

**PONTIFÍCIA UNIVERSIDADE CATÓLICA DO RIO GRANDE DO SUL
ESCOLA POLITÉCNICA
BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO**

RAFT IGUALITÁRIO — ERAFT

ERICK PINTOR

Trabalho de Conclusão IV apresentado como requisito parcial à obtenção do grau de Bacharel em Ciência da Computação na Pontifícia Universidade Católica do Rio Grande do Sul.

Orientador: Prof. Fernando Luís Dotti

**Porto Alegre
2023**

RAFT IGUALITÁRIO — ERAFT

RESUMO

O problema do consenso é um dos desafios fundamentais da área de sistemas distribuídos. Nele, processos independentes devem concordar sobre um determinado valor, proposto entre eles, e que todos adotam. Com uma literatura rica que data do início dos anos 90, o problema do consenso tem sido explorado por diversos ângulos, nos quais sua complexidade é notável. O algoritmo Raft tenta prover uma solução de fácil implementação contando com crescente adoção na indústria. Devido à alta demanda por sistemas distribuídos em redes de longa distância, algoritmos de consenso generalizado e igualitários foram propostos para aumentar a vazão do consenso. Neste trabalho de conclusão de curso introduzimos um novo algoritmo chamado Raft Igualitário, ou Egalitarian Raft (ERaft), que se baseia em Raft para construir um algoritmo de consenso igualitário de fácil implementação.

Palavras-Chave: sistemas distribuídos, consenso, algoritmos de consenso, consenso igualitário, consenso sem líder.

LISTA DE FIGURAS

2.1	Troca de mensagens em uma instância do algoritmo Paxos.	12
2.2	Diagrama de transição de estados do algoritmo Raft.	13
2.3	Troca de mensagens de consenso igualitário com conflito.	16
3.1	Modelo de replicação do <i>log</i>	19
3.2	Diagrama de estados de processos no algoritmo ERaft.	19
3.3	Diagrama de transição de fases de um comando em ERaft.	20
3.4	Exemplo de ciclos de dependências.	20

LISTA DE ALGORITMOS

3.1	Algoritmo que define o estado inicial de um processo p	22
3.2	Algoritmo que adiciona um novo comando ao <i>log</i> do líder.	22
3.3	Algoritmo acionado para a um evento do tipo <i>timeout</i>	23
3.4	Algoritmo que envia comandos do líder para os seguidores.	23
3.5	Algoritmo utilizado pelo seguidor ao receber um comando do líder.	24
3.6	Algoritmo utilizado pelo líder ao receber uma resposta do seguidor.	24
3.7	Algoritmo utilizado pelo líder avalia o progresso do consenso.	25
3.8	Algoritmo utilizado para consolidar logs coletados durante uma eleição. . . .	27

SUMÁRIO

1	INTRODUÇÃO	6
2	CONTEXTUALIZAÇÃO	9
2.1	CONSENSO CLÁSSICO	10
2.1.1	ALGORITMO PAXOS	10
2.1.2	ALGORITMO RAFT	12
2.2	CONSENSO IGUALITÁRIO	14
3	RAFT IGUALITÁRIO	17
3.1	MOTIVAÇÃO	17
3.2	MODELO DO SISTEMA	17
3.3	ALGORITMO	18
3.3.1	IDEIA DO ALGORITMO	18
3.3.2	REPLICAÇÃO DE COMANDOS	21
3.3.3	RECUPERAÇÃO DE FALHAS	26
3.4	INVARIANTES E PROPRIEDADES	28
3.5	CONSIDERAÇÕES PRÁTICAS	28
4	CONCLUSÃO	30
	REFERÊNCIAS	31
	APÊNDICE A – Especificação em TLA⁺	33

1. INTRODUÇÃO

Arquiteturas tradicionais de desenvolvimento de *software* têm sido desafiadas pela crescente demanda por interconectividade. Com a necessidade cada vez mais frequente de expandir seus negócios globalmente, departamentos de tecnologia são forçados a projetar aplicações de alta escalabilidade e tolerantes a falhas. Dentre os vários desafios inerentes a esta categoria de aplicações encontra-se o problema do consenso: como garantir que processos distribuídos e assíncronos concordem sobre um determinado valor proposto. Este é um problema fundamental instanciado nos mais diversos cenários de interesse da Computação Distribuída.

De maneira informal, o consenso consiste em um grupo finito de processos concordarem sobre um valor, proposto entre eles, e que todos adotam. Ou seja, ele representa a capacidade de processos tomarem decisões conjuntas. O problema do consenso deixa de ser trivial quando supomos a possibilidade de falha e que o sistema é assíncrono — nada se pode supor sobre os atrasos de mensagens trocadas e a velocidade relativa dos processos [4]. Visto que a possibilidade de falhas deve inevitavelmente ser considerada em sistemas distribuídos e que na maior parte das vezes o meio de comunicação (rede) apresenta latências não-determinísticas, o estudo e avanço de soluções para o problema do consenso foi e permanece de suma importância no cenário atual de grande conectividade e dependência de sistemas corretos altamente disponíveis.

O trabalho seminal de Leslie Lamport, de 1990, na criação do algoritmo de consenso Paxos [7] (publicado em 1998 e republicado em 2019) deu início a uma nova área de pesquisa que contribui para a criação de sistemas distribuídos tolerantes a falhas. No consenso clássico, existe um processo, normalmente denominado líder, responsável por coordenar as tomadas de decisão que levam ao consenso. Em caso de falha do líder, os demais processos devem eleger um novo processo como líder antes de prosseguir. Intuitivamente, a vazão de tais sistemas é limitada pela capacidade do líder em coordenar o consenso. Ainda, valores propostos devem ser enviados para o líder, o que impõe acréscimo de atraso no tratamento de requisições dos clientes. Esta penalidade aumenta conforme a distribuição geográfica de processos através de redes de longa distância.

Na tentativa de aumentar a vazão em sistemas baseados em consenso, surge na literatura o conceito de consenso generalizado ([14] e [9]). Ao contrário do consenso clássico, que admite que apenas um valor — ou comando — seja coordenado por vez de modo a estabelecer uma ordem total de comandos a serem executados por todos os participantes, o consenso generalizado permite que comandos não-conflitantes sejam coordenados paralelamente. Comandos não-conflitantes são aqueles que do ponto de vista da aplicação podem comutar, ou seja, podem ser executados em qualquer ordem sem alterar o estado final da aplicação. Desta forma, o consenso generalizado adota uma relação de ordem par-

cial entre comandos propostos para permitir maior vazão em cargas de trabalho com baixas taxas de conflito entre comandos.

O próximo passo na evolução da literatura chega em 2013 com o algoritmo Egalitarian Paxos (EPaxos) [12], que introduz o conceito de consenso igualitário. Baseando-se no consenso generalizado, EPaxos permite todos os processos participarem de todas as instâncias de consenso de forma concorrente, aumentando significativamente sua vazão para cargas de trabalho com altas taxas de conflito. Em 2020, o algoritmo Atlas [3], similar a EPaxos, foi publicado trazendo ainda maiores avanços para a vazão neste tipo de sistema.

Outra perspectiva acerca de algoritmos de consenso, tal como documentado pela pesquisa que levou ao desenvolvimento do algoritmo Raft [13], é que a dificuldade de compreensão destes algoritmos representa uma barreira para sua adoção. Na prática, esta dificuldade se traduz em implementações sobre as quais se tem menor confiabilidade acerca de suas propriedades. Sendo assim, Raft foi criado com o intuito de prover um algoritmo legível, de maior facilidade de compreensão, o que resulta em maior probabilidade de implementações aderentes à especificação do protocolo. Como resultado, ele conta hoje com uma rica base de implementações e crescente adoção pela indústria ¹. No entanto, Raft permanece na classe de algoritmos de consenso clássicos com liderança forte, não se beneficiando, assim, dos avanços de vazão e latência alcançados pelos algoritmos de consenso generalizados e igualitários.

Apesar das vantagens descritas na literatura de consenso generalizado e igualitário, tais algoritmos tendem a adicionar funções e regras a algoritmos de complexidade já elevada em termos de compreensão. Não é de nosso conhecimento a tentativa de derivar um algoritmo desta classe, focado na legibilidade de sua especificação, aos moldes de Raft. Supõe-se que tal algoritmo possa acelerar o uso da abordagem de consenso igualitário, na prática, de forma análoga à crescente adoção de Raft. Logo, neste trabalho de conclusão de curso, propomos a pesquisa e o desenvolvimento de um algoritmo de consenso igualitário seguindo os mesmos princípios de projeto deste algoritmo: fácil compreensão e implementação.

Assim, o presente trabalho apresenta as seguintes contribuições:

- Proposta de um algoritmo de consenso igualitário a partir dos princípios de modelagem de Raft;
- Agregação dos principais aspectos de consenso igualitário a partir dos princípios do protocolo Atlas;
- Adequação do funcionamento em fases, do Atlas, para o modelo de replicação de comandos do Raft;

¹Algumas implementações de código aberto de Raft (<https://raft.github.io/#implementations>) são parte de outros sistemas como, por exemplo, Etcd e RethinkDB.

- Proposta de um procedimento recuperação de comandos em caso de falha de processos;
- Especificação do algoritmo proposto utilizando a linguagem formal *Temporal Logic of Actions*² (TLA+).

O capítulo 2 deste trabalho contextualiza o problema do consenso, oferecendo maiores detalhes sobre o consenso clássico e igualitário. O capítulo 3 descreve algoritmo proposto em relação à sua especificação formal incluída no apêndice A. Considerações finais são feitas no capítulo 4.

²<https://lampert.azurewebsites.net/tla/tla.html>

2. CONTEXTUALIZAÇÃO

O problema do consenso pode ser visto como a capacidade de múltiplos processos de concordarem e adotarem um valor proposto entre eles. Este problema se instancia nas mais diversas áreas da Computação Distribuída, como bancos de dados distribuídos ¹, sincronização de configurações em parques de máquinas ² e orquestradores de ambientes ³ [6].

Para ilustrar uma instância possível do problema do consenso, considere o seguinte exemplo: uma instituição financeira oferece saques e depósitos para clientes ao redor do mundo. Suponha dois clientes com acesso à mesma conta bancária, cujo saldo é de R\$100,00 reais. Imagine que ambos os clientes façam, simultaneamente, a solicitação de um saque no valor total do saldo da conta. Intuitivamente, o comportamento esperado é de que apenas um cliente consiga sacar o valor solicitado e que o outro receba uma mensagem informando não haver saldo suficiente para executar o seu pedido. Em uma aplicação não-distribuída, este é um problema trivial, visto que um único processo é responsável por responder às duas solicitações, conseguindo decidir a ordem de execução: autorização da primeira e rejeição da segunda. Por outro lado, em um sistema distribuído, é possível que cada solicitação de saque seja recebida por um processo diferente. Para garantir o mesmo nível de consistência da solução anterior, estes processos devem entrar em consenso sobre quais solicitações estão sendo executadas naquele momento e, principalmente, em que ordem.

Ao generalizar o problema do consenso, obtemos uma classe de algoritmos aplicável a diversos domínios. De maneira genérica, temos que um algoritmo de consenso é projetado para garantir que um conjunto finito de processos distribuídos e assíncronos concordem e adotem um determinado valor proposto — quais saques executar e em qual ordem. Para isso, processos trocam mensagens através de um meio de comunicação (rede). Uma vez determinado o consenso, todos os processos envolvidos devem adotar a decisão tomada, esta sendo imutável. Para considerar um algoritmo de consenso correto, ele deve seguir as seguintes propriedades:

- **Terminação:** todo processo correto adotará a decisão do consenso em algum momento;
- **Validade:** se um processo adota um valor v , então v foi proposto por algum processo;
- **Integridade:** um processo não adota mais de um valor;
- **Acordo:** todo processo correto deve concordar sobre o mesmo valor.

¹Fauna (fauna.com), HBase (hbase.apache.org) e MongoDB (mongodb.com).

²Zookeeper (zookeeper.apache.org.), Consul (www.consul.io) e Etcd (coreos.com/etcd).

³Kubernetes (kubernetes.io), Docker Swarm (github.com/docker/swarm) e Mesos (mesos.apache.org).

Existem diferentes classes de falhas que dificultam preservar tais propriedades, dentre elas, falhas bizantinas tendem a ser particularmente desafiadoras. Em um sistema suscetível a falhas bizantinas, não é possível supor que os processos participantes são corretos. Um processo malicioso pode, intencionalmente, enviar mensagens contraditórias sobre os demais processos — como dizer ter aceitado um saque de R\$300,00 reais mesmo com saldo insuficiente, por exemplo. Este é um problema comum em sistemas de *blockchain* nos quais os participantes do consenso podem não ser entidades confiáveis.

Neste trabalho de conclusão de curso focaremos apenas em algoritmos que toleram falhas não-bizantinas, dentre elas a interrupção abrupta de um processo e sua recuperação (modelo *crash-recovery*), além de falhas de comunicação entre processos. Em ambos os casos, para preservar as propriedades do consenso, deve haver uma etapa de recuperação de um processo falho de forma que o mesmo consiga estar ciente de decisões tomadas durante sua falha. É comum a este tipo de algoritmo suportar a falha de uma minoria de processos sem comprometer a disponibilidade do serviço prestado. No entanto, no caso de falhas catastróficas, nas quais a maioria dos processos falha, ainda se faz necessário preservar as propriedades de segurança e funcionamento (*safety*), garantindo a consistência dos dados do sistema, enquanto a terminação (*liveness*) é violada.

Visto o problema do consenso, nas seções seguintes são apresentados algoritmos de consenso clássicos (seção 2.1) e igualitários (seção 2.2).

2.1 Consenso Clássico

2.1.1 Algoritmo Paxos

Dentre os algoritmos de consenso clássicos, Paxos [7] tem notória importância devido à sua contribuição histórica para a área de sistemas distribuídos. Leslie Lamport estabelece nos anos 90 muitos dos fundamentos utilizados até hoje em algoritmos similares. Em Paxos, um processo líder, também chamado *proposer*, é responsável por coordenar uma instância do algoritmo de consenso através de duas trocas de mensagens com os demais processos, denominados *acceptors* e *learners*. *Acceptors* são processos que participam ativamente do consenso enquanto *learners* apenas recebem, ou aprendem, o resultado do consenso. Paxos é dividido em duas fases, normalmente subdivididas em duas etapas cada:

- Fase 1a:
 - *Prepare(b)*: mensagem enviada pelo *proposer* para *acceptors* onde *b* é o identificador da rodada de troca de mensagens para a instância do consenso que está

sendo executada — Paxos permite múltiplas tentativas, ou rodadas, até que o consenso seja alcançado;

- Fase 1b:
 - *Promise(b, v)*: enviada de *acceptors* para *proposer* em resposta a uma mensagem do tipo *prepare*. Nela, v pode ser nulo se b é desconhecido, ou o valor anteriormente aceito caso contrário;
- Fase 2a:
 - *Accept(b, v)*: enviada pelo *proposer* para *acceptors* após receber um quórum de respostas do tipo *promise*;
- Fase 2b:
 - *Accepted(b, v)* enviada de *acceptors* para o *proposer* e *learners* em resposta a uma mensagem do tipo *accept*.

A figura 2.1 demonstra uma troca de mensagens possível para uma instância de consenso do algoritmo Paxos. Note que, em um ambiente distribuído e assíncrono, quaisquer destas mensagens podem se perder devido à falha do processo emissor, receptor, ou por falha de rede.

Para preservar as propriedades descritas na seção 2, Paxos emprega uma série de regras derivadas do número identificador b e os valores v trocados entre processos. A sobreposição de estados possíveis se torna, especialmente na presença de falhas, um fator de complexidade não-trivial que dificulta a compreensão e reprodução do algoritmo de forma correta.

Múltiplos refinamentos ao algoritmo Paxos surgiram ao longo dos anos tanto na forma do consenso clássico quanto no consenso generalizado. Dentre eles se destacam:

- Consenso clássico:
 - *Paxos Made Simple* [8]
 - *Cheap Paxos* [11]
 - *Paxos Made Moderately Complex* [16]
- Consenso Generalizado:
 - *Generalized Consensus and Paxos* [9]
 - *Fast Paxos* [10]
 - *On Making Generalized Paxos Practical* [15]

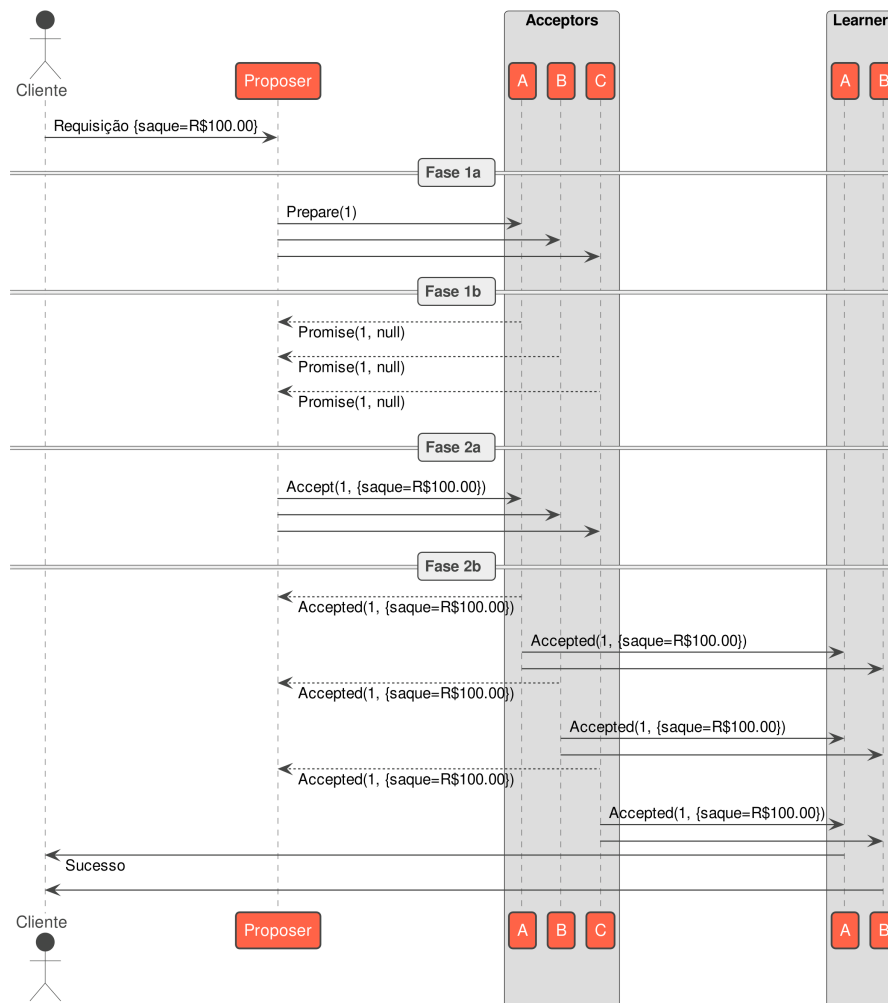


Figura 2.1: Troca de mensagens em uma instância do algoritmo Paxos.

