

1.1 实时监控

1.2 簇点链路

1.3 流控规则

限流阈值类型

QPS

并发线程数

流控模式

快速失败

Warm Up（激增流量）

匀速排队（脉冲流量）

降级规则

慢调用比例

异常比例

异常数

热点参数限流

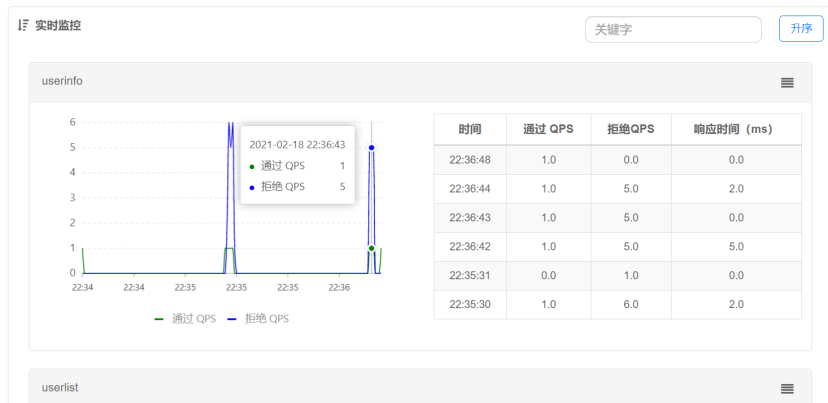
系统规则

授权控制规则

集群规则

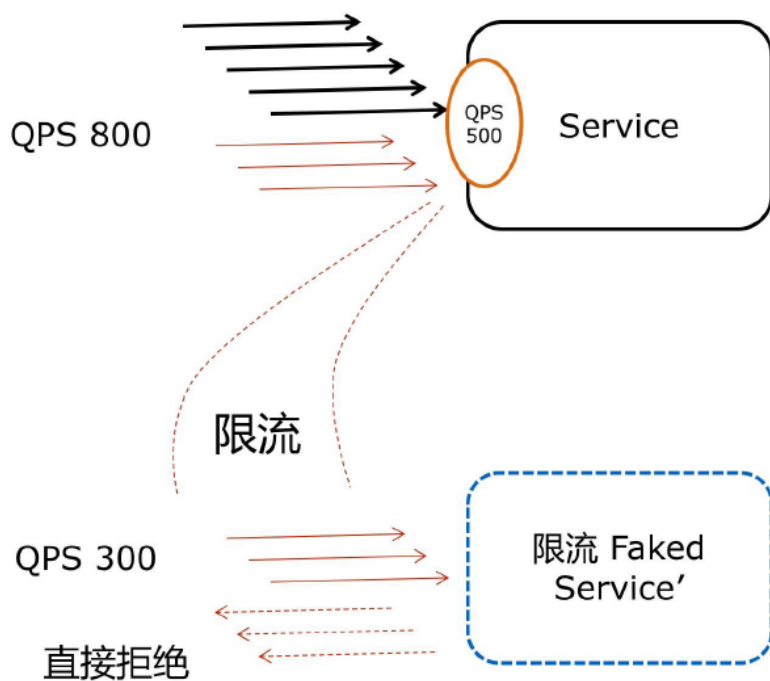
1.1 实时监控

监控接口的通过的QPS和拒绝的QPS



1.2 簇点链路

用来显示微服务的所监控的API



(流控处理逻辑可自定义)

最普适的场景:

- Provider 端控制
- 针对不同调用来源
- Web 接口流控

如何配置规则:

1. 梳理核心接口
2. 通过事前压测评估配置 QPS 阈值

同一个资源可以创建多条限流规则。FlowSlot 会对该资源的所有限流规则依次遍历，直到有规则触发限流或者所有规则遍历完毕。一条限流规则主要由下面几个因素组成，我们可以组合这些元素来实现不同的限流效果。

Field	说明	默认值
resource	资源名，资源名是限流规则的作用对象	
count	限流阈值	
grade	限流阈值类型，QPS 模式 (1) 或并发线程数模式 (0)	QPS 模式
limitApp	流控针对的调用来源	default，代表不区分调用来源
strategy	调用关系限流策略：直接、链路、关联	根据资源本身（直接）
controlBehavior	流控效果（直接拒绝/WarmUp/匀速+排队等待），不支持按调用关系限流	直接拒绝
clusterMode	是否集群限流	否

参考文档：<https://github.com/alibaba/Sentinel/wiki/%E6%B5%81%E9%87%8F%E6%8E%A7%E5%88%B6>

限流阈值类型

QPS (Query Per Second)：每秒请求数，就是说服务器在一秒的时间内处理了多少个请求。

QPS

进入簇点链路选择具体的访问的API，然后点击流控按钮

新增流控规则

资源名

/user/findOrderByUserId/1

针对来源

default

阈值类型

☒ QPS ☐ 线程数

单机阈值

2

QPS大于2就限流

是否集群

☐

流控模式

☒ 直接 ☐ 关联 ☐ 链路

流控效果

☒ 快速失败 ☐ Warm Up ☐ 排队等待

关闭高级选项

新增并继续添加

新增

取消

测试: <http://localhost:8800/user/findOrderByUserId/1>

← → ↺ ⓘ localhost:8800/user/findOrderByUserId/1

Blocked by Sentinel (flow limiting)

并发线程数

并发数控制用于保护业务线程池不被慢调用耗尽。例如，当应用所依赖的下游应用由于某种原因导致服务不稳定、响应延迟增加，对于调用者来说，意味着吞吐量下降和更多的线程数占用，极端情况下甚至导致线程池耗尽。为应对太多线程占用的情况，业内有使用隔离的方案，比如通过不同业务逻辑使用不同线程池来隔离业务自身之间的资源争抢（线程池隔离）。这种隔离方案虽然隔离性比较好，但是代价就是线程数目太多，线程上下文切换的 overhead 比较大，特别是对低延时的调用有比较大的影响。Sentinel 并发控制不负责创建和管理线程池，而是简单统计当前请求上下文的线程数目（正在执行的调用数目），如果超出阈值，新的请求会被立即拒绝，效果类似于信号量隔离。并发数控制通常在调用端进行配置。

