



A Deeper Look into the Inner Workings and Hidden Mechanisms of FICON Performance

- David Lytle, BCAF
- Brocade Communications Inc.
- Monday August 10, 2015 11:15 am to 12:15 pm
- Session Number 17506







Legal Disclosure



Nothing in this presentation shall be deemed to create a warranty of any kind, either express or implied, statutory or otherwise, including but not limited to, any implied warranties of merchantability, fitness for a particular purpose, or non-infringement of third-party rights with respect to any products and services referenced herein.

ADX, Brocade, Brocade Assurance, the B-wing symbol, DCX, Fabric OS, HyperEdge, ICX, MLX, MyBrocade, OpenScript, The Effortless Network, VCS, VDX, Vplane, and Vyatta are registered trademarks, and Fabric Vision and vADX are trademarks of Brocade Communications Systems, Inc., in the United States and/or in other countries. Other brands, products, or service names mentioned may be trademarks of others.



Notes as part of the online handouts



I have saved the PDF files for my presentations in such a way that all of the audience notes are available as you read the PDF file that you download.

If there is a little balloon icon in the upper left hand corner of the slide then take your cursor and put it over the balloon and you will see the notes that I have made concerning the slide that you are viewing.

This will usually give you more information than just what the slide contains.

I hope this helps in your educational efforts!



A deeper look into the Inner Workings and Hidden Mechanisms of FICON Performance



This technical session goes into a fairly deep discussion on some of the design considerations of a FICON infrastructure.

- Among the topics this session will focus on is:
 - Switch Firmware
 - Buffer Credit Flow Control
 - Fabric Congestion, Backpressure and Slow Drain Devices
 - Buffer Credit Starvation
 - Data Encoding, Forward Error Correction and Compression/Encryption

NOTE: Please check for the most recent copy of this presentation at the SHARE website as I make frequent updates.



This Section

- Switch Firmware
- Hardware Support





Switch Firmware

Switching devices run it – and it evolves as technology evolves!

- If a chassis is the body of a switch then the central processing unit (CPU) is the brain and the firmware code is the intelligence residing within that brain on that switch.
- It is the firmware that issues requests out to the rest of the "body" and then monitors and reacts to all of the "sensory information" coming back in to it
- Each switch must have an appropriate release of firmware intelligence in order to run all of the features and functions of the switch:
 - Brocade's firmware is called Fabric Operating System or Fabric OS or more commonly just FOS.
 - Cisco's firmware is called Nexus Operating System (NX-OS) or NxOS.
- Both vendors firmware is used to control the features and functions of their FC and IP SAN products. It is how vendors deliver choice and solutions to customers across multiple generations of products. In a nutshell, firmware helps a switch do what it is supposed to do!
- Current FICON Firmware: Brocade FOS 7.4.0a Cisco – NxOS 6.2(11c)





Previous Generations of IBM Mainframes

End-of-Life For Switch Support Information

• **z9** is not supported on Brocade or Cisco at current firmware releases!

- z9 was supported by Brocade FOS in Sept 2007
 - For Brocade, FOS 5.3 was the first FOS to support z9
 - For Brocade, FOS 7.2.1d is the last FOS to support z9
 - For Cisco, NxOS 6.2(5a) is the last NxOS release to support z9
 - There are now FOS and NxOS releases beyond those above
- **z10** is not supported on Brocade or Cisco at current firmware releases!
- z10 was supported by Brocade FOS in June 2008
 - For Brocade, FOS 6.1.0c was the first FOS to support z10
 - For Brocade, FOS 7.3.1b is the last FOS to support z10
 - For Cisco, NxOS 6.2(5a) is the last NxOS release to support z10
 - There are now FOS and NxOS releases beyond those above

Be very careful of what mainframes are active in your environment and choose the appropriate release of vendor firmware that fully supports all those devices when attached to an I/O infrastructure fabric.

Complete your session evaluations online at www.SHARE.org/Orlando-Eval

MAINERAME

R.I.P.

