



# 容量规划与授权限流降级





# problems

线上机器能承受多大的调用量？

系统需要加机器 or 减机器？

如何控制调用和资源访问？



## Part1 线上压测

---

## Part2 容量规划

---

## Part3 授权限流降级

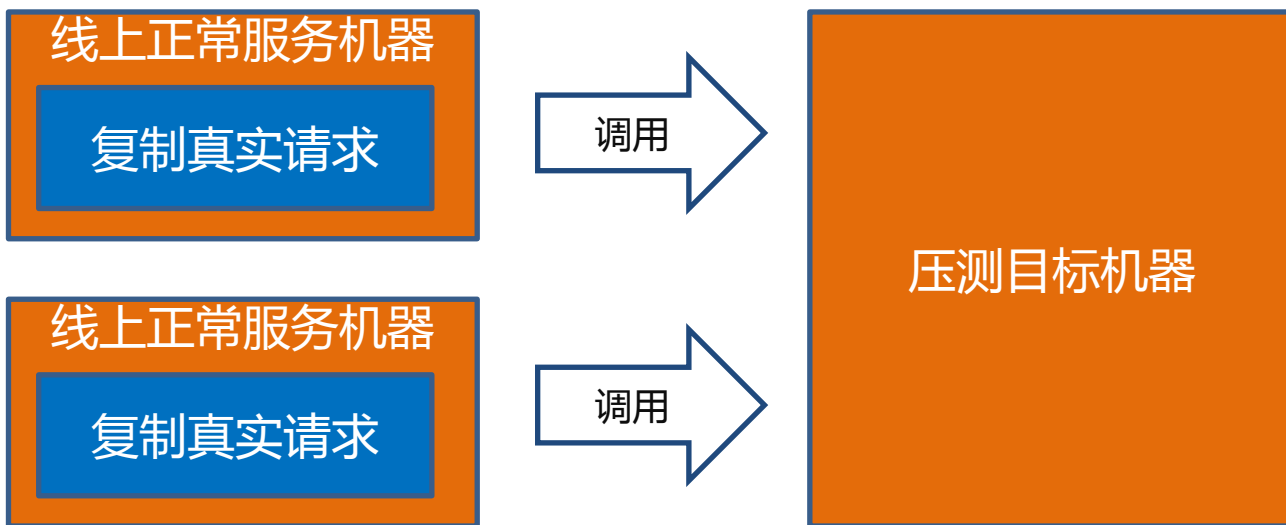
---

- **线上压测方式**
  1. **模拟请求**
