

Architect

SACC

2022 中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2022

· 激发架构性能 点亮业务活力

云上会议 网络直播 | 2022年10月20-22日

IT168.com

ChinaUnix.net

ITPUB

云原生降本之路

技术探索、创新与应用



孟凡杰

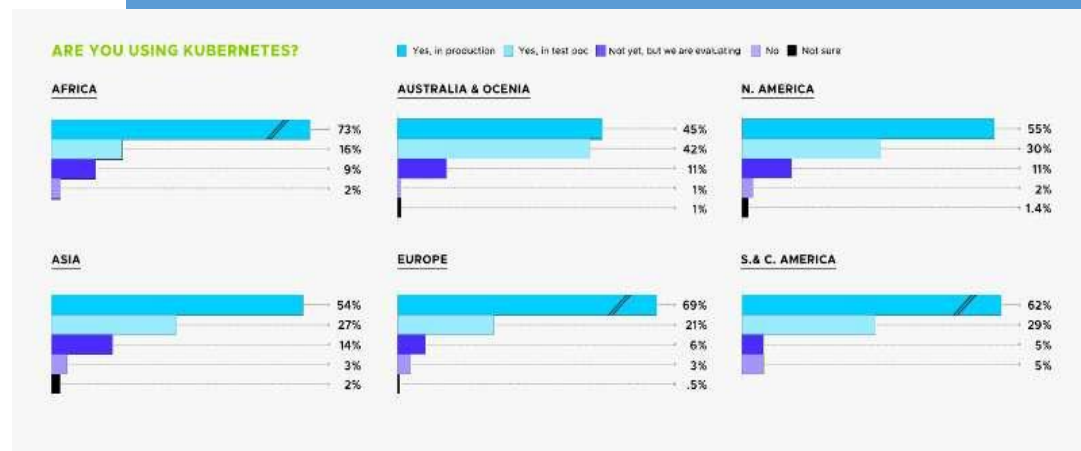
- 腾讯云容器技术专家
- FinOps产品研发负责人、致力于原生成本优化
- 《Kubernetes生产化实践之路》《软件研发效能提升实践》作者
- 极客时间《云原生训练营》讲师
- 技术峰会活跃讲师、出品人

云原生资源利用现状

成本优化成为企业上云的核心关切

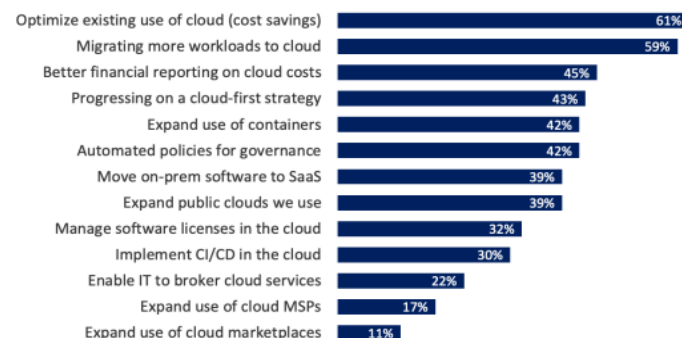


云原生基金会2021年调查显示，云原生的部署率已经达到调查样本的历史性新高
95%的组织已经在调研或使用Kubernetes



Flexera 发布的《2021 云计算市场发展状态报告》
30%-35% 的云支出被浪费了

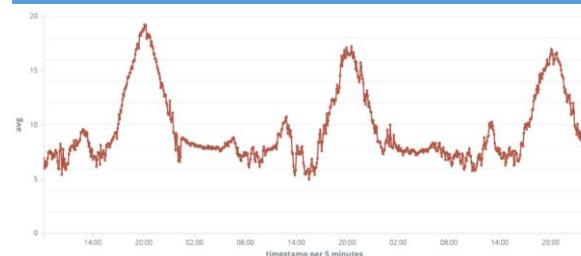
Top Cloud Initiatives for 2021 % of all respondents



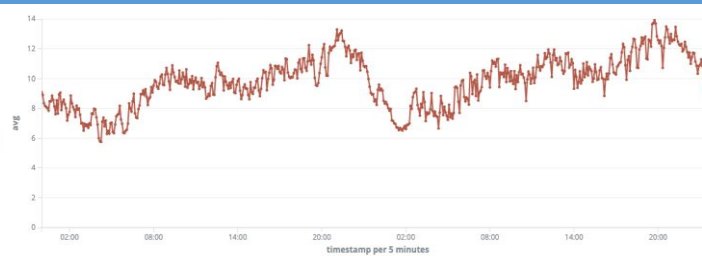
N=750

Source: Flexera 2021 State of the Cloud Report

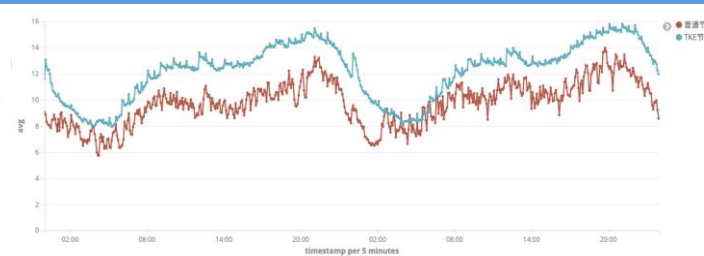
基于腾讯云公有云客户数据分析和调研，客户集群中资源成本浪费非常严重，有众多客户提出关于提高资源利用率的诉求。



