

# 一、数据类型

## 1.1 String

### 介绍

String 是最基本的 key-value 结构，key 是唯一标识是String类型，value 是具体的值，value其实不仅是字符串，也可以是数字（整数或浮点数），value 最多可以容纳的数据长度是 **512M**。

### 内部实现

String 类型的底层的数据结构实现主要是 int 和 SDS（简单动态字符串）。

SDS 相比于 C 的原生字符串：

- **SDS 不仅可以保存文本数据，还可以保存二进制数据。**所以 SDS 不光能存放文本数据，而且能保存图片、音频、视频、压缩文件这样的二进制数据。
- **SDS 获取字符串长度的时间复杂度是  $O(1)$ 。** SDS 结构里用 **len** 属性记录了字符串长度，所以复杂度为  **$O(1)$** 。
- **Redis 的 SDS API 是安全的，拼接字符串不会造成缓冲区溢出。**因为 SDS 在拼接字符串之前会检查 SDS 空间是否满足要求，如果空间不够会自动扩容，所以不会导致缓冲区溢出的问题。

字符串对象的内部编码（encoding）有 3 种：**int**、**raw**和 **embstr**。

如果一个字符串对象保存的是**整数值**，并且这个整数值可以用**long**类型来表示，那么字符串对象会将整数值保存在字符串对象结构的**ptr**属性里面（将**void\***转换成 long），并将字符串对象的编码设置为**int**。

如果字符串对象保存的是一个**字符串**，字符串对象将使用一个简单动态字符串（SDS）来保存这个字符串，如果这个字符串的长度小于等于某个边界值，那么，并将对象的编码设置为**embstr**否则使用**raw**

### 常用指令

```
#设置 key-value 类型的值
set name zs
#根据 key 获得对应的 value
get name
#批量设置 key-value 类型的值
mset name zs age 18
#批量获取多个 key 对应的 value
mget name age
```

```

#判断某个 key 是否存在
exists name
#返回 key 所储存的字符串值的长度
strlen name
#整数值自增
incr age

#删除某个 key 对应的值
del name

#设置key在60秒后过期
expire name 60
#设置 key-value 类型的值, 并设置该key的过期时间(ex/px 秒/毫秒)
set name zs ex 60
#查看数据还有多久过期
ttl name
#插入数据时设置过期时间
set name 10 zs

```

## 应用场景

### 缓存对象

使用 String 来缓存对象有两种方式:

- 直接缓存整个对象的 JSON, 命令例子: `SET user:1 '{"name":"xiaolin", "age":18}'`。
- 采用将 key 进行分离为 user:ID:属性, 采用 MSET 存储, 用 MGET 获取各属性值, 命令例子: `MSET user:1:name xiaolin user:1:age 18 user:2:name xiaomei user:2:age 20`。

### 常规计数

因为 Redis 处理命令是单线程, 所以执行命令的过程是原子的。因此 String 数据类型适合计数场景, 比如计算访问次数、点赞、转发、库存数量等等。

```

#初始化文章的阅读量
> SET aritcle:readcount:1001 0
OK
#阅读量+1
> INCR aritcle:readcount:1001
1
# 获取对应文章的阅读量
> GET aritcle:readcount:1001
"2"

```

## 分布式锁

SET 命令有个 NX 参数可以实现「key不存在才插入」，可以用它来实现分布式锁：

- 如果 key 不存在，则显示插入成功，可以用来表示加锁成功；
- 如果 key 存在，则会显示插入失败，可以用来表示加锁失败。

一般而言，还会对分布式锁加上过期时间，分布式锁的命令如下：

```
SET lock_key unique_value NX PX 10000
```

- lock\_key 就是 key 键；
- unique\_value 是客户端生成的唯一的标识；
- NX 代表只在 lock\_key 不存在时，才对 lock\_key 进行设置操作；
- PX 10000 表示设置 lock\_key 的过期时间为 10s，这是为了避免客户端发生异常而无法释放锁。

而解锁的过程就是将 lock\_key 键删除，但不能乱删，要保证执行操作的客户端就是加锁的客户端。所以，解锁的时候，我们要先判断锁的 unique\_value 是否为加锁客户端，是的话，才将 lock\_key 键删除。

可以看到，解锁是有两个操作，这时就需要 Lua 脚本来保证解锁的原子性，因为 Redis 在执行 Lua 脚本时，可以以原子性的方式执行，保证了锁释放操作的原子性。

```
// 释放锁时，先比较 unique_value 是否相等，避免锁的误释放
if redis.call("get",KEYS[1]) == ARGV[1] then
    return redis.call("del",KEYS[1])
else
    return 0
end
```

这样一来，就通过使用 SET 命令和 Lua 脚本在 Redis 单节点上完成了分布式锁的加锁和解锁。

## 共享-session-信息共享 Session 信息

通常我们在开发后台管理系统时，会使用 Session 来保存用户的会话(登录)状态，这些 Session 信息会被保存在服务器端，但这只适用于单系统应用，如果是分布式系统此模式将不再适用。

例如用户一的 Session 信息被存储在服务器一，但第二次访问时用户一被分配到服务器二，这个时候服务器并没有用户一的 Session 信息，就会出现需要重复登录的问题，问题在于分布式系统每次会把请求随机分配到不同的服务器。

## 1.2 List

### 介绍

List 列表是简单的字符串列表，**按照插入顺序排序**，可以从头部或尾部向 List 列表添加元素。列表的最大长度为  $2^{32} - 1$ ，也即每个列表支持超过 40 亿个元素。

### 内部实现

List 类型的底层数据结构是由**双向链表或压缩列表**实现的：

- 如果列表的元素个数小于 512 个（默认值，可由 `list-max-ziplist-entries` 配置），列表每个元素的值都小于 64 字节（默认值，可由 `list-max-ziplist-value` 配置），Redis 会使用**压缩列表**作为 List 类型的底层数据结构；
- 如果列表的元素不满足上面的条件，Redis 会使用**双向链表**作为 List 类型的底层数据结构；

但是在 Redis 3.2 版本之后，**List 数据类型底层数据结构就只由 quicklist 实现了**，替代了双向链表和压缩列表。

### 常用指令

```
#将多个值插入到key列表
lpush city a b c
rpush city d e f
#移除并返回key列表的头元素/尾元素
lpop city
rpop city

#返回列表指定区间内的元素
lrange city 0 -1

#删除列表
del city

# 从key列表表头弹出一个元素，没有就阻塞timeout秒，如果timeout=0则一直阻塞
BLPOP key [key ...] timeout
# 从key列表表尾弹出一个元素，没有就阻塞timeout秒，如果timeout=0则一直阻塞
BRPOP key [key ...] timeout
```

# 应用场景

## 消息队列

- 消息保序：使用 LPUSH + RPOP;
- 阻塞读取：使用 BRPOP;
- 重复消息处理：生产者自行实现全局唯一 ID;
- 消息的可靠性：使用 BRPOPLPUSH

消息队列在存取消息时，必须要满足三个需求，分别是**消息保序、处理重复的消息和保证消息可靠性**。

Redis 的 List 和 Stream 两种数据类型，就可以满足消息队列的这三个需求。我们先了解下基于 List 的消息队列实现方法，后面在介绍 Stream 数据类型时候，在详细说说 Stream。

### 1、如何满足消息保序需求？

- 生产者使用 `LPUSH key value[value...]` 将消息插入到队列的头部，如果 key 不存在则会创建一个空的队列再插入消息。
- 消费者使用 `RPOP key` 依次读取队列的消息，先进先出。

不过，在消费者读取数据时，有一个潜在的性能风险点。

在生产者往 List 中写入数据时，List 并不会主动地通知消费者有新消息写入，如果消费者想要及时处理消息，就需要在程序中不停地调用 `RPOP` 命令（比如使用一个 `while(1)` 循环）。如果有新消息写入，`RPOP` 命令就会返回结果，否则，`RPOP` 命令返回空值，再继续循环。

所以，即使没有新消息写入 List，消费者也要不停地调用 `RPOP` 命令，这就会导致消费者程序的 CPU 一直消耗在执行 `RPOP` 命令上，带来不必要的性能损失。

为了解决这个问题，Redis 提供了 `BRPOP` 命令。**`BRPOP` 命令也称为阻塞式读取，客户端在没有读到队列数据时，自动阻塞，直到有新的数据写入队列，再开始读取新数据。**和消费者程序自己不停地调用 `RPOP` 命令相比，这种方式能节省 CPU 开销。

### 2、如何处理重复的消息？

消费者要实现重复消息的判断，需要 2 个方面的要求：

- 每个消息都有一个全局的 ID。
- 消费者要记录已经处理过的消息的 ID。当收到一条消息后，消费者程序就可以对比收到的消息 ID 和记录的已处理过的消息 ID，来判断当前收到的消息有没有经过处理。如果已经处理过，那么，消费者程序就不再进行处理了。

但是 List 并不会为每个消息生成 ID 号，所以我们需要自行每个消息生成一个**全局唯一 ID**，生成之后，我们在用 `LPUSH` 命令把消息插入 List 时，需要在消息中包含这个全局唯一 ID。

例如，我们执行以下命令，就把一条全局 ID 为 111000102、库存量为 99 的消息插入了消息队列：

```
> LPUSH mq "111000102:stock:99"
(integer) 1
```

### 3、如何保证消息可靠性？

当消费者程序从 List 中读取一条消息后，List 就不会再留存这条消息了。所以，如果消费者程序在处理消息的过程出现了故障或宕机，就会导致消息没有处理完成，那么，消费者程序再次启动后，就没法再次从 List 中读取消息了。

为了留存消息，List 类型提供了 **BRPOPLPUSH** 命令，这个命令的作用是**让消费者程序从一个 List 中读取消息，同时，Redis 会把这个消息再插入到另一个 List（可以叫作备份 List）留存。**

这样一来，如果消费者程序读了消息但没能正常处理，等它重启后，就可以从备份 List 中重新读取消息并进行处理了。

好了，到这里可以知道基于 List 类型的消息队列，满足消息队列的三大需求（消息保序、处理重复的消息和保证消息可靠性）。

List 作为消息队列有什么缺陷？

**List 不支持多个消费者消费同一条消息**，因为一旦消费者拉取一条消息后，这条消息就从 List 中删除了，无法被其它消费者再次消费。

要实现一条消息可以被多个消费者消费，那么就要将多个消费者组成一个消费组，使得多个消费者可以消费同一条消息，但是 **List 类型并不支持消费组的实现。**

这就要说起 Redis 从 5.0 版本开始提供的 Stream 数据类型了，Stream 同样能够满足消息队列的三大需求，而且它还支持「消费组」形式的消息读取。

## 1.3 Hash

### 介绍

Hash 是一个键值对 (key - value) 集合，其中 value 的形式如：**value= [{field1, value1}, ... {fieldN, valueN}]**。Hash 特别适合用于存储对象。

### 内部实现

Hash 类型的底层数据结构是由**压缩列表或哈希表**实现的：

- 如果哈希类型元素个数小于 **512** 个（默认值，可由 **hash-max-ziplist-entries** 配置），所有值小于 **64** 字节（默认值，可由 **hash-max-ziplist-value** 配置）的话，Redis 会使用**压缩列表**作为 Hash 类型的底层数据结构；

- 如果哈希类型元素不满足上面条件，Redis 会使用**哈希表**作为 Hash 类型的 底层数据结构。

在 Redis 7.0 中，压缩列表数据结构已经废弃了，交由 listpack 数据结构来实现了。

## 常用指令

```
#存储一个键值
hset user1 name zs
hset user1 age 10
hget user1 name
#存储多个键值
hmset user1 name zs age 10
hmget user1 name age
#获取所有键值
hgetall user1
#删除hash表
hdel user1

hlen user1
hexists user1 name
# 为哈希表key中field键的值加上增量n
HINCRBY key field n
```

