

# 微博REDIS定制化之Tribe系统介绍

微博研发中心 张冬洪

2016.11

- ❖ 个人介绍
- ❖ 微博REDIS的使用介绍
- ❖ Tribe系统设计的考量
- ❖ Tribe系统性能测试对比
- ❖ Tribe系统运维点滴

- ❖ Redis中国用户组发起人
- ❖ 微博研发中心DBA，专注于MySQL和NoSQL架构设计与运维以及自动化平台的开发；2016年3月加入微博，目前在微博主要负责Feed系统相关业务的数据库运维和业务保障工作



张冬洪  Lv14

419  
关注

154  
粉丝

86  
微博



Redis2016  Lv4

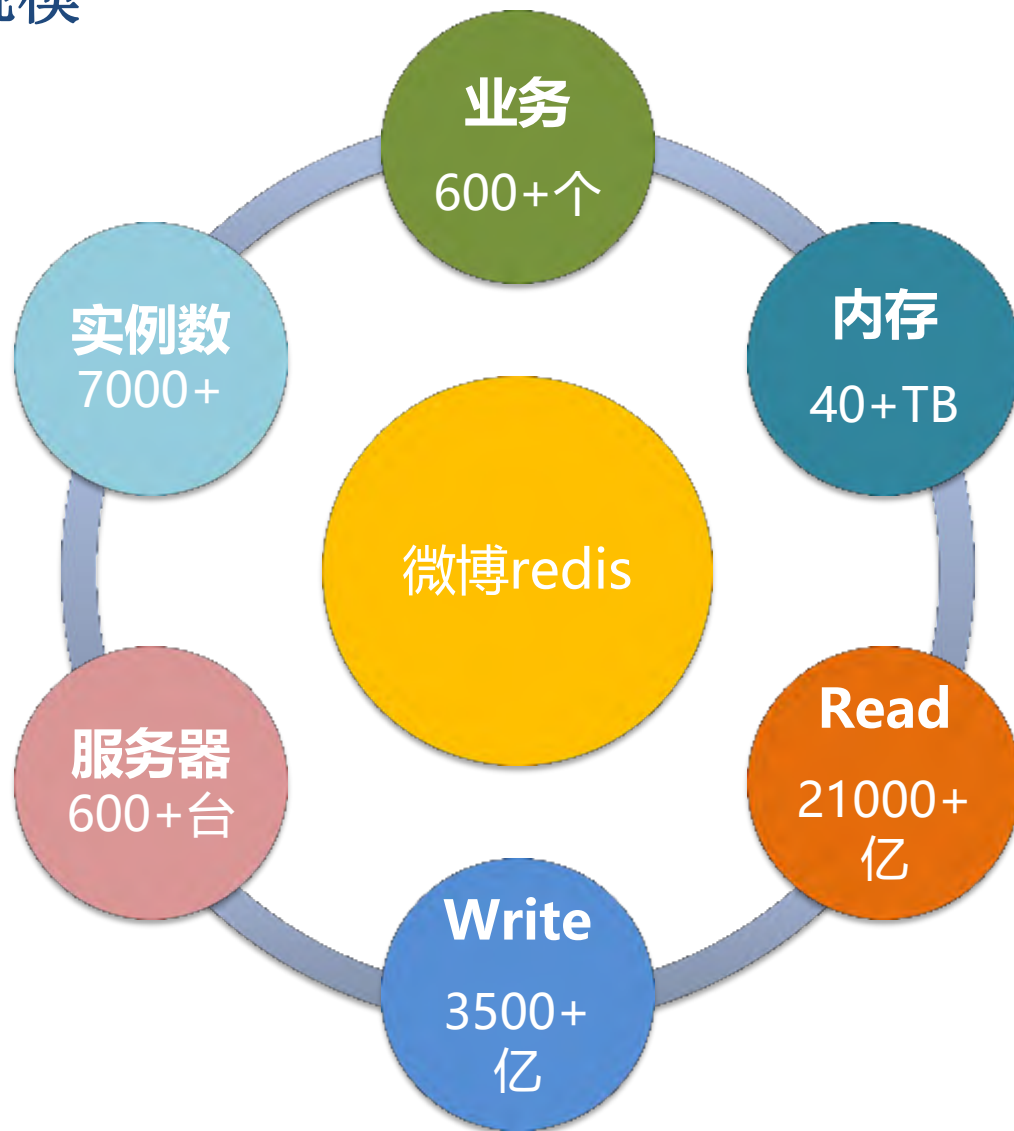
3  
关注

435  
粉丝

15  
微博

- ❖ 个人介绍
- ❖ 微博REDIS的使用介绍
- ❖ Tribe系统设计的考量
- ❖ Tribe系统性能测试对比
