

QCon 全球软件开发大会 【北京站】2016

平台系统的演变

腾讯 孙子荀

自我介绍

个人信息：

Linux内核，在分布式系统，机器学习(并行计算) 上有一定经验。

主要经历：

- 华为/亚信
- 百度
- **2012.** 腾讯

•主导项目：QQ群广告/旋风下载技术2.0/ QQ公众号 后台负责人

1 平台系统

2 防护
接入平面

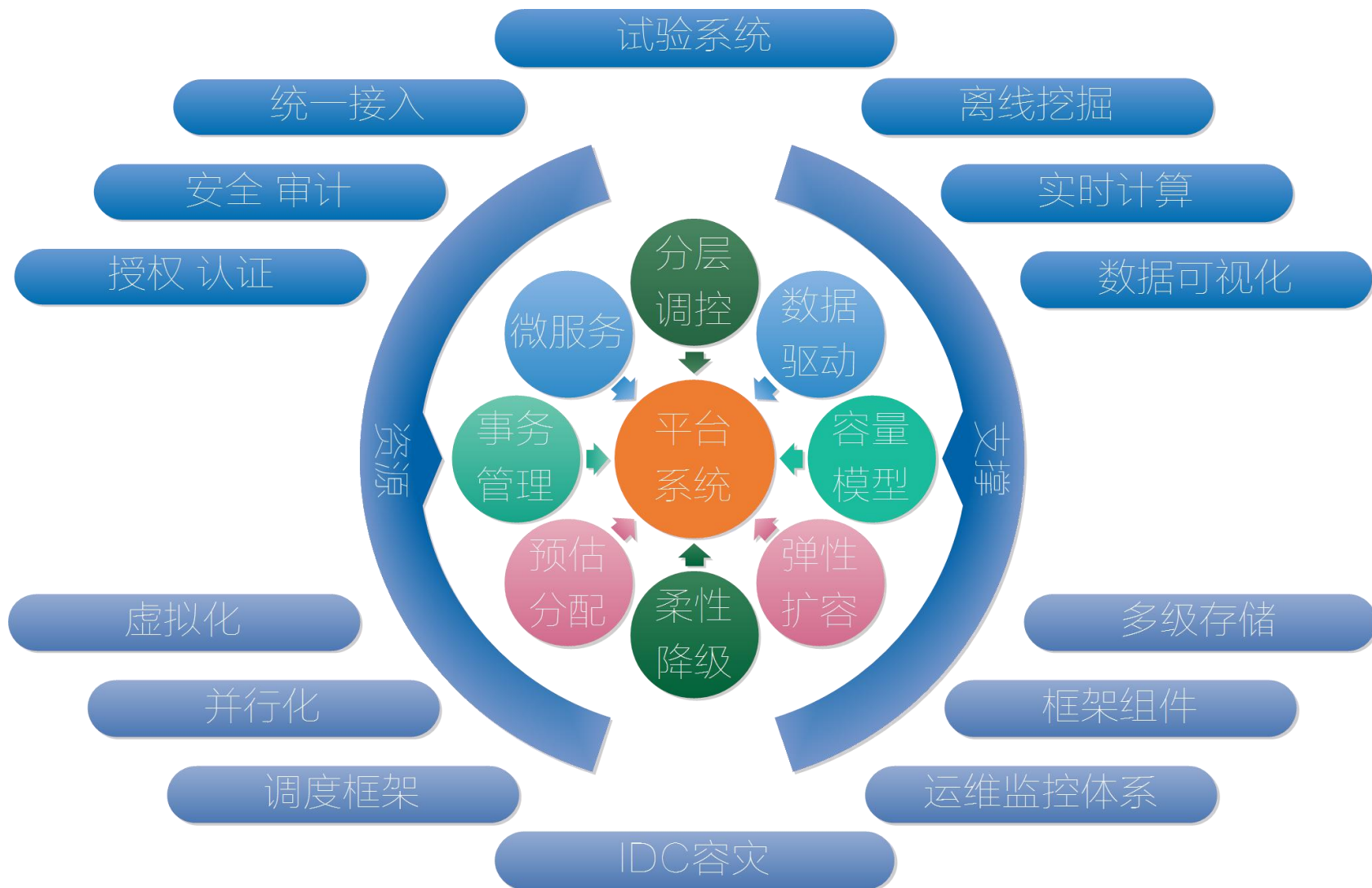
4 对抗
容量模型

6 循环
数据驱动的闭环设计

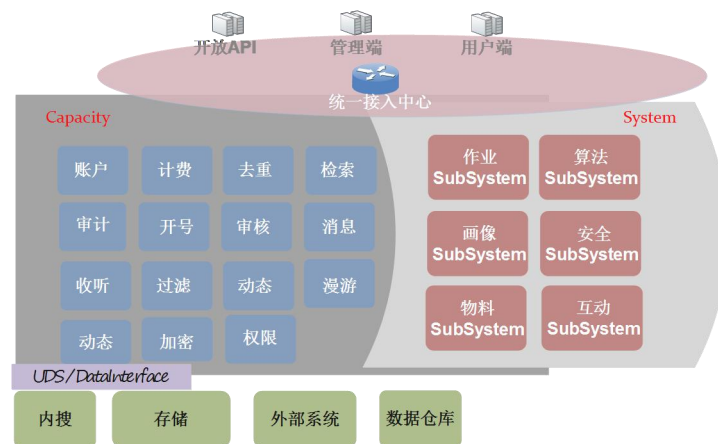
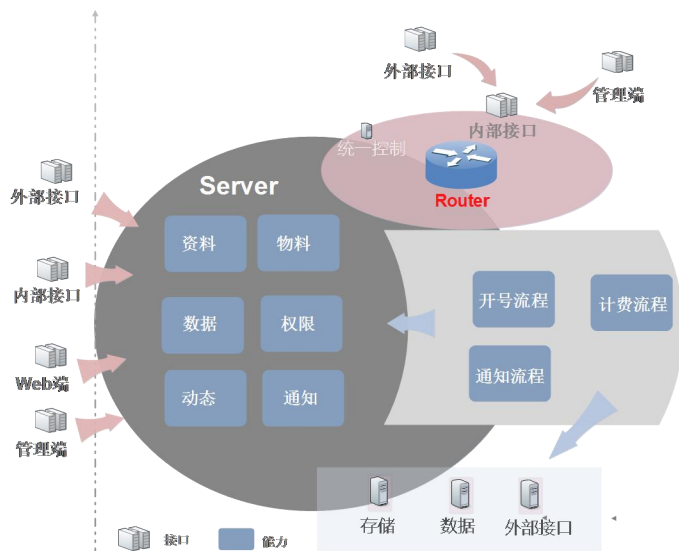
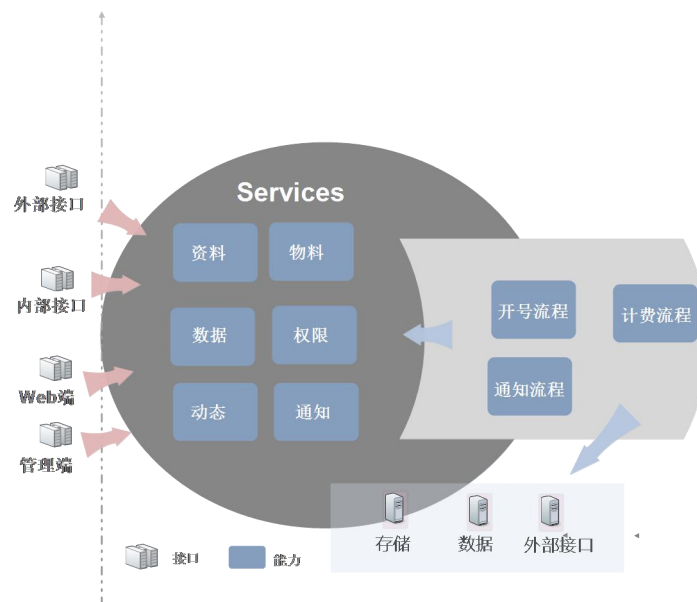
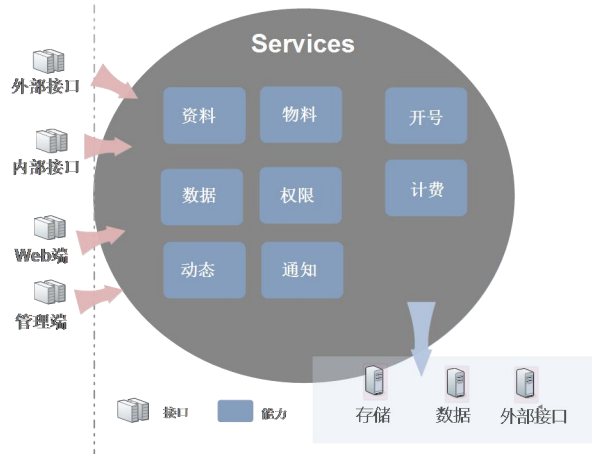
3 建造
框架与调度

