

The Data Center Network L1 Switch Protocol

Thomas L. Rodeheffer
Microsoft Research, Silicon Valley

April 27, 2011

1 Overview

Chuck Thacker *et al.* have proposed a data center network design using FPGAs to implement the switching infrastructure [2]. Inside each switch, the design uses time slot circuit switching of data cells via a distributed crossbar in the data plane and best-effort forwarding of control messages via a high-speed ring in the control plane. The goal is to produce a system with minimum total complexity. Existing networks are extremely complex, have many different failure modes, and are difficult to scale out to the size of a large data center.

The design fits in the context of a data center built in shipping containers. Each container holds a large number of servers and L0 switches, and holds two, top-level L1 switches, called A and B. Cables are run between all the containers to form a complete interconnection graph of the A L1 switches and likewise for the B L1 switches.

This paper presents a TLA+ [1] specification of the L1 switch infrastructure.

1.1 Simplifications

This specification simplifies the full design in a number of respects.

1.1.1 No L0 switches

The full design consists of 64 containers. Each container has 1280 servers, 64 L0 switches, and two L1 switches, called A and B.

This specification concerns itself with the A L1 switches only, ignoring all the servers, L0 switches, and B L1 switches. Where there would be a connection to an L0 switch, the specification connects to a model of an end node that sources and sinks traffic, thus simulating the server nics that would be reached through the L0 switches in the full design.

1.1.2 No link frames

In the full design, data cells and control cells are sent on the links structured inside link frames. Each input unit has a deframer that extracts the cells from the input link frame and each output unit has an enframer that packs the cells into the output link frame.

Furthermore, each input unit uses a double buffer of frame data cells to reorder the data cells between the slot time in their arriving frame and the slot time in which they are sucked through the crossbar into their outgoing frame.

This specification ignores the link frames. In this specification, data cells and control cells are sent on links in any sequence. Each input unit uses a queue per slot to reorder the data cells from their arriving times to their departing times.

1.1.3 Assumed timing constraint

In the full design, worst-case timing assumptions are made in order to guarantee the following timing constraint.

During connection setup, when switch A's down port processes a setup request message and successfully allocates a local slot, two control paths emanate:

- On the first control path, the down port generates a setup answer message which returns around the ring to the upport. The up port converts the setup answer message to a setup xbar message which circles the ring to setup the connection in all slices of the distributed crossbar. Eventually, the distributed crossbar is completely setup and cells that arrive at the down port during the assigned local slot correspond to cells arriving on the new connection at the up port.
- On the second control path, the down port generates a setup request control cell that it sends downstream across the link to the downstream switch or node. If the connection is successfully setup in the downstream component, eventually it replies with a setup ack control cell. This cell arrives at switch A's down port, informing it of the value to use for nextSlot on the downstream link. At this point the down port begins forwarding cells arriving through the crossbar on the new connection.

The timing constraint is that the first control path must complete before the second control path. Otherwise, the down port could forward arbitrary garbage cells on the downstream link, and this would be bad.

The specification guarantees the timing constraint by brute force. A connection identifier is associated with all control cells and crossbar setups, and the setup ack is prevented from being received until the crossbar has been setup.

1.2 Terminology

Terminology in the data center network design document [2] tends to be incomplete and wordy. That is, not everything we need to talk about has a crisp name. So, in this specification, we adopt the following terminology.

The specification has to deal with objects, which often have internal state and actions upon them. Objects can be complicated and contain other objects. We have the following examples of objects.

- A *data cell* object, with type name DataCell (section 9). Data cell objects are sent over links, buffered at input ports, and forwarded through a switch crossbar to output ports.

- A *control cell* object, with type name CtrlCell (section 10). Control cell objects are sent over links from output ports to input ports.
- A *ring message* object, with type name RingMsg (section 12). Ring message objects are sent around the elastic ring inside a switch.
- A *switch* object, with type name Switch (section 18).
- A *port* object (part of a switch), with type name Port (section 14). A port object is called an “i/o unit” in the design document.
- A *gare* object (part of a switch), with type name Gare (section 17). A gare object manages a station on the control ring within a switch on behalf of a number of ports. (“Gare” is the French word for “train station”.) A gare is called a “ring master” in the design document. Unfortunately, “ring master” and “ring message” have the same initials, so in this specification we use a different term for “ring master”.
- A *node* object, with type name Node (section 11). A node object models the network interface card in a server.

In addition to objects, the specification has to deal with references to objects (or to places). In full generality, a reference is a *name*, which can be complicated and contain internal structure. We have the following examples of names.

- A *slot name*, with type name SlotName (section 4). A slot name has an internal structure that identifies the subframe and slottina within the subframe.
- A *port name* (locally within a switch), with type name PortName (section 13). A port name has an internal structure with multiple variants, one variant for a port name of a port that connects externally to a container and another for a port name of a port that connects externally to a halfrow.
- A *gare name* (locally within a switch), with type name GareName (section 15). A gare name has an internal structure which identifies the set of ports it manages.
- A *switch name*, with type name SwitchName (section 16). A switch name has an internal structure which identifies the container of the L1 A switch.

1.3 Modules

The specification is written as a collection of modules, in which each module provides definitions for a particular data type or component of the design. In this way the complexity of the overall design is factored into logical pieces that can be understood one at a time.

First we specify some useful definitions (section 2) and basic system parameters (section 3).

Then we build up basic system modules. This requires specifying the definition of a *chunk sequence number* (section 5), the definition of a *CRC* (section 6), the definition

of a *slot name* (section 4), the definition of *cell type* (section 7), the definition of *control queue type* (section 8), the format of a *data cell* (section 9), the format of a *control cell* (section 10), and the internal state and actions of an *end node* (section 11).

Then we add many modules requires to specify the components of an L1 switching network. This requires specifying the format of a *ring message* (section 12), the definition of a *port name* (section 13), the definition of a *gare name* (section 15), the definition of a *switch name* (section 16), the internal state and actions of a *port* (section 14), the internal state and actions of a *gare* (section 17), the internal state and actions of a *switch* (section 18), and finally a system to connect a number of switches and end nodes (section 19).

2 Useful definitions

MODULE *DCNL1 Useful*

EXTENDS *Naturals, Sequences, FiniteSets, TLC*

Useful TLA constants and operators.

A model constant representing a value that is never (supposed to be) read. A field can be assigned to *NULL* but there should be no checks to see if a field is *NULL*.

CONSTANT *NULL*

Nullable version of type *T*. This is used to define the type of a field that may be assigned *NULL* in some situations.

$Nullable(T) \triangleq T \cup \{NULL\}$

Assert that *x* is in type *T*. Returns TRUE if no error.

$TypeAssert(x, T) \triangleq$
 IF $x \in T$ THEN TRUE ELSE
 $Assert(FALSE, \langle \text{“type error”, } x \rangle)$

The head of sequence *s* if *s* is non-empty, else *n*.

$HeadElse(s, n) \triangleq$ IF $s \neq \langle \rangle$ THEN $Head(s)$ ELSE n

The tail of sequence *s* if *s* is non-empty, else $\langle \rangle$.

$TailElse(s) \triangleq$ IF $s \neq \langle \rangle$ THEN $Tail(s)$ ELSE $\langle \rangle$

Choose an arbitrary sequence of “names” (the elements of a set).

In TLA+ this would more naturally be written as a submodule taking the set $Name$ as a constant and defining all these symbols as operators. The module would then be instantiated using the desired set as a parameter. Unfortunately, TLC fails to recognize that the resulting operators are level 0 constants, and it repeatedly reevaluates them during model checking.

So instead we write this as an operator that returns a record of maps. In this style, TLC does recognize that everything is a level 0 constant, and it evaluates them only once before computing the initial states.

$ArbSeq(Name) \triangleq$

LET

Number of names.

$N \triangleq Cardinality(Name)$

The index set.

$Index \triangleq 1 .. N$

Get name for a given index.

$NameForIndex \triangleq$

LET

RECURSIVE $Recur(-, -)$

$Recur(nn, q) \triangleq$

IF $nn = \{\}$ THEN q ELSE

LET $n \triangleq$ CHOOSE $n \in nn : TRUEIN$

$Recur(nn \setminus \{n\}, Append(q, n))$

IN

$Recur(Name, \langle \rangle)$

Get index for a given name.

$IndexForName \triangleq [n \in Name \mapsto \text{CHOOSE } i \in Index : NameForIndex[i] = n]$

Operations on indexes.

$FirstIndex \triangleq 1$

$IncrIndex \triangleq [i \in Index \mapsto ((i - 1 + N + 1) \% N) + 1]$ cyclic

$DecrIndex \triangleq [i \in Index \mapsto ((i - 1 + N - 1) \% N) + 1]$ cyclic

$LeqIndex \triangleq [i1, i2 \in Index \mapsto i1 \leq i2]$

Operations on names.

$FirstName \triangleq NameForIndex[FirstIndex]$

$IncrName \triangleq [n \in Name \mapsto NameForIndex[IncrIndex[IndexForName[n]]]]$

$DecrName \triangleq [n \in Name \mapsto NameForIndex[DecrIndex[IndexForName[n]]]]$

$LeqName \triangleq [n1, n2 \in Name \mapsto LeqIndex[IndexForName[n1], IndexForName[n2]]]$

IN

[
 $Index \mapsto Index,$

```

NameForIndex ↦ NameForIndex,
IndexForName ↦ IndexForName,
FirstIndex   ↦ FirstIndex,
IncrIndex    ↦ IncrIndex,
DecrIndex    ↦ DecrIndex,
LeqIndex     ↦ LeqIndex,
FirstName    ↦ FirstName,
IncrName     ↦ IncrName,
DecrName     ↦ DecrName,
LeqName     ↦ LeqName
]

```

3 Parameters

MODULE *DCNL1Param*

EXTENDS *DCNL1Useful*

Parameterization constants. Values in the full design are listed in parenthesis.

```

CONSTANT NumContainer    number of containers (64)
CONSTANT NumHalfrow     number of halfrows (64)

```

Note that hardware implementation of the following definitions is zero-based. TLA+ likes it better for natural number ranges to start with one, so we define it the TLA+ way.

```

Container ≜ 1 .. NumContainer
Halfrow   ≜ 1 .. NumHalfrow

```

```

Reason ≜ STRING      teardown reason, opaque
Throttle ≜ BOOLEAN   throttle indication

```

```

CONSTANT Payload    normal data cell payload (32 bytes)

```

For the purpose of the specification, each connection has an identification that is carried by every data cell, control cell, ring message, and status record pertaining to the connection. This identification is not present in the hardware, but is used to specify certain protocol timing assumptions that are guaranteed by the hardware.

```

CONSTANT NumConnection

```

```

Connection ≜ 1 .. NumConnection

```

4 Slot name

MODULE *DCNL1SlotName*

EXTENDS *DCNL1Useful*

Time on the links, and in the data path inside a switch, is divided into recurring frames. Each frame consists of a sequence of subframes. Each *subframe* consists of a sequence of slottinas.

A slot is a particular *slottina* inside a particular *subframe*.

The existence of “*subframe*” in the time hierarchy is partly an implementation artifact in that the hardware design did not want to spend the effort to carry a full bit vector of slot-available status around the ring in the “*SetupReq*” message. Note that the “*SetupReq*” message is already the biggest ring message, being the only two word message in the hardware.

Parameterization constants. Values in the full design are listed in parenthesis.

CONSTANT *NumSubframe* number of subframes (4)

CONSTANT *NumSlottina* number of slottinas (32)

Note that hardware implementation of the following definitions is zero-based. TLA+ likes it better for natural number ranges to start with one, so we define it the TLA+ way.

Subframe \triangleq 1 .. *NumSubframe*

Slottina \triangleq 1 .. *NumSlottina*

Define slot name.

SlotName \triangleq [
 subframe : *Subframe*,
 slottina : *Slottina*
]

Define slot index based on an arbitrary ordering of slot names.

FirstSlotName \triangleq *ArbSeq(SlotName).FirstName*

IncrSlotName \triangleq *ArbSeq(SlotName).IncrName*

DecrSlotName \triangleq *ArbSeq(SlotName).DecrName*

LeqSlotName \triangleq *ArbSeq(SlotName).LeqName*

5 Chunk sequence number

MODULE *DCNL1Chunkseqn*

EXTENDS *DCNL1Useful*

Chunk sequence number.

Parameterization constants. Values in the full design are listed in parenthesis.

CONSTANT *NumChunkseqn* number of chunk sequence numbers (2 ** 16)

Note that hardware implementation of the following definition is zero-based. TLA+ likes it better for natural number ranges to start with one, so we define it the TLA+ way.

$Chunkseqn \triangleq 1 .. NumChunkseqn$

$FirstChunkseqn \triangleq 1$

$IncrChunkseqn(x) \triangleq (x \% NumChunkseqn) + 1$

$DecrChunkseqn(x) \triangleq \text{CHOOSE } x0 \in Chunkseqn : IncrChunkseqn(x0) = x$

6 CRC

MODULE *DCNL1CRC*

EXTENDS *DCNL1Param*

We model a *CRC* as a record containing the sequence of payload values it was computed from.

$CRC \triangleq [crc : Seq(Payload)]$

The initial value of a *CRC*.

$InitCRC \triangleq [crc \mapsto \langle \rangle]$

Update a *crc* by extending it with a payload.

$ExtendCRC(crc, p) \triangleq [crc \text{ EXCEPT } !.crc = Append(@, p)]$

7 Cell type

MODULE *DCNL1CellType*

$CellType_Data \triangleq \text{"CellType_Data"}$

$CellType_Ctrl \triangleq \text{"CellType_Ctrl"}$

```

CellType ≜ {
  CellType_Data,
  CellType_Ctrl
}

```

8 Control queue type

MODULE *DCNL1CtrlQueueType*

```

CtrlQueueType_Req ≜ "CtrlQueueType_Req"
CtrlQueueType_Ups ≜ "CtrlQueueType_Ups"

CtrlQueueType ≜ {
  CtrlQueueType_Req,
  CtrlQueueType_Ups
}

```

9 Data cell

MODULE *DCNL1DataCell*

EXTENDS

```

DCNL1Param,
DCNL1SlotName,
DCNL1Chunkseqn,
DCNL1CRC

```

Data cells on the link.

