

支持乱序执行的 Raft 协议 *

谷晓松¹, 魏恒峰¹, 乔磊², 黄宇¹

¹(计算机软件新技术国家重点实验室(南京大学),江苏 南京 210023)

²(北京控制工程研究所,北京 100190)

通讯作者: 魏恒峰, Email: hfwei@nju.edu.cn

摘要: PolarFS 是阿里巴巴开发的分布式文件系统, 它实现了分布式共识协议 Raft 的一种变体, 称为 ParallelRaft. ParallelRaft 突破了 Raft 中顺序提交、顺序执行的限制, 允许状态机乱序执行用户命令. 然而, 文献表明, ParallelRaft 并未开源, 仅有简短的文字描述, 更缺乏严格的形式化规约. 更进一步, 它的正确性也尚未经过必要的数学论证或形式化检验. 本文旨在为 ParallelRaft 提供严格的形式化规约并证明其正确性. 具体而言, 本文的主要贡献包括: 首先, 为了理清 ParallelRaft 与 Raft 之间的关系, 我们提出了允许乱序提交、顺序执行的 ParallelRaft-SE(Sequential Execution) 协议, 并建立了从 ParallelRaft-SE 到 Multi-Paxos 的精化关系. 其次, 我们发现现有的 ParallelRaft 描述忽略了可能会违反状态一致性的“幽灵日志”问题. 因此, 我们在 ParallelRaft-SE 的基础上提出了 ParallelRaft-CE(Concurrent Execution) 协议. ParallelRaft-CE 限制了 ParallelRaft-SE 在乱序提交阶段的并行度, 从而避免了“幽灵日志”问题. 我们证明了 ParallelRaft-CE 的正确性. 最后, 我们使用 TLA+ 给出了 ParallelRaft-SE 和 ParallelRaft-CE 的形式化规约, 并使用 TLC 模型检验工具验证了从 ParallelRaft-SE 到 Multi-Paxos 的精化关系以及 ParallelRaft-CE 的正确性.

关键词: Raft; ParallelRaft; Multi-Paxos; 共识协议; TLA+; 精化关系; 模型检验

中图法分类号: TP311

中文引用格式: 谷晓松, 魏恒峰, 乔磊, 黄宇. 支持乱序执行的 Raft 协议. 软件学报. <http://www.jos.org.cn/1000-9825/0000.htm>

英文引用格式: Gu XS, Wei HF, Qiao L, Huang Y. Raft with out-of-order executions. Ruan Jian Xue Bao/Journal of Software, 2020 (in Chinese). <http://www.jos.org.cn/1000-9825/0000.htm>

Raft with Out-of-Order Executions

GU Xiao-Song¹, WEI Heng-Feng¹, QIAO Lei², HUANG Yu¹

¹(State Key Laboratory for Novel Software Technology (Nanjing University), Nanjing 210023, China)

²(Beijing Institute of Control Engineering, Beijing 100190, China)

Abstract: PolarFS is a distributed file system developed by Alibaba Inc. with ultra-low latency and high availability. It implemented a variant of the Raft consensus protocol, called ParallelRaft. ParallelRaft breaks Raft's strict serialization restrictions in the commitment and execution of log entries and enables state machines to commit and execute log entries in an out-of-order way. However, ParallelRaft is not open-sourced. It has only a brief description, lacking formal specification. Moreover, the correctness of ParallelRaft has not been manually proven or formally checked. We aim to provide a precise formal specification for ParallelRaft and prove its correctness. Specifically, we make the following main contributions: First, to clarify the relationship between Raft and ParallelRaft, we proposed ParallelRaft-SE (Sequential Execution), which allows out-of-order commitment but prohibits out-of-order executions. We also established a refinement mapping from

* 基金项目: 国家自然科学基金(61702253, 61772258)

Foundation item: National Natural Science Foundation of China (61702253, 61772258)

收稿时间: 0000-00-00; 修改时间: 0000-00-00; 采用时间: 0000-00-00; jos 在线出版时间: 0000-00-00

CNKI 在线出版时间: 0000-00-00

ParallelRaft-SE to Multi-Paxos. Second, we discovered that ParallelRaft, according to its brief description in the literature, neglects the so-called "ghost log entries" phenomenon, which may violate the consistency among state machines. Therefore, based on ParallelRaft-SE, we proposed ParallelRaft-CE (Concurrent Execution). ParallelRaft-CE avoids the "ghost log entries" phenomenon by limiting parallelism in the commitment of log entries. We proved the correctness of ParallelRaft-CE manually. Finally, we provided the formal specifications of ParallelRaft-SE and ParallelRaft-CE using TLA+ (TLA stands for Temporal Logic of Actions), and verified the refinement mapping from ParallelRaft-SE to Multi-Paxos and the correctness of ParallelRaft-CE using the TLC model checker.

Key words: Raft; ParallelRaft; Multi-Paxos; consensus protocols; TLA+; refinement mapping; model checking

分布式共识问题是分布式计算领域的核心问题.它要求多个参与者对某个值或一系列值(也称序列)达成一致^{[2] [3]}.分布式系统通常采用共识协议提供所需的强一致性,如开源 Ceph^[7]、Google 公司的 Spanner^[4]、Oracle 公司的 MySQL^[37]、腾讯的 PaxosStore^[8] 以及阿里巴巴(Alibaba)的 PolarDB^[13] 等分布式数据存储系统.

Multi-Paxos(Paxos)^{[9] [10]} 与 Raft^[11] 是解决分布式共识问题的两种经典协议.它们都基于复制状态机(Replicated State Machine)^[1] 模型,通过投票、日志复制等机制,保证多副本节点上的状态一致性.每个用户命令都将经历"提交"(Commit)与"执行"(Execute)两个阶段:如果某命令收到来自多数派(Majority)服务器的投票,则称该命令被"提交".只有被提交的命令才能被执行.对于多个命令,这两个阶段都可以以顺序或乱序的方式进行.在乱序方式下,较大日志编号所对应的命令可能会在较小日志编号所对应的命令之前被提交或者被执行.比如,Raft 要求顺序提交、顺序执行.Multi-Paxos 则允许乱序提交,但它仍要求顺序执行.

PolarDB 使用了分布式文件系统 PolarFS^[13] .为了提高系统性能,PolarFS 基于 Raft 实现了允许乱序提交、乱序执行的 ParallelRaft 共识协议.然而,文献中并未给出 ParallelRaft 的严格规约,尤其是缺少对乱序执行机制的完整描述.此外,ParallelRaft 尚未经过必要的数学证明或形式化验证.本文旨在为 ParallelRaft 提供严格的形式化规约,并基于精化关系^[20] 以及数学论证证明其正确性.具体而言,本文的主要贡献如下:

1. 首先,为了理清 ParallelRaft 与 Raft 之间的关系,我们基于 Raft 提出了允许乱序提交、但要求顺序执行的 ParallelRaft-SE(Sequential Execution)协议.ParallelRaft 可以看作 ParallelRaft-SE 的乱序执行版本.我们建立了从 ParallelRaft-SE 到 Multi-Paxos 之间的精化关系(表明 ParallelRaft-SE 是 Multi-Paxos 的一种实现),从而证明了 ParallelRaft-SE 的正确性.
2. 其次,我们发现文献中描述的 ParallelRaft 乱序执行机制忽略了可能会破坏状态一致性的"幽灵日志"问题.已有研究表明,在顺序执行方式下(如 Raft 或 Multi-Paxos),"幽灵日志"问题不会破坏副本间的状态一致性.相反地,我们将论证,在乱序执行方式下,它可能造成副本间状态不一致.更进一步,我们指出针对 Raft 和 Multi-Paxos 的解决方案无法解决乱序执行方式下的状态机的一致性问题.
3. 再次,我们在 ParallelRaft-SE 的基础上,提出了支持乱序执行的 ParallelRaft-CE(Concurrent Execution)协议.通过限制 ParallelRaft-SE 在乱序提交阶段的并行度,ParallelRaft-CE 避免了"幽灵日志"问题.我们证明了 ParallelRaft-CE 的正确性.
4. 最后,我们使用 TLA+^{[14] [15] [16]} 语言给出了,ParallelRaft-SE 与 ParallelRaft-CE 的形式化规约,并使用 TLC^[19] 模型检验工具验证了从 ParallelRaft-SE 到 Multi-Paxos 的精化关系以及 ParallelRaft-CE 的正确性.

本文组织如下:第 1 节介绍预备知识,包括 TLA+形式化规约语言、分布式共识问题以及 Multi-Paxos 与 Raft 共识协议.第 2 节介绍 ParallelRaft-SE 协议,并建立从 ParallelRaft-SE 到 Multi-Paxos 的精化关系.第 3 节分析在乱序执行方式下,"幽灵日志"问题对状态一致性的影响,并指出现有解决方案的不足.第 4 节介绍允许乱序执行且避免了"幽灵日志"问题的 ParallelRaft-CE 协议.第 5 节简要证明 ParallelRaft-CE 的正确性.第 6 节使用 TLC 模型检验工具验证从 ParallelRaft-SE 到 Multi-Paxos 的精化关系以及 ParallelRaft-CE 的正确性.第 7 节讨论相关工作.第 8 节总结全文并讨论可能的未来工作.完整的 TLA+规约与模型检验结果参见 GitHub 仓库^[36] .

1 预备知识

1.1 TLA+简介

TLA+是由 Leslie Lamport 基于时序逻辑(Temporal Logic of Actions)^[16] 开发的形式化规约语言^[14] ,尤其适合于描述并发系统或分布式协议.

一个 TLA+规约包含一组变量、初始状态(initial states)和一组动作(actions),通常表示为 $Spec = Init \wedge [Next]_{vars}$,其中, $vars$ 是所有变量的集合.状态(state)是对所有变量的赋值. $Init$ 谓词定义了系统的初始状态, $Next$ 是所有动作的析取式,定义了状态之间的转换关系. $[Next]_{vars}$ 为真当且仅当 $Next$ 为真(某个动作为真,即某个动作被执行)或者所有变量的值保持不变.行为(behavior)是由状态构成的序列. TLA+使用不带撇的变量表示当前状态中变量的值,用带撇的变量表示新状态中的值.这样,动作使用一个包含带撇变量与不带撇的变量的公式描述.例如,动作 $x' = x + 1$ 表示变量 x 在新状态中的值比旧状态中的值大 1.

TLA+支持一阶谓词逻辑以及 ZF(Zermelo-Fraenkel)集合论,可以表达丰富的数据类型.图 1 总结了本文用到的逻辑与集合操作符(operator).文献^[17] 给出了完整的 TLA+操作符列表.

	操作符	含义
逻辑	CHOOSE $x \in S : p$	选择集合 S 中满足条件 p 的元素 x (通常用于符合条件的 x 是唯一的情况下)
集合	SUBSET S	S 的幂集
	$\{e : x \in S\}$	将 e 作用在 S 中的所有元素得到的集合, 如 $\{x^2 : x \in S\}$
	$\{x \in S : p\}$	S 中满足条件 p 的元素构成的集合
函数	$f[e]$	函数 f 作用在参数 e 上
	$[x \in S \mapsto e]$	对于 $x \in S$, 使得 $f(x) = e$ 的函数
	$[f \text{ EXCEPT } !e_1 = e_2]$	函数 $\bar{f} : \bar{f}[e] = \begin{cases} e_2, & \text{if } e = e_1 \\ f[e], & \text{otherwise} \end{cases}$
记录	$e.h$	记录 e 的域 h
	$[h_1 \mapsto e_1, \dots, h_n \mapsto e_n]$	域 h_i 为 e_i 的记录
	$[h_1 : S_1, \dots, h_n : S_n]$	满足域 h_i 属于 S_i 的所有记录构成的集合
	$[r \text{ EXCEPT } !h = e]$	记录 $\bar{r} [h_1 : S_1, \dots, h : e, \dots, h_n : S_n]$
	$[r \text{ EXCEPT } !h = e]$, 其中 e 包含符号 @	e 中的 @ 表示 $r.h$
元组	$e[i]$	元组 e 的第 i 个分量
动作操作符	e'	动作的新状态中 e 的值
	UNCHANGED e	e 不变: $e' = e$
	$[A]_e$	动作 A 成立或者 e 不变: $A \vee (e' = e)$
时序操作符	$\square F$	F 在所有情况下均成立 (\square 表示 “always”)
	$\diamond F$	F 最终成立 (\diamond 表示 “eventually”)

Figure 1 A summary of the TLA+ operators used in this paper

图 1 本文使用的 TLA+操作符

TLA+ 规约以模块(module)的形式组织在一起. 每个模块中, 可以声明常量(CONSTANTS)与变量(VARIABLES), 定义操作符(operator)或者提出定理(THEOREM). 一个模块 M 可以通过扩展(EXTEND)其他模块 M_1, \dots, M_n 引入声明、定义与定理; 在 M 中, 写作 EXTENDS M_1, \dots, M_n , 模块可以被实例化. 比如:

$$IM_1 \triangleq \text{INSTANCE } M_1 \text{ WITH } p_1 \leftarrow e_1, \dots, p_n \leftarrow e_n.$$

其中 p_i 包含 M_1 中的所有常量与变量, e_i 是 M 中的合法表达式. 该语句将 M_1 中的 p_i 替换为相应的 e_i . 我们可以通过 $IM_1!F$ 访问模块 M_1 中的表达式 F . TLA+的隐式替换规则允许我们在 e_j 与 p_j 相同时省略 $p_j \leftarrow e_j$ 子句.

TLC^[19] 是 TLA+的模型检验工具. 它可以遍历所有可能的系统行为, 检查所有的状态, 验证系统是否满足特

定性质.然而有些分布式系统包含无穷的状态,比如 TLA+规约中常常包含自然数,而自然数是无穷的.为了验证这样的系统,TLA+引入模型(model).模型所有的集合都是有穷的,因此系统的状态数目是有穷的.模型检验常常面临组合爆炸的问题,为此,TLC 可以利用模型中的对称性缩减状态空间.例如,假设 CONSTANTS Server 定义了系统中所有进程的集合.在进行模型检验时,我们需要将它实例化为一个有穷的集合,比如 $Server = \{S_1, S_2, S_3\}$.如果 S_1, S_2, S_3 之间的任意排列(比如 S_1 替换 S_2, S_2 替换 S_3, S_3 替换 S_1)都不会影响系统规约满足给定的性质,那么我们可以将 $Server$ 设置为对称集(symmetry set)^{[14] [18]}.