ZSERIES

z9

z10

This Section

- How FICON uses Buffer-to-Buffer Credits
- Determining Buffer Credits Required





FICON Channel Card Buffer Credits (BBC)

FICON Express8



Buffer-to-Buffer Credits (BBC) How they get set - same for FCP and FICON



- At initialization, the 2 ports on each side of a link establish BBC counts:
 - Receiver tells the transmitter the number of frames the receiver can handle
- Each port on different ends of the link can support different BBC values:
 - BBC does not have to be symmetrical (but do make ISLs symmetrical!)
- If a port does not have a buffer credit available, it cannot send a frame:
 - Buffer Credit Count has reached zero on the Transmitter (tx)
- This Mechanism provides the frame's flow control and limits frame drops



- Each and every frame requires up to 2,148 bytes of BBC memory in order to store the frame just before it is sent across the link to the receiver
- Actual frame sizes vary widely but, when stored in BBC memory, always consume a full buffer credit







Buffer-to-Buffer Credits Buffer-to-Buffer flow control



- After initialization, each port knows how many buffers are available in the queue <u>at the other end of the link</u> the transmitter is controlled by this
 This value is known as Transmit (Tx) Credit
 - The purpose is to keep the transmitter from overrunning the receiver



Buffer-to-Buffer Credits Buffer-to-Buffer flow control



- Tx Credit is decremented by one for every frame sent from the CHPID
- No frames may be transmitted after Tx Credit reaches zero
- Tx Credit is incremented by one for each R_RDY received from F_Port



Buffer-to-Buffer Credits (BBC) An Example Showing Several Buffer Credit Deployment Scenarios 16 Gbps Link running more like an 8 Gbps Link Start Of Frame x 🔸 R RDY or VC RDY Acknowledgements Ample BB Credits Well Running 16 Gbps Link (Acks do not use BBCs) SOFx 🕈 R_RDY or VC_RDY Acks 8 x less BB Credits, still ample 2 Gbps Link If 870 bytes. requires SOFx -.82 km/frame R_RDY or VC_RDY Acks 4 Gbps, 870 byte frames, require .41 km of link distance/frame 8 Gbps, 870 byte frames, require .20 km of link distance/frame 16 Gbps, 870 byte frames, require .10 km of link distance/frame

This Section

- Congestion and Backpressure Overview
- Slow Draining Devices





Congestion and Backpressure Overview



These two conditions are not the same thing

- Congestion occurs at the point of restriction
- **Backpressure** is the effect felt by the environment leading up to the point of restriction

I will use an Interstate highway example to demonstrate these concepts



Congestion and Backpressure Overview



No Congestion and No Backpressure

- The highway handles up to 800 cars/trucks per minute and less than 800 cars/trucks per min are arriving
- Time spent in queue (behind slower traffic) is minimal
 - Traffic is flowing very well





Congestion and Backpressure Overview



Congestion

- The highway handles up to 800 cars/trucks per minute and more than 800 cars/trucks per min are arriving
- Trip time (Latency) and vehicle spacing (BBCs consumed) increases
 No more vehicles can be placed onto this road
- Backpressure is experienced by cars slowing down and queuing up



Slow Draining Devices can also jam up a path Think about police cars, side-by-side, going 55 MPH, slowing all traffic

- Slow draining devices are <u>receiving ports</u> that cannot accept the amount of data flowing to them so they actively respond to slow down the link:
 - They do not return Buffer Credit acknowledgements at link rate
 - Data flowing to these devices use up all the BBCs along the path
 - Finally the transmitter runs out of buffer credits cannot transmit the next frame
 - It will only send a frame after an acknowledgement is received gets 1 BBC
 - I/O on the data path is now running only at the speed of the slow draining device
- It is very important to note that this slow down can spread into the fabric and can affect unrelated flows in the fabric trying to use common links

• What causes slow draining devices?

- The most common cause is within the storage device, CHPID or the server itself.
- It often happens because a target port has a slower link rate (e.g. 2G, 4G, 8G) than the I/O source port (e.g. 4G, 8G, 16G) or the Fan-In from the rest of the environment overwhelms the target port.



Slow Draining Devices – 16G DASD

16G DASD is available today – implementing it could cause slow drains!



- This is potentially a very poor performing, infrastructure!
- DASD is about 90% read, 10% write. So, in this case the "drain" of the pipe is the 8Gb CHPID and the "source" of the pipe is the 16Gb storage port.
- The Source can out perform the Drain!
- This can cause congestion and back pressure towards the CHPID. The switch port leading to the CHPID becomes a slow draining device.
- Let us look at this process in more detail!



Congestion and Back Pressure (watch out for ISLs!)

Too many customers will be putting 16G DASD onto 8G CHPIDs (slow drain)



- The is a simple representation of a single CHPID connection
- Of course that won't be true in a real configuration and the results could be much worse and more dire for configuration performance

Physical Tape Drives in a Fabric



BUT... If 8G Tape runs at 500 MBps (2X) it will take 1 GBps to feed it and the best CHPID is <= 620 MBps so compression ratio will be low and maybe Start/Stop!

