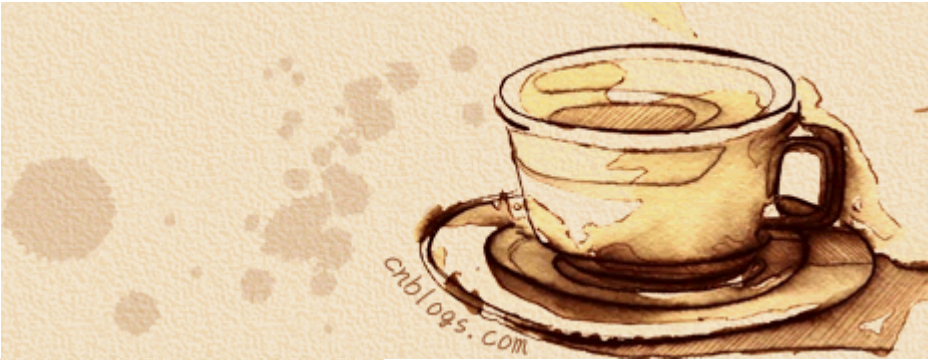


xybaby

竹杖芒鞋轻胜马，谁怕？一蓑烟雨任平生



博客园 首页 新随笔 联系 管理 订阅 XML

随笔- 87 文章- 0 评论- 284

一文搞懂Raft算法

目录

- [raft算法概览](#)
- [leader election](#)
 - [term](#)
 - [选举过程详解](#)
- [log replication](#)
 - [Replicated state machines](#)
 - [请求完整流程](#)
- [safety](#)
- [corner case](#)
 - [stale leader](#)
 - [State Machine Safety](#)
 - [leader crash](#)
- [总结](#)
- [references](#)

正文

raft是工程上使用较为广泛的强一致性、去中心化、高可用的分布式协议。在这里强调了是在工程上，因为在学术理论界，最耀眼的还是大名鼎鼎的Paxos。但Paxos是：少数真正理解的人觉得简单，尚未理解的人觉得很难，大多数人都是一知半解。本人也花了很多时间、看了很多材料也没有真正理解。直到看到raft的论文，两位研究者也提到，他们也花了很长的时间来理解Paxos，他们也觉得很难理解，于是研究出了raft算法。

raft是一个共识算法（consensus algorithm），所谓共识，就是多个节点对某个事情达成一致的看法，即使是在部分节点故障、网络延时、网络分割的情况下。这些年最为火热的加密货币（比特币、区块链）就需要共识算法，而在分布式系统中，共识算法更多用于提高系统的容错性，比如分布式存储中的复制集（replication），在[带着问题学习分布式系统之中心化复制集](#)一文中介绍了中心化复制集的相关知识。raft协议就是一种leader-based的共识算法，与之相应的是leaderless的共识算法。

本文基于论文[In Search of an Understandable Consensus Algorithm](#)对raft协议进行分析，当然，还是建议读者直接看论文。

本文地址：<https://www.cnblogs.com/xybaby/p/10124083.html>

raft算法概览

[回到顶部](#)

Raft算法的头号目标就是容易理解（UnderStandable），这从论文的标题就可以看出来。当然，Raft增强了可理解性，在性能、可靠性、可用性方面是不输于Paxos的。

Raft more understandable than Paxos and also provides a better foundation for building practical systems

为了达到易于理解的目标，raft做了很多努力，其中最主要是两件事情：

- 问题分解
 - 状态简化
- 问题分解是将"复制集中节点一致性"这个复杂的问题划分为数个可以被独立解释、理解、解决的子问题。在raft，子问题包括，*leader election*，*log replication*，*safety*，*membership changes*。而状态简化更好理解，就是对算法做出一些限制，减少需要考虑的状态数，使得算法更加清晰，更少的不确定性（比如，保证新选举出来的leader会包含所有committed log entry）

Raft implements consensus by first electing a distinguished leader, then giving the leader complete responsibility for managing the replicated log. The leader accepts log entries from clients, replicates them on other servers, and tells servers when it is safe to apply log entries to their state machines. A leader can fail or become disconnected from the other servers, in which case a new leader is elected.

上面的引文对raft协议的工作原理进行了高度的概括：raft会先选举出leader，leader完全负责replicated log的管理。leader负责接受所有客户端更新请求，然后复制到follower节点，并在“安全”的时候执行这些请求。如果leader故障，followers会重新选举出新的leader。

这就涉及到raft最新的两个子问题： leader election和log replication

leader election

[回到顶部](#)

raft协议中，一个节点任一时刻处于以下三个状态之一：

- leader
- follower
- candidate

给出状态转移图能很直观的直到这三个状态的区别

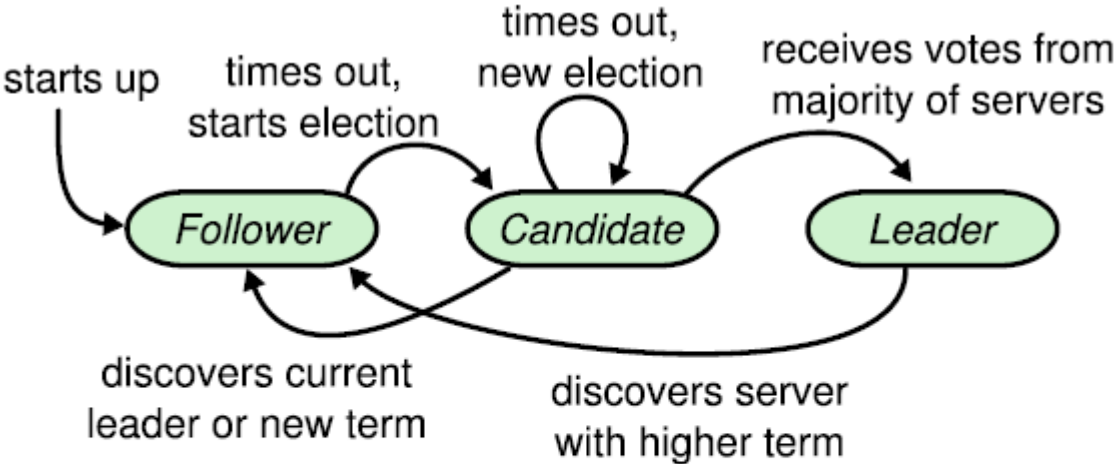


Figure 4: Server states. Followers only respond to requests from other servers. If a follower receives no communication, it becomes a candidate and initiates an election. A candidate that receives votes from a majority of the full cluster becomes the new leader. Leaders typically operate until they fail.

