



QCon 全球软件开发大会
INTERNATIONAL SOFTWARE
DEVELOPMENT CONFERENCE

BEIJING 2018

《MySQL 容器化部署实践》

演讲者 / 王晓波

背景



- 同程旅游早期的数据库都以单库的MySQL。
- MySQL的单库，导致TPS最终还是会成为一个瓶颈。
- MySQL+DB中间件解决水平拆分问题。
- MySQL水平拆分的引入会使数据库实例数量大幅上升,传统运维手段维护成本高,交付能力差。

MySQL数据库为何要Docker化

1. MySQL数据库迅速爆炸式增长后，服务器规模不断增大，快速部署是个问题。
2. 随着业务的发展，扩容数据库的不方便不快捷，也是个问题。
3. 大量数据量小的数据库系统也单独部署在物理机，浪费问题突出。
4. DBA的数据库自动化标准化运维的需求。
5. Docker在同程的大规模使用，应用部署环境100%容器化，有Docker丰富的经验。



让数据库的部署点单化开启

配置	DB架构	硬件选型	机房
2核4G 4核4G	一主一从	SATA-SSD	A机房
4核8G 8核8G	一主多从	PCIE-SSD	B机房
8核16G 16核16G	分片集群	大容量磁盘SAS	C机房
16核64G 32核64G			D机房
32核128G			

容器化之后的MySQL就是一个私有DB云

主从集群创建

分片集群创建

集成高可用方案

巨细无遗的监控项

精美的图形展示

便捷的告警管理

慢日志分析及查看

自动化备份

资源池管理

高可用切换

集群节点管理

扩容缩容

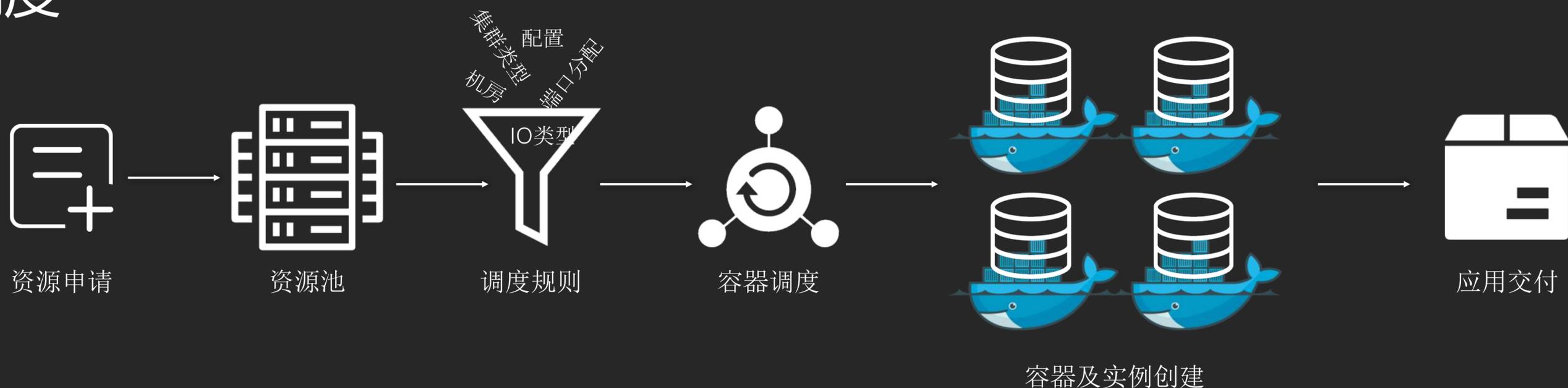
数据库及实例迁移

过载保护机制

总体架构



资源池调度



为了保证MySQL的高可用，需要在Docker容器分配时如何保障主从不在同一宿主机上。我们通过自研Docker容器调度平台管理所有宿主机和容器,自定义Docker容器的分配算法。实现了MySQL的高密度,隔离化,高可用化部署。

调度规则：

- 1.同一复制集群的实例在不同主机上。
- 2.优先分配CPU、内存、磁盘空间资源最空闲的主机。
- 3.根据IO需求调度容器创建在不同IO类型的主机。
- 4.申请新集群时，若IO要求高则按照宿主机的IO情况，优先选择IO最空闲的主机。
- 5.VIP集群必须主从端口一致，Proxy接入的集群端口无需一致。
- 6.VIP集群端口基于网段递增，Proxy集群端口基于IP递增

Docker里放了什么



内核版本

Kernel版本 4.7



操作系统

CentOS 7.2



宿主机

部署服务器监控、容器监控agent容器



容器

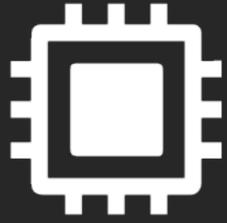
Docker版本 1.12,部署监控及系统服务agent



镜像

MariaDB镜像(按产品)、MySQL5.7镜像(按产品)、监控容器镜像、HA管理系统镜像、实例迁移服务镜像、监控服务端镜像

资源隔离



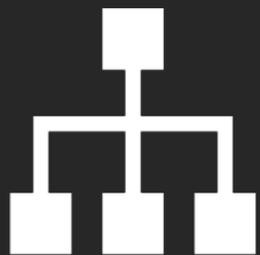
CPU最大超卖3倍，通过cpu-period配合cpu-quota一起使用，来限制容器的CPU的使用量。比cpuset-cpus绑定CPU的方式灵活。



限制容器内存，且内存不超卖。通过—memory限制内存，同时结合MySQL自身参数控制几个内存大户(比如buffer_pool等)，最后配合lxcfs增强隔离性。



IO方面由于我们采用挂载宿主机本地的磁盘设备，还不能做到彻底隔离。所以对于高IO的实例使用的是PCIE-SSD。磁盘空间方面，我们在申请时会预估出一个量，使用超过80%的时候会结合本地磁盘空间评估是否有足够空间扩容，若宿主机剩余空间不足会启动迁移扩容流程。



