



GOPS2018
Shenzhen

GOPS

全球运维大会 2018

2018.4.13-4.14

中国·广东·深圳·南山区 圣淘沙大酒店(翡翠店)





GOPS2018
Shenzhen

阿里巴巴规模化混部技术演进

蒋玲（玲昕） 阿里巴巴-系统软件部-技术专家

个人简介

- 蒋玲，阿里花名玲昕；
- 2012年加入淘宝网，曾负责电商、菜鸟、新零售、阿里影业等业务运维；
- 大促自动化备战产品负责人；
- 电商规模化混部项目负责人；
- 5次参与双11备战，曾荣获“2015年双11技术保障老A（特种兵）称号”。



GOPS2018
Shenzhen



GOPS2018
Shenzhen

目录

- ➔ **1** 阿里巴巴混部探索简介
- 2** 混部方案及架构
- 3** 混部核心技术
- 4** 未来展望



GOPS2018
Shenzhen

阿里巴巴混部探索简介

1. 为什么混部？业务增长 VS 资源成本
2. 何为混部？
3. 混部在阿里的发展
4. 阿里巴巴规模化混部成果

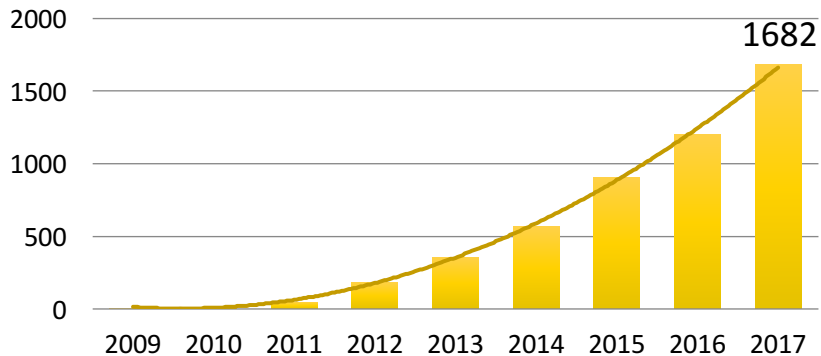


GOPS2018
Shenzhen

为什么混部？

- 业务增长：双11大促日成交额**1682亿**、大数据存储**KPB级**、日均**百万级**任务；
- 不同类型业务部署于**独立的数据中心**，资源体量巨大，**万台**机器；
- 容灾等设计使得在线业务数据中心资源利用率不高，**10%左右**；
- 不同业务特性，表现出**分时复用**的可能性。

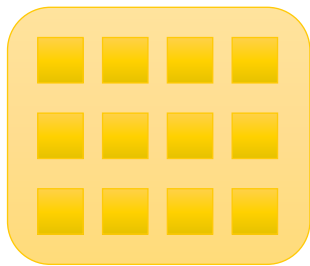
双11交易额（亿）



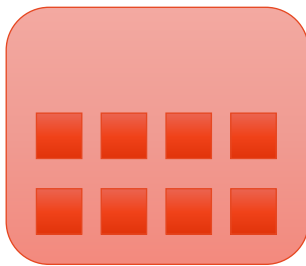


何为混部（Co-location）？

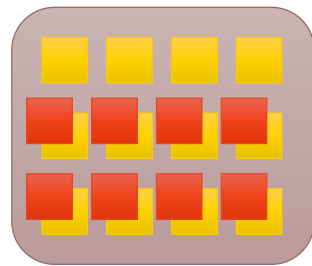
- 不同类型业务混合部署（资源整合）
- 通过调度与隔离的手段进行资源共享与竞争
- 保障不同业务服务质量（服务优先级）



