

小米网redis服务平台 建设历程

徐成选@小米

2017年8月



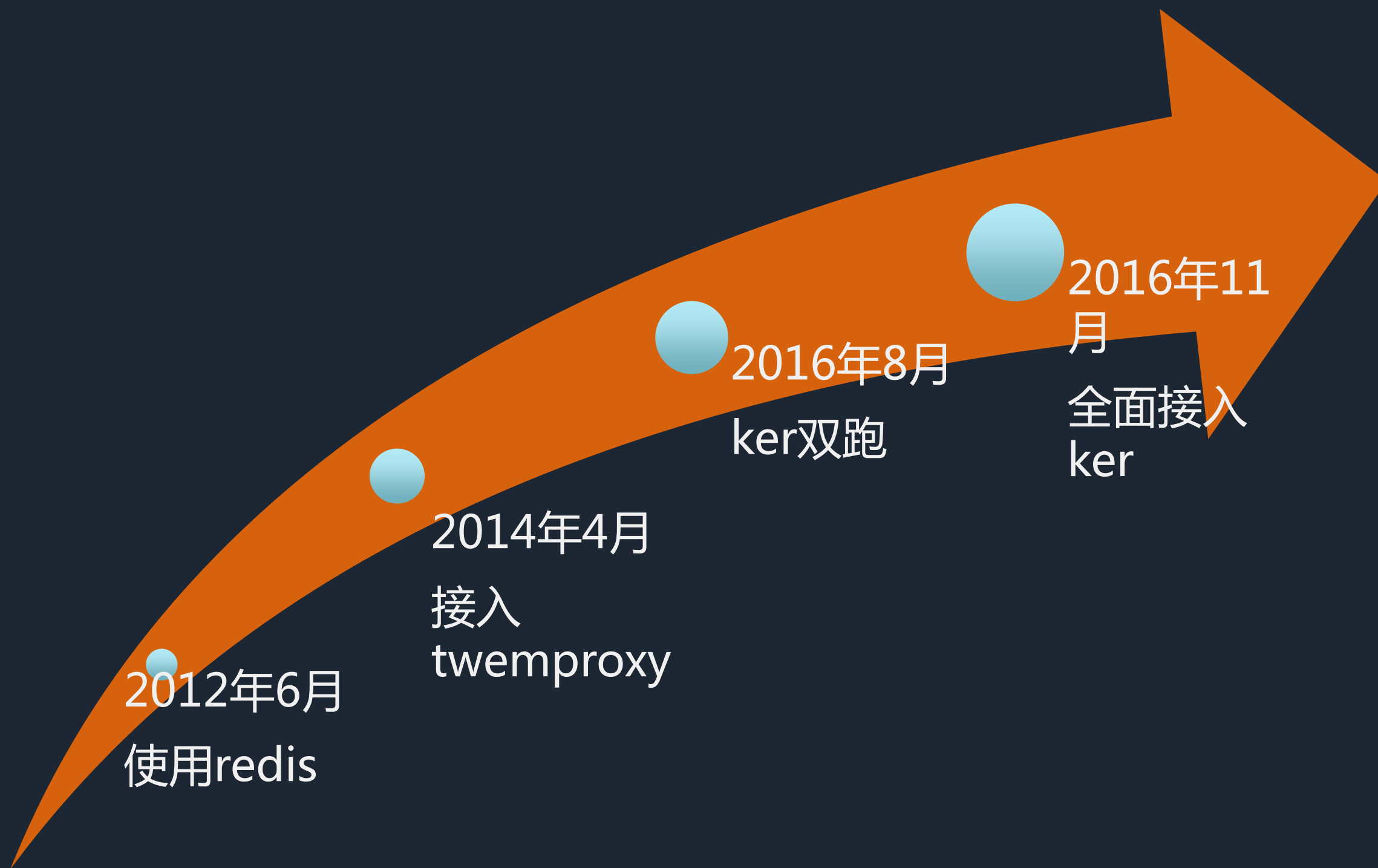
twemproxy面临的问题

KER平台设计考量

KER平台架构

管理平台

监控与报警





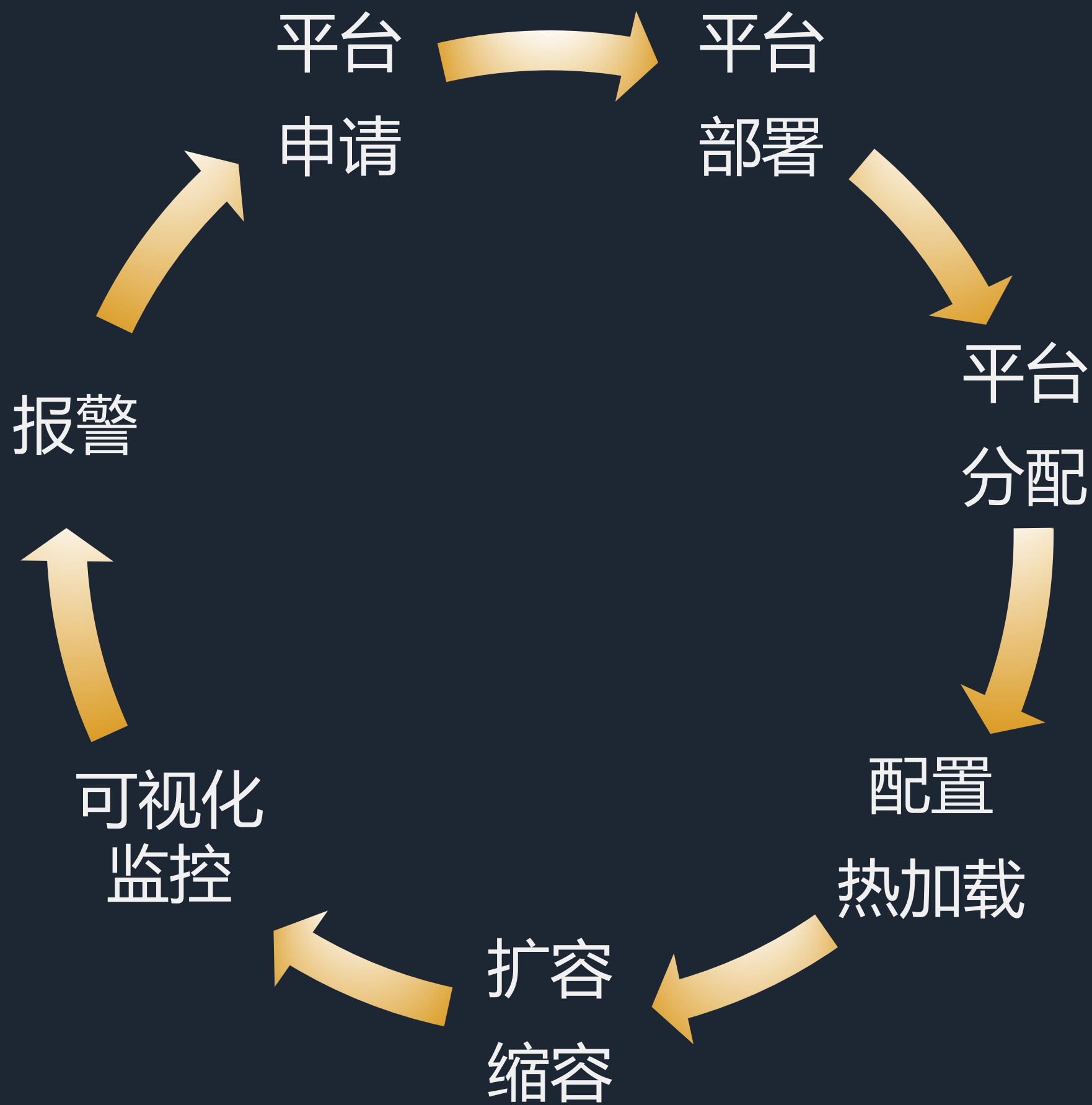
twemproxy面临的问题

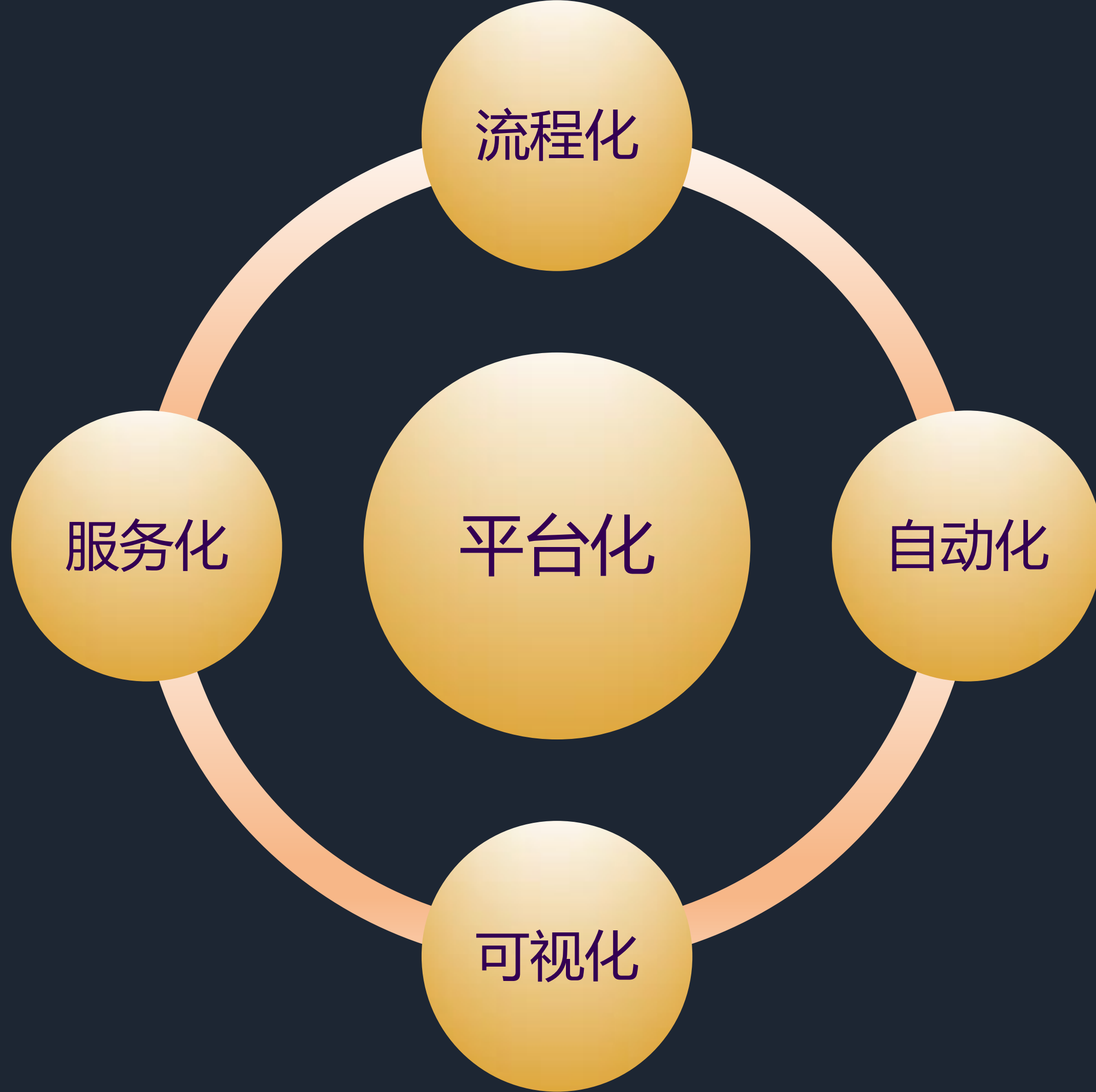
KER平台设计考量

KER平台架构

管理平台

监控与报警





twemproxy面临的问题

KER平台设计考量

KER平台架构

管理平台

监控与报警

