



第十一届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2020

架构革新 高效可控



北京国际会议中心 | 2020/12/21-12/23

新浪微博数据库资源调度平台架构实践

张磊

新浪微博 DBA





个人介绍

- 目前就职于新浪微博基础架构平台，主要负责NoSQL和DB类自动化运维体系建设和方案落地

部门介绍

- 微博数据库平台负责新浪微博所有业务的数据库服务托管
- 覆盖主流关系数据库和NoSQL数据库
- 提供整体的OLTP和OLAP解决方案
- 支撑微博平台、微博主站、微博广告、热门微博和机器学习等公司核心业务
- 推动公司数据库技术创新和落地



资源调度平台shanks的产生背景

- 微博当前资源规模和dba配比
- 其他亟需解决的痛点



资源规模

- 服务器近万台规模
- 近10万的数据库实例数
- PB级关系数据存储
- 万亿级NoSQL访问
- DBA人数个位数，**人均管理实例近万**
- 服务SLA 99.99%

其他痛点

- 数据库资源多样
- 网络环境复杂
- 多云环境
- 热点事件带来的极速峰值流量
- 持续增长的资源访问





期待一个英雄的降临

- 基于资源调度姿势的考量
 - ❖ 主动
 - ❖ 被动
 - ❖ 主动和被动相结合



基于资源调度姿势的考量

➤ 被动姿势

❖ Saltstack

❖ Ansible

❖ 自动化运维平台



基于资源调度姿势的考量

➤ 主动姿势

❖ AIops (智能运维)

❖ SAAS (software as a service)



基于资源调度姿势的考量

- 相对柔和和稳定的方式：主动和被动相结合
 - ❖ 主动感知
 - ❖ 主动处理 + 被动处理
 - ❖ 节奏可控、规模可控

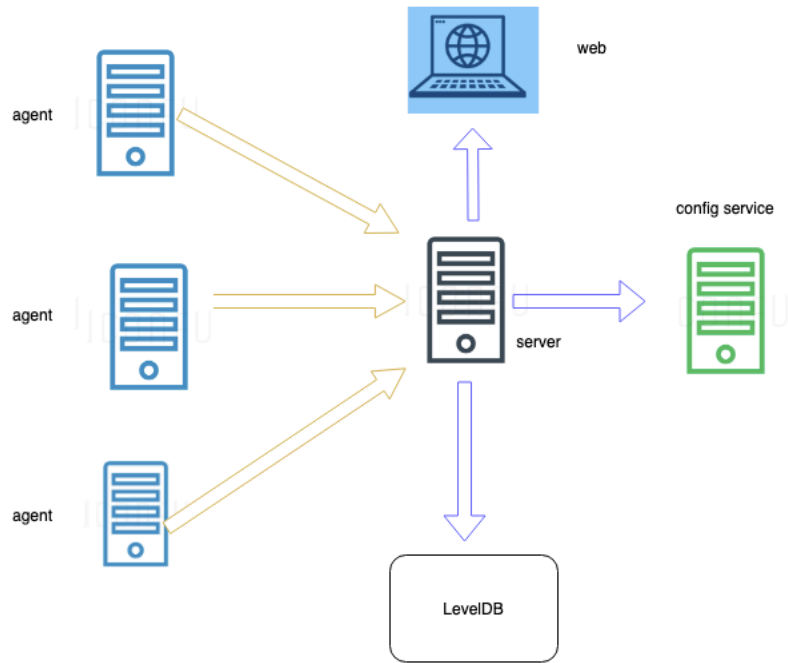
统一资源调度平台 — shanks

- 命名来自漫画：海贼王四皇之一
- 一个资源服务平台(RAAS)
- 万物皆资源，资源即服务
- 为所有支持的资源提供服务化的治理和全生命周期管理
- 提升资源管理效率
- 提升资源稳定性与可靠性
- 基于go开发
- 高性能可扩展
- 弱状态



Shanks架构

- 轻量agent
- 元数据存储基于LevelDB
- 高效的内存访问
- 支持万台规模
- 支持机器和实例层面调度



Shanks架构

➤ Agent干啥

- ❖ metric采集和上报
- ❖ login free
- ❖ 无状态
- ❖ 最低功能原则：别的模块能完成的功能一律不提供
- ❖ 0依赖原则：agent的启动、运行不依赖任何第三方提供的服务
- ❖ 0配置原则：agent的启动、运行不依赖任何静态配置
- ❖ 最小化暴露原则：不额外提供不需要的未知的功能