  2. **复制请求**
  3. **请求引流转发**
  4. **修改负载均衡权重**

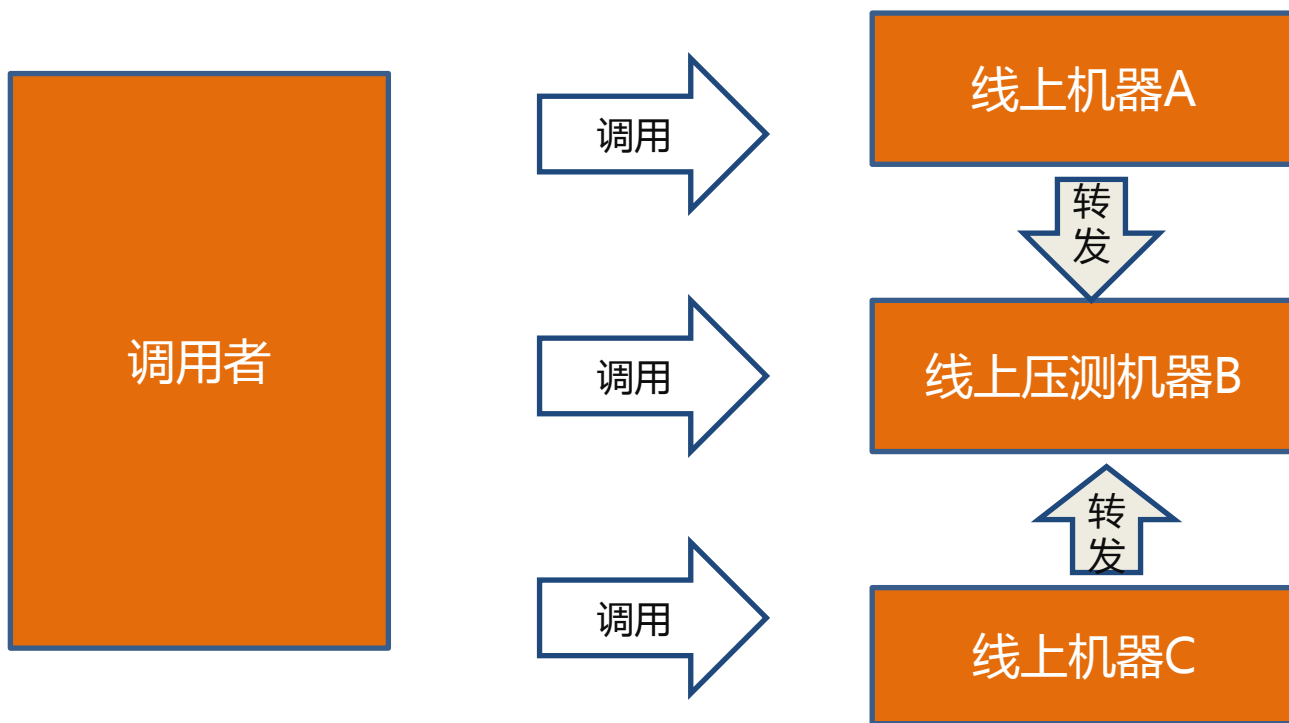
### 模拟请求



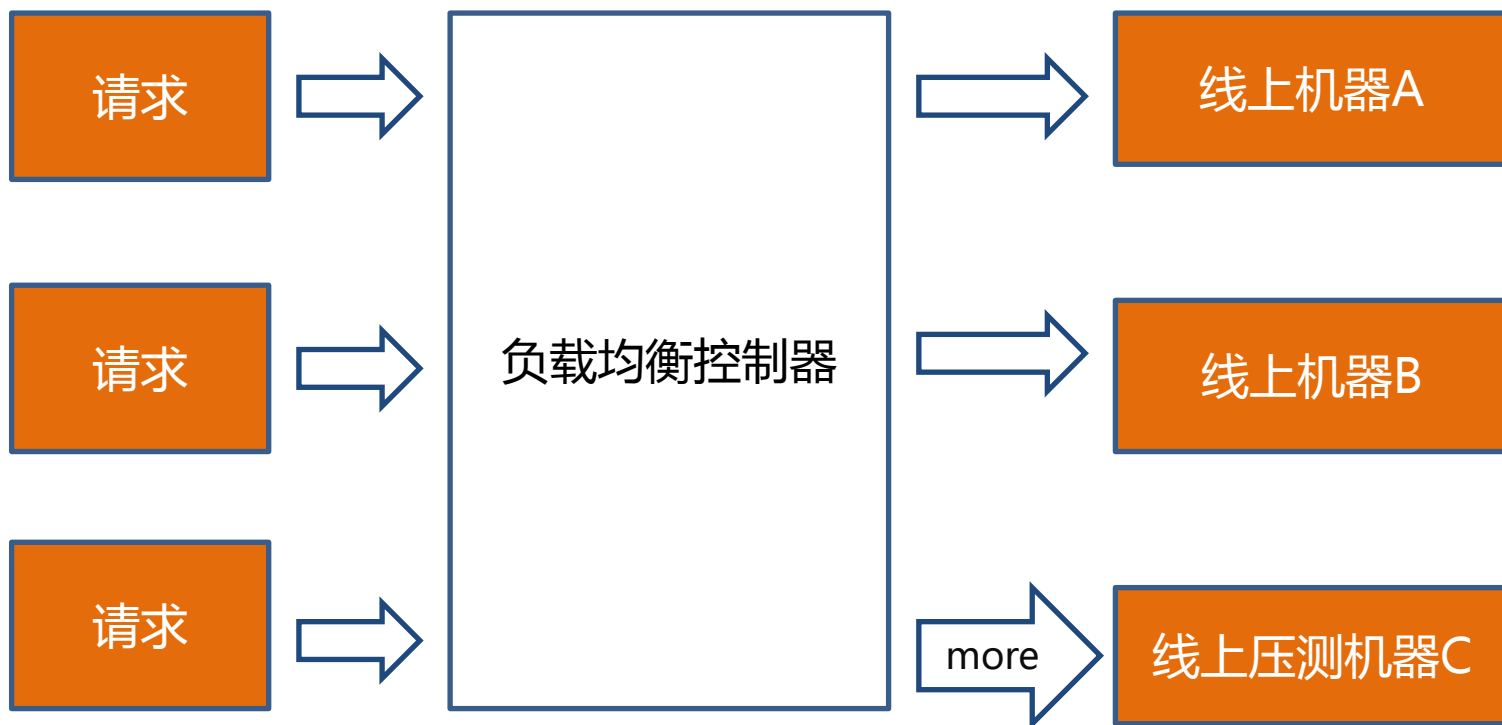
### 复制请求



## 请求引流转发



## 修改负载均衡权重



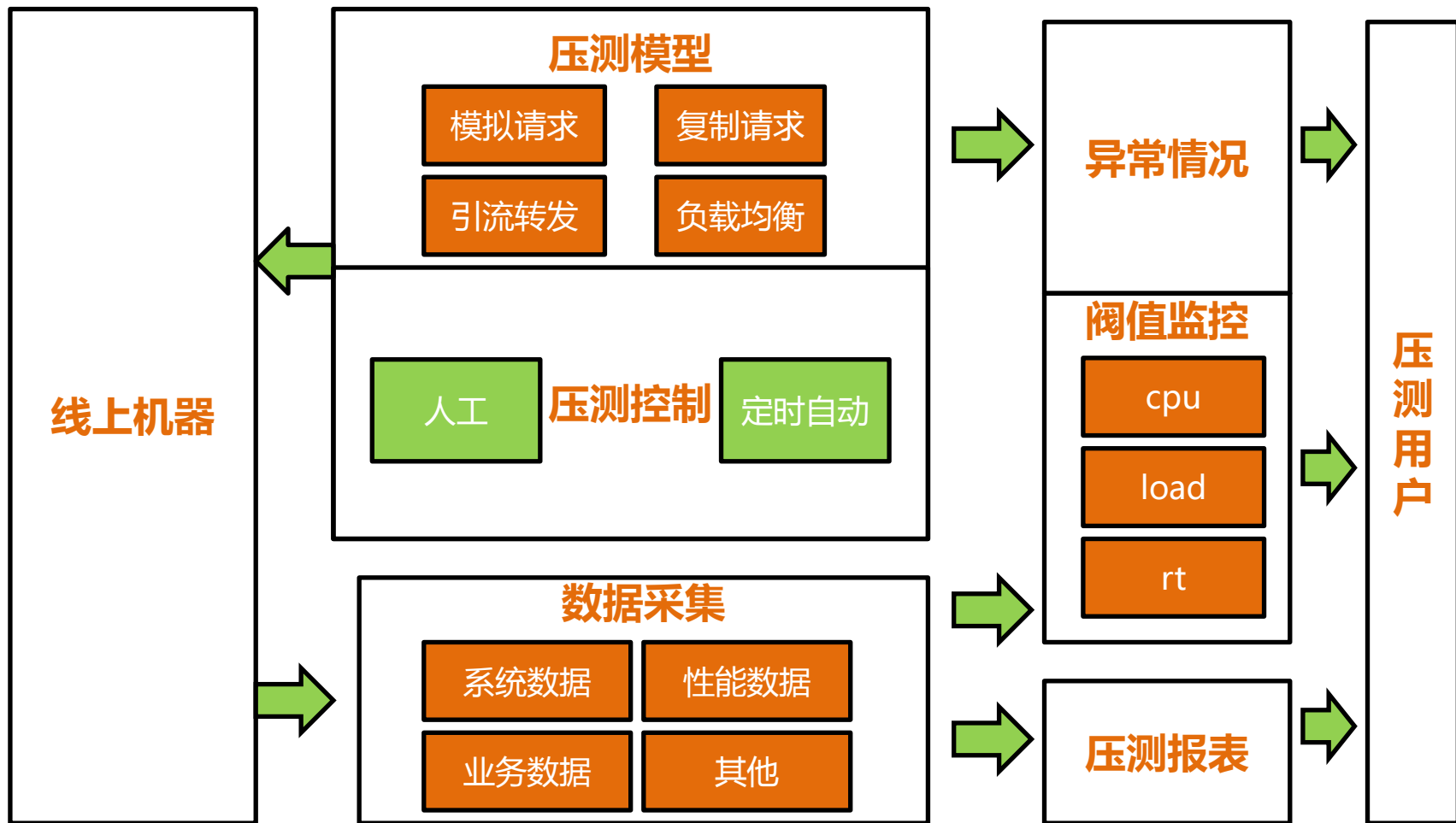
压测方式	优点	缺点
模拟请求	简单应用、对于没什么流量的系统非常好用	模拟请求缺乏真实性，写的请求需特别考虑（脏数据）
复制请求	请求真实，能够放大流量	写请求和响应需特别考虑。需一台不提供服务的机器充当压测目标机
请求引流转发	完全真实的场景，压测数据准确	依赖系统自身的流量、服务类应用不太好转发
修改负载均衡权重	完全真实的场景，压测数据准确	依赖系统自身的流量，需要负载均衡控制器开放接口



