# LeoFS

# Reliable, Scaling and High Performance Storage System

**Yosuke Hara** - @yosukehara

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

with **Masahiro Sanjo**, Coordinator of R.I.T.

楽R天　Ⓡ Rakuten

R.I.T.
Ⓡ Rakuten
Institute of Technology

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

NFS Support

LeoFS Administration at Rakuten

Future Plans
    LeoFS QoS

楽R天 ®Rakuten

# 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**

楽R天 RRakuten

# LeoFS Overview

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

*Load Balancer*

## Keeping High Availability
## Keeping High Performance
## Easy Administration

## Gateway

## Manager

( Erlang RPC)

## Storage

( Erlang RPC)

( TCP/IP,SNMP )

Nagios
ZABBIX

**Monitor**

Storage Engine/Router

Storage Engine/Router

Storage Engine/Router

**GUI Console**

RAM

RAM

RAM

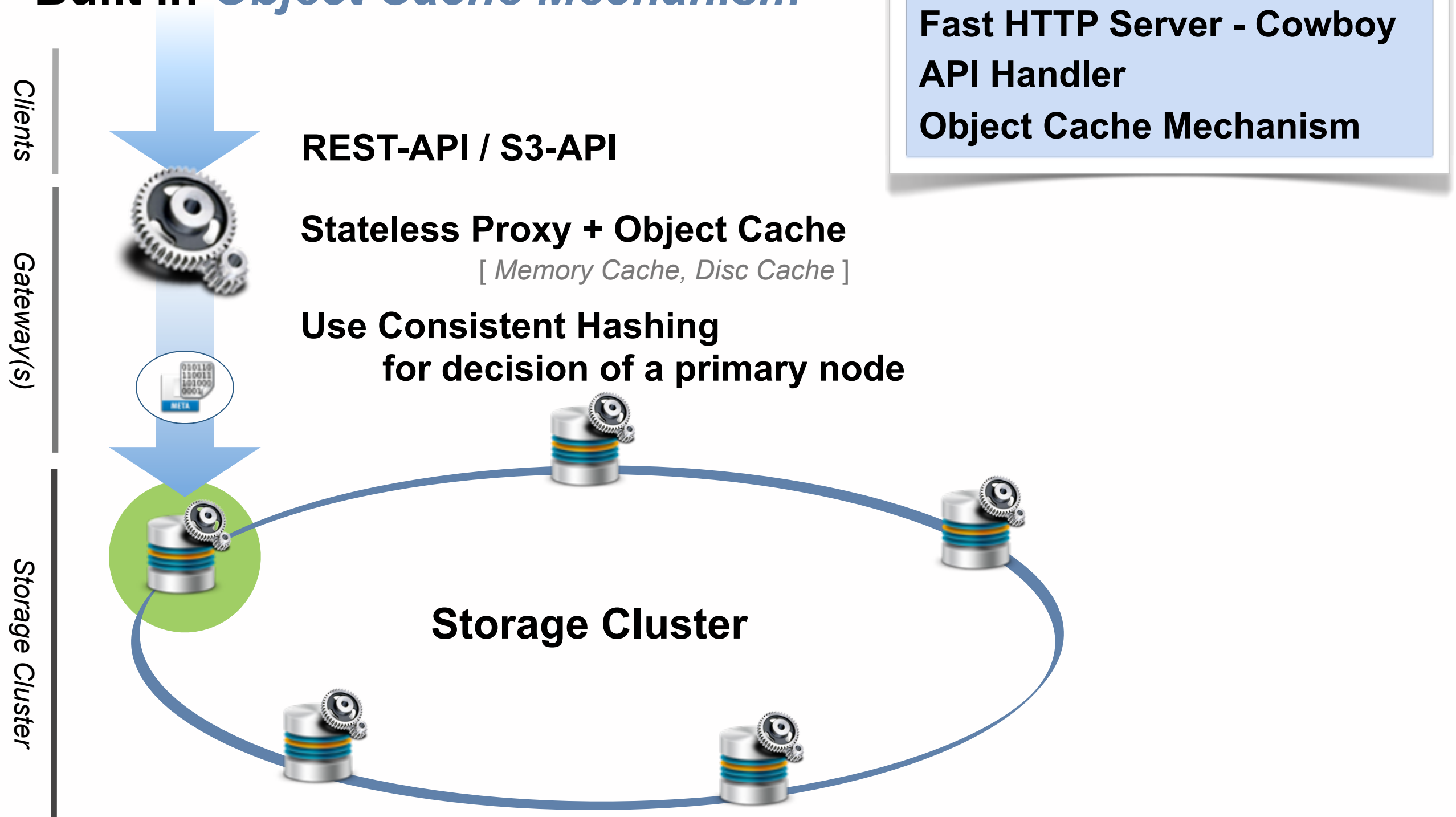Metadata  Object Storage

Metadata  Object Storage

Metadata  Object Storage

楽R天  R Rakuten

# LeoFS Gateway

# LeoFS Overview - Gateway

## HTTP Request and Response

## Built in *Object Cache Mechanism*

Fast HTTP Server - Cowboy
API Handler
Object Cache Mechanism

*Clients*

*Gateway(s)*

REST-API / S3-API

**Stateless Proxy + Object Cache**
*[ Memory Cache, Disc Cache ]*