Paxos continua sendo um dos algoritmos fundamentais da área de sistemas distribuídos. A partir dele, surgem os princípios e propriedades básicas que permitem a construção de algoritmos de consenso consistentes e tolerantes a falhas. No entanto, sua complexidade inspirou pesquisadores a buscar novas formas de projetar algoritmos de consenso, como o algoritmo Raft (seção 2.1.2).

2.1.2 Algoritmo Raft

Raft [13] é um algoritmo de consenso focado em legibilidade e fácil implementação. Ele surge a partir da frustração de seus autores em compreender e reproduzir variações do algoritmo Paxos. Em sua pesquisa, os autores mencionam como até mesmo especialistas na área de sistemas distribuídos tem, por vezes, dificuldade em explicar em detalhes como Paxos resolve o problema do consenso na presença de falhas. Ainda, o nível de abstração em que Paxos é descrito deixa a cargo do programador resolver questões práticas, porém não triviais, de sua implementação. Por exemplo, é comum em sistemas que precisam re-

resolver consenso, o registro (*log*) de decisões passadas para fins de recuperação de falhas. No entanto, Paxos não deixa claro como este *log* deve ser mantido ou quando é seguro remover registros antigos dele.

De modo a prover um algoritmo de consenso de fácil compreensão, Raft toma decisões de projeto pragmáticas com foco na clareza de sua especificação. Como resultado, implementações de Raft tendem a ser mais próximas de sua descrição, dado que o programador tem pouca ou nenhuma questão prática de implementação que não está coberta pela descrição do algoritmo.

Em alto nível, cada processo do Raft pode se encontrar em um dos três estados: “seguidor”, “candidato”, ou “líder”. Um **seguidor** se torna **candidato** se não receber mensagens do **líder** em um determinado período. Ao se tornar candidato, ele promove uma eleição na tentativa de se tornar o líder, ou volta a ser um seguidor caso mal sucedido. A figura 2.2 demonstra as transições possíveis de estado de um processo Raft.

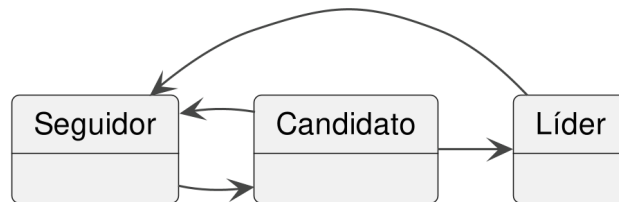


Figura 2.2: Diagrama de transição de estados do algoritmo Raft.

Líderes recebem valores (comandos) de clientes e os adicionam a seu *log*. Periodicamente, o líder envia comandos pendentes para os seguidores, os quais os adicionam a seu *log*. Assim que um comando tenha sido replicado para a maioria dos seguidores, o líder registra o seu índice como persistido (*committed*). Comandos persistidos são executados por todos os processos na ordem do *log* e podem ser removidos se todos os processos já o executaram. Raft emprega apenas um RPC (*Remote Procedure Call*) para replicar comandos entre líder e seguidores:

$$\text{Append}(\text{term}, \text{leaderID}, \text{prevLogIndex}, \text{entries}, \text{leaderCommit}) \rightarrow (\text{term}, \text{success})$$

Note que, ao replicar comandos via *Append*, líder e seguidores trocam toda informação necessária para as tomadas de decisão do algoritmo. Sendo elas:

- Requisição:
 - *term*: identificador único incrementado em cada eleição;
 - *leaderID*: identificador do líder associado ao *term* atual;
 - *prevLogIndex*: posição anterior do *log* onde *entries* se inicia;

- *entries*: quais comandos devem ser adicionados no *log*;
 - *leaderCommit*: qual o maior índice dentre os comandos já persistidos.
- Resposta:
 - *term*: o termo atual do seguidor;
 - *success*: indica se a requisição foi bem sucedida.

Essa estratégia diminui a complexidade do algoritmo visto que a sobreposição possível de estados durante falhas diminui. Ainda, esta decisão de projeto fornece um modelo mental intuitivo para quem implementa o algoritmo: repita o RPC *Append* enquanto o processo for líder; caso contrário, volte ao estado de seguidor.

É importante ressaltar que legibilidade e facilidade de implementação são critérios subjetivos na avaliação de algoritmos. No entanto, Raft traz com sua publicação uma pesquisa qualitativa que sugere uma real vantagem de Raft sobre Paxos na capacidade de programadores não familiarizados em seguirem suas especificações e implementarem uma versão correta de cada algoritmo.

2.2 Consenso Igualitário

O papel de um processo líder em algoritmos de consenso clássicos simplifica o modelo e facilita preservar as propriedades do consenso. Entretanto, ele se torna um fator limitante no desempenho desses algoritmos, principalmente em sistemas distribuídos em longas distâncias geográficas. Em tais sistemas, mesmo que haja um processo geograficamente próximo ao cliente, o mesmo deve enviar seus comandos para o líder, que pode estar localizado em outra região geográfica, causando aumento do tempo de resposta devido à latência de rede.

Alternativas foram propostas para diminuir o papel do líder e aumentar a vazão. Dentre elas, Mencius [1] propõe um algoritmo baseado em Paxos no qual o líder de cada instância de consenso é escolhido deterministicamente de maneira a rotacionar este papel entre os processos participantes. Baseando-se no Paxos Generalizado [9], Mencius consegue concluir instâncias do consenso paralelamente. Este algoritmo promove, então, o balanceamento de carga ao distribuir o papel de liderança entre processos. No entanto, ele ainda é suscetível a alta latência em sistemas geograficamente distribuídos.

O surgimento de Egalitarian Paxos (EPaxos) [12] traz um marco importante na busca por algoritmos de consenso para sistemas distribuídos em redes de longa distância. Esta versão de Paxos permite que cada processo seja líder dos comandos que recebe de seus clientes, reduzindo assim a latência do sistema dado que clientes podem se comunicar com o processo localizado mais próximo de si. Surge, então, uma nova classe de algoritmos

chamados sem-líder (*leaderless*), apesar de o termo “igualitário” ser mais apropriado uma vez que todo processo é líder de seus próprios comandos.

Posteriormente, o artigo “Leaderless State-Machine Replication: Specification, Properties, Limits” [5] reflete sobre esta classe de algoritmos e os generaliza de modo a estabelecer seus componentes e propriedades fundamentais. Conclui-se que algoritmos igualitários requerem uma função de conflito. Dois comandos são ditos conflitantes se, quando executados em diferentes ordens relativas, produzem resultados diferentes. Uma definição mais concreta de conflito estabelece que dois comandos são conflitantes se um escreve no mesmo conjunto de dados acessado (lidos ou escritos) pelo outro comando. Por exemplo, um saque de R\$100,00 reais e um depósito de R\$100,00 reais executados na mesma conta bancária com saldo atual igual a zero produzem resultados diferentes para o cliente se executados em ordens diferentes.

A primeira etapa de um algoritmo de consenso igualitário é coletar os conflitos de seus comandos pendentes, etapa conhecida como caminho rápido (*fast path*). Comandos sem conflitos são persistidos e executados após esta fase. Comandos conflitantes são promovidos para o chamado caminho lento (*slow path*), o qual requer uma nova troca de mensagens para garantir que todos os processos tomem conhecimento de todos os conflitos para estes comandos. Somente ao fim do caminho lento, comandos conflitantes podem ser persistidos e executados para garantir que todos os processos sigam a mesma ordem relativa entre comandos conflitantes. A figura 2.3 demonstra, em alto nível, uma troca de mensagens possível de consenso igualitário com a presença de conflito.

EPaxos emprega função de conflito, caminhos rápidos e lentos. No entanto, a chegada do algoritmo Atlas [3] traz uma contribuição significativa para a área: ao contrário de EPaxos que requer o caminho lento para quaisquer comandos conflitantes, Atlas é capaz de pular esta etapa se determinado que um quórum de processos tem conhecimento dos mesmos conflitos. Com esta otimização, Atlas demonstrou maior desempenho em avaliação experimental tanto sobre EPaxos como sobre Mencius.

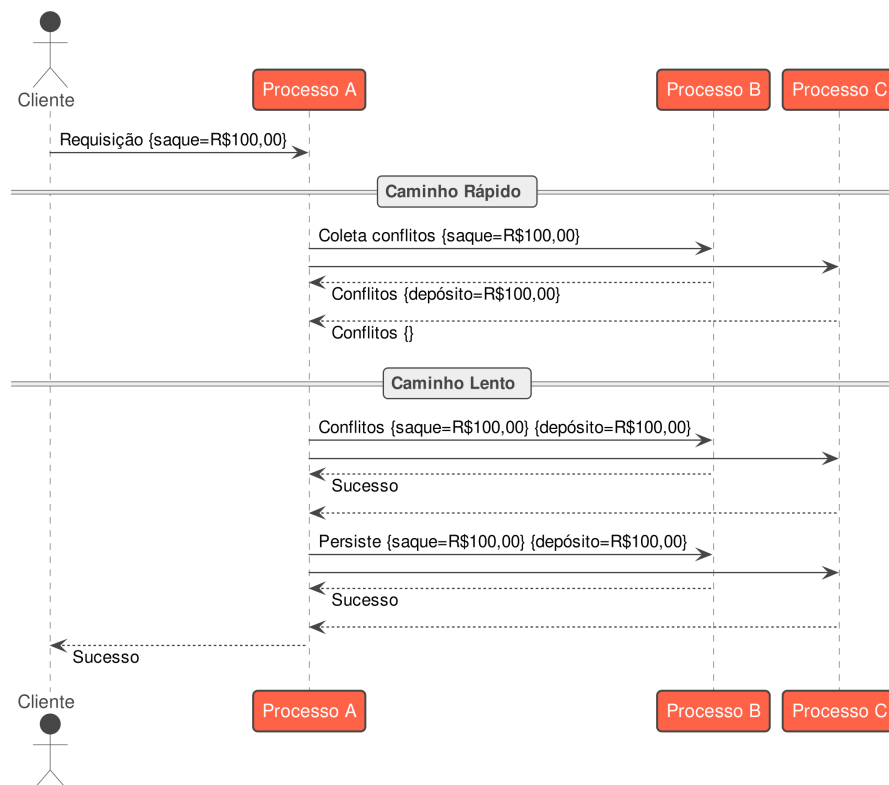


Figura 2.3: Troca de mensagens de consenso igualitário com conflito.

3. RAFT IGUALITÁRIO

3.1 Motivação

Ao analisar as diferentes formas de consenso descritas na literatura, é possível observar a evolução desta área de pesquisa buscando aumentar sua vazão principalmente em sistemas geograficamente distribuídos em redes de longa distância. Todavia, em seu estado atual, Raft permanece incapaz de se beneficiar desta evolução, já que este é um algoritmo clássico com forte papel do processo líder na coordenação do consenso.

Por acreditar que o modelo estabelecido por Raft possa ser adaptado para a classe de algoritmos de consenso igualitário ainda preservando a legibilidade de sua especificação, neste trabalho de conclusão de curso propomos um novo algoritmo de consenso chamado Raft Igualitário, ou Egalitarian Raft (ERaft).

Apesar das mudanças necessárias para adaptar o modelo proposto por Raft ao consenso igualitário, presumimos que, ao adotar Raft como algoritmo base, temos uma plataforma consolidada de fácil compreensão, permitindo tais refinamentos com pouco acréscimo de complexidade. Ademais, o Raft, com todo o seu material disponível, serve como introdução para ERaft da mesma forma que Paxos estabelece a base fundamental para compreender Paxos Generalizado e Paxos Igualitário. Por fim, acreditamos que o modelo criado por Raft pode evoluir e se beneficiar dos avanços documentados na literatura de consenso igualitário com a introdução de ERaft, sendo este possível propulsor da adoção de consenso igualitário pela indústria.

3.2 Modelo do Sistema

Assumimos um sistema distribuído composto por processos interligados. Existe um conjunto ilimitado $C = \{c_1, c_2, \dots\}$ de processos clientes e um conjunto limitado $S = \{s_1, s_2, \dots, s_n\}$ de processos servidores.

O sistema pode se comportar de forma assíncrona, ou seja, não há limites para atrasos de mensagens e velocidades relativas de processo. Entretanto, para garantia de progresso supõe-se que em algum momento o sistema passa a se comportar de forma síncrona. Este momento é chamado de Tempo Global de Estabilização [2] (GST - *Global Stabilization Time*), sendo desconhecido pelos processos.

Assumimos o modelo de falha de colapso (*crash*) e excluimos comportamentos maliciosos e arbitrários, ou seja, sem falhas bizantinas. O processo pode falhar e se recuperar.

Os processos se comunicam somente por troca de mensagens, usando comunicação um-para-um. Não há memória compartilhada. A comunicação é através das primitivas *send(m)* e *receive(m)*, onde *m* é uma mensagem. Se um remetente enviar uma mensagem várias vezes, um destinatário correto acabará recebendo a mensagem.

3.3 Algoritmo

O algoritmo ERaft desenvolvido neste trabalho foi modelado utilizando a linguagem de especificação formal TLA⁺¹ (apêndice A). Sua especificação oferece uma definição precisa para o funcionamento do algoritmo quanto à replicação de comandos e recuperação em caso de falhas. O restante deste capítulo se dedica a apresentar em alto nível o funcionamento do algoritmo de modo a enriquecer o entendimento de sua especificação.

3.3.1 Ideia do Algoritmo

Imagine um sistema com um único processo que implementa o algoritmo Raft. Este processo possui um estado local composto por seus comandos (*log*) e seus índices de replicação, *commit*, e execução. Dado que há apenas um processo, temos que o mesmo será o líder do consenso. Em uma implementação clássica de Raft, ao adicionar novos processos ao sistema, os mesmos se tornam seguidores do líder vigente e passam a replicar o seu *log* (Figura. 3.1a). Já em ERaft, novos processos multiplicam o estado: cada processo tem o seu próprio *log* e replica o *log* dos demais (Figura. 3.1b).

Em Raft, intervalos de tempo fechados (*timeouts*) são utilizados como mecanismo de tomada de decisão nos processos participantes. Ao iniciar o algoritmo, cada processo começa uma contagem de tempo que, ao término de um intervalo predeterminado, dispara uma tomada de decisão antes de iniciar o próximo intervalo de tempo. ERaft utiliza a mesma abordagem, porém, para suportar múltiplos *logs*, é necessário também manter múltiplos *timeouts*, um por *log*.

Assim como em Raft, todos os processos em ERaft iniciam no estado “seguidor”. Ao término de um ciclo de *timeout*, um “seguidor” passa para o estado de “candidato” e inicia uma eleição na tentativa de se tornar líder de seu próprio *log*. Após o término da eleição, em uma execução sem falhas, o processo candidato torna-se líder e passa a receber comandos dos clientes. ERaft traz uma modificação no diagrama de estados de Raft (Figura. 2.2) na qual, por motivos de recuperação de falhas, é permitido que líderes passem para o estado de candidatos sem antes passar para o estado de seguidores (Figura. 3.2) — maiores detalhes são discutidos na seção 3.3.2.

