

Further Details on Cycle Identification

dv-yizhou-cycle-identification-details-1122-v02

Yizhou Li (Huawei)

Guanhua Zhuang (Huawei)

Shoushou Ren (Huawei)

Jiang Li (Huawei)

Jeong-dong Ryoo (ETRI)

Peng Liu (CMCC)

Dong Li (Shenyang Institute of Automation)

Wenbin Dai (Shanghai Jiao Tong University)

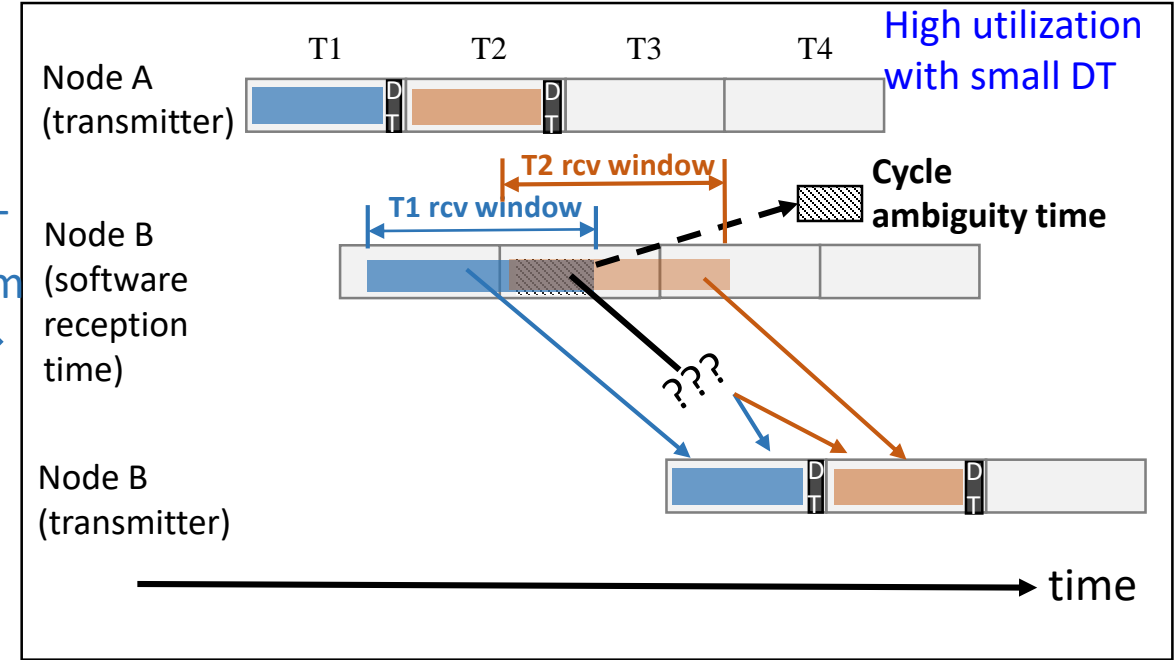
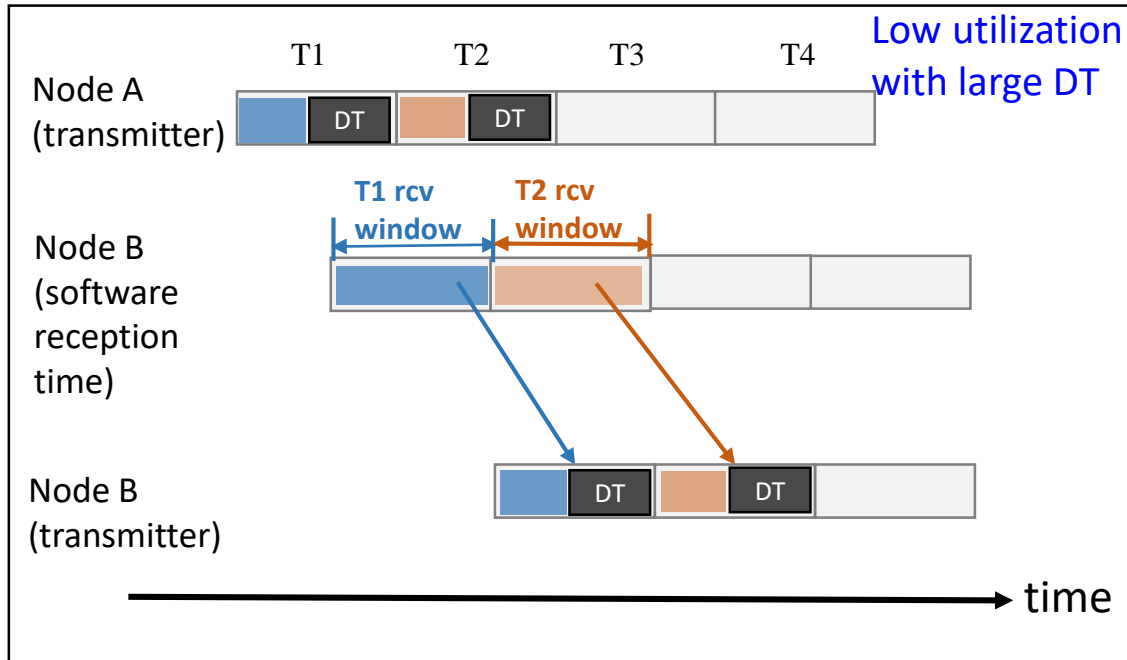
Recap and Goal

