



Buffer-to-Buffer Credits, Exchanges, and Urban Legends

Lou Ricci, EMC

Howard L. Johnson, Brocade

14 August 2013 (8:00am – 9:00am)

Session 14281

QR Code



Legal Stuff

- Notice
 - IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not give you any license to these patents. You can send license inquiries, in writing to: *IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785 U.S.A.*
 - Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.
- Trademarks
 - The following terms are trademarks of the International Business Machines Corporation in the United States, other countries, or both: FICON® IBM® Redbooks™ System z10™ z/OS® zSeries® z10™
 - Other Company, product, or service names may be trademarks or service marks of others.

Abstract

- Performance in a FICON network is influenced by the underlying flow control mechanisms of Fibre Channel. In this session, we examine how Buffer-to-Buffer credits flow from the channel to the control unit. We also look at how exchanges are used in FICON applications and how they change with the introduction of zHPF. During both examinations, we explore the role of the FICON Director in managing Buffer-to-Buffer credits and exchanges over a cascaded network. Throughout the session, we debunk the various FICON “Urban Legends” featuring credits and exchanges. Take the opportunity to learn from two of the FICON industry’s leading experts in channel and fabric development and join our session.

Agenda

- Buffer Credits
 - What are they and how do they work?
 - How do you fill the pipe?
 - What if you can't fill the pipe?
 - What's wrong with multiple senders and one receiver?
 - What happens when the pipes are different sizes?
 - What's it like in the real world?
- Exchanges
 - What are they and how do they work?
 - What's an Exchange?
 - How many exchanges are needed?
 - Can they be "reused?"
 - Can you have too many exchanges?

What are they and how do they work?

BUFFER CREDITS

What is Buffer-to-Buffer Credit?



- The greater the BB Credit....
 - A. The faster frames can be sent
 - B. The farther apart the two ports can be
 - C. The larger the frames can be
 - D. None of the above

What is Buffer-to-Buffer Credit?

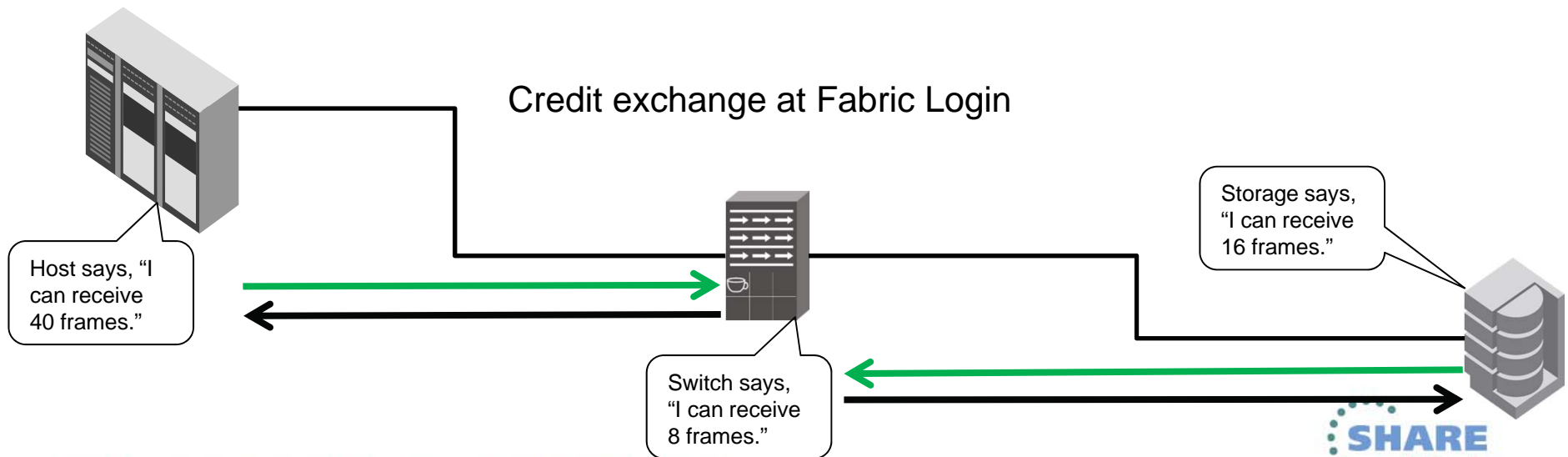
- The greater the BB Credit....
 - A.
 - B. The farther apart the two ports can be
 - C.
 - D.



