



京东虚拟业务系统高可用性设计

京东商城研发部架构师



成为软件技术专家 从 / / / D 全球软件开发大会的必经之路

[北京站] 2018

2018年4月20-22日 北京·国际会议中心

一片三购票中,每张立减2040元

团购享受更多优惠



识别二维码了解更多





下载极客时间App 获取有声IT新闻、技术产品专栏,每日更新



扫一扫下载极客时间App



AICON

全球人工智能与机器学习技术大会

助力人工智能落地

2018.1.13 - 1.14 北京国际会议中心



扫描关注大会官网

TABLE OF

CONTENTS 大纲

- ▶虚拟业务系统
- 》虚拟业务系统高可用性实践
- 〉大促如何确保高可用

虚拟业务系统特点

- 非标品
- 多业务线(20+产品线)
- 》业务复杂程度不一
- 大量第三方供应商交互
- 新业务上线迅速
- 系统维护量大



虚拟业务系统架构v1.0

- 系统抽象度不高
- 重复开发,效率低
- 系统稳定性难保证







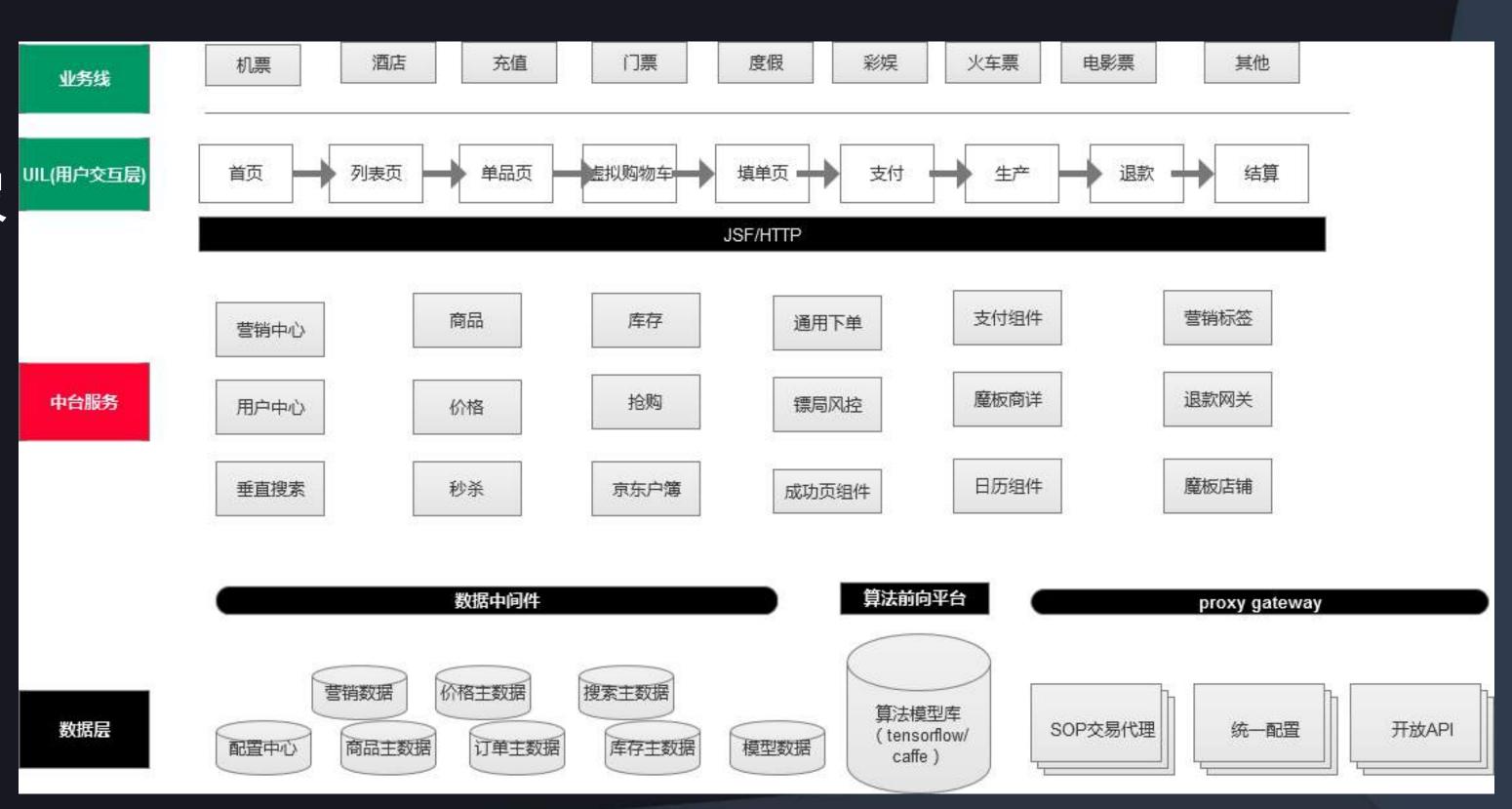






虚拟业务系统架构v2.0

- 基础公共能力下沉
 - 商品、订单等模型抽象
 - 》业务组件提炼
- 业务流程平台化
 - 编排
 - **配置**





部署架构

- 多实例多机房部署
- 同机房不同机柜
- ES集群
- > MySQL—主多从
- 》缓存(JimDB)一主多从

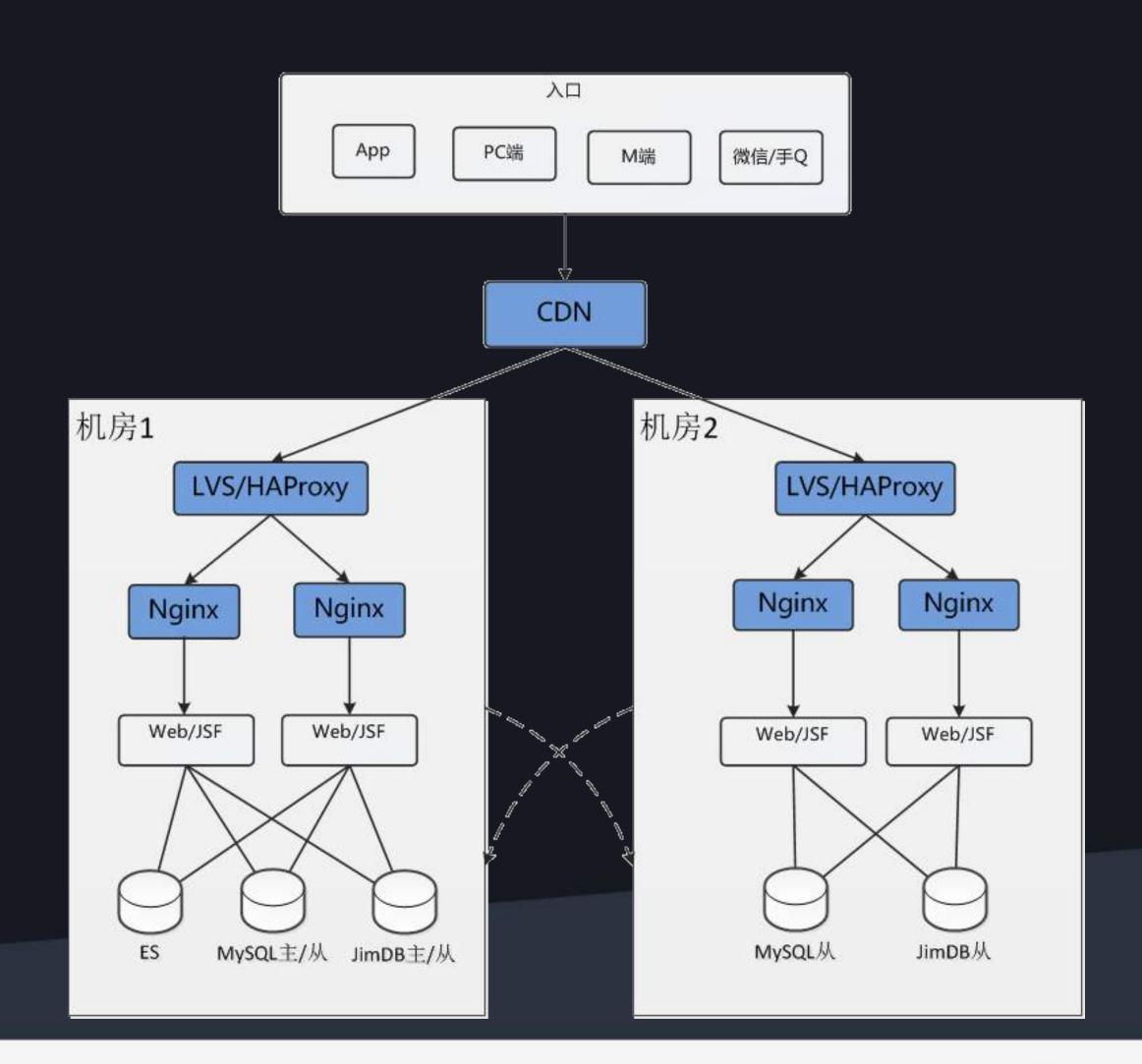




TABLE OF

CONTENTS 大纲

- ▶虚拟业务系统
- > 虚拟业务系统高可用性实践
- 〉大促如何确保高可用



高可用性的度量

A = MTBF / (MTBF + MTTR)

MTBF = 平均故障间隔时间(Mean Time Between Failures)

MTTR = 平均恢复时间 (Mean Time to Recovery)