- [dv-yizhou-cycle-identification-0922-v02](#) in September meeting
 - Showed the goal to improve the bandwidth utilization in small cycle
 - Discussed the cycle ambiguity problem when making dead time (DT) minimum to improve the utilization
 - Proposed to use implicit reception time indicated by cycle id to determine to which bin a frame goes in addition to direct hardware timestamp based one in 802.1Qdv
 - Provided three options to carry cycle id
- Slides 4-8 recaps the basic concept of using cycle id. Please refer to the previous slides for more details if needed.
- Goal of this deck of slides is to answer the questions received:
 - How to initialize the cycle id based system?
 - How to compute cycle and buffer IDs as they have different ID spaces?
 - How the mapping works when the #of bins(buffers) are different on two neighbor nodes?
 - How large should the cycle ID space be? What are the considerations when choosing the field length of it?

Change log

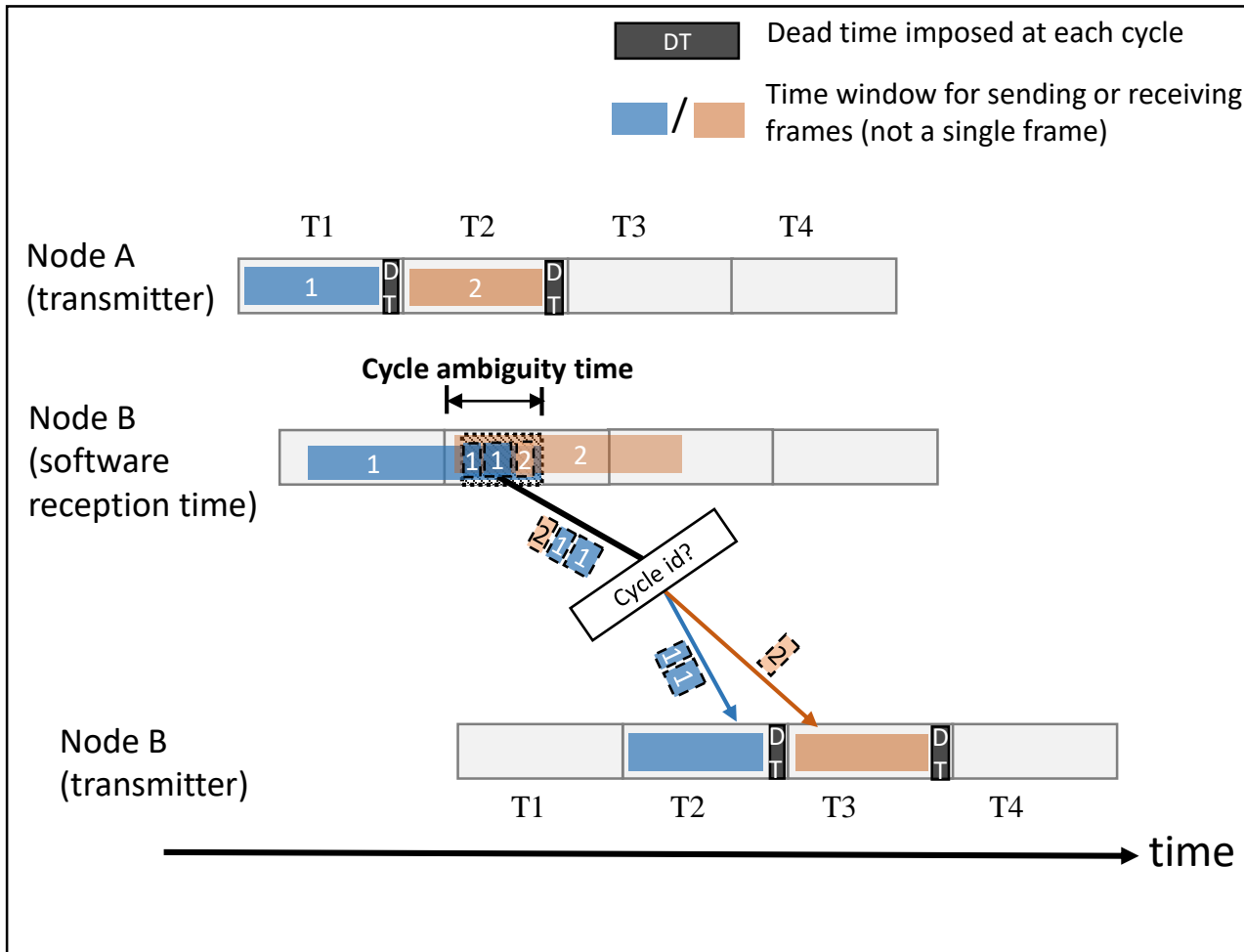
- **dv-yizhou-cycle-identification-details-1122-v02:**
 - Added clarifications for checksum re-compute, multiple cycle time and computation delay.
 - Added explanation on cycle id = implicit pivot reception time
 - Editorial changes from -v01
- **dv-yizhou-cycle-identification-details-1022-v01:**
 - Early dv-yizhou-cycle-identification-0922-v02 showed why and how to use implicit reception time indicated by cycle id to determine to which bin a frame goes
 - Provided further details on that cycle id based implicit time determination approach
 - Focus on the initialization, especially how to establish the cycle id mapping relations for the port pairs.

