

Redis

1. 我们的项目中那些地方用到了redis

应用redis的场景

高性能

假设这么个场景，你有个操作，一个请求过来，吭哧吭哧你各种乱七八糟操作 `mysql`，半天查出来一个结果，耗时 `600ms`。但是这个结果可能接下来几个小时都不会变了，或者变了也可以不用立即反馈给用户。那么此时咋办？

缓存啊，折腾 `600ms` 查出来的结果，扔缓存里，一个 `key` 对应一个 `value`，下次再有人查，别走 `mysql` 折腾 `600ms` 了，直接从缓存里，通过一个 `key` 查出来一个 `value`，`2ms` 搞定。性能提升 `300` 倍。

就是说对于一些需要复杂操作耗时查出来的结果，且确定后面不怎么变化，但是有很多读请求，那么直接将查询出来的结果放在缓存中，后面直接读缓存就好。

高并发

`mysql` 这么重的数据库，压根儿设计不是让你玩儿高并发的，虽然也可以玩儿，但是天然支持不好。`mysql` 单机支撑到 `2000QPS` 也开始容易报警了。

所以要是你有个系统，高峰期一秒钟过来的请求有 `1万`，那一个 `mysql` 单机绝对会死掉。你这个时候就只能上缓存，把很多数据放缓存，别放 `mysql`。缓存功能简单，说白了就是 `key-value` 式操作，单机支撑的并发量轻松一秒几万十几万，支撑高并发 `so easy`。单机承载并发量是 `mysql` 单机的几十倍。

我们的电商项目中用到redis的地方

1. **搜索模块中**：所有的商品的品牌、分类、规格、规格对应的规格选项，这些数据数据量小，且读的频繁（每日去搜索界面都会用到）
2. **购物（订单）车模块中**：订单本身数据量小，而且用户会对订单进行查看、修改、确认
3. **秒杀模块中**：秒杀模块中的商品、以及对应的秒杀订单，这部分数据本身数据量不大且秒杀模块本身被访问的频率就高，完全符合redis 的应用场景

2.redis基础

2.1 redis支持的数据类型

redis 主要有以下几种数据类型：

- string
- hash
- list
- set
- sorted set

string

这是最简单的类型，就是普通的set和get，做简单的KV缓存

set key value

hash

这个是类似 `map` 的一种结构，这个一般就是可以将结构化的数据，比如一个对象（前提是这个对象没嵌套其他的对象）给缓存在 `redis` 里，然后每次读写缓存的时候，可以就操作 `hash` 里的某个字段。

```
hset person name Jenny
hset person age 12
hset person sex female
hset person tel 123123123
```

```
person = {"name":"Jenny","age":12,"sex":"female","tel":"123123123"}
```

list

`list` 是有序列表

比如可以通过 `list` 存储一些列表型的数据结构，类似粉丝列表、文章的评论列表之类的东西。

比如可以通过 `lrange` 命令，读取某个闭区间内的元素，可以基于 `list` 实现分页查询。

比如可以搞个简单的消息队列，从 `list` 头进去，从 `list` 尾巴出来。

set

`set` 是无序集合，自动去重。

直接基于 `set` 将系统里需要去重的数据扔进去，自动就给去重了，如果需要对一些数据进行快速的全局去重，当然也可以基于 `jvm` 内存里的 `HashSet` 进行去重，但是如果某个系统部署在多台机器上呢？得基于 `redis` 进行全局的 `set` 去重。

可以基于 `set` 进行交集、并集、差集的操作。

sorted set

`sorted set` 是排序的 `set`，去重但可以排序，写进去的时候给一个分数，自动根据分数排序。

2.2 redis持久化的方式

`Redis` 提供了两种不同的持久化方法来讲数据存储到硬盘里边：

- `RDB`(基于快照)，将某一时刻的所有数据保存到一个 `RDB` 文件中。
- `AOF`(`append-only-file`)，当 `Redis` 服务器执行写命令的时候，将执行的写命令保存到 `AOF` 文件中。

2.2.1 持久化的选择

`RDB` 和 `AOF` 并不互斥，它俩可以同时使用。我们可以根据他们的优缺点进行选择。

- 不要仅仅使用 `RDB`，因为那样会导致你丢失很多数据；
- 也不要仅仅使用 `AOF`，因为那样有两个问题：第一，你通过 `AOF` 做冷备，没有 `RDB` 做冷备来的恢复速度更快；第二，`RDB` 每次简单粗暴生成数据快照，更加健壮，可以避免 `AOF` 这种复杂的备份和恢复机制的 `bug`；
- `redis` 支持同时开启两种持久化方式，我们可以综合使用 `AOF` 和 `RDB` 两种持久化机制，用 `AOF` 来保证数据不丢失，作为数据恢复的第一选择；用 `RDB` 来做不同程度的冷备，在 `AOF` 文件都丢失或损坏不可用的时候，还可以使用 `RDB` 来进行快速的数据恢复。

如果Redis服务器同时开启了RDB和AOF持久化，服务器会优先使用AOF文件来还原数据(因为AOF更新频率比RDB更新频率要高，还原的数据更完善)。

2.3 过期策略

redis 过期策略是：定期删除+惰性删除；

所谓定期删除，指的是 redis 默认是每隔 100ms 就随机抽取一些设置了过期时间的 key，检查其是否过期，如果过期就删除。

假设 redis 里放了 10w 个 key，都设置了过期时间，你每隔几百毫秒，就检查 10w 个 key，那 redis 基本上就死了，cpu 负载会很高的，消耗在你的检查过期 key 上了。注意，这里可不是每隔 100ms 就遍历所有的设置过期时间的 key，那样就是一场性能上的灾难。实际上 redis 是每隔 100ms 随机抽取一些 key 来检查和删除的。