5 扩张
容灾选型

平台系统



平台系统



接入平面

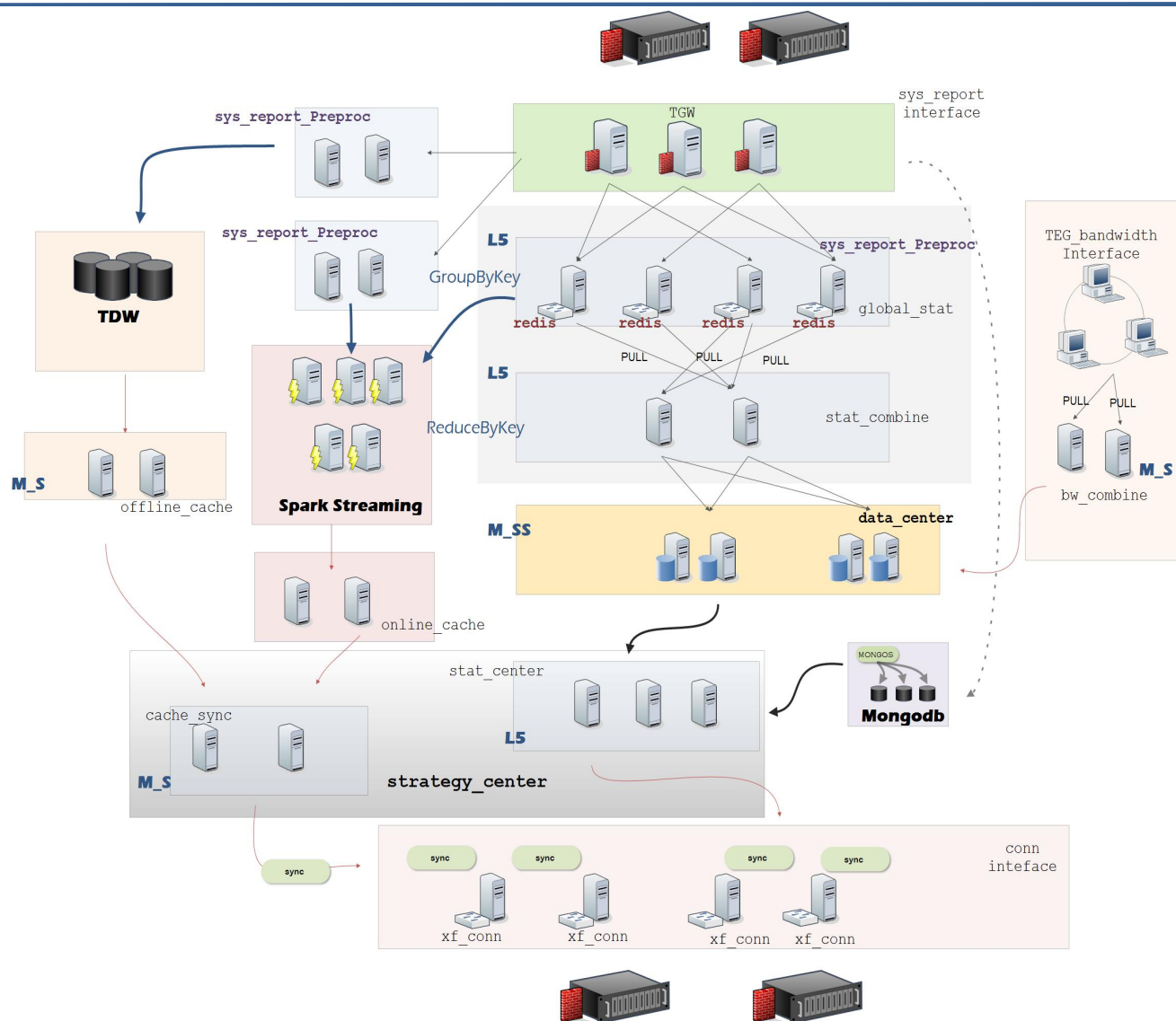
1 接入平面的分类

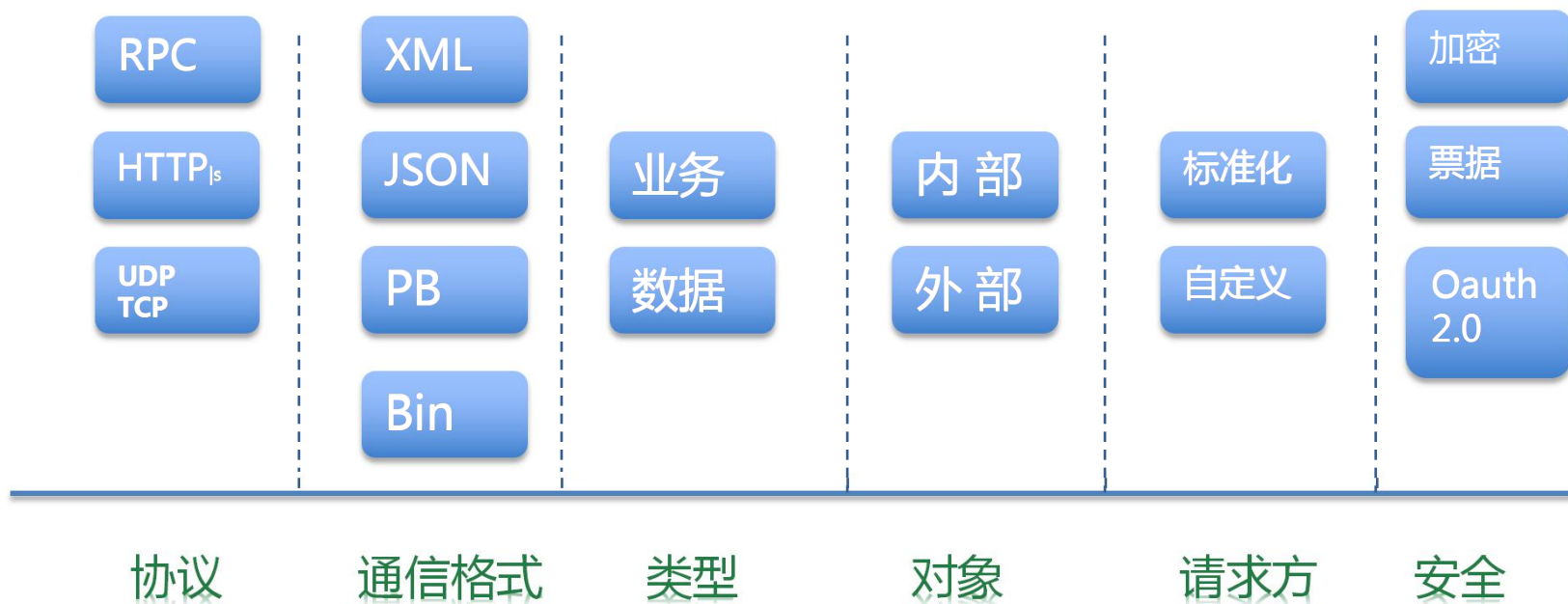
2 微服务路由

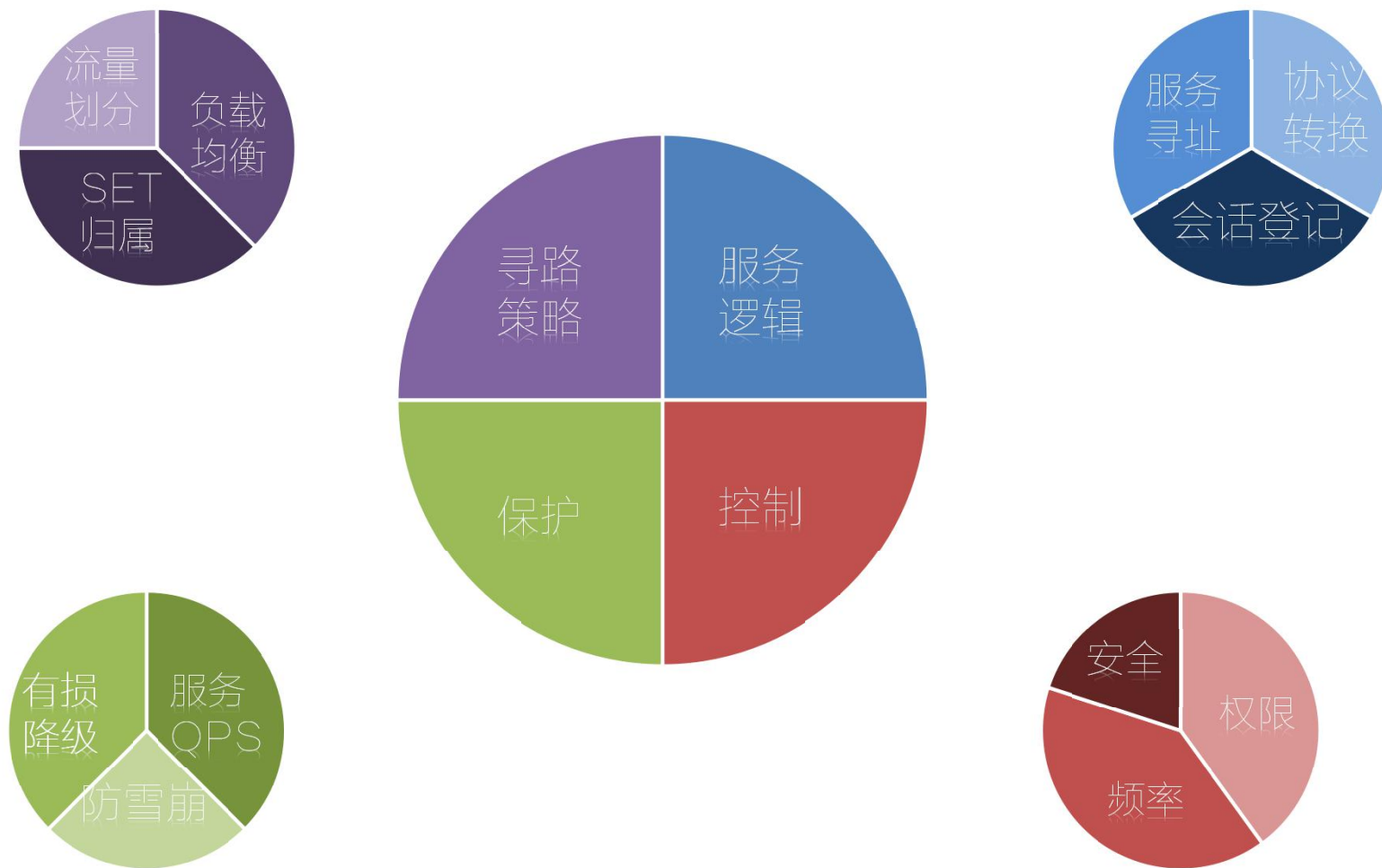
3 权限与频控

4 全网流量调度

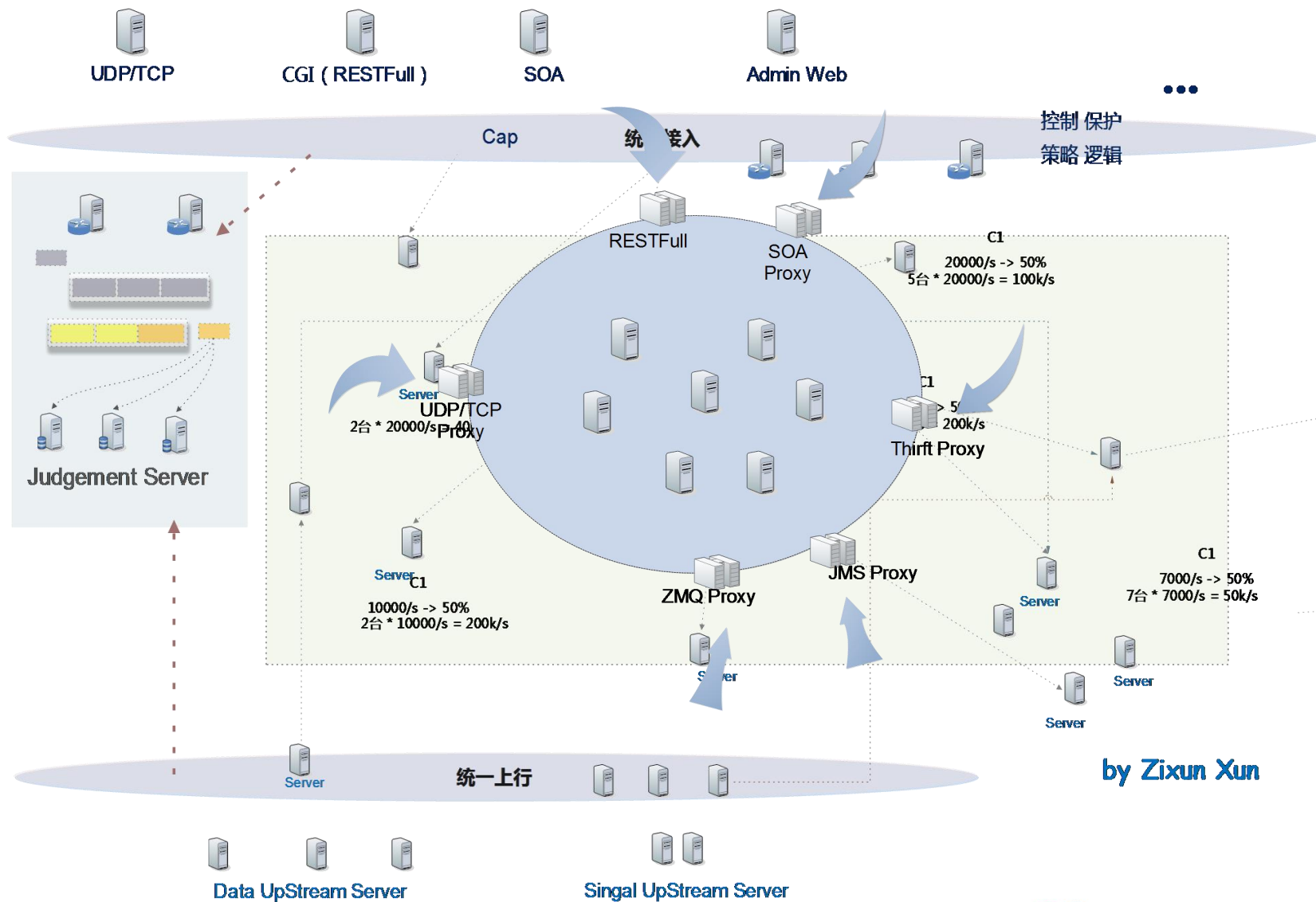
接入平面 | 接入的分类







接入平面 | 权限与频控

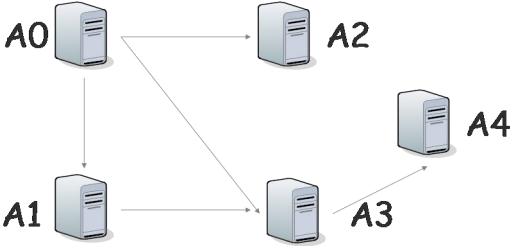
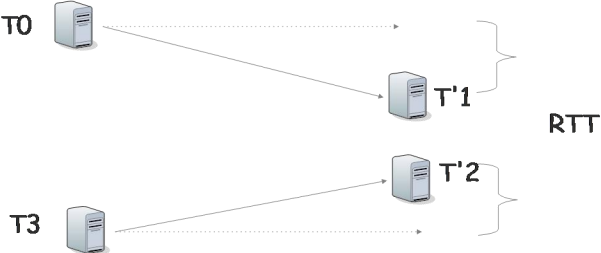


主机时间不同步

滑动平均
自回归模型
神经网络

网络延迟问题

DC分区问题



| 节点 | 0 | 1 | 2 | 3 | 4 |
|----|---|---|---|---|---|
| 0 | 0 | 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 |

| | A0 | A1 | A2 | A3 | A4 |
|----|----|----|----|----|----|
| A0 | 0 | 1 | 1 | 1 | 0 |
| A1 | 1 | 0 | 0 | 1 | 0 |
| A2 | 1 | 0 | 0 | 0 | 0 |
| A3 | 1 | 1 | 0 | 0 | 1 |
| A4 | 0 | 0 | 0 | 1 | 0 |