## 应用场景

### 缓存对象

Hash 类型的 (key, field, value) 的结构与对象的 (对象id, 属性, 值) 的结构相似，也可以用来存储对象。

```
hmset uid:1 name zs age 10
hget all uid:1
"name"
"zs"
"age"
10
```

### 购物车

以用户 id 为 key, 商品 id 为 field, 商品数量为 value, 恰好构成了购物车的3个要素

## 1.4 Set

### 介绍

Set 类型是一个无序并唯一的键值集合，它的存储顺序不会按照插入的先后顺序进行存储。

一个集合最多可以存储  $2^{32}-1$  个元素。概念和数学中个的集合基本类似，可以交集，并集，差集等等，所以 Set 类型除了支持集合内的增删改查，同时还支持多个集合取交集、并集、差集。

Set 类型和 List 类型的区别如下：

- List 可以存储重复元素，Set 只能存储非重复元素；
- List 是按照元素的先后顺序存储元素的，而 Set 则是无序方式存储元素的。

### 内部实现

Set 类型的底层数据结构是由**哈希表或整数集合**实现的：

- 如果集合中的元素都是整数且元素个数小于 512（默认值，`set-maxintset-entries`配置）个，Redis 会使用**整数集合**作为 Set 类型的底层数据结构；
- 如果集合中的元素不满足上面条件，则 Redis 使用**哈希表**作为 Set 类型的底层数据结构。

### 常用指令

```
# 往集合key中存入元素，元素存在则忽略，若key不存在则新建
sadd list a b c
# 从集合key中删除元素
srem list a
# 获取集合key中所有元素
smembers list
# 获取集合key中的元素个数
scard list

# 判断member元素是否存在于集合key中
sismember list a

# 从集合key中随机选出count个元素，元素不从key中删除
srandmember key [count]
# 从集合key中随机选出count个元素，元素从key中删除
spop key [count]

# 交集运算
sinter key [key ...]
```



```
# 将交集结果存入新集合destination中
sinterstore destination key [key ...]

# 并集运算
sunion key [key ...]
# 将并集结果存入新集合destination中
sunionstore destination key [key ...]

# 差集运算
sdiff key [key ...]
# 将差集结果存入新集合destination中
sdiffstore destination key [key ...]
```

## 应用场景

集合的主要几个特性，无序、不可重复、支持并交差等操作。

因此 Set 类型比较适合用来数据去重和保障数据的唯一性，还可以用来统计多个集合的交集、错集和并集等，当我们存储的数据是无序并且需要去重的情况下，比较适合使用集合类型进行存储。

但是要提醒你一下，这里有一个潜在的风险。**Set 的差集、并集和交集的计算复杂度较高，在数据量较大的情况下，如果直接执行这些计算，会导致 Redis 实例阻塞。**

在主从集群中，为了避免主库因为 Set 做聚合计算（交集、差集、并集）时导致主库被阻塞，我们可以选择一个从库完成聚合统计，或者把数据返回给客户端，由客户端来完成聚合统计。

### 点赞（没有重复数据）

Set 类型可以保证一个用户只能点一个赞

```
SADD article:1 uid:1
SADD article:1 uid:2
SADD article:1 uid:3
```

### 共同关注（集合运算）

Set 类型支持交集运算，所以可以用来计算共同关注的好友、公众号等。

key 可以是用户id，value 则是已关注的公众号的id。

```
# uid:1 用户关注公众号 id 为 5、6、7、8、9
> SADD uid:1 5 6 7 8 9
# uid:2 用户关注公众号 id 为 7、8、9、10、11
> SADD uid:2 7 8 9 10 11
# 获取共同关注
> SINTER uid:1 uid:2
1) "7"
2) "8"
3) "9"
# 给uid:2 推荐uid:1关注的公众号
> SDIFF uid:1 uid:2
1) "5"
2) "6"
```

## 抽奖活动

存储某活动中中奖的用户名，Set 类型因为有去重功能，可以保证同一个用户不会中奖两次。

key为抽奖活动名，value为员工名称，把所有员工名称放入抽奖箱：

```
sadd lucky a b c d

#抽取一个人
srandmember lucky 1
```

## 1.5 Zset

### 介绍

Zset 类型（有序集合类型）相比于 Set 类型多了一个排序属性 score（分值），对于有序集合 ZSet 来说，**每个存储元素相当于有两个值组成的，一个是有序集合的元素值，一个是排序值。**

有序集合保留了集合不能有重复成员的特性（分值可以重复），但不同的是，有序集合中的元素可以排序。

### 内部实现

Zset 类型的底层数据结构是由**压缩列表或跳表**实现的：

- 如果有序集合的元素个数小于 128 个，并且每个元素的值小于 64 字节时，Redis 会使用**压缩列表**作为 Zset 类型的底层数据结构；

- 如果有序集合的元素不满足上面的条件，Redis 会使用**跳表**作为 Zset 类型的底层数据结构；

在 Redis 7.0 中，压缩列表数据结构已经废弃了，交由 listpack 数据结构来实现了。

## 常用指令

```
# 往有序集合key中加入带分值元素
zadd list 90 zs 98 wu 33 ll
# 往有序集合key中删除元素
zrem list zs
# 返回有序集合key中元素member的分值
zscore list zs
# 返回有序集合key中元素个数
zcard list

# 为有序集合key中元素member的分值加上increment
zincrby key increment member

# 正序获取有序集合key从start下标到stop下标的元素
zrange key start stop [WITHSCORES]
# 倒序获取有序集合key从start下标到stop下标的元素
zrevrange key start stop [WITHSCORES]

# 返回有序集合中指定分数区间内的成员，分数由低到高排序。
zrangebyscore key min max [WITHSCORES] [LIMIT offset count]

# 返回指定成员区间内的成员，按字典正序排列，分数必须相同。
zrangebylex key min max [LIMIT offset count]
# 返回指定成员区间内的成员，按字典倒序排列，分数必须相同
zrevrangebylex key max min [LIMIT offset count]

# 并集计算(相同元素分值相加)，numberkeys一共多少个key，WEIGHTS每个key对应的分值乘积
ZUNIONSTORE destkey numberkeys key [key...]
# 交集计算(相同元素分值相加)，numberkeys一共多少个key，WEIGHTS每个key对应的分值乘积
ZINTERSTORE destkey numberkeys key [key...]
```

## 应用场景

Zset 类型 (Sorted Set, 有序集合) 可以根据元素的权重来排序, 我们可以自己来决定每个元素的权重值。比如说, 我们可以根据元素插入 Sorted Set 的时间确定权重值, 先插入的元素权重小, 后插入的元素权重大。

在面对需要展示最新列表、排行榜等场景时, 如果数据更新频繁或者需要分页显示, 可以优先考虑使用 Sorted Set。

## 排行榜

小林发表了五篇博文, 分别获得赞为 200、40、100、50、150。

```
# article:1 文章获得了200个赞
> ZADD user:xiaolin:ranking 200 article:1
# article:2 文章获得了40个赞
> ZADD user:xiaolin:ranking 40 article:2
# article:3 文章获得了100个赞
> ZADD user:xiaolin:ranking 100 article:3
# article:4 文章获得了50个赞
> ZADD user:xiaolin:ranking 50 article:4
# article:5 文章获得了150个赞
> ZADD user:xiaolin:ranking 150 article:5
```

可以使用 ZINCRBY 命令为文章5点赞

```
zincrby user:xiaolin:ranking 1 article:5
```

查看某篇文章的赞数, 可以使用 ZSCORE 命令

```
zscore user:xiaolin:ranking article:5
```

获取小林文章赞数最多的 3 篇文章, 可以使用 ZREVRANGE 命令

```
zrevrange user:xiaolin:ranking 0 2
1) "article:1"
2) "200"
3) "article:5"
4) "150"
5) "article:3"
6) "100"
```

获取小林 100 赞到 200 赞的文章, 可以使用 ZRANGEBYSCORE 命令

```
zrangebyscore user:xiaolin:ranking 100 200 withscores
1) "article:3"
2) "100"
3) "article:5"
4) "150"
5) "article:1"
6) "200"
```

## 电话、姓名排序

使用有序集合的 **ZRANGEBYLEX** 或 **ZREVRANGEBYLEX** 可以帮助我们实现电话号码或姓名的排序，我们以 **ZRANGEBYLEX**（返回指定成员区间内的成员，按 key 正序排列，分数必须相同）为例。

**注意：不要在分数不一致的 SortSet 集合中去使用 ZRANGEBYLEX和 ZREVRANGEBYLEX 指令，因为获取的结果会不准确。**

### 1、电话排序

我们可以将电话号码存储到 SortSet 中，然后根据需要来获取号段：

```
> ZADD phone 0 13100111100 0 13110114300 0 13132110901
(integer) 3
> ZADD phone 0 13200111100 0 13210414300 0 13252110901
(integer) 3
> ZADD phone 0 13300111100 0 13310414300 0 13352110901
(integer) 3
```

获取所有号码（按字典正序排列）

```
> ZRANGEBYLEX phone - +
1) "13100111100"
2) "13110114300"
3) "13132110901"
4) "13200111100"
5) "13210414300"
6) "13252110901"
7) "13300111100"
8) "13310414300"
9) "13352110901"
```

获取 132 号段的号码：

```
> ZRANGEBYLEX phone [132 (133
1) "13200111100"
2) "13210414300"
3) "13252110901"
```

获取132、133号段的号码：

```
> ZRANGEBYLEX phone [132 (134
1) "13200111100"
2) "13210414300"
3) "13252110901"
4) "13300111100"
5) "13310414300"
6) "13352110901"
```

## 2、姓名排序

```
> zadd names 0 Toumas 0 Jake 0 Bluetuo 0 Gaodeng 0 Aimini 0
Aidehua
(integer) 6
```

获取所有人的名字：

```
> ZRANGEBYLEX names - +
1) "Aidehua"
2) "Aimini"
3) "Bluetuo"
4) "Gaodeng"
5) "Jake"
6) "Toumas"
```

获取名字中大写字母A开头的所有人：

```
> ZRANGEBYLEX names [A (B
1) "Aidehua"
2) "Aimini"
```

获取名字中大写字母 C 到 Z 的所有人：

```
> ZRANGEBYLEX names [C [Z
1) "Gaodeng"
2) "Jake"
3) "Toumas"
```

## 1.6 BitMap

### 介绍

Bitmap，即位图，是一串连续的二进制数组（0和1），可以通过偏移量（offset）定位元素。BitMap通过最小的单位bit来进行0|1的设置，表示某个元素的值或者状态，时间复杂度为O(1)。

由于 bit 是计算机中最小的单位，使用它进行储存将非常节省空间，特别适合一些数据量大且使用**二值统计的场景**

### 内部实现

Bitmap 本身是用 String 类型作为底层数据结构实现的一种统计二值状态的数据类型。

String 类型是会保存为二进制的字节数组，所以，Redis 就把字节数组的每个 bit 位利用起来，用来表示一个元素的二值状态，你可以把 Bitmap 看作是一个 **bit 数组**。

### 常用指令

bitmap 基本操作：

```
# 设置值，其中value只能是 0 和 1
SETBIT key offset value

# 获取值
GETBIT key offset

# 获取指定范围内值为 1 的个数
# start 和 end 以字节为单位
BITCOUNT key start end
```

bitmap 运算操作：





如何统计这个月首次打卡时间呢？

Redis 提供了 `BITPOS key bitValue [start] [end]` 指令，返回数据表示 Bitmap 中第一个值为 `bitValue` 的 `offset` 位置。

在默认情况下，命令将检测整个位图，用户可以通过可选的 `start` 参数和 `end` 参数指定要检测的范围。所以我们可以执行这条命令来获取 `userID = 100` 在 2022 年 6 月份**首次打卡**日期：

```
BITPOS uid:sign:100:202206 1
```