但是问题是，定期删除可能会导致很多过期 key 到了时间并没有被删除掉，那怎么办呢？所以就是惰性删除了。这就是说，在你获取某个 key 的时候，redis 会检查一下，这个 key 如果设置了过期时间那么是否过期了？如果过期了此时就会删除，不会给你返回任何东西。

但是实际上这还是有问题的，如果定期删除漏掉了很多过期 key，然后你也没及时去查，也就没走惰性删除，此时会怎么样？如果大量过期 key 堆积在内存里，导致 redis 内存块耗尽了，那怎么办呢？这种情况下，我们需要使用内存淘汰机制。

2.3 内存淘汰机制

redis 内存淘汰机制有以下几个：

- **noeviction**：当内存不足以容纳新写入数据时，新写入操作会报错，这个一般没人用吧，实在是太恶心了。
- **allkeys-lru**：当内存不足以容纳新写入数据时，在键空间中，移除最近最少使用的 key（这个是最常用的）。
- **allkeys-random**：当内存不足以容纳新写入数据时，在键空间中，随机移除某个 key，这个一般没人用吧，为啥要随机，肯定是把最近最少使用的 key 给干掉啊。
- **volatile-lru**：当内存不足以容纳新写入数据时，在设置了过期时间的键空间中，移除最近最少使用的 key（这个一般不太合适）。
- **volatile-random**：当内存不足以容纳新写入数据时，在设置了过期时间的键空间中，随机移除某个 key。

volatile-ttl：当内存不足以容纳新写入数据时，在设置了过期时间的键空间中，有更早过期时间的 key 优先移除。

3.redis主从的架构

单机的 redis，能够承载的 QPS 大概就在上万到几万不等。对于缓存来说，一般都是用来支撑读高并发的。因此架构做成主从(master-slave)架构，一主多从，主负责写，并且将数据复制到其它的 slave 节点，从节点负责读。所有的读请求全部走从节点。这样也可以很轻松实现水平扩容，支撑读高并发。

3.1 redis replication 的核心机制

• redis 采用异步方式复制数据到 slave 节点，不过 redis2.8 开始，slave node 会周期性地确认自己每次复制的数据量；

• 一个 master node 是可以配置多个 slave node 的；

• slave node 也可以连接其他的 slave node；

- **slave node** 做复制的时候，不会 **block master node** 的正常工作；
- **slave node** 在做复制的时候，也不会 **block** 对自己的查询操作，它会用旧的数据集来提供服务；但是复制完成的时候，需要删除旧数据集，加载新数据集，这个时候就会暂停对外服务了；
- **slave node** 主要用来进行横向扩容，做读写分离，扩容的 **slave node** 可以提高读的吞吐量。

注意，如果采用了主从架构，那么建议必须**开启 master node** 的持久化，不建议用 **slave node** 作为 **master node** 的数据热备，因为那样的话，如果你关掉 **master** 的持久化，可能在 **master** 宕机重启的时候数据是空的，然后可能一经过复制，**slave node** 的数据也丢了。

另外，**master** 的各种备份方案，也需要做。万一本地的所有文件丢失了，从备份中挑选一份 **rdb** 去恢复 **master**，这样才能**确保启动的时候，是有数据的**，**slave node** 可以自动接管 **master node**，但也可能 **sentinel** 还没检测到 **master failure**，**master node** 就自动重启了，还是可能导致上面所有的 **slave node** 数据被清空。

3.2 redis 主从复制的核心原理

当启动一个 **slave node** 的时候，它会发送一个 **PSYNC** 命令给 **master node**。

如果这是 **slave node** 初次连接到 **master node**，那么会触发一次 **full resynchronization** 全量复制。此时 **master** 会启动一个后台线程，开始生成一份 **RDB** 快照文件，同时还会将从客户端 **client** 新收到的所有写命令缓存在内存中。**RDB** 文件生成完毕后，**master** 会将这个 **RDB** 发送给 **slave**，**slave** 会先**写入本地磁盘，然后再从本地磁盘加载到内存中**，接着 **master** 会将内存中缓存的写命令发送到 **slave**，**slave** 也会同步这些数据。**slave node** 如果跟 **master node** 有网络故障，断开了连接，会自动重连，连接之后 **master node** 仅会复制给 **slave** 部分缺少的数据。

3.3 主从复制的断点续传

从 **redis2.8** 开始，就支持主从复制的断点续传，如果主从复制过程中，网络连接断掉了，那么可以接着上次复制的地方，继续复制下去，而不是从头开始复制一份(增量复制)。

master node 会在内存中维护一个 **backlog**，**master** 和 **slave** 都会保存一个 **replica offset** 还有一个 **master run id**，**offset** 就是保存在 **backlog** 中的。如果 **master** 和 **slave** 网络连接断掉了，**slave** 会让 **master** 从上次 **replica offset** 开始继续复制，如果没有找到对应的 **offset**，那么就会执行一次 **resynchronization**。