¹<https://lampport.azurewebsites.net/tla/tla.html>

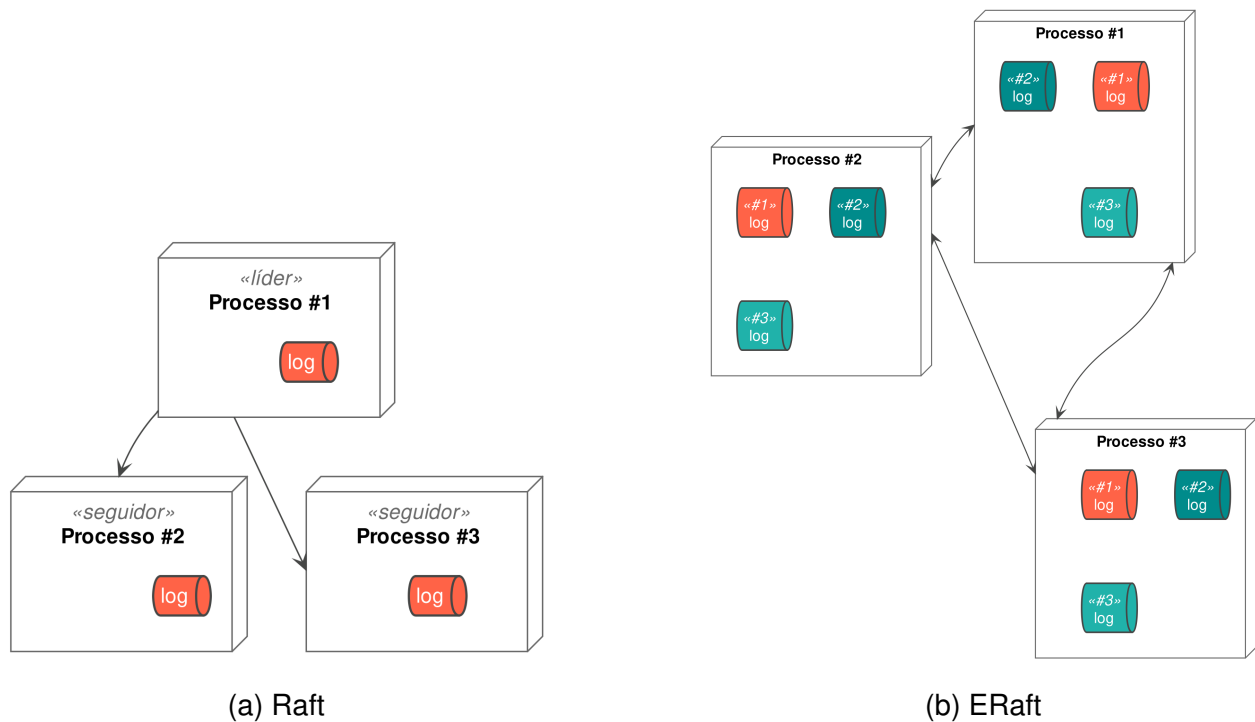


Figura 3.1: Modelo de replicação do *log*. À esquerda, Raft com um único *log* cujo líder é responsável pela replicação para os seguidores. À direita, o algoritmo EReplicação, em que cada processo replica o seu *log* para os demais.

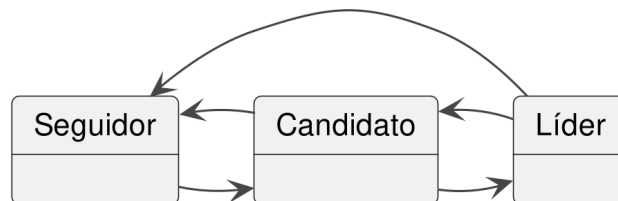


Figura 3.2: Diagrama de estados de processos no algoritmo EReplicação.

Em uma aplicação na qual os comandos não conflitam, o *log* replicado de cada processo pode ser executado em qualquer ordem. No entanto, como comandos conflitantes podem ser submetidos a qualquer processo, é necessário que haja um mecanismo de sincronização entre os mesmos para estabelecer uma ordem parcial de execução que garanta a ordem entre comandos conflitantes e assim a consistência dos dados da aplicação. Para isso, EReplicação adapta os protocolos de troca de mensagens de Raft para a coleta e replicação de conflitos — também chamados dependências. O consenso garante que um comando decidido carregará as mesmas dependências em todas as réplicas. A fase de execução deve garantir que um comando só execute se suas dependências foram satisfeitas (já executadas).

Da mesma forma que Atlas, comandos no *log* são classificados nas fases: “rápida”, “lenta”, e “concluída” (fig. 3.3). Todo comando recebido por um líder tem seu conjunto inicial de dependências calculado e adicionado ao *log* na fase **rápida**. Comandos nesta fase

são enviados a um subgrupo de processos denominado quórum rápido, os quais calculam sua contribuição para o conjunto de dependências e as retornam ao líder. Ao receber um número suficiente de respostas, o líder decide com base nas dependências recebidas se é seguro promover o comando para fase **concluída**, comunicando assim seus pares sobre esta decisão. Caso não seja seguro concluir o consenso, o comando tem seu conjunto de dependências atualizado e promovido para a fase **lenta**, na qual mais uma rodada de trocas de mensagens ocorre antes de promovê-lo para a fase **concluída**.

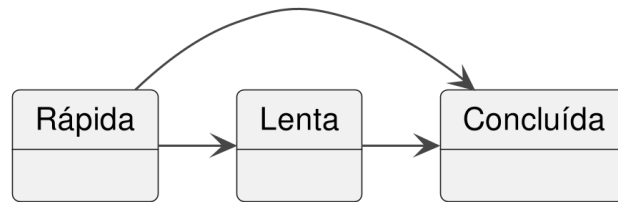


Figura 3.3: Diagrama de transição de fases de um comando em ERep.

Eventualmente cada processo passa a conter uma cópia do *log* de cada um de seus pares. Comandos na fase concluída podem ser executados a qualquer momento desde que executados após suas dependências. Na prática, isso gera um grafo cíclico dirigido, o qual cada processo deve navegar de forma determinística (fig. 3.4). Comandos executados são, por fim, removidos do *log*.

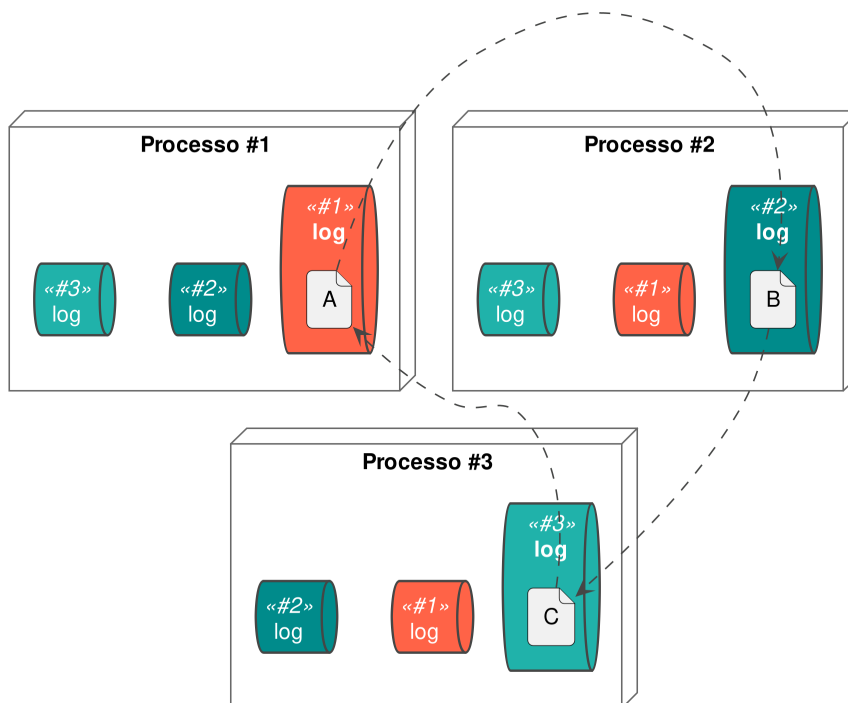


Figura 3.4: Exemplo de ciclos de dependências: comando “A” coordenado pelo processo 1 depende do comando “B” coordenado pelo processo 2 que, por sua vez, depende do comando “C” coordenado pelo processo 3, o qual volta a depender do comando A”.

A recuperação de falhas em algoritmos de consenso é um fator de notável complexidade. Em algoritmos igualitários, é possível que um processo tenha informado seus pares quanto a dependências para seus próprios comandos, os quais não são promovidos para a fase concluída antes de sua falha. Neste caso, execução dos comandos dependentes nos demais processos é comprometida.

Se, ao término de um ciclo de *timeout* um processo não observar mensagens de um determinado par, ele suspeita da falha do mesmo. Caso o mesmo tenha uma dependência não concluída para o par suspeito, ele passa a ser candidato a uma eleição para o *log* deste par. Durante o processo de eleição, o candidato recebe dos demais processos os comandos conhecidos onde o suspeito é líder. Se eleito, o processo conclui o consenso dos comandos recebidos, o que permite que comandos dependentes possam ser executados.

As secções 3.3.2 e 3.3.3 oferecem uma visão em alto nível dos algoritmos de replicação de comandos e recuperação de falhas. Note que os pseudocódigos presentes neste capítulo são um subconjunto mínimo do comportamento do algoritmo, oferecendo uma noção geral de seu funcionamento. A definição precisa dos conceitos aqui introduzidos se encontram na especificação formal no apêndice A.

3.3.2 Replicação de Comandos

ERaft é configurado com base em 3 constantes, sendo elas: um conjunto finito e não vazio de processos conhecidos; uma função booliana e comutativa de conflito definida pela aplicação; e o número f de falhas toleradas pelo sistema (apêndice. A, linhas 113–176). Ainda, de acordo com Atlas, o número de falhas f deve obedecer à relação $1 \leq f \leq \lfloor \frac{n-1}{2} \rfloor$ sendo n o número de processos.

Cada processo possui o estado inicial descrito no algoritmo 3.1. De forma análoga ao estado inicial de Raft, o termo atual para cada processo inicia em zero indicando que não houve eleições de líderes até o momento. Todo processo começa no estado de seguidor e com o *log* vazio. Os índices de replicação de comandos em fase rápida, lenta, e concluída iniciam em zero e o índice do próximo comando a ser enviado para seus pares inicia em 1, dado que não houve comandos submetidos até o momento. Note que, por haver uma sobreposição de estados, todo participante do algoritmo replica o *log* de todos os processos disponíveis e a mesma rotina de inicialização é invocada para cada processo conhecido.

Todo participante do algoritmo escolhe um quórum de processos para coordenar o consenso de comandos na fase rápida. Ao contrário do Atlas, que escolhe um quórum por comando, ERaft escolhe um quórum por processo, simplificando seu algoritmo de recuperação de falhas. Além disso, em conformidade com o Atlas, sua cardinalidade é de $\lfloor \frac{n}{2} \rfloor + f$. A definição formal da escolha de quóruns se encontra no apêndice A (linhas 416–423), sendo equivalente à equação 3.1.

Require: *Servers* ▷ Processos conhecidos

procedure INIT(*p*)

p.currTerm ← 0 ▷ Termo atual

p.currState ← *follower* ▷ Estado atual

p.log ← $\langle \rangle$ ▷ Sequência ordenada de comandos

p.commitIdx ← 0 ▷ Prefixo de comandos concluídos

p.fastMatchIdx ← $[s \mapsto 0 \mid s \in \text{Servers}]$ ▷ Índice de comandos em fase rápida

p.slowMatchIdx ← $[s \mapsto 0 \mid s \in \text{Servers}]$ ▷ Índice de comandos em fase lenta

p.nextIdx ← $[s \mapsto 1 \mid s \in \text{Servers}]$ ▷ Próximo comando a ser enviado a seguidores

p.fastQuorum ← *FastQuorum*(*p*, *Servers*) ▷ Quórum usado para fase rápida

p.votedFor ← $[]$ ▷ Votos concedidos durante eleições

p.recovery ← $[]$ ▷ Estado da recuperação de falhas

end procedure

Algoritmo 3.1: Algoritmo que define o estado inicial de um processo *p*.

$$\text{FastQuorum}(p, \text{Servers}) = Q \subseteq \text{Servers} : |Q| = \lfloor \frac{n}{2} \rfloor + f \wedge p \in Q \quad (3.1)$$

Uma vez definido seu estado inicial, o processo inicia um ciclo fechado de tempo (*timeout*) por processo conhecido. O EReplic é desenhado para funcionar em uma rotina de tratamento de eventos (*event loop*) no qual apenas um evento é processado por vez, o que garante a atomicidade das ações descritas na especificação formal. Os dois principais tipos de eventos são *AppendToLog(cmd)* e *Timeout(p)*. Os demais se referem ao tratamento de mensagens de rede e serão discutidos posteriormente.

Ao receber um novo comando de um cliente, um evento do tipo *AppendToLog(cmd)* é disparado. Ao processá-lo, o líder calcula suas dependências e o adiciona ao seu *log* na fase rápida (algoritmo 3.2, e apêndice A nas linhas 713–729).

procedure APPENDTOLOG(*cmd*)

Deps ← *Dependencies*(*cmd*) ▷ Calcula dependências

cmd ← (*cmd* ↦ *cmd*, *deps* ↦ *Deps*, *phase* ↦ *fast*) ▷ Cria um novo registro

leader.log ← *Append*(*leader.log*, *cmd*) ▷ Adiciona no log

end procedure

Algoritmo 3.2: Algoritmo que adiciona um novo comando ao *log* do líder.

A função de cálculo de dependências é descrita pela equação 3.2 — formalmente definida no apêndice A, linhas 607–627. Este algoritmo utiliza a função de conflito definida pela aplicação para identificar dependências entre o comando proposto e os demais comandos conhecidos.

$$Dependencies(c) = \bigcup \{c' \in leader.log, leader \in Servers \mid Conflict(c, c')\} \quad (3.2)$$

Quando um ciclo de tempo é concluído para um determinado processo, um evento do tipo *Timeout(p)* é gerado, o que dispara a rotina descrita pelo algoritmo 3.3. Nela, caso o processo seja líder, ele verifica se há comandos prontos para mudar de fase antes de disparar a rotina de envio de mensagens para seus pares. Caso contrário, ocorre a suspeita de falha do processo. Em ambos os casos, um novo ciclo de tempo se inicia.

```

procedure ONTIMEOUT(p)
  if p.currState = leader then
    CommitPending(p)
    ReplicateCommands(p)
  else
    SuspectFailure(p)
  end if
  ScheduleTimeout(p)
end procedure

```

- ▷ É líder do consenso
- ▷ Avança comandos pendentes
- ▷ Envia mensagens aos seguidores
- ▷ Suspeita de falha do processo
- ▷ Inicia próximo ciclo de tempo

Algoritmo 3.3: Algoritmo acionado para a um evento do tipo *timeout*.

O algoritmo 3.4 descreve a rotina responsável por enviar comandos de um líder para seus seguidores — modelada formalmente no apêndice A, linhas 749–787. Assim como em Raft, um índice de replicação é utilizado para determinar quais comandos serão enviados para cada seguidor. No entanto, em conformidade com o Atlas, comandos em fase rápida são enviados apenas para seguidores que fazem parte do quórum rápido do líder. Mesmo que o conjunto de comandos a ser enviado seja vazio, o *RPC* de replicação ainda é invocado para que, assim como em Raft, seguidores possam reiniciar o *timeout* do líder, o que impede uma suspeita de falha por inatividade.

```

procedure REPLICATECOMMANDS(leader)
  for all follower ∈ Servers \ {leader} do
    C ← {cmdi ∈ leader.log | i ≥ leader.nextIdx[follower]}
    P ← C' ⊆ C | ∀ c' : c'.phase = fast ⇒ follower ∈ leader.fastQuorum
    AppendRPC(follower, leader.commtIdx, leader.nextIdx[follower] − 1, P)
  end for
end procedure

```

- ▷ Para todo seguidor
- ▷ Comandos a enviar
- ▷ Prefixo de *C*

Algoritmo 3.4: Algoritmo que envia comandos do líder para os seguidores.

Cada *RPC* em Eraft é modelado por um evento de requisição e um de resposta. O algoritmo 3.5 descreve a rotina invocada pelo seguidor ao receber uma sequência de comandos do líder — modelada formalmente no apêndice A (linhas 789–912). Nela, assim como no algoritmo Atlas, o seguidor contribui para o consenso recalculando o conjunto de dependências de comandos na fase rápida à medida que os comandos são adicionados

ao *log* do seguidor. Comandos em fase lenta ou concluído são aceitos desde que não sobrescrevam comandos já concluídos e que o seguidor tenha conhecimento.

```

procedure ONAPPENDREQUEST(leader, follower, leaderCommit, prevLogIndex, Cmds)
  follower.commitIdx  $\leftarrow$  leaderCommit ▷ Aprende índice do líder

  for all  $cmd_i \in Cmds, cmd' \in follower.log_{prevLogIndex+i}$  do
    if  $cmd' = \emptyset$  then ▷ Não há comando neste posição
      if  $cmd.phase = fast$  then ▷ Comando recebido em fase rápida
         $Deps \leftarrow Dependencies(cmd)$  ▷ Recalcula dependências
         $cmd.deps \leftarrow cmd.deps \cup Deps$  ▷ Atualiza comando
      end if
       $follower.log \leftarrow Append(follower.log, cmd)$  ▷ Aceita comando
    else ▷ Comando sobrescreve existente
      if  $cmd'.phase \neq committed$  then ▷ Comando existente não concluído
         $follower.log \leftarrow Replace(follower.log, cmd', cmd)$  ▷ Aceita comando
      end if
    end if
  end for ▷ Devolve dependências para o líder

   $Deps \leftarrow \{cmd_i.deps \in follower.log \mid prevLogIndex < i \leq |Cmds|\}$ 
   $Reply(leader, prevLogIndex, Deps)$ 
end procedure

```

Algoritmo 3.5: Algoritmo utilizado pelo seguidor ao receber um comando do líder.