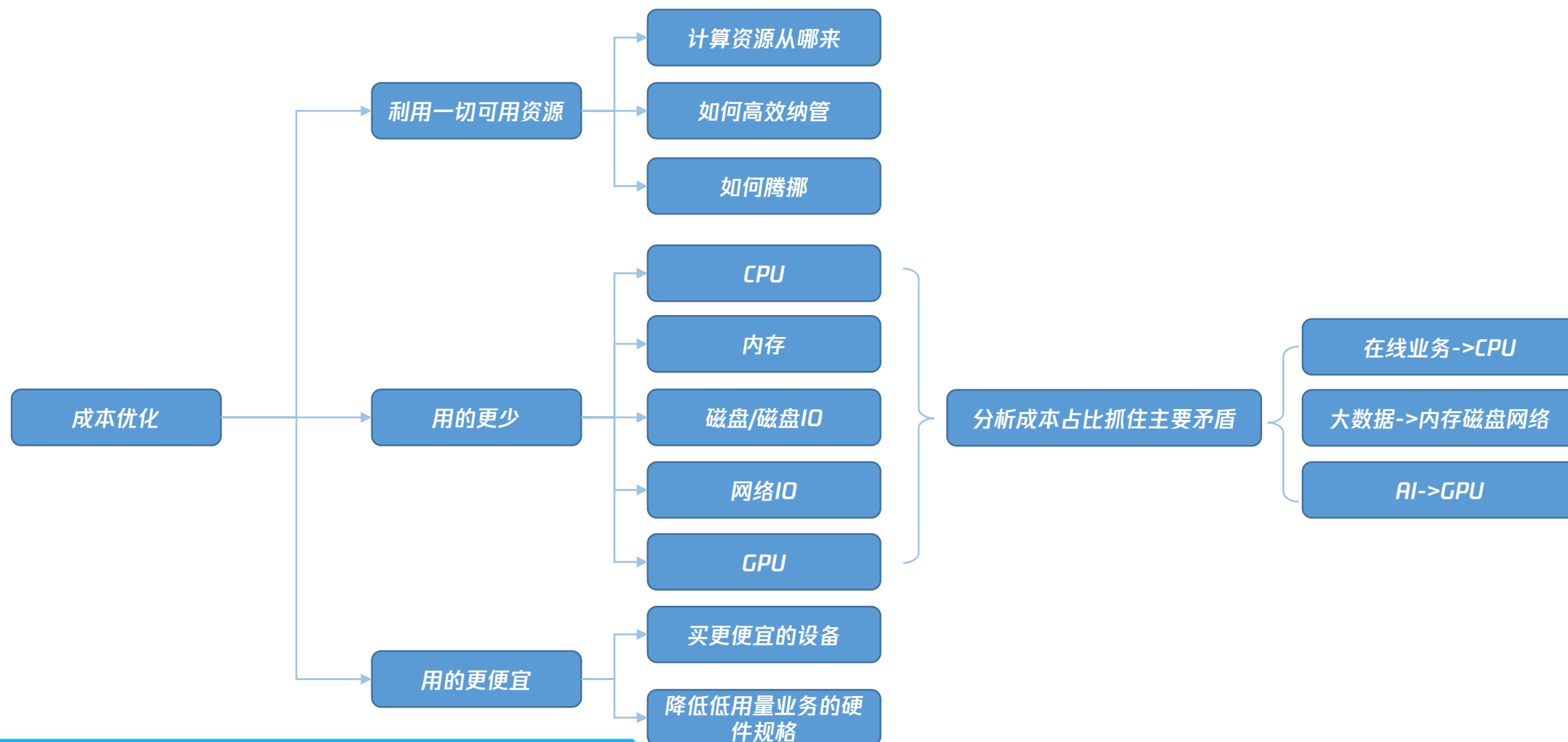
物理机利用率: 10%



虚拟机利用率: 12%



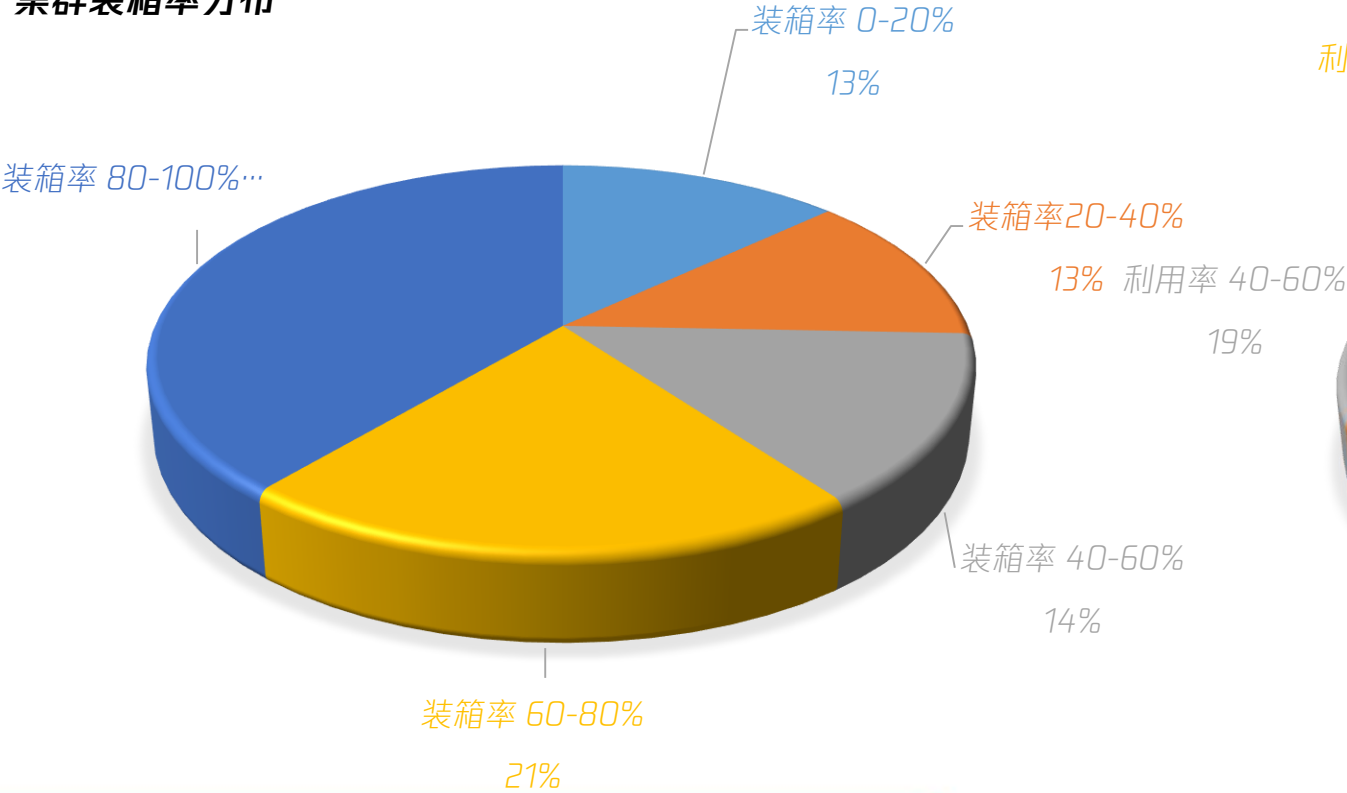
容器化利用率: 14%



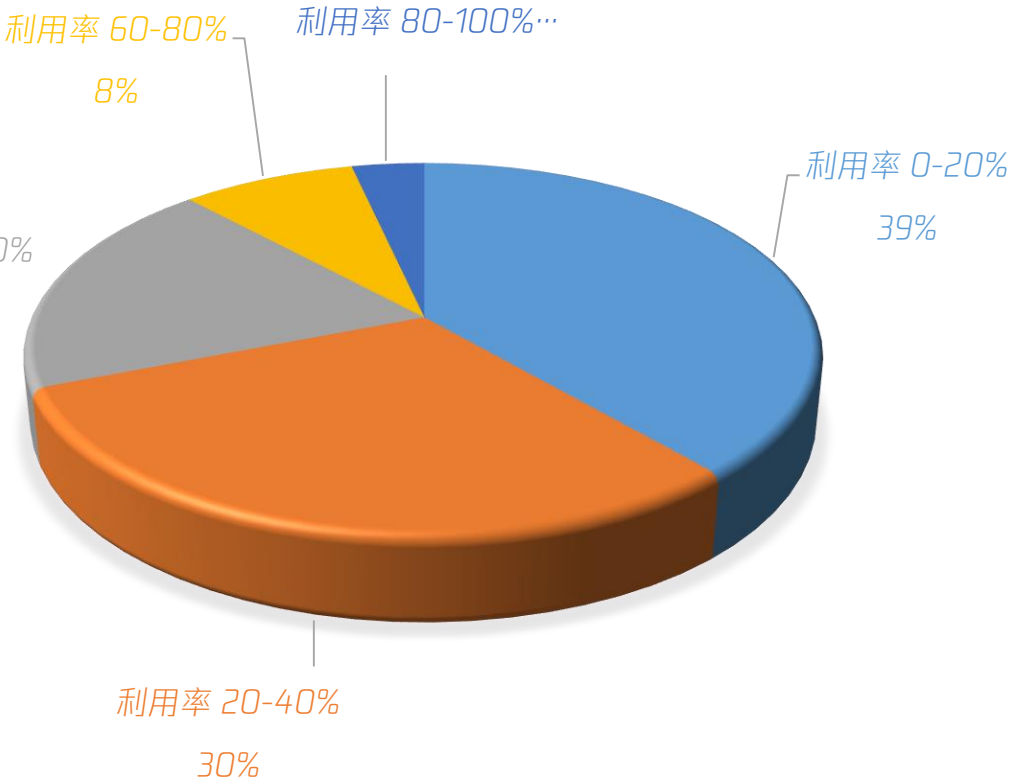
集群装箱率与利用率分布

节点装箱率参差不齐：近一半集群装箱率不足50%
节点利用率低：三分之二集群峰值利用率不到40%

集群装箱率分布

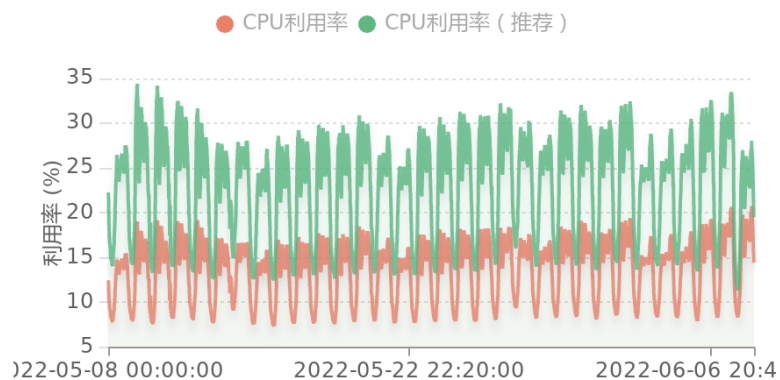


集群峰值利用率分布

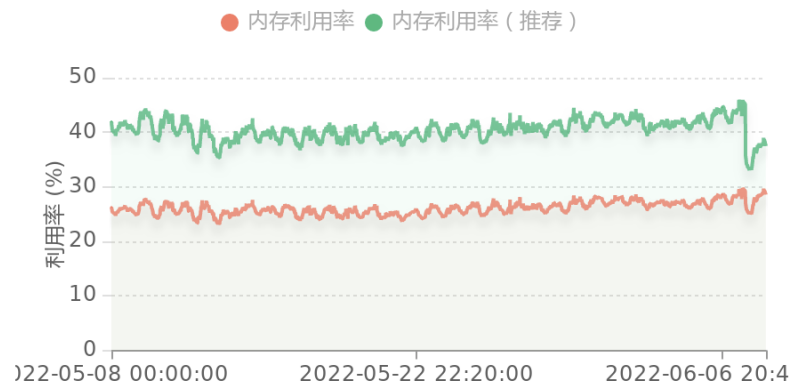


业务资源利用率低：cpu利用率15%，内存利用率25%

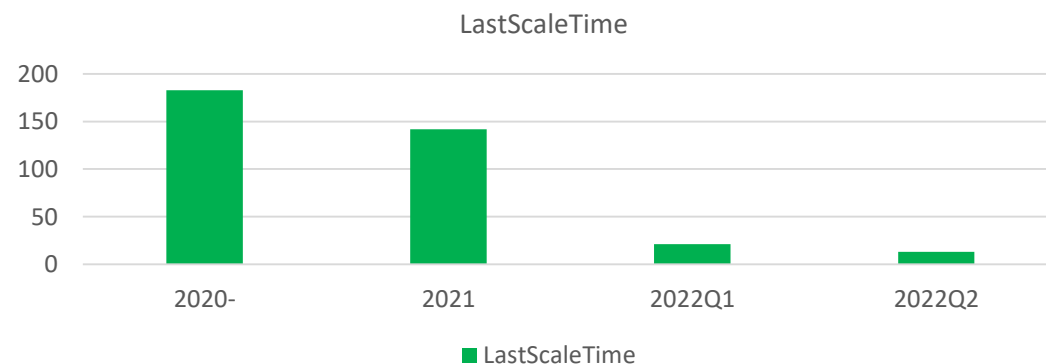
CPU利用率



内存利用率



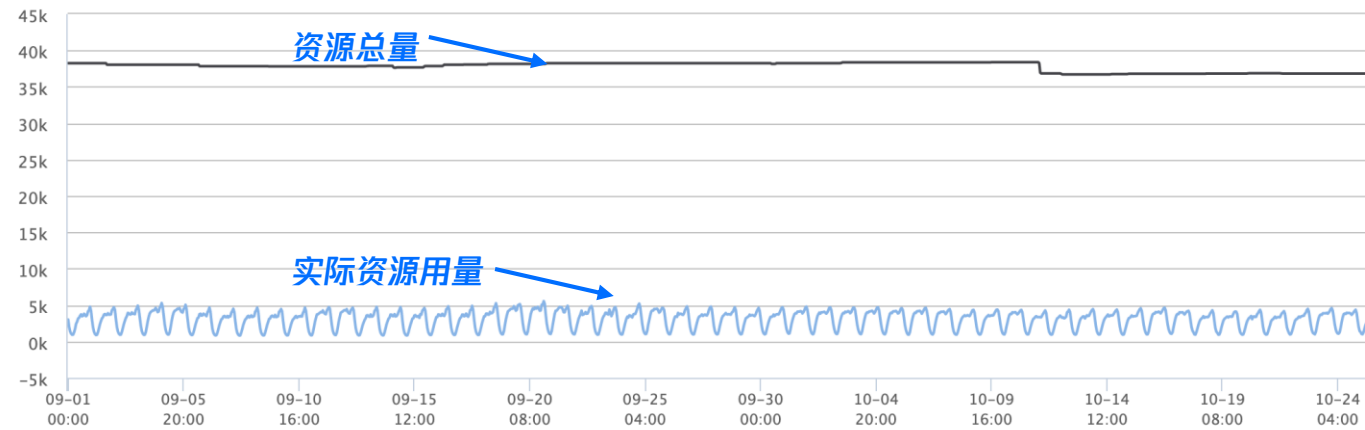
有效弹性占比低：只有10%的HPA在本年度弹出过



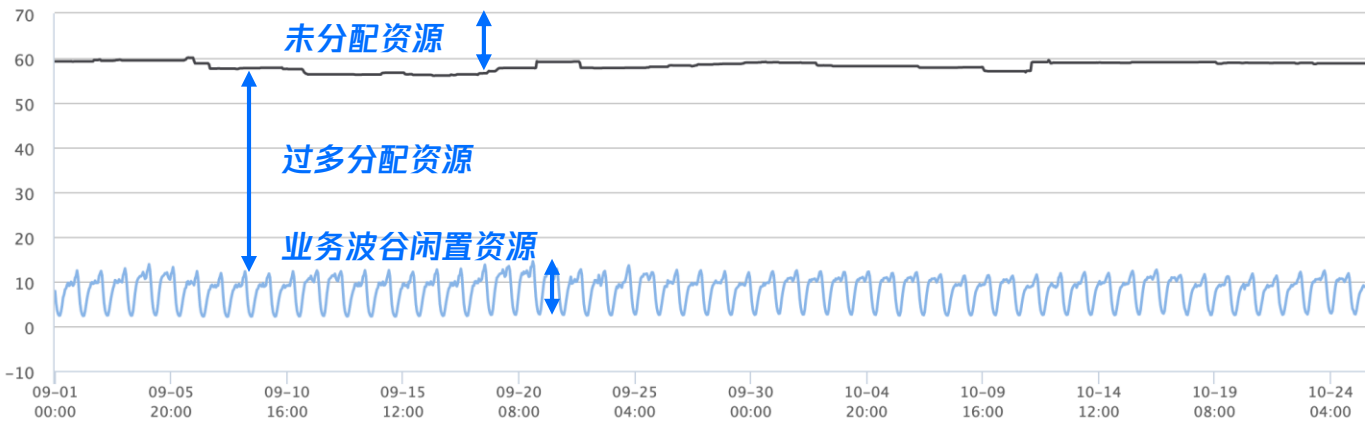
资源浪费原因分析

典型资源利用率-浪费的来源

资源总量以及用量