Shanks架构

➤ Server干啥

- ❖ 接收并分析metrics，如果触发阈值，生成报警事件
- ❖ 监控资源的运行状态
- ❖ 依据配置的策略进行报警，触发不同的exception handler进行处理，生成对应的task list;
- ❖ 接收resource admin提交的任务；
- ❖ 调度并执行任务：寻找合适的agent，将任务转换成相应的command下发给agent，完成任务的执行
- ❖ 提供通用API



核心功能

- 资源操作标准化
- 自动注册、服务发现
- 多维度监控与报警
- 服务自愈
- 弹性调度

设计理念

- 足够简洁、抽象：cs模式、界限清晰
 - ❖ 将redis、mc、mysql、HBase、mcq、qservice、DNS等统一描述成资源
 - ❖ 将资源的变更：包括扩容、缩容、DDL、备份、迁移等等抽象对资源的action，以task的方式提交
 - ❖ Agent安装不依赖任何环境
 - ❖ Agent只做metric上报和login free，Server只做metric分析统计和任务下发
- 足够健壮：应对单点、网络割接、agent挂起
- 足够智能：自动恢复大比例覆盖全网资源异常，降低人为干预
- 足够全面：将日常运维经验逐步反哺到shanks中，让其发挥更大的作用
- 足够灵活：介于主动运维和被动运维之间

资源标准化

- 抽象资源类型和操作
- 提供通用http api
- 支持实例部署、升级、扩缩容和迁移
- 方便和各种运维平台整合
- 批量操作成本更低

自动注册、服务发现

- 服务器初始化后agent进程会自动启动
- Agent会定期上报服务器上面所有服务的相关指标
- agent和server定期通信，如果agent挂掉，server会感知，并会将其拉起
- 上报的资源信息会和统一配置中心结合 供业务使用

多维度监控与报警

- agent上报的相关指标会同步到监控dashboard和指标决策系统
- 指标决策系统 提供多维度的指标聚合计算，生成资源健康检查报告和报警事件
- 决策系统会对报警事件进行判断，触发不同的处理策略

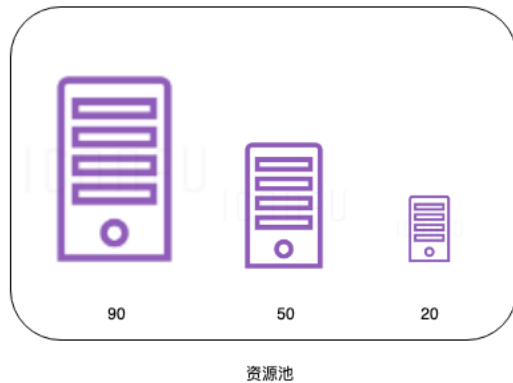


服务自愈

- 支持常规实例资源故障自愈
- 基于标准化api封装
- 自定义配置自愈策略
- 支持多种资源类型故障切换和自动恢复
- 机房级别网络故障切换
- 降低服务故障时间，降低人为干预成本