编辑流控规则

资源名

/user/findOrderByUserId/{id}

针对来源

default

阈值类型

☐ QPS ☒ 线程数

单机阈值

5

是否集群

☐

流控模式

☒ 直接 ☐ 关联 ☐ 链路

关闭高级选项

保存

取消

可以利用jmeter测试

← → ↺ ⓘ localhost:8800/user/findOrderByUserId/1

```
{  
  msg: "接口限流了",  
  code: 100  
}
```

BlockException异常统一处理

springwebmvc接口资源限流入口在HandlerInterceptor的实现类AbstractSentinelInterceptor的preHandle方法中，对异常的处理是BlockExceptionHandler的实现类

| sentinel 1.7.1 引入了sentinel-spring-webmvc-adapter.jar

自定义BlockExceptionHandler 的实现类统一处理BlockException

```
1 @Slf4j
2 @Component
3 public class MyBlockExceptionHandler implements BlockExceptionHandler {
4     @Override
5     public void handle(HttpServletRequest request, HttpServletResponse response, BlockException e) throws Exception {
6         log.info("BlockExceptionHandler BlockException===== "+e.getRule());
7
8         R r = null;
9
10        if (e instanceof FlowException) {
11            r = R.error(100, "接口限流了");
12        }
13        else if (e instanceof DegradeException) {
14            r = R.error(101, "服务降级了");
15        }
16        else if (e instanceof ParamFlowException) {
17            r = R.error(102, "热点参数限流了");
18        }
19        else if (e instanceof SystemBlockException) {
20            r = R.error(103, "触发系统保护规则了");
21        }
22        else if (e instanceof AuthorityException) {
23            r = R.error(104, "授权规则不通过");
24        }
25
26        //返回json数据
27        response.setStatus(500);
28        response.setCharacterEncoding("utf-8");
29        response.setContentType(MediaType.APPLICATION_JSON_VALUE);
30        new ObjectMapper().writeValue(response.getWriter(), r);
31    }
32 }
33 }
```

测试：

← → ↺ ⓘ localhost:8800/user/findOrderByUserId/1#

```
{
  msg: "接口限流了",
  code: 100
}
```

流控模式

基于调用关系的流量控制。调用关系包括调用方、被调用方；一个方法可能会调用其它方法，形成一个调用链路的层次关系。

直接

资源调用达到设置的阈值后直接被流控抛出异常

Blocked by Sentinel (flow limiting)

关联

当两个资源之间具有资源争抢或者依赖关系的时候，这两个资源便具有了关联。比如对数据库同一个字段的读操作和写操作存在争抢，读的速度过高会影响写得速度，写的速度过高会影响读的速度。如果放任读写操作争抢资源，则争抢本身带来的开销会降低整体的吞吐量。可使用关联限流来避免具有关联关系的资源之间过度的争抢，举例来说，`read_db` 和 `write_db` 这两个资源分别代表数据库读写，我们可以给 `read_db` 设置限流规则来达到写优先的目的：设置 `strategy` 为 `RuleConstant.STRATEGY_RELATE` 同时设置 `refResource` 为 `write_db`。这样当写库操作过于频繁时，读数据的请求会被限流。

编辑流控规则

资源名

/order/get

针对来源

default

阈值类型

☒ QPS

☐ 线程数

单机阈值

2

是否集群

☐

流控模式

☐ 直接

☒ 关联

☐ 链路

关联资源

/order/add

流控效果

☒ 快速失败

☐ Warm Up

☐ 排队等待

当/order/add 一秒钟访问三次，将触发/order/get 限流

关闭高级选项

链路

根据调用链路入口限流。

下面中记录了资源之间的调用链路，这些资源通过调用关系，相互之间构成一棵调用树。这棵树的根节点是一个名字为 `getUser` 的虚拟节点，调用链的入口都是这个虚节点子节点。

一棵典型的调用树如下图所示：

```
1  getUser
2  /  \
3  /  \
4  /order/test1 /order/test2
```

上图中来自入口 `/order/test1` 和 `/order/test2` 的请求都调用到了资源 `getUser`，Sentinel 允许只根据某个入口的统计信息对资源限流。

新增流控规则

资源名	getUser <small>/test3&/test4都调用getUser,当getUser达到阈值后针对/test3限流</small>		
针对来源	default		
阈值类型	<input checked="" type="radio"/> QPS <input type="radio"/> 线程数	单机阈值	2
是否集群	<input type="checkbox"/>		
流控模式	<input type="radio"/> 直接 <input type="radio"/> 关联 <input checked="" type="radio"/> 链路		
入口资源	/test3		
流控效果	<input checked="" type="radio"/> 快速失败 <input type="radio"/> Warm Up <input type="radio"/> 排队等待		

测试会发现链路规则不生效

注意，高版本此功能直接使用不生效，如何解决？

从1.6.3 版本开始，Sentinel Web filter默认收敛所有URL的入口context，因此链路限流不生效。
1.7.0 版本开始（对应SCA的2.1.1.RELEASE），官方在CommonFilter 引入了WEB_CONTEXT_UNIFY 参数，用于控制是否收敛context。将其配置为 false 即可根据不同的URL 进行链路限流。
SCA 2.1.1.RELEASE之后的版本,可以通过配置spring.cloud.sentinel.web-context-unify=false即可关闭收敛

```
1 spring.cloud.sentinel.web-context-unify: false
```

测试，此场景拦截不到BlockException，对应@SentinelResource指定的资源必须在@SentinelResource注解中指定blockHandler处理BlockException