Uma vez recebida a resposta de um seguidor, o líder utiliza a rotina descrita pelo algoritmo 3.6 para avançar seu estado. Nesta etapa, comandos na fase rápida acumulam as dependências reportadas pelo seguidor. Note que, assim como em Atlas, o objetivo da fase rápida é tanto replicar comandos quanto coletar dependências, enquanto a fase lenta tem o objetivo de comunicar aos seguidores as dependências consolidadas pelo líder. Portanto, este algoritmo avança o índice rápido para todo comando observado nesta fase, mas por outro lado, somente avança o índice lento se observados comandos nesta fase que possuam o mesmo conjunto de dependências.

```

procedure ONAPPENDREPLY(leader, follower, prevLogIndex, Deps)
  for all  $D_i \in Deps, cmd_{prevLogIndex+i} \in leader.log$  do
    if  $cmd.phase = fast$  then ▷ Comando na fase rápida
       $cmd.deps \leftarrow cmd.deps \cup D$  ▷ Acumula dependências
       $leader.fastMatchIdx[follower] \leftarrow prevLogIndex + i$  ▷ Atualiza o índice rápido
    else if  $cmd.phase = slow \wedge cmd.deps = D$  then ▷ Mesmas dependências
       $leader.slowMatchIdx[follower] \leftarrow prevLogIndex + i$  ▷ Atualiza índice lento
    end if
  end for
end procedure

```

Algoritmo 3.6: Algoritmo utilizado pelo líder ao receber uma resposta do seguidor.

Para cada resposta recebida de um seguidor, o índice do próximo comando a ser enviado a ele é atualizado de acordo com seu estado. Caso o líder observe um comando em fase lenta onde o seguidor reportou um conjunto diferente de dependências, é utilizado o índice deste comando como o próximo a ser enviado a ele de modo a corrigir esta divergência. Caso o líder observe que todos os comandos processados estão em fase rápida, ele define o próximo comando como sendo $leader.fastMatchIdx[follower] + 1$, sinalizando que líder e seguidor estão sincronizados até este ponto (apêndice A, linhas 914–1064).

Em resposta a um evento do tipo *timeout*, líderes invocam a rotina descrita pelo algoritmo 3.7 para avaliar seu progresso. Comandos trocam de fase ao serem replicados para um número suficiente de seguidores conforme sua fase (equação 3.3).

```

procedure COMMITPENDING(leader)
  for all  $cmd_i \in leader.log : leader.commitIdx < i \leq |leader.log|$  do
    if  $cmd_i.phase \in \{slow, fast\} \wedge Replicated(cmd_i)$  then      ▷ Comandos replicados
      if  $cmd_i.phase = slow \vee SkipSlow(cmd_i)$  then          ▷ Fase lenta ou podem pulá-la
         $cmd_i.phase \leftarrow committed$                        ▷ Promove para fase concluída
      else if  $cmd_i.phase = fast$  then                          ▷ Não pode pular fase lenta
         $cmd_i.phase \leftarrow slow$                              ▷ Promove para fase lenta
      end if
    end if
  end for
end procedure

```

Algoritmo 3.7: Algoritmo utilizado pelo líder avalia o progresso do consenso.

$$Replicated(c_i) = |Quorum(c_i)| \geq \begin{cases} \lfloor \frac{n}{2} \rfloor + f & \text{if } c_i.phase = fast \\ f + 1 & \text{if } c_i.phase = slow \end{cases} \quad (3.3)$$

$$Quorum(c_i) = \left\{ s \in Servers \mid i \leq \begin{cases} fastMatchIdx[s] & \text{if } c_i.phase = fast \\ slowMatchIdx[s] & \text{if } c_i.phase = slow \end{cases} \right\} \quad (3.4)$$

De acordo com o Atlas, comandos em fase rápida são replicados para todo o quórum rápido de processos, composto por $\lfloor \frac{n}{2} \rfloor + f$ pares, antes de trocar de fase. Estes comandos são concluídos se todas as suas dependências foram reportadas por ao menos f processos (equação 3.5). Caso contrário, é necessária mais uma rodada de consenso para informar aos seguidores sobre as dependências consolidadas pelo líder. Neste caso, estes comandos são promovidos para a fase lenta. Por outro lado, comandos na fase lenta podem ser considerados concluídos após replicados por pelo menos $f + 1$ processos.

$$\text{SkipSlow}(c) = \forall d \in c.\text{deps} : |\{d' \in c.\text{deps} : \text{SameDep}(d, d')\}| \geq f \quad (3.5)$$

$$\text{SameDep}(d, d') = d.\text{leader} = d'.\text{leader} \wedge d.\text{pos} = d'.\text{pos} \quad (3.6)$$

Note que, ao trocar da fase rápida para lenta ou concluída, um comando pode ter um conjunto de dependências no líder diferente do conjunto replicado para seus seguidores. Neste caso, após processar as mudanças de fase, o líder ajusta o índice do próximo comando a ser enviado para cada seguidor, repetindo comandos promovidos da fase rápida. Para mais detalhes, consulte o apêndice A, linhas 1066–1185.

3.3.3 Recuperação de Falhas

ERaft segue o mesmo princípio de recuperação de falhas que Raft no qual um *RPC* é utilizado para eleger um novo líder do consenso. O mesmo mecanismo é utilizado tanto para eleger um processo como líder de seu próprio *log*, quanto para se eleger líder do log de um par do qual ele deseja recuperar. Ainda, eleições podem ocorrer quando um líder suspeita da falha de um seguidor que pertence ao seu quórum rápido. Neste caso, o líder escolhe um novo quórum no qual o processo suspeito não participa e inicia uma eleição para si mesmo, concluindo seus comandos pendentes antes de utilizar o quórum escolhido — veja apêndice A, linhas 1191–1231.

Toda eleição inicia com o incremento do termo para o processo suspeito de falha e a transição do estado do processo atual para candidato (apêndice A, linhas 651–697). ERaft aplica as mesmas restrições em seus *RPCs* que Raft para garantir que a propagação de novo termo entre líderes e seguidores: o *RPC* de replicação recusa mensagens de termos anteriores, enquanto o *RPC* de eleição se recusa a votar em processos de termos passados. Ainda, qualquer processo volta ao estado de seguidor se receber mensagens de um processo com termo maior que o dele.

Para ser capaz de concluir o consenso dos comandos de um processo suspeito, ERaft modifica o *RPC* de eleição de Raft de modo que além de pedir votos ele também coleta comandos pendentes (apêndice A, linhas 1233–1256 e 1319–1347). Diferente de Raft onde cada processo apenas responde a uma solicitação de voto por termo, em ERaft cada processo vota apenas uma vez por termo, mas toda solicitação é respondida. Esta mudança visa garantir que mesmo que um processo não vote em um candidato, caso ele vença a eleição por maioria simples de votos, o candidato ainda considere os comandos coletados dos pares que não votaram nele.

Um candidato é considerado líder se receber votos de ao menos $\lfloor \frac{n}{2} \rfloor + 1$ processos. No entanto, para garantir que ele observou um conjunto suficiente de comandos para

concluir o consenso, o Atlas exige que $n - f$ respostas sejam coletadas de seus pares. O algoritmo 3.8 descreve como o EReplica consolida os comandos coletados durante a eleição — especificado formalmente no apêndice A, linhas 1349–1513.

```

procedure MERGELOGS(Logs, i)
  if i = 0 then
    return  $\langle \rangle$ 
  else
     $C \leftarrow \{c_i \in \text{Log}, \text{Log} \in \text{Logs}\}$  ▷ Todos os comandos na posição i

     $\text{cmd} \leftarrow c \in C : c.\text{phase} = \text{committed}$  ▷ Escolhe um comando na fase concluída
    if  $\text{cmd} = \emptyset$  then ▷ Se nenhum em fase concluída
       $\text{cmd} \leftarrow c \in C : c.\text{phase} = \text{slow}$  ▷ Escolhe um comando na fase lenta
    end if

    if  $\text{cmd} = \emptyset$  then ▷ Se nenhum em fase lenta
      if  $|C| \geq \lfloor \frac{n}{2} \rfloor$  then ▷ Decide se pode escolher um na fase rápida
         $\text{cmd} \leftarrow c \in C$  ▷ Escolhe um comando na fase rápida
         $\text{cmd}.\text{phase} \leftarrow \text{slow}$  ▷ Promove comando para fase lenta
         $\text{cmd}.\text{deps} \leftarrow \bigcup \{c'.\text{deps} \mid c' \in C\}$  ▷ Consolida dependências
      else
         $\text{cmd} \leftarrow (\text{cmd} \mapsto \emptyset, \text{deps} \mapsto \{\}, \text{phase} \mapsto \text{slow})$  ▷ Anula esta posição
      end if
    end if

    return  $\text{Append}(\text{MergeLogs}(\text{logs}, i - 1), \text{cmd})$  ▷ Consolida restante do log
  end if
end procedure

```

Algoritmo 3.8: Algoritmo utilizado para consolidar logs coletados durante uma eleição.

Caso o candidato observe um comando concluído ou na fase lenta para uma determinada posição do *log*, ele escolhe este comando para concluir o consenso dado que nestas fases seu conjunto de dependências já foi consolidado pelo líder e deve ser preservado. Por outro lado, se apenas comandos na fase rápida forem observados, não é possível determinar se o líder concluiu o consenso deste comando antes de falhar. Neste caso, Atlas prova que a união das dependências reportadas por $\lfloor \frac{n}{2} \rfloor$ processos recriam o mesmo conjunto que o líder observou, estratégia esta também adotada por EReplica. Em último caso, aquela posição do *log* recebe um comando nulo para desbloquear a execução de comandos que dependam dela.

Uma vez consolidados os comandos pendentes para o processo suspeito, o candidato atualiza o *log* do mesmo, passa para o estado de líder e utiliza o *RPC* de replicação discutido na sessão 3.3.2 para comunicar aos demais processos as decisões tomadas por ele. O processo que falhou, ao se recuperar, utiliza o mesmo mecanismo de eleição para tomar conhecimento das decisões feitas durante sua ausência. Para facilitar o retorno do

processo falho, o líder eleito para a recuperação do mesmo abdica da liderança após concluir o consenso do *log* consolidado (apêndice A, linhas 1515–1541).

3.4 Invariantes e Propriedades

De maneira informal é derivado que por utilizar o mesmo mecanismo e quórum de eleições que Raft, ambos provados corretos, as eleições de líderes em ERaft também são corretas. Durante a replicação de comandos e recuperação de falhas, o ERaft utiliza a mesma estratégia e quórums que Atlas, o que o torna também correto nestes quesitos uma vez que Atlas é provado correto.

Para mais detalhes veja a especificação formal do algoritmo que conta com uma seção de invariantes e propriedades verificadas durante a exploração do espaço de estados do mesmo (apêndice A, linhas 1543–1781).

3.5 Considerações Práticas

Uma vez concluído o consenso para uma sequência de comandos, líderes comunicam seu progresso via *RPC* de replicação no qual seu índice de *commit* é propagado. Uma vez observados comandos em posição abaixo deste índice, seguidores podem executá-los em sua máquina de estados local desde que respeitando a restrição de ordem imposta pelo conjunto de dependências de cada comando. Na prática, isso gera um grafo dirigido com a possibilidade de ciclos que deve ser navegado de forma determinística por todos os processos. Métodos tradicionais de navegação de grafos, incluindo os métodos utilizados por EPaxos [12] e Atlas [3] podem ser aplicados a ERaft para implementar a execução.

Para garantir que processos não suspeitem da falha de líderes durante o comportamento normal do sistema, ERaft presume o mesmo comportamento que Raft onde líderes utilizam um intervalo de tempo menor que seguidores ao agendar seus *timeouts* e seguidores reiniciam o *timeout* referente ao líder para toda mensagem recebida dele. Também é esperado que *timeouts* sejam agendados com intervalos aleatórios de tempo para diminuir as chances de múltiplos processos suspeitarem da falha de um par simultaneamente, o que pode resultar em uma competição pela liderança do processo suspeito.

O custo da busca por dependências de um determinado comando cresce linearmente com o tamanho dos *logs* mantidos pelos processos. É possível limitar este custo desde que o limite inferior de busca seja adicionado ao conjunto de dependências do comando. Por exemplo, em um processo onde há um *log* com tamanho 100, se a busca por dependências é configurada para analisar apenas os 10 comandos mais recentes, o índice 89 do *log* deve ser considerado como dependência do comando avaliado. Este limite deve

ser configurado com base na velocidade média de execução de comandos durante o funcionamento normal do sistema para que, na maior parte dos casos, quando o comando concluir seu consenso o índice 89 do qual ele depende já foi executado pelos processos.

É documentada na literatura uma condição em que algoritmos de consenso igualitários podem sofrer de postergação indefinida [5]. Basta que um comando fique pendente e que o processo continue a aceitar novos comandos que dependem dele. Para mitigar este efeito, implementações destes algoritmos devem preferir concluir o consenso de comandos pendentes ao invés de aceitar novos comandos. Esta condição pode ser implementada em ERaft como uma restrição baseada na diferença entre o índice de *commit* e o tamanho do *log* onde o processo passa a rejeitar novos comandos se a diferença atingir um valor limite.

Dado que comandos são promovidos após serem replicados e existem índices de replicação para cada fase, é possível utilizar estes índices para determinar se um seguidor já tem conhecimento de um comando. Então, pode-se reduzir tráfego de rede reenviando apenas as dependências de um comando se promovido da fase rápida para lenta ou concluída.

4. CONCLUSÃO

Neste trabalho de conclusão de curso foi introduzido o algoritmo de consenso igualitário ERaft, inspirado pelos algoritmos Raft e Atlas. Assim como o Raft, ele possui apenas dois *RPCs* de comunicação entre processos visando diminuir a complexidade do algoritmo. Ainda em conformidade com Atlas, todo processo é líder de seus próprios comandos, o que tende a aumentar a vazão do sistema através do consenso e execução de comandos em paralelo, exceto na presença de conflito. A maior contribuição deste trabalho é sua especificação formal escrita em TLA⁺, disponível no apêndice A, que contém a definição precisa do algoritmo proposto.

Trabalhos futuros relacionados a esta pesquisa incluem sua prova via um assistente de provas como Coq¹ ou Isabelle² de forma a aumentar a confiança em sua correção. Além disso, a construção de um simulador visual aos moldes de RaftScope³ pode servir como ferramenta didática e auxiliar na promoção deste algoritmo. Por fim, a construção de uma implementação concreta de ERaft e sua avaliação experimental de desempenho em comparação ao Raft clássico e similares pode evidenciar possíveis otimizações e pontos de melhoria.

¹<https://coq.inria.fr>

²<https://isabelle.in.tum.de>

³<https://raft.github.io>

REFERÊNCIAS BIBLIOGRÁFICAS

- [1] Barcelona, C.-S. “Mencius: building efficient replicated state machines for wans”. In: 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI 08), 2008.
- [2] Dwork, C.; Lynch, N.; Stockmeyer, L. “Consensus in the presence of partial synchrony”, *J. ACM*, vol. 35–2, apr 1988, pp. 288–323.
- [3] Enes, V.; Baquero, C.; Rezende, T. F.; Gotsman, A.; Perrin, M.; Sutra, P. “State-machine replication for planet-scale systems”. In: Proceedings of the Fifteenth European Conference on Computer Systems, 2020, pp. 1–15.
- [4] Fischer, M. J.; Lynch, N. A.; Paterson, M. S. “Impossibility of distributed consensus with one faulty processor”, *Journal of the ACM*, vol. 32–2, 1985, pp. 374–382.
- [5] França Rezende, T.; Sutra, P. “Leaderless state-machine replication: Specification, properties, limits”. In: 34th International Symposium on Distributed Computing (DISC 2020), 2020.
- [6] Howard, H. “Distributed consensus revised”, Tese de Doutorado, University of Cambridge, 2019.
- [7] Lamport, L. “The part-time parliament”, *ACM Trans. Comput. Syst.*, vol. 16–2, may 1998, pp. 133–169.
- [8] Lamport, L. “Paxos made simple”, *ACM SIGACT News (Distributed Computing Column)* 32, 4 (Whole Number 121, December 2001), December 2001, pp. 51–58.
- [9] Lamport, L. “Generalized consensus and paxos”, Relatório Técnico MSR-TR-2005-33, 2005.
- [10] Lamport, L. “Fast paxos”, *Distributed Computing*, vol. 19, October 2006, pp. 79–103.
- [11] Lamport, L.; Massa, M. “Cheap paxos”. In: International Conference on Dependable Systems and Networks, 2004, 2004, pp. 307–314.
- [12] Moraru, I.; Andersen, D. G.; Kaminsky, M. “There is more consensus in egalitarian parliaments”. In: Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles, 2013, pp. 358–372.
- [13] Ongaro, D.; Ousterhout, J. “In search of an understandable consensus algorithm”. In: 2014 {USENIX} Annual Technical Conference ({USENIX}{ATC} 14), 2014, pp. 305–319.

- [14] Pedone, F.; Schiper, A. “Handling Message Semantics with Generic Broadcast Protocols”, *Distrib. Comput.*, vol. 15–2, 2002, pp. 97–107, generic broadcast.
- [15] Rezende, T. F.; Sutra, P.; Saramago, R. Q.; Camargos, L. “On making generalized paxos practical”. In: 2017 IEEE 31st International Conference on Advanced Information Networking and Applications (AINA), 2017, pp. 347–354.
- [16] Van Renesse, R.; Altinbuken, D. “Paxos made moderately complex”, *ACM Comput. Surv.*, vol. 47–3, feb 2015.