3 2 1 3 1

A0

| | A0 | A1 | A2 | A3 | A4 |
|----|----|----|----|----|----|
| A0 | 0 | 1 | 1 | 1 | 0 |
| A1 | 0 | 0 | 0 | 0 | 0 |
| A2 | 0 | 0 | 0 | 0 | 0 |
| A3 | 0 | 0 | 0 | 0 | 0 |
| A4 | 0 | 0 | 0 | 1 | 0 |

0 1 1 2 0

A0 A3

| | A0 | A1 | A2 | A3 | A4 |
|----|----|----|----|----|----|
| A0 | 0 | 0 | 0 | 0 | 0 |
| A1 | 0 | 0 | 0 | 0 | 0 |
| A2 | 0 | 0 | 0 | 0 | 0 |
| A3 | 0 | 0 | 0 | 0 | 0 |
| A4 | 0 | 0 | 0 | 0 | 0 |

0 0 0 0 0

A0 A3

框架与调度

1 服务框架

2 调度

3 虚拟化技术

AN OPERATIONS MODEL FOR MICROSERVICES

C-1 C-2

Legend

- LB = Load Balancer
- CB = Circuit Breaker

Edge server
(Reverse Proxy & Router)

OAuth - SSO - Proxy

CB/LB

OAuth
Authorization
Server

API Services

OAuth Res
API-1
CB/LB

OAuth Res
API-2
CB/LB

OAuth Res
API-3
CB/LB

Service
Discovery

Configuration
Server

Composite Services

MS-3
CB/LB

Monitor
Dashboard

Core Services

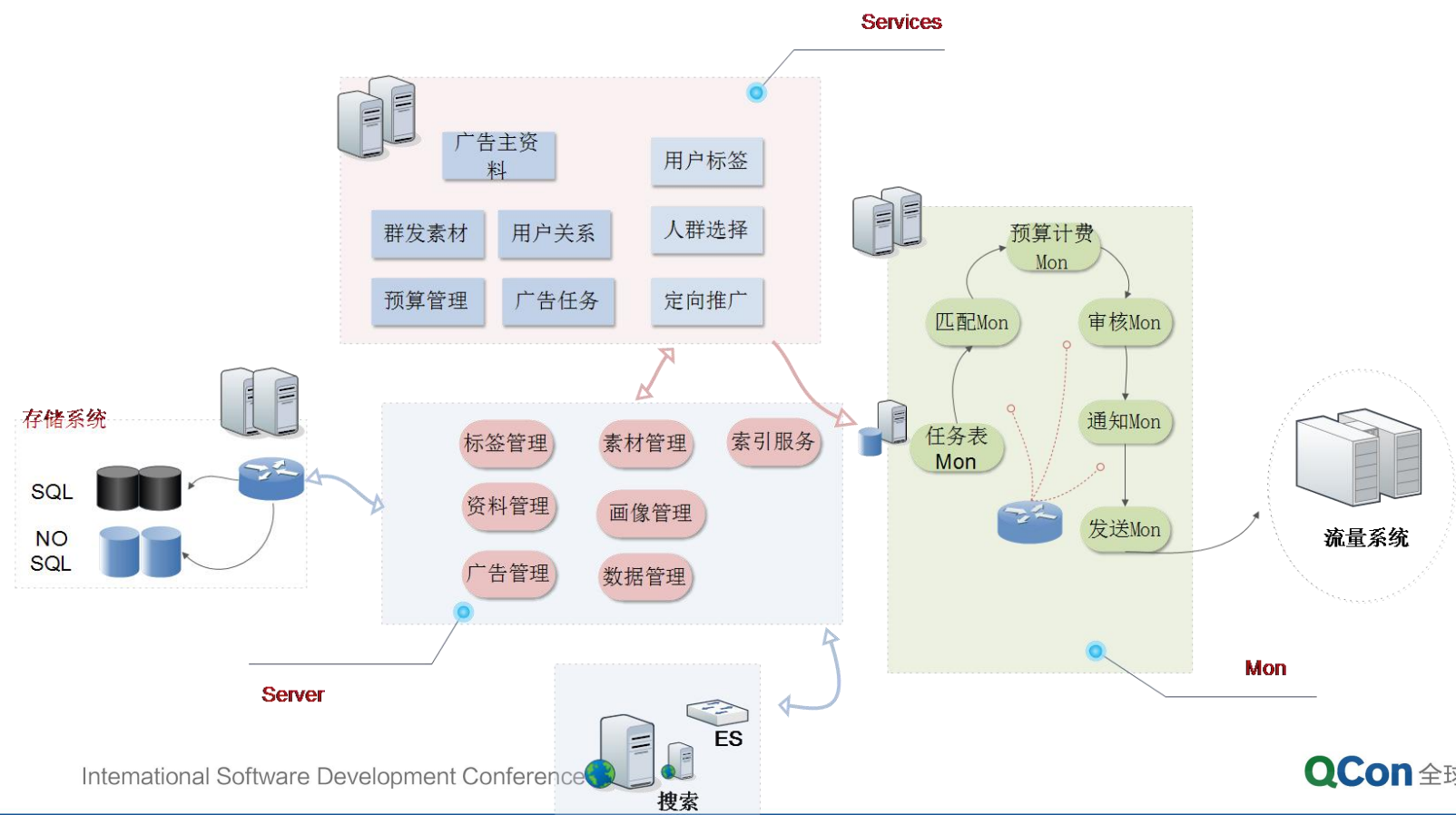
MS-1

MS-2

Logging
Analyses
Dashboard

调度与分配 | 服务框架

| | 特点 | 框架 | 模型 |
|---------|---------|-------------------|------------|
| Service | 对外的能力展现 | Async / Coroutine | RestFull |
| Server | 内部的能力封装 | Leader-Follower | MicoServer |
| Mon | 作业任务 | Reactor/ Proactor | SOA |

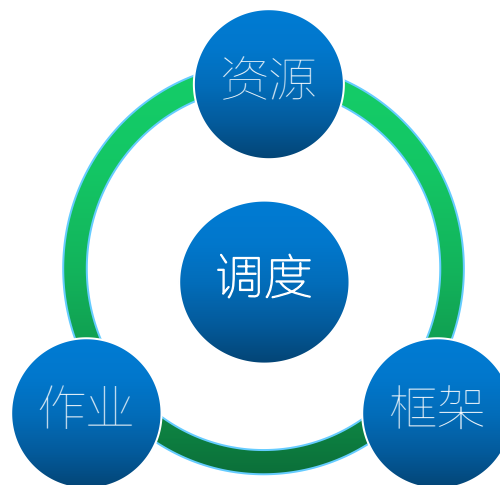


单机调度：

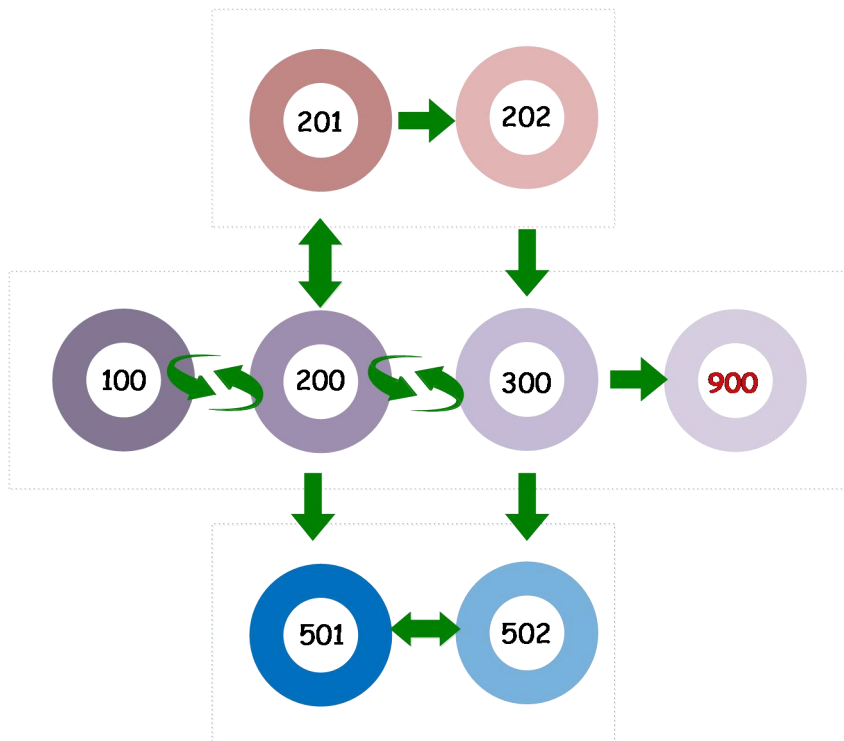
Linux Process Scheduler
Linux IO Scheduler