Flow Control

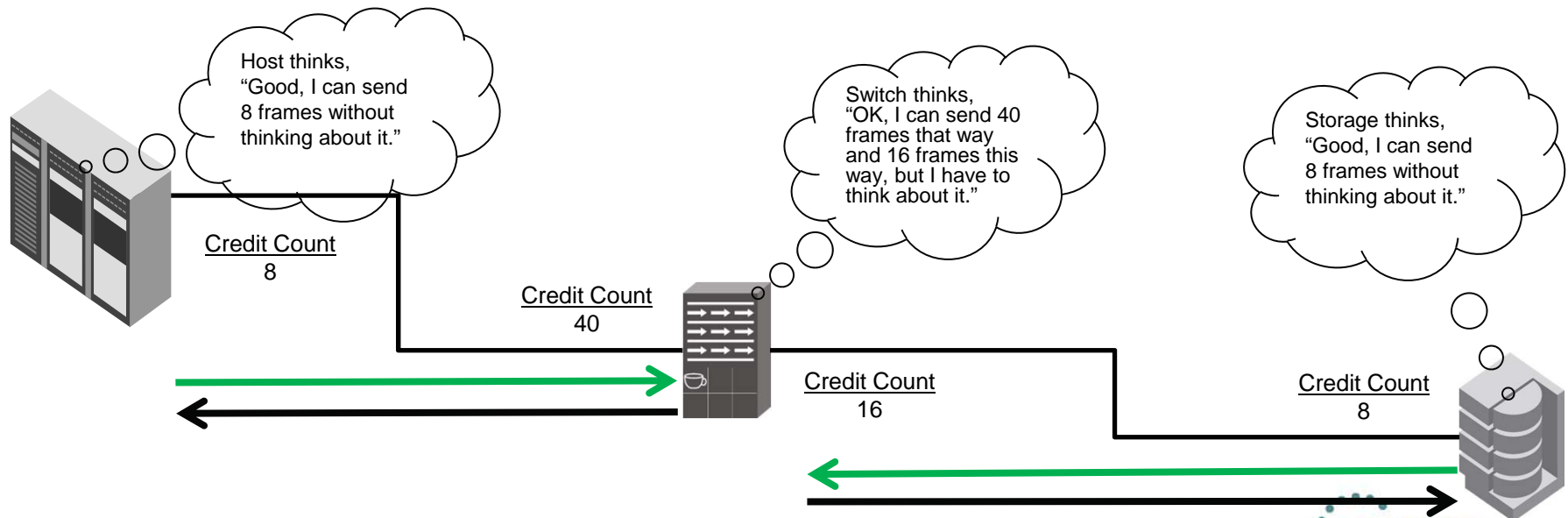
- Related to the devices' ability to receive and process frames
- Manages when frames are coming faster than they can be processed
- Dropped frames occur when frames are arriving too fast to be processed

- Frames can only be transmitted when the receiver is ready
- Credit establishment communicates the number of frames a device can receive at a time
- The credit value is exchanged at login
- Transmission stops when credit runs out
- The receiver indicates when it is ready to receive more frames



Buffer Credit

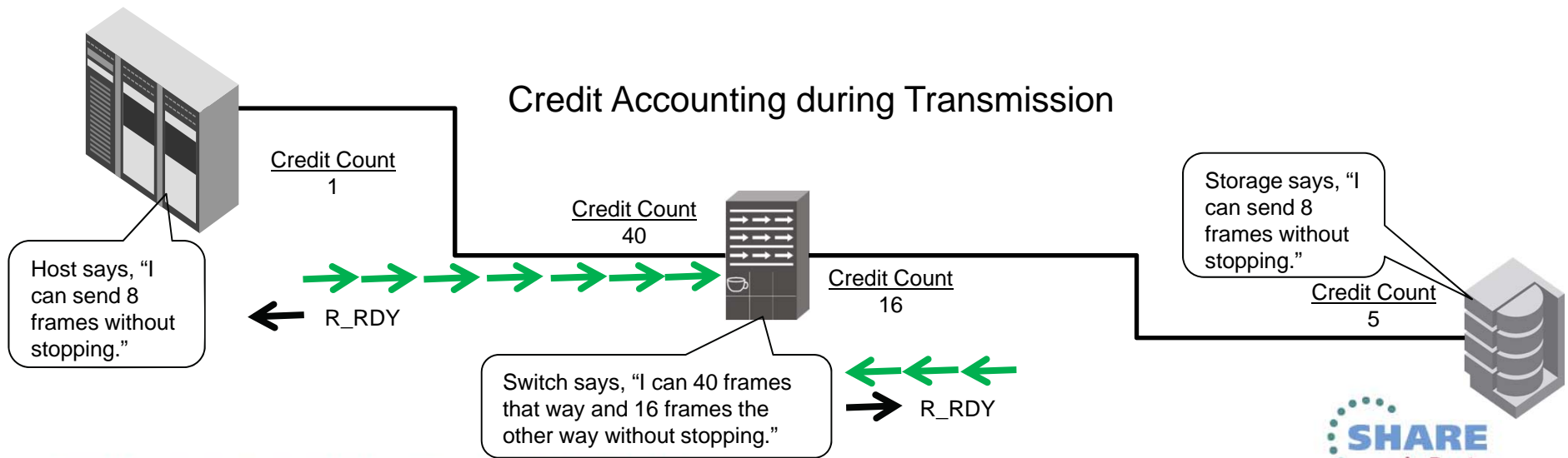
- At initialization, the two ports establish credit
 - Each buffer credit corresponds to a frame (regardless of size)
- Each side can support different values
 - Credit Count
- If a port doesn't have credit, it can't send a frame
 - Credit Count has reached zero
- Mechanism limits frame drops



