

NoSQL集群运营之道

腾讯分布式KV存储实践

赵乐任

关于我

赵乐任(laurentzhao)

腾讯社交网络平台技术运营中心

高级数据运维工程师

前华为DBA(Oracle)

擅长分布式存储、关系数据库

提纲

腾讯分布式KV存储介绍

- CKV实现

运营体系建设

- 鹦鹉螺

运维挑战与实践

- 自动化
- 成本优化
- 问题定位

回顾

腾讯三大KV存储组件

CKV

- 公共业务
- QQ空间、相册、音乐、广点通

Grocery

- QQ

Quorum

- 微信

社交网络业务类型

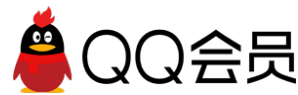
QQ空间：计数、feeds

QQ：关系链、群消息

广点通：用户画像、计费

热点：春节红包、F码

推送：红点、游戏活动



社交网络存储规模

10000+ 存储服务器

四地部署

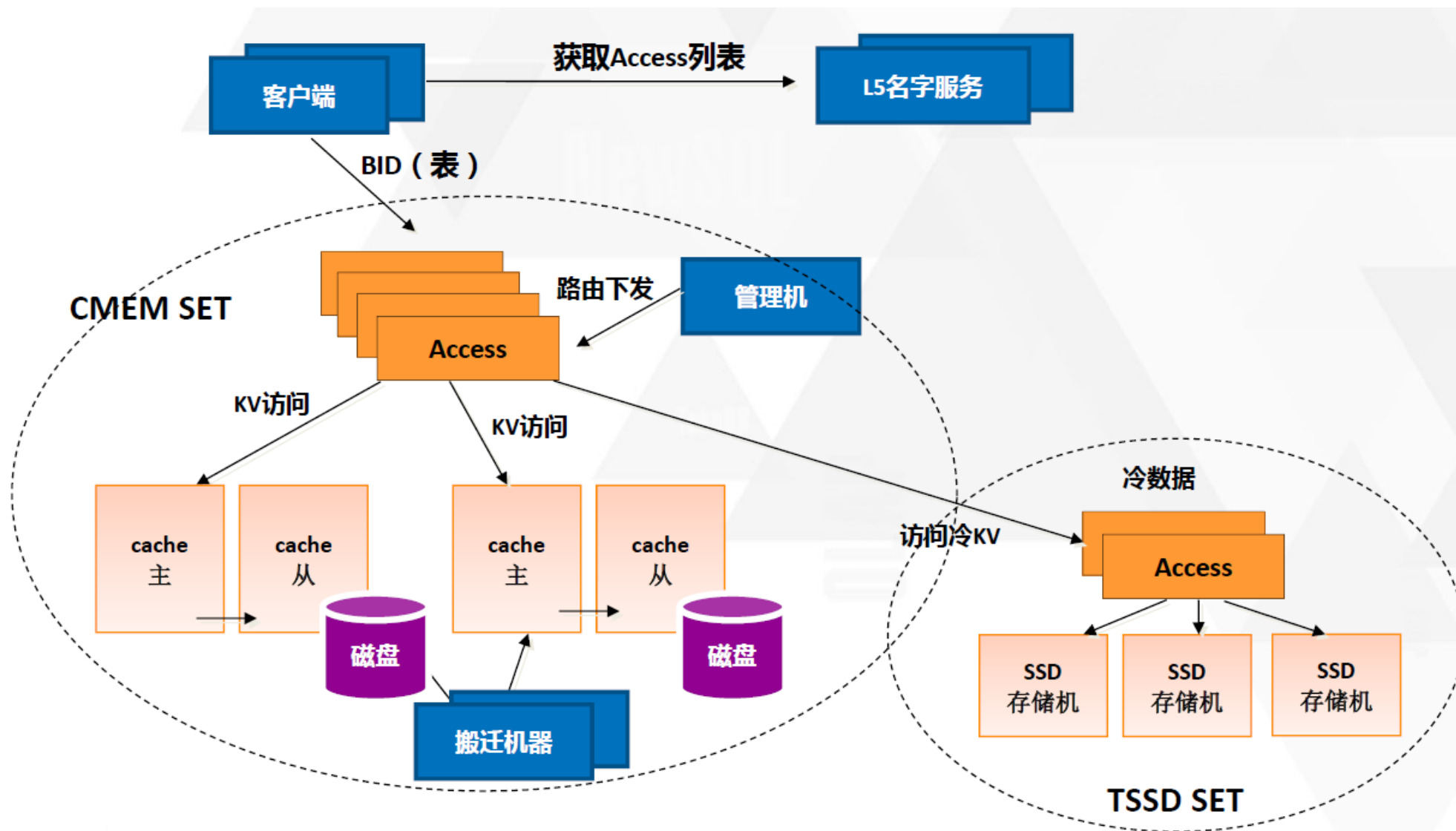
100+ 仓库

140+ TB 内存存储

1亿+ QPS



存储架构(CKV)



存储分片(CKV)