提高系统的可用性的两个途径

• 提高MTBF:减少故障

· 降低MTTR: 快速解决故障



影响可用性的因素与应对方法

	方法	结果
单点故障	冗余	提高MTBF
依赖故障	降级、异步化	提高MTBF
超出承载	限流	提高MTBF
定位、恢复慢	监控报警、应急预案	降低MTTR



冗余设计一消除单点

基础设施层

- 多IDC
- 多DNS
- 多CDN
- 多路由器
- 双网卡

数据层

- 数据库master/ slave
- 缓存一主多从
- ES集群

应用层

- 多实例
- 跨机架
- 服务无状态

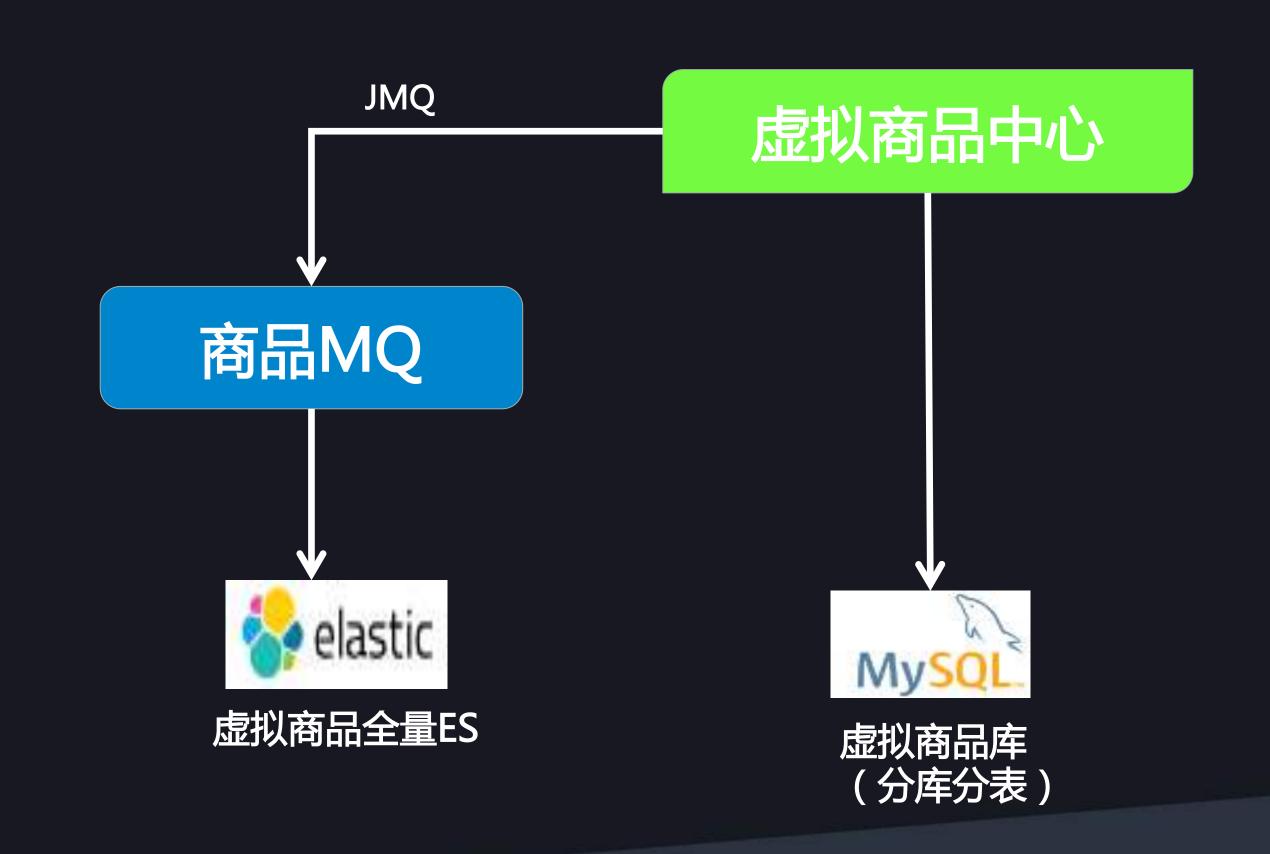
数据冗余实现一复制

应用同步双写

- 一可立刻读
- 一性能差

应用异步双写

- 十性能好
- 有延迟



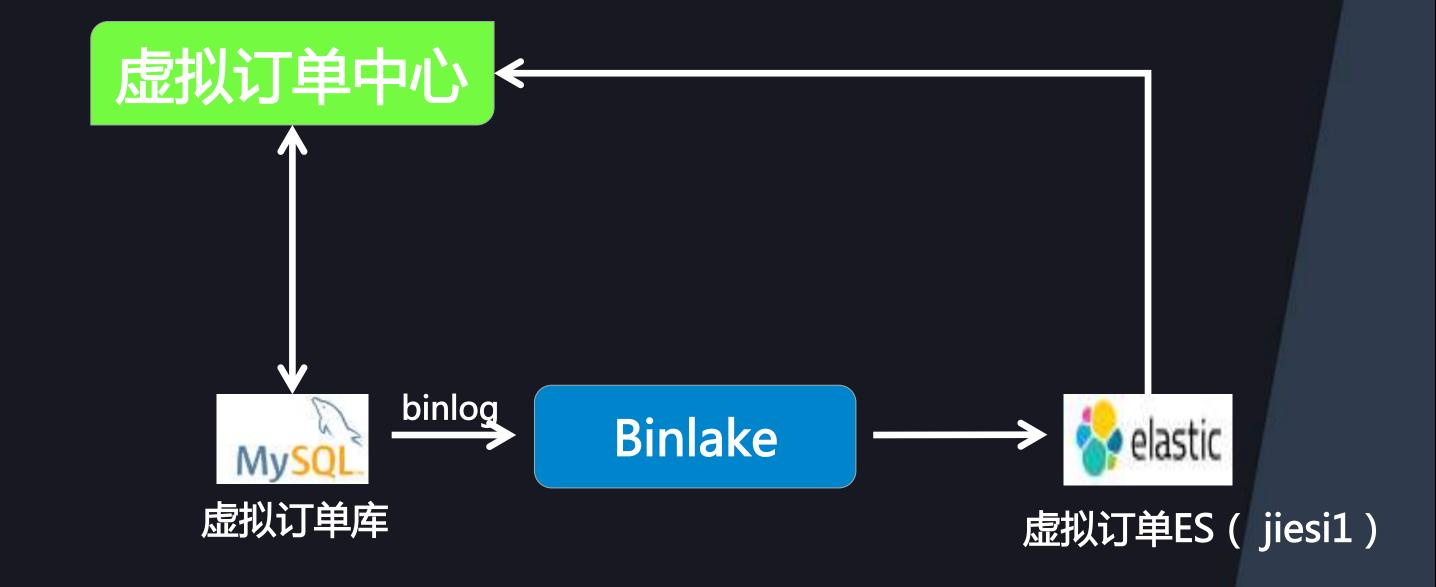
数据冗余实现一复制

利用底层数据存储复制机制

+ 对应用无侵入

开发量少

一延迟较严重





数据本地化

应用场景

- 〉依赖外部数据
- ⇒高频访问
- 》数据变化频率低

实现方式

- 定时主动数据同步
- 第三方通知





降级

场景

依赖的外部系统出现故障

原则

KISS

降级非关键依赖

方法

开关(手动/自动)