- ❖ Tribe系统运维点滴

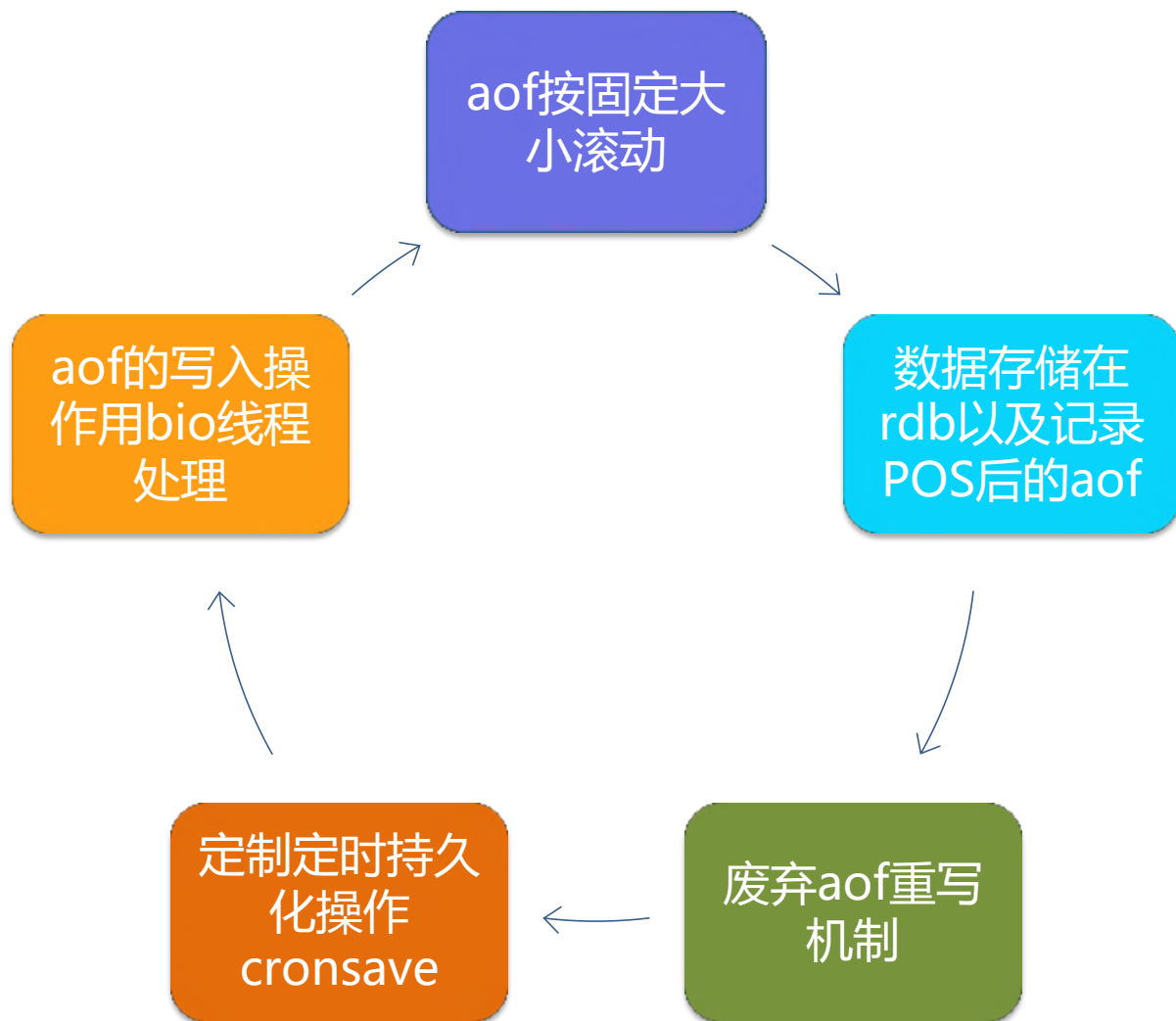
## ❖ 使用规模



## ❖ 定制化演进

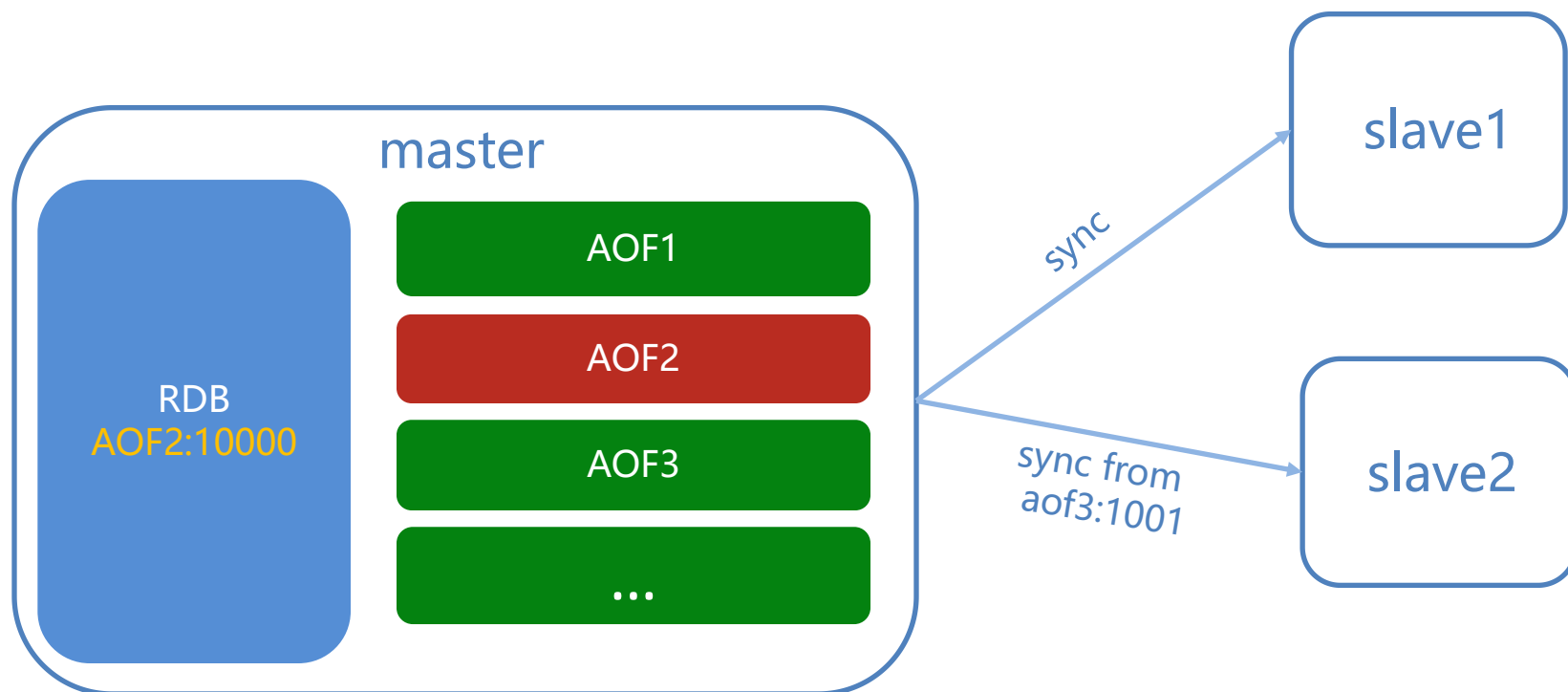


## ❖ 无阻塞落地、数据持久化



## ❖ 主从增量复制的实现

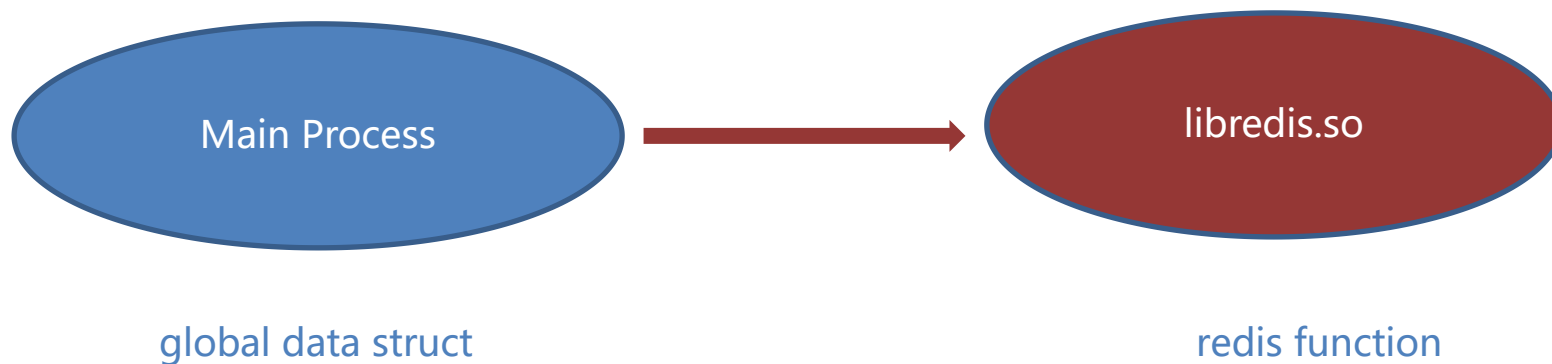
- 使用rdb+aof+aof position方式实现
- 主从只需进行一次全量复制，其余增量，如：主从断开、切主等



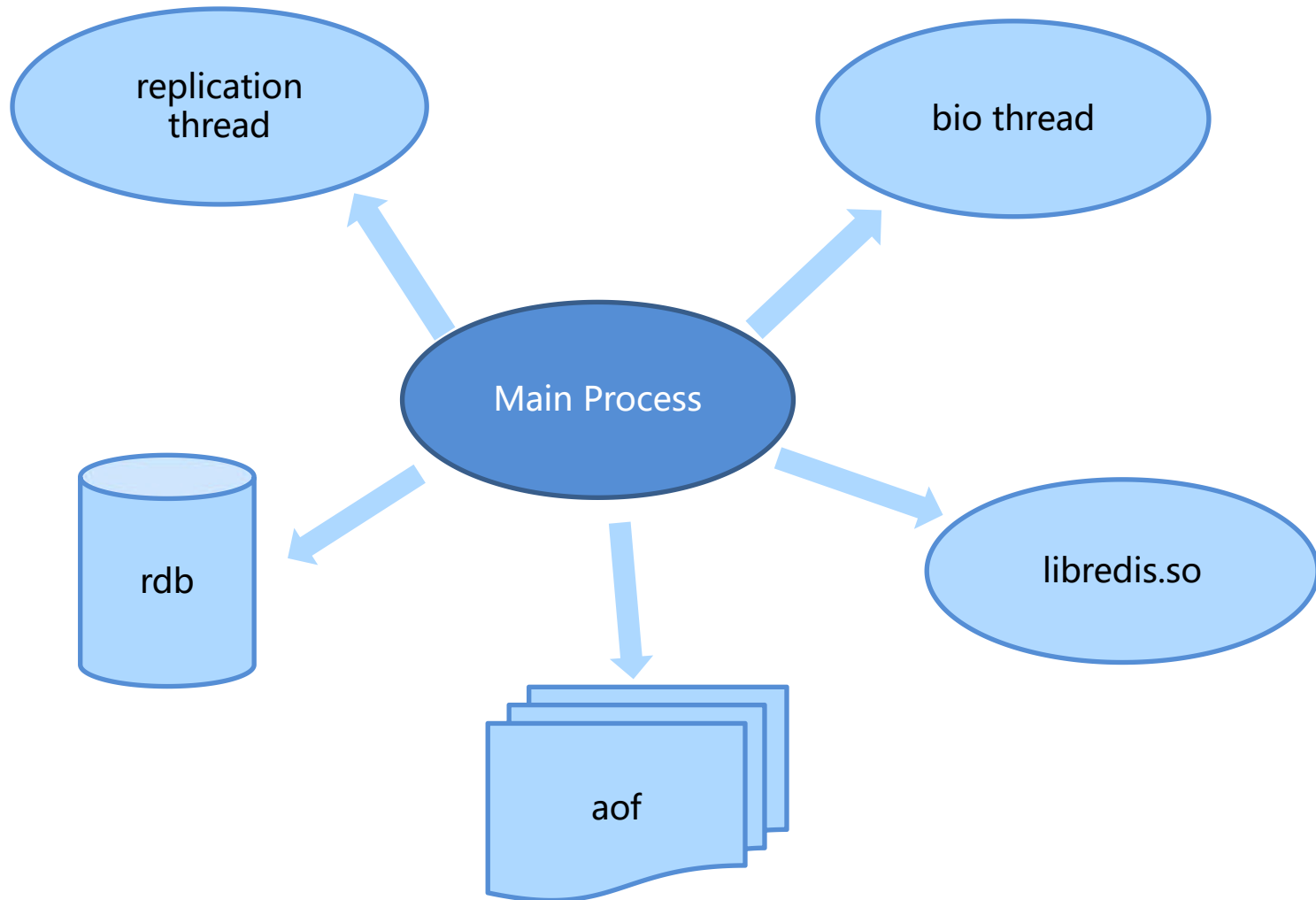


## ❖ 在线热升级

- 主程序存全局数据结构；内存中的数据保存在全局变量里
- 其他核心处理逻辑封装到动态库
- 外部程序通过调用动态库里相应函数来读写数据
- 升级只需替换动态库，无须重新载入数据
- 毫秒级完成升级，无业务影响