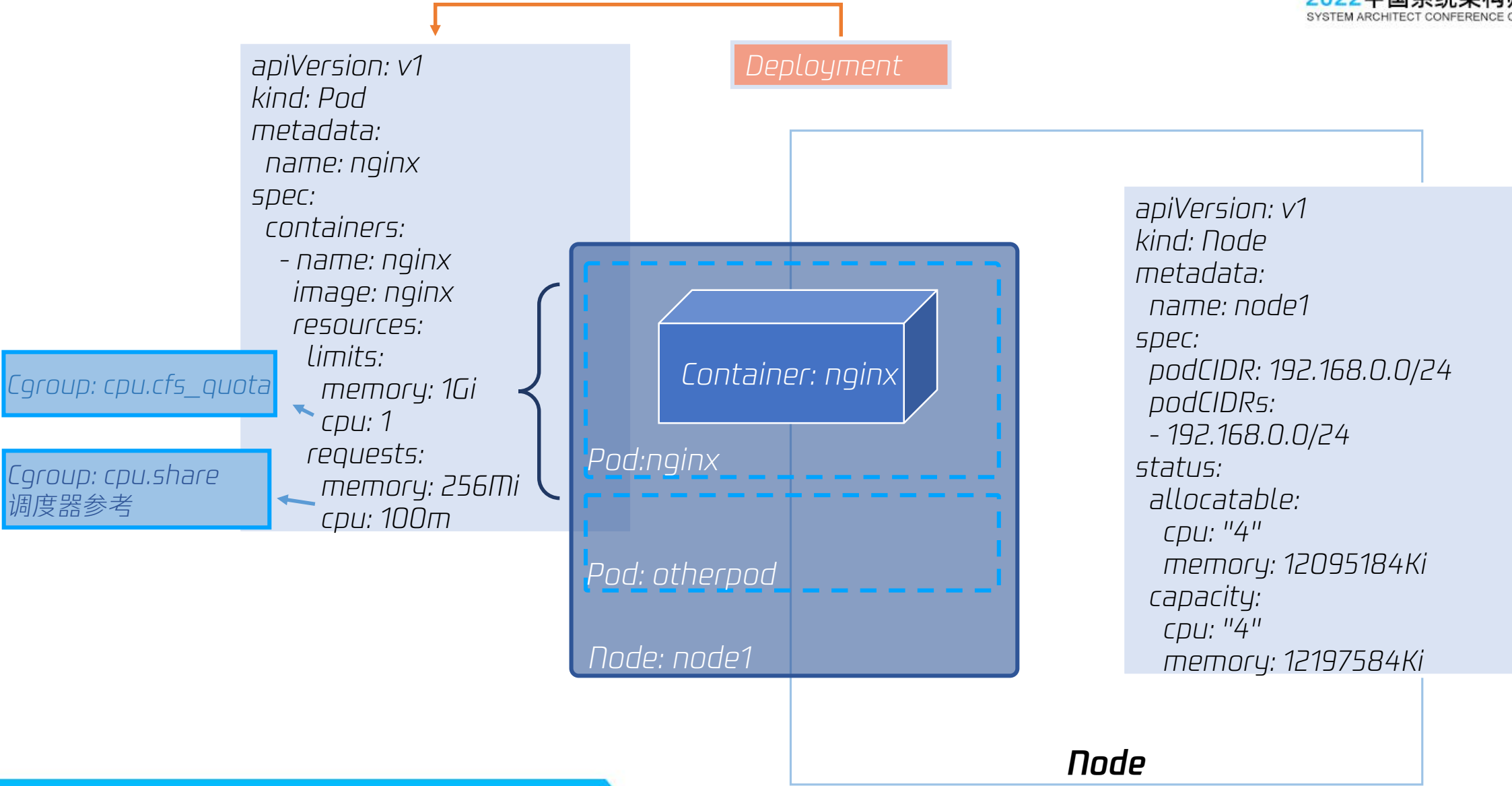


资源分配率和使用率



云成本管理的核心：在保障业务的前提下，最小化资源需求

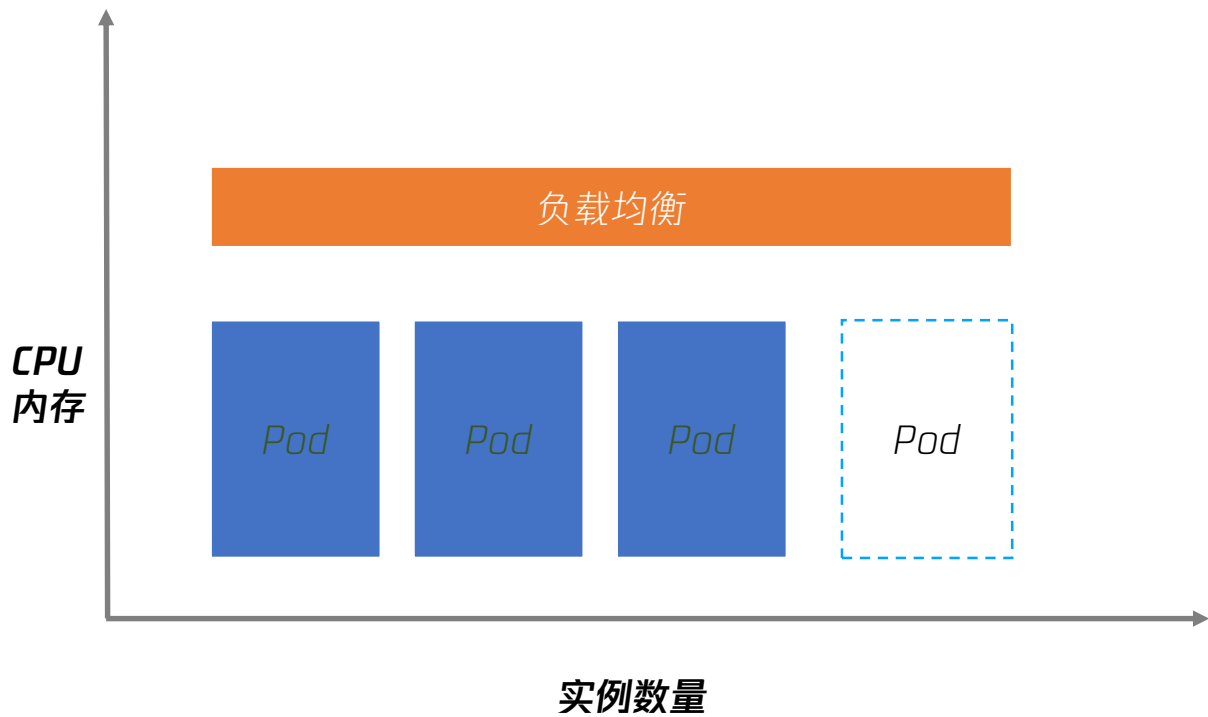
Kubernetes 中的资源分配



Node

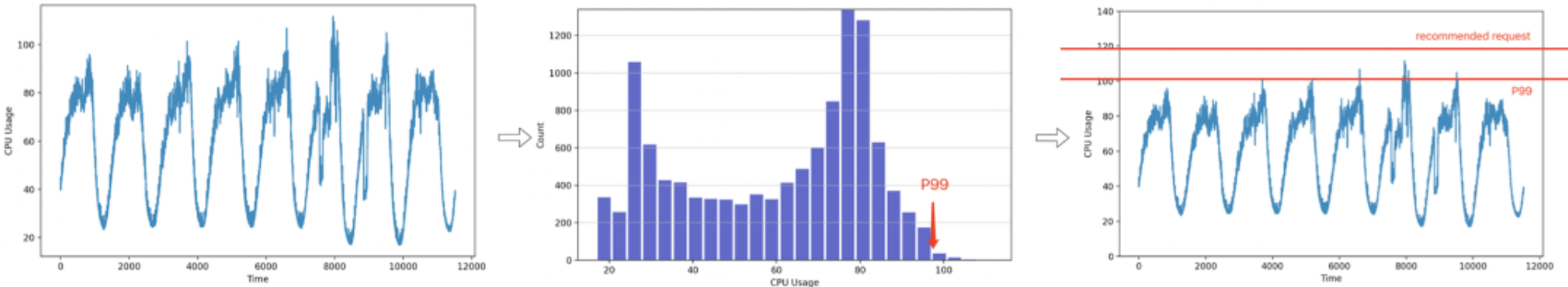
横向伸缩和纵向伸缩

- 应用扩容是指在应用接收到的并发请求已经处于其处理请求极限边界的情形下，扩展处理能力而确保应用高可用的技术手段
- Horizontal Scaling
 - 所谓横向伸缩是指通过增加应用实例数量分担负载的方式来提升应用整体处理能力的方式
- Vertical Scaling
 - 所谓纵向伸缩是指通过增加单个应用实例资源以提升单个实例处理能力，进而提升应用整体处理能力的方式



Kubernetes原生能力的不足

业务现状



配置手段



资源配置

- 基于经验的资源配置不准导致大量浪费



弹性

- 基于阈值的弹性的滞后性导致业务来不及弹



业务稳定性

- CPU是可压缩资源，CPU承压时，不驱逐，所有Pod等比受损
- 发生CPU抢占时，以 `cpu.shares` 公平分配时间片
- 无法确保延迟敏感型业务的稳定性
- 独占式绑核能力造成较大资源浪费