- **压测相关工具**

1. **模拟请求: http\_load, webbench, ab, jmeter, Siege 等**
2. **复制请求: tcpcopy, btrace, nginx post\_action , 自定义 agent 等**
3. **请求引流转发 : apache mod\_jk mod\_proxy, nginx proxy 等**
4. **修改负载均衡权重: F5, LVS, SOA service registration 等负载均衡控制器**

## 淘宝压测平台架构



## 容量名词解释

单机能力 = 单台机器压测阈值qps

单机负荷 = 前一天单台机器最大qps

集群能力 = 单机能力 \* 机器数 ( 机器环境一致 )

集群负荷 = 前一天集群最大qps

水位标准

单机房 ( 70% ) , 双机房 ( 40% ) , 三机房 ( 60% )

## 单个系统容量计算

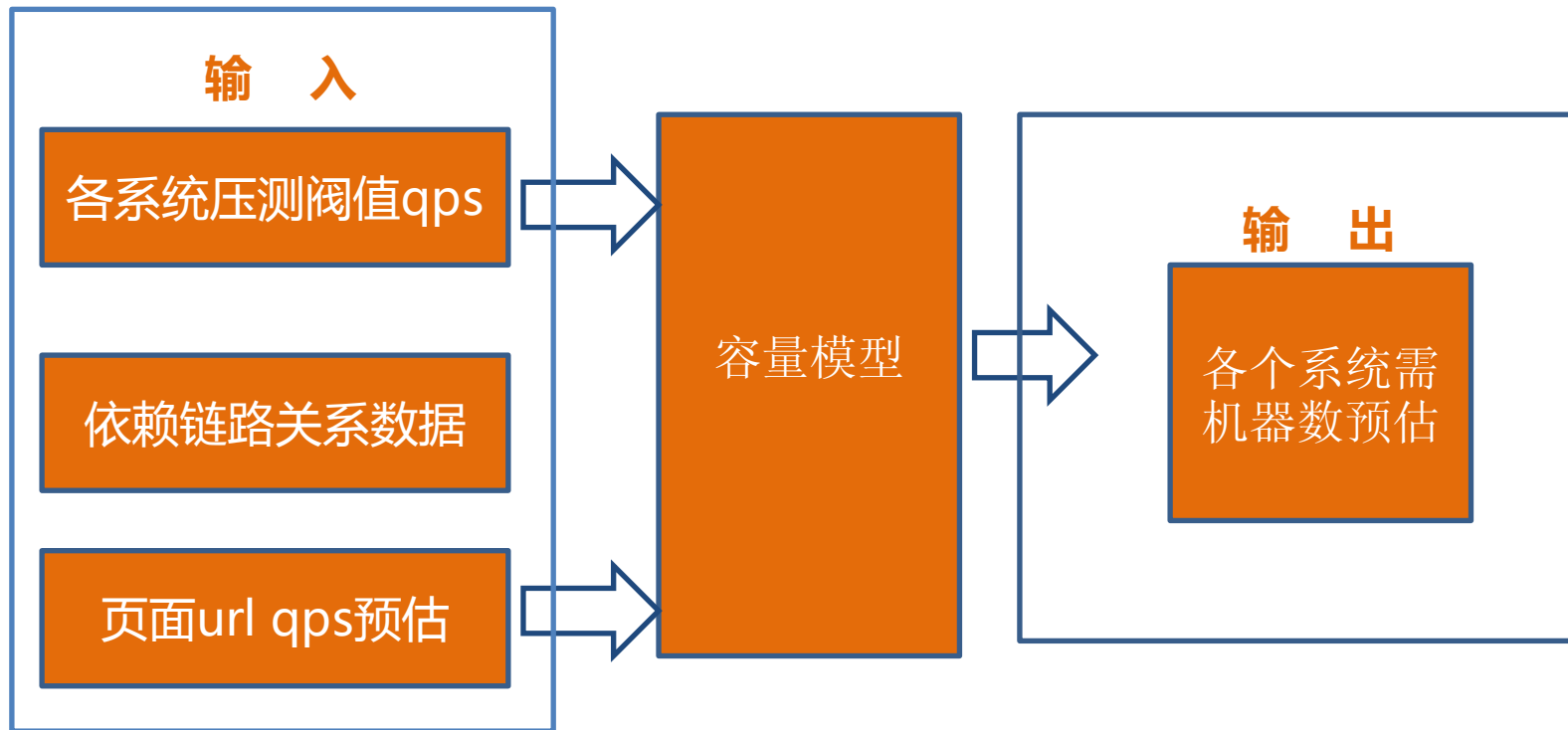
**单机水位 = 单机负荷 / 单机能力 \* 100%**

**集群水位 = 集群负荷 / 集群能力 \* 100%**

**理论机器数 = ( 实际机器数 \* 集群负荷 \* 集群水位 ) / ( 集群能力 \* 水位标准 )**

**机器增减 = 理论机器数 - 实际机器数**

## 链路容量计算



依赖链路关系数据通过eagleeye获取 (类似谷歌的dapper, twitter的zipkin)

