
MODULE *FLAIR*

Formal specification for *FLAIR*.

EXTENDS *Naturals*, *FiniteSets*, *Sequences*, *TLC*, *Reals*

Constants

The set of replicas

CONSTANTS *Replicas*

Replicas states.

CONSTANTS *Follower*, *Leader*

Replicas running state

CONSTANTS *ReplicaUpState*, *ReplicaDownState*

CONSTANTS *SwitchIp*, *SwitchKGroupsNum*, *SwitchStateActive*, *SwitchStateInactive*

Ip address of the switch
Number of key groups

Message Types

CONSTANTS *ClientReadRequest*, *ClientWriteRequest*, *WriteResponse*, *ReadResponse*, *InternalReadRequest*, *InternalWriteRequest*, *AppendEntriesRequest*, *AppendEntriesResponse*

Read request type
Write request type
Write response type
read response type
Read request from the switch to a replica
Write request from the switch to a replica
A request from the leader to the replicas
to append a *log* entry
An append entry resposne from a replica
to the leader

The set of possible keys and values in a client request

CONSTANTS *ValueSpace*, *KeySpace*

Special reserved value

CONSTANTS *Nil*

Variables

Replica vars

VARIABLE <i>state</i> ,	The replica's state (Follower, or <i>Leader</i>).
<i>log</i> ,	Log of all committed and uncommitted operations
<i>commitIndex</i> ,	The index of the highest committed index in the <i>log</i>
<i>currentTerm</i> ,	Current term number
<i>isActive</i> ,	Is the replica up or down
<i>replicaSession</i>	The latest switch id seen by the leader

$replicaVars \triangleq \langle state, log, commitIndex, currentTerm, isActive, replicaSession \rangle$

Leader vars. The following variables are used only on leaders:

VARIABLE <i>nextIndex</i> ,	The next entry to send to each follower.
<i>matchIndex</i> ,	The entry index for which a follower <i>log</i> matches the leader <i>log</i> . This used to calculate commit index
<i>replicaKGroups</i>	A map from a key hash to the last received sequence number for the associated <i>KGroup</i>

$leaderVars \triangleq \langle nextIndex, matchIndex, replicaKGroups \rangle$

Switch vars

VARIABLES <i>switchKGourpArray</i> ,	Array to maintain information about each <i>KGroup</i>
<i>switchSeqNum</i> ,	The last sequence number assigned by the switch
<i>switchTermId</i> ,	The latest term number seen by the switch
<i>switchLeaderId</i> ,	Current leader id as seen by the switch
<i>session</i> ,	Current switch id
	unique value for each session
<i>switchState</i>	Is the switch up or down

$switchVars \triangleq \langle switchKGourpArray, switchState, switchTermId, switchLeaderId, switchSeqNum, session \rangle$

Messages variables

VARIABLE <i>messages</i> ,	messages among replicas
<i>msgsClientSwitch</i> ,	messages from the client to the switch
<i>msgsReplicasSwitch</i>	messages from between replicas and the switch

$msgsVars \triangleq \langle messages, msgsClientSwitch, msgsReplicasSwitch \rangle$

Proof vars, don't appear in implementations

VARIABLES *responsesToClient*

Set of all variables
 $vars \triangleq \langle msgsVars, replicaVars, leaderVars, switchVars, responsesToClient \rangle$

Helpers

Helper to *Append* an element to a set
 $AddToSet(set, element) \triangleq set \cup \{element\}$

Helper to remove an element to a set
 $RemoveFromSet(set, element) \triangleq set \setminus \{element\}$

Helper to return the minimum value from a set,
or undefined if the set is empty.
 $Min(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \leq y$

Helper to return the maximum value from a set,
or undefined if the set is empty.
 $Max(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \geq y$

Helper to choose a random element from a set
 $ChooseRandomly(set) \triangleq \text{CHOOSE } i \in set : \text{TRUE}$

Helper to return the index of the *KGroup* given a hash
 $getIndexFromHash(hash) \triangleq (hash \% SwitchKGroupsNum) + 1$

The set of all quorums. This just calculates simple majorities, but the only
important property is that every quorum overlaps with every other.
 $Quorum \triangleq \{i \in \text{SUBSET}(Replicas) : Cardinality(i) * 2 > Cardinality(Replicas)\}$

Helper to return the term of the last entry in a *log*,
or 0 if the *log* is empty.
 $LastTerm(xlog) \triangleq \text{IF } Len(xlog) = 0 \text{ THEN } 0 \text{ ELSE } xlog[Len(xlog)].term$

Helper to find list of replicas that agree on a *log* index
 $AgreeIndex(index, logs, leaderId) \triangleq$
 $\{leaderId\} \cup \{k \in Replicas :$
 $\quad \wedge k \neq leaderId$
 $\quad \wedge Len(logs[k]) \geq index$
 $\quad \wedge logs[k][index] = logs[leaderId][index]\}$

Helper to return all *log* entries for specific key
 $entriesForKey(s, key) \triangleq \{j \in \text{DOMAIN } log[s] : \wedge log[s][j].key = key\}$

return the max value of a set or *Nil* if the set is empty
 $indexOfLastEntry(entries) \triangleq \text{IF } Cardinality(entries) > 0$
 $\quad \text{THEN } Max(entries) \text{ ELSE } Nil$

helper function to get the *ID* of the leader

if there is a leader that is up, this function returns its *ID*

if there is no current leader, this function finds the next leader based on raft criteria and returns its *ID*

```

getLeaderId  $\triangleq$ 
  LET leadersList  $\triangleq$  { s ∈ DOMAIN state : state[s] =
    Leader ∧ isActive[s] = ReplicaUpState }
  runningReplicas  $\triangleq$  { s ∈ Replicas : isActive[s] = ReplicaUpState }
  replicasWithNonEmptyLog  $\triangleq$  { s ∈ runningReplicas : Len(log[s]) > 0 }
  latestTerm  $\triangleq$  Max({ log[s][Len(log[s])).term :
    s ∈ replicasWithNonEmptyLog } )
  replicasWithLatestTerm  $\triangleq$  { s ∈ replicasWithNonEmptyLog :
    log[s][Len(log[s])).term = latestTerm }
  lengthOfLargestLog  $\triangleq$  Max({ Len(log[s]) : s ∈ replicasWithLatestTerm } )
  replicasWithLongestLog  $\triangleq$  { s ∈ replicasWithLatestTerm :
    Len(log[s]) = lengthOfLargestLog }
  newLeaderId  $\triangleq$  IF Cardinality(replicasWithLongestLog) ≥ 1
    THEN CHOOSE i ∈ replicasWithLongestLog : TRUE
    ELSE IF Cardinality(runningReplicas) > 0
    THEN CHOOSE i ∈ runningReplicas : TRUE
    ELSE Nil
  IN IF Cardinality(leadersList) > 0
    THEN CHOOSE s ∈ leadersList : TRUE
    ELSE newLeaderId

```

Helper to add a message to a set of messages.

Send(*m*) \triangleq *messages'* = *AddToSet*(*messages*, *m*)

Helper to remove a message from a set of messages.

Used when a replica is done processing a *replicaVarsmessage*.