总结：为了解决链路规则引入ComonFilter的方式，除了此处问题，还会导致更多的问题，不建议使用ComonFilter的方式。流控链路模式的问题等待官方后续修复，或者使用AHAS。

快速失败

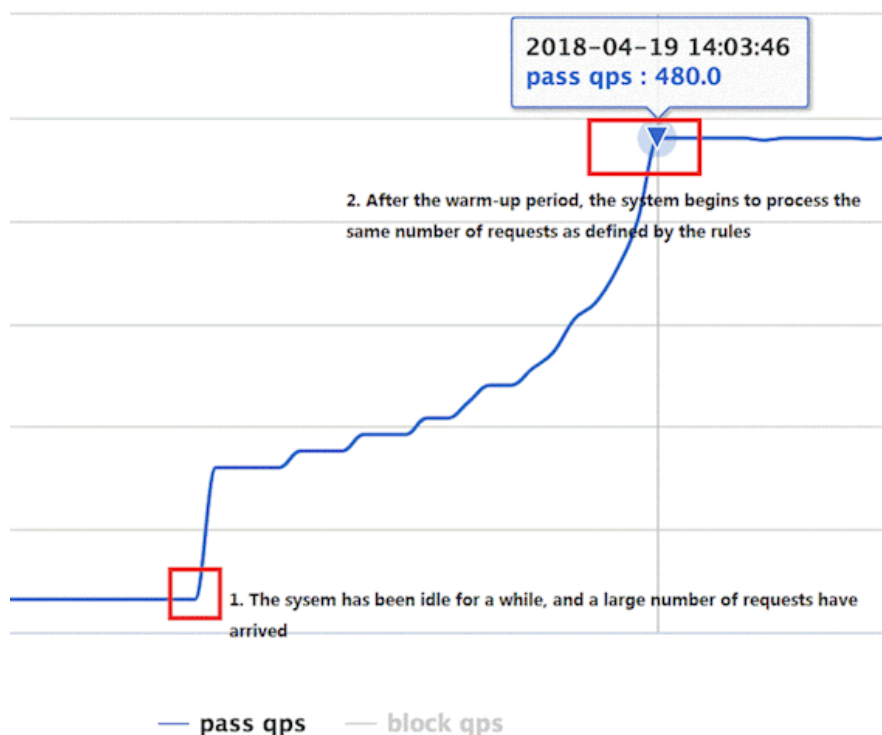
（RuleConstant.CONTROL_BEHAVIOR_DEFAULT）方式是默认的流量控制方式，当QPS超过任意规则的阈值后，新的请求就会被立即拒绝，拒绝方式为抛出FlowException。这种方式适用于对系统处理能力确切已知的情况下，比如通过压测确定了系统的准确水位时。

Warm Up（激增流量）

Warm Up（RuleConstant.CONTROL_BEHAVIOR_WARM_UP）方式，即预热/冷启动方式。当系统长期处于低水位的情况下，当流量突然增加时，直接把系统拉升到高水位可能瞬间把系统压垮。通过“冷启动”，让通过的流量缓慢增加，在一定时间内逐渐增加到阈值上限，给冷系统一个预热的时间，避免冷系统被压垮。

冷加载因子: codeFactor 默认是3，即请求 QPS 从 threshold / 3 开始，经预热时长逐渐升至设定的 QPS 阈值。

通常冷启动的过程系统允许通过的 QPS 曲线如下图所示



测试用例

```
1 @RequestMapping("/test")
2 public String test() {
3     try {
4         Thread.sleep(100);
5     } catch (InterruptedException e) {
6         e.printStackTrace();
7     }
8     return "=====test()=====";
9 }
```

编辑流控规则

编辑流控规则

资源名

/test

针对来源

default

阈值类型

☒ QPS ☐ 线程数

单机阈值

10

是否集群

☐

流控模式

☒ 直接 ☐ 关联 ☐ 链路

流控效果

☐ 快速失败 ☒ Warm Up ☐ 排队等待

预热时长

5 单位s

关闭高级选项

jmeter测试

线程属性

线程数: 200

Ramp-Up时间(秒): 10

循环次数 ☐ 永远 ☐ 1

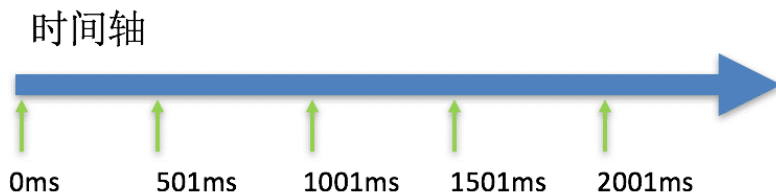
查看实时监控，可以看到通过QPS存在缓慢增加的过程



匀速排队 (脉冲流量)

匀速排队 (`RuleConstant.CONTROL_BEHAVIOR_RATE_LIMITER`) 方式会严格控制请求通过的间隔时间，也即是让请求以均匀的速度通过，对应的是漏桶算法。

该方式的作用如下图所示：



阈值QPS=2时，每隔500ms才允许通过下一个请求

这种方式主要用于处理间隔性突发的流量，例如消息队列。想象一下这样的场景，在某一秒有大量的请求到来，而接下来的几秒则处于空闲状态，我们希望系统能够在接下来的空闲期间逐渐处理这些请求，而不是在第一秒直接拒绝多余请求。