弹性调度

- 支持资源指标维度的容量水位自动调度
- 支持多种资源的弹性扩缩容
- 基于资源池和产品线为基本调度单位
- 支持分钟级扩容百台的规模

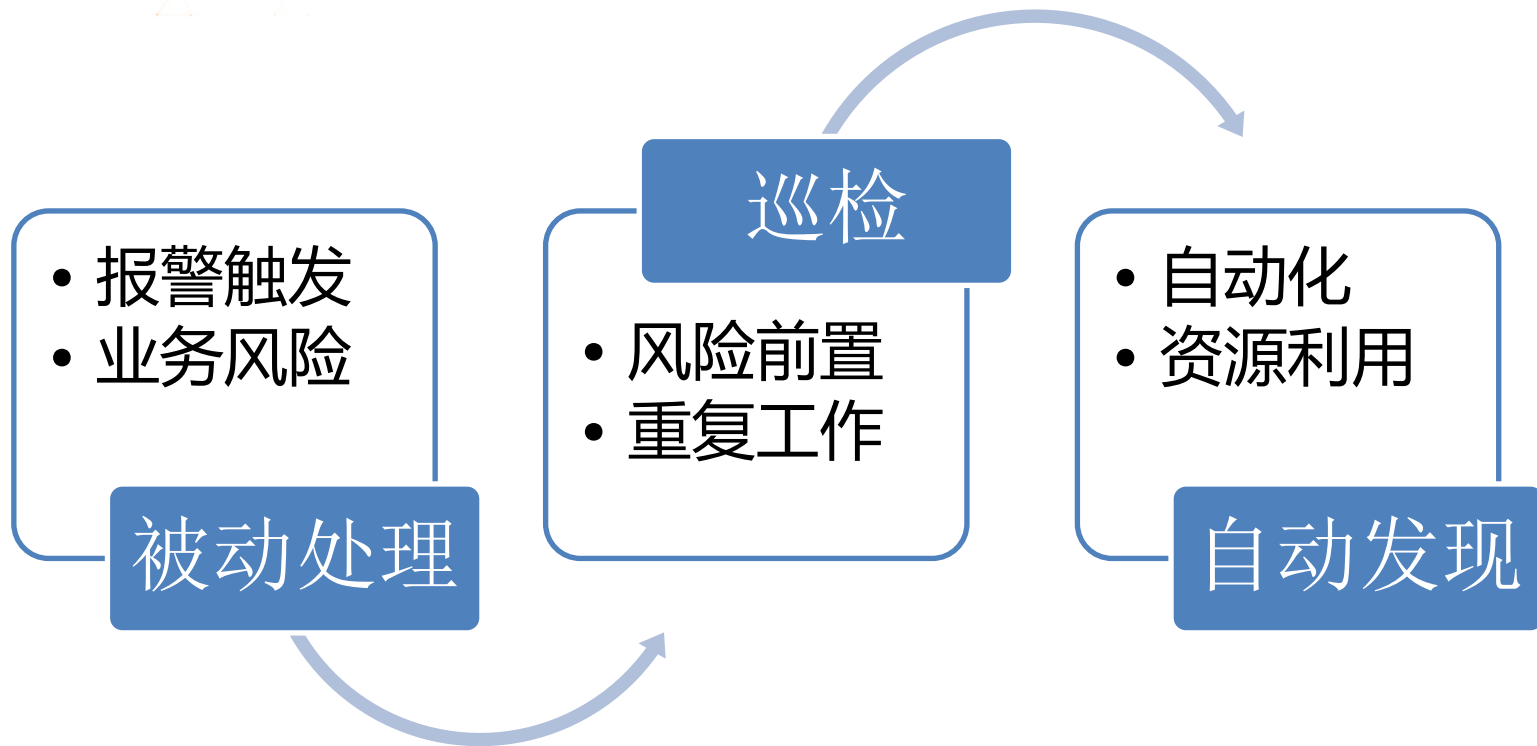




典型应用场景

- 容量问题
- 应对微博热点
- 数据库访问质量自动治理

容量问题





应对微博热点

- 微博典型场景: 新闻热点 明星事件等
- 峰值高
- 资源有状态如何快速扩容
- 成本

自动弹性扩容

- 基于公有云的弹性快速扩容
- 依靠云的弹性来提高资源冗余度
- 降低成本
- 进一步降低弹性扩容时间

```
},  
- elasticspecs: {  
  - : {  
    enabled: false,  
    ecsgroup:  
      - strategy: {  
        hwmk: 2500,  
        lwmark: 1800,  
        maxslaves: 4,  
        minslaves: 1,  
        name: "readload"  
      }  
    },  
  },  
}
```



数据库访问质量治理

- 业务反馈访问资源慢
- 单一实例访问超时
- 但是资源基础监控正常



数据库访问质量治理

- 实时的访问质量监控
- 资源耗时的同比环比报警
- 异常实例自动处理和降级



平台收益

- 资源管理效率提升
- 资源稳定性与可靠性提升
- 自动恢复覆盖全网60%以上资源异常
- 减少资源故障定位和恢复时间
- DBA工作效率提升



未来展望

- 资源服务化、规格化、云化
- 提升DB类快速扩缩容能力
- 资源精细化管理、智能化管理



资源服务化、规格化、云化

- 进一步标准化
- 统一资源使用姿势
- 减少和业务沟通成本，资源开箱即用
- 最大化资源利用率，应对各种复杂场景的资源调度
- 资源和机房解耦



提升DB类快速扩缩容能力

- 数据量大 如何弹性？
- 大DB拆小，DB规格化
- 高效的数据备份恢复体系
- 高效的数据传输体系
- 高效的数据校验体系



资源精细化管理、智能化管理

- 资源全方位的健康检查和优化
- 进一步扩大资源异常自动处理比例
- 资源特征化管理



总结 - 用软件工程的思路做资源运维事

- 找到痛点，定义场景。
- 体验做透、方案优雅。
- 保持克制、体验闭环。
- 小步快跑、快速迭代。

THANKS