昵称：xybaby
园龄：2年7个月
粉丝：578
关注：9
[+加关注](#)

< 2019年8月 >						
日	一	二	三	四	五	六
28	29	30	31	1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	31
1	2	3	4	5	6	7

搜索

找找看

谷歌搜索

我的标签

[distributed system](#)(21)
[python2.7](#)(18)
[python](#)(15)
[programmer](#)(13)
[mongodb](#)(8)
[web](#)(7)
[coroutine](#)(6)
[gunicorn](#)(5)
[programming paradigm](#)(4)
[gevent](#)(4)
[更多](#)

随笔档案

2019年5月 (1)
2019年4月 (1)
2019年1月 (1)
2018年12月 (2)
2018年11月 (1)
2018年10月 (1)
2018年9月 (1)
2018年8月 (3)
2018年7月 (1)
2018年6月 (2)
2018年5月 (3)
2018年4月 (3)
2018年3月 (2)
2018年2月 (3)
2018年1月 (2)
2017年12月 (4)
2017年11月 (5)
2017年10月 (3)
2017年9月 (3)
2017年8月 (4)
2017年7月 (3)
2017年6月 (2)
2017年5月 (2)
2017年4月 (4)
2017年3月 (6)
2017年2月 (7)
2017年1月 (17)

python2. 7

积分与排名

积分 - 208149
排名 - 1801

最新评论

- Re:什么是分布式系统，如何学习分布式系统
@anobscureretreat您过吗？您配吗？ ...
--ExplorerOfJava
- Re:什么是分布式系统，如何学习分布式系统
感谢博主Thanks!(´ω`)/
--rgb-24bit
- Re:什么是分布式系统，如何学习分布式系统
受益匪浅
--super超人
- Re:使用gc、objgraph干掉python内存泄露与循环引用！
好文，刨根问底，又涉及到解决方法和各种工具的使用！
--PegasusWang
- Re:为什么要读源代码，如何阅读源代码
请问作者看代码一般用到哪些工具呢？
--PegasusWang

阅读排行榜

- 什么是分布式系统，如何学习分布式系统(160169)
- gunicorn 简介(36878)
- 从银行转账失败到分布式事务：总结与思考(31165)

19

推荐

0

反对

可以看出所有节点启动时都是follower状态；在一段时间内如果没有收到来自leader的心跳，从follower切换到candidate，发起选举；如果收到majority的造成票（含自己的一票）则切换到leader状态；如果发现其他节点比自己更新，则主动切换到follower。

总之，系统中最多只有一个leader，如果在一段时间里发现没有leader，则大家通过选举-投票选出leader。leader会不停的给follower发心跳消息，表明自己的存活状态。如果leader故障，那么follower会转换成candidate，重新选出leader。

term

从上面可以看出，哪个节点做leader是大家投票选举出来的，每个leader工作一段时间，然后选出新的leader继续负责。这跟民主社会的选举很像，每一届新的履职期称之为**一届任期**，在raft协议中，也是这样的，对应的术语叫**term**。

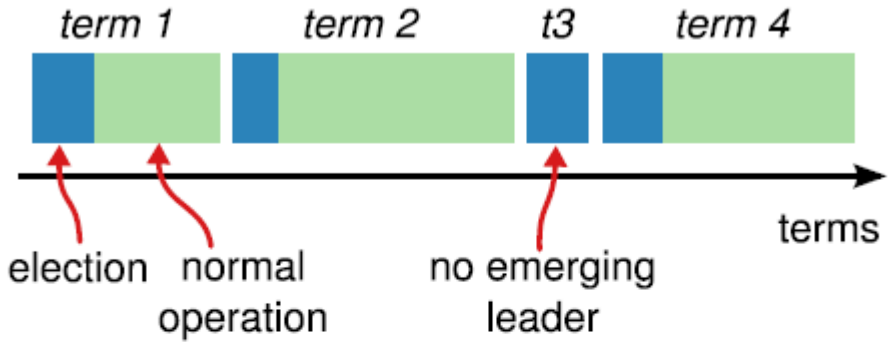


Figure 5: Time is divided into terms, and each term begins with an election. After a successful election, a single leader manages the cluster until the end of the term. Some elections fail, in which case the term ends without choosing a leader. The transitions between terms may be observed at different times on different servers.

term（任期）以选举（election）开始，然后就是一段或长或短的稳定工作期（normal Operation）。从上图可以看到，任期是递增的，这就充当了逻辑时钟的作用；另外，term 3展示了一种情况，就是说没有选举出leader就结束了，然后会发起新的选举，后面会解释这种**split vote**的情况。

选举过程详解

上面已经说过，如果follower在**election timeout**内没有收到来自leader的心跳，（也许此时还没有选出leader，大家都在等；也许leader挂了；也许只是leader与该follower之间网络故障），则会主动发起选举。步骤如下：

- 增加节点本地的 *current term* ，切换到candidate状态

- 投自己一票
- 并行给其他节点发送 *RequestVote RPCs*
- 等待其他节点的回复

在这个过程中，根据来自其他节点的消息，可能出现三种结果

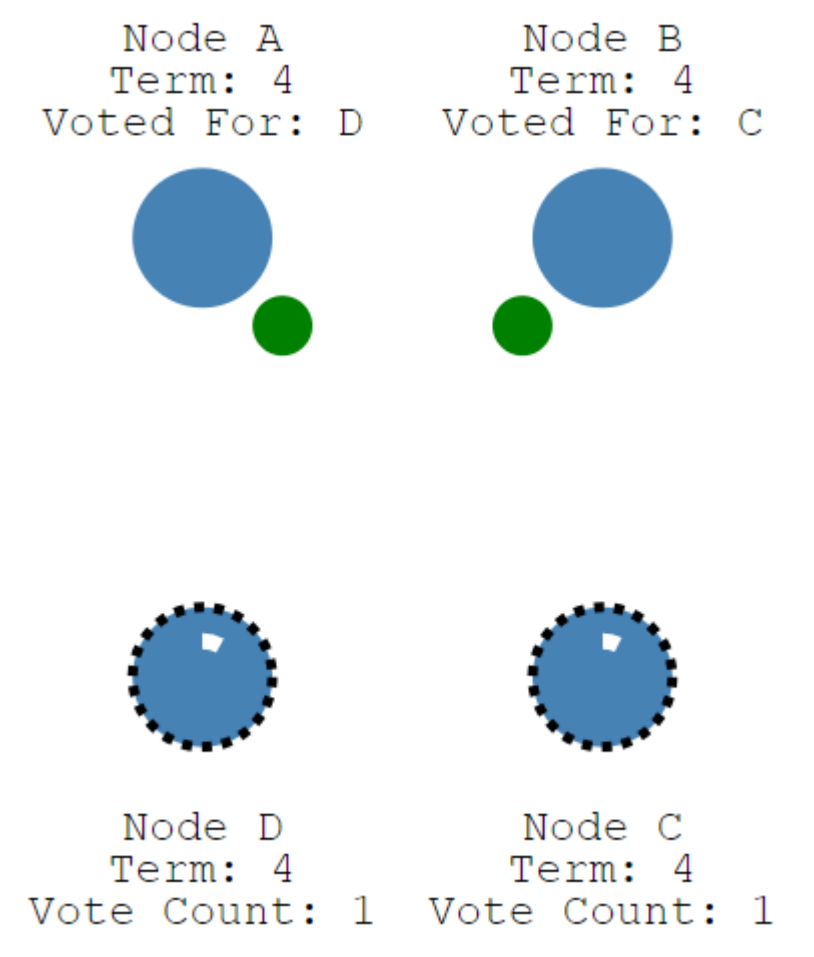
- 收到majority的投票（含自己的一票），则赢得选举，成为leader
- 被告知别人已当选，那么自行切换到follower
- 一段时间内没有收到majority投票，则保持candidate状态，重新发出选举

