



智能弹性容量管理

by 张娟 (希宁)

QCon

全球软件开发大会

成为软件技术专家
的必经之路

[北京站] 2018

2018年4月20-22日 北京·国际会议中心

7折 购票中, 每张立减2040元
团购享受更多优惠



识别二维码了解更多

AiCon

全球人工智能与机器学习技术大会

助力人工智能落地

2018.1.13 - 1.14 北京国际会议中心



扫码关注大会官网



极客时间

重拾极客精神 提升技术认知

下载极客时间App

获取有声IT新闻、技术产品专栏，每日更新



扫一扫下载极客时间App

ABOUT ME

2011年加入阿里

2016年~至今，从事集团弹性资源管理





Contents

01	背景	4
02	智能弹性容量管理	9
03	具体实践	20
04	未来展望	32

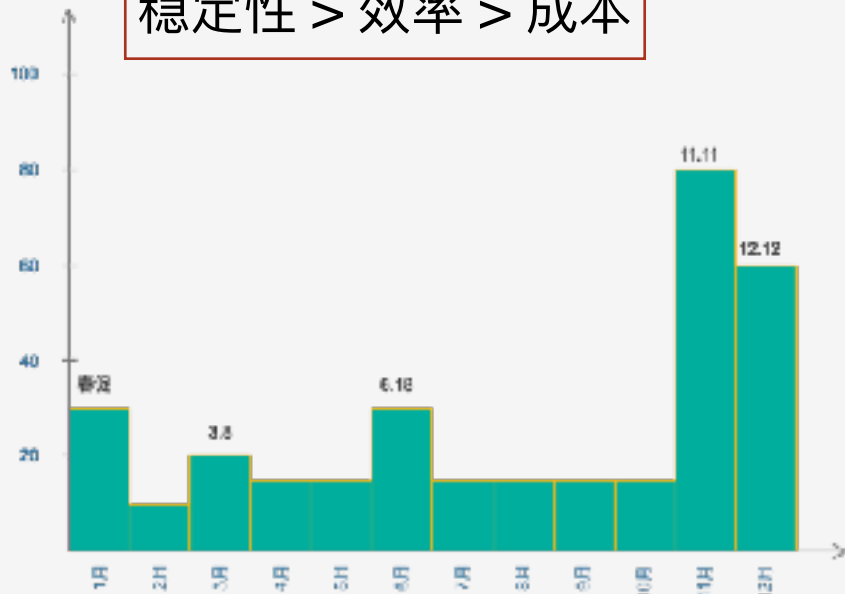
01

背景



成本运维的困境

稳定性 > 效率 > 成本



月度流量趋势图

月度资源保有量趋势图



成本运维的困境

多少合适?

You never know!



传统做法



经验预估



等比预估



压测到目标量级



容量规划的意义

- 用更科学手段做资源运营。



智能弹性容量管理

概述



What

“容量规划”+“弹性伸缩”+“风险评估”

How

“智能决策” + “自动执行”

