# Large scale data near-line loading method and architecture

**FiberHome Telecommunication**
**2017-7-19**

FiberHome

# /usr/bin/whoami

**Shuaifeng Zhou(周帅锋):**

* Big data research and development director ( Fiberhome 2013-)
* Software engineer (Huawei 2007-2013)
* Use and contribute to HBase since 2009
* sfzhou1791@fiberhome.com

# Contents

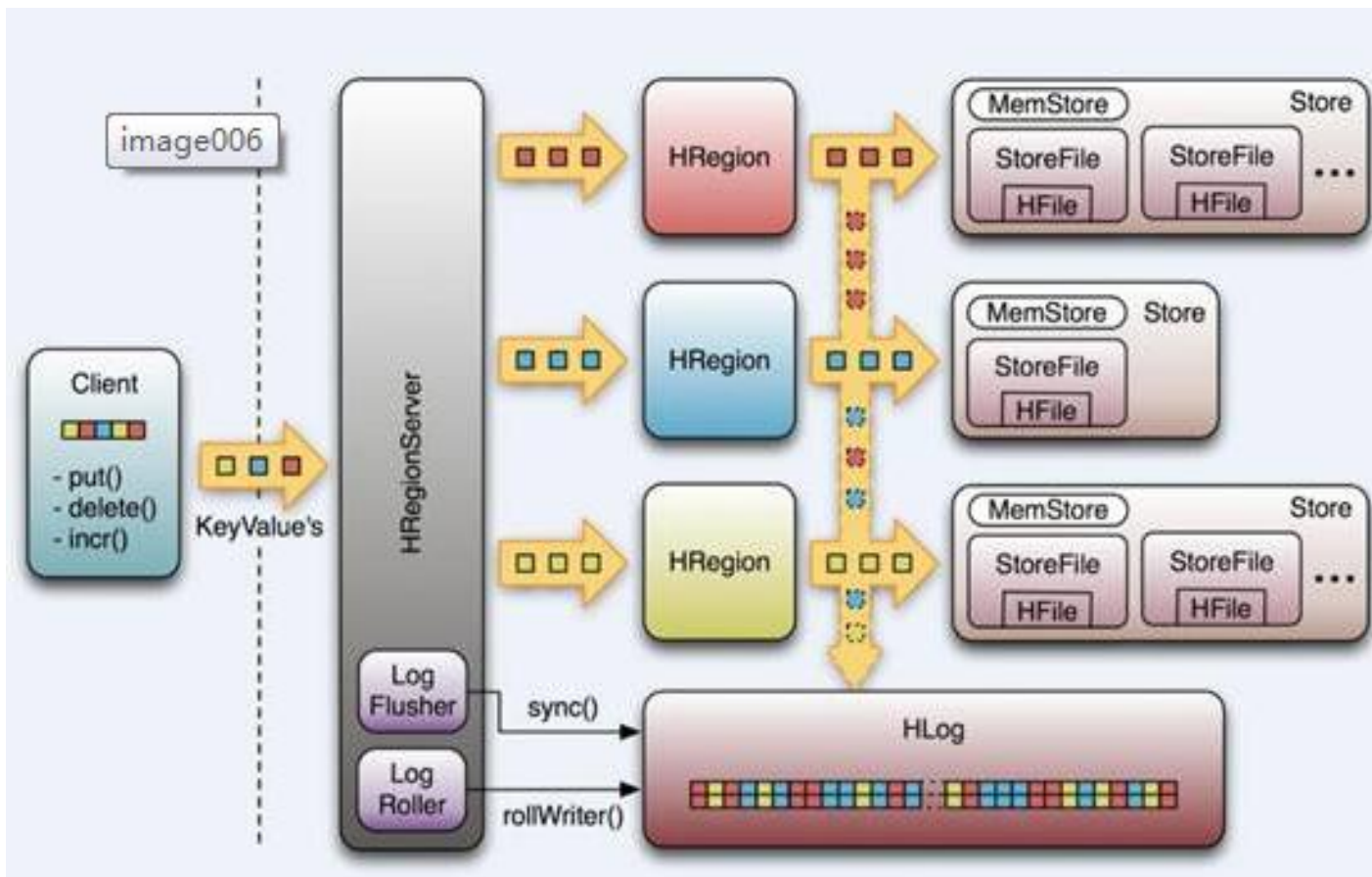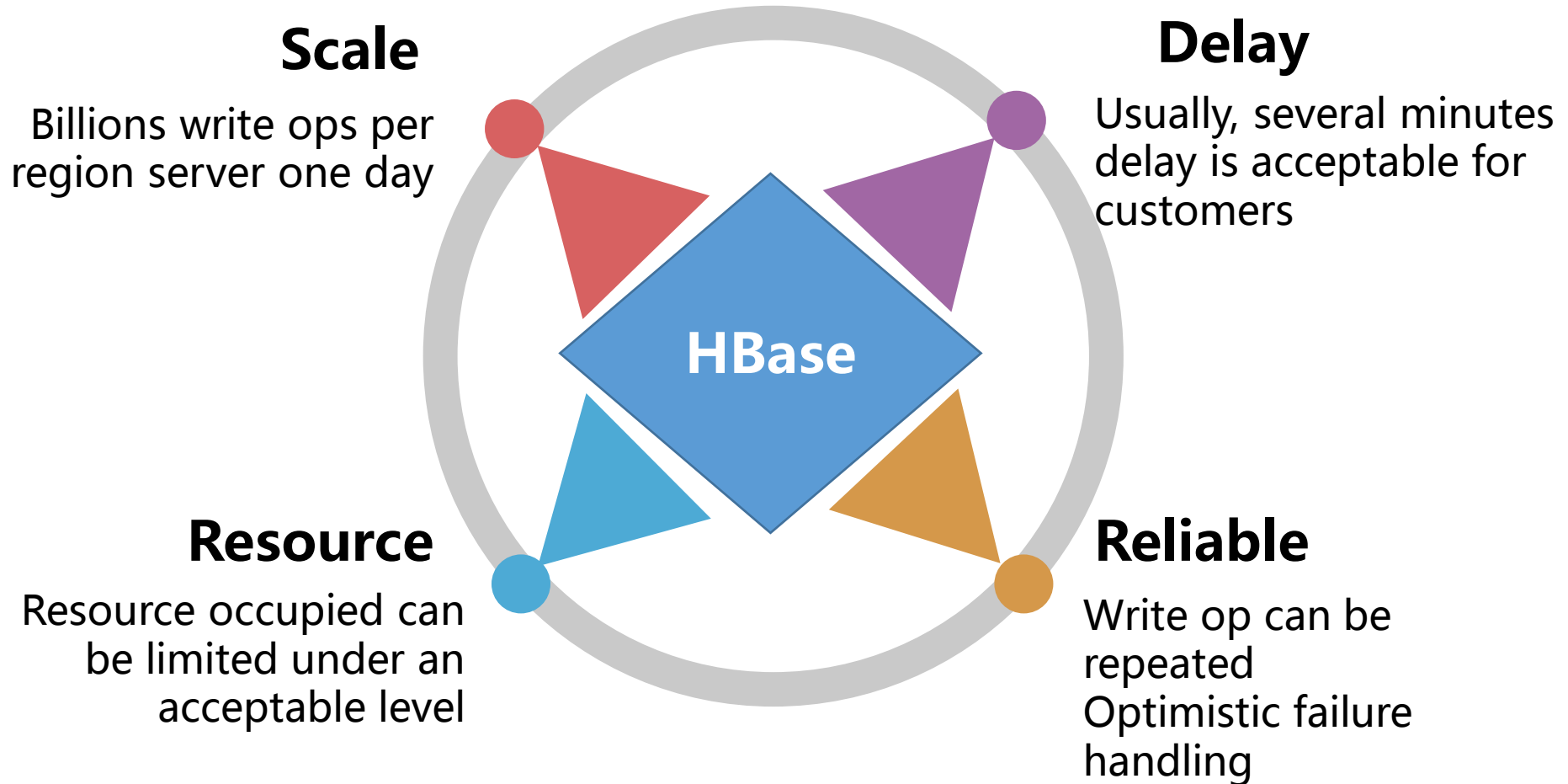# HBase Realtime Data Loading



