

架构实战营直播

接口高可用架构设计

李运华

前阿里资深技术专家（P9）



自我介绍



- 5 年华为，8 年 UC，2 年蚂蚁金服 / 前阿里资深技术专家
- 《编程的逻辑》 / 《从 0 开始学架构》（5w+） / 大厂晋升指南（1w+）



华仔
我已加入学习，邀你一起！

极客时间

从 0 开始学架构

前阿里 P9 技术专家的
实战架构心法

全集

你将获得

- 1. 理解架构设计的本质和目的
- 2. 掌握高性能和高可用架构模式
- 3. 走进 BAT 标准技术架构实战
- 4. 从编程到架构，实现思维跃迁

李运华
前阿里资深技术专家

原价 ¥99
新人仅 **¥19.9**
拼团价 ¥79

华仔
我已加入学习，邀你一起！

极客时间

大厂晋升指南

前阿里 P9 技术专家的升职心法

你将获得

- ☑ 从 P5 到 P9 的升职秘籍
- ☑ 实用的职场晋升技巧
- ☑ 19 个高效工作和学习方法
- ☑ 完整的职场晋升路线

李运华
前阿里资深技术专家 (P9)

2 天倒计时!
原价 ¥129 早鸟优惠 **¥99**
新人专享，仅需 ¥19.9

极客时间 | 训练营

架构实战营

第 4 期

跟着阿里 P9，系统提升你的架构能力

李运华
前阿里资深技术专家 (P9)

你将获得

- ☑ P9 架构师原创的“面向复杂度架构设计”方法论
- ☑ 切合实际业务，详解 4 种核心业务架构设计套路
- ☑ 手把手带你实现业务高性能高可用的中间件系统
- ☑ 业务驱动实战，亿级 IM 系统架构设计与演进

抢占优惠名额
>>>>>>>

教学目标

1. 掌握接口级别高可用设计的架构模式和技巧



架构决定系统质量上限，代码决定系统质量下限！

目录

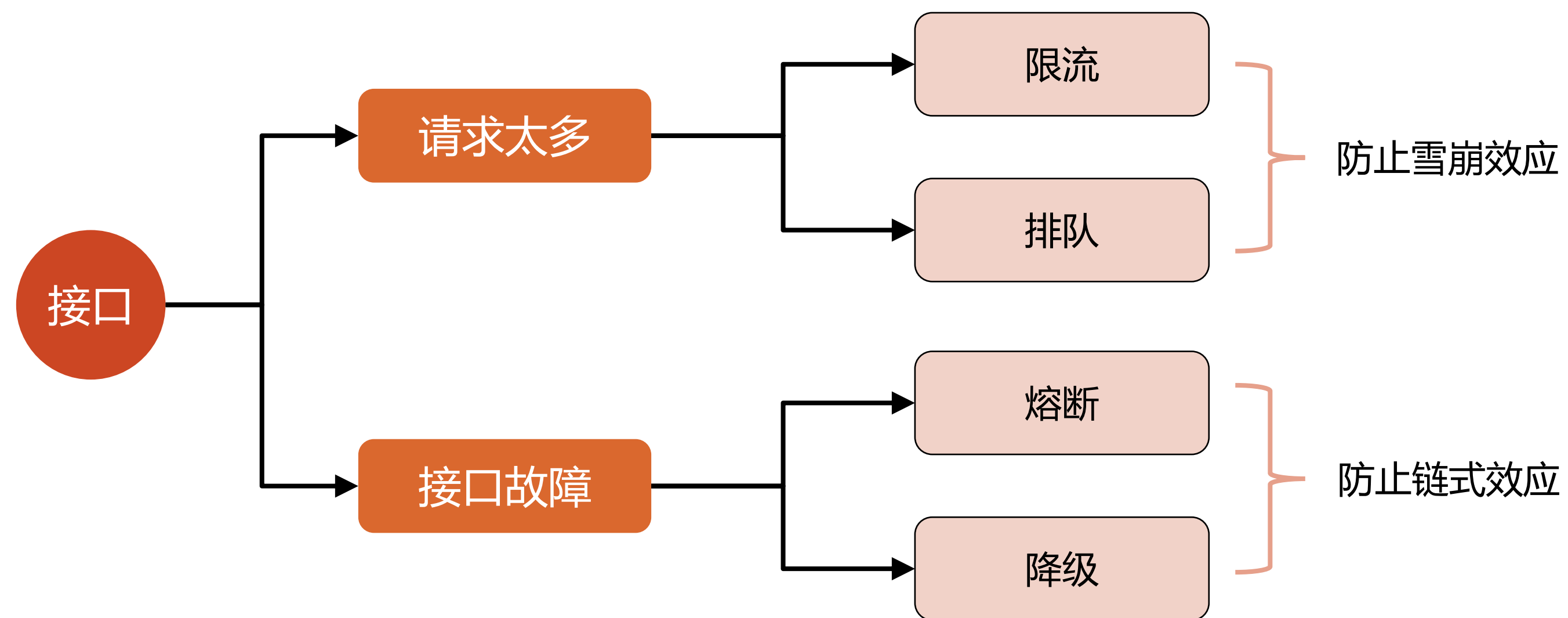
1. 接口高可用整体框架
2. 限流
3. 排队
4. 降级
5. 熔断

1. 接口高可用整体框架

接口高可用整体框架

雪崩效应：请求量超过系统处理能力后导致系统性能螺旋快速下降。

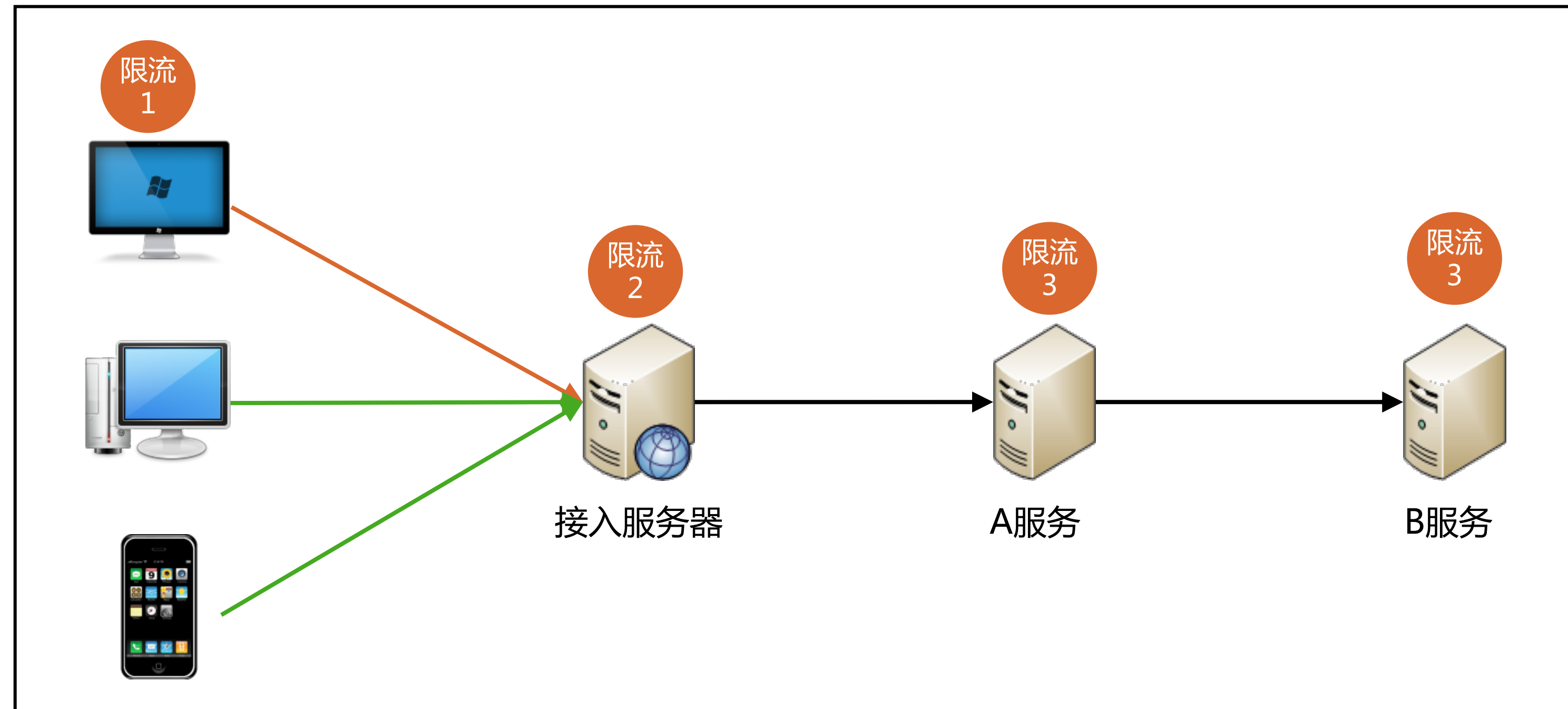
链式效应：某个故障引起后续一连串故障。



接口高可用架构本质上是“丢车保帅”策略，业务或者用户体验会部分有损！

2. 限流

限流



用户请求全流程各个环节都可以限流：

1. **请求端限流**：发起请求的时候就进行限流，被限流的请求实际上并没有发给后端服务器。
2. **接入端限流**：接到业务请求的时候进行限流，避免业务请求进入实际的业务处理流程。
3. **微服务限流**：单个服务的自我保护措施，处理能力不够的时候丢弃新的请求。