注意：匀速排队模式暂时不支持 QPS > 1000 的场景。

资源名: /test

针对来源: default

阈值类型: ☒ QPS ☐ 线程数 单机阈值: 5

是否集群: ☐ 排队等待200ms一个请求

流控模式: ☒ 直接 ☐ 关联 ☐ 链路

流控效果: ☐ 快速失败 ☐ Warm Up ☒ 排队等待

超时时间: 500 请求超时时间是500ms

[关闭高级选项](#)

jemeter压测

线程属性

线程数: 200

Ramp-Up时间(秒): 20

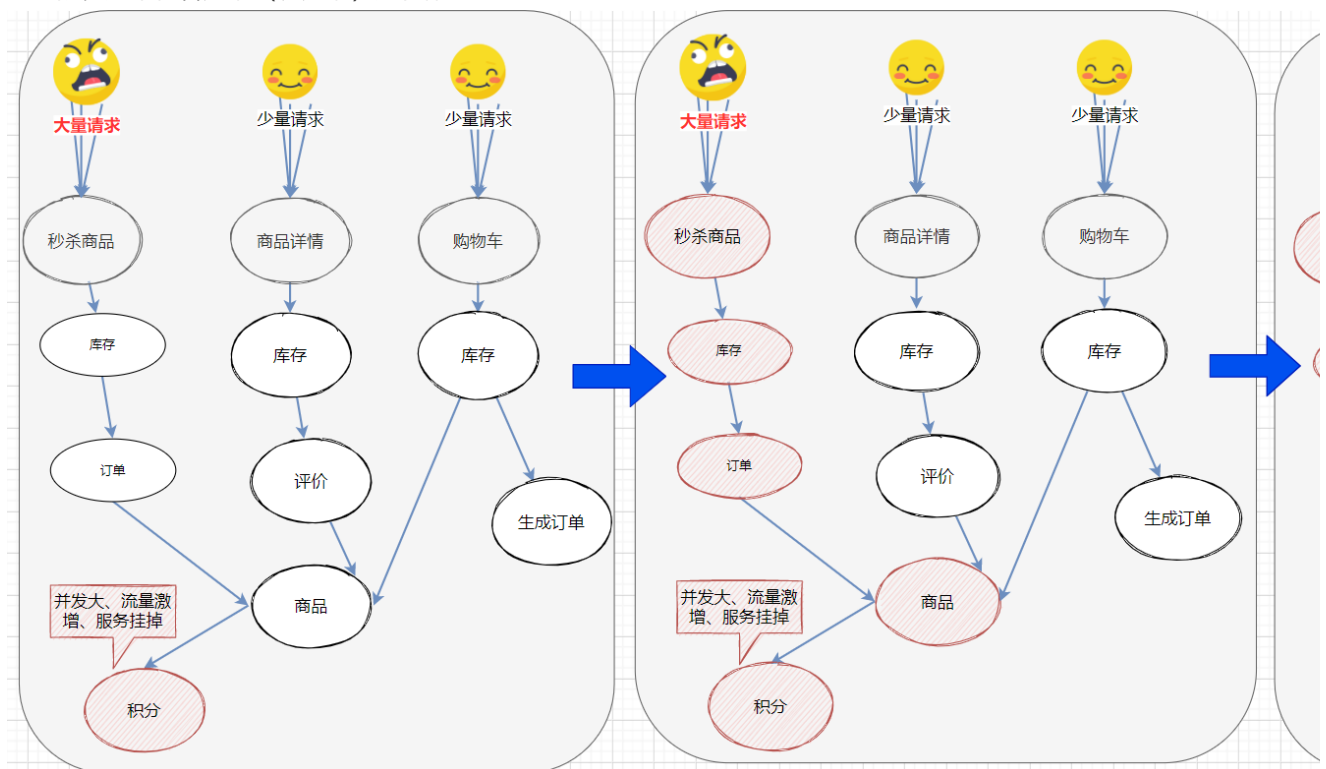
循环次数 ☐ 永远 ☐ 1

查看实时监控，可以看到通过QPS为5，体现了匀速排队效果



降级规则

除了流量控制以外，对调用链路中不稳定的资源进行熔断降级也是保障高可用的重要措施之一。我们需要对不稳定的弱依赖服务调用进行熔断降级，暂时切断不稳定调用，避免局部不稳定因素导致整体的雪崩。熔断降级作为保护自身的手段，通常在客户端（调用端）进行配置。



熔断降级与隔离

保护自身的手段

- 并发控制（信号量隔离）
 - 基于慢调用比例熔断
 - 基于异常比例熔断
- 通常在 **consumer** 端
组合配置

触发熔断后的处理逻辑示例

- 提供 fallback 实现（服务降级）
- 返回错误 result
- 读缓存（DB 访问降级）

熔断降级规则说明

熔断降级规则（DegradeRule）包含下面几个重要的属性：

Field	说明	默认值
resource	资源名，即规则的作用对象	
grade	熔断策略，支持慢调用比例/异常比例/异常数策略	慢调用比例
count	慢调用比例模式下为慢调用临界 RT（超出该值计为慢调用）；异常比例/异常数模式下为对应的阈值	
timeWindow	熔断时长，单位为 s	
minRequestAmount	熔断触发的最小请求数，请求数小于该值时即使异常比率超出阈值也不会熔断（1.7.0 引入）	5
statIntervalMs	统计时长（单位为 ms），如 60*1000 代表分钟级（1.8.0 引入）	1000 ms
slowRatioThreshold	慢调用比例阈值，仅慢调用比例模式有效（1.8.0 引入）	

熔断策略

慢调用比例

慢调用比例 (SLOW_REQUEST_RATIO)：选择以慢调用比例作为阈值，需要设置允许的慢调用 RT（即最大的响应时间），请求的响应时间大于该值则统计为慢调用。当单位统计时长（statIntervalMs）内请求数目大于设置的最小请求数目，并且慢调用的比例大于阈值，则接下来的熔断时长内请求会自动被熔断。经过熔断时长后熔断器会进入探测恢复状态（HALF-OPEN 状态），若接下来的一个请求响应时间小于设置的慢调用 RT 则结束熔断，若大于设置的慢调用 RT 则会再次被熔断。