In the hardware, each data cell has a 2 byte header and a 32 byte *payload*. The header has a parity bit which we ignore.

```

DataCell_Typecode_Null      ≜ "DataCell_Null"
DataCell_Typecode_Normal    ≜ "DataCell_Normal"
DataCell_Typecode_FirstMark ≜ "DataCell_FirstMark"
DataCell_Typecode_ChunkMark ≜ "DataCell_ChunkMark"
DataCell_Typecode_EndMark   ≜ "DataCell_EndMark"

```

$$\text{DataCell_Null} \triangleq [$$

$$\text{typecode} : \{\text{DataCell_Typecode_Null}\}$$

$$]$$

$$\text{DataCell_Normal} \triangleq [$$

$$\text{typecode} : \{\text{DataCell_Typecode_Normal}\},$$

$$\text{connect} : \text{Connection},$$

$$\text{nextSlot} : \text{SlotName},$$

$$\text{payload} : \text{Payload}$$

$$]$$

$$\text{DataCell_FirstMark} \triangleq [$$

$$\text{typecode} : \{\text{DataCell_Typecode_FirstMark}\},$$

$$\text{connect} : \text{Connection},$$

$$\text{nextSlot} : \text{SlotName},$$

$$\text{destC} : \text{Container},$$

$$\text{destH} : \text{Halfrow}$$

$\text{this L1 model terminates at the halfrow, hence there is no } \text{DestN}$

$$]$$

$$\text{DataCell_ChunkMark} \triangleq [$$

$$\text{typecode} : \{\text{DataCell_Typecode_ChunkMark}\},$$

$$\text{connect} : \text{Connection},$$

$$\text{nextSlot} : \text{SlotName},$$

$$\text{chunkseqn} : \text{Chunkseqn},$$

$$\text{crc} : \text{CRC}$$

$$]$$

$$\text{DataCell_EndMark} \triangleq [$$

$$\text{typecode} : \{\text{DataCell_Typecode_EndMark}\},$$

$$\text{connect} : \text{Connection},$$

$$\text{nextSlot} : \text{SlotName}$$

$$]$$

$$\text{DataCell} \triangleq \{ \}$$

$$\cup \text{DataCell_Null}$$

$$\cup \text{DataCell_Normal}$$

$$\cup \text{DataCell_FirstMark}$$

$$\cup \text{DataCell_ChunkMark}$$

$$\cup \text{DataCell_EndMark}$$

$$\text{NullDataCell} \triangleq [\text{typecode} \mapsto \text{DataCell_Typecode_Null}]$$

$$\text{DataCell_Is_Null}(dc) \triangleq dc.\text{typecode} = \text{DataCell_Typecode_Null}$$

$$\text{DataCell_Is_Normal}(dc) \triangleq dc.\text{typecode} = \text{DataCell_Typecode_Normal}$$

$$\text{DataCell_Is_FirstMark}(dc) \triangleq dc.\text{typecode} = \text{DataCell_Typecode_FirstMark}$$

$$\text{DataCell_Is_ChunkMark}(dc) \triangleq dc.\text{typecode} = \text{DataCell_Typecode_ChunkMark}$$

DataCell_Is_EndMark(*dc*) \triangleq *dc.typecode* = *DataCell_Typecode_EndMark*

Nice view for debugging.

```

DataCell_Debug_Nice(dc)  $\triangleq$ 
  LET
    sn  $\triangleq$  dc.nextSlot
    nulldc  $\triangleq$  <"Null">
    normaldc  $\triangleq$  <"DC", sn, dc.payload>
    firstmdc  $\triangleq$  <"FM", sn>
    chunkmdc  $\triangleq$  <"CM", sn, dc.chunkSeqn>
    endmdc  $\triangleq$  <"EM", sn>
  IN
  IF DataCell_Is_Null(dc) THEN nulldc ELSE
  IF DataCell_Is_Normal(dc) THEN normaldc ELSE
  IF DataCell_Is_FirstMark(dc) THEN firstmdc ELSE
  IF DataCell_Is_ChunkMark(dc) THEN chunkmdc ELSE
  IF DataCell_Is_EndMark(dc) THEN endmdc ELSE
  Assert(FALSE, "type error")

```

10 Control cell

MODULE *DCNL1_CtrlCell*

```

EXTENDS
  DCNL1Param,
  DCNL1SlotName,
  DCNL1Chunkseqn

```

Control cells on the link.

In the hardware, each control cell fits in a 32 bit word and contains a parity bit. Control cells received with bad parity are ignored. We do not model the parity bit.

```

CtrlCell_Typecode_SetupReq  $\triangleq$  "CtrlCell_SetupReq"
CtrlCell_Typecode_SetupAck  $\triangleq$  "CtrlCell_SetupAck"
CtrlCell_Typecode_Teardown  $\triangleq$  "CtrlCell_Teardown"
CtrlCell_Typecode_Discard  $\triangleq$  "CtrlCell_Discard"
CtrlCell_Typecode_DataAck  $\triangleq$  "CtrlCell_DataAck"
CtrlCell_Typecode_Restart  $\triangleq$  "CtrlCell_Restart"
CtrlCell_Typecode_Stop  $\triangleq$  "CtrlCell_Stop"
CtrlCell_Typecode_Null  $\triangleq$  "CtrlCell_Null"

```

```

CtrlCell_SetupReq  $\triangleq$  [
  typecode : {CtrlCell_Typecode_SetupReq},
  connect  : Connection,
  DestC    : Container,
  DestH    : Halfrow,
  prevSlot : SlotName
  this L1 model terminates at the halfrow, hence there is no DestN
  this L1 model omits the L0 switches, hence there is no FirstL1
  this L1 model omits VLB routing
]

CtrlCell_SetupAck  $\triangleq$  [
  typecode : {CtrlCell_Typecode_SetupAck},
  connect  : Connection,
  prevSlot : SlotName,
  nextSlot : SlotName
]

CtrlCell_Teardown  $\triangleq$  [
  typecode : {CtrlCell_Typecode_Teardown},
  connect  : Connection,
  prevSlot : SlotName,
  reason   : Reason
]

CtrlCell_Discard  $\triangleq$  [
  typecode : {CtrlCell_Typecode_Discard},
  connect  : Connection,
  prevSlot : SlotName
]

CtrlCell_DataAck  $\triangleq$  [
  typecode : {CtrlCell_Typecode_DataAck},
  connect  : Connection,
  prevSlot : SlotName,
  throttle : BOOLEAN,
  chunkseqn : Chunkseqn
]

CtrlCell_Restart  $\triangleq$  [
  typecode : {CtrlCell_Typecode_Restart},
  connect  : Connection,
  prevSlot : SlotName
]

CtrlCell_Stop  $\triangleq$  [
  typecode : {CtrlCell_Typecode_Stop},
  connect  : Connection,

```

```

    prevSlot : SlotName
]
CtrlCell_Null ≙ [
    typecode : {CtrlCell_Typecode_Null}
]

```

```

CtrlCell ≙ {}
  ∪ CtrlCell_Null
  ∪ CtrlCell_SetupReq
  ∪ CtrlCell_SetupAck
  ∪ CtrlCell_Teardown
  ∪ CtrlCell_Discard
  ∪ CtrlCell_DataAck
  ∪ CtrlCell_Restart
  ∪ CtrlCell_Stop

```

```

NullCtrlCell ≙ [typecode ↦ CtrlCell_Typecode_Null]

```

```

CtrlCell_Is_SetupReq(cc) ≙ cc.typecode = CtrlCell_Typecode_SetupReq
CtrlCell_Is_SetupAck(cc) ≙ cc.typecode = CtrlCell_Typecode_SetupAck
CtrlCell_Is_Teardown(cc) ≙ cc.typecode = CtrlCell_Typecode_Teardown
CtrlCell_Is_Discard(cc) ≙ cc.typecode = CtrlCell_Typecode_Discard
CtrlCell_Is_DataAck(cc) ≙ cc.typecode = CtrlCell_Typecode_DataAck
CtrlCell_Is_Restart(cc) ≙ cc.typecode = CtrlCell_Typecode_Restart
CtrlCell_Is_Stop(cc) ≙ cc.typecode = CtrlCell_Typecode_Stop
CtrlCell_Is_Null(cc) ≙ cc.typecode = CtrlCell_Typecode_Null

```

11 End node

MODULE *DCNL1Node*

```

EXTENDS
  DCNL1Useful,
  DCNL1Param,
  DCNL1SlotName,
  DCNL1Cell,
  DCNL1CellType

```

Emulate an endnode nic that lives at an L1 switch halfrow.

CONSTANT *MaxChunkCells* maximum size of a chunk
BugRecvWrongPacket \triangleq FALSE receive some arbitrary packet
BugRecvWrongCRC \triangleq FALSE receive incorrect *crc*
TraceRecvPacket \triangleq FALSE

Per slot input state.

NodeInSlotState_Avail \triangleq "Avail"
NodeInSlotState_WaitFirstMark \triangleq "WaitFirstMark"
NodeInSlotState_WaitChunk \triangleq "WaitChunk"
NodeInSlotState_End \triangleq "End"

NodeInSlotState \triangleq {
 NodeInSlotState_Avail,
 NodeInSlotState_WaitFirstMark,
 NodeInSlotState_WaitChunk,
 NodeInSlotState_End
}

NodeInSlot \triangleq [
 i_state : *NodeInSlotState*,
 i_connect : *Nullable(Connection)*,
 i_prevSlot : *Nullable(SlotName)*,
 i_chunkseqn : *Chunkseqn*,
 i_chunkcrc : *Nullable(CRC)*,
 i_chunk : *Seq(Payload)*,
 i_packet : *Seq(Payload)*
]

InitNodeInSlot \triangleq [
 i_state \mapsto *NodeInSlotState_Avail*,
 i_connect \mapsto *NULL*,
 i_prevSlot \mapsto *NULL*,
 i_chunkseqn \mapsto *FirstChunkseqn*,
 i_chunkcrc \mapsto *NULL*,
 i_chunk \mapsto $\langle \rangle$,
 i_packet \mapsto $\langle \rangle$
]

Per slot output state.

NodeOutSlotState_Avail \triangleq "Avail"
NodeOutSlotState_SendSetupReq \triangleq "SendSetupReq"
NodeOutSlotState_WaitSetupAck \triangleq "WaitSetupAck"
NodeOutSlotState_SendFirstMark \triangleq "SendFirstMark"
NodeOutSlotState_SendChunkData \triangleq "SendChunkData"
NodeOutSlotState_SendChunkMark \triangleq "SendChunkMark"
NodeOutSlotState_WaitDataAck \triangleq "WaitDataAck"
NodeOutSlotState_SendEndMark \triangleq "SendEndMark"
NodeOutSlotState_WaitTeardown \triangleq "WaitTeardown"
NodeOutSlotState_Abort \triangleq "Abort"
NodeOutSlotState_End \triangleq "End"

NodeOutSlotState \triangleq {
 NodeOutSlotState_Avail,
 NodeOutSlotState_SendSetupReq,
 NodeOutSlotState_WaitSetupAck,
 NodeOutSlotState_SendFirstMark,
 NodeOutSlotState_SendChunkData,
 NodeOutSlotState_SendChunkMark,
 NodeOutSlotState_WaitDataAck,
 NodeOutSlotState_SendEndMark,
 NodeOutSlotState_WaitTeardown,
 NodeOutSlotState_Abort,
 NodeOutSlotState_End
}

NodeOutSlot \triangleq [
 o_state : *NodeOutSlotState*,
 o_connect : *Nullable(Connection)*,
 o_destC : *Nullable(Container)*,
 o_destH : *Nullable(Halfrow)*,
 o_nextSlot : *Nullable(SlotName)*,
 o_chunkseqn : *Chunkseqn*,
 o_chunkcrc : *Nullable(CRC)*,
 o_ixatchunk : *Nat*,
 o_ixinchunk : *Nat*,
 o_packet : *Seq(Payload)*
]

InitNodeOutSlot \triangleq [
 o_state \mapsto *NodeOutSlotState_Avail*,
 o_connect \mapsto *NULL*,
 o_destC \mapsto *NULL*,
 o_destH \mapsto *NULL*,
 o_nextSlot \mapsto *NULL*,

```

o_chunkseqn ↦ FirstChunkseqn,
o_chunkcrc  ↦ NULL,
o_ixatchunk ↦ 1,
o_ixinchunk ↦ 1,
o_packet    ↦ ⟨⟩
]

```

```

Node ≜ [
  outq : [CellType → Seq(Cell)],
  iss  : [SlotName → NodeInSlot],
  oss  : [SlotName → NodeOutSlot],
  slot : SlotName
]

```

```

InitNode ≜ [
  outq ↦ [t ∈ CellType ↦ ⟨⟩],
  iss  ↦ [sn ∈ SlotName ↦ InitNodeInSlot],
  oss  ↦ [sn ∈ SlotName ↦ InitNodeOutSlot],
  slot ↦ FirstSlotName
]

```

Set of received packets at this moment in time.

$Node_ReceivedPackets(n) \triangleq$

LET

slots that have finished receiving a packet

$sns \triangleq \{sn \in SlotName : n.iss[sn].i_state = NodeInSlotState_End\}$

packet received in slot sn

$pkt(sn) \triangleq$ IF *BugRecvWrongPacket* THEN CHOOSE $pkt0 \in Seq(Payload) : TRUE$ ELSE $n.iss[sn].i_packet$

IN

$\{pkt(sn) : sn \in sns\}$

Determine if local slot is available outgoing / incoming.