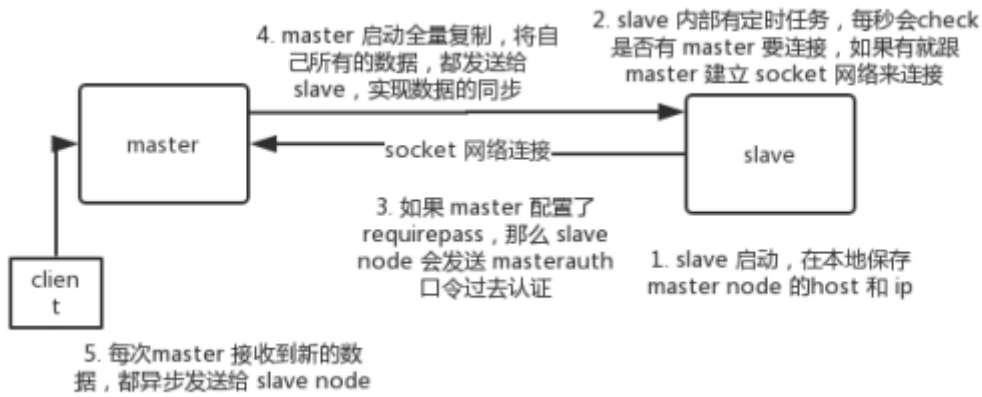
3.4 过期 key 处理

slave 不会过期 **key**，只会等待 **master** 过期 **key**。如果 **master** 过期了一个 **key**，或者通过 **LRU** 淘汰了一个 **key**，那么会模拟一条 **del** 命令发送给 **slave**。

3.5 复制的完整流程

slave node 启动时，会在自己本地保存 **master node** 的信息，包括 **master node** 的 **host** 和 **ip**，但是复制流程没开始。

slave node 内部有个定时任务，每秒检查是否有新的 **master node** 要连接和复制，如果发现，就跟 **master node** 建立 **socket** 网络连接。然后 **slave node** 发送 **ping** 命令给 **master node**。如果 **master** 设置了 **requirepass**，那么 **slave node** 必须发送 **masterauth** 的口令过去进行认证。**master node** **第一次执行全量复制**，将所有数据发给 **slave node**。而在后续，**master node** 持续将写命令，异步复制给 **slave node**。



3.6 全量复制

- master 执行 `bgsave`，在本地生成一份 `rdb` 快照文件。
- master node 将 `rdb` 快照文件发送给 slave node，如果 `rdb` 复制时间超过 60秒（`repl-timeout`），那么 slave node 就会认为复制失败，可以适当调大这个参数。
- master node 在生成 `rdb` 时，会将所有新的写命令缓存在内存中，在 slave node 保存了 `rdb` 之后，再将新的写命令复制给 slave node。
- 如果在复制期间，内存缓冲区持续消耗超过 64MB，或者一次性超过 256MB，那么停止复制，复制失败。
- slave node 接收到 `rdb` 之后，清空自己的旧数据，然后重新加载 `rdb` 到自己的内存中，同时基于旧的数据版本对外提供服务。
- 如果 slave node 开启了 AOF，那么会立即执行 `BGREWRITEAOF`，重写 AOF。

3.7 增量复制

- 如果全量复制过程中，master-slave 网络连接断掉，那么 slave 重新连接 master 时，会触发增量复制。
- master 直接从自己的 backlog 中获取部分丢失的数据，发送给 slave node，默认 backlog 就是 1MB。
- master 就是根据 slave 发送的 `psync` 中的 `offset` 来从 backlog 中获取数据的。

3.8 异步复制

master 每次接收到写命令之后，先在内部写入数据，然后异步发送给 slave node。

4.redis的哨兵高可用

4.1 什么是哨兵

哨兵是 redis 集群机构中非常重要的一个组件，主要有以下功能：

- 集群监控：负责监控 redis master 和 slave 进程是否正常工作。
- 消息通知：如果某个 redis 实例有故障，那么哨兵负责发送消息作为报警通知给管理员。
- 故障转移：如果 master node 挂掉了，会自动转移到 slave node 上。
- 配置中心：如果故障转移发生了，通知 client 客户端新的 master 地址。

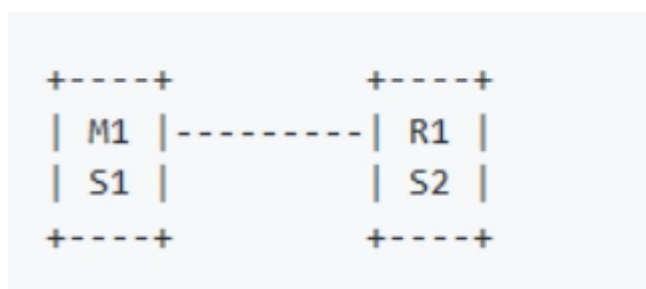
哨兵用于实现 **redis** 集群的高可用，本身也是分布式的，作为一个哨兵集群去运行，互相协同工作。

- 故障转移时，判断一个 **master node** 是否宕机了，需要大部分的哨兵都同意才行，涉及到了分布式选举的问题。
- 即使部分哨兵节点挂掉了，哨兵集群还是能正常工作的。

4.2 哨兵的重点核心

- 哨兵至少需要 3 个实例，来保证自己的健壮性。
- 哨兵 + **redis** 主从的部署架构，是 **不保证数据零丢失**的，只能保证 **redis** 集群的高可用性。