实时型业务



计算型业务



混合部署



GOPS2018
Shenzhen

在线离线混部

在线业务-高优先级



交易、支付、浏览型请求

实时响应，不可降级

延时敏感，不可重试

日常：白天高、夜晚低

大促：脉冲压力高，持续时间短

离线业务-低优先级



计算型任务

离线任务，可短时间降级

延时不敏感，可重试

日常：白天低，夜晚高

大促：部分降级

在线

- 容器化分配资源
- 高优先级
- 在线日常空闲

离线

- 按进程申请资源
- 低优先级
- 离线日常繁忙

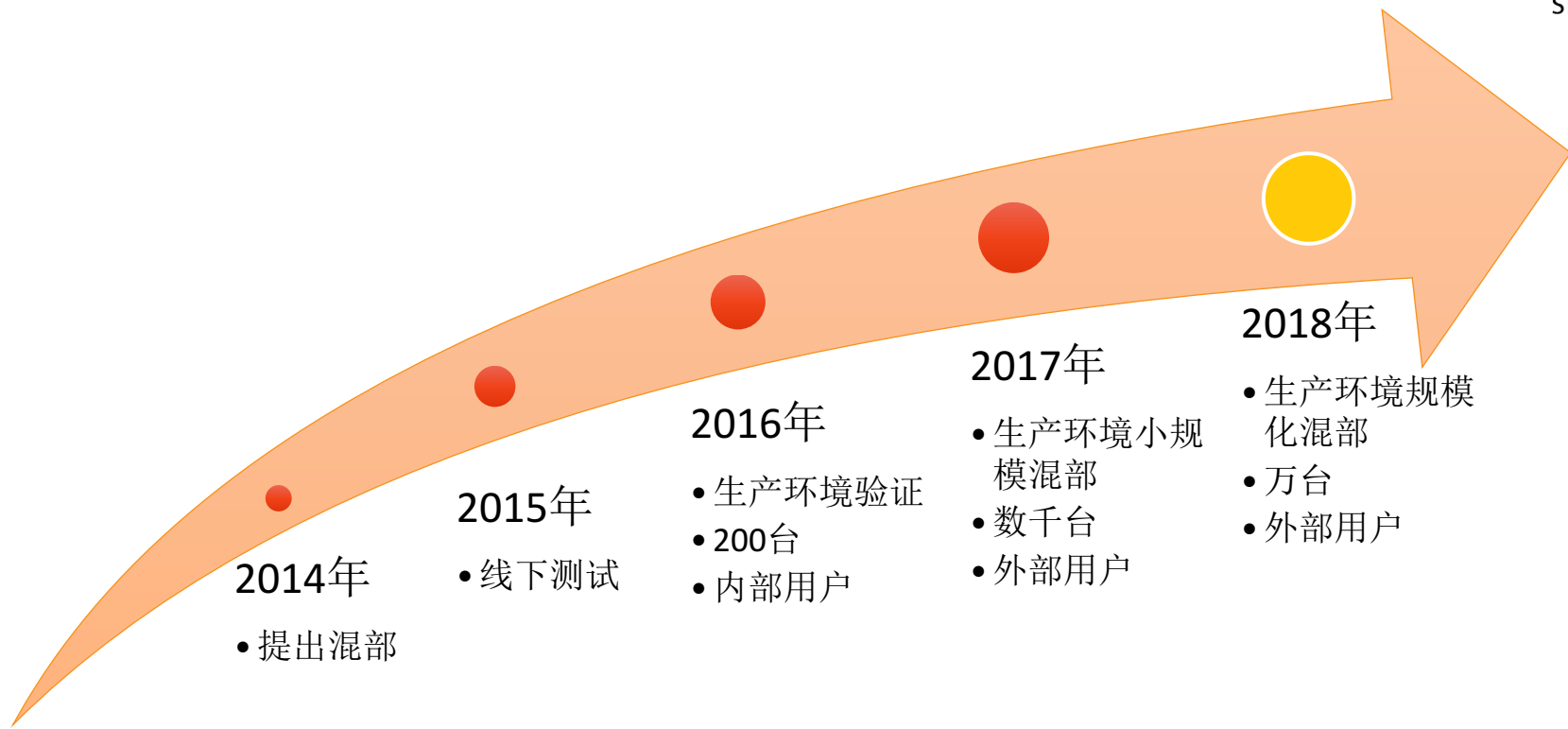
共享与
抢占

- 在、离线各自分配整机资源
- 闲时离线填充
- 竞争时离线退出

阿里巴巴混部探索历程



GOPS2018
Shenzhen





GOPS2018
Shenzhen

阿里巴巴规模化混部成果

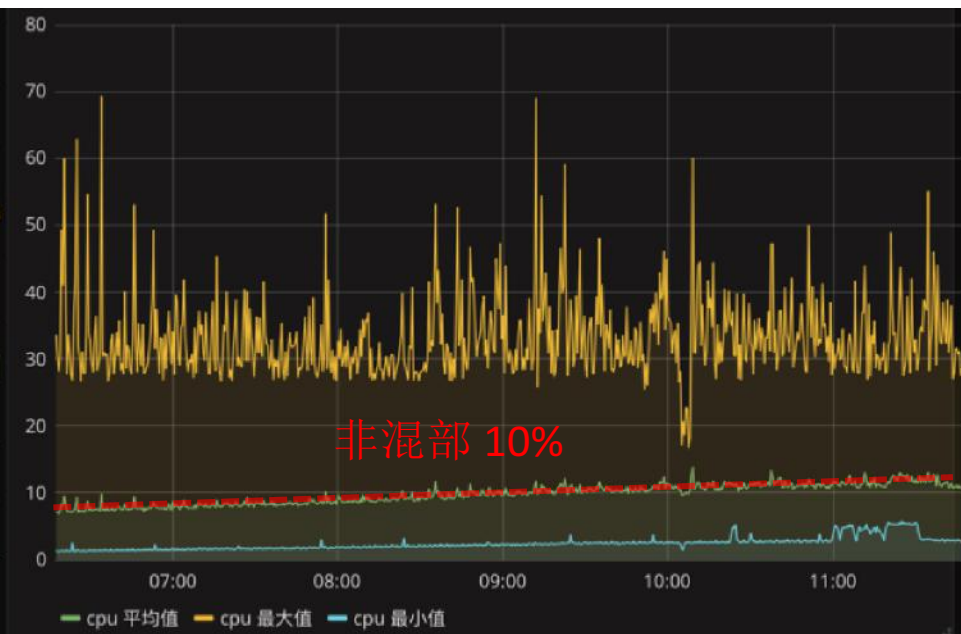
1. 混部规模达**数千台**，经历双11交易核心场景验证；
2. 在线集群上引入离线计算任务（在离线）：**日常CPU利用率 10% -> 40%**；
3. 离线集群上部署在线业务（离在线），支撑**双11大促数W笔/s**交易创建能力；
4. 混部环境下对在线业务服务干扰影响小于**5%**；