$Node_IsAvailOutSlot(n, sn) \triangleq n.oss[sn].o_state = NodeOutSlotState_Avail$

$Node_IsAvailInSlot(n, sn) \triangleq n.iss[sn].i_state = NodeInSlotState_Avail$

Set of available outgoing / incoming local slots.

$Node_AvailOutSlots(n) \triangleq \{sn \in SlotName : Node_IsAvailOutSlot(n, sn)\}$

$Node_AvailInSlots(n) \triangleq \{sn \in SlotName : Node_IsAvailInSlot(n, sn)\}$

Create a packet request.

```
Node_CreateRequest(n, sn, C, H, packet, connect)  $\triangleq$ 
CASE TypeAssert(n, Node) →
CASE TypeAssert(sn, SlotName) →
CASE TypeAssert(C, Container) →
CASE TypeAssert(H, Halfrow) →
CASE TypeAssert(packet, Seq(Payload)) →
CASE TypeAssert(connect, Connection) →

IF n.oss[sn].o_state ≠ NodeOutSlotState_Avail
THEN Assert(FALSE, "protocol error") ELSE
[n EXCEPT
!.oss[sn].o_state = NodeOutSlotState_SendSetupReq,
!.oss[sn].o_connect = connect,
!.oss[sn].o_destC = C,
!.oss[sn].o_destH = H,
!.oss[sn].o_packet = packet
]
```

Update a node output slot state for slot time.

```
Node_Oss_SlotTime(n, sn)  $\triangleq$ 
LET
  oss  $\triangleq$  n.oss[sn]
```

Send a setup request.

Sending a control cell is not logically tied to the slot time, but checking it and queuing the cell at this time makes multiplexing our attention easier.

```
SendSetupReq  $\triangleq$ 
[n EXCEPT
!.outq[CellType_Ctrl] = Append(@, [
  typecode ↦ CtrlCell_Typecode_SetupReq,
  connect ↦ oss.o_connect,
  DestC ↦ oss.o_destC,
  DestH ↦ oss.o_destH,
  prevSlot ↦ sn
]),
!.oss[sn].o_state = NodeOutSlotState_WaitSetupAck,
!.oss[sn].o_ixatchunk = 1,
```

```

    !.oss[sn].o_ixinchunk = 0
  ]

```

Send a first mark.

```

SendFirstMark  $\triangleq$ 
[n EXCEPT
  !.outq[CellType_Data] = Append(@, [
    typecode  $\mapsto$  DataCell_Typecode_FirstMark,
    connect  $\mapsto$  oss.o_connect,
    nextSlot  $\mapsto$  oss.o_nextSlot,
    destC  $\mapsto$  oss.o_destC,
    destH  $\mapsto$  oss.o_destH
  ]),
  !.oss[sn].o_state = NodeOutSlotState_SendChunkData,
  !.oss[sn].o_chunkcrc = InitCRC,
  !.oss[sn].o_ixatchunk = 1,
  !.oss[sn].o_ixinchunk = 0,
  !.oss[sn] =
  IF  $\forall$  @.o_ixatchunk + @.o_ixinchunk > Len(@.o_packet)
     $\vee$  @.o_ixinchunk  $\geq$  MaxChunkCells
  THEN [@ EXCEPT !.o_state = NodeOutSlotState_SendChunkMark]
  ELSE @
]

```

Send chunk data.

```

SendChunkData  $\triangleq$ 
LET
  ix  $\triangleq$  oss.o_ixatchunk + oss.o_ixinchunk
  py  $\triangleq$  oss.o_packet[ix]
IN
[n EXCEPT
  !.outq[CellType_Data] = Append(@, [
    typecode  $\mapsto$  DataCell_Typecode_Normal,
    connect  $\mapsto$  oss.o_connect,
    nextSlot  $\mapsto$  oss.o_nextSlot,
    payload  $\mapsto$  py
  ]),
  !.oss[sn].o_chunkcrc = ExtendCRC(@, py),
  !.oss[sn].o_ixinchunk = @ + 1,
  !.oss[sn] =
  IF  $\forall$  @.o_ixatchunk + @.o_ixinchunk > Len(@.o_packet)
     $\vee$  @.o_ixinchunk  $\geq$  MaxChunkCells
  THEN [@ EXCEPT !.o_state = NodeOutSlotState_SendChunkMark]
  ELSE @
]

```

Send chunk mark.

```

SendChunkMark  $\triangleq$ 
[n EXCEPT
!.outq[CellType_Data] = Append(@, [
    typecode  $\mapsto$  DataCell_Typecode_ChunkMark,
    connect  $\mapsto$  oss.o_connect,
    nextSlot  $\mapsto$  oss.o_nextSlot,
    chunkseqn  $\mapsto$  oss.o_chunkseqn,
    crc  $\mapsto$  oss.o_chunkcrc
]),
!.oss[sn].o_state = NodeOutSlotState_WaitDataAck,
!.oss[sn].o_chunkcrc = NULL
]

```

Send end mark.

```

SendEndMark  $\triangleq$ 
[n EXCEPT
!.outq[CellType_Data] = Append(@, [
    typecode  $\mapsto$  DataCell_Typecode_EndMark,
    connect  $\mapsto$  oss.o_connect,
    nextSlot  $\mapsto$  oss.o_nextSlot
]),
!.oss[sn].o_state = NodeOutSlotState_WaitTeardown
]

```

Recycle output slot.

```

RecycleSlot  $\triangleq$ 
[n EXCEPT
!.oss[sn] = InitNodeOutSlot
]

```

IN

CASE

<i>oss.o_state</i> = <i>NodeOutSlotState_Avail</i>	$\rightarrow n$	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_SendSetupReq</i>	\rightarrow <i>SendSetupReq</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_WaitSetupAck</i>	$\rightarrow n$	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_SendFirstMark</i>	\rightarrow <i>SendFirstMark</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_SendChunkData</i>	\rightarrow <i>SendChunkData</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_SendChunkMark</i>	\rightarrow <i>SendChunkMark</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_WaitDataAck</i>	$\rightarrow n$	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_SendEndMark</i>	\rightarrow <i>SendEndMark</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_WaitTeardown</i>	$\rightarrow n$	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_Abort</i>	\rightarrow <i>RecycleSlot</i>	\square
<i>oss.o_state</i> = <i>NodeOutSlotState_End</i>	\rightarrow <i>RecycleSlot</i>	\square

Update a node for receiving a setup ack control cell.

```

Node_RecvCtrlCell_SetupAck(n, cc)  $\triangleq$ 
LET
  sn  $\triangleq$  cc.prevSlot
  oss  $\triangleq$  n.oss[sn]

  RecvSetupAck  $\triangleq$ 
    CASE Assert(cc.connect = oss.o_connect, "mismatch connection")  $\rightarrow$ 
      [n EXCEPT
        !.oss[sn].o_state = NodeOutSlotState_SendFirstMark,
        !.oss[sn].o_nextSlot = cc.nextSlot,
        !.oss[sn].o_chunkseqn = FirstChunkseqn
      ]
IN
CASE
  oss.o_state = NodeOutSlotState_Avail  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_SendSetupReq  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_WaitSetupAck  $\rightarrow$  RecvSetupAck □
  oss.o_state = NodeOutSlotState_SendFirstMark  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_SendChunkData  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_SendChunkMark  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_WaitDataAck  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_SendEndMark  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_WaitTeardown  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_Abort  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_End  $\rightarrow$  n

```

Update a node for receiving a teardown control cell.

```

Node_RecvCtrlCell_Teardown(n, cc)  $\triangleq$ 
LET
  sn  $\triangleq$  cc.prevSlot
  oss  $\triangleq$  n.oss[sn]

  DoAbort  $\triangleq$ 
    CASE Assert(cc.connect = oss.o_connect, "mismatch connection")  $\rightarrow$ 
      [n EXCEPT !.oss[sn].o_state = NodeOutSlotState_Abort]

  DoEnd  $\triangleq$ 
    CASE Assert(cc.connect = oss.o_connect, "mismatch connection")  $\rightarrow$ 
      [n EXCEPT !.oss[sn].o_state = NodeOutSlotState_End]
IN
CASE
  oss.o_state = NodeOutSlotState_Avail  $\rightarrow$  n □
  oss.o_state = NodeOutSlotState_SendSetupReq  $\rightarrow$  DoAbort □

```

<i>oss.o_state</i> = <i>NodeOutSlotState_WaitSetupAck</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_SendFirstMark</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_SendChunkData</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_SendChunkMark</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_WaitDataAck</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_SendEndMark</i>	→ <i>DoAbort</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_WaitTeardown</i>	→ <i>DoEnd</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_Abort</i>	→ <i>n</i>	□
<i>oss.o_state</i> = <i>NodeOutSlotState_End</i>	→ <i>n</i>	□

Update a node for receiving a data ack control cell.

Node_RecvCtrlCell_DataAck(*n*, *cc*) \triangleq

LET

sn \triangleq *cc.prevSlot*
oss \triangleq *n.oss[sn]*

RecvDataAck \triangleq

CASE *Assert*(*cc.connect* = *oss.o_connect*, "mismatch connection") →

IF *cc.chunkseqn* \neq *oss.o_chunkseqn* THEN *n* ELSE

Got the expected data ack, go on to the next chunk.

[*n* EXCEPT

!.*oss[sn].o_state* = *NodeOutSlotState_SendChunkData*,

!.*oss[sn].o_chunkseqn* = *IncrChunkseqn*(@),

!.*oss[sn].o_chunkcrc* = *InitCRC*,

!.*oss[sn].o_ixatchunk* = @ + *oss.o_ixinchunk*,

!.*oss[sn].o_ixinchunk* = 0,

!.*oss[sn]* =

IF @.*o_ixatchunk* + @.*o_ixinchunk* > *Len*(@.*o_packet*)

THEN [@ EXCEPT !.*o_state* = *NodeOutSlotState_SendEndMark*]

ELSE @

]

IN

CASE

oss.o_state = *NodeOutSlotState_Avail* → *n* □

oss.o_state = *NodeOutSlotState_SendSetupReq* → *n* □

oss.o_state = *NodeOutSlotState_WaitSetupAck* → *n* □

oss.o_state = *NodeOutSlotState_SendFirstMark* → *n* □

oss.o_state = *NodeOutSlotState_SendChunkData* → *n* □

oss.o_state = *NodeOutSlotState_SendChunkMark* → *n* □

oss.o_state = *NodeOutSlotState_WaitDataAck* → *RecvDataAck* □

oss.o_state = *NodeOutSlotState_SendEndMark* → *n* □

oss.o_state = *NodeOutSlotState_WaitTeardown* → *n* □

$oss.o_state = NodeOutSlotState_Abort \quad \rightarrow n \quad \square$
 $oss.o_state = NodeOutSlotState_End \quad \rightarrow n$

Update a node input state for slot sn .

$Node_Iss_SlotTime(n, sn) \triangleq$

LET
 $iss \triangleq n.iss[sn]$

A convenient time to take the finished received packet away and reset the state of the connection.

$AtEnd \triangleq$

CASE IF $TraceRecvPacket$ THEN $PrintT(\langle \text{"recv pkt"}, iss.i_packet \rangle)$ ELSE TRUE \rightarrow
 $[n \text{ EXCEPT } !.iss[sn] = InitNodeInSlot]$

IN

CASE

$iss.i_state = NodeInSlotState_Avail \quad \rightarrow n \quad \square$
 $iss.i_state = NodeInSlotState_WaitFirstMark \rightarrow n \quad \square$
 $iss.i_state = NodeInSlotState_WaitChunk \quad \rightarrow n \quad \square$
 $iss.i_state = NodeInSlotState_End \quad \rightarrow AtEnd$

Advance the slot in a node

$Node_AdvanceSlot(n) \triangleq$

$[[n \mapsto n] \text{ EXCEPT}$
 $!.n = Node_Iss_SlotTime(@, @.slot),$
 $!.n = Node_Oss_SlotTime(@, @.slot),$
 $!.n.slot = IncrSlotName[@]$
 $].n$

Receive a setup request control cell into node n .

$Node_RecvCtrlCell_SetupReq(n, cc) \triangleq$

LET

$prevSlot \triangleq cc.prevSlot$

$sns \triangleq Node_AvailInSlots(n)$

$sn \triangleq \text{CHOOSE } sn1 \in sns : \forall sn2 \in sns : LeqSlotName[sn1, sn2]$

```

iss      ≜ n.iss[sn]

DoNak ≜ [n EXCEPT
        !.outq[CellType_Ctrl] = Append(@, [
            typecode ↦ CtrlCell_Typecode_Teardown,
            connect  ↦ cc.connect,
            prevSlot ↦ prevSlot,
            reason   ↦ "dest full"
        ])
    ]

DoAck ≜ [n EXCEPT
        !.outq[CellType_Ctrl] = Append(@, [
            typecode ↦ CtrlCell_Typecode_SetupAck,
            connect  ↦ cc.connect,
            prevSlot ↦ prevSlot,
            nextSlot ↦ sn
        ]),
        !.iss[sn].i_state = NodeInSlotState_WaitFirstMark,
        !.iss[sn].i_connect = cc.connect,
        !.iss[sn].i_prevSlot = prevSlot,
        !.iss[sn].i_packet = ⟨⟩
    ]

IN
IF sns = {} THEN DoNak ELSE DoAck

```

Receive a first mark data cell into node n .

```
Node_RecvDataCell_FirstMark( $n$ ,  $dc$ ) ≜
```

```
LET
  sn      ≜ dc.nextSlot
  iss     ≜ n.iss[sn]

```

If we are waiting for a first mark . . .

```
ok_iss ≜ [iss EXCEPT
        !.i_state = NodeInSlotState_WaitChunk,
        !.i_chunkseqn = FirstChunkseqn,
        !.i_chunkcrc = InitCRC,
        !.i_chunk = ⟨⟩,
        !.i_packet = ⟨⟩
    ]

ok_node ≜
  CASE Assert( $dc.connect = iss.i_connect$ , "mismatch connection") →
  [n EXCEPT !.iss[sn] = ok_iss]

IN

```

IF $iss.i_state \neq NodeInSlotState_WaitFirstMark$ THEN n ELSE ok_node

Receive a normal data cell into node n .