哨兵集群必须部署 2 个以上节点，如果哨兵集群仅仅部署了 2 个哨兵实例，**quorum = 1**。



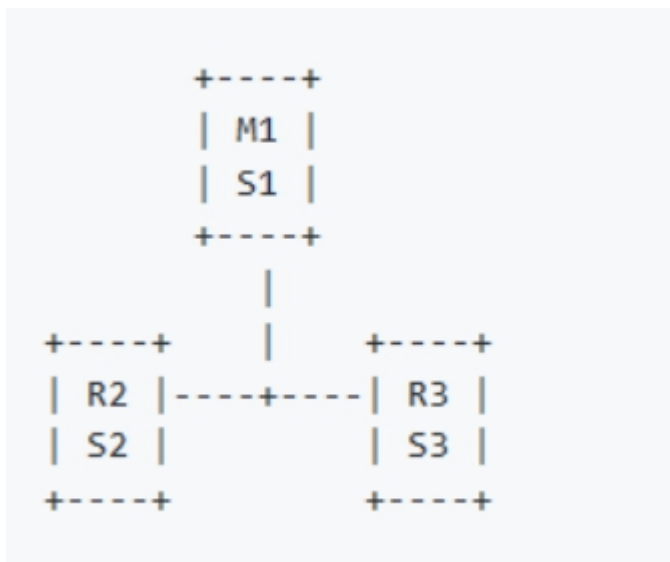
配置 **quorum=1**，如果 **master** 宕机，**s1** 和 **s2** 中只要有 1 个哨兵认为 **master** 宕机了，就可以进行切换，同时 **s1** 和 **s2** 会选举出一个哨兵来执行故障转移。但是同时这个时候，需要 **majority**，也就是大多数哨兵都是运行的。

```

2 个哨兵, majority=2
3 个哨兵, majority=2
4 个哨兵, majority=2
5 个哨兵, majority=3
...
```

如果此时仅仅是 **M1** 进程宕机了，哨兵 **s1** 正常运行，那么故障转移是 **OK** 的。但是如果是整个 **M1** 和 **S1** 运行的机器宕机了，那么哨兵只有 1 个，此时就没有 **majority** 来允许执行故障转移，虽然另外一台机器上还有一个 **R1**，但是故障转移不会执行。

经典的 3 节点哨兵集群是这样的：



配置 `quorum=2`，如果 `M1` 所在机器宕机了，那么三个哨兵还剩下 2 个，`S2` 和 `S3` 可以一致认为 `master` 宕机了，然后选举出一个来执行故障转移，同时 3 个哨兵的 `majority` 是 2，所以还剩下的 2 个哨兵运行着，就可以允许执行故障转移。

4.3 redis 哨兵主备切换的数据丢失问题

4.3.1 两种情况导致数据丢失

主备切换的过程，可能会导致数据丢失：

- 异步复制导致的数据丢失

因为 `master->slave` 的复制是异步的，所以可能有部分数据还没复制到 `slave`，`master` 就宕机了，此时这部分数据就丢失了。

- 脑裂导致的数据丢失

脑裂，也就是说，某个 `master` 所在机器突然脱离了正常的网络，跟其他 `slave` 机器不能连接，但是实际上 `master` 还运行着。此时哨兵可能就会认为 `master` 宕机了，然后开启选举，将其他 `slave` 切换成了 `master`。这个时候，集群里就会有二个 `master`，也就是所谓的脑裂。

此时虽然某个 `slave` 被切换成了 `master`，但是可能 `client` 还没来得及切换到新的 `master`，还继续向旧 `master` 写数据。因此旧 `master` 再次恢复的时候，会被作为一个 `slave` 挂到新的 `master` 上去，自己的数据会清空，重新从新的 `master` 复制数据。而新的 `master` 并没有后来 `client` 写入的数据，因此，这部分数据也就丢失了。

4.3.2 数据丢失问题的解决方案

进行如下配置：

```
min-slaves-to-write 1
min-slaves-max-lag 10
```

表示，要求至少有 1 个 `slave`，数据复制和同步的延迟不能超过 10 秒。

如果说一旦所有的 `slave`，数据复制和同步的延迟都超过了 10 秒钟，那么这个时候，`master` 就不会再接收任何请求了。

- 减少异步复制数据的丢失

有了 `min-slaves-max-lag` 这个配置，就可以确保说，一旦 `slave` 复制数据和 `ack` 延时太长，就认为可能 `master` 宕机后损失的数据太多了，那么就拒绝写请求，这样可以把 `master` 宕机时由于部分数据未同步到 `slave` 导致的数据丢失降低的可控范围内。

- 减少脑裂的数据丢失

如果一个 `master` 出现了脑裂，跟其他 `slave` 丢了连接，那么上面两个配置可以确保说，如果不能继续给指定数量的 `slave` 发送数据，而且 `slave` 超过 10 秒没有给自己 `ack` 消息，那么就直接拒绝客户端的写请求。因此在脑裂场景下，最多就丢失 10 秒的数据。

4.4 sdown 和 odown 转换机制

- `sdown` 是主观宕机，就一个哨兵如果自己觉得一个 `master` 宕机了，那么就是主观宕机
- `odown` 是客观宕机，如果 `quorum` 数量的哨兵都觉得一个 `master` 宕机了，那么就是客观宕机

`sdown` 达成的条件很简单，如果一个哨兵 `ping` 一个 `master`，超过了 `is-master-down-after-milliseconds` 指定的毫秒数之后，就主观认为 `master` 宕机了；如果一个哨兵在规定时间内，收到了 `quorum` 数量的其它哨兵也认为那个 `master` 是 `sdown` 的，那么就认为是 `odown` 了。