Goal - Improve the utilization in small cycle T



- **Why low utilization?** DT is relatively too large to a cycle.
- **A straightforward way to improve utilization:** make DT minimum
 - Absorb only the preemption delay instead of the full time variation, i.e. curve out output delay and processing delay
- **A remaining problem:** cycle ambiguity in software reception time based bin determination

Propose to use the explicit cycle identification



- Carry cycle id and change per hop
- Use cycle id as pivot reception time implicitly for output bin determination instead of use software reception time directly
- Remove the ambiguity
- Achieve the good utilization in small cycles by making DT minimum

How to carry a cycle id

1. R-tag (defined in 802.1CB)

Ethertype (F1-C1)	Reserved (16-bit)	Sequence number (16-bit)
-------------------	-------------------	--------------------------

Define a subtype flag and use the last 4-bit in Reserved field for cycle id.

Ethertype (F1-C1)	Reserved (16-bit)			Sequence number (16-bit)
	flag(1)	Rsvd (11)	Cycle id (4)	

How to carry a cycle id (cont'd)

2. Define a new cycle-tag

Ethertype (cycle-tag)	Subtype (4-bit)	Reserved (4-bit)	Cycle ID (8-bit)
--------------------------	--------------------	---------------------	---------------------

How to carry a cycle id (cont'd)

3. Use vlan stacking + vlan mapping function

- Inner vlan is used for cycle id, use ACL to map from ingress cycle id to egress cycle id
- Outer vlan is used as normal vlan based mac learning and forwarding
- Used in a controlled domain, may not be compatible with some existing s-vlan + c-vlan usage

Ethertype (s-vlan)	vlan-tag (16-bit)	Ethertype (c-vlan)	Cycle id (16-bit)
--------------------	-------------------	--------------------	-------------------

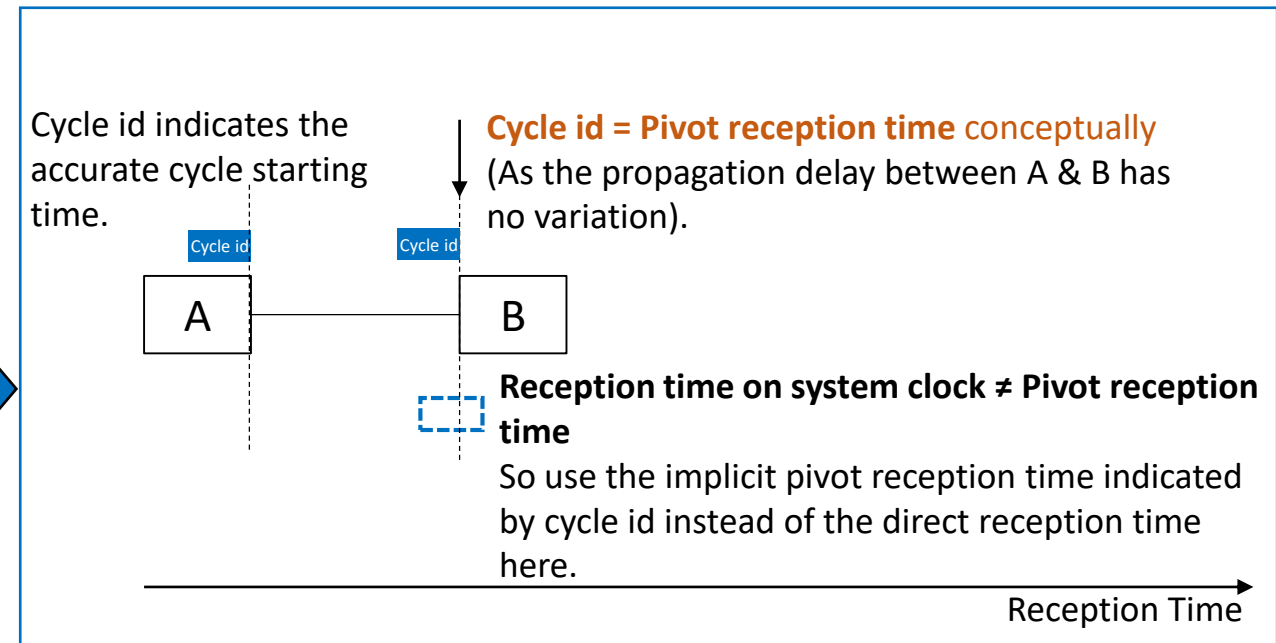
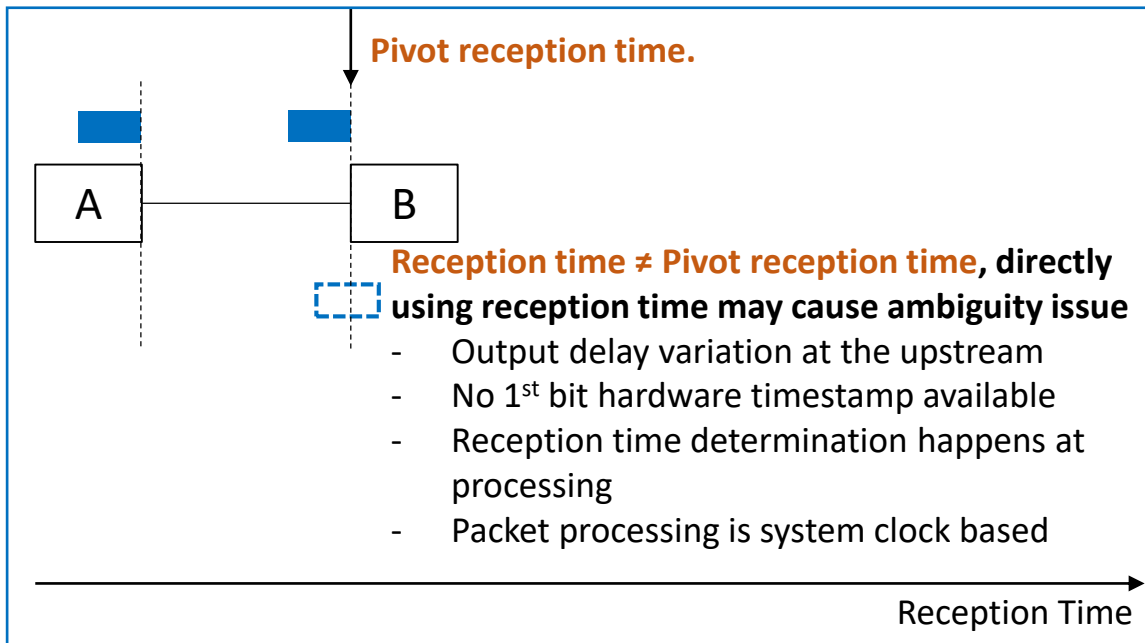
ACL example: `if-match cvlan-id cycle-id-in`

