

Robust Time Delivery – Turning a Wish Into a Reality (How To Use PTP With A Ring Architecture)

David Spencer – WSTS 2014

Ring Architectures



Ring architectures are already common for data transport

- Traffic can pass from any node to any other node through one of two routes
 - Traditionally called East and West (viewed from bottom of ring)
- Provides alternate route in event of single failure
 - As long as ring is geographically separated

For communication between any pair of nodes one path is likely preferred

- Normally the shortest one
- But other path provides a reasonable backup
- Different for different nodes
 - For example SW9 would talk to SW1 anticlockwise while SW5 would talk to SW1 clockwise

But how does this fit in with PTP?

Does a ring of boundary clocks work?





Simple BC Ring

Boundary Clock Ring: Node Architectures A 'normal' Boundary clock



Standard BC:

- A PTP slave port locks to a selected master via a selected port creating a timebase.
- The Master, all egress ports *and* the selected output are tied to the slave timebase.
- The slave port may or may not be able to hitlessly switch on failure.



Let's call this a 'Single BC'

The BMCA – A Quick Refresher



Mechanism to select best master

- When multiple masters are available selects the best one based on number of factors
 - Programmed priorities
 - Characteristics of clock source

□ What about a ring

- Same master will be seen from both sides
- Second stage of BMCA handles this
 - Looks at Steps Removed field
 - Prefers links through least number of BCs
- Path closest to GM will be selected
 - Either east or west depending on where in ring





Boundary Clock Architecture: Conceptual- Single GM, 8-BC node ring



Primary Flow Secondary Flow

6



Each node has two potential masters- selected in Slave by BMCA The Slave's timebase is used for all the egress ports and the output

Boundary Clock Architecture: Single GM, BC nodes, with transmitted 'steps removed'





© 2014 Semtech Corporation

Boundary Clock Architecture: Failure: first action, everything controlled by SR

SEMTECH



© 2014 Semtech Corporation

Boundary Clock Architecture: Failure: first action, everything controlled by SR





Boundary Clock Architecture: Failure: second action





Boundary Clock Architecture: Failure occurs: third action





Boundary Clock Architecture: Failure occurs: Fourth action





© 2014 Semtech Corporation

Boundary Clock Architecture: Failure occurs: Final action







Single BC ring Lab test

Boundary Clock Architecture: Single GM, Lab Setup Initial state





Timing flow in error-free condition

Network Manager Tool: A Display of Networked PTP elements





Initial Lab Setup: GUI screenshot



Grand Master Unknown PTP node 15.88.27.100:1 BC1 15.88.27.30 BC1 BC1 Port 1: Slave 15.88.27.30 BC1 Error: -52,504 ns Port 3: Master Time: 4 June 2014 13:09:37 🗲 ва 15.88.27.31 BC₂ BC2 Port 1: Slave 15.88.27.31 BC₂ Error: -2.699 ns m 62 Port 3: Master Time: 4 June 2014 13:09:37 A BC3 15.88.27.32 BC3 BC3 Port 1: Slave 15.88.27.32 BC3 Error: 20.354 ns Port 3: Master Time: 4 June 2014 13:09:37 SC4 15.88.27.33 BC4 BC4 Port 1: Slave 15.88.27.33 BC4 Error: 19.770 ns C, Port 3: Master Network Cloud Time: 4 June 2014 13:09:37 BC5 15.88.27.34 BC5 BC5 Port 1: Slave 15.88.27.34 BC5 Error: -17.208 ns Port 3: Master Time: 4 June 2014 13:09:37 🗲 воз 15.88.27.35 BC6 BC6 Port 1: Slave 15.88.27.35 BC6 Error: -77.888 ns 6. Port 3: Master Time: 4 June 2014 13:09:37 🚄 вст BC7 15.88.27.36 BC7 Port 1: Slave BC7 15.88.27.36 Error: -80.353 ns 2. Port 3: Master Time: 4 June 2014 13:09:37 BO8 15.88.27.37 BC8 BC8 Port 1: Slave 15.88.27.37 BC8 Error: -2.951 ns Port 3: Master Time: 4 June 2014 13:09:37 Master ports **Slave ports**

Each BOX is a PTP *PORT*

Each BC has two PORTSone is Master & one is Slave

Single BC loop: Initial Lab Setup Semtech GUI screenshot





Single BC loop: Initial Lab Setup Semtech GUI screenshot





Initial Lab Setup: GUI screenshot





Initial Failure: BC1 to BC2 BC2 slave port goes passive





Failure + 5s



© 2014 Semtech Corporation







© 2014 Semtech Corporation



© 2014 Semtech Corporation



© 2014 Semtech Corporation

Single BC Ring- Summary



- **D PTP** Rings can be created
- □ 'Steps removed' is the primary arbitration in each node
- □ Timing passes in two directions from GM to furthest point
- **Temporary timing loops are created**
- □ A failure triggers a chain of events that take time to settle
- Every node that 'switches' may have a 'hit' whilst slave settles
- □ What to do on recovery?



Dual BC Ring

Let's consider other Node Architectures



□ Can we make fault tolerance hitless?

□ Can we remove transient timing loops?

□ Can we improve failover and recovery times?

Boundary Clock Ring: Node Architectures A 'normal' Boundary clock





Single BC: A single slave port locks to a selected master creating a single timebase.

The Master, all egress ports and the selected output are tied to the slave timebase.

The slave port may or may not be able to hitlessly switch on failure.



Dual Slave: Two independent slave ports create two internal timebases.

A single Master selects a slave timebase and uses it for all egress ports and the selected output.