方案审查、演练

例子

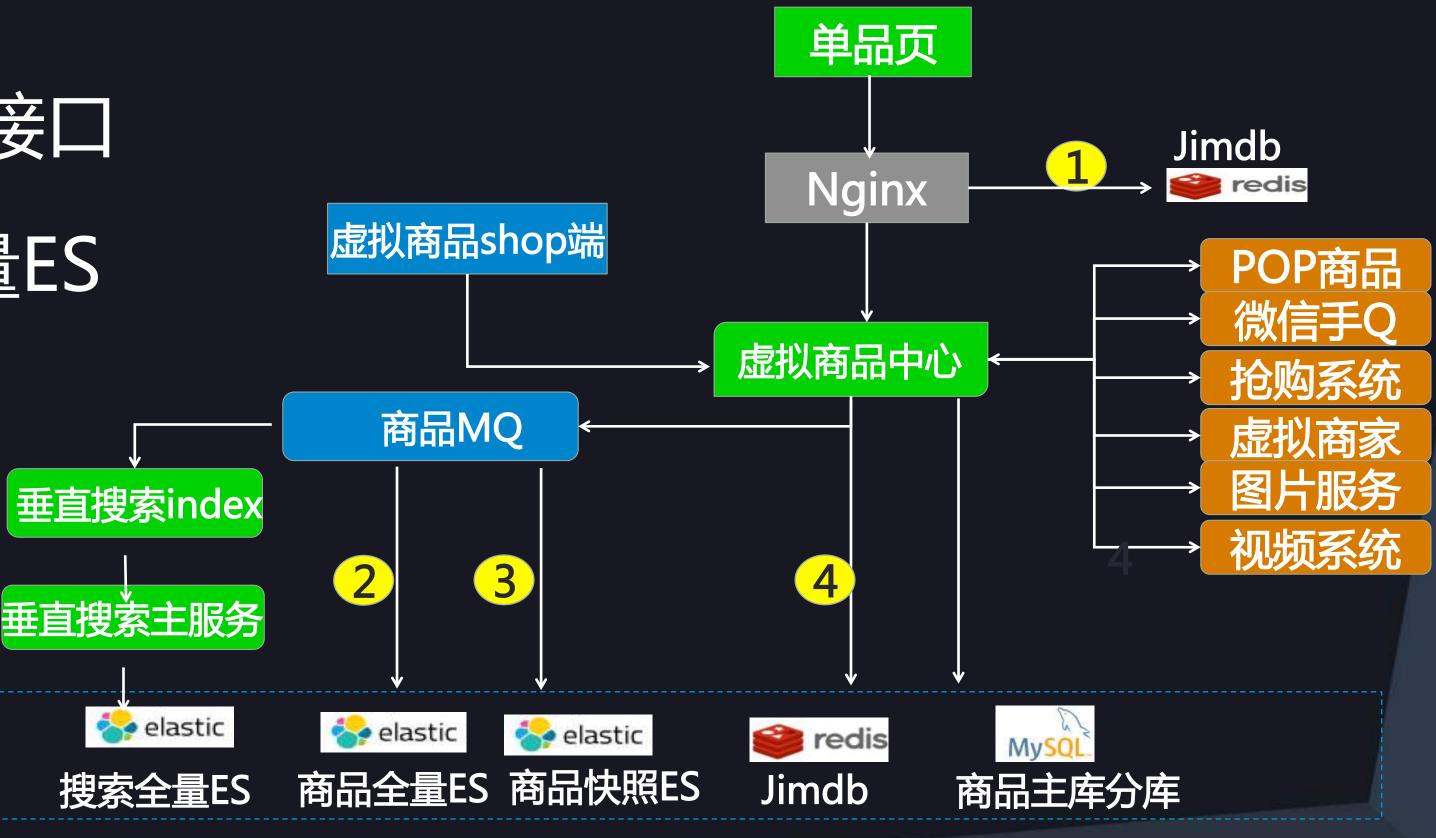
性能降级:比如缓存降级到ES

功能降级:比如优惠券服务故障降级支付方式



降级一虚拟商品系统

- 1 Lua自动降级调用HTTP接口
- 2 开关控制是否实时写全量ES
- 3 开关控制是否写快照
- 4 开关控制是否读Jimdb





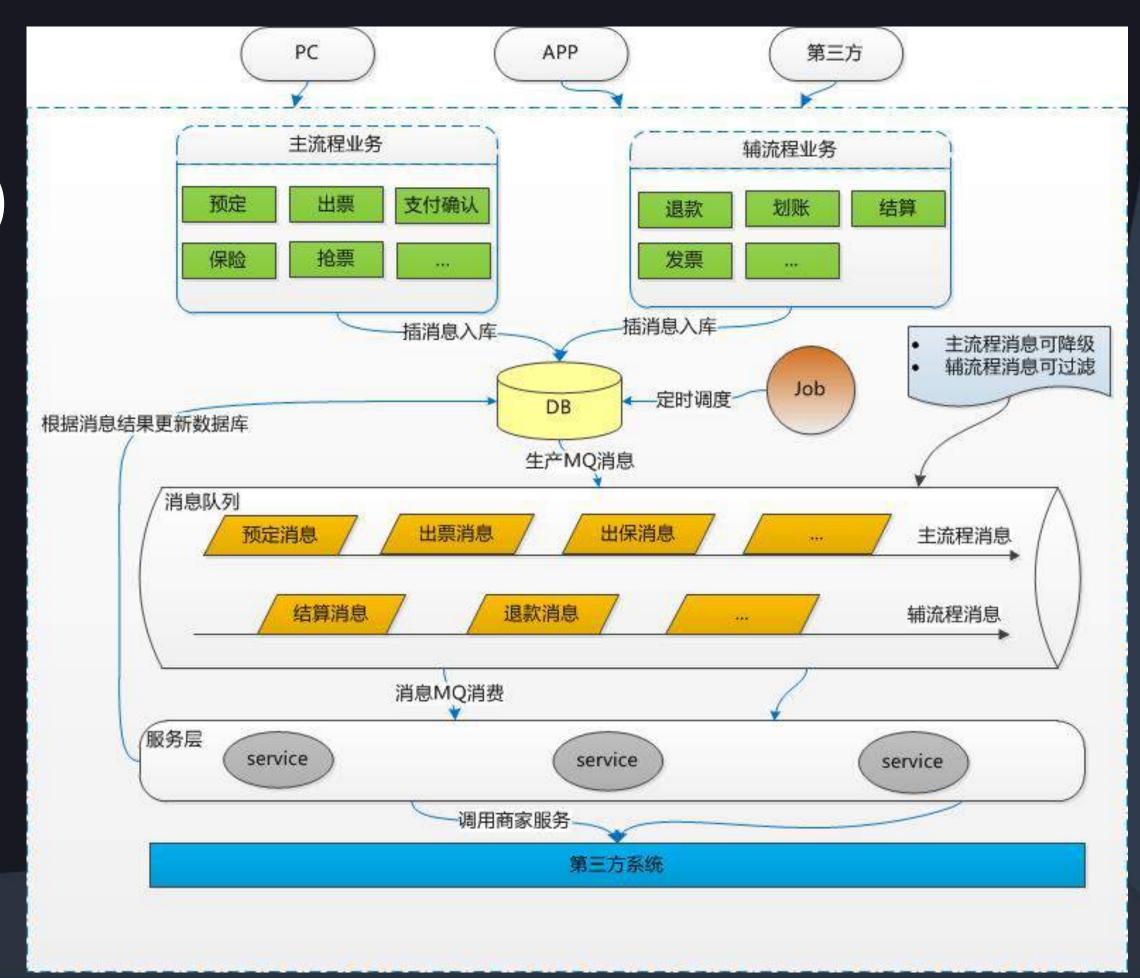
异步化-事件驱动(火车票)

典型场景

>访问第三方系统不稳定(网络等)