在 TLA+中,精化(refinement)关系^[20] 就是逻辑蕴含(logical implication)关系.例如,考虑模块 $AbsModule$ 中的 $AbsSpec$ 规约以及模块 $ImplModule$ 中的 $ImplSpec$ 规约. $AbsSpec$ 包含变量 $x_1, \dots, x_m, y_1, \dots, y_n$, $ImplSpec$ 包括变量 $x_1, \dots, x_m, z_1, \dots, z_p$.令 X 、 Y 与 Z 分别代表变量组 x_1, \dots, x_m 、 y_1, \dots, y_n 与 z_1, \dots, z_p .为了验证 $ImplSpec$ 实现/精化了 $AbsSpec$ (也称 $ImplSpec$ 是 $AbsSpec$ 的精化;即 $ImplSpec \Rightarrow AbsSpec$),我们需要证明,对于每一个满足 $ImplSpec$ 的行为,都存在一种关于变量组 Y 的赋值方式,使得赋值后所得的行为(考虑变量组 X 与 Y)满足 $AbsSpec$ ^[20].具体而言,对于每个 y_i ,我们仅使用 X 与 Z 定义表达式 \bar{y}_i 在 $AbsSpec$ 中做替换 $y_i \leftarrow \bar{y}_i$,得到新的规约 $\overline{AbsSpec}$,然后证明 $ImplSpec$ 实现/蕴含了 $AbsSpec$.这里,替换 $y_i \leftarrow \bar{y}_i$ 被称为精化映射(refinement mapping).为了使用 TLC 检验在精化映射 $y_i \leftarrow \bar{y}_i$ 下, $ImplSpec$ 实现了 $AbsSpec$,我们在模块 $ImplModule$ 中添加定义 $AbsSub \triangleq INSTANCE AbsModule WITH y_1 \leftarrow \bar{y}_1, \dots, y_n \leftarrow \bar{y}_n$,并检验定理 $THEOREM ImplSpec \Rightarrow AbsSub! AbsSpec$.

1.2 分布式共识

分布式共识要求多个副本服务器节点保持状态一致.每个服务器节点可以建模为一个复制状态机,通过执行用户命令进行状态转换.

通常使用复制日志的机制实现复制状态机.每个服务器保存一份日志.日志由按顺序编号(通常是自然数)的日志项组成.每个日志项存储了一条来自用户的命令.服务器每次按顺序从日志中读取下一条已经取得共识的命令,在状态机上执行,并将结果返回给用户.由于服务器按照编号顺序执行一条命令,我们称此为顺序执行.传统的分布式共识协议(Multi-Paxos,Raft)都是顺序执行.在此条件下,我们假设副本服务器上的复制状态机具有相同的初始状态,那么只需要保证日志之间的一致性,就可以保证服务器之间的状态一致.因此,分布式共识问题可以转化为保证不同副本服务器节点上日志之间的一致性.具体而言,它要求

- 非平凡性(Nontriviality):只能对用户发出的命令取得共识.
- 一致性(Consistency):每个位置最多只能对一条命令取得共识.

服务器通过运行共识协议保证日志的一致性.本文考虑异步消息传递系统中的共识协议,其故障模型(failure model)如下:

- 服务器可能终止执行(fail by stop),但不会出现拜占庭错误^[9].
- 消息可能被延时、乱序到达、丢失或重复,但消息的内容不会被篡改.

1.3 Paxos协议

Paxos 协议是解决分布式共识问题的经典协议,它允许一组服务器对单个值(即单个日志项)取得共识,是 Multi-Paxos 的基础.Paxos 中定义了三种角色:提议者(Proposer),从节点(Acceptor)和学习者(Learner).提议者提出值(value),接受者选择值,学习者学习被选中的值.Paxos 包括两个阶段,每个阶段又包括两个子阶段^[9]:

- 准备(Prepare)阶段(也称第一阶段,Phase1)
 - 子阶段 Phase1a:提议者选择一个全局唯一的编号 b (通常是自然数),向所有接受者发送编号为 b 的 Prepare 请求.
 - 子阶段 Phase1b:接受者收到编号为 b 的 Prepare 请求.如果它之前收到过编号大于 b 的 Prepare 请求,则忽略该编号为 b 的请求;否则,接受者将自己接受过的提议(包括编号与值)中编号最大的回复给提议者.
- 接受(Accept)阶段(也称第二阶段,Phase2)

- 子阶段 Phase2a: 提议者收到来自超过半数的接受者发送的针对编号为 b 的 Prepare 请求的回复。如果回复中不包含任何提议,则提议者可以任选一个值(通常是该提议者收到的用户命令);否则,提议者选择编号最大的回复对应的值。设提议者选择的值为 v ,则提议者向所有接受者发送内容为 (b, v) 的 Accept 请求。
- 子阶段 Phase2b: 接受者收到编号为 b 的 Accept 请求。如果接受者没有收到过编号大于 b 的请求,则可以接受该 Accept 请求。否则,接受者忽略该 Accept 请求。

1.4 Multi-Paxos 协议

Multi-Paxos 针对每个日志项运行一个独立的 Paxos 实例,从而支持副本服务器对日志(日志项的序列)达成共识。由于各个 Paxos 实例是并行独立运行的,Multi-Paxos 允许乱序提交用户命令,也就是说,它允许副本服务器先就编号较大的日志项对应的用户命令达成共识,而无需等待之前的日志项完成共识。

在实际实现时,常常通过将 phase1 的消息分批(batching)的方式提高性能。此时,系统包含一个恢复阶段,恢复阶段本质上是将不同实例的 phase1 消息集中处理的过程。

模块 *MultiPaxos* 描述了 Multi-Paxos 协议,其中包括三个常量:

- *Acceptors*: 所有接受者构成的集合。
- *Value*: 所有可能的提议值构成的集合。
- *Nil*: 特殊记号,不属于 *Value*

我们使用自然数表示可能的提议编号以及每个实例的编号,即 $\text{Ballots} \triangleq \text{Nat}$, $\text{Instances} \triangleq \text{Nat}.\text{Quorums}$ 定义了一种特殊的议会系统(quorum system):它的每个元素都是由超过半数的接受者构成的集合。

Multi-Paxos 的规约中包括六个变量:

- *ballot*: $\text{ballot}[a]$ 表示接受者 a 记录的最大的提议编号,也是 a 可以接受的最小的提议编号。
- *vote*: $\text{vote}[a][i][b]$ 表示接受者 a 对编号为 i 的实例在 a 的提议编号为 b 时接受的提议值。如果 a 在提议编号为 b 时,没有对编号为 i 的实例接受任何值,则 $\text{vote}[a][i][b] = \text{Nil}$ 。
- *leaderVote*: $\text{leaderVote}[b][i]$ 为提议者在提议编号为 b 时,为编号为 i 的实例提出的值与提议编号 b 组成的二元组。初始状态下,对所有的提议编号 b 与实例编号 i , $\text{leaderVote}[b][i] = \langle -1, \text{Nil} \rangle$ 。
- *1amsgs, 1bmsgs, 2amsgs*: 不同类型的消息集合。

```
MODULE MultiPaxos
EXTENDS Integers, FiniteSets
CONSTANTS Acceptors, Nil, Value
Ballots  $\triangleq$  Nat
Instances  $\triangleq$  Nat
Quorums  $\triangleq$  { $Q \in \text{SUBSET } \text{Acceptors} : \text{Cardinality}(Q) > \text{Cardinality}(\text{Acceptors}) \setminus 2$ }
Max( $s$ )  $\triangleq$  CHOOSE  $x \in s : \forall y \in s : x \geq y$ 
VARIABLES ballot, vote, leaderVote, 1amsgs, 1bmsgs, 2amsgs
```

协议定义的主要动作包括:

- *Phase1a(b)*: 对应 Paxos 的 *Phase1a* 阶段。提议者选取提议编号 b , 并且向其他节点发送 Prepare 请求, 提议的编号为 b 。
- *Phase1b(a, b)*: 对应 Paxos 的 *Phase1b* 阶段。节点 a 收到编号为 b 的投票请求。若 $b > \text{ballot}[a]$, a 设置 $\text{ballot}[a]$ 为 b , 并且对每一个实例编号 i , a 将它接受过的提议编号最大的提议值与提议编号回复给提议者。

$$\begin{aligned}
\vdash & \\
Phase1a(b) &\triangleq \\
&\wedge \textit{1ams} = \textit{1ams} \cup \{\langle b \rangle\} \\
&\wedge \text{UNCHANGED } \langle \textit{ballot}, \textit{vote}, \textit{leaderVote}, \textit{1bmsgs}, \textit{2ams} \rangle \\
MaxAcceptorVote(a, i) &\triangleq \\
&\text{LET } \textit{maxBallot} \triangleq \textit{Max}(\{b \in \textit{Ballots} : \textit{vote}[a][i][b] \neq \textit{Nil}\} \cup \{-1\}) \\
&\quad v \triangleq \text{IF } \textit{maxBallot} > -1 \text{ THEN } \textit{vote}[a][i][\textit{maxBallot}] \text{ ELSE } \textit{Nil} \\
&\text{IN } \langle \textit{maxBallot}, v \rangle \\
Phase1b(a, b) &\triangleq \\
&\wedge \textit{ballot}[a] < b \\
&\wedge \langle b \rangle \in \textit{1ams} \\
&\wedge \textit{ballot}' = [\textit{ballot} \text{ EXCEPT } ![a] = b] \\
&\wedge \textit{1bmsgs}' = \textit{1bmsgs} \cup \\
&\quad \{\langle b, \{\langle i, \text{MaxAcceptorVote}(a, i) \rangle : i \in \textit{Instances}\}, a \rangle\} \\
&\wedge \text{UNCHANGED } \langle \textit{vote}, \textit{leaderVote}, \textit{1ams}, \textit{2ams} \rangle \\
IncreaseBallot(a, b) &\triangleq \\
&\wedge \textit{ballot}[a] < b \\
&\wedge \textit{ballot}' = [\textit{ballot} \text{ EXCEPT } ![a] = b] \\
&\wedge \text{UNCHANGED } \langle \textit{vote}, \textit{leaderVote}, \textit{1ams}, \textit{1bmsgs}, \textit{2ams} \rangle
\end{aligned}$$

- *Merge(b)*:与*Propose(b, i)*,*Phase2a(b, i)*动作一起对应于 Paxos 的*Phase2a*阶段.在*Merge*动作中,提议者根据收到的投票回复更新自己的日志,但是不发起 Accept 请求.当发起提议编号为**b**的 Prepare 请求的提议者收到一个多数派的回复时,对每一个实例编号*i*,根据回复(方法如 Paxos 协议的 Phase2a 阶段描述)选取提议值并保存在*leaderVote[b][i]*中.值得注意的是,只有当*leaderVote[b][i]*之前没有被修改过时(*leaderVote[b][i] = (-1, Nil)*)时才能更新.因为 Multi-Paxos 中每个提议编号只能有一个提议者发起投票,且对于每个实例编号,它只能提出一个提议值.
- *Propose(b, i)*:提议者对实例编号*i*提出提议值.若*leaderVote[b][i] = (-1, Nil)*,那么提议者选取一个合法的提议值*v*,更新*leaderVote[b][i] = (b, v)*,并发起请求.否则提议者直接将*leaderVote[b][i]*发送给接受者.
- *Phase2a(b, i)*:提议者根据*Merge(b)*和*Propose(b, i)*中对实例编号*i*确定的提议值,发起 Accept 请求.
- *Vote(a, b, i)*:对应于 Paxos 的*Phase2b*阶段.接受者*a*收到提议编号为*b*,实例编号为*i*的 Accept 请求时,若*b*不小于*a*的提议编号(*b ≥ ballot[a]*),则*a*接受 Accept 请求,并更新*vote[a][i][b]*为对应的提议值.

$$\begin{aligned}
& \text{---} \\
& 1bMsgs(b, Q) \triangleq \{m \in 1bmsgs : m[3] \in Q \wedge m[1] = b\} \\
& MaxVote(b, i, Q) \triangleq \\
& \quad \text{LET } entries \triangleq \text{UNION } \{m[2] : m \in 1bMsgs(b, Q)\} \\
& \quad ientries \triangleq \{e \in entries : e[1] = i\} \\
& \quad maxBal \triangleq \text{Max}(\{e[2][1] : e \in ientries\}) \\
& \text{IN CHOOSE } v \in Value \cup \{Nil\} : \exists e \in ientries : \\
& \quad \wedge e[2][1] = maxBal \wedge e[2][2] = v \\
& lastInstance(b, Q) \triangleq \text{LET } entries \triangleq \text{UNION } \{m[2] : m \in 1bMsgs(b, Q)\} \\
& \quad valid \triangleq \{e \in entries : e[2][1] \neq -1\} \\
& \quad \text{IN IF } valid = \{\} \text{ THEN } -1 \text{ ELSE } \text{Max}(\{e[1] : e \in valid\}) \\
& Merge(b) \triangleq \wedge \exists Q \in Quorums : \\
& \quad \wedge \forall a \in Q : \exists m \in 1bMsgs(b, Q) : m[3] = a \\
& \quad \wedge leaderVote' = [leaderVote \text{ EXCEPT } !(b) = [i \in Instances \mapsto \\
& \quad \quad \text{IF } (\wedge i \in 0 .. lastInstance(b, Q) \\
& \quad \quad \wedge leaderVote[b][i][1] = -1) \\
& \quad \quad \text{THEN } \langle b, MaxVote(b, i, Q) \rangle \\
& \quad \quad \text{ELSE } leaderVote[b][i])] \\
& \quad \wedge \text{UNCHANGED } \langle vote, ballot, 1amsgs, 1bmsgs, 2amsgs \rangle \\
& Propose(b, i) \triangleq \wedge leaderVote[b][i][1] = -1 \\
& \quad \wedge \exists Q \in Quorums : \\
& \quad \wedge \forall a \in Q : \exists m \in 1bMsgs(b, Q) : m[3] = a \\
& \quad \wedge \text{LET } maxV \triangleq MaxVote(b, i, Q) \\
& \quad \quad safe \triangleq \text{IF } maxV \neq Nil \text{ THEN } \{maxV\} \text{ ELSE } \\
& \quad \quad \quad Value \cup \{Nil\} \\
& \quad \text{IN } \exists v \in safe : leaderVote' = [leaderVote \text{ EXCEPT } \\
& \quad \quad !(b)[i] = \langle b, v \rangle] \\
& \quad \wedge \text{UNCHANGED } \langle vote, ballot, 1amsgs, 1bmsgs, 2amsgs \rangle \\
& Phase2a(b, i) \triangleq \\
& \quad \wedge leaderVote[b][i][1] = b \\
& \quad \wedge 2amsgs' = 2amsgs \cup \{\langle b, i, leaderVote[b][i] \rangle\} \\
& \quad \wedge \text{UNCHANGED } \langle ballot, vote, leaderVote, 1amsgs, 1bmsgs \rangle \\
& Vote(a, b, i) \triangleq \\
& \quad \wedge ballot[a] \leq b \\
& \quad \wedge ballot' = [ballot \text{ EXCEPT } !(a) = b] \\
& \quad \wedge \exists m \in 2amsgs : \\
& \quad \quad \wedge m[2] = i \wedge m[1] = b \\
& \quad \quad \wedge vote' = [vote \text{ EXCEPT } !(a)[i][b] = m[3][2]] \\
& \quad \wedge \text{UNCHANGED } \langle leaderVote, 1amsgs, 1bmsgs, 2amsgs \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{---}
\end{aligned}$$

