



第九届中国数据库技术大会
DATABASE TECHNOLOGY CONFERENCE CHINA 2018

从零到壹 我们是这样打造高可用公有云 Redis服务

腾讯 高级工程师 冯伟源

DTCC
2018

2018.05.10 – 12 北京国际会议中心



IT168.com

ChinaUnix

ITPUB

个人经历



平安

数据库工程师 (Oracle)
保险银行投资等金融业务支持
SQL与实例优化、迁移备份监控

2012

唯品会



数据库工程师 (MySQL)
订单商品网站业务DB支持
主流开源技术与架构

唯品会
vip.com

2014



腾讯

高级数据库工程师 (Redis)
海量NoSQL运营保障
自动化运维开发、运营规划

2015



目录 contents

01

方案简介
PRODUCT DESCRIPTION

02

架构设计
ARCHITECTURE DESIGN

03

运营系统
SUPPORT SYSTEM

04

运营思考
OPERATE THINK

PART 01

方案简介

DTCC
2018

数领先机 智赢未来 (9)

IT168.com

ChinaUnix

ITPUB

Redis外部环境

人类活动

海量并发

NoSQL

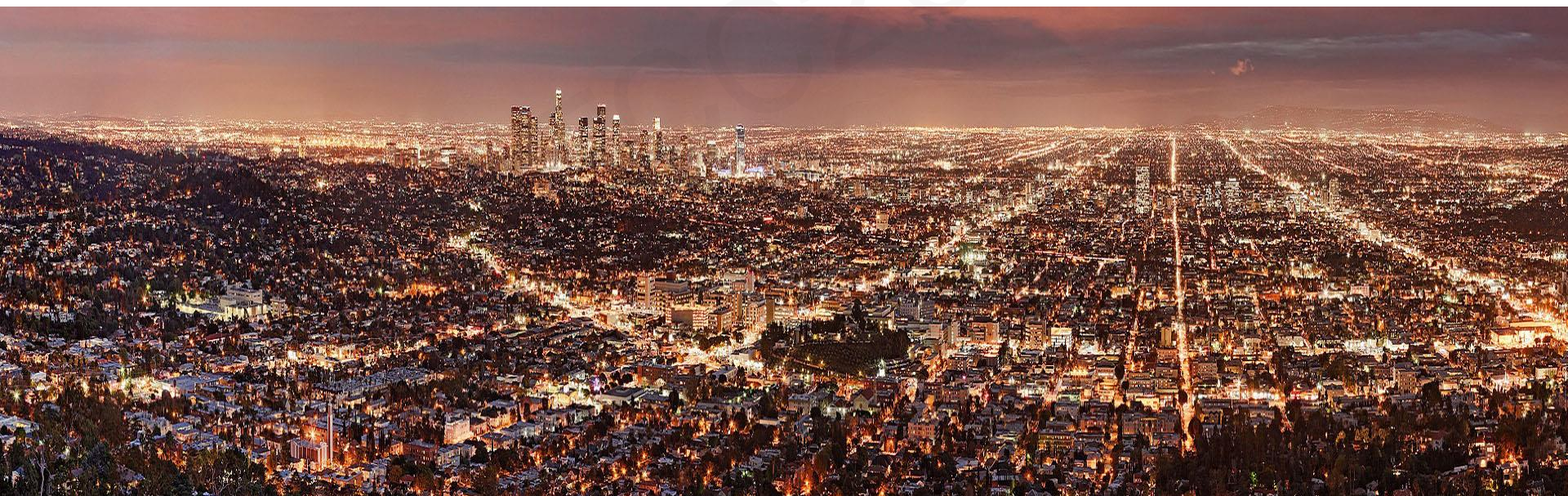
Redis

远程字典服
务器

简洁、极
致，高效、
开源

结构丰富、
生态活跃

2018 DB-
Engine KV
类第二





托管部署



数据保障



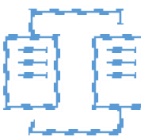
平滑拓展



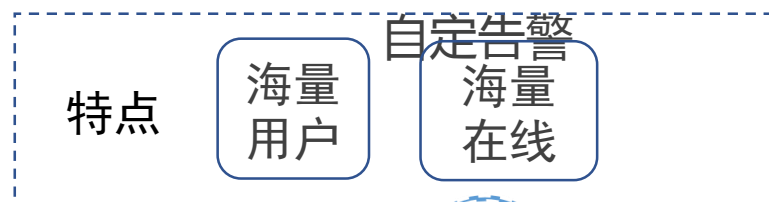
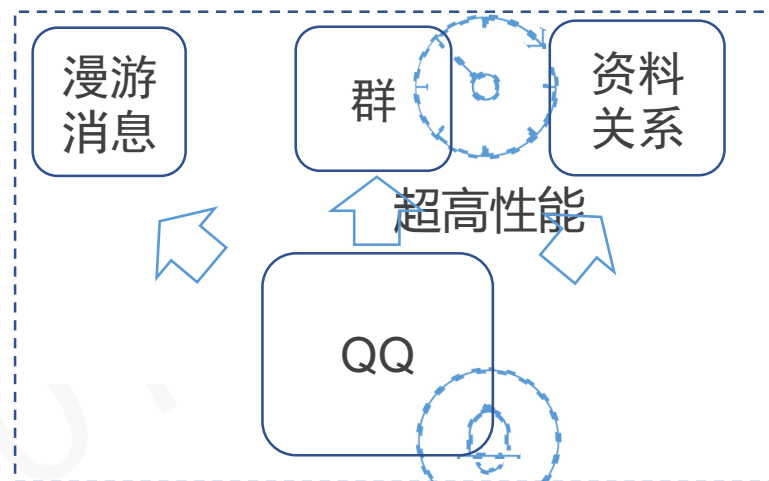
全面监控