- WAL/Flush/Compact Triple IO pressure
- Read/Write operations share resource:
  - Cpu
  - Network
  - Disk IO
  - Handler
- Read performance decrease too much when write load is heavy

# Why near-line data loading?



**Scale**
Billions write ops per region server one day

**Delay**
Usually, several minutes delay is acceptable for customers

**Resource**
Resource occupied can be limited under an acceptable level

**Reliable**
Write op can be repeated
Optimistic failure handling

**HBase**

**Large scale data loading reliably with acceptable time delay and resource occupation**
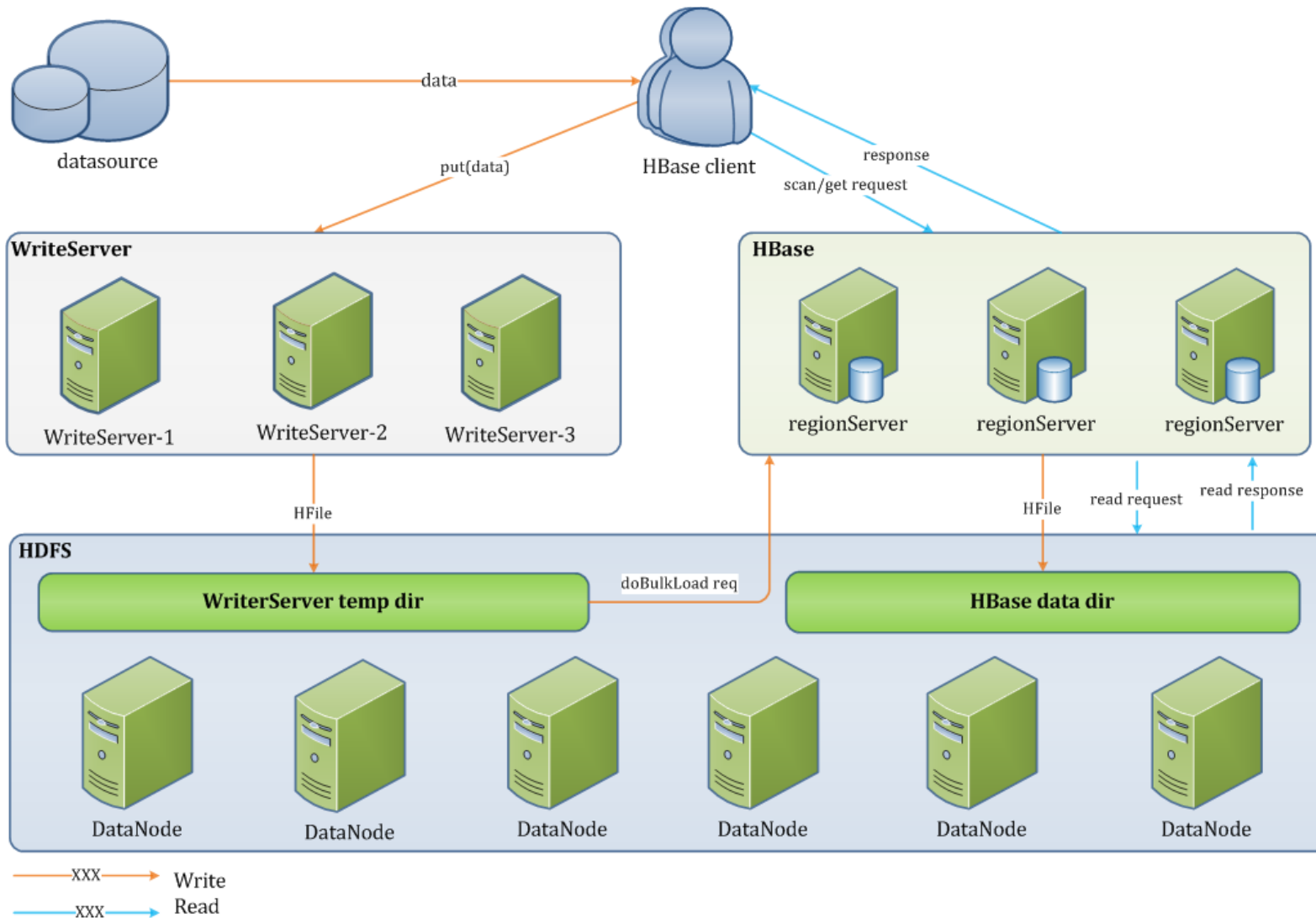
# Contents
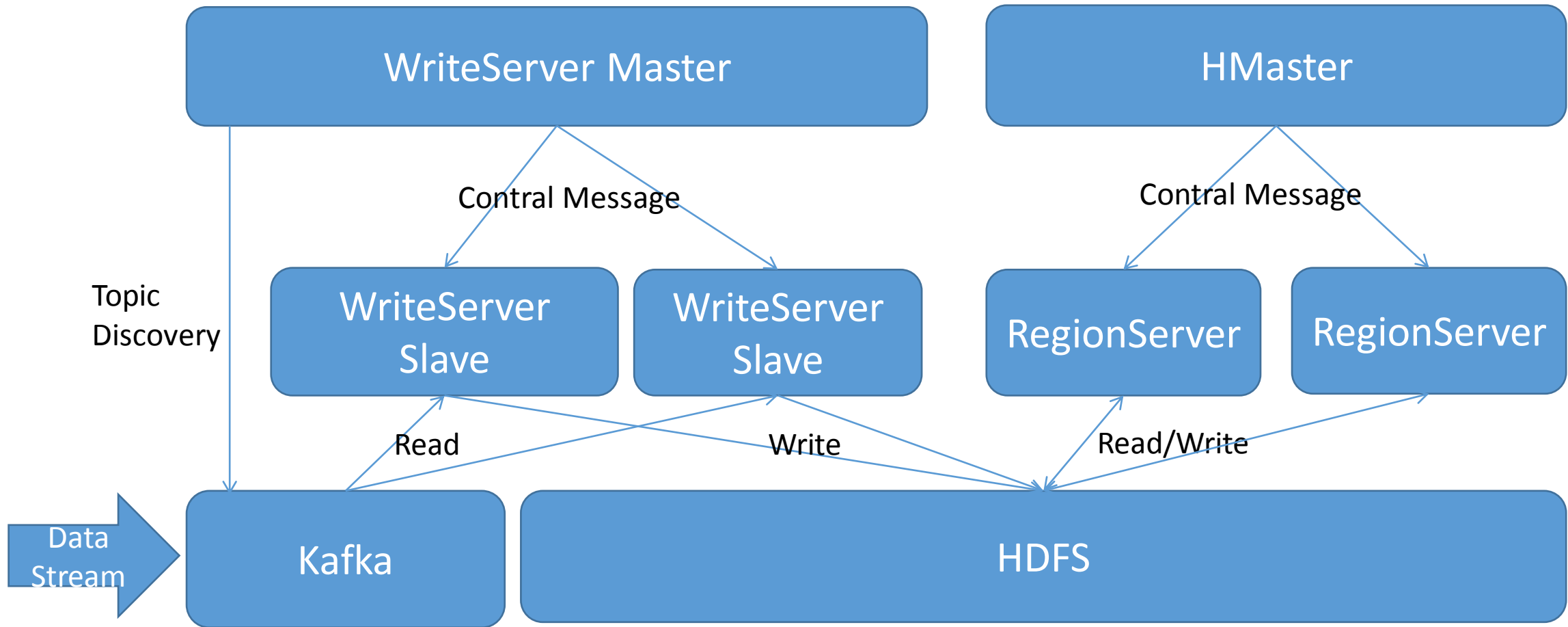
# Read-Write split data loading



- Independent WriterServer to handle put request
- RegionServer only handle read request
- WriteServer write HFile on HDFS, send do-bulkload operation.
- Several minutes delay between put and data readable.

