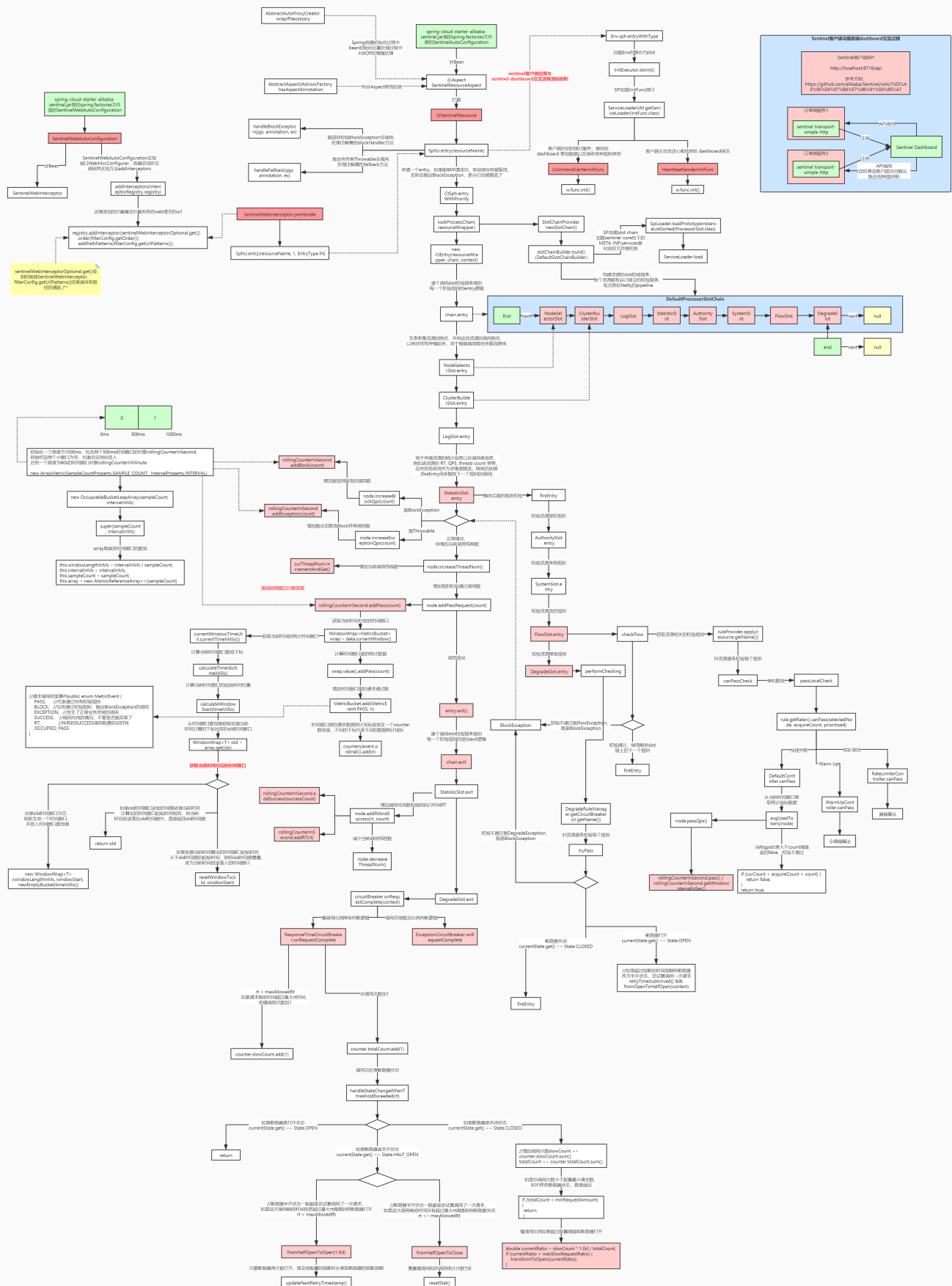


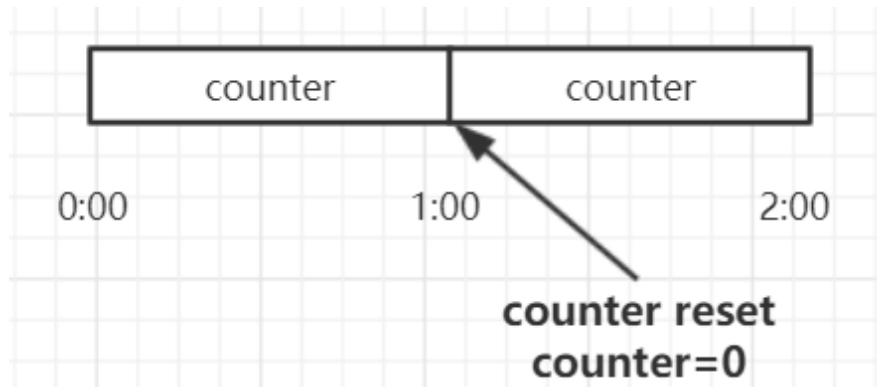
Sentinel限流、熔断降级源码架构图



常见限流算法精讲

计数器法

计数器法是限流算法里最简单也是最容易实现的一种算法。比如我们规定，对于A接口来说，我们1分钟的访问次数不能超过100个。那么我们可以这么做：在一开始的时候，我们可以设置一个计数器counter，每当一个请求过来的时候，counter就加1，如果counter的值大于100并且该请求与第一个请求的间隔时间还在1分钟之内，那么说明请求数过多；如果该请求与第一个请求的间隔时间大于1分钟，且counter的值还在限流范围内，那么就重置 counter。