Next 定义了次态关系。*Spec* 定义了完整的行为规约。

$$\begin{aligned}
 Next &\triangleq \\
 &\vee \exists a \in \text{Acceptors}, b \in \text{Ballots} : \text{IncreaseBallot}(a, b) \\
 &\vee \exists b \in \text{Ballots} : \text{Phase1a}(b) \\
 &\vee \exists a \in \text{Acceptors}, b \in \text{Ballots} : \text{Phase1b}(a, b) \\
 &\vee \exists b \in \text{Ballots} : \text{Merge}(b) \\
 &\vee \exists b \in \text{Ballots}, i \in \text{Instances} : \text{Propose}(b, i) \\
 &\vee \exists b \in \text{Ballots}, i \in \text{Instances} : \text{Phase2a}(b, i) \\
 &\vee \exists a \in \text{Acceptors}, b \in \text{Ballots}, i \in \text{Instances} : \text{Vote}(a, b, i)
 \end{aligned}$$

$$Spec \triangleq Init \wedge \square[Next]_{\langle leaderVote, ballot, vote, 1amsgs, 1bmsgs, 2amsgs \rangle}$$

1.5 Raft协议

Raft 是一种更易于理解的分布式共识协议^[11]。它加强了日志项之间的串行性,简化了协议的设计。Raft 中的每个节点都维护一个递增的变量,称为任期(term)。任期本质上是节点共同维护的逻辑时钟,通过任期,节点可以发现过时的消息。具体而言,节点间发送消息时,需携带自身当前的任期。如果节点收到的消息携带的任期值小于该节点自身当前的任期值,则拒绝该消息。否则,则更新自身的任期值。节点向日志添加新的日志项时,会将自身当前的任期值也保存在日志项中,这成为该日志项的任期。

Raft 中节点有三种角色:主节点(Leader)从节点(Follower)和候选节点(Candidate)。初始状态下,所有的节点都是从节点。Raft 协议主要包括两部分:选举主节点以及主节点向从节点同步日志。正常情况下,主节点通过定时向其它节点发送心跳信号来维护自己的主节点身份。当从节点一段时间内没有收到心跳信号时,便转变为候选节点。候选节点首先增大自身任期,然后将任期发送给所有节点,以发起选举。收到选举请求后,节点将比较自身任期和选举请求携带的任期。如果自身任期较大,或者自己在本任期内已为其它候选节点投过票,则拒绝此次选举请求。收到多数派投票的候选节点将成为新的主节点。Raft 的多数派投票机制确保选举的安全性(Election Safety)^[11]:每个任期内最多只有一个主节点。为了保证新任主节点的完全性(Leader Completeness),即它的日志应包含所有已被提交的日志项,Raft 引入了如下规则:如果候选节点的日志比自身的日志旧^[11],那么节点将拒绝该选举请求。节点间可以通过比较日志中编号最大的日志项的编号与任期(分别称为lastIndex 和 lastTerm)判断日志的新旧^[11]。

主节点按照日志编号顺序向从节点同步日志项。从节点需要按照编号顺序接受主节点的日志项。在收到编号更小的日志项之前,从节点不能接受编号更大的日志项。主节点与从节点间通过确认(ack)机制维护下一个可以接受的编号。与 Multi-Paxos 中不同实例独立发送和接受日志项不同,Raft 在所有的日志项之间通过编号建立了全序关系。通过这样的限制,Raft 中节点的日志不会出现空洞,同时保证了节点间日志的一致性(Log Matching Property)^[11]:如果两个节点上的日志在同一位置拥有相同日志项,那么之前所有位置上的日志项也必定相同。根据日志的一致性,如果某个日志项已被提交,那么日志中所有编号更小的日志项也已被提交。

2 ParallelRaft-SE 协议及精化

Raft 要求顺序提交、顺序执行用户命令。这种限制使得它不适用于高并发系统^[13]。因此,阿里巴巴公司基于 Raft 提出了 ParallelRaft,它允许乱序提交、乱序执行用户命令^[13]。为了理清 ParallelRaft 与 Raft 之间的关系,我们提出了 ParallelRaft-SE。ParallelRaft-SE 允许乱序提交,但仍要求顺序执行用户命令。因此,可以将 ParallelRaft-SE 视为顺序提交版本的 ParallelRaft。为了证明 ParallelRaft-SE 的正确性,我们建立了从它到 Multi-Paxos 的精化关系。

2.1 ParallelRaft-SE协议

ParallelRaft-SE 沿用了 Raft 的角色(主节点、从节点、候选节点)、任期的概念以及选主机制.ParallelRaft-SE 协议主要包括三部分:主节点向从节点同步日志,主节点的选取和日志恢复.

在 ParallelRaft-SE 中,主节点可以并发地向从节点发送多个日志项.从节点收到日志项后立即接受并确认,而无需等待编号更小的日志项.因此,ParallelRaft-SE 支持日志的乱序接受和提交.

ParallelRaft-SE 的选举机制与 Raft 基本相同,都需要保证选举的安全性.不同的是,由于 ParallelRaft-SE 不具有 Raft 的串行性,节点的日志中可能有空洞,因此在选举过程中,ParallelRaft-SE 无法通过节点比较日志的新旧保证新任主节点的完全性.因此,ParallelRaft-SE 加入了日志恢复阶段.

ParallelRaft-SE 的日志恢复过程的基本思想是新任主节点收集其它节点的日志,并运行 Paxos 协议恢复那些可能已被提交但新任主节点缺失的日志项.经过恢复,新任主节点满足完全性,此时可以向从节点同步日志项.

ParallelRaft-SE 的规约中使用了如下的常量与变量(仅介绍相对于 Multi-Paxos 模块新引入的部分):

- *Server*是所有参与共识的节点的集合.
- *Follower, Candidate, Leader*是服务器的三种不同状态.
- *r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, r3amsgs, negMsgs*是不同类型的消息的集合.
- *currentTerm[i]*为节点*i*所记录的最大的任期.
- *currentState[i]*为节点*i*的状态,在任何时候都是*Follower, Candidate, Leader*之一.
- *vote[i][n][t]*表示节点*i*在任期*t*接受的序号为*n*的日志项.引入*vote*是为了构建 ParallelRaft-SE 到 Multi-Paxos 的精化,实际协议中不需要*vote*.
- *leaderLog*记录了每个任期的主节点的日志.*leaderLog*同样是为了构建 ParallelRaft-SE 到 Multi-Paxos 的精化,实际协议中不需要.*leadrLog[t][n]*表示任期*t*的主节点的日志中的编号为*n*的日志项.每个日志项是形如(*t', v, b*)的三元组.其中*t'*是日志项的任期,*v*是提议值,*b*是一个布尔值,当且仅当日志项被提交时,*b*为真.
- *log[i][n]*为节点*i*的日志中编号为*n*的日志项.

MODULE *ParallelRaft - SE*

```

EXTENDS Integers, FiniteSets, Sequences, TLC
CONSTANTS Server, Follower, Candidate, Leader, Nil, Value

Quorums  $\triangleq$  { $i \in \text{SUBSET } (\text{Server}) : \text{Cardinality}(i) * 2 > \text{Cardinality}(\text{Server})$ }
Index  $\triangleq$  Nat
Term  $\triangleq$  Nat

VARIABLE r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, r3amsgs, negMsgs,
       currentTerm, currentState, vote, leaderLog, log

serverVars  $\triangleq$  (currentTerm, currentState)
vars  $\triangleq$  (r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, r3amsgs, negMsgs, log,
           serverVars, leaderLog, vote)

```

主要的动作包括:

- *Timeout(i)*: 从节点或候选节点*i*因未收到来自主节点的消息而超时,则增大自身任期 ($\text{currentTerm}'[i] = \text{currentTerm}[i] + 1$),将状态转变为选举者 ($\text{currentState}'[i] = \text{Candidate}$).
- *RequestVote(i)*: 候选节点*i*向其它节点发送自身任期,发起选举请求.
- *HandleRequestVoteRequest(i)*: 节点*i*收到选举请求*m*,如果*m*携带的任期大于*i*的任期 ($m[1] > \text{currentTerm}[i]$),则*i*升级任期,接受该选举请求,并将自己的日志发送给候选节点.否则,*i*拒绝该选举.

请求.

$$\begin{aligned}
 \text{Timeout}(i) &\triangleq \\
 &\wedge \text{currentState}[i] \in \{\text{Follower}, \text{Candidate}\} \\
 &\wedge \text{currentTerm}' = [\text{currentTerm EXCEPT } ![i] = \text{currentTerm}[i] + 1] \\
 &\wedge \text{currentState}' = [\text{currentState EXCEPT } ![i] = \text{Candidate}] \\
 &\wedge \text{UNCHANGED } \langle r1amsgs, r1bmsgs, log, r2amsgs, r2bmsgs, r3amsgs, negMsgs, leaderLog, vote \rangle
 \end{aligned}$$

$$\begin{aligned}
 \text{RequestVote}(i) &\triangleq \\
 &\wedge \text{currentState}[i] = \text{Candidate} \\
 &\wedge r1amsgs' = r1amsgs \cup \{\langle \text{currentTerm}[i], i \rangle\} \\
 &\wedge \text{UNCHANGED } \langle \text{serverVars}, r1bmsgs, log, r2amsgs, r2bmsgs, r3amsgs, negMsgs, leaderLog, vote \rangle
 \end{aligned}$$

$$\begin{aligned}
 \text{HandleRequestVoteRequest}(i) &\triangleq \\
 &\wedge \exists m \in r1amsgs : \\
 &\quad \text{LET } j \triangleq m[2] \\
 &\quad \text{grant} \triangleq m[1] > \text{currentTerm}[i] \\
 &\quad \text{entries} \triangleq \{\langle n, \log[i][n] \rangle : n \in \text{Index}\} \\
 &\quad \text{IN} \\
 &\quad \vee \wedge \text{grant} \\
 &\quad \wedge \text{UpdateTerm}(i, m[1]) \\
 &\quad \wedge r1bmsgs' = r1bmsgs \cup \{\langle m[1], \text{entries}, i, j \rangle\} \\
 &\quad \wedge \text{UNCHANGED } negMsgs \\
 &\quad \vee \wedge \neg \text{grant} \\
 &\quad \wedge negMsgs' = negMsgs \cup \{\langle \text{currentTerm}[i], j \rangle\} \\
 &\quad \wedge \text{UNCHANGED } \langle \text{currentState}, \text{currentTerm}, r1bmsgs \rangle \\
 &\wedge \text{UNCHANGED } \langle log, r1amsgs, r2amsgs, r2bmsgs, r3amsgs, vote, leaderLog \rangle
 \end{aligned}$$

- *BecomeLeader(i)*: 收到多数派投票的候选节点*i*成为主节点并根据收到的日志恢复可能缺失的日志项. 对每个日志编号*n*, 考察它收到的所有编号为*n*的日志项, 选择其中任期最大的日志项, 将该日志项的任期修改为*i*自身的任期. 为了构建从 *ParallelRaft-SE* 到 *Multi-Paxos* 的精化关系, 恢复过程将修改 *leaderLog*, 而非直接修改 *log[i]*. 之后, *i*可以向自己发送 *RequestSync* 请求, 完成日志更新.
- *RequestSync(i)*: 主节点*i*向其它节点同步日志项.

$Merge(entries, term) \triangleq$
 LET $committed \triangleq \{e \in entries : e[3] = \text{TRUE}\}$
 $chosen \triangleq$
 CASE $committed = \{\} \rightarrow \text{CHOOSE } x \in entries :$
 $\forall y \in entries : x[1] \geq y[1]$
 $\square \quad committed \neq \{\} \rightarrow \text{CHOOSE } x \in committed : \text{TRUE}$
 $safe \triangleq chosen[2]$
 IN $\langle term, safe, chosen[3] \rangle$

$BecomeLeader(i) \triangleq$
 $\wedge currentState[i] = Candidate$
 $\wedge \exists Q \in Quorums :$
 LET $voteGranted \triangleq \{m \in r1bmsgs : m[4] = i \wedge m[3] \in Q$
 $\wedge m[1] = currentTerm[i]\}$
 $allLog \triangleq \text{UNION } \{m[2] : m \in voteGranted\}$
 $valid \triangleq \{e \in allLog : e[2][1] \neq -1\}$
 $end \triangleq \text{IF } valid = \{\} \text{ THEN } -1 \text{ ELSE } \text{Max}(\{e[1] : e \in valid\})$
 IN
 $\wedge \forall q \in Q : \exists m \in voteGranted : m[3] = q$
 $\wedge leaderLog' = [leaderLog \text{ EXCEPT } !currentTerm[i]] =$
 $[n \in Index \mapsto \text{IF } n \in 0..end \text{ THEN}$
 $Merge(\{l[2] : l \in \{t \in allLog : t[1] = n\}, currentTerm[i]\})$
 $\text{ELSE } \langle -1, Nil, \text{FALSE} \rangle]$
 $\wedge currentState' = [currentState \text{ EXCEPT } !i = Leader]$
 $\wedge \text{UNCHANGED } \langle currentTerm, r1amsgs, r2amsgs, r1bmsgs, r2bmsgs, r3amsgs,$
 $\quad negMsgs, log, vote \rangle$

$RequestSync(i) \triangleq$
 $\wedge currentState[i] = Leader$
 $\wedge \text{LET } sync \triangleq \{n \in Index : leaderLog[currentTerm[i]][n][1] \neq -1\}$
 IN
 $\exists n \in sync : r2amsgs' = r2amsgs \cup$
 $\{\langle currentTerm[i], n, leaderLog[currentTerm[i]][n], i \rangle\}$
 $\wedge \text{UNCHANGED } \langle serverVars, log, r1amsgs, r1bmsgs, r2bmsgs,$
 $\quad r3amsgs, negMsgs, leaderLog, vote \rangle$

-
- $HandleRequestSyncRequest(i)$: i 收到 $RequestSync$ 请求 m , 如果 m 携带的任期不小于 i 自身的任期 ($m[1] \geq currentTerm[i]$), 则 i 接受该请求, 升级自己的任期, 更新 $log[i]$ 与 $vote$, 并回复确认.
 - $CommitEntry(i)$: 收到多数派节点对某个日志项的确认后, 主节点 i 将该日志项标记为已提交 ($committed$).
 - $RequestCommit(i)$: 主节点将已提交的日志项发送给其它节点.

