

## Zookeeper 的设计目标？

- 1、简单的数据结构，Zookeeper 使得分布式程序能够通过一个共享的树形结构的名称空间来进行相互协调，即 Zookeeper 服务器内存中的数据模型由一系列被称为ZNode的数据节点组成，Zookeeper 将全量的数据存储在内存中，以此来提高服务器吞吐、减少延迟的目的。
- 2、可以构建集群 Zookeeper 集群通常由一组机器构成，组成 Zookeeper 集群的而每台机器都会在内存中维护当前服务器状态，并且每台机器之间都相互通信。
- 3、顺序访问，对于来自客户端的每个更新请求，Zookeeper都会分配一个全局唯一的递增编号，这个编号反映了所有事务操作的先后顺序。
- 4、高性能，Zookeeper和Redis一样全量数据存储在内存中，100%读请求压测 QPS 12-13W。

## Zookeeper 有哪些应用场景？

Zookeeper 的功能很强大，应用场景很多，结合我们实际工作中使用 Dubbo 框架的情况，Zookeeper 主要是做注册中心用。

- 基于 Dubbo 框架开发的提供者、消费者都向 Zookeeper 注册自己的 URL，消费者还能拿到并订阅提供者的注册URL，以便在后续程序的执行中去调用提供者。
- 而提供者发生了变动，也会通过 Zookeeper 向订阅的消费者发送通知。
- 统一命名服务。

命名服务是指通过指定的名字来获取资源或服务的地址，利用zk创建一个全局的路径，即时唯一的路径，这个路径就可以作为一个名字，指向集群中机器或者提供服务的地址，又或者一个远程的对象等。

- 分布式锁服务。

这个比较好理解，Zookeeper 实现的分布式锁的可靠性会比 Redis 实现的分布式锁高，当然相对来说，性能会低。

- 配置管理。

例如说，Spring Cloud **Config** Zookeeper，就实现了基于 Zookeeper 的 Spring Cloud **Config** 的实现，提供配置中心的服务。

- 注册与发现

是否有机器加入或退出

所有机器约定在父目录下创建临时目录节点，然后监听父目录节点下的子节点变化。一旦有机器挂掉，该机器与 ZooKeeper 的连接断开，其所创建的临时目录节点也被删除，所有其他机器都收到通知：某个节点被删除了。

- Master 选举。

基于 Zookeeper 实现分布式协调，从而实现主从的选举。这个在 Kafka、Elastic-Job 等等中间件，都有所使用到。

- 分布式锁。

有了 ZooKeeper 的一致性文件系统，锁的问题变得容易。锁服务可以分成两类，一个是保持独占，另一个是控制时序。

- 1、保持独占，我们把 znode 看作是一把锁，通过 createZnode 的方式来实现。所有客户端都去创建 **/distribute\_lock** 节点，最终成功创建的那个客户端也即拥有了这把锁。用完删除掉自己创建的 **/distribute\_lock** 节点就释放出锁。
- 2、控制时序，**/distribute\_lock** 已经预先存在，所有客户端在它下面创建临时顺序编号目录节点，和 Master 一样，编号最小的获得锁，用完删除，依次方便。

- 队列管理

- 1、同步队列，当一个队列的成员都聚齐时，这个队列才可用，否则一直等待。在约定的目录下创建临时目录节点，监听节点数目是否是我们要求的数目。
- 2、队列按照 FIFO 方式进行入队和出队操作。和分布式锁服务中的控制时序的场景基本原理一致，入列有编号，出列按编号。创建 PERSISTENT\_SEQUENTIAL 节点，创建成功时 Watcher 通知等待的队列，队列删除序列号最小的节点以消费。此场景下，znode 用于消息存储，znode 存储的数据就是消息队列中的消息内容，SEQUENTIAL 序列号就是消息的编号，按序取出即可。由于创建的节点是持久化的，所以不必担心队列消息丢失的问题。