$Node_RecvDataCell_Normal(n, dc) \triangleq$

LET

$sn \triangleq dc.nextSlot$
 $iss \triangleq n.iss[sn]$

If we are waiting for a normal data cell . . .

$ok_iss \triangleq [iss \text{ EXCEPT}$
 $\quad !.i_chunkcrc = ExtendCRC(@, dc.payload),$
 $\quad !.i_chunk = Append(@, dc.payload)$
 $\quad]$

$ok_node \triangleq$
CASE $Assert(dc.connect = iss.i_connect, \text{"mismatch connection"}) \rightarrow$
 $[n \text{ EXCEPT } !.iss[sn] = ok_iss]$

IN

IF $iss.i_state \neq NodeInSlotState_WaitChunk$ THEN n ELSE ok_node

Receive a chunk mark data cell into node n .

$Node_RecvDataCell_ChunkMark(n, dc) \triangleq$

LET

$sn \triangleq dc.nextSlot$
 $iss \triangleq n.iss[sn]$

Send data ack control cell corresponding to this chunk mark.

$SendAck(n0) \triangleq$
 $[n0 \text{ EXCEPT}$
 $\quad !.outq[CellType_Ctrl] = Append(@, [$
 $\quad \quad typecode \mapsto CtrlCell_Typecode_DataAck,$
 $\quad \quad connect \mapsto iss.i_connect,$
 $\quad \quad prevSlot \mapsto iss.i_prevSlot,$
 $\quad \quad throttle \mapsto FALSE,$
 $\quad \quad chunkseqn \mapsto dc.chunkseqn$
 $\quad])$
 $\quad]$

Accept the current chunk. Append it to the packet being received and increment the chunk sequence number.

$AcceptChunk(n0) \triangleq$

$[n0 \text{ EXCEPT}$
 $\quad !.iss[sn].i_packet = @ \circ iss.i_chunk,$

```

    !.iss[sn].i_chunkseqn = IncrChunkseqn(@)
  ]
  Start receiving the next chunk.
  StartChunk(n0) ≜
  [n0 EXCEPT
    !.iss[sn].i_chunkcrc = IF BugRecvWrongCRC THEN @ ELSE
      InitCRC,
    !.iss[sn].i_chunk = ⟨⟩
  ]
IN
IF iss.i_state ≠ NodeInSlotState_WaitChunk THEN n ELSE
CASE
  (iss.i_chunkseqn = dc.chunkseqn ∧ iss.i_chunkcrc ≠ dc.crc)
  ⇒ PrintT("crc error")
  →
CASE Assert(dc.connect = iss.i_connect, "mismatch connection") →
IF iss.i_chunkseqn = dc.chunkseqn ∧ iss.i_chunkcrc = dc.crc
THEN
  Chunk mark has the sequence number of the next chunk we are expecting and the crc is correct.
  [[n ↦ n] EXCEPT
    !.n = SendAck(@),
    !.n = AcceptChunk(@),
    !.n = StartChunk(@)
  ].n
ELSE
IF iss.i_chunkseqn = IncrChunkseqn(dc.chunkseqn)
THEN
  Chunk mark repeats the sequence number of the last chunk we accepted.
  [[n ↦ n] EXCEPT
    !.n = SendAck(@),
    !.n = StartChunk(@)
  ].n
ELSE
  Something else.
  [[n ↦ n] EXCEPT
    !.n = StartChunk(@)
  ].n

```

Receive an end mark data cell into node n .

$Node_RecvDataCell_EndMark(n, dc) \triangleq$

```

LET
  sn       $\triangleq$  dc.nextSlot
  iss      $\triangleq$  n.iss[sn]

  endcc   $\triangleq$  [typecode  $\mapsto$  CtrlCell_Typecode_Teardown,
                connect   $\mapsto$  iss.i_connect,
                prevSlot  $\mapsto$  iss.i_prevSlot,
                reason    $\mapsto$  "end mark"
              ]

```

If we accept the end mark . . .

```

ok_iss   $\triangleq$  [iss EXCEPT
                !.i_state = NodeInSlotState_End,
                !.i_chunk =  $\langle \rangle$ ,
                !.i_packet = @
              ]
ok_node  $\triangleq$ 
  CASE Assert(dc.connect = iss.i_connect, "mismatch connection")  $\rightarrow$ 
  [n EXCEPT
   !.outq[CellType_Ctrl] = Append(@, endcc),
   !.iss[sn] = ok_iss
  ]
IN
IF iss.i_state  $\neq$  NodeInSlotState_WaitChunk THEN n ELSE ok_node

```

Receive a data cell from the link into node n .

$Node_RecvDataCell(n, dc) \triangleq$

```

CASE
  DataCell_Is_FirstMark(dc)  $\rightarrow$  Node_RecvDataCell_FirstMark(n, dc)  $\square$ 
  DataCell_Is_Normal(dc)    $\rightarrow$  Node_RecvDataCell_Normal(n, dc)    $\square$ 
  DataCell_Is_ChunkMark(dc)  $\rightarrow$  Node_RecvDataCell_ChunkMark(n, dc)  $\square$ 
  DataCell_Is_EndMark(dc)   $\rightarrow$  Node_RecvDataCell_EndMark(n, dc)   $\square$ 
  DataCell_Is_Null(dc)      $\rightarrow$  n

```

Receive a control cell from the link into node n .

$Node_RecvCtrlCell(n, cc) \triangleq$

```

CASE
  CtrlCell_Is_SetupReq(cc)  $\rightarrow$  Node_RecvCtrlCell_SetupReq(n, cc)  $\square$ 
  CtrlCell_Is_SetupAck(cc)  $\rightarrow$  Node_RecvCtrlCell_SetupAck(n, cc)  $\square$ 
  CtrlCell_Is_DataAck(cc)   $\rightarrow$  Node_RecvCtrlCell_DataAck(n, cc)  $\square$ 
  CtrlCell_Is_Teardown(cc)  $\rightarrow$  Node_RecvCtrlCell_Teardown(n, cc)  $\square$ 
  CtrlCell_Is_Null(cc)      $\rightarrow$  n

```

Interface to link.

Receive a cell of the indicated type from the link into node n .

$Node_RecvCell(n, t, cell) \triangleq$
LET
 $ready \triangleq TRUE$
 $update \triangleq$
 CASE
 $t = CellType_Data \rightarrow Node_RecvDataCell(n, cell) \square$
 $t = CellType_Ctrl \rightarrow Node_RecvCtrlCell(n, cell)$
IN
IF $\neg ready$ THEN [$ready \mapsto ready$] ELSE
[
 $ready \mapsto ready,$
 $update \mapsto update$
]

Extract a cell of the indicated type onto the link from node n .

$Node_TakeCell(n, t) \triangleq$
LET
 $ready \triangleq n.outq[t] \neq \langle \rangle$
 $cell \triangleq Head(n.outq[t])$
 $update \triangleq [n \text{ EXCEPT } !.outq[t] = Tail(@)]$
IN
IF $\neg ready$ THEN [$ready \mapsto ready$] ELSE
[
 $ready \mapsto ready,$
 $cell \mapsto cell,$
 $update \mapsto update$
]

12 Ring message

MODULE *DCNL1RingMsg*

EXTENDS
 DCNL1Param,

DCNL1SlotName,
DCNL1Chunkseqn,
DCNL1PortName

Messages on the ring.

In the hardware, each ring message except *SetupReq* fits in one 40 bit word. The *SetupReq* message takes two 40 bit words. However, we ignore this difference.

<i>RingMsg_TypeCode_RedSetupReq</i>	\triangleq	“RingMsg_RedSetupReq”	up port to down port
<i>RingMsg_TypeCode_RedSetupAns</i>	\triangleq	“RingMsg_RedSetupAns”	down port to up port
<i>RingMsg_TypeCode_RedSetupXbar</i>	\triangleq	“RingMsg_RedSetupXbar”	up port to up port
<i>RingMsg_TypeCode_RedNull</i>	\triangleq	“RingMsg_RedNull”	
<i>RingMsg_TypeCode_GreenTeardown</i>	\triangleq	“RingMsg_GreenTeardown”	down port to up port
<i>RingMsg_TypeCode_GreenGCClear</i>	\triangleq	“RingMsg_GreenGCClear”	down port to up port
<i>RingMsg_TypeCode_GreenDiscard</i>	\triangleq	“RingMsg_GreenDiscard”	down port to up port
<i>RingMsg_TypeCode_GreenDataAck</i>	\triangleq	“RingMsg_GreenDataAck”	down port to up port
<i>RingMsg_TypeCode_GreenRestart</i>	\triangleq	“RingMsg_GreenRestart”	down port to up port
<i>RingMsg_TypeCode_GreenStop</i>	\triangleq	“RingMsg_GreenStop”	down port to up port
<i>RingMsg_TypeCode_GreenNull</i>	\triangleq	“RingMsg_GreenNull”	

RingMsg_RedSetupReq \triangleq [
typecode : { *RingMsg_TypeCode_RedSetupReq* },
connect : *Connection*,
DestC : *Container*,
DestH : *Halfrow*,
prevSlot : *SlotName*,
upPort : *PortName*,
downPort : *PortName*,
subframe : *Subframe*,
occupied : [*Slottina* \rightarrow BOOLEAN]
 this L1 model terminates at the halfrow, hence there is no *DestN*
 this L1 model omits the L0 switches, hence there is no *FirstL1*
 this L1 model omits VLB routing
]

RingMsg_RedSetupAns \triangleq [
 The hardware encodes a setup answer ring message as a setup request with *processed* = TRUE. For the specification it is cleaner to use a separate ring message for the setup answer.
typecode : { *RingMsg_TypeCode_RedSetupAns* },
connect : *Connection*,
DestC : *Container*,
DestH : *Halfrow*,
prevSlot : *SlotName*,
upPort : *PortName*,
downPort : *PortName*,
]

```

    subframe : Subframe,
    occupied  : [Slottina → BOOLEAN ]
]

```

```

RingMsg_RedSetupXbar  $\triangleq$  [

```

The hardware encodes a setup xbar ring message as a setup request with processed = TRUE and xDone = TRUE. For the specification it is cleaner to use a separate ring message for the setup xbar.

```

    typecode : { RingMsg_Typecode_RedSetupXbar },
    connect  : Connection,
    DestC    : Container,
    DestH    : Halfrow,
    prevSlot : SlotName,
    upPort   : PortName,
    downPort : PortName,
    subframe : Subframe,
    occupied  : [Slottina → BOOLEAN ]
]

```

```

RingMsg_RedNull  $\triangleq$  [
    typecode : { RingMsg_Typecode_RedNull }
]

```

```

RingMsg_GreenTeardown  $\triangleq$  [
    typecode : { RingMsg_Typecode_GreenTeardown },
    connect  : Connection,
    prevSlot : SlotName,
    localSlot : SlotName,
    upPort   : PortName,
    reason   : Reason
]

```

```

RingMsg_GreenGCClear  $\triangleq$  [
    typecode : { RingMsg_Typecode_GreenGCClear },
    connect  : Connection,
    upPort   : PortName,
    localSlot : SlotName
]

```

```

RingMsg_GreenDiscard  $\triangleq$  [
    typecode : { RingMsg_Typecode_GreenDiscard },
    connect  : Connection,
    prevSlot : SlotName,
    upPort   : PortName
]

```

```

RingMsg_GreenDataAck  $\triangleq$  [
    typecode : { RingMsg_Typecode_GreenDataAck },

```

connect : *Connection*,
prevSlot : *SlotName*,
upPort : *PortName*,
chunkseqn : *Chunkseqn*,
throttle : *Throttle*
]

RingMsg_GreenRestart \triangleq [
typecode : {*RingMsg_Typecode_GreenRestart*},
connect : *Connection*,
prevSlot : *SlotName*,
upPort : *PortName*
]

RingMsg_GreenStop \triangleq [
typecode : {*RingMsg_Typecode_GreenStop*},
connect : *Connection*,
prevSlot : *SlotName*,
upPort : *PortName*
]

RingMsg_GreenNull \triangleq [
typecode : {*RingMsg_Typecode_GreenNull*}
]

RingMsg \triangleq {
 \cup *RingMsg_RedSetupReq*
 \cup *RingMsg_RedSetupAns*
 \cup *RingMsg_RedSetupXbar*
 \cup *RingMsg_RedNull*
 \cup *RingMsg_GreenTeardown*
 \cup *RingMsg_GreenGCClear*
 \cup *RingMsg_GreenDiscard*
 \cup *RingMsg_GreenDataAck*
 \cup *RingMsg_GreenRestart*
 \cup *RingMsg_GreenStop*
 \cup *RingMsg_GreenNull*

RedNullRingMsg \triangleq [*typecode* \mapsto *RingMsg_Typecode_RedNull*]

GreenNullRingMsg \triangleq [*typecode* \mapsto *RingMsg_Typecode_GreenNull*]

RingMsg_Is_RedSetupReq(*rm*) \triangleq *rm.typecode* = *RingMsg_Typecode_RedSetupReq*

RingMsg_Is_RedSetupAns(*rm*) \triangleq *rm.typecode* = *RingMsg_Typecode_RedSetupAns*

RingMsg_Is_RedSetupXbar(*rm*) \triangleq *rm.typecode* = *RingMsg_Typecode_RedSetupXbar*

