
MODULE *P4RuntimeRevision*

EXTENDS *Naturals, FiniteSets, Sequences, TLC*

The set of all *ONOS* nodes

CONSTANTS *Nodes*

Stream states

CONSTANTS *Open, Closed*

Write message types

CONSTANTS *WriteRequest, WriteResponse*

Response status constants

CONSTANTS *Ok, PermissionDenied*

Empty value

CONSTANT *Nil*

The current state of mastership elections

VARIABLES *term, master, backups*

The current mastership event queue for each node

VARIABLE *events*

The current mastership state for each node

VARIABLE *masterships*

The state of all streams and their requests and responses

VARIABLE *streams, requests, responses*

The term of the last successful write to the device

VARIABLE *lastTerm*

Counting variables used to enforce state constraints

VARIABLES *mastershipChanges, streamChanges, messageCount*

A history of successful writes to the switch used for model checking

VARIABLE *history*

Mastership/consensus related variables

$mastershipVars \triangleq \langle term, master, backups, mastershipChanges \rangle$

Node related variables

$nodeVars \triangleq \langle events, masterships \rangle$

Stream related variables

$streamVars \triangleq \langle streams, streamChanges \rangle$

Message related variables
 $messageVars \triangleq \langle requests, responses, messageCount \rangle$

Device related variables
 $deviceVars \triangleq \langle lastTerm, history \rangle$

A sequence of all variables
 $vars \triangleq \langle mastershipVars, nodeVars, streamVars, messageVars, deviceVars \rangle$

Helpers

Returns a sequence with the head removed
 $Pop(q) \triangleq SubSeq(q, 2, Len(q))$

Returns a sequences with the element at the given index removed
 $Drop(q, i) \triangleq SubSeq(q, 1, i - 1) \circ SubSeq(q, i + 1, Len(q))$

Returns the set of values in f
 $Range(f) \triangleq \{f[x] : x \in DOMAIN\ f\}$

Returns the maximum value from a set or undefined if the set is empty
 $Max(s) \triangleq CHOOSE\ x \in s : \forall y \in s : x \geq y$

This section models the messaging between controller nodes and the device. Messaging is modelled on *TCP*, providing strict ordering between controller and device via sequences. The 'requests' sequence represents the messages from controller to device for each node, and the 'responses' sequence represents the messages from device to each node. Requests and responses are always received from the head of the queue and are never duplicated or reordered.

Sends request 'm' on the stream for node 'n'
 $SendRequest(n, m) \triangleq$
 $\wedge requests' = [requests\ EXCEPT\ ![n] = Append(requests[n], m)]$
 $\wedge messageCount' = messageCount + 1$

Indicates whether a request of type 't' is at the head of the queue for node 'n'
 $HasRequest(n, t) \triangleq Len(requests[n]) > 0 \wedge requests[n][1].type = t$

Returns the next request in the queue for node 'n'
 $NextRequest(n) \triangleq requests[n][1]$

Discards the request at the head of the queue for node 'n'
 $DiscardRequest(n) \triangleq requests' = [requests\ EXCEPT\ ![n] = Pop(requests[n])]$

Sends response 'm' on the stream for node 'n'
 $SendResponse(n, m) \triangleq$
 $\wedge responses' = [responses\ EXCEPT\ ![n] = Append(responses[n], m)]$
 $\wedge messageCount' = messageCount + 1$

Indicates whether a response of type 't' is at the head of the queue for node 'n'
 $HasResponse(n, t) \triangleq Len(responses[n]) > 0 \wedge responses[n][1].type = t$

Returns the next response in the queue for node 'n'
 $NextResponse(n) \triangleq responses[n][1]$

Discards the response at the head of the queue for node 'n'
 $DiscardResponse(n) \triangleq responses' = [responses \text{ EXCEPT } ![n] = Pop(responses[n])]$

Indicates whether the stream for node 'n' is *Open*
 $IsStreamOpen(n) \triangleq streams[n].state = Open$

Indicates whether the stream for node 'n' is *Closed*
 $IsStreamClosed(n) \triangleq streams[n].state = Closed$

This section models the mastership election service used by the controller to elect masters. Mastership changes through join and leave steps. Mastership is done through a consensus service, so these steps are atomic. When a node joins or leaves the mastership election, events are queued to notify nodes of the mastership change. Nodes learn of mastership changes independently of the state change in the consensus service.