Access

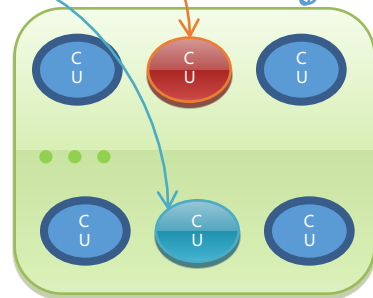
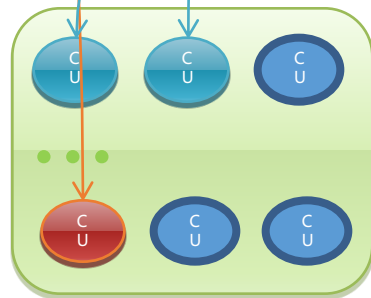
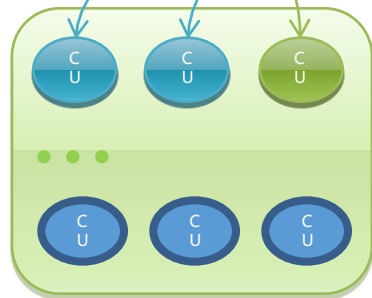
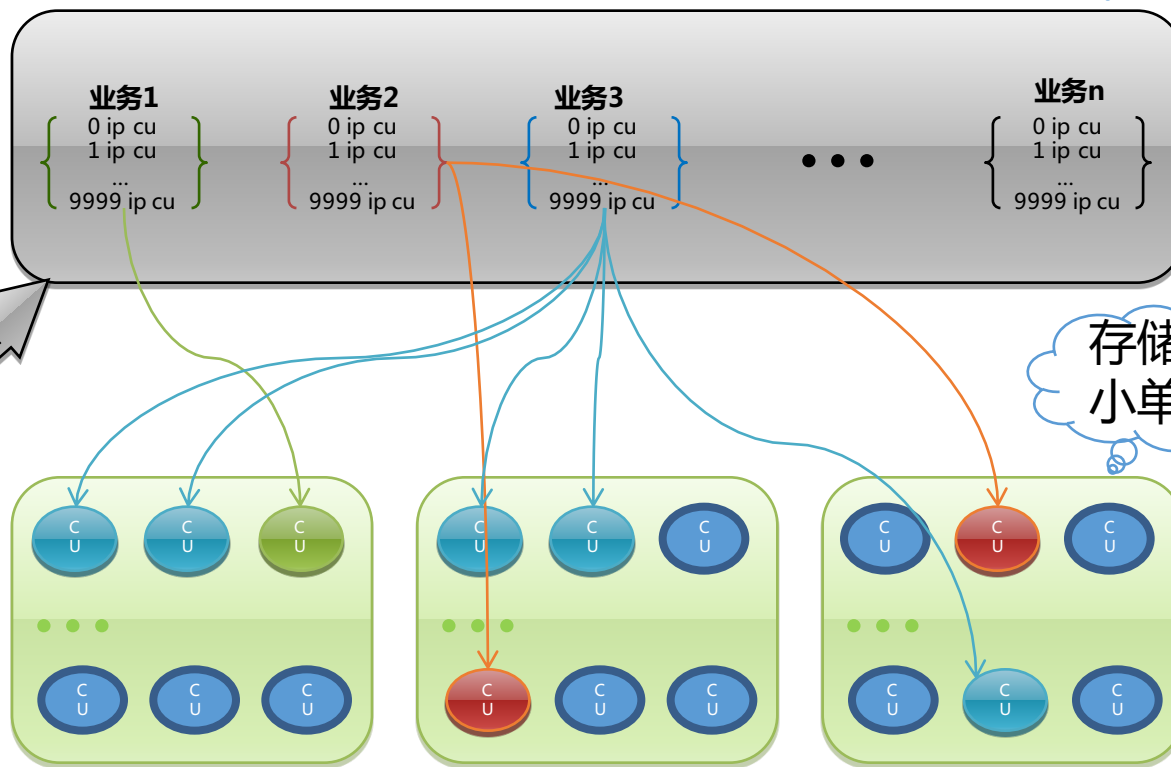
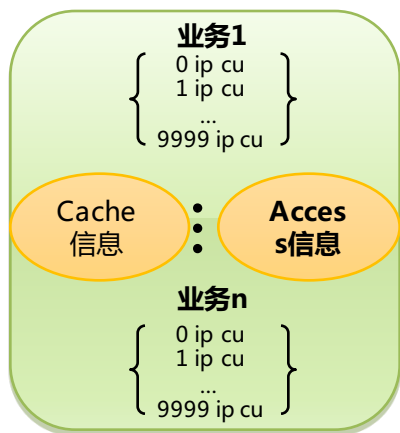
业务数据hash
到1W桶个