实现方式

- 事件驱动业务流程
- >失败利用MQ重试机制
- 注意幂等





限流

典型场景

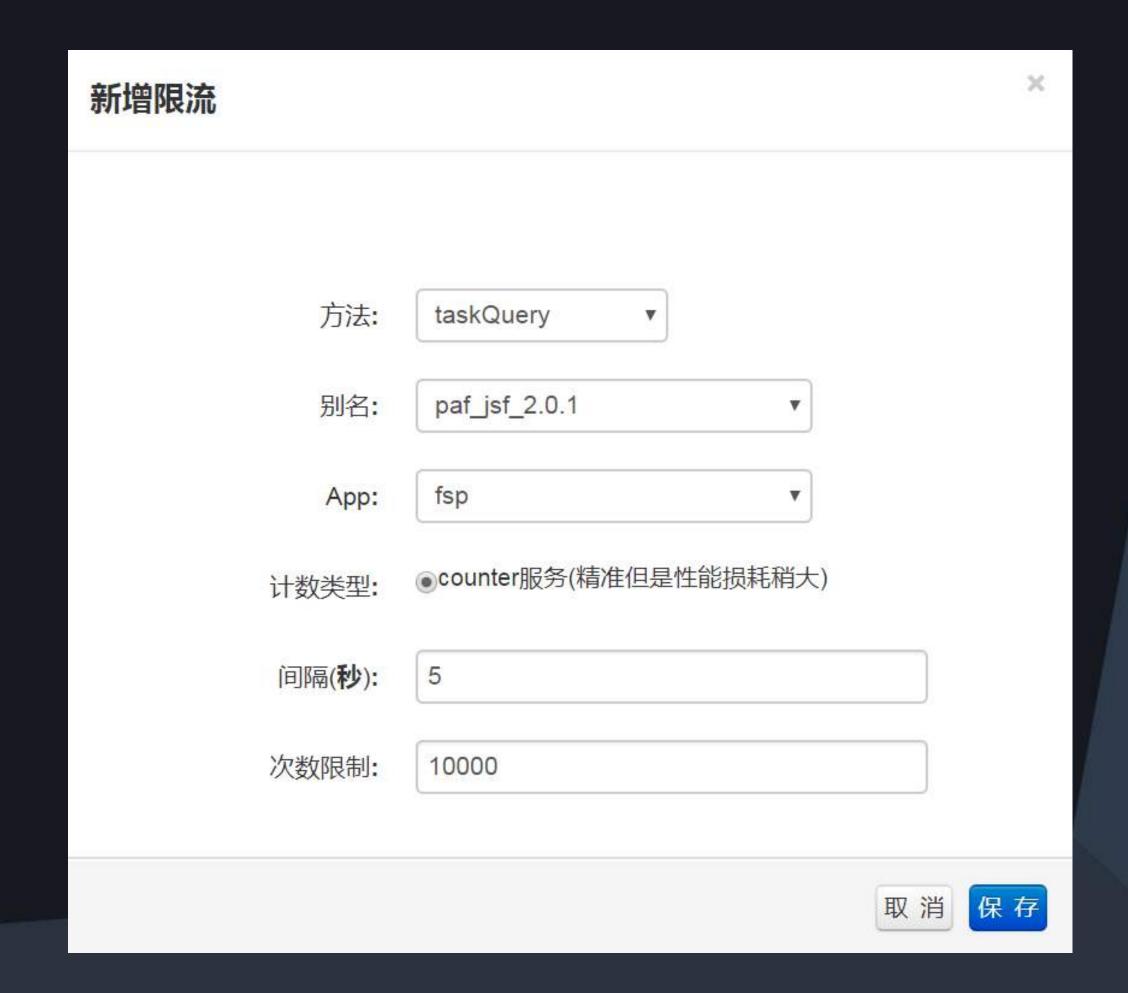
> 机票防刷;抢火车票

接入层限流

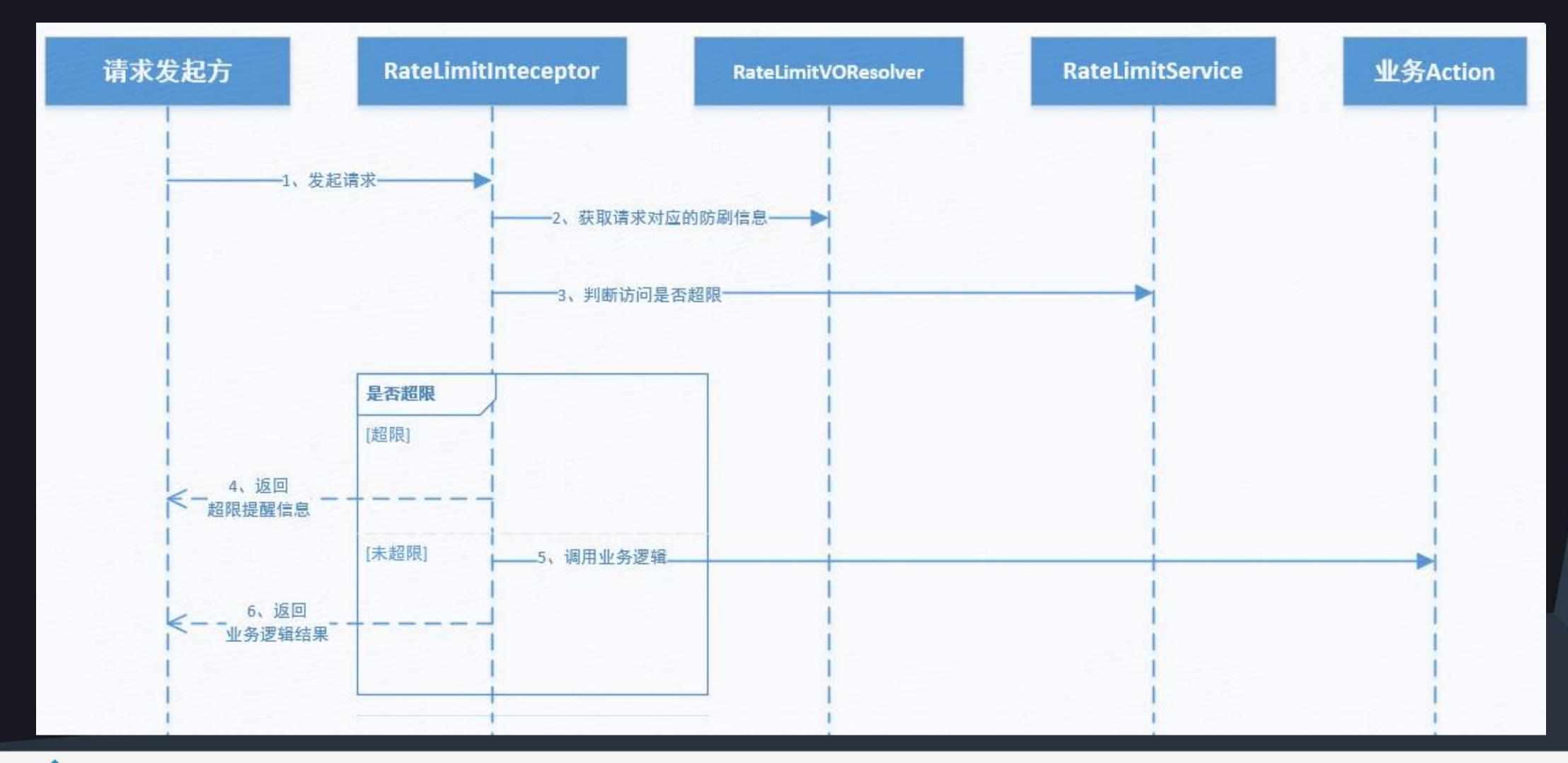
- Nginx+Lua+Redis
- > IP、账户、地区等

应用限流

> 服务框架设置(方法级)



限流





监控与报警

- > 系统监控(CPU,内存、负载等)
- > 网络监控(网络流入量,流出量等
- > 磁盘监控(使用率、读写速度等)
- > 容器监控 (线程数、SWAP使用等
- » 服务监控(TP99,调用量等)
- > 业务监控(订单状态等)