Node 'n' joins the mastership election

If the current 'master' is *Nil*, set the master to node 'n', increment the 'term', and send a mastership change event to each node. If the current 'master' is non-*Nil*, append node 'n' to the sequence of 'backups'.

$$\begin{aligned}
JoinMastershipElection(n) \triangleq & \\
& \wedge \vee \wedge master = Nil \\
& \quad \wedge term' = term + 1 \\
& \quad \wedge master' = n \\
& \quad \wedge backups' = \langle \rangle \\
& \quad \wedge events' = [i \in Nodes \mapsto Append(events[i], [\\
& \quad \quad \quad term \mapsto term', \\
& \quad \quad \quad master \mapsto master', \\
& \quad \quad \quad backups \mapsto backups'])] \\
& \vee \wedge master \neq Nil \\
& \quad \wedge master \neq n \\
& \quad \wedge n \notin Range(backups) \\
& \quad \wedge backups' = Append(backups, n) \\
& \quad \wedge events' = [i \in Nodes \mapsto Append(events[i], [\\
& \quad \quad \quad term \mapsto term, \\
& \quad \quad \quad master \mapsto master, \\
& \quad \quad \quad backups \mapsto backups'])] \\
& \quad \wedge UNCHANGED \langle term, master \rangle \\
& \quad \wedge mastershipChanges' = mastershipChanges + 1 \\
& \quad \wedge UNCHANGED \langle masterships, streamVars, messageVars, deviceVars \rangle
\end{aligned}$$

Node 'n' leaves the mastership election

If node 'n' is the current 'master' and a backup exists, increment the 'term', promote the first backup to master, and send a mastership change event to each node. If node 'n' is the current 'master' and no backups exist, set the 'master' to *Nil*. If node 'n' is in the sequence of 'backups', simply remove it.

$$\begin{aligned}
\text{LeaveMastershipElection}(n) &\triangleq \\
&\wedge \vee \wedge \text{master} = n \\
&\quad \wedge \vee \wedge \text{Len}(\text{backups}) > 0 \\
&\quad \quad \wedge \text{term}' = \text{term} + 1 \\
&\quad \quad \wedge \text{master}' = \text{backups}[1] \\
&\quad \quad \wedge \text{backups}' = \text{Pop}(\text{backups}) \\
&\quad \quad \wedge \text{events}' = [i \in \text{Nodes} \mapsto \text{Append}(\text{events}[i], [\\
&\quad \quad \quad \text{term} \mapsto \text{term}', \\
&\quad \quad \quad \text{master} \mapsto \text{master}', \\
&\quad \quad \quad \text{backups} \mapsto \text{backups'}])] \\
&\quad \vee \wedge \text{Len}(\text{backups}) = 0 \\
&\quad \quad \wedge \text{master}' = \text{Nil} \\
&\quad \quad \wedge \text{UNCHANGED } \langle \text{term}, \text{backups}, \text{events} \rangle \\
&\vee \wedge n \in \text{Range}(\text{backups}) \\
&\quad \wedge \text{backups}' = \text{Drop}(\text{backups}, \text{CHOOSE } j \in \text{DOMAIN } \text{backups} : \text{backups}[j] = n) \\
&\quad \wedge \text{UNCHANGED } \langle \text{term}, \text{master}, \text{events} \rangle \\
&\wedge \text{mastershipChanges}' = \text{mastershipChanges} + 1 \\
&\wedge \text{UNCHANGED } \langle \text{masterships}, \text{streamVars}, \text{messageVars}, \text{deviceVars} \rangle
\end{aligned}$$

This section models controller-side mastership arbitration. The controller nodes receive mastership change events from the mastership service and send master arbitration requests to the device. Additionally, master nodes can send write requests to the device.