## Zookeeper 的通知机制是什么？

Zookeeper 允许客户端向服务端的某个 znode 注册一个 Watcher 监听，当服务端的一些指定事件触发了这个 Watcher，服务端会向指定客户端发送一个事件通知来实现分布式的通知功能，然后客户端根据 Watcher 通知状态和事件类型做出业务上的改变。


整个流程如下：

- 第一步，客户端注册 Watcher。
- 第二步，服务端处理 Watcher。
- 第三步，客户端回调 Watcher。

Watcher 的特性总结：

- 1、一次性。

无论是服务端还是客户端，一旦一个 Watcher 被触发，Zookeeper 都会将其从相应的存储中移除。这样的设计有效的减轻了服务端的压力，不然对于更新非常频繁的节点，服务端会不断的向客户端发送事件通知，无论对于网络还是服务端的压力都非常大。

 注意哟，这个特性可以变成一个面试题「Zookeeper 对节点的 watch 监听通知是永久的吗？」。

如果我们使用 Apache Curator 作为操作 Zookeeper 的客户端，它可以帮我们自动透明的实现持续的 watch 操作，非常方便。

- 2、客户端串行执行。

客户端 Watcher 回调的过程是一个串行同步的过程。

- 3、轻量级 Watch 机制。

1. Watcher 通知非常简单，只会告诉客户端发生了事件，而不会说明事件的具体内容。
2. 客户端向服务端注册 Watcher 的时候，并不会把客户端真实的 Watcher 对象实体传递到服务端，仅仅是在客户端请求中使用 boolean 类型属性进行了标记。

- 4、Watcher event异步发送Watcher的通知事件从Server发送到Client 是异步的，这就存在一个问题，不同的客户端和服务端之间通过Socket 进行通信，由于网络延迟或其他因素导致客户端在不通的时刻监听到事件，由于 Zookeeper本身提供了orderingguarantee，即客户端监听事件后，才会感知它所监视znode 发生了变化。所以我们使用Zookeeper不能期望能够监控到节点每次的变化。Zookeeper只能保证最终的一致性，而无法保证强一致性。
- 5、可以注册 Watcher 的操作：getData、exists、getChildren。
- 6、可以触发 Watcher 的操作：create、delete、setData。
- 7、当一个 Client连接到一个新的服务器上时，watch将会被以任意会话事件触发。当与一个服务器失去连接的时候，是无法接收到 watch 的。而当 Client 重新连接时，如果需要的话，所有先前注册过的watch，都会被重新注册。通常这是完全透明的。只有在一个特殊情况下，watch 可能会丢失：对于一个未创建的znode 的 exists watch，如果在客户端断开连接期间被创建了，并且随后在客户端连接上之前又删除了，这种情况下，这个 watch 事件可能会被丢失。

## Zookeeper 采用什么权限控制机制？

对于 Zookeeper，它采用 ACL ( Access Control List ) 访问控制列表。包括三个方面：

- 权限模式 ( Scheme )
  - IP : 从 IP 地址粒度进行权限控制
  - 【常用】Digest : 最常用, 用类似于 username:password 的权限标识来进行权限配置, 便于区分不同应用来进行权限控制。
  - World : 最开放的权限控制方式, 是一种特殊的 digest 模式, 只有一个权限标识 “world:anyone” 。
  - Super : 超级用户。
- 授权对象
  - 授权对象指的是权限赋予的用户或一个指定实体, 例如 IP 地址或是机器等。
- 权限 Permission
  - CREATE : 数据节点创建权限, 允许授权对象在该 znode 下创建子节点。
  - DELETE : 子节点删除权限, 允许授权对象删除该数据节点的子节点。
  - READ : 数据节点的读取权限, 允许授权对象访问该数据节点并读取其数据内容或子节点列表等。
  - WRITE : 数据节点更新权限, 允许授权对象对该数据节点进行更新操作。
  - ADMIN : 数据节点管理权限, 允许授权对象对该数据节点进行 ACL 相关设置操作。