GOPS2018
Shenzhen

阿里巴巴规模化混部成果

在线集群上引入离线计算任务：日常CPU利用率 10% -> 40%

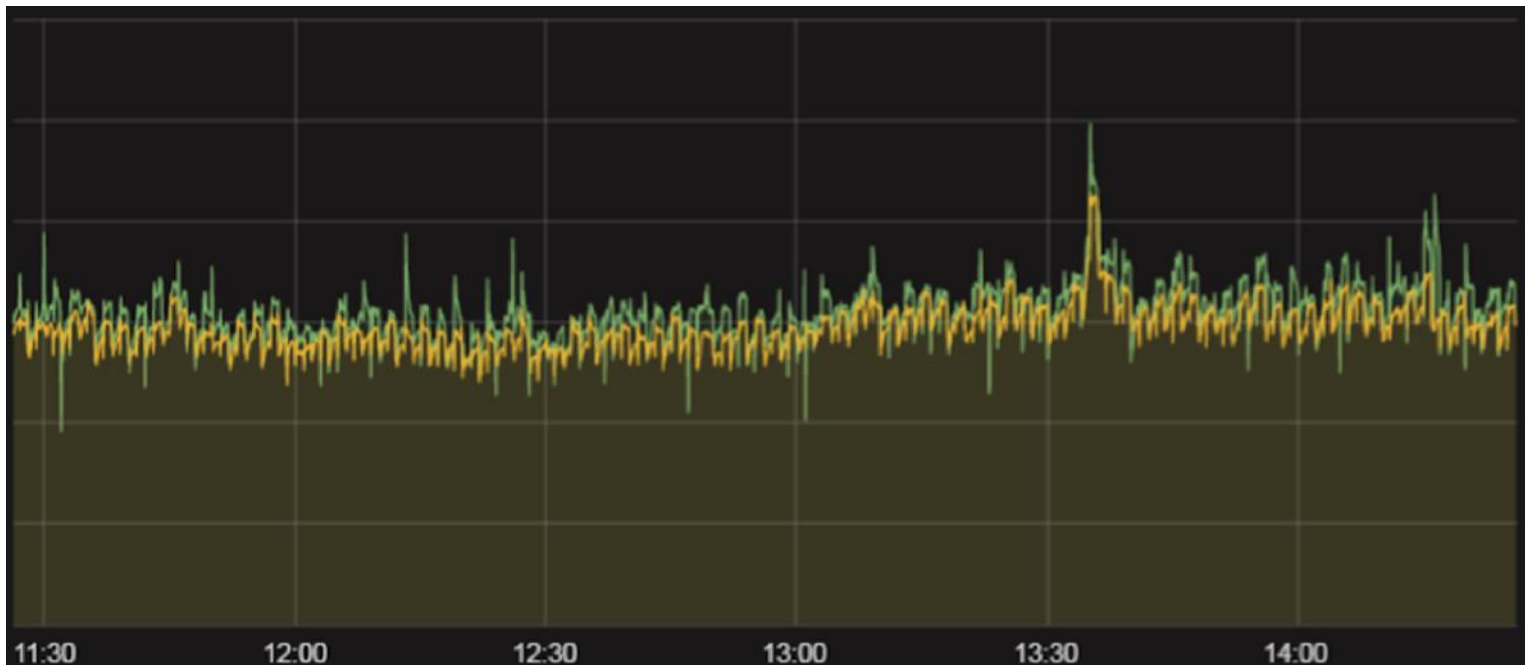




GOPS2018
Shenzhen

阿里巴巴规模化混部成果

混部环境在线交易服务RT表现：**混部干扰 < 5%**





GOPS2018
Shenzhen

目录

1 阿里巴巴混部探索简介

➔ 2 混部方案及架构

3 混部核心技术

4 未来展望



GOPS2018
Shenzhen

混部方案及架构

1. 混部整体架构
2. 混部场景业务部署策略
 - 计算存储分离技术
 - 无中生有：资源共享
3. 混部场景业务运行策略
 - 大促：站点快上快下
 - 日常：分时复用



混部整体架构

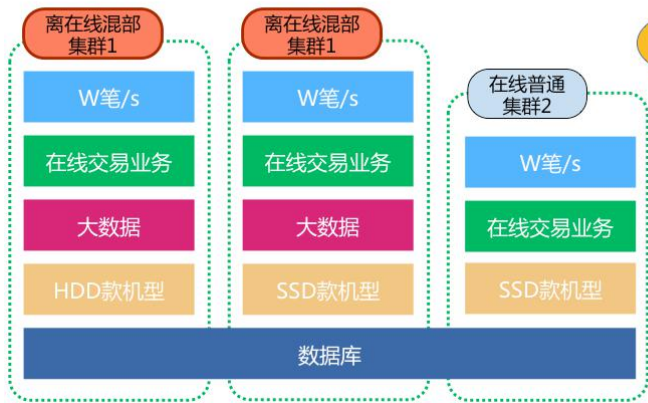
- 合并资源，分时复用；
- 业务资源调度：sigma、Fuxi；
- 0层：协调一层调度及资源分配；
- 内核级资源隔离；
- 支持业务服务优先级，保障SLA。