分布式：

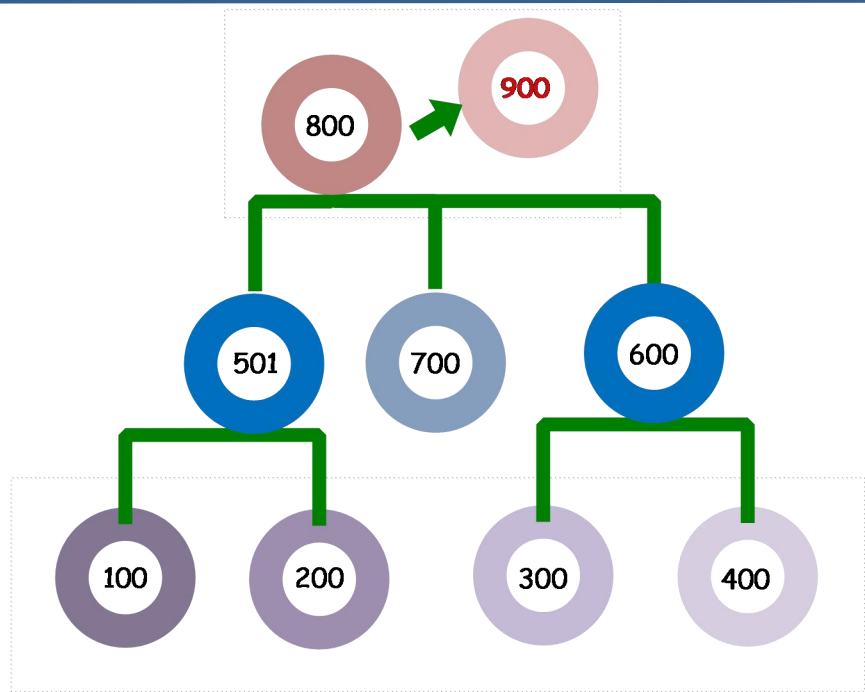
maui torque
mesos
yarn



| | 资源 | 框架 | 作业单元 | 调度 |
|-------|------------|-----------|---------|------|
| Linux | CPU时间 | Scheduler | Process | CFS |
| Mesos | 集群硬件资源 | MR/Spark | 并行化算子 | DRF |
| Ads | empression | RTB | Request | eCPM |



Filter模式状态机
串行



Switch模式 状态机
并行

资源:

硬件(CPU,内存,外存,网络) -> Cloud Computation

带宽 -> CDN

流量 -> A D X

类别:

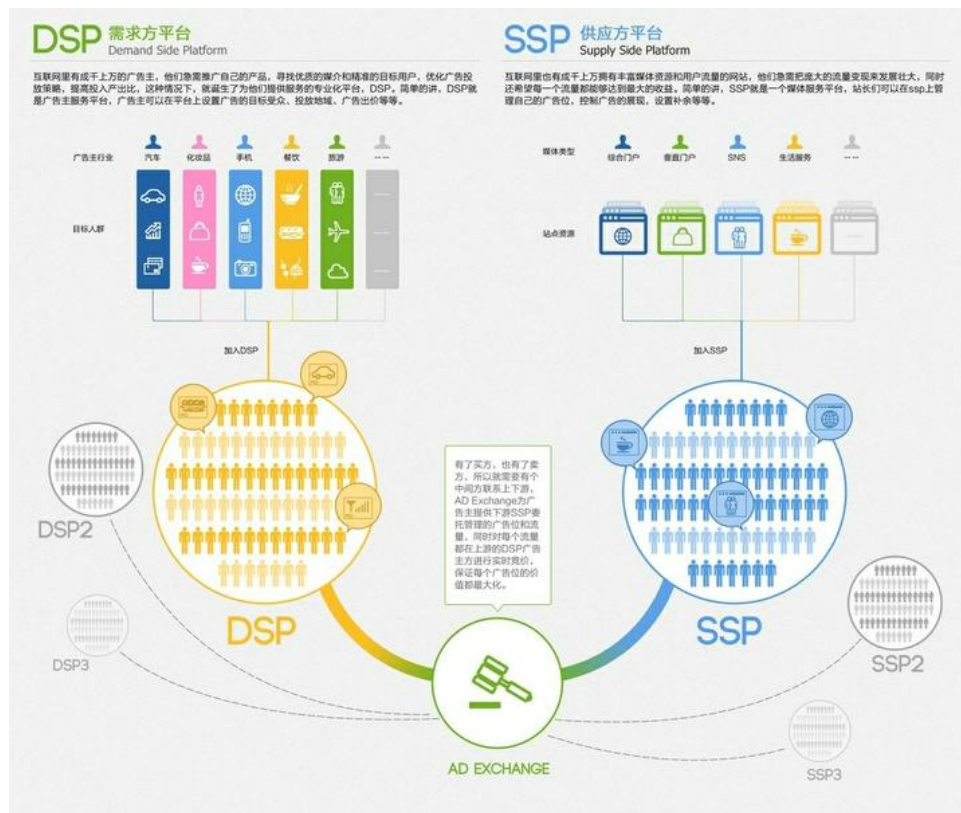
实时决策 (L2R HWM 贪心)

离线决策 (运筹学)

分类, 预测

决策依据:

三方最优
平台, 需求方, 资源供给方



容量模型

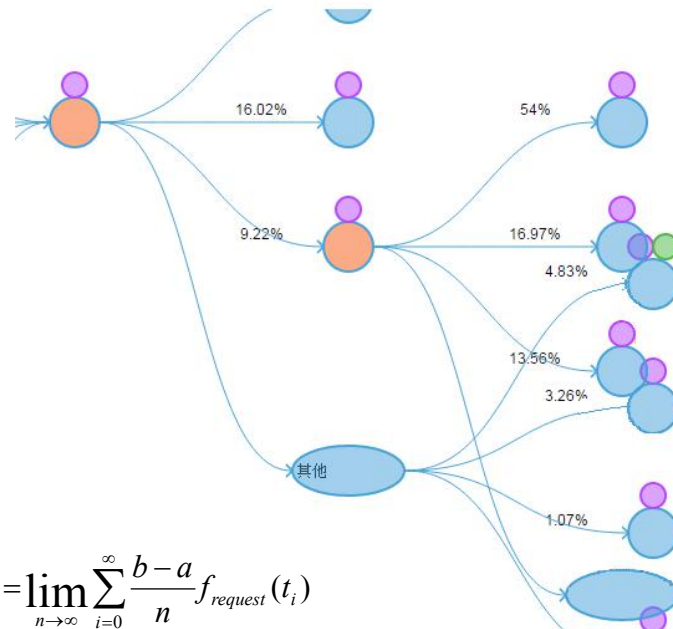
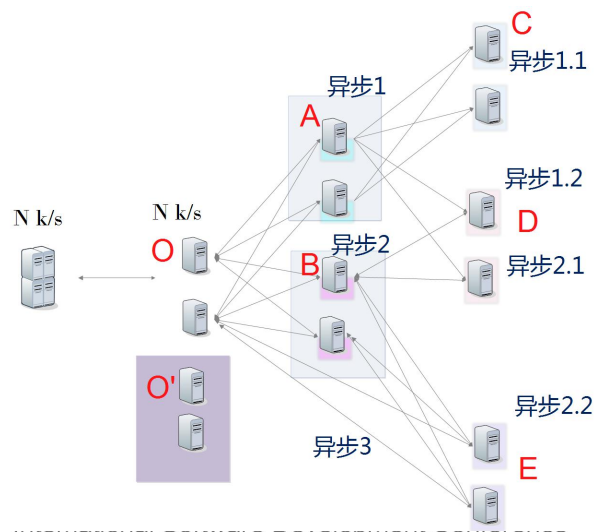
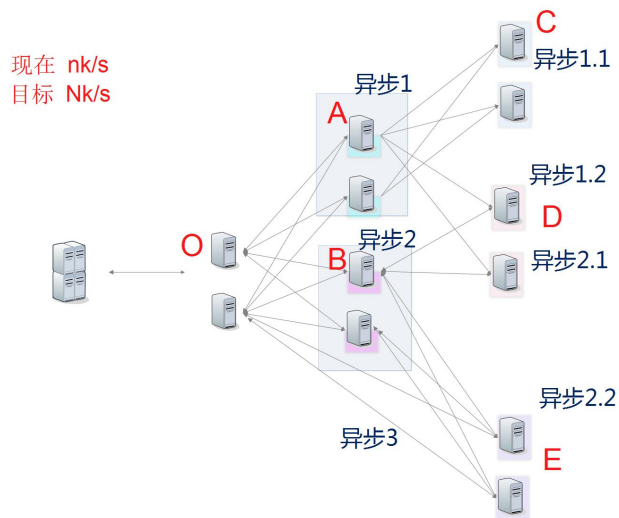
1 一切尽在掌控