限流具体实现方式

请求端 限流

【常见手段】

1. 限制请求次数，例如按钮变灰）；
2. 嵌入简单业务逻辑，例如生成随机数。

【优缺点】

1. 实现简单；
2. 流量本地就控制住了；
3. 防君子不防小人（脚本）。

接入端 限流

【常见手段】

1. 限制同一用户请求频率；
2. 随机抛弃无状态请求，例如限流浏览请求，不限流下单请求。

【优缺点】

1. 实现复杂；
2. 可以防刷；
3. 限流阈值可能需要人工判断。

微服务 限流

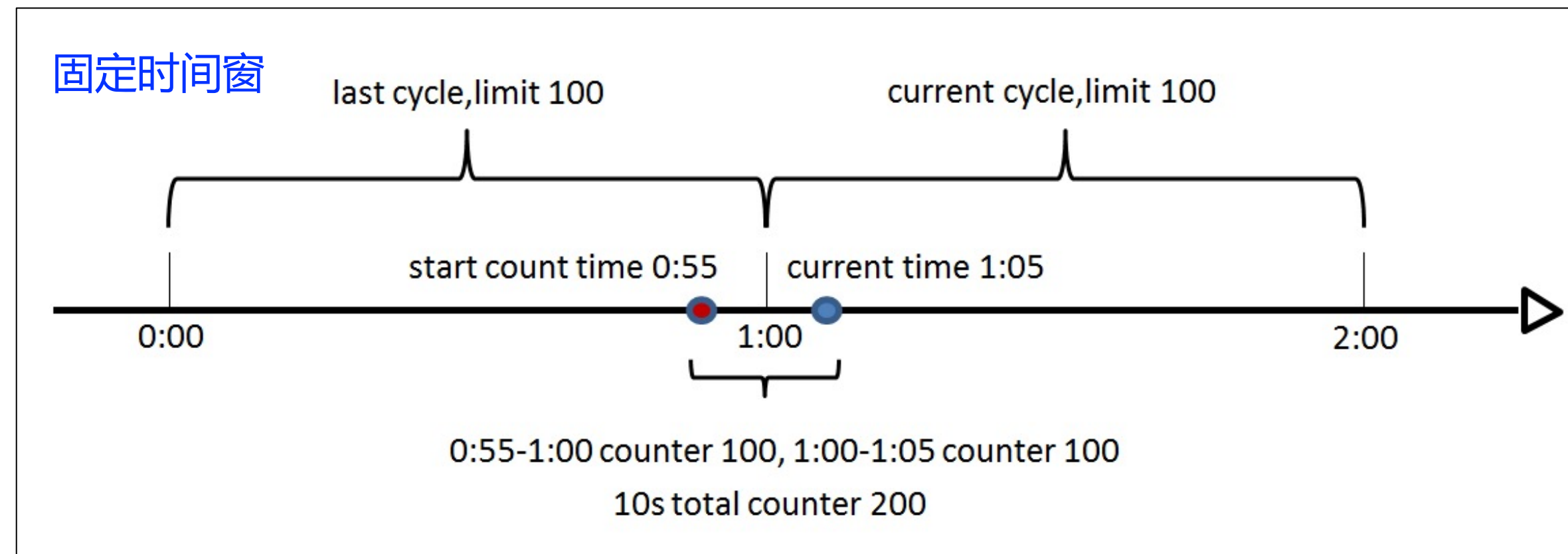
【常见手段】

根据处理能力，丢弃无法处理的请求。

【优缺点】

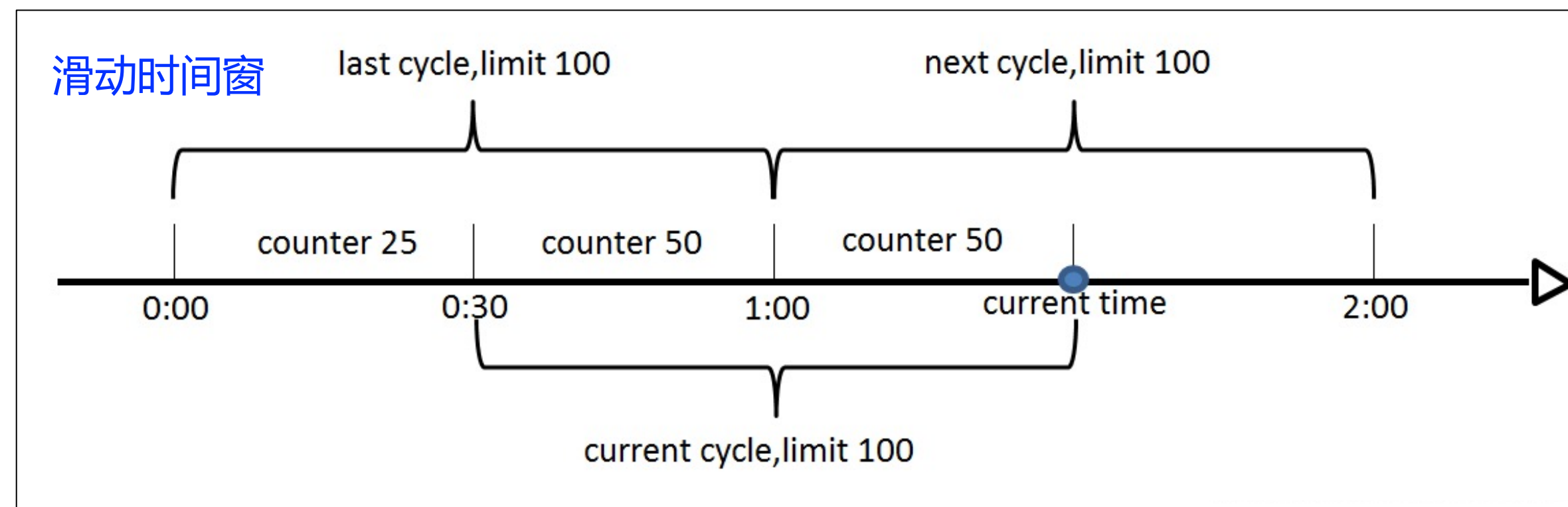
1. 实现简单；
2. 处理能力难以精准配置。

限流算法1 - 固定 & 滑动 时间窗



【设计原理】

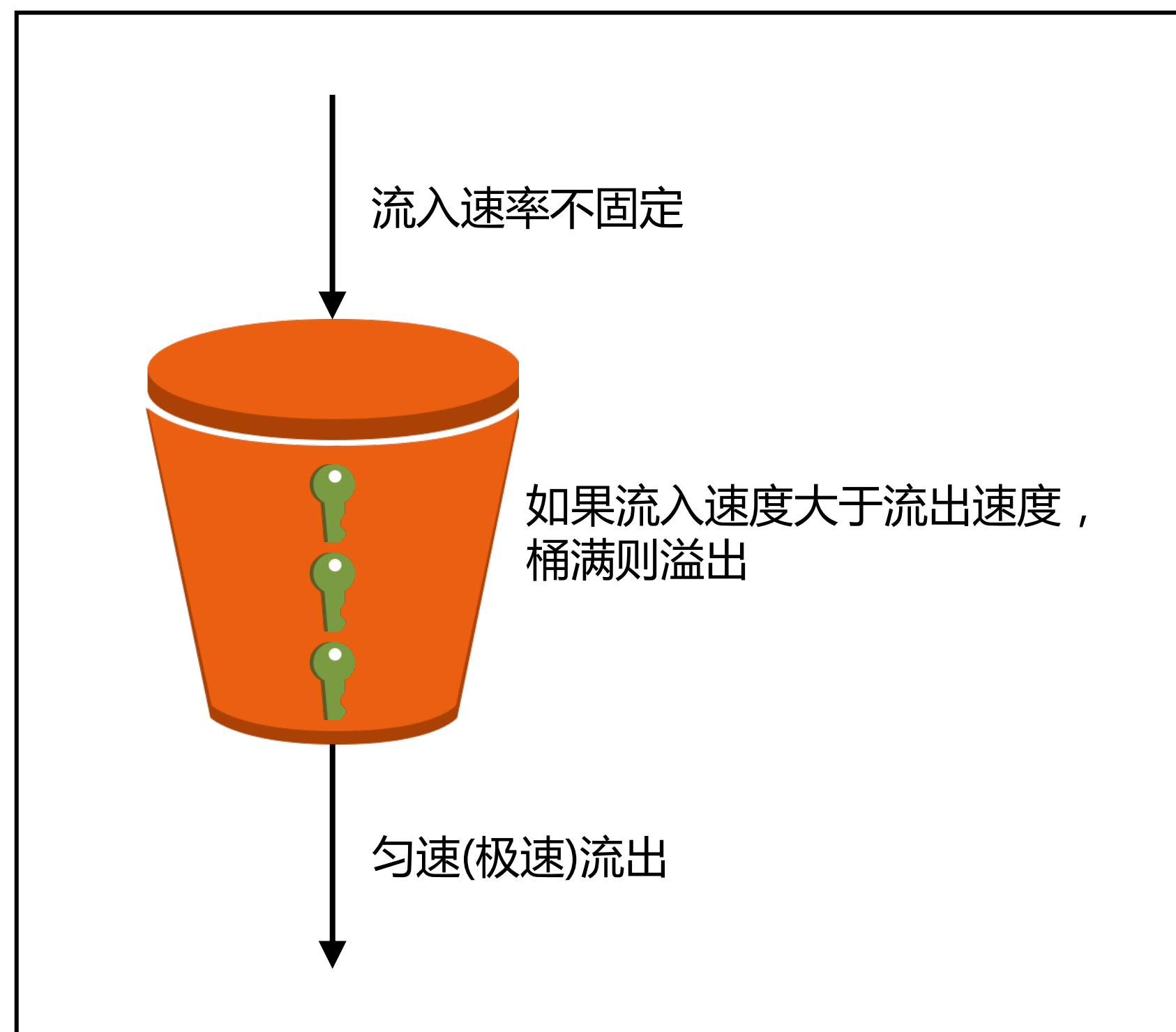
1. 统计**固定时间周期内**的请求量，超过阈值则限流；
2. 存在**临界点**问题，如图中的红蓝两点对应的时间范围。