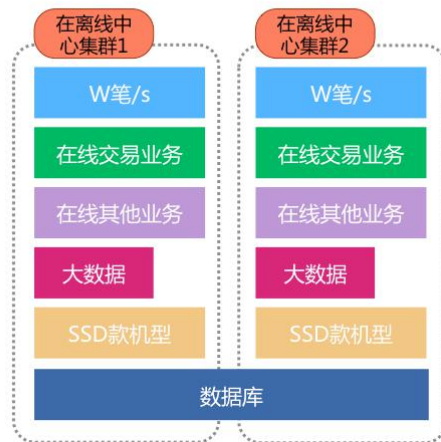
混部场景在线业务部署策略

在线业务：三地多单元；中心同城双机房



离在线集群

在离线集群

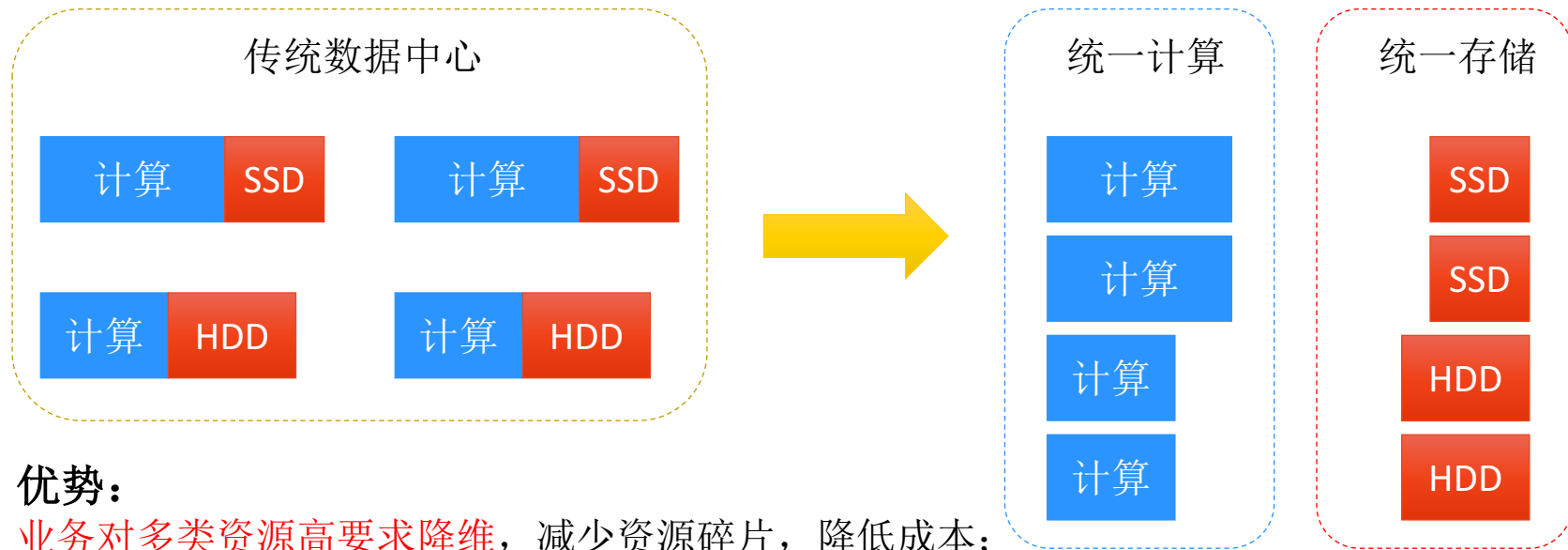


- 资源整合：网络环境、机型约束
- 单元化架构：
 - 单元内交易闭环；
 - 单元流量根据userid进行分流；



GOPS2018
Shenzhen

计算存储分离技术



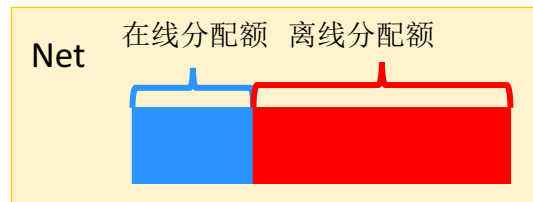
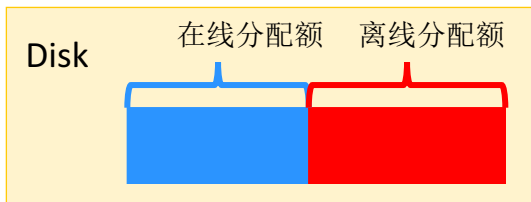
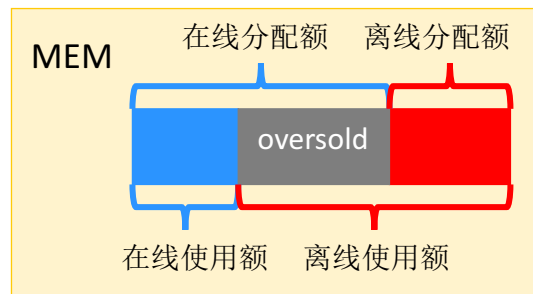
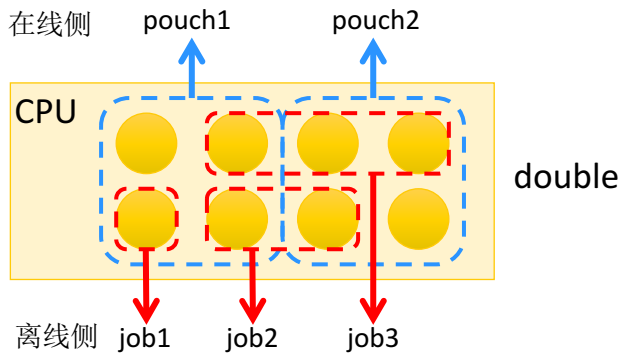
优势:

业务对多类资源高要求降维，减少资源碎片，降低成本；
调度复杂度降级；
充分利用高带宽网络红利；



混部集群资源分配

无中生有；充分共享；竞争隔离。



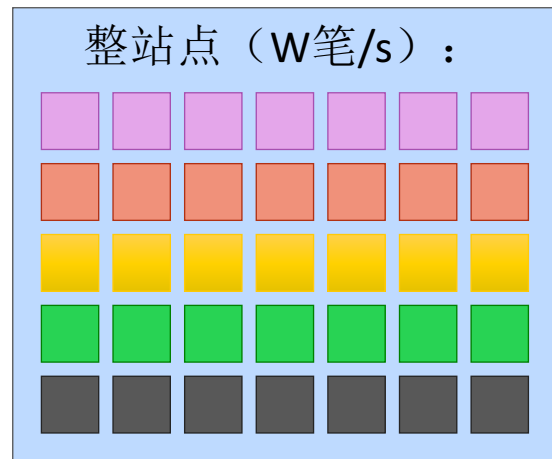
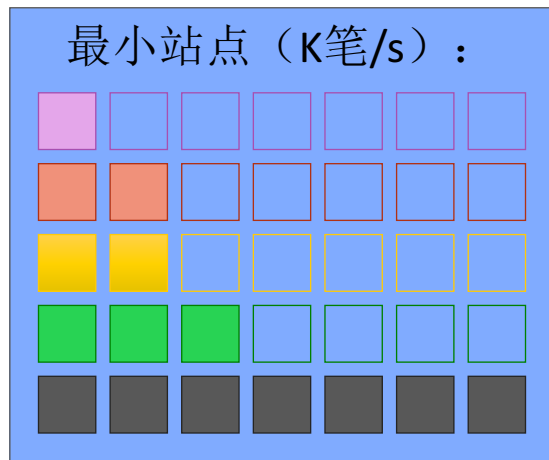
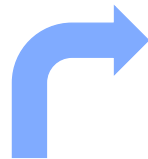


GOPS2018
Shenzhen

大促资源退让机制：站点快上快下

- 在线站点维度容量管控
 - 精细化容量模型建模；
 - 独占型及可伸缩型应用；
- 快速站点容量伸缩能力；
- 在线、离线资源调配机制；
- 装箱调度能力；
- 整站点运行时，离线有部分资源损失，业务降级；