1w个桶与存储单元间的映射关系

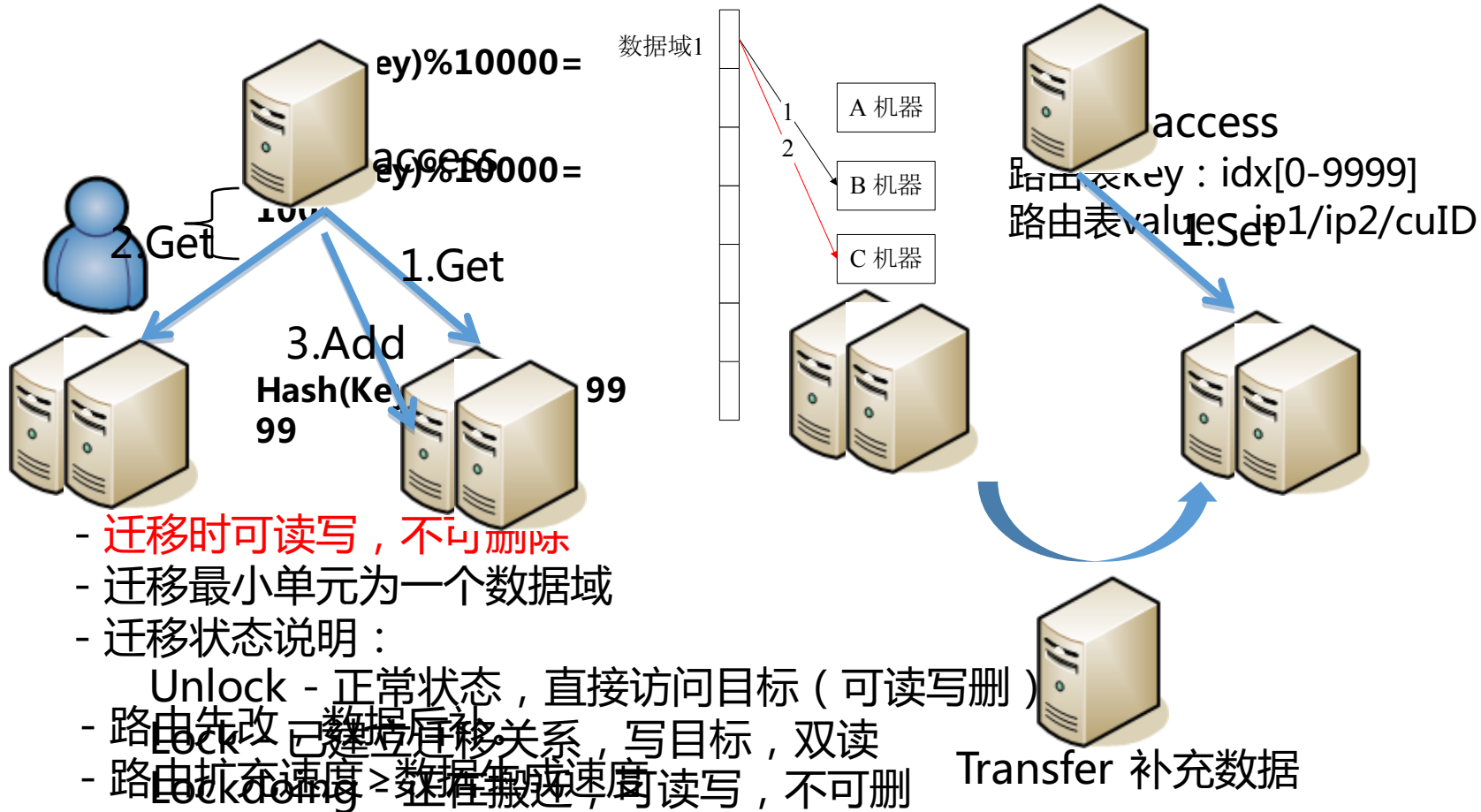
存储最小单元

Master

推送路由



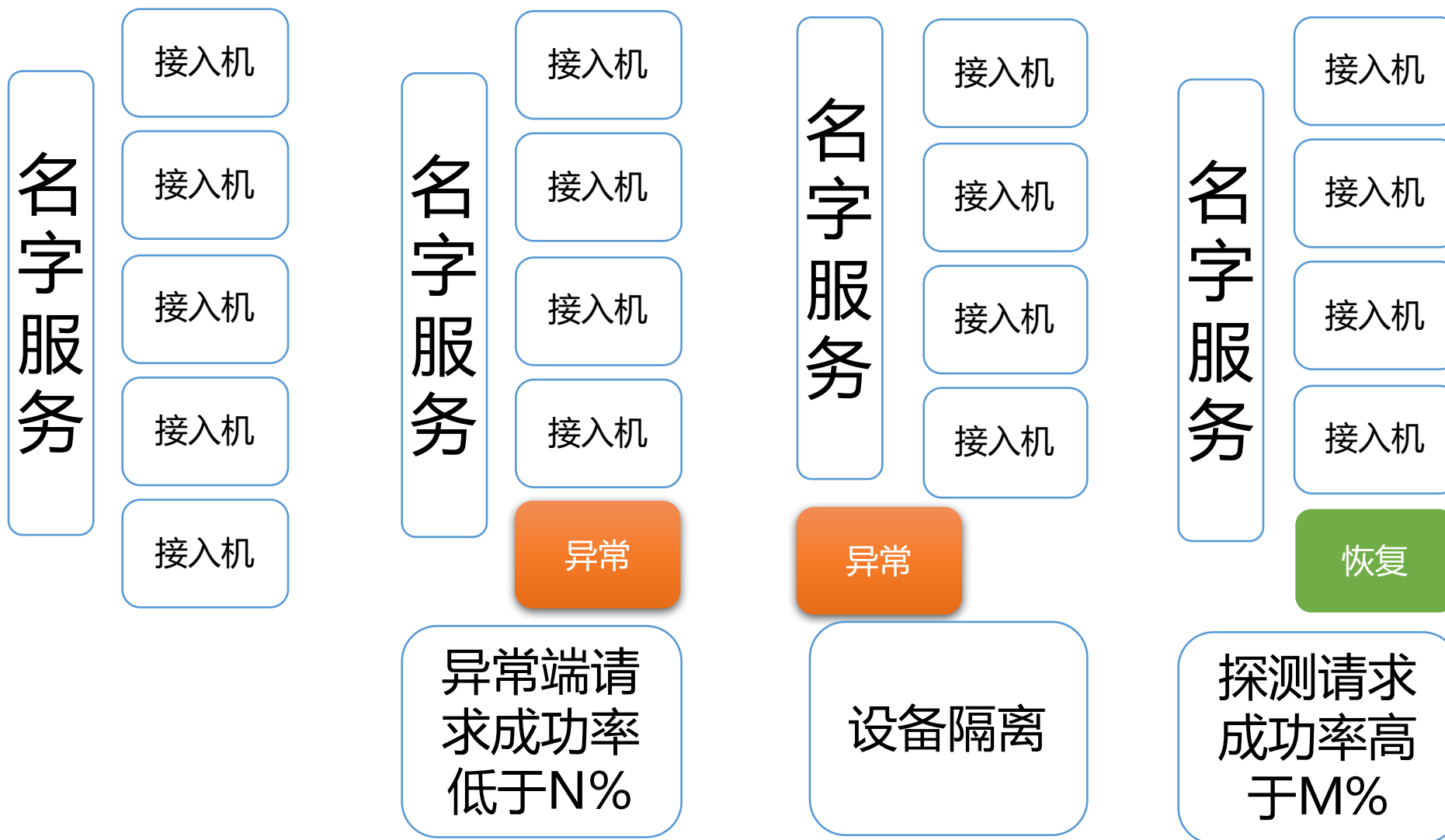
存储搬迁(CKV)



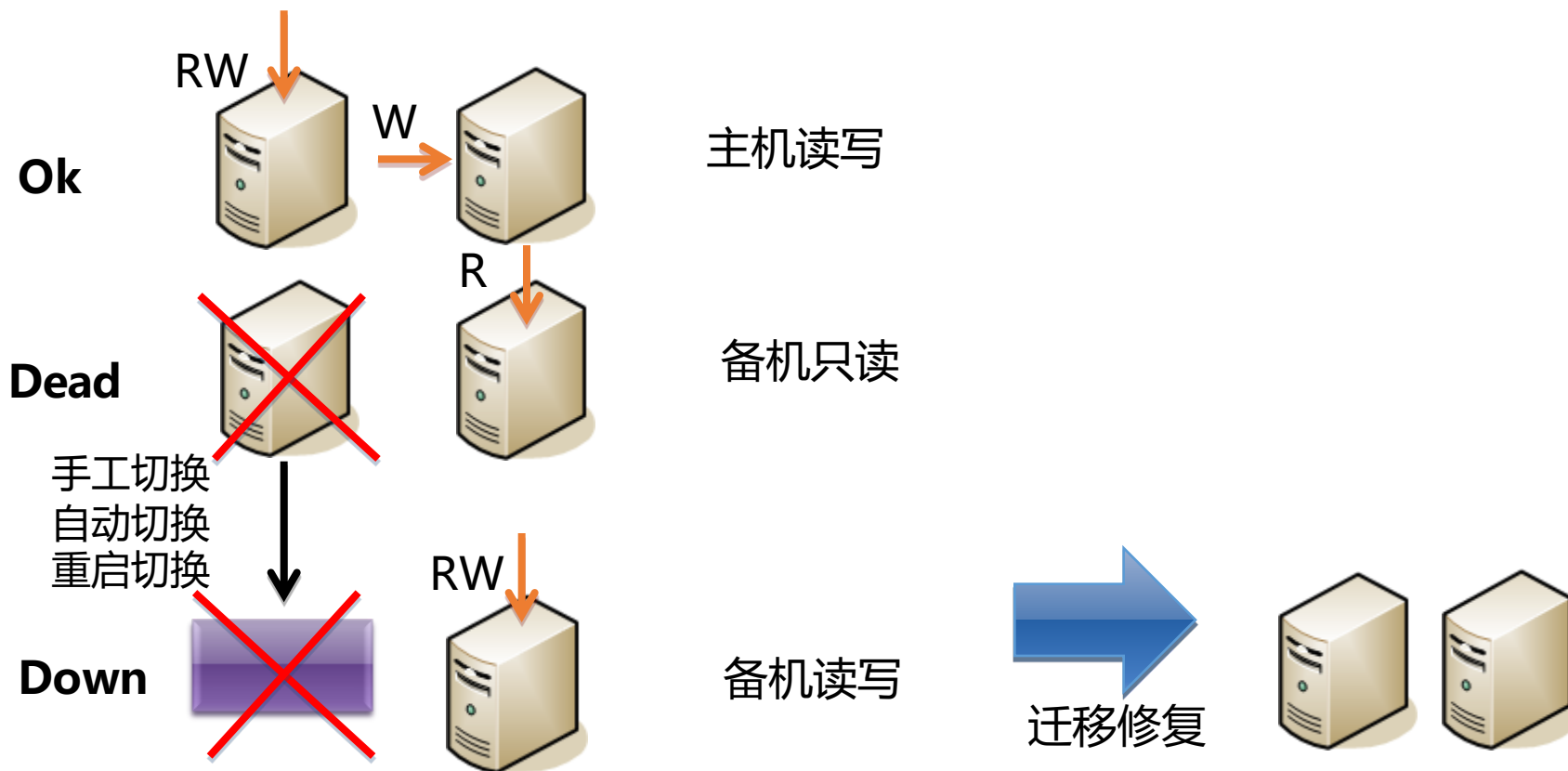
数据搬迁(扩缩容)
对业务无感知
高效(30min/台)
可控(1/w)