目前使用的host模式，无法隔离网络。但是考虑到10G接入，且单机密度可控的情况下，网络消耗不会过载。另外目前我们已经在线下尝试结合Ovs+Dpdk的方案实现网络隔离。

PS：容器虚拟化带来轻量高效，快速部署的同时，docker容器在隔离性方面也存在一些缺陷。例如，在容器内部proc文件中可以看到Host宿主机上的proc信息。这样就导致了一些问题，比如监控信息不准确、限制内存会导致应用程序OOM等。我们基于lxcfs组件来增强容器的隔离性。

容器的调度



- 提供两种API：
- Docker API的封装，用于创建,删除,修改。
- 集群管理API，用于集群管理。
- Scheduler调度：选择最优节点创建容器。
- Agent：用于连接监控模块，上报系统运行状况。

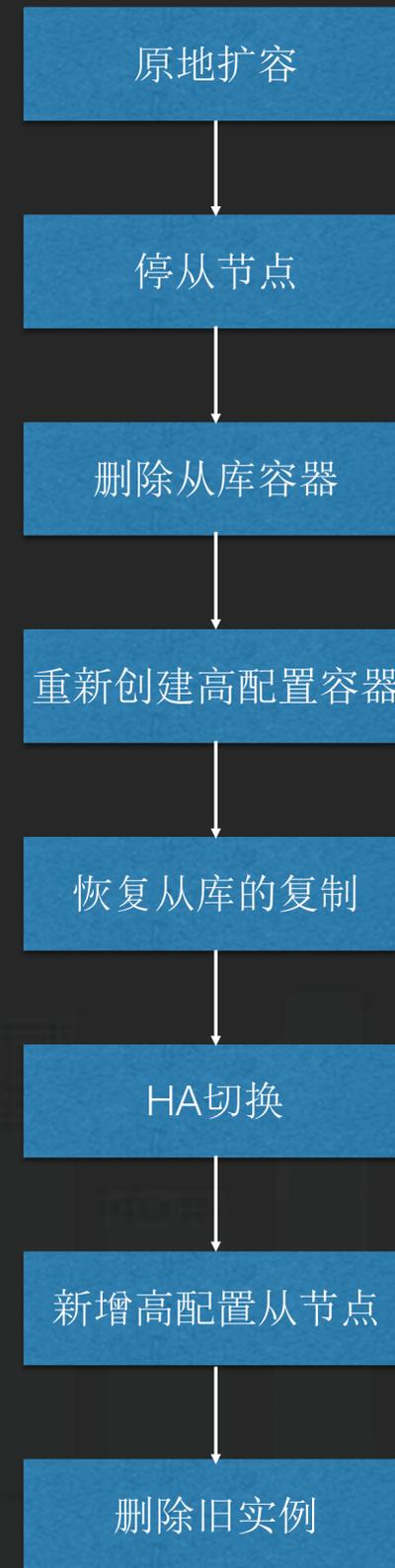
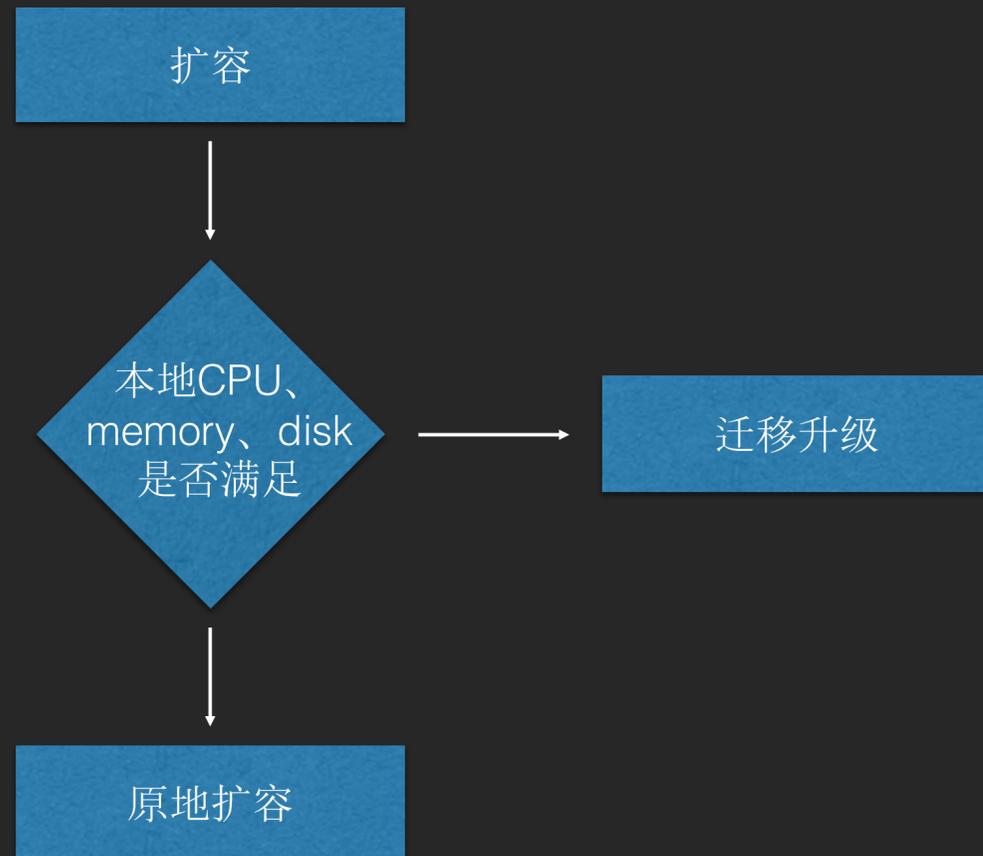
磁盘挂载