面临挑战

不会配

不敢配

不能配

全链路降本的思考和价值主张

问题

1. 运维侧无法从全局视角管控装箱率, 每个业务单独配置超卖比

1. request配置不合理
2. K8s原生调度策略是默认均衡优先
3. 在提升节点使用率时, 负载过高可能影响业务稳定

1. 业务存在波峰波谷
2. 存在复杂任务类型时 [高优在线任务、低优离线任务], 难以动态的完成资源管理和调度, 确保在线业务不受影响

步骤

1. 利用一切可用资源、提升集群装箱率

2. 提升资源利用率峰值

3. 提升资源利用率均值

方案

- 多种集群和节点形态
- 节点容量缩放

1. 基于负载画像调度
2. 负载感知调度
3. 重调度器进行实时负载再平衡

1. 削峰填谷
 1. 弹性
 2. 混布
 3. 服务质量保证

多种集群形态

利用一切可利用计算资源

Serverless 集群

集群特点

- 集群控制面托管到 TKE Meta Cluster 集群
- Nodeless, 直接部署应用
- 轻量级虚拟机, 快速弹出资源

适用场景

- 希望简化集群运维、提升集群稳定性
- 对容器隔离性有强需求
- 对容器弹性速度有强需求
- 希望降低集群使用成本

TKE 托管集群

集群特点

- 集群控制面托管到TKE Meta Cluster集群
- 用户直接添加Slave节点
- 控制面资源免费
- 支持自定义组件参数

适用场景

- 对集群中存在 Node 节点有强需求
- 对操作系统定制化有强需求
- 希望使用 TKE 提供的 qGPU 方案

TKE 独立集群

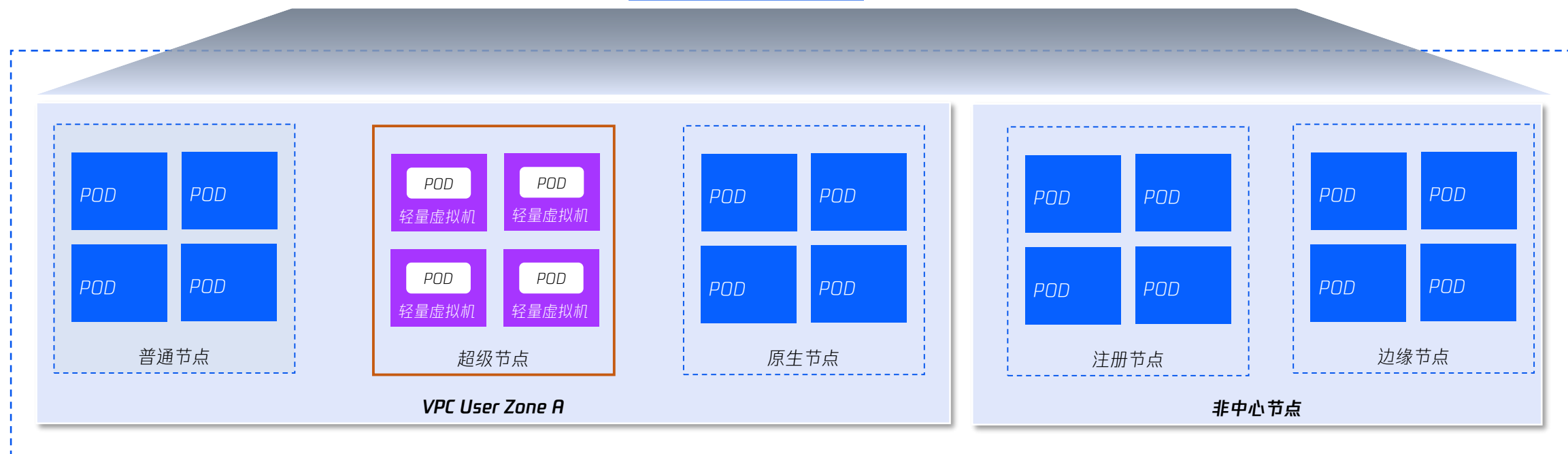
集群特点

- 用户创建 CVM 节点运行 k8s 控制面组件
- 根据集群规模选择 Master 节点配置
- Master 按节点配置收费

适用场景

- 对控制面监控数据和日志有要求
- 测试集群, 有自定义 CRD 需要依赖控制面日志做调试
- 自建 K8S, 有 K8S 定制需求

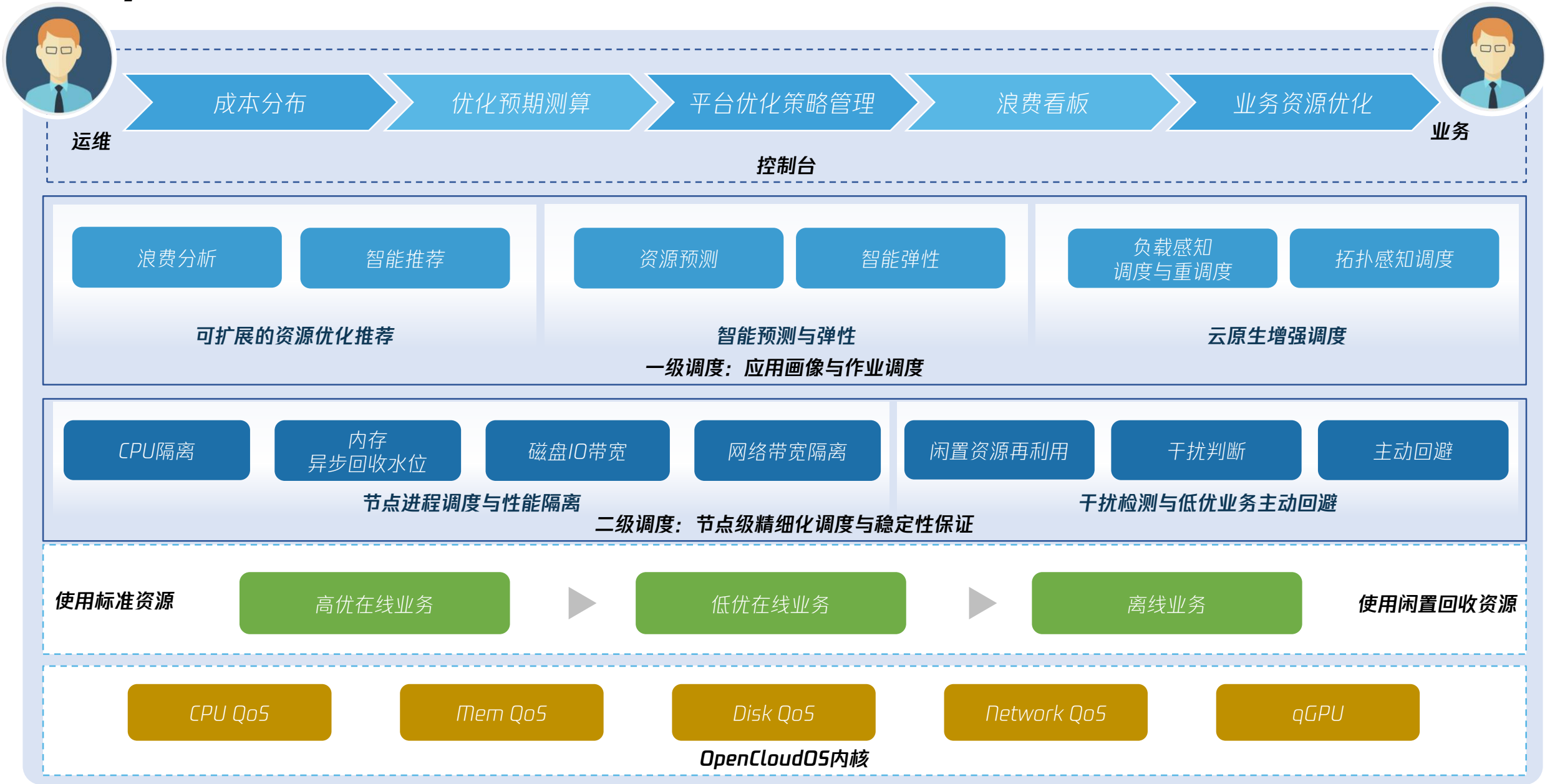
托管集群



- 按照集群资源付费
- 需要运维集群节点
- 容器共享OS内核
- 安全隔离性弱
- 全托管
- 按照Pod资源付费
- 无需运维集群节点、极致弹性
- 虚拟机级别隔离性
- 全托管
- Yaml管理资源
- 轻运维
- Finops
- 需要专线
- 无需本地部署Master
- 旧代次设备低成本接入
- 弱网支持
- 设备分钟
- 分布式健康检查

Crane - 智能调度系统助力 资源优化

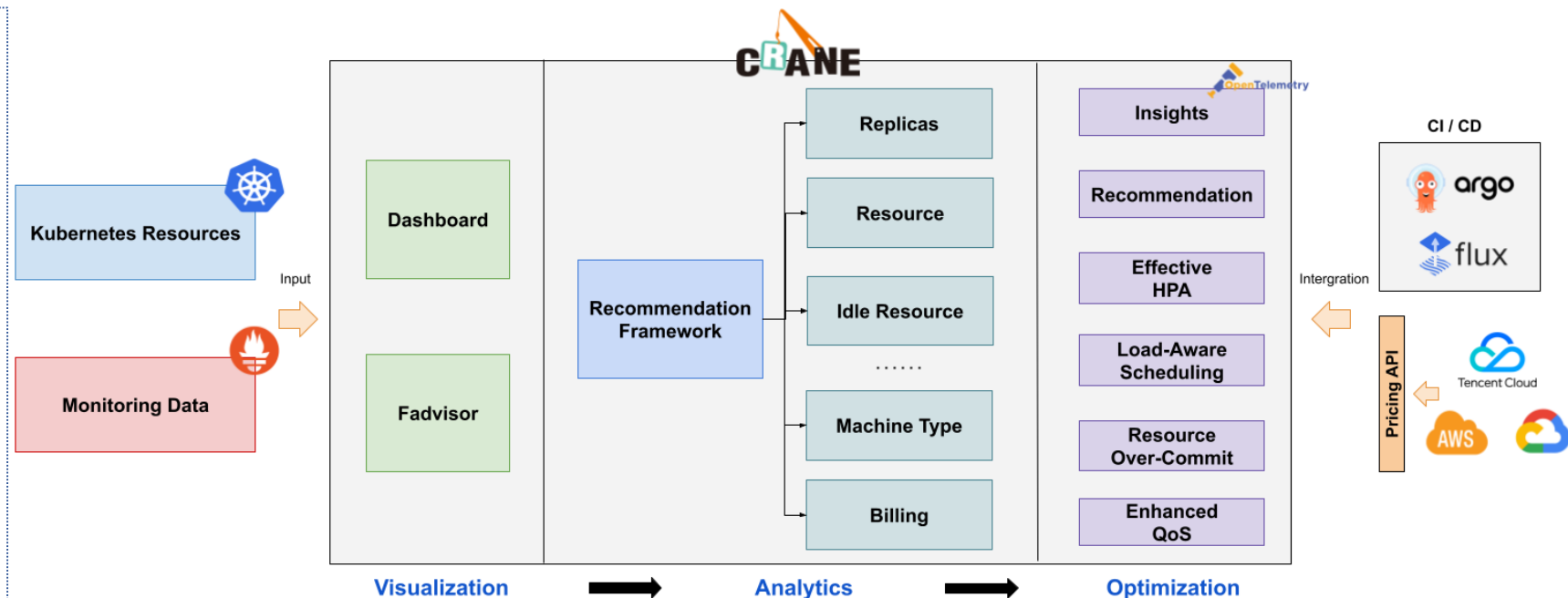
FinOps理念指导腾讯内部海量集群落地经验系统性输出



统一推荐框架

A FinOps framework for Cloud Resource Analytics and Economics

```
apiVersion: analysis.crane.io/v1alpha1
kind: RecommendationRule
metadata:
  name: workload
spec:
  namespaceSelector:
    any: true
  recommenders:
    - name: WorkloadReplicas
    - name: WorkloadResource
  resourceSelectors:
    - apiVersion: apps/v1
      kind: Deployment
  runInterval: 1h
```