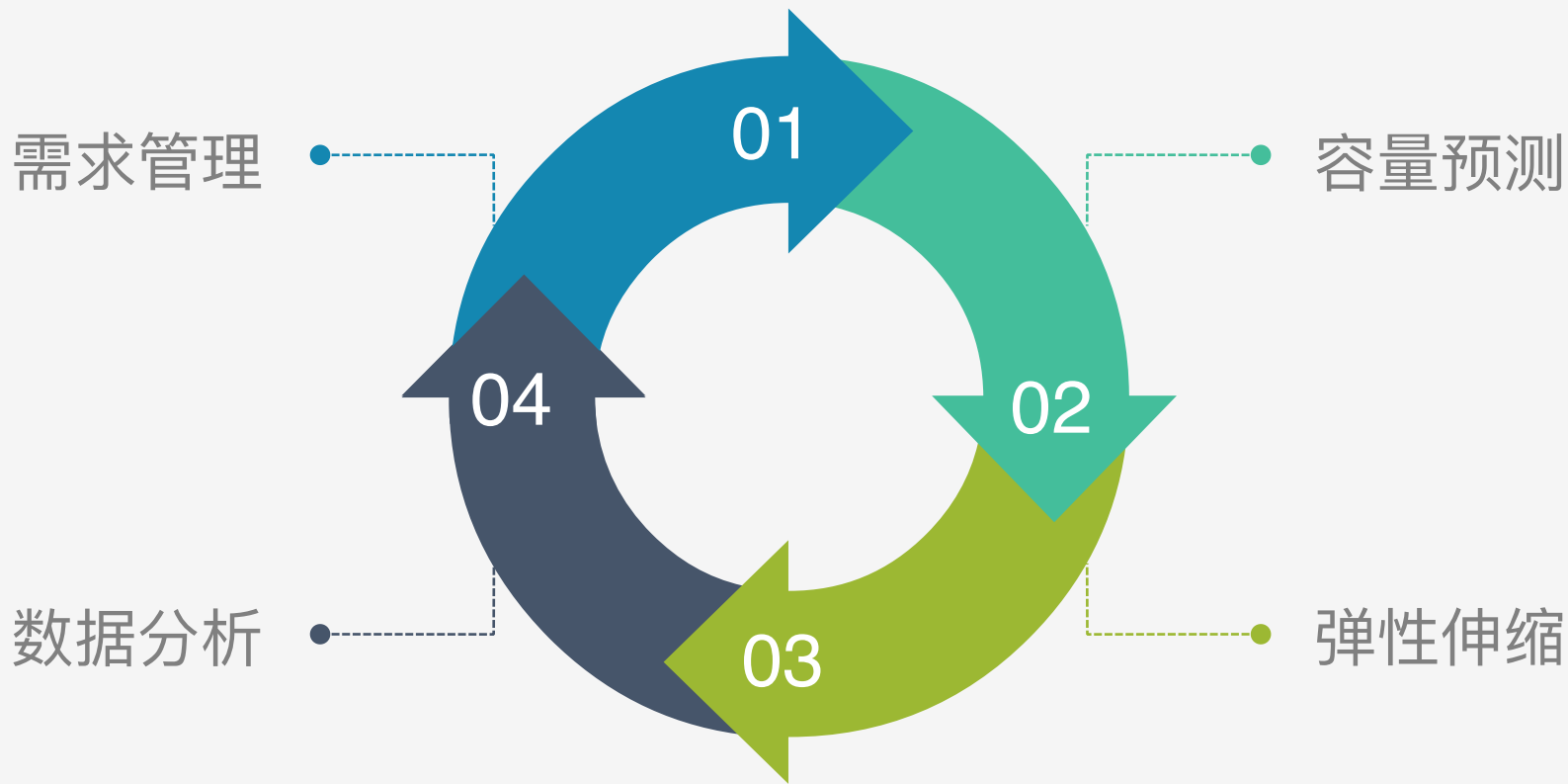
Target

“最优利用” + “容量自治”

Who

“Dev+Ops” + “业务”