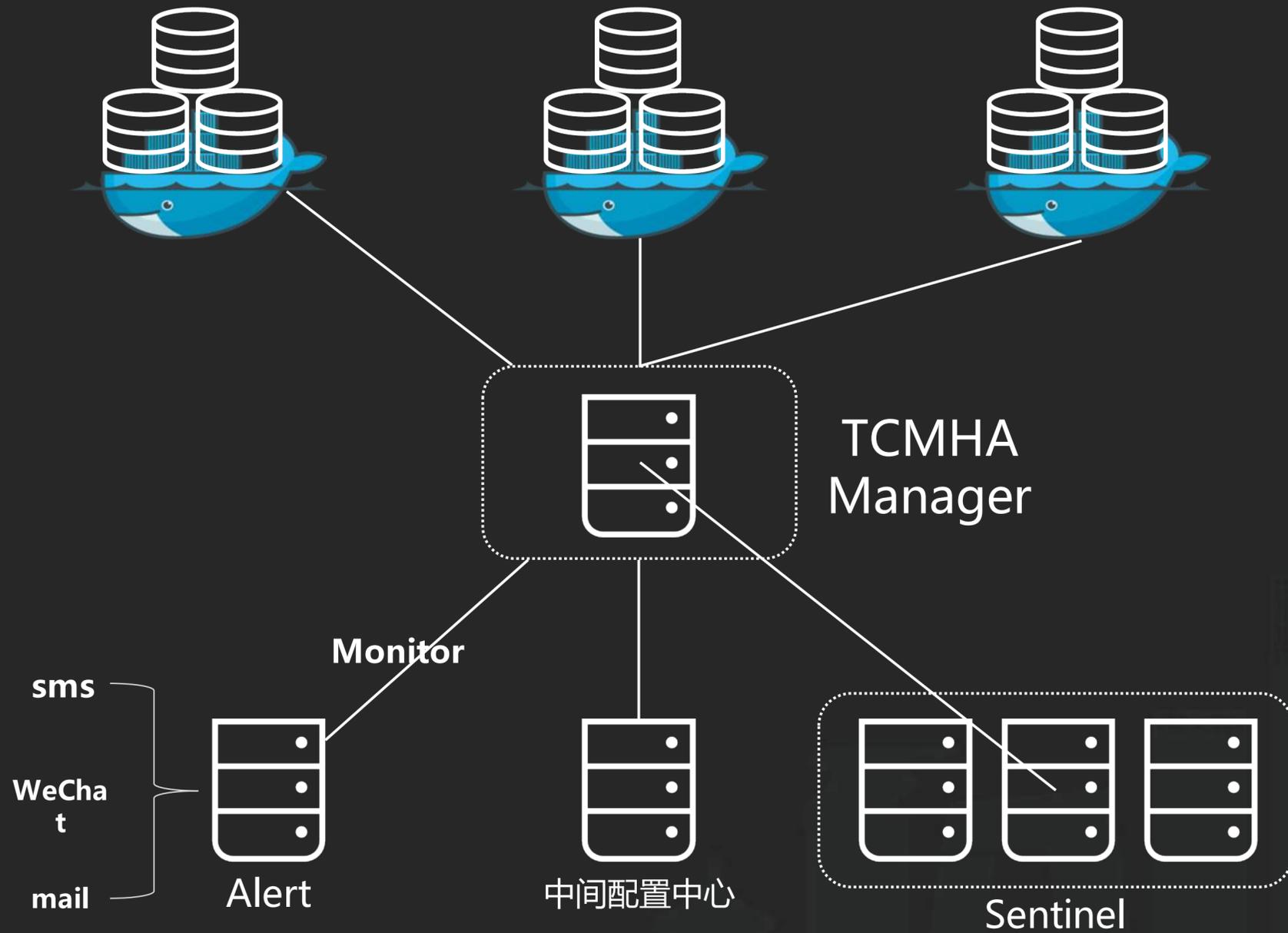
- 为了保证容器内的MySQL实例有更好的磁盘IO性能。采用了本地宿主机磁盘挂载到容器内的方式(每个实例对应一个文件夹)。这种方式的优点是IO性能最佳，随之而来的缺点是磁盘容量不好估算，有可能在使用了一段时间后出现磁盘空间不足的问题，这个时候则会启动迁移扩容的流程。已经在和提供高密度IO分布式存储解决方案厂商接触，计划测试平台接入分布式存储的方案。

集群扩容

扩容逻辑



高可用管理



TCMHA 管理分为两个部分：

- MariaDB、MySQL5.6的高可用管理是基于开源工具MHA定制开发后的工具完成的，该工具支持了MariaDB的gtid、更完善的哨兵检测机制、对接了DB中间件等定制化功能。
- MySQL5.7的MGR复制集群是我们自己写的一套高可用组件配合DB中间，实现无感知的高可用切换。