KER迭代过程

KER1.0

- 兼容twemproxy
- 配置热加载
- 自助接入、平台运维
- 可视化监控

KER2.0

- 基于规则报警
- 动态扩缩容
- SSD
- 一键主从切换

KER3.0

- 自动化部署、配置
- 服务化



整体架构



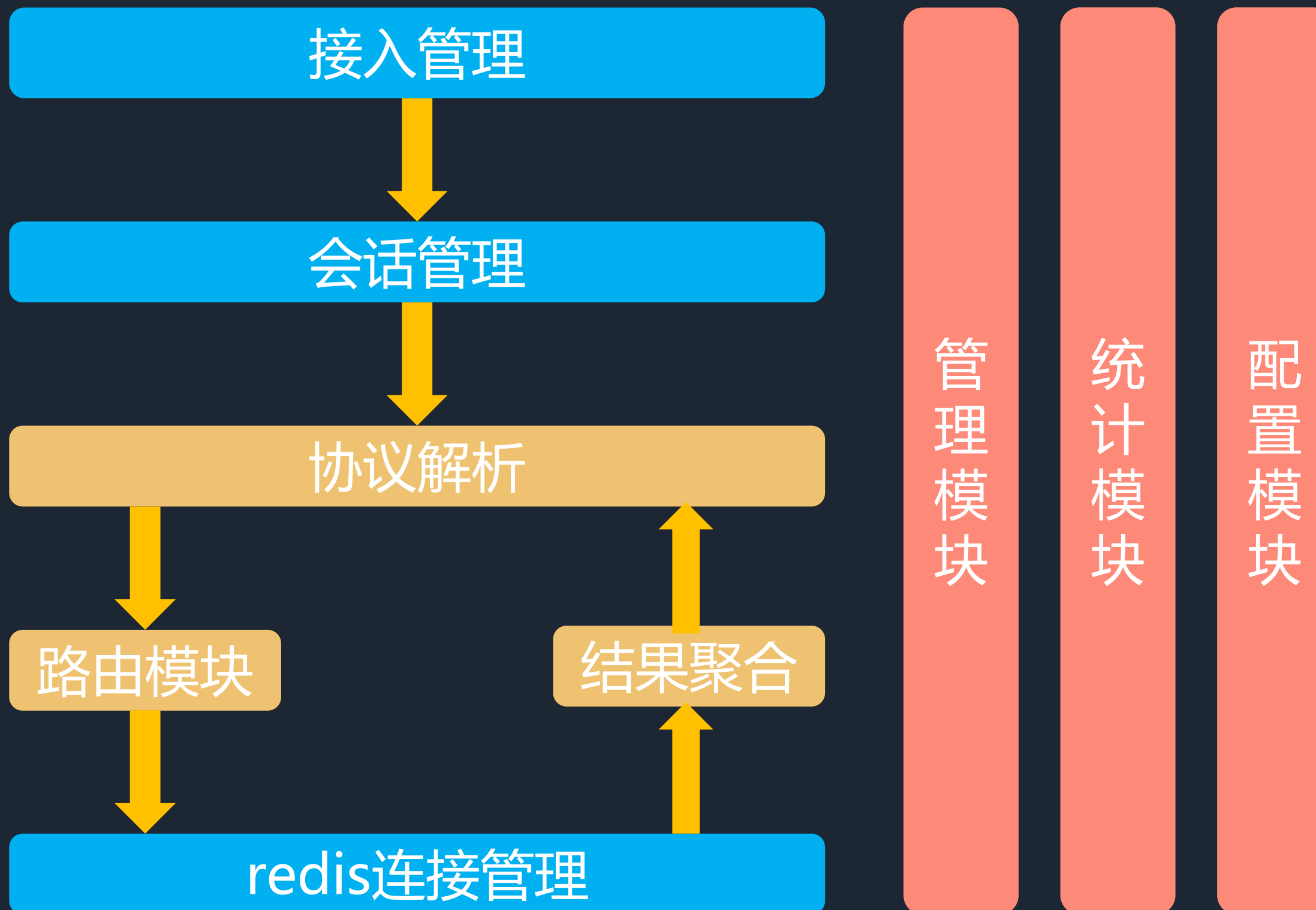
ker-proxy



只做最核心的功能



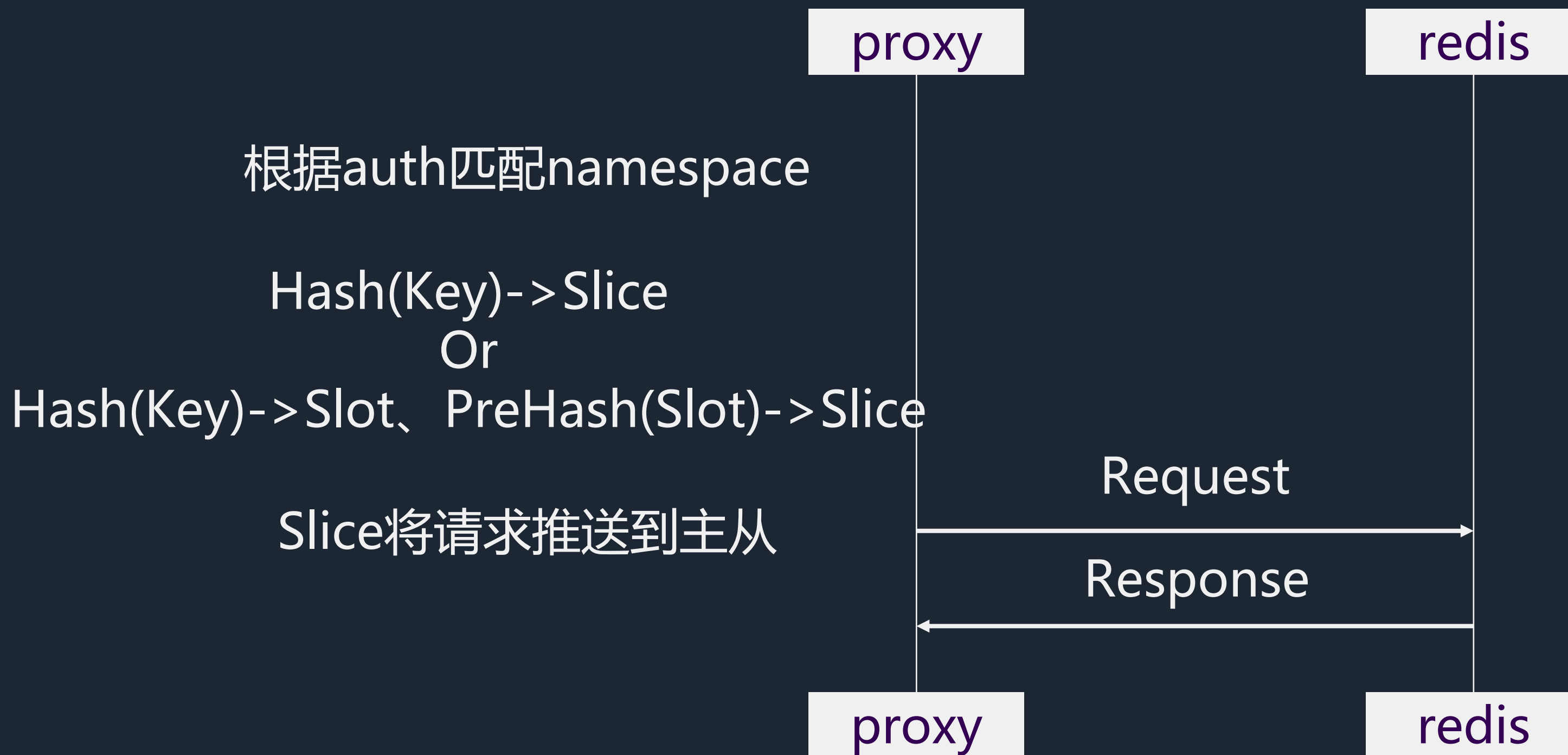
ker-proxy



ker-proxy



ker-proxy处理流程



ker-cc



全局掌控，干所有脏活、累活



