

# 飞天 互联网规模的操作系统

Apsara: an Internet-scale Operating System

唐 洪 阿里云首席架构师



# 互联网规模的基础设施

## Internet-scale Infrastructure

大陆五个区域：华北（两个）、华南、华东（两个）

大陆以外十个区域：香港、新加坡、美西、美东、日本、印度、迪拜、德国、悉尼、台湾

边际网络：全球600多POP节点、带宽储备超过20T





## PC 操作系统

应用商店

输入 / 输出

系统调用

账号管理  
权限控制

OS 内核

PC 硬件

## 飞天

云市场

接入层：数据传输、内容发布、网络接入

云服务层：云服务的 Web API

云服务层：账号、认证授权、计量、结算

内核层：数据中心级别的集群计算系统

物理层：互联网规模的基础设施

## 云市场

接入层：数据传输、内容分发、网络接入

云服务层：云服务的Web API

云服务层：账号、认证授权、计量、结算

内核层：数据中心级别的集群计算系统

物理层：互联网规模的基础设施

自动化运维  
天基

云市场：VM镜像、容器镜像、编排模版、API服务

数据传输

内容分发

网络接入

数据智能：商业智能、数据开发、人工智能

安全服务：密钥管理、云盾

连接编排服务：弹性伸缩、资源编排、通知队列、分布式事务管理

计算

存储

数据库

网络

账号

认证授权

计量

结算

盘古：分布式存储管理

伏羲：分布式资源调度

分布式协同

安全管理

日志采集

监控报警、跟踪诊断