APÊNDICE A – ESPECIFICAÇÃO EM TLA⁺

1 ————— MODULE *ERaft* —————

Egalitarian Raft (*ERaft*) is an egalitarian – also called leaderless – version of the well known Raft consensus algorithm. *ERaft* heavily draws inspiration from the Atlas consensus protocol to achieve optimal commit latency in the presence of conflicting commands.

The following section offers a high level overview of the *ERaft* algorithm, whereas the specification itself provides a precise definition of each component mentioned.

— HIGH LEVEL OVERVIEW —

ERaft works similarly to Raft: every server process starts at the “follower” state; followers can become a “candidate” due to inactivity (timeout); a round of election happens to which only one candidate wins by a majority vote – split votes are disregarded and a new election is called. The winning candidate transitions to the “leader” state and begins managing the *log* of commands submitted by clients. *ERaft* slightly modifies this model by introducing another dimension: each server process has its own *log*.

A server can be simultaneously at the “leader” state for its own *log* while being at the “follower” state for every other server process in the system during normal operation. For example:

Server 1:

(leader) Server 1 → [a, b, c]

(follower) Server 2 → [d]

(follower) Server 3 → [e, f]

Server 2:

(follower) Server 1 → [a, b, c]

(leader) Server 2 → [d]

(follower) Server 3 → [e, f]

Server 3:

(follower) Server 1 → [a, b, c]

(follower) Server 2 → [d]

(leader) Server 3 → [e, f]

Commands are submitted by clients to one of the server processes available. Upon receiving a command, the server will compute the set of known dependencies for the command before appending it to its own *log*. The commands’ set of dependencies are used to determine the *log*’s execution order in which commands can only be executed after their dependencies so that logical data consistency is preserved. *ERaft* detect dependencies via an abstract conflict function provided by the application.

As per Raft, a single *RPC* – called “append” – is used to replicate commands between leaders and followers. *ERaft* extends this *RPC* to exchange dependencies discovered during command replication. As per Atlas, commands can be in different phases. *ERaft* establishes that every new command appended to the leader’s *log* starts at the “fast” phase and is sent via broadcast to a quorum of server processes which reply with their own set of computed dependencies for it. If the leader determines that the command can be safely committed based on the reported dependencies, the command is moved to the “committed” phase. Otherwise, commands are promoted to the “slow” phase and another call to the append *RPC* is required before committing it.

Committed commands are immutable and can be executed by the application which must respect the ordering constraints established by their dependency set. In practice, this results in a cyclic graph that servers have to traverse and solve cycles deterministically.

Upon suspecting of failures, servers can become candidates for some other server's *log*. The election process derived from Raft – here called “recovery” *RPC* – combines election with command discovery giving the elected recovery leader the necessary view of the failing server's *log* so that it can complete the pending commands or replace them with *NULL* if necessary. A recovery is also triggered when a leader needs to reconfigure its quorums when detecting failed member.

Recovery leaders, once elected, proceed with standard append *RPC* calls in order to replicate and commit commands discovered during its election. Recovery leaders never append new commands to the failed

leader's *log*. Once all known commands are committed by the recovery leader, it abdicates its leadership by going back to the “follower” state.

This specification is meant as a comprehensive model of both command replication and failure recovery portion of the algorithm while committed commands execution and *log* truncation are left to a separate specification.

— REFERENCES —

- Ongaro, D. and Ousterhout, J. (2014). In search of an understandable consensus algorithm. In 2014 *{USENIX}* Annual Technical Conference (*{USENIX}**{ATC}* 14), pages 305–319.
- Enes, V., Baquero, C., Rezende, T. F., Gotsman, A., Perrin, M., and Sutra, P. (2020). State-machine replication for planet-scale systems. In Proceedings of the Fifteenth *European Conference on Computer Systems*, pages 1–15.

101 EXTENDS *FiniteSets, Naturals, Sequences, TLC, Utils*

103

A finite set of interconnected server processes.

E.g.: {*s*1, *s*2, *s*3}

109 CONSTANT *Servers*

A finite set of commands to be proposed by servers.

E.g.: {*a*, *b*, *c*}

116 CONSTANT *Commands*

A commutative function that takes 2 commands and return TRUE iff they conflict with each other. Returns FALSE otherwise.

E.g.: LAMBDA *x, y* : *x* = *y*

124 CONSTANT *Conflict*(-, -)

A constant representing the absence of a value.

129 CONSTANT *NULL*

A natural number indicating the amount of tolerated failed processes (as per Atlas).

135 CONSTANT *F*

The number of server processes in the system.

140 $N \triangleq \text{Cardinality}(\text{Servers})$

The size of a process fast quorum (as per Atlas).

$$145 \quad FQ \triangleq (N \div 2) + F$$

The size of a process slow quorum (as per Atlas).

$$150 \quad SQ \triangleq F + 1$$

The size of a recovery quorum (as per Atlas).

$$155 \quad RQ \triangleq N - F$$

The size of an election quorum (as per Raft).

$$160 \quad EQ \triangleq (N \div 2) + 1$$

162 ASSUME

Assume *Servers* and *Commands* to be non-empty finite sets, *NULL* values are non-conflicting, the conflict function to be commutative, and *F* to be properly configured (as per Atlas).

$$\begin{aligned}
 168 \quad & \wedge \text{IsFiniteSet}(\text{Servers}) \\
 169 \quad & \wedge \text{Cardinality}(\text{Servers}) > 0 \\
 170 \quad & \wedge \text{IsFiniteSet}(\text{Commands}) \\
 171 \quad & \wedge \text{Cardinality}(\text{Commands}) > 0 \\
 172 \quad & \wedge \forall \text{cmd} \in \text{Commands} : \\
 173 \quad & \quad \text{Conflict}(\text{cmd}, \text{NULL}) = \text{FALSE} \\
 174 \quad & \wedge \forall a, b \in \text{Commands} \cup \{\text{NULL}\} : \\
 175 \quad & \quad \text{Conflict}(a, b) = \text{Conflict}(b, a) \\
 176 \quad & \wedge 1 \leq F \wedge F \leq (N - 1) \div 2
 \end{aligned}$$

178 |

— STATE VARIABLES —

This section describes the state variables used throughout the spec. Notably, the spec state is multi-dimensional. Server variables are defined as a function from “ $s1 \rightarrow s2 \rightarrow \text{value}$ ” where “ $s1$ ” defines the server process in which the state variable resides, and “ $s2$ ” represents a server from which “ $s1$ ” has knowledge about (possibly the same as “ $s1$ ”). For example, a possible slice of the *log* variable is:

```
(s1:> (s1:> ([cmd : a, ...], ...))
  @@ (s2:> ([cmd : b, ...], ...)))
```

In the example above, both sequences of commands resides in the server $s1$. The server $s1$ has command “ a ” appended to its own *log*. The server $s2$ has replicated command “ b ” to server “ $s1$ ”.

Variables used during recovery contain a third dimension as they are used to track both the state of the leader processes as well as the state of the recovering process. For example, take a possible slice of the *nextIdx* variable:

```
(s1:> (s1:> (s1:> 2)
  @@ (s2:> 1)
  @@ (s3:> 1))
  @@ (s2:> (s1:> 3)
  @@ (s2:> 3)
  @@ (s3:> 1))
  @@ ..)
```

Again, the first dimension defines the server process where the state resides. In the example above, $s1$ has its own $nextIdx$ set to 1 for both servers $s2$ and $s3$. Moreover, while coordinating the recovery the process $s2$, it has $s2$'s $nextIdx$ set to 1 for server $s3$.

A function from a server process to a function from leader to its current term – defined as a natural number. Similarly to Raft, term act as a synchronization mechanism when a leader is suspected to have failed and another server process attempts to recover its log .

221 VARIABLE $currTerm$

223 $CurrTermOK \triangleq$

For every known server process, its term is a natural number.

227 $\forall server, leader \in Servers : currTerm[server][leader] \in Nat$

229 $InitCurrTerm \triangleq$

All servers start at the same term 0.

233 $currTerm = [server \in Servers \mapsto [leader \in Servers \mapsto 0]]$

A function from server to its current state.

238 VARIABLE $currState$

240 $CurrStateOK \triangleq$

Current state transitions are similar to Raft's:

- Servers start as followers;
- Followers transition to candidates;
- Candidates transition to follower or leader;
- Leaders transition to follower or candidate.

252 $\forall server, leader \in Servers :$

253 $currState[server][leader] \in \{ \text{"follower"}, \text{"candidate"}, \text{"leader"} \}$

255 $InitCurrState \triangleq$

Every leader starts as a follower. It must call for a recovery – analogous to a Raft election – to claim ownership of its log .

260 $currState = [server \in Servers \mapsto [leader \in Servers \mapsto \text{"follower"}]]$

A function from a server to the leader's ordered sequence of commands agreed upon consensus.

266 VARIABLE log

268 $Dependencies \triangleq$

Dependencies is a set of records containing the reporting peer, the command's leader, and its position in the leader's log .

273 $[reporter : Servers, leader : Servers, pos : Nat]$

275 $LogEntries \triangleq$

A log entry is composed of its current phase, a single command, and its set of dependencies. Although not necessary for correctness, the log entry also carries its position in the leader's log to simplify specification.

282 $[phase : \{\text{"fast"}, \text{"slow"}, \text{"committed"}\},$
 283 $cmd : Commands \cup \{NULL\},$
 284 $deps : \text{SUBSET } Dependencies,$
 285 $pos : Nat]$

287 $LogOK \triangleq$

For every server and leader's log , each log entry is correctly indexed, and contained within the set of possible log entries.

Moreover $NULL$ commands have no dependencies. During failure recovery, $NULL$ values can replace commands from which peers depend on but wasn't yet replicated by the leader before failing.

296 $\forall server, leader \in Servers :$
 297 $\forall i \in \text{DOMAIN } log[server][leader] :$
 298 LET
 299 $entry \triangleq log[server][leader][i]$
 300 IN
 301 $\wedge entry.pos = i$
 302 $\wedge entry \in LogEntries$
 303 $\wedge entry.cmd = NULL \Rightarrow entry.deps = \{\}$

305 $InitLog \triangleq$

Servers start with an empty sequence of commands for every known leader.

310 $log = [server \in Servers \mapsto [leader \in Servers \mapsto \langle \rangle]]$

A function from server process to a function from leaders to the followers' next index: a natural number indicating the position of the next command in the leader's log to be sent to the associated follower (as per Raft).

Note that the $nextIdx$ can move ahead of the log length. If $nextIdx \leq Len(log)$, there are commands to be sent to the follower. Otherwise, leader has already sent all its log entries to the follower, thus resulting in $nextIdx$ being set to $Len(log) + 1$.

323 VARIABLE $nextIdx$

325 $NextIdxOK \triangleq$

For every server, every leader's next index with regards to its follower is a natural number greater than zero.

330 $\forall server, leader, follower \in Servers :$
 331 $nextIdx[server][leader][follower] \in Nat \setminus \{0\}$

333 $InitNextIdx \triangleq$

Since logs start empty, $nextIdx$ starts as 1 for every known leader which indicates that at start there are no commands to be replicated.

338 $nextIdx = [$
 339 $server \in Servers \mapsto [$
 340 $leader \in Servers \mapsto [$
 341 $follower \in Servers \mapsto 1$
 342 $]]]$

Functions from server processes to a function from leaders to follower's match index: a natural number indicating the position where the follower matches the *log* from the leader with regards to fast and slow phases, respectively.

This is an adaptation of Raft's match index concept to allow for Atlas fast and slow phase commits.

353 VARIABLES *fastMatchIdx*, *slowMatchIdx*

355 *MatchIndexesOK* \triangleq

For every known server, the leader's match indexes with regards to each follower is be a natural number.

360 $\forall server, leader, follower \in Servers :$

361 $\wedge fastMatchIdx[server][leader][follower] \in Nat$

362 $\wedge slowMatchIdx[server][leader][follower] \in Nat$

364 *InitMatchIndexes* \triangleq

Since logs start empty, match indexes start at 0 for every known follower indicating no commands were replicated yet.

369 $\wedge fastMatchIdx = [$
 370 $server \in Servers \mapsto [$
 371 $leader \in Servers \mapsto [$
 372 $follower \in Servers \mapsto 0$

373 $]]]$

374 $\wedge slowMatchIdx = [$
 375 $server \in Servers \mapsto [$
 376 $leader \in Servers \mapsto [$
 377 $follower \in Servers \mapsto 0$

378 $]]]$

A function from servers to leader's commit index: a natural number indicating the highest command position committed for a given leader, or zero if unknown (as per Raft).

385 VARIABLE *commitIdx*

387 *CommitIdxOK* \triangleq

For every known leader, its commit index is a natural number.

391 $\forall server, leader \in Servers : commitIdx[server][leader] \in Nat$

393 *InitCommitIdx* \triangleq

Start with zero as the commit index for every known server process.

397 $commitIdx = [server \in Servers \mapsto [leader \in Servers \mapsto 0]]$

A function from leaders to their fast quorum of server processes, including itself, from which non-diverging responses can be safely committed without a second round of consensus (as per Atlas).

404 VARIABLES *fastQuorum*

406 *FastQuorumOK* \triangleq

A leader's fast quorum is correct if it's a subset of server processes, is properly sized, and contains the leader itself.

411 $\forall server \in Servers :$
412 $\wedge fastQuorum[server] \subseteq Servers$
413 $\wedge Cardinality(fastQuorum[server]) = FQ$
414 $\wedge server \in fastQuorum[server]$

416 $ChooseFastQuorum(leader, peers) \triangleq$
Arbitrarily chooses a fast quorum for the given leader among its available peers.

421 $CHOOSE quorum \in SUBSET (peers \cup \{leader\}) :$
422 $\wedge leader \in quorum$
423 $\wedge Cardinality(quorum) = FQ$

425 $InitFastQuorum \triangleq$
Arbitrarily chooses a fast quorum per leader among all server processes at start.

430 $fastQuorum = [leader \in Servers \mapsto ChooseFastQuorum(leader, Servers)]$

A function from server process to a function from term to server tracking for which server the current process voted in the associated term.
As per Raft, a server only votes for a single peer for a given election.

439 VARIABLE $votedFor$

441 $VotedForOK \triangleq$
For every leader known by server, known terms are natural numbers, and the associated cast vote is a server process.

446 $\forall server, leader \in Servers :$
447 $\forall term \in DOMAIN votedFor[server][leader] :$
448 $\wedge term \in Nat$
449 $\wedge votedFor[server][leader][term] \in Servers$

451 $InitVotedFor \triangleq$
Start with an empty mapping function as no votes have been cast.

455 $votedFor = [$
456 $server \in Servers \mapsto [$
457 $leader \in Servers \mapsto [$
458 $term \in \{\} \mapsto \{\}$
459 $]]]$

A function from a server process to a leader from which the server is attempting to recover.

465 VARIABLE $recovery$

467 $Recovery \triangleq$
Recovery is a temporary state maintained by candidates during the recovery process. It tracks the highest known command index for the failing leader.
Also, it tracks replies collected during election, which are composed of a reporting peer, the commit index for the failed leader known by the peer, a flag indicating if the vote has been granted, and a slice of the *log* from the failing leader according to the reporting peer.

```

478  [replies : SUBSET [
479    reporter      : Servers,
480    commitIdx     : Nat,
481    voteGranted   : BOOLEAN ,
482    slice         : Seq(LogEntries)
483  ],
484  maxIdx : Nat]

```

```

486 RecoveryOK  $\triangleq$ 

```

For every server, the state of each leader where server attempts recovery belongs to the set of possible recovery states.

```

491  $\forall$  server  $\in$  Servers :
492    $\forall$  leader  $\in$  DOMAIN recovery[server] :
493      $\wedge$  leader  $\in$  Servers
494      $\wedge$  recovery[server][leader]  $\in$  Recovery

```

```

496 InitRecovery  $\triangleq$ 

```

Start with no recovery in progress.

```

500 recovery = [server  $\in$  Servers  $\mapsto$  [leader  $\in$  {}  $\mapsto$  {}]]

```

A set of messages in the system's network. Note that the spec makes no assumption about message ordering. The usage of a set implies message uniqueness but, the spec makes no attempt to prevent duplicated message handling, thus modeling message loss and repeated message delivery from which the algorithm is safe against.

Note that three fields are present in every message: type, indicating the message type; src representing the message's source server; and dest, modeling the message's destination server. These fields are meant to model the network stack. Concrete implementations of the algorithm can omit them.

```

515 VARIABLE messages

```

```

517 AppendRPC  $\triangleq$ 

```

As per Raft, the *AppendRPC* is used to replicate commands from leaders to followers. Since a recovery leader can temporarily assume the leadership of a peer's *log*, the leader field indicates which leader's *log* is being replicated at the time. A message with *src* \neq *leader* can be seen as the recovery leader *src* replicating *log* on behalf of the failed leader.