- 4G tape WILL BE a slow drain Tape is about 90% write and 10% read on average
- The FX8S/16S using Command Mode FICON can transfer (write) ~620MBps
- BUT... a 4G Tape interface handles only about 380MBps (400MBps * .95 = 380MBps)
- AND...the maximum bandwidth 4G tape needs to accept and compress is 504 MBps (at 2:1 compression) for Oracle (T10000D – 252 MBps head speed) and about 500 MBps (at 2:1 compression) for IBM (TS1140 – 250 MBps head speed)
- So a single CHPID attached to a 380 MBps, 4G physical tape interface:
 - Can stream ONLY 1 IBM tape/CHPID (620 / 380 = 1.63) poor compression on 1
 - Can stream ONLY 1 Oracle tape/CHPID (620 / 380 = 1.63) poor compression on 1

This Section

- Determining How Many Buffer Credits Are Required
- RMF Reports for Switched-FICON



Buffer Credits – required on every FC link!

Why FICON Never Averages a Full Frame Size

- There are three metrics that are required to determine the number of buffer credits required for a link – absolutely necessary for ISL links of long length:
 - The speed of the link (pretty easy to figure out)
 - The cable distance of the link (DWDM can make this harder to obtain)
 - The average frame size (historically, the most difficult metric to determine)
- Average frame size is the hardest to determine compression can even change it!
 - RMF 74-7 records or the *portbuffershow* CLI command find Avg. Frame Size (Tx)
 - You will find that FICON just never averages full frame size
 - Below is a simple FICON 4K write that demonstrates average frame size



Average = (76+2016+2084+100+68) / 5 = 869 Bytes

Buffer Credits really needed for ISLs over distance!

The Impact of Average Frame Size on Buffer Credits – example is 50 km

A distance	of 50KM wi	th 100% link	utilization	Additional	2Gbps	4Gbps	8Gbps	10Gbps	16Gbps
SOF, Header, CRC, EOF	Payload	Total Frame Bytes (Average Frame Size)	Smaller than full frame by xx%	Compression and a little head room	Buffer Credis Required 8b10b	Buffer Credis Required 8b10b	Buffer Credis Required 8b10b	Buffer Credis Required 64b66b	Buffer Credis Required 64b66b
36	2112	2148	0.00%	NO 2.2:1	56 N/A	105 N/A	204 442	254 551	402 877
36	1824	1860	13.41%	NO 2.2:1	64 N/A	121 N/A	235 509	292 635	464 1012
36	1114	1150	46.46%	NO 2.2:1	99 N/A	191 N/A	376 820	469 1023	746 1633
36	964	1000	53.45%	NO 2.2:1	113 N/A	219 N/A	432 942	538 1175	857 1877
36	834	870	59.50%	NO 2.2:1	129 N/A	251 N/A	495 1081	617 1350	984 2156
36	476	512	76.16%	NO 2.2:1	214 N/A	422 N/A	837 1833	1044 2290	1667 3660
36	164	200	90.69%	NO 2.2:1	538 N/A	1069 N/A	2132 4682	2663 5851	4257 9358

Obtaining the best FICON reporting

Obtain RMF reports on the switched-FICON infrastructure





Control Unit Port (CUP) provides users with an ability to monitor and manage FICON fabrics and cascaded fabrics. RMF, Systems Automation.

Also, some z/OS functions use CUP for communications and control.

- Long ago IBM made a decision that any resource that could impact performance and I/O predictability must be reported on by RMF.
- One such impactful resource is buffer credits.
- However, IBM chose to only report on buffer credits within RMF when switches are deployed to handle the I/O infrastructure.
- So a user must deploy switched-FICON and then that user can deploy and enable the CUP feature in order to get this data into RMF:
 - RMF 74-7 records



Switched-FICON Lets You Use CUP

Control Unit Port Is A Great Tool for FICON



- CUP code is part of a switches firmware:
 - Used to provide FICON Director Reports in RMF (buffer credits, frame size, etc.)
 - Used for in-band communication with Systems Automation
 - Supports Prohibit Dynamic Connectivity Mask (PDCMs)
 - Supports Dynamic Channel Path Management (DCM)
 - Provides Service Information Messages (SIM) to the z/OS console for FICON device hardware failures
 - Provides CUP Diagnostics support
 - Supports IBM's z/OS Health Checker
- Since each vendor writes their own CUP code they have complete freedom to provide capabilities such as RMF Reporting and CUP Diagnostics as well as providing enriched support for IBM's zHealth Check
- It is recommended that every FICON switch should have CUP enabled