## 授权限流降级场景 ???

我提供的某个接口或者资源我只想被A应用访问

我提供的某个接口或者资源我不想被B应用访问

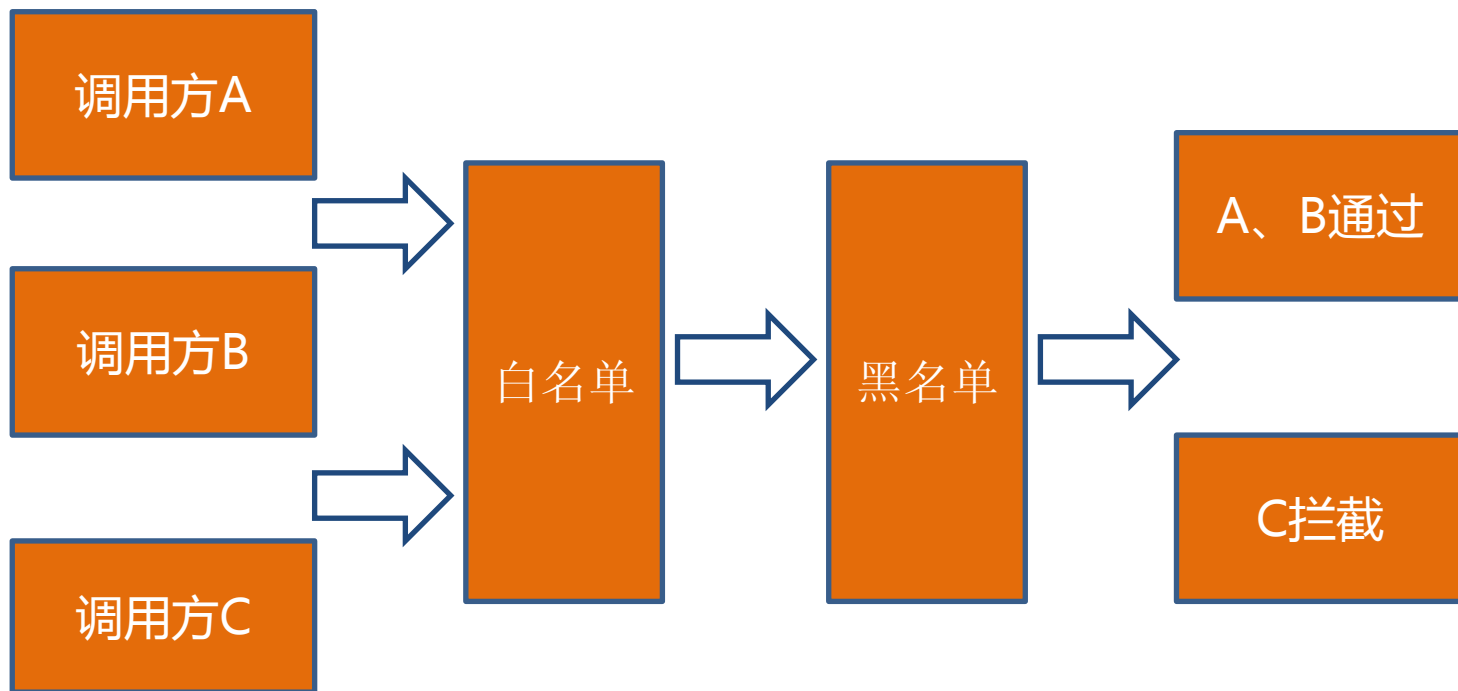