Discard(*m*) \triangleq *messages'* = *RemoveFromSet*(*messages*, *m*)

Messages schemes and data structures

```

createClnetReadRequest(key)  $\triangleq$ 
  [mtype      ↦ ClientReadRequest,
   mkey       ↦ key,
   mhash      ↦ key]

```

```

createClnetWriteRequest(key, value)  $\triangleq$ 
  [mtype      ↦ ClientWriteRequest,
   mkey       ↦ key,
   mvalue     ↦ value,
   mhash      ↦ key]

```

When the switch forwards a read request, it adds:

- 1 – *logIndex*: The replica should *executereplicaVars* this index before serving the request
- 2 – *seqNum*: Which is used to ensure the linearizability of the request response

$createInternalReadRequest(msg, KGroup, forwardTo) \triangleq$

<i>mtype</i>	$\mapsto InternalReadRequest,$
<i>mkey</i>	$\mapsto msg.mkey,$
<i>mhash</i>	$\mapsto msg.mhash,$
<i>msession</i>	$\mapsto session,$
<i>mterm</i>	$\mapsto switchTermId,$
<i>mleaderId</i>	$\mapsto switchLeaderId,$
<i>mlogIndex</i>	$\mapsto KGroup.logIndex,$
<i>mkGroupSeqNum</i>	$\mapsto KGroup.seqNum,$
<i>msource</i>	$\mapsto SwitchIp,$
<i>mdest</i>	$\mapsto forwardTo]$

When the switch forwards a write request, it adds a sequence number that is used to order write requests.

The leader processes the requests in order.

$createInternalWriteRequest(msg, KGroup) \triangleq$

<i>mtype</i>	$\mapsto InternalWriteRequest,$
<i>mkey</i>	$\mapsto msg.mkey,$
<i>mvalue</i>	$\mapsto msg.mvalue,$
<i>mhash</i>	$\mapsto msg.mhash,$
<i>msession</i>	$\mapsto session,$
<i>mterm</i>	$\mapsto switchTermId,$
<i>mleaderId</i>	$\mapsto switchLeaderId,$
<i>mkGroupSeqNum</i>	$\mapsto KGroup.seqNum,$
<i>msource</i>	$\mapsto SwitchIp,$
<i>mdest</i>	$\mapsto switchLeaderId]$

$createReadResponse(m, i, leaderId, value, status, logIndex) \triangleq$

<i>mtype</i>	$\mapsto ReadResponse,$
<i>mkey</i>	$\mapsto m.mkey,$
<i>mvalue</i>	$\mapsto value,$
<i>mhash</i>	$\mapsto m.mhash,$
<i>mstatus</i>	$\mapsto status,$
<i>mlogIndex</i>	$\mapsto logIndex,$
<i>mkGroupSeqNum</i>	$\mapsto m.mkGroupSeqNum,$
<i>mterm</i>	$\mapsto currentTerm[i],$
<i>mleaderId</i>	$\mapsto leaderId,$
<i>mllLogs</i>	$\mapsto log,$ for correctness check only
<i>mcommitIndex</i>	$\mapsto commitIndex,$
<i>msession</i>	$\mapsto m.msession,$

$$\begin{array}{ll} msource & \mapsto i, \\ mdest & \mapsto SwitchIp] \end{array}$$

$$\begin{array}{l} createWriteResponse(i, logEntry, index, status, replicaIds) \triangleq \\ [mtype \mapsto WriteResponse, \\ mkey \mapsto logEntry.key, \\ mvalue \mapsto logEntry.value, \\ mhash \mapsto logEntry.hash, \\ mstatus \mapsto status, \\ mlogIndex \mapsto index, \\ mkGroupSeqNum \mapsto logEntry.seqNum, \\ msession \mapsto logEntry.switchId, \\ mreplicaIds \mapsto replicaIds, \\ mterm \mapsto currentTerm[i], \\ mallLogs \mapsto log, \text{ for correctness check only} \\ mcommitIndex \mapsto commitIndex, \\ msource \mapsto i, \\ mdest \mapsto SwitchIp] \end{array}$$

$$\begin{array}{l} createKGroup(leaderAcked, replicasIds, seqNum, logIndex) \triangleq \\ [leaderAcked \mapsto leaderAcked, \\ replicasIds \mapsto replicasIds, \\ seqNum \mapsto seqNum, \\ logIndex \mapsto logIndex] \end{array}$$

$$\begin{array}{l} createLogEntry(i, msg) \triangleq \\ [term \mapsto currentTerm[i], \\ key \mapsto msg.mkey, \\ value \mapsto msg.mvalue, \\ hash \mapsto msg.mhash, \\ seqNum \mapsto msg.mkGroupSeqNum, \\ switchId \mapsto msg.msession] \end{array}$$

$$\begin{array}{l} createResHistoryEntry(msg, tag) \triangleq \\ [msg \mapsto msg, \\ switchKGroupEntry \mapsto switchKGourpArray[getIndexFromHash(msg.mhash)], \\ tag \mapsto tag] \end{array}$$

Variables initialization

$$\begin{array}{l} InitSwitchVars \triangleq \\ \text{Initially the switch is inactive and} \end{array}$$

$KGroupArray$ is stable
 $\wedge switchKGroupArray = [i \in 1 \dots SwitchKGroupsNum$
 $\quad \mapsto createKGroup(TRUE, \{\}, Nil, Nil)]$
 $\wedge switchState = SwitchStateInactive$
 $\wedge switchTermId = 0$
 $\wedge switchLeaderId = Nil$
 $\wedge switchSeqNum = 0$
 $\wedge session = 0$

$InitMsgsSets \triangleq$
 $\wedge msgsClientSwitch = \{\}$
 $\wedge msgsReplicasSwitch = \{\}$
 $\wedge messages = \{\}$

$InitReplicaVars \triangleq$
 $\wedge currentTerm = [i \in Replicas \mapsto 1]$
 $\wedge state = [i \in Replicas \mapsto Follower]$
 $\wedge isActive = [i \in Replicas \mapsto ReplicaUpState]$
 $\wedge replicaSession = [i \in Replicas \mapsto 0]$
 $\wedge log = [i \in Replicas \mapsto \langle \rangle]$
 $\wedge commitIndex = [i \in Replicas \mapsto 0]$

$InitLeaderVars \triangleq$
 $\wedge nextIndex = [i \in Replicas \mapsto [j \in Replicas \mapsto 1]]$
 $\wedge matchIndex = [i \in Replicas \mapsto [j \in Replicas \mapsto 0]]$
 $\wedge replicaKGroups = [i \in Replicas \mapsto$
 $\quad [j \in 1 \dots SwitchKGroupsNum \mapsto 0]]$

$Init \triangleq \wedge InitReplicaVars$
 $\wedge InitLeaderVars$
 $\wedge InitSwitchVars$
 $\wedge InitMsgsSets$
 $\quad \text{Needed to prove safety}$
 $\wedge responsesToClient = \{\}$

Variables actions

Client actions

$\text{client sends a read request to read key } k$
 $IssueReadRequest(key) \triangleq$
 $\wedge LET request \triangleq createClientReadRequest(key)$
 $\quad IN msgsClientSwitch' = AddToSet(msgsClientSwitch, request)$
 $\wedge UNCHANGED \langle messages, replicaVars, leaderVars,$
 $\quad switchVars, msgsReplicasSwitch, responsesToClient \rangle$

client issues a write request to update key k
 $IssueWriteRequest(key, value) \triangleq$
 $\wedge LET request \triangleq createClientWriteRequest(key, value)$
 $IN msgsClientSwitch' = AddToSet(msgsClientSwitch, request)$
 $\wedge UNCHANGED \langle messages, replicaVars, leaderVars,$
 $switchVars, msgsReplicasSwitch, responsesToClient \rangle$

Switch actions

Switch state changes from Active to inactive
 $switchFails \triangleq$
 $\wedge switchState = SwitchStateActive$
 $\wedge switchState' = SwitchStateInactive$
 $\wedge UNCHANGED \langle replicaVars, leaderVars, msgsVars, responsesToClient,$
 $switchSeqNum, session, switchKGourpArray,$
 $switchLeaderId, switchTermId \rangle$

Switch handles a read request from a client.
The request has a hash that is mapped to a
 $KGroup$. If the $KGroup$ is stable, the request will
be forwarded to one of the replicas, otherwise it
will be forwarded to the leader.