遍布全球的几十个数据中心、数百个POP节点



# 数据湖 Data Lake

轻量

弹性

在线

# Future Data Lake

离线

重量

开源托管产品：Redcap MapReduce、HBase、GreenShark

单机：RDS 支持SQL、流计算、图计算

分布式：AnalyticDB、PetaData





2016 杭州·云栖大会  
THE COMPUTING CONFERENCE



# 低扰的高精度监控跟踪

High-resolution Monitoring and Tracing with Minimum Overhead

默认开启、7x24、秒级采样

全精度用户请求跟踪



# 指令级别的性能优化

Instruction-level Optimization

基于共享内存的数据采集：多进程并发无锁写入，zero-copy

时间戳获取：13ns vs 40ns

随机标识生成：4ns vs 16ns

支持每秒上百万事件的采集





## 幽霊審点升级、迁移变成常态

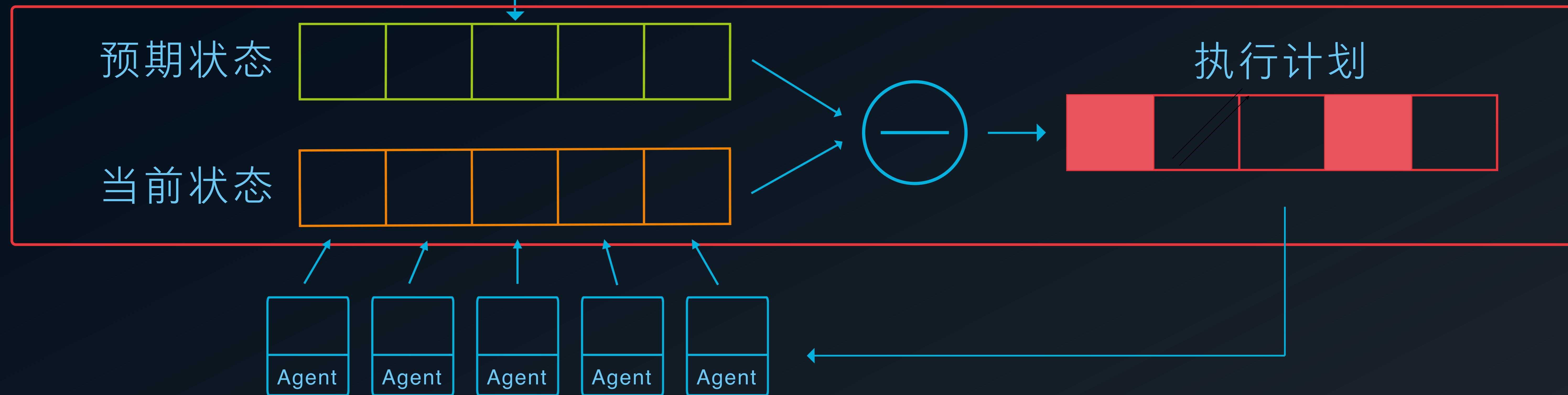
# 精瘦拨带系统运行状态

## 线上调试



天基Master

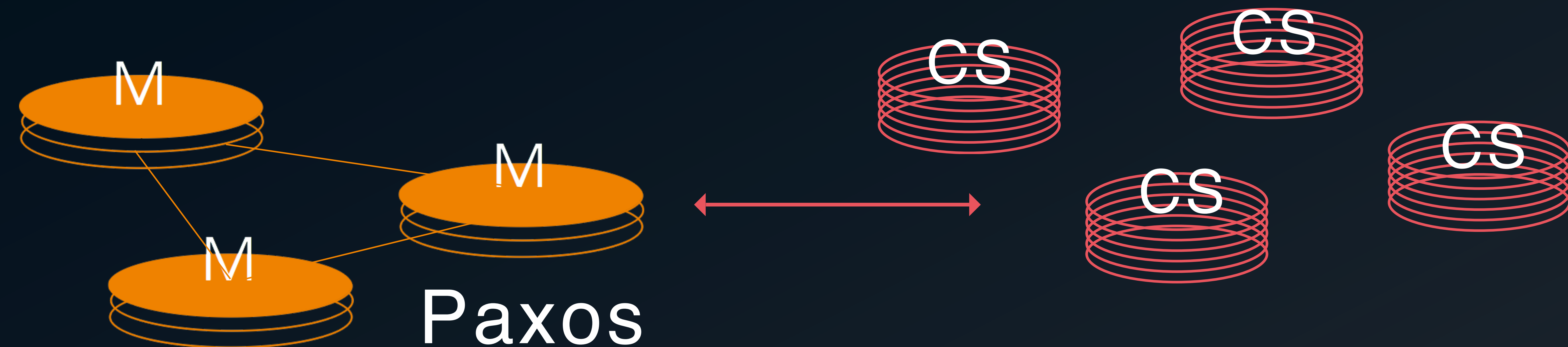
部署、升级、扩容、下线、配置变更



# 面向恢复的计算模型

Recovery-oriented Computing Model





# 盘古：分布式存储管理

Pangu: Distributed Storage Management

