

---

MODULE *P4RuntimeElection*

---

EXTENDS *Naturals, FiniteSets, Sequences, TLC*

The set of all *ONOS* nodes

CONSTANTS *Nodes*

Stream states

CONSTANTS *Open, Closed*

Master arbitration message types

CONSTANTS *MasterArbitrationUpdate*

Write message types

CONSTANTS *WriteRequest, WriteResponse*

Response status constants

CONSTANTS *Ok, AlreadyExists, 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 current set of elections for the switch, the greatest of which is the current master

VARIABLE *elections*

Counting variables used to enforce state constraints

VARIABLES *mastershipChanges, streamChanges, messageCount*

---

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  
 $deviceVars \triangleq \langle elections \rangle$

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  $f$   
 $Range(f) \triangleq \{f[x] : x \in \text{DOMAIN } f\}$

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

Messaging between the *Nodes* and the device are modelled on *TCP*. For each node, a request and response sequence provides ordered messaging between the two points. Requests and responses are always received from the head of the queue and are never duplicated or reordered, and request and response queues only last the lifetime of the stream. When a stream is closed, all that stream's requests and responses are lost.

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

Indicates whether any requests are in 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'  
 $DiscardRequest(n) \triangleq requests' = [requests \text{ EXCEPT } ![n] = Pop(requests[n])]$

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

Indicates whether any responses are in 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'  
 $DiscardResponse(n) \triangleq responses' = [responses \text{ EXCEPT } ![n] = Pop(responses[n])]$

This section models mastership arbitration on the controller side. Mastership election occurs in two distinct types of state changes. One state change occurs to change the mastership in the consensus layer, and the other occurs when a node actually learns of the mastership change. Nodes will always learn of mastership changes in the order in which they occur, and nodes will always learn of a mastership change. This, of course, is not representative of practice but is sufficient for modelling the mastership election algorithm.

Adds a node to the mastership election  
 $JoinMastershipElection(n) \triangleq$   
 $\wedge \vee \wedge master = Nil$   
 $\wedge term' = term + 1$   
 $\wedge master' = n$   
 $\wedge backups' = \langle \rangle$   
 $\wedge events' = [i \in Nodes \mapsto Append(events[i], [$   
 $term \mapsto term',$   
 $master \mapsto master',$   
 $backups \mapsto backups'])]$   
 $\vee \wedge master \neq Nil$   
 $\wedge master \neq n$   
 $\wedge n \notin Range(backups)$   
 $\wedge backups' = Append(backups, n)$   
 $\wedge events' = [i \in Nodes \mapsto Append(events[i], [$   
 $term \mapsto term,$   
 $master \mapsto master,$   
 $backups \mapsto backups'])]$   
 $\wedge UNCHANGED \langle term, master \rangle$   
 $\wedge mastershipChanges' = mastershipChanges + 1$   
 $\wedge UNCHANGED \langle masterships, streamVars, messageVars, deviceVars \rangle$

Removes a node from the mastership election  
 $LeaveMastershipElection(n) \triangleq$   
 $\wedge \vee \wedge master = n$   
 $\wedge \vee \wedge Len(backups) > 0$   
 $\wedge term' = term + 1$   
 $\wedge master' = backups[1]$   
 $\wedge backups' = Pop(backups)$   
 $\wedge events' = [i \in Nodes \mapsto Append(events[i], [$

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

Sets the current master to node 'n' if it's not already set

$$\begin{aligned}
\text{SetMastership}(n) & \triangleq \\
& \vee \wedge \text{master} = n \\
& \quad \wedge \text{UNCHANGED } \langle \text{mastershipVars} \rangle \\
& \vee \wedge \text{master} \neq n \\
& \quad \wedge \text{term}' = \text{term} + 1 \\
& \quad \wedge \text{master}' = n \\
& \quad \wedge \vee \wedge n \in \text{Range}(\text{backups}) \\
& \quad \quad \wedge \text{backups}' = \text{Drop}(\text{backups}, \text{CHOOSE } j \in \text{DOMAIN } \text{backups} : \text{backups}[j] = n) \\
& \quad \vee \wedge n \notin \text{Range}(\text{backups}) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{backups} \rangle \\
& \quad \wedge \text{mastershipChanges}' = \text{mastershipChanges} + 1
\end{aligned}$$