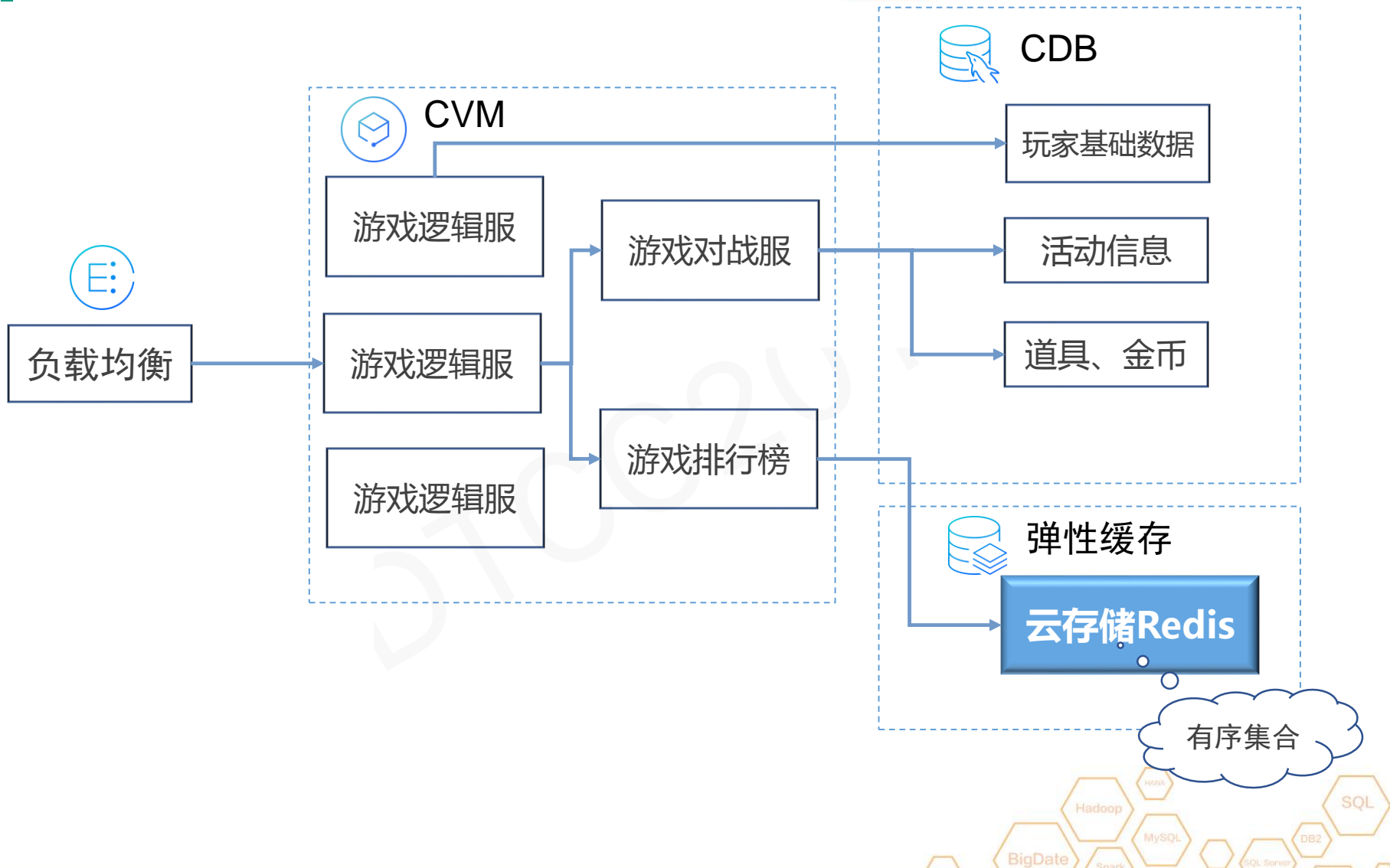
数据迁移



自动容灾



应用场景-游戏





小红书

mobike
摩拜单车

富途证券
FUTU5.COM



WeBank

微影时代
北京微影时代科技有限公司

招联金融
MUCFC.COM



Kingdee

猎豹移动

ANSWERN
安心保险



QQ音乐

QQ空间
分享生活 留住感动

微云

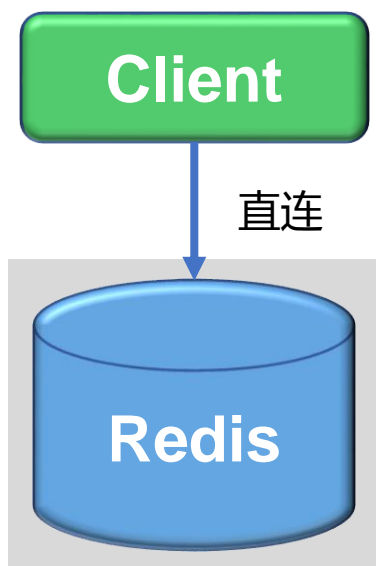
PART 02

架构设计



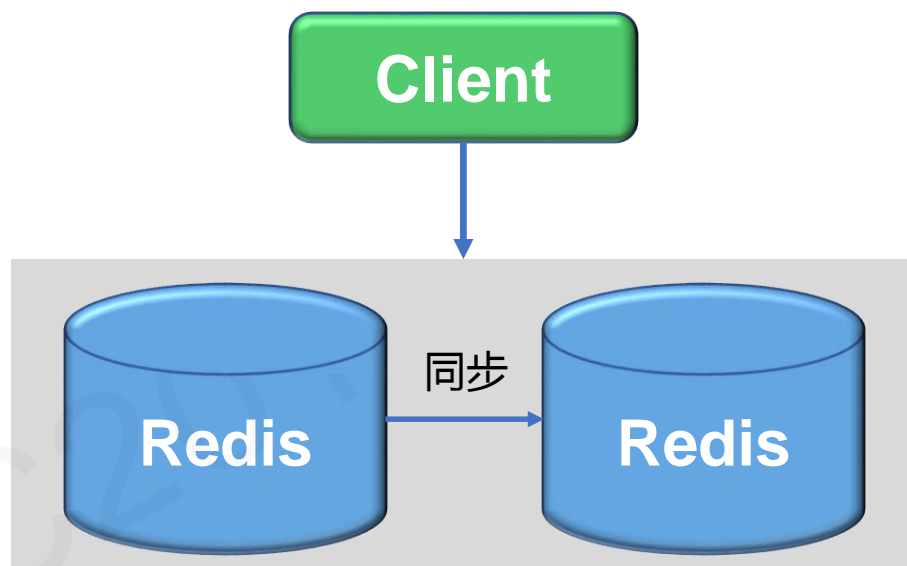
有哪些传统的Redis使用方法？





可靠性

- 数据易丢失
- 服务不可靠



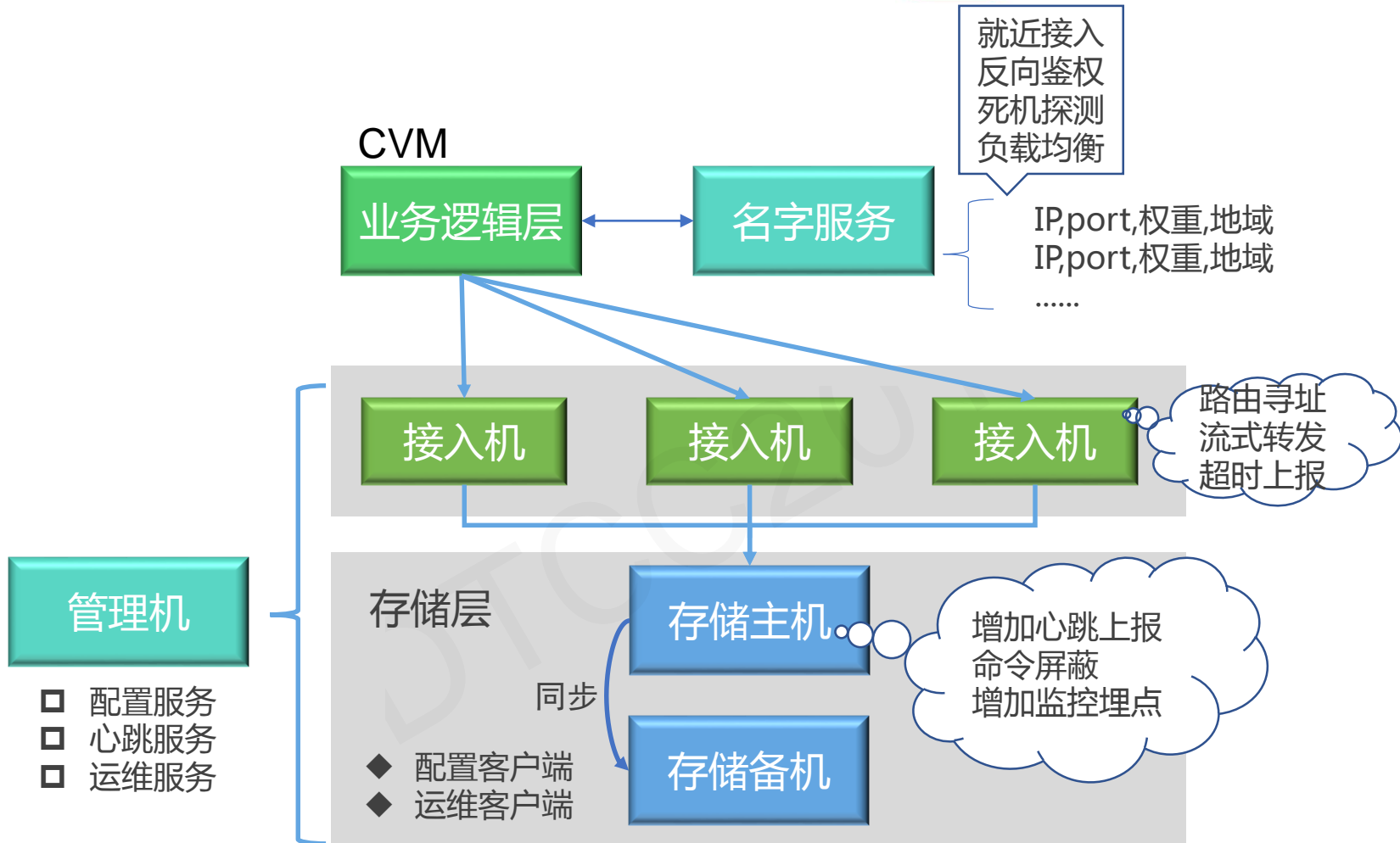
自动化

- 无法自动容灾切换



自动容灾切换-CRS主从版