4.5 哨兵集群的自动发现机制

哨兵互相之间的发现，是通过 `redis` 的 `pub/sub`(发布/订阅) 系统实现的，每个哨兵都会往 `sentinel:hello` 这个 `channel` 里发送一个消息，这时候所有其他哨兵都可以消费到这个消息，并感知到其他的哨兵的存在。

每隔两秒钟，每个哨兵都会往自己监控的某个 `master+slaves` 对应的 `sentinel:hello` `channel` 里发送一个消息，内容是自己的 `host`、`ip` 和 `runid` 还有对这个 `master` 的监控配置。

每个哨兵也会去监听自己监控的每个 `master+slaves` 对应的 `sentinel:hello` `channel`，然后去感知到同样在监听这个 `master+slaves` 的其他哨兵的存在。

每个哨兵还会跟其他哨兵交换对 `master` 的监控配置，互相进行监控配置的同步。

4.5 slave 配置的自动纠正

哨兵会负责自动纠正 `slave` 的一些配置，比如 `slave` 如果要成为潜在的 `master` 候选人，哨兵会确保 `slave` 复制现有 `master` 的数据；如果 `slave` 连接到了一个错误的 `master` 上，比如故障转移之后，那么哨兵会确保它们连接到正确的 `master` 上。

4.6 slave->master 选举算法

如果一个 `master` 被认为 `odown` 了，而且 `majority` 数量的哨兵都允许主备切换，那么某个哨兵就会执行主备切换操作，此时首先要选举一个 `slave` 来，会考虑 `slave` 的一些信息：

- 跟 `master` 断开连接的时长
- `slave` 优先级
- 复制 `offset`
- `run id`

如果一个 `slave` 跟 `master` 断开连接的时间已经超过了 `down-after-milliseconds` 的 10 倍，外加 `master` 宕机的时长，那么 `slave` 就被认为不适合选举为 `master`。

$(\text{down-after-milliseconds} * 10) + \text{milliseconds_since_master_is_in_SDOWN_state}$

接下来会对 `slave` 进行排序：

- 按照 `slave` 优先级进行排序，`slave priority` 越低，优先级就越高。
- 如果 `slave priority` 相同，那么看 `replica offset`，哪个 `slave` 复制了越多的数据，`offset` 越靠后，优先级就越高。
- 如果上面两个条件都相同，那么选择一个 `run id` 比较小的那个 `slave`。

4.7 quorum 和 majority

quorum: 确认odown的最少的哨兵数量

majority: 授权进行主从切换的最少的哨兵数量

每次一个哨兵要做主备切换，首先需要 `quorum` 数量的哨兵认为 `odown`，然后选举出一个哨兵来做切换，这个哨兵还需要得到 `majority` 哨兵的授权，才能正式执行切换。

如果 `quorum < majority`，比如 5 个哨兵，`majority` 就是 3，`quorum` 设置为 2，那么就 3 个哨兵授权就可以执行切换。

但是如果 `quorum >= majority`，那么必须 `quorum` 数量的哨兵都授权，比如 5 个哨兵，`quorum` 是 5，那么必须 5 个哨兵都同意授权，才能执行切换。

5.Redis Cluster

5.1 redis cluster 介绍

- 自动将数据进行分片，每个 `master` 上放一部分数据
- 提供内置的高可用支持，部分 `master` 不可用时，还是可以继续工作的

在 `redis cluster` 架构下，每个 `redis` 要放开两个端口号，比如一个是 6379，另外一个就是加1w的端口号，比如 16379。

16379 端口号是用来进行节点间通信的，也就是 `cluster bus` 的东西，`cluster bus` 的通信，用来进行故障检测、配置更新、故障转移授权。`cluster bus` 用了另外一种二进制的协议，`gossip` 协议，用于节点间进行高效的数据交换，占用更少的网络带宽和处理时间。

5.2 使用场景

`redis cluster`，主要是针对**海量数据+高并发+高可用**的场景。`redis cluster` 支撑 N 个 `redis master node`，每个 `master node` 都可以挂载多个 `slave node`。这样整个 `redis` 就可以横向扩容了。如果你要支撑更大数据量的缓存，那就横向扩容更多的 `master` 节点，每个 `master` 节点就能存放更多的数据了。

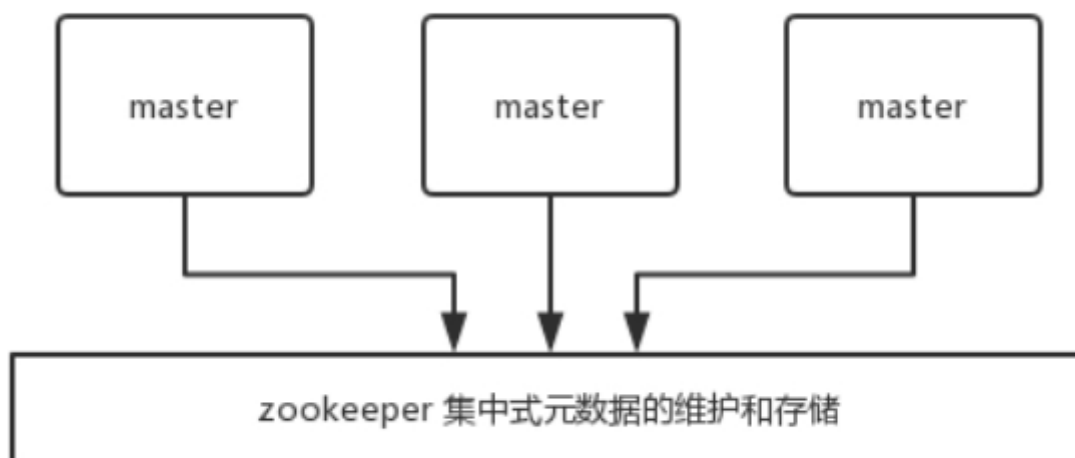
如果你的数据量很少，主要是承载高并发高性能的场景，比如你的缓存一般就几个 G，单机就足够了，可以使用 `replication`，一个 `master` 多个 `slaves`，要几个 `slave` 跟你要求的读吞吐量有关，然后自己搭建一个 `sentinel` 集群去保证 `redis` 主从架构的高可用性。