智能容量管理反馈环



容量预测公式



抽象成简化公式：

$$\text{目标机器数} = \text{预测流量} / \text{应用单机能力 (预测)}$$

流量预测

自然态流量预测	利用集群流量时序特征回归
非自然态流量预测	全链路流量模型，基于业务目标各应用集群流量预测，线性回归

单机能力预测



First step

应用特征分析

1

Second step

提取关键性能指标，
建立算法模型

2

Third Step

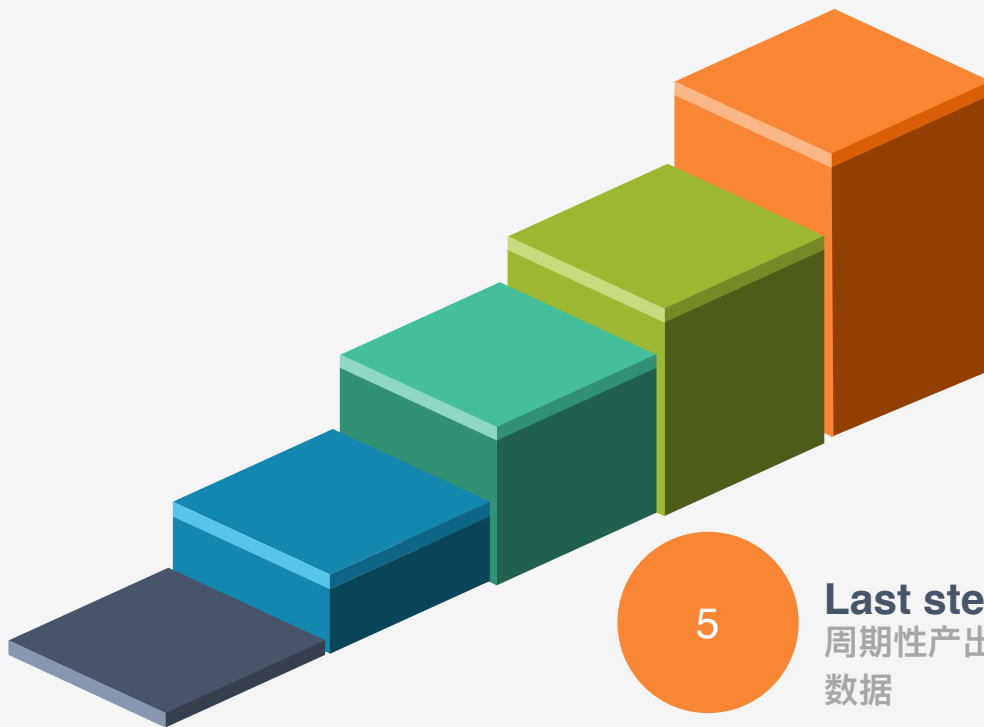
通过日常样本数据回
归预测

3