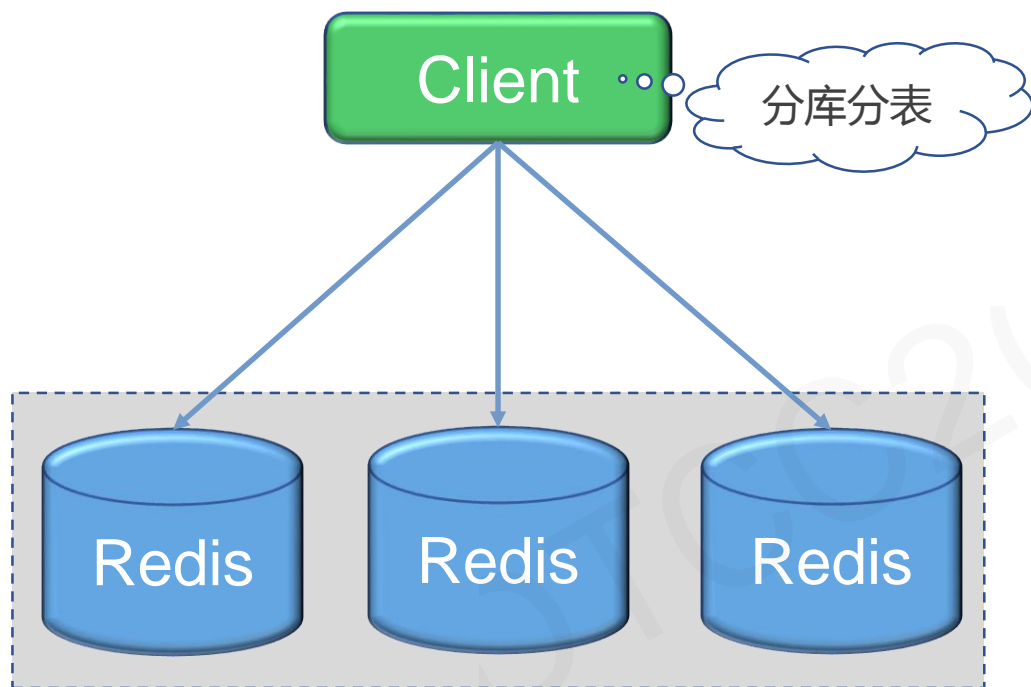






容量与性能的扩展-CRS集群版





- 特点
 - ✓ 多套实例
 - ✓ 客户端实现分片逻辑

- 不足
 - ✓ 开发者分片麻烦
 - ✓ 静态分片，缺乏动态扩容能力
 - ✓ 可运维性差
 - ✓ 不支持跨Key命令和事务

其他集群实现方案

Twemproxy

Codis

Redis Cluster

持久化时质量抖动

内存碎片内存溢出

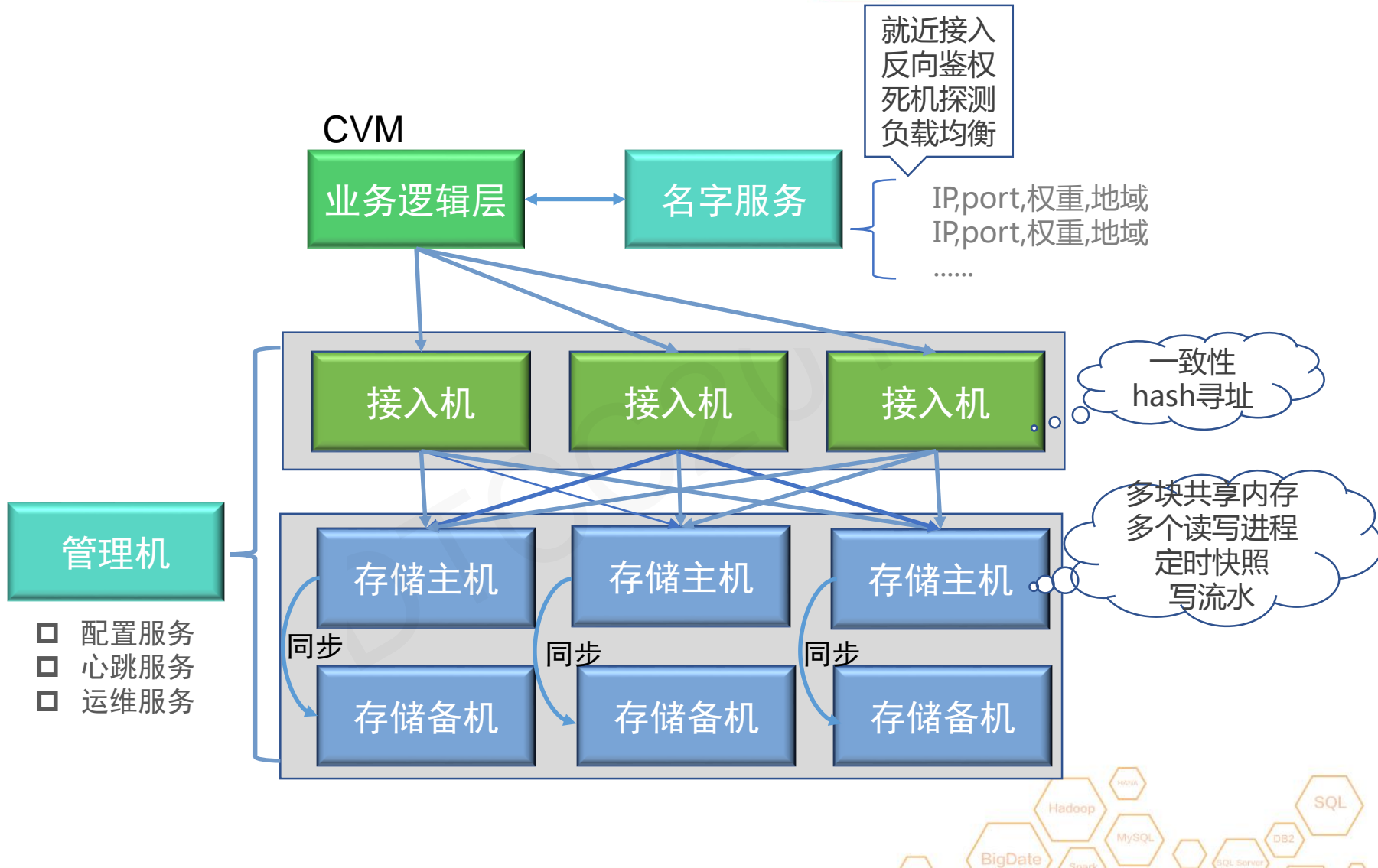
存活误判监控告警

备份回档流水审计

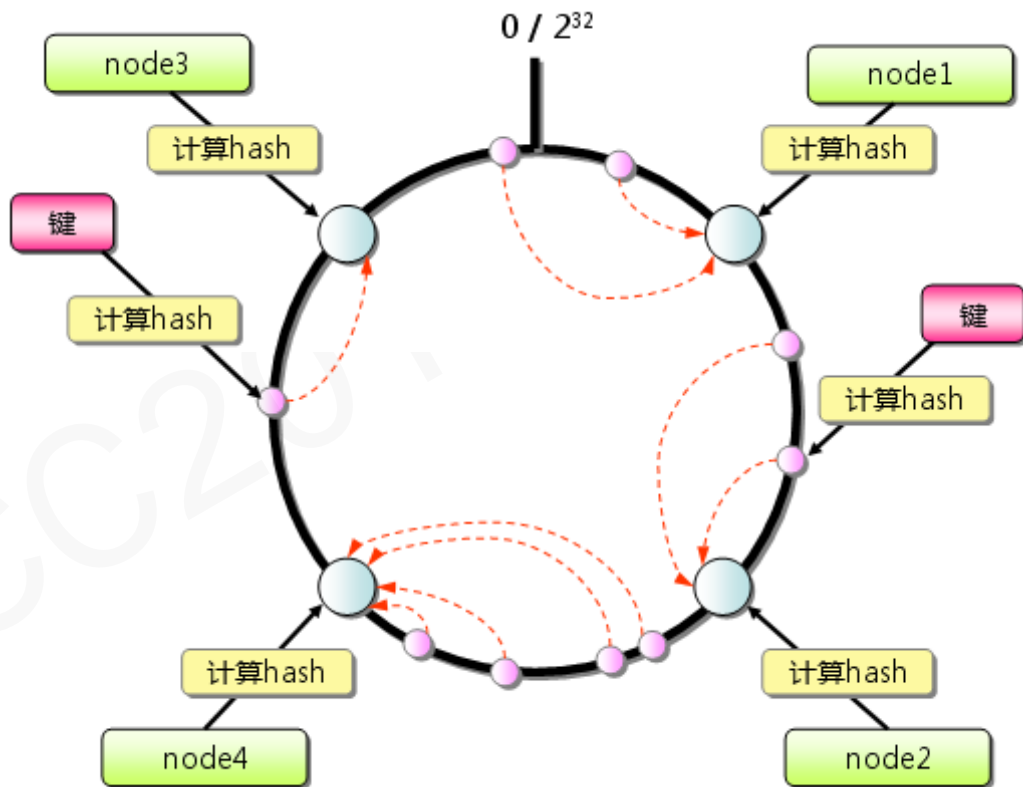
数据安全权限分级

多租户资源隔离差

CRS集群版



- 物理机->虚拟节点
- 按hash(key)寻址
- 一致性hash
 - 自动计算配置
 - 无限平行扩展
 - 打散热点
 - 迁移最小化
 - 数据均匀分布($\pm 3\%$)