RingMsg_Is_RedNull(*rm*) \triangleq *rm.typecode* = *RingMsg_Typecode_RedNull*

RingMsg_Is_GreenTeardown(*rm*) \triangleq *rm.typecode* = *RingMsg_Typecode_GreenTeardown*

$$\begin{aligned}
RingMsg_Is_GreenGCClear(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenGCClear \\
RingMsg_Is_GreenDiscard(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenDiscard \\
RingMsg_Is_GreenDataAck(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenDataAck \\
RingMsg_Is_GreenRestart(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenRestart \\
RingMsg_Is_GreenStop(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenStop \\
RingMsg_Is_GreenNull(rm) &\triangleq rm.typecode = RingMsg_Typecode_GreenNull
\end{aligned}$$

The setup request, setup answer, and setup xbar messages all use a *subframe* number and occupied map that need to be evaluated in the same way, to see if there is a *localSlot* available and if so, what the least such slot is. Here is how we compute this.

$$\begin{aligned}
RingMsg_Eval_Occupied(subframe, occupied) &\triangleq \\
\text{LET} & \\
sns &\triangleq \{sn \in SlotName : \\
&\quad \wedge sn.subframe = subframe \\
&\quad \wedge \neg occupied[sn.slottina] \\
&\quad \} \\
success &\triangleq sns \neq \{\} \\
localSlot &\triangleq \text{CHOOSE } sn1 \in sns : \forall sn2 \in sns : LeqSlotName[sn1, sn2] \\
\text{IN} & \\
\text{IF } \neg success &\text{ THEN } [success \mapsto success] \text{ ELSE} \\
[& \\
success &\mapsto success, \\
localSlot &\mapsto localSlot \\
] &
\end{aligned}$$

13 Port name

MODULE *DCNL1PortName*

EXTENDS *DCNL1Param*

Port names inside an *L1* switch.

$$\begin{aligned}
PortName_Typecode_Halfrow &\triangleq \text{"PortName_Halfrow"} \\
PortName_Typecode_Container &\triangleq \text{"PortName_Container"}
\end{aligned}$$

$$\begin{aligned}
PortName_Halfrow &\triangleq [\\
typecode &: \{PortName_Typecode_Halfrow\},
\end{aligned}$$

```

    to      : Halfrow
  ]
PortName_Container  $\triangleq$  [
  typecode : {PortName_Typecode_Container},
  to      : Container
]
PortName  $\triangleq$  { }
   $\cup$  PortName_Halfrow
   $\cup$  PortName_Container

PortName_Is_Halfrow(pn)  $\triangleq$  pn.typecode = PortName_Typecode_Halfrow
PortName_Is_Container(pn)  $\triangleq$  pn.typecode = PortName_Typecode_Container

```

Define an arbitrary ordering of port names.

```

FirstPortName  $\triangleq$  ArbSeq(PortName).FirstName
IncrPortName  $\triangleq$  ArbSeq(PortName).IncrName
DecrPortName  $\triangleq$  ArbSeq(PortName).DecrName
LeqPortName  $\triangleq$  ArbSeq(PortName).LeqName

```

14 Port

MODULE *DCNL1Port*

EXTENDS

```

DCNL1Useful,
DCNL1Param,
DCNL1Cell,
DCNL1CellType,
DCNL1DataCell,
DCNL1CtrlCell,
DCNL1RingMsg,
DCNL1SwitchName,
DCNL1CtrlQueueType

```

Definition of the state maintained in a port.

The input state per slot of a port.

```

PortInSlot  $\triangleq$  [
  i_assign : BOOLEAN,           the connection is assigned
  i_connect : Nullable(Connection), connection identifier (spec only)
  i_buf : Seq(DataCell)        arriving data cells
]

```

```

InitPortInSlot  $\triangleq$  [
  i_assign  $\mapsto$  FALSE,
  i_connect  $\mapsto$  NULL,
  i_buf  $\mapsto$   $\langle \rangle$ 
]

```

The output state per slot of a port.

Some of the output related fields are declared as *Nullable* so that they can accept a *NULL* value. When the field is set to *NULL* it should never be read. Using a *NULL* value where it is not expected should lead to a type error during model checking.

```

PortOutSlot  $\triangleq$  [
  o_assign : BOOLEAN,           the connection is assigned
  o_connect : Nullable(Connection), connection identifier (spec only)
  o_nextset : BOOLEAN,         the next slot is set up
  o_prevSlot : Nullable(SlotName), upstream switch's localSlot
  o_nextSlot : Nullable(SlotName), downstream switch's localSlot
  o_upPort : Nullable(PortName) upstream port in this switch
]

```

```

InitPortOutSlot  $\triangleq$  [
  o_assign  $\mapsto$  FALSE,
  o_connect  $\mapsto$  NULL,
  o_nextset  $\mapsto$  FALSE,
  o_prevSlot  $\mapsto$  NULL,
  o_nextSlot  $\mapsto$  NULL,
  o_upPort  $\mapsto$  NULL
]

```

State of a port.

```

Port  $\triangleq$  [
  switchname : SwitchName,
  portname : PortName,
  link_outq : [CellType  $\rightarrow$  Seq(Cell)], departing cells by type
  locked : [Subframe  $\rightarrow$  BOOLEAN], which subframes are locked
  locked_connect : [Subframe  $\rightarrow$  Nullable(Connection)], for what connection
]

```

```

ccinq      : [CtrlQueueType → Seq(CtrlCell)],    ctrl cells in from link
iss        : [SlotName → PortInSlot],
oss        : [SlotName → PortOutSlot],

toxbars_in : DataCell      most recent cell sent into crossbar
]

```

```

InitPort(swn, pn) ≜ [
  switchname  ↦ swn,
  portname    ↦ pn,

  link_outq   ↦ [t ∈ CellType ↦ ⟨⟩],

  locked      ↦ [sf ∈ Subframe ↦ FALSE],
  locked_connect ↦ [sf ∈ Subframe ↦ NULL],

  ccinq       ↦ [t ∈ CtrlQueueType ↦ ⟨⟩],

  iss         ↦ [sn ∈ SlotName ↦ InitPortInSlot],
  oss         ↦ [sn ∈ SlotName ↦ InitPortOutSlot],

  toxbars_in  ↦ NullDataCell
]

```

Get the container for port p .

```
Port_MyContainer(p) ≜ p.switchname.container
```

Interface to link.

Receive a data cell from the link into port p . Compute the updated state of the port.

```

Port_RecvDataCell(p, dc) ≜
CASE TypeAssert(dc, DataCell) →

IF dc.typecode = DataCell_Typecode_Null
THEN
  ignore arriving null data cell
  p
ELSE
LET
  sn ≜ dc.nextSlot
  iss ≜ p.iss[sn]
IN

```

```

IF  $\neg iss.i\_assign$  THEN
  ignore arriving data cell whose next slot is not assigned for a connection
  p
ELSE
  enqueue arriving data cell based on its next slot
  CASE Assert(dc.connect = iss.i_connect,
    (“connection mismatch”, dc, p.portname, iss) →
    [p EXCEPT  $!.iss[sn].i\_buf = Append(@, dc)$ ]

```

Receive a control cell from the link into port *p*. Compute the updated state of the port.

```

Port_RecvCtrlCell(p, cc)  $\triangleq$ 
  LET
    pack  $\triangleq$  LET
      sn  $\triangleq$  cc.prevSlot
      oss  $\triangleq$  p.oss[sn]
    IN
      CASE Assert(cc.connect = oss.o_connect, “connection mismatch”) →
      IF  $\neg oss.o\_assign$  THEN p ELSE
      [p EXCEPT
         $!.oss[sn].o\_nextSlot = cc.nextSlot,$ 
         $!.oss[sn].o\_nextset = TRUE$ 
      ]
    pinq(t)  $\triangleq$  [p EXCEPT  $!.ccinq[t] = Append(@, cc)$ ]
  IN
  CASE
    CtrlCell_Is_SetupReq(cc) → pinq(CtrlQueueType_Req) □
    CtrlCell_Is_SetupAck(cc) → pack □
    CtrlCell_Is_Teardown(cc) → pinq(CtrlQueueType_Ups) □
    CtrlCell_Is_Discard(cc) → pinq(CtrlQueueType_Ups) □
    CtrlCell_Is_DataAck(cc) → pinq(CtrlQueueType_Ups) □
    CtrlCell_Is_Restart(cc) → pinq(CtrlQueueType_Ups) □
    CtrlCell_Is_Stop(cc) → pinq(CtrlQueueType_Ups) □
    CtrlCell_Is_Null(cc) → p

```

Receive a cell of the indicated type.

```

Port_RecvCell(p, t, cell)  $\triangleq$ 
  LET
    ready  $\triangleq$  TRUE
    update  $\triangleq$ 
    CASE
      t = CellType_Data → Port_RecvDataCell(p, cell) □

```

```

     $t = \text{CellType\_Ctrl} \rightarrow \text{Port\_RecvCtrlCell}(p, \text{cell})$ 
  IN
  IF  $\neg \text{ready}$  THEN  $[\text{ready} \mapsto \text{ready}]$  ELSE
  [
     $\text{ready} \mapsto \text{ready}$ ,
     $\text{update} \mapsto \text{update}$ 
  ]

```

Extract a cell of the indicated type onto the link from port p .

```

 $\text{Port\_TakeCell}(p, t) \triangleq$ 
  LET
     $\text{ready} \triangleq p.\text{link\_outq}[t] \neq \langle \rangle$ 
     $\text{cell} \triangleq \text{Head}(p.\text{link\_outq}[t])$ 
     $\text{update} \triangleq [p \text{ EXCEPT } !.\text{link\_outq}[t] = \text{Tail}(@)]$ 
  IN
  IF  $\neg \text{ready}$  THEN  $[\text{ready} \mapsto \text{ready}]$  ELSE
  [
     $\text{ready} \mapsto \text{ready}$ ,
     $\text{cell} \mapsto \text{cell}$ ,
     $\text{update} \mapsto \text{update}$ 
  ]

```

Send data cell dc from port p . Compute the updated state of the port.

Each port sends one data cell per slot time. In the hardware, this feeds into the framer which constructs frames for transmission on the link.

We do not model link frames in the specification, so we just queue data cells for transmission on the link.

```

 $\text{Port\_SendDataCell}(p, dc) \triangleq$ 
  CASE  $\text{TypeAssert}(dc, \text{DataCell}) \rightarrow$ 
  [
     $[p \text{ EXCEPT } !.\text{link\_outq}[\text{CellType\_Data}] = \text{Append}(@, dc)]$ 
  ]

```

Send ctrl cell cc from port p . Compute the updated state of the port.

In the hardware, this feeds into a queue which the framer draws upon to construct frames for transmission on the link. In the hardware, there are 128 data cells and 64 control cells in each frame.

We do not model link frames in the specification, so we just queue control cells for transmission on the link.

```

 $\text{Port\_SendCtrlCell}(p, cc) \triangleq$ 
  CASE  $\text{TypeAssert}(cc, \text{CtrlCell}) \rightarrow$ 

```

$[p \text{ EXCEPT } !.link_outq[CellType_Ctrl] = Append(@, cc)]$

Interface to crossbar.

Send a data cell to the crossbar from port p in slot name sn . Compute a record with the following fields:

p = the updated state of the port

dc = the data cell sent into the crossbar

$Port_Xbar_TakeSlot(p, sn) \triangleq$

```

LET
   $dc \triangleq HeadElse(p.iss[sn].i\_buf, NullDataCell)$ 
IN
  [
     $p \mapsto [p \text{ EXCEPT}$ 
       $!.to\_xbar\_in = dc,$ 
       $!.iss[sn].i\_buf = TailElse(@)$ 
    ],
     $dc \mapsto dc$ 
  ]

```

Receive a data cell from the crossbar into port p during slot name sn . Compute the updated state of the port.

$Port_Xbar_RecvSlot(p, sn, dc) \triangleq$

```

LET
   $oss \triangleq p.oss[sn]$ 
IN
  IF  $dc.typecode = DataCell\_Typecode\_Null$  THEN
    a null data cell coming out of the crossbar – ignore it
  ELSE
    IF
       $\vee \neg oss.o\_assign$ 
       $\vee \neg oss.o\_nextset$ 
    THEN
      this slot is not assigned or we do not yet know the next slot – ignore the data cell
    ELSE

```

this slot is connected – twiddle the *nextSlot* in the data cell and send it

```
Port_SendDataCell(p, [dc EXCEPT !.nextSlot = oss.o_nextSlot])
```

15 Gare name

MODULE *DCNL1GareName*

EXTENDS

DCNL1Param,
DCNL1PortName

Index set and assignments of the port gares in an *L1* switch.

CONSTANT *MpxPortGare* how many ports are multiplexed to each port gare

Construct a set of partitions of the set of port names. The partitions are disjoint subsets that when combined include all port names. As many partitions as possible must be of size *MpxPortGare*.

PortNamePartition \triangleq

LET

RECURSIVE *Recur*(-, -, -)

Recur(*parts*, *mkpart*, *pns*) \triangleq

IF

\vee *Cardinality*(*mkpart*) = *MpxPortGare* subset is required size

\vee *mkpart* \neq {} \wedge *pns* = {} non-empty subset and done

THEN

Include {*mkpart*} as a subset in the partition.

Recur(*parts* \cup {*mkpart*}, {}, *pns*)

ELSE

IF *pns* = {}

THEN

Construction of the partition is complete.

parts

ELSE

Take an arbitrary port name from *pns* and include it in the subset being constructed.

LET

pn \triangleq CHOOSE *pn* \in *pns* : TRUE

IN

Recur(*parts*, *mkpart* \cup {*pn*}, *pns* \setminus {*pn*})

IN

$Recur(\{\}, \{\}, PortName)$

Gare names in an L1 switch.