站点快上：
1小时内快速拉起整站点容量；



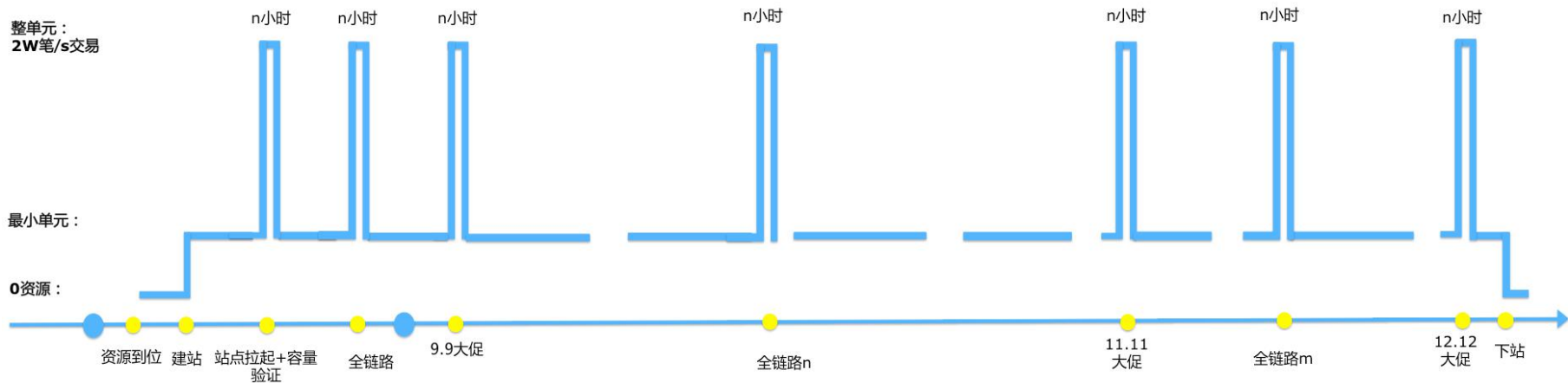
站点快下：
半小时内快速缩容到最小站点容量，释放资源给离线；



GOPS2018
Shenzhen

大促站点快上快下：站点运行计划

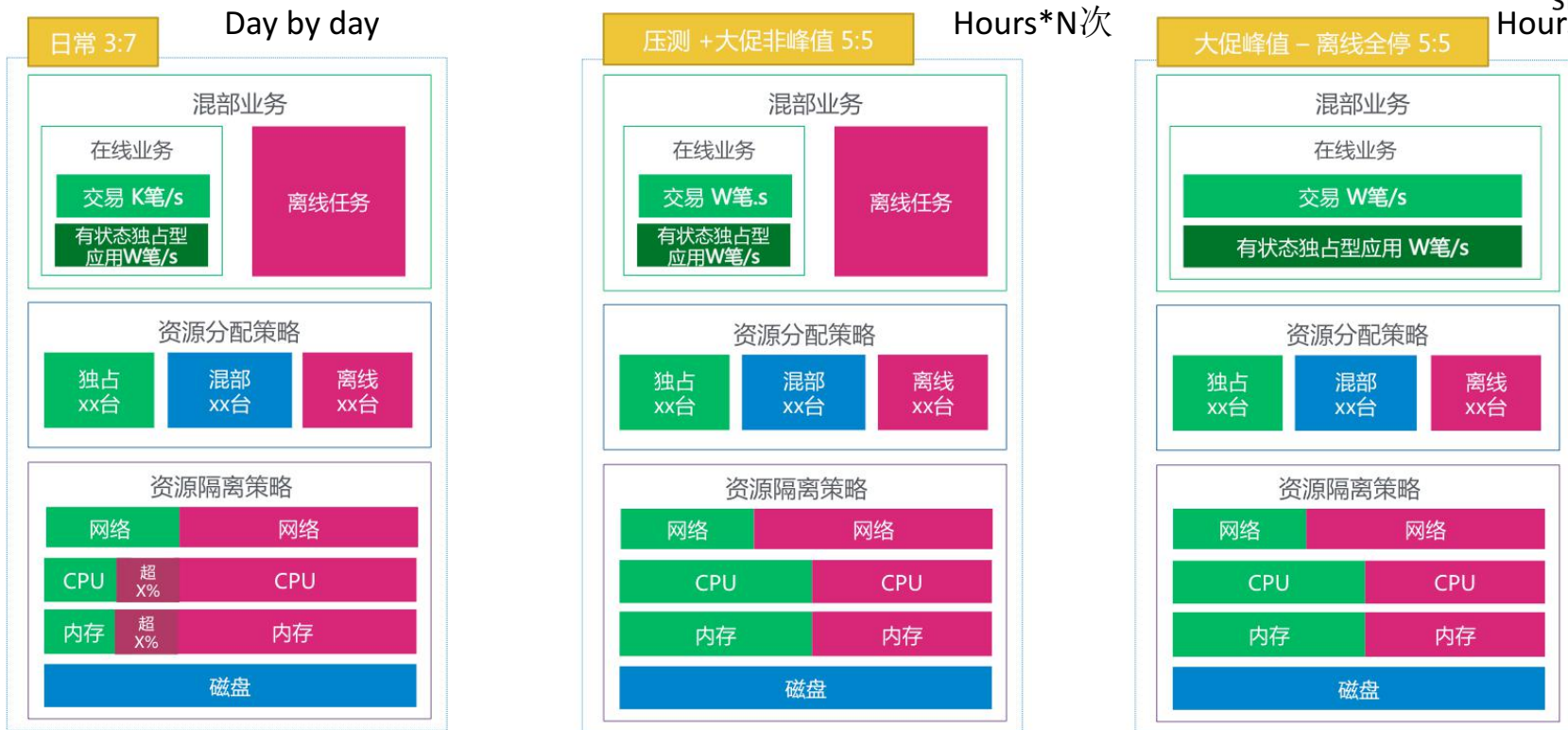
- 离在线集群支撑99大促、双11大促、双12大促；
- 快上1小时、快下半小时，万笔交易能力伸缩，涉及近千个在线容器；
- 离线业务平滑顺延，业务无损，用户无感知（半小时全停）；
- 充分共享物理资源。





GOPS2018
Shenzhen

快上快下：不同场景下的资源分配





GOPS2018
Shenzhen

日常资源退让机制：分时复用

在线日常流量曲线：



- 时间维度优化资源分配；
- 弹性伸缩，分时复用；
- 平均CPU利用率提升至60%+；

在线扩容
离线扩容

在线扩容
离线扩容



GOPS2018
Shenzhen

目录

1 阿里巴巴混部探索简介

2 混部方案及架构

→ 3 混部核心技术

4 未来展望



GOPS2018
Shenzhen

混部核心技术

1. 内核资源隔离技术

- CPU HT隔离：Noise Clean，解决超线程资源争抢问题，一堆HT核不会同时跑离线、在线任务
- CPU调度隔离：CFS基础上增加Task Preempt特性，提高在线服务调度优先级
- CPU缓存隔离：CAT，三级缓存（LLC）通道隔离（Broadwell及以上）
- 内存隔离：Cgroup内存用量隔离、Bandwidth Control内存带宽隔离、OOM优先级；
- **内存弹性**：在线闲时，离线调度突破mem Cgroup limit
- IO隔离：IO带宽隔离
- 网络QoS隔离：单机TC增强管控；金银铜牌业务等级定义（在线银、离线铜），全网络分等级带宽保障；