**Use Consistent Hashing**
**for decision of a primary node**
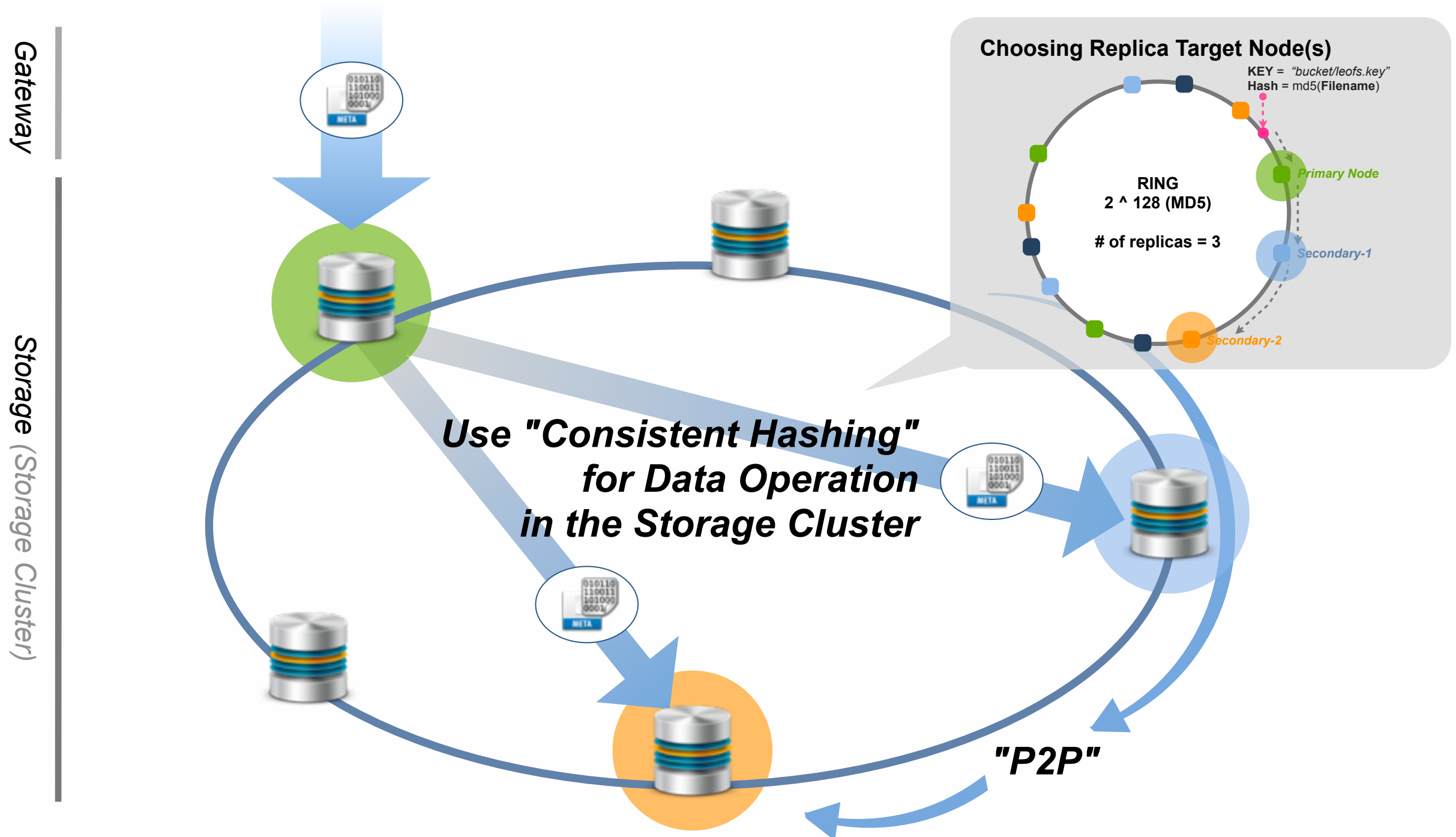
*Storage Cluster*

**Storage Cluster**

楽R天 ®Rakuten

# LeoFS Storage
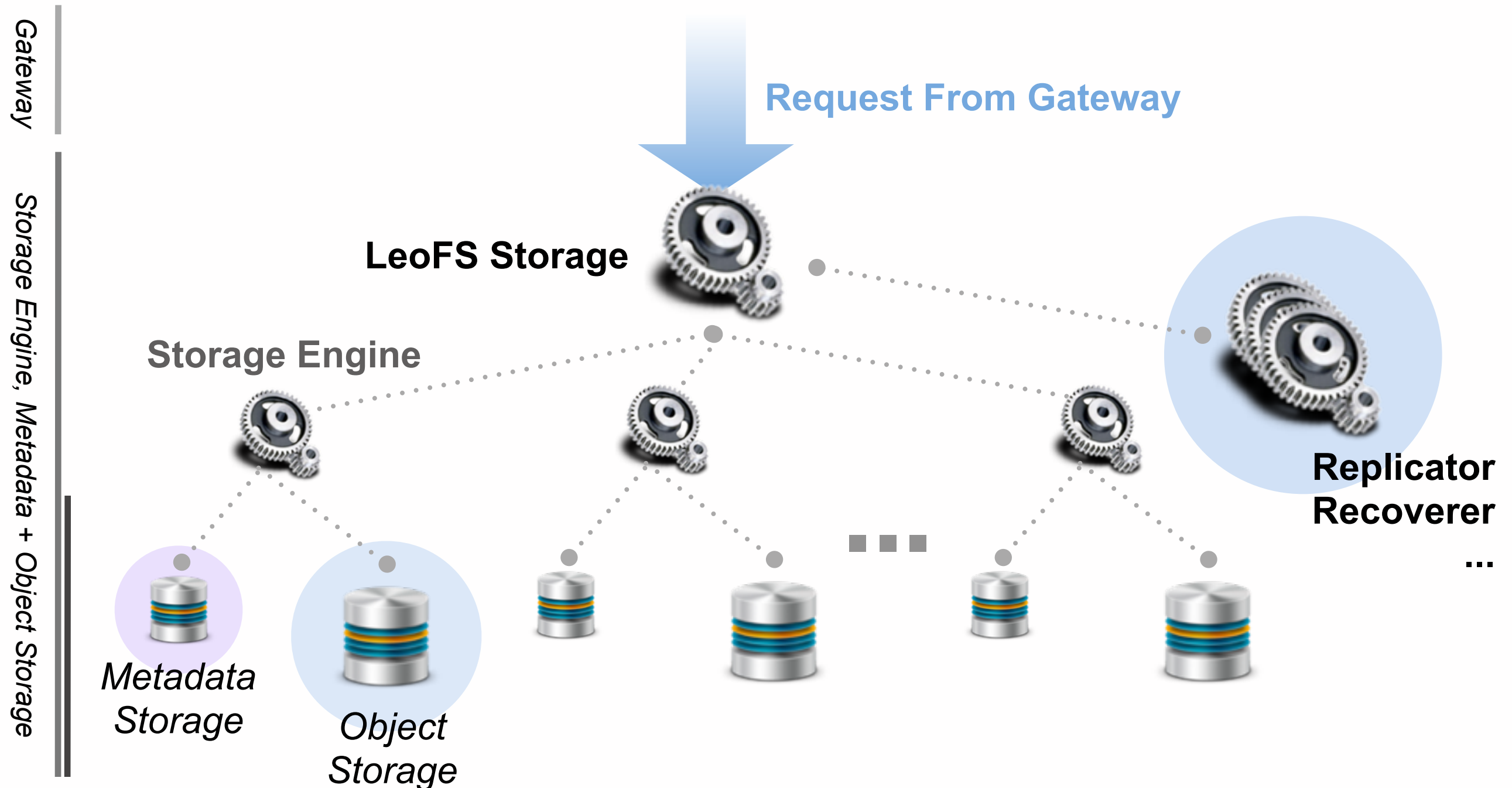
# LeoFS Overview - Storage

## WRITE: Auto Replication
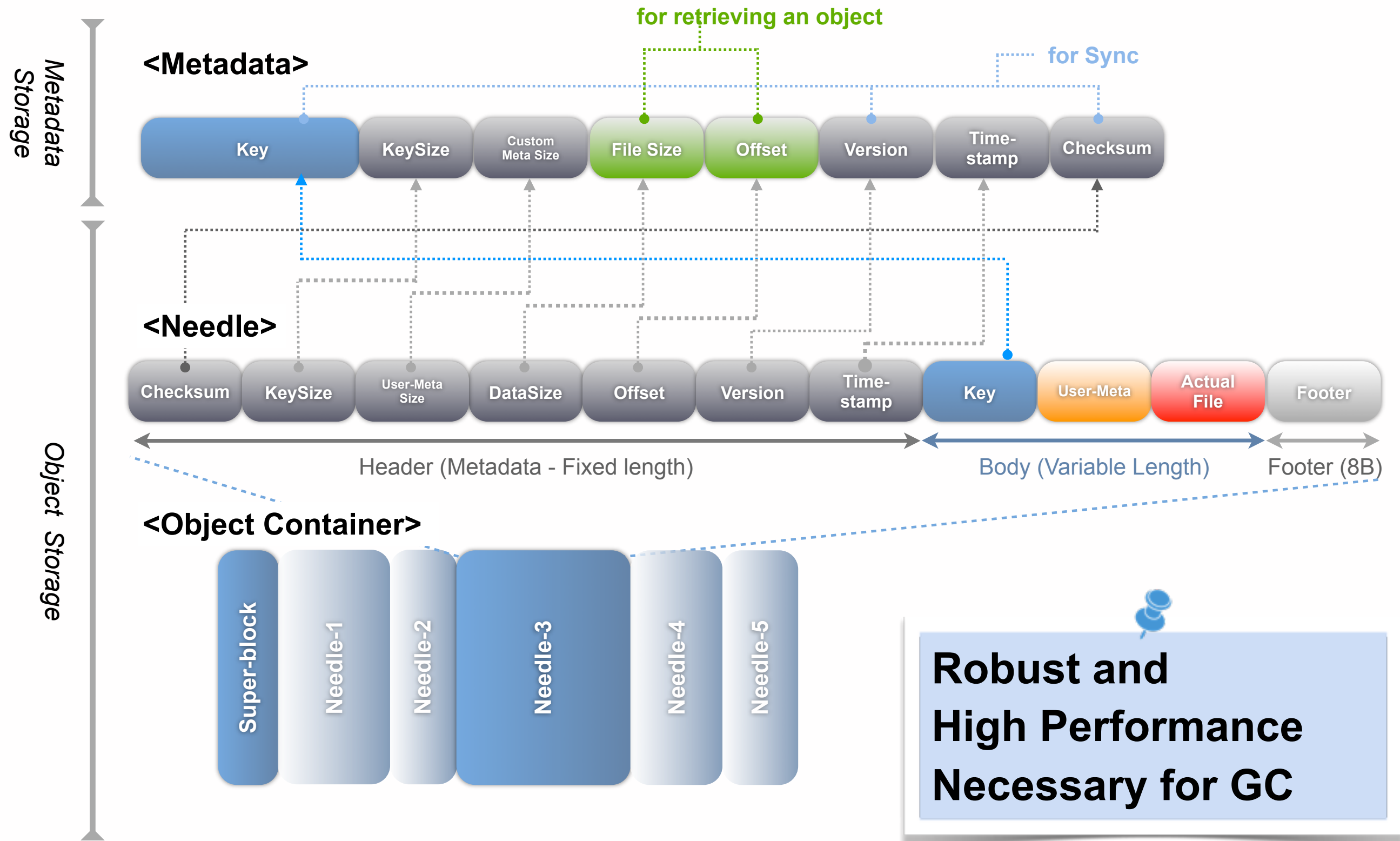## READ : Auto Repair of an Inconsistent Object with Async

# LeoFS Overview - Storage

**Storage consists of *Object Storage* and *Metadata Storage***

**Includes *Replicator* and *Recoverer* for the eventual consistency**



Gateway

Storage Engine, Metadata + Object Storage

**Request From Gateway**

**LeoFS Storage**

**Storage Engine**

**Replicator**
**Recoverer**

...