## ❖ 定制化版本总结



- ❖ 个人介绍
- ❖ 微博REDIS的使用介绍
- ❖ Tribe系统设计的考量
- ❖ Tribe系统性能测试对比
- ❖ Tribe系统运维点滴

## ❖ 当前结构及问题

业务端访问逻辑较重，每个业务需要实现一套



- 1. 数据迁移通过脚本完成，繁琐
- 2. 扩容/缩容（拆分）需要业务端更改配置，高耦合

健康检查，流量切换需要人工干预；HA已比较成熟，自动扩容系统也较完善



## ❖ Tribe系统设计背景

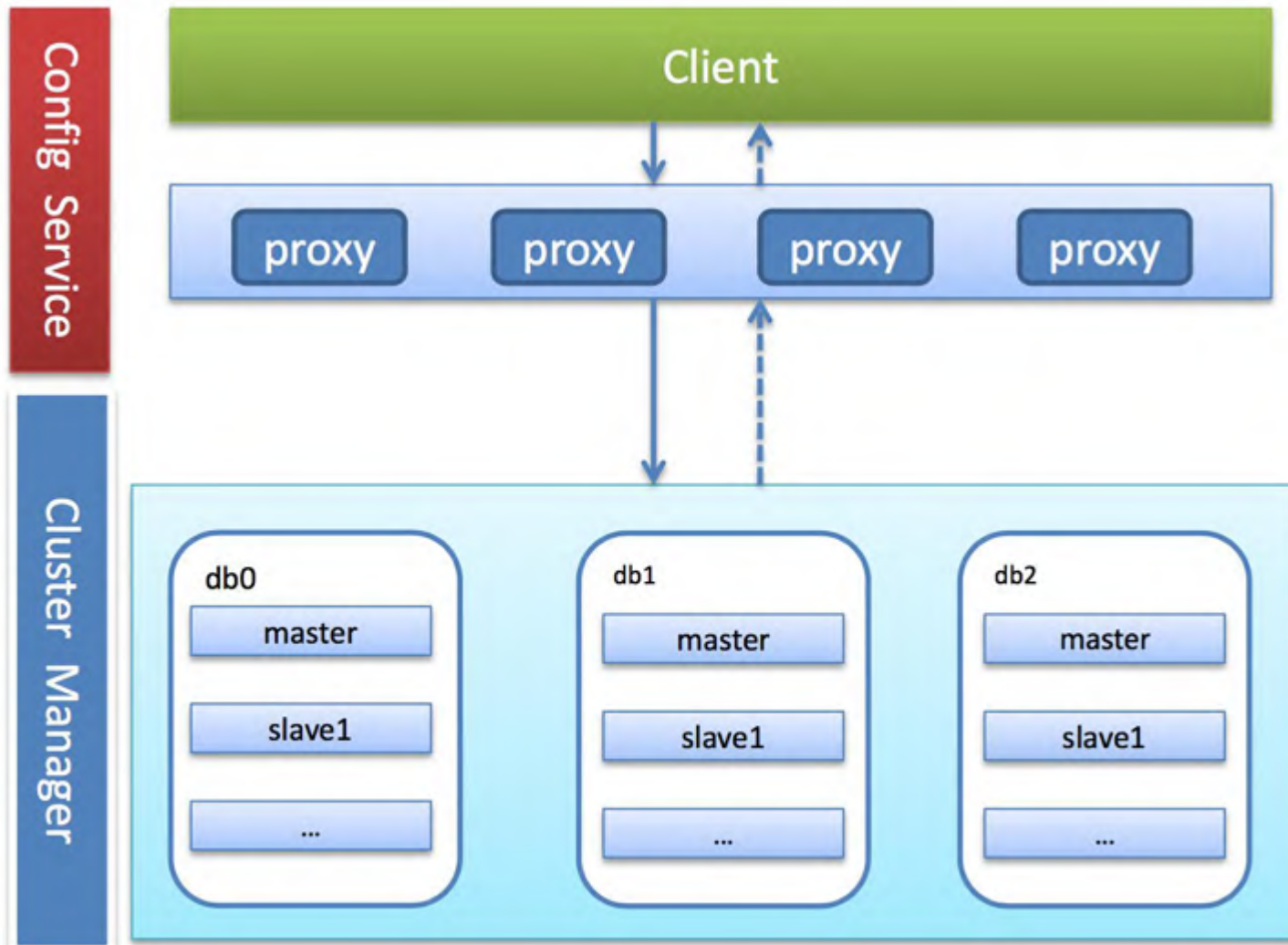
### 开发痛点

- 1、重复开发拆库，读写分离、负载均衡、跨机房部署调配、复杂度高
- 2、短连接访问redis，吞吐量=1/3长连接，RT=2倍吞吐量
- 3、拆库扩容需要业务进行双写，工作量大