Receives a mastership change event from the consensus layer on node 'n'

$$\begin{aligned}
\text{LearnMastership}(n) & \triangleq \\
& \wedge \text{Len}(\text{events}[n]) > 0 \\
& \wedge \text{LET } e \triangleq \text{events}[n][1] \\
& \quad m \triangleq \text{masterships}[n] \\
& \text{IN} \\
& \quad \vee \wedge e.\text{term} > m.\text{term} \\
& \quad \quad \wedge \text{masterships}' = [\text{masterships} \text{ EXCEPT } ![n] = [ \\
& \quad \quad \quad \text{term} \mapsto e.\text{term}, \\
& \quad \quad \quad \text{master} \mapsto e.\text{master}, \\
& \quad \quad \quad \text{backups} \mapsto e.\text{backups}, \\
& \quad \quad \quad \text{sent} \mapsto \text{FALSE}]] \\
& \quad \vee \wedge e.\text{term} = m.\text{term} \\
& \quad \quad \wedge \text{masterships}' = [\text{masterships} \text{ EXCEPT } ![n] = [ \\
& \quad \quad \quad \text{term} \mapsto e.\text{term}, \\
& \quad \quad \quad \text{master} \mapsto e.\text{master}, \\
& \quad \quad \quad \text{backups} \mapsto e.\text{backups}, \\
& \quad \quad \quad \text{sent} \mapsto m.\text{sent}] \\
& \wedge \text{events}' = [\text{events} \text{ EXCEPT } ![n] = \text{Pop}(\text{events}[n])] \\
& \wedge \text{UNCHANGED } \langle \text{mastershipVars}, \text{streamVars}, \text{messageVars}, \text{deviceVars} \rangle
\end{aligned}$$

Notifies the device of node 'n' mastership info if it hasn't already been sent

$SendMasterArbitrationUpdateRequest(n) \triangleq$

$\wedge streams[n] = Open$

$\wedge LET\ m \triangleq masterships[n]$

IN

$\wedge m.term > 0$

$\wedge \neg m.sent$

$\wedge \vee \wedge m.master = n$

$\wedge SendRequest(n, [$   
 $\quad type \mapsto MasterArbitrationUpdate,$   
 $\quad election\_id \mapsto m.term + Cardinality(Nodes)])$

$\vee \wedge m.master \neq n$

$\wedge n \in Range(m.backups)$

$\wedge SendRequest(n, [$   
 $\quad type \mapsto MasterArbitrationUpdate,$   
 $\quad election\_id \mapsto m.term + Cardinality(Nodes) - \text{CHOOSE } i \in \text{DOMAIN } m.backups : m.backups[i]$

$\wedge masterships' = [masterships \text{ EXCEPT } ![n].sent = \text{TRUE}]$

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

Receives a master arbitration update response on node 'n'

$ReceiveMasterArbitrationUpdateResponse(n) \triangleq$

$\wedge streams[n] = Open$

$\wedge HasResponse(n, MasterArbitrationUpdate)$

$\wedge LET\ m \triangleq NextResponse(n)$

IN

$\vee \wedge m.status = Ok$

$\wedge SetMastership(n)$

$\vee \wedge m.status = AlreadyExists$

$\wedge \text{UNCHANGED } \langle mastershipVars \rangle$

$\wedge DiscardResponse(n)$

$\wedge \text{UNCHANGED } \langle nodeVars, deviceVars, streamVars, requests, messageCount \rangle$

Sends a write request to the device from node 'n'

$SendWriteRequest(n) \triangleq$

$\wedge streams[n] = Open$

$\wedge LET\ m \triangleq masterships[n]$

IN

$\wedge m.term > 0$

$\wedge m.master = n$

$\wedge SendRequest(n, [$   
 $\quad type \mapsto WriteRequest,$   
 $\quad election\_id \mapsto m.term + Cardinality(Nodes)])$

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

Receives a write response on node 'n'

$ReceiveWriteResponse(n) \triangleq$

```

 $\wedge$  streams[n] = Open
 $\wedge$  HasResponse(n, WriteResponse)
 $\wedge$  LET m  $\triangleq$  NextResponse(n)
IN
  TODO: This should be used to determine whether writes from old masters are allowed
   $\vee$  m.status = Ok
   $\vee$  m.status = PermissionDenied
 $\wedge$  DiscardResponse(n)
 $\wedge$  UNCHANGED  $\langle$ mastershipVars, nodeVars, deviceVars, streamVars, requests, messageCount $\rangle$ 

```

---

This section models the *P4* switch. The switch side manages stream states between the device and the controller. Streams are opened and closed in a single state transition for the purposes of this model. Switches can handle two types of messages from the controller nodes: *MasterArbitrationUpdate* and *Write*.