$GareName \triangleq [$
 $part : PortNamePartition$
]

$FirstGareName \triangleq ArbSeq(GareName).FirstName$
 $IncrGareName \triangleq ArbSeq(GareName).IncrName$
 $DecrGareName \triangleq ArbSeq(GareName).DecrName$
 $GareIndexForName \triangleq ArbSeq(GareName).IndexForName$

Gare name assigned to a given port name.

$GareName_ForPortName \triangleq$
 $[pn \in PortName \mapsto \text{CHOOSE } gn \in GareName : pn \in gn.part]$

16 Switch name

MODULE $DCNL1SwitchName$

EXTENDS $DCNL1Param$

Switch names.

$SwitchName \triangleq [$
 $container : Container$
]

$SwitchName_ForContainer(c) \triangleq [$
 $container \mapsto c$
]

17 Gare

MODULE *DCNL1Gare*

EXTENDS

DCNL1Param,
DCNL1GareName,
DCNL1Port,
DCNL1RingMsg

Gare.

Gare \triangleq [
 garename : *GareName*,
 ring_out : *RingMsg*,

 port : [*PortName* \rightarrow *Nullable(Port)*],
 priopn : [*CtrlQueueType* \rightarrow *PortName*]
]

InitGare(*swn*, *gn*) \triangleq [
 garename \mapsto *gn*,
 ring_out \mapsto LET *x* \triangleq *GareIndexForName*[*gn*]%2IN
 CASE
 x = 0 \rightarrow *GreenNullRingMsg* \square
 x = 1 \rightarrow *RedNullRingMsg*,

 port \mapsto [*pn* \in *PortName* \mapsto
 IF *pn* \notin *gn*.*part* THEN *NULL* ELSE *InitPort*(*swn*, *pn*)
],
 priopn \mapsto [*t* \in *CtrlQueueType* \mapsto *FirstPortName*]
]

Set of port names assigned to gare *g*.

Gare_AssignedPortNames(*g*) \triangleq
g.garename.part

Update a gare to send a ring message.

$$\begin{aligned}
 \text{Gare_SendRingMsg}(g, rm) &\triangleq \\
 &\text{CASE TypeAssert}(g, \text{Gare}) \rightarrow \\
 &\text{CASE TypeAssert}(rm, \text{RingMsg}) \rightarrow \\
 &[g \text{ EXCEPT } !.ring_out = rm]
 \end{aligned}$$

Process a request control cell in port name pn at gare g . Compute the updated gare.

$$\begin{aligned}
 \text{Gare_Work_PortRequest}(g, pn) &\triangleq \\
 \text{LET} & \\
 cc &\triangleq \text{Head}(g.port[pn].ccinq[\text{CtrlQueueType_Req}]) \\
 p &\triangleq [g.port[pn] \text{ EXCEPT } !.ccinq[\text{CtrlQueueType_Req}] = \text{Tail}(@)]
 \end{aligned}$$

$$\text{Exec}(\text{downPort}, \text{DestC}, \text{DestH}, \text{prevSlot}) \triangleq$$

LET

set of unlocked subframes for port p

$$usfs \triangleq \{sf \in \text{Subframe} : \neg p.locked[sf]\}$$

number of free slots in subframe sf

$$\begin{aligned}
 cntf(sf) &\triangleq \text{Cardinality}(\{sn \in \text{SlotName} : \\
 &\quad \wedge sn.subframe = sf \\
 &\quad \wedge \neg p.iss[sn].i_assign\})
 \end{aligned}$$

unlocked subframe with most number of free slots

$$\begin{aligned}
 subframe &\triangleq \text{CHOOSE } sf1 \in usfs : \\
 &\quad \forall sf2 \in usfs : \\
 &\quad \quad cntf(sf1) \geq cntf(sf2)
 \end{aligned}$$

map of which slottinas are currently occupied in that subframe

$$\begin{aligned}
 occupied &\triangleq [sti \in \text{Slottina} \mapsto \\
 &\quad \text{LET} \\
 &\quad \quad sn \triangleq [subframe \mapsto subframe, slottina \mapsto sti] \\
 &\quad \text{IN} \\
 &\quad \quad p.iss[sn].i_assign \\
 &\quad]
 \end{aligned}$$

launch the setup request ring message

$$\begin{aligned}
 \text{SendG} &\triangleq [\text{Gare_SendRingMsg}(g, [\\
 &\quad \text{typecode} \mapsto \text{RingMsg_Typecode_RedSetupReq}, \\
 &\quad \text{connect} \mapsto cc.connect, \\
 &\quad \text{DestC} \mapsto \text{DestC},
 \end{aligned}$$

```

    DestH    ↦ DestH,
    prevSlot ↦ prevSlot,
    upPort   ↦ pn,
    downPort ↦ downPort,
    subframe ↦ subframe,
    occupied ↦ occupied
  ]) EXCEPT

  !.priopn[CtrlQueueType_Req] = IncrPortName[pn], rotate the priority

  !.port[pn] = [p EXCEPT
    !.locked[subframe] = TRUE, lock the subframe
    !.locked_connect[subframe] = cc.connect for what connection
  ]
]

the request cannot be serviced now
WaitG ≜ g

fail this connection attempt
FailG ≜ [g EXCEPT

  !.priopn[CtrlQueueType_Req] = IncrPortName[pn], rotate the priority

  !.port[pn] = Port_SendCtrlCell(p, [
    typecode ↦ CtrlCell_Typecode_Teardown,
    connect  ↦ cc.connect,
    prevSlot ↦ prevSlot,
    reason   ↦ "full"
  ])
]

IN
IF ∀ sf ∈ Subframe : cntf(sf) = 0 THEN FailG ELSE
IF usfs ≠ {} ∧ cntf(subframe) > 0 THEN SendG ELSE
WaitG

```

```

downPortToC(DestC) ≜ [
  typecode ↦ PortName_Typecode_Container,
  to       ↦ DestC
]

```

```

downPortToH(DestH) ≜ [
  typecode ↦ PortName_Typecode_Halfrow,
  to       ↦ DestH
]

```

```

IN
IF  $cc.DestC \neq Port\_MyContainer(p)$ 
THEN
  make a setup request to connect to container  $DestC$ 
  Exec( $downPortToC(cc.DestC)$ ,  $cc.DestC$ ,  $cc.DestH$ ,  $cc.prevSlot$ )
ELSE
  make a setup request to connect to halfrow  $DestH$ 
  Exec( $downPortToH(cc.DestH)$ ,  $cc.DestC$ ,  $cc.DestH$ ,  $cc.prevSlot$ )

```

Process an upstream control cell in port name pn at gare g . Compute the updated gare.

```

Gare_Work_PortUpstream( $g$ ,  $pn$ )  $\triangleq$ 
LET
   $cc$   $\triangleq$   $Head(g.port[pn].ccinq[CtrlQueueType-Ups])$ 
   $p$   $\triangleq$  [ $g.port[pn]$  EXCEPT  $!.ccinq[CtrlQueueType-Ups] = Tail(@)$ ]

   $sn$   $\triangleq$   $cc.prevSlot$ 
   $oss$   $\triangleq$   $p.oss[sn]$ 

  UsePermitG  $\triangleq$  [ $g$  EXCEPT
     $!.priopn[CtrlQueueType-Ups] = IncrPortName[pn]$ ,
     $!.port[pn] = p$ 
  ]

  UpstreamG( $rm$ )  $\triangleq$  Gare_SendRingMsg(UsePermitG,  $rm$ )

  matchconn  $\triangleq$  Assert( $cc.connect = oss.o\_connect$ ,
    ("connection mismatch",  $cc$ ,  $pn$ ,  $oss$ ))

  TeardownG  $\triangleq$ 
  CASE matchconn  $\rightarrow$ 
  IF  $\neg oss.o\_assign$  THEN UsePermitG ELSE
  [UpstreamG([
    typecode  $\mapsto$  RingMsg_TypeCode_GreenTeardown,
    connect  $\mapsto$   $oss.o\_connect$ ,
    localSlot  $\mapsto$   $sn$ ,
    prevSlot  $\mapsto$   $oss.o\_prevSlot$ ,
    upPort  $\mapsto$   $oss.o\_upPort$ ,
    reason  $\mapsto$   $cc.reason$ 
  ]) EXCEPT
   $!.port[pn].oss[sn] = InitPortOutSlot$ 
  ]

```

```

DiscardG  $\triangleq$ 
CASE matchconn  $\rightarrow$ 
IF  $\neg$ oss.o_assign THEN UsePermitG ELSE
UpstreamG([
  typecode  $\mapsto$  RingMsg_Typecode_GreenDiscard,
  connect  $\mapsto$  oss.o_connect,
  prevSlot  $\mapsto$  oss.o_prevSlot,
  upPort  $\mapsto$  oss.o_upPort
])

DataAckG  $\triangleq$ 
CASE matchconn  $\rightarrow$ 
IF  $\neg$ oss.o_assign THEN UsePermitG ELSE
UpstreamG([
  typecode  $\mapsto$  RingMsg_Typecode_GreenDataAck,
  connect  $\mapsto$  oss.o_connect,
  prevSlot  $\mapsto$  oss.o_prevSlot,
  upPort  $\mapsto$  oss.o_upPort,
  chunkseqn  $\mapsto$  cc.chunkseqn,
  throttle  $\mapsto$  cc.throttle
])

RestartG  $\triangleq$ 
CASE matchconn  $\rightarrow$ 
IF  $\neg$ oss.o_assign THEN UsePermitG ELSE
UpstreamG([
  typecode  $\mapsto$  RingMsg_Typecode_GreenRestart,
  connect  $\mapsto$  oss.o_connect,
  prevSlot  $\mapsto$  oss.o_prevSlot,
  upPort  $\mapsto$  oss.o_upPort
])

StopG  $\triangleq$ 
CASE matchconn  $\rightarrow$ 
IF  $\neg$ oss.o_assign THEN UsePermitG ELSE
UpstreamG([
  typecode  $\mapsto$  RingMsg_Typecode_GreenStop,
  connect  $\mapsto$  oss.o_connect,
  prevSlot  $\mapsto$  oss.o_prevSlot,
  upPort  $\mapsto$  oss.o_upPort
])

IN
CASE
CtrlCell_Is_Teardown(cc)  $\rightarrow$  TeardownG  $\square$ 
CtrlCell_Is_Discard(cc)  $\rightarrow$  DiscardG  $\square$ 
CtrlCell_Is_DataAck(cc)  $\rightarrow$  DataAckG  $\square$ 
CtrlCell_Is_Restart(cc)  $\rightarrow$  RestartG  $\square$ 

```

$CtrlCell_Is_Stop(cc) \rightarrow StopG$

Process a control cell from the indicated type queue of one of our assigned input ports, and generate a ring message. Compute the updated gare.

$Gare_Work(g, t) \triangleq$

LET

Given a set of port names, find the one that is least higher than $priopn$. Compute a record with fields $ok =$ there is a highest priority port name, $pn =$ that port name.

$Prio(pns, priopn) \triangleq$

LET

$elevpns \triangleq \{pn \in pns : LeqPortName[priopn, pn]\}$

$priopns \triangleq \text{IF } elevpns \neq \{\} \text{ THEN } elevpns \text{ ELSE } pns$

$pn \triangleq \text{CHOOSE } pn1 \in priopns : \forall pn2 \in priopns : LeqPortName[pn1, pn2]$

$ok \triangleq priopns \neq \{\}$

IN

IF $\neg ok$ THEN $[ok \mapsto ok]$ ELSE

[

$ok \mapsto ok,$

$pn \mapsto pn$

]

$pns \triangleq Gare_AssignedPortNames(g)$

$q_pns \triangleq \{pn \in pns : g.port[pn].ccinq[t] \neq \langle \rangle\}$

$q_prio \triangleq Prio(q_pns, g.priopn[t])$

IN

IF $q_prio.ok$ THEN

CASE

$t = CtrlQueueType_Req \rightarrow Gare_Work_PortRequest(g, q_prio.pn) \square$

$t = CtrlQueueType_Ups \rightarrow Gare_Work_PortUpstream(g, q_prio.pn)$

ELSE

g

Compute the updated state of a gare after it processes a setup request ring message.