The append *RPC* has been modified to include in its response a sequence of dependencies set for each replicated command, thus allowing for dependency collection like in the Atlas protocol.

```

530 LET

```

```

531 Request  $\triangleq$ 
532   [type      : {"append-request"},
533    src       : Servers,
534    dest      : Servers,
535    leader    : Servers,
536    term      : Nat,
537    leaderCommit : Nat,
538    prevLogIndex : Nat,
539    entries   : Seq(LogEntries)]

```



```

541   Reply  $\triangleq$ 
542     [type      : { "append-reply" },
543     src       : Servers,
544     dest      : Servers,
545     leader    : Servers,
546     success   : BOOLEAN ,
547     term      : Nat,
548     prevLogIndex : Nat,
549     dependencies : Seq(SUBSET Dependencies)]
550   IN
551     Request  $\cup$  Reply
553 RecoveryRPC  $\triangleq$ 
    The RecoveryRPC coordinates an election analogous to Raft. Moreover, it allows for participat-
    ing servers to report back on the known leader's log so that the elected leader can consolidate
    and make progress on pending commands on behalf of the failing leader.
560   LET
561     Request  $\triangleq$ 
562       [type      : { "recovery-request" },
563       src       : Servers,
564       dest      : Servers,
565       leader    : Servers,
566       term      : Nat,
567       commitIdx : Nat]
569     Reply  $\triangleq$ 
570       [type      : { "recovery-reply" },
571       src       : Servers,
572       dest      : Servers,
573       leader    : Servers,
574       term      : Nat,
575       maxIdx    : Nat,
576       commitIdx : Nat,
577       voteGranted : BOOLEAN ,
578       slice     : Seq(LogEntries)]
579   IN
580     Request  $\cup$  Reply
582 MessagesOK  $\triangleq$ 
    A network message is valid if it's part of the either append or recovery RPCs.
583   messages  $\subseteq$  AppendRPC  $\cup$  RecoveryRPC
589 InitMessages  $\triangleq$ 
    The network starts with an empty set of messages.
593   messages = {}

```

The sequence of all variables defined in the spec.

```

598 vars  $\triangleq$  ⟨
599   currTerm, currState, log, nextIdx, fastMatchIdx, slowMatchIdx,
600   commitIdx, fastQuorum, votedFor, recovery, messages
601 ⟩
602 |

```

— Helper Functions —

```

607 ComputeDeps(peer, command)  $\triangleq$ 
    Compute dependencies for the given command according to the peer's view of all leaders' logs.
612 LET
613   allLogs  $\triangleq$  {
614     ⟨leader, SetOf(log[peer][leader])⟩ : leader ∈ Servers
615   }
616   allEntries  $\triangleq$  {
617     ⟨leader, entry⟩ : entry ∈ llog} : ⟨leader, llog⟩ ∈ allLogs
618   }
619   conflicts  $\triangleq$  {
620     ⟨leader, entry⟩ ∈ UNION allEntries : Conflict(command, entry.cmd)
621   }
622 IN {
623   [reporter ↦ peer,
624     leader ↦ leader,
625     pos ↦ entry.pos
626   ] : ⟨leader, entry⟩ ∈ conflicts
627 }

```

```

629 MaxDepPos(server, leader)  $\triangleq$ 
    Returns the highest uncommitted index from which any command known by server depend on
    the given leader.
634 LET
635   allLogs  $\triangleq$  UNION {
636     SetOf(log[server][peer]) : peer ∈ Servers
637   }
638   uncommittedDeps  $\triangleq$  {
639     entry ∈ allLogs :
640     ∃ dep ∈ entry.deps :
641       ∧ dep.leader = leader
642       ∧ dep.pos > commitIdx[server][leader]
643   }
644   uncommittedDepsPos  $\triangleq$  {
645     entry.pos : entry ∈ uncommittedDeps
646   }
647 IN
648   CHOOSE i ∈ uncommittedDepsPos ∪ {0} :

```

649 $\forall j \in uncommittedDepsPos : i \geq j$

651 $StartRecovery(server, leader) \triangleq$

Set all variables needed to start a recovery for the given leader with server as a candidate.

Similarly to Raft, elections start by incrementing the leader's current term, changing the server's state to candidate with regards to the leader, and voting for itself. Additionally, the temporary state of recovery is reset to account for the candidate server vote.

```

661  $\wedge currTerm' = [$ 
662    $currTerm \text{ EXCEPT } ![server][leader] = @ + 1$ 
663  $]$ 
664  $\wedge currState' = [$ 
665    $currState \text{ EXCEPT } ![server][leader] = \text{"candidate"}$ 
666  $]$ 
667  $\wedge votedFor' = [$ 
668    $votedFor \text{ EXCEPT } ![server][leader] =$ 
669      $(currTerm'[server][leader] :> server) @@@$ 
670  $]$ 

```

The maximum index for the leader's *log* is the highest command position known by the candidate, rather the *log* length or a dependency reported by the leader but not replicated to the candidate. Moreover, considering that a candidate always vote for itself, it fills the replies field with an equivalent reply record accounting for its own.

```

679  $\wedge recovery' = [$ 
680    $recovery \text{ EXCEPT } ![server] =$ 
681      $(leader :> [$ 
682        $maxIdx \mapsto Max($ 
683          $MaxDepPos(server, leader),$ 
684          $Len(log[server][leader])$ 
685        $),$ 
686        $replies \mapsto \{[$ 
687          $voteGranted \mapsto \text{TRUE},$ 
688          $reporter \mapsto server,$ 
689          $commitIdx \mapsto commitIdx[server][leader],$ 
690          $slice \mapsto SubSeq($ 
691            $log[server][leader],$ 
692            $commitIdx[server][leader] + 1,$ 
693            $Len(log[server][leader])$ 
694          $),$ 
695        $]\}$ 
696      $)] @@@$ 
697  $]$ 

```

699 — Actions —

Like Raft, *ERaft* is design to work on a timeout loop. Since *ERaft* maintains multiple logs, Raft's original timeout design must also be adapted to work as one timeout per *log*. Although time isn't modeled by the spec, the following actions describe operations ran in response to timeouts.

Note that, although any interleaving among defined actions is possible, each action is itself an atomic operation.

713 $AppendCommandToLog(leader, command) \triangleq$

This action models a client submitting a command to given leader.

Given a leader and a command, compute the command's set of dependencies and append it to the leader's *log* in the fast phase.

720 $\wedge currState[leader][leader] = \text{"leader"}$
 721 $\wedge log' = [$
 722 $log \text{ EXCEPT } ![leader][leader] =$
 723 $Append(@, [$
 724 $phase \mapsto \text{"fast"},$
 725 $cmd \mapsto command,$
 726 $deps \mapsto ComputeDeps(leader, command),$
 727 $pos \mapsto Len(log[leader][leader]) + 1$
 728 $])$
 729 $]$

The leader won't attempt to replicate its own *log* to itself, thus, the following indexes won't be updated in response to the append *RPC* call. Therefore, adjusting their values here will simplify decisions made later in the algorithm as the leader won't have to be treated as a special server elsewhere.

737 $\wedge nextIdx' = [nextIdx \text{ EXCEPT } ![leader][leader][leader] = @ + 1]$
 738 $\wedge fastMatchIdx' = [fastMatchIdx \text{ EXCEPT } ![leader][leader][leader] = @ + 1]$
 739 $\wedge slowMatchIdx' = [slowMatchIdx \text{ EXCEPT } ![leader][leader][leader] = @ + 1]$
 740 $\wedge \text{UNCHANGED } \langle$
 741 $currTerm, currState, commitIdx, fastQuorum,$
 742 $votedFor, recovery, messages$
 743 \rangle

— REPLICATION ACTIONS —

749 $SendAppendRequest(server, leader) \triangleq$

This action models a leader timeout, either original ($server = leader$) or recovery leader ($server \neq leader$).

Upon a timeout, initiates the append *RPC* from server on behalf of the leader by sending an "append-request" message in order to replicate pending commands and/or its latest commit index.

As per Raft, *log* entries sent are a slice of the leader's *log* from *nextIdx* to $Len(log)$. Note that, as per Atlas, commands in the "fast" phase are only sent to the leader's fast quorum. Otherwise, commands are sent via broadcast to all servers.

763 $\wedge currState[server][leader] = \text{"leader"}$
 764 $\wedge messages' = messages \cup \{$
 765 $[type \mapsto \text{"append-request"},$
 766 $src \mapsto server,$
 767 $dest \mapsto follower,$

```

768     leader      ↦ leader,
769     term        ↦ currTerm[server][leader],
770     leaderCommit ↦ commitIdx[server][leader],
771     prevLogIndex ↦ nextIdx[server][leader][follower] - 1,
772     entries     ↦ TakeWhile(
773         SubSeq(
774             log[server][leader],
775             nextIdx[server][leader][follower],
776             Len(log[server][leader])
777         ),
778     LAMBDA entry :
779         entry.phase = "fast" ⇒
780         follower ∈ fastQuorum[leader]
781     )
782 ] : follower ∈ Servers \ {server}
783 }
784 ∧ UNCHANGED ⟨
785     currTerm, currState, log, nextIdx, fastMatchIdx, slowMatchIdx,
786     commitIdx, fastQuorum, votedFor, recovery
787 ⟩

```

789 $HandleAppendRequest(msg) \triangleq$

This action models a server receiving an "append-request" message from an alleged leader.

794 $\wedge msg.type = \text{"append-request"}$

Discard old/repeated messages with outdated commit indexes.

798 $\wedge msg.leaderCommit \geq commitIdx[msg.dest][msg.leader]$

As per Raft, leaders may attempt to proceed with *log* replication after failing. If noticing an outdated term, reject the append and reply with the most recent term for the leader.

Moreover, if the follower is behind the leader's *log*, which can occur when the elected recovery leader has more commands in its *log* than the follower, reply with the *Len(log)* so that the leader can adjust its *nextIdx* for the next append request.

809 $\wedge \text{IF } \vee msg.term < currTerm[msg.dest][msg.leader]$

810 $\vee msg.prevLogIndex > Len(log[msg.dest][msg.leader])$

811 THEN

812 $\wedge messages' = messages \cup \{$

813 $type \quad \mapsto \text{"append-reply"},$

814 $success \quad \mapsto \text{FALSE},$

815 $src \quad \mapsto msg.dest,$

816 $dest \quad \mapsto msg.src,$

817 $leader \quad \mapsto msg.leader,$

818 $term \quad \mapsto currTerm[msg.dest][msg.leader],$

819 $prevLogIndex \mapsto Len(log[msg.dest][msg.leader]),$

820 $dependencies \mapsto \langle \rangle$

821 $\}$

```

822     }
823     ^ UNCHANGED <
824         currTerm, currState, log, nextIdx, fastMatchIdx, slowMatchIdx,
825         commitIdx, fastQuorum, votedFor, recovery
826     >
827 ELSE
828     ^ LET

```

Append behaves differently depending on the command's current phase. Commands in the slow or committed phase are simply accepted as long as they don't override an existing committed command.

New commands in the fast phase are accepted with an updated set of dependencies that takes the follower's view of the commands; dependencies set into account.

On the other hand, existing commands in the fast phase are preserved since repeated calls to the append *RPC* requires the set of dependencies to remain the same until the leader makes a decision to either commit or promote the command to the slow phase (as per Atlas).

```

845     entries  $\triangleq$  Map(msg.entries, LAMBDA entry :
846         LET
847             llog  $\triangleq$  log[msg.dest][msg.leader]
848             curr  $\triangleq$  llog[entry.pos]
849         IN
850             IF entry.phase  $\neq$  "fast"
851                 THEN
852                     IF ^ entry.pos  $\leq$  Len(llog)
853                         ^ curr.phase = "committed"
854                         THEN curr
855                         ELSE entry
856                 ELSE
857                     IF entry.pos  $\leq$  Len(llog)
858                         THEN curr
859                     ELSE [
860                         entry EXCEPT !.deps =
861                             @  $\cup$  ComputeDeps(msg.dest, entry.cmd)
862                     ]
863             )
864     IN

```

If receiving an append from a leader with higher term, abdicate your current state (candidate or leader) and start following the new leader (as per Raft).

```

870     ^ currTerm' = [
871         currTerm EXCEPT ![msg.dest][msg.leader] =
872         msg.term
873     ]
874     ^ currState' = [
875         currState EXCEPT ![msg.dest][msg.leader] =
876         "follower"

```

```

877 ]
      Update the follower's log according to the updated log entries and reply back to
      leader if there are any dependencies to report. Note that the leader might send
      empty append messages to restate its leadership and prevent recoveries from
      happening, as well as communicate progress on its commit index (as per Raft).
886  $\wedge \text{log}' = [$ 
887    $\text{log EXCEPT ![msg.dest][msg.leader]} =$ 
888      $\text{ReplaceSlice}(@, \text{msg.prevLogIndex}, \text{entries})$ 
889 ]
890  $\wedge \text{commitIdx}' = [$ 
891    $\text{commitIdx EXCEPT ![msg.dest][msg.leader]} =$ 
892      $\text{msg.leaderCommit}$ 
893 ]
894  $\wedge \text{IF Len}(\text{entries}) > 0$ 
895   THEN
896      $\text{messages}' = \text{messages} \cup \{$ 
897        $\text{[type} \quad \mapsto \text{"append-reply"},$ 
898        $\text{success} \quad \mapsto \text{TRUE},$ 
899        $\text{src} \quad \mapsto \text{msg.dest},$ 
900        $\text{dest} \quad \mapsto \text{msg.src},$ 
901        $\text{leader} \quad \mapsto \text{msg.leader},$ 
902        $\text{term} \quad \mapsto \text{msg.term},$ 
903        $\text{prevLogIndex} \mapsto \text{msg.prevLogIndex},$ 
904        $\text{dependencies} \mapsto \text{Map}(\text{entries}, \text{LAMBDA } e : e.\text{deps})$ 
905     ]
906   }
907   ELSE
908     UNCHANGED messages
909  $\wedge$  UNCHANGED (
910    $\text{nextIdx}, \text{fastMatchIdx}, \text{slowMatchIdx},$ 
911    $\text{fastQuorum}, \text{votedFor}, \text{recovery}$ 
912 )

```

914 $\text{HandleAppendReply}(\text{msg}) \triangleq$

This action models a server receiving an "append-reply" message from a follower.

919 $\wedge \text{msg.type} = \text{"append-reply"}$

Reject replies when the destination is not the *log*'s leader or from outdated terms (as per Raft).

924 $\wedge \text{currState}[\text{msg.dest}][\text{msg.leader}] = \text{"leader"}$

925 $\wedge \text{msg.term} \geq \text{currTerm}[\text{msg.dest}][\text{msg.leader}]$

Reject old/repeated messages with *log* slices starting at an outdated replication position.

If the *log* slice starts ahead the fast match index, this is a slice of new commands being replicated to the follower. As a result of the append *RPC*, the fast match index gets updated to the maximum *log* position known to be replicated to the follower.

If the *log* slice starts behind of the fast match index, this is either: a slice containing commands promoted to the slow phase and are, therefore, being replicated again to the follower, to which case the *log* slice starts ahead the slow match index; or else it's an old message and must be discarded.

941 $\wedge \text{msg.prevLogIndex} < \text{fastMatchIdx}[\text{msg.dest}][\text{msg.leader}][\text{msg.src}] \Rightarrow$
942 $\text{msg.prevLogIndex} \geq \text{slowMatchIdx}[\text{msg.dest}][\text{msg.leader}][\text{msg.src}]$

On a failed append request, consider outdated follower or leadership loss.

947 $\wedge \text{IF } \neg \text{msg.success}$
948 THEN
949 IF $\text{msg.term} = \text{currTerm}[\text{msg.dest}][\text{msg.leader}]$
950 THEN

If still the leader, the follower must be behind in *log* replication and has set the *prevLogIndex* in the reply to its known *log* length. Reset the next index at the leader so that next append *RPC* can succeed. This is slight different than traditional Raft design where it simply decrements the next index by one.

959 $\wedge \text{nextIdx}' = [$
960 nextIdx EXCEPT $![\text{msg.dest}][\text{msg.leader}][\text{msg.src}] =$
961 $\text{msg.prevLogIndex} + 1$
962 $]$
963 \wedge UNCHANGED \langle
964 $\text{currTerm}, \text{currState}, \text{log}, \text{fastMatchIdx}, \text{slowMatchIdx},$
965 $\text{commitIdx}, \text{fastQuorum}, \text{votedFor}, \text{recovery}, \text{messages}$
966 \rangle
967 ELSE

If receiving a reply from a peer with a higher term, fallback to being a follower (as per Raft).

972 $\wedge \text{currState}' = [$
973 currState EXCEPT $![\text{msg.dest}][\text{msg.leader}] =$
974 "follower"
975 $]$
976 $\wedge \text{currTerm}' = [$
977 currTerm EXCEPT $![\text{msg.dest}][\text{msg.leader}] =$
978 msg.term
979 $]$
980 \wedge UNCHANGED \langle
981 $\text{log}, \text{nextIdx}, \text{fastMatchIdx}, \text{slowMatchIdx}, \text{commitIdx},$
982 $\text{fastQuorum}, \text{votedFor}, \text{recovery}, \text{messages}$
983 \rangle
984 ELSE
985 \wedge LET

When handling a reply, map over the commands in the same *log* slice and update their dependencies set depending on their current phase.

Commands in the slow or committed phase are kept as they are, while commands in the fast phase accumulate dependencies reported by their followers.