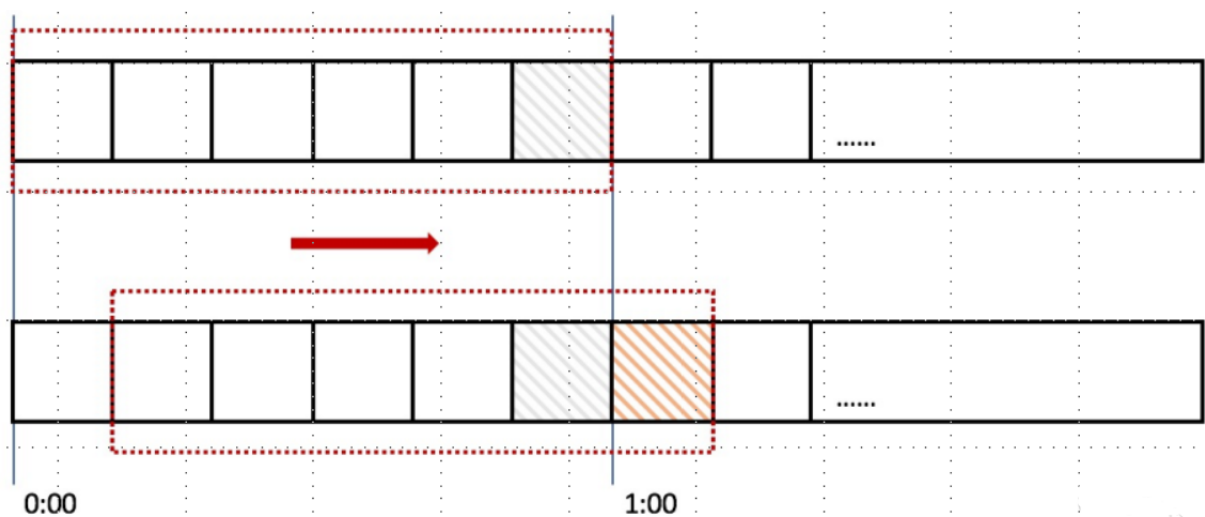


具体算法的伪代码：

```
1  /**
2   * 最简单的计数器限流算法
3   */
4  public class Counter {
5      public long timeStamp = System.currentTimeMillis(); // 当前时间
6      public int reqCount = 0; // 初始化计数器
7      public final int limit = 100; // 时间窗口内最大请求数
8      public final long interval = 1000 * 60; // 时间窗口ms
9
10     public boolean limit() {
11         long now = System.currentTimeMillis();
12         if (now < timeStamp + interval) {
13             // 在时间窗口内
14             reqCount++;
15             // 判断当前时间窗口内是否超过最大请求控制数
16             return reqCount <= limit;
17         } else {
18             timeStamp = now;
19             // 超时后重置
20             reqCount = 1;
21             return true;
22         }
23     }
24 }
```

滑动时间窗口算法

滑动时间窗口，又称rolling window。为了解决计数器法统计精度太低的问题，引入了滑动窗口算法。下面这张图，很好地解释了滑动窗口算法：



在上图中，整个红色的矩形框表示一个时间窗口，在我们的例子中，一个时间窗口就是一分钟。然后将时间窗口进行划分，比如图中，我们就将滑动窗口划成了6格，所以每格代表的是10秒钟。每过10秒钟，我们的时间窗口就会往右滑动一格。每一个格子都有自己独立的计数器counter，比如当一个请求在0:35秒的时候到达，那么0:30~0:39对应的counter就会加1。计数器算法其实就是滑动窗口算法。只是它没有对时间窗口做进一步地划分，所以只有1格。由此可见，当滑动窗口的格子划分的越多，那么滑动窗口的滚动就越平滑，限流的统计就会越精确。

具体算法的伪代码：

```
1  /**
2   * 滑动时间窗口限流实现
3   * 假设某个服务最多只能每秒钟处理100个请求，我们可以设置一个1秒钟的滑动时间窗口，
4   * 窗口中有10个格子，每个格子100毫秒，每100毫秒移动一次，每次移动都需要记录当前服务请求的次数
5   */