### 运维痛点

- 1、扩容迁移需业务配合完成，沟通成本高
- 2、拆分、导入、校验、双写，每一步都需人工完成，出错率大
- 3、拆库困难、导致单端口数据量大，难于备份和恢复

## ❖ Tribe系统架构



- 请求路由
- 读写分离
- 负载均衡
- 配置更新
- 数据聚集
- 动态扩容
- Failover
- 数据迁移
- 服务成本
- 集群管理

## ❖ Tribe系统功能

### • 请求路由

- 数据根据配置路由表，自动进行切分
- 读写请求根据配置路由表，自动路由

### • 读写分离

- 支持配置读写分离以及只读主库

### • 负载均衡

- 支持配置机房优先策略
- 支持round-robin负载均衡策略

### • 配置更新

- Proxy自动从configservice抓取配置、并在线更新并保存本地

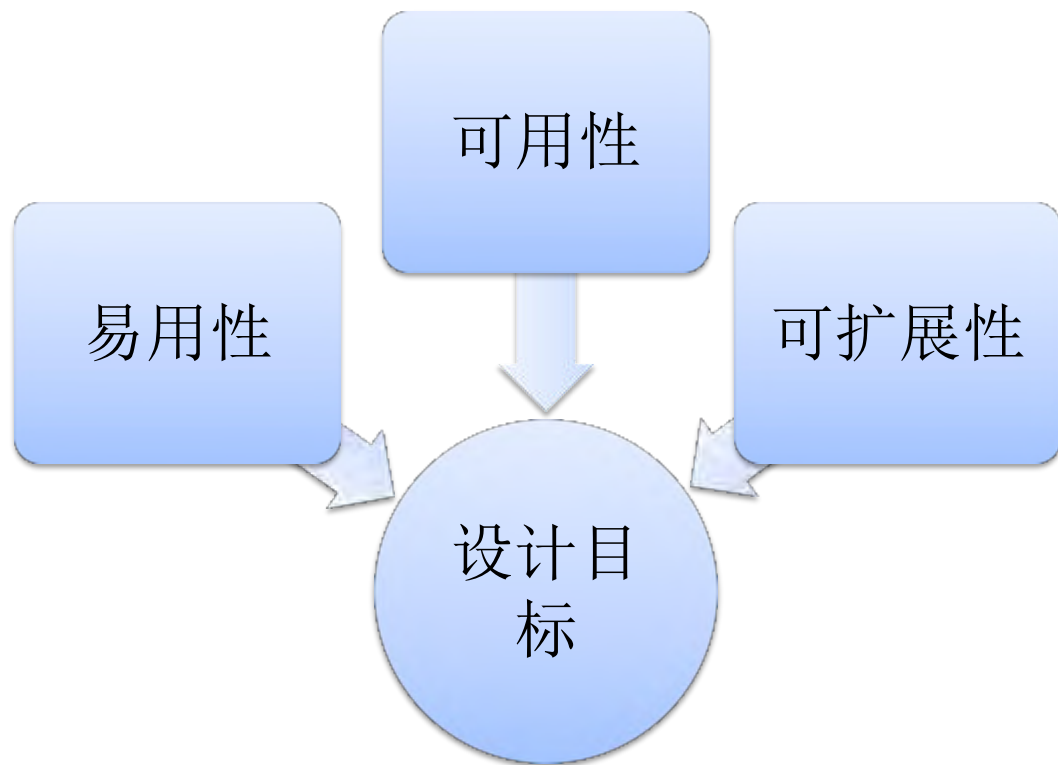
### • 数据聚集

- 支持配置前缀，将相同前缀的数据落到同一个slot中

### • 动态扩容

- 支持在线拆分端口
- 支持自动切换流量

## ❖ Tribe系统设计目标



### 易用性

- 简化业务访问逻辑，让业务更专注代码
- 减少前后端间耦合，减少变更成本

### 可用性

- 健康检查
- 自动流量切换
- 自动故障恢复

