

EXTENDS *Integers*,
 FiniteSets

CONSTANTS *H_NODES*,
 H_MAX_VERSION

VARIABLES *msgs*,
 nodeTS,
 nodeState,
 nodeRcvdAcks,
 nodeLastWriter,
 nodeLastWriteTS,
 nodeWriteEpochID,
 aliveNodes,
 epochID

The consistent invariant: all alive nodes in valid state should have the same value / *TS*

$HConsistent \triangleq$

$$\forall k, s \in aliveNodes : \begin{aligned} & \vee nodeState[k] \neq \text{"valid"} \\ & \vee nodeState[s] \neq \text{"valid"} \\ & \vee nodeTS[k] = nodeTS[s] \end{aligned}$$

$HMessage \triangleq$ Messages exchanged by the Protocol

$$[type : \{ \text{"INV"}, \text{"ACK"} \}, sender : H_NODES, \\ epochID : 0 \dots (Cardinality(H_NODES) - 1), \\ version : 0 \dots H_MAX_VERSION, \\ tieBreaker : H_NODES]$$

∪

$$[type : \{ \text{"VAL"} \}, version : 0 \dots H_MAX_VERSION, \\ tieBreaker : H_NODES]$$

$HTypeOK \triangleq$ The type correctness invariant

$$\begin{aligned} & \wedge msgs \subseteq HMessage \\ & \wedge \forall n \in H_NODES : nodeRcvdAcks[n] \subseteq (H_NODES \setminus \{n\}) \\ & \wedge nodeLastWriter \in [H_NODES \rightarrow H_NODES] \\ & \wedge nodeLastWriteTS \in [H_NODES \rightarrow [version : 0 \dots H_MAX_VERSION, \\ & \hspace{10em} tieBreaker : H_NODES]] \\ & \wedge nodeTS \in [H_NODES \rightarrow [version : 0 \dots H_MAX_VERSION, \\ & \hspace{10em} tieBreaker : H_NODES]] \\ & \wedge nodeState \in [H_NODES \rightarrow \{ \text{"valid"}, \text{"invalid"}, \text{"invalid_write"}, \\ & \hspace{10em} \text{"write"}, \text{"replay"} \}] \end{aligned}$$

membership and epoch id related

$$\begin{aligned} & \wedge aliveNodes \subseteq H_NODES \\ & \wedge epochID \in 0 \dots (Cardinality(H_NODES) - 1) \end{aligned}$$

$$\wedge \text{ nodeWriteEpochID} \in [H_NODES \rightarrow 0 \dots (\text{Cardinality}(H_NODES) - 1)]$$

$$\begin{aligned}
HInit &\triangleq \text{The initial predicate} \\
&\wedge \text{ msgs} = \{\} \\
&\quad \text{membership and epoch id related} \\
&\wedge \text{ epochID} = 0 \\
&\wedge \text{ aliveNodes} = H_NODES \\
&\wedge \text{ nodeWriteEpochID} = [n \in H_NODES \mapsto 0] \\
&\quad \text{Init rest per node replica metadata} \\
&\wedge \text{ nodeRcvdAcks} = [n \in H_NODES \mapsto \{\}] \\
&\wedge \text{ nodeState} = [n \in H_NODES \mapsto \text{"valid"}] \\
&\wedge \text{ nodeLastWriter} = [n \in H_NODES \mapsto \text{CHOOSE } k \in H_NODES : \\
&\quad \quad \quad \forall m \in H_NODES : k \leq m] \\
&\wedge \text{ nodeTS} = [n \in H_NODES \mapsto [\text{version} \mapsto 0, \\
&\quad \quad \quad \text{tieBreaker} \mapsto \\
&\quad \quad \quad \text{CHOOSE } k \in H_NODES : \\
&\quad \quad \quad \forall m \in H_NODES : k \leq m]] \\
&\wedge \text{ nodeLastWriteTS} = [n \in H_NODES \mapsto [\text{version} \mapsto 0, \\
&\quad \quad \quad \text{tieBreaker} \mapsto \\
&\quad \quad \quad \text{CHOOSE } k \in H_NODES : \\
&\quad \quad \quad \forall m \in H_NODES : k \leq m]]
\end{aligned}$$