【设计原理】

1. 统计**滑动时间周期内**的请求量，超过阈值则限流；
2. 判断比较准确，但**实现稍微复杂**。

限流算法2 - 漏桶



【基本原理】

请求放入“桶”（消息队列等），业务处理单元（线程/进程/服务）从桶里拿请求处理，桶满则丢弃新的请求。

【技术本质】

总量控制，桶大小是设计关键。

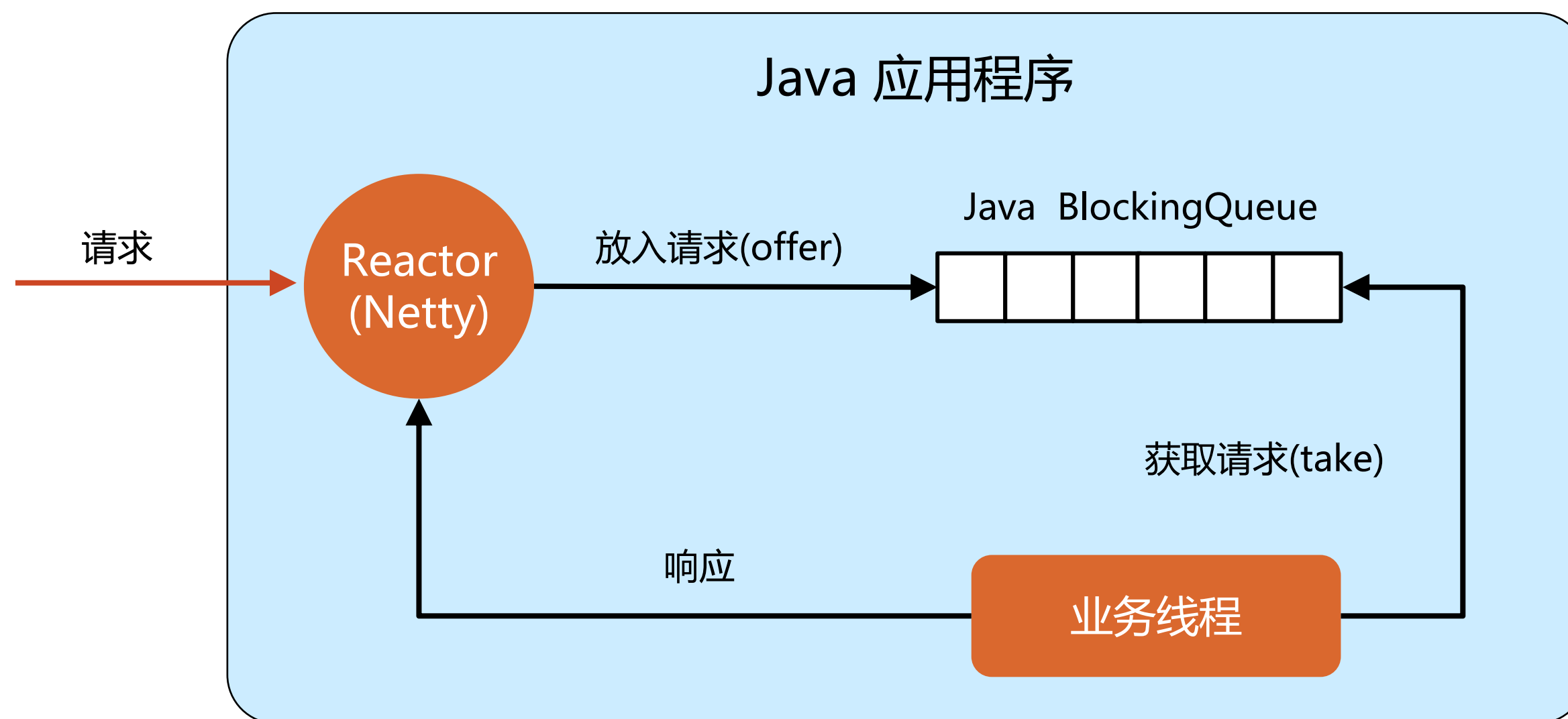
【优缺点】

1. 桶大小动态调整比较困难，例如 Java BlockingQueue；
2. 无法精确控制流出速度（处理速度）；
3. 突发流量时丢弃的请求较少。

【应用场景】

瞬时高并发流量，例如0点签到，整点秒杀。

Java 限流的漏桶算法简单示例



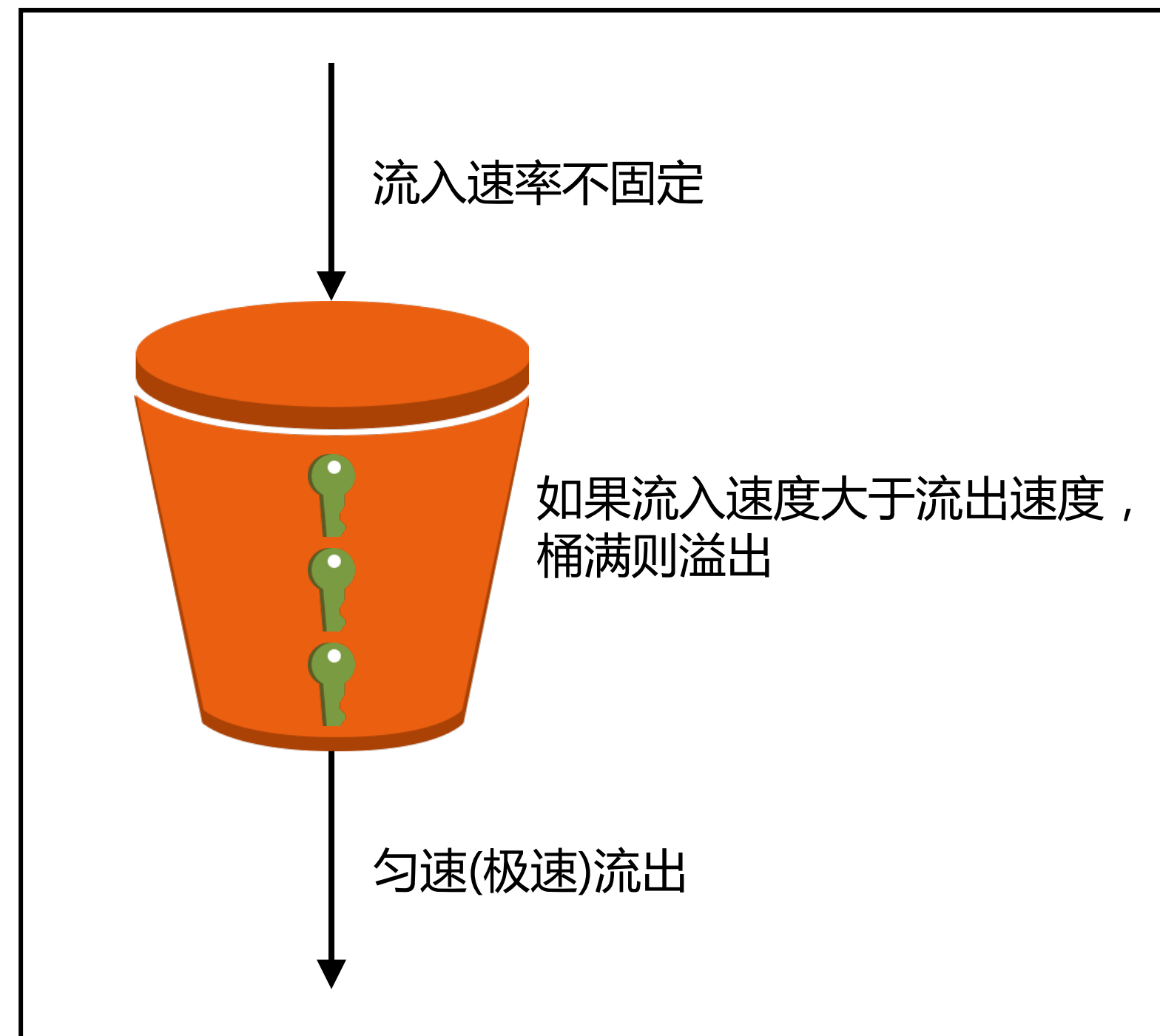
【设计关键】

1. 业务线程和 IO 线程分离，通过队列传递请求；
2. BlockingQueue 的长度配置，太长没作用，太短浪费。



如果是 Tomcat、SpringBoot 这类框架怎么办？

漏桶算法变种 - 写缓冲(Buffer)