etcd

实例管理

配置管理

后台任务



ker-agent



物理机实例部署、配置、启停



配置热加载

twemproxy

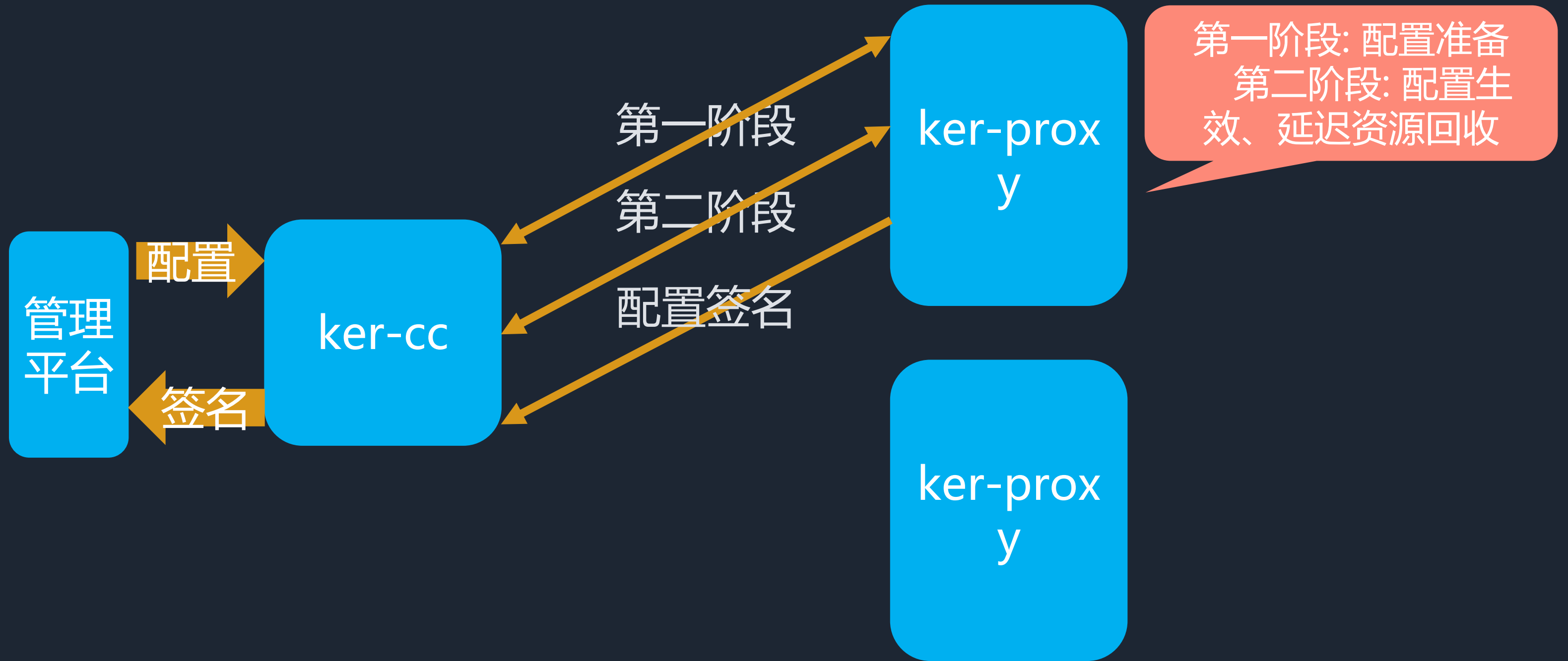
- 3000行的yaml文件
- 批量替换、批量重启
- 易出错

KER

- 配置按照namespace分别存储到mysql、etcd、本地缓存文件
- ker-proxy上报当前配置签名，校验配置变更是否生效、一致
- ker-cc校验配置准确性，通过两阶段提交的方式变更ker-proxy配置
- COW策略
- 第一阶段全量复制配置
- 第二阶段以namespace为单位，生效配置、延迟回收动态资源