$$\begin{aligned}
HandleRequestSyncRequest(i) &\triangleq \\
&\wedge \exists m \in r2amsgs : \\
&\quad \text{LET } j \triangleq m[4] \\
&\quad grant \triangleq m[1] \geq currentTerm[i] \\
&\quad \text{IN} \\
&\quad \wedge \vee \wedge m[1] > currentTerm[i] \\
&\quad \wedge UpdateTerm(i, m[1]) \\
&\quad \vee \wedge m[1] \leq currentTerm[i] \\
&\quad \wedge \text{UNCHANGED } \langle currentTerm, currentState \rangle \\
&\quad \wedge \vee \wedge grant \\
&\quad \wedge log' = [log \text{ EXCEPT } ![i][m[2]] = m[3]] \\
&\quad \wedge vote' = [vote \text{ EXCEPT } ![i][m[2]][m[1]] = m[3][2]] \\
&\quad \wedge r2bmsgs' = r2bmsgs \cup \{ \langle m[1], m[2], i, j \rangle \} \\
&\quad \wedge \text{UNCHANGED } negMsgs \\
&\quad \vee \wedge \neg grant \\
&\quad \wedge negMsgs' = negMsgs \cup \{ \langle currentTerm[i], j \rangle \} \\
&\quad \wedge \text{UNCHANGED } \langle vote, r2bmsgs, log \rangle \\
&\wedge \text{UNCHANGED } \langle r1amsgs, r1bmsgs, r2amsgs, r3amsgs, leaderLog \rangle \\
CommitEntry(i) &\triangleq \\
&\wedge \exists index \in Index, Q \in Quorums : \\
&\quad \text{LET } syncSuccess \triangleq \{ m \in r2bmsgs : \\
&\quad \quad m[4] = i \wedge m[3] \in Q \\
&\quad \quad \wedge m[1] = currentTerm[i] \wedge m[2] = index \} \\
&\quad \text{IN} \\
&\quad \wedge currentState[i] = Leader \\
&\quad \wedge \forall q \in Q : \exists m \in syncSuccess : m[3] = q \\
&\quad \wedge leaderLog' = [leaderLog \text{ EXCEPT } !(currentTerm[i])[index][3] = \text{TRUE}] \\
&\wedge \text{UNCHANGED } \langle serverVars, log, r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, \\
&\quad r3amsgs, negMsgs, vote \rangle \\
RequestCommit(i) &\triangleq \\
&\wedge currentState[i] = Leader \\
&\wedge \text{LET } committed \triangleq \{ n \in Index : leaderLog[currentTerm[i]][n][3] = \text{TRUE} \} \text{ IN} \\
&\quad \exists n \in committed : r3amsgs' = r3amsgs \cup \{ \langle currentTerm[i], n, i \rangle \} \\
&\wedge \text{UNCHANGED } \langle serverVars, log, r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, \\
&\quad negMsgs, leaderLog, vote \rangle
\end{aligned}$$

- *HandleRequestCommit(i)*: 节点*i*收到来自主节点的*ReuestCommit*请求*m*,如果*m*携带的任期不小于*i*自身的任期,则*i*将相应的日志项标记为已提交.
- *ClientRequest(i)*: 收到用户的命令*v*后,主节点*i*将*v*作为新的日志项添加到日志中.

$$\begin{aligned}
 & HandleRequestCommitRequest(i) \triangleq \\
 & \wedge \exists m \in r3amsgs : \\
 & \quad \text{LET } grant \triangleq currentTerm[i] \leq m[1] \\
 & \quad j \triangleq m[3] \\
 & \quad \text{IN} \\
 & \quad \wedge \vee \wedge m[1] > currentTerm[i] \\
 & \quad \wedge UpdateTerm(i, m[1]) \\
 & \quad \vee \wedge m[1] \leq currentTerm[i] \\
 & \quad \wedge \text{UNCHANGED } \langle currentTerm, currentState \rangle \\
 & \quad \wedge \vee \wedge grant \\
 & \quad \wedge log[i][m[2]][1] = m[1] \\
 & \quad \wedge log' = [log \text{ EXCEPT } ![i][m[2]][3] = \text{TRUE}] \\
 & \quad \wedge \text{UNCHANGED } negMsgs \\
 & \quad \vee \wedge \neg grant \\
 & \quad \wedge negMsgs' = negMsgs \cup \{\langle currentTerm[i], j \rangle\} \\
 & \quad \wedge \text{UNCHANGED } log \\
 & \wedge \text{UNCHANGED } \langle serverVars, r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, \\
 & \quad r3amsgs, leaderLog, vote \rangle
 \end{aligned}$$

$$\begin{aligned}
 & ClientRequest(i) \triangleq \\
 & \text{LET } ind \triangleq \{b \in Index : leaderLog[currentTerm[i]][b][1] \neq -1\} \\
 & nextIndex \triangleq \text{IF } ind = \{\} \\
 & \quad \text{THEN } 0 \\
 & \quad \text{ELSE } \text{Max}(ind) + 1 \\
 & \text{IN} \\
 & \wedge currentState[i] = Leader \\
 & \wedge \exists v \in Value : leaderLog' = [leaderLog \text{ EXCEPT } !(currentTerm[i])[nextIndex] = \\
 & \quad \langle currentTerm[i], v, \text{FALSE} \rangle] \\
 & \wedge \text{UNCHANGED } \langle serverVars, log, r1amsgs, r1bmsgs, r2amsgs, r2bmsgs, r3amsgs, \\
 & \quad negMsgs, vote \rangle
 \end{aligned}$$

Next 定义了次态关系。*Spec* 定义了完整的行为规约。

$$\begin{aligned}
 Next \triangleq & \vee \exists i \in Server : \text{Timeout}(i) \\
 & \vee \exists i \in Server : \text{RequestVote}(i) \\
 & \vee \exists i \in Server : \text{HandleRequestVoteRequest}(i) \\
 & \vee \exists i \in Server : \text{BecomeLeader}(i) \\
 & \vee \exists i \in Server : \text{CommitEntry}(i) \\
 & \vee \exists i \in Server : \text{ClientRequest}(i) \\
 & \vee \exists i, j \in Server : \text{RequestCommit}(i) \\
 & \vee \exists i \in Server : \text{HandleRequestCommitRequest}(i) \\
 & \vee \exists i, j \in Server : \text{RequestSync}(i) \\
 & \vee \exists i \in Server : \text{HandleRequestSyncRequest}(i)
 \end{aligned}$$

$$Spec \triangleq Init \wedge \square[Next]_{vars}$$

2.2 精化ParallelRaft-SE到Multi-Paxos

ParallelRaft-SE 支持乱序提交、顺序执行用户命令.这与 Multi-Paxos 相同.此外,ParallelRaft-SE 的日志恢复阶段本质上是使用 Paxos 对可能缺失的日志项重确认,而 Multi-Paxos 也通过 Paxos 发起提议.实际上,我们可以建立从 ParallelRaft-SE 到 Multi-Paxos 的精化关系,从而也证明了 ParallelRaft-SE 的正确性.这种精化关系基于 ParallelRaft-SE 与 Multi-Paxos 的以下相似之处:

- *RequestVote*对应于*Phase1a*.ParallelRaft-SE 中的任期对应于 Multi-Paxos 中的提议编号.
- *HandleRequestVote*对应于*Phase1b*.二者都需要通过比较任期/提议编号来决定是否同意选举/接受提议请求,且都需要在回复中包含自己的日志.
- *BecomeLeader*对应于*Merge*.收到多数派回复的主节点/提议者运行 Paxos 恢复可能缺失的日志项.
- *RequestSync*对应于*Phase2a*. 在 ParallelRaft-SE 中,主节点完成日志恢复后,向从节点同步日志项.在 Multi-Paxos 中,提议者通过运行 Paxos 选出提议值后,向接受者同步.
- *HandleRequestSync*对应于*Vote*.在 ParallelRaft-SE 中,从节点收到同步请求后更新自己的日志.在 Multi-Paxos 中,接受者收到提议者的提议后,接受提议并记录在本地.
- *ClientRequest*对应于*Propose*. 在 ParallelRaft-SE 中,主节点收到用户命令,将其作为日志项添加到日志的末尾.在 Multi-Paxos 中,若对于某个编号,提议者没有收到提议值,可以提出任意合法的值.因此可以向日志末尾追加合法的日志项.

ParallelRaft-SE 规约给出了 *ParallelRaft-SE* 和 *Multi-Paxos* 在常量、变量之间的精化映射:

$$\begin{aligned}
 & \vdash \\
 & \begin{aligned}
 & \text{Acceptors} \triangleq \text{Server} \\
 & \text{Ballots} \triangleq \text{Term} \\
 & \text{Instances} \triangleq \text{Index} \\
 & \text{ballot} \triangleq \text{currentTerm} \\
 & \text{leaderVote} \triangleq [i \in \text{Ballots} \mapsto [j \in \text{Index} \mapsto \langle \text{leaderLog}[i][j][1], \text{leaderLog}[i][j][2] \rangle]] \\
 & \text{1amsgs} \triangleq \{\langle m[1] \rangle : m \in r1amsgs\} \\
 & \text{1bmsgs} \triangleq \{\langle \langle e[1], \langle e[2][1], e[2][2] \rangle \rangle : e \in m[2] \}, m[3] \rangle : m \in r1bmsgs\} \\
 & \text{2amsgs} \triangleq \{\langle m[1], m[2], \langle m[3][1], m[3][2] \rangle \rangle : m \in r2amsgs\} \\
 & \text{Spec} \triangleq \text{Init} \wedge \square[\text{Next}]_{vars} \\
 & \text{MP} \triangleq \text{INSTANCE MultiPaxos} \\
 & \text{THEOREM Refinement} \triangleq \text{Spec} \Rightarrow \text{MP}! \text{Spec}
 \end{aligned}
 \end{aligned}$$

3 乱序执行模型与"幽灵日志"问题

相对于 Raft 而言,ParallelRaft-SE 支持乱序提交.但是,它仍然要求顺序执行用户命令,不适合于高并发系统.在 PolarFS 文件系统应用场景中,ParallelRaft 需要支持乱序执行,也就是允许状态机先执行编号较大的已提交日志项,而不必等待编号较小的日志项被提交或被执行^[13].为了在乱序执行情况下,仍能满足状态一致性,需要保证被乱序执行的命令是无冲突的.因此,本节先介绍 ParallelRaft 所采用的乱序执行模型.在分析 ParallelRaft 协议的正确性时,我们发现,文献^[13] 中关于 ParallelRaft 的描述忽略了可能会违反状态一致性的"幽灵日志"问题.本节将分析该问题给乱序执行机制带来的挑战.

3.1 乱序执行的模型

ParallelRaft 的乱序执行模型给出了用户命令之间的冲突判断规则^[13].在 PolarFS 文件系统应用场景中,每

个用户命令都包含该命令所访问数据的逻辑区块地址(Logic Block Address; LBA).LBA 不重叠的命令不存在冲突,可以乱序执行;相反,LBA 有重叠的命令需要按日志编号顺序执行.

在该模型下,每个命令在被执行前,首先检查该命令是否与编号更小的命令存在冲突.由于日志项是乱序接受(可能尚未被提交)的,日志中可能存在"空洞".为了在存在空洞的情况下,仍能判断当前命令是否可被执行,ParallelRaft 要求每个日志项记录它之前的 K (K 是待定参数)个日志项的 LBA,称为"向后看缓冲区"(Look Behind Buffer;LBF).因此,只要日志中不存在长度超过 K 的"空洞",就可以判断任意两个日志项/命令之间是否存在冲突;否则,"空洞"之后的日志项都需要等待.

3.2 "幽灵日志"问题

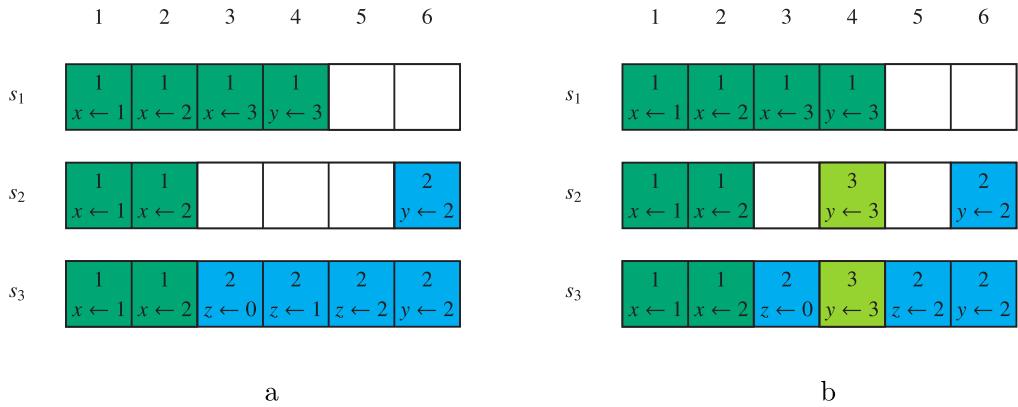


Figure 2 The "Ghost Log Entries" Phenomenon Violates Consistency

图 2 "幽灵日志"破坏状态一致性

经过日志恢复后,ParallelRaft 中主节点的日志不存在"空洞",因此它可以为每个日志项计算 LBF.之后,主节点可以将 LBF 随日志项一起发送给从节点.正确的冲突判断要求每个日志项的 LBF 都准确地记录了它之前 K 个日志项的 LBA. 然而,"幽灵日志"现象可能会导致主节点计算出的 LBF 与实际不符.图 2 描述了"幽灵日志"现象,它包含三个阶段(图中每个方格表示一个日志项,它记录了该日志项的任期以及所携带的用户命令.日志项从 1 开始编号):

- 如图 2 a,第一阶段中, s_1 是主节点. s_1 向日志中添加了编号为 1,2,3,4 的 4 个日志项.其中编号为 1 和 2 的日志项被 s_2 和 s_3 接受,这两个日志项被提交.编号为 3 和 4 的日志项还没有被 s_2 和 s_3 接受. s_1 失效.
- 第二阶段中, s_3 成为主节点. s_3 的主节点在恢复过程中没有收到 s_1 未提交的项(注意, s_3 只需要从多数派节点收集日志项).因此完成恢复完成后, s_3 的日志中不包含 s_1 未提交的日志项.如图 2 a,之后 s_3 向日志中添加了编号为 3,4,5,6 的新的日志项,并将编号为 6 的日志项发送给了 s_2 ,该项被提交.由于这一项是对 y 的修改($y \leftarrow 2$),与之前的日志项无冲突.因此 s_3 执行 $y \leftarrow 2$.(乱序执行).之后 s_3 也失效,系统选出新的主节点 s_2 .
- 如图 2 b,第三阶段中, s_3 在恢复过程中的收到了 s_1 (此时 s_1 恢复正常)未提交的日志项(编号为 4),并将它发送给 s_3 .这一项被提交后, s_3 需要执行.它也是对变量 y 的修改($y \leftarrow 3$),根据乱序执行的规则,应该先执行.但是 s_3 已经执行 $y \leftarrow 2$,因此导致了执行顺序的不一致.

s_1 在第一阶段未提交的日志项不在第二阶段的主节点 s_3 的日志中,但是之后又出现在了第三阶段的主节点 s_2 的日志中.我们称这种现象为"幽灵日志".这使得 s_3 计算出的 LBF 与实际不符,可能会导致错误的冲突判断,进

而影响 ParallelRaft 在乱序执行下的正确性.避免出现"幽灵日志"的关键在于,在日志恢复的过程中,第二阶段的主节点 s_3 需要明确 s_1 中未提交的日志项的状态:要么将它们重新提交,要么将它们删除并保证之后不会被新的主节点提交."幽灵日志"现象在顺序执行(如 Multi-Paxos 与 Raft)下也可能出现,但不会影响协议的正确性^[12].在顺序执行下,"幽灵日志"问题很容易解决.接下来我们将简要论述 Multi-Paxos 与 Raft 解决"幽灵日志"问题的方案并论证这些方案无法解决乱序执行下的"幽灵日志"问题.