在线业务高优先级保障



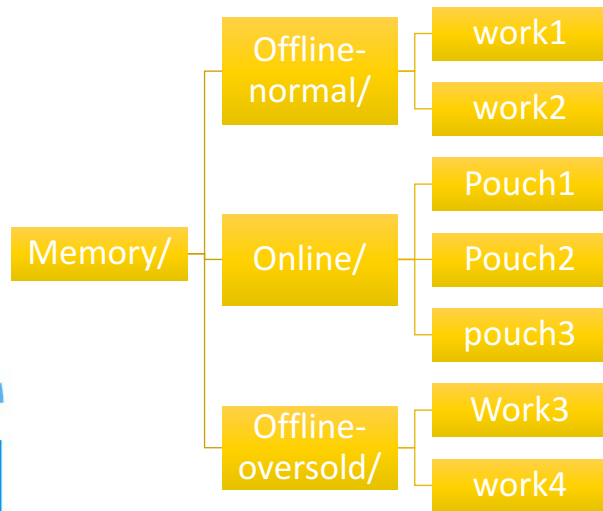
混部核心技术—Memory动态超卖

□ 内存容量隔离与抢占—动态内存分配

Cgroup分组控制、在线内存初始配额；

新增超卖cgroup，将实际未消耗的物理内存分配给离线超卖任务；

在、离线实际使用内存间保留buffer值，用于满足在线内存使用增量；



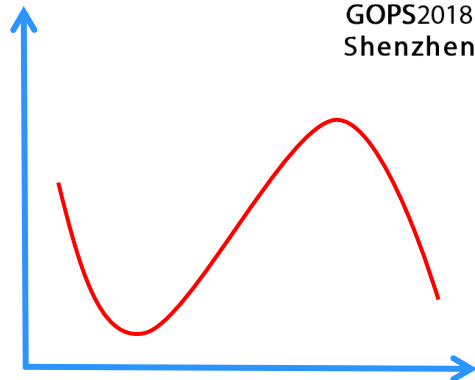
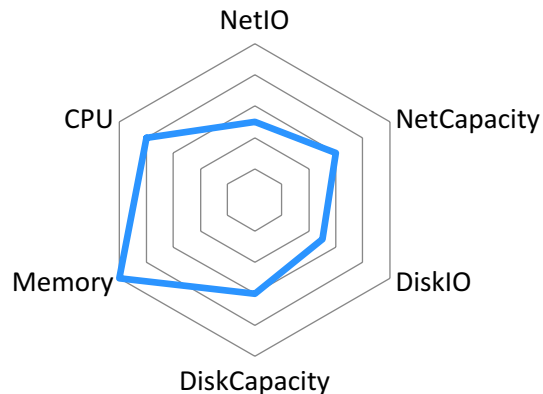