私有内存

重启数据释放

COW机制

预留两倍内存

页表复制带来抖动

VS

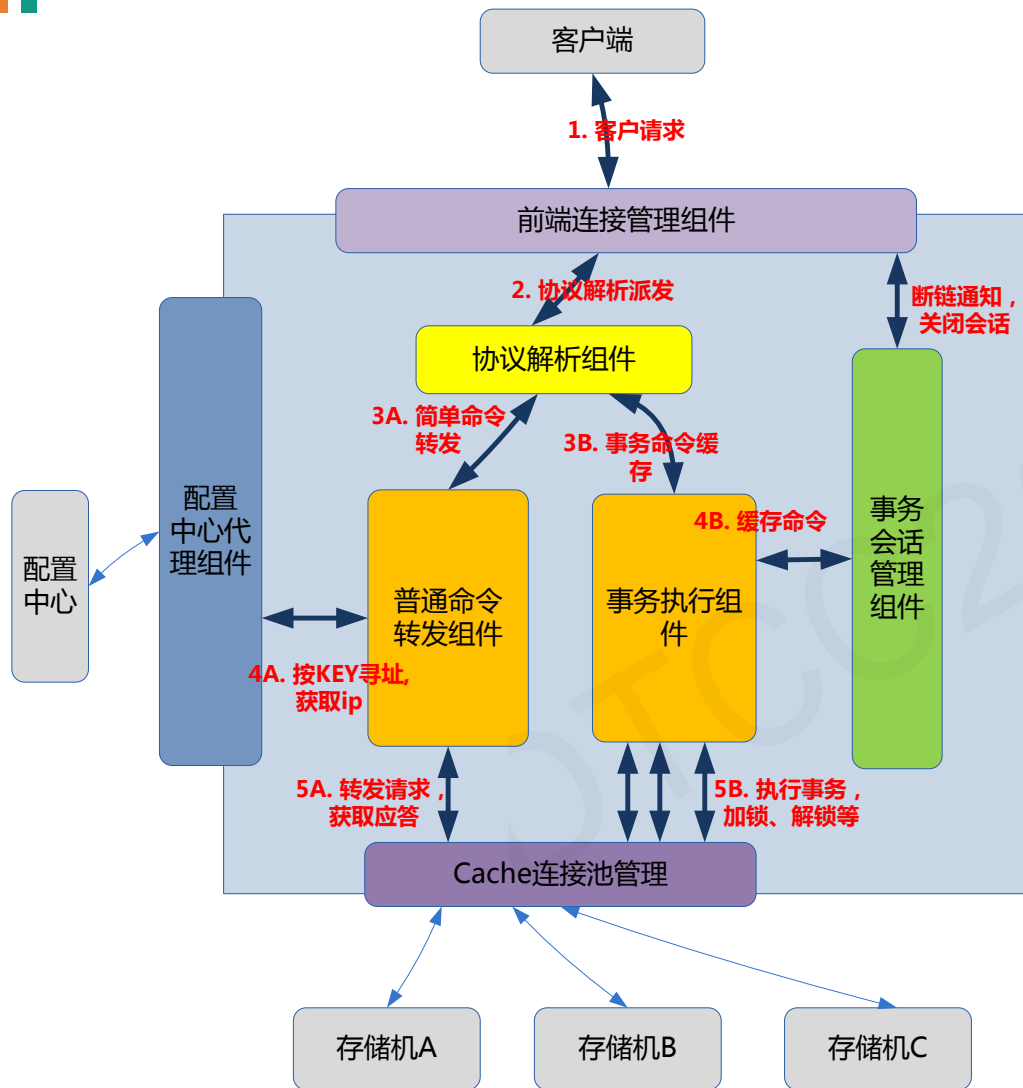
共享内存

数据安全与版本发布更容易

遍历+流水

占用内存更少

不存在页表复制



- 命令队列
- 路由功能
- 流式转发
- 链接池
- 命令统计

基本

- Auth鉴权
- 命令过滤
- 死机探测

安全

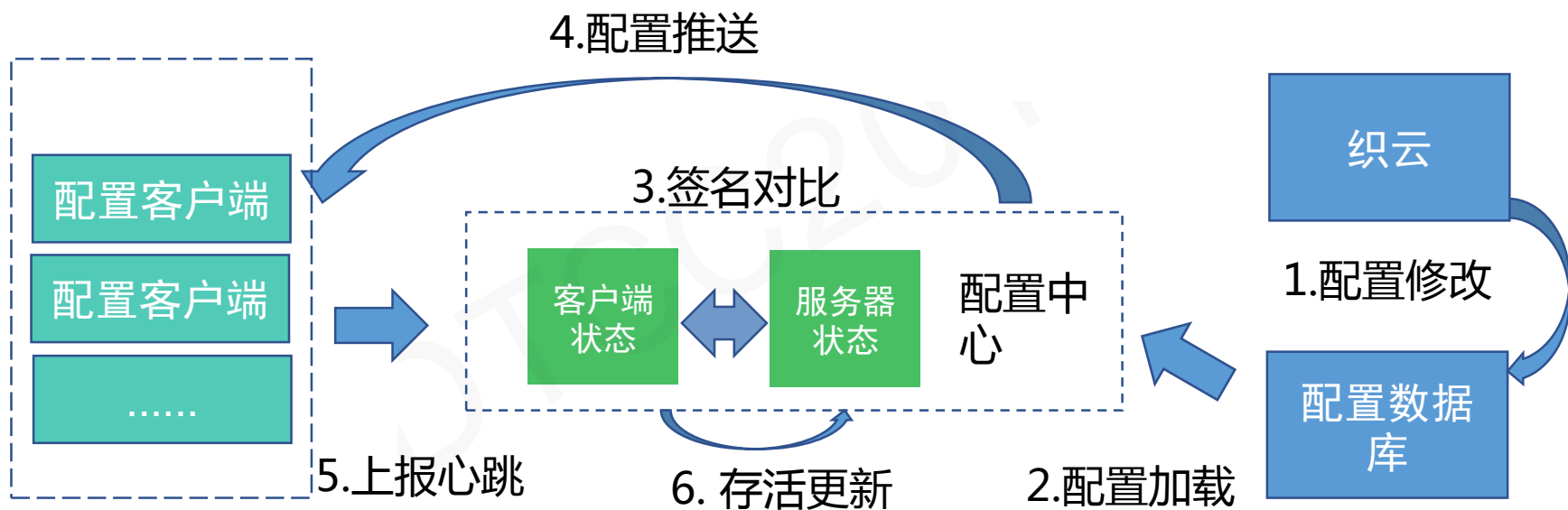
- 写远程流水
- 事务会话
- 防雪崩的超时机制
- 读写端口 (实现中)
- 热Key导流 (实现中)

进阶

CRS配置中心和心跳管理



- ❑ 配置加载：配置表-> 服务端配置
- ❑ 签名对比：筛选出需要接收配置的机器
- ❑ 配置推送：将配置推送到指定机器
- ❑ 上报心跳：上报心跳并更新客户端状态
- ❑ 存活更新：将服务器记载的各机器状态做更改





心跳管理

烽火台



运维模块

工兵



信息外报

对外信使

配置模块

行政中枢



巡检模块

巡逻部队



远程日志系统

ULS，统一远程日志服务

配额管理系统

根据配置对资源限额限流

监控系统

Monitor监控
多维监控

01 04

02 05

03 06

流水中心

记录写流水，快照回档
与审计

名字服务

CMLB，带负载均衡、死机
探测、回包统计、反向鉴权
功能的名字服务系统

冷备系统

微服务定义：微服务是一种软件架构风格，它是以专注于单一责任与功能的小型功能区块为基础，利用模组化的方式组合出复杂的大型应用程序，各功能区块使用与语言无关的 API 集相互通讯。