统一的存储管理

基于Paxos的高可用架构

单集群一万台服务器，十亿级文件数，EB级别存储空间



# 数据可靠是最高优先级

Data Reliability is the Highest Priority

默认三副本数据冗余，分布在不同机架，数据可靠性达到10个9  
支持跨数据中心的副本分布



# 数据可靠是最高优先级

Data Reliability is the Highest Priority

纠删码模式数据冗余：同样数据可靠性，存储开销从3x降低到1.375x



# 数据可靠是最高优先级

Data Reliability is the Highest Priority

端到端的数据校验：防止数据读写链路上的任何环节的问题导致数据发生错误

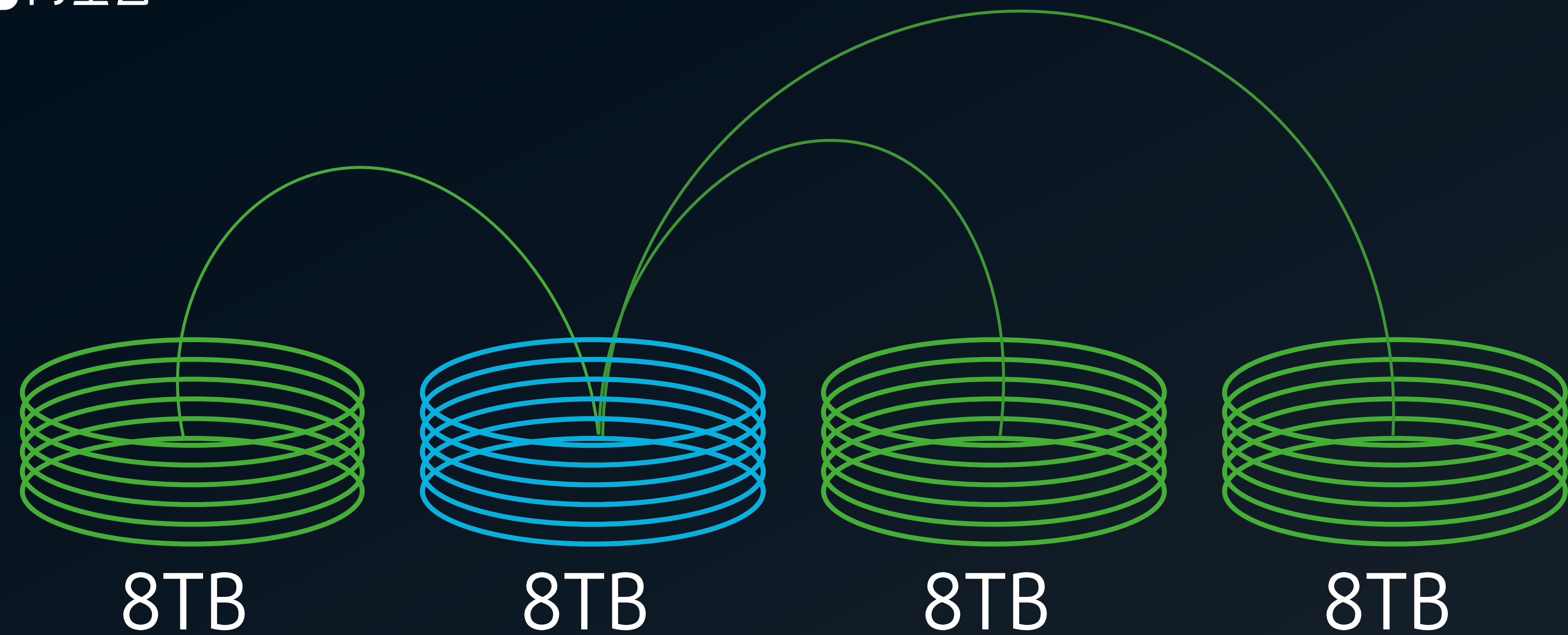


# 数据可靠是最高优先级

Data Reliability is the Highest Priority

并发冗余恢复：确保故障后数据即刻恢复冗余





# 硬件替换 + 原地恢复

Hardware Replacement with in-place Recovery

恢复速度受限于硬盘写入带宽

最快27小时恢复冗余

前台应用无法读写





# 盘古：并发冗余恢复

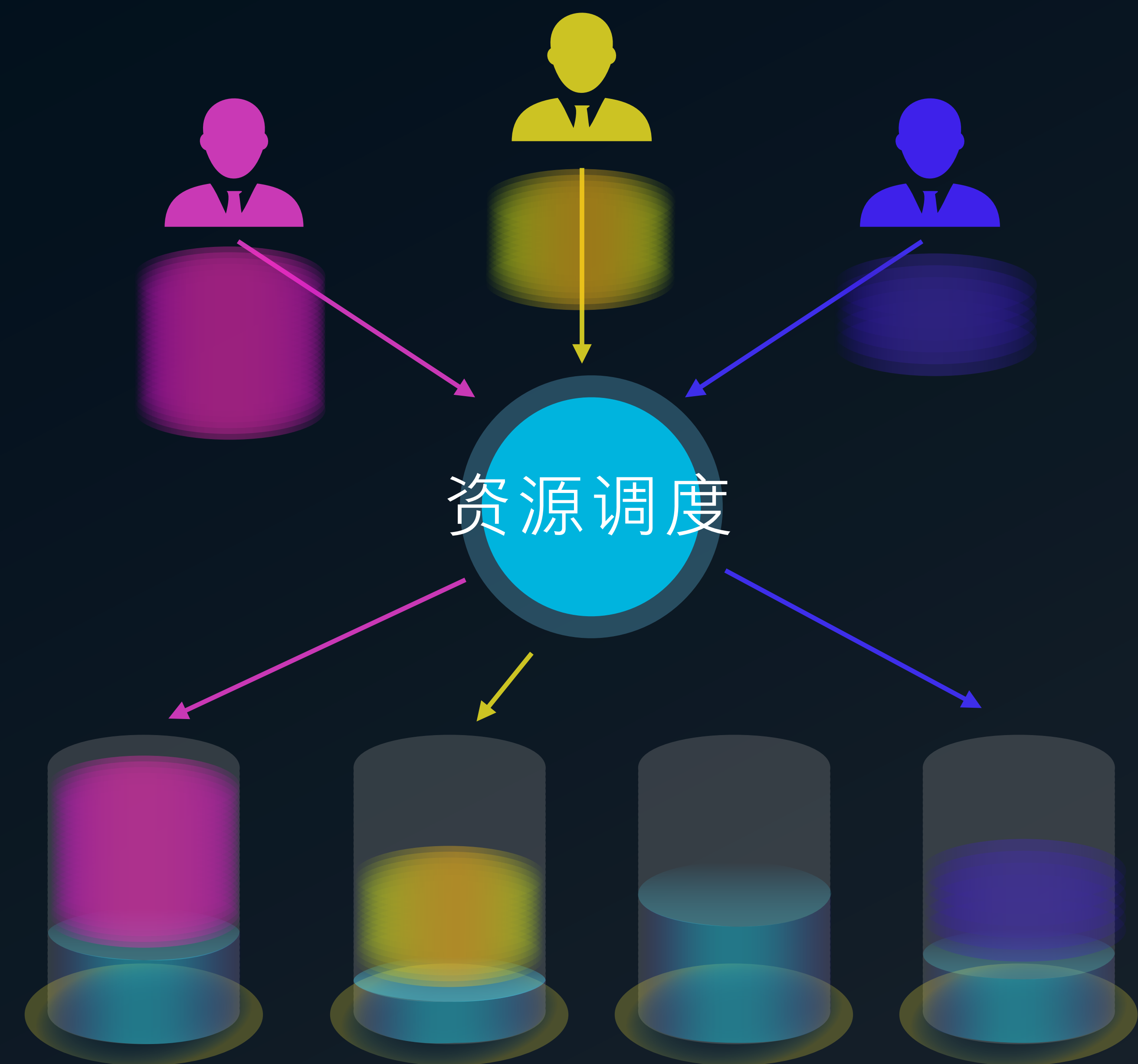
Pangu: Parallel Redundancy Recovery with Traffic Throttling