第一种情况，赢得了选举之后，新的leader会立刻给所有节点发消息，广而告之，避免其余节点触发新的选举。在这里，先回到投票者的视角，投票者如何决定是否给一个选举请求投票呢，有以下约束：

- 在任一任期内，单个节点最多只能投一票
- 候选人知道的信息不能比自己的少（这一部分，后面介绍log replication和safety的时候会详细介绍）
- first-come-first-served 先来先得

第二种情况，比如有三个节点A B C。A B同时发起选举，而A的选举消息先到达C，C给A投了一票，当B的消息到达C时，已经不能满足上面提到的第一个约束，即C不会给B投票，而A和B显然都不会给对方投票。A胜出之后，会给B,C发心跳消息，节点B发现节点A的term不低于自己的term，知道有已经有Leader了，于是转换成follower。

第三种情况，没有任何节点获得majority投票，比如下图这种情况：



总共有四个节点，Node C、Node D同时成为了candidate，进入了term 4，但Node A投了NodeD一票，NodeB投了Node C一票，这就出现了平票 split vote的情况。这个时候大家都在等啊等，直到超时后重新发起选举。如果出现平票的情况，那么就延长了系统不可用的时间（没有leader是不能处理客户端写请求的），因此raft引入了randomized election timeouts来尽量避免平票情况。同时，leader-based 共识算法中，节点的数目都是奇数个，尽量保证majority的出现。

log replication

[回到顶部](#)

当有了leader，系统应该进入对外工作期了。客户端的一切请求来发送到leader，leader来调度这些并发请求的顺序，并且保证leader与followers状态的一致性。raft中的做法是，将这些请求以及执行顺序告知followers。leader和followers以相同的顺序来执行这些请求，保证状态一致。

Replicated state machines

共识算法的实现一般是基于复制状态机（Replicated state machines），何为复制状态机：

If two identical, **deterministic** processes begin in the same state and get the same inputs in the same order, they will produce the same output and end in the same state.

简单来说：**相同的初识状态 + 相同的输入 = 相同的结束状态**。引文中有一个很重要的词 `deterministic`，就是说不同节点要以相同且确定性的函数来处理输入，而不要引入一下不确定的值，比如本地时间等。如何保证所有节点 `get the same inputs in the same order`，使用replicated log是一个很不错的注意，log具有持久化、保序的特点，是大多数分布式系统的基石。

因此，可以这么说，在raft中，leader将客户端请求（command）封装到一个log entry，将这些log entries复制（replicate）到所有follower节点，然后大家按相同顺序应用（apply）log entry中的command，则状态肯定是一致的。

19

推荐

0

反对

- python __getattr__ 巧妙应用(26761)
- 使用gc、objgraph干掉python内存泄露与循环引用! (22850)
- 分布式学习最佳实践：从分布式系统的特征开始（附思维导图）(20930)
- python性能优化(19174)
- 带着问题学习分布式系统之数据分片(16761)
- CAP理论与MongoDB—一致性、可用性的一些思考(15164)
- python yield generator 详解(15136)
- python bottle 简介(14767)
- 关于负载均衡的一切：总结与思考(12056)
- 带着问题学习分布式系统(11411)
- 一文搞懂Raft算法(11111)
- greenlet 详解(10954)
- 同步与异步，回调与协程(10866)
- python属性查找 深入理解（attribute lookup）(10069)
- 用信号量和读写锁解决读者写者问题(9984)
- 日志的艺术（The art of logging）(9651)
- 关于metaclass，我原以为我是懂的(9424)
- gevent调度流程解析(9405)
- 并发与同步、信号量与管理、生产者消费者问题(9395)
- 我的进程去哪儿了，谁杀了我的进程(9318)
- 什么是Mixin模式：带实现的协议(8806)
- 性能优化指南：性能优化的一般性原则与方法(8043)

评论排行榜

- 什么是分布式系统，如何学习分布式系统(37)
- 从银行转账失败到分布式事务：总结与思考(28)
- 带着问题学习分布式系统之数据分片(19)
- 带着问题学习分布式系统(14)
- 分布式学习最佳实践：从分布式系统的特征开始（附思维导图）(13)
- 怎样才算得上合格的程序员(12)
- 作为程序员，再也不想和PM干架了(10)
- 我在博客园的这一年(10)
- 关于负载均衡的一切：总结与思考(10)
- Hey, man, are you ok? - - 关于心跳、故障监测、lease机制(8)
- 想要升职加薪？先管理好时间与目标！(7)
- 我的进程去哪儿了，谁杀了我的进程(7)
- 一文搞懂Raft算法(7)
- 分布式系统中生成全局ID的总结与思考(7)
- 技术领导（Technical Leader）画像(5)
- 啊，我的程序为啥卡住啦(5)
- 并发与同步、信号量与管理、生产者消费者问题(5)
- python yield generator 详解(5)
- 带着SMART原则重新出发(5)
- 性能优化指南：性能优化的一般性原则与方法(5)

推荐排行榜

- 什么是分布式系统，如何学习分布式系统(228)
- 从银行转账失败到分布式事务：总结与思考(96)
- 分布式学习最佳实践：从分布式系统的特征开始（附思维导图）(71)
- 日志的艺术（The art of logging）(45)
- 关于负载均衡的一切：总结与思考(31)
- 带着问题学习分布式系统(24)
- 怎样才算得上合格的程序员(23)
- 想要升职加薪？先管理好时间与目标！(22)
- 我在博客园的这一年(20)
- 一文搞懂Raft算法(19)
- 啊，我的程序为啥卡住啦(18)
- 我的进程去哪儿了，谁杀了我的进程(15)
- 分布式系统中生成全局ID的总结与思考(15)
- 技术领导（Technical Leader）画像(14)
- CAP理论与MongoDB—一致性、可用性的一些思考(14)
- 带着问题学习分布式系统之中心化复制集(14)
- 带着问题学习分布式系统之数据分片(14)
- 使用gc、objgraph干掉python内存泄露与循环引用! (12)
- python yield generator 详解(12)
- 性能优化指南：性能优化的一般性原则与方法(11)
- python性能优化(11)
- 大型网站技术架构：摘要与读书笔记(10)
- Hey, man, are you ok? - - 关于心跳、故障监测、lease机制(10)
- 不想再被鄙视？那就看进来！ 一文搞懂Python2字符编码(9)
- 并发与同步、信号量与管理、生产者消费者问题(9)