- 迁移时可读写，不可删除
- 迁移最小单元为一个数据域
- 迁移状态说明：
 - Unlock - 正常状态，直接访问目标（可读写删）
 - Lock - 已建立迁移关系，写目标，双读
 - Lock doing - 还在搬运，可读写，不可删

存储容灾(CKV)

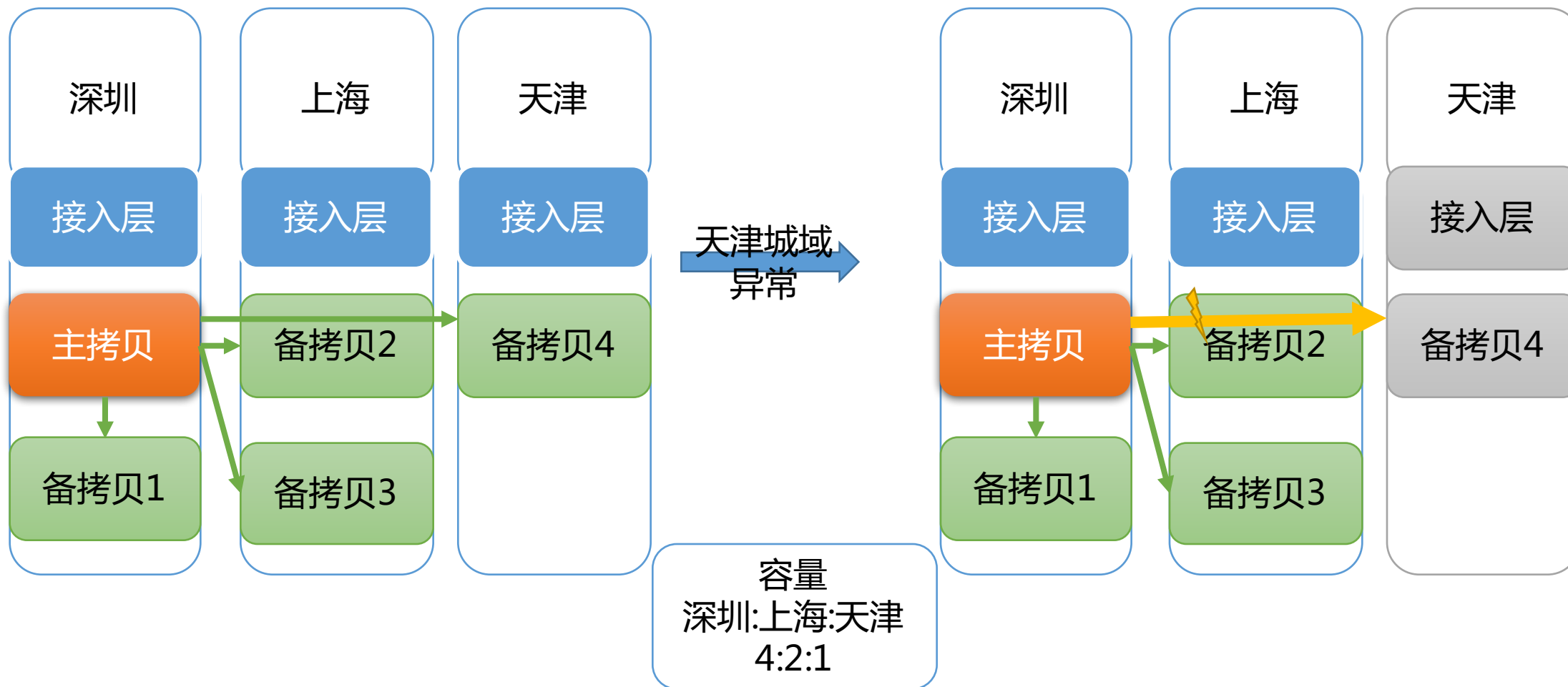


存储容灾(CKV)