恢复速度与集群规模成反比关系

万台规模集群，不到一分钟恢复冗余

前端应用完全无感知





# 资源调度的挑战

Challenges of Resource Scheduling

万台集群规模

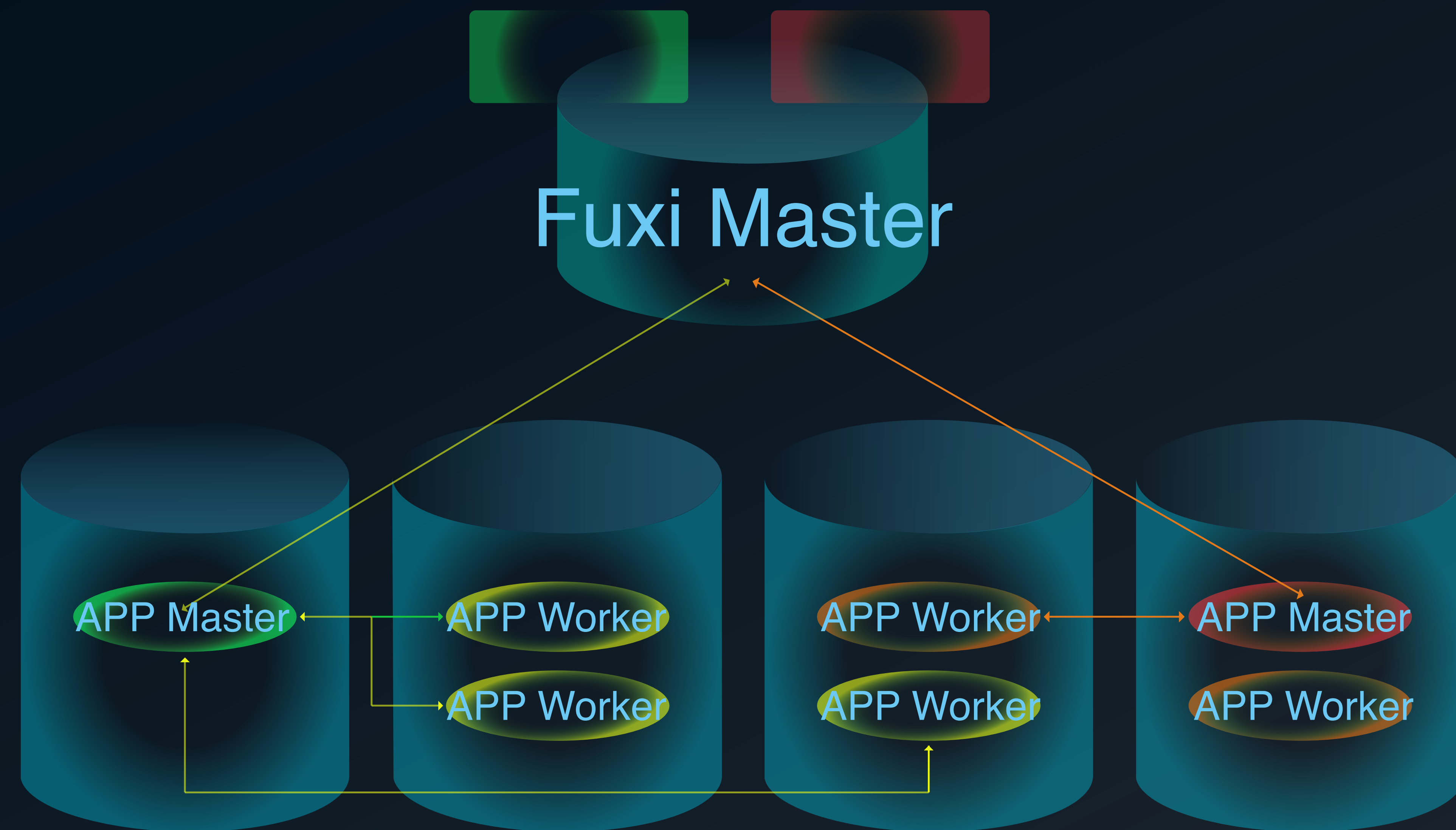
负载均衡

多维度资源请求

复杂的调度约束

额度控制

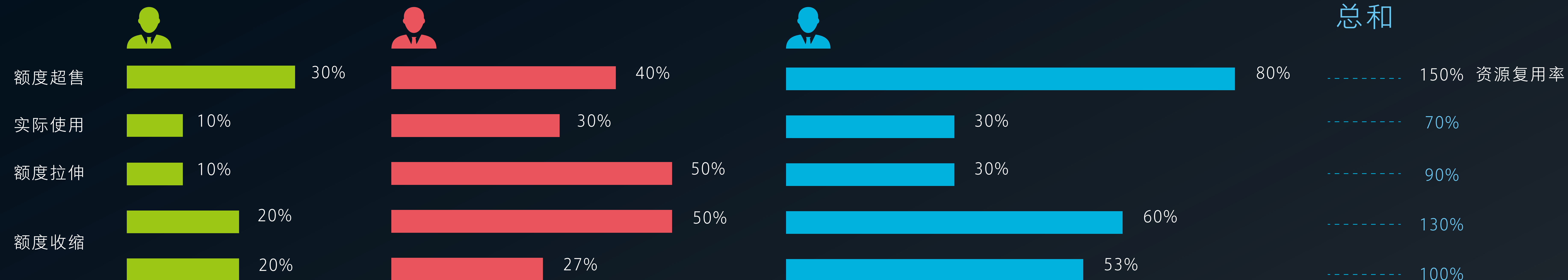




两级调度：批发 + 零售  
增量调度：一次请求、多次分配

单集群规模一万台，10万个进程，毫秒级响应  
2015年排序竞赛四项冠军，100TB排序377秒完成





# 兼顾效率与公平

Achieving Efficiency and Fairness

弹性额度

离线在线混合调度

日常利用率: 5% vs 54%

峰值利用率: 22% vs 64%



规模

Scale

高可靠

Reliability

性能

Performance

高可用

High Availability

效率

Efficiency

开放

Openness

自动升级部署API接口修复  
管理上云迁移基础设施  
云上数据迁移系统  
数据中心级灾备服务模式  
2015年排序竞赛四组数据排序377秒完成  
云市场兼容的灾备服务模式



# 为了无法计算的价值

Alibaba Cloud, More than just cloud