下图形象展示了这种log-based replicated state machine

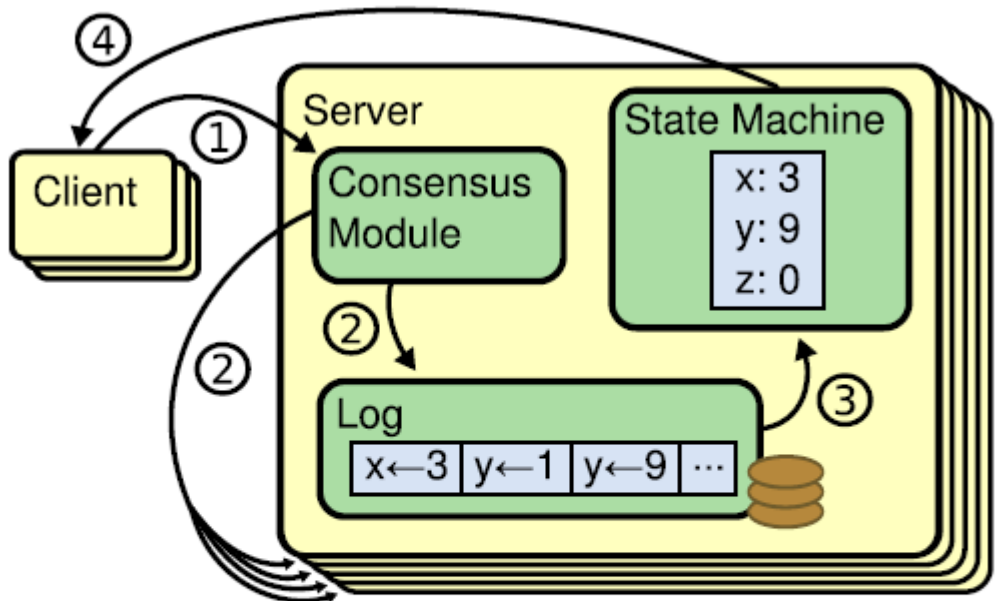


Figure 1: Replicated state machine architecture. The consensus algorithm manages a replicated log containing state machine commands from clients. The state machines process identical sequences of commands from the logs, so they produce the same outputs.

请求完整流程

当系统（leader）收到一个来自客户端的写请求，到返回给客户端，整个过程从leader的视角来看会经历以下步骤：

- leader append log entry
- leader issue AppendEntries RPC in parallel
- leader wait for majority response
- leader apply entry to state machine
- leader reply to client
- leader notify follower apply log

可以看到日志的提交过程有点类似两阶段提交(2PC)，不过与2PC的区别在于，leader只需要大多数（majority）节点的回复即可，这样只要超过一半节点处于工作状态则系统就是可用的。

那么日志在每个节点上是什么样子的呢

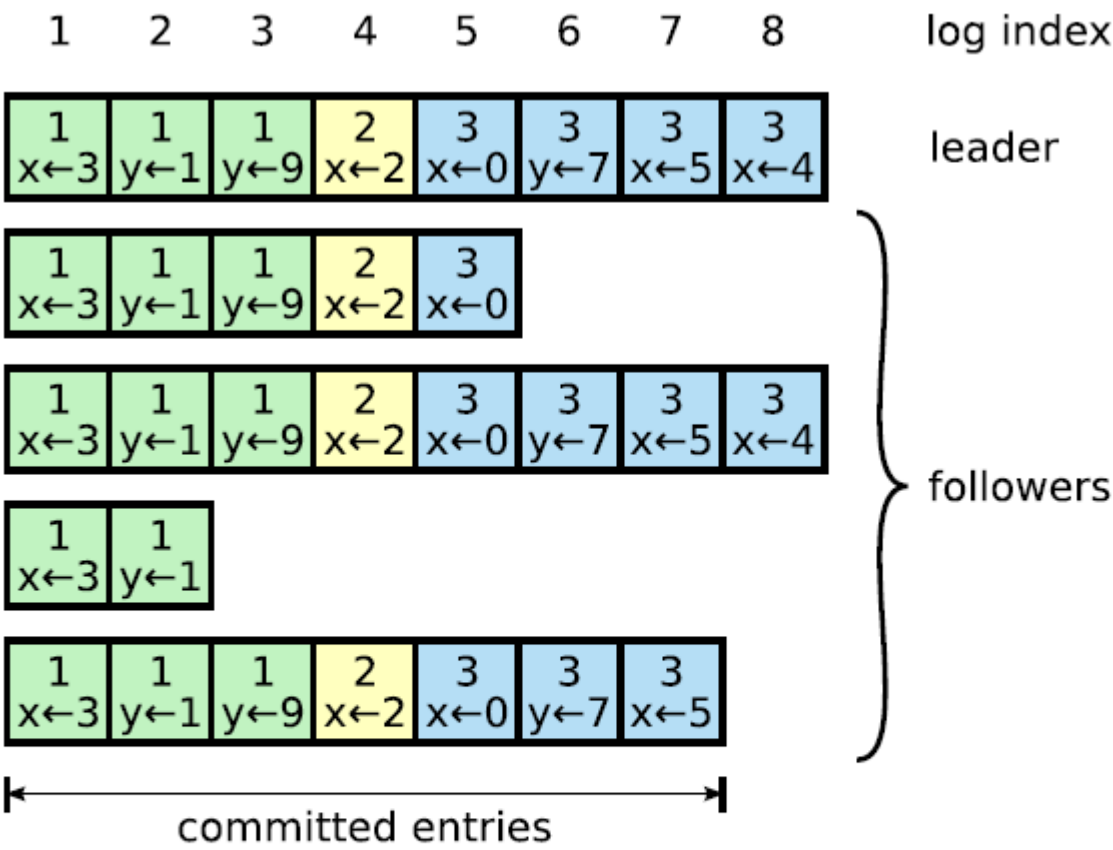


Figure 6: Logs are composed of entries, which are numbered sequentially. Each entry contains the term in which it was created (the number in each box) and a command for the state machine. An entry is considered *committed* if it is safe for that entry to be applied to state machines.

不难看到，logs由顺序编号的log entry组成，每个log entry除了包含command，还包含产生该log entry时的leader term。从上图可以看到，五个节点的日志并不完全一致，raft算法为了保证高可用，并不是强一致性，而是最终一致性，leader会不断尝试给follower发log entries，直到所有节点的log entries都相同。

在上面的流程中，leader只需要日志被复制到大多数节点即可向客户端返回，一旦向客户端返回成功消息，那么系统就必须保证log（其实是log所包含的command）在任何异常的情况下都不会发生回滚。这里有两个词：commit（committed），apply(applied)，前者是指日志被复制到了大多数节点后日志的状态；而后者则是节点将日志应用到状态机，真正影响到节点状态。

The leader decides when it is safe to apply a log entry to the state machines; such an entry is called committed. Raft guarantees that committed entries are durable and will eventually be executed by all of the available state machines. A log entry is committed once the leader that created the entry has replicated it on a majority of the servers

safety

[回到顶部](#)

在上面提到只要日志被复制到majority节点，就能保证不会被回滚，即使在各种异常情况下，这根leader election提到的选举约束有关。在这一部分，主要讨论raft协议在各种各样的异常情况下如何工作的。

衡量一个分布式算法，有许多属性，如

- safety: nothing bad happens,
- liveness: something good eventually happens.