`remark cvlan-id cycle-id-out`

Further details on Cycle id

Conceptually cycle id = pivot reception time

- Conceptually Cycle id = pivot reception time of the frame.
- It is a direct timestamp for cycle start time of the frame's transmission bin on upstream node's output port
- Consequently it indicates the pivot reception time of the frame with no/little time variation experienced on the downstream node's input port



Some assumptions in this deck to simplify the illustration

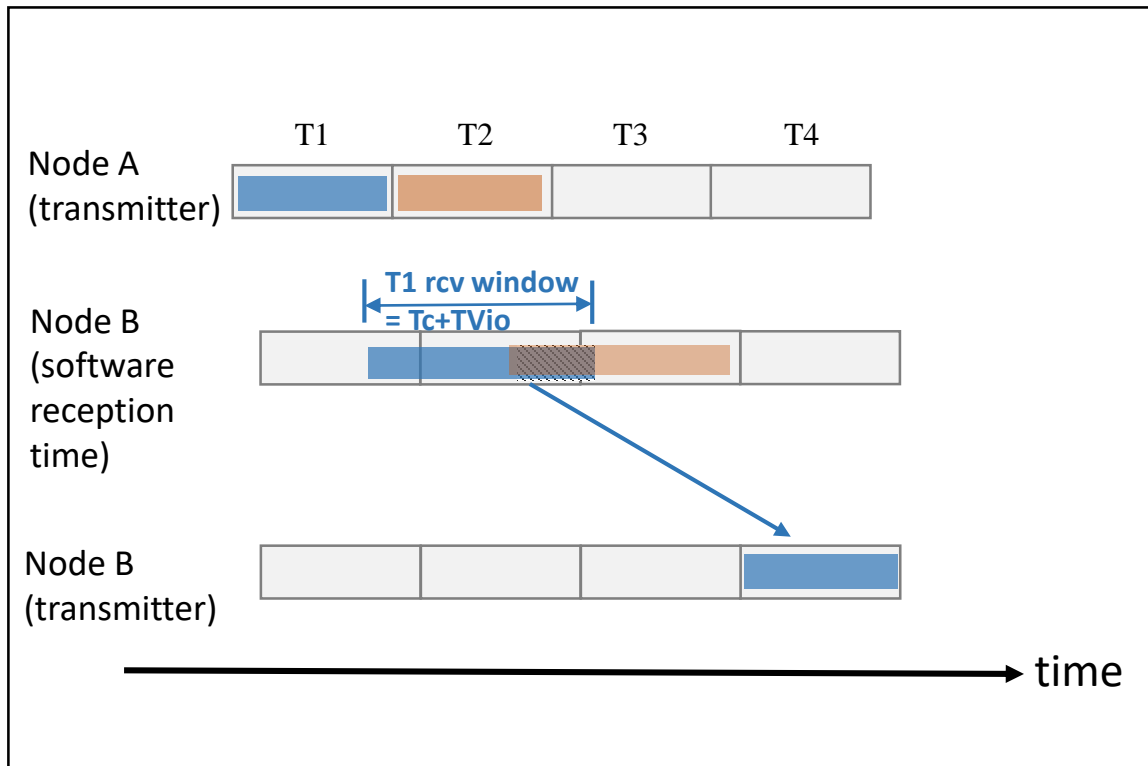
- Only one priority class uses Qdv mechanisms at a physical port.
- The configuration and calculation is for that priority.
- Cycle time T_c remains the same for the same priority in the domain.
- The slides provide a practical way in implementation as example, but it is not the only way.
- Try to reuse terms defined in new-finn-multiple-CQF-0921-v02 as much as possible