$SwitchHandleClientRead(msg) \triangleq$
 $LET kGroup \triangleq switchKGourpArray[getIndexFromHash(msg.mhash)]$
 $forwardTo \triangleq IF kGroup.leaderAcked \wedge$
 $kGroup.seqNum \neq Nil$
 $THEN ChooseRandomly(\{x \in kGroup.replicasIds :$
 $x \neq switchLeaderId\})$
 $ELSE IF kGroup.leaderAcked \wedge$
 $kGroup.seqNum = Nil$
 $THEN switchLeaderId$
 $ELSE switchLeaderId$
 $internalMsg \triangleq createInternalReadRequest(msg, kGroup, forwardTo)$
 IN
 $\wedge msgsReplicasSwitch' = AddToSet(msgsReplicasSwitch, internalMsg)$
 $\wedge UNCHANGED \langle replicaVars, leaderVars, switchVars,$
 $responsesToClient, messages, msgsClientSwitch \rangle$

Switch handles a write request from a client
 $SwitchHandleClientWrite(msg) \triangleq$
 $Mark the KGroup associated with the key as unstable$
 $LET kGroup \triangleq switchKGourpArray[getIndexFromHash(msg.mhash)]$

$$\begin{aligned}
\text{updatedKGroup} &\triangleq [kGroup \text{ EXCEPT } \text{!.leaderAcked} = \text{FALSE}, \\
&\quad \text{!.seqNum} = \text{switchSeqNum} + 1, \\
&\quad \text{!.replicasIds} = \{\}, \\
&\quad \text{!.logIndex} = \text{Nil}] \\
\text{internalMsg} &\triangleq \text{createInternalWriteRequest}(\text{msg}, \text{updatedKGroup})
\end{aligned}$$

IN

$$\begin{aligned}
&\wedge \text{msgsReplicasSwitch}' = \text{AddToSet}(\text{msgsReplicasSwitch}, \text{internalMsg}) \\
&\wedge \text{switchKGourpArray}' = [\text{switchKGourpArray} \text{ EXCEPT} \\
&\quad \text{![getIndexFromHash}(\text{msg.mhash})] = \text{updatedKGroup}] \\
&\wedge \text{switchSeqNum}' = \text{switchSeqNum} + 1 \quad \text{Increments the sequence number} \\
&\wedge \text{UNCHANGED } \langle \text{replicaVars}, \text{leaderVars}, \text{responsesToClient}, \\
&\quad \text{messages}, \text{msgsClientSwitch}, \text{switchState}, \\
&\quad \text{switchTermId}, \text{switchLeaderId}, \text{session} \rangle
\end{aligned}$$

Switch receives a read or write request from a client

$$\begin{aligned}
\text{SwitchReceiveFromClient} &\triangleq \\
&\wedge \text{switchState} = \text{SwitchStateActive} \\
&\wedge \text{Cardinality}(\text{msgsClientSwitch}) > 0 \\
&\wedge \text{LET } \text{msg} \triangleq \text{ChooseRandomly}(\text{msgsClientSwitch}) \\
&\quad \text{type} \triangleq \text{msg.mtype} \\
\text{IN } &\vee \wedge \text{type} = \text{ClientReadRequest} \\
&\quad \wedge \text{SwitchHandleCleintRead}(\text{msg}) \\
&\vee \wedge \text{type} = \text{ClientWriteRequest} \\
&\quad \wedge \text{SwitchHandleCleintWrite}(\text{msg})
\end{aligned}$$

Switch handles a read response

$$\begin{aligned}
\text{SwitchHandleReadResponse}(\text{msg}) &\triangleq \\
&\wedge \vee \wedge \text{msg.msource} \neq \text{switchLeaderId} \quad \text{msg is from follower} \\
&\wedge \text{msg.mterm} = \text{switchTermId} \quad \text{msg.term} = \text{switch term} \\
&\wedge \text{msg.mstatus} = \text{TRUE} \quad \text{The operation was succeeded} \\
&\quad \text{Map the key to a KGroup based on hash.} \\
&\quad \text{The response will be sent to teh client if} \\
&\quad 1 - \text{msg.seqNum} = \text{KGroup.seqNum} \\
&\quad 2 - \text{KGroup is stable} \\
&\quad \text{otherwise the request will be dropped} \\
&\wedge \text{LET } \text{KGroup} \triangleq \text{switchKGourpArray}[\text{getIndexFromHash}(\text{msg.mhash})] \\
&\quad \text{isSeqOk} \triangleq \text{KGroup.seqNum} = \text{msg.mkGroupSeqNum} \\
\text{IN } &\vee \wedge \text{isSeqOk} \\
&\quad \wedge \text{KGroup.leaderAcked} \\
&\quad \wedge \text{responsesToClient}' = \text{AddToSet}(\text{responsesToClient}, \\
&\quad \quad \text{createResHistoryEntry}(\text{msg}, \text{Nil}))
\end{aligned}$$

$$\begin{aligned}
& \text{msg.seqNum!} = \text{KGroup.seqNum} \\
\vee \wedge \neg \text{isSeqOk} \\
& \wedge \text{responsesToClient}' = \text{AddToSet}(\text{responsesToClient}, \\
& \quad \text{createResHistoryEntry}(\text{msg}, \text{Nil})) \\
& \wedge \text{UNCHANGED} \langle \text{switchKGourpArray} \rangle \\
\vee \wedge \text{msg.msource} \neq \text{switchLeaderId} \\
& \wedge \text{UNCHANGED} \langle \text{switchKGourpArray}, \text{responsesToClient} \rangle \\
\vee \wedge \text{msg.mstatus} = \text{FALSE} \\
& \wedge \text{UNCHANGED} \langle \text{switchKGourpArray}, \text{responsesToClient} \rangle \\
\wedge \text{UNCHANGED} \langle \text{replicaVars}, \text{leaderVars}, \text{msgsVars}, \\
& \quad \text{switchSeqNum}, \text{session}, \text{switchLeaderId}, \\
& \quad \text{switchTermId}, \text{switchState} \rangle
\end{aligned}$$

Switch receives a message from a replica
 $\text{SwitchReceiveFromReplica}(\text{msg}) \triangleq$
 msg term id is larger than switch term id
 \Rightarrow switch stops processing request by setting its status to inactive

$$\begin{aligned}
\vee \wedge \text{switchState} &= \text{SwitchStateActive} \\
& \wedge \text{msg.msession} = \text{session} \\
& \wedge \text{msg.mterm} > \text{switchTermId} \\
& \wedge \text{switchState}' = \text{SwitchStateInactive} \\
& \wedge \text{UNCHANGED} \langle \text{replicaVars}, \text{leaderVars}, \text{msgsVars}, \\
& \quad \text{msgsClientSwitch}, \text{switchVars}, \text{responsesToClient} \rangle
\end{aligned}$$

msg is coming from an old leader
 \Rightarrow switch just ignore the message

$$\begin{aligned}
\vee \wedge \text{switchState} &= \text{SwitchStateActive} \\
& \wedge \text{msg.msession} = \text{session} \\
& \wedge \text{msg.mterm} < \text{switchTermId} \\
& \wedge \text{UNCHANGED} \langle \text{replicaVars}, \text{leaderVars}, \text{msgsVars}, \\
& \quad \text{responsesToClient}, \text{switchVars} \rangle
\end{aligned}$$

msg.switchId does not match switchId
 \Rightarrow switch just ignore the message

$$\begin{aligned}
\vee \wedge \text{switchState} &= \text{SwitchStateActive} \\
& \wedge \text{msg.msession} \neq \text{session} \\
& \wedge \text{UNCHANGED} \langle \text{replicaVars}, \text{leaderVars}, \text{msgsVars}, \\
& \quad \text{responsesToClient}, \text{switchVars} \rangle
\end{aligned}$$

Switch is active and the read response passes the safety check
 \Rightarrow switch processes the read response

$$\begin{aligned}
\vee \wedge \text{switchState} &= \text{SwitchStateActive} \\
& \wedge \text{msg.msession} = \text{session} \\
& \wedge \text{msg.mtype} = \text{ReadResponse} \\
& \wedge \text{SwitchHandleReadResponse}(\text{msg})
\end{aligned}$$

Switch is active and the write response passes the safety check
 \Rightarrow switch processes the write response
 $\vee \wedge \text{switchState} = \text{SwitchStateActive}$
 $\wedge \text{msg.msession} = \text{session}$
 $\wedge \text{msg.mtype} = \text{WriteResponse}$
 $\wedge \text{SwitchHandleWriteResponse}(\text{msg})$

Replica actions

Replica i fails and stops processing msgs .
 It loses everything but its currentTerm and log .