需要注意的是，因为 `offset` 从 0 开始的，所以我们需要将返回的 `value + 1`。

## 判断用户登陆状态

Bitmap 提供了 `GETBIT`、`SETBIT` 操作，通过一个偏移值 `offset` 对 `bit` 数组的 `offset` 位置的 `bit` 位进行读写操作，需要注意的是 `offset` 从 0 开始。

只需要一个 `key = login_status` 表示存储用户登陆状态集合数据，将用户 ID 作为 `offset`，在线就设置为 1，下线设置 0。通过 `GETBIT` 判断对应的用户是否在线。5000 万用户只需要 6 MB 的空间。

假如我们要判断 `ID = 10086` 的用户的登陆情况：

第一步，执行以下指令，表示用户已登录。

```
SETBIT login_status 10086 1
```

第二步，检查该用户是否登陆，返回值 1 表示已登录。

```
GETBIT login_status 10086
```

第三步，登出，将 `offset` 对应的 `value` 设置成 0。

```
SETBIT login_status 10086 0
```

## 连续签到用户总数

如何统计出这连续 7 天连续打卡用户总数呢？

我们把每天的日期作为 Bitmap 的 `key`，`userID` 作为 `offset`，若是打卡则将 `offset` 位置的 `bit` 设置成 1。

`key` 对应的集合的每个 `bit` 位的数据则是一个用户在该日期的打卡记录。

一共有 7 个这样的 Bitmap，如果我们能对这 7 个 Bitmap 的对应的 `bit` 位做『与』运算。同样的 `UserID` `offset` 都是一样的，当一个 `userID` 在 7 个 Bitmap 对应对应的 `offset` 位置的 `bit = 1` 就说明该用户 7 天连续打卡。

结果保存到一个新 Bitmap 中，我们再通过 BITCOUNT 统计 bit = 1 的个数便得到了连续打卡 7 天的用户总数了。

Redis 提供了 BITOP operation destkey key [key ...] 这个指令用于对一个或者多个 key 的 Bitmap 进行位元操作。

- operation 可以是 and、OR、NOT、XOR。当 BITOP 处理不同长度的字符串时，较短的那个字符串所缺少的部分会被看作 0。空的 key 也被看作是包含 0 的字符串序列。

假设要统计 3 天连续打卡的用户数，则是将三个 bitmap 进行 AND 操作，并将结果保存到 destmap 中，接着对 destmap 执行 BITCOUNT 统计，如下命令：

```
# 与操作
BITOP AND destmap bitmap:01 bitmap:02 bitmap:03
# 统计 bit 位 = 1 的个数
BITCOUNT destmap
```

即使一天产生一个亿的数据，Bitmap 占用的内存也不大，大约占 12 MB 的内存 ( $10^8/8/1024/1024$ )，7 天的 Bitmap 的内存开销约为 84 MB。同时我们最好给 Bitmap 设置过期时间，让 Redis 删除过期的打卡数据，节省内存。

## 1.7 HyperLogLog

### 介绍

Redis HyperLogLog 是 Redis 2.8.9 版本新增的数据类型，是一种用于「统计基数」的数据集合类型，**基数统计就是指统计一个集合中不重复的元素个数**。但要注意，HyperLogLog 是统计规则是基于概率完成的，不是非常准确，标准误差率是 0.81%。

所以，简单来说 HyperLogLog **提供不精确的去重计数**。

HyperLogLog 的优点是，在输入元素的数量或者体积非常非常大时，计算基数所需的内存空间总是固定的、并且是很小的。

在 Redis 里面，**每个 HyperLogLog 键只需要花费 12 KB 内存，就可以计算接近  $2^{64}$  个不同元素的基数**，和元素越多就越耗费内存的 Set 和 Hash 类型相比，HyperLogLog 就非常节省空间。

这什么概念？举个例子给大家对比一下。

用 Java 语言来说，一般 long 类型占用 8 字节，而 1 字节有 8 位，即：1 byte = 8 bit，即 long 数据类型最大可以表示的数是： $2^{63}-1$ 。对应上面的  $2^{64}$  个数，假设此时有  $2^{63}-1$  这么多个数，从 0 ~  $2^{63}-1$ ，按照 long 以及 1k = 1024 字节的规则来计算内存总数，就是： $((2^{63}-1) * 8/1024)K$ ，这是很庞大的一个数，存储空间远远超过 12K，而 HyperLogLog 却可以用 12K 就能统计完。

## 常用指令

HyperLogLog 命令很少，就三个。

```
# 添加指定元素到 HyperLogLog 中
PFADD key element [element ...]

# 返回给定 HyperLogLog 的基数估算值。
PFCOUNT key [key ...]

# 将多个 HyperLogLog 合并为一个 HyperLogLog
PFMERGE destkey sourcekey [sourcekey ...]
```

## 应用场景

### 百万级网页 UV 计数

Redis HyperLogLog 优势在于只需要花费 12 KB 内存，就可以计算接近  $2^{64}$  个元素的基数，和元素越多就越耗费内存的 Set 和 Hash 类型相比，HyperLogLog 就非常节省空间。所以，非常适合统计百万级以上的网页 UV 的场景。

在统计 UV 时，你可以用 PFADD 命令（用于向 HyperLogLog 中添加新元素）把访问页面的每个用户都添加到 HyperLogLog 中。

```
PFADD page1:uv user1 user2 user3 user4 user5
```

接下来，就可以用 PFCOUNT 命令直接获得 page1 的 UV 值了，这个命令的作用就是返回 HyperLogLog 的统计结果。

```
PFCOUNT page1:uv
```

不过，有一点需要你注意一下，HyperLogLog 的统计规则是基于概率完成的，所以它给出的统计结果是有一定误差的，标准误差率是 0.81%。

这也就意味着，你使用 HyperLogLog 统计的 UV 是 100 万，但实际的 UV 可能是 101 万。虽然误差率不算大，但是，如果你需要精确统计结果的话，最好还是继续用 Set 或 Hash 类型。

## 1.8 GEO

### 介绍

Redis GEO 是 Redis 3.2 版本新增的数据类型，主要用于存储地理位置信息，并对存储的信息进行操作。

在日常生活中，我们越来越依赖搜索“附近的餐馆”、在打车软件上叫车，这些都离不开基于位置信息服务（Location-Based Service, LBS）的应用。LBS 应用访问的数据是和人或物关联的一组经纬度信息，而且要能查询相邻的经纬度范围，GEO 就非常适合应用在 LBS 服务的场景中。

## 内部实现

GEO 本身并没有设计新的底层数据结构，而是直接使用了 Sorted Set 集合类型。

GEO 类型使用 GeoHash 编码方法实现了经纬度到 Sorted Set 中元素权重分数的转换，这其中的两个关键机制就是「对二维地图做区间划分」和「对区间进行编码」。一组经纬度落在某个区间后，就用区间的编码值来表示，并把编码值作为 Sorted Set 元素的权重分数。

这样一来，我们就可以把经纬度保存到 Sorted Set 中，利用 Sorted Set 提供的“按权重进行有序范围查找”的特性，实现 LBS 服务中频繁使用的“搜索附近”的需求。

## 常用指令

```
# 存储指定的地理空间位置，可以将一个或多个经度(longitude)、纬度(latitude)、位置名称(member)添加到指定的 key 中。
GEOADD key longitude latitude member [longitude latitude member ...]

# 从给定的 key 里返回所有指定名称(member)的位置（经度和纬度），不存在的返回 nil。
GEOPOS key member [member ...]

# 返回两个给定位置之间的距离。
GEODIST key member1 member2 [m|km|ft|mi]

# 根据用户给定的经纬度坐标来获取指定范围内的地理位置集合。
GEORADIUS key longitude latitude radius m|km|ft|mi [WITHCOORD]
[WITHDIST] [WITHHASH] [COUNT count] [ASC|DESC] [STORE key]
[STOREDIST key]
```

## 应用场景

### 滴滴叫车

这里以滴滴叫车的场景为例，介绍下具体如何使用 GEO 命令：GEOADD 和 GEORADIUS 这两个命令。

假设车辆 ID 是 33，经纬度位置是 (116.034579, 39.030452)，我们可以用一个 GEO 集合保存所有车辆的经纬度，集合 key 是 cars:locations。

执行下面的这个命令，就可以把 ID 号为 33 的车辆当前经纬度位置存入 GEO 集合中：

```
GEOADD cars:locations 116.034579 39.030452 33
```

BS 应用执行下面的命令时，Redis 会根据输入的用户经纬度信息（116.054579, 39.030452），查找以这个经纬度为中心的 5 公里内的车辆信息，并返回给 LBS 应用。

```
GEOADIUS cars:locations 116.054579 39.030452 5 km ASC COUNT 10
```

## 1.9 Stream

### 介绍

Redis Stream 是 Redis 5.0 版本新增加的数据类型，Redis 专门为消息队列设计的数据类型。

在 Redis 5.0 Stream 没出来之前，消息队列的实现方式都有着各自的缺陷，例如：

- 发布订阅模式，不能持久化也就无法可靠的保存消息，并且对于离线重连的客户端不能读取历史消息的缺陷；
- List 实现消息队列的方式不能重复消费，一个消息消费完就会被删除，而且生产者需要自行实现全局唯一 ID。

基于以上问题，Redis 5.0 便推出了 Stream 类型也是此版本最重要的功能，用于完美地实现消息队列，它支持消息的持久化、支持自动生成全局唯一 ID、支持 ack 确认消息的模式、支持消费组模式等，让消息队列更加的稳定和可靠。

### 常用指令

Stream 消息队列操作命令：

- XADD：插入消息，保证有序，可以自动生成全局唯一 ID；
- XLEN：查询消息长度；
- XREAD：用于读取消息，可以按 ID 读取数据；
- XDEL：根据消息 ID 删除消息；
- DEL：删除整个 Stream；
- XRANGE：读取区间消息
- XREADGROUP：按消费组形式读取消息；
- XPENDING 和 XACK：
  - XPENDING 命令可以用来查询每个消费组内所有消费者「已读取、但尚未确认」的消息；
  - XACK 命令用于向消息队列确认消息处理已完成；

# 应用场景

## 消息队列

生产者通过 **XADD 命令插入一条消息**：

```
# * 表示让 Redis 为插入的数据自动生成一个全局唯一的 ID
# 往名称为 mymq 的消息队列中插入一条消息，消息的键是 name，值是 xiaolin
> XADD mymq * name xiaolin
"1654254953808-0"
```

插入成功后会返回全局唯一的 ID："1654254953808-0"。消息的全局唯一 ID 由两部分组成：

- 第一部分“1654254953808”是数据插入时，以毫秒为单位计算的当前服务器时间；
- 第二部分表示插入消息在当前毫秒内的消息序号，这是从 0 开始编号的。例如，“1654254953808-0”就表示在“1654254953808”毫秒内的第 1 条消息。

消费者通过 **XREAD 命令从消息队列中读取消息**时，可以指定一个消息 ID，并从这个消息 ID 的下一条消息开始进行读取（注意是输入消息 ID 的下一条信息开始读取，不是查询输入ID的消息）。

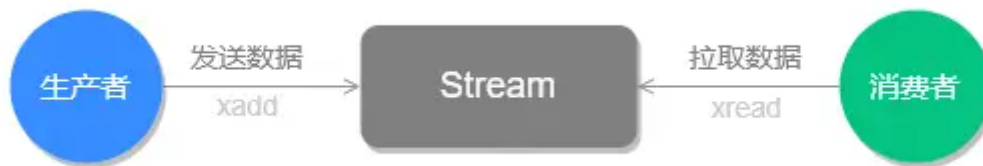
```
# 从 ID 号为 1654254953807-0 的消息开始，读取后续的所有消息（示例中一共 1 条）。
> XREAD STREAMS mymq 1654254953807-0
1) 1) "mymq"
   2) 1) 1) "1654254953808-0"
      2) 1) "name"
      2) "xiaolin"
```