统一推荐框架

- Filter
- Prepare
- Recommend
- Observe

灵活扩展的推荐插件

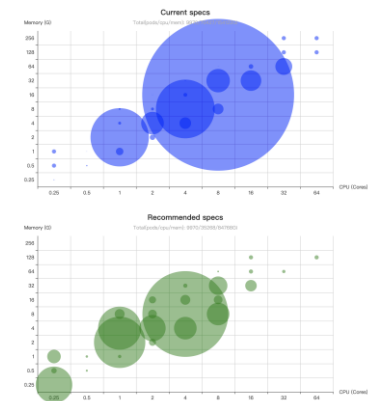
- Request推荐
- 副本数推荐
- 弹性推荐
- 闲置PVC优化建议

多种结果展示方式

- 自动化Action
- 成本展示
- 功耗与碳排放量展示

生态集成

- Kubesphere成本优化套件
- 与ArgoCD等流水线集成



FFT预测与智能弹性

提前扩容

- 时间序列算法: FFT快速傅里叶变换
- 取预测窗口最大值: 提前扩容
- 基于 Custom Metric
- metric 兜底保护

减少无效缩容

- 预测未来可以减少不必要的缩容

支持 Cron 配置

- 应对大促节假日等有规律的流量洪峰

易于使用

- 完全兼容社区 HPA
- 支持 Dryrun 观测
- 指标支持 Prometheus Metric

```
apiVersion: autoscaling.crane.io/v1alpha1
kind: EffectiveHorizontalPodAutoscaler
metadata:
```

```
  name: php-apache
```

```
spec:
```

```
  scaleTargetRef:
```

```
    apiVersion: apps/v1
```

```
    kind: Deployment
```

```
    name: php-apache
```

```
  minReplicas: 1
```

```
  maxReplicas: 10
```

```
  scaleStrategy: Auto
```

```
  metrics:
```

```
    - type: Resource
```

```
      resource:
```

```
        name: cpu
```

```
      target:
```

```
        type: Utilization
```

```
        averageUtilization: 50
```

```
prediction:
```

```
  predictionWindowSeconds: 3600
```

```
  predictionAlgorithm:
```

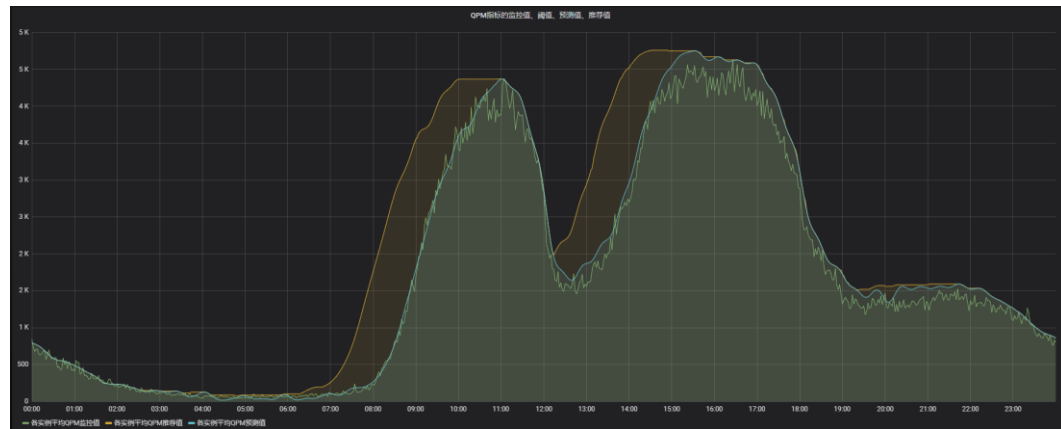
```
    algorithmType: dsp
```

```
  dsp:
```

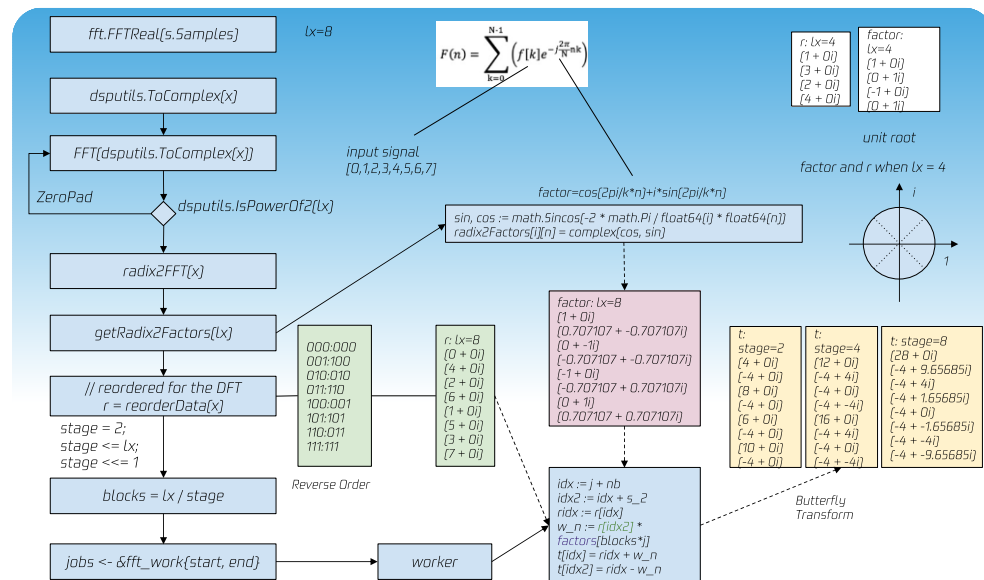
```
    sampleInterval: "60s"
```

```
    historyLength: "3d"
```

预测弹性效果



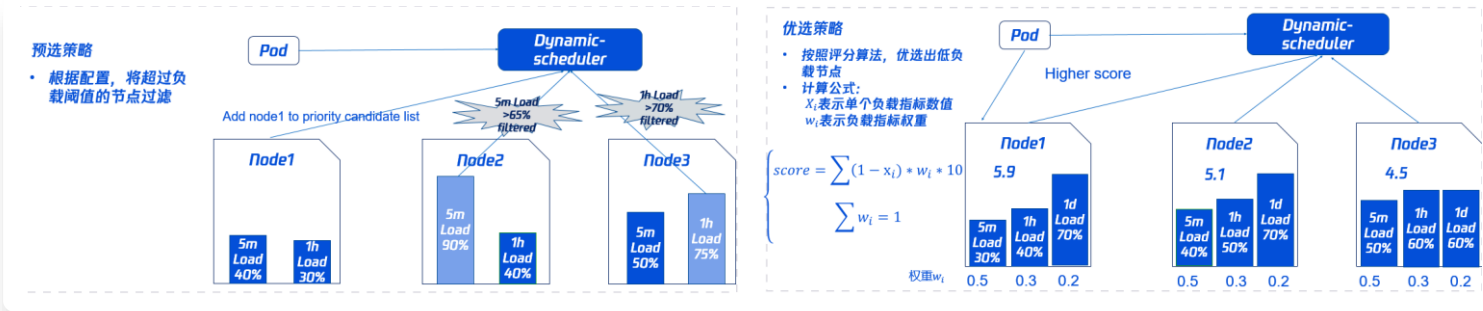
—— 实时Metric —— 预测Metric —— 弹性计算Metric



增强调度

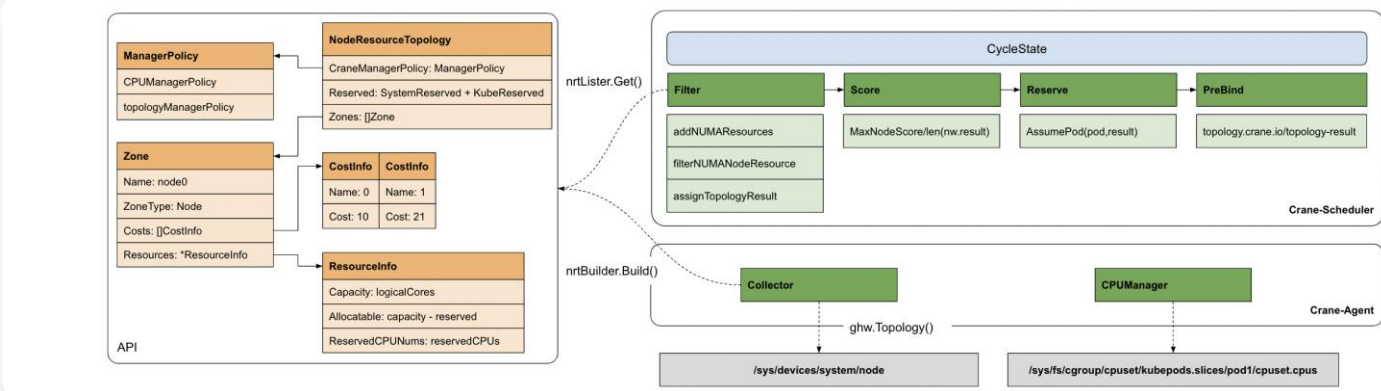
负载感知调度

- 基于真实负载
- 引入多维负载指标，感知业务波峰
 - 1h内最大利用率
 - 1天内最大利用率



拓扑感知调度

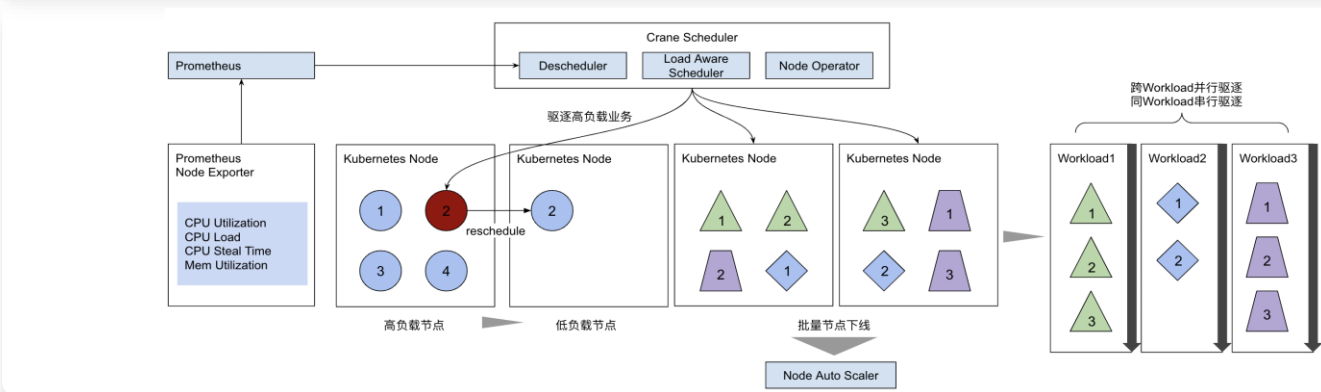
- 性能优先，Numa亲和
- 负载均衡，选择空闲的Numa



来自TEG星辰算力团队的技术沉淀

重调度器

- 驱逐负载过高的节点上的低优Pod
- 支持固定时间窗口腾空节点
- 支持基于特定标签批量驱逐Pod



基于多优先级业务的混布

NodeQOS: 节点指标水位

rules:

- actionName: eviction
- avoidanceThreshold: 2
- metricRule:
 - name: cpu_total_utilization
 - value: 80
 - name: cpu-usage