Returns the highest election *ID* for the given elections

*ElectionId*(*e*)  $\triangleq$  *Max*(*Range*(*e*))

Returns the master for the given elections

*Master*(*e*)  $\triangleq$

```

IF Cardinality( $\{i \in \text{Range}(e) : i > 0\}$ ) > 0 THEN
  CHOOSE n  $\in$  DOMAIN e : e[n] = ElectionId(e)
ELSE
  Nil

```

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

When a new stream is opened, the '*requests*' and '*responses*' queues for the node are cleared and the '*streams*' state is set to '*Open*'.

```

ConnectStream(n)  $\triangleq$ 
 $\wedge$  streams[n] = Closed
 $\wedge$  streams' = [streams EXCEPT ![n] = Open]
 $\wedge$  requests' = [requests EXCEPT ![n] =  $\langle \rangle$ ]
 $\wedge$  responses' = [responses EXCEPT ![n] =  $\langle \rangle$ ]
 $\wedge$  streamChanges' = streamChanges + 1
 $\wedge$  UNCHANGED  $\langle$ mastershipVars, nodeVars, deviceVars, messageCount $\rangle$ 

```

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

When the stream is closed, the '*requests*' and '*responses*' queues for the node are cleared and a '*MasterArbitrationUpdate*' is sent to all remaining connected nodes to notify them of a mastership change if necessary.

```

CloseStream(n)  $\triangleq$ 
 $\wedge$  streams[n] = Open
 $\wedge$  elections' = [elections EXCEPT ![n] = 0]
 $\wedge$  streams' = [streams EXCEPT ![n] = Closed]
 $\wedge$  requests' = [requests EXCEPT ![n] =  $\langle \rangle$ ]
 $\wedge$  LET oldMaster  $\triangleq$  Master(elections)

```

$$\begin{aligned}
& newMaster \triangleq Master(elections') \\
\text{IN} \\
& \vee \wedge oldMaster \neq newMaster \\
& \wedge responses' = [i \in \text{DOMAIN } streams' \mapsto \\
& \quad \text{IF } i = newMaster \text{ THEN} \\
& \quad \quad Append(responses[i], [ \\
& \quad \quad \quad type \mapsto MasterArbitrationUpdate, \\
& \quad \quad \quad status \mapsto Ok, \\
& \quad \quad \quad election\_id \mapsto ElectionId(elections')]) \\
& \quad \text{ELSE} \\
& \quad \quad Append(responses[i], [ \\
& \quad \quad \quad type \mapsto MasterArbitrationUpdate, \\
& \quad \quad \quad status \mapsto AlreadyExists, \\
& \quad \quad \quad election\_id \mapsto ElectionId(elections')])]) \\
& \wedge messageCount' = messageCount + 1 \\
& \vee \wedge oldMaster = newMaster \\
& \wedge responses' = [responses \text{ EXCEPT } ![n] = \langle \rangle] \\
& \wedge \text{UNCHANGED } \langle messageCount \rangle \\
& \wedge streamChanges' = streamChanges + 1 \\
& \wedge \text{UNCHANGED } \langle mastershipVars, nodeVars \rangle
\end{aligned}$$

Handles a master arbitration update on the device

If the *election\_id* is already present in the 'elections', send an 'AlreadyExists' response to the node. Otherwise,

$$\begin{aligned}
& HandleMasterArbitrationUpdate(n) \triangleq \\
& \wedge streams[n] = Open \\
& \wedge HasRequest(n, MasterArbitrationUpdate) \\
& \wedge \text{LET } m \triangleq NextRequest(n) \\
\text{IN} \\
& \vee \wedge m.election\_id \in Range(elections) \\
& \wedge SendResponse(n, [ \\
& \quad type \mapsto MasterArbitrationUpdate, \\
& \quad election\_id \mapsto m.election\_id, \\
& \quad status \mapsto AlreadyExists]) \\
& \wedge \text{UNCHANGED } \langle deviceVars \rangle \\
& \vee \wedge m.election\_id \notin Range(elections) \\
& \wedge elections' = [elections \text{ EXCEPT } ![n] = m.election\_id] \\
& \wedge \text{LET } oldMaster \triangleq Master(elections) \\
& \quad newMaster \triangleq Master(elections') \\
\text{IN} \\
& \vee \wedge oldMaster \neq newMaster \\
& \wedge responses' = [i \in \text{DOMAIN } streams \mapsto \\
& \quad \text{IF } i = newMaster \text{ THEN} \\
& \quad \quad Append(responses[i], [ \\
& \quad \quad \quad type \mapsto MasterArbitrationUpdate,
\end{aligned}$$

$$\begin{aligned}
& \text{status} \mapsto \text{Ok}, \\
& \text{election\_id} \mapsto \text{ElectionId}(\text{elections}') \\
& \text{ELSE} \\
& \quad \text{Append}(\text{responses}[i], [ \\
& \quad \quad \text{type} \mapsto \text{MasterArbitrationUpdate}, \\
& \quad \quad \text{status} \mapsto \text{AlreadyExists}, \\
& \quad \quad \text{election\_id} \mapsto \text{ElectionId}(\text{elections}') \\
& \quad \quad \wedge \text{messageCount}' = \text{messageCount} + 1 \\
& \quad \vee \wedge \text{oldMaster} = \text{newMaster} \\
& \quad \wedge \text{SendResponse}(n, [ \\
& \quad \quad \text{type} \mapsto \text{MasterArbitrationUpdate}, \\
& \quad \quad \text{status} \mapsto \text{Ok}, \\
& \quad \quad \text{election\_id} \mapsto \text{ElectionId}(\text{elections}') \\
& \wedge \text{DiscardRequest}(n) \\
& \wedge \text{UNCHANGED } \langle \text{mastershipVars}, \text{nodeVars}, \text{streamVars} \rangle
\end{aligned}$$