配置热加载



动态扩缩容

Slice路由

- 翻倍扩容
- 提供清理工具，线下清理脏数据

Slots路由

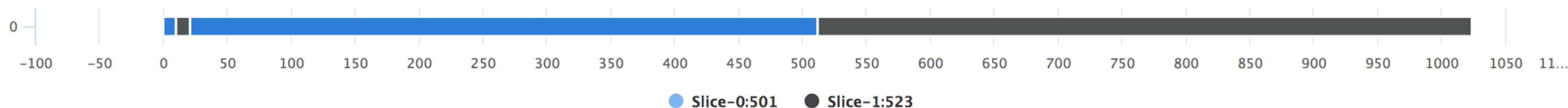
- 每个Namespace预分配1024个槽
 - DBA以槽为单位发起数据迁移
 - ker-cc以后台任务的方式不断迁移槽内数据
 - ker-proxy主动迁移槽内数据

动态扩缩容-平台

迁移Slots

设置参数 ~ to 迁移Slots

Slots分布



Slots信息

Show entries

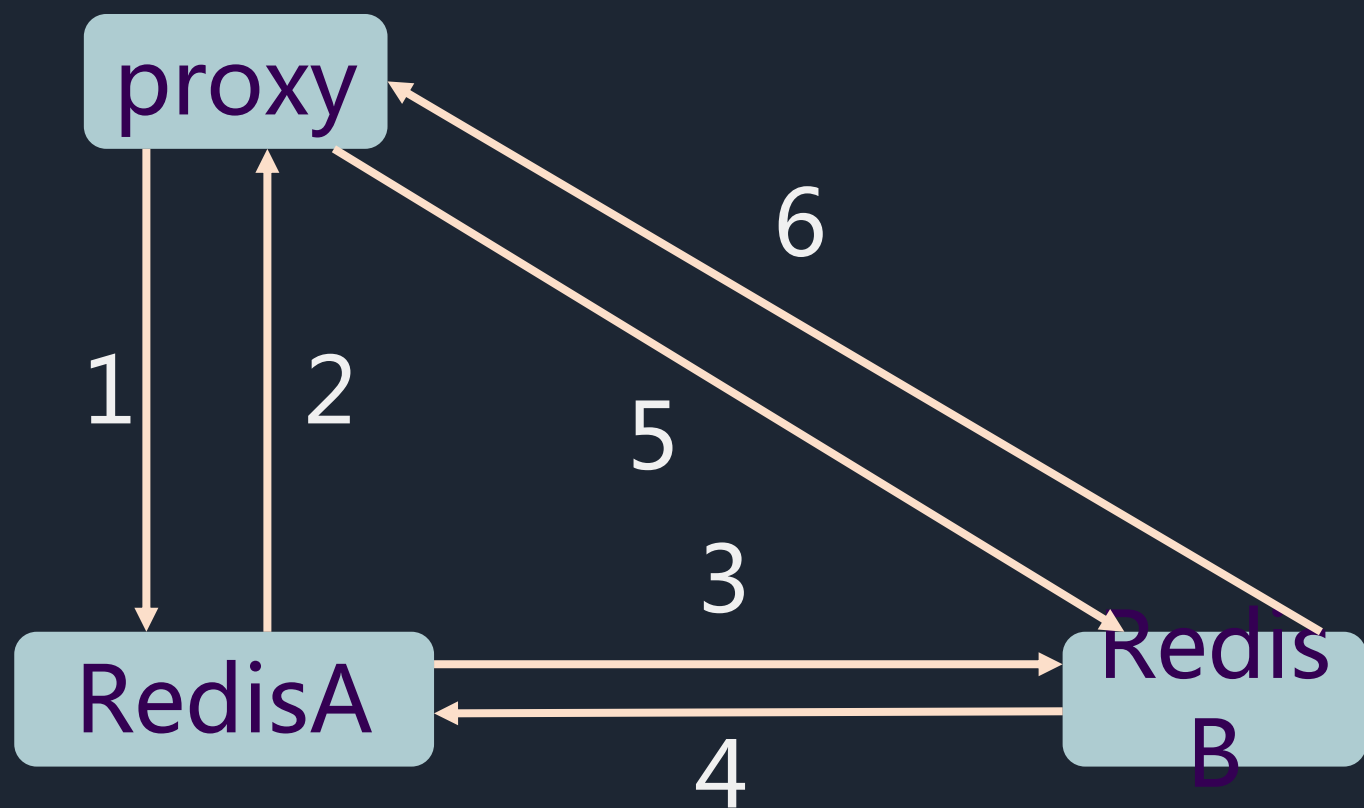
Search:

SlotId	SliceId	State	TargetSliceId
0	0	-1	-1
1	0	-1	-1

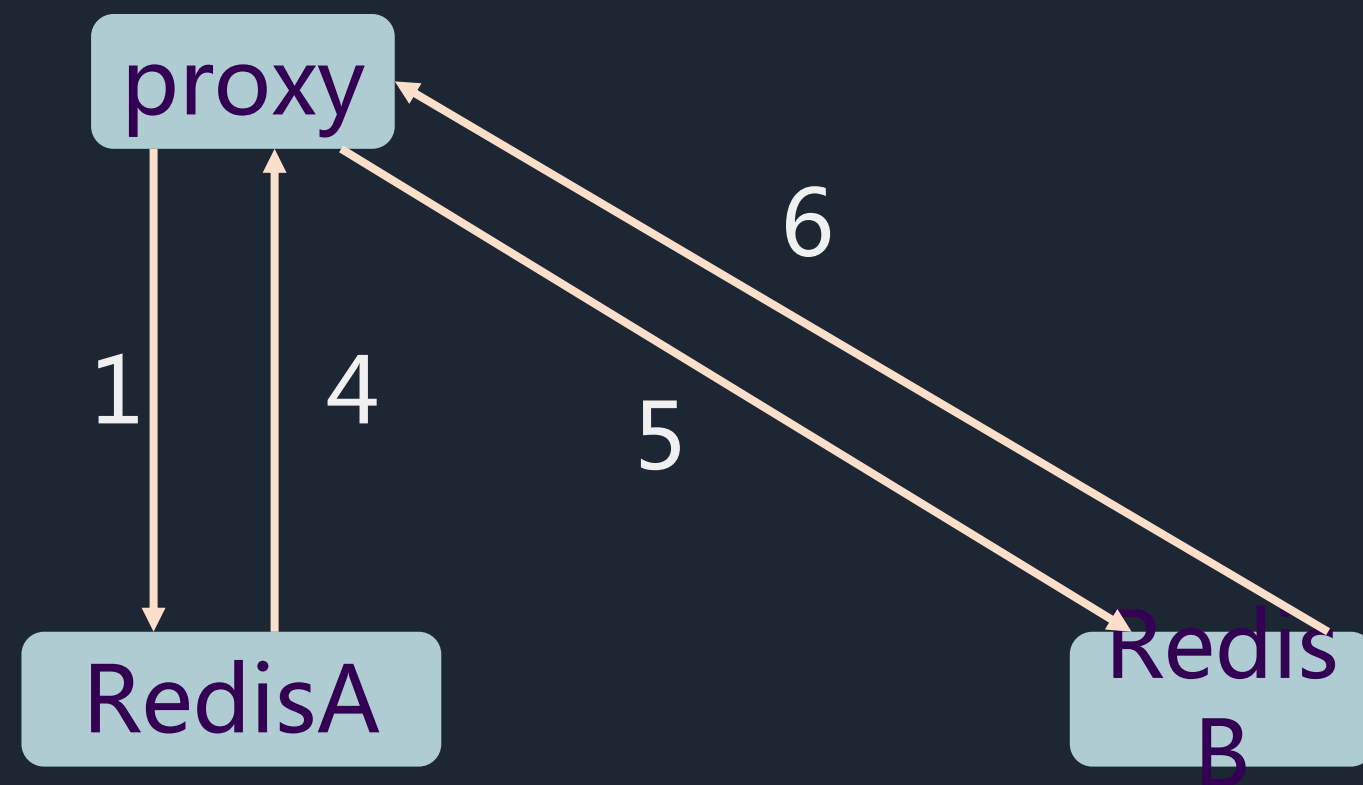
迁移信息

Slot	Slice	Target	Status
32	0	1	正在迁移
33	0	1	等待迁移
34	0	1	等待迁移

动态扩缩容-迁移



- 1 - SLOTSMGR TONE
- 2 - OK
- 3 - SLOTSRESTORE
- 4 - OK
- 5 - redis cmd
- 6 - OK



- 1 - SLOTSMGR TONE
- 4 - OK
- 5 - redis cmd
- 6 - OK

服务化

需求

- 基础服务接入微服务体系
- 应用系统跨过LVS直连ker-proxy

方案

- ker-proxy启动后注册信息到etcd、定时更新
- soa服务只需要配置auth即可通过xredis库访问redis
- 多个IP负载均衡、IP对应连接池延迟创建

服务化

首页 → 服务管理 → systech → ker_proxy → 实例列表

ker_proxy 实例列表 注册中心主机异常状态 绑定服务树/添加主机

IP	主机名称	机房区域	机房名	环境	程序md5	内存占用	日志级别	权限认证	启动时间	服务主机管理操作
10.236.120.9:10514	10.236.120.9	cn		测试环境				已关闭		重启 下线 cpu mem 火焰图 删除

共 1

```
//lvs
if p.kerDirect == 0 {
    reply, err = p.do(p.pool, cmd, args...)
    return
} else { //ker
    s, err := xconfig.GetService(KER_SERVICE)
    if err != nil {
        return reply, err
    }

    ip, port, err := xbalance.DoBalance(s)
    if err != nil {
        return reply, err
    }
}
```

```
p.kerLock.RLock()
pool, ok := p.kerPool[k]
if ok {
    p.kerLock.RUnlock()
    reply, err = p.do(pool, cmd, args...)
    return reply, err
}
p.kerLock.RUnlock()
p.kerLock.Lock()
```

```
pool, err = NewXRedisPool(p.name, ip, port, p.auth, uint(p.connTimeout),
    uint(p.readTimeout), uint(p.writeTimeout), p.maxOpen)
p.kerPool[k] = pool
p.kerLock.Unlock()
```



twemproxy面临的问题

KER平台设计考量

KER平台架构

管理平台

监控与报警

我的业务

认领redis

申请

业务方接入Ker-Redis代理

服务 (APP_ID)

请选择一个服务 (app_id)

机房

请选择一个机房(cluster)

类型

缓存

是否读写分离

使用读写分离:性能高, 但是牺牲一致性

命名空间*

命名空间

前缀:

所需存储空间

10

MB

一年后存储

20

MB

读 Qps

1000

次数

写 Qps

1000

次数

key-value对平均大小

100

Byte

申请描述信息



管理平台-DBA分配

236	系统组	auth_test_one	金山云测试	js_systech_test_slots12	申请中	上线	存储	slots	miredis	是	yangsi	2017-06-30 16:56:37	dba管理	删除	下线	修改	查看
242	数据组	redis_cost_sharing	大陆测试	cn_midata_cost_sharing	申请中	上线	缓存	slots	miredis	是	liuyifan3	2017-07-14 16:33:28	dba管理	删除	下线	修改	查看
248	门店组	xmstore_it	金山云测试	js_xmstore_test	申请中		缓存	slots	miredis	是	hanliguo	2017-07-27 11:39:47	dba管理	删除	下线	修改	查看
3	运维组	redis_test.damn	大陆测试	redis_test.damn	申请完成	上线	缓存	slice		是	zhanghua	2016-11-10 09:57:13	dba管理	删除	下线	修改	查看
4	运维组	redis_test.damn	大陆测试	cn_redis_test.damn	申请完成	上线	缓存	slice		是	zhangxionghu	2016-10-12 23:14:54	dba管理	删除	下线	修改	查看
6	系统组	ker_redis_all_test	大陆测试	xm_stock_aws	申请完成	上线	缓存	slice		是	zhangxionghu	2016-10-17 09:35:23	dba管理	删除	下线	修改	查看

服务 (APP_ID)	xuchengxuan_test_hello---系统组
机房	大陆测试
管理用户	xuchengxuan---徐成选
类型	缓存
所需存储空间	100 MB

[保存到db](#)
[添加slice文本](#)
[添加slice列表](#)
[重置表单](#)
[发送到ker中控](#)

命名空间:

Proxy Auth:

Redis Auth:

redis代理-管理信息

字段	命名空间	ProxyAuth	RedisAuth	HashTag	Readonly	读写分离
DB	cn_systech_cn_systech_bigpyer_test_1	cn_systech_cn_systech_bigpyer_test_1_1wwwVfn7nbp1N	6621		否	是
REMOTE	cn_systech_cn_systech_bigpyer_test_1	cn_systech_cn_systech_bigpyer_test_1_1wwwVfn7nbp1N	6621		否	是

redis代理-Slice列表

字段	slice_id	master	slave	slots	操作
DB	0	10.236.120.9:7480		[0,9] [21,29] [41,511]	修改 切换 配置 查看
REMOTE		10.236.120.9:7480			
DB	1	10.236.120.9:8480		[10,20] [30,40] [512,1023]	修改 切换 配置 查看
REMOTE		10.236.120.9:8480			





twemproxy面临的问题

KER平台设计考量

KER平台架构

管理平台

监控与报警

监控大盘



ker-proxy



Zoom Out

Last 30 minutes



集群 MAX QPS

733.74 K

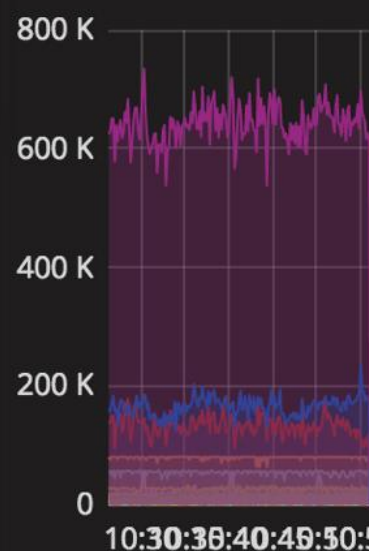
集群网络负载MAX

276.4 MBps

集群平均响应时间

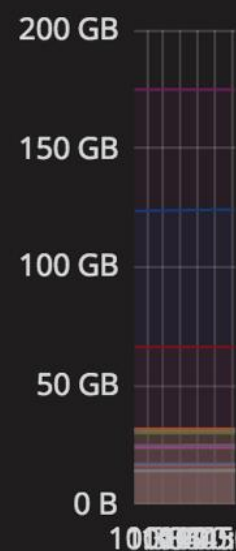
245 μs

业务QPS



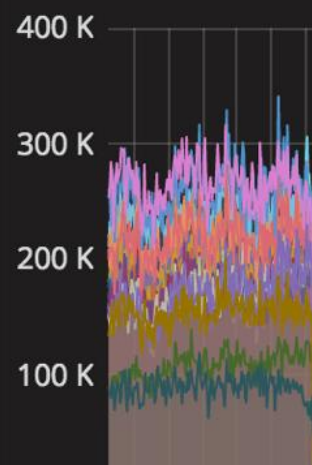
	max	avg
总 QPS	734 K	639 K
cn_b2c_calc_uservisit	232 K	165 K
b2ccalc_016_tmp	178 K	132 K
systemtech_union_cache	81 K	79 K
cn_shield	58 K	54 K
xm_stock_cn	56 K	17 K
kfs_mic	32 K	26 K
b2ccalc_incl_storage	32 K	23 K
snodex_common_session	25 K	7 K

业务内存使用情况



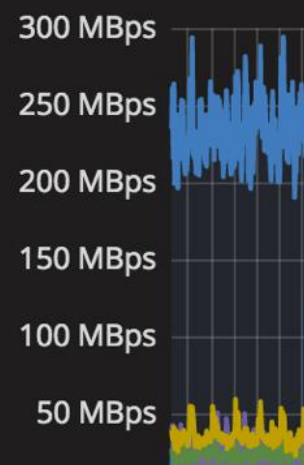
	max	avg
pmc_cart	174.7 GB	174.7 GB
cn_b2c_calc_uservisit	124.3 GB	124.0 GB
order_risk_ipdata	66.3 GB	66.3 GB
weixin_platform	32.0 GB	32.0 GB
misite_userset_damiao	31.6 GB	31.6 GB
delivery_address	31.1 GB	31.1 GB
xmpro_new	30.1 GB	30.1 GB
b2ccalc_016_tmp	25.0 GB	24.9 GB
systemtech_nginxspark_streaminglogstat	24.0 GB	24.0 GB