3.3 Multi-Paxos 的解决方案

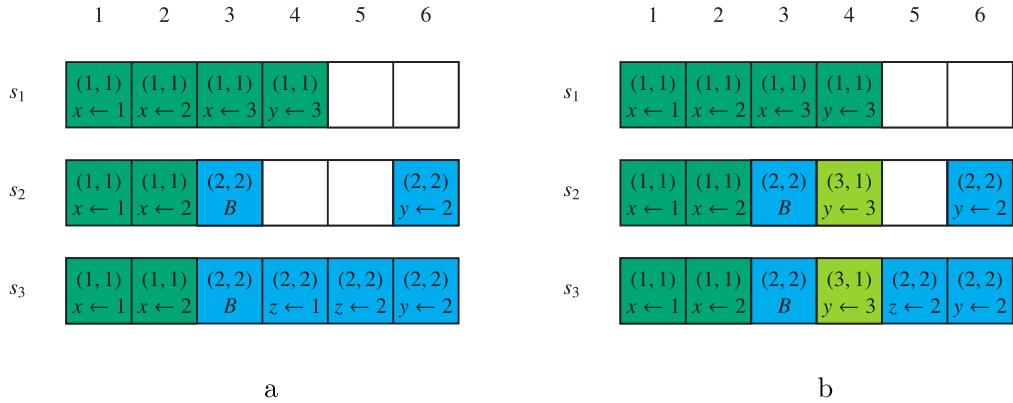


Figure 3 Passive Way to Address the "Ghost Log Entries" Phenomenon

图 3 "幽灵日志"被动式解决方案

Multi-Paxos 采用了被动式的解决方案.提议者在添加日志项时保存它的生成时间(即,生成该日志项的提议者的提议编号).节点可以通过日志项的生成时间检测出"幽灵日志",并忽略它们.为此,新的提议者完成恢复过程后,向日志中写入一条特殊的空操作日志项,称为栅栏(barrier).只有等栅栏日志项被提交后,提议者才能向日志中添加新的日志项.因此,栅栏日志项标记了日志恢复的终点以及新的提议者添加日志项的起点.如果之后某个日志项的生成时间小于栅栏日志项的生成时间,则可断定它为"幽灵日志".

如图 3 a 所示,每个日志项中的二元组分别记录了该日志项的提议编号(后续过程中可能更改)与生成时间. s_1 是提议编号 1 的提议者,它向日志中添加了日志编号为 1,2,3,4 的日志项,它们的提议编号与生成时间均为 1.其中日志编号为 3 和 4 的日志项没有被提交.之后, s_1 失效. s_3 成为提议编号 2 的提议者. s_3 收到 s_2 的日志项,完成恢复过程后,向日志中写入一条栅栏日志项(日志编号为 3,在图中用 *B* 标记),其提议编号与生成时间为 2. s_3 首先将栅栏日志项发送给 s_2 ,并将它提交.之后, s_3 响应用户请求,并向日志中添加了日志编号为 4,5,6 的日志项,其中日志编号为 6 的日志项被 s_2 接受(因此被提交). s_3 执行该日志项后失效.

s_2 成为提议编号 3 的提议者.如图 3 b,在恢复过程中, s_2 从 s_1 的日志中找到了日志编号为 4($y \leftarrow 3$)的日志项并将它添加到日志中.这一项的提议编号为 3,但是它是 s_1 写入日志的,生成时间为 1.而 s_2 的日志中编号为 3 的栅栏日志项的生成时间为 2.由此, s_2 可以判断日志编号为 4 的日志项是"幽灵日志". s_2 从日志中删除该项,从而解决了"幽灵日志"问题.

3.4 Raft的解决方案

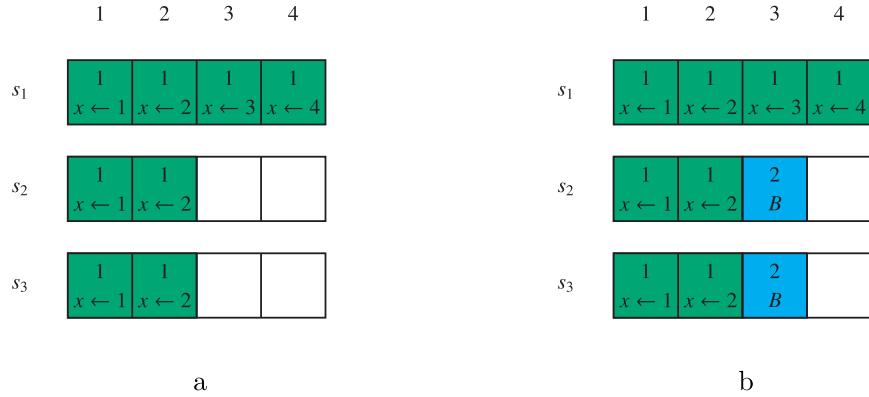


Figure 4 Proactive Way to Address the "Ghost Log Entries" Phenomenon

图 4 “幽灵日志”主动式解决方案

Raft 采用了主动式的解决方案^[12] :利用选主机制中的限制条件,使得拥有“幽灵日志”的节点无法成为新的主节点.为此,Raft 中选出的新主节点向日志中添加一条栅栏日志项.只有等栅栏日志项被提交后,新的主节点才能响应用户的请求.根据 Raft 的选举规则,此时,没有栅栏日志项的节点将无法成为主节点(日志较旧,无法获得多数派投票).另一方面,在日志项同步阶段,接收到栅栏日志项的节点将会删除栅栏日志项后的所有日志项,故不存在“幽灵日志”.

如图 4 a, s_1 是任期 1 的主节点,它向日志中写入了编号为 1,2,3,4 的四个日志项.其中编号为 3 和 4 的日志项还没有提交,主节点变更为 s_2 .如图 4 b, s_2 成为任期 2 的主节点后,先向日志中写入一条栅栏日志项并提交,之后 s_2 才能响应用户请求.即使 s_2 之后失效,根据 Raft 的选主规则, s_1 若想成为主节点,需要先删除多余的日志项(编号为 3 和 4),并添加栅栏日志项.因此,编号为 3 和 4 的两个日志项不会出现在新的主节点的日志中.

3.5 乱序执行下的挑战

ParallelRaft-SE 是 Multi-Paxos 的一种精化,也会出现“幽灵日志”现象.但是,直接应用 Multi-Paxos 的解决方案并不能保证 ParallelRaft-SE(或 ParallelRaft)在乱序执行下的正确性.换句话说,在乱序执行下,“幽灵日志”现象并非导致状态不一致的必要条件.比如,主节点 i 的日志中包含未提交的日志项.当主节点发生变更后, i 乱序地收到了新主节点的日志项. i 并不知道日志中未提交的日志项已经过时,从而误以为不存在冲突,导致不一致的行为.

另一方面,Raft 利用它的串行性在选举过程中消除了“幽灵日志”.然而,ParallelRaft-SE 支持乱序提交,日志中可能存在“空洞”.这使得 ParallelRaft-SE 也不能直接应用 Raft 的解决方案.

4 ParallelRaft-CE 协议与规约

为了解决乱序执行下的“幽灵日志”问题,我们基于 ParallelRaft-SE 提出了 ParallelRaft-CE 协议.ParallelRaft-CE 通过限制 ParallelRaft-SE 在乱序提交阶段的并行度,避免了“幽灵日志”问题.ParallelRaft-CE 主要包括三个部分:日志同步机制、选主机制及其日志恢复机制.

我们先给出 ParallelRaft-CE 的几条关键性质,其中“同步号”与“主节点候选者”的概念将在接下来的协议描述中介绍.下一节将简要证明这些关键性质以及 ParallelRaft-CE 的正确性.

1. 选举的安全性:每个任期最多选出一个主节点.
2. 节点的任期与同步号单调递增.
3. 设主节点候选者的任期为 t ,同步号为 s ,则 $s < t$ 且系统内不存在任期大于 s 的日志项.
4. 主节点完全性:若主节点的任期为 t ,那么它的日志中包括所有任期小于 t 且被提交的日志项.
5. 一致性:对任意两个节点,如果它们的同步号都大于某个任期值 t ,则它们的日志中包含相同的任期为 t 的日志项.

4.1 日志同步机制

在 ParallelRaft-CE 中,日志项的发送与接受是部分乱序的:任期相同的日子项可以并发地发送与接受,任期不同的日志项则需按序发送与接受.为此,每个节点维护一个同步号(sync),表示目前节点仅接受任期等于同步号的日志项.当从节点收到主节点的日志项时,会检查该日志项的任期.如果日志项的任期与从节点的同步号相同,则从节点接受该日志项,并回复确认消息;否则,从节点拒绝该日志项,并将自己的同步号发送给主节点.主节点额外维护每个从节点的同步号.当收到从节点的同步号时,主节点升级自己记录的从节点同步号.当从节点对当前任期的每一条日志项都完成确认后,主节点通知从节点将同步号设置为下一任期.这种日志同步机制限制了提交阶段的并行度.

4.2 选主机制

ParallelRaft-CE 的选主机制与 Raft 相似.不同的是,ParallelRaft-CE 通过同步号判断日志的新旧:节点的同步号越大,它的日志就越新.候选节点发起选举时,将自己的同步号也发送给其它节点.只有当接收节点的同步号不大于候选节点的同步号时,才同意该选举请求.

4.3 日志恢复机制

在 ParallelRaft-CE 中,任期相同的日子项可以并发发送与接受.因此,日志中可能存在"空洞",新的主节点需要对上一任期的日子项进行恢复.恢复的目标是使得主节点的同步号升级为它的任期,此时系统便已对所有任期小于主节点任期的日子项达成了共识.我们将尚未完成日志恢复的主节点称为主节点候选者(Leader Candidate).恢复过程完成后,所有任期小于主节点任期的日子项的状态是确定的:已经被提交(或执行)的日子项仍在新主节点的日志中;之前未提交的日子项要么被提交,要么被丢弃且之后不会再被提交.

假设主节点候选者 l 的任期为 t ,同步号为 s (根据性质 3, $s < t$).根据性质 5 系统已对任期小于 s 的日子项达成共识.再根据性质 3,系统中不存在任期超过 s 的日子项.因此 l 仅需要恢复任期等于 s 的日子项.恢复日志项的方法与 ParallelRaft-SE(或 Multi-Paxos)相同, l 需要从多数派节点收集日志项,并从中选择最"新"的日子项.由于日志项的任期是固定的,ParallelRaft-CE 为每个日志项增加一个新的变量:日期(date).日期相当于 Paxos 中日志项的提议编号.日志编号相同,日期越大的日志项就越"新".

l 需要分三阶段将恢复得到的日子项同步给从节点.设 l 记录某从节点的同步号为 n .如果 $n > s$,则直接进入第三阶段;如果 $n = s$,则直接进入第二阶段.否则,先进入第一阶段.

第一阶段($n < s$): l 先向该从节点发送任期为 n 的日子项.当这些项被从节点全部确认后, l 通知从节点升级到满足以下条件的下一个同步号 k :

1. $k > n$;
2. l 的日子项中含有任期为 k 的日子项;
3. k 是满足条件 1,2 的最小值.

收到升级同步号请求后,从节点删除日志中所有未确认的任期为 n 的日子项,并将同步号升级为 k .之后,主节点候选者从节点发送任期为 k 的日子项.持续以上过程,直到从节点的同步号升级为 s .进入第二阶段.

第二阶段($n = s$): l 持续向从节点同步日子项,直到任期为 s 的所有日子项都被相同的多数派提交,记这个多数派为 Q .进入第三阶段.

第三阶段($n \geq s$): l 将它的同步号升级为自身的任期 t (根据性质 3,这是安全的).然后, l 通知 Q 中的节点将自

身同步号升级为 t .收到某多数派节点(均来自 Q)的确认消息后,恢复过程结束, i 正式成为主节点.

4.4 ParallelRaft-CE的TLA+规约

ParallelRaft-CE 使用常量 *LeaderCandidate* 表示主节点候选者角色.其中的变量包括:

- *messages*: 节点发送的消息的集合.
- *currentTerm[i]*: 节点*i*记录的最大的任期.
- *currentState[i]*: 节点*i*的状态,为*Leader*,*LeaderCandidate*,*Follower*,*Candidate*之一.
- *votedFor*: 每个节点一个任期内只能为一个候选节点投票.*votedFor[i]*表示*i*在本任期(*currentTerm[i]*)内投票的候选节点.若*i*没有为任何节点投票,则*votedFor[i] = Nil*.
- *sync[i]*: 节点*i*的同步号.
- *end[i][t]*: 节点*i*记录的可接受的任期为 t 的日志项的最大编号.*ParallelRaft-CE* 使用 Paxos 的投票机制确定每一个任期的最大编号.
- *log[i]*: 节点*i*的日志.*log[i][k]*为*i*的日志中编号为 k 的日志项.每个日志项是 $\langle t, d, v, b \rangle$ 的四元组.其中 t 是日志项的任期,不可修改; d 是日期,用于判断日志项的新旧; v 是提议值; b 是布尔值,当且仅当日志项被提交时, b 为真.
- *syncTrack*: 状态为*Leader*或*LeaderCandidate*的节点用来记录其它节点的同步号.*syncTrack[i][a]*为节点*i*记录的节点的同步号.
- *elections, halfElections*为历史变量,记录了选举信息.