在任何系统模型下，都需要满足safety属性，即在任何情况下，系统都不能出现不可逆的错误，也不能向客户端返回错误的内容。比如，raft保证被复制到大多数节点的日志不会被回滚，那么就是safety属性。而raft最终会让所有节点状态一致，这属于liveness属性。

raft协议会保证以下属性

- Election Safety:** at most one leader can be elected in a given term. §5.2
- Leader Append-Only:** a leader never overwrites or deletes entries in its log; it only appends new entries. §5.3
- Log Matching:** if two logs contain an entry with the same index and term, then the logs are identical in all entries up through the given index. §5.3
- Leader Completeness:** if a log entry is committed in a given term, then that entry will be present in the logs of the leaders for all higher-numbered terms. §5.4
- State Machine Safety:** if a server has applied a log entry at a given index to its state machine, no other server will ever apply a different log entry for the same index. §5.4.3

Figure 3: Raft guarantees that each of these properties is true at all times. The section numbers indicate where each property is discussed.

Election safety

选举安全性，即任一任期内最多一个leader被选出。这一点非常重要，在一个复制集中任何时刻只能有一个leader。系统中同时有多余一个leader，被称之为脑裂（brain split），这是非常严重的问题，会导致数据的覆盖丢失。在raft中，两点保证了这个属性：

- 一个节点某一任期内最多只能投一票；
- 只有获得majority投票的节点才会成为leader。

因此，**某一任期内一定只有一个leader。**

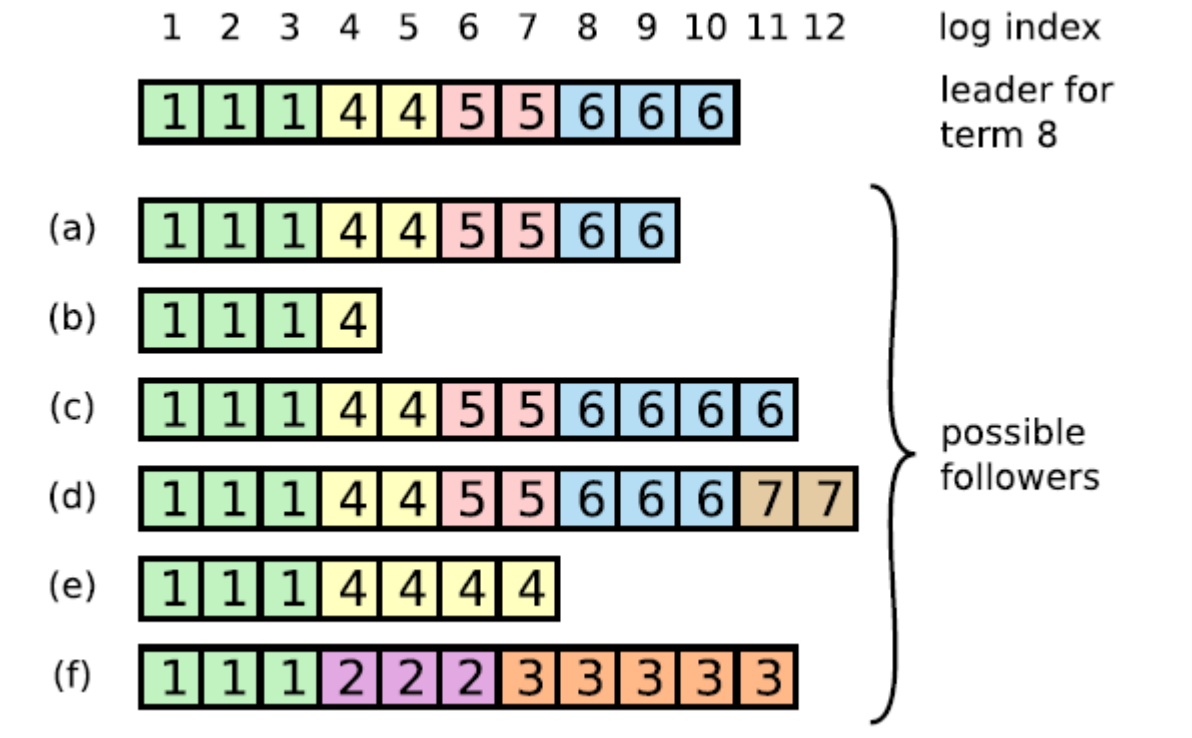
log matching

很有意思，log匹配特性，就是说如果两个节点上的某个log entry的log index相同且term相同，那么在该index之前的所有log entry应该都是相同的。如何做到的？依赖于以下两点

- If two entries in different logs have the same index and term, then they store the same command.
- If two entries in different logs have the same index and term, then the logs are identical in all preceding entries.

首先，leader在某一term的任一位置只会创建一个log entry，且log entry是append-only。其次，consistency check。leader在AppendEntries中包含最新log entry之前的一个log 的term和index，如果follower在对应的term index找不到日志，那么就会告知leader不一致。

在没有异常的情况下，log matching是很容易满足的，但如果出现了node crash，情况就会变得负责。比如下图



注意：上图的a-f不是6个follower，而是某个follower可能存在的六个状态

leader、follower都可能crash，那么follower维护的日志与leader相比可能出现以下情况

- 比leader日志少，如上图中的ab
- 比leader日志多，如上图中的cd
- 某些位置比leader多，某些日志比leader少，如ef（多少是针对某一任期而言）

当出现了leader与follower不一致的情况，leader强制follower复制自己的log

To bring a follower’s log into consistency with its own, the leader must find the latest log entry where the two logs agree, delete any entries in the follower’s log after that point, and send the follower all of the leader’s entries after that point.

leader会维护一个nextIndex[]数组，记录了leader可以发送每一个follower的log index，初始化为eader最后一个log index加1，前面也提到，leader选举成功之后会立即给所有follower发送AppendEntries RPC（不包含任何log entry，也充当心跳消息），那么流程总结为：

- s1 leader 初始化nextIndex[x]为 leader最后一个log index + 1
- s2 AppendEntries里prevLogTerm prevLogIndex来自 logs[nextIndex[x] - 1]
- s3 如果follower判断prevLogIndex位置的log term不等于prevLogTerm，那么返回 false，否则返回True
- s4 leader收到follower的恢复，如果返回值是True，则nextIndex[x] -= 1，跳转到s2. 否则
- s5 同步nextIndex[x]后的所有log entries

leader completeness vs elction restriction

leader完整性：如果一个log entry在某个任期被提交（committed），那么这条日志一定会出现在所有更高term的leader的日志里面。这个跟leader election、log replication都有关。

- 一个日志被复制到majority节点才算committed
- 一个节点得到majority的投票才能成为leader，而节点A给节点B投票的其中一个前提是，B的日志不能比A的日志旧。下面的引文指处了如何判断日志的新旧

voter denies its vote if its own log is more up-to-date than that of the candidate.

If the logs have last entries with different terms, then the log with the later term is more up-to-date. If the logs end with the same term, then whichever log is longer is more up-to-date.