Terms – buffer related

- n : the number of physical ports. Range: $0 \sim n-1$
- $B_{i,o}$: the number of buffers required for the input and output port pair (i, o)
- B_o : the number of physical output buffers on port o .
 $B_o = \text{Max}(B_{i,o}), i = 0 \sim n-1$ on a given port o
- buf_id :
 - id assigned to the physical buffer in which a frame should be placed at output port o when the frame is processed.
 - $\text{buf_id} \in [0, B_o - 1]$

Terms – latency related

- TV: sum of the frame-based time variations
 - Include output delay variation, processing delay variation and link delay variation
 - Link delay is the first-bit-out to last-bit-in delay that varies with packet size
- $TV_{i,o}$: TV value for the input and output port pair (i, o).
- Calculate $B_{i,o}$ as follows:



$$B_{i,o} = \text{floor}(TV_{i,o} / T_c) + 4 \quad \text{--- baseline}$$

- Max receiving window time interval $T_c + TV_{i,o}$
- # of receiving buffer to accommodate the max receiving window time:

$$\text{floor}((T_c + TV_{i,o}) / T_c) + 2 = \text{floor}(TV_{i,o} / T_c) + 3$$
- # of sending buffer = 1
- Hence, $B_{i,o} = \text{receiving buffer} + \text{sending buffer}$

$$= \text{floor}(TV_{i,o} / T_c) + 4$$
- Special cases:
 - use one less buffer, i.e. $B_{i,o} = \text{floor}(TV_{i,o} / T_c) + 3$
 - rev window spans over 2 cycles in case of lucky or intentional cycle phase shift between A & B. We assume A&B independently run cycles here, hence always use at least the baseline calculation
 - Use more buffers, i.e. $B_{i,o} = \text{floor}(TV_{i,o} / T_c) + 5/6/./n$
 - Introduce intentional extra delay, e.g. to balance replicated frames over diff paths

Terms – cycle related

- **cycle_id**: cycle id assigned to a frame when the frame is placed to a buffer with **buf_id** at an output port **o**
 - Max value of **cycle_id** **C** is limited by the field length, e.g. if length is 4, then $C = 16$, so that $\text{cycle_id} \in [0, 15]$
 - **C** is the same for all ports on all nodes in a domain
 - increments (modulo **C**) each cycle time
- **N**: the least common multiple over all B_o and **C** in a node. Value of **N** can be different on the nodes.
- **S_i**: Each input port **i** assigns each received frame a logical buffer selector **S_i**, which is an integer in the range 0 through **N-1**, and which increments (modulo **N**) each cycle. The value of **S_i** is **s_i**.
 - **s_i** is directly used by data frames, not for (cycle) mapping determination frames.

4 steps in initialization and determination in the following slides

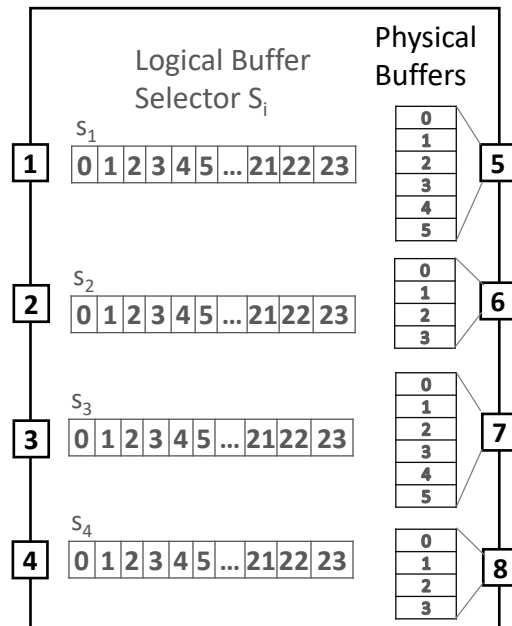
- Step 1: Base parameter provisioning
- Step 2: Initialization of buffer selector and output buffers
- Step 3: Determination of cycle id mapping relation for port pair (i,o)

One-time static calculation and provisioning

-
- Step 4: Determination of output cycle id and buffer id for data frame
 - This step is for real data traffic, not part of initialization and provisioning.