995 $\text{entries} \triangleq \text{Map}(\text{$


```

996     SubSeq(
997         log[msg.dest][msg.leader],
998         msg.prevLogIndex + 1,
999         msg.prevLogIndex + Len(msg.dependencies)
1000     ),
1001     LAMBDA entry :
1002         IF entry.phase ≠ "fast"
1003         THEN entry
1004         ELSE [
1005             entry EXCEPT !.deps =
1006             @ ∪ msg.dependencies[entry.pos - msg.prevLogIndex]
1007         ]
1008     )

```

For commands in the slow phase, identify the prefix of commands with matching dependencies with the replying follower.

```

1014     slow ≜ SelectSeq(entries, LAMBDA entry :
1015         entry.phase = "slow"
1016     )
1017     matching ≜ TakeWhile(slow, LAMBDA entry :
1018         entry.deps = msg.dependencies[entry.pos - msg.prevLogIndex]
1019     )
1020 IN

```

Update the log to account for new dependencies learned.

```

1025     ∧ log' = [
1026         log EXCEPT ![msg.leader][msg.leader] =
1027         ReplaceSlice(@, msg.prevLogIndex, entries)
1028     ]

```

Fast match index is updated to match the highest command position known to be replicated at the replying follower.

```

1034     ∧ fastMatchIdx' = [
1035         fastMatchIdx EXCEPT ![msg.dest][msg.leader][msg.src] =
1036         Max(@, Last(entries).pos)
1037     ]

```

Slow match index is updated to the highest contiguous command known to be have matching dependencies with the replying follower.

```

1043     ∧ slowMatchIdx' = [
1044         slowMatchIdx EXCEPT ![msg.dest][msg.leader][msg.src] =
1045         IF Len(matching) > 0 THEN Last(matching).pos ELSE @
1046     ]

```

Next index is set to fast match + 1 if no slow commands are being handled in this reply, slow match + 1 if there are matching commands in the slow phase, or else remains unchanged.

```

1053     ∧ nextIdx' = [
1054         nextIdx EXCEPT ![msg.dest][msg.leader][msg.src] =
1055         CASE Len(slow) = 0 →

```

```

1056         fastMatchIdx'[msg.dest][msg.leader][msg.src] + 1
1057     □ Len(matching) > 0 →
1058         slowMatchIdx'[msg.dest][msg.leader][msg.src] + 1
1059     □ OTHER → @
1060     ]
1061     ∧ UNCHANGED ⟨
1062         currTerm, currState, commitIdx, fastQuorum,
1063         votedFor, recovery, messages
1064     ⟩
1066 CommitReplicatedCommands(server, leader) ≜
    This action models the part of the leader's timeout loop where it decides which commands to
    commit (as per Raft).
1071     ∧ currState[server][leader] = "leader"
1072     ∧ LET
    Select commands above the leader's commit index which are known to be replicated to a
    sufficient number of followers according to their phase.
1078     replicated ≜ TakeWhile(
1079         SubSeq(
1080             log[server][leader],
1081             commitIdx[server][leader] + 1,
1082             Len(log[server][leader])
1083         ),
1084         LAMBDA entry :
1085             CASE entry.phase = "committed" → TRUE
1086             □ entry.phase = "slow" →
1087                 LET matching ≜ {
1088                     follower ∈ Servers :
1089                         entry.pos ≤ slowMatchIdx[server][leader][follower]
1090                 }
1091                 IN
1092                     Cardinality(matching) ≥ SQ
1093             □ entry.phase = "fast" →
1094                 LET matching ≜ {
1095                     follower ∈ Servers :
1096                         entry.pos ≤ fastMatchIdx[server][leader][follower]
1097                 }
1098                 IN
1099                     Cardinality(matching) ≥ FQ
1100     )

```

Update replicated commands: commands committed are preserved; commands in the slow phase are promoted to the committed phase; commands in the fast phase are promoted to either the committed or slow phases.

As per Atlas, commands in the fast phase can skip the slow phase if at least F nodes, excluding itself, have reported the same dependency. Otherwise commands must be promoted to the slow phase for another append *RPC* call.

```

1112 updated  $\triangleq$  Map(replicated, LAMBDA entry :
1113   CASE entry.phase = "committed"  $\rightarrow$  entry
1114     □ entry.phase = "slow"  $\rightarrow$  [entry EXCEPT !.phase = "committed"]
1115     □ entry.phase = "fast"  $\rightarrow$ 
1116       LET
1117         depsReportedByPeers  $\triangleq$  {
1118           dep  $\in$  entry.deps :
1119             dep.reporter  $\neq$  leader
1120         }
1121         canSkipSlowPhase  $\triangleq$ 
1122            $\forall$  dep  $\in$  depsReportedByPeers :
1123             LET
1124               sameDepReports  $\triangleq$  {
1125                 d  $\in$  depsReportedByPeers :
1126                    $\wedge$  d.leader = dep.leader
1127                    $\wedge$  d.pos = dep.pos
1128               }
1129             IN
1130               Cardinality(sameDepReports)  $\geq$  F
1131         IN [
1132           entry EXCEPT !.phase =
1133             IF canSkipSlowPhase
1134               THEN "committed"
1135               ELSE "slow"
1136         ]
1137   )
1138 promoted  $\triangleq$  SelectSeq(updated, LAMBDA entry :
1139    $\wedge$  log[server][leader][entry.pos].phase = "fast"
1140    $\wedge$  entry.phase  $\in$  {"slow", "committed"}
1141 )
1142 committed  $\triangleq$  TakeWhile(updated, LAMBDA entry :
1143   entry.phase = "committed"
1144 )
1145 IN
  Only executes this action if promoting or committing commands.
1150  $\wedge$   $\forall$  Len(promoted)  $>$  0
1151  $\forall$  Len(committed)  $>$  0
  Updates the log to account for new command phases.
1155  $\wedge$  log' = [
1156   log EXCEPT ![server][leader] =
1157     ReplaceSlice(@, commitIdx[server][leader], updated)
1158 ]
  Commit index is updated to the highest contiguous log position that is committed at
  the leader.

```

```

1163    $\wedge$   $commitIdx' = [$ 
1164        $commitIdx$  EXCEPT  $![server][leader] =$ 
1165       IF  $Len(committed) = 0$  THEN  $@$  ELSE  $Last(committed).pos$ 
1166   ]
```

When promoting an entry from fast phase to either slow or committed, it's possible that not all followers are aware of all the dependencies collected by the leader during the initial append *RPC* call. Therefore, adjust the next index for all followers to ensure that receive an updated set of dependencies from the leader.

```

1175    $\wedge$   $nextIdx' = [$ 
1176        $nextIdx$  EXCEPT  $![server][leader] =$ 
1177       IF  $Len(promoted) = 0$  THEN  $@$  ELSE
1178          $(server := @[server]) @@ [$ 
1179            $other \in Servers \mapsto Head(promoted).pos$ 
1180         ]
1181   ]
```

```

1182    $\wedge$  UNCHANGED  $\langle$ 
1183      $currTerm, currState, fastMatchIdx, slowMatchIdx,$ 
1184      $fastQuorum, votedFor, recovery, messages$ 
1185    $\rangle$ 

```

— FAILURE RECOVERY ACTIONS —

```

1191 SuspectFailure( $server, leader$ )  $\triangleq$ 
```

This action models a follower timeout while waiting for leader messages. Start a recovery process if suspecting of its own failure ($server = leader$), or suspecting of some other process failure from which server has a dependency on.

```

1198    $\wedge currState[server][leader] = \text{"follower"}$ 
1199    $\wedge \vee server = leader$ 
1200      $\vee MaxDepPos(server, leader) > 0$ 
1201    $\wedge StartRecovery(server, leader)$ 
1202    $\wedge$  UNCHANGED  $\langle$ 
1203      $log, nextIdx, fastMatchIdx, slowMatchIdx,$ 
1204      $fastQuorum, commitIdx, messages$ 
1205    $\rangle$ 
```

```

1207 ReconfigureFastQuorum( $server, leader$ )  $\triangleq$ 
```

This action models a part of the follower timeout where it determines if the suspected failed leader is part of the server's fast quorum. Since Atlas require a fast quorum of responses before making progress on the fast phase, this quorum must be reconfigured upon failures to restore the server's progress.

Upon identifying a failure of a leader in its fast quorum, the server chooses a different fast quorum that does not contain the failing process. Additionally, it calls for its own recovery to complement pending its pending commands before start using the newly chosen quorum.

```

1221    $\wedge server \neq leader$ 
1222    $\wedge currState[server][server] = \text{"leader"}$ 
1223    $\wedge leader \in fastQuorum[server]$ 
1224    $\wedge fastQuorum' = [$ 

```

```

1225     fastQuorum EXCEPT ![server] =
1226         ChooseFastQuorum(server, Servers \ {leader})
1227     ]
1228     ∧ StartRecovery(server, server)
1229     ∧ UNCHANGED ⟨
1230         log, nextIdx, fastMatchIdx, slowMatchIdx, commitIdx, messages
1231     ⟩

```

1233 *SendRecoveryRequest*(server, leader) \triangleq

As per Raft, elections happen as part of the server's timeout loop. As long as server is a candidate for the leader's *log*, send a "recovery-request" message via broadcast to all peers.

Note that since recoveries start a new term, and that appends from previous terms are discarded, the leader's commit index will remain the same until the end of the recovery process.

```

1243     ∧ currState[server][leader] = "candidate"
1244     ∧ messages' = messages ∪ {
1245         [type      ↦ "recovery-request",
1246          src       ↦ server,
1247          dest      ↦ follower,
1248          leader    ↦ leader,
1249          term      ↦ currTerm[server][leader],
1250          commitIdx ↦ commitIdx[server][leader]
1251         ] : follower ∈ Servers \ {server}
1252     }
1253     ∧ UNCHANGED ⟨
1254         currTerm, currState, log, nextIdx, fastMatchIdx, slowMatchIdx,
1255         commitIdx, fastQuorum, votedFor, recovery
1256     ⟩

```

1258 *HandleRecoveryRequest*(msg) \triangleq

This action models a server receiving a recovery request from a candidate peer for a given leader's *log*.

Messages from older terms are discarded. Otherwise, if receiving a recovery request from a higher term, it abdicates its current state and fallback to the follower state (as per Raft election rules).

```

1267     ∧ msg.term ≥ currTerm[msg.dest][msg.leader]
1268     ∧ currTerm' = [
1269         currTerm EXCEPT ![msg.dest][msg.leader] =
1270             msg.term
1271     ]
1272     ∧ currState' = [
1273         currState EXCEPT ![msg.dest][msg.leader] =
1274             IF msg.term = currTerm[msg.dest][msg.leader]
1275             THEN @ ELSE "follower"
1276     ]

```

Only vote for a single candidate per term. Note that the semantics of $f @ @ g$ preserves the domain of f , meaning that if $msg.term$ already belongs to $votedFor[msg.dest][msg.leader]$, the association $(msg.term :> msg.src)$ is discarded.

```

1283  $\wedge votedFor' = [$ 
1284    $votedFor \text{ EXCEPT } ![msg.dest][msg.leader] =$ 
1285      $@ @ @ (msg.term :> msg.src)$ 
1286 ]
```

Reply back with the server's commit index for the failed leader, and a *log* slice starting from the candidate's commit index. This mechanism allows the candidate to become aware of committed commands it wasn't aware at the time it initiated the recovery process. Also, replies with the voter's maximum known command index for the failed leader.

```

1295  $\wedge messages' = messages \cup \{$ 
1296    $[type \quad \mapsto \text{"recovery-reply"},$ 
1297      $src \quad \mapsto msg.dest,$ 
1298      $dest \quad \mapsto msg.src,$ 
1299      $leader \mapsto msg.leader,$ 
1300      $term \quad \mapsto msg.term,$ 
1301      $commitIdx \mapsto commitIdx[msg.dest][msg.leader],$ 
1302      $voteGranted \mapsto votedFor'[msg.dest][msg.leader][msg.term] = msg.src,$ 
1303      $maxIdx \quad \mapsto Max($ 
1304        $MaxDepPos(msg.dest, msg.leader),$ 
1305        $Len(log[msg.dest][msg.leader])$ 
1306     ),
1307      $slice \quad \mapsto SubSeq($ 
1308        $log[msg.dest][msg.leader],$ 
1309        $msg.commitIdx + 1,$ 
1310        $Len(log[msg.dest][msg.leader])$ 
1311     )
1312   ]
1313 }
1314  $\wedge \text{UNCHANGED } \langle$ 
1315    $log, nextIdx, fastMatchIdx, slowMatchIdx,$ 
1316    $commitIdx, fastQuorum, recovery$ 
1317  $\rangle$ 
```

```

1319  $HandleRecoveryReply(msg) \triangleq$ 
```

This action models a candidate receiving "recovery-reply" messages from its peers. Replies from different terms are discarded, as well as already registered votes.

```

1325  $\wedge currState[msg.dest][msg.leader] = \text{"candidate"}$ 
1326  $\wedge msg.term = currTerm[msg.dest][msg.leader]$ 
1327  $\wedge msg.src \notin \{$ 
1328    $reply.reporter : reply \in recovery[msg.dest][msg.leader].replies$ 
1329  $\}$ 
```

Keep track of replies received as well as the maximum command index known among all peers.

```

1334  $\wedge$  recovery' = [
1335   recovery EXCEPT
1336     ![msg.dest][msg.leader].maxIdx = Max(@, msg.maxIdx),
1337     ![msg.dest][msg.leader].replies = @  $\cup$  {
1338       [reporter       $\mapsto$  msg.src,
1339        commitIdx     $\mapsto$  msg.commitIdx,
1340        voteGranted  $\mapsto$  msg.voteGranted,
1341        slice         $\mapsto$  msg.slice]
1342     }
1343 ]
1344  $\wedge$  UNCHANGED  $\langle$ 
1345   currTerm, currState, log, nextIdx, fastMatchIdx, slowMatchIdx,
1346   commitIdx, fastQuorum, votedFor, messages
1347  $\rangle$ 

```

1349 *ClaimLeadership*(*server*, *peer*) \triangleq

This action models the part of the server's timeout loop where it decides to claim ownership of the leader's *log* if sufficient votes were gathered.

```

1355  $\wedge$  currState[server][peer] = "candidate"
1356  $\wedge$  LET
1357   replies       $\triangleq$  recovery[server][peer].replies
1358   maxIdx       $\triangleq$  recovery[server][peer].maxIdx
1359   votes        $\triangleq$  {r  $\in$  replies : r.voteGranted}
1360   commitIdxs  $\triangleq$  {r.commitIdx : r  $\in$  replies}

```

If elected, the recovery leader computes the current state of the failed leader's *log* by merging the *log* slices received by the voters.

1367 *MergedSlice*[*n* \in *Nat*] \triangleq

1368 IF *n* = 0 THEN \langle ELSE

1369 LET

Define the set of relevant replies for the *n*th position on the merged *log* slice as the set of replies with slice length greater or equal to *n*. Commands in the fast, slow, and committed phase are derived from relevant replies.

```

1376   relevant     $\triangleq$  {r  $\in$  replies : n  $\leq$  Len(r.slice)}
1377   fast        $\triangleq$  {r  $\in$  relevant : r.slice[n].phase = "fast"}
1378   slow        $\triangleq$  {r  $\in$  relevant : r.slice[n].phase = "slow"}
1379   committed  $\triangleq$  {r  $\in$  relevant :
1380      $\vee$  r.slice[n].phase = "committed"
1381      $\vee$  r.slice[n].pos  $\leq$  r.commitIdx
1382   }

```

The pick function chooses an arbitrary command for the *n*th position of the merged *log* slice.

1387 *Pick*(*replySet*) \triangleq (CHOOSE *r* \in *replySet* : TRUE).*slice*[*n*]

1388 IN

To ensure consistency, the following rules must be respected:

- *Commands* in the committed phase MUST be preserved as consensus have been achieved and some processes might already have executed them;
- *Commands* in the slow phase MUST be preserved as their dependency set have been defined by the leader, which may have committed and executed it already.
- *Commands* in the fast phase reported by at least $\text{floor}(N/2)$ MUST be preserved;

As per *Atlas*, if $\text{floor}(N/2)$ replies contains a command in the fast phase, the recovery leader can't determine if the original leader had received all $\text{floor}(N/2)+F$ replies required to promote the command to the slow or committed phases. Therefore, the command MUST be preserved with the SAME set of dependencies.

Atlas proves that the union of the dependencies reported by $\text{floor}(N/2)$ servers is guaranteed to be the same as the original leader's set of dependencies if it had decided to skip the slow phase for a given command. Otherwise, the leader promoted the command to slow phase but didn't replicated it to enough servers since the candidate didn't see it. In either case, preserve the command and promote it to the slow phase for another round of consensus.

– Otherwise, nullify the command.

If the command in the n th position is in the fast phase but not yet replicated to $\text{floor}(N/2)$ servers, the recovery leader is certain that the original leader didn't promote it to the slow or committed phases yet.

On the other hand, if the n th position is unknown by any of the replying hosts, it means that the leader had reported n as a dependency to one of its peers but it hasn't replicated it anywhere yet.