上面两点都提到了majority：commit majority and vote majority，根据Quorum，这两个majority一定是有重合的，因此被选举出的leader一定包含了最新的committed log。raft与其他协议（Viewstamped Replication、mongodb）不同，raft始终保证leade包含最新的已提交的日志，因此leader不会从follower catchup日志，这也大大简化了

19

推荐

0

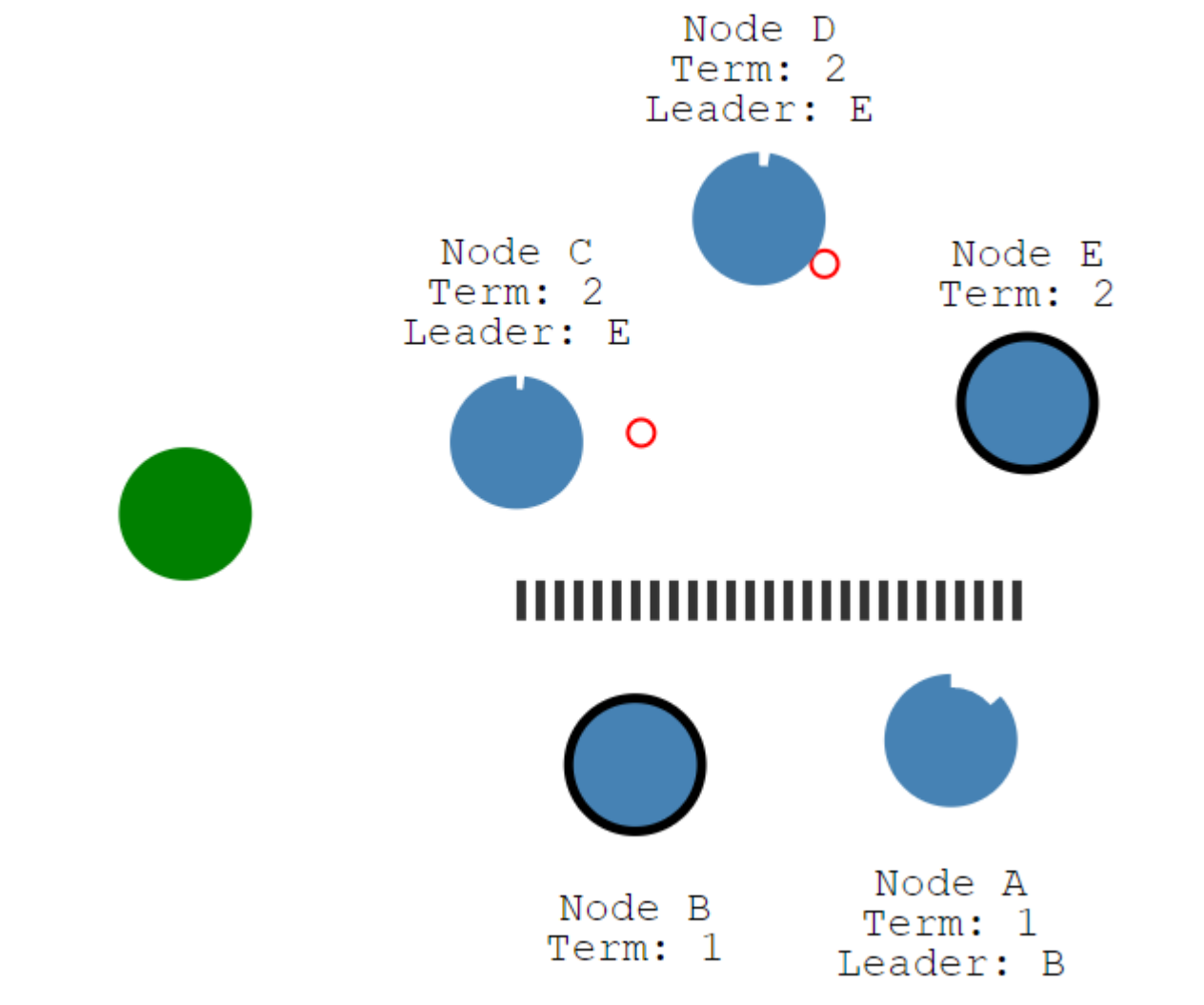
反对

corner case

[回到顶部](#)

stale leader

raft保证Election safety，即一个任期内最多只有一个leader，但在网络分割（network partition）的情况下，**可能会出现两个leader，但两个leader所处的任期是不同的**。如下图所示



系统有5个节点ABCDE组成，在term1，Node B是leader，但Node A、B和Node C、D、E之间出现了网络分割，因此Node C、D、E无法收到来自leader（Node B）的消息，在election time之后，Node C、D、E会分期选举，由于满足majority条件，Node E成为了term 2的leader。因此，在系统中貌似出现了两个leader：term 1的Node B，term 2的Node E，Node B的term 更旧，但由于无法与Majority节点通信，NodeB仍然会认为自己是leader。

在这样的情况下，我们来考虑读写。

首先，如果客户端将请求发送到了NodeB，NodeB无法将log entry 复制到majority节点，因此不会告诉客户端写入成功，这就不会有问题。

对于读请求，stale leader可能返回stale data，比如在read-after-write的一致性要求下，客户端写入到了term2任期的leader Node E，但读请求发送到了Node B。如果要保证不返回stale data，leader需要check自己是否过时了，办法就是与大多数节点通信一次，这个可能会出现效率问题。另一种方式是使用lease，但这就会依赖物理时钟。

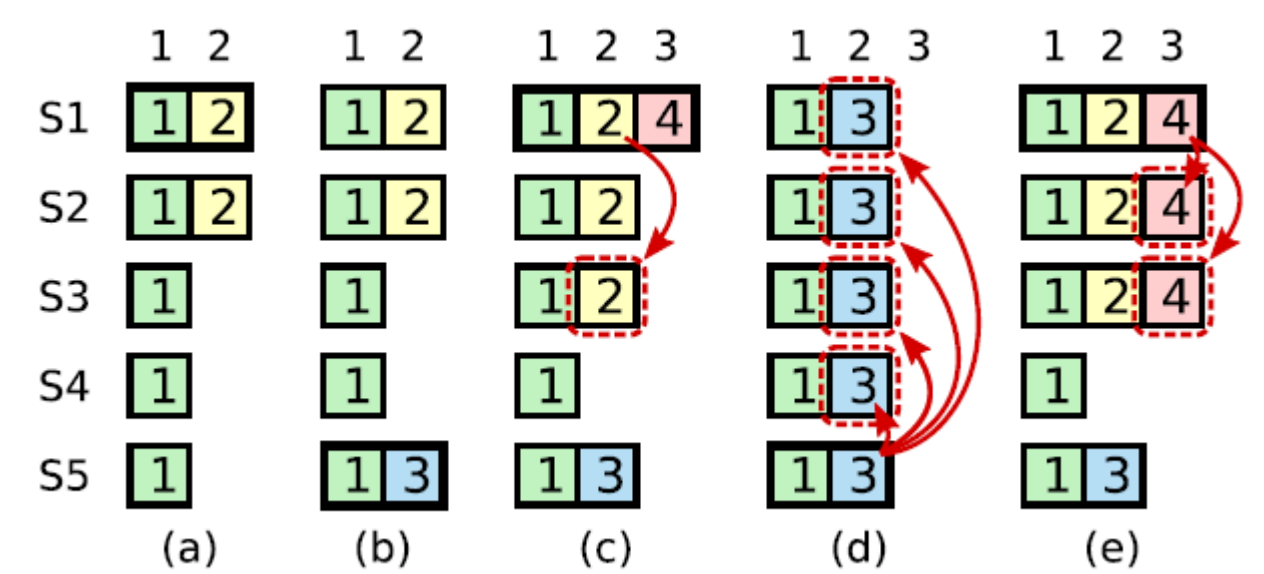
从raft的论文中可以看到，leader转换成follower的条件是收到来自更高term的消息，如果网络分割一直持续，那么stale leader就会一直存在。而在raft的一些实现或者raft-like协议中，leader如果收不到majority节点的消息，那么可以自己step down，自行转换到follower状态。