*Metadata Storage*

*Object Storage*

...
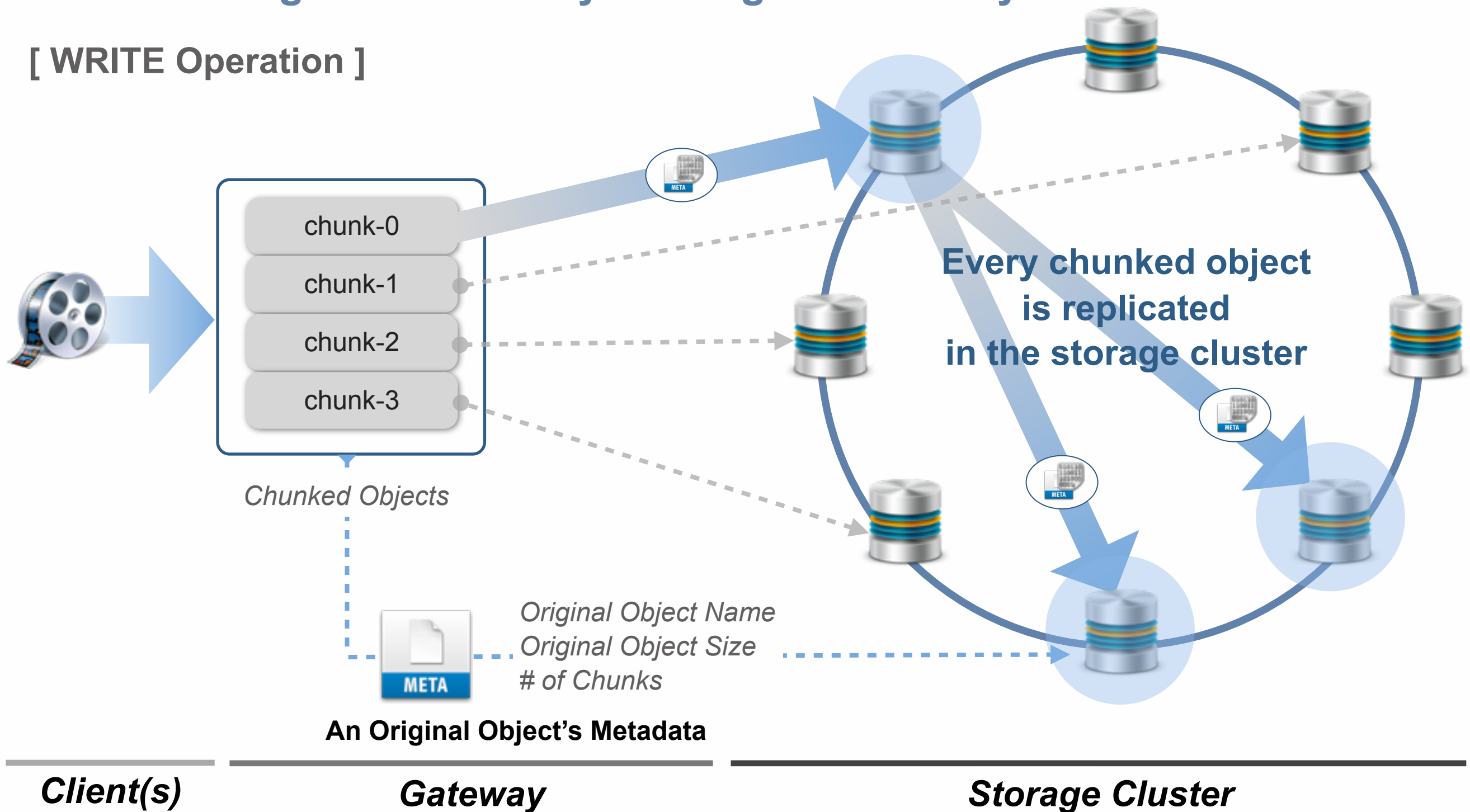
楽R天 ®Rakuten

# LeoFS Overview - Storage - Data Structure

# LeoFS Overview - Storage - Large Object Support

**To Equalize Disk Usage in Every Storage Node**
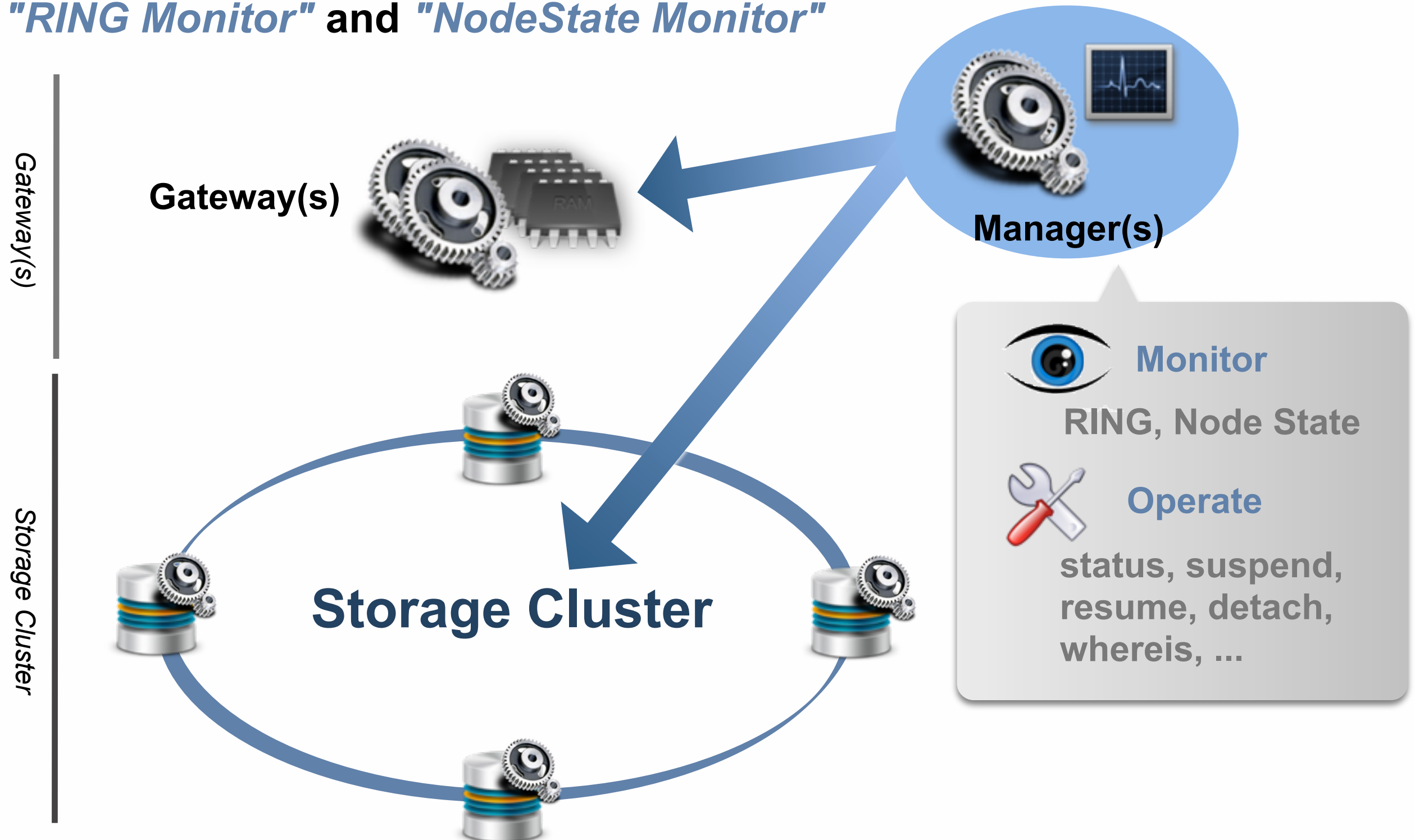**To Realise High I/O efficiency and High Availability**

[ WRITE Operation ]



chunk-0
chunk-1
chunk-2
chunk-3

*Chunked Objects*

**Every chunked object is replicated in the storage cluster**

*Original Object Name*
*Original Object Size*
*# of Chunks*

**META**

**An Original Object's Metadata**

*Client(s)*     *Gateway*     *Storage Cluster*

楽R天 Ⓡ Rakuten

# LeoFS Manager

# LeoFS Overview - Manager

**Operate LeoFS - Gateway and Storage Cluster**

*"RING Monitor"* **and** *"NodeState Monitor"*

Gateway(s)

**Gateway(s)**

**Manager(s)**

Storage Cluster

**Storage Cluster**

👁 **Monitor**

RING, Node State

🔧 **Operate**

status, suspend, resume, detach, whereis, ...