2 压测与扩容

3 试验系统

4 SET化部署

容量模型 | 压测与扩容



$$\int_a^b f_{request}(x)dx = \lim_{n \rightarrow \infty} \sum_{i=0}^{\infty} \frac{b-a}{n} f_{request}(t_i)$$

$$\int_{a+t}^{b+t} f_{response}(x)dx = \lim_{n \rightarrow \infty} \sum_{i=0}^{\infty} \frac{b-a}{n} f_{response}(t_i)$$

$$1s \Rightarrow L$$

$$1s = \sum T_{Brust}$$

$$J \Rightarrow T_{Brust}$$

$$All_{except} = J \cdot \frac{1}{T_{Brust}} \sum_i^k Num_i = J$$

$$L = All_{reality} = J \cdot \frac{1}{(T_{Brust} + 2T_{rtt/2} + T_1 + T_2 + T_{..})}$$

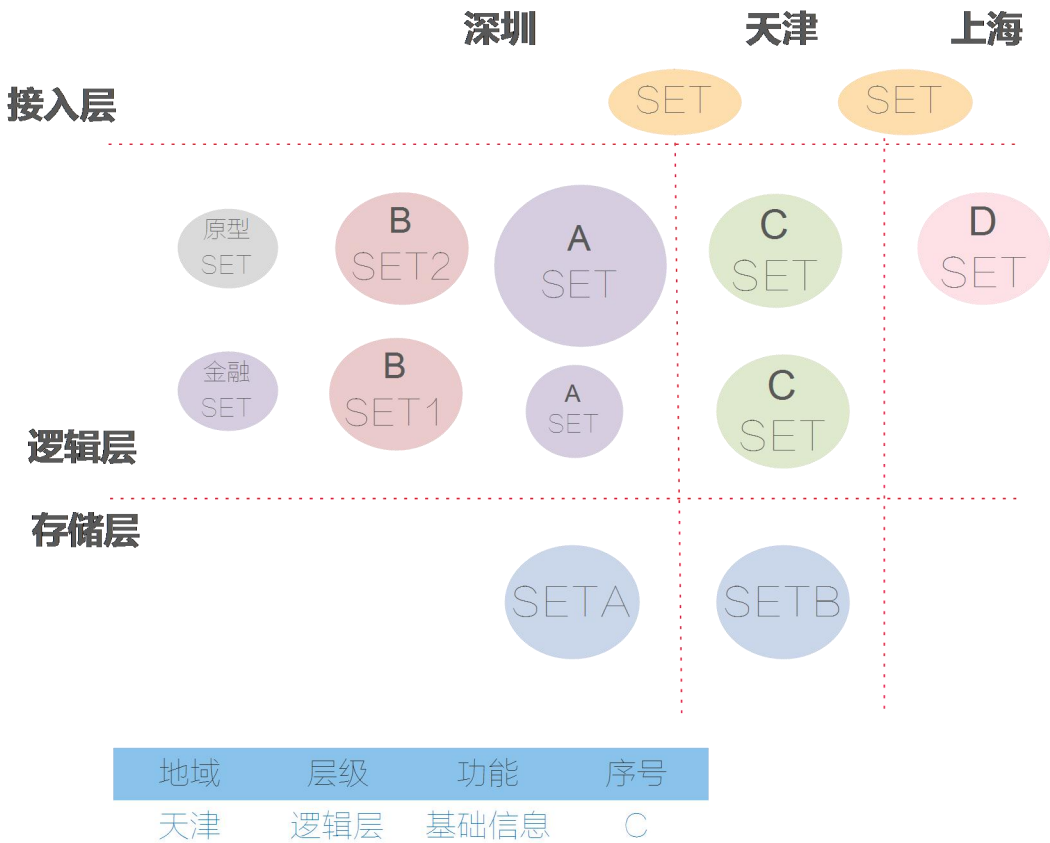
$$T_1 = f(Num_1, \frac{1}{P_1}, T_{rtt})$$

$$Num_i = f(J / NodeNum, k)$$

$$Num_i \propto k$$

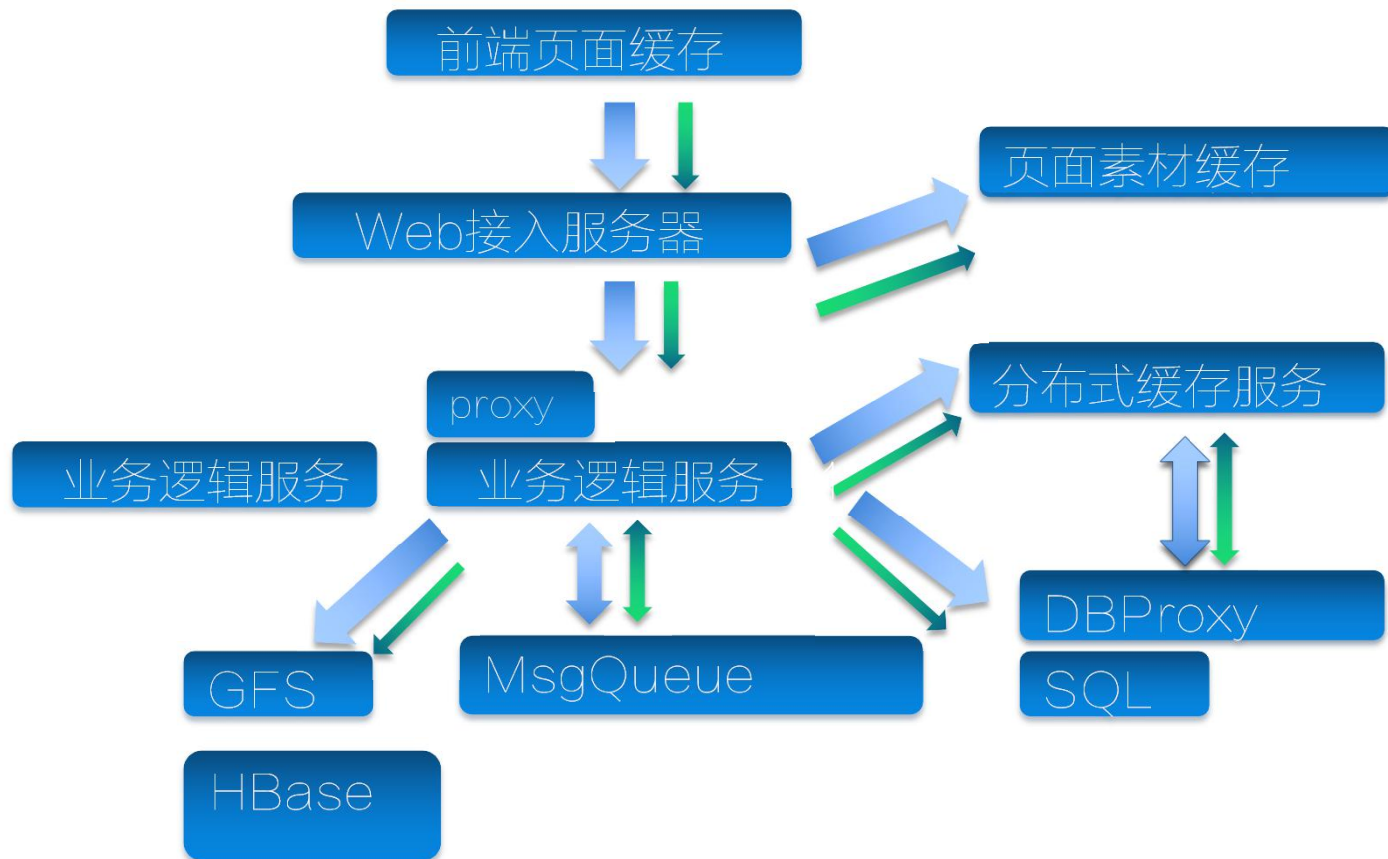
要求：
标准化、规模化、模块化

- 粒度：
- 1) 承载 容量
 - 2) 打包 业务
 - 3) 独立子系统



步骤:

- 灰度策略
- 流量切分
- 部署升级
- 观察验证
- 全量上线

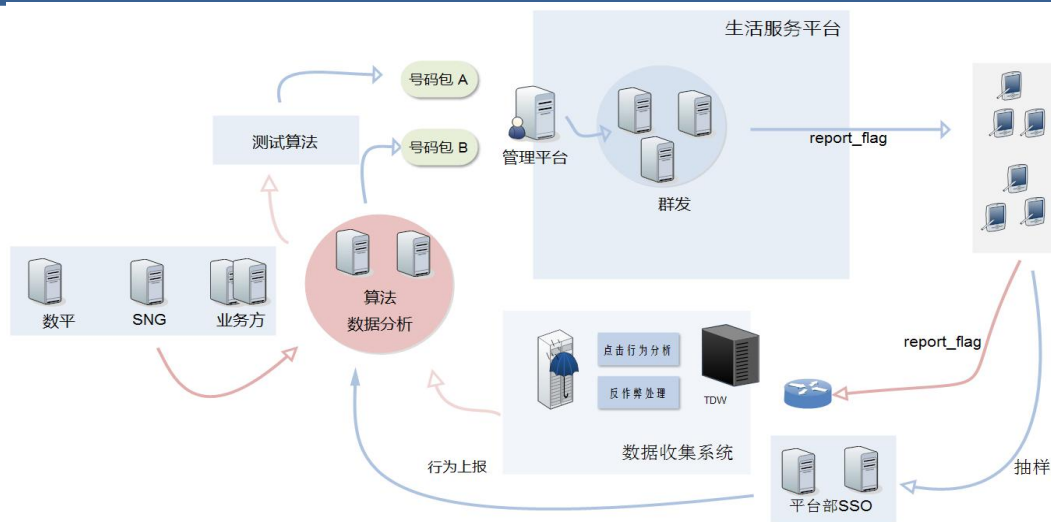


系统层面:

分层流量标记
动态策略

数据层面:

统计学效果检验
多因素方法分析
置信度抽样数
辛普森悖论



容灾



方案比较

容灾选型 | 方案比较

| 系统模型 | 接入/逻辑层 | 存储层 | 方法 | 标准 |
|--------|-------------|------------------|---|--|
| 作业平台系统 | Check Point | 分布式文件系统 P2P网络 | 回退恢复： (backward recovery) 系统从当前错误的 状态回到先前正确 的状态 | 提交作业成功执行 |
| 存储平台系统 | 接入层容灾 | 多副本机制 | 一致性模型 | C A P |
| 推送平台系统 | 分层可靠协议 | 多副本机制 | 协议+备份 | 短信平台Mo Mt 入网标准 72小时 内有接受条件 到达率： 99.99% 及时率 >= 95% |
| 金融平台系统 | 事务性操作 幂等 | 可靠存储 | 多级流水对账 严格事务性 | 可溯源 回退 有据可查 ACID |