MODULE *ParallelRaft - CE*

CONSTANTS *Server*, *Follower*, *Candidate*, *Leader*, *LeaderCandidate*, *Nil*, *Value*,
RequestVoteRequest, *RequestVoteResponse*,
RequestCommitRequest, *RequestCommitResponse*,
RequestSyncRequest, *RequestSyncResponse*,
UpdateSyncRequest, *UpdateSyncResponse*

Quorums $\triangleq \{i \in \text{SUBSET } (\text{Server}) : \text{Cardinality}(i) * 2 > \text{Cardinality}(\text{Server})\}$

VARIABLE *messages*, *currentTerm*, *currentState*, *votedFor*, *sync*, *end*, *log*, *syncTrack*

serverVars $\triangleq \langle \text{currentTerm}, \text{currentState}, \text{votedFor}, \text{sync}, \text{end} \rangle$

VARIABLE *halfElections*, *elections*
electionVars $\triangleq \langle \text{halfElections}, \text{elections} \rangle$

VARIABLE *allLogs*, *allEntries*, *allSynced*

vars $\triangleq \langle \text{messages}, \text{allLogs}, \text{allEntries}, \text{log}, \text{serverVars},$
syncTrack, *allSynced*, *electionVars* \rangle

ParallelRaft-CE 的主要动作包括:

- *UpdateTerm(i)*: s 收到消息 m , 若 $m.mterm > \text{currentTerm}[i]$, 那么 i 升级任期 ($\text{currentTerm}'[i] = m.mterm$), 转变为从节点 ($\text{currrentState}'[i] = \text{Follower}$), 并将投票的记录置空 ($\text{voteFor}'[i] = \text{Nil}$), 因为 i 在新的任期还没有为任何节点投票.在 *ParallelRaft-CE* 中, 节点收到任何请求都需要检查任期.当消息的任期大于节点的任期时, 节点执行 *UpdateTerm*, 之后再响应请求; 当消息的任期与节点任期相同时, 节点直接处理请求; 当消息的任期小于节点任期时, 节点拒绝任何请求, 并将自己的任期回复给发送者.
- *RequestVote(i)*: i 发起选举.与 *ParallelRaft-SE(Multi-Paxos)* 不同的是, i 需要将自己的同步号 ($\text{sync}[i]$) 发送给其它节点.

- *HandleRequestVoteRequest(i)*: 收到选举请求 m 的节点*i*比较任期($currentTerm[i]$ 和 $m.mterm$)与同步号($sync[i]$ 和 $m.msSync$)。当 $m.mterm = currentTerm[i]$ 且 $sync[i] \geq m.msSync$ 且*i*在本任期内还没有同意其它节点的选举请求,*i*同意选举请求,并将自己的日志中任期为 $m.msSync$ 的项以及*i*可以接受的任期为 $m.msSync$ 的日志项的最大编号($end[i][m.msSync]$)发送给选举发起者。

$UpdateTerm(i) \triangleq$
 $\wedge \exists m \in messages :$
 $\wedge m.mterm > currentTerm[i]$
 $\wedge \vee m.mdest = i$
 $\quad \vee m.mdest = Nil$
 $\wedge currentTerm' = [currentTerm \text{ EXCEPT } ![i] = m.mterm]$
 $\wedge currentState' = [currentState \text{ EXCEPT } ![i] = Follower]$
 $\wedge votedFor' = [votedFor \text{ EXCEPT } ![i] = Nil]$
 $\wedge \text{UNCHANGED } \langle messages, sync, log, syncTrack, electionVars, allSynced, end \rangle$

$RequestVote(i) \triangleq$
 $\wedge currentState[i] = Candidate$
 $\wedge Send([mtype \mapsto RequestVoteRequest,$
 $\quad mterm \mapsto currentTerm[i],$
 $\quad msync \mapsto sync[i],$
 $\quad msources \mapsto i,$
 $\quad mdest \mapsto Nil])$
 $\wedge \text{UNCHANGED } \langle serverVars, syncTrack, log, electionVars, allSynced \rangle$

$HandleRequestVoteRequest(i) \triangleq$
 $\wedge \exists m \in messages :$
 $\text{LET } j \triangleq m.msources$
 $\quad syncOK \triangleq \wedge m.msSync \geq sync[i]$
 $\quad grant \triangleq \wedge syncOK$
 $\quad \wedge votedFor[i] \in \{Nil, j\}$
 $\quad \wedge currentTerm[i] = m.mterm$
 IN
 $\quad \wedge m.mterm \leq currentTerm[i]$
 $\quad \wedge m.mtype = RequestVoteRequest$
 $\quad \wedge \vee grant \wedge votedFor' = [votedFor \text{ EXCEPT } ![i] = j]$
 $\quad \quad \vee \neg grant \wedge \text{UNCHANGED } votedFor$
 $\quad \wedge Send([mtype \mapsto RequestVoteResponse,$
 $\quad \quad mterm \mapsto currentTerm[i],$
 $\quad \quad mvoteGranted \mapsto grant,$
 $\quad \quad mlog \mapsto \text{LET } C \triangleq \{n \in Index : log[i][n].term = sync[i]\}$
 $\quad \quad \quad \text{IN } \{\langle n, log[i][n] \rangle : n \in C\},$
 $\quad \quad mend \mapsto end[i][m.msSync],$
 $\quad \quad msources \mapsto i,$
 $\quad \quad mdest \mapsto j])$
 $\quad \wedge \text{UNCHANGED } \langle currentTerm, currentState, sync, log, syncTrack, electionVars, allSynced, end \rangle$

- *BecomeLeaderCandidate(i)*: 若候选节点*i*从一个多数派(*voteGranted*)的节点收到同意选举的消息,*i*转变为主节点候选者。*i*要恢复任期等于自己同步号($sync[i]$)的日志项。因此*i*首先确定能接受的任期为 $sync[i]$ 的日志项的最大编号,将它保存在 $end[i][sync[i]]$ 中。方法为在收到的消息的 $m.mend$ 中选择最新

的.之后对每个可接受的编号, i 在收到的对应编号的日志项中选择最新的,并写入日志.对每个节点 p,i 初始化 $syncTrack[i][p] = sync[i]$.

```

Merge(entries, term, date)  $\triangleq$  IF entries = {} THEN [term  $\mapsto$  term,
                                                       date  $\mapsto$  date,
                                                       value  $\mapsto$  Nil,
                                                       committed  $\mapsto$  FALSE]
                                         ELSE
                                         LET
                                         committed  $\triangleq$  {e  $\in$  entries : e.committed = TRUE}
                                         chosen  $\triangleq$ 
                                         CASE committed = {}  $\rightarrow$  CHOOSE x  $\in$  entries :
                                          $\forall y \in entries : x.date \geq y.date$ 
                                          $\square$  committed  $\neq$  {}  $\rightarrow$  CHOOSE x  $\in$  committed : TRUE
                                         IN
                                         [term  $\mapsto$  chosen.term,
                                         date  $\mapsto$  date,
                                         value  $\mapsto$  chosen.value,
                                         committed  $\mapsto$  chosen.committed]

BecomeLeaderCandidate(i)  $\triangleq$ 
 $\wedge currentState[i] = Candidate$ 
 $\wedge \exists P \in Quorums :$ 
  LET voteGranted  $\triangleq$  {m  $\in$  messages :  $\wedge m.mtype = RequestVoteResponse$ 
                            $\wedge m.mdest = i$ 
                            $\wedge m.msource \in P$ 
                            $\wedge m.mterm = currentTerm[i]$ 
                            $\wedge m.mvoteGranted = TRUE$ }
  allLog  $\triangleq$  UNION {m.mlog : m  $\in$  voteGranted}
  endLine  $\triangleq$  LET allPoint  $\triangleq$  {m.mend : m  $\in$  voteResponded}
  IN e  $\triangleq$  CHOOSE e1  $\in$  allPoint :
   $(\forall e2 \in allPoint : e1[1] \geq e2[1])$ 
  toRecover  $\triangleq$  {n  $\in$  0 .. endLine : log[i][n].committed = FALSE}
  toSync  $\triangleq$  {⟨n, Merge({l[2] : l  $\in$  {t  $\in$  allLog : t[1] = n}}), sync[i], currentTerm[i]⟩
                : n  $\in$  toRecover}
  IN
   $\wedge \forall p \in P : \exists m \in voteGranted : m.msource = p$ 
   $\wedge log' = [log \text{ EXCEPT } ![i] = [n \in Index \mapsto \text{IF } n \in toRecover \text{ THEN}$ 
  (CHOOSE e  $\in$  toSync : e[1] = n)[2]
  ELSE log[i][n]]]
   $\wedge end' = [end \text{ EXCEPT } ![i][sync[i]] = \langle currentTerm[i], end \rangle]$ 
 $\wedge currentState' = [currentState \text{ EXCEPT } ![i] = LeaderCandidate]$ 
 $\wedge syncTrack' = [syncTrack \text{ EXCEPT } ![i] = [j \in Server \mapsto sync[i]]]$ 
 $\wedge \text{UNCHANGED } \langle messages, currentTerm, votedFor, sync, elections, allSynced \rangle$ 

```

- **RequestSync(i):**当*i*为主节点或主节点候选者时,对其它每个节点*p*,*i*向*p*发送任期为 $syncTrack[i][p]$ 的日志项.这些日志项可以乱序地发送和接受.
- **HandleRequestSyncRequest(i):**节点*i*收到同步请求*m*.当*m.msSync < sync[i]*或*m.msSync > sync[i]*时,*i*拒绝请求并回复*RequestSyncResponse*消息,将*sync[i]*告知发送者.当*m.msSync = sync[i]*时,*i*同意请求,回复确认并根据*m.mend*(可以接受的任期为*m.msSync*的日志项的最大编号)和*m.mentries*(任期为*m.msSync*的日志项)修改*end[i]*和*log[i]*.对日志的修改为:删除日志中所有编号大于*m.mend*的日志项,并将编号不超过*m.mend*的日志项替换为*m.mentries*中相应的项.

\vdash

$$\begin{aligned}
 & RequestSync(i) \triangleq \\
 & \quad \wedge currentState[i] \in \{LeaderCandidate, Leader\} \\
 & \quad \wedge \exists s \in 0..sync[i] : \\
 & \quad \quad \text{LET } start \triangleq \min(\{n \in Index : log[i][n].term = s\}) \\
 & \quad \quad \text{end} \triangleq \max(\{n \in Index : log[i][n].term = s\}) \\
 & \quad \text{IN} \\
 & \quad \quad \wedge Send([mtype \mapsto RequestSyncRequest, \\
 & \quad \quad \quad mterm \mapsto currentTerm[i], msync \mapsto s, \\
 & \quad \quad \quad mstart \mapsto start, mend \mapsto end, \\
 & \quad \quad \quad mentries \mapsto \text{IF } start = -1 \text{ THEN Nil ELSE} \\
 & \quad \quad \quad \quad [n \in start..end \mapsto log[i][n]], \\
 & \quad \quad \quad \quad msource \mapsto i, mdest \mapsto Nil]) \\
 & \quad \wedge \text{UNCHANGED } \langle serverVars, logVars, electionVars, syncTrack, allSynced \rangle \\
 & HandleRequestSyncRequest(i) \triangleq \\
 & \quad \wedge \exists m \in messages : \\
 & \quad \quad \text{LET } j \triangleq m.msource \\
 & \quad \quad grant \triangleq \wedge m.mterm = currentTerm[i] \\
 & \quad \quad \wedge m.msSync = sync[i] \\
 & \quad \text{IN} \\
 & \quad \quad \wedge m.mtype = RequestSyncRequest \\
 & \quad \wedge m.mterm \leq currentTerm[i] \\
 & \quad \wedge j \neq i \\
 & \quad \wedge \vee \wedge grant \\
 & \quad \quad \wedge log' = [log \text{ EXCEPT } ![i] = [n \in Index \mapsto \\
 & \quad \quad \quad \text{IF } n < m.mstart \text{ THEN } log[i][n] \\
 & \quad \quad \quad \text{ELSE IF } n \in m.mstart..m.mend \\
 & \quad \quad \quad \quad \text{THEN } m.mentries[n] \\
 & \quad \quad \quad \text{ELSE } [term \mapsto -1, date \mapsto -1, \\
 & \quad \quad \quad \quad value \mapsto Nil, committed \mapsto \text{FALSE}]]] \\
 & \quad \wedge endPoint' = [endPoint \text{ EXCEPT } ![i][sync[i]] = \langle currentTerm[i], m.mend \rangle] \\
 & \quad \vee \wedge \neg grant \\
 & \quad \wedge \text{UNCHANGED } \langle log, endPoint \rangle \\
 & \quad \wedge Send([mtype \mapsto RequestSyncResponse, mterm \mapsto currentTerm[i], \\
 & \quad \quad msyncGranted \mapsto grant, msync \mapsto sync[i], \\
 & \quad \quad mstart \mapsto m.mstart, mend \mapsto m.mend, \\
 & \quad \quad msource \mapsto i, mdest \mapsto j]) \\
 & \quad \wedge \text{UNCHANGED } \langle currentTerm, currentState, sync, votedFor,
 \end{aligned}$$

\vdash

- $HandleRequestSyncResponse(i)$: i 为主节点或主节点候选者时,收到节点 p 的 $RequestSyncResponse$ 消息 m .若 m 是确认消息($m.msSyncGranted = \text{TRUE}$),且 p 的同步号小于 i 的同步号($m.msSync < sync[i]$),则 s 向 p 发送 $UpdateSyncRequest$ 请求,通知 p 升级到下一个同步号.下一个同步号的选举方法如上文所述.若 m 是否定消息,则 i 修改 $syncTrack[i][p]$ ($syncTrack'[i][p] = m.msSync$).
- $UpdateSync(i)$:状态为主节点候选者的节点 i ,如果存在一个多数派的节点完成对任期为 $sync[i]$ 的日志项的同步,且 i 收到它们的确认($RequestSyncResponse$ 消息),那么 s 向这些节点发送 $UpdateSyncRequest$ 消息,通知它们升级同步号为 i 的任期($currentTerm[i]$).

$$\begin{aligned}
 & \text{HandleRequestSyncResponse}(i) \triangleq \\
 & \quad \wedge \exists m \in \text{messages} : \\
 & \quad \quad \text{LET } j \triangleq m.\text{msourceIN} \\
 & \quad \quad \wedge m.\text{mtype} = \text{RequestSyncResponse} \\
 & \quad \quad \wedge m.\text{mdest} = i \\
 & \quad \quad \wedge \text{currentTerm}[i] = m.\text{mterm} \\
 & \quad \quad \wedge \text{currentState}[i] \in \{\text{Leader}, \text{LeaderCandidate}\} \\
 & \quad \quad \wedge \text{syncTrack}' = [\text{syncTrack EXCEPT } ![i][j] = m.\text{msync}] \\
 & \quad \quad \wedge \vee \wedge m.\text{msyncGranted} \\
 & \quad \quad \quad \wedge m.\text{msync} < \text{sync}[i] \\
 & \quad \quad \quad \wedge \text{Send}([\text{mtype} \mapsto \text{UpdateSyncRequest}, \\
 & \quad \quad \quad \quad \text{mterm} \mapsto \text{currentTerm}[i], \\
 & \quad \quad \quad \quad \text{msync} \mapsto \text{Min}(\{\text{sync}[i]\} \cup \{k \in \text{Nat} : k > m.\text{msync} \wedge \\
 & \quad \quad \quad \quad \quad \text{Cardinality}(\{n \in \text{Index} : \log[i][n].\text{term} = k\}) > 0\}), \\
 & \quad \quad \quad \quad \text{msource} \mapsto i, \\
 & \quad \quad \quad \quad \text{mdest} \mapsto \{j\})] \\
 & \quad \vee \wedge \neg m.\text{msyncGranted} \\
 & \quad \quad \wedge \text{UNCHANGED messages} \\
 & \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \log, \text{electionVars} \rangle \\
 \\
 & \text{UpdateSync}(i) \triangleq \\
 & \quad \wedge \text{currentState}[i] = \text{LeaderCandidate} \\
 & \quad \wedge \exists Q \in \text{Quorums} : \\
 & \quad \quad \text{LET } \text{syncUpdated} \triangleq \{m \in \text{messages} : \wedge m.\text{mtype} = \text{RequestSyncResponse} \\
 & \quad \quad \quad \wedge m.\text{mterm} = \text{currentTerm}[i] \\
 & \quad \quad \quad \wedge m.\text{msyncGranted} = \text{TRUE} \\
 & \quad \quad \quad \wedge m.\text{msync} = \text{sync}[i] \\
 & \quad \quad \quad \wedge m.\text{msource} \in Q \\
 & \quad \quad \quad \wedge m.\text{mdest} = i\} \\
 & \quad \quad \text{IN} \\
 & \quad \quad \wedge \forall q \in Q : (\exists m \in \text{syncUpdated} : m.\text{msource} = q) \vee q = i \\
 & \quad \quad \wedge \text{Send}([\text{mtype} \mapsto \text{UpdateSyncRequest}, \\
 & \quad \quad \quad \text{mterm} \mapsto \text{currentTerm}[i], \\
 & \quad \quad \quad \text{msync} \mapsto \text{currentTerm}[i], \\
 & \quad \quad \quad \text{msource} \mapsto i, \\
 & \quad \quad \quad \text{mdest} \mapsto Q]) \\
 & \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \log, \text{syncTrack}, \text{electionVars} \rangle
 \end{aligned}$$