楽R天 🅡Rakuten

# 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**

楽❿天 ⓇRakuten

# Brief Benchmark Report

## 1st Case:

### Group of Value Ranges
*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
*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 (R:W=9:1)

**Environment:**

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

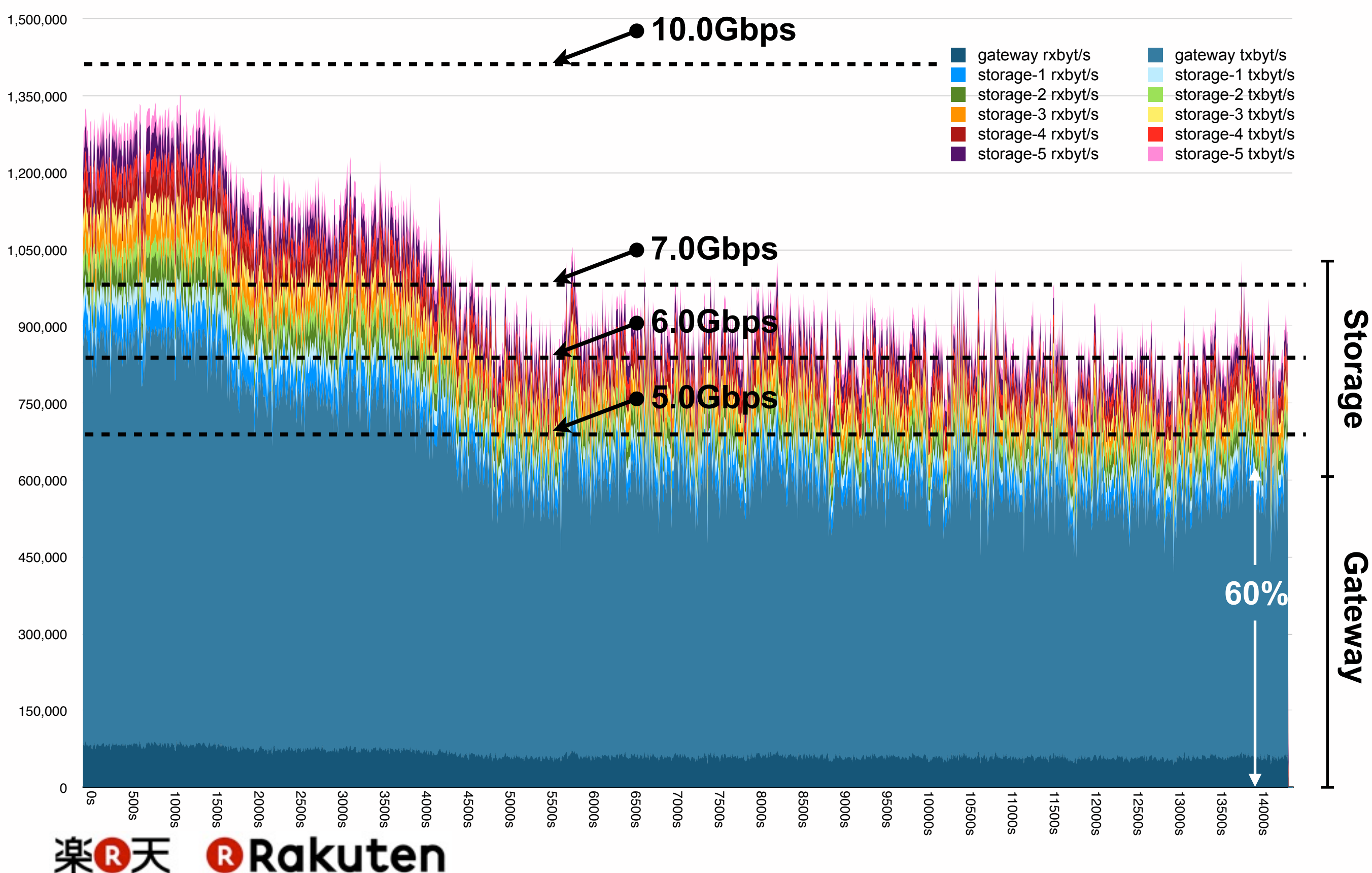**System Consistency Level:** [ N:3, W:2, R:1, D:2 ]

**Benchmark Configuration:**

| | |
|---|---|
| Duration | 4.0h |
| **R:W** | **9:1** |
| # of Concurrent Processes | 64 |
| # of Keys | 100,000 |
| Value Size | |

| Range (byte) | | Percentage |
|---|---|---|
| 1024 | 10240 | 24.00% |
| 10241 | 102400 | 30.00% |
| 10241 | 819200 | 30.00% |
| 819201 | 1572864 | 16.00% |

楽R天 ⓇRakuten

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



**1,500ops**

**No Errors**

**OPS**

**Latency**

**50ms**

**50ms**

# Brief Benchmark Report  - 1st Case / Network Traffic



10.0Gbps

7.0Gbps

6.0Gbps

5.0Gbps

Storage

Gateway

60%

Legend:
- gateway rxbyt/s
- gateway txbyt/s
- storage-1 rxbyt/s
- storage-1 txbyt/s
- storage-2 rxbyt/s
- storage-2 txbyt/s
- storage-3 rxbyt/s
- storage-3 txbyt/s
- storage-4 rxbyt/s
- storage-4 txbyt/s
- storage-5 rxbyt/s
- storage-5 txbyt/s

# Brief Benchmark Report - 1st Case / Memory and CPU



Memory Usage

gateway
storage-1
storage-2
storage-3
storage-4
storage-5

CPU Load 5min

1.0

楽R天 ®Rakuten

# Brief Benchmark Report  - 2nd Case (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 | **Range (byte)** | | **Percentage** |
| | 1024 | 10240 | 24.00% |
| | 10241 | 102400 | 30.00% |
| | 10241 | 819200 | 30.00% |
| | 819201 | 1572864 | 16.00% |

楽R天 Rakuten

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



**1,000ops**

**No Errors**

**OPS**

**60-70ms**

**80-90ms**

**Latency**

楽R天 R Rakuten

# Compare 1st case with 2nd case

# Brief Benchmark Report

## 1st Case - Network Traffic



**7.0Gbps**
**6.0Gbps**

Legend:
- gateway rxbyt/s
- storage-1 rxbyt/s
- storage-2 rxbyt/s
- storage-3 rxbyt/s
- storage-4 rxbyt/s
- storage-5 rxbyt/s
- gateway txbyt/s
- storage-1 txbyt/s
- storage-2 txbyt/s
- storage-3 txbyt/s
- storage-4 txbyt/s
- storage-5 txbyt/s

