

Raft(01)

基本概念

定义: Raft是一个共识算法,即集群中的多个节点对某某一件事能够达成一致,及时出现了部分的节点故障、网络延迟,甚至是网络分区的情况下

应用: Raft简单来说就是一切以领导者为主,来实现一系列数据的共识,属于单领导者的分布式算法

问题: 在如图的网络拓扑中,假设出现了网络分区,那么Raft算法如何保证在同一时间有且只有一个领导者处理写请求的?

节点分类: 领导者(Leader), 候选人(Candidate), 追随者(Follower)

成员身份: 领导者(Leader) 集群节点的领军人物,其工作是处理写请求,管理日志复制以及定时向其余节点发送心跳包

候选人(Candidate) 心跳包的内容就是:我还活着,你们这群小崽子们老实一点,不要随随便便选举新的领导者上位

追随者(Follower) 当领导者因故障下线时,候选人通过RPC调用向其它节点发送请求投票信息,通知其它节点进行投票,如果多数节点通过,将称为新的领导者

普通打工人,接收领导者的一切处理需求。当领导者的心跳包发送超时时,就主动站出来报称自己当领导者,追随者(Follower) 此时身份转换成候选人

解决一致性以及高可用性问题和集团的董事长以及董事会非常相似

选举领导过程

节点初始化: 在集群节点初始化时,所有节点的角色均为追随者,并且,每个节点会随机地分配一个超时时间,该超时时间就是接收到领导者心跳包的超时时间

投票阶段: 每个节点在发起投票前都会给自己投一票,并且更新自己的任期编号

任期: 领导者是有任期的,当领导者任期到期时,会重新开始选举,任期编号会随着选举过程而变化

变化时机: 当追随者接收领导者的请求心跳包超时,则转变角色,由追随者转变成候选人,并给自己的任期加1

关于投票: 首先,每一轮的投票是有超时时间限制的,不可能永久地进行投票,该投票的超时时间为随机值,一旦超过这个值,本轮投票无效,进行下一轮投票

split vote: 假设集群节点中两个节点同时发起了投票请求,并且获得的票数相同,会发生什么?

日志复制

在Raft算法中,副本数据是以日志的形式存在的,当领导者接收到客户端的写操作时,处理写请求的过程就是将数据转换成日志项并对其进行复制

在Raft算法中,副本数据是以日志的形式存在的,而日志则由日志项所构成

日志项: 指令由客户端所制定,表示对某个key进行怎样的操作,例如 SET

Raft日志格式: 索引值,任期编号

MySQL的两阶段提交: redo log prepare, binlog, redo log prepare, binlog, flush, redo log commit

Raft两阶段提交: 写操作就不算真正的完成,写操作就不算真正的完成

一致性保证: 当领导者和追随者的日志出现不一致时,领导者将强制地要求追随者复制自己的日志项,从而达到日志的一致性

日志覆盖过程: G2在收到RPC日志复制信息后,首先判断索引为6,任期是否为6,任期不为6,则返回给领导者一个错误