后端 QPS



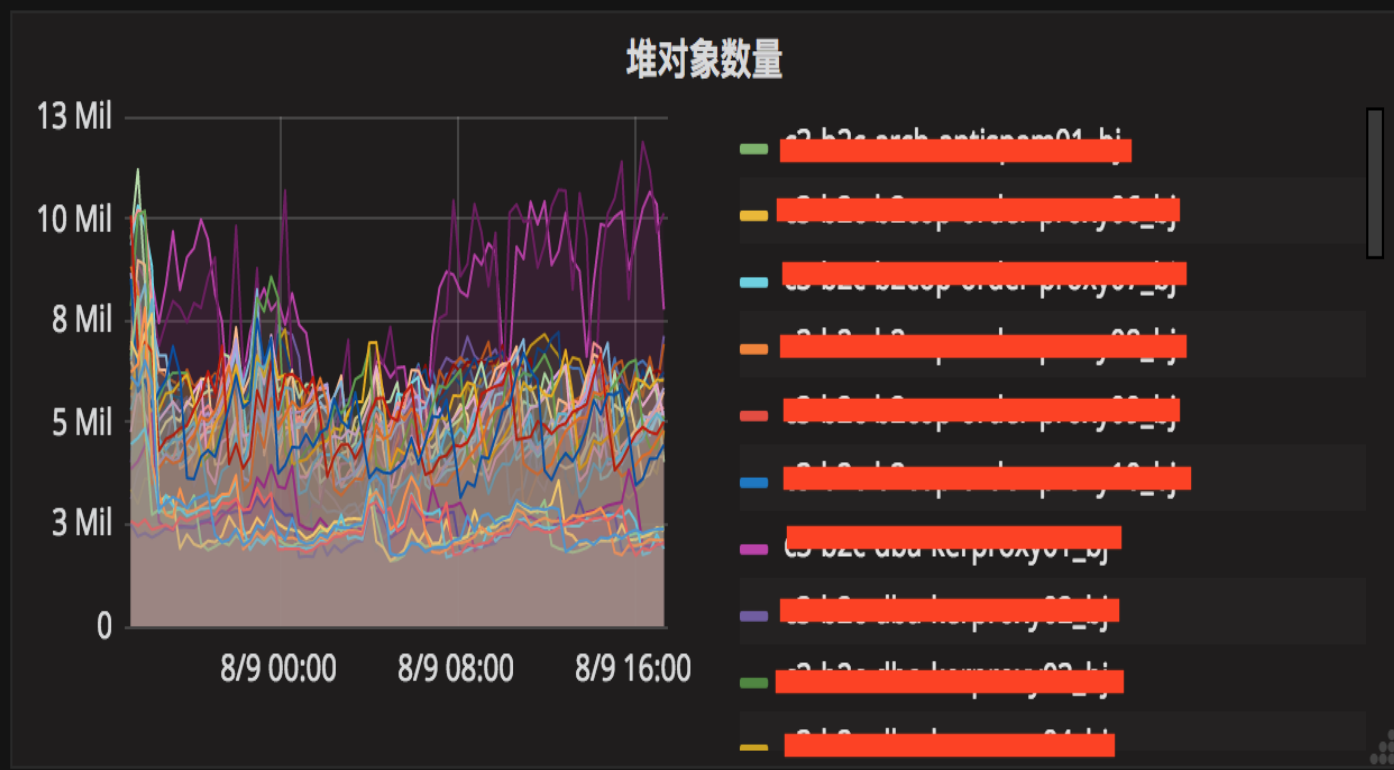
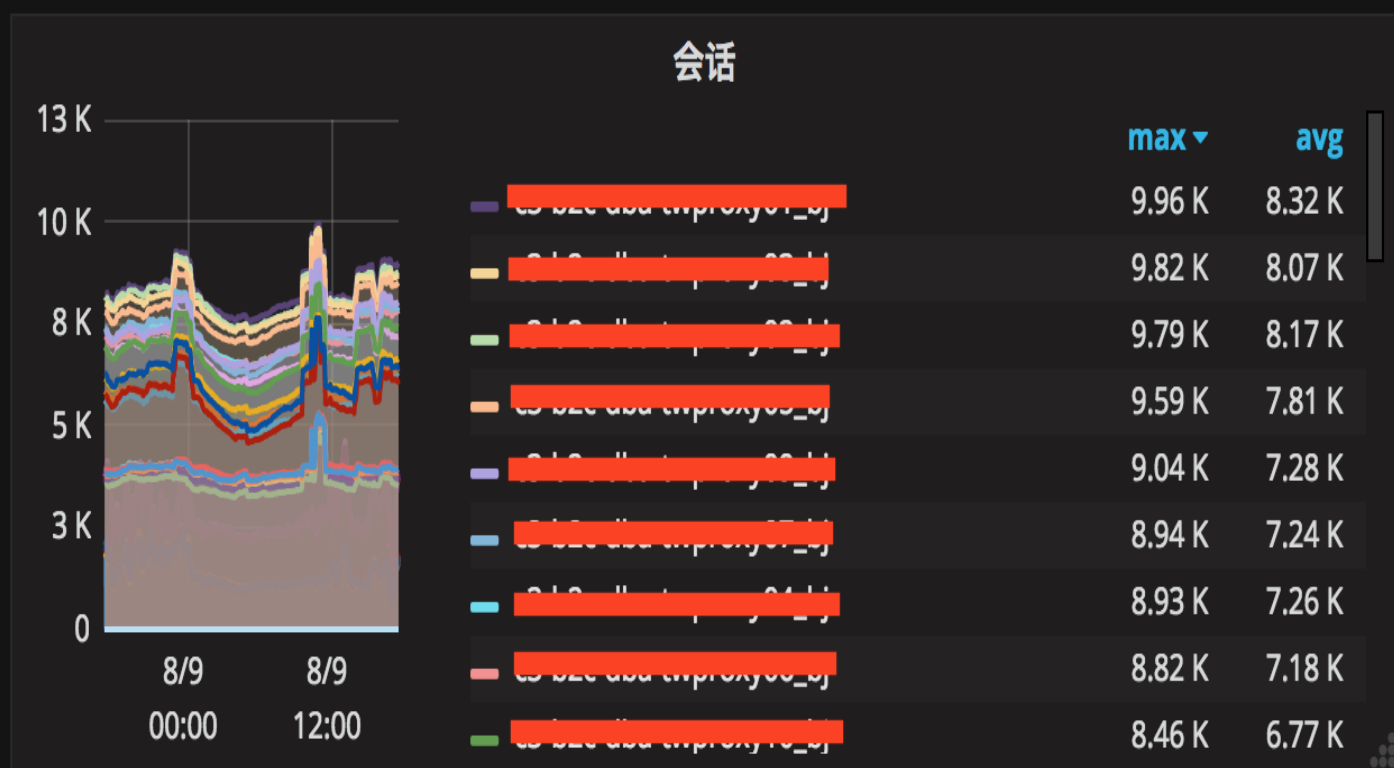
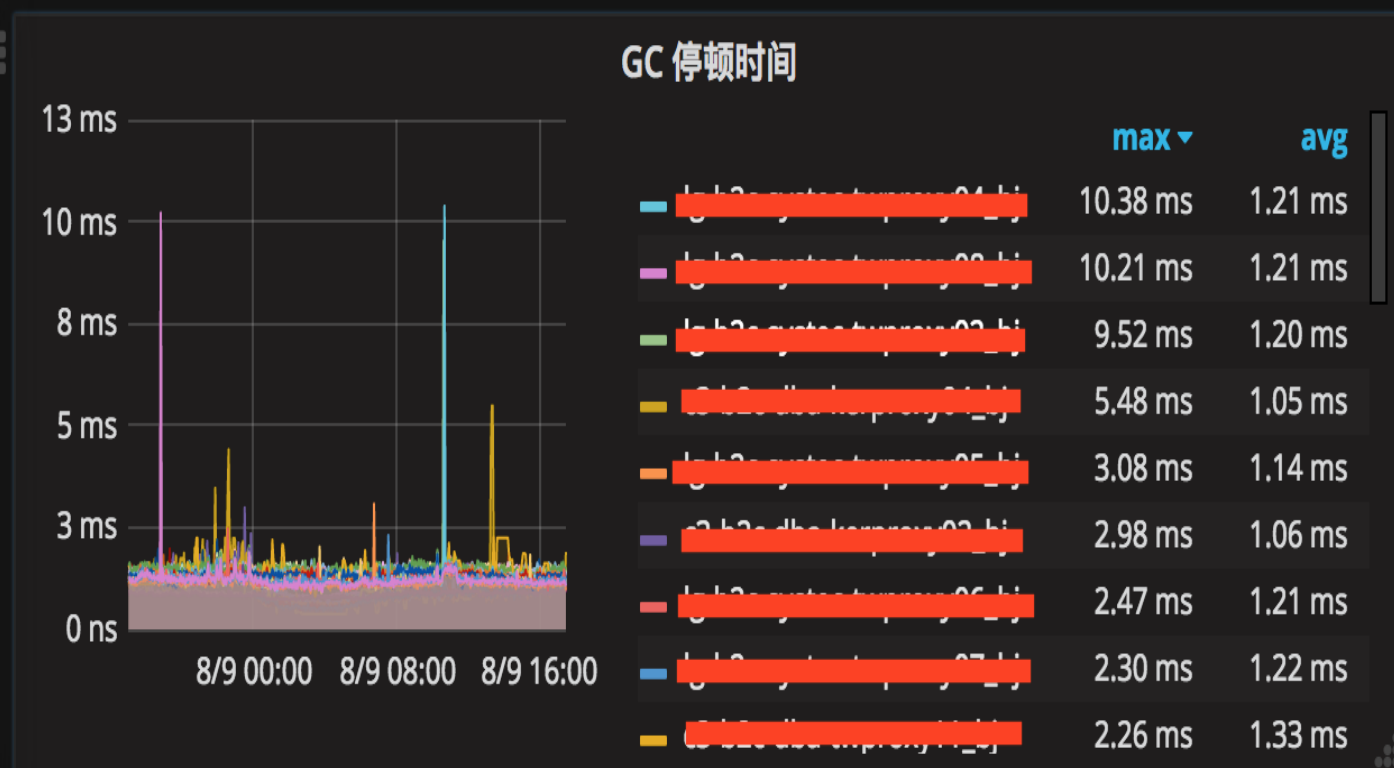
	max	avg
sgpaws_b2c_systemtech_cached_redis06_miao		
sgpaws_b2c_systemtech_cached_redis05_miao		
sgpaws_b2c_systemtech_cached_redis04_miao		
sgpaws_b2c_systemtech_cached_redis03_miao		
sgpaws_b2c_systemtech_cached_redis02_miao		
sgpaws_b2c_systemtech_cached_redis01_miao		
CS-b2c-systemtech-cache20_bj		