## minus 0.7Gbps

## 2nd Case - Network Traffic



**7.0Gbps**
**6.0Gbps**

楽天 Rakuten

# Brief Benchmark Report

## 1st Case - Disk util%



## 2nd Case - Disk util%



1.8x high

# 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**

楽R天 ®Rakuten

# Multi Data Center Replication

# Multi Data Center Replication

**HIGH-Scalability**
**HIGH-Availability** **+** **Easy Operation for Admins**

Europe
US
Tokyo
Singapore

**NO SPOF**

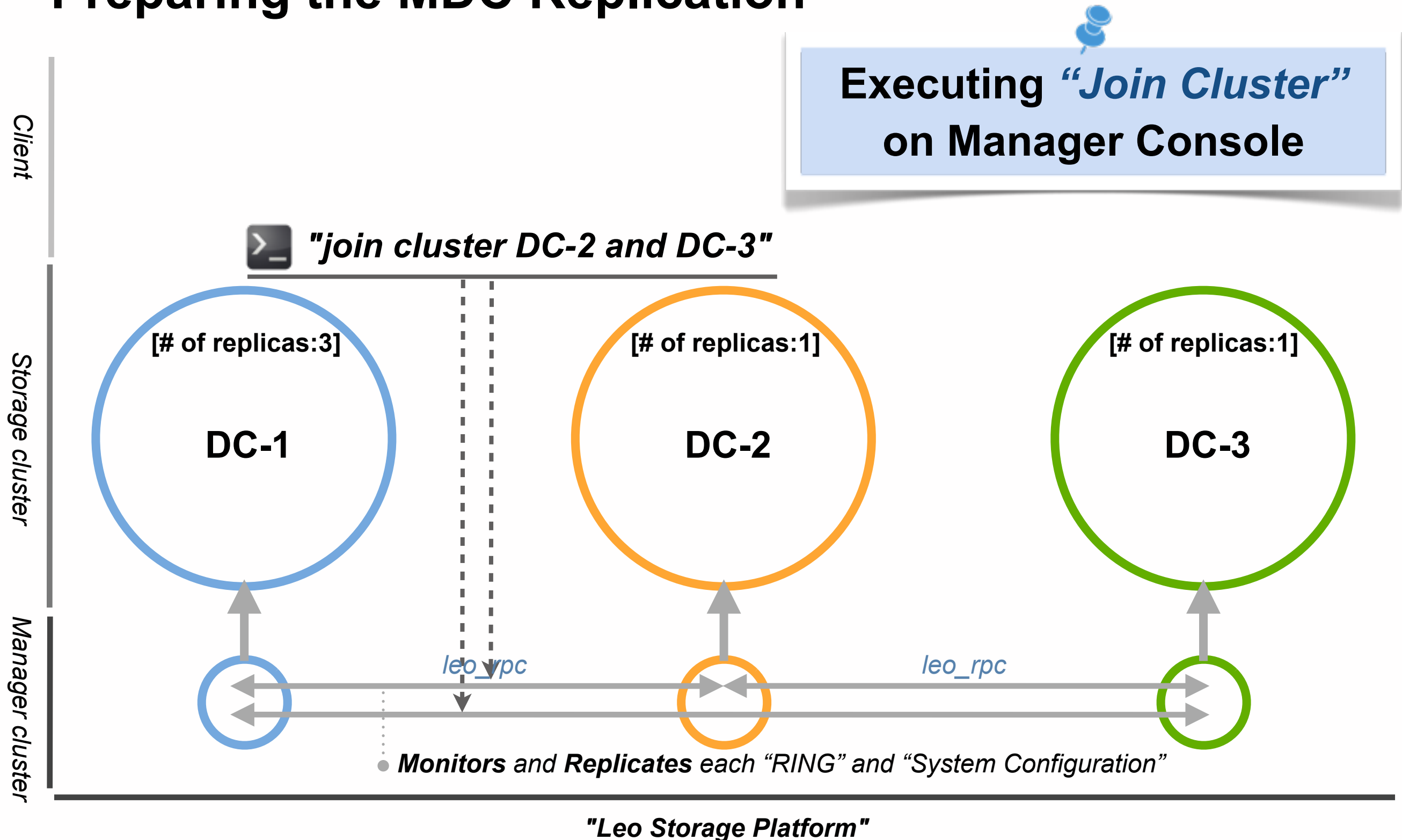**NO Performance Degradation**

楽R天 **R** Rakuten

# 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**

楽R天  ®Rakuten

# Multi Data Center Replication

## Preparing the MDC Replication

Executing *"Join Cluster"* on Manager Console

*Client*

*Storage cluster*

*Manager cluster*

*"join cluster DC-2 and DC-3"*

[# of replicas:3]

**DC-1**

[# of replicas:1]

**DC-2**

[# of replicas:1]

**DC-3**

*leo_rpc*

*leo_rpc*

● **Monitors** and **Replicates** each "RING" and "System Configuration"

*"Leo Storage Platform"*

楽天 Rakuten

# Multi Data Center Replication

## Stacking objects



**Client** | Application(s)

Request to
the Target Region

*Temporally Stacking objects*
 *- One container's capacity is* *32MB*
 *- When capacity is full,*
   *send it to remote cluster(s)*

*\* 32MB: default capacity - able to set optional value*

**Storage cluster**

[# of replicas:3]

**DC-1**

[# of replicas:1]

**DC-2**

[# of replicas:1]

**DC-3**

**Manager cluster**

*leo_rpc*

*leo_rpc*

● ***Monitors*** *and* ***Replicates*** *each "RING" and "System Configuration"*

*"Leo Storage Platform"*

楽R天  ®Rakuten

# Multi Data Center Replication

# Transferring stacked objects

Client

Storage cluster

Manager cluster

Application(s)

Request to
the Target Region

**Stacked an object with a metadata**

**Compress it with LZ4**

**Stacked objects**

leo_rpc

leo_rpc

DC-1

DC-2

DC-3

**Replicated an object**

leo_rpc

*Monitors* and *Replicates* each "RING" and "System Configuration"

*"Leo Storage Platform"*

Rakuten

# Multi Data Center Replication

## Investigating stored objects



Application(s)

Request to
the Target Region

1) Receive metadata of stored objects
2) Compare them at the local cluster
3) Fix inconsistent objects

Client

Storage cluster

Manager cluster

DC-1

DC-2

DC-3

leo_rpc

leo_rpc

leo_rpc

leo_rpc

*Monitor* and *Replicate* each "RING" and "System Configuration"

"Leo Storage Platform"