6  public class SlidingTimeWindow {
7      //服务访问次数，可以放在Redis中，实现分布式系统的访问计数
8      Long counter = 0L;
9      //使用LinkedList来记录滑动窗口的10个格子。
10     LinkedList<Long> slots = new LinkedList<Long>();
11
12     public static void main(String[] args) throws InterruptedException {
13         SlidingTimeWindow timeWindow = new SlidingTimeWindow();
14
15         new Thread(new Runnable() {
16             @Override
17             public void run() {
18                 try {
19                     timeWindow.doCheck();
```

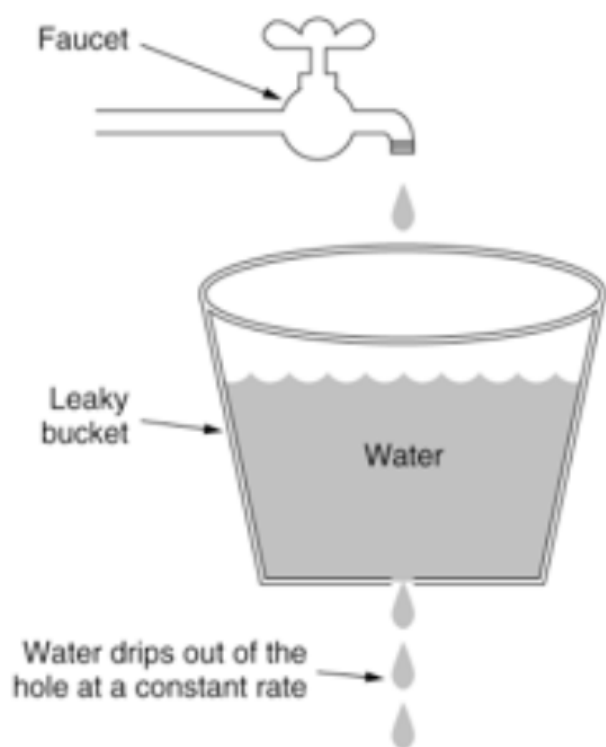
```

20 } catch (InterruptedException e) {
21     e.printStackTrace();
22 }
23 }
24 }).start();
25
26 while (true){
27     //TODO 判断限流标记
28     timeWindow.counter++;
29     Thread.sleep(new Random().nextInt(15));
30 }
31 }
32
33 private void doCheck() throws InterruptedException {
34     while (true) {
35         slots.addLast(counter);
36         if (slots.size() > 10) {
37             slots.removeFirst();
38         }
39         //比较最后一个和第一个，两者相差100以上就限流
40         if ((slots.peekLast() - slots.peekFirst()) > 100) {
41             System.out.println("限流了。。");
42             //TODO 修改限流标记为true
43         }else {
44             //TODO 修改限流标记为false
45         }
46
47         Thread.sleep(100);
48     }
49 }
50 }