- $\text{HandleUpdateSyncRequest}(i)$: i 收到 UpdateSyncRequest 请求 m , 升级自己的同步号 ($\text{sync}'[i] = m.\text{msync}$), 并将日志中的所有日志项标记为提交. 然后回复确认, 将升级后的同步号通知发送者. 收到 i 的确认的主节点(主节点候选者)根据回复修改对 i 的同步号的记录 (syncTrack).
- $\text{HandleUpdateSyncResponse}(i)$: 主节点或主节点候选者 i 收到 UpdateSyncReponse 回复后更新对发送者同步号的记录.

$\text{HandleUpdateSyncRequest}(i) \triangleq$

$$\begin{aligned} & \exists m \in \text{messages} : \\ & \quad \text{LET } grant \triangleq \wedge \text{currentTerm}[i] = m.mterm \\ & \quad \wedge m.msSync > sync[i] \\ & \quad j \triangleq m.msource \\ & \text{IN} \\ & \quad \wedge m.mtype = \text{UpdateSyncRequest} \\ & \quad \wedge i \in m.mdest \\ & \quad \wedge m.mterm \leq \text{currentTerm}[i] \\ & \quad \wedge \vee \wedge grant \\ & \quad \wedge sync' = [sync \text{ EXCEPT } !i = m.msSync] \\ & \quad \wedge log' = [log \text{ EXCEPT } !i = [n \in \text{Index} \mapsto \\ & \quad \quad \text{IF } log[i][n].term = sync[i] \text{ THEN} \\ & \quad \quad \quad log[i][n].committed = \text{TRUE} \\ & \quad \quad \text{ELSE } log[i][n]]] \\ & \vee \wedge \neg grant \\ & \quad \wedge \text{UNCHANGED } \langle log, sync \rangle \\ & \wedge \text{Send}([mtype \mapsto \text{UpdateSyncResponse}, mterm \mapsto \text{currentTerm}[i], \\ & \quad mupdateSyncGranted \mapsto grant, msync \mapsto sync'[i], \\ & \quad msource \mapsto i, mdest \mapsto j]) \\ & \wedge \text{UNCHANGED } \langle \text{currentTerm}, \text{currentState}, \text{votedFor}, \text{end}, \text{syncTrack}, \text{electionVars} \rangle \end{aligned}$$

$\text{HandleUpdateSyncResponse}(i) \triangleq$

$$\begin{aligned} & \wedge \exists m \in \text{messages} : \\ & \quad \text{LET } j \triangleq m.msource \text{ IN} \\ & \quad \wedge m.mtype = \text{UpdateSyncResponse} \\ & \quad \wedge m.mdest = i \\ & \quad \wedge \text{currentTerm}[i] = m.mterm \\ & \quad \wedge \text{currentState}[i] \in \{\text{Leader}, \text{LeaderCandidate}\} \\ & \quad \wedge \vee \wedge m.mupdateSyncGranted \\ & \quad \wedge syncTrack' = [syncTrack \text{ EXCEPT } !i[j] = m.msSync] \\ & \vee \wedge \neg m.mupdateSyncGranted \\ & \quad \wedge \text{UNCHANGED } syncTrack \\ & \wedge \text{UNCHANGED } \langle \text{messages}, \text{serverVars}, log, \text{electionVars} \rangle \end{aligned}$$

- $\text{BecomeLeader}(i)$: 主节点候选者 i 执行 UpdateSync 后, 若有一个多数派 Q 的节点的同步号升级为 $\text{currentTerm}[i]$, 且 i 收到它们的确认 ($\forall q \in Q: syncTrack[i][q] = currrentTerm[i]$), 那么 i 转变为主节点, 并提交日志中的所有日志项.
- $\text{ClientRequest}(i)$: i 为主节点时, 可以响应用户的请求, 将用户的命令插入到日志中.

$$\begin{aligned}
& \text{BecomeLeader}(i) \triangleq \\
& \quad \wedge \text{currentState}[i] = \text{LeaderCandidate} \\
& \quad \wedge \exists Q \in \text{Quorums} : \forall q \in Q : (q = i \vee \text{syncTrack}[i][q] = \text{currentTerm}[i]) \\
& \quad \quad \wedge \text{elections}' = \text{elections} \cup \\
& \quad \quad \{[\text{eterm} \mapsto \text{currentTerm}[i], \\
& \quad \quad \quad \text{esync} \mapsto \text{sync}[i], \\
& \quad \quad \quad \text{eleader} \mapsto i, \\
& \quad \quad \quad \text{evotes} \mapsto Q, \\
& \quad \quad \quad \text{evoterLog} \mapsto \{\log[k] : k \in Q\}, \\
& \quad \quad \quad \text{elog} \mapsto \log[i]\}] \\
& \quad \wedge \text{sync}' = [\text{sync EXCEPT } !(i) = \text{currentTerm}[i]] \\
& \quad \wedge \text{currentState}' = [\text{currentState EXCEPT } !(i) = \text{Leader}] \\
& \quad \wedge \log' = [\log \text{ EXCEPT } !(i) = [n \in \text{Index} \mapsto \\
& \quad \quad \quad \text{IF } \log[i][n].\text{term} = \text{sync}[i] \text{ THEN} \\
& \quad \quad \quad \quad \log[i][n].\text{committed} \mapsto \text{TRUE}] \\
& \quad \quad \quad \text{ELSE } \log[i][n]]] \\
& \quad \wedge \text{UNCHANGED } \langle \text{messages}, \text{currentTerm}, \text{votedFor}, \text{end}, \text{syncTrack}, \text{halfElections} \rangle \\
& \text{ClientRequest}(i, v) \triangleq \\
& \quad \text{LET } \text{nextIndex} \triangleq \log\text{Tail}(\log[i]) + 1 \\
& \quad \text{entry} \triangleq [\text{term} \mapsto \text{currentTerm}[i], \\
& \quad \quad \text{date} \mapsto \text{currentTerm}[i], \\
& \quad \quad \text{value} \mapsto v, \\
& \quad \quad \text{committed} \mapsto \text{FALSE}] \\
& \quad \text{IN} \\
& \quad \wedge \text{currentState}[i] = \text{Leader} \\
& \quad \wedge \log' = [\log \text{ EXCEPT } !(i)[\text{nextIndex}] = \text{entry}] \\
& \quad \wedge \text{UNCHANGED } \langle \text{messages}, \text{serverVars}, \text{electionVars}, \text{syncTrack} \rangle
\end{aligned}$$

Next 定义了次态关系.Spec 定义了完整的行为规约.

$$\begin{aligned}
\text{Next} \triangleq & \quad \vee \exists i \in \text{Server} : \text{Restart}(i) \\
& \vee \exists i \in \text{Server} : \text{Timeout}(i) \\
& \vee \exists i \in \text{Server} : \text{UpdateTerm}(i) \\
& \vee \exists i \in \text{Server} : \text{RequestVote}(i) \\
& \vee \exists i \in \text{Server} : \text{HandleRequestVoteRequest}(i) \\
& \vee \exists i \in \text{Server} : \text{BecomeLeaderCandidate}(i) \\
& \vee \exists i \in \text{Server} : \text{BecomeLeader}(i) \\
& \vee \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequest}(i, v) \\
& \vee \exists i, j \in \text{Server} : \text{RequestSync}(i) \\
& \vee \exists i \in \text{Server} : \text{HandleRequestSyncRequest}(i) \\
& \vee \exists i \in \text{Server} : \text{HandleRequestSyncResponse}(i) \\
& \vee \exists i, j \in \text{Server} : \text{UpdateSync}(i) \\
& \vee \exists i \in \text{Server} : \text{HandleUpdateSyncRequest}(i) \\
& \vee \exists i \in \text{Server} : \text{HandleUpdateSyncResponse}(i)
\end{aligned}$$

$$\text{Spec} \triangleq \text{Init} \wedge \square[\text{Next}]_{\text{vars}}$$

5 ParallelRaft-CE 的正确性证明

本节我们简要论证 ParallelRaft-CE 协议的正确性.ParallelRaft-CE 的正确性包括两部分:一是共识协议的一

致性;二是 ParallelRaft-CE 中不会出现"幽灵日志"现象,由此可以得到复制状态机的安全性:节点间执行每个日志项的顺序没有冲突,状态机的状态一致.

- 节点的任期与同步号单调递增:根据规约容易得到,任何节点的任期和同步号是单调递增的.
- 根据规约,节点只会修改任期等于同步号的日志项,而不会增加,删除或修改日志中任期小于同步号的日志项.
- 选举的安全性:每个任期最多选出一个主节点. ParallelRaft-CE 中从节点成为主节点需要两个阶段:一是从节点通过选举成为主节点候选者;二是主节点候选者转变为主节点.其中,选取主节点候选者的方式与 Raft 中选取主节点的方式相同,根据 Raft 选主的安全性,可以得到 ParallelRaft-CE 中每个任期最多选出一个主节点候选者.主节点候选者可能成功转变为主节点,可能由于失效,或是任期更大的节点发起选举等原因没有转变为主节点,在这种情况下,系统进入下一任期.因此每个任期最多选出一个主节点.
- 任期相同,编号相同的日志项一定包含相同的命令.根据选举的安全性,以及对任意编号,每个主节点最多添加一个日志项可以得到.这个性质与 Raft 相同.
- 性质 3:首先任何节点的同步号不大于任期.由于节点在发起选举需要增大任期(同步号不变),且在选举过程中任期与同步号不变.因此主节点候选者的同步号小于任期.设主节点候选者 i 的同步号为 s ,任期为 $t (s < t)$.
 - 对任意 $s < k < t$,不存在任期为 k 的主节点.使用反证法,假设存在节点 j 为任期 k 的主节点.根据规约,存在多数派节点 Q_1 为 i 投票. j 由主节点候选者转变为主节点,因此存在多数派节点 Q_2 升级同步号为 k .根据多数派的性质, $Q_1 \cap Q_2 \neq \emptyset$. 设 $r \in Q_1 \cap Q_2$. 若 r 先为 i 投票,那么 r 升级任期为 t . 由于 $t > k$,因此 r 之后会拒绝 j 升级同步号的请求,这与 $r \in Q_2$ 矛盾. 若 r 先升级同步号为 k ,那么由于 $k > s$,之后 r 会拒绝 i 的选举请求(注意选举投票要比较同步号),这与 $r \in Q_1$ 矛盾. 由此,可知不存在这样的 k .
 - 对 $k \geq t$,不存在任期为 k 的主节点.若系统中存在任期为 k 的主节点,那么 k 在选举中获得多数派节点 Q 的投票,因此 Q 中的节点的任期至少为 k .由于节点的任期递增,因此 i 发起选举时, Q 中的节点的任期要么大于 t ,要么已经为任期为 t 的其它节点投过票,因此不会为 i 投票.这与 i 成为任期为 t 的主节点候选者矛盾.
 - 综上,系统中没有出现过任期大于 s 的主节点,因此不存在任期大于 s 的日志项.
- 主节点的完全性:ParallelRaft-CE 中,任期为 t_1 主节点失效后,新的主节点候选者(任期为 t_2)会对任期为 t_1 的日志项进行重确认,在重确认的过程中,所有上一任期的日志项的状态被确定:被提交或者丢弃.重确认过程结束后,主节点候选者成为新的主节点.重确认的过程实质上是对每个位置执行 Paxos 协议的过程,因此可以保证取得共识,且之前取得共识的日志项不会丢失.根据规约,当完成重确认后,一个多数派的节点升级同步号为 t_2 ,且由于节点的同步号递增,之后,同步号小于 t_2 的节点无法通过选举成为主节点候选者,因此已经取得共识的任期为 t_1 的日志项不会改变.
- 一致性:ParallelRaft-CE 的主节点根据从节点的同步号同步日志项,以及恢复阶段的正确性(性质 3 与主节点的完全性),保证了节点间日志的一致性.
- ParallelRaft-CE 中不会出现"幽灵日志"现象.在 ParallelRaft-CE 中,新的主节点候选者通过重确认过程对上一主节点的日志项的进行重确认.而当重确认过程结束后,主节点候选者的同步号升级为任期,且由于节点的同步号递增,以及选举机制保证了同步号落后的节点无法成为主节点.因此系统不会重新提交之前节点的日志项.

6 模型检验

本节使用 TLC 模型检验工具验证 ParallelRaft-CE 协议的正确性,并验证 ParallelRaft-SE 到 Multi-Paxos 的

精化关系.

6.1 实验设置

在所有的实验中,我们调整参与者集合 Proposer(Server),提议者集合 Value,提议编号集合 Ballot(Term)的大小,并将前两者设置为对称集^[14],以提高 TLC 的验证效率.我们使用 10 个线程进行实验,以下是我们的实验统计结果:已遍历(BFS 方式遍历)的系统状态图的直径,已遍历的状态的数量,已发现的不同状态的数量以及检验花费的时间(单位是 hh:mm:ss).在验证 ParallelRaft-SE 精化 Multi-Paxos 时,我们设置 TLC 已检验的不同状态数量超过 1 亿时,人为停止实验.在验证 ParallelRaft-CE 满足一致性时,我们设置 TLC 已检验的不同状态数超过 2 亿时,人为停止实验.