楽R天  R Rakuten

# NFS Support

# Future Plans

# NFS Support

## Data-HUB: Centralize unstructured data in LeoFS

# LeoFS Administration at Rakuten

*Presented by Masahiro Sanjo*

*Rakuten Institute of Technology*

楽R天　Ⓡ Rakuten

**Storage Platform**

**File Sharing Service**

**Others**

*Portal Site*

*Photo Storage*

*Background Storage of OpenStack*

楽R天 ®Rakuten

# 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/600TB

# of Files: 600Million

Daily Growth: 100GB

Daily Reqs: 13Million

E-Commerce

Photo share

Review

Portal & Contents
*(Movie)*

Blog

Storage Platform

Recruiting

Bookmark

Insurance

Calendar

楽R天 ®Rakuten

# Storage Platform - System Layout

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

*Load Balancer / Cache Servers*

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

**Gateway x 4**

**Manager x 2**

**Nagios**
**Ganglia**

*( Erlang RPC)*

*( Erlang RPC)*

*( TCP/IP,SNMP )*

**Monitor**

**Storage x 14**

**GUI Console**

楽R天  ®Rakuten

# Storage Platform - Monitor

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

**Ganglia Agent**

**Send Mail Alert**

Gateway x 4

Storage x 14

Manager x 2

*( Erlang RPC)*

*( Erlang RPC)*

*( TCP/IP,SNMP )*

Nagios
Ganglia

Monitor

GUI Console

楽R天  Rakuten

# Covering All Services
# with Multi DC Replication

# File Sharing Service



**+**



*https://owncloud.com/*

# File Sharing Service - Required Targets



**Reduce Costs**
**Handle Confidential Files**
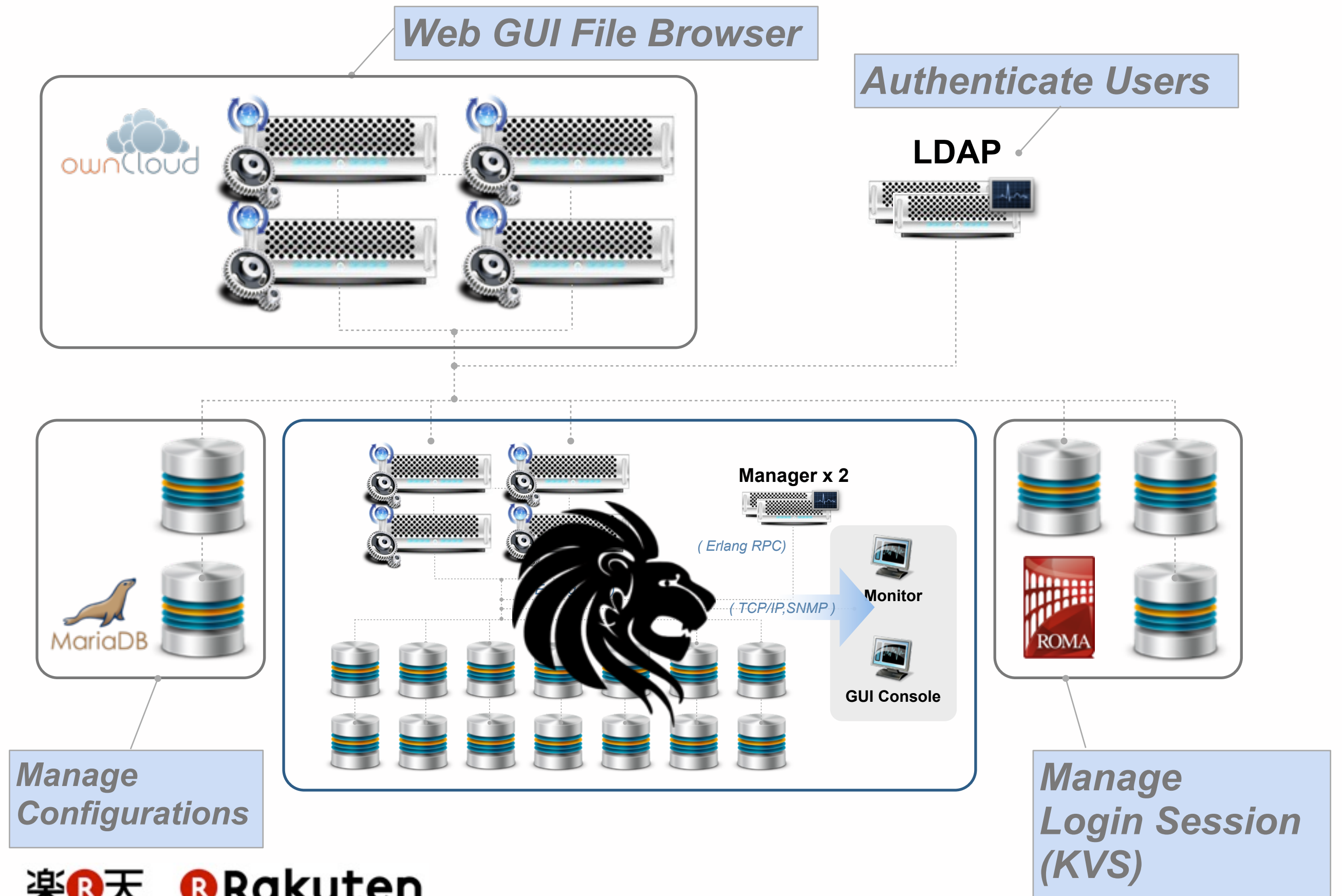**Store Large Files**
**Scale Easily**