Handles a write request on the device

$$\begin{aligned}
\text{HandleWrite}(n) & \triangleq \\
& \wedge \text{streams}[n] = \text{Open} \\
& \wedge \text{HasRequest}(n, \text{WriteRequest}) \\
& \wedge \text{LET } m \triangleq \text{NextRequest}(n) \\
& \text{IN} \\
& \quad \vee \wedge \text{Cardinality}(\text{DOMAIN elections}) = 0 \\
& \quad \quad \wedge \text{SendResponse}(n, [ \\
& \quad \quad \quad \text{type} \mapsto \text{WriteResponse}, \\
& \quad \quad \quad \text{status} \mapsto \text{PermissionDenied}] \\
& \quad \vee \wedge \text{ElectionId}(\text{elections}) \neq m.\text{election\_id} \\
& \quad \quad \wedge \text{SendResponse}(n, [ \\
& \quad \quad \quad \text{type} \mapsto \text{WriteResponse}, \\
& \quad \quad \quad \text{status} \mapsto \text{PermissionDenied}] \\
& \quad \vee \wedge m.\text{election\_id} \notin \text{Range}(\text{elections}) \\
& \quad \quad \wedge \text{elections}[n] = m.\text{election\_id} \\
& \quad \quad \wedge \text{SendResponse}(n, [ \\
& \quad \quad \quad \text{type} \mapsto \text{WriteResponse}, \\
& \quad \quad \quad \text{status} \mapsto \text{Ok}] \\
& \wedge \text{DiscardRequest}(n) \\
& \wedge \text{UNCHANGED } \langle \text{mastershipVars}, \text{nodeVars}, \text{deviceVars}, \text{streamVars} \rangle
\end{aligned}$$


---


$$\begin{aligned}
\text{Init} & \triangleq \\
& \wedge \text{term} = 0 \\
& \wedge \text{master} = \text{Nil} \\
& \wedge \text{backups} = \langle \rangle \\
& \wedge \text{events} = [n \in \text{Nodes} \mapsto \langle \rangle] \\
& \wedge \text{masterships} = [n \in \text{Nodes} \mapsto [\text{term} \mapsto 0, \text{master} \mapsto 0, \text{backups} \mapsto \langle \rangle, \text{sent} \mapsto \text{FALSE}]]
\end{aligned}$$



$$\begin{aligned}
&\wedge \text{streams} = [n \in \text{Nodes} \mapsto \text{Closed}] \\
&\wedge \text{requests} = [n \in \text{Nodes} \mapsto \langle \rangle] \\
&\wedge \text{responses} = [n \in \text{Nodes} \mapsto \langle \rangle] \\
&\wedge \text{elections} = [n \in \text{Nodes} \mapsto 0] \\
&\wedge \text{mastershipChanges} = 0 \\
&\wedge \text{streamChanges} = 0 \\
&\wedge \text{messageCount} = 0
\end{aligned}$$

$$\begin{aligned}
\text{Next} &\triangleq \\
&\vee \exists n \in \text{Nodes} : \text{ConnectStream}(n) \\
&\vee \exists n \in \text{Nodes} : \text{CloseStream}(n) \\
&\vee \exists n \in \text{Nodes} : \text{JoinMastershipElection}(n) \\
&\vee \exists n \in \text{Nodes} : \text{LeaveMastershipElection}(n) \\
&\vee \exists n \in \text{Nodes} : \text{LearnMastership}(n) \\
&\vee \exists n \in \text{Nodes} : \text{SendMasterArbitrationUpdateRequest}(n) \\
&\vee \exists n \in \text{Nodes} : \text{HandleMasterArbitrationUpdate}(n) \\
&\vee \exists n \in \text{Nodes} : \text{ReceiveMasterArbitrationUpdateResponse}(n) \\
&\vee \exists n \in \text{Nodes} : \text{SendWriteRequest}(n) \\
&\vee \exists n \in \text{Nodes} : \text{HandleWrite}(n) \\
&\vee \exists n \in \text{Nodes} : \text{ReceiveWriteResponse}(n)
\end{aligned}$$

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


---

\ \* Modification History  
\ \* Last modified Sun Feb 17 00:35:14 PST 2019 by jordanhalterman  
\ \* Created Thu Feb 14 11:33:03 PST 2019 by jordanhalterman