# Architecture

# WriteServer Master

## Topic Management

- Discover new kafka topics
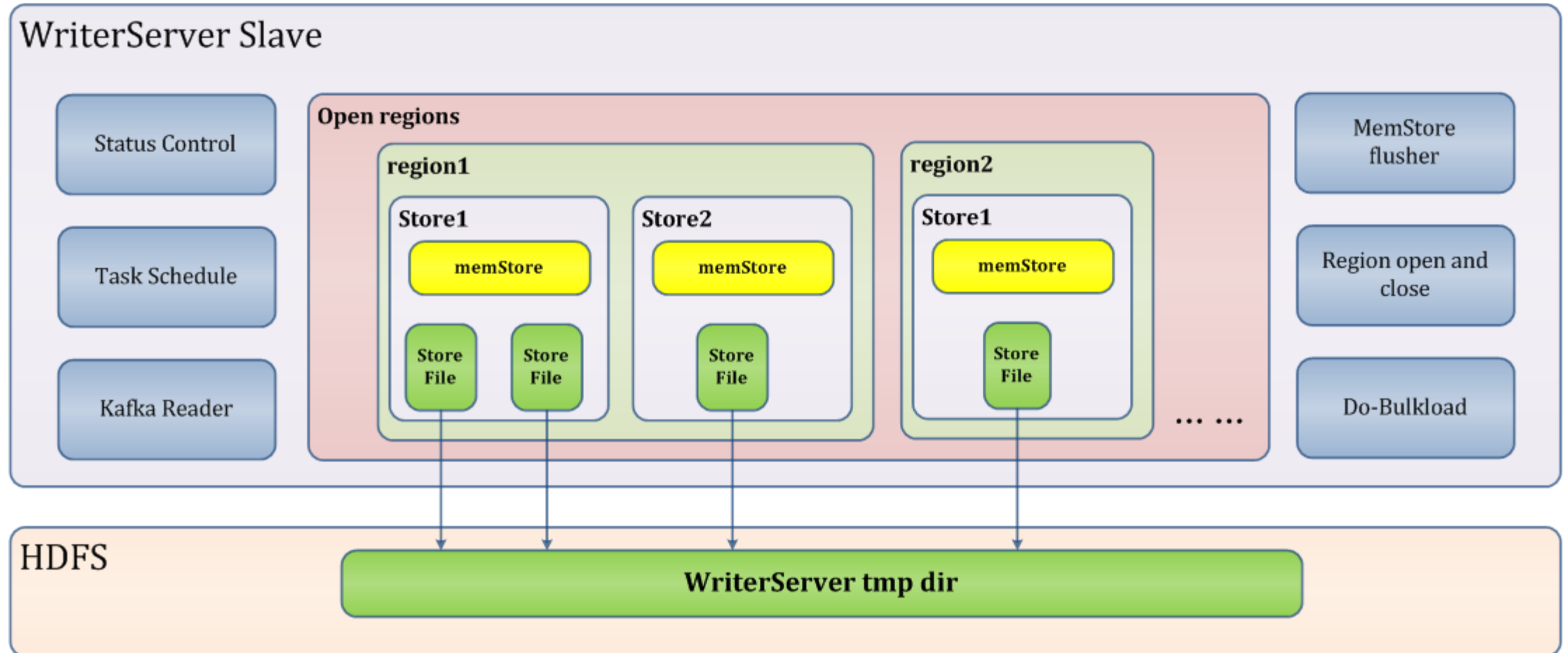- Receive loading request
- Loading records statistic

## Task Management

- Create new loading tasks every five minutes or every 10,000 records
- Find a slave to load the task
- Task status control

## Slaves Management

- Slave status report to master
- Balance
- failover

# WriteServer Slave

# Failure Handling

**Meta Data based Failure Handling**

**Task Meta Data** is the descricption info of a task, include the topic, partitions, start and end offset, status. Stored on disk.

**Task Meta Data** is constructed when a task is created by master, and change status to succeed when slave finish the task.