实验使用的机器配置为: 2.40 GHz 10 核 CPU 以及 64GB 内存,TLC 版本号为 1.7.0

6.2 验证结果

6.2.1 ParallelRaft-SE 精化 Multi-Paxos

TLC 模型 (参与者数量, 提议值个数, 投票轮数)	状态图直径	状态数	不同状态数	检验时间 (hh : mm : ss)
(3,2,2)	22	1183766512	104836664	10:05:25
(3,2,3)	23	899846293	102806000	09:54:49
(3,3,2)	22	950017774	100064549	12:49:10
(3,3,3)	22	828085252	100020363	12:40:07
(4,2,2)	24	1045345827	100093827	12:52:38
(4,2,3)	22	1483023804	100044759	09:42:07
(4,3,2)	24	1150819236	100101884	23:21:09
(4,3,3)	21	1452458568	100054689	20:04:08

Figure 5 Model Checking Results of Verifying the Refinement from ParallelRaft-SE to Multi-Paxos

图 5 ParallelRaft-SE 精化 Multi-Paxos 的验证结果

图 5 给出了在不同配置下验证 ParallelRaft-SE 精化 Multi-Paxos 的验证结果.精化关系在 ParallelRaft-SE 的 TLA+ 规约中给出.从实验数据可以发现,参与者数量与提议值个数对实验的规模以及检验时间的影响较大,而投票轮数对实验的规模影响相对较小.

6.2.2 ParallelRaft-CE 满足一致性

TLC 模型 (参与者数量, 提议值个数, 投票轮数)	状态图直径	状态数	不同状态数	检验时间 (hh : mm : ss)
(2,2,2)	26	1274709286	201161468	01:51:02
(2,2,3)	24	1104433959	202169934	01:31:11
(2,3,2)	26	1218607009	200208445	02:27:13
(2,3,3)	24	1108870367	200916171	02:22:13
(3,2,2)	22	1370489381	200340341	03:02:17
(3,2,3)	19	1221006004	200913649	02:40:23
(3,3,2)	22	1369865234	200265437	05:44:01
(3,3,3)	19	1236091812	200041489	05:02:33

Figure 6 Model Checking Results of Verifying that ParallelRaft-CE Satisfies Consistency

图 6 ParallelRaft-CE 满足一致性的验证结果

图 6 给出了在不同配置下验证 ParallelRaft-CE 满足一致性的验证结果. ParallelRaft-CE 的一致性形式化定义为:

$$\begin{aligned} AllEntries(i) &\triangleq \{\langle n, \log[i][n] \rangle : n \in Index\} \\ Consistency &\triangleq \forall i, j \in Server : \\ &\{e \in AllEntries(i) : e[2].term \geq 0 \wedge e[2].term < \text{Min}(\{sync[i], sync[j]\})\} = \\ &\{e \in AllEntries(j) : e[2].term \geq 0 \wedge e[2].term < \text{Min}(\{sync[i], sync[j]\})\} \end{aligned}$$

其中 $AllEntries(i)$ 为节点 i 的所有日志项(日志编号与日志的二元组). $Consistency$ 定义了 ParallelRaft-CE 的一致性: 即任意两个节点, 若它们的同步号都大于 t , 那么它们的日志中有相同的任期为 t 的日志项.

ParallelRaft-CE 协议更加复杂, 不同状态的数量较多, 因此我们设置不同状态数达到 2 亿时停止实验.

7 相关工作

经典分布式共识协议 Multi-Paxos(Paxos)^[9] [10] 衍生出了多种变体^[25] [26]. 如 Disk Paxos^[27]、 Cheap Paxos^[28]、 Fast Paxos^[22]、 Generalized Paxos^[23]、 Stoppable Paxos^[29]、 Vertical Paxos^[30]、 Byzantine Paxos^[31]、 EPaxos^[32]、 CASPaxos^[33]、 TPaxos^[35] 等. Multi-Paxos 支持乱序提交、顺序执行用户命令. Raft^[11] 也可以看作 Multi-Paxos(Paxos) 的变体^[21], 它限制了乱序提交行为, 通过顺序提交、顺序执行用户命令简化了协议设计. 本文关注分布式文件系统 PolarFS 使用的 ParallelRaft 协议^[13]. 它基于 Raft 实现了乱序提交、乱序执行, 更适用于高并发系统. 为了理清 ParallelRaft 与 Raft 之间的关系, 我们提出了 ParallelRaft-SE. ParallelRaft-SE 是 Multi-Paxos 的一种变体, 支持乱序提交、顺序执行用户命令.

利用无冲突命令之间的可交换性可以提高共识协议的性能. 在 Generalized Paxos^[23] 中, 如果命令无冲突, 用户可以绕过主节点, 直接将命令广播给所有节点, 减少通信开销、提高性能. 微软公司提出的分布式复制框架 Tribble^[24] 利用多线程技术并发执行无冲突的命令项, 同时保证状态一致性. ParallelRaft^[13] 根据命令的逻辑区块地址进行冲突判断, 无冲突的命令可以乱序执行. 我们分析了 ParallelRaft 的乱序执行机制, 发现文献^[13] 中的描述忽略了可能会违反状态一致性的“幽灵日志”问题. 因此, 我们提出了 ParallelRaft-CE, 它通过限制乱序提交阶段的并行度, 避免了“幽灵日志”问题.

使用形式化规约语言描述分布式协议, 并使用模型检验工具验证预期的性质, 可以有效地增强人们对协议可靠性的信心. Lamport 使用 TLA+[14] [16] 描述了 Paxos^[38]、Fast Paxos^[22]、Byzantine Paxos^[31] 等协议, 并使用模型检验工具 TLC 验证了这些协议(在受限规模下的)正确性. Ongaro 给出了 Raft 的 TLA+ 规约^[12]. 本文给出了 ParallelRaft-SE 和 ParallelRaft-CE 的 TLA+ 规约, 并验证了它们(在受限规模下的)正确性.

精化技术^[20] 有助于理解各种协议之间的关系. 在研究 Paxos 时, Lamport 提出了共识问题的抽象描述 Consensus, 给出了分布式共识问题的集中式解决方案 Voting^[38], 并构建了从 Paxos 到 Voting 以及从 Voting 到 Consensus 的精化关系^[38]. Yi 等人在研究腾讯公司 PaxosStore 系统中的 TPaxos 时, 提出了 Voting 的一种变体 EagerVoting, 并构建了从 TPaxos 到 EagerVoting 以及从 EagerVoting 到 Consensus 的精化关系^[35]. 本文构建了从 ParallelRaft-SE 到 Multi-Paxos 的精化关系, 证明了 ParallelRaft-SE 的正确性.

8 总结与未来工作

本文的目标是为 PolarFS 文件系统中使用的 ParallelRaft 协议^[13] (它支持乱序提交、乱序执行用户命令) 提供严格的形式化规约并证明其正确性. 首先, 为了理清 ParallelRaft 与 Raft 之间的关系, 我们提出了允许乱序提交、顺序执行的 ParallelRaft-SE 协议, 并建立了从 ParallelRaft-SE 到 Multi-Paxos 的精化关系. 其次, 我们发现现有的 ParallelRaft 描述忽略了可能会违反状态一致性的“幽灵日志”问题, 并提出了 ParallelRaft-CE 协议. 通过限制 ParallelRaft-SE 在乱序提交阶段的并行度, ParallelRaft-CE 避免了“幽灵日志”问题. 最后, 我们给出了 ParallelRaft-SE 和 ParallelRaft-CE 的 TLA+ 规约, 并使用 TLC 模型检验工具验证了从 ParallelRaft-SE 到 Multi-Paxos 的精化

关系以及 ParallelRaft-CE 的正确性.

目前,我们正在使用 TLAPS(TLA+ Proof System)^{[15] [34]} 定理证明系统为 ParallelRaft-CE(以及 Raft)开发机械化正确性证明.此外,我们还将从顺序/乱序提交、顺序/乱序执行的角度对已知的分布式共识协议进行综述,并利用形式化方法研究它们之间的关系.

References:

- [1] Schneider, Fred B . Implementing fault-tolerant services using the state machine approach: a tutorial. *Acm Computing Surveys*, 1990, 22(4):299-319.
- [2] Fischer MJ, Lynch NA, Paterson MS. Impossibility of distributed consensus with one faulty process. *Journal of the ACM (JACM)*. 1985 Apr 1;32(2):374-82.
- [3] Herlihy M. Wait-free synchronization[J]. *ACM Transactions on Programming Languages and Systems (TOPLAS)*. 1991 Jan 1;13(1):124-49.
- [4] Corbett JC, Dean J, Epstein M, Fikes A, Frost C, Furman JJ, Ghemawat S, Gubarev A, Heiser C, Hochschild P, Hsieh W. Spanner: Google's globally distributed database. *ACM Transactions on Computer Systems (TOCS)*. 2013 Aug 1;31(3):1-22.
- [5] Cao W, Liu Y, Cheng Z, Zheng N, Li W, Wu W, Ouyang L, Wang P, Wang Y, Kuan R, Liu Z. POLARDB Meets Computational Storage: Efficiently Support Analytical Workloads in Cloud-Native Relational Database. In 18th USENIX Conference on File and Storage Technologies (FAST 2020) 2020 (pp. 29-41).
- [6] Isard M. Autopilot: automatic data center management. *ACM SIGOPS Operating Systems Review*. 2007 Apr 1;41(2):60-7.
- [7] Weil SA. Ceph: reliable, scalable, and high-performance distributed storage (Doctoral dissertation, University of California, Santa Cruz).
- [8] Zheng J, Lin Q, Xu J, Wei C, Zeng C, Yang P, Zhang Y. PaxosStore: high-availability storage made practical in WeChat. *Proceedings of the VLDB Endowment*. 2017 Aug 1;10(12):1730-41.
- [9] Lamport L. Paxos made simple. *ACM Sigact News*. 2001 Dec 4;32(4):18-25.
- [10] Lamport L. The part-time parliament. *ACM Transactions on Computer Systems (TOCS)*, 1998, 16(2): 133-169.
- [11] Ongaro D, Ousterhout J. In search of an understandable consensus algorithm. In 2014 USENIX Annual Technical Conference ({USENIX} {ATC} 14) 2014 (pp. 305-319).
- [12] Ongaro D. Consensus: Bridging theory and practice (Doctoral dissertation, Stanford University).
- [13] Cao W, Liu Z, Wang P, Chen S, Zhu C, Zheng S, Wang Y, Ma G. PolarFS: an ultra-low latency and failure resilient distributed file system for shared storage cloud database. *Proceedings of the VLDB Endowment*. 2018 Aug 1;11(12):1849-62.
- [14] Lamport L. Specifying Systems: The TLA+ Language and Tools for Hardware and Software Engineers. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2002.
- [15] Lamport L. The TLA+ hyperbook. <http://research.microsoft.com/enus/um/people/lamport/tla/hyperbook.html> [31. 8. 2015]. 2015.
- [16] Lamport L. The temporal logic of actions. *ACM Transactions on Programming Languages and Systems (TOPLAS)*. 1994 May 1;16(3):872-923.
- [17] Lamport L. Summary of tla+. <http://lamport.azurewebsites.net/tla/summary-standalone.pdf>
- [18] Yuan D, Luo Y, Zhuang X, Rodrigues GR, Zhao X, Zhang Y, Jain PU, Stumm M. Simple testing can prevent most critical failures: An analysis of production failures in distributed data-intensive systems. In 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 2014) 2014 (pp. 249-265).
- [19] Yu Y, Manolios P, Lamport L. Model checking TLA+ specifications. In Advanced Research Working Conference on Correct Hardware Design and Verification Methods 1999 Sep 27 (pp. 54-66). Springer, Berlin, Heidelberg.
- [20] Abadi M, Lamport L. The existence of refinement mappings. *Theoretical Computer Science*. 1991 May 31;82(2):253-84.
- [21] Wang Z, Zhao C, Mu S, Chen H, Li J. On the Parallels between Paxos and Raft, and how to Port Optimizations. In Proceedings of the 2019 ACM Symposium on Principles of Distributed Computing. 2019 Jul 16 (pp. 445-454).
- [22] Lamport L. Fast paxos. *Distributed Computing*. 2006 Oct 1;19(2):79-103.
- [23] Lamport L. Generalized consensus and paxos. Tech. Rep. MSR-TR-2005-33, Microsoft Research, 2005.
- [24] Guo Z, Hong C, Yang M, Zhou D, Zhou L, Zhuang L. Paxos made parallel. Microsoft Research Asia, Tech. Rep. 2012 Nov 29;118.

- [25] Butler Lampson. 2001. The ABCD's of Paxos. In Proceedings of the twentieth annual ACM symposium on Principles of distributed computing (PODC '2001). Association for Computing Machinery, New York, NY, USA, 13.
- [26] Van Renesse R, Altinbuken D. Paxos made moderately complex. ACM Computing Surveys (CSUR). 2015 Feb 17;47(3):1-36.
- [27] Gafni E, Lamport L. Disk paxos. Distributed Computing. 2003 Feb 1;16(1):1-20.
- [28] Lamport L, Massa M. Cheap paxos. InInternational Conference on Dependable Systems and Networks, 2004 Jun 28 (pp. 307-314). IEEE.
- [29] Lamport L, Malkhi D, Zhou L. Stoppable paxos. TechReport, Microsoft Research. 2008 Apr 28.
- [30] Lamport L, Malkhi D, Zhou L. Vertical paxos and primary-backup replication. In Proceedings of the 28th ACM symposium on Principles of distributed computing 2009 Aug 10 (pp. 312-313).
- [31] Lamport L. Byzantizing Paxos by refinement. InInternational Symposium on Distributed Computing 2011 Sep 20 (pp. 211-224). Springer, Berlin, Heidelberg.
- [32] Moraru I, Andersen DG, Kaminsky M. There is more consensus in egalitarian parliaments. In Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles 2013 Nov 3 (pp. 358-372).
- [33] Rystsov D. CASPaxos: Replicated state machines without logs. arXiv preprint arXiv:1802.07000. 2018 Feb 20.
- [34] Chaudhuri, K, Doligez, D, Lamport, L, Merz, SA TLA+ Proof System. In Proceedings of the LPAR Workshops, CEUR Workshop 2008 (pp. 17-37).
- [35] Yi XC, Wei HF, Huang Y, Qiao L, Lu J. TPaxos in PaxosStore: derivation, specification, and refinement. Ruan Jian Xue Bao/Journal of Software, 2020, 31(8) (in Chinese).
- [36] <https://github.com/HappyCS-Gu/Parallel-Raft-tla>
- [37] <https://www.mysql.com/>
- [38] Lamport L, Merz S, Doligez D. A TLA + specification of Paxos and its refinement.
<https://github.com/tlaplus/Examples/tree/master/specifications/Paxos> (Accessed Sep 5, 2019).

附中文参考文献:

- [35] 易星辰,魏恒峰,黄宇,乔磊,吕建.PaxosStore 中共识算法 TPaxos 的推导、规约与精化.软件学报,2020,31(8)