DB中间件

兼容mysql协议

支持SELECT/INSERT/UPDATE/DELETE语句

支持单DB实例上的inner join

支持单DB实例上的事务

支持聚合函数：max、min、sum、avg、count

支持：distinct、order by、group by、limit、

top:definition text

支持多种拆分方式

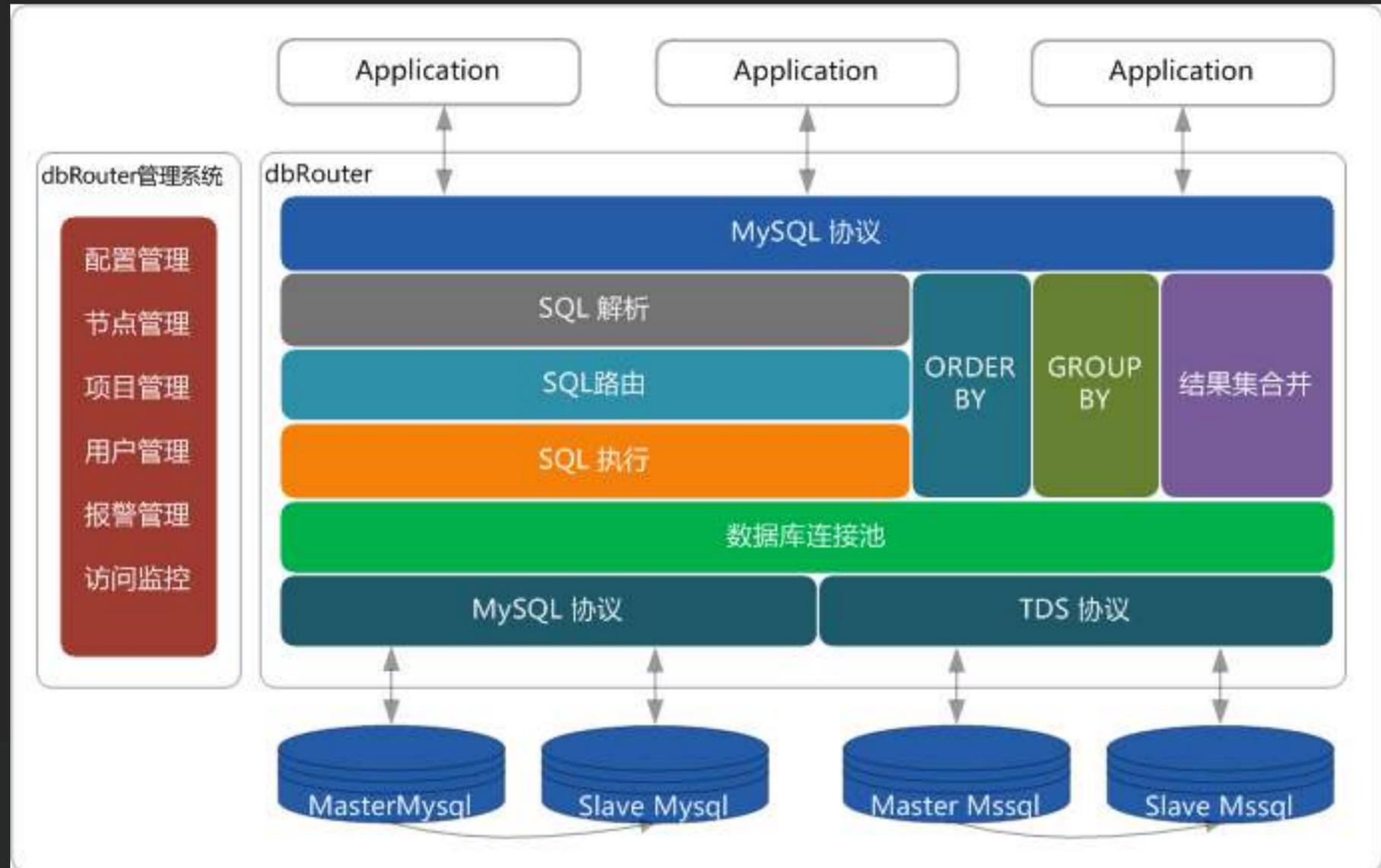
不分区

根据关键字段，进行hash分区

根据时间字段，进行时间分区

根据关键字段，进行区间分区

读写分离模式



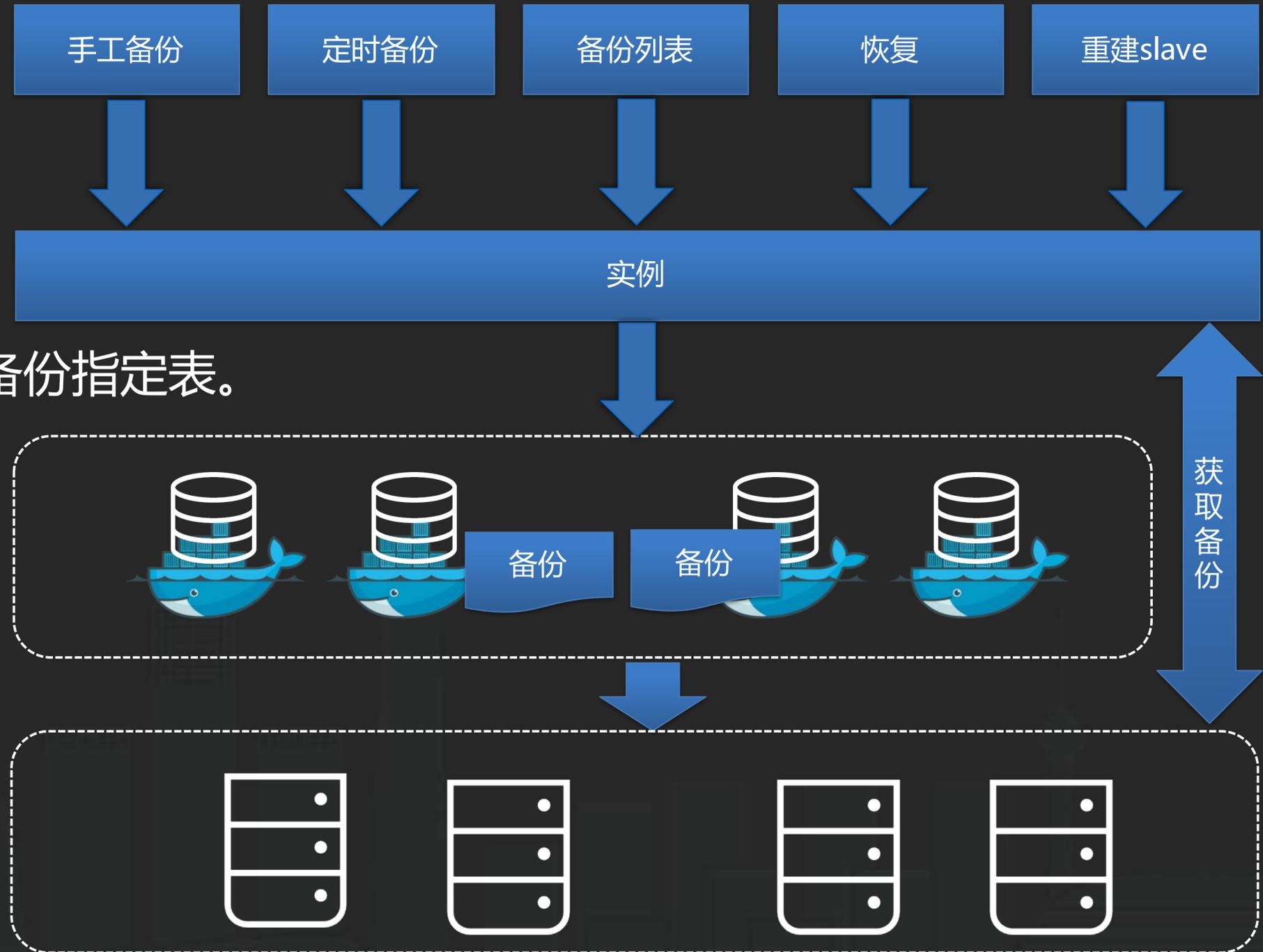
集群管理

- HA切换
- Slave增加删除
- 创建新集群
- 查看监控
- 慢日志分析
- 慢日志查看
- 手动备份

The screenshot displays the MySQL Cluster Management web interface. On the left is a dark sidebar with navigation options: 仪表盘 (Dashboard), 管理 (Management), 项目列表 (Project List), MySQL集群管理 (MySQL Cluster Management), 服务器管理 (Server Management), 告警管理 (Alert Management), 任务队列 (Task Queue), VIP管理 (VIP Management), 综合查询 (General Query), 参数模板 (Parameter Template), 脚本管理 (Script Management), 工具 (Tools), 数据校验 (Data Verification), 实例迁移 (Instance Migration), 日志 (Logs), 操作日志 (Operation Log), 备份日志 (Backup Log), and 告警日志 (Alert Log). The main content area is titled 'MySQL集群管理' and shows a cluster overview with a dropdown menu and a 'PRODUCT' button. Below this, it displays cluster details: 集群名称: TCC... (Cluster Name), 生产线: ... (Production Line), 项目名称: ... (Project Name), and 版本: ... (Version). A '集群维护' (Cluster Maintenance) section includes buttons for SLAVE管理, HA管理, Proxy管理, 监控查看, and 信息修改. The central part of the interface is a table listing nodes with columns for Proxy, 节点ID (Node ID), 高可用 (High Availability), VIP, 主机 (Host), 类型 (Type), 