## Zookeeper 的会话管理是怎么样 的？

ZooKeeper 的每个客户端都维护一组服务端信息, 在创建连接时由应用指定, 客户端随机选择一个服务端进行连接, 连接成功后, 服务端为每个连接分配一个唯一标识。

- 客户端在创建连接时可以指定溢出时间, 客户端会周期性的向服务端发送 PING 请求来保持连接。
- 当客户端检测到与服务端断开连接后, 客户端将自动选择服务端列表中的另一个服务端进行重连。客户端允许应用修改服务端列表, 但修改可能导致客户端与服务端的重连。

## Zookeeper 的部署方式？

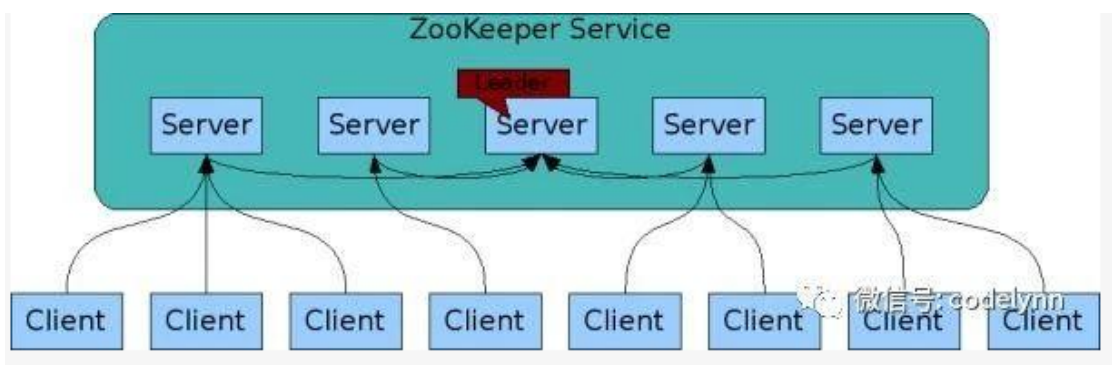
Zookeeper 有两种部署方式：

- 1、单机
- 2、集群

Zookeeper 集群，是一个由多个 Server 组成，一个 Leader，多个 Follower。（这个不同于我们常见的 Master/Slave 模式）Leader 为客户端服务器提供读写服务，除了 Leader 外其他的机器只能提供读服务。

每个 Server 保存一份数据副本全数据一致，分布式读 Follower，写由 Leader 实施更新请求转发，由 Leader 实施更新请求顺序进行，来自同一个 Client 的更新请求按其发送顺序依次执行数据更新原子性，一次数据更新要么成功，要么失败。

全局唯一数据视图，Client 无论连接到哪个 Server，数据视图都是一致的实时性，在一定事件范围内，Client 能读到最新数据。



一般来说，测试环境部署单机，而生产环境必须必须必须部署集群。

## 集群如果有 3 台机器，挂掉 1 台集群还能工作吗？挂掉 2 台呢？

记住一个原则：过半存活即可用。所以挂掉 1 台可以继续工作，挂掉 2 台不可以工作。

## 集群支持动态添加机器吗？

在 3.5 版本开始，支持动态扩容。

而在 3.5 版本之前，Zookeeper 在这方面不太好。所以需要如下两种方式：

- 全部重启：关闭所有Zookeeper服务，修改配置之后启动。不影响之前客户端的会话。
- 逐个重启：顾名思义。这是比较常用的方式

## ZooKeeper 的工作原理？

ZooKeeper 的核心是原子广播，这个机制保证了各个 Server 之间的同步。实现这个机制的协议叫做 Zab 协议。Zab 协议有两种模式，它们分别是恢复模式（选主）和广播模式（同步）：