防止我的某个接口或者某种资源被过度访问

防止我对某个接口或者某种资源的过度访问

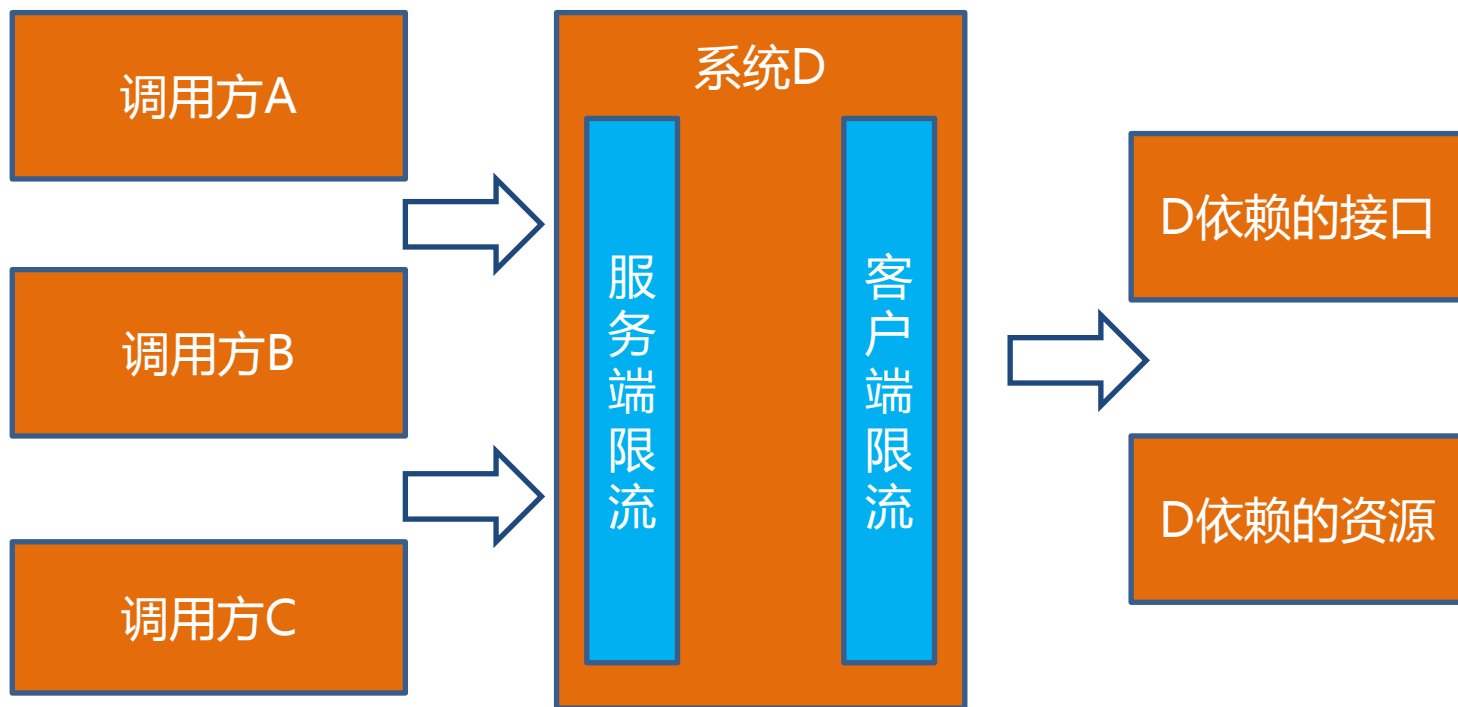