角色 (Role), Data_Size, Log_Size, 状态 (Status), 创建人 (Creator), and 管理 (Management). The table shows six nodes (ID 1-6) with their respective configurations and status indicators. Each node has a '详情' (Details) and '操作' (Action) button. At the bottom right, there is a 'QCon 2018·北京站' logo.

Proxy	节点ID	高可用	VIP	主机	类型	角色	Data_Size	Log_Size	状态	创建人	管理	
<input type="checkbox"/>	1	●	8	无	172.16.1.101	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.102	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>					172.16.1.103	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>	2	●	8	无	172.16.1.104	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.105	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>	3	●	8	无	172.16.1.106	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.107	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>	4	●	8	无	172.16.1.108	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.109	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>	5	●	8	无	172.16.1.110	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.111	16C64G	Slave	100 GB	200 GB	●		详情 操作
<input type="checkbox"/>	6	●	8	无	172.16.1.112	16C64G	Master	100 GB	200 GB	●	...	详情 操作
<input type="checkbox"/>					172.16.1.113	16C64G	Slave	100 GB	200 GB	●		详情 操作

备份系统



多样化备份：

1.物理备份：

通过xtrabackup备份20G以上的实例。

2.逻辑备份：

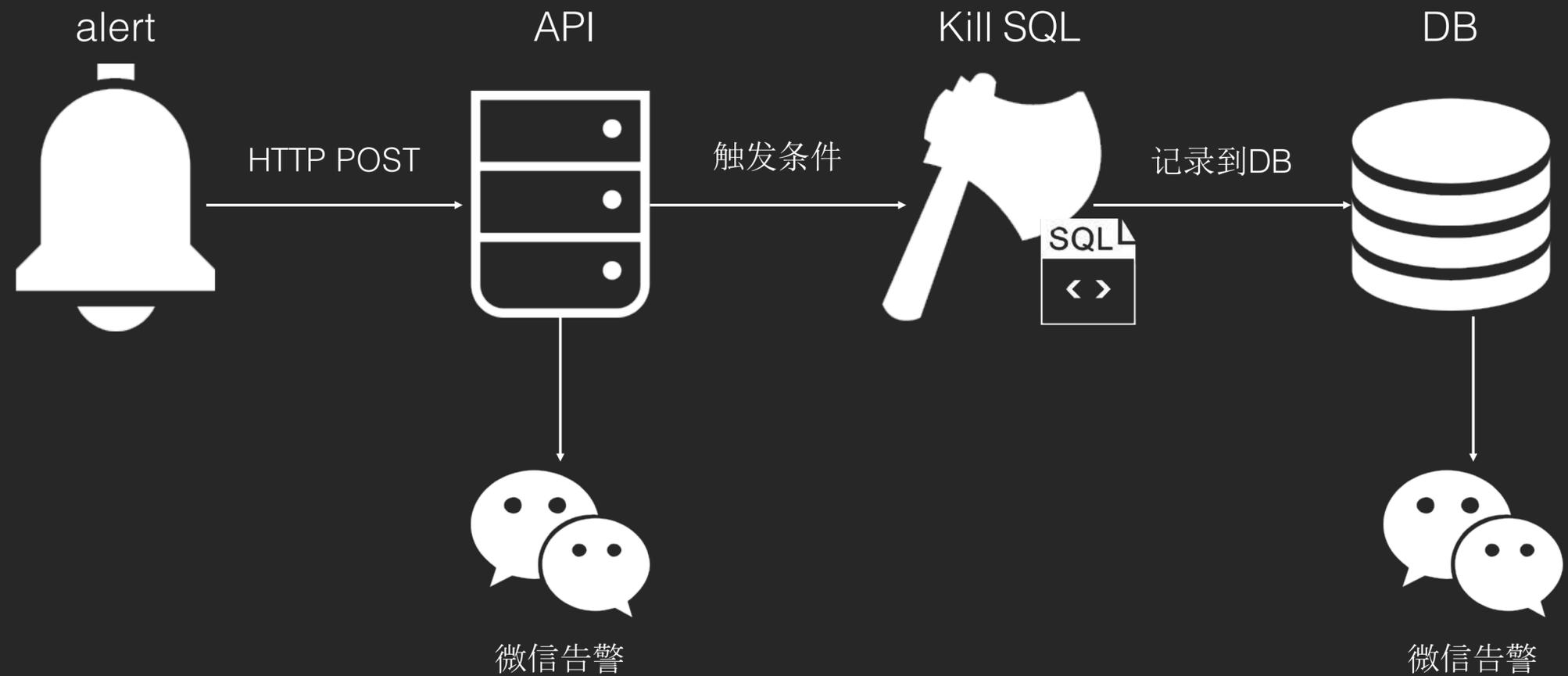
通过mydumper备份20G以下的实例或备份指定表。

备份策略：

1.每日自动化备份。

2.DBA临时手工备份。

过载保护



目的:

- 当系统数据库较高时确保大多数请求能够正常访问。

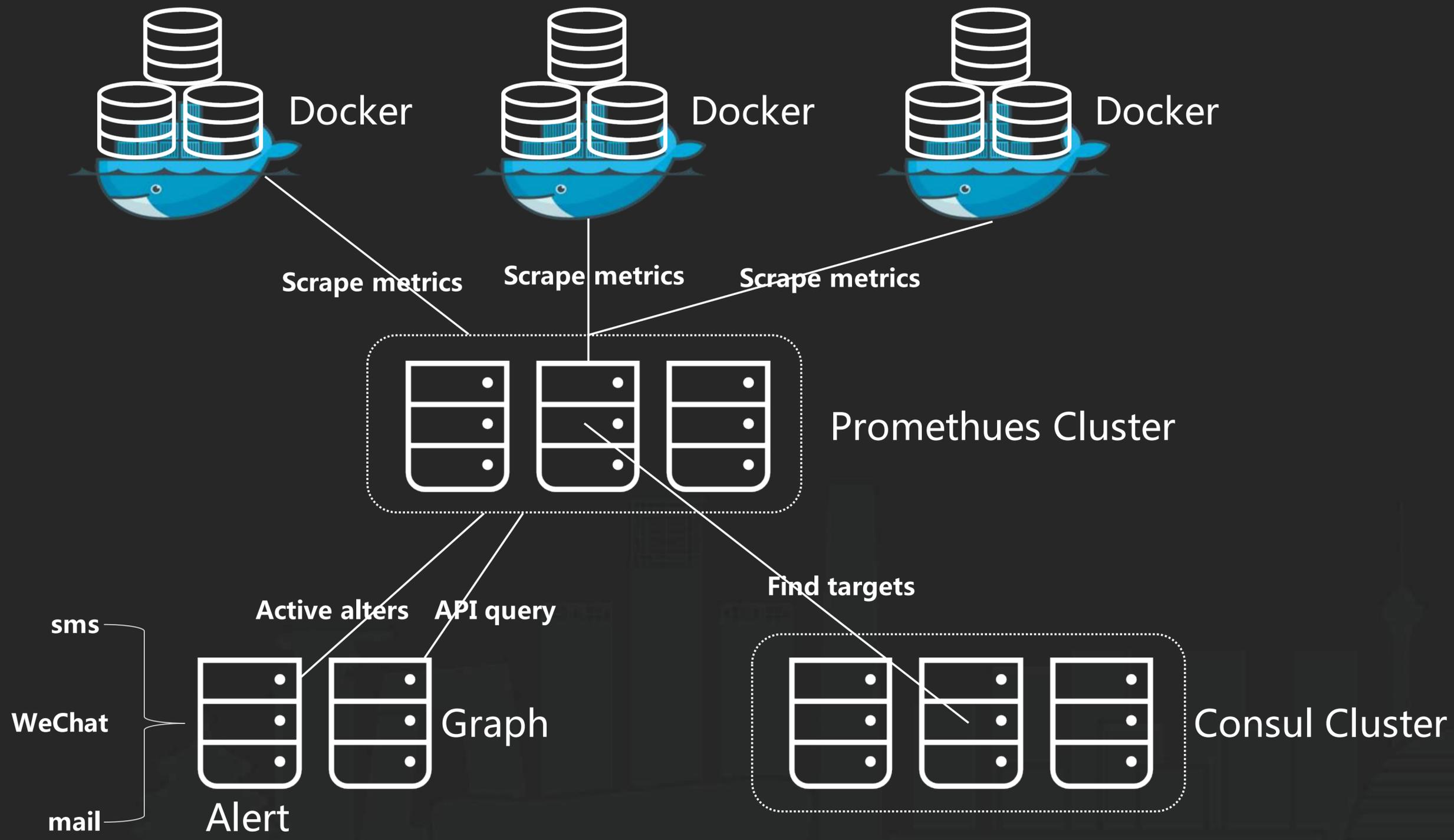
触发条件:

- 依赖监控告警系统,thread_running > 30(可自定义)

kill哪些语句:

- select语句
- 非特定系统账号(可自定义)
- 执行时间超过10S
- 只Kill sql请求, 不kill连接

监控告警系统



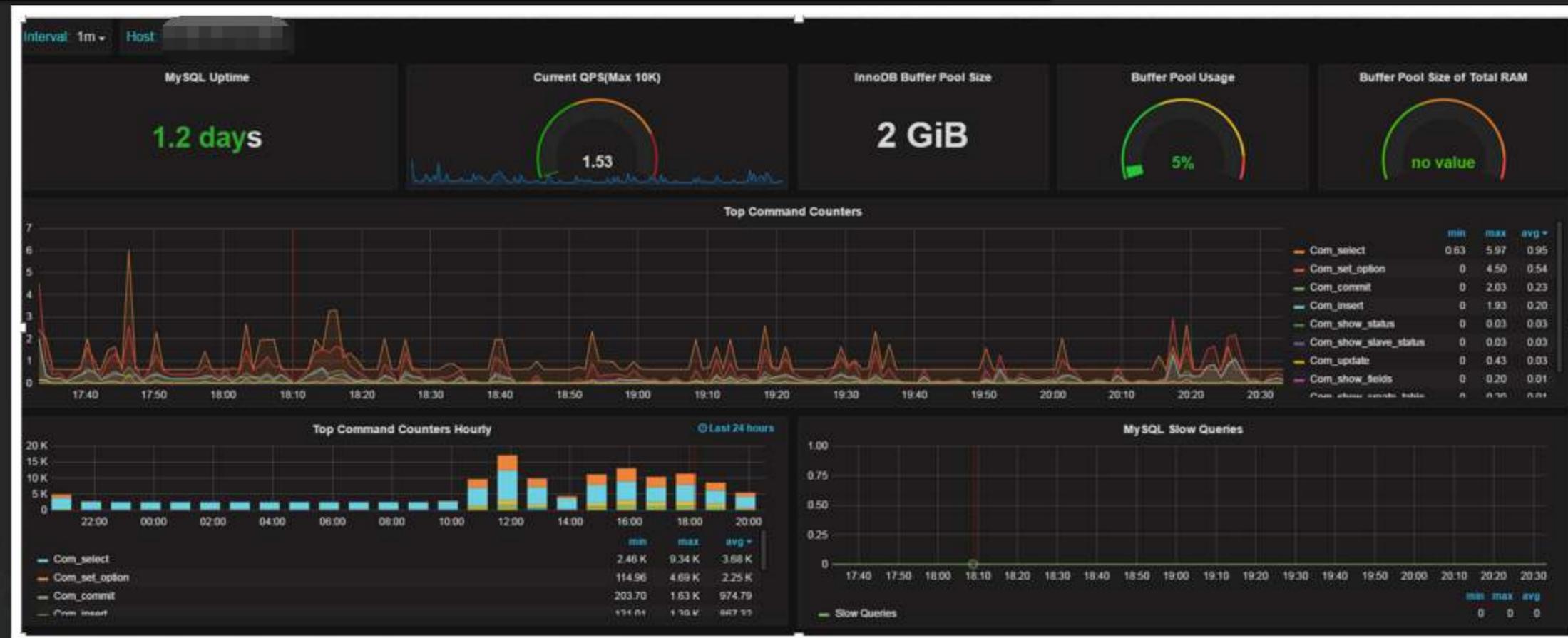
MySQL监控

多达上百项的详细的监控信息：

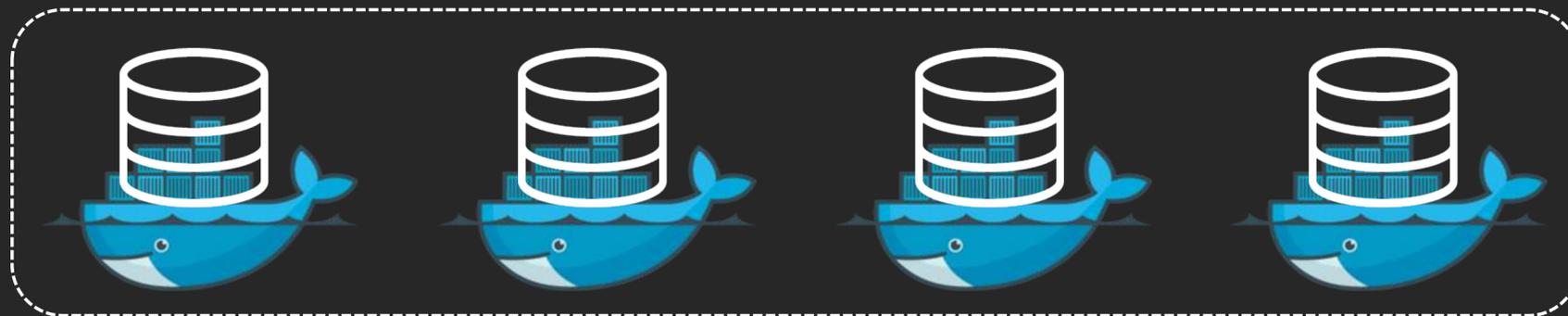
集群名称: [模糊] 生产线: [模糊] 项目名称: [模糊]

集群维护: SLAVE管理 HA管理 **监控查看** 信息修改

Proxy	节点ID	高可用	VIP	MySQL	主机	类型	角色	Data_Size	Log_Size	状态	创建人	管理
	1	●		MySQL 容器	[模糊]	1C1G	Master	[模糊]	[模糊]	●	[模糊]	详情 操作
				接收人管理	[模糊]	1C1G	Slave	[模糊]	[模糊]	●	[模糊]	详情 操作



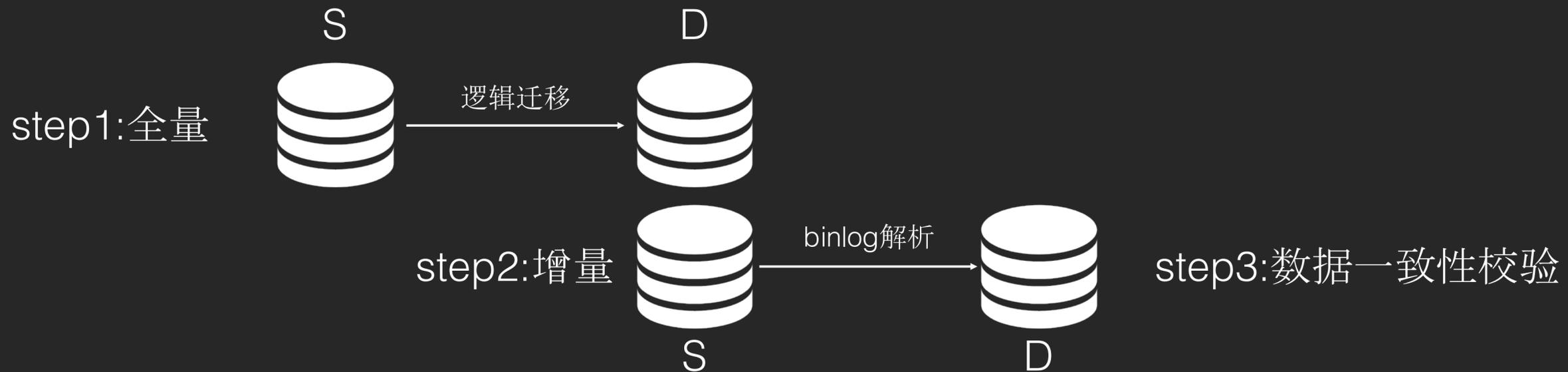
慢日志分析



告警：
当慢SQL每秒产生的数量超过阈值的时触发告警通知DBA及相关人员及时关注问题。

分析：
定期将实例端的慢SQL分析后录入数据库，然后通过系统页面查看慢SQL的执行时长、每日次数、每日平均耗时等多个维度的指标。同时也支持ui上查看执行计划和表结构信息。