GOPS2018
Shenzhen

混部核心技术

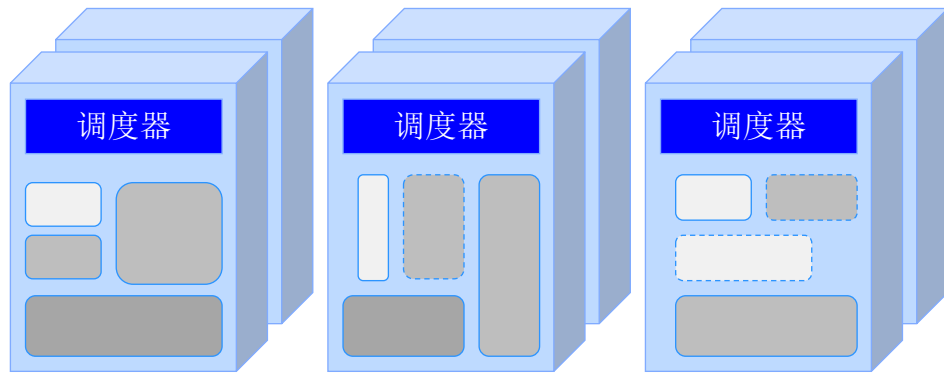
2. 在线集群调度

- 应用画像
- 资源调度：装箱、亲和互斥
- 应用自动伸缩、分时复用
- 站点快上快下



3. 离线集群调度

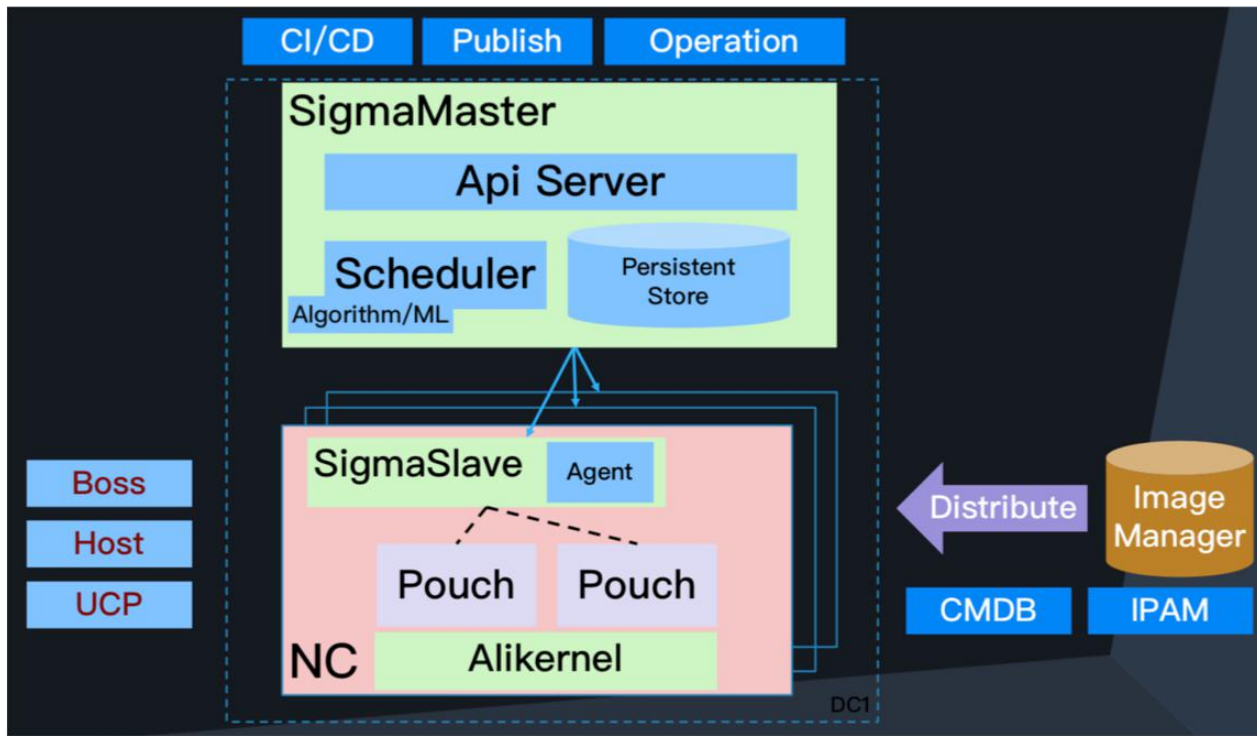
- 分等级任务调度
- 动态内存超卖
- 无损降级、有损降级





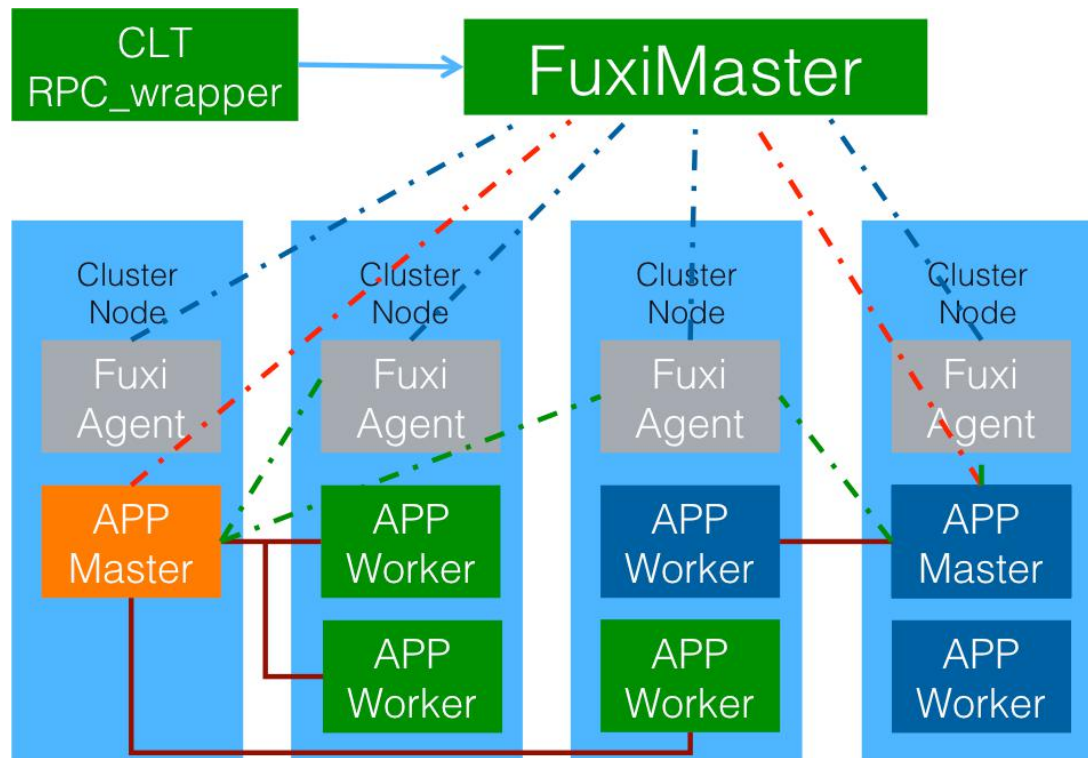
混部核心技术—在线资源调度sigma

- 兼容Kubernetes API, 和开源社区共建
- 采用阿里Pouch容器 (兼容OCI标准)
- 通过阿里多年大规模及双11验证





混部核心技术—离线资源调度

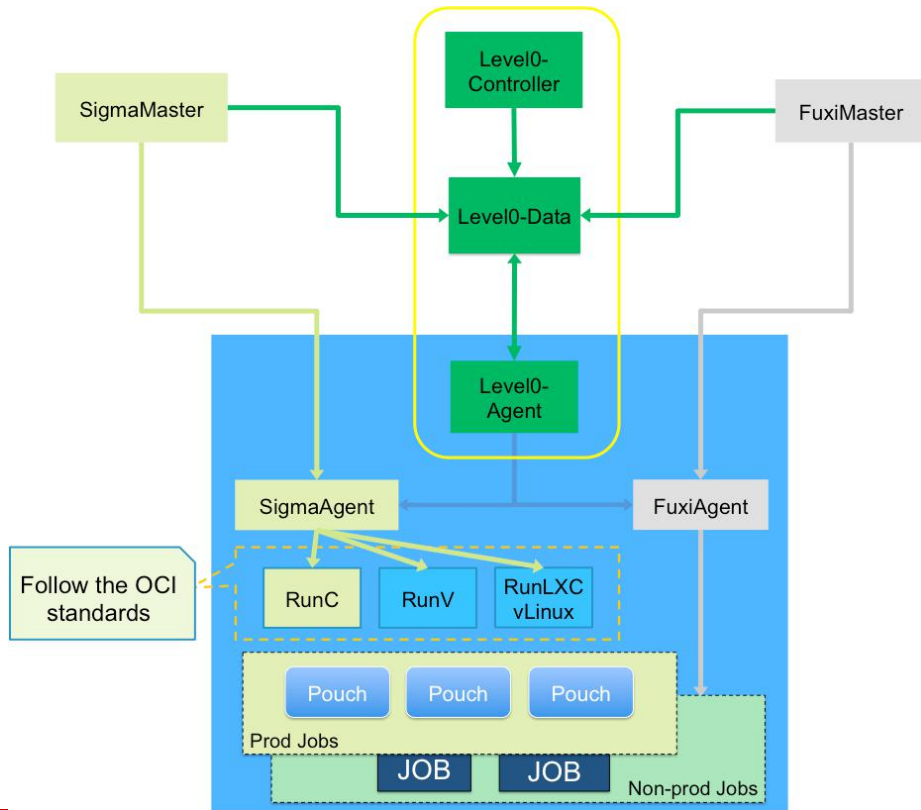


- 面向海量数据处理和大规模计算类型的复杂应用
- 提供了一个数据驱动的多级流水线并行计算框架，在表述能力上兼容MapReduce, Map-Reduce-Merge, Cascading, FlumeJava 等多种编程模式。
- 高可扩展性，支持十万以上级的并行任务调度，能根据数据分布优化网络开销。自动检测故障和系统热点，重试失败任务，保证作业稳定可靠运行完成。



混部核心技术—统一资源调度0层