$\text{Stop}(i) \triangleq$
 $\wedge \text{isActive}[i] = \text{ReplicaUpState}$
 $\wedge \text{isActive}' = [\text{isActive} \text{ EXCEPT } ![i] = \text{ReplicaDownState}]$
 $\wedge \text{UNCHANGED } \langle \text{leaderVars}, \text{msgsVars}, \text{switchVars}, \text{responsesToClient},$
 $\text{currentTerm}, \text{log}, \text{commitIndex},$
 $\text{state}, \text{replicaSession} \rangle$

Replica i becomes active. The replica starts as follower

$\text{Start}(i) \triangleq$
 $\wedge \text{isActive}[i] = \text{ReplicaDownState}$
 $\wedge \text{isActive}' = [\text{isActive} \text{ EXCEPT } ![i] = \text{ReplicaUpState}]$
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}]$
 $\wedge \text{nextIndex}' = [\text{nextIndex} \text{ EXCEPT } ![i] = [j \in \text{Replicas} \mapsto 1]]$
 $\wedge \text{matchIndex}' = [\text{matchIndex} \text{ EXCEPT } ![i] = [j \in \text{Replicas} \mapsto 0]]$
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = 0]$
 $\wedge \text{replicaKGroups}' = [\text{replicaKGroups} \text{ EXCEPT } ![i] =$
 $\quad [j \in 1 \dots \text{SwitchKGroupsNum} \mapsto 0]]$
 $\wedge \text{UNCHANGED } \langle \text{msgsVars}, \text{switchVars}, \text{responsesToClient},$
 $\text{currentTerm}, \text{log}, \text{replicaSession} \rangle$

Select a new leader if non of the replicas is a leader

This helper selects the new leader based on raft criteria.

That is, the node with the log with the highest term id and
 longest log , if mutiple nodes have the same log id.

$\text{ElectLeader} \triangleq$
 $\wedge \text{LET } \text{runningReplicas} \triangleq \{s \in \text{Replicas} : \text{isActive}[s] = \text{ReplicaUpState}\}$
 $\text{replicasWithNonEmptyLog} \triangleq \{s \in \text{runningReplicas} : \text{Len}(\text{log}[s]) > 0\}$
 $\text{latestTerm} \triangleq \text{Max}(\{\text{log}[s][\text{Len}(\text{log}[s])].\text{term} :$
 $\quad s \in \text{replicasWithNonEmptyLog}\})$
 $\text{replicasWithLatestTerm} \triangleq \{s \in \text{replicasWithNonEmptyLog} :$
 $\quad \text{log}[s][\text{Len}(\text{log}[s])].\text{term} = \text{latestTerm}\}$
 $\text{lengthOfLargestLog} \triangleq \text{Max}(\{\text{Len}(\text{log}[s]) :$
 $\quad s \in \text{replicasWithLatestTerm}\})$

$$\begin{aligned}
& \text{replicasWithLongestLog} \triangleq \{s \in \text{replicasWithLatestTerm} : \\
& \quad \text{Len}(\log[s]) = \text{lengthOfLargestLog}\} \\
& \text{newLeaderId} \triangleq \text{IF } \text{Cardinality}(\text{replicasWithLongestLog}) \geq 1 \\
& \quad \text{THEN CHOOSE } i \in \text{replicasWithLongestLog} : \text{TRUE} \\
& \quad \text{ELSE IF } \text{Cardinality}(\text{runningReplicas}) > 0 \\
& \quad \text{THEN CHOOSE } i \in \text{runningReplicas} : \text{TRUE} \\
& \quad \text{ELSE Nil} \\
& \text{newTerm} \triangleq \text{IF } \text{newLeaderId} \neq \text{Nil} \\
& \quad \text{THEN } \text{currentTerm}[\text{newLeaderId}] + 1 \text{ ELSE Nil} \\
& \text{majority} \triangleq \text{IF } \text{newLeaderId} \neq \text{Nil} \\
& \quad \text{THEN CHOOSE } g \in \text{Quorum} : \text{newLeaderId} \in g \\
& \quad \text{ELSE Nil} \\
& \text{newCurrentTerm} \triangleq [j \in \text{Replicas} \mapsto \text{IF } j \in \text{majority} \\
& \quad \text{THEN } \text{newTerm} \\
& \quad \text{ELSE } \text{currentTerm}[j]] \\
& \text{IN } \wedge \text{Cardinality}(\text{runningReplicas}) * 2 > \text{Cardinality}(\text{Replicas}) \\
& \wedge \forall i \in \text{runningReplicas} : \text{state}[i] \in \{\text{Follower}\} \\
& \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![\text{newLeaderId}] = \text{Leader}] \\
& \wedge \text{nextIndex}' = [\text{nextIndex} \text{ EXCEPT } ![\text{newLeaderId}] = \\
& \quad [j \in \text{Replicas} \mapsto \text{Len}(\log[\text{newLeaderId}]) + 1]] \\
& \wedge \text{matchIndex}' = [\text{matchIndex} \text{ EXCEPT } ![\text{newLeaderId}] = \\
& \quad [j \in \text{Replicas} \mapsto 0]] \\
& \wedge \text{currentTerm}' = \text{newCurrentTerm} \\
& \wedge \text{UNCHANGED } \langle \text{switchVars}, \text{msgsVars}, \text{responsesToClient}, \\
& \quad \text{replicaKGroups}, \text{commitIndex}, \log, \\
& \quad \text{isActive}, \text{replicaSession} \rangle
\end{aligned}$$

leader i populates the switch $KGroup$ array
with information about unstable $KGroups$

$$\begin{aligned}
& \text{fillSwitchKGroup}(i) \triangleq \\
& \text{LET } \text{startIndex} \triangleq 1 \\
& \quad \text{endIndex} \triangleq \text{Len}(\log[i]) \\
& \quad \text{Scan the } \log \text{ and map each key in the } \log \\
& \quad \text{to its associated } KGroup \\
& \quad \text{keysToKGroups} \triangleq [j \in \text{startIndex} \dots \text{endIndex} \mapsto \\
& \quad \quad \text{getIndexFromHash}(\log[i][j].\text{hash})] \\
& \quad \text{For each } KGroup, \text{ find the last } \log \text{ entry} \\
& \quad \text{that should be used to update the switch} \\
& \quad KGroup \\
& \text{mapLogEntryToKGroupIndex} \triangleq \\
& \quad [j \in 1 \dots \text{SwitchKGroupsNum} \mapsto \\
& \quad \text{IF } \text{Cardinality}(\{k \in \text{DOMAIN } \text{keysToKGroups} : \text{keysToKGroups}[k] = j\}) > 0 \\
& \quad \text{THEN } \text{Max}(\{k \in \text{DOMAIN } \text{keysToKGroups} : \text{keysToKGroups}[k] = j\}) \\
& \quad \text{ELSE Nil}]
\end{aligned}$$