PART 03

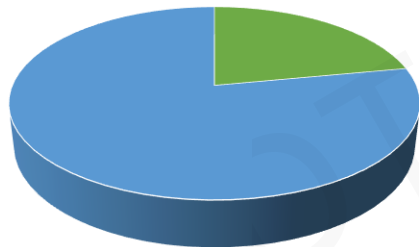
运营系统

运维本质



系统运维的本质是人与计算机共同参与的一项系统工程

成本投入



■ 上线前 ■ 上线后

40-90%的开销在运营阶段

1年研发

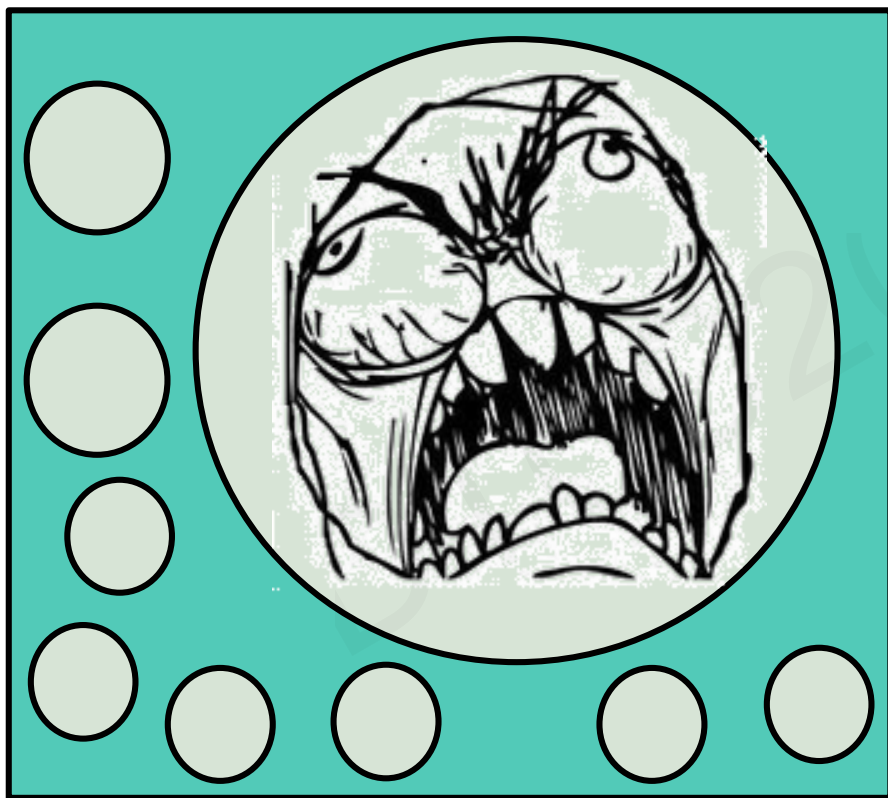


3-5年运营



- 挑战一：元信息的一致性管理？
- 挑战二：万台设备的作业方式？
- 挑战三：如何实现智能调度？

元信息混乱带来的不一致性



地域？机房？交换机？

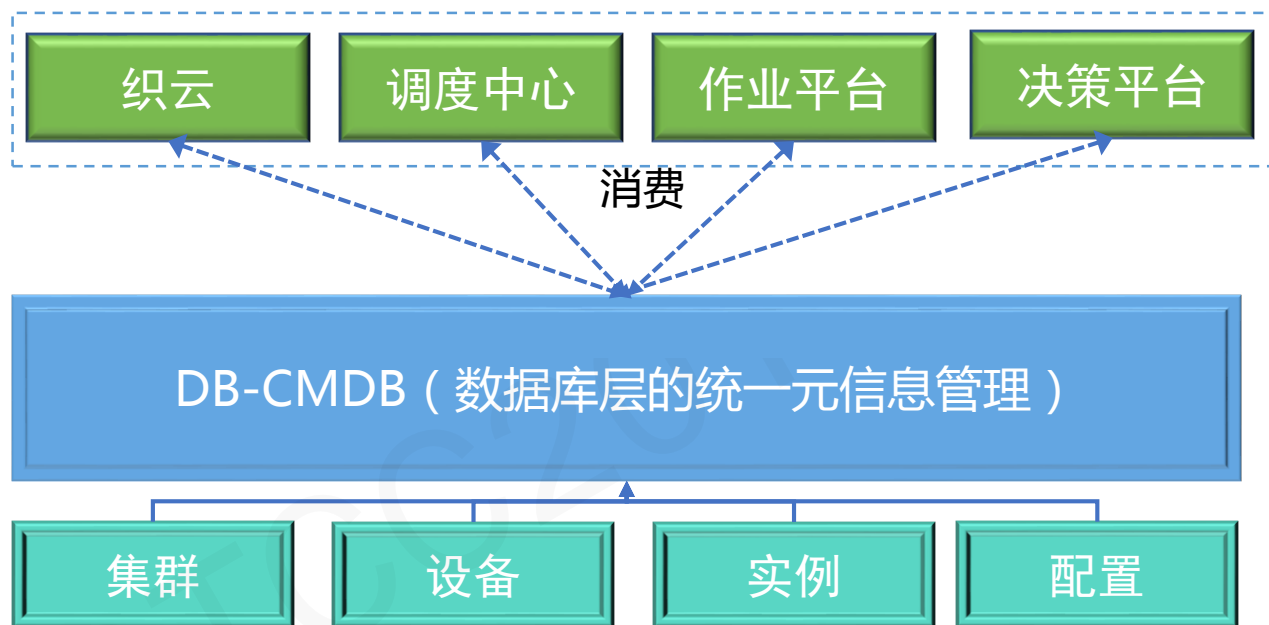
集群，IP，端口？

实例，容量，配置？

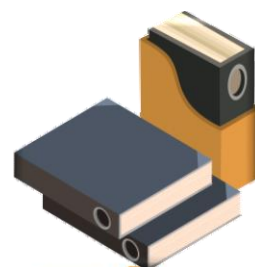
挑战一：元信息管理

- ✓ 信息互联，消除孤岛
- ✓ 数据同步和一致性保障
- ✓ 统一配置与元信息
- ✓ 基础信息的统一管控

全 准 入



- 梳理元信息，公共特征，分类建模，
- 抽象属性与方法，定义数据结构，
- 设计API，定制同步。



挑战一：元信息管理

鸚鵡螺2.0

Welcome, varianfeng

DB CMD8

请选择

显示字段

选择显示字段

输入搜索字段

列表页

Qedis系统列表

Show 10 entries

id	系统id	系统名称	用途	状态	用途	创建者	创建时间	所属中心	备注
3	110001	广州二区-测试	测试	使用中	测试	varianfeng	2018-11-15 17:49:38	10	
4	110001	广州二区-测试	测试	使用中	测试	varianfeng	2018-11-15 17:49:38	10	
6	110001	上海一区-测试	测试	使用中	测试	varianfeng	2018-11-15 17:49:38	10	
7	110001	北京一区-测试	测试	使用中	测试	varianfeng	2018-11-15 17:49:38	10	
8	110001	北京一区-测试	测试	使用中	测试	varianfeng	2018-11-15 17:49:38	10	

挑战二：万台设备的作业方式



腾讯云

织云



大家是怎么做批量作业的？

DTCC
2018

数领先机 智赢未来 (9)