RMF 74-7 FICON Director Activity Report



28

FICON DIRECTOR ACTIVITY

								_	PAGE 1
	z/0S	V1R8		SYSTEM ID	ABCD	START (04/12/2009-04	.30.00 IN	TERVAL 000.15.00
				RPT VERSIO	ON V1R8 H	AMF END (04/12/2009-04	.45.00 Ci	CLE 1.000 SECONDS
IODF =	A2 CR-	-DATE:	03/27/2009	CR-TIME: 18	8.43.51	ACT: ACTIVATE	2		
SWITCH	DEVICE:	032B	SWITCH ID: 21	B TYPE	006140	MODEL 001 N	AN MCD PI	ANT: 01	SERIAL: 00000HIJKLMN
PORT	-CONNE(CTION-	AVG FRAME	AVG FRAM	Æ SIZE	PORT BANDWID	TH (MB/SEC)	ERROR	
ADDR	UNIT	ID	PACING	READ	WRITE	READ	WRITE	COUNT	
05	CHP-H	05		849	1436	8.63	17.34	0	
07	CHP	6B		1681	1395	50.87	10.32	0	
09	CHP	15	0	833	1429	11.96	20.49	0	
0C	CHP-H	64	0	939	1099	0.39	0.50	0	
0D	CHP	6B	0	1328	1823	3.56	12.73	0	
OF	CHP-H	66	0	1496	1675	1.85	2.61	0	
10	CHP	64	0	644	1380	0.03	0.13	0	
13	CHP-H	19	0	907	885	0.58	0.45	0	
16	CU	C800	0	1241	738	20.97	5.72	0	Eabria with
	CU	CA00				70.10	3.82	0	
1A	CHP	15	0	1144	1664	0.65	1.18	0	zHPF Enabled
1B	CHP	0D	0	510	1759	0.12	1.72	0	
1E	CHP-H	05	0	918	894	0.59	0.45	0	
1F	CHP	21	0	1243	1736	0.97	1.70	0	
20	CU	E900	0	1429	849	17.66	8.85	0	
	CU	E800							
	CU	E700							
22	CHP	10	0	923	1753	0.55	2.78	0	
23	CHP	54	0	1805	69	20.80	7.30	0	
24	CHP	64	0	89	1345	0.00	0.00	0	
27	CHP	6B		1619	82	0.01	0.00	0	
28	SWITCH	95	270	998	789	50.32	10.56	0	
2B	CHP	70	U	69	2022	0.00	0.71	0	
	Indicator of		Average	e Frame	Indicator of h	now relevant			
			cizo Dy	and Ty	Decing Dela	wand Ava			
Starvation			SIZERX	(buffer credits) Frame Size is to performance					
			(buffer					SHARE®	
							-		in Orlando 2015

Complete your session evaluations online at www.SHARE.org/Orlando-Eval

Ē

Buffer Credit Starvation

Detecting Problems such as Buffer Credit Starvation

FICON Director

FICON Switch

Produce the FICON Director Activity Report by creating the RMF 74-7 records, but this is only available when utilizing switched-FICON!

 Running Control Unit Port (CUP) per switching device allows RMF to access the CUP port – always x"FE" – to provide switch statistics back to RMF at its request

IODF = A2

PORT

ADDR

05

07

09

0C

OD

OF

13

16

17

1A

1B

1E

1F

20

22

23

24

27

28

28

SWITCH DEVICE:

UNIT

CHP-H

CHP-H

CHP-H

CHP-H

CHP

CHP

CHP

CHP

CHP

CHP

CHP

CHP

CHP

CU

CU

CU

CHP

CHP

CHP

CHP

CHP

CHP

CHP-H

2/05 V1R8

-CONNECTION-

CR-DATE:

032B

05

6B

15

64

6B

66

64

19

12

0B

15

OD

05

21

E900

E800

E700

10

54

64

6B

95

03/27/2009

WITCH ID:

AVG FRAME

PACING



FICON Path to x"FE"

z System

longer in its port but the BBC count was at zero so the frame could not be sent

INTERVAL 000.15.00

PAGE

FICON DIRECTOR ACTIVITY

ACT: ACTIVATE

-- READ ---

8.63

0.87

11.96

0.39

3.56

1.85

0.03

0.58

0.97

0.10

0.65

0.12

0.59

0.97

0.55

0.80

0.00

0.01

10.32

0.00

17.66

MODEL: 001

START 04/12/2009-04.30.00

MAN: MCD

-- WRITE --

PORT BANDWIDTH (MB/SEC)