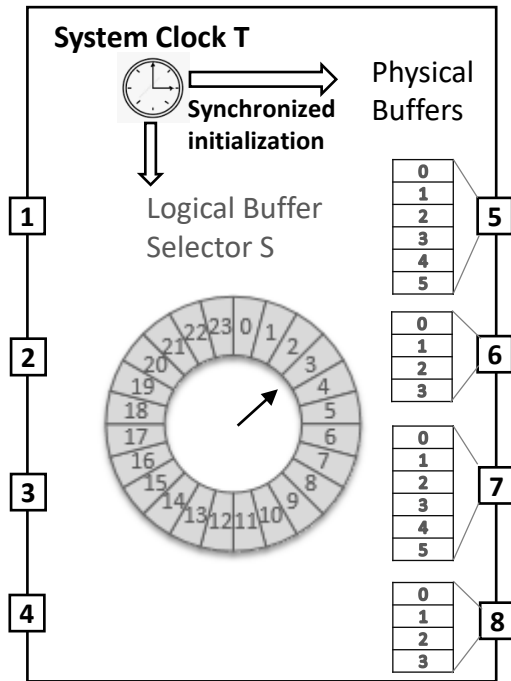
for each data frame

Example Step 1– base parameter provisioning



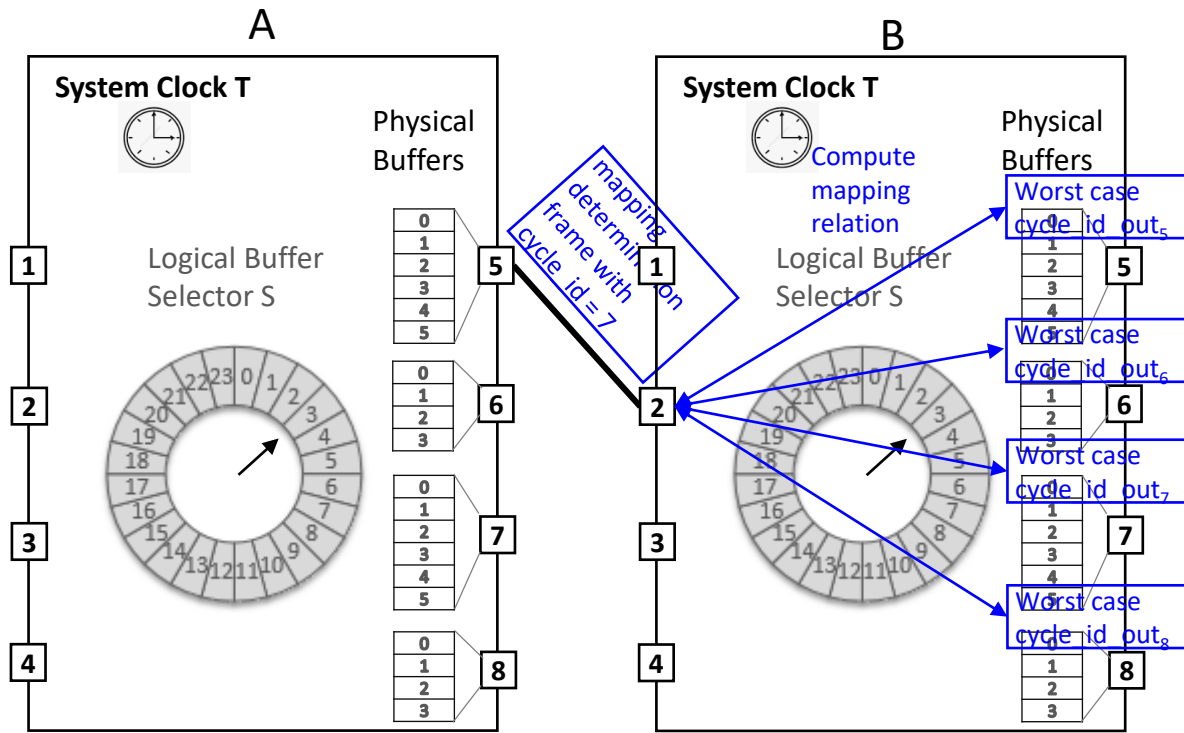
- The node has $n=8$ port.
- Cycle time T_c is set and time variation $TV_{i,o}$ for port pair (i,o) is known.
- Calculate $B_{i,o} = \text{floor}(TV_{i,o} / T_c) + 4$ and B_o is the $\max(B_{i,o})$ for a port o .
- Assume the calculation results of the number of physical buffers B_o on each port:
 - $B_5 = B_7 = 6, B_6 = B_8 = 4$
 - $B_1 - B_4$ not shown, assume all are 4
- Assume Max # of cycle_id $C = 8, \text{cycle_id} \in [0, 7]$
 - i.e. length of cycle_id = 3 bit. (This is for picture simplicity. In reality, it should be larger)
- Then $N = 24$ (the least common multiple of 4,6,8)
- Logical buffer selector S_i makes value s_i rotate between 0 – 23.
 - Note: It is the index assignment for calculation simplicity. No real buffer attached.

Example Step 2 – initialization of buffer selector and output buffers



- All the ports of a node use the same system time T . T is the time elapsed (in nanosecond) from the latest system startup. All the components in the node can use T as a time source.
- **System Synchronized Initialization:**
 - Initialize the starting time for S_i and buf_id on all ports to be multiple of $N \cdot T_c$.
 - Logically make the position and shifting for all S_i the same, thus a single selector S can be used to simplify the implementation. Let s be the value of the selector S .
 - Facilitate halfway port enabling
 - Initialize $s = 0$ and the transmitting buffer id = 0 on all ports
 - increments (modulo N) each cycle
- At any time T , the following computations hold
 - $s = (T \bmod (N \cdot T_c)) / T_c$ ---- “/” is floor division
 - buffer id currently transmitting frames on port o :
Tx-buf-id = $s \bmod B_o$
 - cycle id used by currently transmitting buffer:
Tx-cycle-id = $s \bmod C$

Example Step 3 – Determination of cycle id mapping relation for port pair (i,o)



Pre-requisite:

- Node A & B have run their cycles independently with cycle T_c

Purpose: For the input and output port pair (i, o) on B, determine mapping parameter $M_{i,o}$ so that a stable cycle mapping equation $\text{cycle_id_out} = (\text{cycle_id_in} + M_{i,o}) \bmod C$ can be used for future data frames over port pair (i,o) no matter what time variation are experienced by that frame.