5.3 节点间的内部通信机制

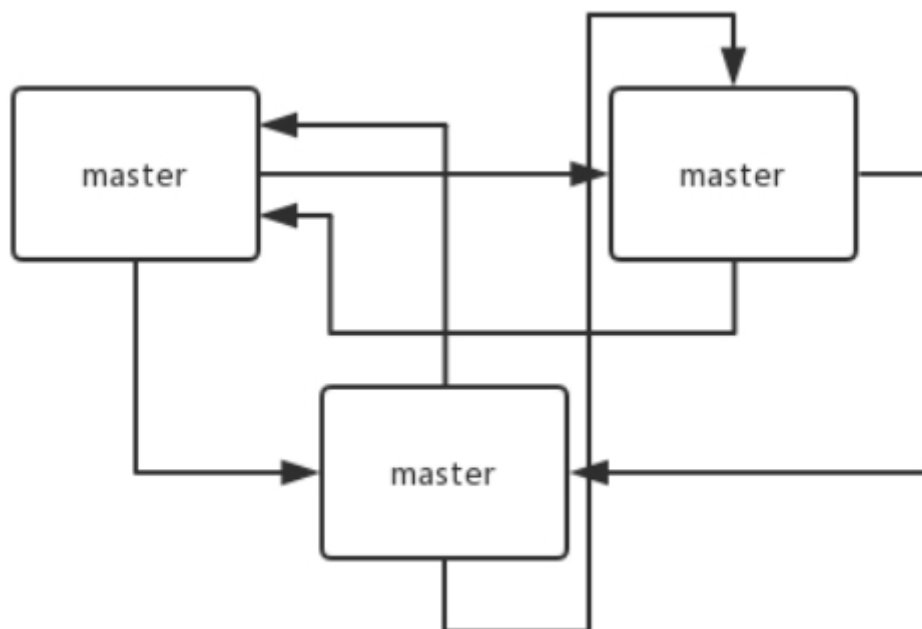
5.3.1 基本通信原理

集群元数据的维护有两种方式：集中式、`Gossip` 协议。`redis cluster` 节点间采用 `gossip` 协议进行通信。

集中式是将集群元数据（节点信息、故障等等）集中存储在某个节点上。集中式元数据集中存储的一个典型代表，就是大数据领域的 **storm**。它是分布式的大数据实时计算引擎，是集中式的元数据存储的结构，底层基于 **zookeeper**（分布式协调的中间件）对所有元数据进行存储维护。



redis 维护集群元数据采用另一个方式，**gossip** 协议，所有节点都持有一份元数据，不同的节点如果出现了元数据的变更，就不断将元数据发送给其它的节点，让其它节点也进行元数据的变更。



5.3.1.1 gossip 协议

gossip 协议包含多种消息，包含 **ping,pong,meet,fail** 等等。

- **meet**: 某个节点发送 **meet** 给新加入的节点，让新节点加入集群中，然后新节点就会开始与其它节点进行通信。

redis-trib.rb add-node

其实内部就是发送了一个 `gossip meet` 消息给新加入的节点，通知那个节点去加入我们的集群。

- **ping**: 每个节点都会频繁给其它节点发送 **ping**，其中包含自己的状态还有自己维护的集群元数据，互相通过 **ping** 交换元数据。
- **pong**: 返回 **ping** 和 **meet**，包含自己的状态和其它信息，也用于信息广播和更新。
- **fail**: 某个节点判断另一个节点 **fail** 之后，就发送 **fail** 给其它节点，通知其它节点说，某个节点宕机了。

5.3.1.2 ping 消息深入

ping 时要携带一些元数据，如果很频繁，可能会加重网络负担。

每个节点每秒会执行 **10** 次 **ping**，每次会选择 **5** 个最久没有通信的其它节点。当然如果发现某个节点通信延时达到了 `cluster_node_timeout / 2`，那么立即发送 **ping**，避免数据交换延时过长，落后的时间太长了。比如说，两个节点之间都 **10** 分钟没有交换数据了，那么整个集群处于严重的元数据不一致的情况，就会有问题。所以 `cluster_node_timeout` 可以调节，如果调得比较大，那么会降低 **ping** 的频率。