如果想要实现阻塞读（当没有数据时，阻塞住），可以调用 **XREAD 时设定 BLOCK 配置项**，实现类似于 BRPOP 的阻塞读取操作。

比如，下面这命令，设置了 BLOCK 10000 的配置项，10000 的单位是毫秒，表明 XREAD 在读取最新消息时，如果没有消息到来，XREAD 将阻塞 10000 毫秒（即 10 秒），然后再返回。

```
# 命令最后的"$"符号表示读取最新的消息
> XREAD BLOCK 10000 STREAMS mymq $
(nil)
(10.00s)
```

Stream 的基础方法，使用 xadd 存入消息和 xread 循环阻塞读取消息的方式可以实现简易版的消息队列，交互流程如下图所示：



前面介绍的这些操作 List 也支持的，接下来看看 Stream 特有的功能。

Stream 可以使用 **XGROUP 创建消费组**，创建消费组之后，Stream 可以使用 XREADGROUP 命令让消费组内的消费者读取消息。

创建两个消费组，这两个消费组消费的消息队列是 mymq，都指定从第一条消息开始读取：

```
# 创建一个名为 group1 的消费组，0-0 表示从第一条消息开始读取。
> XGROUP CREATE mymq group1 0-0
OK
# 创建一个名为 group2 的消费组，0-0 表示从第一条消息开始读取。
> XGROUP CREATE mymq group2 0-0
OK
```

消费组 group1 内的消费者 consumer1 从 mymq 消息队列中读取所有消息的命令如下：

```
# 命令最后的参数">"，表示从第一条尚未被消费的消息开始读取。
> XREADGROUP GROUP group1 consumer1 STREAMS mymq >
1) 1) "mymq"
   2) 1) 1) "1654254953808-0"
      2) 1) "name"
      2) "xiaolin"
```

**消息队列中的消息一旦被消费组里的一个消费者读取了，就不能再被该消费组内的其他消费者读取了，即同一个消费组里的消费者不能消费同一条消息。**

比如说，我们执行完刚才的 XREADGROUP 命令后，再执行一次同样的命令，此时读到的就是空值了：

```
> XREADGROUP GROUP group1 consumer1 STREAMS mymq >
(nil)
```

但是，**不同消费组的消费者可以消费同一条消息（但是有前提条件，创建消息组的时候，不同消费组指定了相同位置开始读取消息）。**



比如说，刚才 group1 消费组里的 consumer1 消费者消费了一条 id 为 1654254953808-0 的消息，现在用 group2 消费组里的 consumer1 消费者消费消息：

```
> XREADGROUP GROUP group2 consumer1 STREAMS mymq >
1) 1) "mymq"
   2) 1) 1) "1654254953808-0"
      2) 1) "name"
         2) "xiaolin"
```

因为我创建两组的消费组都是从第一条消息开始读取，所以可以看到第二组的消费者依然可以消费 id 为 1654254953808-0 的这条消息。因此，不同的消费组的消费者可以消费同一条消息。

使用消费组的目的是让组内的多个消费者共同分担读取消息，所以，我们通常会让每个消费者读取部分消息，从而实现消息读取负载在多个消费者间是均衡分布的。

例如，我们执行下列命令，让 group2 中的 consumer1、2、3 各自读取一条消息。

```
# 让 group2 中的 consumer1 从 mymq 消息队列中消费一条消息
> XREADGROUP GROUP group2 consumer1 COUNT 1 STREAMS mymq >
1) 1) "mymq"
   2) 1) 1) "1654254953808-0"
      2) 1) "name"
         2) "xiaolin"

# 让 group2 中的 consumer2 从 mymq 消息队列中消费一条消息
> XREADGROUP GROUP group2 consumer2 COUNT 1 STREAMS mymq >
1) 1) "mymq"
   2) 1) 1) "1654256265584-0"
      2) 1) "name"
         2) "xiaolincoding"

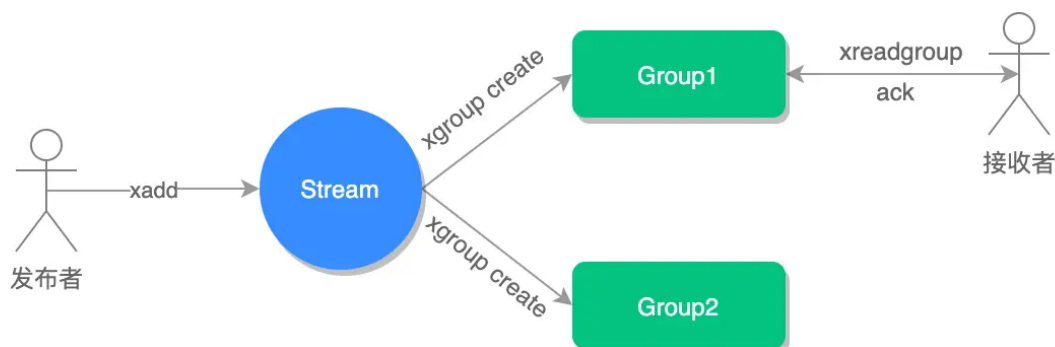
# 让 group2 中的 consumer3 从 mymq 消息队列中消费一条消息
> XREADGROUP GROUP group2 consumer3 COUNT 1 STREAMS mymq >
1) 1) "mymq"
   2) 1) 1) "1654256271337-0"
      2) 1) "name"
         2) "Tom"
```

基于 Stream 实现的消息队列，如何保证消费者在发生故障或宕机再次重启后，仍然可以读取未处理完的消息？

Streams 会自动使用内部队列（也称为 PENDING List）留存消费组里每个消费者读取的消息，直到消费者使用 XACK 命令通知 Streams“消息已经处理完成”。



消费确认增加了消息的可靠性，一般在业务处理完成之后，需要执行 XACK 命令确认消息已经被消费完成，整个流程的执行如下图所示：



如果消费者没有成功处理消息，它就不会给 Streams 发送 XACK 命令，消息仍然会留存。此时，**消费者可以在重启后，用 XPENDING 命令查看已读取、但尚未确认处理完成的消息。**

例如，我们来查看一下 group2 中各个消费者已读取、但尚未确认的消息个数，命令如下：

```
127.0.0.1:6379> XPENDING mymq group2
1) (integer) 3
2) "1654254953808-0" # 表示 group2 中所有消费者读取的消息最小 ID
3) "1654256271337-0" # 表示 group2 中所有消费者读取的消息最大 ID
4) 1) 1) "consumer1"
    2) "1"
    2) 1) "consumer2"
    2) "1"
    3) 1) "consumer3"
    2) "1"
```

如果想查看某个消费者具体读取了哪些数据，可以执行下面的命令：

```
# 查看 group2 里 consumer2 已从 mymq 消息队列中读取了哪些消息
> XPENDING mymq group2 - + 10 consumer2
1) 1) "1654256265584-0"
    2) "consumer2"
    3) (integer) 410700
    4) (integer) 1
```

可以看到，consumer2 已读取的消息的 ID 是 1654256265584-0。

**一旦消息 1654256265584-0 被 consumer2 处理了，consumer2 就可以使用 XACK 命令通知 Streams，然后这条消息就会被删除。**

```
> XACK mymq group2 1654256265584-0
(integer) 1
```

当我们再使用 XPENDING 命令查看时，就可以看到，consumer2 已经没有已读取、但尚未确认处理的消息了。

```
> XPENDING mymq group2 - + 10 consumer2
(empty array)
```

好了，基于 Stream 实现的消息队列就说到这里了，小结一下：

- 消息保序：XADD/XREAD
- 阻塞读取：XREAD block
- 重复消息处理：Stream 在使用 XADD 命令，会自动生成全局唯一 ID；
- 消息可靠性：内部使用 PENDING List 自动保存消息，使用 XPENDING 命令查看消费组已经读取但是未被确认的消息，消费者使用 XACK 确认消息；
- 支持消费组形式消费数据

Redis 基于 Stream 消息队列与专业的消息队列有哪些差距？

一个专业的消息队列，必须要做到两大块：

- 消息不丢。
- 消息可堆积。

### 1、Redis Stream 消息会丢失吗？

使用一个消息队列，其实就分为三大块：**生产者、队列中间件、消费者**，所以为了保证消息就是保证三个环节都不能丢失数据。



Redis Stream 消息队列能不能保证三个环节都不丢失数据？

- Redis 生产者会不会丢消息？生产者会不会丢消息，取决于生产者对于异常情况的处理是否合理。从消息被生产出来，然后提交给 MQ 的过程中，只要能正常收到（MQ 中间件）的 ack 确认响应，就表示发送成功，所以只要处理好返回值和异常，如果返回异常则进行消息重发，那么这个阶段是不会出现消息丢失的。

- Redis 消费者会不会丢消息？不会，因为 Stream（MQ 中间件）会自动使用内部队列（也称为 PENDING List）留存消费组里每个消费者读取的消息，但是未被确认的消息。消费者可以在重启后，用 X\_PENDING 命令查看已读取、但尚未确认处理完成的消息。等到消费者执行完业务逻辑后，再发送消费确认 XACK 命令，也能保证消息的不丢失。
- Redis 消息中间件会不会丢消息？

会

，Redis 在以下 2 个场景下，都会导致数据丢失：

- AOF 持久化配置为每秒写盘，但这个写盘过程是异步的，Redis 宕机时会存在数据丢失的可能
- 主从复制也是异步的，主从切换时，也存在丢失数据的可能（opens new window）。

可以看到，Redis 在队列中间件环节无法保证消息不丢。像 RabbitMQ 或 Kafka 这类专业的队列中间件，在使用时是部署一个集群，生产者在发布消息时，队列中间件通常会写「多个节点」，也就是有多个副本，这样一来，即便其中一个节点挂了，也能保证集群的数据不丢失。

## 2、Redis Stream 消息可堆积吗？

Redis 的数据都存储在内存中，这就意味着一旦发生消息积压，则会导致 Redis 的内存持续增长，如果超过机器内存上限，就会面临被 OOM 的风险。

所以 Redis 的 Stream 提供了可以指定队列最大长度的功能，就是为了避免这种情况发生。

当指定队列最大长度时，队列长度超过上限后，旧消息会被删除，只保留固定长度的新消息。这么来看，Stream 在消息积压时，如果指定了最大长度，还是有可能丢失消息的。

但 Kafka、RabbitMQ 专业的消息队列它们的数据都是存储在磁盘上，当消息积压时，无非就是多占用一些磁盘空间。

因此，把 Redis 当作队列来使用时，会面临的 2 个问题：

- Redis 本身可能会丢数据；
- 面对消息挤压，内存资源会紧张；

所以，能不能将 Redis 作为消息队列来使用，关键看你的业务场景：

- 如果你的业务场景足够简单，对于数据丢失不敏感，而且消息积压概率比较小的情况下，把 Redis 当作队列是完全可以的。
- 如果你的业务有海量消息，消息积压的概率比较大，并且不能接受数据丢失，那么还是用专业的消息队列中间件吧。

补充：Redis 发布/订阅机制为什么不可以作为消息队列？

发布订阅机制存在以下缺点，都是跟丢失数据有关：

1. 发布/订阅机制没有基于任何数据类型实现，所以不具备「数据持久化」的能力，也就是发布/订阅机制的相关操作，不会写入到 RDB 和 AOF 中，当 Redis 宕机重启，发布/订阅机制的数据也会全部丢失。
2. 发布订阅模式是“发后既忘”的工作模式，如果有订阅者离线重连之后不能消费之前的历史消息。
3. 当消费端有一定的消息积压时，也就是生产者发送的消息，消费者消费不过来时，如果超过 32M 或者是 60s 内持续保持在 8M 以上，消费端会被强行断开，这个参数是在配置文件中设置的，默认值是 `client-output-buffer-limit pubsub 32mb 8mb 60`。