### 可扩展性

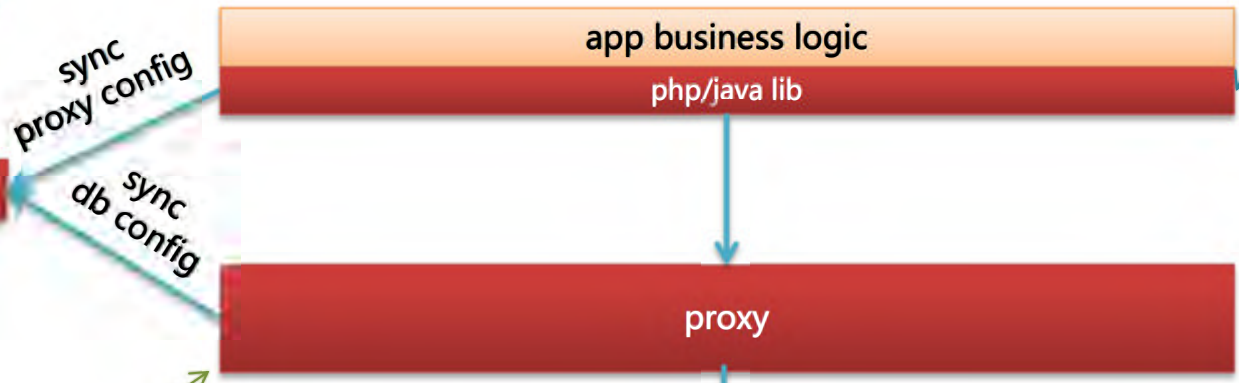
- 简单、扩展即加机器
- 扩展的策略，延迟分配
- 扩展后也能收缩



## Tribe系统设计一期

- 1. 配置管理
- 2. db信息
- 3. proxy信息

configservice

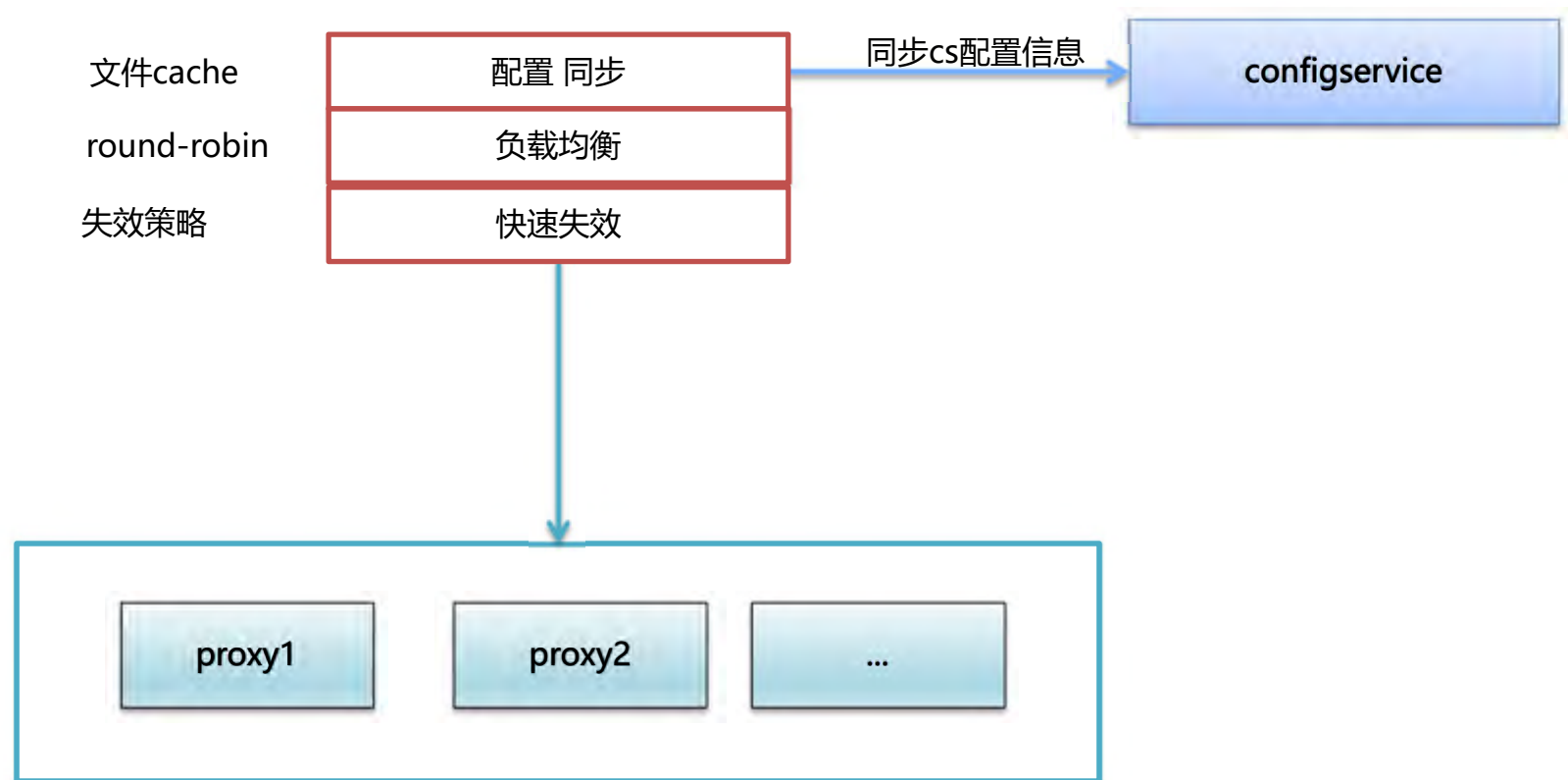


- 1. 配置信息同步
- 2. 负载均衡
- 3. 快速失效, 慢恢复

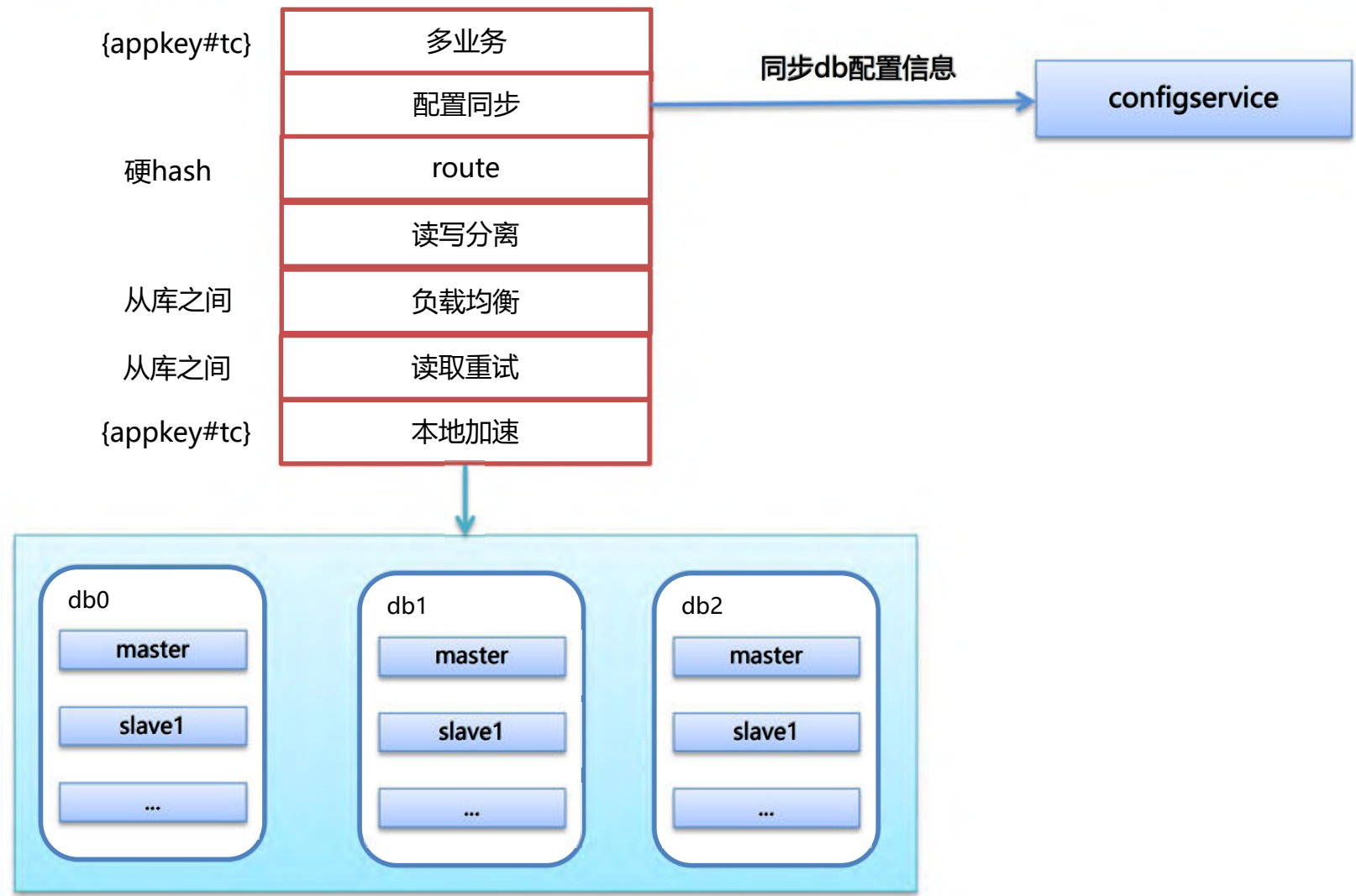
- 1. modula路由
- 2. 读写分离
- 3. 从库读取重试
- 4. 从库负载均衡, 支持指定机房策略
- 5. 支持多业务



## ❖ Java/php lib

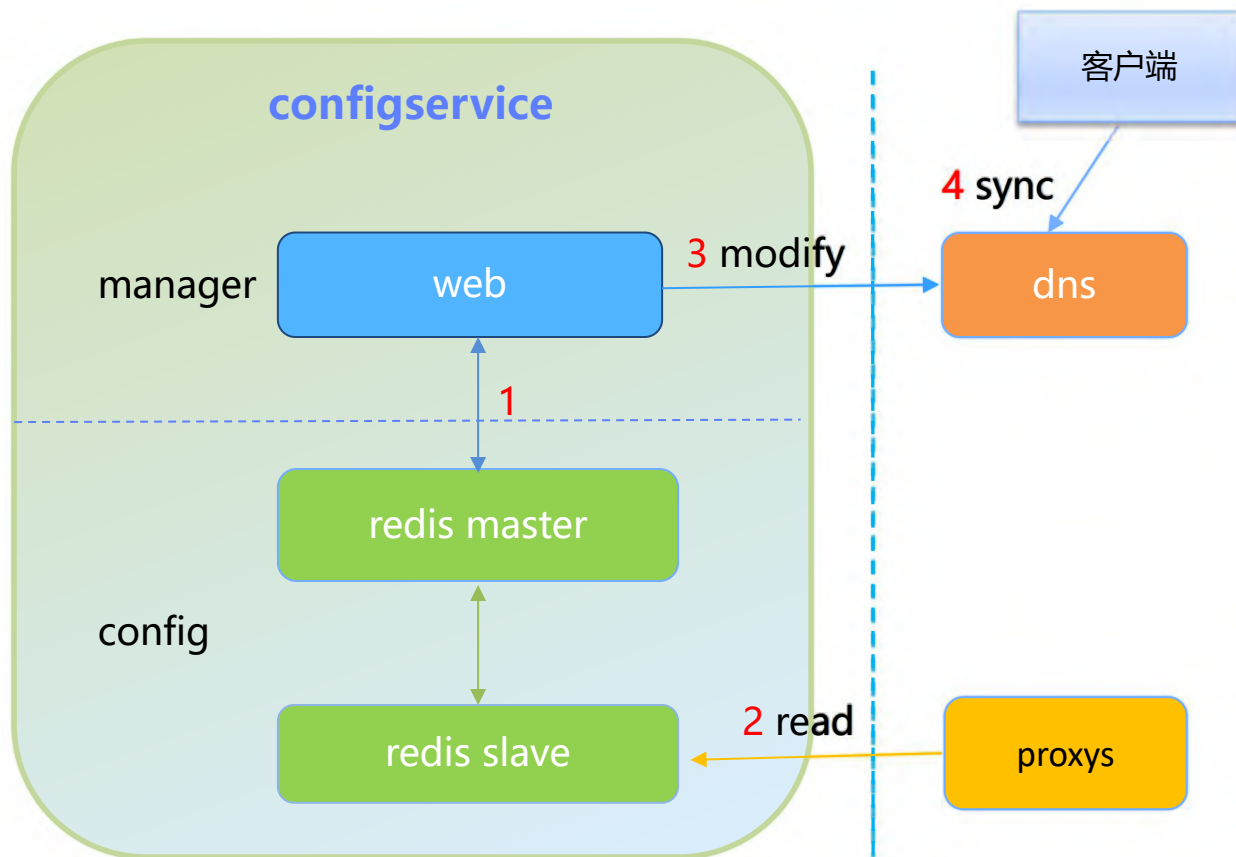


## ❖ proxy功能



## ❖ configService功能

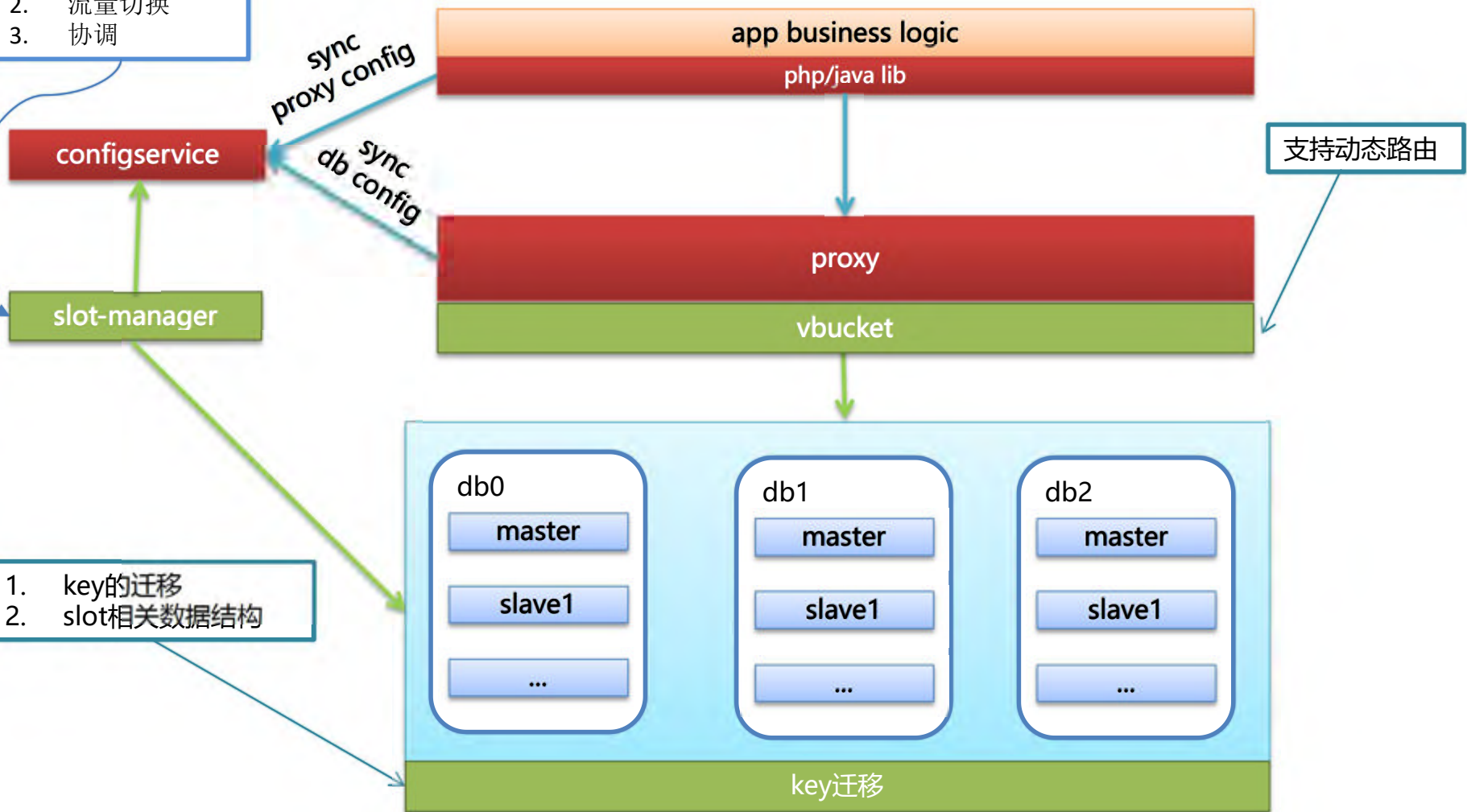
配置信息查询/修改
配置信息存储



1. 前端web将配置信息更新到redis
2. 每个proxy从redis中读取配置信息
3. web通过修改DNS来完成proxy的上下线
4. 客户端会通过DNS来sync proxy信息