$Gare_ProcessSetupReq(g, rm) \triangleq$

CASE $TypeAssert(g, Gare) \rightarrow$

CASE $TypeAssert(rm, RingMsg_RedSetupReq) \rightarrow$

LET

```

occupied  $\triangleq$  LET dp  $\triangleq$  g.port[rm.downPort]IN
  [sti  $\in$  Slottina  $\mapsto$ 
     $\vee$  rm.occupied[sti]
     $\vee$  LET sn  $\triangleq$  [slottina  $\mapsto$  sti, subframe  $\mapsto$  rm.subframe]IN
      dp.oss[sn].o_assign
  ]

eval  $\triangleq$  RingMsg-Eval-Occupied(rm.subframe, occupied)

OutG  $\triangleq$  Gare-SendRingMsg(g, [
  typecode  $\mapsto$  RingMsg-Typecode-RedSetupAns,
  connect  $\mapsto$  rm.connect,
  DestC  $\mapsto$  rm.DestC,
  DestH  $\mapsto$  rm.DestH,
  prevSlot  $\mapsto$  rm.prevSlot,
  upPort  $\mapsto$  rm.upPort,
  downPort  $\mapsto$  rm.downPort,
  subframe  $\mapsto$  rm.subframe,
  occupied  $\mapsto$  occupied
])

IN
IF  $\neg$ eval.success THEN OutG ELSE
[OutG EXCEPT !.port[rm.downPort] =
[Port-SendCtrlCell(@, [
  typecode  $\mapsto$  CtrlCell-Typecode-SetupReq,
  connect  $\mapsto$  rm.connect,
  DestC  $\mapsto$  rm.DestC,
  DestH  $\mapsto$  rm.DestH,
  prevSlot  $\mapsto$  eval.localSlot
])]
EXCEPT !.oss[eval.localSlot] =
[
  o_assign  $\mapsto$  TRUE,
  o_connect  $\mapsto$  rm.connect,
  o_nextset  $\mapsto$  FALSE,
  o_prevSlot  $\mapsto$  rm.prevSlot,
  o_nextSlot  $\mapsto$  NULL,
  o_upPort  $\mapsto$  rm.upPort
]
]
]
]

```

Compute the updated state of a gare after it processes a setup answer ring message.

```
Gare_ProcessSetupAns(g, rm)  $\triangleq$ 
CASE TypeAssert(g, Gare)  $\rightarrow$ 
CASE TypeAssert(rm, RingMsg_RedSetupAns)  $\rightarrow$ 
```

```
LET
  eval  $\triangleq$  RingMsg_Eval_Occupied(rm.subframe, rm.occupied)
```

Send output ring message and unlock the up port.

```
UnlkG  $\triangleq$  [Gare_SendRingMsg(g, [
  typecode  $\mapsto$  RingMsg_Typecode_RedSetupXbar,
  connect  $\mapsto$  rm.connect,
  DestC  $\mapsto$  rm.DestC,
  DestH  $\mapsto$  rm.DestH,
  prevSlot  $\mapsto$  rm.prevSlot,
  upPort  $\mapsto$  rm.upPort,
  downPort  $\mapsto$  rm.downPort,
  subframe  $\mapsto$  rm.subframe,
  occupied  $\mapsto$  rm.occupied
])] EXCEPT

!.port[rm.upPort] = [ @ EXCEPT
!.locked[rm.subframe] =
  CASE Assert(@, "already unlocked subframe on setup ans")  $\rightarrow$ 
  FALSE,
!.locked_connect[rm.subframe] =
  CASE Assert(@ = rm.connect, "connection mismatch")  $\rightarrow$ 
  NULL
]
]

FailG  $\triangleq$  [UnlkG EXCEPT
!.port[rm.upPort] = Port_SendCtrlCell(@, [
  typecode  $\mapsto$  CtrlCell_Typecode_Teardown,
  connect  $\mapsto$  rm.connect,
  prevSlot  $\mapsto$  rm.prevSlot,
  reason  $\mapsto$  "setup fail"
])]
]

AssnG  $\triangleq$  [UnlkG EXCEPT
!.port[rm.upPort].iss[eval.localSlot] = [ @ EXCEPT
!.i_assign =
  CASE Assert( $\neg$ @, "already assigned slot on setup ans")  $\rightarrow$ 
  TRUE,
!.i_connect = rm.connect
```

```

]
]
PassG  $\triangleq$  [AssnG EXCEPT
  !.port[rm.upPort] = Port_SendCtrlCell(@, [
    typecode  $\mapsto$  CtrlCell_TypeCode_SetupAck,
    connect  $\mapsto$  rm.connect,
    prevSlot  $\mapsto$  rm.prevSlot,
    nextSlot  $\mapsto$  eval.localSlot
  ])
]
IN
IF  $\neg$ eval.success THEN FailG ELSE PassG

```

Receive ring message rm at gare g . Compute the new state of the gare.

```

Gare_RecvRingMsg(g, rm)  $\triangleq$ 
CASE TypeAssert(g, Gare)  $\rightarrow$ 
CASE TypeAssert(rm, RingMsg)  $\rightarrow$ 

LET
  MyU  $\triangleq$  rm.upPort  $\in$  Gare_AssignedPortNames(g)
  MyD  $\triangleq$  rm.downPort  $\in$  Gare_AssignedPortNames(g)

  OutCopy  $\triangleq$  Gare_SendRingMsg(g, rm)
  OutRedNull  $\triangleq$  Gare_SendRingMsg(g, RedNullRingMsg)
  OutGreenNull  $\triangleq$  Gare_SendRingMsg(g, GreenNullRingMsg)

  Upstream(cc)  $\triangleq$  [OutGreenNull EXCEPT
    !.port[rm.upPort] = Port_SendCtrlCell(@, cc)
  ]

  UpInit(g1)  $\triangleq$  [g1 EXCEPT !.port[rm.upPort] =
    CASE Assert(@.iss[rm.localSlot].i_assign, "protocol error")  $\rightarrow$ 
    [@ EXCEPT !.iss[rm.localSlot] = InitPortInSlot]
  ]

  matchconn  $\triangleq$  LET
    iss  $\triangleq$  g.port[rm.upPort].iss
  IN
    Assert(rm.connect = iss[rm.localSlot].i_connect,
      ("up port connection mismatch", rm, iss))

  DoTeardown  $\triangleq$ 
CASE matchconn  $\rightarrow$ 
UpInit(Upstream([

```

```

        typecode ↦ CtrlCell_Typecode_TearDown,
        connect ↦ rm.connect,
        prevSlot ↦ rm.prevSlot,
        reason ↦ rm.reason
    )))
DoGCClear ≜
    CASE matchconn →
        UpInit(OutGreen.Null)
DoDiscard ≜
    cannot matchconn as no localSlot in a discard ring message
    CASE matchconn →
        Upstream([
            typecode ↦ CtrlCell_Typecode_Discard,
            connect ↦ rm.connect,
            prevSlot ↦ rm.prevSlot
        ])
DoDataAck ≜
    cannot matchconn as no localSlot in a data ack ring message
    CASE matchconn →
        Upstream([
            typecode ↦ CtrlCell_Typecode_DataAck,
            connect ↦ rm.connect,
            prevSlot ↦ rm.prevSlot,
            chunkseqn ↦ rm.chunkseqn,
            throttle ↦ rm.throttle
        ])
DoRestart ≜
    cannot matchconn as no localSlot in a restart ring message
    CASE matchconn →
        Upstream([
            typecode ↦ CtrlCell_Typecode_Restart,
            connect ↦ rm.connect,
            prevSlot ↦ rm.prevSlot
        ])
DoStop ≜
    cannot matchconn as no localSlot in a stop ring message
    CASE matchconn →
        Upstream([
            typecode ↦ CtrlCell_Typecode_Stop,
            connect ↦ rm.connect,
            prevSlot ↦ rm.prevSlot
        ])

```

```

IN
CASE
  RingMsg_Is_RedSetupReq(rm)     $\wedge$  MyD  $\rightarrow$  Gare_ProcessSetupReq(g, rm)            $\square$ 
  RingMsg_Is_RedSetupAns(rm)     $\wedge$  MyU  $\rightarrow$  Gare_ProcessSetupAns(g, rm)            $\square$ 
  RingMsg_Is_RedSetupXbar(rm)    $\wedge$  MyU  $\rightarrow$  OutRedNull                            $\square$ 
  RingMsg_Is_GreenTeardown(rm)  $\wedge$  MyU  $\rightarrow$  DoTeardown                            $\square$ 
  RingMsg_Is_GreenGCClear(rm)   $\wedge$  MyU  $\rightarrow$  DoGCClear                            $\square$ 
  RingMsg_Is_GreenDiscard(rm)   $\wedge$  MyU  $\rightarrow$  DoDiscard                            $\square$ 
  RingMsg_Is_GreenDataAck(rm)   $\wedge$  MyU  $\rightarrow$  DoDataAck                            $\square$ 
  RingMsg_Is_GreenRestart(rm)   $\wedge$  MyU  $\rightarrow$  DoRestart                            $\square$ 
  RingMsg_Is_GreenStop(rm)      $\wedge$  MyU  $\rightarrow$  DoStop                                $\square$ 
  RingMsg_Is_RedNull(rm)         $\rightarrow$  Gare_Work(OutCopy, CtrlQueueType_Req)  $\square$ 
  RingMsg_Is_GreenNull(rm)       $\rightarrow$  Gare_Work(OutCopy, CtrlQueueType_Ups)  $\square$ 
OTHER
   $\rightarrow$  OutCopy

```

18 Switch

MODULE *DCNL1Switch*

EXTENDS

```

  DCNL1Param,
  DCNL1SlotName,
  DCNL1CtrlCell,
  DCNL1DataCell,
  DCNL1RingMsg,
  DCNL1SwitchName,
  DCNL1Port,
  DCNL1Gare

```

BugNoTimingAssumption \triangleq FALSE

State of a *L1* switch.

```

SwitchXbarInfo  $\triangleq$  [
  upPort : PortName,
  connect : Nullable(Connection)
]

```

```

InitSwitchXbarInfo  $\triangleq$  [

```

```

    upPort ↦ FirstPortName,
    connect ↦ NULL
]

```

```

Switch ≜ [
    switchname : SwitchName,
    gare       : [GareName → Gare],
    xbarinfo   : [SlotName → [PortName → SwitchXbarInfo]],
    slot       : SlotName
]

```

```

InitSwitch(swn) ≜ [
    switchname ↦ swn,
    gare       ↦ [gn ∈ GareName ↦ InitGare(swn, gn)],
    xbarinfo   ↦ [sn ∈ SlotName ↦ [pn ∈ PortName ↦ InitSwitchXbarInfo]],
    slot       ↦ FirstSlotName
]

```

Get a port of a switch.

```

Switch_Port(sw, pn) ≜
    LET
        gn ≜ GareName_ForPortName[pn]
    IN
        sw.gare[gn].port[pn]

```

Update a port of a switch.

```

Switch_Port_Update(sw, p) ≜
    LET
        pn ≜ p.portname
        gn ≜ GareName_ForPortName[pn]
    IN
        [sw EXCEPT !.gare[gn].port[pn] = p]

```

In switch *sw*, on port name *pn*, determine if the switch is ready to receive data cell *dc*.

$Switch_Port_RecvCell_Ready(sw, pn, t, cell) \triangleq$

CASE

$t = CellType_Data \rightarrow \text{TRUE}$

□

$t = CellType_Ctrl \rightarrow$

IF *BugNoTimingAssumption* THEN TRUE ELSE

In switch *sw*, on port name *pn*, determine if the switch is ready to receive a given control cell. There is a timing requirement that the crossbar must be set up for a connection before a setup ack control cell for that connection arrives at the down port input unit.

LET

$cc \triangleq cell$

$match(sn) \triangleq sw.xbarinfo[sn][pn].connect = cc.connect$

IN

CASE

$CtrlCell_Is_SetupReq(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_SetupAck(cc) \rightarrow match(cc.prevSlot) \square$

$CtrlCell_Is_Teardown(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_Discard(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_DataAck(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_Restart(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_Stop(cc) \rightarrow \text{TRUE} \quad \square$

$CtrlCell_Is_Null(cc) \rightarrow \text{TRUE}$

In switch *sw*, on port name *pn*, receive a given cell.

$Switch_Port_RecvCell(sw, pn, t, cell) \triangleq$

LET

$p \triangleq Switch_Port(sw, pn)$

$recv \triangleq Port_RecvCell(p, t, cell)$

$ready \triangleq Switch_Port_RecvCell_Ready(sw, pn, t, cell) \wedge recv.ready$

$update \triangleq Switch_Port_Update(sw, recv.update)$

IN

IF $\neg ready$ THEN $[ready \mapsto ready]$ ELSE

[

$ready \mapsto ready,$

$update \mapsto update$

]

In switch *sw*, on port name *pn*, extract a cell.

$Switch_Port_TakeCell(sw, pn, t) \triangleq$

LET

```

p       $\triangleq$  Switch_Port(sw, pn)
take    $\triangleq$  Port_TakeCell(p, t)

ready   $\triangleq$  take.ready
cell    $\triangleq$  take.cell
update  $\triangleq$  Switch_Port_Update(sw, take.update)
IN
IF  $\neg$ ready THEN [ready  $\mapsto$  ready] ELSE
[
  ready  $\mapsto$  ready,
  cell    $\mapsto$  cell,
  update  $\mapsto$  update
]

```

Get a sequence of port outputs (debugging).

```

Switch_Debug_LinkOut(sw)  $\triangleq$ 
[pn  $\in$  PortName  $\mapsto$  Switch_Port(sw, pn).dcoutq]

```

Advance a crossbar slot time. Compute the updated state of the switch.

```

Switch_AdvanceSlot(sw)  $\triangleq$ 
LET
  pmap0    $\triangleq$  [pn  $\in$  PortName  $\mapsto$  Switch_Port(sw, pn)]
  sn       $\triangleq$  sw.slot
  xbar(pn)  $\triangleq$  sw.xbarinfo[sn][pn].upPort

  from(p)  $\triangleq$  pmap0[xbar(p.portname)]   port sending to p

  ip(p)    $\triangleq$  Port_Xbar_TakeSlot(p, sn).p   updated input port p
  dc(p)    $\triangleq$  Port_Xbar_TakeSlot(p, sn).dc   data cell sent to xbar by p
  xb(p)    $\triangleq$  Port_Xbar_RecvSlot(ip(p), sn, dc(from(p)))   full update of p

  pmap1    $\triangleq$  [pn  $\in$  PortName  $\mapsto$  xb(pmap0[pn])]
IN
[sw EXCEPT
  !.gare = [gn  $\in$  GareName  $\mapsto$ 
    [@[gn] EXCEPT
      !.port = [pn  $\in$  PortName  $\mapsto$ 
        IF pn  $\notin$  gn.part THEN NULL ELSE pmap1[pn]
      ]
    ]
  ]
]

```

```

    ],
    !.slot = IncrSlotName[@]
  ]

```

Compute the updated state of the switch after pushing the out ring message from each *gare* into the next *gare* around the ring.