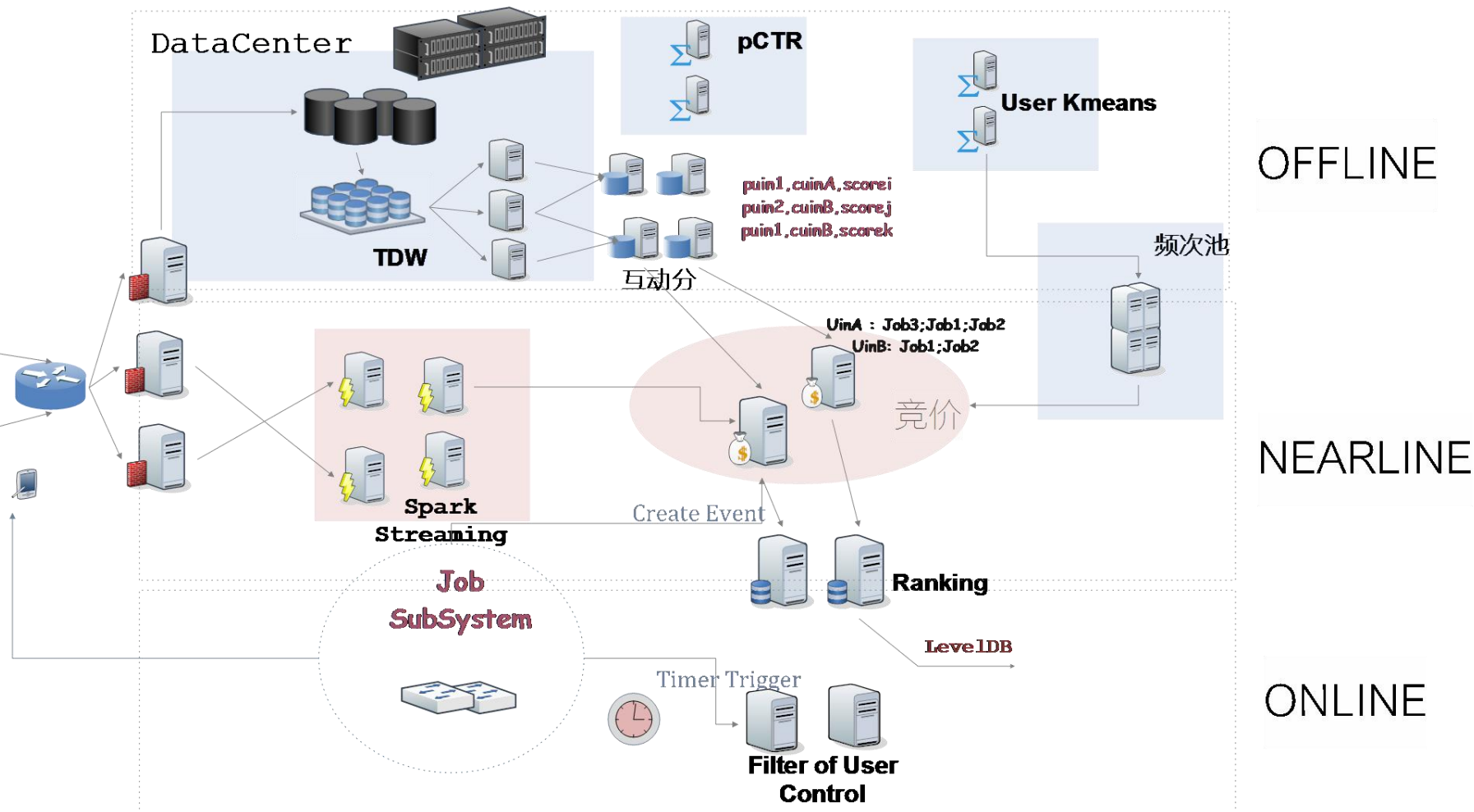
数据闭环

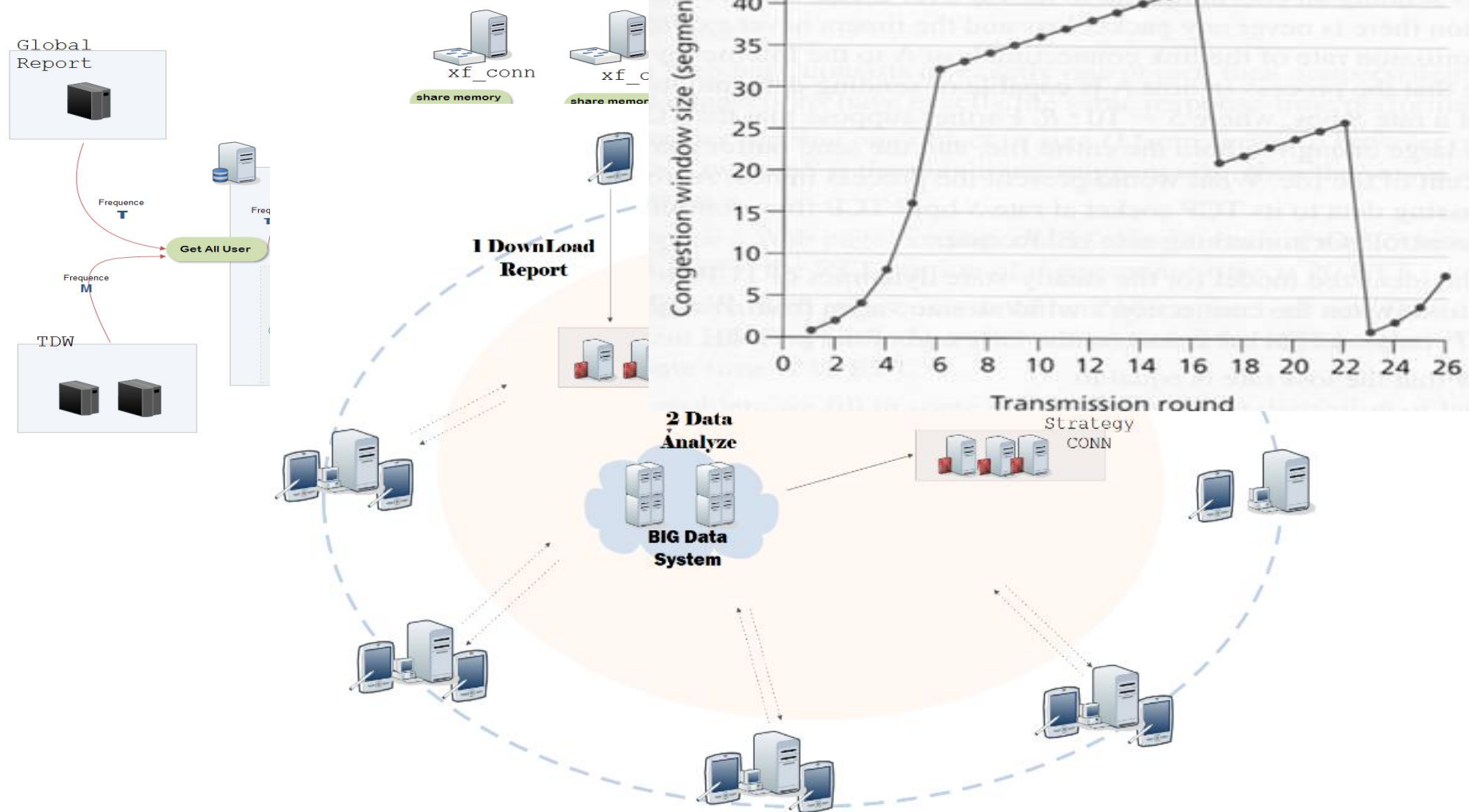
1

数据分层

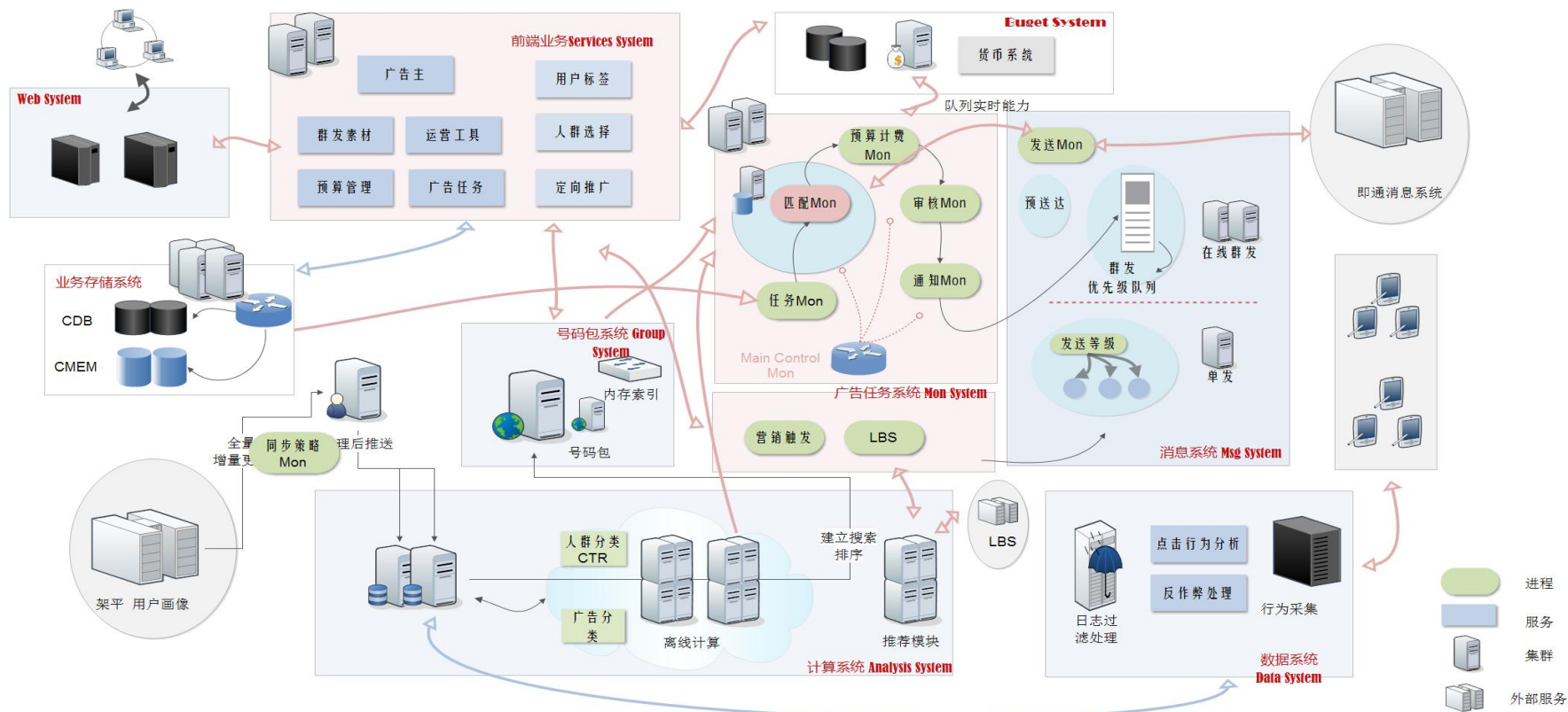
2

数据驱动

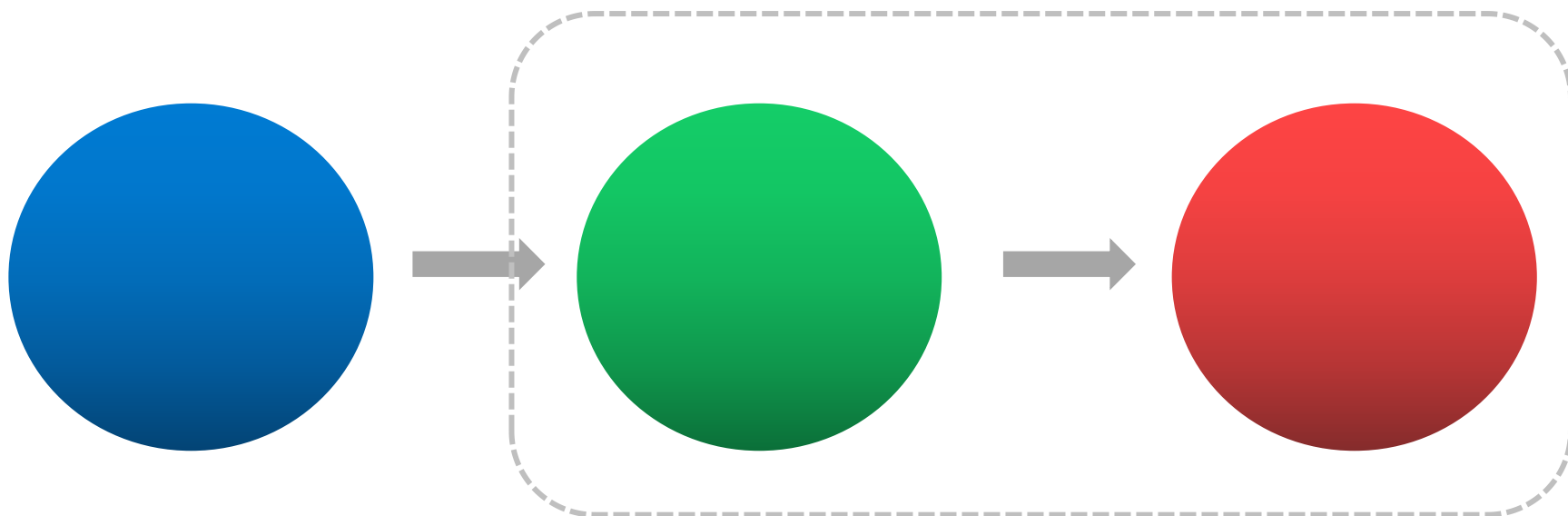




数据闭环 | 数据驱动



平台系统





THANKS!