Fourth step

压测验证，效果反馈

4



5

Last step

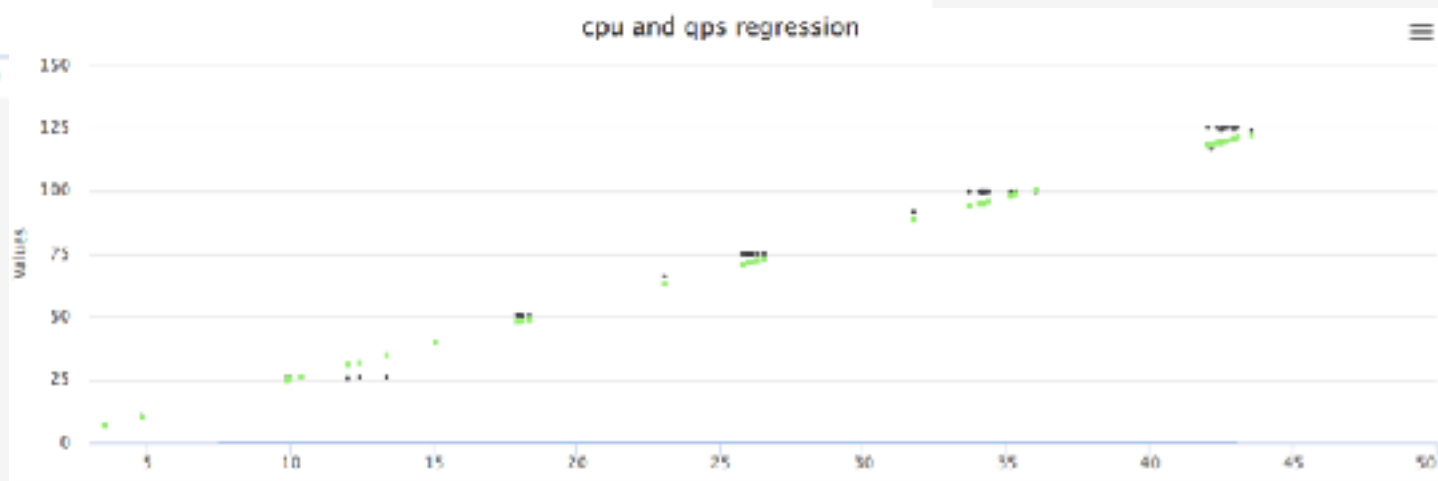
周期性产出单机性能
数据

线性回归模型

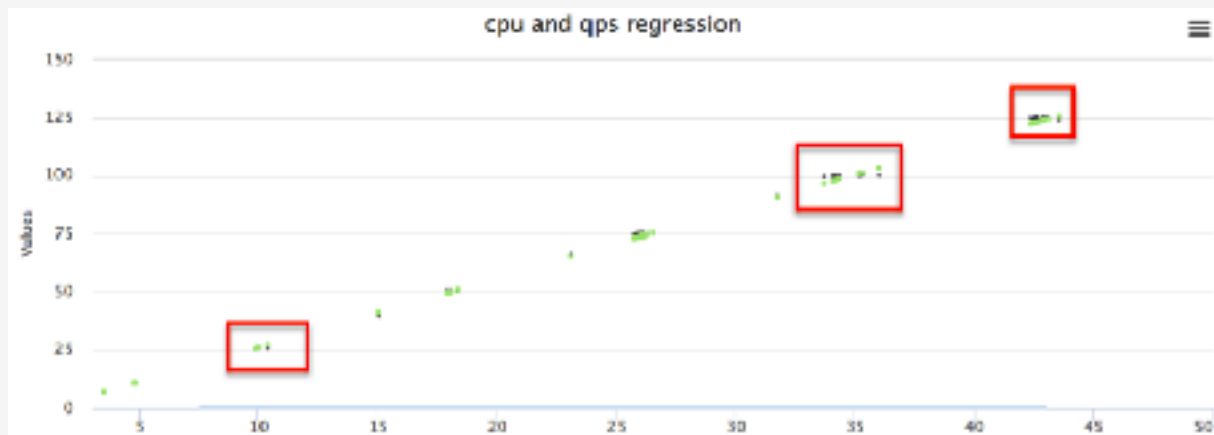


“.”为真实样本点

“.”为拟合后的点

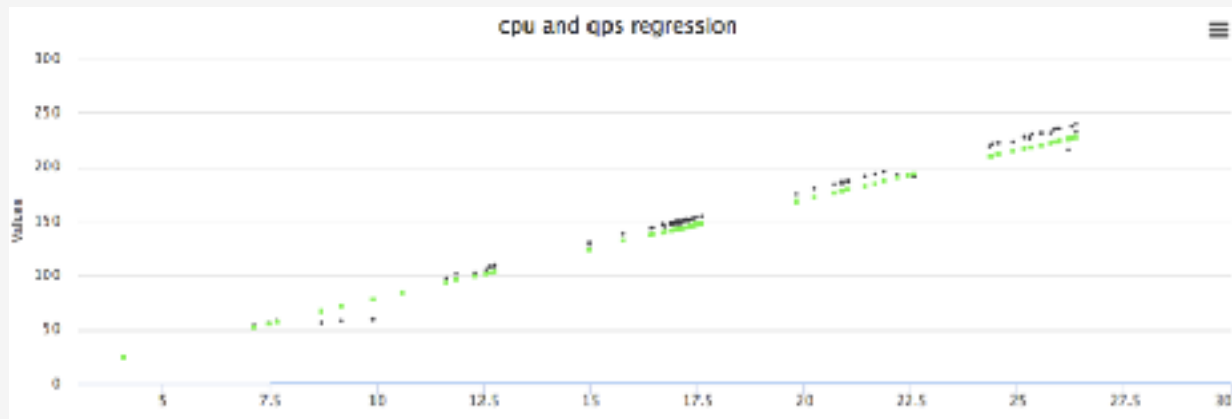


模型校正



残差降噪

局部加权



效果评估

理论评估

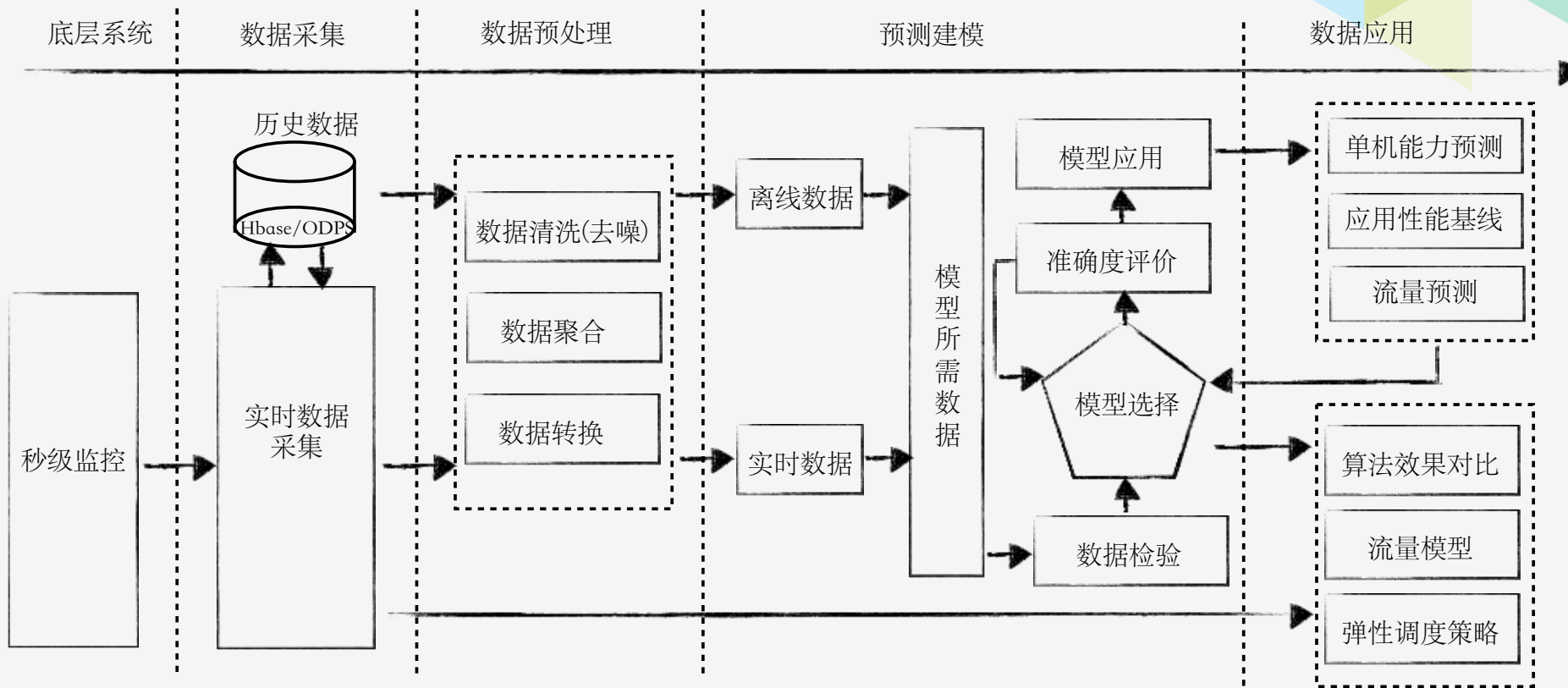