PodQOS: 业务分级与规则

allowedActions:

- eviction
- scopeSelector:
 - labelSelector:
 - matchLabels:
 - app-type: offline

AvoidanceAction: 回避动作

coolDownSeconds: 300

description: evict low priority pods

eviction:

- terminationGracePeriodSeconds: 30



CPU QoS

进程优先级CPU Burst

CPUSet管理超线程隔离

节点弹性资源水位控制

Mem QoS

异步回收全局分级水位

pagecache上限全局分级水位

节点全局pagecache管理

DiskIO QoS

Direct IOPSBuffered IOPS

Direct BPSBuffered BPS

NetIO QoS

入站流量限速出站流量限速

Pod分级限速端口白名单

节点网卡限速

全维度性能隔离

业务指标采集

节点指标采集

资源预测

闲置资源回收

干扰判断

主动回避

干扰检测与主动回避

实践案例 - 腾讯内部自研落地

FinOps Framework

FinOps定义了一系列云财务管理规则和最佳实践，通过助力工程和财务团队、技术和业务团队彼此合作，进行数据驱动的成本决策，使组织能够获得最大收益。

原则

团队协作
成本节省，人人有责
中心化团队驱动FinOps
实时报表助力决策
业务价值驱动决策
灵活利用云上成本模型

角色



能力

理解用量和成本

绩效跟踪和展示

实时决策

费率优化

用量优化

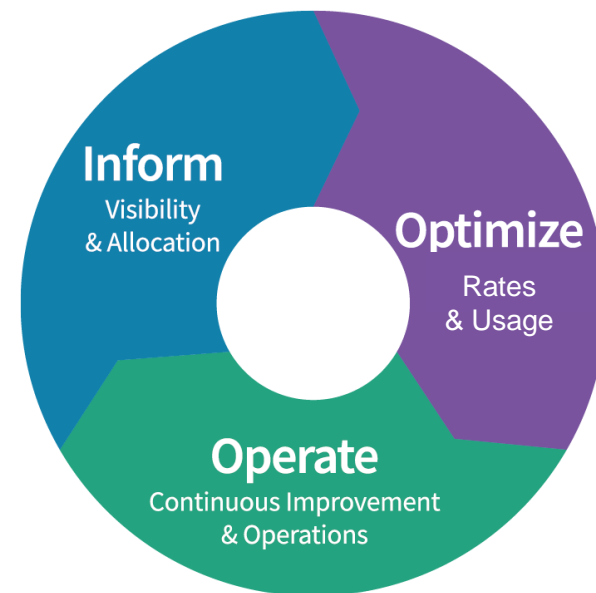
组织支撑

成熟度

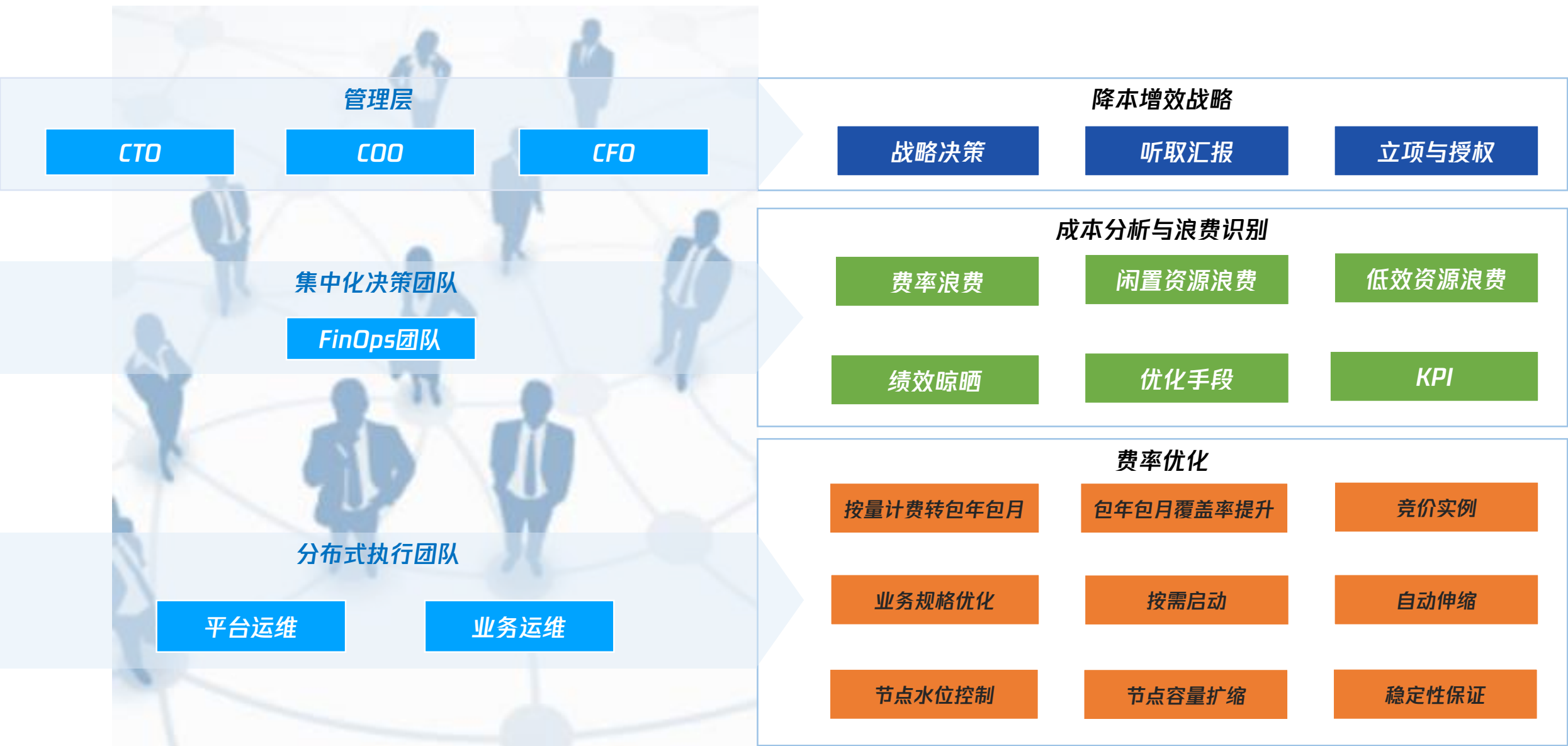


阶段

Crawl Walk Run

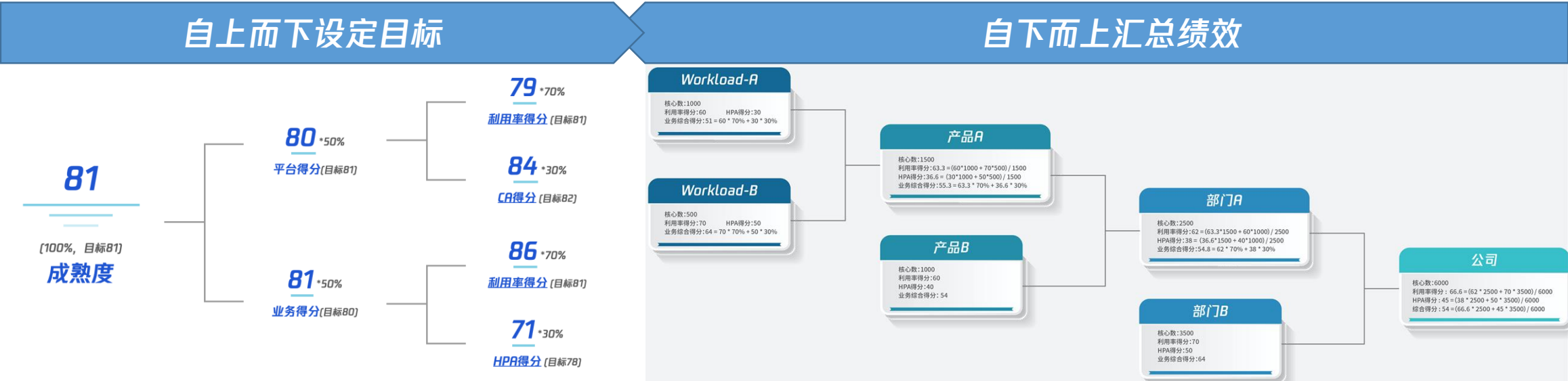


如何有效推动降本



目标设定与绩效晾晒

- 为优化腾讯内部业务云资源，腾讯定义了云成熟度模型
- 该模型从平台侧以及业务侧考核各个BG的云资源使用情况
- 总成熟度得分 = 业务侧得分 * 50% + 平台侧得分 * 50%
- 得分情况从作业，产品，部门维度层层汇总，并以该结果作为考核参考指标



数据驱动的成本分析与成果测算

• 离线数仓

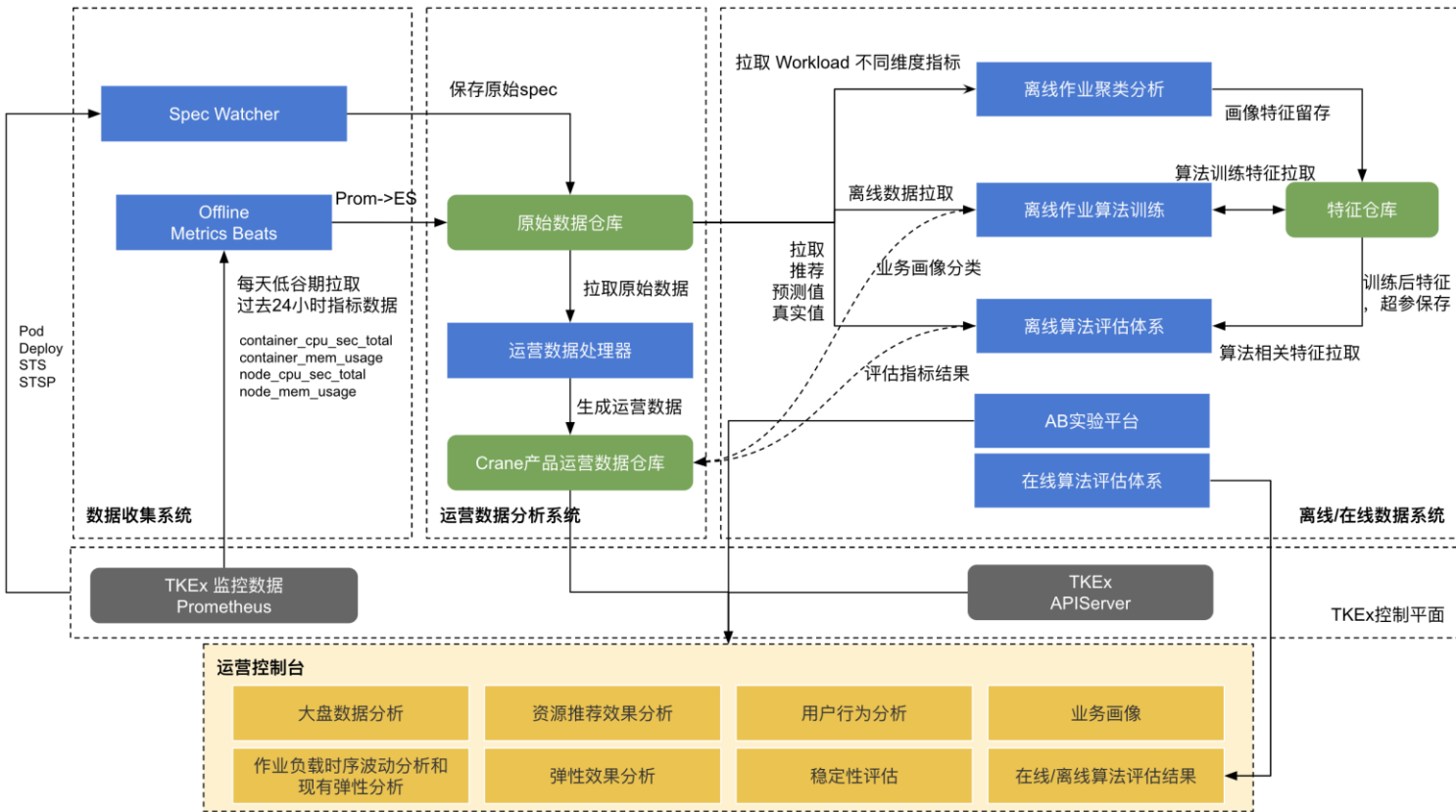
- 自定义Spec Watcher捕获workload变动
- 基于Prometheus Metrics Beats每日凌晨拉取当天业务指标
- 针对Metrics Beats做了自定义存储优化，存储空间降低数个量级