系统负载太高可以降级掉不重要的应用对我的调用

依赖的非关键调用长时间没有响应可以对其进行降级

## 授权

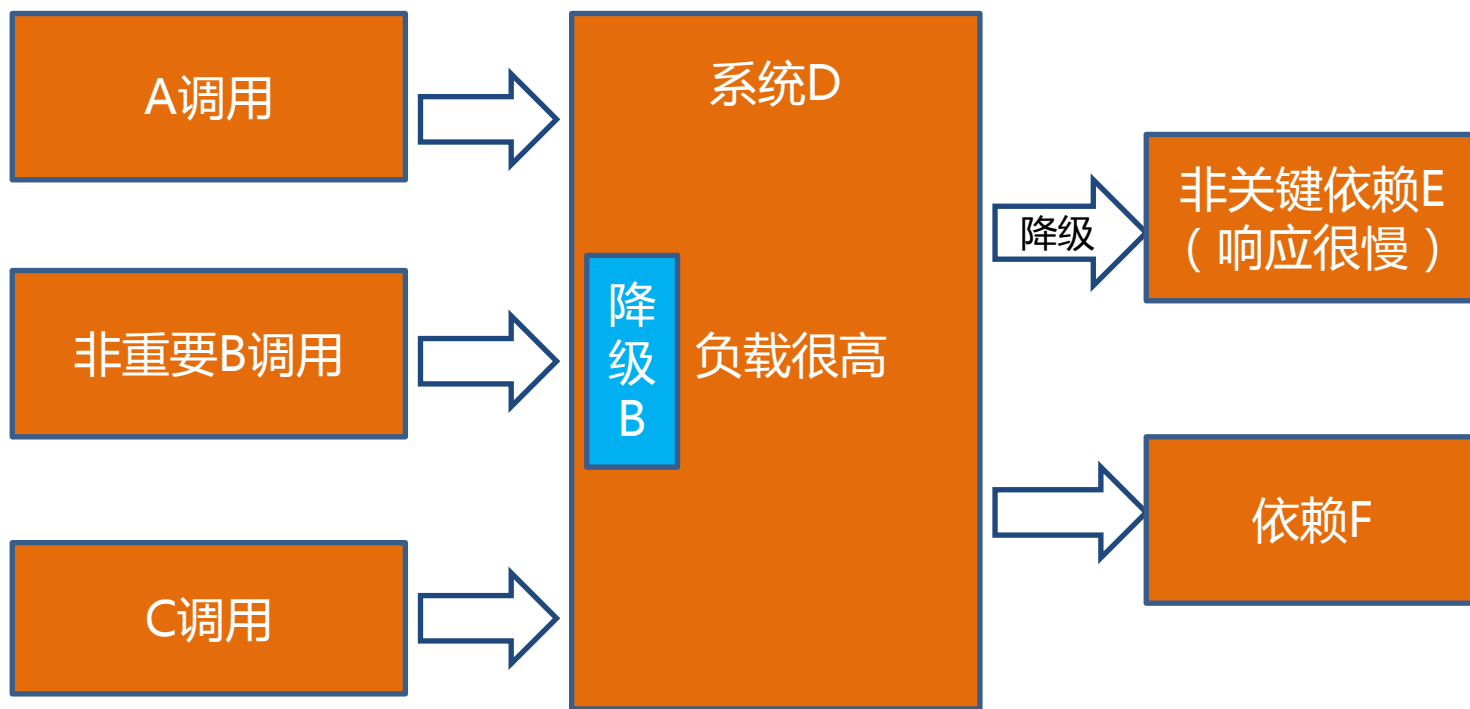