Node 'n' receives a mastership change event from the mastership service

When a mastership change event is received, the node's local mastership state is updated. If the mastership term has changed, the node will set a flag to push the mastership change to the device in the master arbitration step.

$$\begin{aligned}
\text{LearnMastership}(n) &\triangleq \\
&\wedge \text{Len}(\text{events}[n]) > 0 \\
&\wedge \text{LET } e \triangleq \text{events}[n][1] \\
&\quad m \triangleq \text{masterships}[n] \\
&\text{IN} \\
&\quad \vee \wedge e.\text{term} > m.\text{term} \\
&\quad \quad \wedge \text{masterships}' = [\text{masterships} \text{ EXCEPT } ![n] = [\\
&\quad \quad \quad \text{term} \mapsto e.\text{term}, \\
&\quad \quad \quad \text{master} \mapsto e.\text{master}, \\
&\quad \quad \quad \text{backups} \mapsto e.\text{backups}] \\
&\quad \vee \wedge e.\text{term} = m.\text{term} \\
&\quad \quad \wedge \text{masterships}' = [\text{masterships} \text{ EXCEPT } ![n] = [\\
&\quad \quad \quad \text{term} \mapsto e.\text{term}, \\
&\quad \quad \quad \text{master} \mapsto e.\text{master}, \\
&\quad \quad \quad \text{backups} \mapsto e.\text{backups}]
\end{aligned}$$

$$\wedge events' = [events \text{ EXCEPT } ![n] = Pop(events[n])]$$

$$\wedge \text{UNCHANGED } \langle mastershipVars, streamVars, messageVars, deviceVars \rangle$$

Master node 'n' sends a *WriteRequest* to the device

To write to the device, the node must have an open stream, must have received a mastership change event from the mastership service (stored in 'masterships') indicating it is the master, and must have received a *MasterArbitrationUpdate* from the switch indicating it is the master (stored in 'isMaster') for the same term as was indicated by the mastership service. The term is sent with the *WriteRequest* for model checking.

$$\begin{aligned} & SendWriteRequest(n) \triangleq \\ & \quad \wedge IsStreamOpen(n) \\ & \quad \wedge LET \ m \triangleq masterships[n] \\ & \quad IN \\ & \quad \quad \wedge m.term > 0 \\ & \quad \quad \wedge m.master = n \\ & \quad \quad \wedge SendRequest(n, [\\ & \quad \quad \quad type \quad \quad \mapsto WriteRequest, \\ & \quad \quad \quad term \quad \quad \mapsto m.term]) \\ & \quad \wedge \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, streamVars, responses \rangle \end{aligned}$$

Node 'n' receives a write response from the device

$$\begin{aligned} & ReceiveWriteResponse(n) \triangleq \\ & \quad \wedge IsStreamOpen(n) \\ & \quad \wedge HasResponse(n, WriteResponse) \\ & \quad \wedge LET \ m \triangleq NextResponse(n) \\ & \quad IN \\ & \quad \quad \vee m.status = Ok \\ & \quad \quad \vee m.status = PermissionDenied \\ & \quad \wedge DiscardResponse(n) \\ & \quad \wedge \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, streamVars, requests, messageCount \rangle \end{aligned}$$

This section models a *P4 Runtime* device. In this spec, the device's only role is to accept writes from controller nodes. Mastership order is maintained using a simple fencing token stored in 'lastTerm'. The device ensures only new masters can write to it by rejecting writes from older masters according to their provided term.

Opens a new stream between node 'n' and the device

When a stream is opened, the 'streams' state for node 'n' is set to *Open*. Stream creation is modelled as a single step to reduce the state space.

$$\begin{aligned} & ConnectStream(n) \triangleq \\ & \quad \wedge IsStreamClosed(n) \\ & \quad \wedge streams' = [streams \text{ EXCEPT } ![n].state = Open] \\ & \quad \wedge streamChanges' = streamChanges + 1 \\ & \quad \wedge \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, messageVars \rangle \end{aligned}$$

Closes an open stream between node 'n' and the device

When a stream is closed, the 'streams' state for node 'n' is set to *Closed*, and the 'requests' and 'responses' queues for the node are cleared.

$$\begin{aligned}
\text{CloseStream}(n) &\triangleq \\
&\wedge \text{IsStreamOpen}(n) \\
&\wedge \text{streams}' = [\text{streams} \text{ EXCEPT } ![n] = [\text{state} \mapsto \text{Closed}, \text{term} \mapsto 0]] \\
&\wedge \text{requests}' = [\text{requests} \text{ EXCEPT } ![n] = \langle \rangle] \\
&\wedge \text{responses}' = [\text{responses} \text{ EXCEPT } ![n] = \langle \rangle] \\
&\wedge \text{streamChanges}' = \text{streamChanges} + 1 \\
&\wedge \text{UNCHANGED } \langle \text{mastershipVars}, \text{nodeVars}, \text{deviceVars}, \text{messageCount} \rangle
\end{aligned}$$

The device receives a *WriteRequest* from node 'n'

