users



LDAP

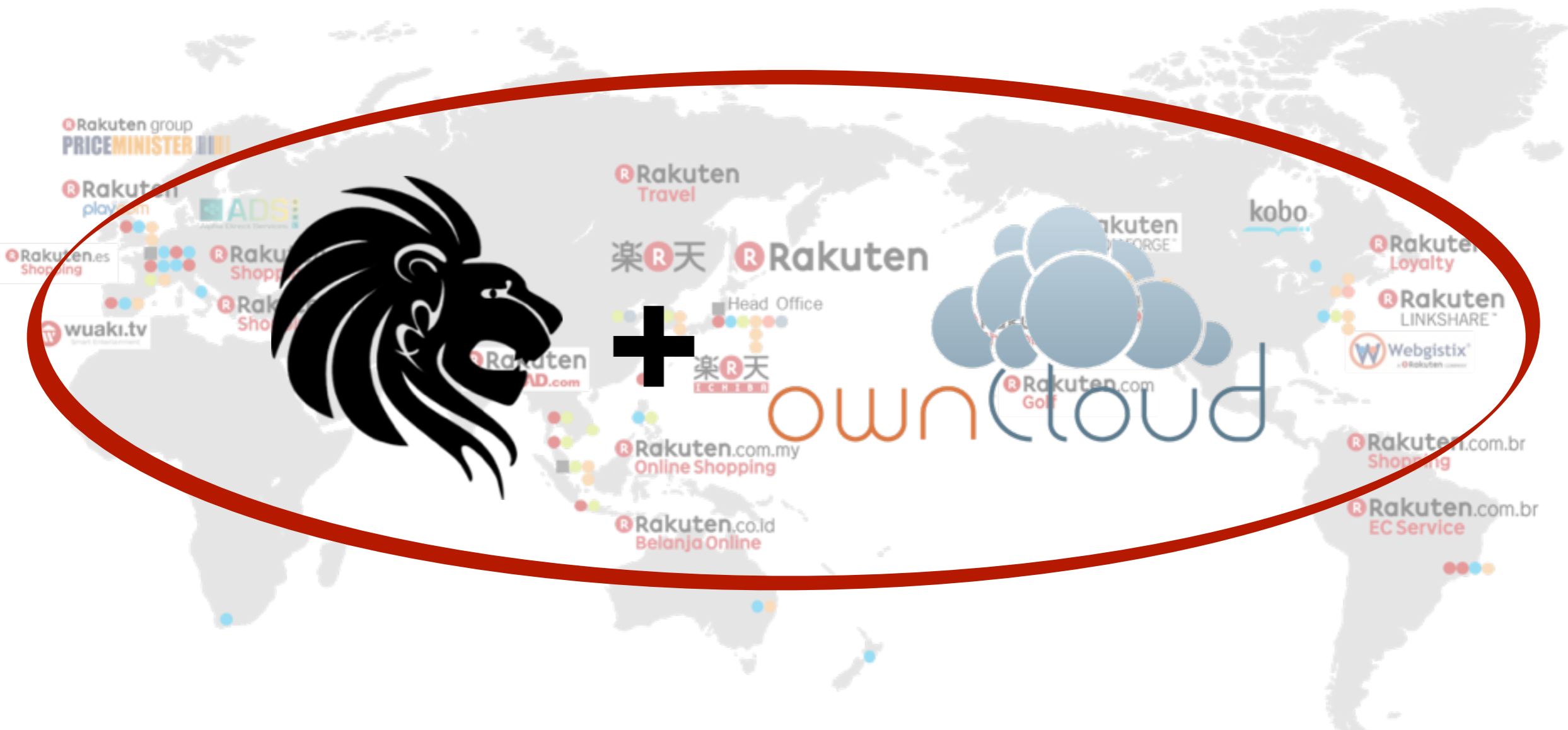


Manage Configurations

Manage Login Session (KVS)