A boolean array that indicates whether
a *log* entry is committed or not

$$leaderAckedFlag \triangleq [j \in startIndex .. endIndex \mapsto j \leq commitIndex[i]]$$

Get the set of replicas that
acknowledged each *log* entry

$$keysToReplicas \triangleq [j \in startIndex .. endIndex \mapsto \begin{array}{l} \text{IF } leaderAckedFlag[j] \\ \text{THEN } AgreeIndex(j, log, i) \\ \text{ELSE } \{i\} \end{array}]$$

IN Update the switch *KGroup* Array to match the leader's
KGroup. If the leader do not have any writes for a
KGroup, then the *KGroup* is stable

$$\begin{aligned} \wedge switchKGroupArray' = \\ [j \in 1 .. SwitchKGroupsNum \mapsto \\ \text{IF } mapLogEntryToKGroupIndex[j] \neq Nil \\ \text{THEN } createKGroup(leaderAckedFlag[mapLogEntryToKGroupIndex[j]], \\ keysToReplicas[mapLogEntryToKGroupIndex[j]], \\ 0, mapLogEntryToKGroupIndex[j]) \\ \text{ELSE } createKGroup(TRUE, Replicas, Nil, Nil)] \end{aligned}$$

Leader *i* activates the switch

$$LeaderActivateSwitch(i) \triangleq$$

$$\begin{aligned} \wedge state[i] = Leader \quad & \text{Replica is the leader} \\ \wedge isActive[i] = ReplicaUpState \quad & \text{Replica is active} \\ \wedge switchState = SwitchStateInactive \quad & \text{Switch is inactive} \\ \wedge replicaSession' = [replicaSession \text{ EXCEPT } ![i] = session + 1] \\ \wedge session' = session + 1 \quad & \text{Update the replica and switch sessions} \\ \wedge switchTermId' = currentTerm[i] \quad & \text{Update switch term number} \\ \wedge switchLeaderId' = i \quad & \text{Update leader id of the switch} \\ \wedge switchSeqNum' = 0 \quad & \text{Reset the sequence number} \\ \wedge fillSwitchKGroup(i) \quad & \text{Populate switch } KGroup \text{ array} \\ \wedge switchState' = SwitchStateActive \quad & \text{activate the switch} \\ \wedge \text{UNCHANGED } \langle leaderVars, msgsVars, responsesToClient, \\ state, log, commitIndex, currentTerm, \\ isActive \rangle \end{aligned}$$

Any *RPC* with a newer term causes the recipient to advance its term first.

$$UpdateTerm(i, j, m) \triangleq$$

$$\begin{aligned} \wedge m.mterm > currentTerm[i] \\ \wedge currentTerm' &= [currentTerm \text{ EXCEPT } ![i] = m.mterm] \\ \wedge state' &= [state \text{ EXCEPT } ![i] = Follower] \\ &\text{messages is unchanged so } m \text{ can be processed further.} \end{aligned}$$

\wedge UNCHANGED $\langle \text{switchVars}, \text{msgsVars}, \text{leaderVars},$
 $\text{responsesToClient}, \text{isActive}, \text{log},$
 $\text{commitIndex}, \text{replicaSession} \rangle$

Responses with stale terms are ignored.

$\text{DropStaleResponse}(i, j, m) \triangleq$
 $\wedge m.\text{mterm} < \text{currentTerm}[i]$
 \wedge UNCHANGED $\langle \text{replicaVars}, \text{leaderVars}, \text{switchVars},$
 $\text{responsesToClient}, \text{msgsVars} \rangle$

Replica i receives a read request (m) from a client.

$\text{ReplicaReceiveReadRequest}(m, i) \triangleq$

The request received by the leader

$\vee \wedge \text{state}[i] = \text{Leader}$
 Check the message term numebr
 $\wedge m.\text{mterm} = \text{currentTerm}[i]$
 Get the index of the last committed entry
 $\wedge \text{LET } \text{committedEntries} \triangleq \{j \in \text{DOMAIN } \text{log}[i] : \wedge \text{log}[i][j].\text{key} = m.\text{mkey}$
 $\wedge j \leq \text{commitIndex}[i]\}$
 $\text{lastEntryIndex} \triangleq \text{IF } \text{Cardinality}(\text{committedEntries}) > 0$
 $\text{THEN } \text{Max}(\text{committedEntries}) \text{ ELSE } \text{Nil}$
 $\text{success} \triangleq \text{IF } \text{lastEntryIndex} = \text{Nil} \text{ THEN FALSE ELSE TRUE}$
 $\text{value} \triangleq \text{IF } \text{success} \text{ THEN } \text{log}[i][\text{lastEntryIndex}].\text{value} \text{ ELSE } \text{Nil}$
 IN $\wedge \text{msgsReplicasSwitch}' = \text{AddToSet}(\text{msgsReplicasSwitch},$
 $\text{createReadResponse}(m, i, \text{getLeaderId},$
 $\text{value}, \text{success}, \text{lastEntryIndex}))$
 \wedge UNCHANGED $\langle \text{replicaVars}, \text{leaderVars}, \text{switchVars},$
 $\text{msgsClientSwitch}, \text{messages},$
 $\text{responsesToClient} \rangle$

The request received by a follower and $\text{msg.mlogIndex} > 0$

i.e, the switch processed a write associated with the same

KGroup of the key

$\vee \wedge \text{state}[i] = \text{Follower}$
 $\wedge m.\text{mterm} = \text{currentTerm}[i]$ $\text{msg.term} = \text{replica term}$
 $\wedge m.\text{mlogIndex} \leq \text{Len}(\text{log}[i])$ The replica has the index
 $\wedge m.\text{mlogIndex} > 0$ Switch Kgroup entry is not empty
 Get the last committed log entry for the requested key
 $\wedge \text{LET } \text{logEntriesForKey} \triangleq \text{entriesForKey}(i, m.\text{mkey})$
 $\text{filteredEntries} \triangleq \{j \in \text{logEntriesForKey} : j \leq m.\text{mlogIndex}\}$
 $\text{lastEntryIndex} \triangleq \text{IF } \text{Cardinality}(\text{filteredEntries}) > 0$
 $\text{THEN } \text{Max}(\text{filteredEntries}) \text{ ELSE } \text{Nil}$
 $\text{requestedEntry} \triangleq \text{IF } \text{Cardinality}(\text{filteredEntries}) > 0$
 $\text{THEN } \text{log}[i][\text{lastEntryIndex}] \text{ ELSE } \text{Nil}$

$$\begin{aligned}
& isCommitted \triangleq m.mlogIndex \leq commitIndex[i] \\
& success \triangleq \text{IF } Cardinality(filteredEntries) > 0 \\
& \quad \text{THEN } requestedEntry.key = m.mkey \text{ ELSE } FALSE \\
& value \triangleq \text{IF } success \text{ THEN } requestedEntry.value \text{ ELSE } Nil \\
\text{IN } & \wedge msgsReplicasSwitch' = AddToSet(msgsReplicasSwitch, \\
& \quad createReadResponse(m, i, getLeaderId, \\
& \quad \quad value, success, lastEntryIndex)) \\
& \wedge \vee \wedge isCommitted \\
& \quad \wedge \text{UNCHANGED } \langle commitIndex \rangle \\
& \vee \wedge \neg isCommitted \\
& \quad \wedge success \\
& \quad \wedge commitIndex' = [commitIndex \text{ EXCEPT } ![i] = m.mlogIndex] \\
& \wedge \text{UNCHANGED } \langle leaderVars, switchVars, responsesToClient, \\
& \quad log, state, currentTerm, isActive, \\
& \quad msgsClientSwitch, messages, replicaSession \rangle
\end{aligned}$$