• 离线算法评估

- 针对大量离线指标数据评估预测算法准确性
- 超参调优

• 运营数据分析

- 聚类与业务画像
- 大盘现状与走势
- 优化进展
- 用户行为分析



产品化控制台助力业务优化

成本可视化

- 用量
- 基于预测的趋势分析

成本和浪费识别

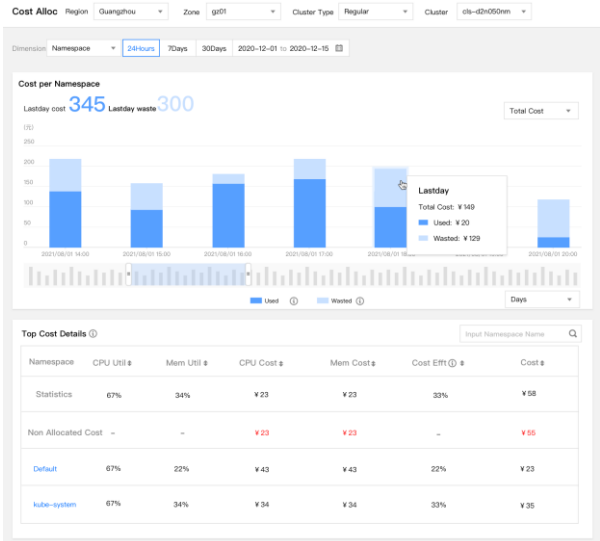
- 与计费API整合的费用展示

灵活的汇聚维度

- 按部门
- 按项目
- 按应用类型
- 按自定义标签

方便可信的优化能力

- 可支持原地升降配的规格优化
- 弹性推荐、定时弹性和预测弹性
- 三条黄金曲线展示推荐值的来源



Pod 更新策略

配置实例镜像

根据对本工作负载的历史监控, 建议使用推荐值更改当前工作负载的 CPU 与内存配置。

请查看当前工作负载资源推荐量, 当前配置量, 与实际用量的关系图。

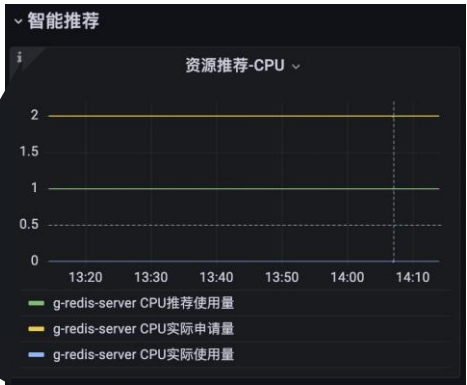
实例容器 redis 配置

容器名: redis

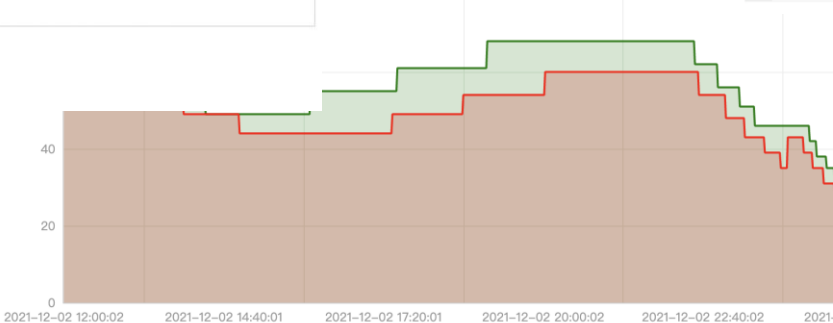
镜像地址: ccr.ccs.tencentyun.com/tkeimages/redis 选择镜像

镜像版本: 5.0.8

CPU 与内存限制: 2 核 2Gi 推荐值: 1 核 / 1Gi

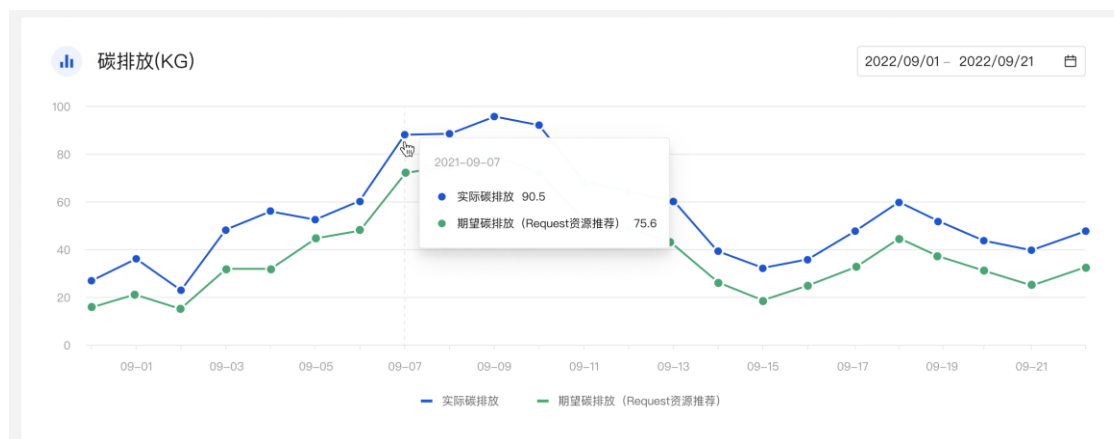


	metrics_name	score
0	accuracy_score	0.979167
1	precision_score	0.979651
2	recall_score	0.979167
3	f1-score	0.978342
4	Missing detection rate	0.000000
5	roc-auc-score	0.916667





$$\text{Average Watts} = \text{Min Watts} + \text{Avg vCPU Utilization} * (\text{Max Watts} - \text{Min Watts})$$



中国生态环境部2022年发布的《电力行业温室气体核算指南修订版》中明确提出，电网排放因子采用0.5810tCO₂/MWh

$$\text{Average Carbon Emissions} = \text{Average Watts} * \text{Emission Factor} = \text{Average Watts} * 0.5810$$

平台侧优化 - 节点容量缩放和水位管理

集群大盘可视化

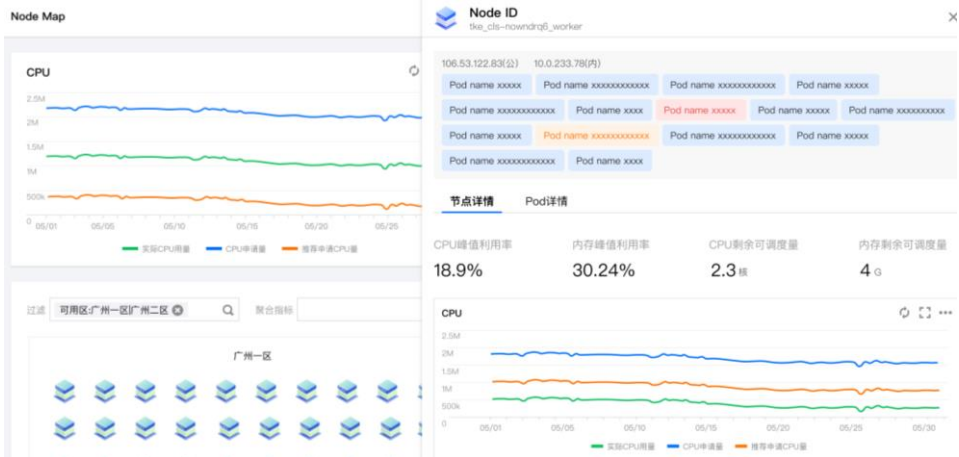
- 集群总体利用率
- 节点利用率热力图

节点容量缩放

- 可定义节点容量缩放比例，放大节点可分配资源，提升装箱率

节点水位控制

- 可自定义节点水位，控制节点利用率上限
- 动态调度器依据利用率装箱，确保真实利用率与目标利用率一致
- 可配置紧缩优先调度策略，方便退还闲置节点



托管节点专用调度器

托管节点专用调度器支持通过配置参数，在保障托管节点稳定性的同时，提升托管节点资源的利用率。

作用的节点范围

节点类型	作用节点	CPU 目标利用率	内存目标利用率	CPU 规格放大系数	内存规格放大系数
托管节点/节点池	托管节点池下拉列表	推荐30~60 %	推荐40~70 %	推荐1~3	推荐1~2