- 选主：当服务启动或者 Leader 崩溃后，Zab 就进入了恢复模式，当新的 Leader 被选举出来，且大多数 Server 完成了和 Leader 的状态同步以后，恢复模式就结束了。
  - 首先，选举产生新的 Leader 服务器。
  - 然后，集群中 Follower 服务器开始与新的 Leader 服务器进行数据同步。
  - 当集群中超过半数机器与该 Leader 服务器完成数据同步之后，退出恢复模式进入消息广播模式，
- 同步：状态同步保证了 Leader 和 Server 具有相同的系统状态。
  - Leader 服务器开始接收客户端的事务请求，生成事务提案来进行事务请求处理。

## ZooKeeper 是如何保证事务的顺序一致性的？

ZooKeeper 采用了递增的事务 id 来识别，所有的 proposal（提议）都在被提出的时候加上了 zxid。zxid 实际上是一个 64 位数字。

- 高 32 位是 epoch 用来标识 Leader 是否发生了改变，如果有新的 Leader 产生出来，epoch 会自增。
- 低 32 位用来递增计数。

当新产生的 proposal 的时候，会依据数据库的两阶段过程，首先会向其他的 Server 发出事务执行请求，如果超过半数的机器都能执行并且能够成功，那么就会开始执行。

## ZooKeeper 集群中个服务器之间是怎样通信的？

Leader 服务器会和每一个 Follower/Observer 服务器都建立 TCP 连接，同时为每个 Follower/Observer 都创建一个叫做 LearnerHandler 的实体。

- LearnerHandler 主要负责 Leader 和 Follower/Observer 之间的网络通讯，包括数据同步，请求转发和 Proposal 提议的投票等。
- Leader 服务器保存了所有 Follower/Observer 的 LearnerHandler。

## ZAB 和 Paxos 算法的联系与区别？



Paxos 算法是分布式选举算法，Zookeeper 使用的 ZAB 协议（Zookeeper 原子广播）。

二者有相同的地方：

- 都有一个 Leader，用来协调 N 个 Follower 的运行
- Leader 要等待超半数的 Follower 做出正确反馈之后才进行提案。
- 二者都有一个值来代表 Leader 的周期。ZAB 协议中，每个 Proposal 中都包含一个 epoch 值来代表当前的 Leader 周期，Paxos 中名字为 Ballot。

不同的地方在于：

- ZAB 用来构建高可用的分布式数据主备系统（Zookeeper），Paxos 是用来构建分布式一致性状态机系统。

Paxos 算法、ZAB 协议要想讲清楚可不是一时半会的事儿，自 1990 年莱斯利·兰伯特提出 Paxos 算法以来，因为晦涩难懂并没有受到重视。后续几年，兰伯特通过好几篇论文对其进行更进一步地解释，也直到 06 年谷歌发表了三篇论文，选择 Paxos 作为 Chubby cell 的一致性算法，Paxos 才真正流行起来。

## Zookeeper 的选举过程？

当 Leader 崩溃，或者 Leader 失去大多数的 Follower，这时 Zookeeper 进入恢复模式，恢复模式需要重新选举出一个新的 Leader，让所有的 Server 都恢复到一个正确的状态。

Zookeeper 的选举算法有两种：一种是基于 basic paxos 实现的，另外一种是基于 fast paxos 算法实现的。系统默认的选举算法为 fast paxos。

相对详细的，胖友可以看看《[【分布式】Zookeeper的Leader选举](#)》和《[Zookeeper 源码分析 —— Zookeeper Leader 选举算法](#)》。



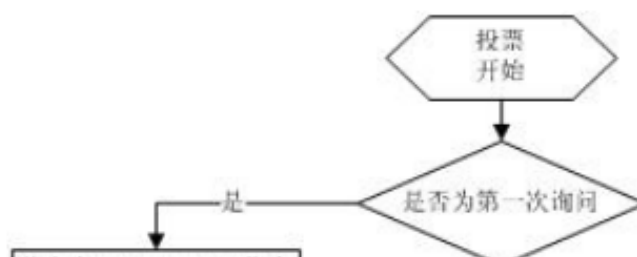
- 不同阶段的选举流程
  - 服务器启动时期的 Leader 选举。
  - 服务器运行时期的 Leader 选举。
- 三种选举算法
  - LeaderElection ：使用 basic paxos 算法。
  - FastLeaderElection ：使用 fast paxos 算法。
  - AuthFastLeaderElection ：在 FastLeaderElection 的基础上，增加认证。
  - 最终在 Zookeeper 3.4.0 版本之后，只保留 FastLeaderElection 版本。