所以，发布/订阅机制只适合即时通讯的场景，比如构建哨兵集群（`opens new window`）的场景采用了发布/订阅机制。

## 二、持久化

Redis 的读写操作都是在内存中，所以 Redis 性能才会高，但是当 Redis 重启后，内存中的数据就会丢失，那为了保证内存中的数据不会丢失，Redis 实现了数据持久化的机制，这个机制会把数据存储到磁盘，这样在 Redis 重启就能够从磁盘中恢复原有的数据。

Redis 共有三种数据持久化的方式：

- **AOF 日志**：每执行一条写操作命令，就把该命令以追加的方式写入到一个文件里；
- **RDB 快照**：将某一时刻的内存数据，以二进制的方式写入磁盘；
- **混合持久化方式**：Redis 4.0 新增的方式，集成了 AOF 和 RDB 的优点；

### 2.1 AOF持久化

#### 2.1.1 AOF日志

Redis 每执行一条写操作命令，就把该命令以追加的方式写入到AOF日志中

在 Redis 中 AOF 持久化功能默认是不开启的，需要我们修改 `redis.conf` 配置文件中的以下参数：

```
// redis.conf
appendonly yes // 表示是否开启AOF持久化(默认 no，关闭):
appendfilename "appendonly.aof" // AOF持久化文件的名称
```

AOF日志格式

set name xiaolin 日志格式为

```
*3
$3
set
$4
name
$7
xiaolin
```

「\*3」表示当前命令有三个部分，每部分都是以「\$+数字」开头，后面紧跟着具体的命令、键或值。然后，这里的「数字」表示这部分中的命令、键或值一共有多少字节。例如，「\$3 set」表示这部分有 3 个字节，也就是「set」命令这个字符串的长度。

先执行命令，再写AOF日志的好处

第一个好处，**避免额外的检查开销。**

因为如果先将写操作命令记录到 AOF 日志里，再执行该命令的话，如果当前的命令语法有问题，那么如果不进行命令语法检查，该错误的命令记录到 AOF 日志里后，Redis 在使用日志恢复数据时，就可能会出错。

第二个好处，**不会阻塞当前写操作命令的执行**，因为当写操作命令执行成功后，才会将命令记录到 AOF 日志。

AOF日志的潜在风险

第一个风险，执行写操作命令和记录日志是两个过程，那当 Redis 在还没来得及将命令写入到硬盘时，服务器发生宕机了，这个数据就会有**丢失的风险**。

第二个风险，前面说道，由于写操作命令执行成功后才记录到 AOF 日志，所以不会阻塞当前写操作命令的执行，但是**可能会给「下一个」命令带来阻塞风险**。

## 2.1.2 三种写回策略

Redis写入AOF日志的过程

1. Redis 执行完写操作命令后，会将命令追加到 `server.aof_buf` 缓冲区；
2. 然后通过 `write()` 系统调用，将 `aof_buf` 缓冲区的数据写入到 AOF 文件，此时数据并没有写入到硬盘，而是拷贝到了**内核缓冲区 page cache**，等待内核将数据写入硬盘；
3. 具体内核缓冲区的数据什么时候写入到硬盘，由内核决定。

三种写回策略

在 `redis.conf` 配置文件中的 `appendfsync` 配置项可以有以下 3 种参数可填：

- **Always**，这个单词的意思是「总是」，所以它的意思是每次**写操作命令执行完后，同步将 AOF 日志数据写回硬盘；主进程会阻塞**
- **Everysec**，这个单词的意思是「每秒」，所以它的意思是每次写操作命令执行完后，先将命令**写入到 AOF 文件的内核缓冲区，然后子线程异步执行每隔一秒将缓冲区里的内容写回到硬盘；**
- **No**，意味着不由 Redis 控制写回硬盘的时机，转交给操作系统控制写回的时机，也就是每次写操作命令执行完后，先将命令**写入到 AOF 文件的内核缓冲区，再由操作系统决定何时将缓冲区内容写回硬盘。**

这 3 种写回策略都无法完美解决「主进程阻塞」和「减少数据丢失」的问题，因为两个问题是对立的，偏向于一边的话，就会要牺牲另外一边，原因如下：

- Always 策略的话，可以最大程度保证数据不丢失，但是由于它每执行一条写操作命令就同步将 AOF 内容写回硬盘，所以是不可避免会影响主进程的性能；
- No 策略的话，是交由操作系统来决定何时将 AOF 日志内容写回硬盘，相比于 Always 策略性能较好，但是操作系统写回硬盘的时机是不可预知的，如果 AOF 日志内容没有写回硬盘，一旦服务器宕机，就会丢失不定数量的数据。
- Everysec 策略的话，是折中的一种方式，避免了 Always 策略的性能开销，也比 No 策略更能避免数据丢失，当然如果上一秒的写操作命令日志没有写回到硬盘，发生了宕机，这一秒内的数据自然也会丢失。

写回策略	写回时机	优点	缺点
Always	同步写回	可靠性高、最大程度保证数不丢失	每个写命令都要写回硬盘，性能开销大
Everysec	每秒写回	性能适中	宕机时会丢失1秒内的数据
No	由操作系统控制写回	性能好	宕机时丢失的数据可能会很多

### 2.1.3 AOF重写机制

AOF 日志是一个文件，随着执行的写操作命令越来越多，**AOF文件的大小会越来越大**。Redis 为了避免 AOF 文件越写越大，提供了 **AOF 重写机制**，当 AOF 文件的大小超过所设定的阈值后，Redis 就会**启用 AOF 重写机制，来压缩 AOF 文件**。

假设前后执行了「`set name xiaolin`」和「`set name xiaolincoding`」这两个命令的话，就会将这两个命令记录到 AOF 文件。但是**在使用重写机制后，就会读取 name 最新的 value（键值对），然后用一条「`set name xiaolincoding`」命令记录到新的 AOF 文件**。等到全部记录完后，就将新的 AOF 文件替换掉现有的 AOF 文件。这样一来，一个键值对在重写日志中只用一条命令就行了。

重写机制的妙处在于，尽管某个键值对被多条写命令反复修改，**最终也只需要根据这个「键值对」当前的最新状态，然后用一条命令去记录键值对**



为什么重写 AOF 的时候，不直接复用现有的 AOF 文件，而是先写到新的 AOF 文件再覆盖过去。

因为如果 AOF 重写过程中失败了，现有的 AOF 文件就会造成污染，可能无法用于恢复使用。

## 2.1.4 后台重写

重写比较耗时，所以重写的操作不能放在主进程里。

所以，Redis 的**重写AOF过程是由后台子进程 bgrewriteaof 来完成的**，这么做可以达到两个好处：

- 子进程进行 AOF 重写期间，主进程可以继续处理命令请求，从而避免阻塞主进程；
- 子进程带有主进程的数据副本（数据副本怎么产生的后面会说），这里使用子进程而不是线程，因为如果是使用线程，多线程之间会共享内存，那么在修改共享内存数据的时候，需要通过加锁来保证数据的安全，而这样就会降低性能。而使用子进程，创建子进程时，父子进程是共享内存数据的，不过这个共享的内存只能以只读的方式，而当父子进程任意一方修改了该共享内存，就会发生「写时复制」，于是父子进程就有了独立的数据副本，就不用加锁来保证数据安全。

### 写时复制

**父进程创建子进程时只会复制页表而共享同一块内存，只有在发生修改内存数据的情况时，物理内存才会被复制一份。**这时CPU就会触发**写保护中断**，这个写保护中断是由于违反权限导致的，然后操作系统会在「写保护中断处理函数」里进行**物理内存的复制**，并重新设置其内存映射关系，将父子进程的内存读写权限设置为**可读写**，最后才会对内存进行写操作，这个过程被称为「**写时复制(\*Copy On Write\*)**」。

**这样的目的是为了减少创建子进程时的性能损耗，从而加快创建子进程的速度，毕竟创建子进程的过程中，是会阻塞主线程的。**

### AOF 重写缓冲区

重写 AOF 日志过程中，如果主进程修改了已经存在 key-value，此时这个 key-value 数据在子进程的内存数据就跟主进程的内存数据不一致了，这时要怎么办呢？

为了解决这种数据不一致问题，Redis 设置了一个 **AOF 重写缓冲区**，这个缓冲区在创建bgrewriteaof 子进程之后开始使用。在重写 AOF 期间，当 Redis 执行完一个写命令之后，它会**同时将这个写命令写入到「AOF 缓冲区」和「AOF 重写缓冲区」**。

### 后台重写过程

主线程创建一个bgrewriteaof子进程执行AOF重写，创建子进程会复制页表阻塞一会。在发生写时复制时也会阻塞。

在 bgrewriteaof 子进程执行 AOF 重写期间，主进程需要执行以下三份工作：

- 执行客户端发来的命令；
- 将执行后的写命令追加到「AOF 缓冲区」；
- 将执行后的写命令追加到「AOF 重写缓冲区」；

子线程完成 AOF 重写工作（**扫描数据库中所有数据，逐一把内存数据的键值对转换成一条命令，再将命令记录到重写日志**）后，会向主进程发送一条信号，信号是进程间通讯的一种方式，且是异步的。

主进程收到该信号后，会调用一个**信号处理函数**，该函数主要做以下工作：

- 将 AOF 重写缓冲区中的所有内容追加到新的 AOF 的文件中，使得新旧两个 AOF 文件所保存的数据库状态一致；
- 新的 AOF 的文件进行改名，覆盖现有的 AOF 文件。

## 2.2 RDB快照

快照就是将redis中的数据写入到RDB文件中来实现持久化的，Redis 的快照是**全量快照**，也就是说每次执行快照，都是把内存中的「所有数据」都记录到磁盘中。**所以执行快照是一个比较重的操作，如果频率太频繁，可能会对 Redis 性能产生影响。如果频率太低，服务器故障时，丢失的数据会更多。**

Redis 提供了两个命令来生成 RDB 文件，分别是 **save** 和 **bgsave**，他们的区别就在于是否在「主线程」里执行：

- 执行了 **save** 命令，就会在主线程生成 RDB 文件，由于和执行操作命令在同一个线程，所以如果写入 RDB 文件的时间太长，**会阻塞主线程**；
- 执行了 **bgsave** 命令，**会创建一个子进程来生成 RDB 文件**，这样可以**避免主线程的阻塞**；

### 写时复制的情况

bgsave 快照过程中，如果主线程修改了共享数据，**发生了写时复制后，RDB 快照保存的是原本的内存数据**，而主线程刚修改的数据，是没办法在这一时间写入 RDB 文件的，只能交由下一次的 bgsave 快照。

## 2.3 RDB 和 AOF合体

### RDB的问题

尽管 RDB 比 AOF 的数据恢复速度快，但是快照的频率不好把握：

- 如果频率太低，两次快照间一旦服务器发生宕机，就可能会比较多的数据丢失；
- 如果频率太高，频繁写入磁盘和创建子进程会带来额外的性能开销。

如果想要开启混合持久化功能，可以在 Redis 配置文件将下面这个配置项设置成 yes：



```
aof-use-rdb-preamble yes
```

混合持久化工作在 **AOF 日志重写过程**。

当开启了混合持久化时，在 AOF 重写日志时，**fork** 出来的重写子进程会先将与主线程共享的内存数据以 RDB 方式写入到 AOF 文件，然后主线程处理的操作命令会被记录在重写缓冲区里，重写缓冲区里的增量命令会以 AOF 方式写入到 AOF 文件，写入完成后通知主进程将新的含有 RDB 格式和 AOF 格式的 AOF 文件替换旧的 AOF 文件。