【基本原理】

如果漏桶的容量无限（例如用 Kafka 消息队列），则漏桶可以用来做写缓冲。

【技术本质】

同步改异步，缓冲所有请求，慢慢处理。

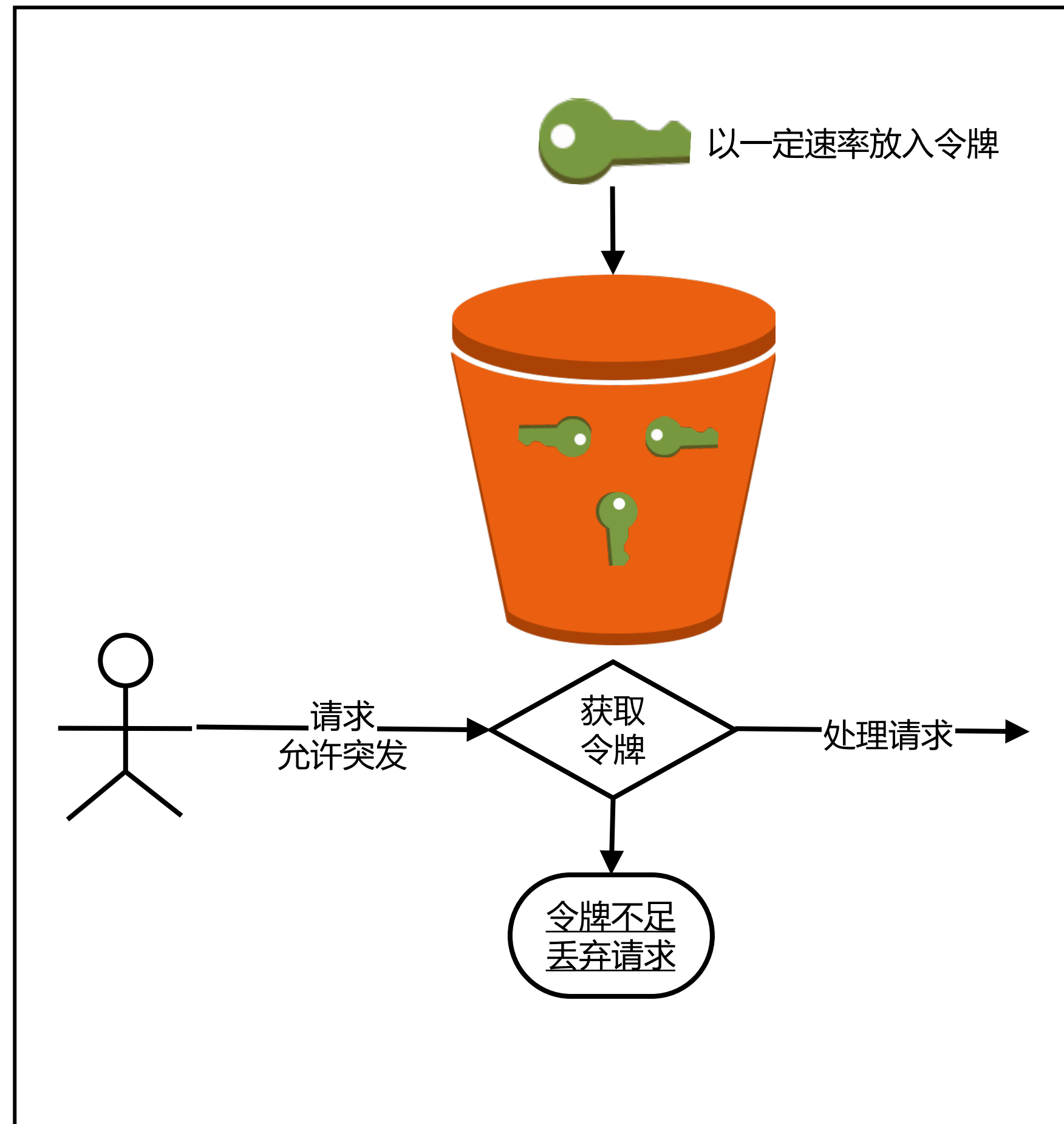
【应用场景】

高并发写入请求，例如热门微博评论。



为什么看微博的请求可以丢弃，而写评论请求却全部缓冲起来？

限流算法3 - 令牌桶



【基本原理】

某个处理单元按照**指定速率**将令牌放入“桶”（消息队列等），业务处理单元收到请求后需要获取令牌，获取不到就丢弃请求。

【技术本质】

速率控制，令牌产生的速度是设计关键。

【优缺点】

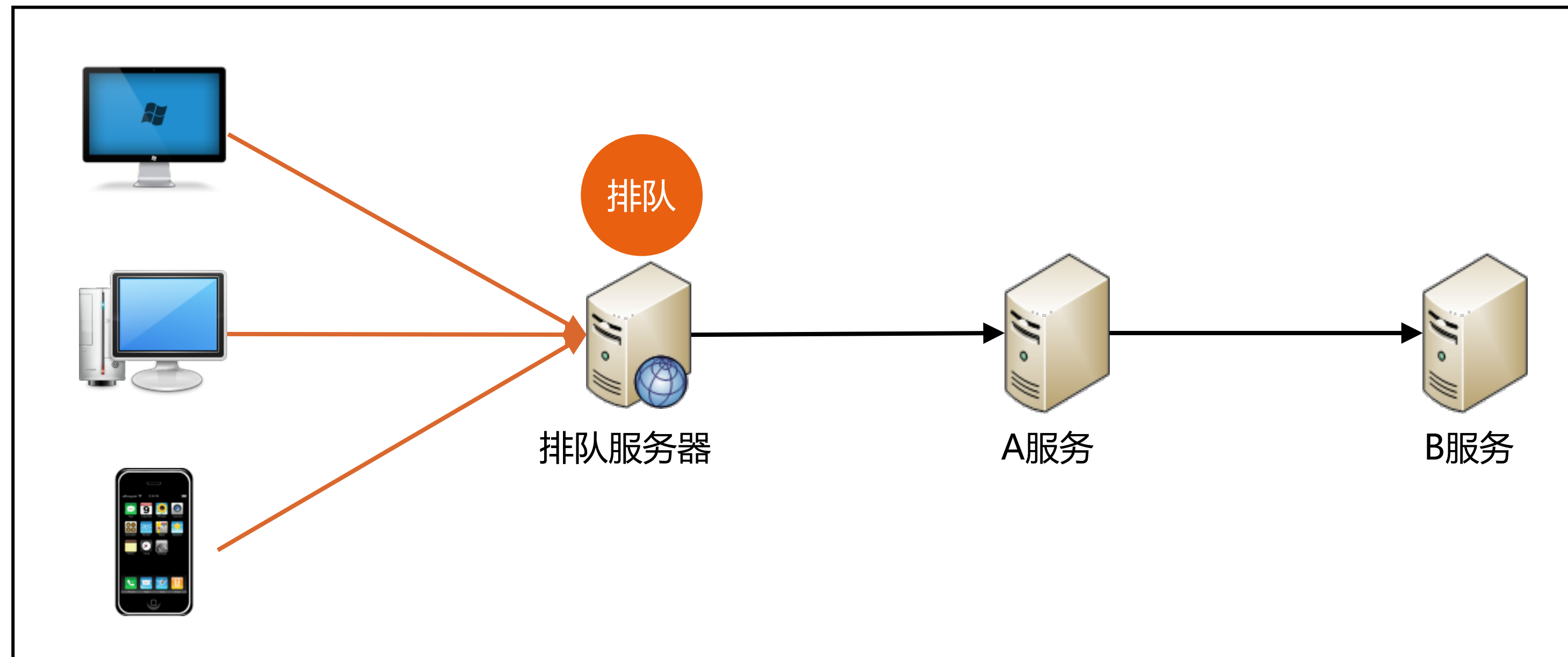
1. 可以**动态调整**处理速度；
2. 突发流量的时候可能丢弃很多请求；
3. 实现相对复杂。

