─────────────── 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*

  A sequence of successful writes to the switch used for model checking
VARIABLE *writes*

─────────────────────────────────────────────────────────

  Mastership/consensus related variables
$mastershipVars \triangleq \langle term, master, backups, mastershipChanges \rangle$

  Node related variables
$nodeVars \triangleq \langle events, masterships \rangle$

1

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

$messageVars \triangleq \langle requests,\ responses,\ messageCount \rangle$

$deviceVars \triangleq \langle elections,\ writes \rangle$

$vars \triangleq \langle mastershipVars,\ nodeVars,\ streamVars,\ messageVars,\ deviceVars \rangle$

---

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'
$SendRequest(n,\ m) \triangleq$
$\quad \wedge requests' = [requests \text{ EXCEPT } ![n] = Append(requests[n],\ m)]$
$\quad \wedge messageCount' = messageCount + 1$

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'
$SendResponse(n,\ m) \triangleq$

$$\land\ responses' = [responses \text{ EXCEPT } ![n] = Append(responses[n],\ m)]$$
$$\land\ messageCount' = messageCount + 1$$

Indicates whether any responses are in the queue for node 'n'
$$HasResponse(n,\ t)\ \triangleq\ Len(responses[n]) > 0 \land responses[n][1].type = t$$

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 disctinct 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
$$JoinMastershipElection(n)\ \triangleq$$
$$\land\ \lor\ \land\ master = Nil$$
$$\land\ term' = term + 1$$
$$\land\ master' = n$$
$$\land\ backups' = \langle\rangle$$
$$\land\ events' = [i \in Nodes \mapsto Append(events[i],\ [$$
$$term \mapsto term',$$
$$master \mapsto master',$$
$$backups \mapsto backups'])]$$
$$\lor\ \land\ master \neq Nil$$
$$\land\ master \neq n$$
$$\land\ n \notin Range(backups)$$
$$\land\ backups' = Append(backups,\ n)$$
$$\land\ events' = [i \in Nodes \mapsto Append(events[i],\ [$$
$$term \mapsto term,$$
$$master \mapsto master,$$
$$backups \mapsto backups'])]$$
$$\land\ \text{UNCHANGED } \langle term,\ master \rangle$$
$$\land\ mastershipChanges' = mastershipChanges + 1$$
$$\land\ \text{UNCHANGED } \langle masterships,\ streamVars,\ messageVars,\ deviceVars \rangle$$

Removes a node from the mastership election
$$LeaveMastershipElection(n)\ \triangleq$$
$$\land\ \lor\ \land\ master = n$$
$$\land\ \lor\ \land\ Len(backups) > 0$$
$$\land\ term' = term + 1$$
$$\land\ master' = backups[1]$$

$$\land\ backups' = Pop(backups)$$
$$\land\ events' = [i\quad \in Nodes \mapsto Append(events[i], [$$
$$term \mapsto term',$$
$$master \mapsto master',$$
$$backups \mapsto backups'])]$$
$$\lor\ \land\ Len(backups) = 0$$
$$\land\ master' = Nil$$
$$\land\ \textsc{unchanged}\ \langle term,\ backups,\ events\rangle$$
$$\lor\ \land\ n \in Range(backups)$$
$$\land\ backups' = Drop(backups, \textsc{choose}\ j \in \textsc{domain}\ backups : backups[j] = n)$$
$$\land\ \textsc{unchanged}\ \langle term,\ master,\ events\rangle$$
$$\land\ mastershipChanges' = mastershipChanges + 1$$
$$\land\ \textsc{unchanged}\ \langle masterships,\ streamVars,\ messageVars,\ deviceVars\rangle$$