QPS {CPU(Max) / 2} 预测 QPS {CPU(Max)}

实际评估

压测验证

$$1 - \frac{\text{abs}(\text{qps真实} - \text{qps预测})}{\text{qps真实}}$$

数据处理框架



APM应用性能管理

维护应用性能基线。



04

具体实践





实践场景

01 日常弹性

02 分时复用

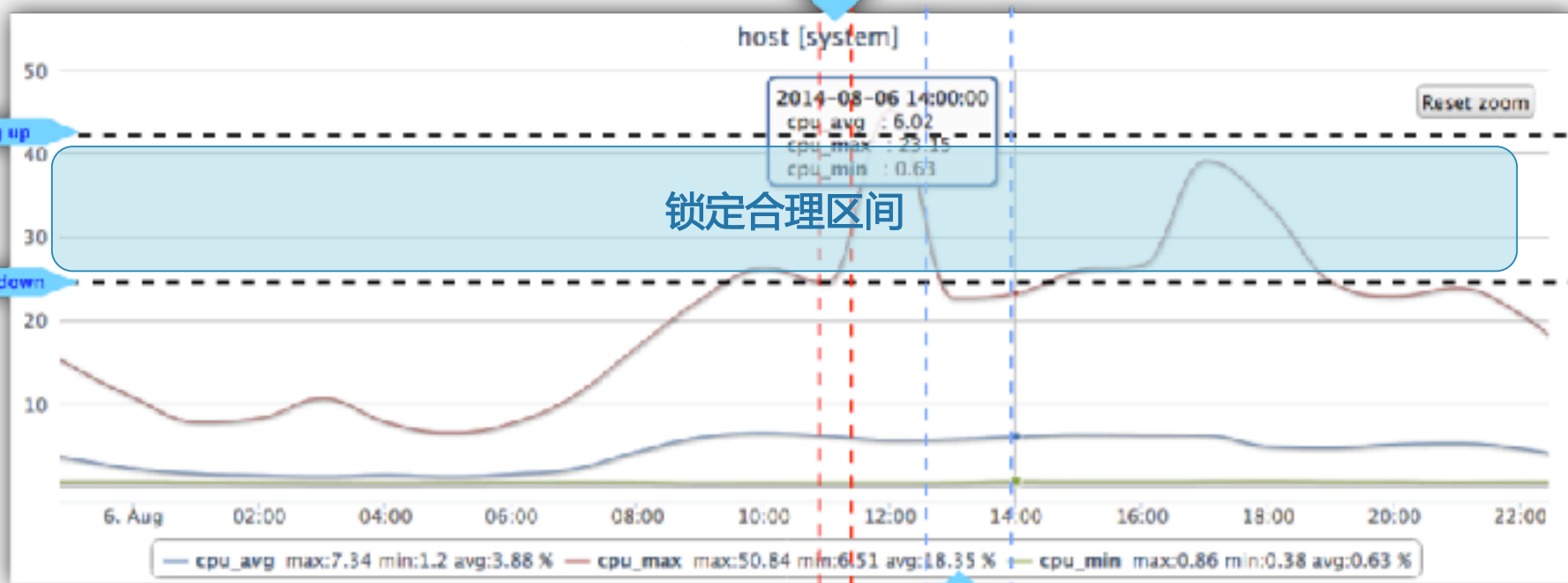
03 边压边弹



04 IDC引流评估

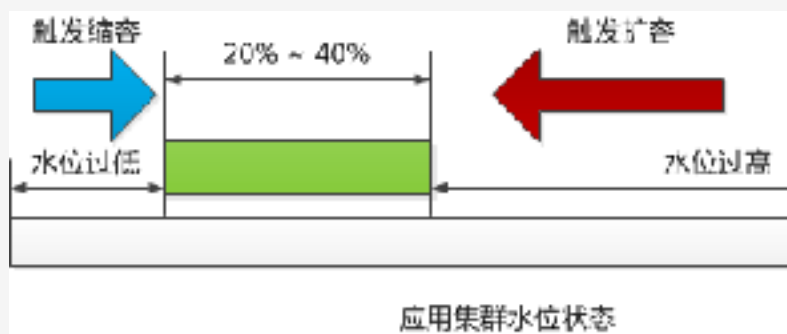