【典型应用场景】

1. 控制访问第三方服务的速度，防止把下游压垮；
2. 控制自己的处理速度，防止过载

3. 排队

排队

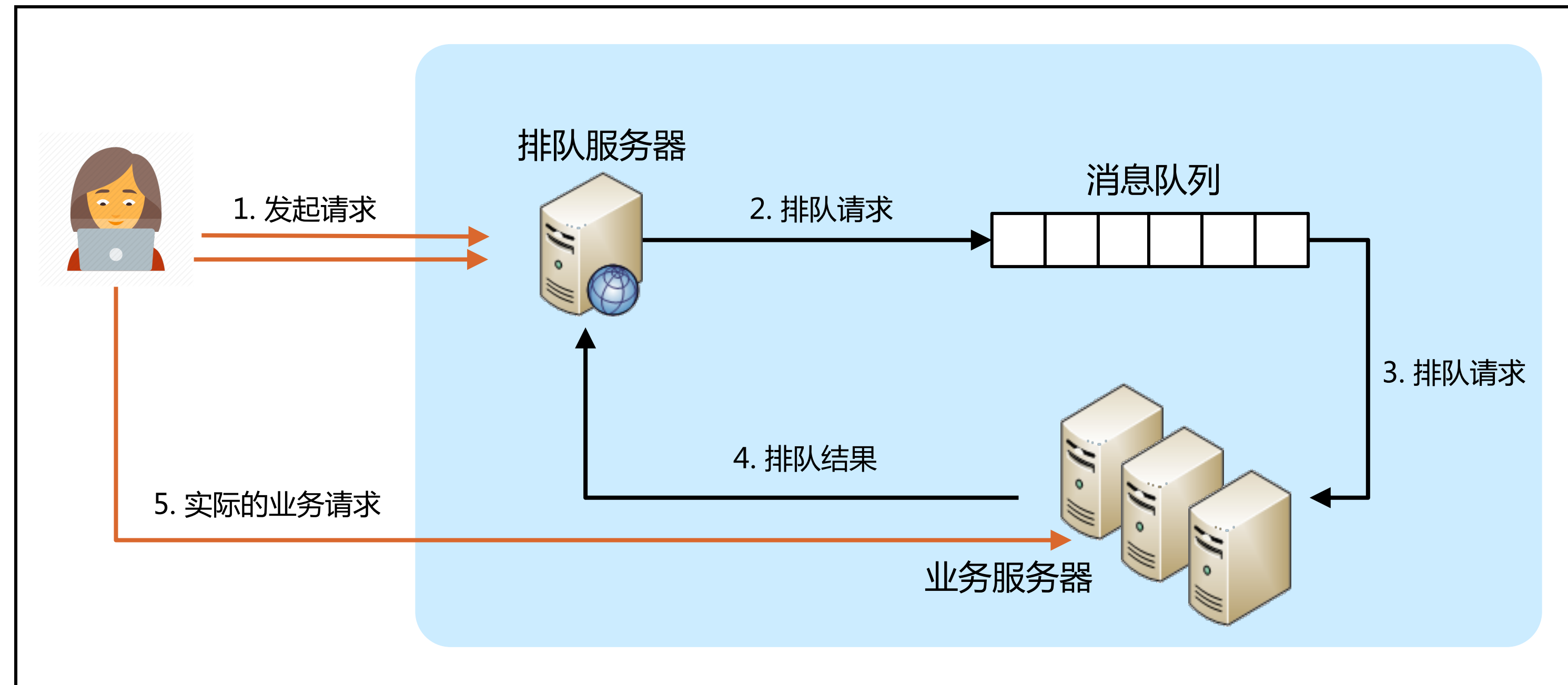


基本原理：收到请求后并不同步处理，而是将请求放入队列，系统根据能力异步处理。

技术本质：请求缓存 + 同步改异步 + 请求端轮询。

应用场景：秒杀、抢购。

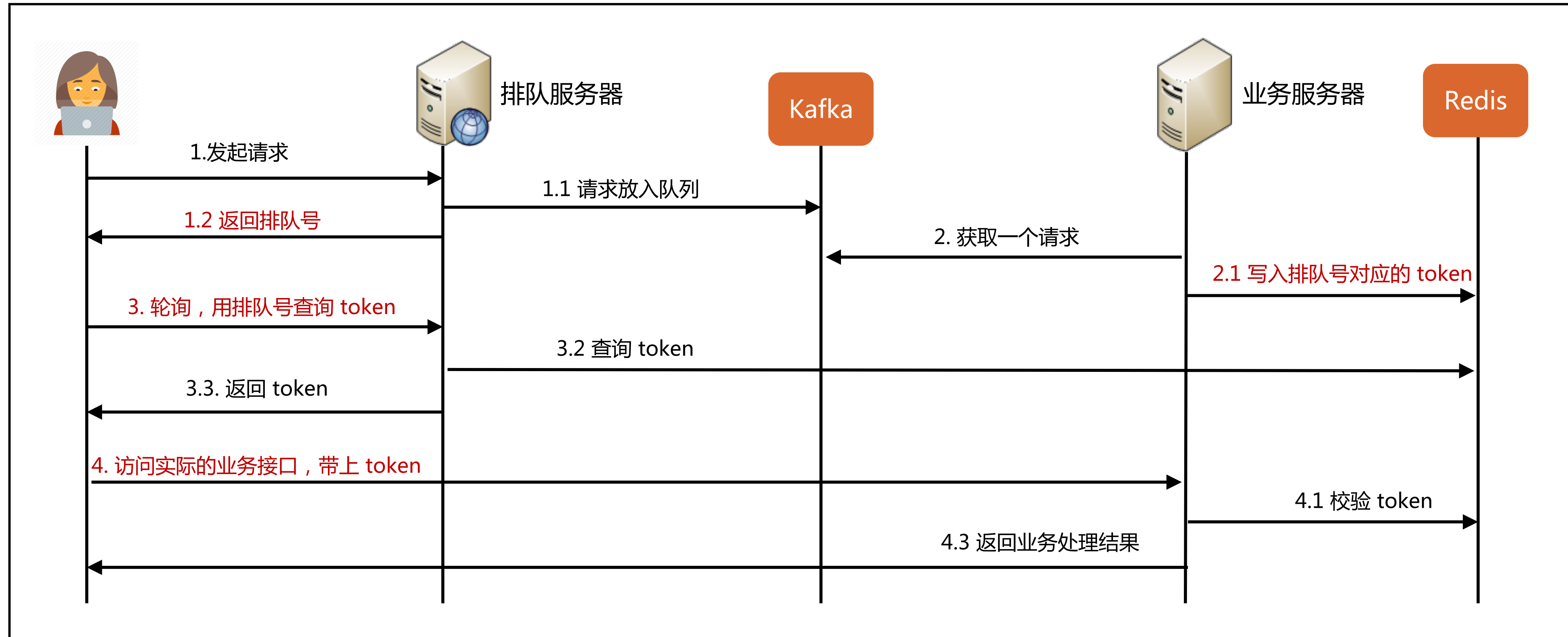
排队的架构示意图



【设计关键】

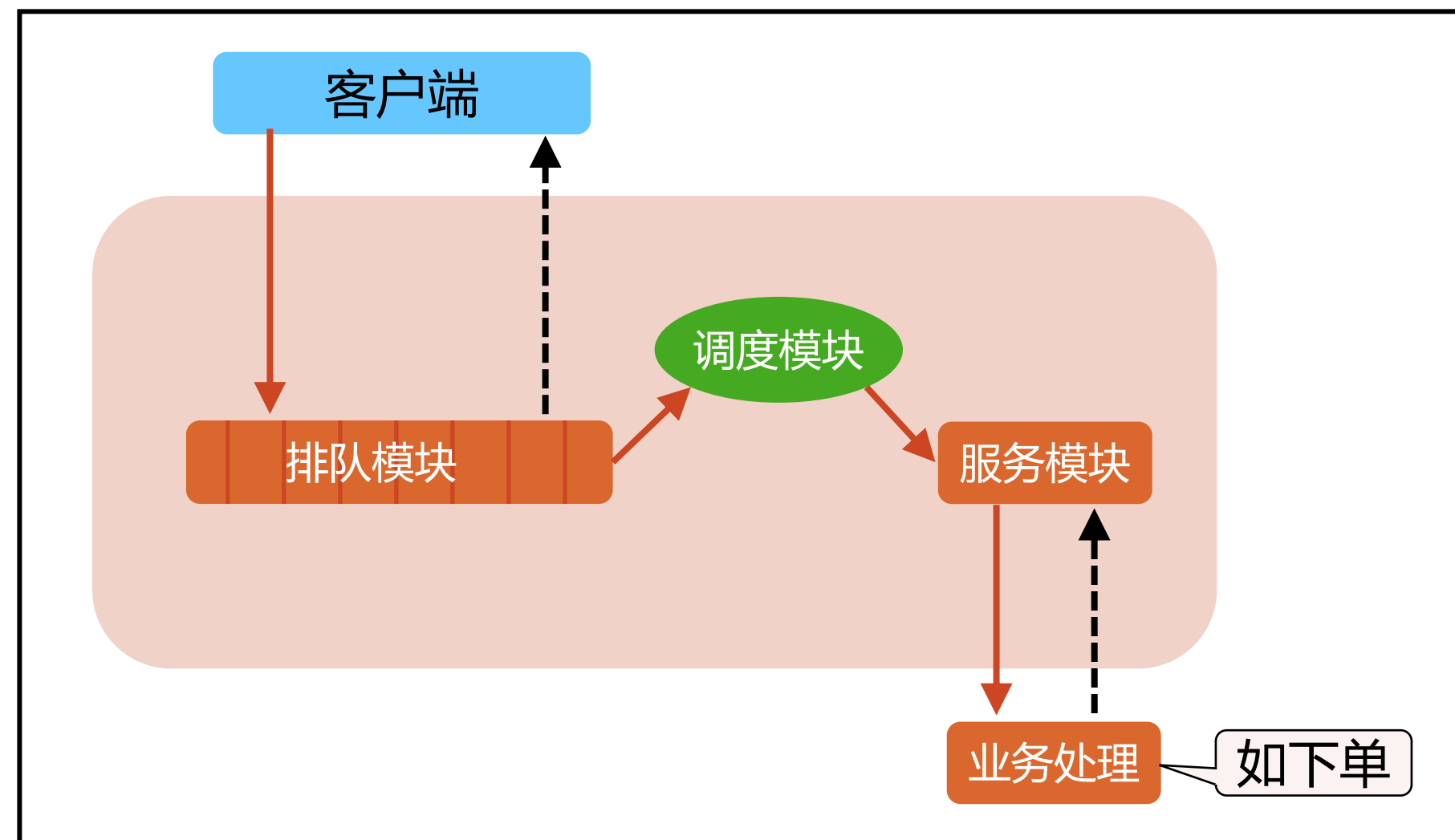
1. 如何设计异步处理流程；
2. 如何保证用户体验（前端、客户端交互）。

排队的具体实现方案示例



为什么要用 token，直接用排队号不就可以了么？

1号店双十一秒杀排队



排队模块：

负责接收用户的抢购请求，将请求以先入先出的方式保存下来。每一个参加秒杀活动的商品保存一个队列，队列的大小可以根据参与秒杀的商品数量（或加点余量）自行定义。

调度模块：

负责排队模块到服务模块的动态调度，不断检查服务模块，一旦处理能力有空闲，就从排队队列头上把用户访问请求调入服务模块。

服务模块：

是负责调用真正业务处理服务，并返回处理结果，并调用排队模块的接口回写业务处理结果。 [参考链接](#)