也就是说，使用了混合持久化，AOF 文件的**前半部分是 RDB 格式的全量数据，后半部分是 AOF 格式的增量数据**。

这样的好处在于，重启 Redis 加载数据的时候，由于前半部分是 RDB 内容，这样**加载的时候速度会很快**。

加载完 RDB 的内容后，才会加载后半部分的 AOF 内容，这里的内容是 Redis 后台子进程重写 AOF 期间，主线程处理的操作命令，可以使得**数据更少的丢失**。

## 2.4 Redis大Key对持久化有什么影响？

当 AOF 写回策略配置了 Always 策略，如果写入是一个大 Key，主线程在执行 `fsync()` 函数的时候，阻塞的时间会比较久，因为当**写入的数据量很大的时候，数据同步到硬盘这个过程是很耗时的**。

AOF 重写机制和 RDB 快照 (`bgsave` 命令) 的过程，都会分别通过 `fork()` 函数创建一个子进程来处理任务。会有两个阶段会导致阻塞父进程（主线程）：

- 创建子进程的途中，由于要复制父进程的页表等数据结构，阻塞的时间跟页表的大小有关，**页表越大，阻塞的时间也越长**；
- 发生写时复制，这期间会拷贝物理内存，**由于大 Key 占用的物理内存会很大，那么在复制物理内存这一过程，就会比较耗时，所以有可能会阻塞父进程**。

大 key 除了会影响持久化之外，还会有以下的影响。

- 客户端超时阻塞。由于 Redis 执行命令是单线程处理，然后在操作大 key 时会比较耗时，那么就会阻塞 Redis，从客户端这一视角看，就是很久很久都没有响应。
- 引发网络阻塞。每次获取大 key 产生的网络流量较大，如果一个 key 的大小是 1 MB，每秒访问量为 1000，那么每秒会产生 1000MB 的流量，这对于普通千兆网卡的服务器来说是灾难性的。
- 阻塞工作线程。如果使用 `del` 删除大 key 时，会阻塞工作线程，这样就没办法处理后续的命令。
- 内存分布不均。集群模型在 slot 分片均匀情况下，会出现数据和查询倾斜情况，部分有大 key 的 Redis 节点占用内存多，QPS 也会比较大。

如何避免大 Key 呢？

最好在设计阶段，就把大 key 拆分成一个一个小 key。或者，定时检查 Redis 是否存在大 key，如果该大 key 是可以删除的，不要使用 DEL 命令删除，因为该命令删除过程会阻塞主线程，而是用 unlink 命令 (Redis 4.0+) 删除大 key，因为该命令的删除过程是异步的，不会阻塞主线程。

## 三、功能篇

### 3.1 过期删除策略

#### 3.1.1 过期时间命令

先说一下对 key 设置过期时间的命令。设置 key 过期时间的命令一共有 4 个：

- `expire <key> <n>`：设置 key 在 n 秒后过期，比如 `expire key 100` 表示设置 key 在 100 秒后过期；
- `pexpire <key> <n>`：设置 key 在 n 毫秒后过期，比如 `pexpire key2 100000` 表示设置 key2 在 100000 毫秒 (100 秒) 后过期。
- `expireat <key> <n>`：设置 key 在某个时间戳 (精确到秒) 之后过期，比如 `expireat key3 1655654400` 表示 key3 在时间戳 1655654400 后过期 (精确到秒)；
- `pexpireat <key> <n>`：设置 key 在某个时间戳 (精确到毫秒) 之后过期，比如 `pexpireat key4 1655654400000` 表示 key4 在时间戳 1655654400000 后过期 (精确到毫秒)

当然，在设置字符串时，也可以同时对 key 设置过期时间，共有 3 种命令：

- `set <key> <value> ex <n>`：设置键值对的时候，同时指定过期时间 (精确到秒)；
- `set <key> <value> px <n>`：设置键值对的时候，同时指定过期时间 (精确到毫秒)；
- `setex <key> <n> <value>`：设置键值对的时候，同时指定过期时间 (精确到秒)。

如果你想查看某个 key 剩余的存活时间，可以使用 `TTL <key>` 命令。

```
# 设置键值对的时候，同时指定过期时间位 60 秒
> setex key1 60 value1
OK

# 查看 key1 过期时间还剩多少
> ttl key1
(integer) 56
> ttl key1
(integer) 52
```

如果突然反悔，取消 key 的过期时间，则可以使用 `PERSIST <key>` 命令。

```
# 取消 key1 的过期时间
> persist key1
(integer) 1

# 使用完 persist 命令之后,
# 查下 key1 的存活时间结果是 -1, 表明 key1 永不过期
> ttl key1
(integer) -1
```

### 3.1.2 如何判定 key 已过期了?

Redis 会把该 key 带上过期时间存储到一个**过期字典** (expires dict) 中, 也就是说「过期字典」保存了数据库中所有 key 的过期时间。

字典实际上是哈希表, 哈希表的最大好处就是让我们可以用  $O(1)$  的时间复杂度来快速查找。当我们查询一个 key 时, Redis 首先检查该 key 是否存在于过期字典中:

- 如果不在, 则正常读取键值;
- 如果存在, 则会获取该 key 的过期时间, 然后与当前系统时间进行比对, 如果比系统时间大, 那就没有过期, 否则判定该 key 已过期

### 3.1.3 过期删除策略

#### 定时删除

定时删除策略的做法是, **在设置 key 的过期时间时, 同时创建一个定时事件, 当时间到达时, 由事件处理器自动执行 key 的删除操作。**

定时删除策略的**优点**:

- 可以保证过期 key 会被尽快删除, 也就是内存可以被尽快地释放。因此, 定时删除对内存是最友好的。

定时删除策略的**缺点**:

- 在过期 key 比较多的情况下, 删除过期 key 可能会占用相当一部分 CPU 时间, 在内存不紧张但 CPU 时间紧张的情况下, 将 CPU 时间用于删除和当前任务无关的过期键上, 无疑会对服务器的响应时间和吞吐量造成影响。所以, 定时删除策略对 CPU 不友好。

惰性删除策略的做法是, **不主动删除过期键, 每次从数据库访问 key 时, 都检测 key 是否过期, 如果过期则删除该 key。**

#### 惰性删除

惰性删除策略的**优点**:

- 因为每次访问时，才会检查 key 是否过期，所以此策略只会使用很少的系统资源，因此，惰性删除策略对 CPU 时间最友好。

惰性删除策略的**缺点**：

- 如果一个 key 已经过期，而这个 key 又仍然保留在数据库中，那么只要这个过期 key 一直没有被访问，它所占用的内存就不会释放，造成了一定的内存空间浪费。所以，惰性删除策略对内存不友好。

### 定期删除

定期删除策略的做法是，每隔一段时间「随机」从数据库中取出一定数量的 key 进行检查，并删除其中的过期key。

定期删除策略的**优点**：

- 通过限制删除操作执行的时长和频率，来减少删除操作对 CPU 的影响，同时也能删除一部分过期的数据减少了过期键对空间的无效占用。

定期删除策略的**缺点**：

- 内存清理方面没有定时删除效果好，同时没有惰性删除使用的系统资源少。
- 难以确定删除操作执行的时长和频率。如果执行的太频繁，定期删除策略变得和定时删除策略一样，对CPU不友好；如果执行的太少，那又和惰性删除一样了，过期 key 占用的内存不会及时得到释放。

### redis的过期删除策略

Redis 选择「惰性删除+定期删除」这两种策略配和使用，以求在合理使用 CPU 时间和避免内存浪费之间取得平衡。

## 3.2 内存淘汰策略

当 Redis 的运行内存已经超过 Redis 设置的最大内存之后，则会使用内存淘汰策略删除符合条件的 key，以此来保障 Redis 高效的运行。

### 3.2.1 如何设置 Redis 最大运行内存？

在配置文件 redis.conf 中，可以通过参数 `maxmemory <bytes>` 来设定最大运行内存，只有在 Redis 的运行内存达到了我们设置的最大运行内存，才会触发内存淘汰策略。不同位数的操作系统，maxmemory 的默认值是不同的：

- 在 64 位操作系统中，maxmemory 的默认值是 0，表示没有内存大小限制，那么不管用户存放多少数据到 Redis 中，Redis 也不会对可用内存进行检查，直到 Redis 实例因内存不足而崩溃也无作为。
- 在 32 位操作系统中，maxmemory 的默认值是 3G，因为 32 位的机器最大只支持 4GB 的内存，而系统本身就需要一定的内存资源来支持运行，所以 32 位操作系

统限制最大 3 GB 的可用内存是非常合理的，这样可以避免因为内存不足而导致 Redis 实例崩溃。

### 3.2.2 内存淘汰策略

Redis 内存淘汰策略共有八种，这八种策略大体分为「不进行数据淘汰」和「进行数据淘汰」两类策略。

#### 1、不进行数据淘汰的策略

**noeviction** (Redis 3.0 之后，默认的内存淘汰策略)：它表示当运行内存超过最大设置内存时，不淘汰任何数据，这时如果有新的数据写入，会报错通知禁止写入，不淘汰任何数据，但是如果没用数据写入的话，只是单纯的查询或者删除操作的话，还是可以正常工作。

#### 2、进行数据淘汰的策略

针对「进行数据淘汰」这一类策略，又可以细分为「在设置了过期时间的数据中进行淘汰」和「在所有数据范围内进行淘汰」这两类策略。

**在设置了过期时间的数据中进行淘汰：**

- **volatile-random**：随机淘汰设置了过期时间的任意键值；
- **volatile-ttl**：优先淘汰更早过期的键值。
- **volatile-lru** (Redis 3.0 之前，默认的内存淘汰策略)：淘汰所有设置了过期时间的键值中，最久未使用的键值；
- **volatile-lfu** (Redis 4.0 后新增的内存淘汰策略)：淘汰所有设置了过期时间的键值中，最少使用的键值；

**在所有数据范围内进行淘汰：**

- **allkeys-random**：随机淘汰任意键值；
- **allkeys-lru**：淘汰整个键值中最久未使用的键值；
- **allkeys-lfu** (Redis 4.0 后新增的内存淘汰策略)：淘汰整个键值中最少使用的键值。

使用 `config get maxmemory-policy` 命令，来查看当前 Redis 的内存淘汰策略，命令如下：

```
127.0.0.1:6379> config get maxmemory-policy
1) "maxmemory-policy"
2) "noeviction"
```

### 3.2.3 LRU和LFU算法

LRU 全称是 Least Recently Used 翻译为**最近最少使用**，会选择淘汰最近最少使用的数据。

传统 LRU 算法的实现是基于「链表」结构，链表中的元素按照操作顺序从前往后排列，最新操作的键会被移动到表头，当需要内存淘汰时，只需要删除链表尾部的元素即可，因为链表尾部的元素就代表最久未被使用的元素。

Redis 并没有使用这样的方式实现 LRU 算法，因为**传统的 LRU 算法存在两个问题**：

- 需要用链表管理所有的缓存数据，这会带来额外的空间开销；
- 当有数据被访问时，需要在链表上把该数据移动到头端，如果有大量数据被访问，就会带来很多链表移动操作，会很耗时，进而会降低 Redis 缓存性能。

Redis 实现的是一种**近似 LRU 算法**，目的是为了更好的节约内存，它的**实现方式是在 Redis 的对象结构体中添加一个额外的字段，用于记录此数据的最后一次访问时间**。

#### Redis实现的LRU算法的优缺点

不用为所有的数据维护一个大链表，节省了空间占用；

不用在每次数据访问时都移动链表项，提升了缓存的性能；

但是 LRU 算法有一个问题，**无法解决缓存污染问题**，比如应用一次读取了大量的数据，而这些数据只会被读取这一次，那么这些数据会留存在 Redis 缓存中很长一段时间，造成缓存污染