数据库及实例迁移



- 实例迁移我们支持了全实例、库、表级别的迁移,并且同时兼容各种MySQL版本。全量迁移基于开源工具mydumper和myloader做了定制化开发,解决了一些问题,比如utf8mb4导致乱码的问题等。增量复制使用的是公司自研的基于binlog复制的产品。该工具支持指定库、表级别复制,另外也支持源与目标命名不一致(结构一致)的复制。

实例迁移

实例迁移 实例、库、表迁移

:

平台外的实例迁移至平台内

平台内数据迁移

实例迁移

首页 / 实例迁移

源信息

IP Port

账号

密码

是否使用所有库 否 是

库名

目标信息

集群名称

项目名称

生产线

创建人

节点ID

开始迁移

迁移

INTERNATIONAL SOFTWARE DEVELOPMENT CONFERENCE

QCon 2018·北京站

成果

- 屏蔽底层物理资源，降低决策时间
- 资源利用率提升30倍
- 90%的基础运维工作完全自动化
- 交付能力提升70倍
- 生产环境近3500个实例在平台上工作

效率提升

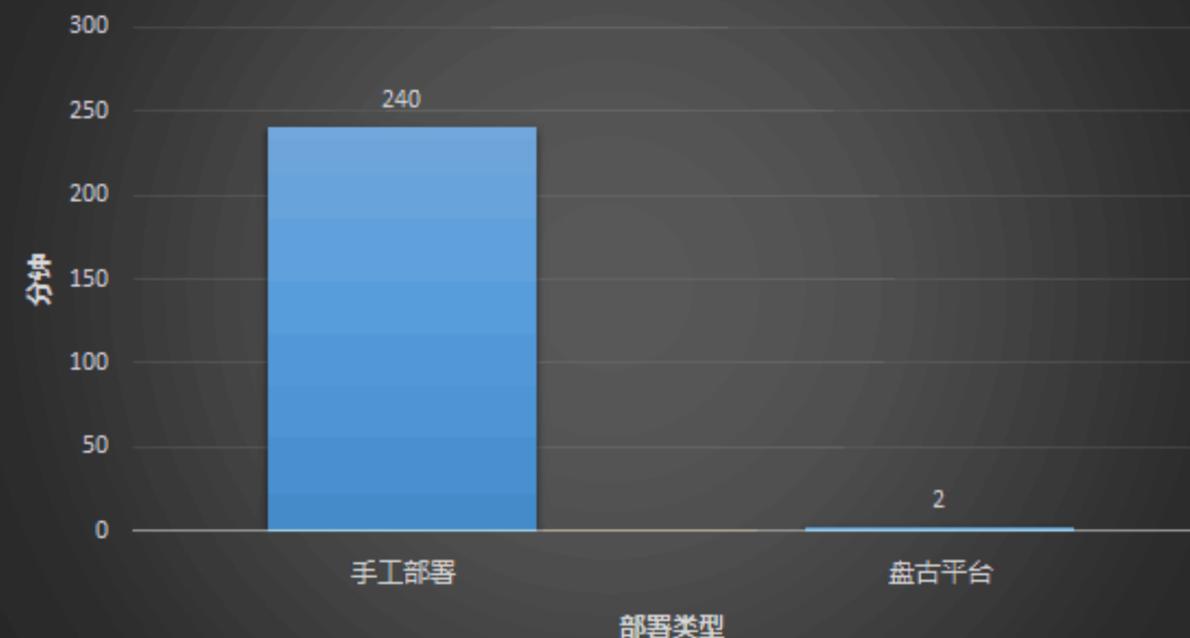
手工部署：

部署一套高可用集群+备份，配置监控。至少30分钟。部署32个节点的分片的集群，至少一个上午。无系统化管理，资源分配情况，无法统一调配，服务器资源利用率低。

MySQL容器平台：

部署一套高可用集群+自动化备份+慢日志分析+监控。用时1-2分钟。部署32个节点的分片集群，只需5分钟。标准化的系统管理，部署环境统一、配置文件统一。系统化的操作降低人为失误和重复劳动。资源使用集中管理，有效利用服务器资源。

32节点分片部署时间对比



HA管理 监控查看 信息修改

节点ID	高可用	VIP	主机	类型	角色	Data_Size	Log_Size	状态
6	●			8C32G	Slave	1 GB	GB	●
				8C32G	Master	1 GB	GB	●
				8C32G	Slave	1 GB	GB	●

数据校验

首页 / 数据校验

源

IP

Port

库名

排除的表

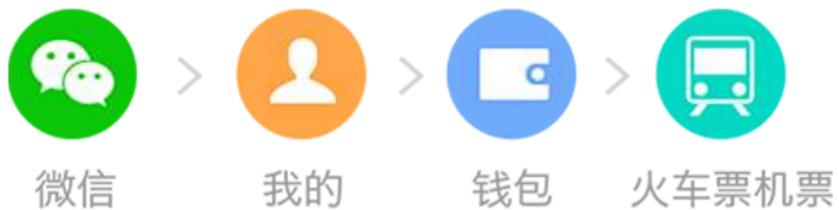
验证的表

目标

IP

Port

开始校验



晓波

江苏 苏州

扫一扫上面的二维码图案，加我微信

Thanks!



主办方 **Geekbang** > **InfoQ**
极客邦科技