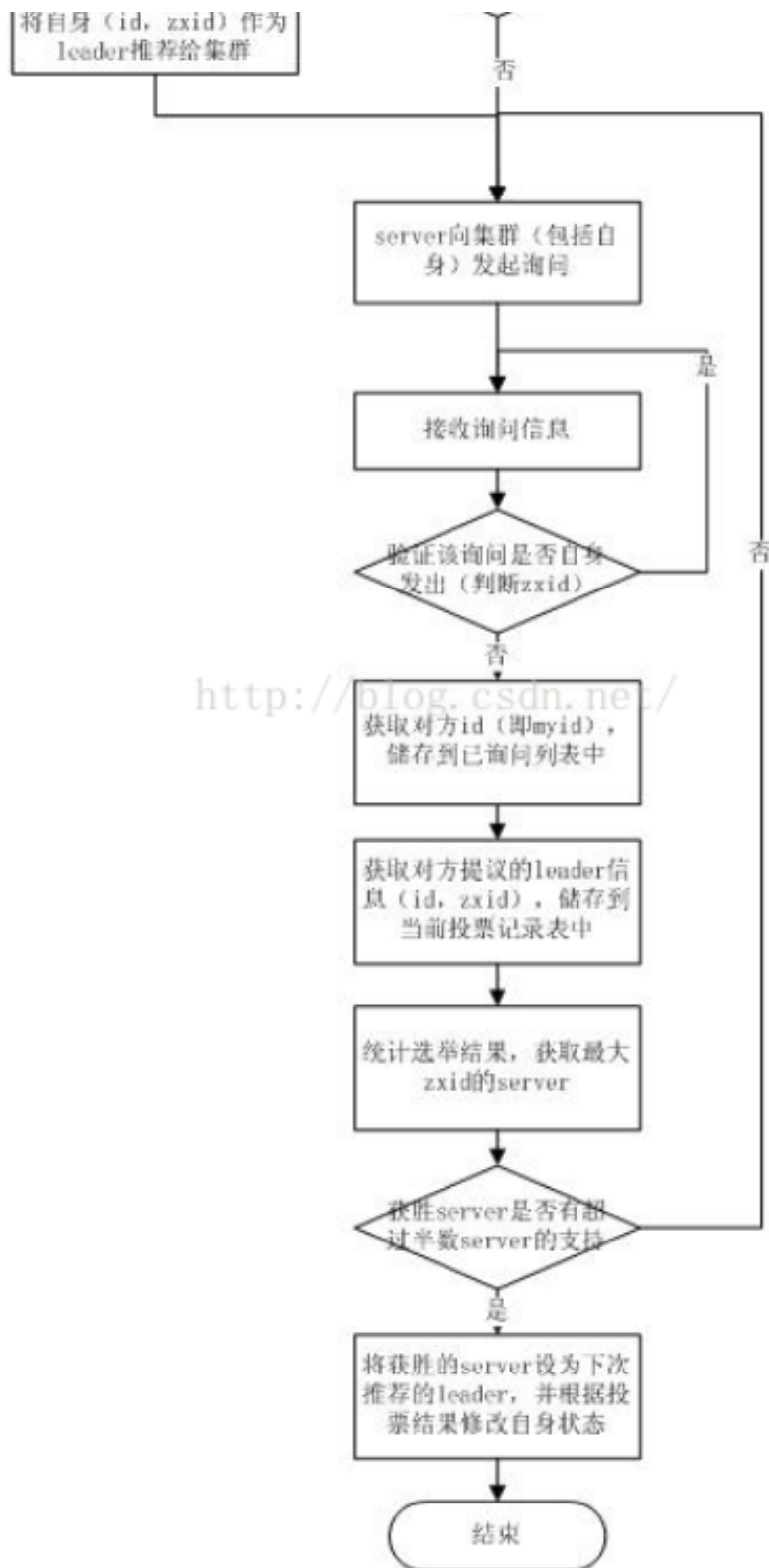
## Zookeeper 选主流程(basic paxos) ?