**Recover:** Redo failed tasks when slave down or master restart.
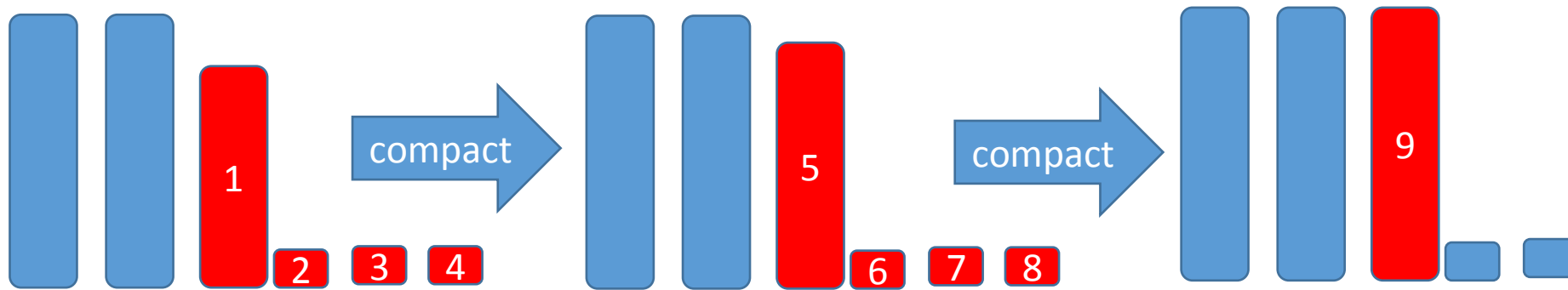
# Contents

# Balance

**Load balance according tasks:**

- Send new tasks to slaves with less tasks on handling

- Try to send tasks of one topic to a few fixed slaves
  - avoid one region open everywhere
  - Less region open, less small files

- Keep region opened for a while, even there are no tasks
  - avoid region open/close too frequently

# Compact

- Small files with higher priority

- Avoid one large file together with many small files compact again and again
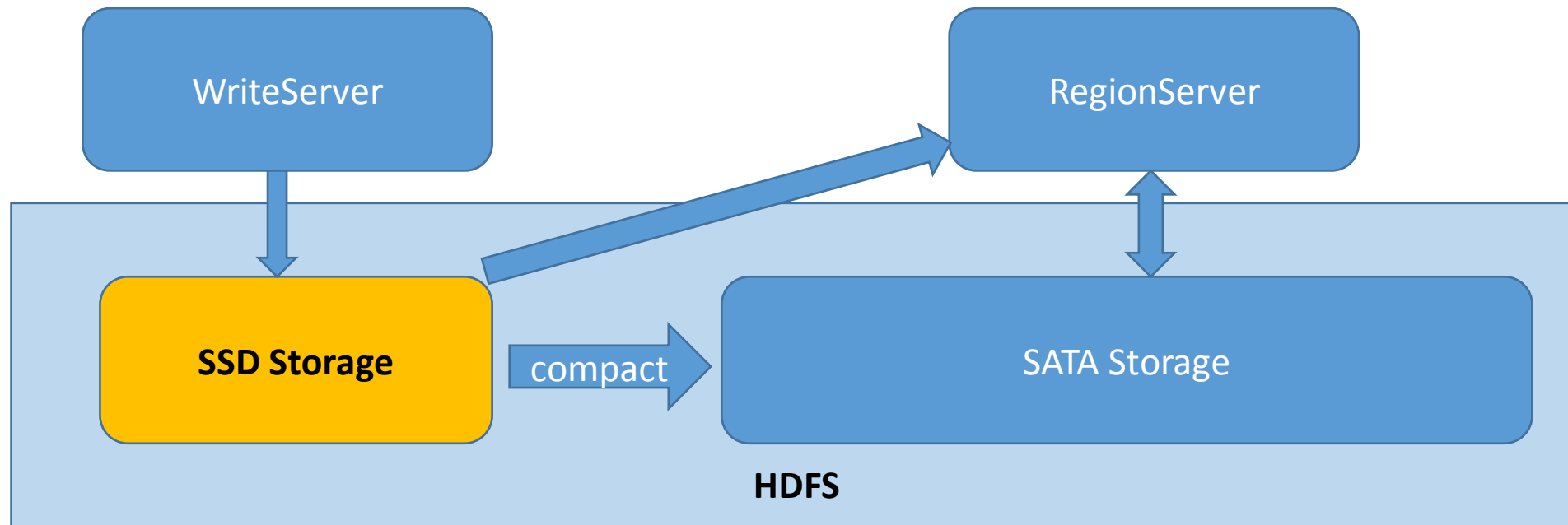
# StoreEngine