编辑降级规则

资源名

/test

熔断策略

☒ 慢调用比例

☐ 异常比例

☐ 异常数

最大 RT

101

比例阈值

0.2

熔断时长

3

s

最小请求数

5

单位ms，当前允许的最大响应时间是101ms,大于当前值就是慢调用

单位时间内请求数目大于5，并且慢调用比例大于0.2，在接下来的3s内会自动熔断

保存

取消

测试用例

```
1 @RequestMapping("/test")
2 public String test() {
3     try {
4         Thread.sleep(100);
5     } catch (InterruptedException e) {
6         e.printStackTrace();
7     }
8     return "=====test()=====";
9 }
```

jemeter压测/test接口，保证每秒请求数超过配置的最小请求数

线程属性

线程数

200

Ramp-Up时间（秒）

20

循环次数

☐ 永远

☒ 1

查看实时监控，可以看到断路器熔断效果



此时浏览器访问会出现服务降级结果

← → ↺

localhost:8800/test

{

msg: "服务降级了",

code: 101

}

异常比例

异常比例 (ERROR_RATIO): 当单位统计时长 (statIntervalMs) 内请求数目大于设置的最小请求数目，并且异常的比例大于阈值，则接下来的熔断时长内请求会自动被熔断。经过熔断时长后熔断器会进入探测恢复状态 (HALF-OPEN 状态)，若接下来的一个请求成功完成 (没有错误) 则结束熔断，否则会再次被熔断。异常比率的阈值范围是 [0.0, 1.0]，代表 0% - 100%。

测试用例

```
1 @RequestMapping("/test2")
```

```

2 public String test2() {
3     atomicInteger.getAndIncrement();
4     if (atomicInteger.get() % 2 == 0){
5         //模拟异常和异常比率
6         int i = 1/0;
7     }
8
9     return "=====test2()=====";
10 }

```

配置降级规则

编辑降级规则

资源名: /test2

熔断策略: ☐ 慢调用比例 ☒ 异常比例 ☐ 异常数

比例阈值: 0.4

熔断时长: 3 s 最小请求数: 5

1s内请求数大于5，并且异常比例大于0.6，接下来熔断时长内请求就会自动熔断

查看实时监控，可以看到断路器熔断效果



异常数

异常数 (ERROR_COUNT): 当单位统计时长内的异常数目超过阈值之后会自动进行熔断。经过熔断时长后熔断器会进入探测恢复状态 (HALF-OPEN 状态)，若接下来的一个请求成功完成 (没有错误) 则结束熔断，否则会再次被熔断。

注意: 异常降级仅针对业务异常，对 Sentinel 限流降级本身的异常 (BlockException) 不生效。

配置降级规则

资源名: /test2

熔断策略: ☐ 慢调用比例 ☐ 异常比例 ☒ 异常数

异常数: 2

熔断时长: 3 s 最小请求数: 6

当单位统计时长1s内的异常数目超过阈值2之后会自动进行熔断

jmeter测试

线程属性

线程数: 60

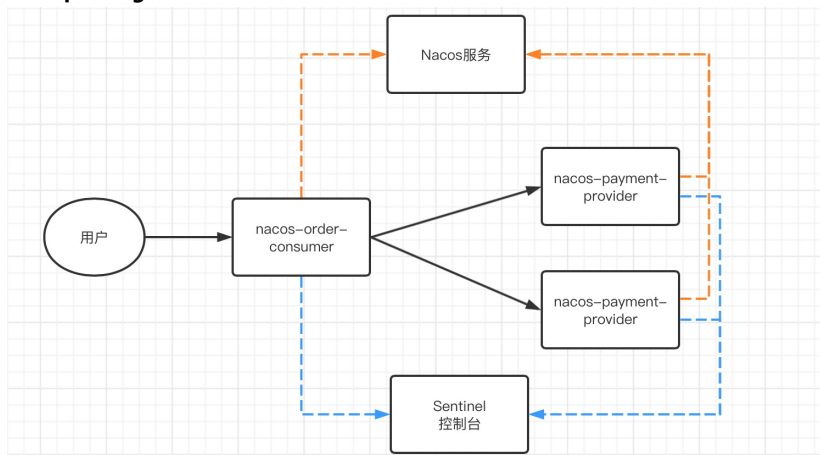
Ramp-Up时间 (秒): 10

循环次数: 永远 1

查看实时监控，可以看到断路器熔断效果



整合openfeign进行降级



引入

```

1 <dependency>
2 <groupId>com.alibaba.cloud</groupId>
3 <artifactId>spring-cloud-starter-alibaba-sentinel</artifactId>
4 </dependency>
5 <dependency>
6 <groupId>org.springframework.cloud</groupId>
7 <artifactId>spring-cloud-starter-openfeign</artifactId>
8 </dependency>

```

application.yml

```

1 #对Feign的支持
2 feign:
3   sentinel:
4     enabled: true # 添加feign对sentinel的支持
5

```

openfeign接口

```

1 @FeignClient(value = "nacos-payment-provider", fallback = ConsumerFallBackService.class )
2 public interface ConsumerService {
3
4   @GetMapping(value = "/paymentSQL/{id}")
5   public CommonResult<Payment> paymentSQL(@PathVariable("id") Long id);
6 }