$Switch_AdvanceRing(sw) \triangleq$

```

LET
  ring message entering gare gn
  rmin(gn)  $\triangleq$  sw.gare[DecrGareName[gn]].ring_out
IN
[sw EXCEPT
  Update each gare according to the ring message entering it.
  !.gare = [gn  $\in$  GareName  $\mapsto$  Gare_RecvRingMsg(@[gn], rmin(gn))],

```

Update *xbar* info according to any setup *xbar* ring message entering the first *gare*. In the actual hardware, each *gare* handles a bit slice of the crossbar, and the crossbar setup happens for each bit slice as the setup *xbar* ring message enters the respective *gare*. We simplify this by doing all the crossbar setup at the first *gare*.

```

!.xbarinfo =
  LET
    rm  $\triangleq$  rmin(FirstGareName)
    eval  $\triangleq$  RingMsg_Eval_Occupied(rm.subframe, rm.occupied)
  IN
  IF  $\neg$ RingMsg_Is_RedSetupXbar(rm) THEN @ ELSE
  IF  $\neg$ eval.success THEN @ ELSE
  [@ EXCEPT ![eval.localSlot][rm.downPort] =
  [@ EXCEPT
    !.upPort = rm.upPort,
    !.connect = rm.connect
  ]
  ]
]

```

19 Full L1 system

MODULE DCNL1

EXTENDS

DCNL1Param,
DCNL1SlotName,
DCNL1CtrlCell,
DCNL1DataCell,
DCNL1RingMsg,
DCNL1SwitchName,
DCNL1Port,
DCNL1Gare,
DCNL1Switch,
DCNL1Node

Specification of the L1 switch infrastructure in Chuck *Thacker*'s Data Center Network version 4.4.

CONSTANT *MaxPacketCells*
 CONSTANT *MinPacketCells*

State of the system.

nodeName \triangleq [
 swName : *SwitchName*,
 hr : *Halfrow*
]

Request \triangleq [
 connect : *Connection*,
 srcx : *nodeName*,
 dstx : *nodeName*,
 pkt : *Seq(Payload)*
]

State \triangleq [
 node : [*nodeName* \rightarrow *Node*], nodes in the system
 switch : [*SwitchName* \rightarrow *Switch*], switches in the system
 request_cnt : *Nat*, number of requests launched
 request_set : SUBSET *Request*, set of requests lauched
 debug : *Any*, debugging info
]

InitState \triangleq [
 node \mapsto [*nn* \in *nodeName* \mapsto *InitNode*],
]

```

switch      ↦ [swn ∈ SwitchName ↦ InitSwitch(swn)],
request_cnt ↦ 0,
request_set ↦ {},
debug      ↦ NULL
]

```

Link definitions

A link end is either a node or a port.

$LinkEnd_Typecode_Node \triangleq \text{"LinkEnd_Node"}$

$LinkEnd_Typecode_Port \triangleq \text{"LinkEnd_Port"}$

$LinkEnd_Node \triangleq [$
 $typecode : \{LinkEnd_Typecode_Node\},$
 $nn : NodeName$
 $]$

$LinkEnd_Port \triangleq [$
 $typecode : \{LinkEnd_Typecode_Port\},$
 $swn : SwitchName,$
 $pn : PortName$
 $]$

$LinkEnd \triangleq \{ \}$
 $\cup LinkEnd_Node$
 $\cup LinkEnd_Port$

The set of all port-to-port links.

$Link_PortPort \triangleq$
 LET
 $slkes \triangleq$
 $\{lke \in LinkEnd_Port :$
 $PortName_Is_Container(lke.pn)$
 $\}$
 $dlke(slke) \triangleq$
 CHOOSE $lke \in LinkEnd_Port :$
 LET
 $spn \triangleq slke.pn$
 $dpn \triangleq lke.pn$
 IN
 $\wedge PortName_Is_Container(dpn)$

$$\begin{aligned} & \wedge dpn.to = slke.swn.container \\ & \wedge lke.swn.container = spn.to \\ \text{IN} \\ & \{[srcx \mapsto slke, dstx \mapsto dlke(slke)] : slke \in slkes\} \end{aligned}$$

The set of all port-to-node and node-to-port links.

$$\begin{aligned} Link_PortNode & \triangleq \\ \text{LET} \\ & slkes \triangleq \\ & \quad \{lke \in LinkEnd_Port : \\ & \quad \quad PortName_Is_Halfrow(lke.pn) \\ & \quad \} \\ & dlke(slke) \triangleq \\ & \quad \text{CHOOSE } lke \in LinkEnd_Node : \\ & \quad \text{LET} \\ & \quad \quad spn \triangleq slke.pn \\ & \quad \quad dnn \triangleq lke.nn \\ & \quad \text{IN} \\ & \quad \quad \wedge dnn.swn = slke.swn \\ & \quad \quad \wedge dnn.hr = spn.to \\ \text{IN} \\ & \{ \} \\ & \cup \{[srcx \mapsto slke, dstx \mapsto dlke(slke)] : slke \in slkes\} \\ & \cup \{[dstx \mapsto slke, srcx \mapsto dlke(slke)] : slke \in slkes\} \end{aligned}$$

The set of all links.

$$\begin{aligned} Link & \triangleq \{ \} \\ & \cup Link_PortPort \\ & \cup Link_PortNode \end{aligned}$$

In state s , on link lk , receive a cell of the indicated type.

The receiver might not be ready for this particular cell yet. The hardware avoids this situation by using timing assumptions that are carefully checked in the hardware design. In the specification, we model it explicitly.

$$\begin{aligned} State_Link_RecvCell(s, lk, t, cell) & \triangleq \\ \text{LET} \\ & DoNode \triangleq \\ & \quad \text{LET} \\ & \quad \quad nn \triangleq lk.dstx.nn \\ & \quad \quad recv \triangleq Node_RecvCell(s.node[nn], t, cell) \end{aligned}$$

```

    ready  $\triangleq$  recv.ready
    update  $\triangleq$  [s EXCEPT !.node[nn] = recv.update]
  IN
  IF  $\neg$ ready THEN [ready  $\mapsto$  ready] ELSE
  [
    ready  $\mapsto$  ready,
    update  $\mapsto$  update
  ]
DoPort  $\triangleq$ 
  LET
    swn  $\triangleq$  lk.dstx.swn
    pn  $\triangleq$  lk.dstx.pn
    recv  $\triangleq$  Switch_Port_RecvCell(s.switch[swn], pn, t, cell)

    ready  $\triangleq$  recv.ready
    update  $\triangleq$  [s EXCEPT !.switch[swn] = recv.update]
  IN
  IF  $\neg$ ready THEN [ready  $\mapsto$  ready] ELSE
  [
    ready  $\mapsto$  ready,
    update  $\mapsto$  update
  ]
IN
CASE
lk.dstx.typecode = LinkEnd_Typecode_Node  $\rightarrow$  DoNode  $\square$ 
lk.dstx.typecode = LinkEnd_Typecode_Port  $\rightarrow$  DoPort

```

In state *s*, on link *lk*, extract a cell of the indicated type from the sender.

```

State_Link_TakeCell(s, lk, t)  $\triangleq$ 
  LET
    DoNode  $\triangleq$ 
      LET
        nn  $\triangleq$  lk.srcx.nn
        take  $\triangleq$  Node_TakeCell(s.node[nn], t)

        ready  $\triangleq$  take.ready
        cell  $\triangleq$  take.cell
        update  $\triangleq$  [s EXCEPT !.node[nn] = take.update]
      IN
      IF  $\neg$ ready THEN [ready  $\mapsto$  ready] ELSE
      [
        ready  $\mapsto$  ready,
        cell  $\mapsto$  cell,

```

```

    update ↦ update
  ]
DoPort ≜
LET
  swn ≜ lk.srcx.swn
  pn ≜ lk.srcx.pn
  take ≜ Switch_Port_TakeCell(s.switch[swn], pn, t)

  ready ≜ take.ready
  cell ≜ take.cell
  update ≜ [s EXCEPT !.switch[swn] = take.update]
IN
IF ¬ready THEN [ready ↦ ready] ELSE
[
  ready ↦ ready,
  cell ↦ cell,
  update ↦ update
]
IN
CASE
lk.srcx.typecode = LinkEnd_Typecode_Node → DoNode□
lk.srcx.typecode = LinkEnd_Typecode_Port → DoPort

```

VARIABLE *state*

NEXT STATE RELATIONS

Advance the slot time everywhere.

```

NextSlot ≜
state' = [state EXCEPT
  !.node = [nn ∈ NodeName ↦ Node_AdvanceSlot(@[nn])],
  !.switch = [swn ∈ SwitchName ↦ Switch_AdvanceSlot(@[swn])],
  !.debug = ⟨"adv slot"⟩
]

```

Advance the ring in switch *swn*.

```

NextRingSwitch(swn) ≜
LET
  sw0 ≜ state.switch[swn]

```

```

    sw1  $\triangleq$  Switch_AdvanceRing(sw0)
  IN
   $\wedge$  state' = [state EXCEPT
                !.switch[swn] = sw1,
                !.debug = ⟨"sw adv ring", swn⟩
              ]

```

Move a cell of the indicated type across link *lk*.

```

NextLinkCell(lk, t)  $\triangleq$ 
  LET
    take  $\triangleq$  State_Link_TakeCell(state, lk, t)
    recv  $\triangleq$  State_Link_RecvCell(take.update, lk, t, take.cell)
  IN
   $\wedge$  take.ready
   $\wedge$  recv.ready
   $\wedge$  state' = [recv.update EXCEPT
                !.debug = ⟨"link", t, lk.srcx, lk.dstx, take.cell⟩
              ]

```

Create a request to send a packet. With back-to-back nodes, the destination is always the other node than the source. However, we still have to pick *destV*, *destC*, and *destH* for the setup request control cell.

```

NextCreateRequest  $\triangleq$ 
   $\exists$  srcx  $\in$  NodeName :
  LET
    cnt  $\triangleq$  state.request_cnt + 1
    srcn  $\triangleq$  state.node[srcx]
  IN
   $\wedge$  cnt  $\leq$  NumConnection
   $\wedge$ 
     $\exists$  s  $\in$  Node_AvailOutSlots(srcn) :
     $\exists$  C  $\in$  Container :
     $\exists$  H  $\in$  Halfrow :
     $\exists$  len  $\in$  MinPacketCells .. MaxPacketCells :
     $\exists$  pkt  $\in$  [1 .. len  $\rightarrow$  Payload] :
  LET
    dstx  $\triangleq$  [swn  $\mapsto$  SwitchName_ForContainer(C), hr  $\mapsto$  H]
  IN
  state' =
    [state EXCEPT
     !.node[srcx] = Node_CreateRequest(@, s, C, H, pkt, cnt),
     !.request_cnt = cnt,

```

```

!.request_set = @ ∪ {[
    connect ↦ cnt,
    srcx    ↦ srcx,
    dstx    ↦ dstx,
    pkt     ↦ pkt
  ]},
!.debug = ⟨“create req”, srcx, dstx, pkt⟩
]

```

INVARIANTS

The state must always be of the proper type.

$$InvType \triangleq \wedge state \in State$$

Any packets received must have been sent.

$$InvRecvPacket \triangleq \forall dstx \in NodeName : \text{LET } \begin{array}{l} dstn \triangleq state.node[dstx] \\ pkts \triangleq Node_ReceivedPackets(dstn) \end{array} \text{IN } \begin{array}{l} \forall pkt \in pkts : \\ \exists r \in state.request_set : \\ \wedge r.dstx = dstx \\ \wedge r.pkt = pkt \end{array}$$

TEMPORAL PROPERTIES

All node output slots always eventually available.

$$AllNodeOutSlotEvAvail \triangleq \begin{array}{l} \forall nn \in NodeName : \\ \forall sn \in SlotName : \end{array}$$

$\Box \Diamond \text{Node_IsAvailOutSlot}(\text{state.node}[nn], sn)$

All node input slots always eventually available.

$\text{AllNodeInSlotEvAvail} \triangleq$
 $\forall nn \in \text{NodeName} :$
 $\forall sn \in \text{SlotName} :$
 $\Box \Diamond \text{Node_IsAvailInSlot}(\text{state.node}[nn], sn)$

All links are eventually quiet.

$\text{AllLinkEvQuiet} \triangleq$
 $\forall lk \in \text{Link} :$
 $\forall t \in \text{CellType} :$
 $\Box \Diamond \neg \text{State_Link_TakeCell}(\text{state}, lk, t).ready$

Initial state.

$\text{Init} \triangleq$
 $\wedge \text{FirstGareName} \neq \text{NULL}$
 $\wedge \text{state} = \text{InitState}$

Next state relation.

$\text{Next} \triangleq$
 $\vee \exists swn \in \text{SwitchName} : \text{NextRingSwitch}(swn)$

 $\vee \exists lk \in \text{Link} :$
 $\quad \exists t \in \text{CellType} :$
 $\quad \text{NextLinkCell}(lk, t)$

 $\vee \text{NextSlot}$

 $\vee \text{NextCreateRequest}$

Liveness.

$\text{Liveness} \triangleq$
 $\wedge \forall swn \in \text{SwitchName} :$
 $\quad \text{WF}_{\text{state}}(\text{NextRingSwitch}(swn))$

 $\wedge \forall lk \in \text{Link} :$
 $\quad \forall t \in \text{CellType} :$
 $\quad \text{WF}_{\text{state}}(\text{NextLinkCell}(lk, t))$

$\wedge \text{WF}_{state}(NextSlot)$

System specification.

$Spec \triangleq$
 $\wedge Init$
 $\wedge \square[Next]_{state}$
 $\wedge Liveness$

View removing debugging info.

$View \triangleq [state \text{ EXCEPT } !.debug = NULL]$

References

- [1] L. Lamport. *Specifying Systems: The TLA+ Language and Tools for Hardware and Software Engineers*. Addison-Wesley, 2002.
- [2] C. Thacker, A. Nowatzky, T. Rodeheffer, and F. Yu. A data center network using FPGAs (v4.5), April 2011. Unpublished.