LFU 全称是 Least Frequently Used 翻译为**最近最不常用**，LFU 算法是根据数据访问次数来淘汰数据的，它的核心思想是“如果数据过去被访问多次，那么将来被访问的频率也更高”。**实现方式是在 Redis 的对象结构体中添加一个额外的字段，用于记录此数据的访问次数**。

## 四、高可用

### 4.1 主从复制

由于数据都是存储在一台服务器上，如果出事就完犊子了，比如：

- 如果服务器发生了宕机，由于数据恢复是需要点时间，那么这个期间是无法服务新的请求的；
- 如果这台服务器的**硬盘出现了故障，可能数据就都丢失了**。

要避免这种单点故障，最好的办法是将数据备份到其他服务器上，让这些服务器也可以对外提供服务，这样即使有一台服务器出现了故障，其他服务器依然可以继续提供服务。

**主从复制**解决了多台服务器之间的数据一致性问题

#### 数据的读写操作是否每台服务器都可以处理等问题

**主服务器可以读写，而从服务器只能读**。所有的数据修改只在主服务器上进行，然后将最新的数据同步给从服务器，这样就使得主从服务器的数据是一致的。



## 4.1.2 第一次同步

我们在服务器 B 上执行下面这条命令：

```
# 服务器 B 执行这条命令
replicaof <服务器 A 的 IP 地址> <服务器 A 的 Redis 端口号>
```

接着，服务器 B 就会变成服务器 A 的「从服务器」，然后与主服务器进行第一次同步。

第一次同步包括**三个阶段**：

- 第一阶段是建立链接、协商同步；
- 第二阶段是主服务器同步数据给从服务器；
- 第三阶段是主服务器发送新写操作命令给从服务器。

### 第一阶段：建立链接、协商同步

从服务器执行了replicaof命令后，从服务器就会给主服务器发送 **psync** 命令，表示要进行数据同步psync 命令包含两个参数，分别是**主服务器的runID**和**复制进度offset**。第一次分别为？ -1。每个 Redis 服务器在启动时都会自动生产一个随机的runID 来唯一标识自己。

主服务器收到 psync 命令后，会用 **FULLRESYNC** 作为响应命令返回给对方。包含**主服务器的 runID** 和**主服务器目前的复制进度 offset**

FULLRESYNC 响应命令的意图是采用**全量复制**的方式，也就是主服务器会把所有的数据都同步给从服务器。

### 第二阶段：主服务器同步数据给从服务器

接着，主服务器会执行 **bgsave** 命令来生成 **RDB 文件**，然后把文件发送给从服务器。进行数据同步

但是，**这期间的写操作命令并没有记录到刚刚生成的 RDB 文件中**，这时主从服务器间的数据就不一致了。那么为了保证主从服务器的数据一致性，**主服务器在下面这三个时间间隙中将收到的写操作命令，写入到 replication buffer 缓冲区里：**

- 主服务器生成 RDB 文件期间；
- 主服务器发送 RDB 文件给从服务器期间；
- 「从服务器」加载 RDB 文件期间；

### 第三阶段：主服务器发送新写命令给从服务器

完成 RDB 的载入后，会回复一个确认消息给主服务器。

接着，主服务器将 replication buffer 缓冲区里所记录的写操作命令发送给从服务器，从服务器执行来自主服务器 replication buffer 缓冲区里发来的命令，**这时主从服务器的数据就一致了。**

### 4.1.3 命令传播

主从服务器在完成第一次同步后，双方之间就会维护一个 TCP 连接。

主服务器可以通过这个连接继续将写操作命令传播给从服务器，然后从服务器执行该命令，使得与主服务器的数据库状态相同。**写操作命令会写入到 replication buffer 缓冲区里**

而且这个连接是长连接的，目的是避免频繁的 TCP 连接和断开带来的性能开销。

上面的这个过程被称为**基于长连接的命令传播**，通过这种方式来保证第一次同步后的主从服务器的数据一致性。

### 4.1.4 分摊主服务器压力

主从复制过程中，主服务器会做两件耗时的操作：生成 RDB 文件和传输 RDB 文件。如果从服务器数量非常多，而且都与主服务器进行全量同步的话，就会带来两个问题：

- 由于是通过 bgsave 命令来生成 RDB 文件的，那么主服务器就会忙于使用 fork() 创建子进程，如果主服务器的内存数据非大，在执行 fork() 函数时是会阻塞主线程的，从而使得 Redis 无法正常处理请求；
- 传输 RDB 文件会占用主服务器的网络带宽，会对主服务器响应命令请求产生影响。

**因此从服务器可以自己建立从服务器来分摊主服务器压力。**

我们在「从服务器」上执行下面这条命令，使其作为目标服务器的从服务器：

```
replicaof <目标服务器的IP> 6379
```

### 4.1.5 增量复制

在一个长连接中，如果网络故障断开，从服务器就没法收到主服务器的写命令，导致网络重新连接后主从数据不一致。

从主从服务器会采用**增量复制**的方式继续同步，**也就是只会把网络断开期间主服务器接收到的写操作命令，同步给从服务器**。主要有三个步骤：

- 从服务器在恢复网络后，会发送 psync 命令给主服务器，此时的 psync 命令里的 offset 参数不是 -1；
- 主服务器收到该命令后，然后用 CONTINUE 响应命令告诉从服务器接下来采用增量复制的方式同步数据；
- 然后主服务将主从服务器断线期间，所执行的写命令发送给从服务器，然后从服务器执行这些命令。

**主服务器怎么知道要将哪些增量数据发送给从服务器呢？**

答案藏在这两个东西里：



- `repl_backlog_buffer`，是一个「**环形**」缓冲区，用于主从服务器断连后，从中找到差异的数据；
- `replication offset`，标记上面那个缓冲区的同步进度，主从服务器都有各自的偏移量，主服务器使用 `master_repl_offset` 来记录自己「写」到的位置，从服务器使用 `slave_repl_offset` 来记录自己「读」到的位置。

网络断开后，当从服务器重新连上主服务器时，主服务器根据自己的 `master_repl_offset` 和 `slave_repl_offset` 之间的差距，然后来决定对从服务器执行哪种同步操作：

- 如果判断出从服务器要读取的数据还在 `repl_backlog_buffer` 缓冲区里，那么主服务器将采用**增量同步**的方式；
- 相反，如果判断出从服务器要读取的数据已经不存在 `repl_backlog_buffer` 缓冲区里，那么主服务器将采用**全量同步**的方式。

**为了避免在网络恢复时，主服务器频繁地使用全量同步的方式，我们应该调整下 `repl_backlog_buffer` 缓冲区大小，尽可能的大一些，减少出现从服务器要读取的数据被覆盖的概率，从而使得主服务器采用增量同步的方式。**

## 4.2 哨兵

在 Redis 的主从架构中，由于主从模式是读写分离的，如果主节点 (master) 挂了，那么将没有主节点来服务客户端的写操作请求，也没有主节点给从节点 (slave) 进行数据同步了。

**哨兵 (\*Sentinel\*) 机制**，它的作用是实现**主从节点故障转移**。它会监测主节点是否存活，如果发现主节点挂了，它就会选举一个从节点切换为主节点，并且把新主节点的相关信息通知给从节点和客户端。

### 哨兵机制是如何工作的？

**哨兵**其实是一个运行在特殊模式下的 Redis 进程，所以它也是一个节点。哨兵节点主要负责三件事情：**监控、选主、通知**。

哨兵会每隔 1 秒**给所有主从节点发送 PING 命令**，当主从节点收到 PING 命令后，会发送一个响应命令给哨兵，这样就可以判断它们是否在正常运行。如果主节点或者从节点没有在规定时间内响应哨兵的 PING 命令，哨兵就会将它们标记为**主观下线**。

当一个哨兵判断主节点为「主观下线」后，就会向其他哨兵发起命令，其他哨兵收到这个命令后，就会根据自身和主节点的网络状况，做出赞成投票或者拒绝投票的响应。**当这个哨兵的赞同票数达到哨兵配置文件中的 `quorum` 配置项设定的值后，这时主节点就会被该哨兵标记为客观下线**。

哪个哨兵节点判断主节点为客观下线，这个哨兵节点就是**候选者**，所谓的候选者就是想当 Leader 的哨兵。

候选者会向其他哨兵发送命令，表明希望成为 **Leader** 来执行主从切换，并让所有其他哨兵对它进行投票。每个哨兵只有一次投票机会，如果用完后就不能参与投票了，可以投给自己或投给别人，但是只有候选者才能把票投给自己。那么在投票过程中，任何一个「候选者」，要满足两个条件：

- 第一，拿到半数以上的赞成票；
- 第二，拿到的票数同时还需要大于等于哨兵配置文件中的 `quorum` 值。

## 主从故障转移的过程是怎样的？

选举出了哨兵 `leader` 后，就可以进行主从故障转移了

- 第一步：在已下线主节点（旧主节点）属下的所有「从节点」里面，挑选出一个从节点，并将其转换为主节点。
- 第二步：让已下线主节点属下的所有「从节点」修改复制目标，修改为复制「新主节点」；
- 第三步：将新主节点的 IP 地址和信息，通过「发布者/订阅者机制」通知给客户端；
- 第四步：继续监视旧主节点，当这个旧主节点重新上线时，将它设置为新主节点的从节点；

### 步骤一：选出新主节点

首先要把网络状态不好的从节点过滤掉，接下来要对所有从节点进行三轮考察：**优先级、复制进度、ID 号**。在进行每一轮考察的时候，哪个从节点优先胜出，就选择其作为新主节点。

- 第一轮考察：哨兵首先会根据从节点的优先级来进行排序，优先级越小排名越靠前，
- 第二轮考察：如果优先级相同，则查看复制的下标，哪个从「主节点」接收的复制数据多，哪个就靠前。
- 第三轮考察：如果优先级和下标都相同，就选择从节点 ID 较小的那个。

### 步骤二：将从节点指向新主节点

兵 `leader` 下一步要做的就是，让已下线主节点属下的所有「从节点」指向「新主节点」，这一动作可以通过向「从节点」发送 **SLAVEOF** 命令来实现。

### 步骤三：通知客户的主节点已更换

这主要通过 **Redis** 的发布者/订阅者机制来实现的。每个哨兵节点提供发布者/订阅者机制，客户端可以从哨兵订阅消息。

主从切换完成后，哨兵就会向 **+switch-master** 频道发布新主节点的 IP 地址和端口的消息，这个时候客户端就可以收到这条信息，然后用这里面的新主节点的 IP 地址和端口进行通信了。

## 步骤四：将旧主节点变为从节点

继续监视旧主节点，当旧主节点重新上线时，哨兵集群就会向它发送 `SLAVEOF` 命令，让它成为新主节点的从节点

## 4.3 集群

# 五、缓存篇

用户的数据一般都是存储于数据库，数据库的数据是落在磁盘上的，磁盘的读写速度可以说是计算机里最慢的硬件了。

因为 Redis 是内存数据库，我们可以将数据库的数据缓存在 Redis 里，相当于数据缓存在内存，内存的读写速度比硬盘快好几个数量级，这样大大提高了系统性能。

## 5.1 什么是缓存雪崩、击穿、穿透？

### 5.1.1 缓存雪崩

当大量缓存数据在同一时间过期 或者 Redis 故障宕机时，如果此时有大量的用户请求，都无法在 Redis 中处理，于是全部请求都直接访问数据库，从而导致数据库的压力骤增，严重的会造成数据库宕机，从而形成一系列连锁反应，造成整个系统崩溃，这就是缓存雪崩的问题。