- 通过sigma和fuxi完成在线离线的各自调度
- 通过零层相互协调资源配比





GOPS2018
Shenzhen

目录

1 阿里巴巴混部探索简介

2 混部方案及架构

3 混部核心技术

➔ 4 未来展望



GOPS2018
Shenzhen

未来展望

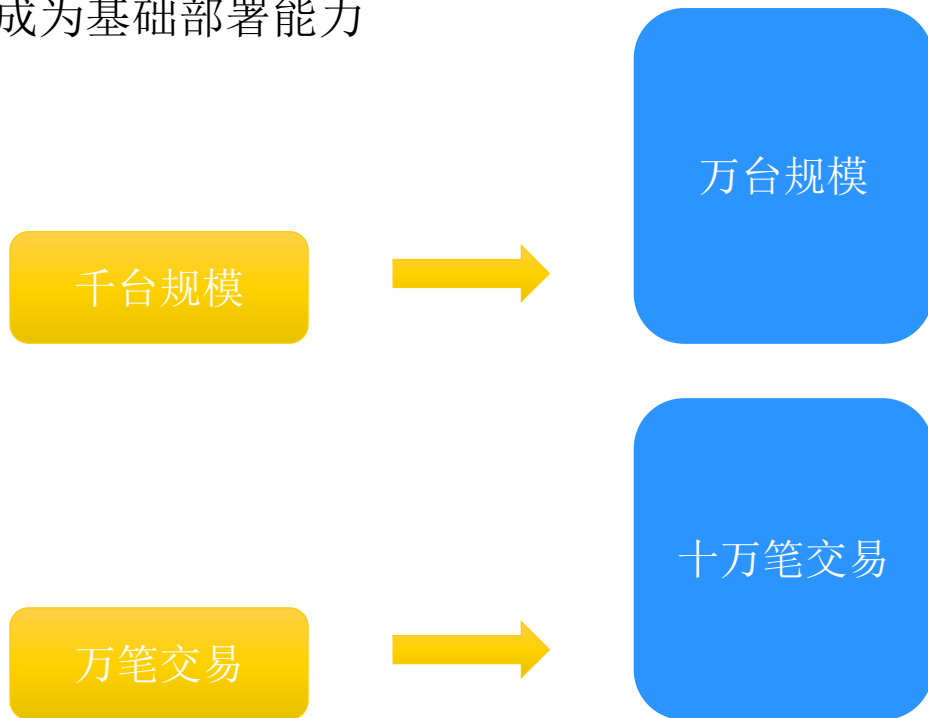
1. 规模化
2. 多元化
3. 精细化



GOPS2018
Shenzhen

规模化：混部成为基础技术能力

全面混部，将成为基础部署能力

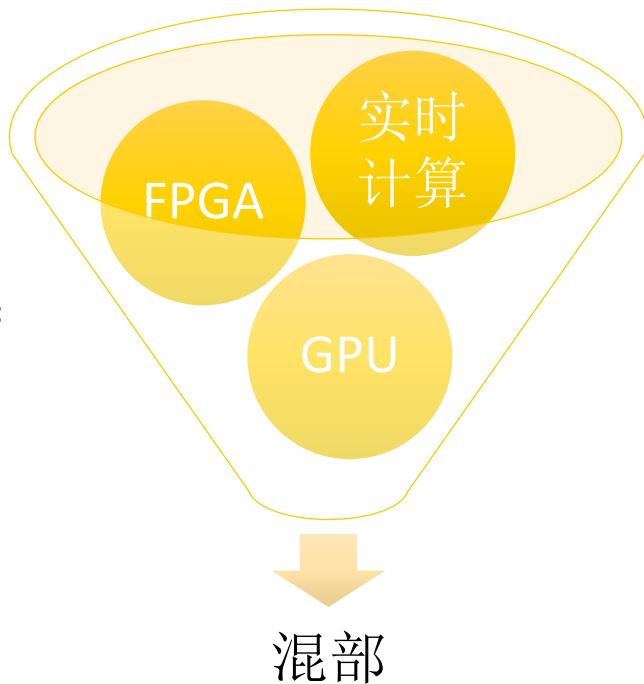




GOPS2018
Shenzhen

多元化：不同类型的业务及资源混部

- 不同业务类型；
- 不同硬件设备；
- 不同基础设施（云上、云下）；



精细化



GOPS2018
Shenzhen

精细化业务资源画像；

精细化容量管理；

精细化资源调度；

精细化内核隔离；

精细化监控及运维配套；

THE END



GOPS2018
Shenzhen

Q & A

微信号: smallfishxx





GOPS2018
Shenzhen



Thanks

高效运维社区
开放运维联盟

荣誉出品



GOPS2018
Shenzhen

想第一时间看到高效运维社区
的新动态吗？