## 限流

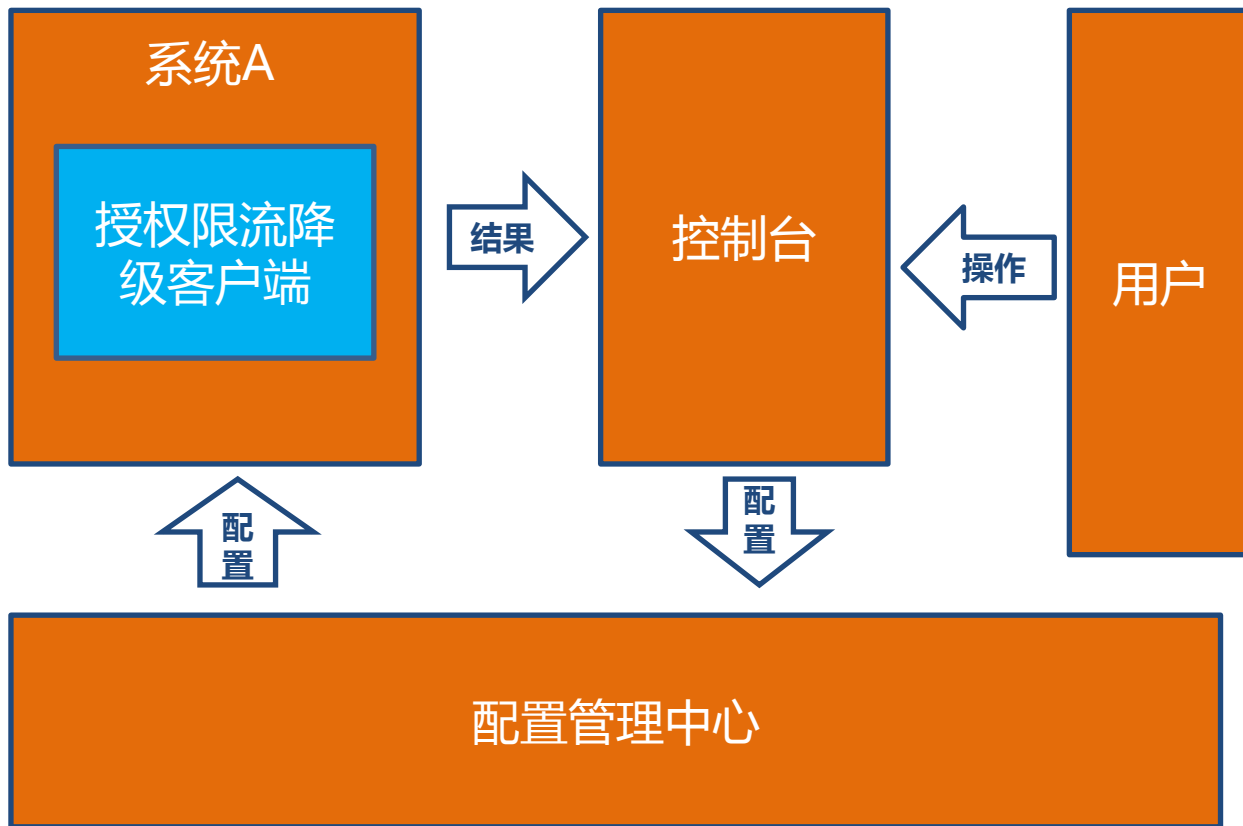




## 降级



## 授权限流降级部署结构图



THANK  
YOU