A buffer maintaining all network messages. Messages are only appended to this variable (not removed once delivered) intentionally to check protocols tolerance in duplicates and reorderings

$$\text{send}(m) \triangleq \text{msgs}' = \text{msgs} \cup \{m\}$$

Check if all acknowledgments for a write have been received

$$\text{receivedAllAcks}(n) \triangleq \text{nodeRcvdAcks}[n] = H_NODES \setminus \{n\}$$

$$\text{receivedAllAcks}(n) \triangleq (\text{aliveNodes} \setminus \{n\}) \subseteq \text{nodeRcvdAcks}[n]$$

$\text{equalTS}(v1, tb1, v2, tb2) \triangleq$ Timestamp equality

$$\begin{aligned}
&\wedge v1 = v2 \\
&\wedge tb1 = tb2
\end{aligned}$$

$\text{greaterTS}(v1, tb1, v2, tb2) \triangleq$ Timestamp comparison

$$\begin{aligned}
&\vee v1 > v2 \\
&\vee \wedge v1 = v2 \\
&\quad \wedge tb1 > tb2
\end{aligned}$$

$\text{isAlive}(n) \triangleq n \in \text{aliveNodes}$

$\text{nodeFailure}(n) \triangleq$ Emulate a node failure

Make sure that there are atleast 3 alive nodes before killing a node

$$\begin{aligned}
&\wedge \exists k, m \in \text{aliveNodes} : \wedge k \neq n \\
&\quad \wedge m \neq n
\end{aligned}$$

$HRead(n) \triangleq$ Execute a read
 $\wedge nodeState[n] = \text{"valid"}$
 $\wedge h_upd_nothing$

$HWrite(n) \triangleq$ Execute a write
 $\wedge nodeState[n] \in \{\text{"valid"}, \text{"invalid"}\}$
 writes in invalid state are also supported as an optimization
 $\wedge nodeState[n] \in \{\text{"valid"}\}$
 $\wedge nodeTS[n].version < H_MAX_VERSION$
 $\wedge h_upd_actions(n, nodeTS[n].version + 1, n, \text{"write"}, \{\})$

$HCoordWriteReplay(n) \triangleq$ Execute a write-replay after a membership re-config
 $\wedge nodeState[n] \in \{\text{"write"}, \text{"replay"}\}$
 $\wedge nodeWriteEpochID[n] < epochID$
 $\wedge \neg receivedAllAcks(n)$ optimization to not replay when we have gathered *acks* from all alive
 $\wedge h_upd_replay_actions(n, nodeRcvdAcks[n])$

$HRcvAck(n) \triangleq$ Process a received acknowledgment
 $\exists m \in msgs :$
 $\wedge m.type = \text{"ACK"}$
 $\wedge m.epochID = epochID$
 $\wedge m.sender \neq n$
 $\wedge m.sender \notin nodeRcvdAcks[n]$
 $\wedge equalTS(m.version,$
 $\quad m.tieBreaker,$
 $\quad nodeLastWriteTS[n].version,$
 $\quad nodeLastWriteTS[n].tieBreaker)$
 $\wedge nodeState[n] \in \{\text{"write"}, \text{"invalid_write"}, \text{"replay"}\}$
 $\wedge nodeRcvdAcks' = [nodeRcvdAcks \text{ EXCEPT } ![n] =$
 $\quad nodeRcvdAcks[n] \cup \{m.sender\}]$
 $\wedge \text{UNCHANGED } \langle msgs, nodeLastWriter, nodeLastWriteTS,$
 $\quad aliveNodes, nodeTS, nodeState, epochID, nodeWriteEpochID \rangle$

$HSendVals(n) \triangleq$ Send validations once received acknowledgments from all alive nodes
 $\wedge nodeState[n] \in \{\text{"write"}, \text{"replay"}\}$
 $\wedge receivedAllAcks(n)$
 $\wedge nodeState' = [nodeState \text{ EXCEPT } ![n] = \text{"valid"}]$
 $\wedge send([type \mapsto \text{"VAL"},$
 $\quad version \mapsto nodeTS[n].version,$
 $\quad tieBreaker \mapsto nodeTS[n].tieBreaker])$
 $\wedge \text{UNCHANGED } \langle nodeTS, nodeLastWriter, nodeLastWriteTS,$
 $\quad aliveNodes, nodeRcvdAcks, epochID, nodeWriteEpochID \rangle$