If the *WriteRequest* 'term' is greater than or equal to the highest term received by the device, the write is accepted, the highest term is updated, and the write is recorded in history for model checking. Otherwise, the write was sent by an old master and is rejected with a *PermissionDenied* error.

$$\begin{aligned}
\text{HandleWrite}(n) &\triangleq \\
&\wedge \text{IsStreamOpen}(n) \\
&\wedge \text{HasRequest}(n, \text{WriteRequest}) \\
&\wedge \text{LET } m \triangleq \text{NextRequest}(n) \\
&\text{IN} \\
&\quad \vee \wedge \text{lastTerm} \leq m.\text{term} \\
&\quad \quad \wedge \text{lastTerm}' = m.\text{term} \\
&\quad \quad \wedge \text{history}' = \text{Append}(\text{history}, [\text{node} \mapsto n, \text{term} \mapsto m.\text{term}]) \\
&\quad \quad \wedge \text{SendResponse}(n, [\\
&\quad \quad \quad \text{type} \mapsto \text{WriteResponse}, \\
&\quad \quad \quad \text{status} \mapsto \text{Ok}]) \\
&\quad \vee \wedge \text{lastTerm} > m.\text{term} \\
&\quad \quad \wedge \text{SendResponse}(n, [\\
&\quad \quad \quad \text{type} \mapsto \text{WriteResponse}, \\
&\quad \quad \quad \text{status} \mapsto \text{PermissionDenied}]) \\
&\quad \quad \wedge \text{UNCHANGED } \langle \text{lastTerm}, \text{history} \rangle \\
&\wedge \text{DiscardRequest}(n) \\
&\wedge \text{UNCHANGED } \langle \text{mastershipVars}, \text{nodeVars}, \text{streamVars} \rangle
\end{aligned}$$

The invariant asserts that the device will not allow a write from an older master if it has already accepted a write from a newer master. This is determined by comparing the mastership terms of accepted writes. For this invariant to hold, terms may only increase in the history of writes.

$$\begin{aligned}
\text{TypeInvariant} &\triangleq \\
&\wedge \forall x \in 1 \dots \text{Len}(\text{history}) : \\
&\quad \forall y \in x \dots \text{Len}(\text{history}) : \\
&\quad \quad \text{history}[x].\text{term} \leq \text{history}[y].\text{term} \\
&\wedge \forall x \in 1 \dots \text{Len}(\text{history}) : \\
&\quad \forall y \in x \dots \text{Len}(\text{history}) : \\
&\quad \quad \text{history}[x].\text{term} = \text{history}[y].\text{term} \Rightarrow \text{history}[x].\text{node} = \text{history}[y].\text{node}
\end{aligned}$$

$Init \triangleq$
 $\wedge term = 0$
 $\wedge master = Nil$
 $\wedge backups = \langle \rangle$
 $\wedge events = [n \in Nodes \mapsto \langle \rangle]$
 $\wedge masterships = [n \in Nodes \mapsto [term \mapsto 0, master \mapsto Nil, backups \mapsto \langle \rangle]]$
 $\wedge streams = [n \in Nodes \mapsto [state \mapsto Closed, term \mapsto 0]]$
 $\wedge requests = [n \in Nodes \mapsto \langle \rangle]$
 $\wedge responses = [n \in Nodes \mapsto \langle \rangle]$
 $\wedge lastTerm = 0$
 $\wedge mastershipChanges = 0$
 $\wedge streamChanges = 0$
 $\wedge messageCount = 0$
 $\wedge history = \langle \rangle$

$Next \triangleq$
 $\vee \exists n \in Nodes : ConnectStream(n)$
 $\vee \exists n \in Nodes : CloseStream(n)$
 $\vee \exists n \in Nodes : JoinMastershipElection(n)$
 $\vee \exists n \in Nodes : LeaveMastershipElection(n)$
 $\vee \exists n \in Nodes : LearnMastership(n)$
 $\vee \exists n \in Nodes : SendWriteRequest(n)$
 $\vee \exists n \in Nodes : HandleWrite(n)$
 $\vee \exists n \in Nodes : ReceiveWriteResponse(n)$

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

\backslash * Modification History
 \backslash * Last modified *Wed Feb 20 16:09:21 PST 2019* by *jordanhalterman*
 \backslash * Created *Thu Feb 14 11:33:03 PST 2019* by *jordanhalterman*