```

openfeign的fallback实现类

```

1
2 @Component
3 public class ConsumerFallBackService implements ConsumerService {
4   @Override
5   public CommonResult<Payment> paymentSQL(Long id) {
6     return new CommonResult<Payment>("500", "进入兜底方法---ConsumerFallBackService", null);
7   }
8 }

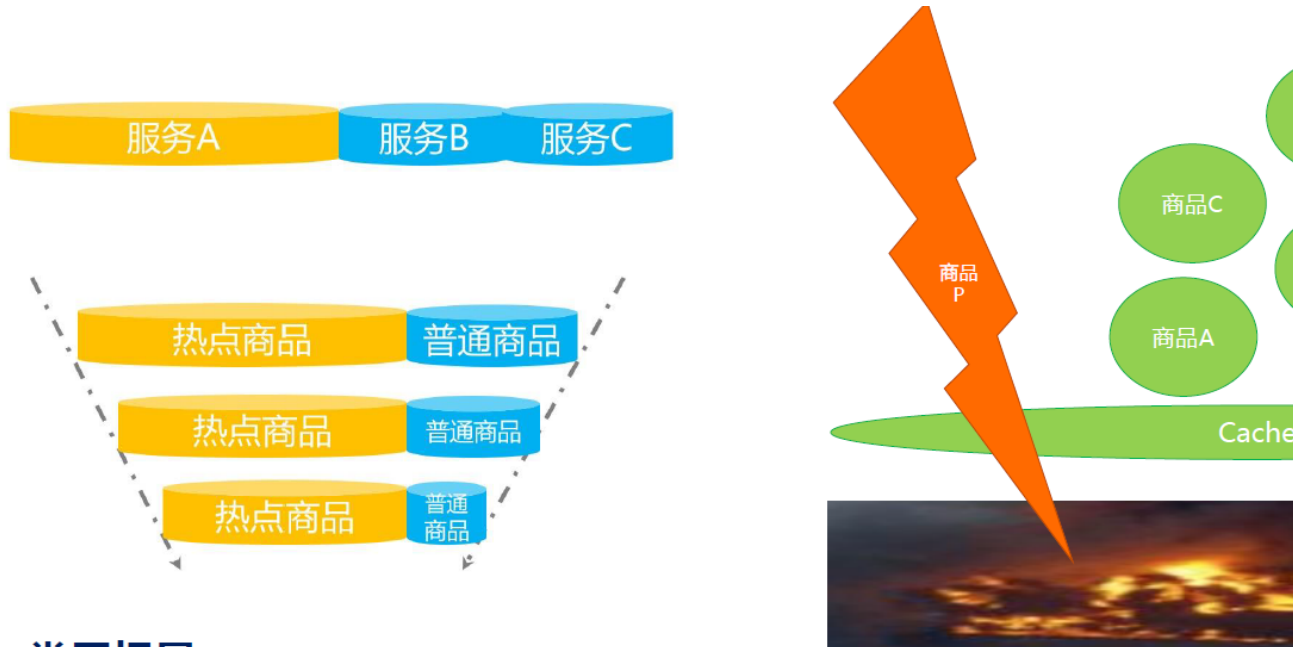
```

```
7 }  
8 }
```

热点参数限流

热点识别流控

何为热点？热点即经常访问的数据。很多时候我们希望统计某个热点数据中访问频次最高的数据，并对其访问进行限制。比如：

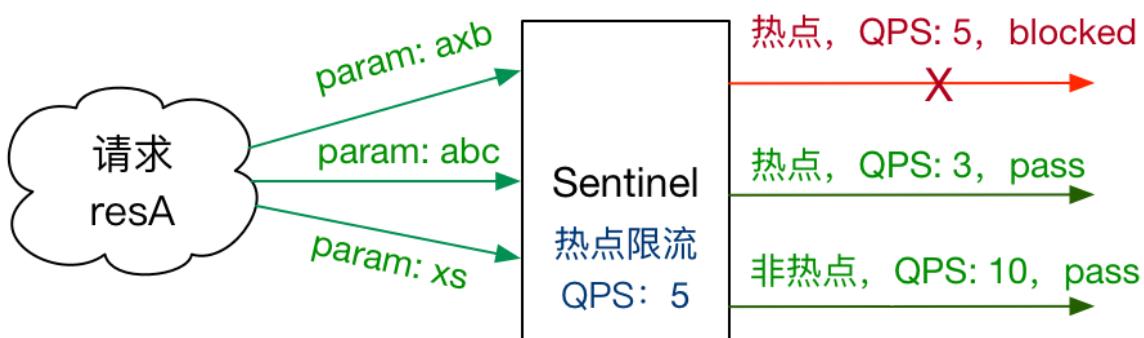


常用场景：

- 热点商品访问/操作控制
- 用户/IP 防刷

实现原理：热点淘汰策略 (LRU)

热点参数限流会统计传入参数中的热点参数，并根据配置的限流阈值与模式，对包含热点参数的资源调用进行限流。热点参数限流可以看做是一种特殊的流量控制，仅对包含热点参数的资源调用生效。



注意：

1. 热点规则需要使用@SentinelResource("resourceName")注解，否则不生效
2. 参数必须是7种基本数据类型才会生效

测试用例

```
1 @RequestMapping("/info/{id}")  
2 @SentinelResource(value = "userinfo",  
3   blockHandlerClass = CommonBlockHandler.class,  
4   blockHandler = "handleException2",  
5   fallbackClass = CommonFallback.class,  
6   fallback = "fallback")
```

```
7 )
8 public R info(@PathVariable("id") Integer id){
9     UserEntity user = userService.getById(id);
10    return R.ok().put("user", user);
11 }
```

单机阈值：针对所有参数的值进行设置的一个公共的阈值

1. 假设当前 参数 大部分的值都是热点流量， 单机阈值就是针对热点流量进行设置， 额外针对普通流量进行参数值流控
2. 假设当前 参数 大部分的值都是普通流量， 单机阈值就是针对普通流量进行设置， 额外针对热点流量进行参数值流控

配置热点参数规则

注意：资源名必须是@SentinelResource(value="资源名")中 配置的资源名，热点规则依赖于注解

编辑热点规则

资源名	userinfo <small>必须是@SentinelResource注解配置的资源</small>		
限流模式	QPS 模式		
参数索引	<input type="text" value="0"/>		
单机阈值	<input type="text" value="2"/>	统计窗口时长	<input type="text" value="1"/> 秒
是否集群	<input type="checkbox"/>		