Sets the current master to node 'n' if it's not already set
$$SetMastership(n)\ \triangleq$$
$$\lor\ \land\ master = n$$
$$\land\ \textsc{unchanged}\ \langle mastershipVars\rangle$$
$$\lor\ \land\ master \neq n$$
$$\land\ term' = term + 1$$
$$\land\ master' = n$$
$$\land\ \lor\ \land\ n\ \in Range(backups)$$
$$\land\ backups' = Drop(backups, \textsc{choose}\ j \in \textsc{domain}\ backups : backups[j] = n)$$
$$\lor\ \land\ n \notin Range(backups)$$
$$\land\ \textsc{unchanged}\ \langle backups\rangle$$
$$\land\ mastershipChanges' = mastershipChanges + 1$$

Receives a mastership change event from the consensus layer on node 'n'
$$LearnMastership(n)\ \triangleq$$
$$\land\ Len(events[n]) > 0$$
$$\land\ \textsc{let}\ e\ \triangleq\ events[n][1]$$
$$m\ \triangleq\ masterships[n]$$
$$\textsc{in}$$
$$\lor\ \land\ e.term > m.term$$
$$\land\ masterships' = [masterships\ \textsc{except}\ ![n] = [$$
$$term\quad \mapsto e.term,$$
$$master \mapsto e.master,$$
$$backups \mapsto e.backups,$$
$$sent\quad \mapsto \textsc{false}]]$$
$$\lor\ \land\ e.term = m.term$$
$$\land\ masterships' = [masterships\ \textsc{except}\ ![n] = [$$
$$term\quad \mapsto e.term,$$
$$master \mapsto e.master,$$
$$backups \mapsto e.backups,$$
$$sent\quad \mapsto m.sent]]$$

4

$\wedge\ events' = [events\ \text{EXCEPT}\ ![n] = Pop(events[n])]$
$\wedge\ \text{UNCHANGED}\ \langle mastershipVars,\ streamVars,\ messageVars,\ deviceVars\rangle$

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

$SendMasterArbitrationUpdateRequest(n)\ \triangleq$
    $\wedge\ streams[n] = Open$
    $\wedge\ \text{LET}\ m\ \triangleq\ masterships[n]$
       $\text{IN}$
          $\wedge\ m.term > 0$
          $\wedge\ \neg m.sent$
          $\wedge\ \vee\ \wedge\ m.master = n$
                $\wedge\ SendRequest(n, [$
                    $type\ \ \ \ \ \ \ \ \ \ \ \ \mapsto MasterArbitrationUpdate,$
                    $election\_id \mapsto m.term + Cardinality(Nodes),$
                    $term\ \ \ \ \ \ \ \ \ \ \ \mapsto m.term])$
             $\vee\ \wedge\ m.master \neq n$
                $\wedge\ n \in Range(m.backups)$
                $\wedge\ SendRequest(n, [$
                    $type\ \ \ \ \ \ \ \ \ \ \ \ \mapsto MasterArbitrationUpdate,$
                    $election\_id \mapsto m.term + Cardinality(Nodes) - \text{CHOOSE}\ i \in \text{DOMAIN}\ m.backups : m.back\cdots$
                    $term\ \ \ \ \ \ \ \ \ \ \ \mapsto m.term])$
       $\wedge\ masterships' = [masterships\ \text{EXCEPT}\ ![n].sent = \text{TRUE}]$
       $\wedge\ \text{UNCHANGED}\ \langle mastershipVars,\ events,\ deviceVars,\ streamVars,\ responses\rangle$

Receives a master arbitration update response on node 'n'

$ReceiveMasterArbitrationUpdateResponse(n)\ \triangleq$
    $\wedge\ streams[n] = Open$
    $\wedge\ HasResponse(n, MasterArbitrationUpdate)$
    $\wedge\ \text{LET}\ m\ \triangleq\ NextResponse(n)$
       $\text{IN}$
          $\vee\ \wedge\ m.status = Ok$
             $\wedge\ SetMastership(n)$
          $\vee\ \wedge\ m.status = AlreadyExists$
             $\wedge\ \text{UNCHANGED}\ \langle mastershipVars\rangle$
    $\wedge\ DiscardResponse(n)$
    $\wedge\ \text{UNCHANGED}\ \langle nodeVars,\ deviceVars,\ streamVars,\ requests,\ messageCount\rangle$