IT168.com

ChinaUnix

ITPUB

挑战二：万台设备的作业方式



腾讯云

织云



手工式运维



堆人力扛业务



难积累与传承

万台设备
亿级QPS

规模上来后，
原始的运维方式无法应对



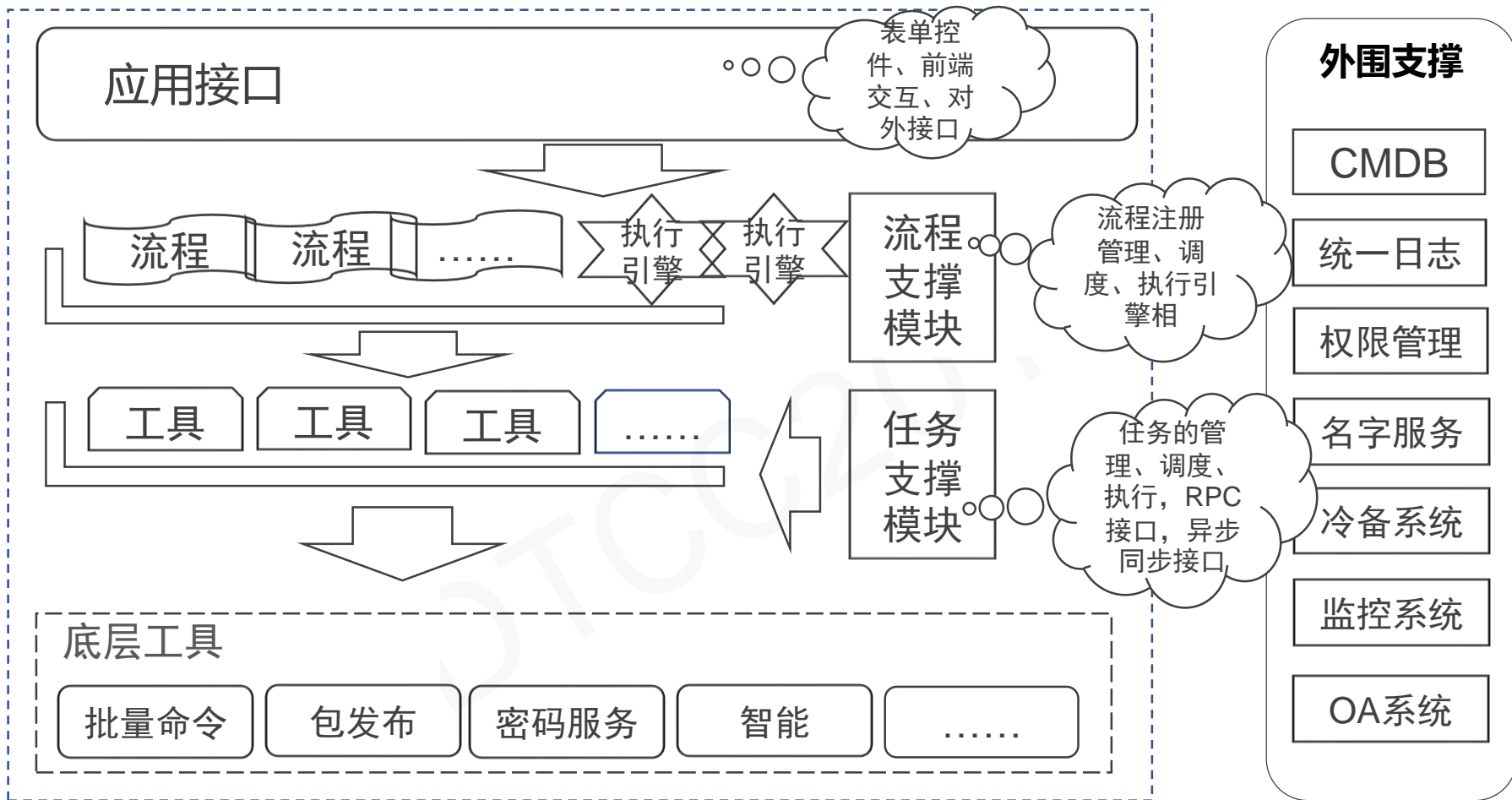
DTCC
2018

数领先机 智赢未来 (9)

IT168

ChinaUnix

ITPUB



作业平台

任务名: 输出口自动下线流程

流程引擎: crsinstallEngine

参数

```
1 {
2   "dist_ips": [],
3   "systemid": 0,
4   "cmibenv": "",
5   "dst_module_id": 0,
6   "disable_seconds": 0,
7   "minutes": 3,
8   "check_cmib_flag": 1
9 }
```

返回值

```
1 {
2   "code": -1,
3   "msg": "",
4   "succ_ips": [],
5   "fail_ips": [],
6   "detail": {}
7 }
```

备注

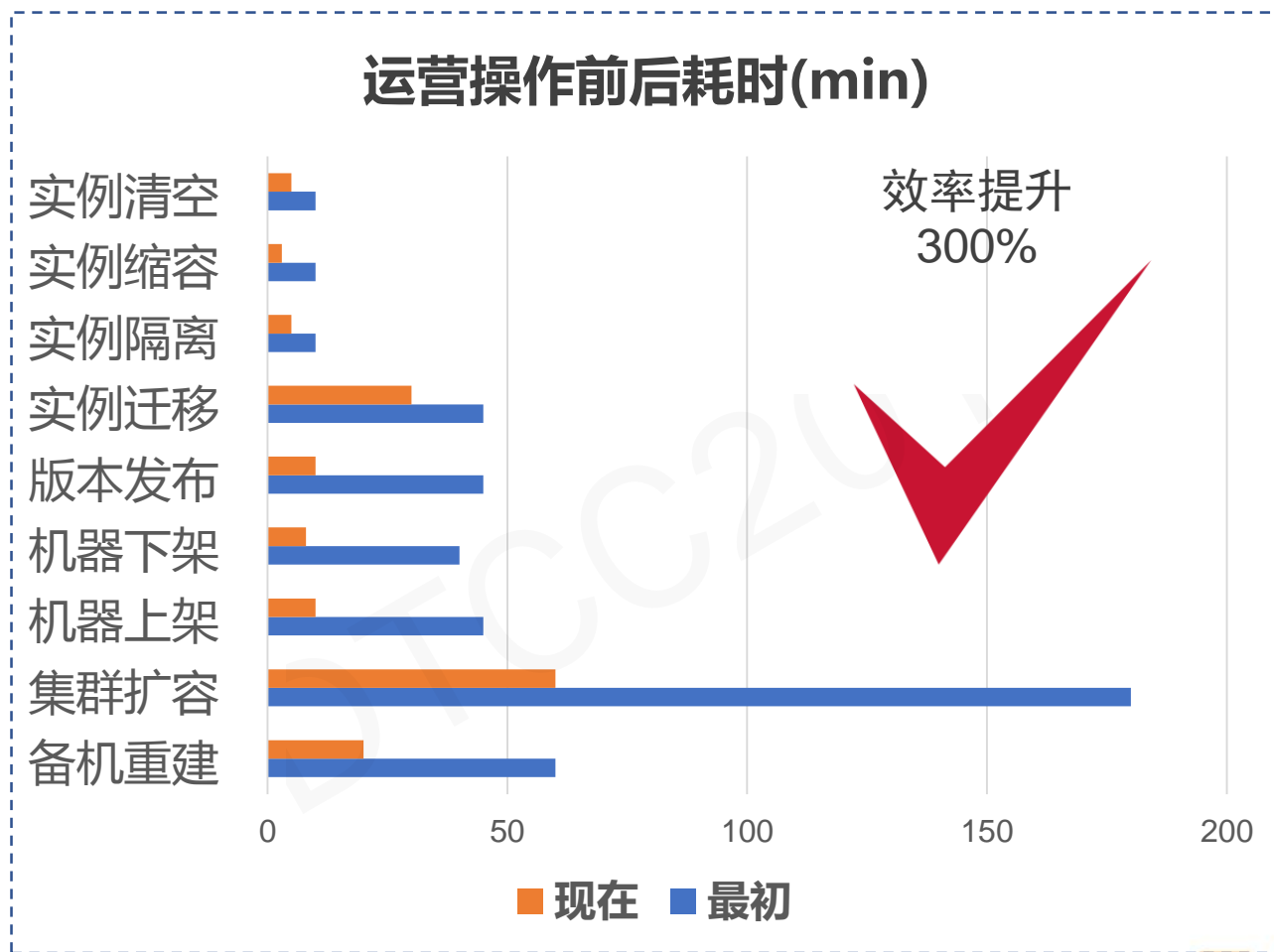
"dst_module_id": 0 下线模块: 824739 日常下线模块 912011 下线隔离模块
"disable_seconds": 0 隔离时间 (s)
"minutes": 3 检测流量时间 (min)
check_cmib_flag 检查 ip 列表是否属于 systemid, 0 不检查, 1 检查

任务列表

添加任务

序号	任务id	任务名	操作
1	260	根据数据库登记的信息判断ip列表中是否有...	上移 下移 删除
2	271	判断ip列表是否在对应的appid中	上移 下移 删除
3	256	在名字服务 (cmib) 中禁用ip列表	上移 下移 删除
4	261	等待	上移 下移 删除
5	179	检查机器是否有流量	上移 下移 删除
6	257	清理 ip 列表的黑名单配置	上移 下移 删除
7	259	在 CMIB 黑名单数据库中删除	上移 下移 删除
8	258	批量转移黑名单	上移 下移 删除