## Customized store engine:

- organize store files in two queues
  - one can be read and compact
  - The other can only be compact
  - If there are too many files, new file will not be readable until they are compact
- **Some new files discovered later better than all files can not be read before time out**
  - Occasionally data explosion can be handled
  - Region need split
  - "Hot key" should be handled
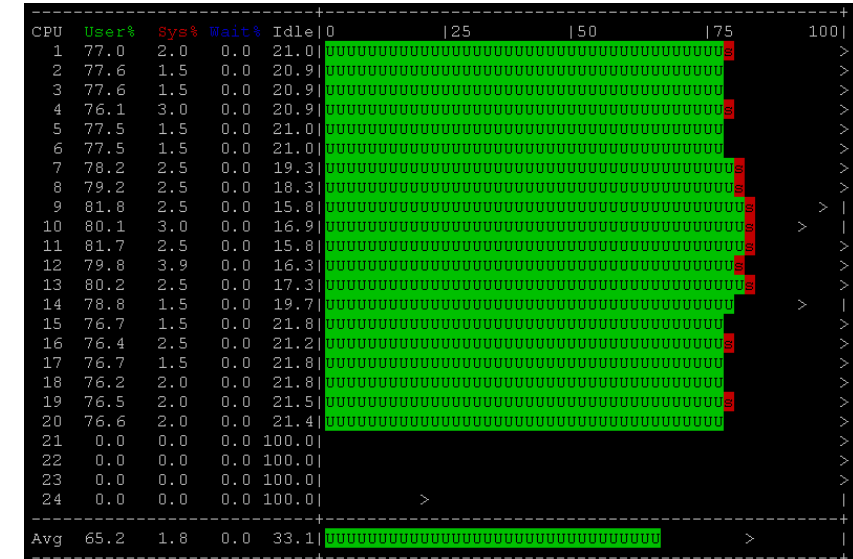
# HDFS Heterogeneous Storage Usage

- **Use SSD storage as WriteServer tmp dir**
- **Use SATA as HBase data dir storage**
  - WriteServer write HFile on SSD
  - Load HFile to HBase(Only move)
  - Change to SATA storage after compact by regionServer

# Resource Control

## Resource used by WriteServer should be controllable:



- Memory:
  - JVM parameters 30~50GB memory
  - Large Memory Store will avoid small files
  - Too Large memory store will cause gc problems

- **CPU:**
  - Slave can use 80% cpu cores at most
  - Compare to real-time data load, a big optimize is we can control the cpu occupation by write operations.
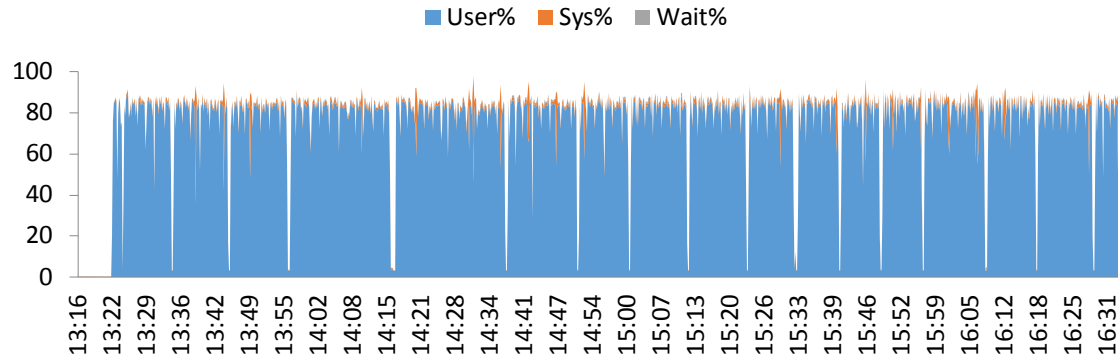
# Contents