具体到参数值限流，配置参数值为3,限流阈值为1

参数例外项			
参数类型	<input type="text" value="int"/>		
参数值	<input type="text" value="3"/>	限流阈值	<input type="text" value="1"/> + 添加
参数值	参数类型	限流阈值	操作

测试：

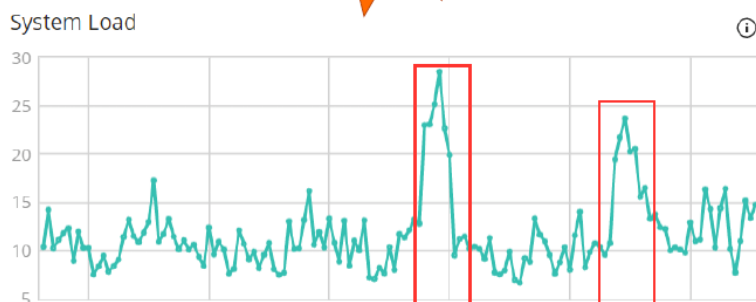
<http://localhost:8800/user/info/1> 限流的阈值为3

<http://localhost:8800/user/info/3> 限流的阈值为1

```
<  →  ↻  localhost:8800/user/info/3
{
  msg: "===被限流啦===com.alibaba.csp.sentinel.slots.block.flow.param.ParamFlowException: 3",
  code: -1
}
```

系统规则

CPU 迅速打满!



- 容量评估不到位, 某个大流量接口没有配置, 导致系统崩溃, 才
- 突然发现机器的 Load 和 CPU us 高, 但却没有办法很快的确认到问题的, 也来不及处理。
- 当其中一台机器挂了之后, 本该目的流量被负载均衡到另外的机器上也被打挂了, 引起系统雪崩。
- 希望有个全局的兜底防护, 即使也希望有一定的保护机制。

规则难配?

依赖难梳理?

We NEI
结合系统指标和服务容量,

Sentinel 系统自适应限流从整体维度对应用入口流量进行控制, 结合应用的 Load、CPU 使用率、总体平均 RT、入口 QPS 和并发线程数等几个维度的监控指标, 通过自适应的流控策略, 让系统的入口流量和系统的负载达到一个平衡, 让系统尽可能跑在最大吞吐量的同时保证系统整体的稳定性。

- **Load 自适应** (仅对 Linux/Unix-like 机器生效): 系统的 load1 作为启发指标, 进行自适应系统保护。当系统 load1 超过设定的启发值, 且系统当前的并发线程数超过估算的系统容量时才会触发系统保护 (BBR 阶段)。系统容量由系统的 $\text{maxQps} * \text{minRt}$ 估算得出。设定参考值一般是 $\text{CPU cores} * 2.5$ 。

<https://www.cnblogs.com/gentlemanhai/p/8484839.html>

- **CPU usage** (1.5.0+ 版本): 当系统 CPU 使用率超过阈值即触发系统保护 (取值范围 0.0-1.0), 比较灵敏。
- **平均 RT**: 当单台机器上所有入口流量的平均 RT 达到阈值即触发系统保护, 单位是毫秒。
- **并发线程数**: 当单台机器上所有入口流量的并发线程数达到阈值即触发系统保护。
- **入口 QPS**: 当单台机器上所有入口流量的 QPS 达到阈值即触发系统保护。

编写系统规则

阈值类型: ☐ LOAD ☐ RT ☐ 线程数 ☒ 入口 QPS ☐ CPU 使用率

阈值: 5 当前机器所有入口流量>5,触发系统保护

jemeter配置

线程属性

线程数: 60

Ramp-Up时间 (秒): 10

循环次数: ☒ 永远 1

测试结果

```
localhost:8800/user/info/1

{
  msg: "触发系统保护规则了",
  code: 103
}
```

授权控制规则

很多时候，我们需要根据调用来源来判断该次请求是否允许放行，这时候可以使用 Sentinel 的来源访问控制（黑白名单控制）的功能。来源访问控制根据资源的请求来源（origin）限制资源是否通过，若配置白名单则只有请求来源位于白名单内时才可通过；若配置黑名单则请求来源位于黑名单时不通过，其余的请求通过。

来源访问控制规则（AuthorityRule）非常简单，主要有以下配置项：

- resource：资源名，即限流规则的作用对象。
- limitApp：对应的黑名单/白名单，不同 origin 用，分隔，如 appA, appB。
- strategy：限制模式，AUTHORITY_WHITE 为白名单模式，AUTHORITY_BLACK 为黑名单模式，默认为白名单模式。

配置授权规则

编辑授权规则



第一步：实现 `com.alibaba.csp.sentinel.adapter.spring.webmvc.callback.RequestOriginParser` 接口，在 `parseOrigin` 方法中区分来源，并交给 spring 管理

注意：如果引入 `CommonFilter`，此处会多出一个

```
RequestOriginParser.java (sentinel-web-servlet-1.8.0-sources.jar\com\alibaba\csp\sentinel\adapter\servlet\callback)
RequestOriginParser.java (sentinel-spring-webmvc-adapter-1.8.0-sources.jar\com\alibaba\csp\sentinel\adapter\spring\webmvc\c

1 import com.alibaba.csp.sentinel.adapter.spring.webmvc.callback.RequestOriginParser;
2 import org.springframework.stereotype.Component;
3
4 import javax.servlet.http.HttpServletRequest;
5
6 /**
7  * @author Fox
8  */
9 @Component
10 public class MyRequestOriginParser implements RequestOriginParser {
11     /**
12      * 通过request获取来源标识，交给授权规则进行匹配
13      * @param request
14      * @return
15      */
16     @Override
17     public String parseOrigin(HttpServletRequest request) {
18         // 标识字段名称可以自定义
19         String origin = request.getParameter("serviceName");
20         // if (StringUtil.isBlank(origin)){
21         // throw new IllegalArgumentException("serviceName参数未指定");
22         // }
23         return origin;
24     }
25 }
```