业务流量

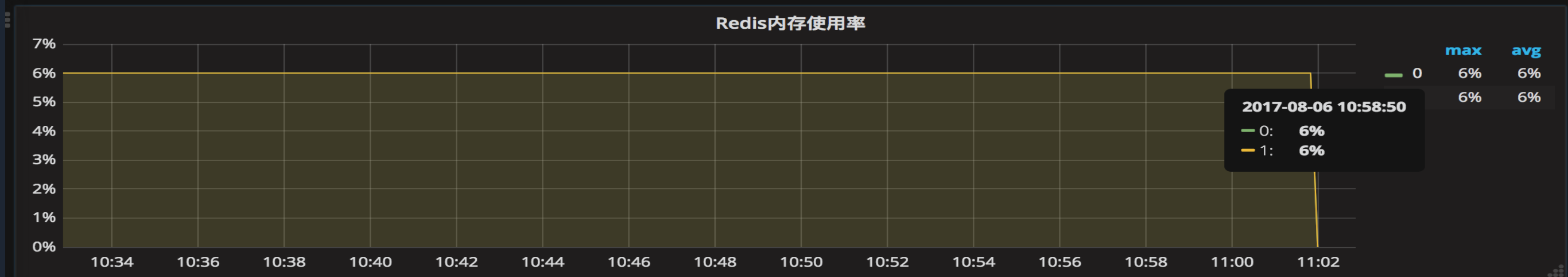
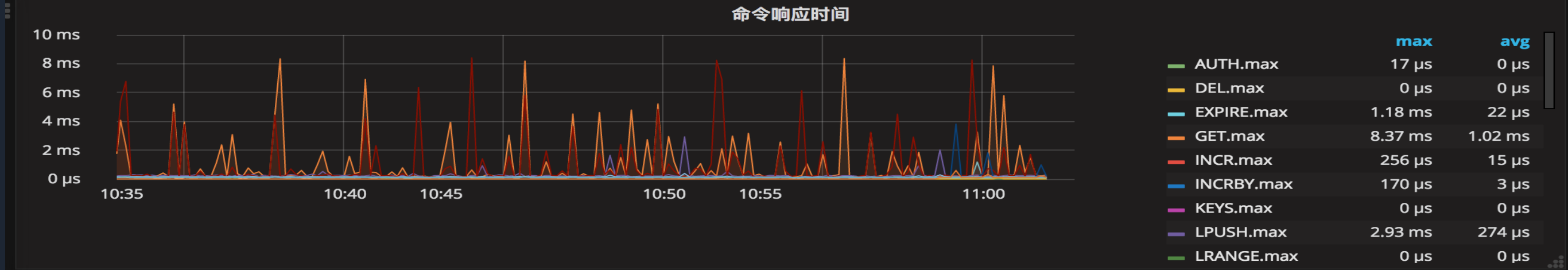
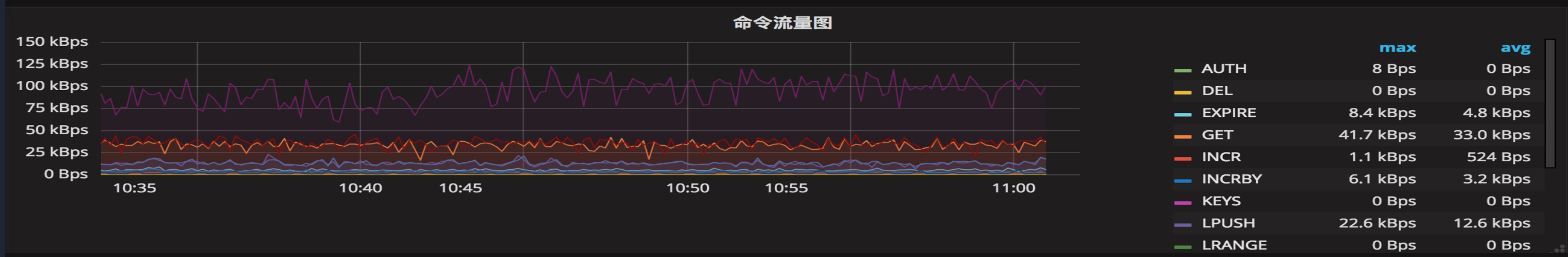
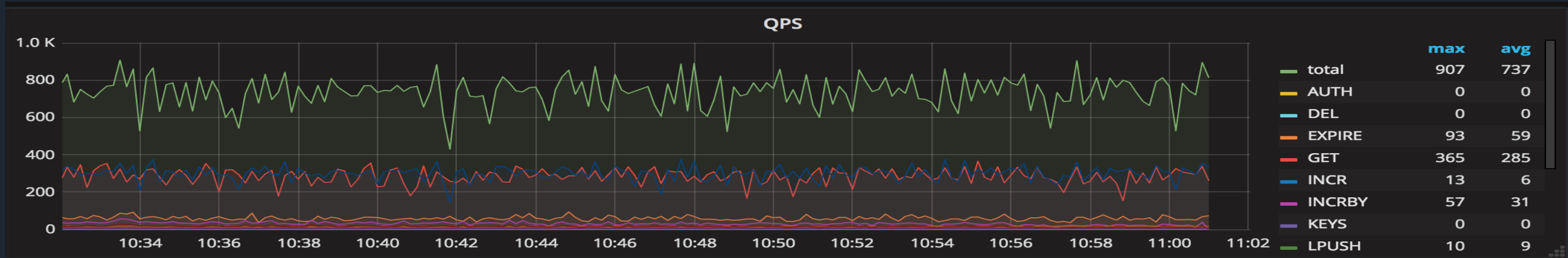


	max	avg
总流量	294.1 MBps	231.9 MBps
mifa_chopapi_miao_new	59.2 MBps	35.6 MBps
mibbs_cache	50.4 MBps	22.4 MBps
xms	40.1 MBps	23.2 MBps
cn_b2c_back_api_d_b2c_srv	33.5 MBps	10.9 MBps
midrop_oms	29.0 MBps	16.5 MBps
deconnectivity	21.1 MBps	19.7 MBps





监控-某业务



规则ID	规则名称	数据源	Topic	Grafana名称	Grafana图表	合并策略	发送策略	优先级	逐级上报	状态	操作
1	内存阈值-xm_stock- [REDACTED]	KAFKA	b2c_ker_cc_applogs_error			endpoint_metric	notify1	P2	否	启用	禁用 修改 删除 异常信息
3	实例存活检测-proxy- [REDACTED]	KAFKA	b2c_ker_cc_applogs_error			endpoint_metric_priority	notify7	P0	否	启用	禁用 修改 删除 异常信息

数据源: KAFKA

Topic: b2c_ker_cc_applogs_error

合并策略: endpoint_metric

发送策略: notify1

优先级: P2

报警接收人: xuchengxuan---徐成选

逐级上报: 否

+ 添加规则

Namespace: =(字符串) xm_stock [删除](#)

and Host: =(字符串) [REDACTED] [删除](#)

and MemRate: >=(数字) 98 [删除](#)

元字段: x Namespace=xm_stock x Host=[REDACTED]

备注: 内存阈值-xm_stock-[REDACTED]

[提交](#) [重置](#)



商城系统组

- 分布式缓存
- 消息中间件
- 数据库中间件
- 微服务平台、框架
- 日志处理、挖掘
- 监控、报警等

Thanks