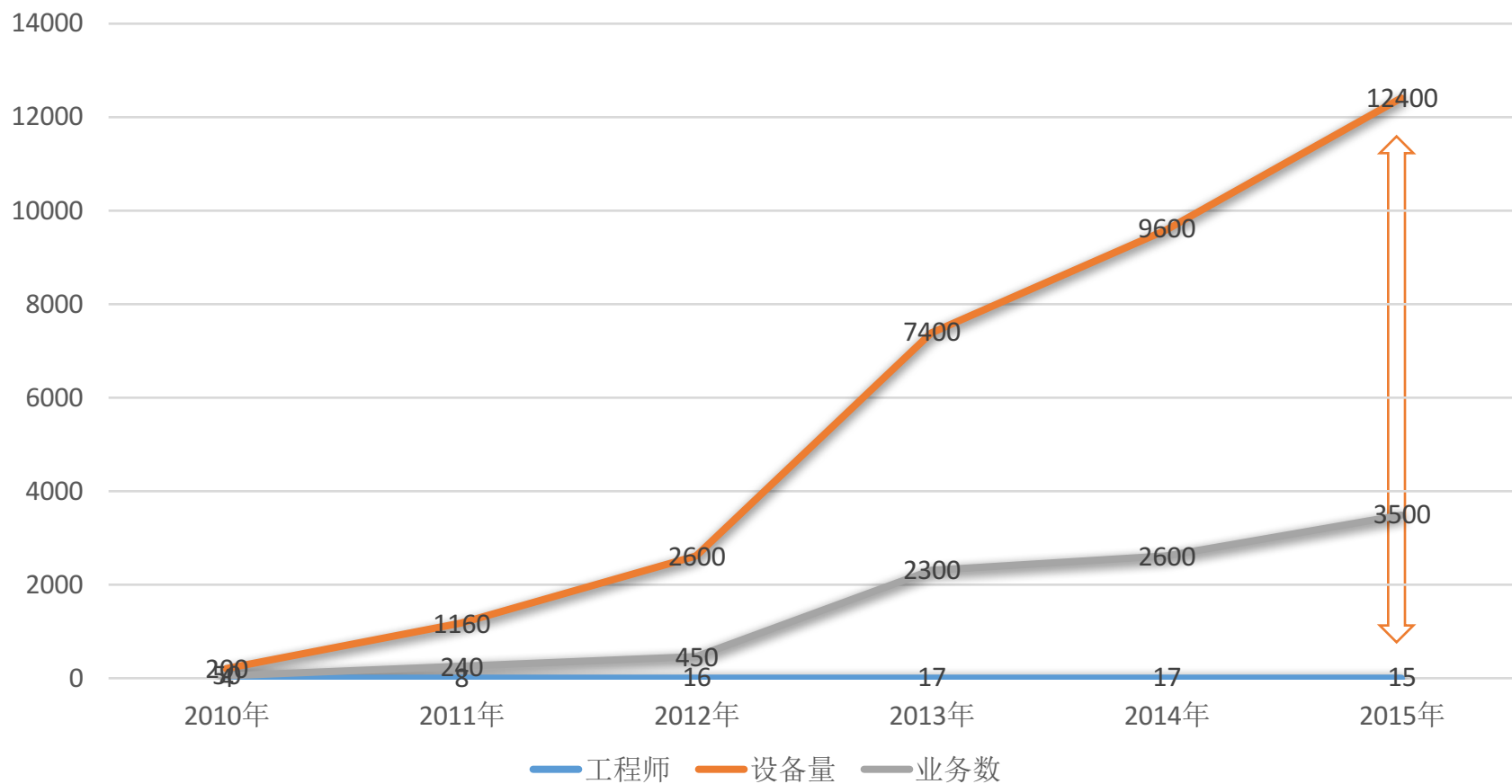
要及时地把数据迁移到新的一套机器上面，以免备机也出问题影响业务。

存储容灾(CKV)



野蛮生长

业务增长vs运维工程师



业务
2010年 10个
2015年 3500个

人均设备维护量
2010年 50台
2015年 1030台

存储运维角色

业务规划



存储运维

工具建设



运营平台建设

存储运维

业务/开发

运营平台

平台开发

开发框架

yii

django

前端展示

highcharts

bootstrap

队列存储

Celery

Redis

流程策略

上线

扩容

屏蔽

下线

数据分析

其他

扩容

收敛

其他

基础组件

存储组件

路由

备机接管

过期删除

接入层

资源调度

备份

存储层

数据下沉

其他

公共组件

名字服务

包管理

远程命令

告警通道

网管系统

统一登录

配置管理

接入网关

其他

运营平台建设

新增流程 修改流程 删除流程 使用选定流程模版建单 搜索关键字:

#	<input type="checkbox"/>	流程模版ID	流程模版名	流程类型	更新人
1	<input type="checkbox"/>	36	Grocery_SSDCacheInstall	上线流程	vivowan
2	<input type="checkbox"/>	35	Grocery_SSDCcInterfaceInstall	上线流程	vivowan
3	<input type="checkbox"/>	23	GroceryInterfaceInstall	上线流程	summer
4	<input type="checkbox"/>	22	GroceryCcInterfaceInstall	上线流程	vivowan
5	<input type="checkbox"/>	21	GroceryCacheInstall	上线流程	amfreez