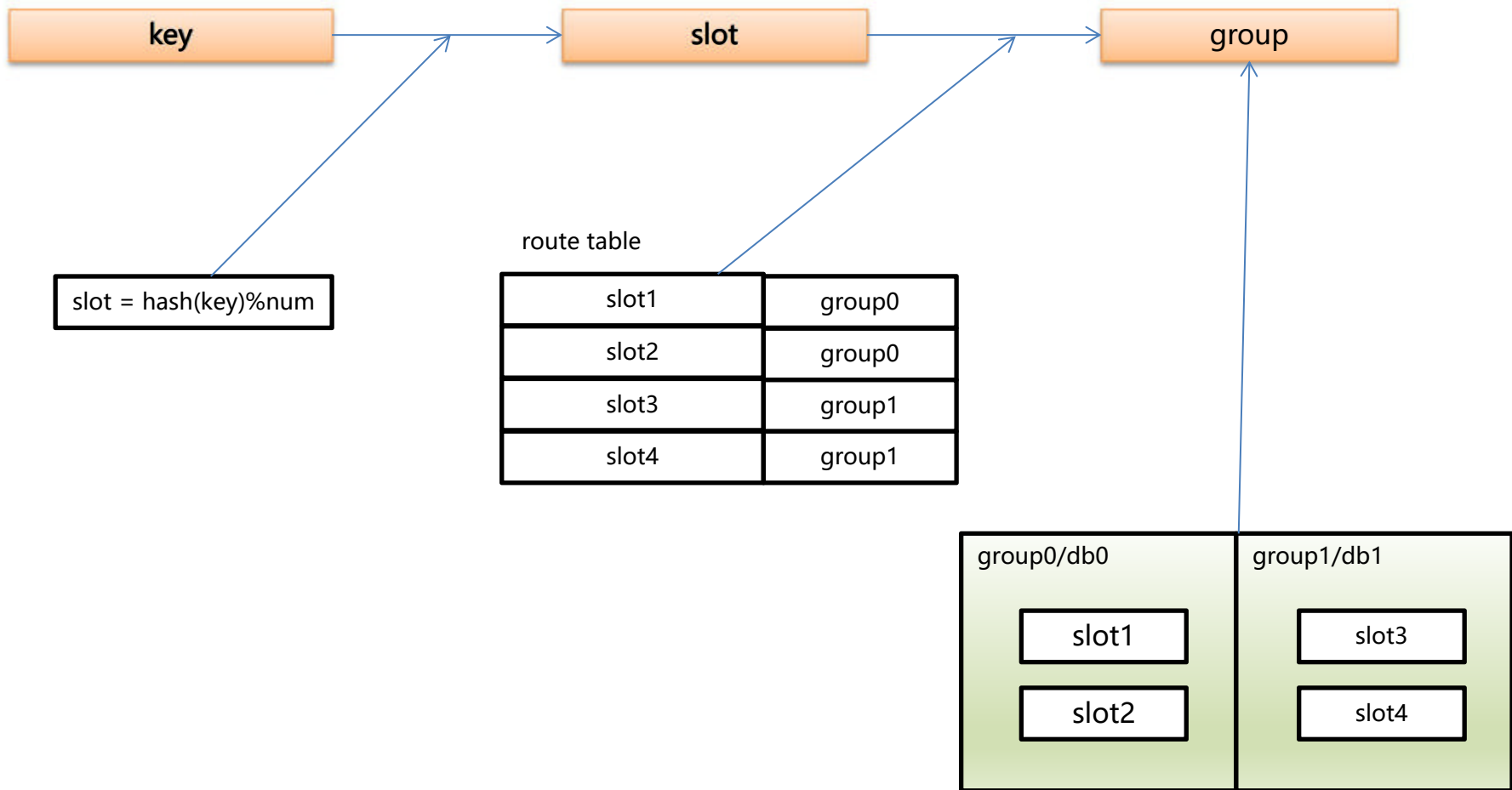
## ❖ Tribe系统设计二期

1. 在线slot迁移
2. 流量切换
3. 协调

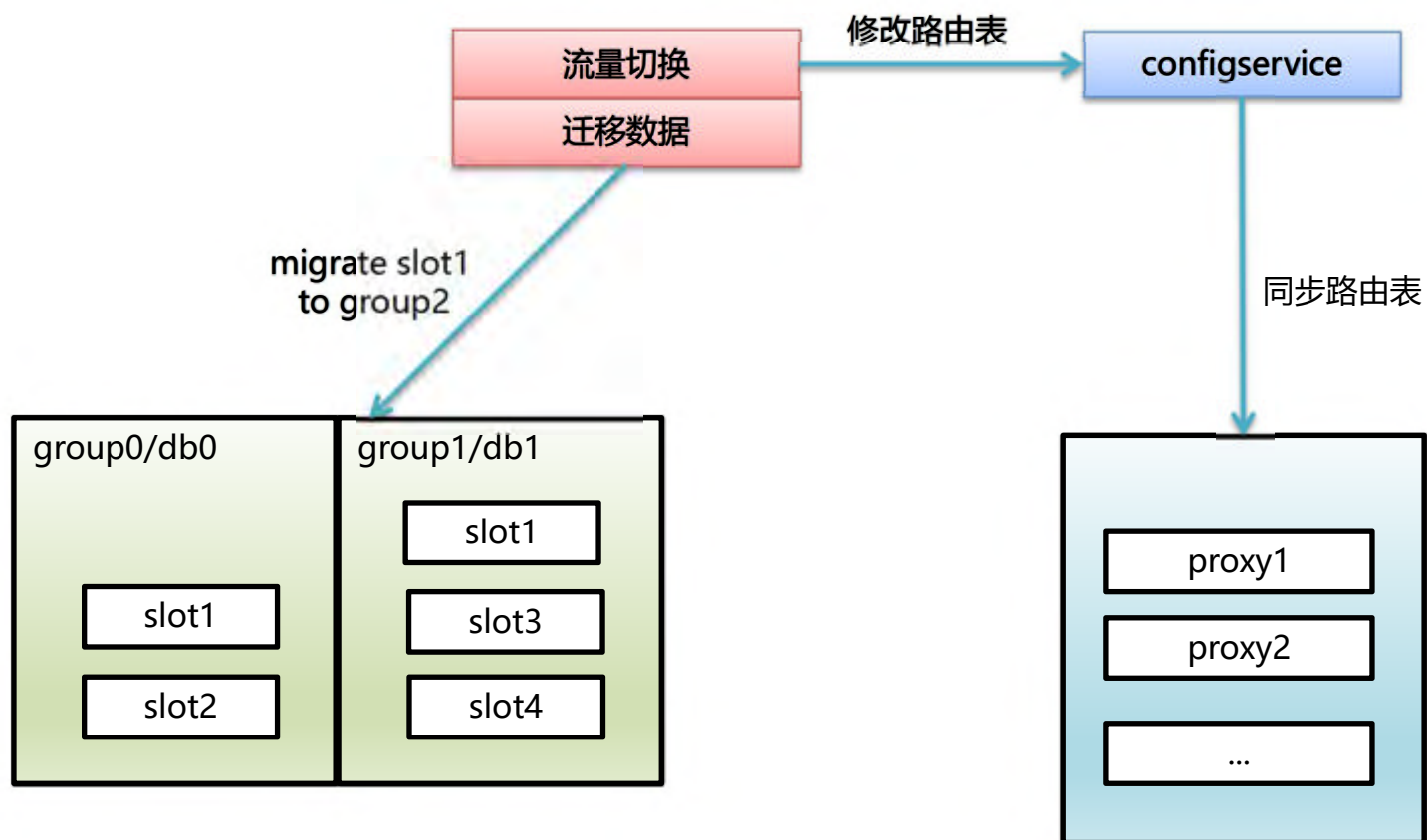


1. key的迁移
2. slot相关数据结构

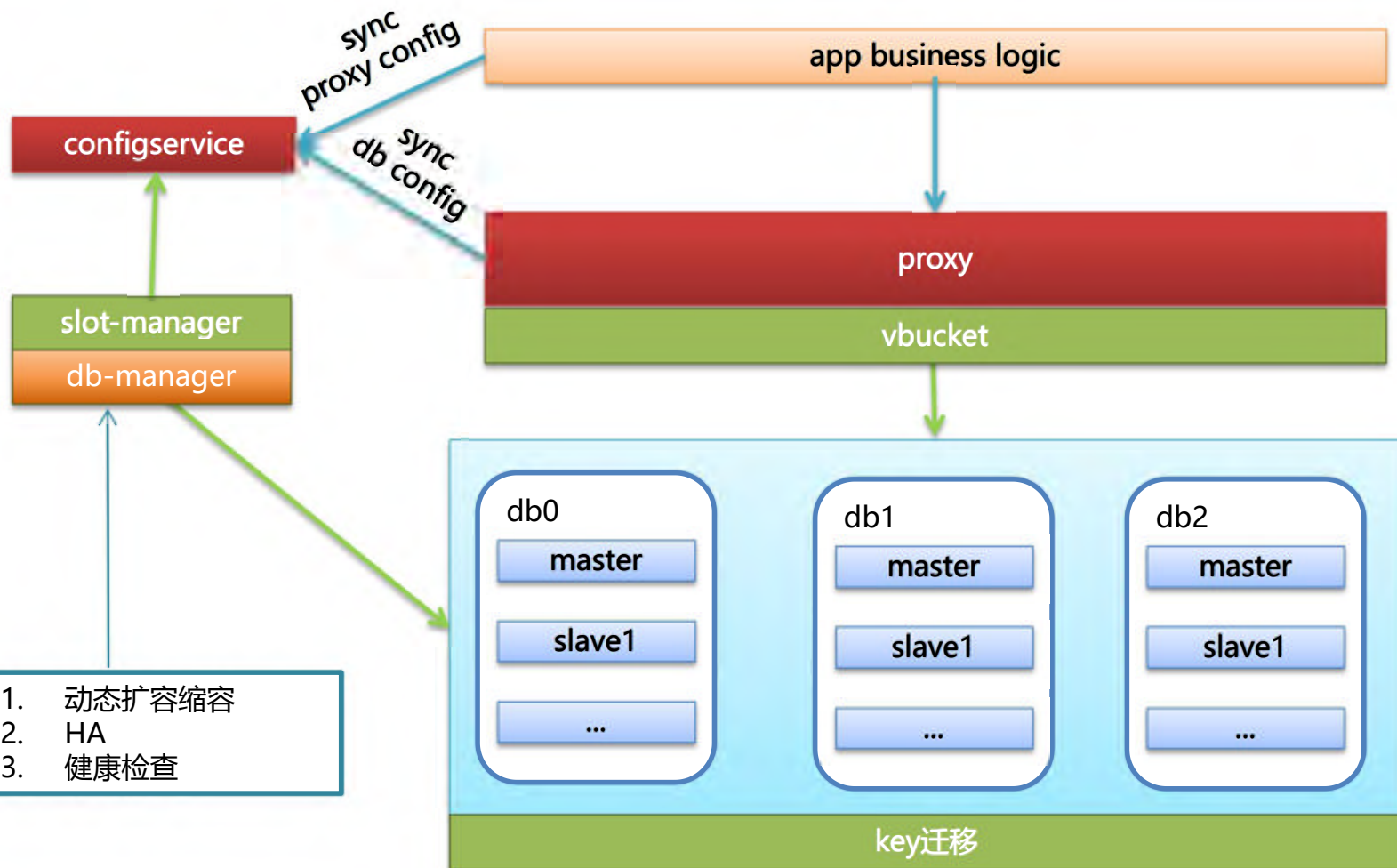
## ❖ 动态路由



## ❖ 在线slot迁移



## ❖ Tribe系统设计三期

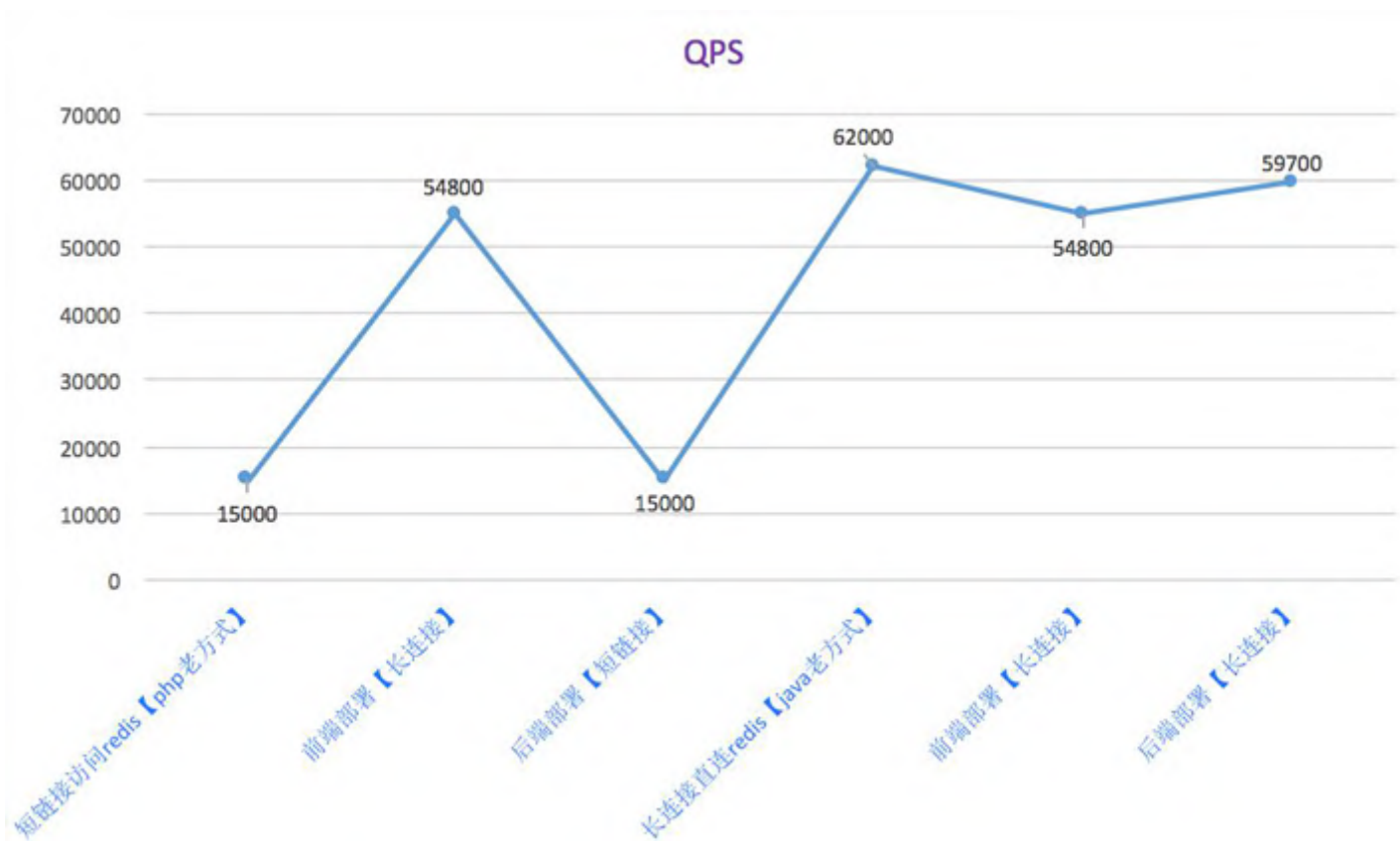




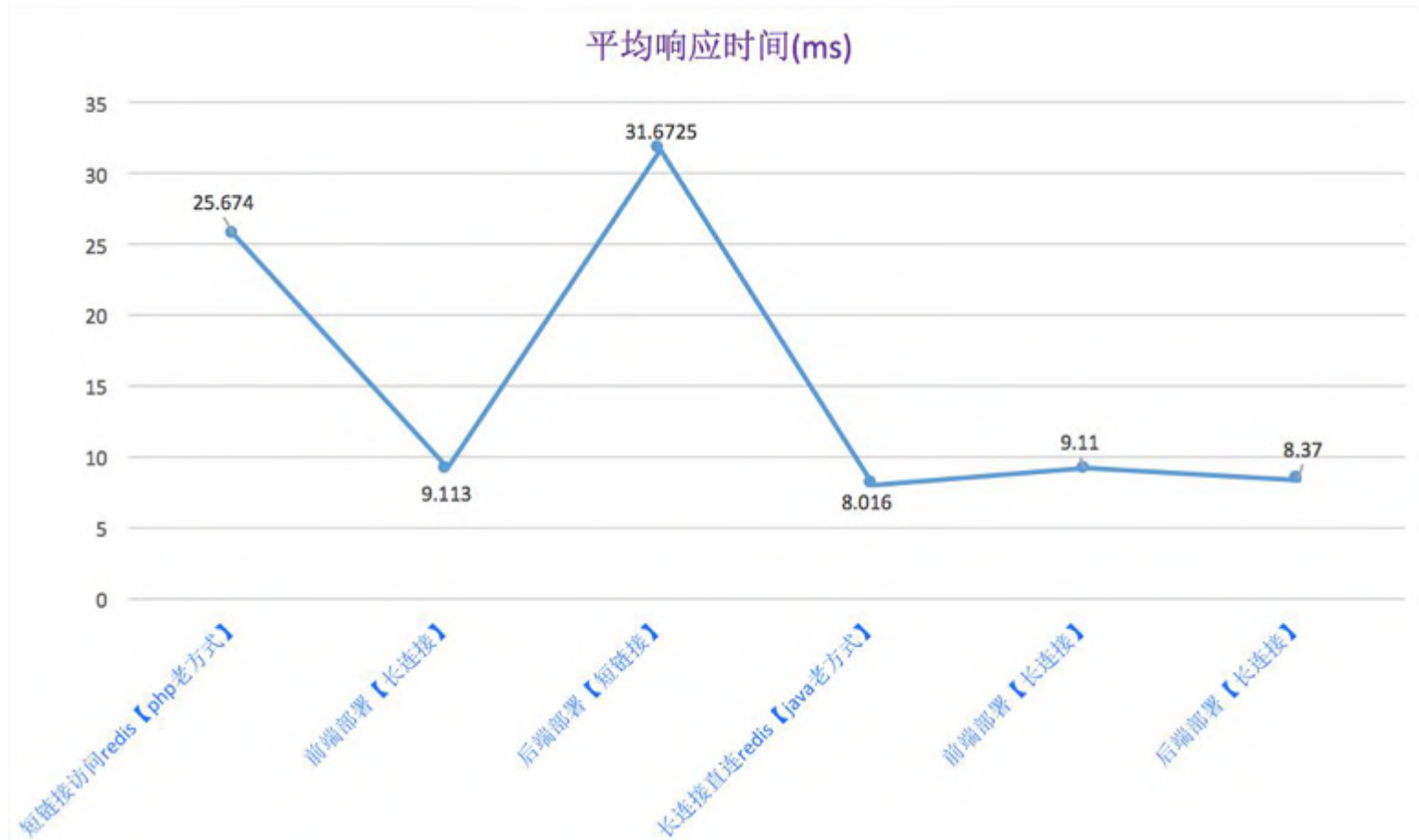
- ❖ 个人介绍
- ❖ 微博REDIS的使用介绍
- ❖ Tribe系统设计的考量
- ❖ Tribe系统性能测试对比
- ❖ Tribe系统运维点滴

## ❖ Tribe系统性能测试

- 300w/set/get
- 500线程并发
- Key:21B / value:10B



## ❖ Tribe系统性能测试



- ❖ 个人介绍
- ❖ 微博REDIS的使用介绍
- ❖ Tribe系统设计的考量
- ❖ Tribe系统性能测试对比
- ❖ Tribe系统运维点滴

## ❖ 添加业务

### ☰ 业务群列表

业务群:

添加

名称	操作
sso_db1	<a href="#">删除</a>
sso_db2	<a href="#">删除</a>

## ❖ 添加Cluster

### ☰ Cluster列表

Cluster:

添加

Cluster	操作1	操作2
sso_part_1	<a href="#">打开failover</a>	<a href="#">删除</a>
sso_part_2	<a href="#">打开failover</a>	<a href="#">删除</a>
test_db	<a href="#">打开failover</a>	<a href="#">删除</a>
weibo_healthy_7679	<a href="#">关闭failover</a>	<a href="#">删除</a>

## ❖ 基础配置

### 基本选项

Cluster Key	<input type="text" value="weibo_healthy_7679"/>
hash	<input type="text" value="crc32a"/>
distribution	<input type="text" value="slot"/>
auto_eject_hosts	<input type="text" value="true"/>
local_first	<input type="text" value="true"/>
master_emrg_read	<input type="text" value="true"/>