抢购的人可真多
排在您前面的人超过 **500** 位
别急，我们的速度比火箭快~

剩余时间
约2分钟



抢购成功，订单正在处理
中。。。

即将跳转至支付页面，优惠马上领回家~

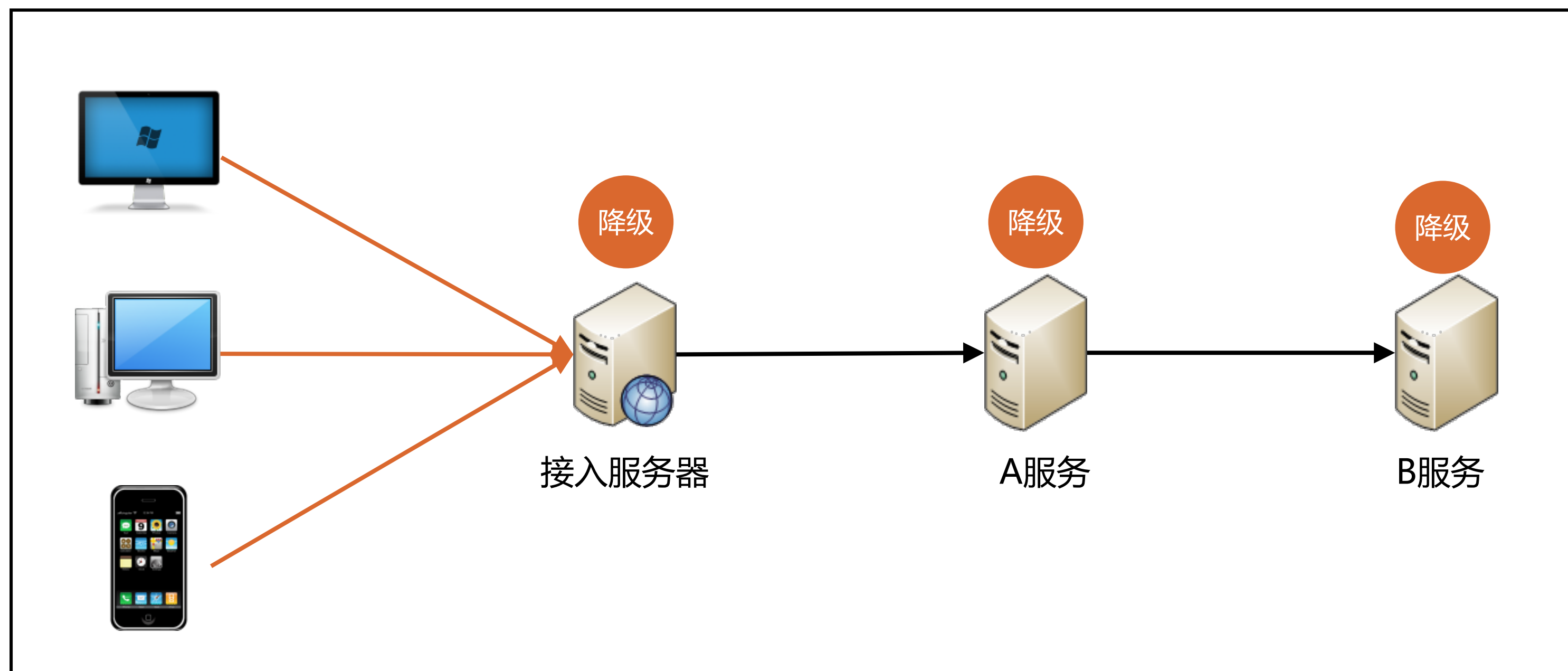


抢购成功



4. 降级

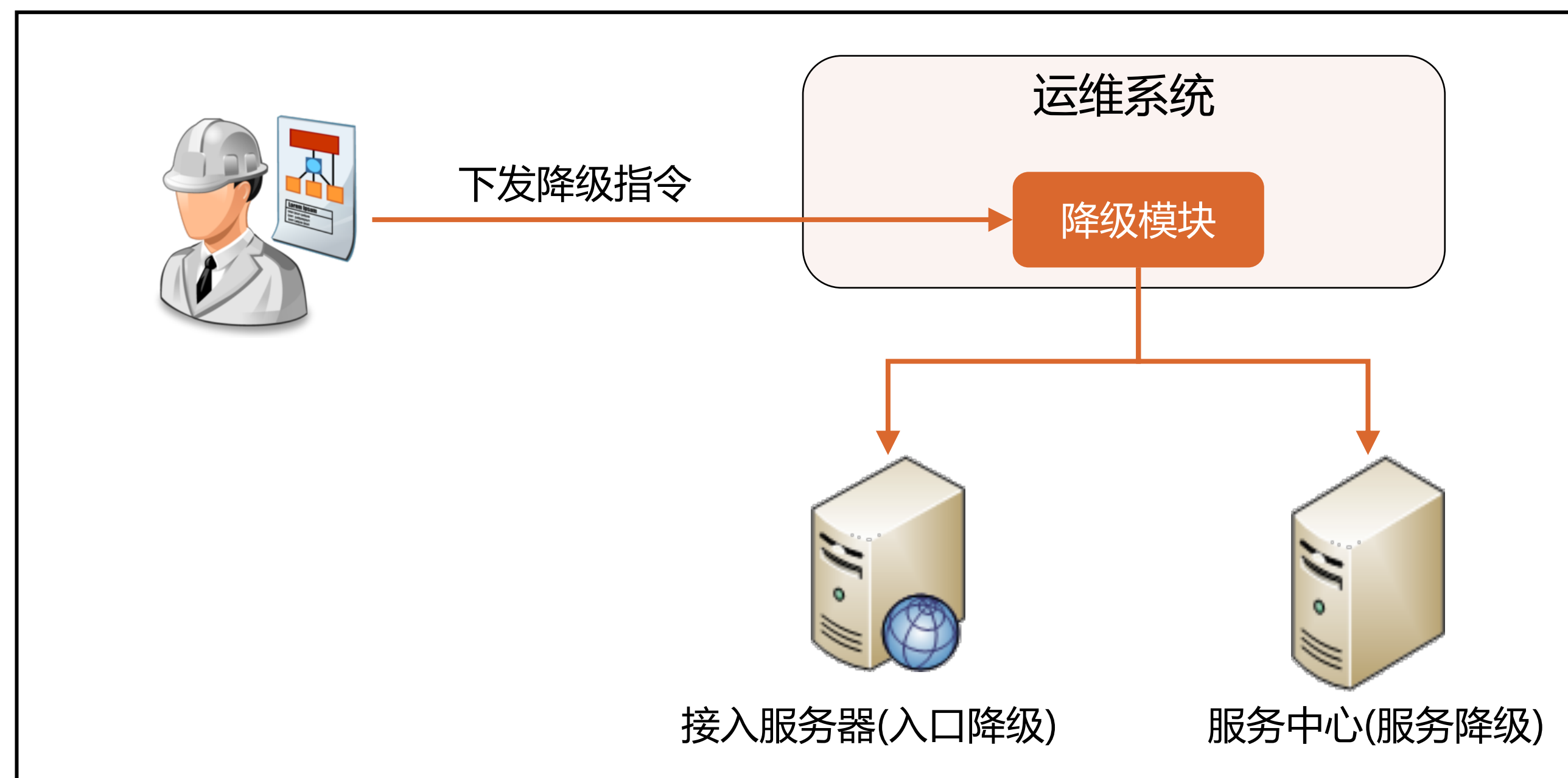
降级



基本原理：直接停用某个接口或者 URL，收到请求后直接返回错误（例如 HTTP 503）。

应用场景：故障应急，通常将**非核心业务降级**，保住核心业务，例如降级日志服务、升级服务等。

降级架构实现

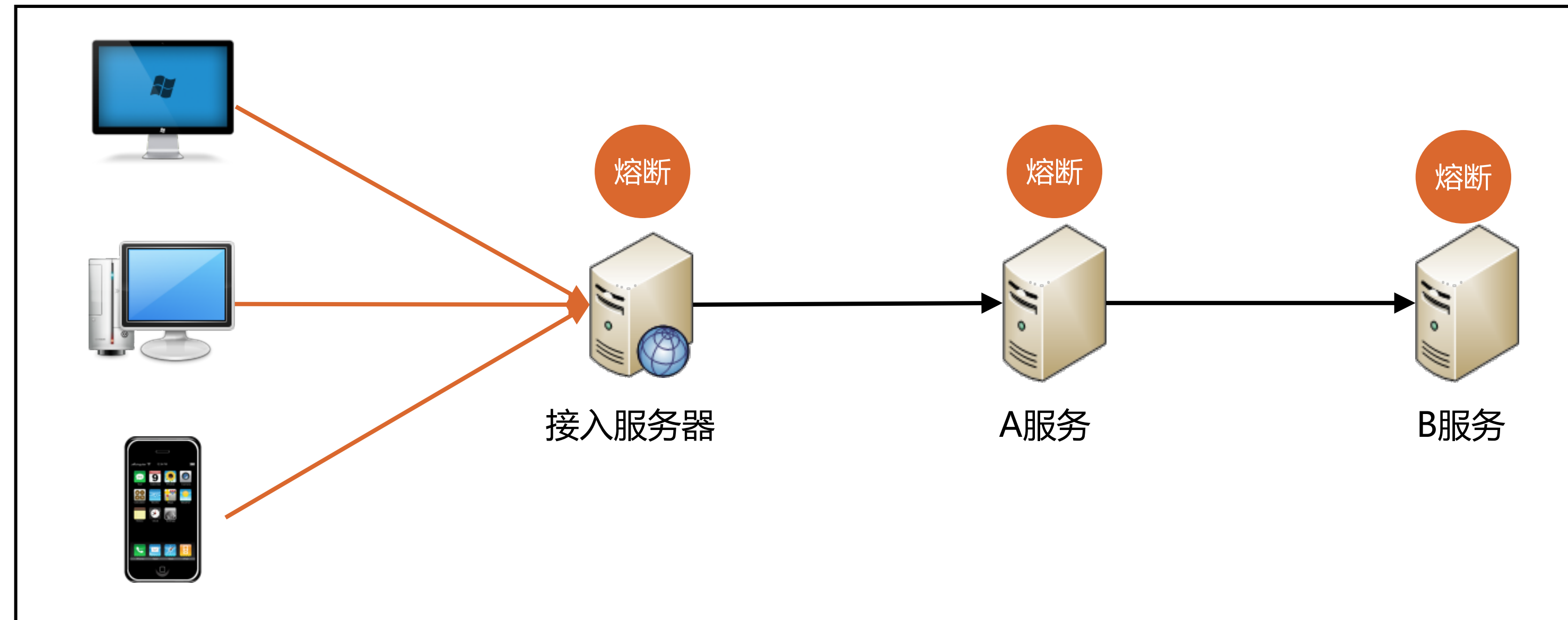


设计要点：

1. 独立系统操作降级，可以是独立的降级系统，也可以是嵌入到其它系统的降级功能。