#### 解决方法

- 均匀设置过期时间；
- 如果发现访问的数据不在 Redis 里，就加个互斥锁，保证同一时间内只有一个请求来构建缓存
- 缓存不设置有效期，而是让缓存“永久有效”，并将更新缓存的工作交由后台线程定时更新。
- 启动服务熔断机制，暂停业务应用对缓存服务的访问，直接返回错误

### 5.1.2 缓存击穿

如果缓存中的某个热点数据过期了，此时大量的请求访问了该热点数据，就无法从缓存中读取，直接访问数据库，数据库很容易就被高并发的请求冲垮，这就是缓存击穿的问题。

#### 解决方法

- 启动服务熔断机制，暂停业务应用对缓存服务的访问，直接返回错误
- 不给热点数据设置过期时间

### 5.1.3 缓存穿透

当用户访问的数据，**既不在缓存中，也不在数据库中**，导致请求在访问缓存时，发现缓存缺失，再去访问数据库时，发现数据库中也没有要访问的数据，没办法构建缓存数据，来服务后续的请求。那么当有大量这样的请求到来时，数据库的压力骤增，这就是**缓存穿透**的问题。

#### 解决方法

- **缓存空值或者默认值**，可以针对查询的数据，在缓存中设置一个空值或者默认值，这样**后续请求就可以从缓存中读取到空值或者默认值**，返回给应用，而不会继续查询数据库。
- 非法请求的限制，在 API 入口处我们要判断请求参数是否合理
- **使用布隆过滤器快速判断数据是否存在**，避免通过查询数据库来判断数据是否存在。

布隆过滤器会通过 3 个操作完成标记：

- 第一步，使用 N 个哈希函数分别对数据做哈希计算，得到 N 个哈希值；
- 第二步，将第一步得到的 N 个哈希值对位图数组的长度取模，得到每个哈希值在位图数组的对应位置。
- 第三步，将每个哈希值在位图数组的对应位置的值设置为 1；

例如当应用要查询数据 x 是否数据库时，通过布隆过滤器的3个哈希函数得到3个哈希值，再得到位图数组的3个位置，只要查到位图数组的这3个位置的值是否全为 1，只要有一个为 0，就认为数据 x 不在数据库中。

缓存异常	产生原因	应对方案
缓存雪崩	大量数据同时过期	- 均匀设置过期时间，避免同一时间过期 - 互斥锁，保证同一时间只有一个应用在构建缓存 - 双 key 策略，主 key 设置过期时间，备 key 永久，主 key 过期时，返回备 key 的内容 - 后台更新缓存，定时更新、消息队列通知更新
	Redis 故障宕机	- 服务熔断 - 请求限流 - 构建 Redis 缓存高可靠集群
缓存击穿	频繁访问的热点数据过期	- 互斥锁 - 不给热点数据设置过期时间，由后台更新缓存
缓存穿透	访问的数据既不在缓存，也不在数据库	- 非法请求的限制； - 缓存空值或者默认值； - 使用布隆过滤器快速判断数据是否存在；

## 5.2 数据库和缓存如何保证一致性？

无论是「先更新数据库，再更新缓存」，还是「先更新缓存，再更新数据库」，这两个方案都存在并发问题，当两个请求并发更新同一条数据的时候，可能会出现**缓存和数据库中的数据不一致的现象**。

**不更新缓存，而是删除缓存中的数据。然后，到读取数据时，发现缓存中没了数据之后，再从数据库中读取数据，更新到缓存中。**

该策略又可以细分为「读策略」和「写策略」。

#### **写策略的步骤：**

- 更新数据库中的数据；
- 删除缓存中的数据。

#### **读策略的步骤：**

- 如果读取的数据命中了缓存，则直接返回数据；
- 如果读取的数据没有命中缓存，则从数据库中读取数据，然后将数据写入到缓存，并且返回给用户。

#### **5.2.1 先删后写还是先写后删？**

写策略有两种顺序：

##### **先删缓存后写 DB**

产生脏数据的概率较大（若出现脏数据，则意味着再不更新的情况下，查询得到的数据均为旧的数据）。

比如两个并发操作，一个是更新操作，另一个是查询操作，更新操作删除缓存后，查询操作没有命中缓存，先把老数据读出来后放到缓存中，然后更新操作更新了数据库。于是，在缓存中的数据还是老的数据，导致缓存中的数据是脏的，而且还一直这样脏下去了。

##### **先写 DB 再删缓存**

产生脏数据的概率较小，但是会出现一致性的问题；若更新操作的时候，同时进行查询操作并命中，则查询得到的数据是旧的数据。但是不会影响后面的查询。

比如一个是读操作，但是没有命中缓存，然后就到数据库中取数据，此时来了一个写操作，写完数据库后，让缓存失效，然后之前的那个读操作再把老的数据放进去，所以会造成脏数据。

#### **5.2.2 如何保证两个操作都能执行成功？**

在删除缓存（第二个操作）的时候可能会失败，导致缓存中的数据是旧值，而数据库是最新值。

解决方法：

- 重试机制。
- 订阅 MySQL binlog，再操作缓存。

## 重试机制

我们可以引入**消息队列**，将第二个操作（删除缓存）要操作的数据加入到消息队列，由消费者来操作数据。

- 如果应用**删除缓存失败**，可以从消息队列中重新读取数据，然后再次删除缓存，这个就是**重试机制**。当然，如果重试超过的一定次数，还是没有成功，我们就需要向业务层发送报错信息了。
- 如果**删除缓存成功**，就要把数据从消息队列中移除，避免重复操作，否则就继续重试。

## 订阅 MySQL binlog，再操作缓存

第一步更新数据库，那么更新数据库成功就会产生一条变更日志，记录在 binlog 里。于是我们就可以通过订阅 binlog 日志，拿到具体要操作的数据，然后再执行缓存删除

# 六、常见面试题

## 认识Redis

### 1. 什么是Redis?

Redis是一个开源的**内存数据结构存储**。用做数据库、缓存、消息代理。

提供**多种数据结构体**

**存储在内存中**，有较高的读写性能

Redis提供数据**持久化**功能

支持**主从复制**，可以实现数据冗余、负载均衡和故障恢复

通过**哨兵模式**和**Redis集群**，可以实现自动故障转移和扩展读写能力。

### 2. Redis和Memcached有什么区别?

Redis 与 Memcached **共同点**:

1. 都是基于内存的数据库，一般都用来当做缓存使用。
2. 都有过期策略。
3. 两者的性能都非常高。

Redis 与 Memcached **区别**:

- Redis 支持的数据类型更丰富 (String、Hash、List、Set、ZSet) , 而 Memcached 只支持最简单的 key-value 数据类型;
- Redis 支持数据的持久化, 可以将内存中的数据保持在磁盘中, 重启的时候可以再次加载进行使用, 而 Memcached 没有持久化功能, 数据全部存在内存之中, Memcached 重启或者挂掉后, 数据就没了;



- Redis 原生支持集群模式, Memcached 没有原生的集群模式, 需要依靠客户端来实现往集群中分片写入数据;
- Redis 支持发布订阅模型、Lua 脚本、事务等功能, 而 Memcached 不支持;

## 3.为什么用Redis作为MySQL的缓冲?

高性能

高并发

单台设备的 Redis 的 QPS (Query Per Second, 每秒钟处理完请求的次数) 是 MySQL 的 10 倍, Redis 单机的 QPS 能轻松破 10w, 而 MySQL 单机的 QPS 很难破 1w。

其它优点

## Redis数据结构

### 1.redis各种数据类型应用场景?

- String 类型的应用场景: 缓存对象、常规计数、分布式锁、共享 session 信息等。
- List 类型的应用场景: 消息队列 (但是有两个问题: 1. 生产者需要自行实现全局唯一 ID; 2. 不能以消费组形式消费数据) 等。
- Hash 类型: 缓存对象、购物车等。
- Set 类型: 聚合计算 (并集、交集、差集) 场景, 比如点赞、共同关注、抽奖活动等。
- Zset 类型: 排序场景, 比如排行榜、电话和姓名排序等。

### 2.五种常见的 Redis 数据类型是怎么实现?

string: SDS

list: 双向列表、压缩列表 (3.0) , quicklist (7.0)

hash: 哈希表、压缩列表 (3.0) , 哈希表、listpack (7.0)

set: 哈希表、整数集

zset: 跳表、压缩列表 (3.0) , 跳表、listpack (7.0)

## Redis线程模型

## 1.Redis是单线程吗？

Redis单线程指的是「接收客户端请求→解析请求 →进行数据读写等操作→发送数据给客户端」这个过程是由一个线程（主线程）来完成的，这也是我们常说 Redis 是单线程的原因。

但是，Redis 程序并不是单线程的，Redis 在启动的时候，是会启动后台线程（BI0）的：

- Redis 在 2.6 版本，会启动 2 个后台线程，分别处理关闭文件、AOF刷盘这两个任务；
- Redis 在 4.0 版本之后，新增了一个新的后台线程，用来异步释放 Redis 内存，也就是 lazyfree 线程。例如执行 unlink key / flushdb async / flushall async 等命令，会把这些删除操作交给后台线程来执行，好处是会导致 Redis 主线程卡顿。因此，当我们要删除一个大 key 的时候，不要使用 del 命令删除，因为 del 是在主线程处理的，这样会导致 Redis 主线程卡顿，因此我们应该使用 unlink 命令来异步删除大key。

## 2.Redis 单线程模式是怎样的？

## 3.Redis 采用单线程为什么还这么快？

- Redis 的大部分操作都在内存中完成，并且采用了高效的数据结构，因此 Redis 瓶颈可能是机器的内存或者网络带宽，而并非 CPU，既然 CPU 不是瓶颈，那么自然就采用单线程的解决方案了；
- Redis 采用单线程模型可以避免多线程之间的竞争，省去了多线程切换带来的时间和性能上的开销，而且也不会导致死锁问题。
- Redis 采用了 I/O 多路复用机制处理大量的客户端 Socket 请求，IO 多路复用机制是指一个线程处理多个 IO 流，就是我们经常听到的 select/epoll 机制。简单来说，在 Redis 只运行单线程的情况下，该机制允许内核中，同时存在多个监听 Socket 和已连接 Socket。内核会一直监听这些 Socket 上的连接请求或数据请求。一旦有请求到达，就会交给 Redis 线程处理，这就实现了一个 Redis 线程处理多个 IO 流的效果。

## 4.Redis 6.0 之前为什么使用单线程？

使用了单线程后，可维护性高。多线程模型虽然在某些方面表现优异，但是它却引入了程序执行顺序的不确定性，带来了并发读写的一系列问题，增加了系统复杂度、同时可能存在线程切换、甚至加锁解锁、死锁造成的性能损耗。

Redis 通过 AE 事件模型以及 IO 多路复用等技术，即使单线程处理性能也非常高，因此没有必要使用多线程。



## 5.Redis 6.0 之后为什么引入了多线程？

在 Redis 6.0 版本之后，也采用了多个 I/O 线程来处理网络请求，这是因为随着网络硬件的性能提升，Redis 的性能瓶颈有时会出现网络 I/O 的处理上。所以为了提高网络 I/O 的并行度，Redis 6.0 对于网络 I/O 采用多线程来处理。但是对于命令的执行，Redis 仍然使用单线程来处理

## Redis持久化

### 1.AOF 日志是如何实现的？

### 2.RDB 快照是如何实现的呢？

### 3.为什么会有混合持久化？

## Redis集群

### 1.Redis 如何实现服务高可用？

主从复制、哨兵模式、切片集群。

### 2.集群脑裂导致数据丢失怎么办？

## Redis过期删除与内存淘汰

### 1.Redis 使用的过期删除策略是什么？

## Redis缓存设计

# Redis 实战

1.Redis 如何实现延迟队列？

2.Redis 的大 key 如何处理？

3.Redis 管道有什么用？

4.Redis 事务支持回滚吗？

5.如何用 Redis 实现分布式锁的？