05 其他

日常弹性



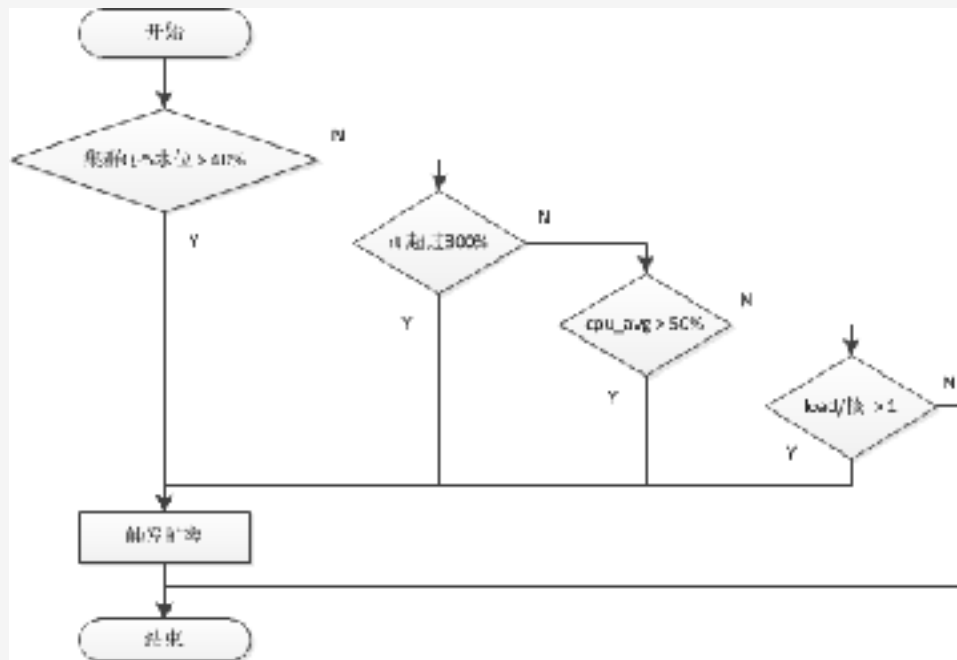
缩容
持续时间

日常弹性



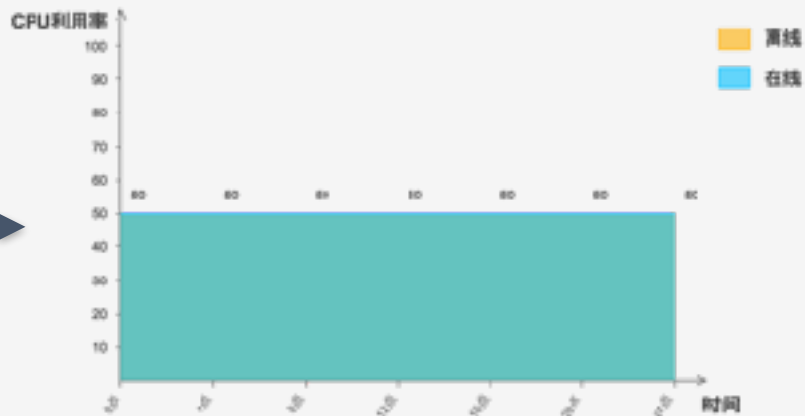
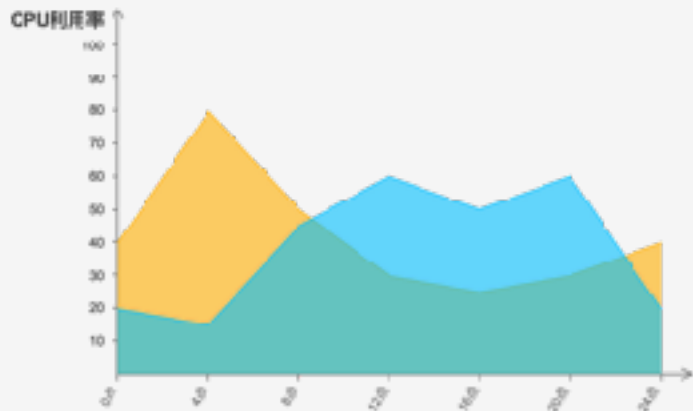
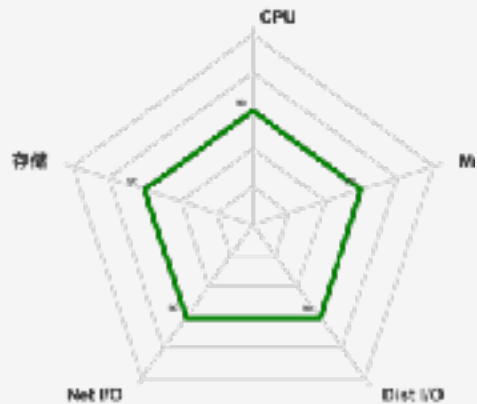
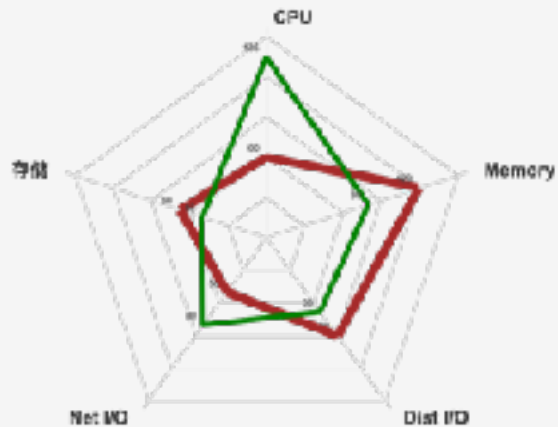
备注：集群水位 = 集群QPS / (单机QPS
极限值 * 机器数) 单机QPS极限能力