The request received by a follower and $msg.mlogIndex = -1$.

i.e., the switch did not process any write associated

$$\begin{aligned}
& \vee \wedge state[i] = Follower \\
& \wedge m.mterm = currentTerm[i] \quad msg.term = \text{replica term} \\
& \wedge m.mlogIndex = Nil \quad \text{Switch Kgroup entry is empty} \\
& \quad \text{Get the last committed log entry for the requested key} \\
& \wedge \text{LET } committedEntries \triangleq \{j \in \text{DOMAIN } log[i] : \wedge log[i][j].key = m.mkey \\
& \quad \wedge j \leq commitIndex[i]\} \\
& lastEntryIndex \triangleq \text{IF } Cardinality(committedEntries) > 0 \\
& \quad \text{THEN } Max(committedEntries) \text{ ELSE } Nil \\
& success \triangleq \text{IF } lastEntryIndex = Nil \text{ THEN } FALSE \text{ ELSE } TRUE \\
& value \triangleq \text{IF } success \text{ THEN } log[i][lastEntryIndex].value \text{ ELSE } Nil \\
\text{IN } & \wedge msgsReplicasSwitch' = AddToSet(msgsReplicasSwitch, \\
& \quad createReadResponse(m, i, \\
& \quad \quad getLeaderId, value, success, \\
& \quad \quad lastEntryIndex)) \\
& \wedge \text{UNCHANGED } \langle replicaVars, leaderVars, switchVars, \\
& \quad responsesToClient, msgsClientSwitch, \\
& \quad messages \rangle
\end{aligned}$$

Leader i receives a write request.

$ReplicaReceiveWriteRequest(m, i) \triangleq$

Safety checks

$$\begin{aligned}
& \vee \wedge m.mterm = currentTerm[i] \quad msg.term = \text{replica.term} \\
& \wedge m.mleaderId = i \quad \text{The leaderId field in the message is the replica} \\
& \wedge state[i] = Leader \quad \text{The replica is the leader} \\
& \wedge m.msession = replicaSession[i] \quad msg.session = \text{replica.session} \\
& \quad \text{Get the highest sequence number received for his KGroup}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{LET } latestSeenSeqNum \triangleq replicaKGroups[i][getIndexFromHash(m.mhash)] \\
& \quad entry \triangleq createLogEntry(i, m) \\
& \quad newLog \triangleq Append(log[i], entry) \\
& \quad \quad msg.seqNum > KGroup.seqNum \\
\text{IN } & \vee \wedge m.mkGroupSeqNum > latestSeenSeqNum \\
& \quad \quad \text{Update the KGroup seqNum} \\
& \quad \wedge replicaKGroups' = [replicaKGroups \text{ EXCEPT } ![i] = \\
& \quad \quad \quad [@ \text{ EXCEPT } ![getIndexFromHash(m.mhash)] = \\
& \quad \quad \quad \quad m.mkGroupSeqNum]] \\
& \quad \wedge log' = [log \text{ EXCEPT } ![i] = newLog] \quad \text{Append to log} \\
& \quad \wedge \text{UNCHANGED } \langle switchVars, msgsVars, responsesToClient, \\
& \quad \quad \quad nextIndex, matchIndex, commitIndex, \\
& \quad \quad \quad replicaSession, state, currentTerm, isActive \rangle \\
& \vee \wedge m.mkGroupSeqNum \leq latestSeenSeqNum \\
& \quad \wedge \text{UNCHANGED } \langle replicaVars, leaderVars, switchVars, \\
& \quad \quad \quad responsesToClient, msgsVars \rangle \\
& \vee \wedge \vee m.mterm \neq currentTerm[i] \\
& \quad \vee m.mleaderId \neq i \\
& \quad \vee state[i] \neq Leader \\
& \quad \vee m.msession \neq replicaSession[i] \\
& \wedge \text{UNCHANGED } \langle replicaVars, leaderVars, switchVars, \\
& \quad \quad \quad responsesToClient, msgsVars \rangle
\end{aligned}$$

Leader i sends follower j an *AppendEntries* request containing up to 1 entry.
While implementations may want to send more than 1 at a time, this spec uses
just 1 because it minimizes atomic regions without loss of generality.

$$\begin{aligned}
AppendEntries(i, j) & \triangleq \\
& \wedge i \neq j \quad \text{Avoid sending to itself} \\
& \wedge state[i] = Leader \quad \text{The sender is the leader} \\
& \wedge isActive[i] = ReplicaUpState \quad \text{The sender is active} \\
& \quad \text{Get the index of the next entry that must} \\
& \quad \text{be sent to the follower } j \\
& \wedge \text{LET } prevLogIndex \triangleq nextIndex[i][j] - 1 \\
& \quad prevLogTerm \triangleq \text{IF } prevLogIndex > 0 \text{ THEN} \\
& \quad \quad log[i][prevLogIndex].term \\
& \quad \quad \text{ELSE} \\
& \quad \quad 0 \\
& \quad \text{Send up to 1 entry, constrained by the end of the log.} \\
& \quad lastEntry \triangleq Min(\{Len(log[i]), nextIndex[i][j]\}) \\
& \quad entries \triangleq SubSeq(log[i], nextIndex[i][j], lastEntry) \\
& \quad m \triangleq \begin{array}{ll} [mtype & \mapsto AppendEntriesRequest, \\ mterm & \mapsto currentTerm[i], \\ mprevLogIndex & \mapsto prevLogIndex, \\ mprevLogTerm & \mapsto prevLogTerm, \end{array}
\end{aligned}$$

$mentries \mapsto entries,$
 $mlog$ is used as a history variable for the proof.
 It would not exist in a real implementation.
 $mlog \mapsto log[i],$
 $mcommitIndex \mapsto Min(\{commitIndex[i], lastEntry\}),$
 $msource \mapsto i,$
 $mdest \mapsto j]$

IN $\wedge Len(entries) > 0$
 $\wedge Send(m)$
 $\wedge UNCHANGED \langle replicaVars, leaderVars, switchVars,$
 $msgsClientSwitch, msgsReplicasSwitch,$
 $responsesToClient \rangle$

Leader i advances its $commitIndex$.
 This is done as a separate step from handling $AppendEntries$ responses,
 in part to minimize atomic regions, and in part so that leaders of
 single-server clusters are able to mark entries committed.

$AdvanceCommitIndex(i) \triangleq$
 $\wedge state[i] = Leader$
 $\wedge isActive[i] = ReplicaUpState$
 LET $\text{The set of replicas that agree up through an index.}$
 $Agree(index) \triangleq \{i\} \cup \{k \in Replicas : matchIndex[i][k] \geq index\}$
 $\text{The maximum indexes for which a quorum agrees}$
 $agreeIndexes \triangleq \{index \in 1 \dots Len(log[i]) : Agree(index) \in Quorum\}$
 $\text{Entries that are replicated to majorities but not yet committed}$
 $IndicesToCommit \triangleq \{index \in agreeIndexes : index > commitIndex[i]\}$
 $majoritiesPerIndex \triangleq [index \in agreeIndexes \mapsto Agree(index)]$
 $\text{New value for } commitIndex'[i]$
 $newCommitIndex \triangleq$
 IF $\wedge agreeIndexes \neq \{\}$
 $\wedge log[i][Max(agreeIndexes)].term = currentTerm[i]$
 THEN $Max(agreeIndexes)$
 ELSE $commitIndex[i]$
 IN $\wedge newCommitIndex > commitIndex[i]$
 $\text{For each entry that will be committed, send a write response to clients}$
 LET $indices \triangleq commitIndex[i] + 1 \dots newCommitIndex$
 $msgs \triangleq [x \in indices \mapsto createWriteResponse(i, log[i][x],$

$$\begin{aligned}
& x, \text{TRUE}, \text{majoritiesPerIndex}[x]] \\
\text{msgsAsSet} \triangleq & \{ \text{createWriteResponse}(i, \text{log}[i][x], \\
& x, \text{TRUE}, \text{majoritiesPerIndex}[x]) \\
& : x \in \text{indices} \}
\end{aligned}$$

$$\begin{aligned}
& \text{IN} \quad \wedge \text{msgsReplicasSwitch}' = \text{msgsReplicasSwitch} \cup \text{msgsAsSet} \\
& \text{Update the commit index at the replica} \\
& \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{newCommitIndex}] \\
& \wedge \text{UNCHANGED} \langle \text{leaderVars}, \text{switchVars}, \text{responsesToClient}, \\
& \quad \text{messages}, \text{msgsClientSwitch}, \text{log}, \text{currentTerm}, \\
& \quad \text{isActive}, \text{replicaSession}, \text{state} \rangle
\end{aligned}$$