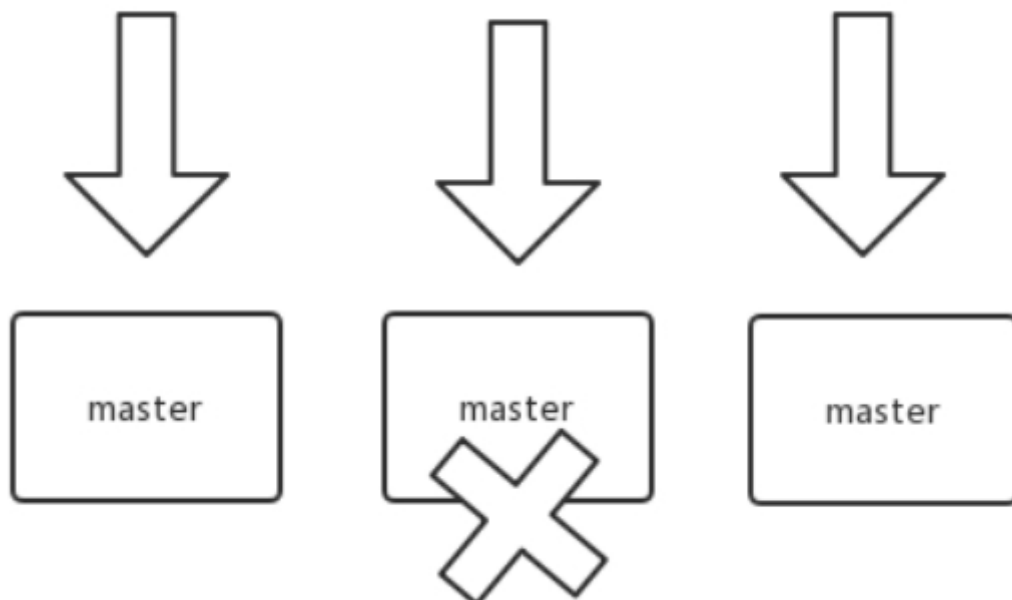
每次 **ping**，会带上自己节点的信息，还有就是带上 **1/10** 其它节点的信息，发送出去，进行交换。至少包含 **3** 个其它节点的信息，最多包含 总节点数减 **2** 个其它节点的信息。

5.3.2 分布式寻址算法

- **hash** 算法（大量缓存重建）
- 一致性 **hash** 算法（自动缓存迁移）+ 虚拟节点（自动负载均衡）
- **redis cluster** 的 **hash slot** 算法

5.3.2.1 hash 算法

来了一个 **key**，首先计算 **hash** 值，然后对节点数取模。然后打在不同的 **master** 节点上。一旦某一个 **master** 节点宕机，所有请求过来，都会基于最新的剩余 **master** 节点数去取模，尝试去取数据。这会导致 *大部分请求过来，全部无法拿到有效的缓存*，导致大量的流量涌入数据库。