基础标准化
任务原子化
流程可视化

[MemSsdInterface_1.1]执行流程



刷新步骤

开始执行

暂停执行

终止执行

强制执行
选定步骤

强制从选
定步骤开
始执行

下面两个按钮
仅调试时使用，
请慎用！

重置任务
为未执行

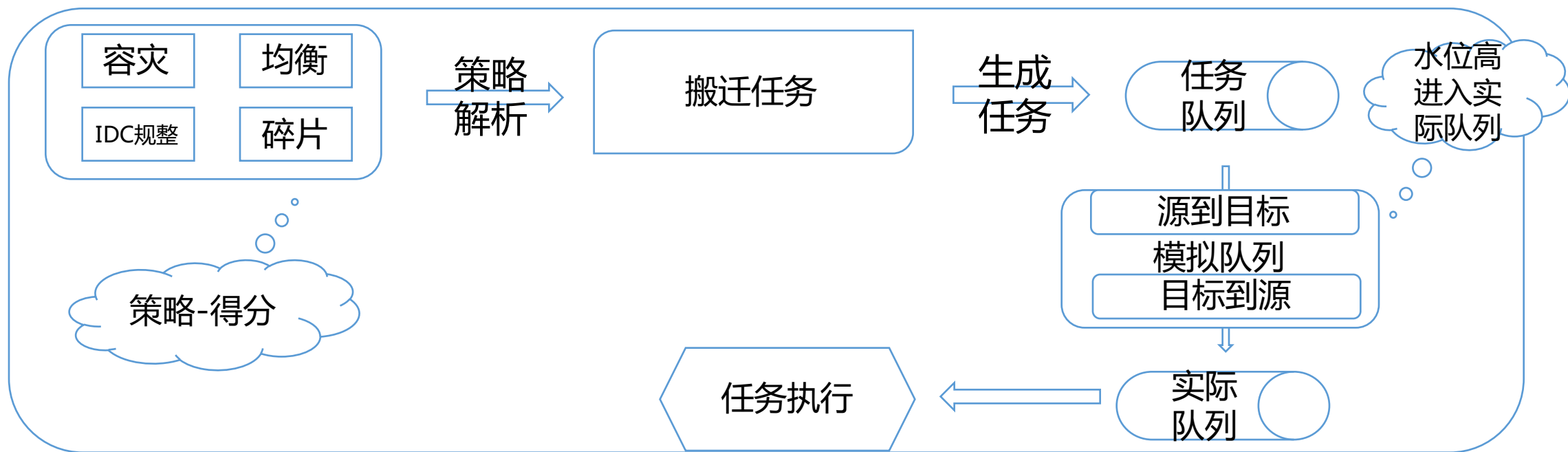
自动化建设

- 调度能力

- 名字服务实现接入层弹性
- 数据搬迁实现存储层弹性

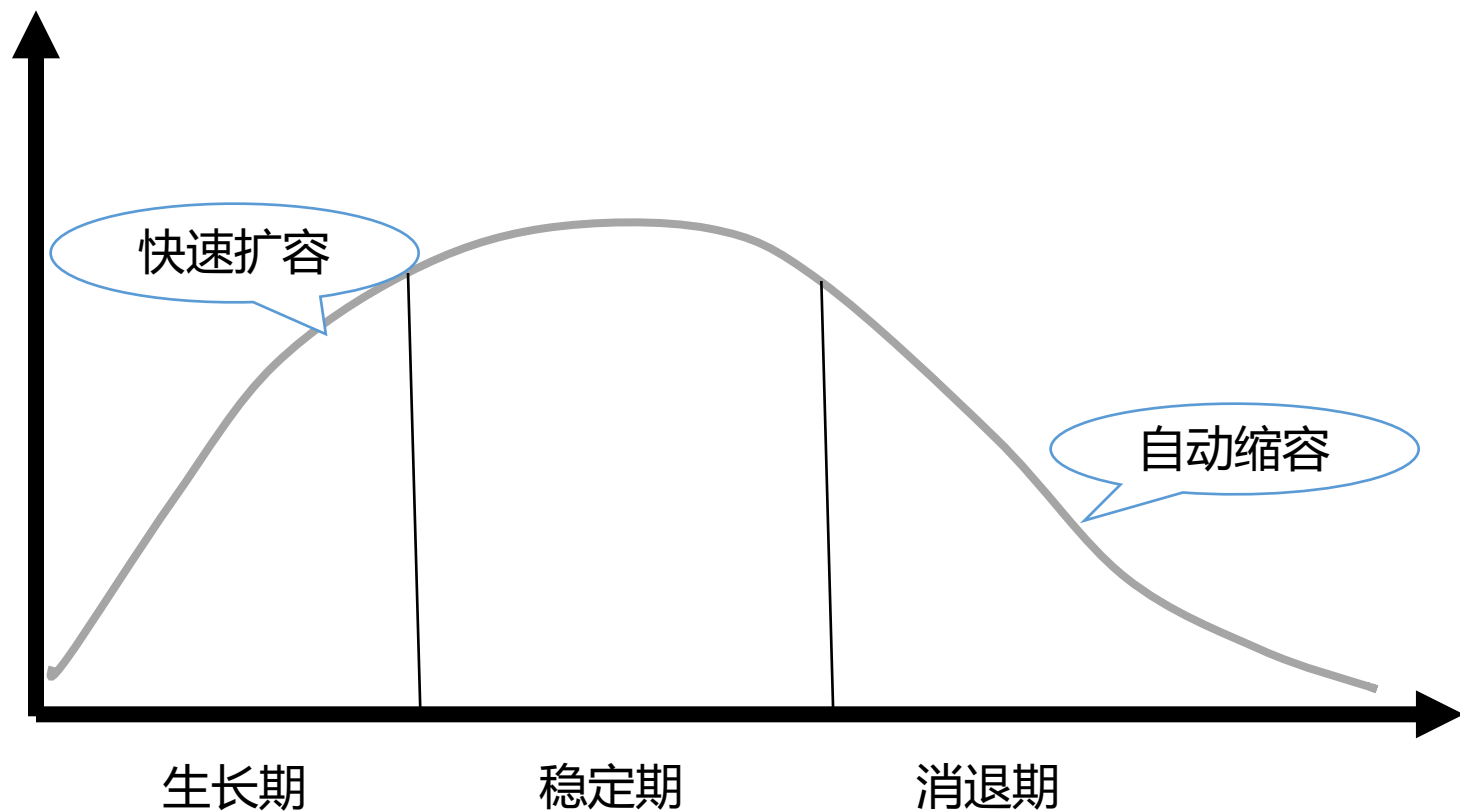
- 调度策略

- 负载均衡
- 高可用保障、容灾等



自动化建设

- 生命周期管理



自动化扩容

自动扩容

机房	SET	Bid	扩容cu个数	迁移后使用率	扩容次数
深圳机房	仓库	101020086	12	71%	2
深圳机房	ayC	20500026	3	63%	1
深圳机房	ayC	20500025	1	0%	1
深圳机房	ayC	20500024	10	77%	2
深圳机房	ayC	20500021	16	79%	3
深圳机房	ayC	20500020	3	9%	3
深圳机房	ayC	20500019	5	70%	3
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				
深圳机房	ayC				

自动扩容				
机房	完成个数	失败个数	涉及cu个数	成功率
机房	32	2	93	94%
机房	11	2	13	84%
自营	9	0	53	100%
机房	18	8	18	69%
SD机房	11	1	68	91%
总计	81	13	245	86%

每周两百多起实例自动扩容

自动化缩容

鸚鵡螺系统邮件

2015年12月3日

CMEM缩容报告：

BID	Region	Depot	可用总量	使用率	操作	
101021091	深圳自 营	[ISD]深圳_自 营 仓库	2 TB	67.00%	趋势图 缩容	
20050003	深圳自 营	[ISD]深圳_自 营 仓库	.02 GB	74.00%	趋势图 缩容	
101060004	深圳自 营	[ISD]深圳_自 营 仓库	.19 GB	12.00%	趋势图 缩容	
101021061	深圳自 营	[ISD]深圳_自 营 仓库	.69 GB	44.00%	趋势图 缩容	
101021474	深圳自 营	[ISD]深圳_自 营 仓库	.1 GB	52.00%	趋势图 缩容	
20050004	深圳自 营	[ISD]深圳_自 营 仓库	10.185.11.78 _SET3	10.129.130.216 103 30.66 GB	63.19 GB 82.00%	趋势图 缩容

hi vis :

bid 四级模块

102020034[N][腾讯

在 2015-12-02 发生了缩容。

缩容可能由于满足如下策略:

- 1、容量首次发现低于70%以下，持续3天
- 2、缩容到容量占用80%以上或者cu 个数为1为止
- 3、内部业务bid上线30天后，才会触发自动缩容

请关注该业务bid是否申请容量过大或者业务尚未上线而造成
需要恢复容量，烦请联系ckv_sng_helper

[CMEM]

邮件由系统自动发送，请勿直接回复！有疑问请联系: nerrissaniu

鸚鵡螺由 SNG - 社交网络运营部 - 数据运维组 开发维护

成本优化

提升访问密度的方法

- ◆ 业务逻辑优化
- ◆ 数据压缩
- ◆ 过期淘汰
- ◆ 备机复用
- ◆ 冷热分离
- ◆ 碎片整理

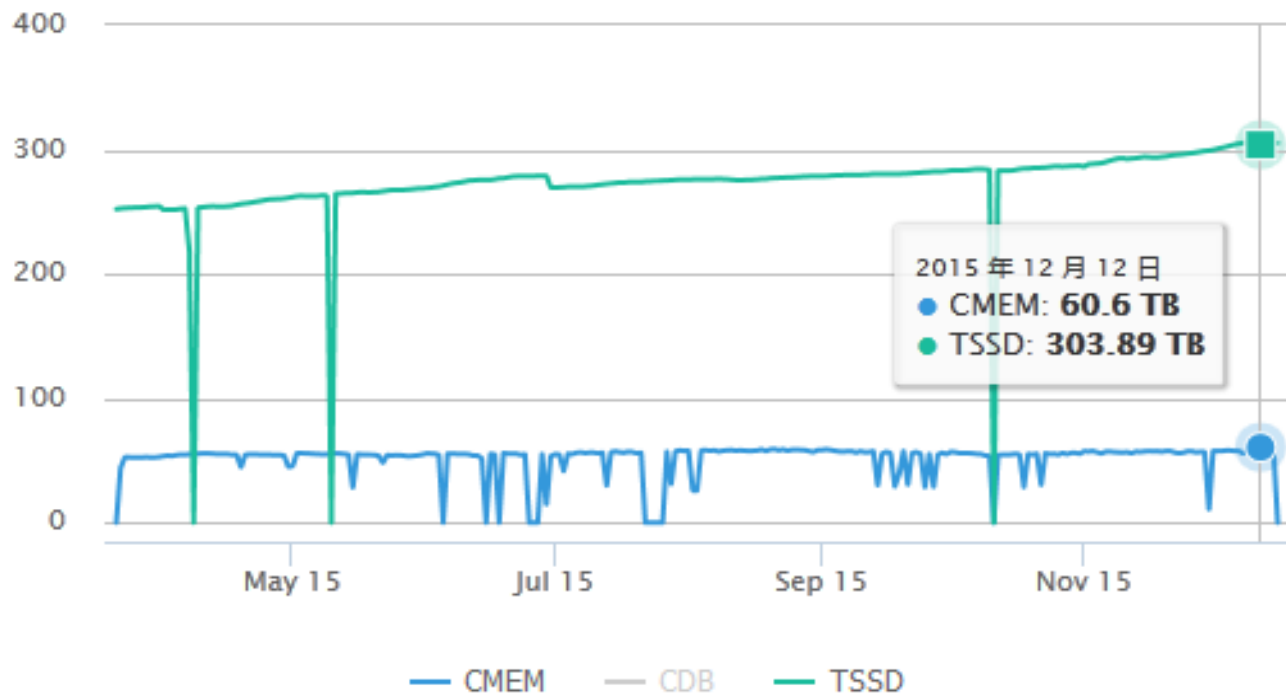
业务QPS/存储量



冷热分离(CKV)

整体淘汰比(内存/总量)84%

数据存储量趋势 (单位: TB)

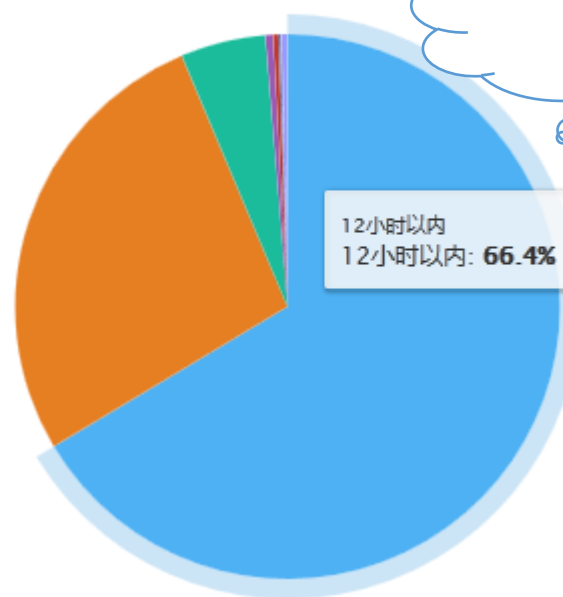


元数据
分析

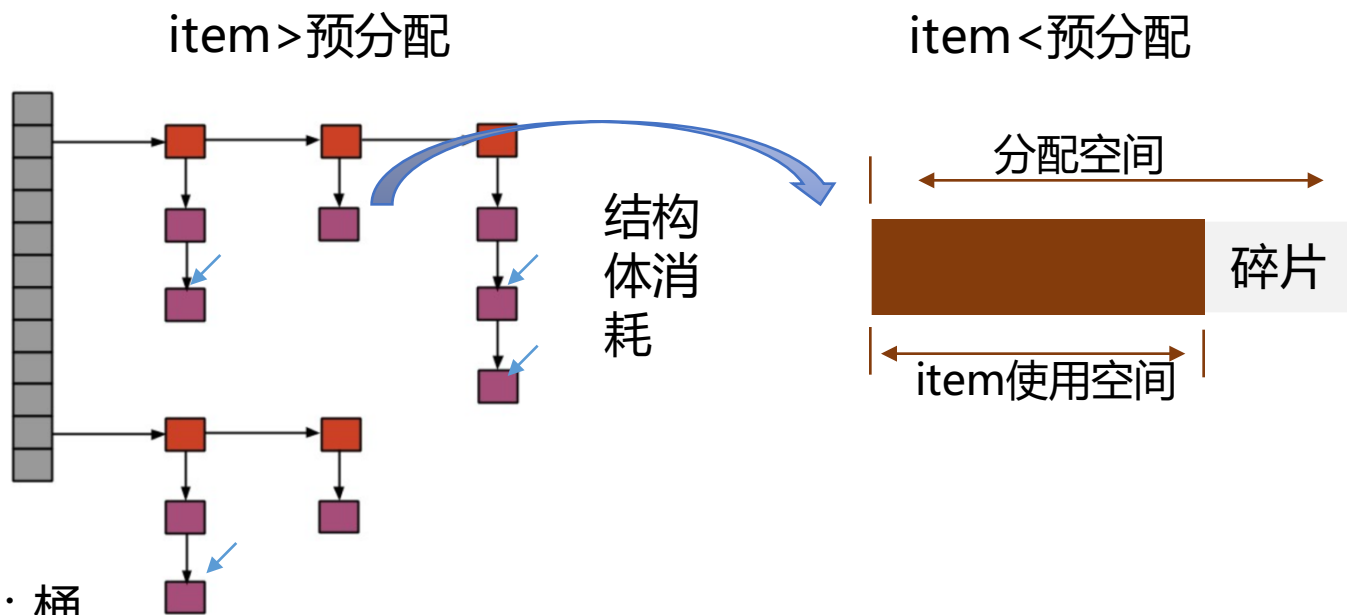
ret	key	stamp	lastread	lastupdate	expire
0:Success	31333034333832383532	90	1450270800 [2015-12-16 21:00:00]	1448089200 [2015-11-21 15:00:00]	NULL