timeout	<input type="text" value="1000"/>
server_retry_timeout	<input type="text" value="10000"/>
server_failure_limit	<input type="text" value="5"/>

提交更改

```
// 从库挂掉，是否自动剔除
auto_eject_hosts: 'true'
// 数据分布方式，支持slot和modula
// modula = hash(key) % [port number]
distribution: slot
// hash函数选择
hash: crc32a
// 优先访问本地从库
local_first: 'true'
// 如果从库全部挂掉，允许临时读取主库
master_emrg_read: 'true'
// 主库ip配置
masters:
- 192.168.1.11:6070:tc
- 192.168.1.12:6080:tc
// 机器连续失败多少次，被认定为failure
server_failure_limit: '6'
// 如果机器failure后，多少ms以后重试是否恢复
server_retry_timeout: '10000'
// 从库配置
slaves:
- - 192.168.1.21:6050:tc
- - ""
// proxy和后端redis的超时时间
timeout: '10'
```

## ❖ 配置group

Group列表

ID	Master	Slaves	操作	操作1
4	ip:port:dc	ip:port:dc ip:port:dc ip:port:dc	添加	
17679	[REDACTED]	[REDACTED]	域名管理	自删除

## ❖ 更新路由表

路由列表 同步路由表

vbucket start	vbucket end	group id	stat	keynum
0	254	17680	consistent	29964500
255	500	17681	consistent	28853395
501	768	17682	consistent	31609336
769	1023	17679	consistent	29969000

## ❖ 添加

DNS信息

IP:  IDC: 请选择 添加Proxy

机房	域名	ip	版本	配置签名	路由签名	黑名单MD5	操作1	操作2
mas	[REDACTED]	[REDACTED]	2.0.0.35	ced33b8c376777e300892f2582d109f5	bce3471942b6b106c04d98184c5f2258	67d1ed85baaa97a8456d0511653797b	下线	停止报警
		[REDACTED]	2.0.0.35	ced33b8c376777e300892f2582d109f5	bce3471942b6b106c04d98184c5f2258	67d1ed85baaa97a8456d0511653797b	下线	停止报警