```

漏桶算法

漏桶算法，又称leaky bucket。



从图中我们可以看到，整个算法其实十分简单。首先，我们有一个固定容量的桶，有水流进来，也有水流出去。对于流进来的水来说，我们无法预计一共有多少水会流进来，也无法预计水流的速度。但是对于流出去的水来说，这个桶可以固定水流出的速率。而且，当桶满了之后，多余的水将会溢出。

我们将算法中的水换成实际应用中的请求，我们可以看到漏桶算法天生就限制了请求的速度。当使用了漏桶算法，我们可以保证接口会以一个常速速率来处理请求。所以漏桶算法天生不会出现临界问题。

具体的伪代码如下：

```
1  /**
2   * 漏桶限流算法
3   */
4  public class LeakyBucket {
5      public long timeStamp = System.currentTimeMillis(); // 当前时间
6      public long capacity; // 桶的容量
7      public long rate; // 水漏出的速度(每秒系统能处理的请求数)
8      public long water; // 当前水量(当前累积请求数)
9
10     public boolean limit() {
11         long now = System.currentTimeMillis();
12         water = Math.max(0, water - ((now - timeStamp)/1000) * rate); // 先执行漏水，
            计算剩余水量
13         timeStamp = now;
14         if ((water + 1) < capacity) {
15             // 尝试加水,并且水还未满
16             water += 1;
```