The master can hitlessly switch between slaves on failure.



Dual BC: Two completely independent boundary clocks.

Two slaves create a timebase for two dedicated masters.

One of the two timebases is selected as the output.

Master ports never switch on failure, merely enter holdover.

All switching is performed hitlessly in the output selection.

Node Architectures Summary



□ 'Single BC'

- Traditional in the sense that it selects a timing ingress port and locks to it
- 'Crosses that bridge...'
 - ie. When a failure occurs, then press the mechanisms into action, not before

'Dual Slave'

- A BC with two independent slave ports.
- Each slave is permanently locked to the selected master
- Switching performance is hitless, as slaves are pre-locked but fault management is the same as the 'Single BC'

'Dual BC'

- Two completely independent boundary clocks
- Can maintain a timing flow in both directions
- Permanently ready for an instant and hitless response to a failure

Dual Boundary Clock Architecture: Conceptual- Single GM, 2 contra-rotating 'sync rings'





Each node has two independent boundary clocks Each slave has only ONE acceptable master- the upstream node Each slave is dedicated to a master Node phase output can switch hitlessly between timebases based on BMCA

33

Primary Flow Secondary Flow

Dual BC Test setup 1: Lab setup (add a second GM for clarity)



CC=Clock Class SR=Steps Removed



Dual BC Test setup 1: Failure between BC1 & BC2







Dual BC ring Lab test

How the dual BC-Ring looks in the GUI





Each cloud represents a domain

- Clockwise ring domain 1
- Anticlockwise ring domain 10

□ Two grandmasters

One per domain

□ Four PTP ports per node

- Two slaves
- Two masters

<u>M</u> anaged De	vices			
Device/Port	Address	State	Domain	Sh
⊟ BC1	15.88.27.30			
Port 1	15.88.27.30	Slave	1	~
Port 2	15.88.27.30	Slave	10	
Port 3	15.88.27.30	Master	1	
Port 4	15.88.27.30	Master	10	~
🖃 BC2	15.88.27.31			
Port 1	15.88.27.31	Slave	1	
Port 2	15.88.27.31	Slave	10	
Port 3	15.88.27.31	Master	1	
Port 4	15.88.27.31	Master	10	
🖃 BC3	15.88.27.32			
Port 1	15.88.27.32	Slave	1	
Port 2	15.88.27.32	Slave	10	
Port 3	15.88.27.32	Master	1	
Port 4	15.88.27.32	Master	10	
🖃 BC4	15.88.27.33			

Prior to the failure:





Prior to the failure:







Failure +4 seconds: only change is BC2 port 1 enters holdover



SEMTECH

No further change: as long as the failure persists



				CI	
				3	
lan	aged De	vices			
De	vice/Port	Address	State	Domain	Sh
-	BC1	15.88.27.30			
	Port 1	15.88.27.30	Slave	1	~
	Port 2	15.88.27.30	Slave	10	~
	Port 3	15.88.27.30	Master	1	~
	Port 4	15.88.27.30	Master	10	~
-	BC2	15.88.27.31			
	Port 1	15.88.27.31	Passive	1	
	Port 2	15.88.27.31	Slave	10	
	Port 3	15.88.27.31	Master	1	
	Port 4	15.88.27.31	Master	10	
-	BC3	15.88.27.32			_
	Port 1	15.88.27.32	Slave	1	
	Port 2	15.88.27.32	Slave	10	
	Port 3	15.88.27.32	Master	1	
	Port 4	15.88.27.32	Master	10	
=	BC4	15.88.27.33			
	Port 1	15.88.27.33	Slave	1	
	Port 2	15.88.27.33	Slave	10	
	Port 3	15.88.27.33	Master	1	
	Port 4	15.88.27.33	Master	- 10	
-	BC5	15.88.27.34			
	Port 1	15.88.27.34	Slave	1	
	Port 2	15.88.27.34	Slave	- 10	
	Port 3	15 88 27 34	Master	1	
	Port 4	15 88 27 34	Master	10	
=	BC6	15 88 27 35	Mascel	10	•
	Dort 1	15,88,27,25	Slava	1	
	Port 2	15,00,27,35	Slave	10	
	Port 2	15.00.27.33	Diave	10	
	PURUS Dort 4	15.00.27.35	Macher	10	
_	FUIL 4	15.00.27.33	master	10	•
	Dort 1	15.00.27.30	Slave	1	
	Port 1	15.00.27.30	Slave	10	
	PULC	10,00,27,30	DIAVE	10	*

Managed Device:

Port 4 BC5

Port 4 BC7

Port 3

Port 4

Port 1

Port 2

Port 3

Port 4

E BC8

15.88.27.36

15.88.27.36

15.88.27.37

15.88.27.37

15.88.27.37

15.88.27.37

15.88.27.37

Master

Master

Slave

Slave

Master

Master

1

10

1

1

10

10

🖃 BC6

🖃 BC1

🖃 BC2

BC3

🖃 BC4

6/5/2014

~

~

~

~

~

~

Conclusion



Rings of boundary clocks work

Simple boundary clock

- Single slave instance
- Has to relock on network break
- Will recover in time

Dual-slave boundary clocks

- Two slave instances so backup path is pre-locked
- Can failover immediately on network break
- Much shorter recovery time

Dual-path boundary clock

- Two boundary clocks in one
- Dual redundant flows at all times
- No rearrangement needed on network break

And If It Goes Wrong ... ?



□ We all have to work to an impossible hour to fix it



Thank you to Rich Lansdowne and Malcolm Green for planning and running these tests