Replica i receives an *AppendEntries* request from Replica j .
This just handles $m.\text{entries}$ of length 0 or 1, but
implementations could safely accept more by treating
them the same as multiple independent requests of 1 entry.

$$\begin{aligned}
& \text{HandleAppendEntriesRequest}(i, j, m) \triangleq \\
& \text{LET } \text{logOk} \triangleq \quad \vee m.\text{mprevLogIndex} = 0 \\
& \quad \vee \wedge m.\text{mprevLogIndex} > 0 \\
& \quad \quad \wedge m.\text{mprevLogIndex} \leq \text{Len}(\text{log}[i]) \\
& \quad \quad \wedge m.\text{mprevLogTerm} = \text{log}[i][m.\text{mprevLogIndex}].\text{term} \\
& \text{IN} \quad \wedge m.\text{mterm} \leq \text{currentTerm}[i] \\
& \quad \wedge \vee \wedge \text{reject request} \\
& \quad \quad \vee m.\text{mterm} < \text{currentTerm}[i] \\
& \quad \quad \vee \wedge m.\text{mterm} = \text{currentTerm}[i] \\
& \quad \quad \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \quad \quad \wedge \neg \text{logOk} \\
& \quad \wedge \text{Send}([\text{mtype} \quad \mapsto \text{AppendEntriesResponse}, \\
& \quad \quad \text{mterm} \quad \mapsto \text{currentTerm}[i], \\
& \quad \quad \text{msuccess} \quad \mapsto \text{FALSE}, \\
& \quad \quad \text{mmatchIndex} \quad \mapsto 0, \\
& \quad \quad \text{msource} \quad \mapsto i, \\
& \quad \quad \text{mdest} \quad \mapsto j]) \\
& \quad \wedge \text{UNCHANGED} \langle \text{replicaVars} \rangle \\
& \quad \text{process the request} \\
& \quad \vee \wedge m.\text{mterm} = \text{currentTerm}[i] \\
& \quad \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \quad \wedge \text{logOk} \\
& \quad \wedge \text{LET } \text{index} \triangleq m.\text{mprevLogIndex} + 1 \\
& \quad \quad \text{IN} \quad \vee \text{already done with request} \\
& \quad \quad \quad \wedge \vee m.\text{mentries} = \langle \rangle \\
& \quad \quad \quad \vee \wedge \text{Len}(\text{log}[i]) \geq \text{index} \\
& \quad \quad \quad \quad \wedge m.\text{mentries} \neq \langle \rangle \\
& \quad \quad \quad \quad \wedge \text{log}[i][\text{index}].\text{term} = m.\text{mentries}[1].\text{term} \\
& \quad \quad \quad \text{This could make our } \text{commitIndex} \text{ decrease (for}
\end{aligned}$$

example if we process an old, duplicated request),
but that doesn't really affect anything.

$$\begin{aligned}
& \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \\
& \quad \quad \quad m.m\text{commitIndex}] \\
& \wedge \text{Send}([mtype \quad \mapsto \text{AppendEntriesResponse}, \\
& \quad \quad mterm \quad \mapsto \text{currentTerm}[i], \\
& \quad \quad msuccess \quad \mapsto \text{TRUE}, \\
& \quad \quad mmatchIndex \mapsto m.mprevLogIndex + \\
& \quad \quad \quad \quad \quad \text{Len}(m.mentries), \\
& \quad \quad msource \quad \mapsto i, \\
& \quad \quad mdest \quad \mapsto j]) \\
& \wedge \text{UNCHANGED } \langle \log \rangle \\
\vee & \text{conflict: remove 1 entry} \\
& \wedge m.mentries \neq \langle \rangle \\
& \wedge \text{Len}(\log[i]) \geq \text{index} \\
& \wedge \log[i][\text{index}].term \neq m.mentries[1].term \\
& \wedge \text{LET } new \triangleq [\text{index2} \in 1 \dots (\text{Len}(\log[i]) - 1) \mapsto \\
& \quad \quad \quad \log[i][\text{index2}]] \\
& \quad \text{IN } \log' = [\log \text{ EXCEPT } ![i] = new] \\
& \wedge \text{Send}([mtype \quad \mapsto \text{AppendEntriesResponse}, \\
& \quad \quad mterm \quad \mapsto \text{currentTerm}[i], \\
& \quad \quad msuccess \quad \mapsto \text{FALSE}, \\
& \quad \quad mmatchIndex \mapsto 0, \\
& \quad \quad msource \quad \mapsto i, \\
& \quad \quad mdest \quad \mapsto j]) \\
& \wedge \text{UNCHANGED } \langle \text{commitIndex} \rangle \\
\vee & \text{no conflict: append entry} \\
& \wedge m.mentries \neq \langle \rangle \\
& \wedge \text{Len}(\log[i]) = m.mprevLogIndex \\
& \wedge \log' = [\log \text{ EXCEPT } ![i] = \\
& \quad \quad \quad \text{Append}(\log[i], m.mentries[1])] \\
& \wedge \text{Send}([mtype \quad \mapsto \text{AppendEntriesResponse}, \\
& \quad \quad mterm \quad \mapsto \text{currentTerm}[i], \\
& \quad \quad msuccess \quad \mapsto \text{TRUE}, \\
& \quad \quad mmatchIndex \mapsto m.mprevLogIndex + \\
& \quad \quad \quad \quad \quad \text{Len}(m.mentries), \\
& \quad \quad msource \quad \mapsto i, \\
& \quad \quad mdest \quad \mapsto j]) \\
& \wedge \text{UNCHANGED } \langle \text{commitIndex} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{leaderVars}, \text{switchVars}, \text{responsesToClient}, \\
& \quad \text{msgsClientSwitch}, \text{msgsReplicasSwitch}, \\
& \quad \text{isActive}, \text{currentTerm}, \text{state}, \text{replicaSession} \rangle
\end{aligned}$$