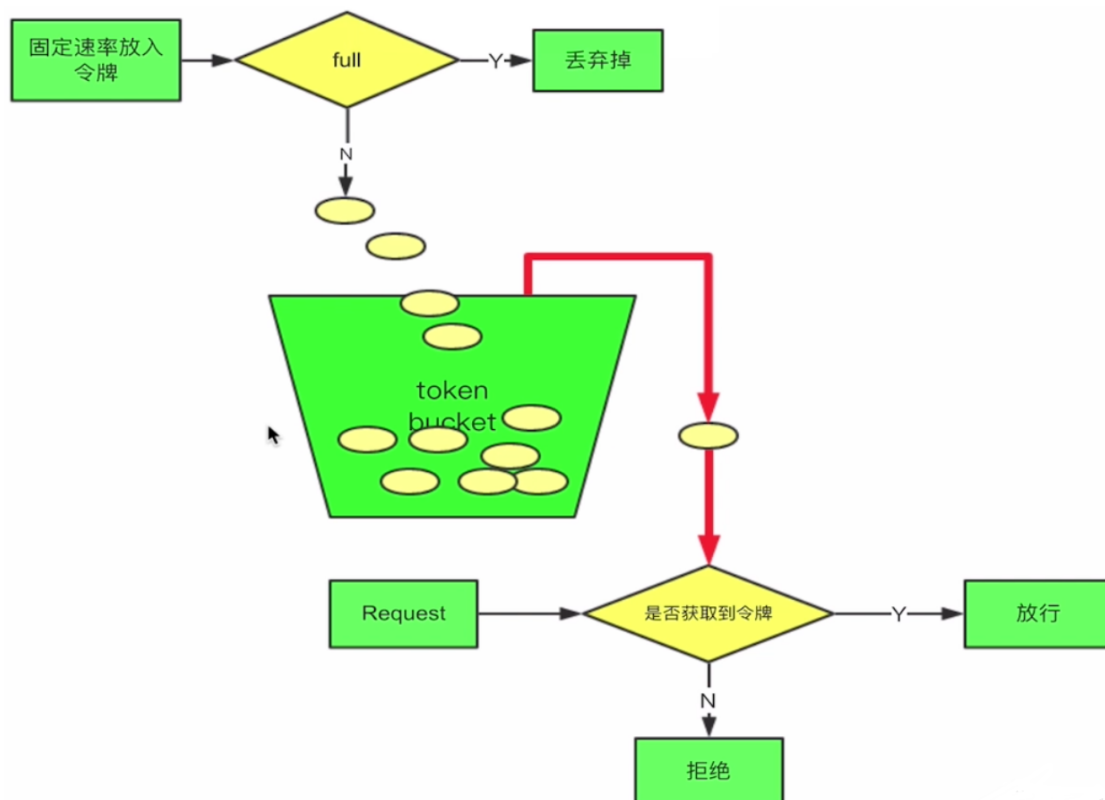
```

17 return true;
18 } else {
19 // 水满，拒绝加水
20 return false;
21 }
22 }
23 }

```

令牌桶算法

令牌桶算法，又称token bucket。同样为了理解该算法，我们来看一下该算法的示意图：



从图中我们可以看到，令牌桶算法比漏桶算法稍显复杂。首先，我们有一个固定容量的桶，桶里存放着令牌（token）。桶一开始是空的，token以一个固定的速率 r 往桶里填充，直到达到桶的容量，多余的令牌将会被丢弃。每当一个请求过来时，就会尝试从桶里移除一个令牌，如果没有令牌的话，请求无法通过。

具体的伪代码如下：

```

1 /**
2  * 令牌桶限流算法
3  */
4 public class TokenBucket {
5     public long timeStamp = System.currentTimeMillis(); // 当前时间
6     public long capacity; // 桶的容量
7     public long rate; // 令牌放入速度
8     public long tokens; // 当前令牌数量
9

```

```
10 public boolean grant() {
11     long now = System.currentTimeMillis();
12     // 先添加令牌
13     tokens = Math.min(capacity, tokens + (now - timeStamp) * rate);
14     timeStamp = now;
15     if (tokens < 1) {
16         // 若不到1个令牌,则拒绝
17         return false;
18     } else {
19         // 还有令牌, 领取令牌
20         tokens -= 1;
21         return true;
22     }
23 }
24 }
```

限流算法小结

计数器 VS 滑动窗口：

- 1 计数器算法是最简单的算法，可以看成是滑动窗口的低精度实现。
- 2 滑动窗口由于需要存储多份的计数器（每一个格子存一份），所以滑动窗口在实现上需要更多的存储空间。
- 3 也就是说，如果滑动窗口的精度越高，需要的存储空间就越大。

漏桶算法 VS 令牌桶算法：

- 1 漏桶算法和令牌桶算法最明显的区别是令牌桶算法允许流量一定程度的突发。
- 2 因为默认的令牌桶算法，取走token是不需要耗费时间的，也就是说，假设桶内有100个token时，那么可以瞬间允许100个请求通过。
- 3 当然我们需要具体情况具体分析，只有最合适的算法，没有最优的算法。

- 1 文档：10-VIP-Sentinel限流、熔断降级源码剖析
- 2 链接：<http://note.youdao.com/noteshare?id=19206aab917935d3a101fa94d2b68700&sub=f2FAAB34BDBA4CA5AA8FF90217407055>