# File Sharing Service - Usage

ownCloud

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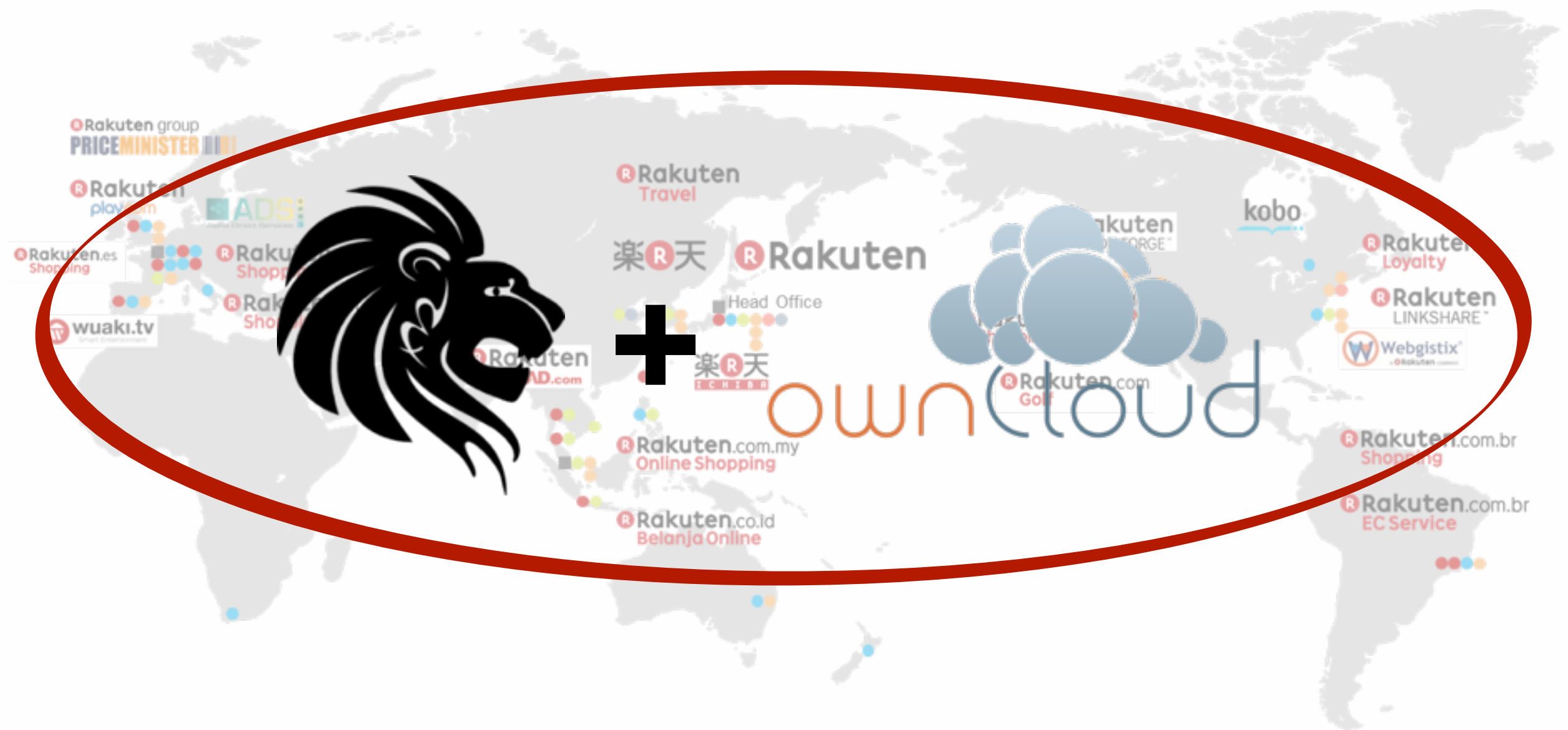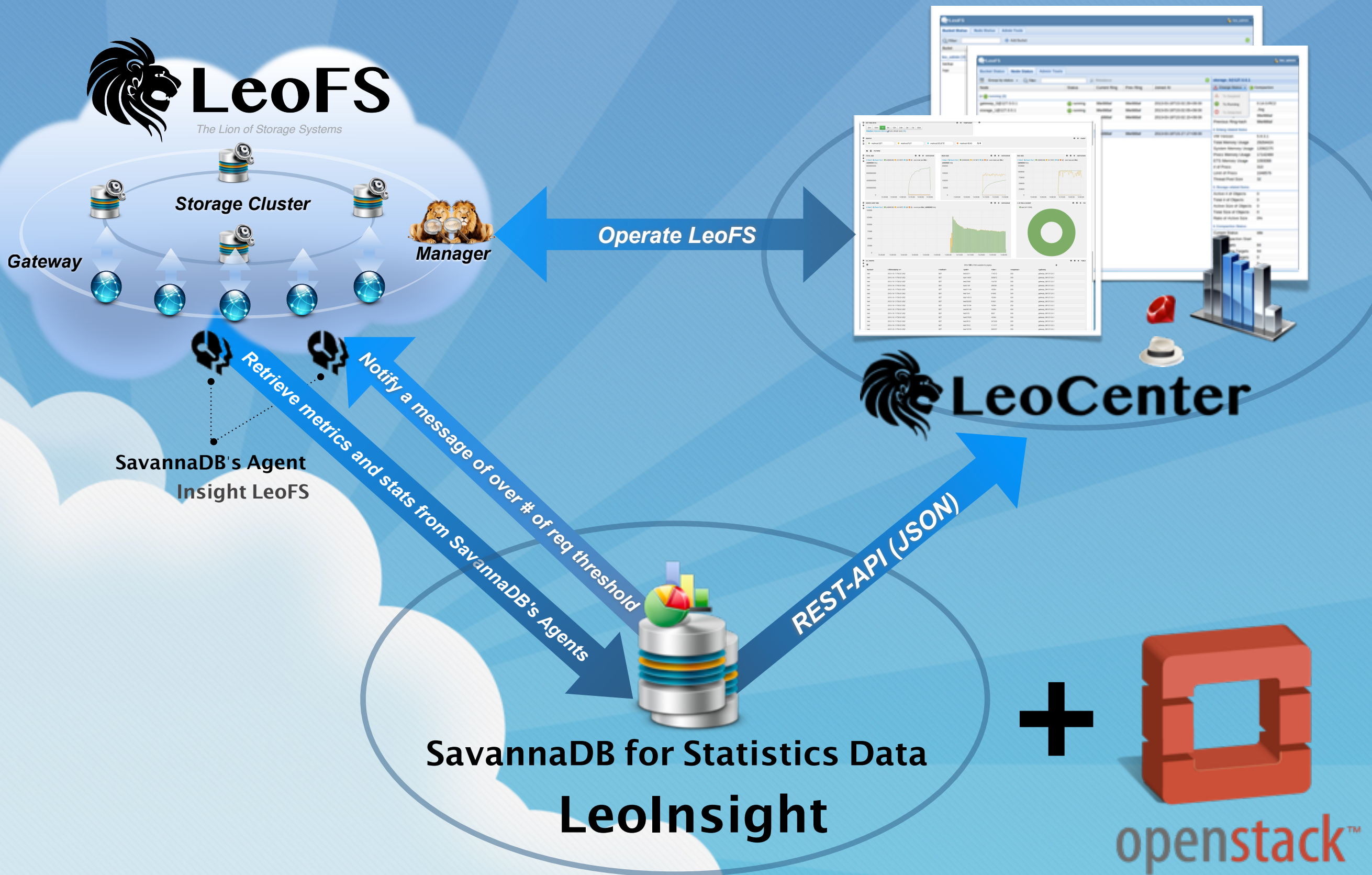Over 4,000 Users, Over 10,000 Teams**

楽R天 ❘ Ⓡ Rakuten

# Empowering the Services and the Users Through the Cloud Storage

# Future Plans

# Set Sail for "Cloud Storage"

**Website: leo-project.net**
**Twitter: @LeoFastStorage**