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

$SendWriteRequest(n)\ \triangleq$
    $\wedge\ streams[n] = Open$
    $\wedge\ \text{LET}\ m\ \triangleq\ masterships[n]$
       $\text{IN}$
          $\wedge\ m.term > 0$
          $\wedge\ m.master = n$
          $\wedge\ SendRequest(n, [$
             $type\ \ \ \ \ \ \ \ \ \ \ \mapsto WriteRequest,$

5

$$election\_id \mapsto m.term + Cardinality(Nodes),$$
$$term \qquad \mapsto m.term])$$
$$\land \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, streamVars, responses \rangle$$

Receives a write response on node 'n'
$ReceiveWriteResponse(n) \triangleq$
    $\land streams[n] = Open$
    $\land HasResponse(n, WriteResponse)$
    $\land \text{LET } m \triangleq NextResponse(n)$
       IN
            *TODO*: This should be used to determine whether writes from old masters are allowed
            $\lor m.status = Ok$
            $\lor m.status = PermissionDenied$
    $\land DiscardResponse(n)$
    $\land \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, streamVars, requests, messageCount \rangle$

---

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 $Cardinality(\{i \in Range(e) : i > 0\}) > 0$ THEN
        CHOOSE $n \in \text{DOMAIN } e : e[n] = ElectionId(e)$
      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'.
$ConnectStream(n) \triangleq$
    $\land streams[n] = Closed$
    $\land streams' = [streams \text{ EXCEPT } ![n] = Open]$
    $\land requests' = [requests \text{ EXCEPT } ![n] = \langle\rangle]$
    $\land responses' = [responses \text{ EXCEPT } ![n] = \langle\rangle]$
    $\land streamChanges' = streamChanges + 1$
    $\land \text{UNCHANGED } \langle mastershipVars, nodeVars, deviceVars, messageCount \rangle$

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.

$CloseStream(n) \triangleq$
 $\wedge\ streams[n] = Open$
 $\wedge\ elections'\ = [elections\ \text{EXCEPT}\ ![n] = 0]$
 $\wedge\ streams' = [streams\ \text{EXCEPT}\ ![n] = Closed]$
 $\wedge\ requests' = [requests\ \text{EXCEPT}\ ![n] = \langle\rangle]$
 $\wedge\ \text{LET}\ oldMaster\ \triangleq\ Master(elections)$
    $newMaster\ \triangleq\ Master(elections')$
  $\text{IN}$
   $\vee\ \wedge\ oldMaster \neq newMaster$
    $\wedge\ responses' = [i \in \text{DOMAIN}\ streams' \mapsto$
         $\text{IF}\ i = newMaster\ \text{THEN}$
           $Append(responses[i], [$
             $type\quad\quad\ \mapsto MasterArbitrationUpdate,$
             $status\quad\ \mapsto Ok,$
             $election\_id \mapsto ElectionId(elections')])$
         $\text{ELSE}$
           $Append(responses[i], [$
             $type\quad\quad\ \mapsto MasterArbitrationUpdate,$
             $status\quad\ \mapsto AlreadyExists,$
             $election\_id \mapsto ElectionId(elections')])]$
    $\wedge\ messageCount' = messageCount + 1$
   $\vee\ \wedge\ oldMaster = newMaster$
    $\wedge\ responses' = [responses\ \text{EXCEPT}\ ![n] = \langle\rangle]$
    $\wedge\ \text{UNCHANGED}\ \langle messageCount\rangle$
 $\wedge\ streamChanges' = streamChanges + 1$
 $\wedge\ \text{UNCHANGED}\ \langle mastershipVars, nodeVars, writes\rangle$