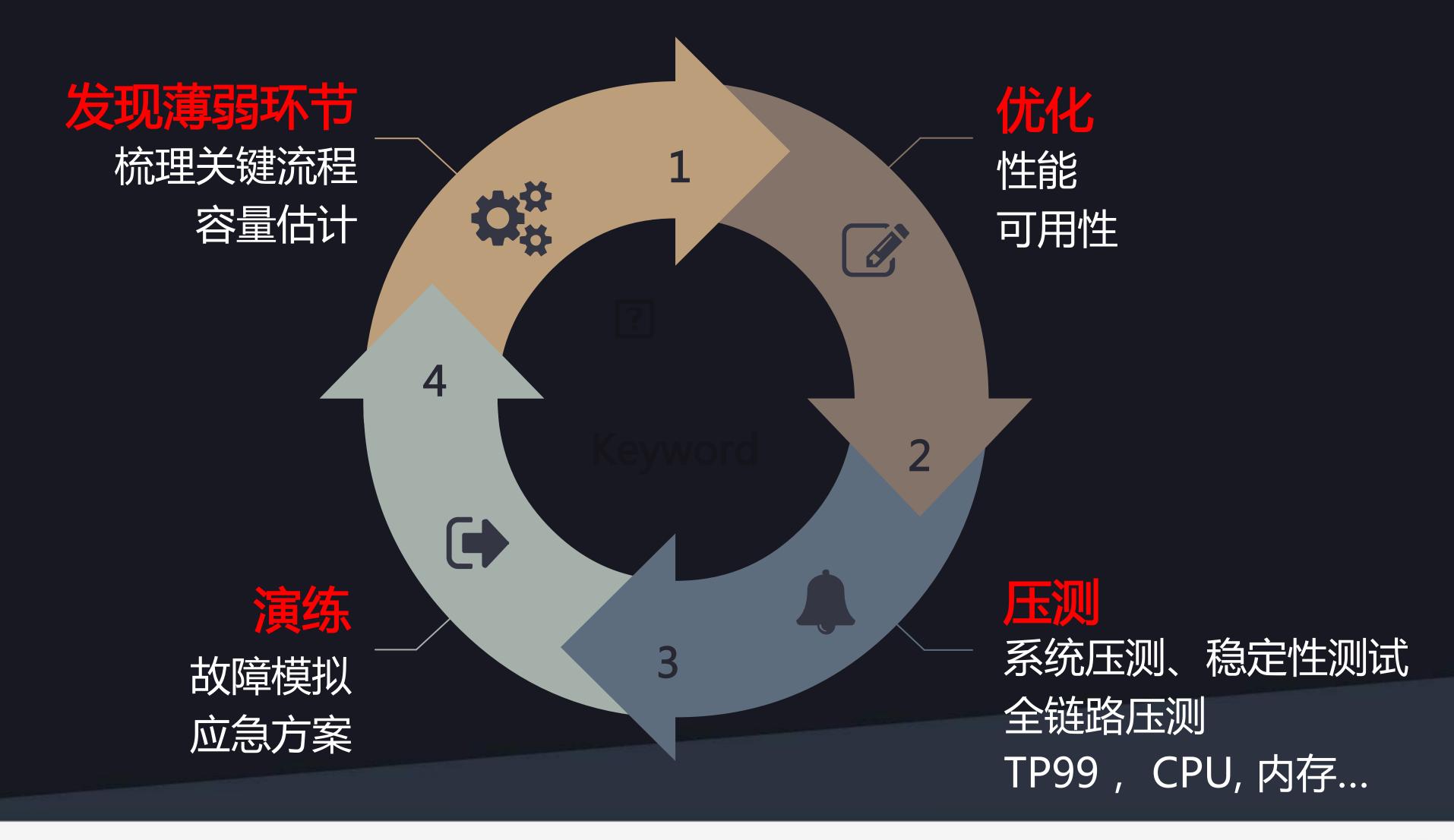


TABLE OF

CONTENTS 大纲

- ▶虚拟业务系统
- 》虚拟业务系统高可用性实践
- > 大促如何确保高可用

大促备战





系统压测

- 》测试环境/预发环境
- 关注吞吐量、性能和稳定性
- 常态化



场景编号₽	系统名称□	方法₽	服务器	线程数₽	Avg (ms)	tp99₽	实际 TPS (笔/秒)₽	PV(笔/ 分)₽	错误率(%)
20171127_01	com.jd.jmi.bean.client.service.JmiBeanTranQueryService	queryBeanTranDetails.	単机₽	20₽	14₽	32 ₽	1 <mark>4</mark> 94 ₽	89640₽	0.11‰

场景编号₽	keyName <i>⊷</i>	并发用户数₽	IP₽	CPU (%)₽	MEM (%) ₽	网络流入 (MB/s) ₽	网络流出 (MB/s)₽	load(%)	Тср₽	UMP 平均响应时间。	UMPTP99 响应时间₽	UMP 可用率₽
20171127_01	jmi_bean_soa_jsf_tran_detail_query •	20₽	2	16.53₽	70.18₽	2.62₽	3.46₽	0 ₽	490₽	1 3₽	29₽	99.85





全链路压测

- 在公网模拟真实用户行为
- > 从CDN节点压测
- > 使用Forcebot执行压测
- ▶ 监控访问量、TP99、CPU等



容量规划

规划

- > 基于历史数据、促销计划预测
- > 根据压测分配资源

扩容方法

- ➤ 应用水平扩容(增加docker)
- > 硬件垂直扩容(内存、升级SSD等)
- > 基础设施扩容(数据库分库分表、ES分片扩容、缓存分片扩容等)

项目	机	房1	机	房2			2016.1			017.1
	实例数	服务能力(万/ 秒)	实例数	服务能力(万/ 秒)	实例数	服务能 力(万/ 秒)	1.11流 量(万/ 秒))	6.18流 量(万/ 秒)	计	11预 流量 万/秒)
XXXX	24	1.6	24	1.6	12	0.8	1.6	2		3

全场景故障模拟演练

- > 开发人员自助服务
- ▶故障模拟种类
 - 系统故障
 - 网络故障
 - 》服务故障
- 可视化效果数据