如何实现咖啡运维？

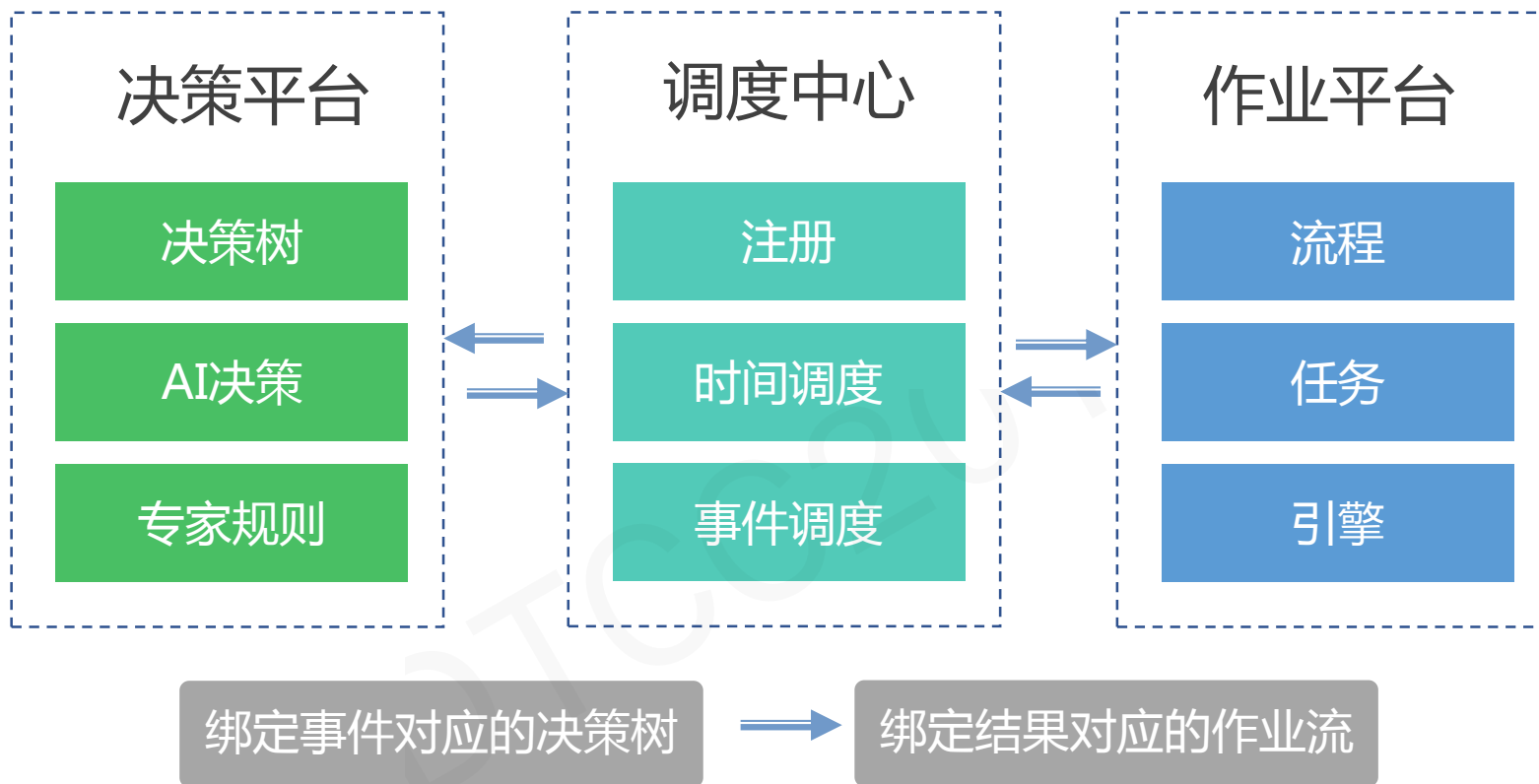
运维闭环的构造

挑战三：如何实现智能调度？

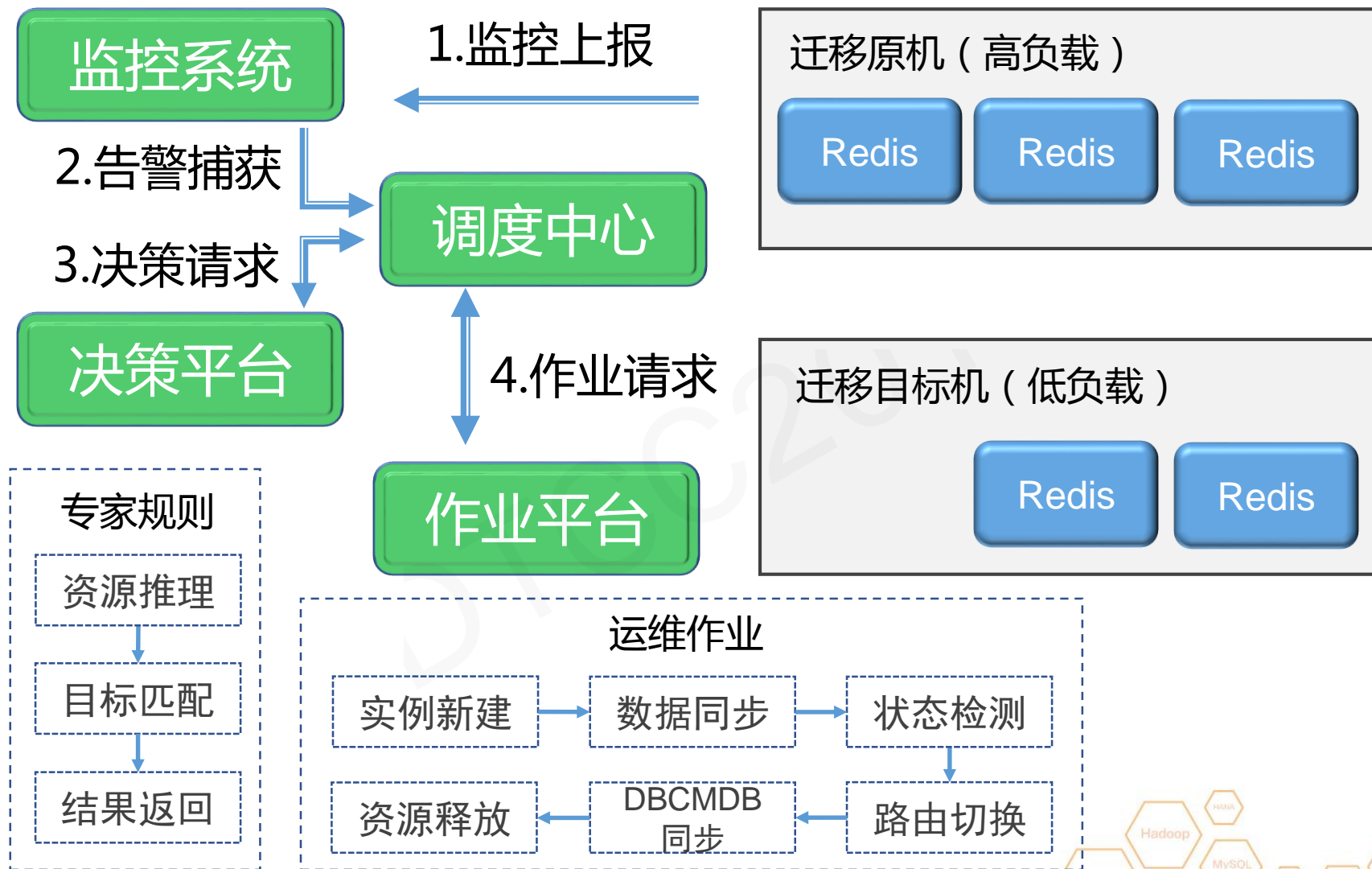


腾讯云

织云



智能调度全流程



注册事件并关联决策

设计：

注册

Show 10 entries Search:

事件名	中文名	描述	关联流程	修改	删除
singlePoint	单点		singlePoint	修改	删除



触发单点事件后的处理：

设计：

Show 10 entries Search:

实例ID	事件名	事件参数	决策流程	决策结果	运维操作ID	接收时间
111	singlePoint	systemId serverNameList:["server_000","server_001"]	singlePoint	extendCopy		2018-03-28 09:30:28.000000