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,

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

$$\text{IN}$$

$$\lor \land \text{oldMaster} \neq \text{newMaster}$$
$$\land \text{responses}' = [i \in \text{DOMAIN } \text{streams} \mapsto$$
$$\text{IF } i = \text{newMaster } \text{THEN}$$
$$Append(\text{responses}[i], [$$
$$type \qquad \mapsto MasterArbitrationUpdate,$$
$$status \qquad \mapsto Ok,$$
$$election\_id \mapsto ElectionId(\text{elections}')])$$
$$\text{ELSE}$$
$$Append(\text{responses}[i], [$$
$$type \qquad \mapsto MasterArbitrationUpdate,$$
$$status \qquad \mapsto AlreadyExists,$$
$$election\_id \mapsto ElectionId(\text{elections}')])]$$
$$\land \text{messageCount}' = \text{messageCount} + 1$$
$$\lor \land \text{oldMaster} = \text{newMaster}$$
$$\land SendResponse(n, [$$
$$type \qquad \mapsto MasterArbitrationUpdate,$$
$$status \qquad \mapsto Ok,$$
$$election\_id \mapsto ElectionId(\text{elections}')])$$

$$\land DiscardRequest(n)$$
$$\land \text{UNCHANGED } \langle mastershipVars, nodeVars, streamVars, writes \rangle$$

Handles a write request on the device
$$HandleWrite(n) \triangleq$$
$$\land \text{streams}[n] = Open$$
$$\land HasRequest(n, WriteRequest)$$
$$\land \text{LET } m \triangleq NextRequest(n)$$
$$\text{IN}$$
$$\lor \land \text{elections}[n] = m.election\_id$$
$$\land Master(\text{elections}) = n$$
$$\land \text{writes}' = Append(\text{writes}, [node \mapsto n, term \mapsto m.term])$$
$$\land SendResponse(n, [$$
$$type \quad \mapsto WriteResponse,$$
$$status \mapsto Ok])$$
$$\lor \land \lor \text{elections}[n] \neq m.election\_id$$
$$\lor Master(\text{elections}) \neq n$$
$$\land SendResponse(n, [$$
$$type \quad \mapsto WriteResponse,$$
$$status \mapsto PermissionDenied])$$
$$\land \text{UNCHANGED } \langle writes \rangle$$
$$\land DiscardRequest(n)$$
$$\land \text{UNCHANGED } \langle mastershipVars, nodeVars, elections, streamVars \rangle$$

---

The invariant asserts that no master can write to the switch after the switch

has been notified of a newer master

$TypeInvariant \triangleq \forall\, i \in \text{DOMAIN}\ writes : i = 1 \vee writes[i-1].term \leq writes[i].term$

---

$Init \triangleq$
    $\wedge\ term = 0$
    $\wedge\ master = Nil$
    $\wedge\ backups = \langle\rangle$
    $\wedge\ events = [n \in Nodes \mapsto \langle\rangle]$
    $\wedge\ masterships = [n \in Nodes \mapsto [term \mapsto 0,\ master \mapsto 0,\ backups \mapsto \langle\rangle,\ sent \mapsto \text{FALSE}]]$
    $\wedge\ streams = [n \in Nodes \mapsto Closed]$
    $\wedge\ requests = [n \in Nodes \mapsto \langle\rangle]$
    $\wedge\ responses = [n \in Nodes \mapsto \langle\rangle]$
    $\wedge\ elections\ = [n \in Nodes \mapsto 0]$
    $\wedge\ mastershipChanges = 0$
    $\wedge\ streamChanges = 0$
    $\wedge\ messageCount = 0$
    $\wedge\ writes = \langle\rangle$

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

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

---

\ * Modification History
\ * Last modified Sun *Feb* 17 03:10:44 *PST* 2019 by *jordanhalterman*
\ * Created *Thu Feb* 14 11:33:03 *PST* 2019 by *jordanhalterman*