Mapping determination frame: special frame to determine the cycle mapping relation between two neighbors during system initialization and auditing. It experiences the least time variation

- Shortest frame – 64B
- Highest priority

Cycle mapping relation computation example:

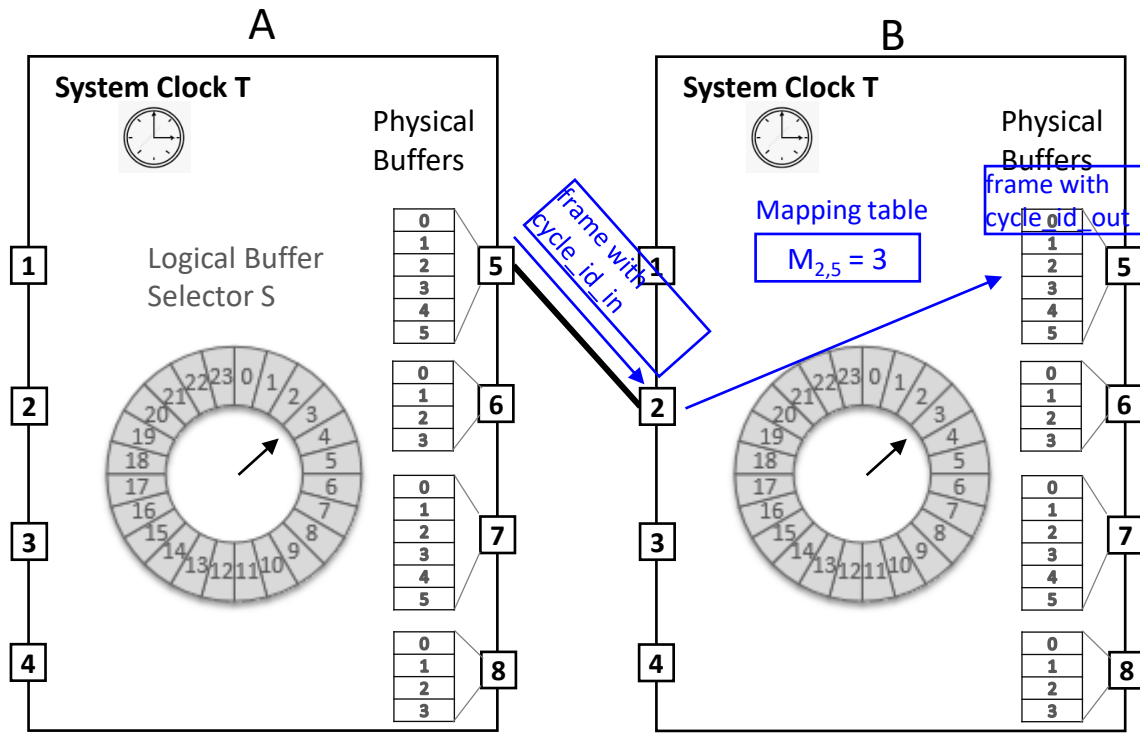
- B receives the **mapping determination frame** from port 2
- Assume B's system time is T. Compute the potential worst case s for the frame over each output port o as:

$$s_o = ((T + TV_{2,o}) \bmod (N * T_c)) / T_c$$
- Compute cycle_id_out for the potential worst case frame over each output port o as:

$$\text{cycle_id_out}_o = (s_o + 1) \bmod C$$
- Compute $M_{i,o} = (\text{cycle_id_out}_o - \text{cycle_id_in}_i + C) \bmod C$
 - Assume the frame carries cycle_id_in = 7 from the incoming port 2 and the computed cycle_id_out is 2 on output port 5, then

$$M_{2,5} = (2 - 7 + 8) \bmod 8 = 3$$
- $M_{i,o}$ can be simplified to M_i for a port i if M_i takes $\text{Max}(M_{i,o})$ for a specific i. That will make the # of mapping parameters $O(n)$ instead of $O(n^2)$, $n = \#$ of ports.