		[REDACTED]	2.0.0.35	ced33b8c376777e300892f2582d109f5	bce3471942b6b106c04d98184c5f2258	67d1ed85baaa97a8456d0511653797b	下线	停止报警

## ❖ 添加产品线信息

### 产品线信息

所属产品:

proxy端口:

部署模式:

报警接收人:

## ❖ 数据迁移

### 数据迁移

jobid	source group	target group	vbucket start	vbucket end	计划时间	job status	clean status	开始时间	结束时间	操作
	<input type="text" value="17679"/>	<input type="text" value="17679"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>					<input type="button" value="添加计划"/>
389	17679	17682	501	768	2016-05-06 17:40:00	迁移成功	清理成功	2016-05-06 17:40:05	2016-05-06 18:18:27	<input type="button" value="查看日志"/>
388	17679	17681	255	500	2016-05-06 16:20:00	迁移成功	清理成功	2016-05-06 16:20:09	2016-05-06 16:54:44	<input type="button" value="查看日志"/>
387	17679	17680	201	254	2016-05-06 14:50:00	迁移成功	清理成功	2016-05-06 14:50:05	2016-05-06 14:57:23	<input type="button" value="查看日志"/>
386	17679	17680	101	200	2016-05-06 13:30:00	迁移成功	清理成功	2016-05-06 13:30:01	2016-05-06 13:44:08	<input type="button" value="查看日志"/>
385	17679	17680	51	100	2016-05-06 12:00:00	迁移成功	清理成功	2016-05-06 12:00:03	2016-05-06 12:06:59	<input type="button" value="查看日志"/>
384	17679	17680	0	50	2016-05-06 11:00:00	迁移成功	清理成功	2016-05-06 11:00:04	2016-05-06 11:07:07	<input type="button" value="查看日志"/>



以微博之力 让世界更美！

*weibo.com*