触发模式：手动、自动、定时



触发策略

分时复用背景



分时复用挑战

I 服务SLO保证

事件模型、QoS监控
和熔断机制

III 精细化分时调度

全时段精细化削峰填
谷



II

资源边界最大化

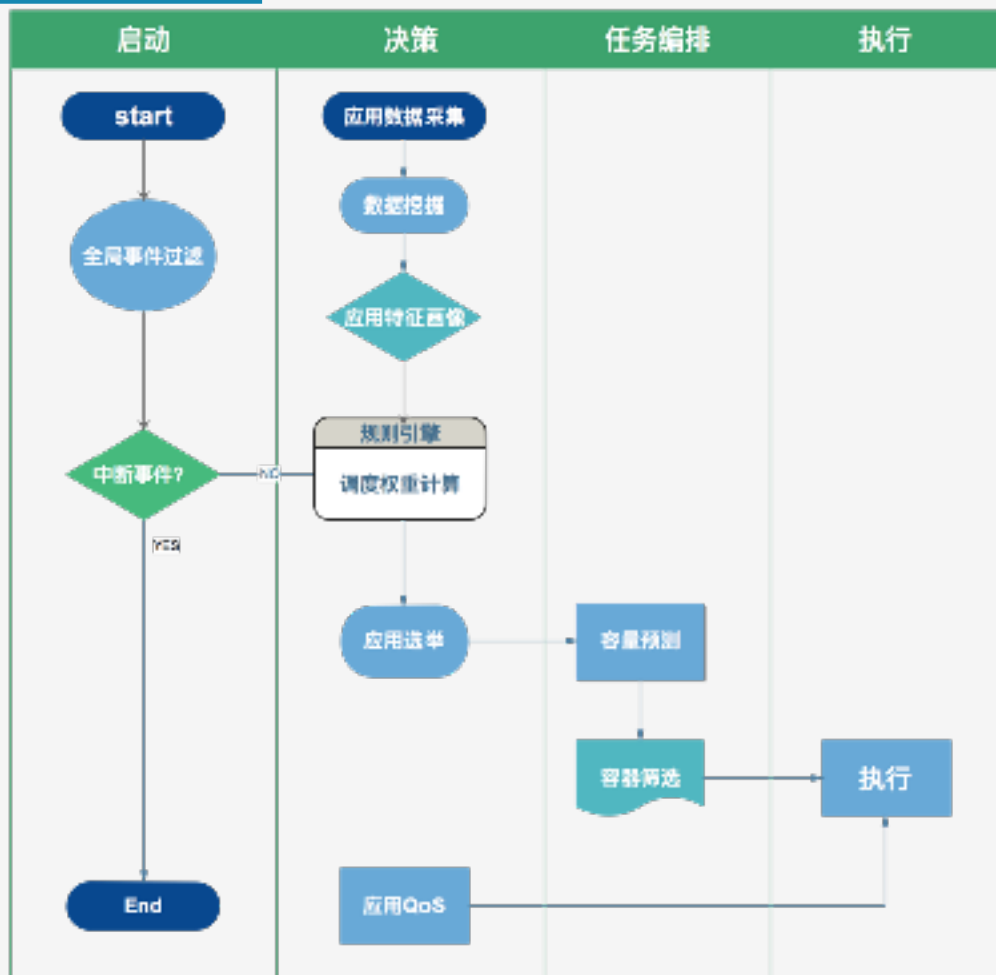
精准权重调度和容量预
测

IV

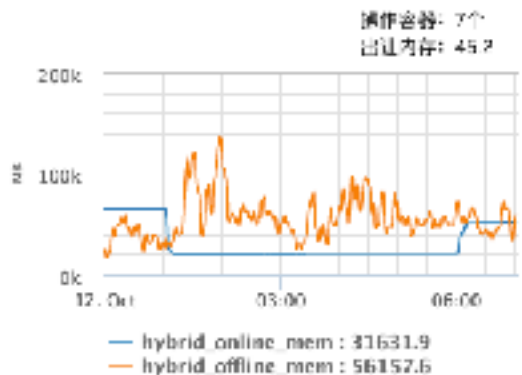
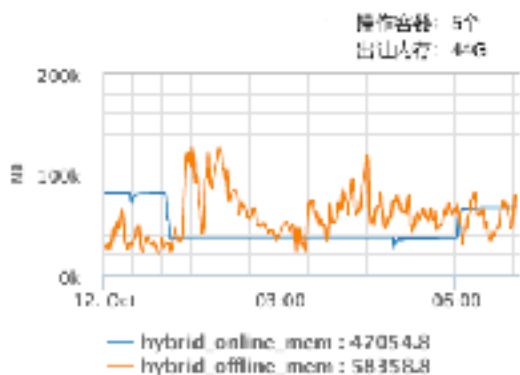
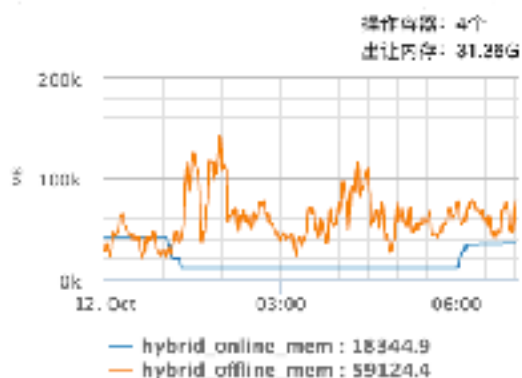
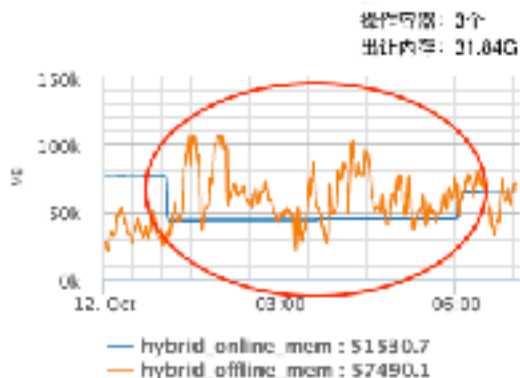
高效执行，快速 恢复