2. 人工判断，人工执行，不要信 AIOps 之类的噱头。

5. 熔断

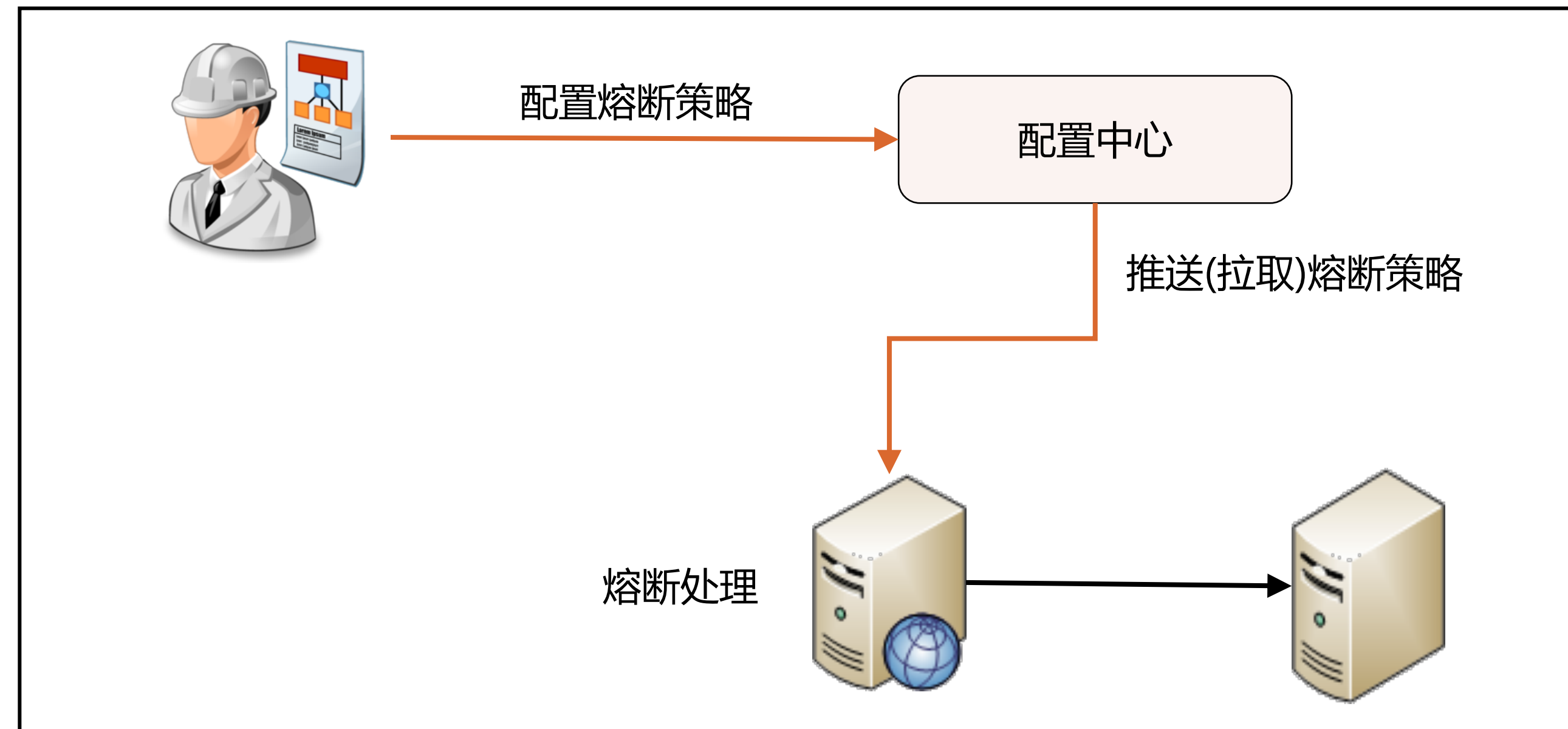
熔断



基本原理：下游接口故障的时候，一定时期内不再调用。

应用场景：服务自我保护，防止故障链式效应。

熔断架构实现

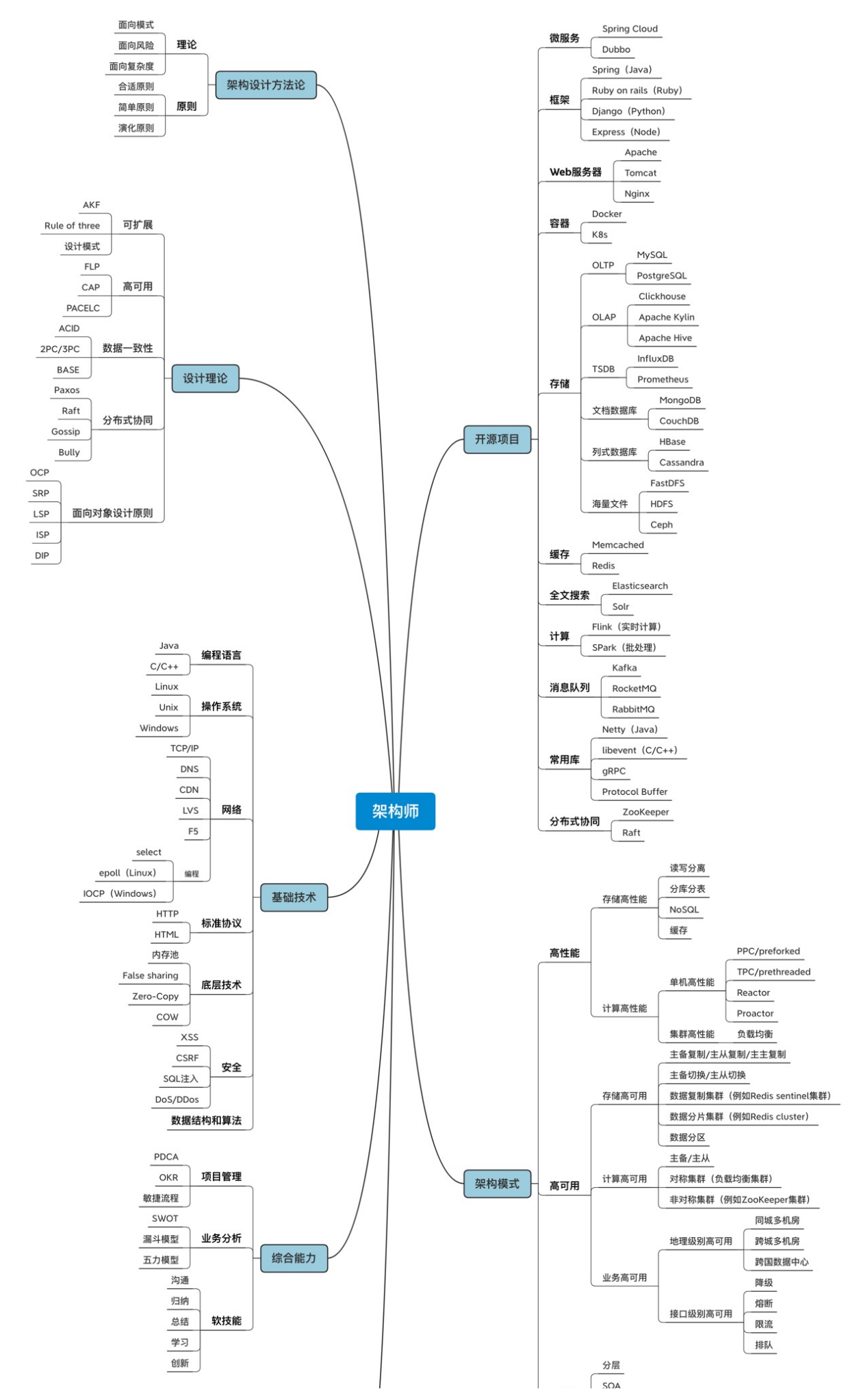


【实现细节】

1. 可以通过配置中心，也可以通过配置文件来配置熔断策略；
2. 熔断处理一般由框架或者 SDK 提供，例如 [Dubbo + Hystrix](#)；
3. 熔断策略一般按照失败次数、失败比例、响应时长等来确定。

Q&A

大厂面经	拼多多面经.pdf
技术栈面试题	快手面经.pdf
Nginx面试题.pdf	腾讯面经.pdf
Spring面试题.pdf	美团面经.pdf
MongoDB面试题.pdf	蚂蚁金服面经.pdf
Spring Cloud面试题.pdf	Memcached面试题.pdf
Elasticsearch面试题.pdf	Docker面试题.pdf
MySQL面试题.pdf	Redis面试题.pdf
Netty 面试题.pdf	ClickHouse 面试题.pdf
Dubbo面试题.pdf	RocketMQ面试题.pdf



直播三重福利免费领

200 道架构面试题 + 架构师技能图谱 + 本次直播分享 PPT



扫码领取

你有过类似的想法吗？

1. 做架构设计离我太远了，现在应该还用不到
2. 只有做到架构师才要学习架构设计
3. 只有写代码才是技术提升，架构设计太空了
4. 架构要学的东西太多了，不知道怎么学
5. 学架构一定要自己设计并研发一个系统，要跟随大厂的架构设计的方法来



- 以阿里职级为例，什么级别开始要求架构设计能力？

实战营主要内容



THANKS

 极客时间 | 训练营