State Machine Safety

前面在介绍safety的时候有一条属性没有详细介绍，那就是State Machine Safety：

State Machine Safety: if a server has applied a log entry at a given index to its state machine, no other server will ever apply a different log entry for the same index.

如果节点将某一位置的log entry应用到了状态机，那么其他节点在同一位置不能应用不同的日志。简单点来说，所有节点在同一位置（index in log entries）应该应用同样的日志。但是似乎有某些情况会违背这个原则：



上图是一个较为复杂的情况。在时刻(a)，s1是leader，在term2提交的日志只赋值到了s1 s2两个节点就crash了。在时刻（b），s5成为了term 3的leader，日志只赋值到了s5，然后crash。然后在(c)时刻，s1又成为了term 4的leader，开始赋值日志，于是把term2的日志复制到了s3，此刻，可以看出term2对应的日志已经被复制到了majority，因此是committed，可以被状态机应用。不幸的是，接下来（d）时刻，s1又crash了，s5重新当选，然后将term3的日志复制到所有节点，这就出现了一种奇怪的现象：被复制到大多数节点（或者说可能已经应用）的日志被回滚。

究其根本，是因为term4时的leader s1在（C）时刻提交了之前term2任期的日志。为了杜绝这种情况的发生：

Raft never commits log entries from previous terms by counting replicas.
Only log entries from the leader’s current term are committed by counting replicas; once an entry from the current term has been committed in this way, then all prior entries are committed indirectly because of the Log Matching Property.

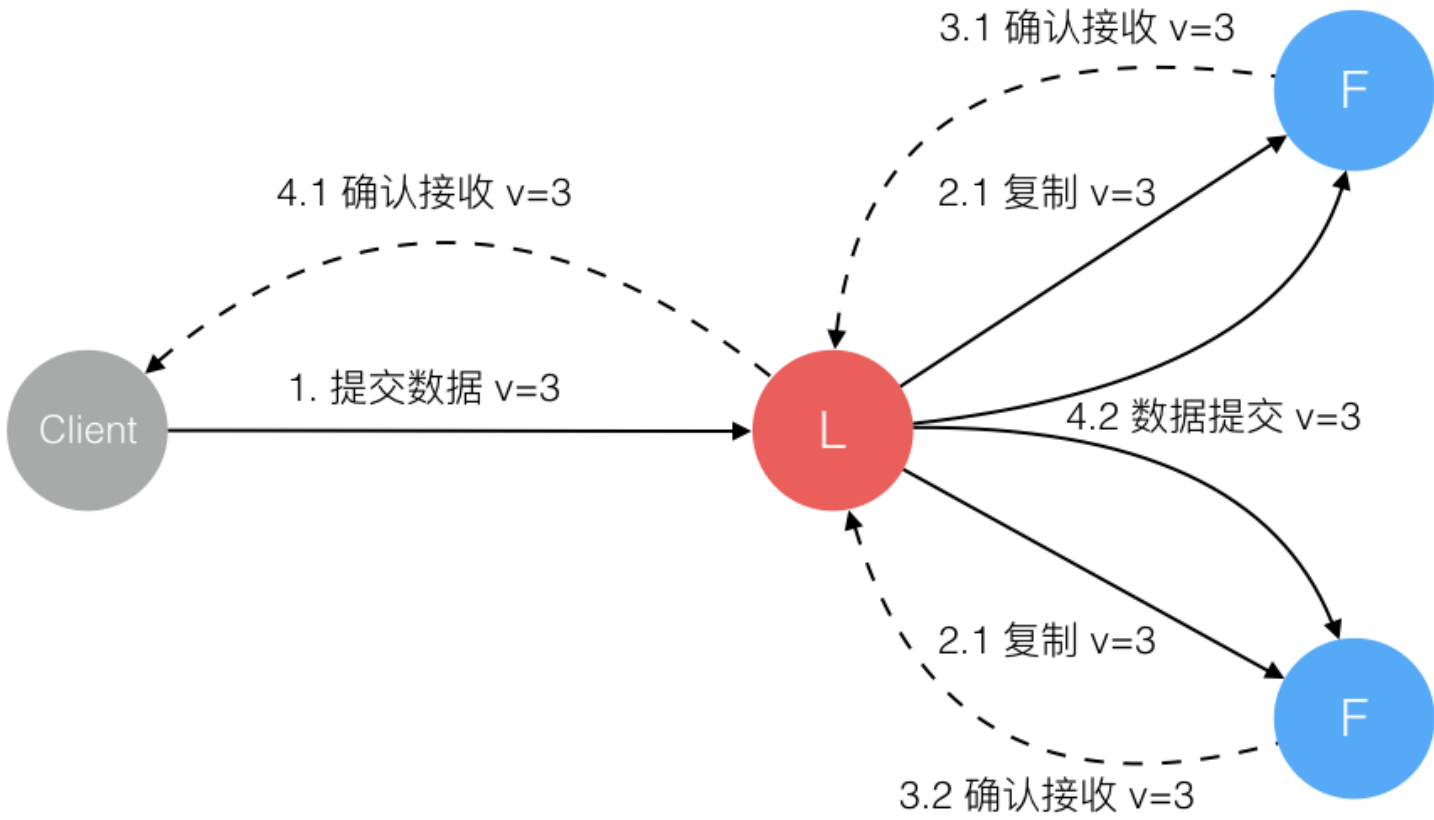
也就是说，某个leader选举成功之后，不会直接提交前任leader时期的日志，而是通过提交当前任期的日志的时候“顺手”把之前的日志也提交了，具体怎么实现了，在log matching部分有详细介绍。那么问题来了，如果leader被选举后没有收到客户端的请求呢，论文中有提到，在任期开始的时候发立即尝试复制、提交一条空的log

Raft handles this by having each leader commit a blank no-op entry into the log at the start of its term.