Showing 1 to 1 of 1 entries

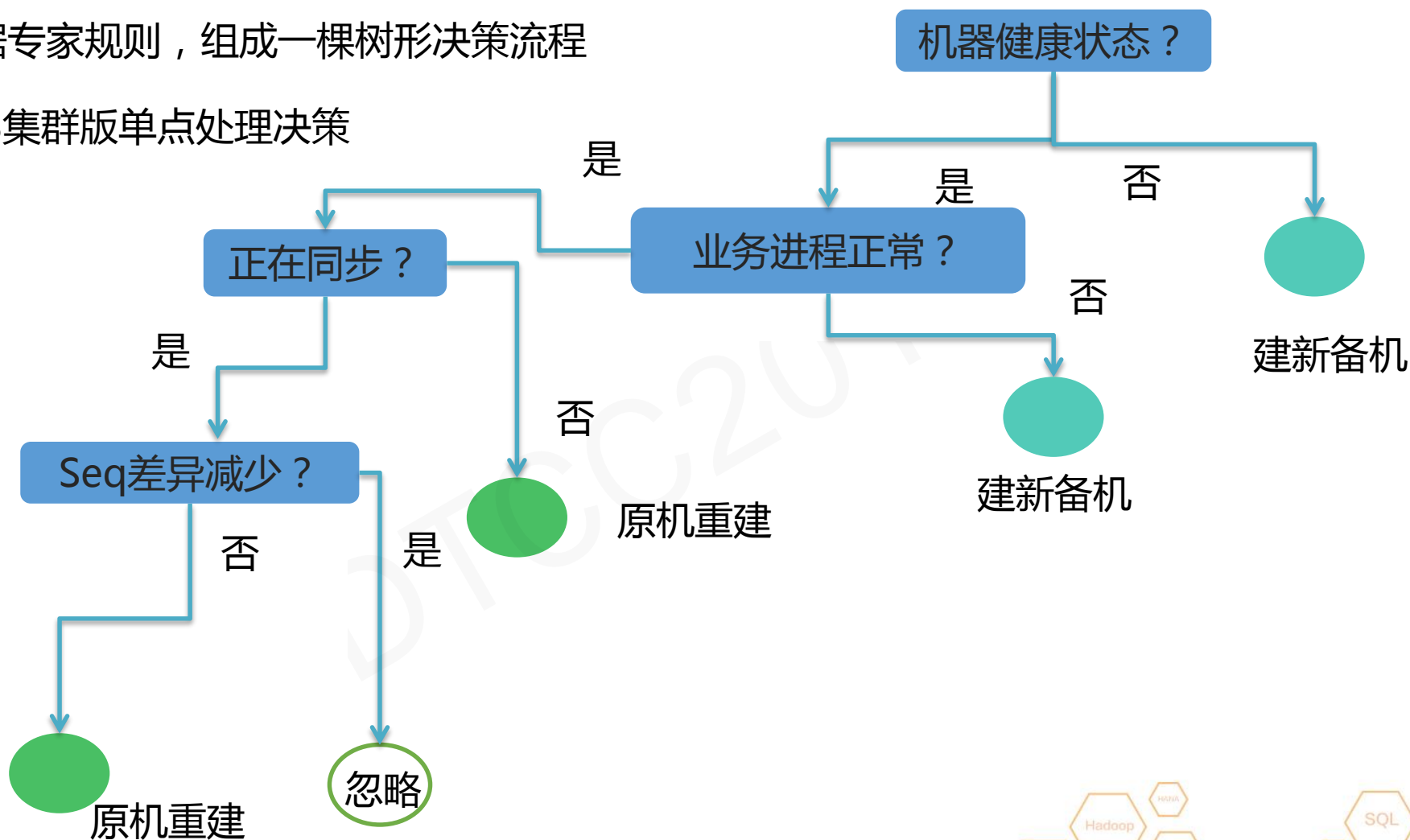
Previous Next

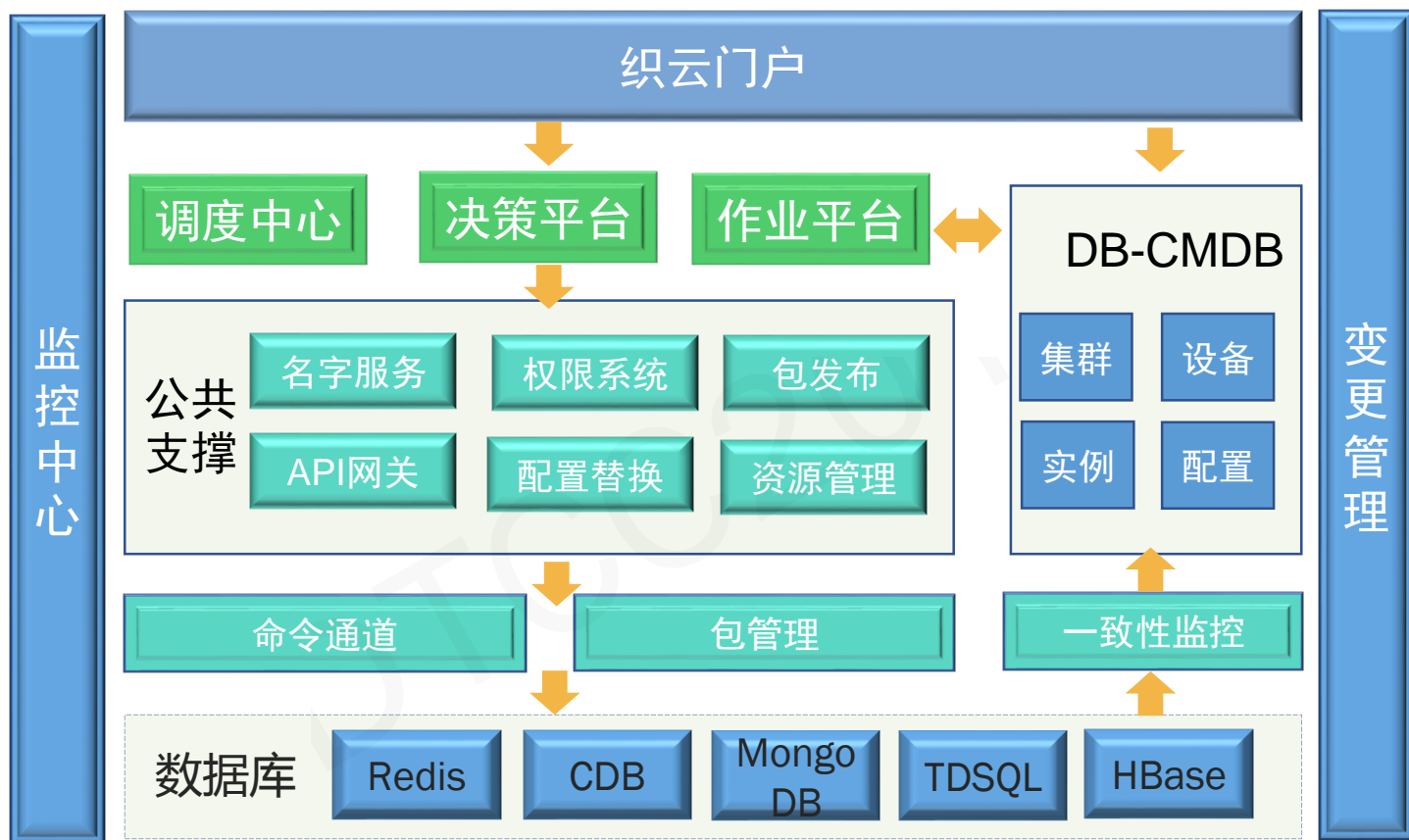
调用的决策流程

执行的运维操作及实例ID

依据专家规则，组成一棵树形决策流程

CRS集群版单点处理决策





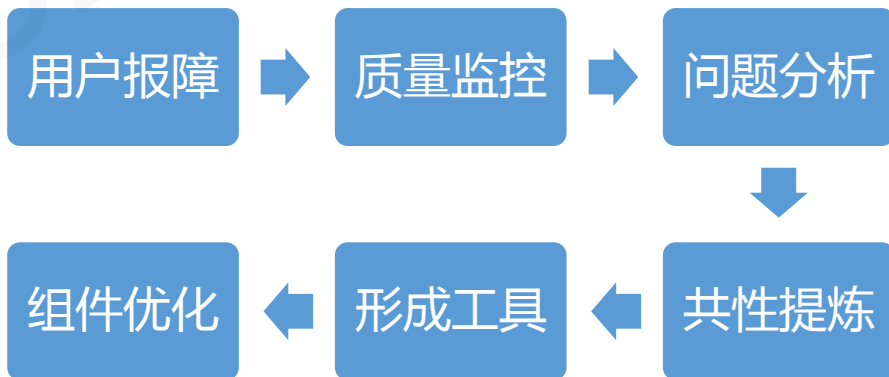
PART 04

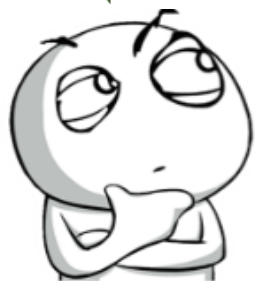
运营思考



海量运营之道

能力、质量、效率、成本





支持



推动



引领





@织云

企业级智能一体化
运维平台

謝

@我





讲师申请

联系电话（微信号）：18612470168

关注“ITPUB”更多
技术干货等你来拿~

与百度外卖、京东、魅族等先后合作系列分享活动



让学习更简单

微学堂是以ChinaUnix、ITPUB所组建的微信群为载体，定期邀请嘉宾对热点话题、技术难题、新产品发布等进行移动端的在线直播活动。

截至目前，累计举办活动期数60+，参与人次40000+。

ITPUB学院

ITPUB学院是盛拓传媒IT168企业事业部（ITPUB）旗下
企业级在线学习咨询平台
历经18年技术社区平台发展
汇聚5000万技术用户
紧随企业一线IT技术需求
打造全方式技术培训与技术咨询服务
提供包括企业应用方案培训咨询（包括企业内训）
个人实战技能培训（包括认证培训）
在内的全方位IT技术培训咨询服务

ITPUB学院讲师均来自于企业
一些工程师、架构师、技术经理和CTO
大会演讲专家1800+
社区版主和博客专家500+

培训特色

无限次免费播放
随时随地在线观看
碎片化时间集中学习
聚焦知识点详细解读
讲师在线答疑
强大的技术人脉圈

八大课程体系

基础架构设计与建设
大数据平台
应用架构设计与开发
系统运维与数据库
传统企业数字化转型
人工智能
区块链
移动开发与SEO



联系我们

联系人：黄老师
电话：010-59127187
邮箱：edu@itpub.net
网址：edu.itpub.net
培训微信号：18500940168