Replica i receives an *AppendEntries* response from Replica j

$$\begin{aligned}
& \text{HandleAppendEntriesResponse}(i, j, m) \triangleq \\
& \wedge m.mterm = \text{currentTerm}[i] \\
& \wedge \vee \wedge m.msucccess \text{ successful} \\
& \quad \wedge \text{nextIndex}' = [\text{nextIndex} \text{ EXCEPT } ![i][j] = m.mmatchIndex + 1] \\
& \quad \wedge \text{matchIndex}' = [\text{matchIndex} \text{ EXCEPT } ![i][j] = m.mmatchIndex] \\
& \vee \wedge \neg m.msucccess \text{ not successful} \\
& \quad \wedge \text{nextIndex}' = [\text{nextIndex} \text{ EXCEPT } ![i][j] = \\
& \quad \quad \quad \text{Max}(\{\text{nextIndex}[i][j] - 1, 1\})] \\
& \quad \wedge \text{UNCHANGED } \langle \text{matchIndex} \rangle \\
& \wedge \text{Discard}(m) \\
& \wedge \text{UNCHANGED } \langle \text{replicaVars}, \text{switchVars}, \text{responsesToClient}, \\
& \quad \text{msgsClientSwitch}, \text{msgsReplicasSwitch}, \text{replicaKGroups} \rangle
\end{aligned}$$

process a message. The message will be processed

by its recipient based on its type

$$\begin{aligned}
& \text{Receive}(m) \triangleq \\
& \text{LET } i \triangleq m.mdest \\
& \quad j \triangleq m.msource \\
& \text{IN} \quad \text{Any RPC with a newer term causes the recipient to advance} \\
& \quad \text{its term first. Responses with stale terms are ignored.} \\
& \wedge \text{isActive}[i] = \text{ReplicaUpState} \\
& \wedge \vee \wedge m.mtype \in \{\text{AppendEntriesRequest}, \text{AppendEntriesResponse}\} \\
& \quad \wedge \text{UpdateTerm}(i, j, m) \\
& \quad \text{Append entry request from replica } j \text{ to replica } i \\
& \quad \vee \wedge m.mtype = \text{AppendEntriesRequest} \\
& \quad \quad \wedge \text{HandleAppendEntriesRequest}(i, j, m) \\
& \quad \text{Append entry response from replica } j \text{ to replica } i \\
& \quad \vee \wedge m.mtype = \text{AppendEntriesResponse} \\
& \quad \quad \wedge \vee \text{DropStaleResponse}(i, j, m) \\
& \quad \quad \vee \text{HandleAppendEntriesResponse}(i, j, m) \\
& \quad \text{Read request from the switch to replica } i \\
& \quad \vee \wedge m.mtype = \text{InternalReadRequest} \\
& \quad \quad \wedge \text{ReplicaReceiveReadRequest}(m, i) \\
& \quad \text{Write request from the switch to replica } i \\
& \quad \vee \wedge m.mtype = \text{InternalWriteRequest} \\
& \quad \quad \wedge \text{ReplicaReceiveWriteRequest}(m, i)
\end{aligned}$$

$$\begin{aligned}
& \text{ReplicasReceiveRaftInternalMsgs} \triangleq \\
& \wedge \exists m \in \text{messages} : \text{Receive}(m)
\end{aligned}$$

$$\begin{aligned}
& \text{ReplicasReceiveFromSwitch} \triangleq \\
& \wedge \exists m \in \text{msgsReplicasSwitch} : \\
& \quad \wedge m.mdest \neq \text{SwitchIp} \wedge \text{Receive}(m)
\end{aligned}$$

Defines how all *variables* (*Replicas*, *Client*, *Switch*) may transition.

$Next \triangleq$

client transitions

$\wedge \vee \exists key \in KeySpace : IssueReadRequest(key)$

$\vee \exists key \in KeySpace : \exists value \in ValueSpace :$

$IssueWriteRequest(key, value)$

Switch transitions

$\vee switchFails$

$\vee SwitchReceiveFromClient$

$\vee \wedge Cardinality(msgsReplicasSwitch) > 0$

$\wedge LET messagesToSwitch \triangleq \{x \in msgsReplicasSwitch : x.mdest = SwitchIp\}$

$IN \quad \wedge \exists msg \in messagesToSwitch : SwitchReceiveFromReplica(msg)$

Replica transitions

$\vee \exists i \in Replicas : Stop(i)$

$\vee \exists i \in Replicas : Start(i)$

$\vee ElectLeader$

$\vee ReplicasReceiveRaftInternalMsgs$

$\vee ReplicasReceiveFromSwitch$

Leader transitions

$\vee \exists i \in Replicas : LeaderActivateSwitch(i)$

$\vee \exists i, j \in Replicas : AppendEntries(i, j)$

$\vee \exists i \in Replicas : AdvanceCommitIndex(i)$

Safety Invariants

In the following statements, each response has the logs of all replicas at the time it was served by a replica (This is only for safety check and should not be the case for real implementations)

Invariant that defines that read responses that passes the switch to the client are linearizable. Check the logs of all replica to get the last committed *log* index that update the requested key. The returned value “*msg.mvalue*” should equal the value in that *log* index.

$isForwardedToClientReadSafe(response) \triangleq$

$LET msg \triangleq response.msg$

$logEntriesForKey \triangleq entriesForKey(msg.mleaderId, msg.mkey)$

$$\begin{aligned}
committedEntries &\triangleq \{j \in logEntriesForKey : \\
&\quad AgreeIndex(j, msg.mallLogs, msg.mleaderId) \\
&\quad \in Quorum \wedge j \leq msg.mlogIndex\} \\
lastCommittedIndex &\triangleq indexOfLastEntry(committedEntries) \\
lastEntry &\triangleq log[msg.mleaderId][lastCommittedIndex]
\end{aligned}$$

IN

$$\begin{aligned}
&\vee \wedge msg.mstatus = \text{TRUE} \\
&\quad \wedge lastCommittedIndex = msg.mlogIndex \\
&\quad \wedge lastEntry.key = msg.mkey \\
&\quad \wedge lastEntry.value = msg.mvalue
\end{aligned}$$

Invariant that defines that write responses that
passes the switch to the client are committed
on majority of the replicas.

$$\begin{aligned}
isForwardedToClientWriteCorrect(response) &\triangleq \\
&\quad \text{LET } msg \triangleq response.msg \\
&\quad \quad isOnMajority \triangleq AgreeIndex(msg.mlogIndex, \\
&\quad \quad \quad msg.mallLogs, msg.msource) \in Quorum \\
\text{IN } isOnMajority
\end{aligned}$$

Invariant that defines the correctness
of all responses forwarded to the client

$$\begin{aligned}
InvResponsesToClientCorrectness &\triangleq \\
&\quad \forall res \in responsesToClient : \\
&\quad \quad \vee \wedge res.msg.mtype = ReadResponse \\
&\quad \quad \quad \wedge isForwardedToClientReadSafe(res) \\
&\quad \quad \vee \wedge res.msg.mtype = WriteResponse \\
&\quad \quad \quad \wedge isForwardedToClientWriteCorrect(res)
\end{aligned}$$

Invariant that defines that no two entries have
the same sequence number or log index.

$$\begin{aligned}
InvSwitchRegisterCorrectness &\triangleq \\
&\quad \wedge \forall i, j \in 1 \dots SwitchKGroupsNum : \\
&\quad \quad \vee \wedge i \neq j \\
&\quad \quad \quad \wedge switchKGourpArray[i].leaderAcked = \text{TRUE} \\
&\quad \quad \quad \wedge switchKGourpArray[j].leaderAcked = \text{TRUE} \\
&\quad \quad \quad \wedge switchKGourpArray[i].seqNum \neq switchKGourpArray[j].seqNum \\
&\quad \quad \quad \wedge switchKGourpArray[i].logIndex \neq switchKGourpArray[j].logIndex \\
&\quad \quad \vee i = j \\
&\quad \quad \vee switchKGourpArray[i].leaderAcked = \text{FALSE} \\
&\quad \quad \vee switchKGourpArray[j].leaderAcked = \text{FALSE}
\end{aligned}$$

invariant that defines that no two leaders exist
at the same time

$$InvLeaderElectionSafety \triangleq$$

LET $runningReplicas \triangleq \{i \in Replicas : isActive[i] = ReplicaUpState\}$
 $leaders \triangleq \{i \in runningReplicas : state[i] = Leader\}$
 IN $Cardinality(leaders) \leq 1$