演练记录													
任务名称	任务ID	创建用户	演练选择方式	选择对象	任务总数	已完成	成功	异常	失败	创建时间	执行状态	任务状态	操作
tomcat进程关闭	332152486		应用		4	4	4	0	0	2017年10月31日 17:34	执行完成	正在进行	Q查看详情
CPU_100%消耗(触发全部cpu核数)	575561449		应用		1	1	1	0	0	2017年10月31日 17:17	执行完成	正在进行	Q查看详情
磁盘打满	877921284		应用		1	1	1	0	0	2017年10月31日 17:14	执行完成	正在进行	Q查看详情
JSF随机provider调用异常故障	713108904		应用		3	3	3	0	0	2017年10月30日 21:08	执行完成	正在进行	Q查看详情
限制访问JFS	983108570		应用		2	2	2	0	0	2017年10月30日 21:02	执行完成	正在进行	Q查看详情
CPU_100%消耗(触发全部cpu核数)	161108268		应用		2	2	2	0	0	2017年10月30日 20:57	执行完成	正在进行	Q查看详情
tomcat进程关闭	256667964		应用		2	2	2	0	0	2017年10月30日 20:52	执行完成	正在进行	Q查斯斯
tomcat进程关闭	578136886		机房		1	1	1	0	0	2017年10月30日 17:48	执行完成	正在进行	Q查看详情



应急预案

> 数量: 虚拟业务200+

> 类型:应用/依赖

演练:与故障模拟结合

京东集团-CMO体系-商城研发部-虚拟平台-			
1 基本信息 类型:依赖服务预案 所属系统: 系统负责人: 相关联系人:	状态:有效 所属应用: 应用负责人:		生效时间:60s
依赖系统&应用: 1. 所属系统: 余额系统 所属应用: 余额服务 依赖接口: 接口地址:		系统负责人: 应用负责人:	
2 详细信息 应急场景:下单页余额无法使用 启动条件:1.通过ump报警方查询余额接口方法可用率低于80% 2.tps 启动影响:1.业务线无法使用余额虚拟资产 执行步骤:1.通知产品,联系业务人员关闭余额支付方式。2.联系余额 其 他: 附 件:		H.	





THANKYOU

如有需求,欢迎至[讲师交流会议室]与我们的讲师进一步交流