因此，在上图中，不会出现（C）时刻的情况，即term4任期的leader s1不会复制term2的日志到s3。而是如同(e)描述的情况，通过复制-提交 term4的日志顺便提交term2的日志。如果term4的日志提交成功，那么term2的日志也一定提交成功，此时即使s1crash，s5也不会重新当选。

leader crash

follower的crash处理方式相对简单，leader只要不停的给follower发消息即可。当leader crash的时候，事情就会变得复杂。在[这篇文章](#)中，作者就给出了一个更新请求的流程图。



我们可以分析leader在任意时刻crash的情况，有助于理解raft算法的容错性。

总结

[回到顶部](#)

raft将共识问题分解成两个相对独立的问题，leader election，log replication。流程是先选举出leader，然后leader负责复制、提交log（log中包含command）

为了在任何异常情况下系统不出错，即满足safety属性，对leader election，log replication两个子问题有诸多约束

leader election约束：

- 同一任期内最多只能投一票，先来先得
- 选举人必须比自己知道的更多（比较term，log index）

log replication约束：

- 一个log被复制到大多数节点，就是committed，保证不会回滚
- leader一定包含最新的committed log，因此leader只会追加日志，不会删除覆盖日志
- 不同节点，某个位置上日志相同，那么这个位置之前的所有日志一定是相同的
- Raft never commits log entries from previous terms by counting replicas.

本文是在看完raft论文后自己的总结，不一定全面。个人觉得，如果只是相对raft协议有一个简单了解，看这个[动画演示](#)就足够了，如果想深入了解，还是要看论文，论文中Figure 2对raft算法进行了概括。最后，还是找一个实现了raft算法的系统来看看更好。

references

[回到顶部](#)

<https://web.stanford.edu/~ouster/cgi-bin/papers/raft-atc14>
<https://raft.github.io/>
<http://thesecretlivesofdata.com/raft/>

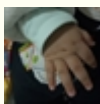
本文版权归作者xybaby（博文地址：<http://www.cnblogs.com/xybaby/>）所有，欢迎转载和商用，请在文章页面明显位置给出原文链接并保留此段声明，否则保留追究法律责任的权利，其他事项，可留言咨询。

标签: [distributed system](#)

好文要顶

关注我

收藏该文



[xybaby](#)
[关注 - 9](#)
[粉丝 - 578](#)
[+加关注](#)

« 上一篇: [设计数据密集型应用第三部分：派生数据](#)
» 下一篇: [Raft与MongoDB复制集协议比较](#)

posted @ 2018-12-17 09:35 [xybaby](#) 阅读(11111) 评论(7) 编辑 收藏

评论

1楼 2018-12-19 00:21 | 幻天芒



搬出我理解Raft的神器：<http://thesecretlivesofdata.com/raft/>

支持(4) 反对(0)

2楼[楼主] 2018-12-19 11:07 | xybaby



@ 幻天芒
是的啊，这个确实牛逼

支持(0) 反对(0)

3楼 2018-12-19 14:56 | PowerShell免费软件



这个算法有大漏洞!!!
貌似mysql5.7同步机遇这个算法，引发了2月前的github惨案。

3主分布式，或者5主分布式理论，大家知道吧？
我升级完善（玩傻^_^）了它。
我发明，或者说完善了分布式原理。大家看看：
估计3年后，有人意识到，应该这么做!!!

---【我发明，完善的，分布式理论】---
1 只有3主，或5主，架构。以下都假设为3主。

2 主只写，从只读。

19

推荐

0

反对

4 集群新建时，都是从。自动投票选主。想人为让某台机子成为主，可以作弊，人为提升其cpu，内存权重。使其成为“高富帅”

4 3主选出后，其他机子都是从。
比如新建1个集群，10台从。经过选主后，成了3主7从。此时并不是性能最好的机子为主。

5 3主中的数据，正常来讲是同步的。若有已同步的从，并且从的内存大，cpu好，并且根据统计现在不忙。
则主从切换，踢出最差的主，从变成主。如此经过几次迭代，cpu+内存+网络最好的机子，成了主。

6 3主中的数据，正常来讲是同步的。若有已同步的从。此时某台主，确认已挂了。就从 从机 中挑个好的，主从切换，从变成主。
挂掉的“高富帅”主回来后，变成从，开始同步。同步完成后，因为你人为提升了它的权重。那么它还能优先成主。
挂掉的“屌丝”主回来后，变成从，开始同步。同步完成后，老老实实排队，想成主，遥遥无期。。。

我们知道，3主可以坏1台，5主可以坏2台。
现有的分布式理论，存在若主少了，抗风险下降。主多了，同步复杂的问题。
用我这个理论，可以补充主。而且实现起来，不难。

7 所有主同时全挂的，世界末日，我不想谈。我想说说，主缓慢的都挂了的事。
若主缓慢的都挂了，从都变成主了。最后只剩3个主了，则集群不能读了。
这里可以设定一个开关，名字叫做【3台时，是否开启读写分离】。
若只有3台主，集群是不能读，还是可以读。

8
若又挂了一个主，剩2个主，还是可以写。
若又挂了一个主，剩1个主，则只读。
老理就是这么设定的，对吧。

欢迎嘲笑登场

支持(0) 反对(0)

4楼 2018-12-19 22:41 | 杰哥很忙

第二种情况，比如有三个节点A B C。A B同时发起选举，而A的选举消息先到达C，C给A投了一票，当B的消息到达C时，已经不能满足上面提到的第一个约束，即C不会给A投票

这里是不是写错了，C不会给B投票吧

支持(0) 反对(0)

5楼[楼主] 2018-12-19 22:52 | xybaby

@ 杰哥很忙
嗯，多谢，已经修正

支持(0) 反对(0)

6楼 2019-04-18 09:51 | 大耳朵石头

请教两个问题：
1) 除了candidate发起选举，哪些场景下follower term值会更新？
2) 网络分区，一个follower被隔离。这个follower的后续行为会怎么样？会切换为candidate不停地发起选举吧？当这个candidate的term大于主分区后整个网络连通了，这个时候会怎么样？

支持(0) 反对(0)

7楼[楼主] 2019-04-20 15:25 | xybaby

@ 大耳朵石头
1) term的更新一定是源于一轮新的选举，不管选举是否成功。对于follower，在其发起选举的时候会自增term；同时，如果收到的消息term大于本地，也会跟新本地term
2) 这个问题在 Four modifications for the Raft consensus algorithm 文章中有详细分析，也可以参考我写的“Raft与MongoDB复制集协议比较”这篇文章

支持(0) 反对(0)

刷新评论 刷新页面 返回顶部

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问网站首页](#)。

【推荐】超50万C++/C#源码：大型实时仿真组态图形源码
【推荐】ALIYUN90%，大陆-港澳专线直连，创业者上云首选。
【推荐】程序员问答平台，解决您开发中遇到的技术难题

相关博文：
· 一文搞懂Raft算法
· 一文搞懂Raft算法
· 一文搞懂Raft算法
· Paxos, Raft, Zab一致性协议-Raft篇
· 解读Raft（二选举和日志复制）