File Sharing Service - Future Plans

Cover 25 Countries/Regions Over 20,000 Users



Empowering the Services and the Users Through the Cloud Storage

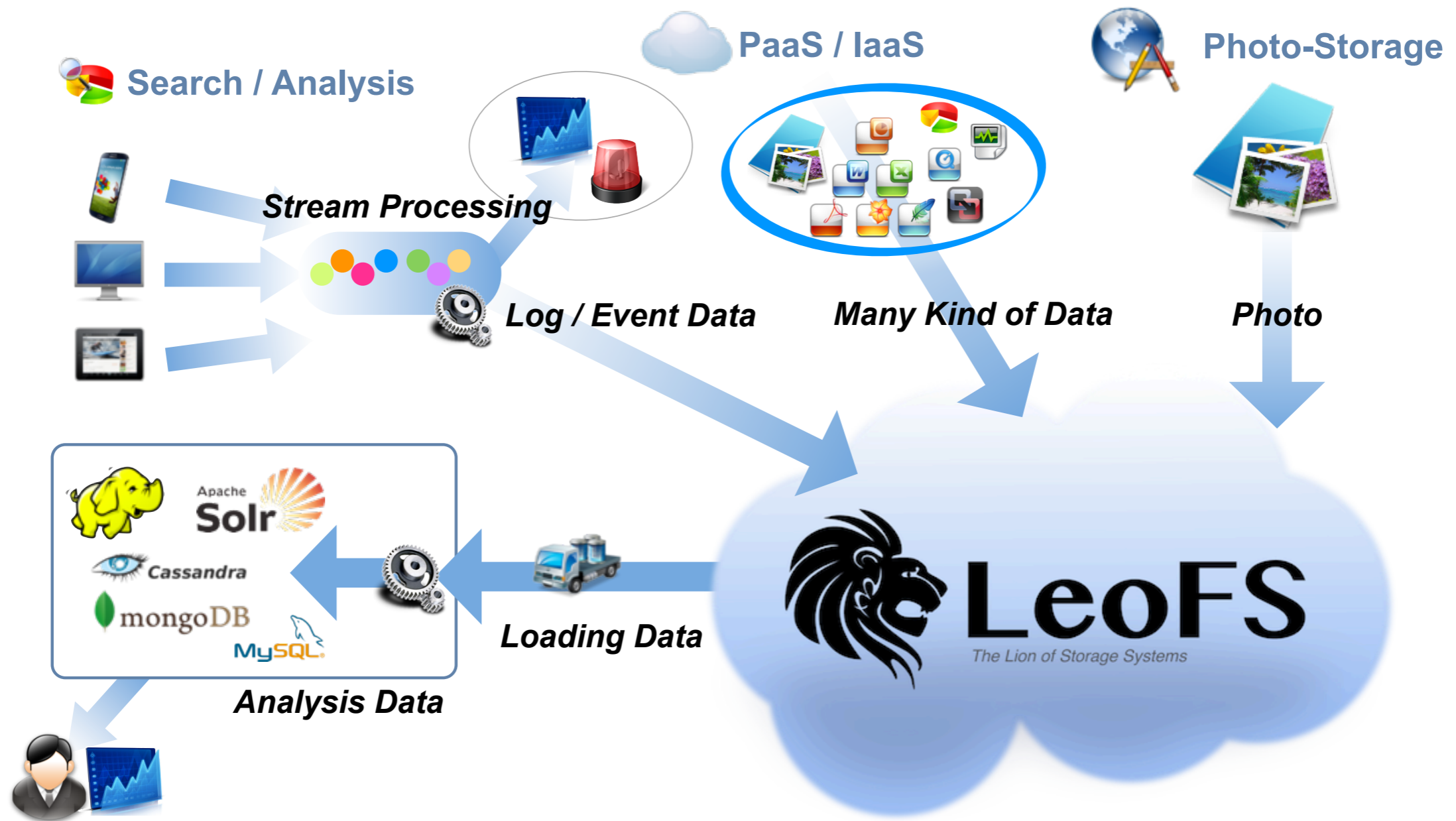


Future Plans

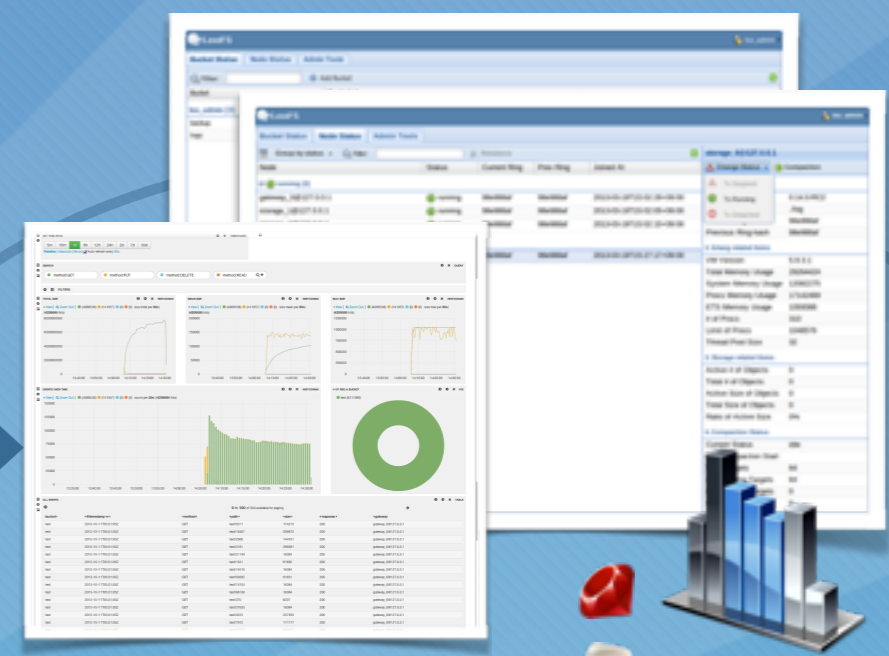
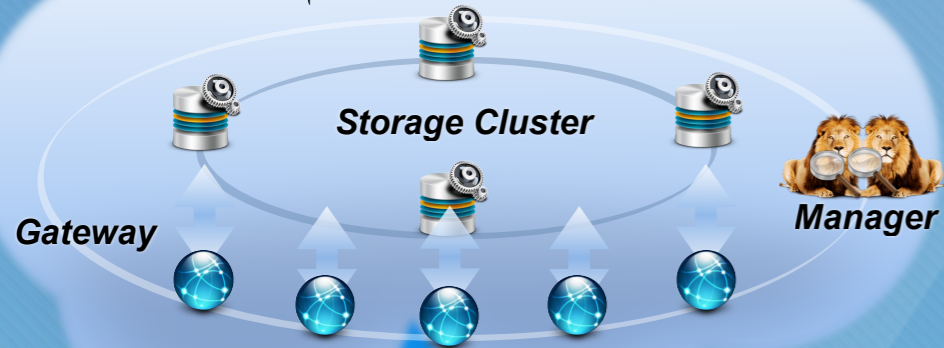
Future Plans

NFS Support

Data-HUB: Centralize unstructured data in LeoFS



Future Plans



SavannaDB's Agent
Insight LeoFS



REST-API (JSON)

**SavannaDB for Statistics Data
LeoInsight**



openstack™



Set Sail for “Cloud Storage”

Website: leo-project.net

Twitter: [@LeoFastStorage](https://twitter.com/LeoFastStorage)

Facebook: www.facebook.com/org.leofs