Example Step 4 – Data frames: determination of output cycle id and buffer id



Output cycle id:

- $\text{cycle_id_out}_o = (\text{cycle_id_in}_i + M_{i,o}) \bmod C$
- e.g. as $M_{2,5} = 3$ for port pair (2,5), the cycle_id determination is defined as

$$\text{cycle_id_out} = (\text{cycle_id_in} + 3) \bmod 8$$

Output buffer id:

- Cycle id used by currently transmitting buffer:

$$\text{Tx-cycle-id} = s \bmod C$$
- Compute the id offset between targeting placement buffer and current transmitting buffer as:

$$\text{offset} = (\text{cycle_id_out}_o - \text{Tx-cycle-id} + C) \bmod C$$
- Compute the buf_id on port o to place the frame as:

$$\text{buf_id}_o = (s + \text{offset}) \bmod B_o$$

Summary of calculations – on each node, for each priority

Step 1: Base parameter provisioning

1. Set cycle time T_c and time variation $TV_{i,o}$ for port pair (i,o)

2. Calculate number of physical buffers B_o on each port o as:

- $B_{i,o} = \text{floor}(TV_{i,o} / T_c) + 4$
- $B_o = \max(B_{i,o})$ for port o

One-time static calculation and provisioning

Step 2: Initialization of buffer selector and output buffers

3. Calculate N = the least common multiple of B_o and C

4. Start system synchronized initialization on every node for buffer selector S and ports

Step 3: Determination of cycle id mapping relation

5. Use Mapping Determination Frame to calculate the mapping relation value $M_{i,o}$ for port pair (i,o) as follows:

- $s_o = [((T + TV_{i,o}) \bmod (N * T_c)) / T_c]$
- $\text{cycle_id_out}_o = (s_o + 1) \bmod C$
- $M_{i,o} = (\text{cycle_id_out}_o - \text{cycle_id_in}_i + C) \bmod C$

Step 4: Data frame mapping

6. For an incoming frame for port pair (i,o), calculate the output cycle id **cycle_id_out_o** and output buffer id **buf_id_o**

- output cycle id: $\text{cycle_id_out}_o = (\text{cycle_id_in}_i + M_{i,o}) \bmod C$
- output buffer id:
 - $\text{Tx-cycle-id} = s \bmod C$
 - $\text{offset} = (\text{cycle_id_out}_o - \text{Tx-cycle-id} + C) \bmod C$
 - $\text{buf_id}_o = (s + \text{offset}) \bmod B_o$

for each data frame

Choice of cycle id length

- Number of cycle id $C = 2^L$ where L is the length in bits of cycle id
- C should be larger than the number of physical buffer B_o on any port o in step 1 base parameter provisioning
- Discuss: L to be at least 5?
 - $C = 32$ ($L = 5$) should be sufficiently large as it can support a 32-buffer port.
 - If we make $L=8$, that would allow 256-buffer on any port.
- Preferred way of tagging (back to slide 6-8)?

Checksum re-compute

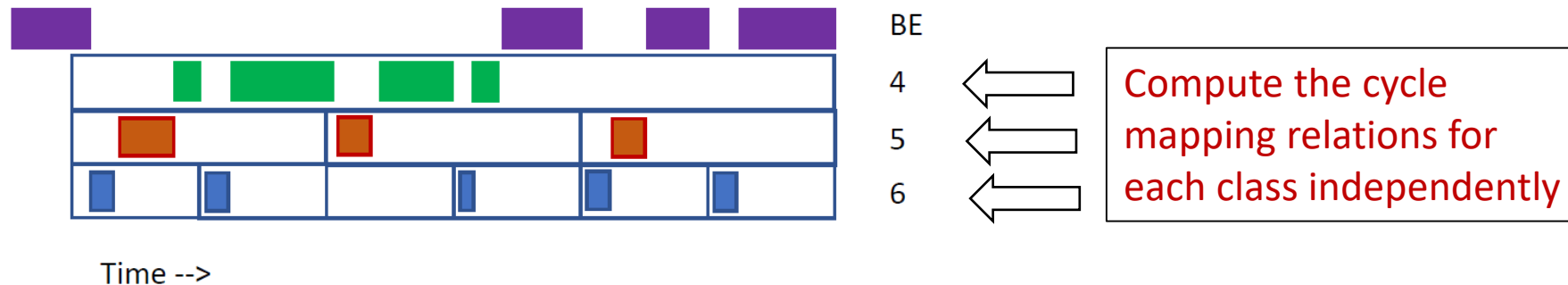
Will Checksum be required to re-compute per hop?

- Yes as the cycle id is changed every hop, but it does not introduce additional cost.
- Checksum re-compute is always performed in normal forwarding at the intermediate nodes regardless of whether any header value changes.
- Hence no additional checksum re-compute cost caused by cycle id change per hop.
- If VLAN variant tagging is used, it is equivalent to normal VLAN mapping procedures in switch implementation.

Support of multiple time cycles

The slides showed the example for a single priority traffic class which has a single time cycle. It can also support multiple time cycles if CQF enhancement mechanisms are enabled on multiple priority traffic classes.

- Use the multiple cycle time picture in new-finn-multiple-CQF-0921-v02 as example.



- Each priority traffic class to establish the cycle mapping relations independently.

Computation delay considerations

- Most computation is one-time thing at the initialization. Not contribute to the data plane delay.
- For each data frame, it is simply a mapping from `cycle_id_in` to `cycle_id_out` and determination of buffer id using add operations.
- Additional computation delay for each data frame is negligible.

Takeaways and Discussions

- Cycle id conceptually is an implicit pivot reception time of the frame. Not impacted by the time variation experienced by the frame.
 - It helps to remove the ambiguity of cycle identification and improve the bandwidth utilizations when software reception time is used.
 - The calculations for initialization and mapping relation establishment are feasible.
 - Per-frame data plane computation is simple.
- Good or bad idea for inclusion in Qdv?