作用的 Workload 范围 请选择NS

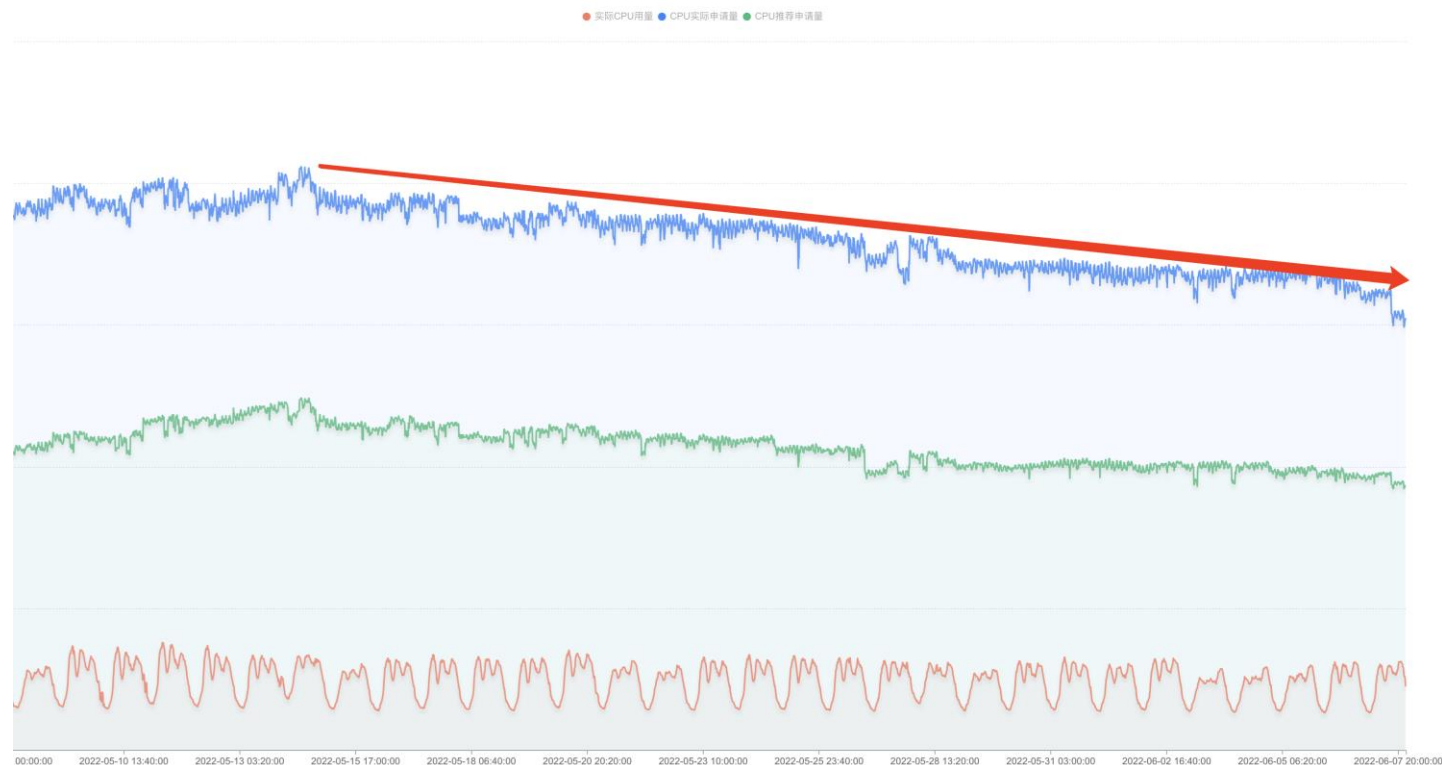
调度策略 紧缩优先

紧缩优先，会让 Pod 先调度到利用率相对较高、但没有超过上述设置的目标利用率的节点上，提升节点利用率，详情请[查看](#)。

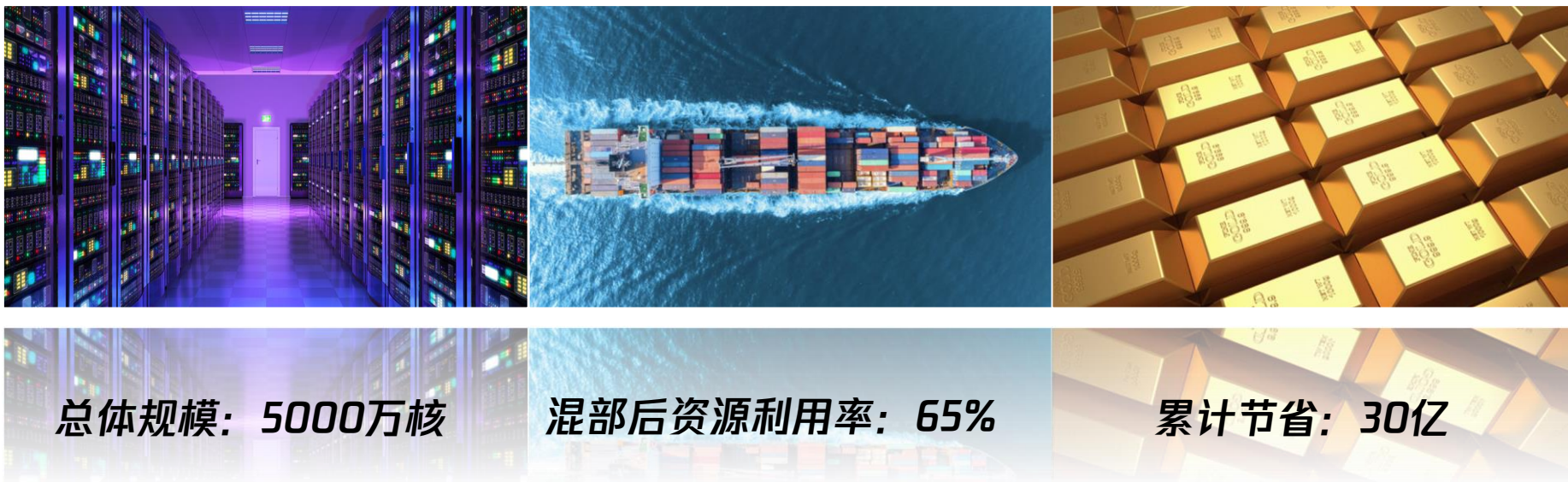
保存 取消

内部大规模落地的成效

- 在腾讯内部自研业务大规模落地
- 部署至数百个Kubernetes集群
- 管控数百万CPU核
- 全面上线一个月内，大盘总核数缩减25%



腾讯内部海量自研业务云上成本优化成效

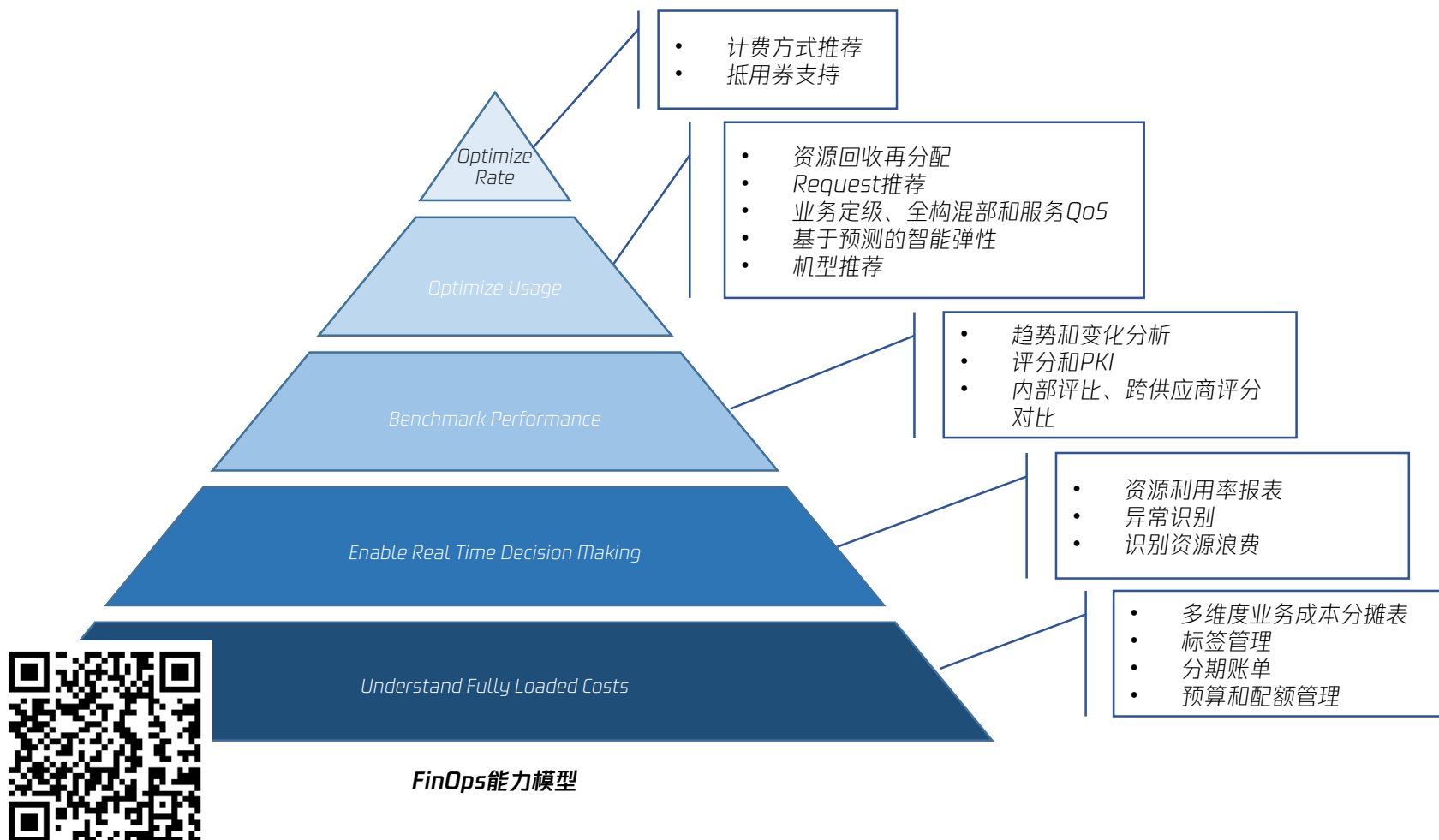


Crane

Cloud Resource Analytics and Economics

期待与您共建，推动科技进步，惠及更多受众!

<https://github.com/gocrane/crane>



开展 FinOps 布道

《降本之源-云原生成本管理白皮书》

《云成本优化节能减排白皮书》

FinOps “双降”讲坛

FinOps Community 公众号

FinOps 开源项目 Crane

制定 FinOps 标准

FinOps 国内首家顶级会员

信通院云管优秀案例

牵头《云原生 FinOps 能力成熟度模型》

参与《云成本优化工具技术要求》

参与《企业云成本优化能力与效果成熟度模型》

参与《云财务运营成熟度模型》

成立 FinOps 产业联盟

腾讯云牵头、40 家企业共同发起成立“FinOps 产业标准工作组”形成国内首个 FinOps 产业标准生态联盟，整合产、学、研、用各方力量，推动 FinOps 落地实践。

荣誉

国家级科学卓越奖《云计算中心科技卓越奖》

信通院云原生产业联盟《2022年度云原生技术创新领航者》

云计算标准和开源推进委员会《2022年度云优化优秀案例》

推荐插件与控制台

- 统一推荐插件开发
- 基于T的产品化控制台
<http://dashboard.gocrane.io/>

算法增强

- IEG蓝鲸产品中心顶会论文加持:

A Transferable Time Series Forecasting Service Using Deep Transformer Model for Online Systems

A Semi-Supervised VAE Based Active Anomaly Detection Framework in Multivariate Time Series for Online Systems

装箱优化

- IEG Eunomia项目组:
基于遗传算法的装箱优化
- TEG星辰算力团队:
更多的重调度指标和策略支持

多集群的资源优化

- 多集群HPA副本数调节

期待专家共建

- 前端工程师
- 算法工程师
- 调度、内核专家
- 多集群管理专家

欢迎试用落地

- 业务优化: 规格优化、副本数优化、智能弹性
- 平台优化: 装箱率提升、利用率提升、资源腾挪



THANKS

Architect