key时间分布

业务访问时间
画像



碎片整理



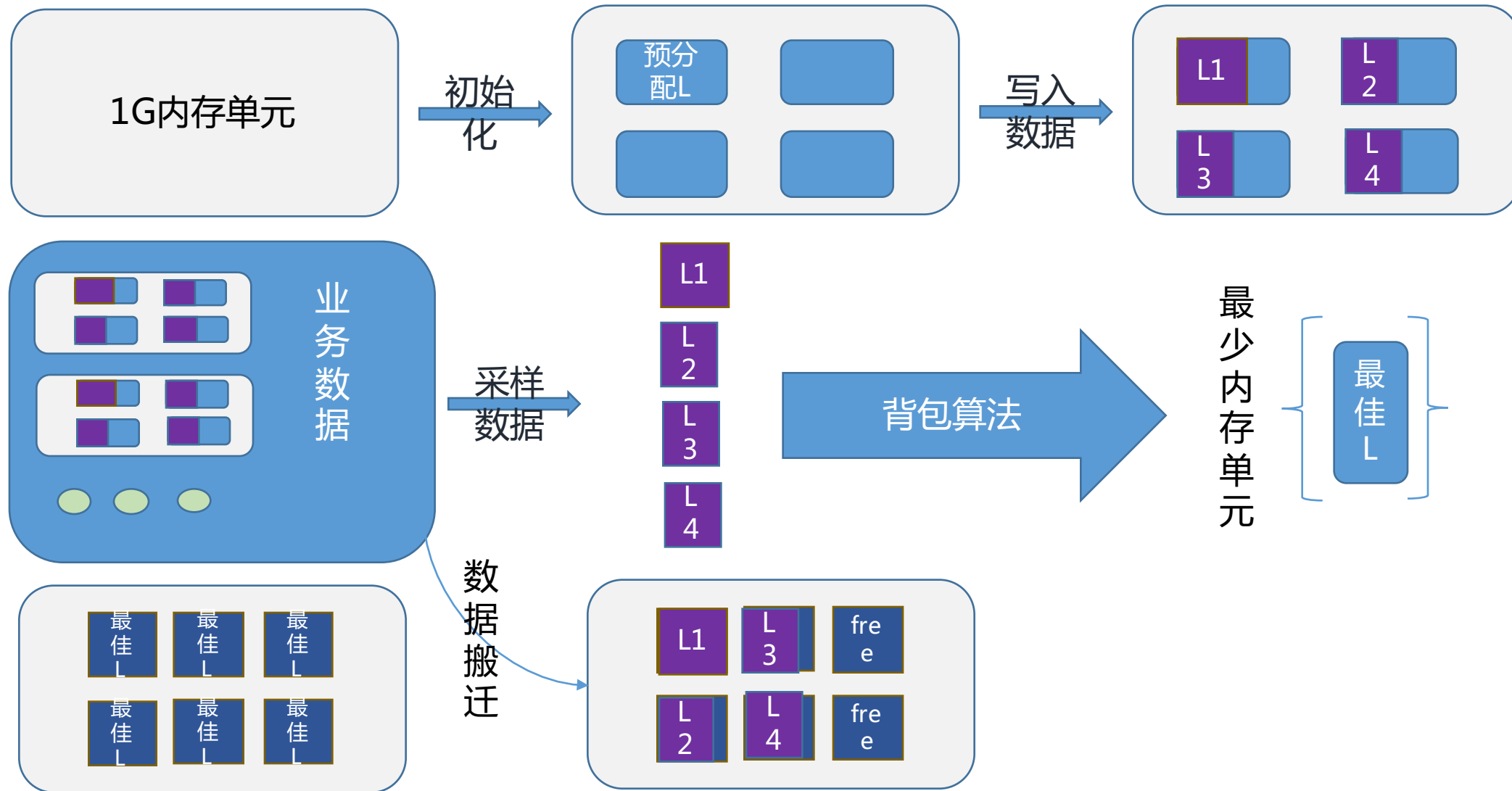
Bucket : 桶

Node : 节点, 对应一个key-value

Obj : 内存块, 大小固定的一块内存

- a)、灰色的对应hash的每个桶bucket
- b)、橙色的对应node
- c)、紫色的对应obj

碎片整理



问题定位

分钟粒度数据上报

Bid 101020401 2015-12-16 21:00:00 2015-12-16 23:00:00

- 注: 1. 下面3个表分别给出了单点源ip异常、单点目标ip异常以及源&目标ip联合异常
2. 若返回码为正均用1表示
3. 若数据库中查不到对应ip的出/入流量, 出/入包量用-1表示

请求端数据展示

源IP	返回码	最大延迟	异常数目	>50ms请求数
214	1	201.803	13	1302
217	1	315.472	13	1227
75	-13200	270.586	3	360
215	1	177.003	4	94
151	1	396.197	4	86

聚合结论引导

ip	数目	占比
最大时延ip(dst_ip)分析		
!5	516	36.86%

目标端数据展示

源IP	返回码	最大延迟	异常数目	>50ms请求数
220	1	186.381	2	1071
58	1	122.660	1	953
4	1	435.892	18	434
75	-13200	270.586	2	359
217	1	315.472	4	274

回顾

腾讯分布式KV存储介绍

- CKV实现

运营平台建设

- 鸚鵡螺

运维挑战与实践

- 自动化
- 成本优化
- 问题定位