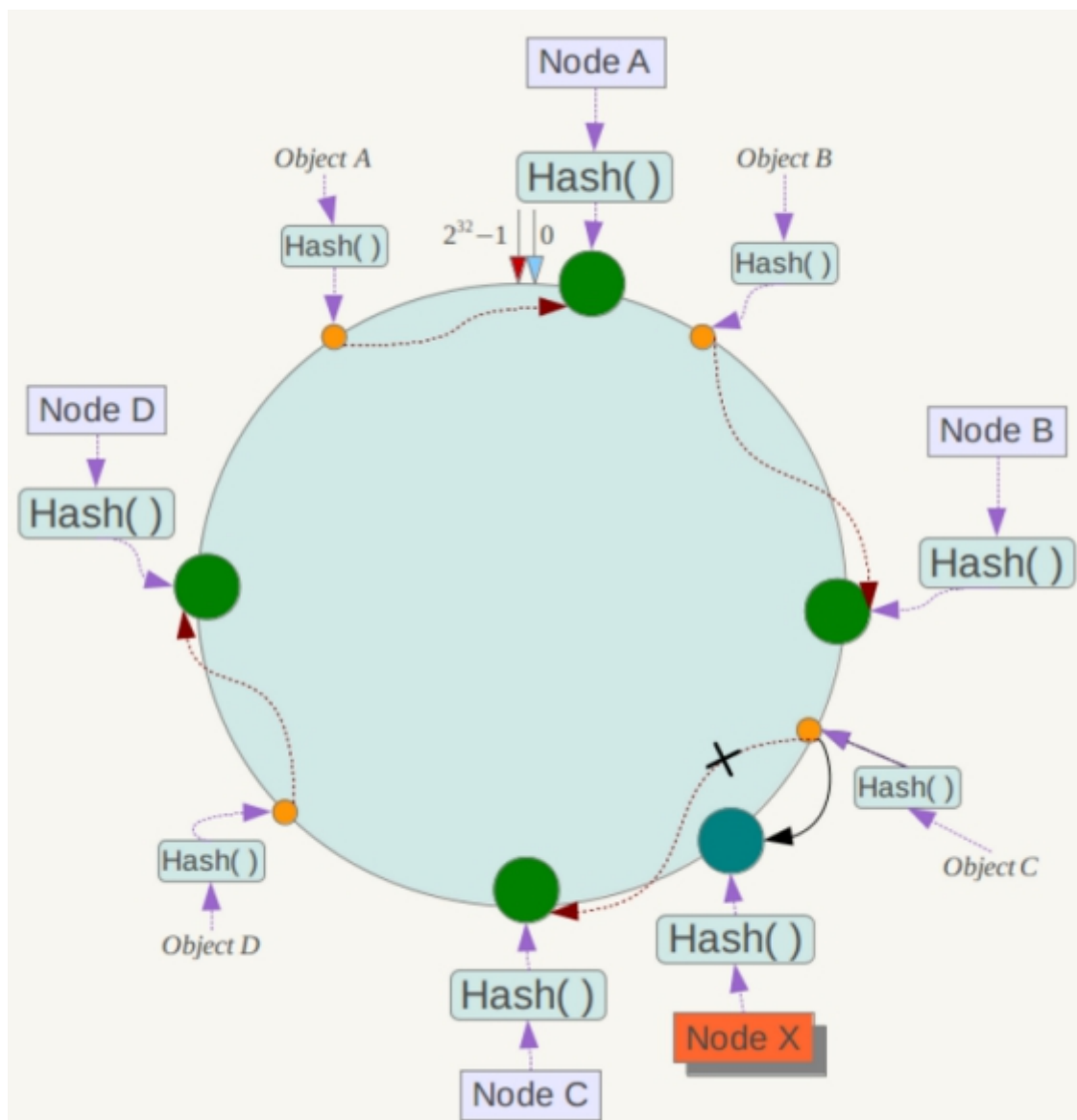
5.3.2.2 一致性 hash 算法

一致性 hash 算法将整个 hash 值空间组织成一个虚拟的圆环，整个空间按顺时针方向组织，下一步将各个 master 节点（使用服务器的 ip 或主机名）进行 hash。这样就能确定每个节点在其哈希环上的位置。

来了一个 key，首先计算 hash 值，并确定此数据在环上的位置，从此位置沿环顺时针“行走”，遇到的第一个 master 节点就是 key 所在位置。

在一致性哈希算法中，如果一个节点挂了，受影响的数据仅仅是此节点到环空间前一个节点（沿着逆时针方向行走遇到的第一个节点）之间的数据，其它不受影响。增加一个节点也同理。

然而，一致性哈希算法在节点太少时，容易因为节点分布不均匀而造成缓存热点的问题。为了解决这种热点问题，一致性 hash 算法引入了虚拟节点机制，即对每一个节点计算多个 hash，每个计算结果位置都放置一个虚拟节点。这样就实现了数据的均匀分布，负载均衡。

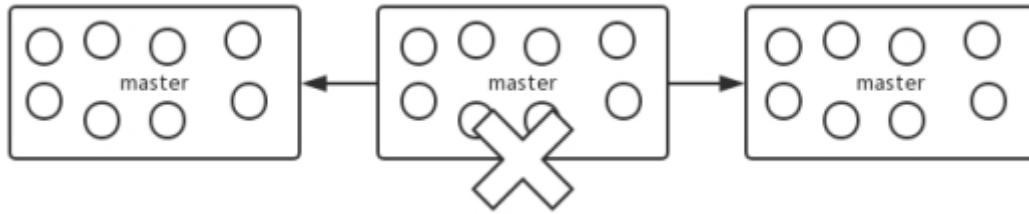


5.3.2.3 redis cluster 的 hash slot 算法

redis cluster 有固定的 16384 个 hash slot，对每个 key 计算 CRC16 值，然后对 16384 取模，可以获得 key 对应的 hash slot。

redis cluster 中每个 master 都会持有部分 slot，比如有 3 个 master，那么可能每个 master 持有 5000 多个 hash slot。hash slot 让 node 的增加和移除很简单，增加一个 master，就将其他 master 的 hash slot 移动部分过去，减少一个 master，就将它的 hash slot 移动到其他 master 上去。移动 hash slot 的成本是非常低的。客户端的 api，可以对指定的数据，让他们走同一个 hash slot，通过 hash tag 来实现。

任何一台机器宕机，另外两个节点，不影响的。因为 key 找的是 hash slot，不是机器。



5.4 redis cluster 的高可用与主备切换原理

redis cluster 的高可用的原理，几乎跟哨兵是类似的。

5.4.1 判断节点宕机

如果一个节点认为另外一个节点宕机，那么就是 **pfail**，**主观宕机**。如果多个节点都认为另外一个节点宕机了，那么就是 **fail**，**客观宕机**，跟哨兵的原理几乎一样，**sdown**，**odown**。

在 **cluster-node-timeout** 内，某个节点一直没有返回 **pong**，那么就被认为 **pfail**。

如果一个节点认为某个节点 **pfail** 了，那么会在 **gossip ping** 消息中，**ping** 给其他节点，如果**超过半数**的节点都认为 **pfail** 了，那么就会变成 **fail**。

5.4.2 从节点过滤

对宕机的 master node，从其所有的 slave node 中，选择一个切换成 master node。

检查每个 slave node 与 master node 断开连接的时间，如果超过了 **cluster-node-timeout * cluster-slave-validity-factor**，那么就**没有资格**切换成 master。

5.4.3 从节点选举

每个从节点，都根据自己对 master 复制数据的 **offset**，来设置一个选举时间，**offset** 越大（复制数据越多）的从节点，选举时间越靠前，优先进行选举。

所有的 master node 开始 slave 选举投票，给要进行选举的 slave 进行投票，如果大部分 master node ($N/2 + 1$) 都投票给了某个从节点，那么选举通过，那个从节点可以切换成 master。

从节点执行主备切换，从节点切换为主节点。

5.4.4 与哨兵比较

整个流程跟哨兵相比，非常类似，所以说，redis cluster 功能强大，直接集成了 replication 和 sentinel 的功能。