- 1、选举线程由当前 Server 发起选举的线程担任，其主要功能是对投票结果进行统计，并选出推荐的 Server 。
- 2、选举线程首先向所有 Server 发起一次询问(包括自己)。
- 3、选举线程收到回复后，验证是否是自己发起的询问(验证 zxid 是否一致)，然后获取对方的 id(myid)，并存储到当前询问对象列表中，最后获取对方提议的 Leader 相关信息(id, zxid)，并将这些信息存储到当次选举的投票记录表中。
- 4、收到所有 Server 回复以后，就计算出 zxid 最大的那个 Server，并将这个 Server 相关信息设置成下一次要投票的 Server 。
- 5、线程将当前 zxid 最大的 Server 设置为当前 Server 要推荐的 Leader，如果此时获胜的 Server 获得  $n/2+1$  的 Server 票数，设置当前推荐的 Leader 为获胜的 Server，将根据获胜的 Server 相关信息设置自己的状态，否则，继续这个过程，直到 Leader 被选举出来。

通过流程分析我们可以得出：要使 Leader 获得多数 Server 的支持，则 Server 总数必须是奇数  $2n+1$ ，且存活的 Server 的数目不得少于  $n+1$ 。每个 Server 启动后都会重复以上流程。

在恢复模式下，如果是刚从崩溃状态恢复的或者刚启动的 Server 还会从磁盘快照中恢复数据和会话信息，Zookeeper 会记录事务日志并定期进行快照，方便在恢复时进行状态恢复。