In both cases, replace the command with *NULL* so that, after committing it, servers can make progress on dependent commands' execution.

```

1437 MergedSlice[n - 1] ◦ ⟨
1438   CASE committed ≠ {} → Pick(committed)
1439   □ slow ≠ {} → Pick(slow)
1440   □ Cardinality(fast) ≥ (N ÷ 2) → [
1441     phase ↦ "slow",
1442     pos ↦ Pick(fast).pos,
1443     cmd ↦ Pick(fast).cmd,
1444     deps ↦ UNION {r.slice[n].deps : r ∈ fast}
1445   ]
1446   □ OTHER → [
1447     phase ↦ "slow",
1448     cmd ↦ NULL,
1449     pos ↦ commitIdx[server][peer] + n,
1450     deps ↦ {}
1451   ]
1452   ⟩
1453 IN

```

As per *Atlas*, a recovery can only occur after $N - F$ (*RQ*) nodes reply. As per *Raft*, a candidate can only get elected with a majority vote of $\text{floor}(N/2) + 1$ (*EQ*).


```

1459  $\wedge \text{Cardinality}(\text{replies}) \geq RQ$ 
1460  $\wedge \text{Cardinality}(\text{votes}) \geq EQ$ 
    Provided enough replies and votes, the candidate becomes the leader for the failed
    server's log.
1465  $\wedge \text{currState}' = [$ 
1466      $\text{currState}$  EXCEPT ![server][peer] =
1467     "leader"
1468 ]
    The leader's log get replaced by the merged slice computed among all the received
    replies.
1473  $\wedge \text{log}' = [$ 
1474      $\text{log}$  EXCEPT ![server][peer] =
1475      $\text{ReplaceSlice}(@,$ 
1476          $\text{commitIdx}[\text{server}][\text{peer}],$ 
1477          $\text{MergedSlice}[\text{maxIdx}]$ 
1478     )
1479 ]
    Commit index is set to the highest commit index reported by the peers.
1484  $\wedge \text{commitIdx}' = [$ 
1485      $\text{commitIdx}$  EXCEPT ![server][peer] =
1486     CHOOSE  $i \in \text{commitIdxs} \cup \{@\}$  :
1487      $\forall j \in \text{commitIdxs} : i \geq j$ 
1488 ]
    Replication indexes are reset based on the previous commit index known by the recovery
    leader, therefore, commencing new rounds of the append RPC in order to replicate
    the merged log slice to its peers.
1495  $\wedge \text{nextIdx}' = [$ 
1496      $\text{nextIdx}$  EXCEPT ![server][peer] =
1497     ( $\text{server} :> \text{Len}(\text{log}[\text{server}][\text{peer}]) + 1$ ) @@ [
1498      $\text{other} \in \text{Servers} \mapsto \text{commitIdx}[\text{server}][\text{peer}] + 1$ 
1499 ]
1500 ]
1501  $\wedge \text{fastMatchIdx}' = [$ 
1502      $\text{fastMatchIdx}$  EXCEPT ![server][peer] =
1503     ( $\text{server} :> \text{Len}(\text{log}[\text{server}][\text{peer}])$ ) @@ [
1504      $\text{other} \in \text{Servers} \mapsto \text{commitIdx}[\text{server}][\text{peer}]$ 
1505 ]
1506 ]
1507  $\wedge \text{slowMatchIdx}' = [$ 
1508      $\text{slowMatchIdx}$  EXCEPT ![server][peer] =
1509     ( $\text{server} :> \text{Len}(\text{log}[\text{server}][\text{peer}])$ ) @@ [
1510      $\text{other} \in \text{Servers} \mapsto \text{commitIdx}[\text{server}][\text{peer}]$ 
1511 ]
1512 ]

```

1513 \wedge UNCHANGED $\langle currTerm, fastQuorum, votedFor, recovery, messages \rangle$

1515 $AbdicateLeadership(server, peer) \triangleq$

This action models a part of the recovery leader's timeout loop where it determines if the recovery has been successful (commit index = log length).

Considering that the recovery leader never appends new commands to the failed leader's log, eventually in a correct system, it will commit the log slice collected during its election. At that point, the recovery leader can abdicate its leadership of the failed leader's log by resetting its current state to "follower".

The failed leader promotes a new recovery for its own log once it comes back, which lets it learn about the consensus made by the recovery leader in its absence.

1531 $\wedge server \neq peer$

1532 $\wedge currState[server][peer] = \text{"leader"}$

1533 $\wedge commitIdx[server][peer] = Len(log[server][peer])$

1534 $\wedge currState' = [$

1535 $currState \text{ EXCEPT } ![server][peer] =$

1536 "follower"

1537 $]$

1538 \wedge UNCHANGED \langle

1539 $currTerm, log, nextIdx, fastMatchIdx, slowMatchIdx,$

1540 $commitIdx, fastQuorum, votedFor, recovery, messages$

1541 \rangle

1543

— INVARIANTS, ACTION PROPERTIES, AND LIVENESS PROPERTIES —

1548 $TypeOK \triangleq$

Invariant: types must stay valid for every state reached.

1552 $\wedge CurrTermOK$

1553 $\wedge CurrStateOK$

1554 $\wedge LogOK$

1555 $\wedge NextIdxOK$

1556 $\wedge MatchIndexesOK$

1557 $\wedge CommitIdxOK$

1558 $\wedge FastQuorumOK$

1559 $\wedge VotedForOK$

1560 $\wedge RecoveryOK$

1561 $\wedge MessagesOK$

1563 $CurrStateTransitions \triangleq$

Action Property: state transitions are valid.

Rules:

- Followers stay as followers or transition to candidates;

- Candidates can stay as candidates, transition to follower when receiving an append or recovery reply message with higher term, or transition to leader after winning an election;

- Leader's can stay as leaders, become candidates if suspecting of a failure in its fast quorum, or fallback to follower if receiving an append or recovery request message with higher term.

```

1579 □[∀ server, leader ∈ Servers :
1580   LET
1581     curr ≜ currState[server][leader]
1582     next ≜ currState'[server][leader]
1583   IN
1584     CASE curr = "follower" → next ∈ {"follower", "candidate"}
1585     □ curr = "candidate" → next ∈ {"follower", "candidate", "leader"}
1586     □ curr = "leader" → next ∈ {"follower", "candidate", "leader"}
1587   ]currState

```

— REPLICATION INVARIANTS AND ACTION PROPERTIES —

```

1593 NextIdxIsBoundedByTheLog ≜
    Invariant: next index is always  $1 \leq next \text{ index} \leq Len(log) + 1$ .
1597   ∀ server, leader, follower ∈ Servers :
1598     nextIdx[server][leader][follower] ∈ 1 .. Len(log[server][leader]) + 1
1600 MatchIdxsFollowEachOther ≜
    Invariant: fast match index is always greater or equal to the slow match index.
1605   ∀ server, leader, follower ∈ Servers :
1606     LET
1607       fastIdx ≜ fastMatchIdx[server][leader][follower]
1608       slowIdx ≜ slowMatchIdx[server][leader][follower]
1609     IN
1610       fastIdx ≥ slowIdx
1612 CommitIdxMovesForward ≜
    Action Property: commit index is always moving forward.
1616   □[∀ server, leader ∈ Servers :
1617     LET
1618       curr ≜ commitIdx[server][leader]
1619       next ≜ commitIdx'[server][leader]
1620     IN
1621       next ≥ curr
1622   ]commitIdx
1624 LogAlwaysGrow ≜
    Action Property: log length is either unchanged or it grows.
1628   □[∀ server, leader ∈ Servers :
1629     Len(log'[server][leader]) ≥ Len(log[server][leader])
1630   ]log
1632 EntryPhaseTransitions ≜

```

Action Property: command phase transitions are valid.

Rules:

- *Commands* in the fast phase can be promoted to the slow or committed phase;
- *Commands* in the slow phase can only be promoted to the committed phase;
- *Commands* in the committed phase remains in the same phase.

```

1646 □[∀ server, leader ∈ Servers :
1647   ∀ curr ∈ SetOf(log[server][leader]) :
1648     LET next ≜ log'[server][leader][curr.pos]
1649     IN CASE curr.phase = "fast"      → next.phase ∈ {"fast", "slow", "committed"}
1650         □ curr.phase = "slow"      → next.phase ∈ {"slow", "committed"}
1651         □ curr.phase = "committed" → next.phase = "committed"
1652   ]log

```

1654 *CommittedEntriesMatch* ≜

A *log* entry committed by a leader, if known to be committed by a follower, must contain the same command and set of dependencies.

```

1659 ∀ server, leader ∈ Servers :
1660   LET
1661     committed ≜ {
1662       entry ∈ SetOf(log[server][leader]) :
1663         ∨ entry.phase = "committed"
1664         ∨ entry.pos ≤ commitIdx[server][leader]
1665     }
1666   IN
1667     ∀ entry ∈ committed :
1668       ∀ follower ∈ Servers \ {server} :
1669         ( ∧ entry.pos ∈ DOMAIN log[follower][leader]
1670           ∧ ∨ entry.pos ≤ commitIdx[follower][leader]
1671           ∨ log[follower][leader][entry.pos].phase = "committed"
1672         ) ⇒
1673         ∧ log[follower][leader][entry.pos].cmd = entry.cmd
1674         ∧ log[follower][leader][entry.pos].deps = entry.deps

```

1676 *CommittedEntriesDoNotChange* ≜

Action Property: commands in the committed phase don't change.

```

1680 □[∀ server, leader ∈ Servers :
1681   LET
1682     committed ≜ {
1683       entry ∈ SetOf(log[server][leader]) :
1684         ∨ entry.phase = "committed"
1685         ∨ entry.pos ≤ commitIdx[server][leader]
1686     }
1687   IN
1688     ∀ curr ∈ committed :

```

1689 $curr = log'[server][leader][curr.pos]$
 1690 $]_{log}$

— RECOVERY INVARIANTS AND ACTION PROPERTIES —

1696 *CurrTermMovesForward* \triangleq

Action Property: a leader's current term can only move forward.

1700 $\square[\forall server, leader \in Servers :$
 1701 $currTerm'[server][leader] \geq currTerm[server][leader]$
 1702 $]_{currTerm}$

1704 *OnlyVoteOncePerTerm* \triangleq

Action Property: A server only votes for a single candidate on a given term.

1709 $\square[\forall server, leader \in Servers :$
 1710 $\forall term \in \text{DOMAIN } votedFor[server][leader] :$
 1711 LET
 1712 $curr \triangleq votedFor[server][leader][term]$
 1713 $next \triangleq votedFor'[server][leader][term]$
 1714 IN
 1715 $curr = next$
 1716 $]_{votedFor}$

1718 *NoDoubleLeader* \triangleq

Invariant: for a given term, only one leader gets elected.

1722 LET
 1723 $terms \triangleq \{$
 1724 $currTerm[server][peer] :$
 1725 $server, peer \in Servers$
 1726 $\}$
 1727 IN
 1728 $\forall \langle leader, term \rangle \in Servers \times terms :$
 1729 LET
 1730 $leaders \triangleq \{$
 1731 $server \in Servers :$
 1732 $\wedge currTerm[server][leader] = term$
 1733 $\wedge currState[server][leader] = \text{"leader"}$
 1734 $\}$
 1735 IN
 1736 $Cardinality(leaders) \leq 1$

1738 *RecoveryLeaderDoNotUseFastPhase* \triangleq

Invariant: recovery leader never submit replicate commands in the fast phase. Commands in the fast phase discovered during its election are kept or nullified but, in either case, they are promoted to the slow phase for the next round of consensus.

1745 $\forall server, leader \in Servers :$
 1746 $server \neq leader \wedge currState[server][leader] = \text{"leader"} \Rightarrow$

1747 $\forall \text{entry} \in \text{SetOf}(\text{log}[\text{server}][\text{leader}]) :$
 1748 $\text{entry.phase} \neq \text{"fast"}$

1750 $\text{NoDivergingRecoveryReplies} \triangleq$

Invariant: candidates only consider one vote per election participant.

1755 $\forall \text{server}, \text{leader} \in \text{Servers} :$
 1756 $\text{currState}[\text{server}][\text{leader}] = \text{"leader"} \Rightarrow$
 1757 $\forall \text{peer} \in \text{Servers} :$
 1758 LET
 1759 $\text{replies} \triangleq \{$
 1760 $\text{reply} \in \text{recovery}[\text{server}][\text{leader}].\text{replies} :$
 1761 $\text{reply.reporter} = \text{peer}$
 1762 $\}$
 1763 IN
 1764 $\text{Cardinality}(\text{replies}) \leq 1$

LIVENESS PROPERTIES

1770 $\text{LeaderReplicatesItsLog} \triangleq$

In a scenario with bounded terms and log lengths, leaders eventually finish consensus on every command.

1775 $\diamond \square \forall \text{leader} \in \text{Servers} :$
 1776 $\wedge \text{commitIdx}[\text{leader}][\text{leader}] = \text{Len}(\text{log}[\text{leader}][\text{leader}])$
 1777 $\wedge \forall \text{follower} \in \text{Servers} :$
 1778 $\wedge \text{commitIdx}[\text{follower}][\text{leader}] = \text{commitIdx}[\text{leader}][\text{leader}]$
 1779 $\wedge \forall \text{entry} \in \text{SetOf}(\text{log}[\text{leader}][\text{leader}]) :$
 1780 $\wedge \text{log}[\text{leader}][\text{follower}][\text{entry.pos}].\text{cmd} = \text{entry.cmd}$
 1781 $\wedge \text{log}[\text{leader}][\text{follower}][\text{entry.pos}].\text{deps} = \text{entry.deps}$

1783

— SPEC BEHAVIOR —

1788 $\text{Init} \triangleq$

1789 $\wedge \text{InitCurrTerm}$
 1790 $\wedge \text{InitCurrState}$
 1791 $\wedge \text{InitLog}$
 1792 $\wedge \text{InitNextIdx}$
 1793 $\wedge \text{InitMatchIndexes}$
 1794 $\wedge \text{InitCommitIdx}$
 1795 $\wedge \text{InitFastQuorum}$
 1796 $\wedge \text{InitVotedFor}$
 1797 $\wedge \text{InitRecovery}$
 1798 $\wedge \text{InitMessages}$

1800 $\text{AppendCommands} \triangleq$

1801 $\exists \langle \text{leader}, \text{cmd} \rangle \in \text{Servers} \times \text{Commands} :$

```

1802     AppendCommandToLog(leader, cmd)

1804 ReplicateCommands  $\triangleq$ 
1805    $\exists$  server, leader  $\in$  Servers :
1806      $\vee$  SendAppendRequest(server, leader)
1807      $\vee$  CommitReplicatedCommands(server, leader)

1809 RecoverFailures  $\triangleq$ 
1810    $\exists$  server, leader  $\in$  Servers :
1811      $\vee$  SuspectFailure(server, leader)
1812      $\vee$  ReconfigureFastQuorum(server, leader)
1813      $\vee$  SendRecoveryRequest(server, leader)
1814      $\vee$  ClaimLeadership(server, leader)
1815      $\vee$  AbdicateLeadership(server, leader)

1817 HandleNetworkMessages  $\triangleq$ 
1818    $\exists$  msg  $\in$  messages :
1819     CASE msg.type = "append-request"  $\rightarrow$  HandleAppendRequest(msg)
1820      $\square$  msg.type = "append-reply"  $\rightarrow$  HandleAppendReply(msg)
1821      $\square$  msg.type = "recovery-request"  $\rightarrow$  HandleRecoveryRequest(msg)
1822      $\square$  msg.type = "recovery-reply"  $\rightarrow$  HandleRecoveryReply(msg)

1824 Next  $\triangleq$ 
1825    $\vee$  AppendCommands
1826    $\vee$  ReplicateCommands
1827    $\vee$  RecoverFailures
1828    $\vee$  HandleNetworkMessages
1829    $\vee$  UNCHANGED vars

1831 Fairness  $\triangleq$ 
1832    $\wedge$  WFvars(ReplicateCommands)
1833    $\wedge$  WFvars(HandleNetworkMessages)

1835 Spec  $\triangleq$ 
1836    $\wedge$  Init
1837    $\wedge$   $\square$ [Next]vars
1838    $\wedge$  Fairness
1839 ]

```

* Modification History
* Last modified Wed Jun 07 21:20:36 BRT 2023 by erickpintor
* Created Thu Aug 12 10:51:52 BRT 2021 by erickpintor

```

1  |----- MODULE Utils -----|
   | Utils offer common functionality that is not part of TLA+ std libs. |
5  LOCAL INSTANCE Naturals
6  LOCAL INSTANCE Sequences
8  Max(a, b)  $\triangleq$ 
   | Given two numbers a and b, returns a iff a succeeds b. |
12 IF a > b THEN a ELSE b
14 SetOf(seq)  $\triangleq$ 
   | Given a sequence seq, return a set of all its elements. |
18 {seq[x] : x ∈ DOMAIN seq}
20 Last(seq)  $\triangleq$ 
   | Given a non-empty sequence seq, returns its last element. |
24 seq[Len(seq)]
26 TakeWhile(seq, P(-))  $\triangleq$ 
   | Given a sequence seq, returns the prefix of elements satisfying P. |
30 LET
31   RECURSIVE Take(-)
32   Take(s)  $\triangleq$ 
33     IF Len(s) > 0 ∧ P(Head(s))
34     THEN Head(s) ∘ Take(Tail(s))
35     ELSE ⟨ ⟩
36 IN
37   Take(seq)
39 Map(seq, F(-))  $\triangleq$ 
   | Given a sequence of elements, return a sequence of elements mapped to the function F. |
44 [i ∈ DOMAIN seq ↦ F(seq[i])]
46 ReplaceSlice(seq, prev, replacement)  $\triangleq$ 
   | Given a sequence and a previous index, replaces the sub sequence delimited from (prev, |
   | replacement |] with the given replacement. |
51 LET prefix  $\triangleq$  SubSeq(seq, 1, prev)
52     suffix  $\triangleq$  SubSeq(seq, prev + Len(replacement) + 1, Len(seq))
53 IN prefix ∘ replacement ∘ suffix
54 |-----|
   | * Modification History
   | * Last modified Tue Mar 21 12:30:14 BRT 2023 by erickpintor
   | * Created Mon Oct 04 07:27:56 BRT 2021 by erickpintor

```