测试：origin是order的请求不通过。

```
localhost:8800/test3?serviceName=order
```

```
{
  msg: "授权规则不通过",
  code: 104
}
```

```
localhost:8800/test3?serviceName=order2
```

```
{
  id: 1,
  username: "fox",
  age: 31
}
```

集群规则

为什么要使用集群流控呢？假设我们希望给某个用户限制调用某个 API 的总 QPS 为 50，但机器数可能很多（比如有 100 台）。这时候我们很自然地就想到，找一个 server 来专门来统计总的调用量，其它的实例都与这台 server 通信来判断是否可以调用。这就是最基础的集群流控的方式。

另外集群流控还可以解决流量不均匀导致总体限流效果不佳的问题。假设集群中有 10 台机器，我们给每台机器设置单机限流阈值为 10 QPS，理想情况下整个集群的限流阈值就为 100 QPS。不过实际情况下流量到每台机器可能会不均匀，会导致总量没有到的情况下某些机器就开始限流。因此仅靠单机维度去限制的话会无法精确地限制总体流量。而集群流控可以精确地控制整个集群的调用总量，结合单机限流兜底，可以更好地发挥流量控制的效果。

<https://github.com/alibaba/Sentinel/wiki/%E9%9B%86%E7%BE%A4%E6%B5%81%E6%8E%A7>

集群流控中共有两种身份：

- Token Client：集群流控客户端，用于向所属 Token Server 通信请求 token。集群限流服务端会返回给客户端结果，决定是否限流。
- Token Server：即集群流控服务端，处理来自 Token Client 的请求，根据配置的集群规则判断是否应该发放 token（是否允许通过）。

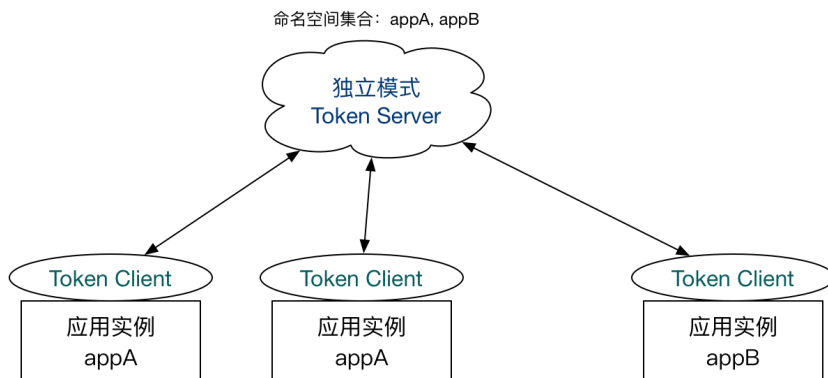
Sentinel 集群流控支持限流规则和热点规则两种规则，并支持两种形式的阈值计算方式：

- **集群总体模式**：即限制整个集群内的某个资源的总体 qps 不超过此阈值。
- **单机均摊模式**：单机均摊模式下配置的阈值等同于单机能够承受的限额，token server 会根据连接数来计算总的阈值（比如独立模式下有 3 个 client 连接到了 token server，然后配的单机均摊阈值为 10，则计算出的集群总量就为 30），按照计算出的总的阈值来进行限制。这种方式根据当前的连接数实时计算总的阈值，对于机器经常进行变更的环境非常适合。

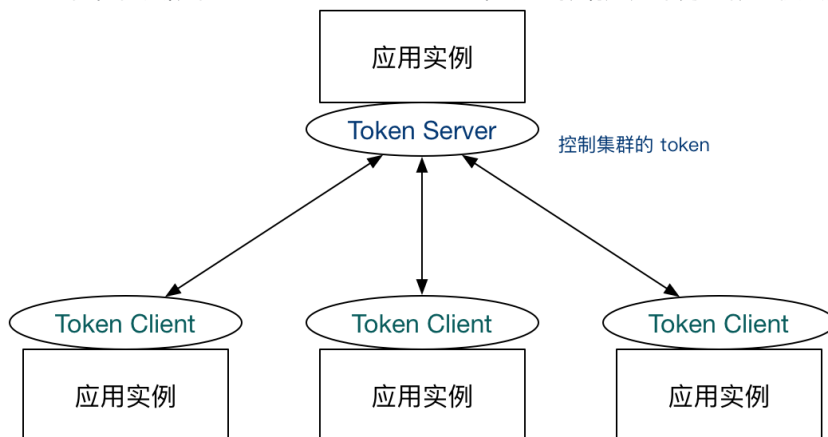
启动方式

Sentinel 集群限流服务端有两种启动方式：

- 独立模式（Alone），即作为独立的 token server 进程启动，独立部署，隔离性好，但是需要额外的部署操作。独立模式适合作为 Global Rate Limiter 给集群提供流控服务。



- 嵌入模式 (Embedded) , 即作为内置的 token server 与服务在同一进程中启动。在此模式下, 集群中各个实例都是对等的, token server 和 client 可以随时进行转变, 因此无需单独部署, 灵活性比较好。但是隔离性不佳, 需要限制 token server 的总 QPS, 防止影响应用本身。嵌入模式适合某个应用集群内部的流控。



云上版本 AHAS Sentinel 提供**开箱即用的全自动托管集群流控能力**, 无需手动指定/分配 token server 以及管理连接状态, 同时支持分钟小时级别流控、大流量低延时场景流控场景, 同时支持 Istio/Envoy 场景的 Mesh 流控能力。

文档: 08-1 Sentinel控制台规则配置详解.not...

链接: [http://note.youdao.com/noteshare?](http://note.youdao.com/noteshare?id=fc7d801d1a8213fc4d1b691302e82a62&sub=675ACF22A59841AFAECD58A0A9B4E151)

id=fc7d801d1a8213fc4d1b691302e82a62&sub=675ACF22A59841AFAECD58A0A9B4E151