数据分析为前提，多
种执行策略

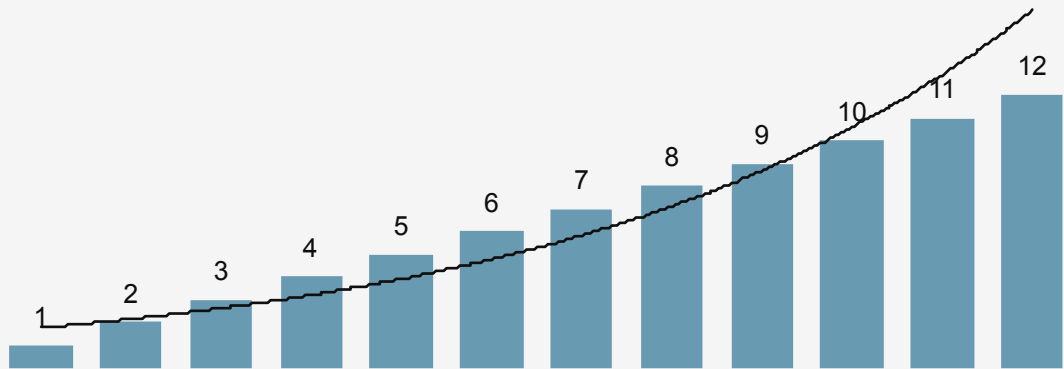
分时复用核心模块



分时复用效果

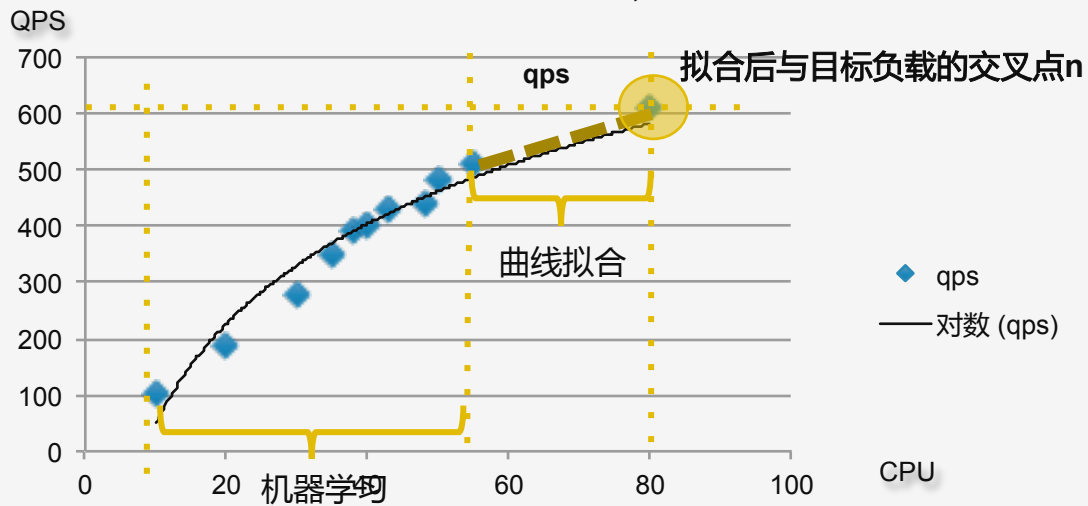


边压边弹



压测交易量增长趋势

负载预测,容量评估



◆ qps
— 对数 (qps)

边压边弹目标



压测无人值守



05

未来展望



- 智能化容量自治。
- 整体集群资源0冗余。





扫一扫上面的二维码图案，加我微信

F&Q

Please Join Us!



Thanks!