# Loading Performance

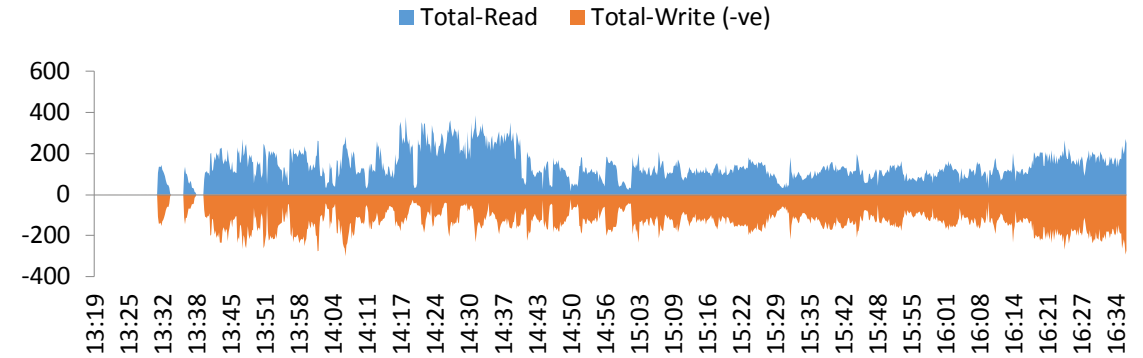| WriterServer Slave | |
|---|---|
| **CPU** | Intel(R) Xeon(R) CPU E5-2640 v2 @ 2.00GHz |
| **Memory** | 128G |
| **Disk** | 1TB SSD * 4 |
| **Network** | 10GE |
| Record size | 1KB |
| Compress | Snappy |
| **Performance** | **300,000 records/s** |

**One WriteServer slave can match 5 RegionServer's loading requirements before RegionServer reach compact limitation.**
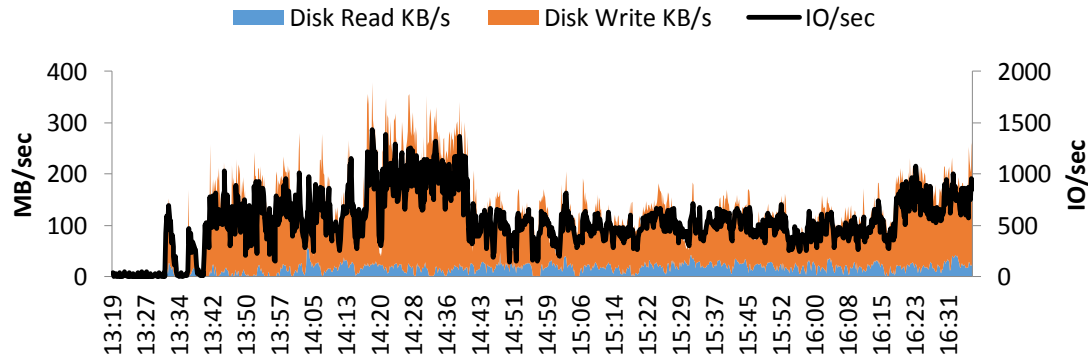
# Resource Performance

**CPU Total WS-Slave5 – 2017/2/17**



**Network I/O WS-Slave5 (MB/s) - 2017/2/17**



**Disk total MB/s WS-Slave5- 2017/2/17**



**Memory**

JVM: aways use memory as much as assigned

GC: config gc policy to avoid full gc.

# Contents

1 Motivation

2 Solution

3 Optimization

4 Tests

5 Summarize

# Summarize

## We proposed an read-write split near-line loading method and architecture:

- Increase loading performance

- Control resource used by write operation, make sure read operation can not be starved

- Provide an architecture  corresponding with kafka and hdfs

- Provide some optimize method, eg: compact, balance, etc.

- Provide test result