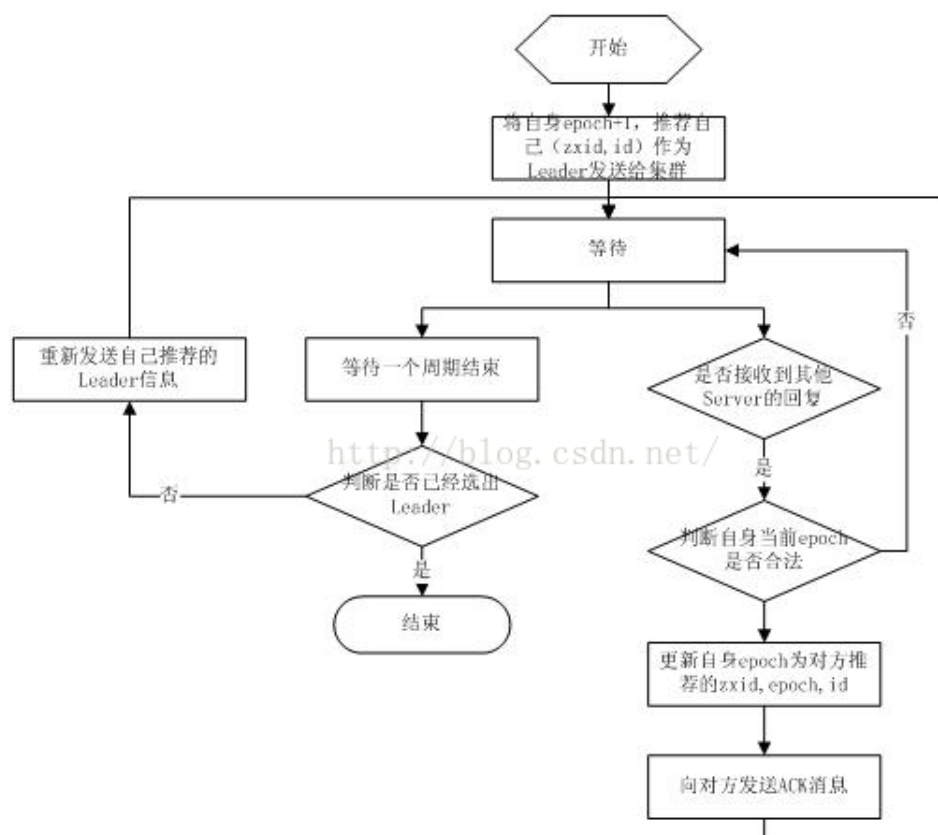


Zookeeper 选主流程(fast paxos) ?

由于 LeaderElection 收敛速度较慢，所以 Zookeeper 引入了 FastLeaderElection 选举算法，FastLeaderElection 也成了 Zookeeper 默认的 Leader 选举算法。

FastLeaderElection 是标准的 Fast Paxos 的实现。它首先向所有 Server 提议自己要成为 Leader，当其它 Server 收到提议以后，解决 epoch 和 zxid 的冲突，并接受对方的提议，然后向对方发送接受提议完成的消息。重复这个流程，最后一定能选举出 Leader。

FastLeaderElection 算法通过异步的通信方式来收集其它节点的选票，同时在分析选票时又根据投票者的当前状态来作不同的处理，以加快 Leader 的选举进程。



## 为什么 Zookeeper 集群推荐节点数是单数？

在统计投票时，有个过半的概念，大于集群机器数量的一半，即大于或等于 $(n/2+1)$ 。那么我们来看看如下的统计：

| 集群数量 | 至少正常运行数量             | 允许挂掉的数量 |
|------|----------------------|---------|
| 2    | 2 的半数为 1，半数以上最少为 2   | 0       |
| 3    | 3 的半数为 1.5，半数以上最少为 2 | 1       |
| 4    | 4 的半数为 2，半数以上最少为 3   | 1       |
| 5    | 5 的半数为 2.5，半数以上最少为 3 | 2       |
| 6    | 6 的半数为 3，半数以上最少为 4   | 2       |

通过以上可以发现：

- 3 台服务器和 4 台服务器都最多允许 1 台服务器挂掉，5 台服务器和 6 台服务器都最多允许 2 台服务器挂掉，明显 4 台服务器成本高于 3 台服务器成本，6 台服务器成本高于 5 服务器成本。
- 这是由于半数以上投票通过决定的。所以，Zookeeper 集群推荐节点数是单数。

## Zookeeper 是否需存在脑裂？

按道理说，Zookeeper 选举不会存在脑裂问题，因为需要  $n / 2 + 1$  投票通过，才能执行对应的写操作。但是听朋友说，实际场景下，貌似发生过脑裂问题。关于这块，芳苒心里也不太有底，欢迎在星球一起讨论。

- [《Zookeeper 已经分布式环境中的假死脑裂》](#)
- [《zookeeper \(二\) 常见问题汇总》](#)

## 机器中为什么会有 Leader？

在分布式环境中，有些业务逻辑只需要集群中的某一台机器进行执行，其他的机器可以共享这个结果，这样可以大大减少重复计算，提高性能，于是就需要进行 Leader 选举。

## Zookeeper 的同步流程？

选完 Leader 以后，Zookeeper 就进入状态同步过程。

- 1、Leader 等待 Server 连接。
  - 2、Follower 连接 Leader ，将最大的 zxid 发送给 Leader 。
  - 3、Leader 根据 Follower 的 zxid 确定同步点。
  - 4、完成同步后通知 Follower 已经成为 update 状态。
  - 5、Follower 收到 update 消息后，又可以重新接受 Client 的请求进行服务了。
- 当然，同步流程并不是像上述描述的这么简单，具体的，还是得看看[《Zookeeper Leader 和 Learner 的数据同步》](#)。