$HCoordinatorActions(n) \triangleq$ Actions of a read/write coordinator

$\vee HRead(n)$
 $\vee HCoordWriteReplay(n)$
 $\vee HWrite(n)$
 $\vee HRcvAck(n)$
 $\vee HSendVals(n)$

$HRcvInv(n) \triangleq$ Process a received invalidation
 $\exists m \in msgs :$
 $\wedge m.type = \text{"INV"}$
 $\wedge m.epochID = epochID$
 $\wedge m.sender \neq n$
always acknowledge a received invalidation (irrelevant to the timestamp)
 $\wedge send([type \mapsto \text{"ACK"},$
 $\quad sender \mapsto n,$
 $\quad epochID \mapsto epochID,$
 $\quad version \mapsto m.version,$
 $\quad tieBreaker \mapsto m.tieBreaker])$
 $\wedge \vee \wedge greaterTS(m.version,$
 $\quad m.tieBreaker,$
 $\quad nodeTS[n].version,$
 $\quad nodeTS[n].tieBreaker)$
 $\wedge nodeLastWriter' = [nodeLastWriter \text{ EXCEPT } ![n] = m.sender]$
 $\wedge nodeTS' = [nodeTS \text{ EXCEPT } ![n].version = m.version,$
 $\quad ![n].tieBreaker = m.tieBreaker]$
 $\wedge \vee \wedge nodeState[n] \in \{\text{"valid"}, \text{"invalid"}, \text{"replay"}\}$
 $\wedge nodeState' = [nodeState \text{ EXCEPT } ![n] = \text{"invalid"}]$
 $\vee \wedge nodeState[n] \in \{\text{"write"}, \text{"invalid_write"}\}$
 $\wedge nodeState' = [nodeState \text{ EXCEPT } ![n] = \text{"invalid_write"}]$
 $\vee \wedge \neg greaterTS(m.version,$
 $\quad m.tieBreaker,$
 $\quad nodeTS[n].version,$
 $\quad nodeTS[n].tieBreaker)$
 $\wedge \text{UNCHANGED } \langle nodeState, nodeTS, nodeLastWriter, nodeWriteEpochID \rangle$
 $\wedge \text{UNCHANGED } \langle nodeLastWriteTS, aliveNodes, nodeRcvdAcks, epochID, nodeWriteEpochID \rangle$

$HRcvVal(n) \triangleq$ Process a received validation
 $\exists m \in msgs :$
 $\wedge nodeState[n] \neq \text{"valid"}$
 $\wedge m.type = \text{"VAL"}$
 $\wedge equalTS(m.version,$
 $\quad m.tieBreaker,$
 $\quad nodeTS[n].version,$
 $\quad nodeTS[n].tieBreaker)$

$$\begin{aligned}
& \wedge nodeState' = [nodeState \text{ EXCEPT } ![n] = \text{"valid"}] \\
& \wedge \text{UNCHANGED } \langle msgs, nodeTS, nodeLastWriter, nodeLastWriteTS, \\
& \quad aliveNodes, nodeRcvdAcks, epochID, nodeWriteEpochID \rangle
\end{aligned}$$

$$\begin{aligned}
HFollowerWriteReplay(n) & \triangleq \text{Execute a write-replay when coordinator failed} \\
& \wedge nodeState[n] = \text{"invalid"} \\
& \wedge \neg isAlive(nodeLastWriter[n]) \\
& \wedge h_upd_replay_actions(n, \{\})
\end{aligned}$$

$$\begin{aligned}
HFollowerActions(n) & \triangleq \text{Actions of a write follower} \\
& \vee HRcvInv(n) \\
& \vee HFollowerWriteReplay(n) \\
& \vee HRcvVal(n)
\end{aligned}$$

$$\begin{aligned}
HNext & \triangleq \text{Modeling Hermes protocol (Coordinator and Follower actions while emulating failures)} \\
& \exists n \in aliveNodes : \\
& \quad \vee HFollowerActions(n) \\
& \quad \vee HCoordinatorActions(n) \\
& \quad \vee nodeFailure(n) \text{ emulate node failures}
\end{aligned}$$