04/12/2009-04.45.00

17.34

0.32

0.50

2.61

0.13

0.45

1.72

0.82

1.18

1.72

0.45

1.70

8.85

2.78

0.00

0.00

30.56

0.71

12.73

20.49

PLANT: 01

ERROP

COUNT

0

0

0

0

0

SYSTEM ID PRD1

CR-TIME: 16.43.51

READ

849

1681

833

939

1328

1496

644

907

1241

1144

510

918

1243

1429

923

89

1805

1619

918

69

0

0

27

685

RPT VERSION VIR8 RMF

AVG FRAME SIZE

TYPE: 006140

WRITE

1436

1395

1429

1099

1823

1675

1380

885

1738

1688

1664

1759

894

849

1753

1345

1589

2022

69

82

1736



FICON Director Activity Report With Frame Delay



Frame Delay Does Not Mean There Is A Problem

Check for C3 Discards and Frame Delay

- When experiencing performance problems that might be BBC related:
 - Correlate "time at zero buffer credits" along with Class 3 (C3) discards
 - C3 discard: a port cannot deliver a frame before the Hold Time is exceeded
 - Often as long as 500 milliseconds (ms) or half a second
 - When the hold time value is exceeded then that frame is dropped from the port by the switch port and the I/O must be re-driven – too many of these is BAD!
 - One of the reasons for C3 discards to occur is running out of buffer credits.
 - Fabric congestion could also result in this situation.
- Other CLI commands are available to provide additional information that will be useful when determining if real buffer credit issues exist. For Brocade:
 - portstatsshow
 - framelog
 - porterrshow
- So a user can pretty easily figure out if there is a potential problem (e.g. BBC at zero for longer than 2.5µs) or a really big problem (e.g. C3 discards occurring at 200-500 milliseconds) but not too much in between! That "in-between" 2.5 µs and 500 ms takes some real troubleshooting!

This Section

- Data Encoding
- Forward Error Correction (FEC)
- Compression and Encryption on ISLs



Data Encoding – Patented by IBM 8b/10b compared to 64b/66b





10Gbps and 16Gbps will always use 64b/66b data encoding



Forward Error Correction (FEC)

Enhances transmission reliability and thus performance

- 16 Gbps devices can recover from bit errors on a 10 Gb or a 16 Gb data stream (64b/66b encoding)
- Corrects up to 11 bit errors per each 2,112 <u>bits</u> in a frame transmission:
 - Makes as many as 11 bit corrections per each 264 bytes of frame data
- Almost like adding additional power (db) to a link
- Each vendor of 16 Gbps devices can decide how they want, or do not want, to support FEC on their interfaces







Forward Error Correction (FEC)

For 16G ISLs, FICON Express16S (FX16S) CHPIDS and 16G DASD



- FEC allows FICON channels to operate at higher speeds, over longer distances, with reduced power and higher throughput
- Those paths retain the same reliability and robustness that FICON has historically been known to provide at slower data rates
- FEC is for ISLs, 16G channels and also for IBM DS8870 DASD currently:
 EMC, HDS and Oracle have yet to announce their plans to support 16G FEC

Forward Error Correction (FEC)

And FICON Express16S CHPIDS and IBM DS8870



Here is a of diagram to show how FEC is deployed with z13 and 16G switching.

z13 using FX16S



Compression and Encryption on ISL Links



Better site-to-site security across ISLs – for In-Flight Data Only

This is in-flight switch-to-switch encryption and compression

- Secure Transfers with Encryption:
 - Encrypts data, at any speed, for 16 Gbps
 FC ports on Gen 5 Directors:
 - Uses AES-GCM algorithm for both authentication and encryption
 - Uses 256-bit encryption key
- Compression for Network Efficiency:
 - 2:1 compression on ISL data
 - Uses LZO algorithm
- When compression is enabled, the number of buffer credits required will at least double on that ISL link
- Requires no license and can be used together or separately or not at all









My Next Presentation:

z13 and Brocade Resilient, Intelligent and Synergistic I/O Processing



Wednesday August 12, 2015 - 4:30pm to 5:30pm -- Session 17503

Please Fill Out Your Online Evaluations!!

our

Eval

Thank You For Attending Today!

This was session:

17506

And Please Indicate in the evaluation if there are other presentations you would like to see us present in this track at SHARE!



My Reaction!





QR Code