Credit accounting after Fabric Login

Receiver Ready (R_RDY)

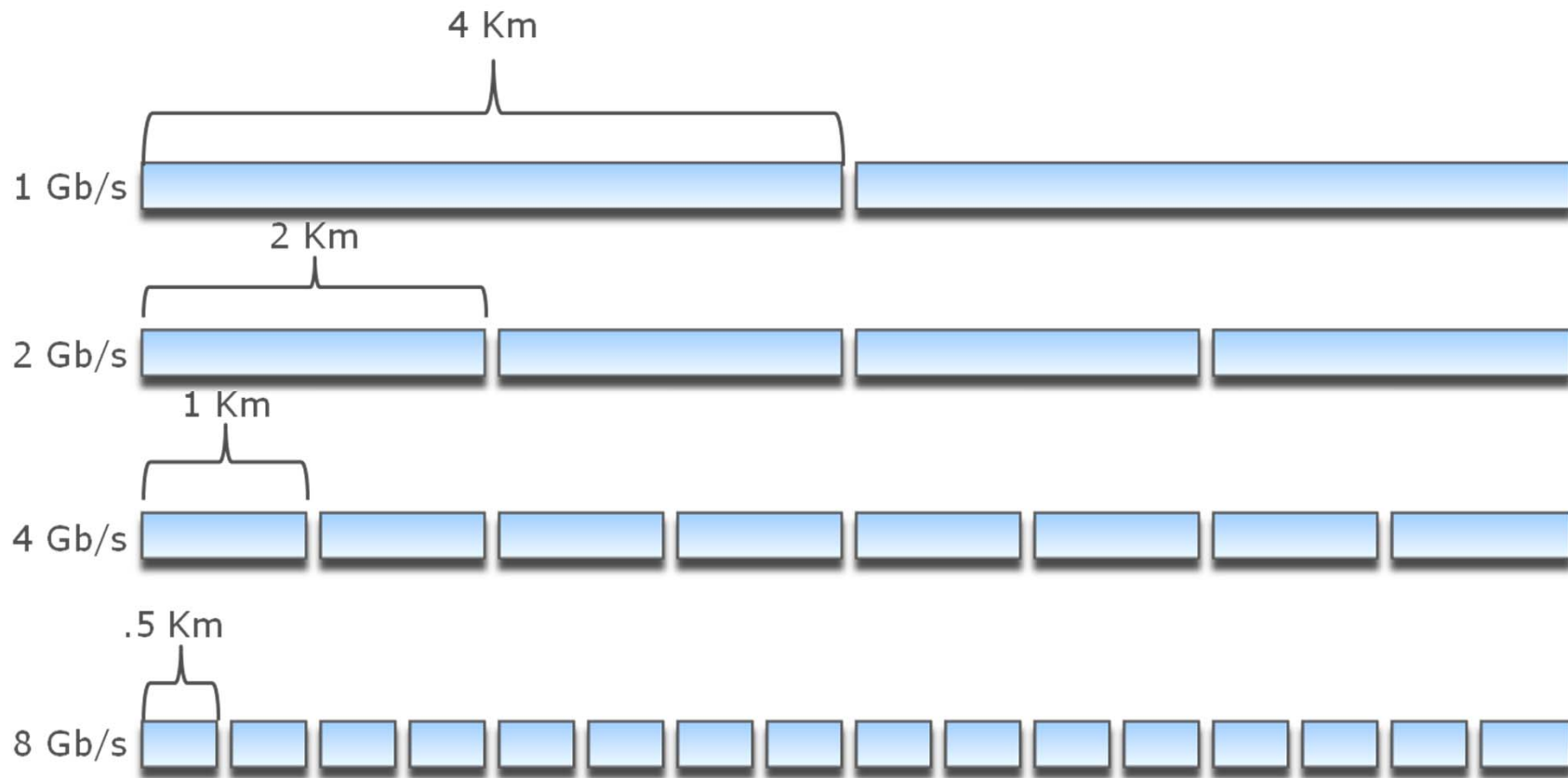
- R_RDY
 - Used for link level flow control
 - Called buffer-to-buffer credit (BB Credit)
- R_RDY is not a frame
 - It is a “primitive” so it doesn’t consume a buffer
- Frame transmission
 - BB Credit is decremented
 - Once for each frame transmitted
 - When BB Credit = 0
 - Transmission stops
- Frame received
 - R_RDY is sent
 - Causes transmitter to increment BB Credit



Urban Legend: Buffer Credits at Zero are a Problem

- Buffer credit determines DISTANCE
 - The distance two nodes can be apart and still maintain full link frame rate
- Buffer credit is the number of FRAME buffers
 - A port provides for it's NEAREST neighbor for RECEIVING frames
 - Does NOT have to be symmetrical
- Buffer credit is a FRAME count
 - Not a data SIZE
 - A 1 byte frame consumes 1 buffer credit
 - A 2K byte frame consumes 1 buffer credit
- Number of credits needed is determined by:
 - Raw Link Speed
 - Speed of light thru a fiber
 - Distance between two adjacent nodes

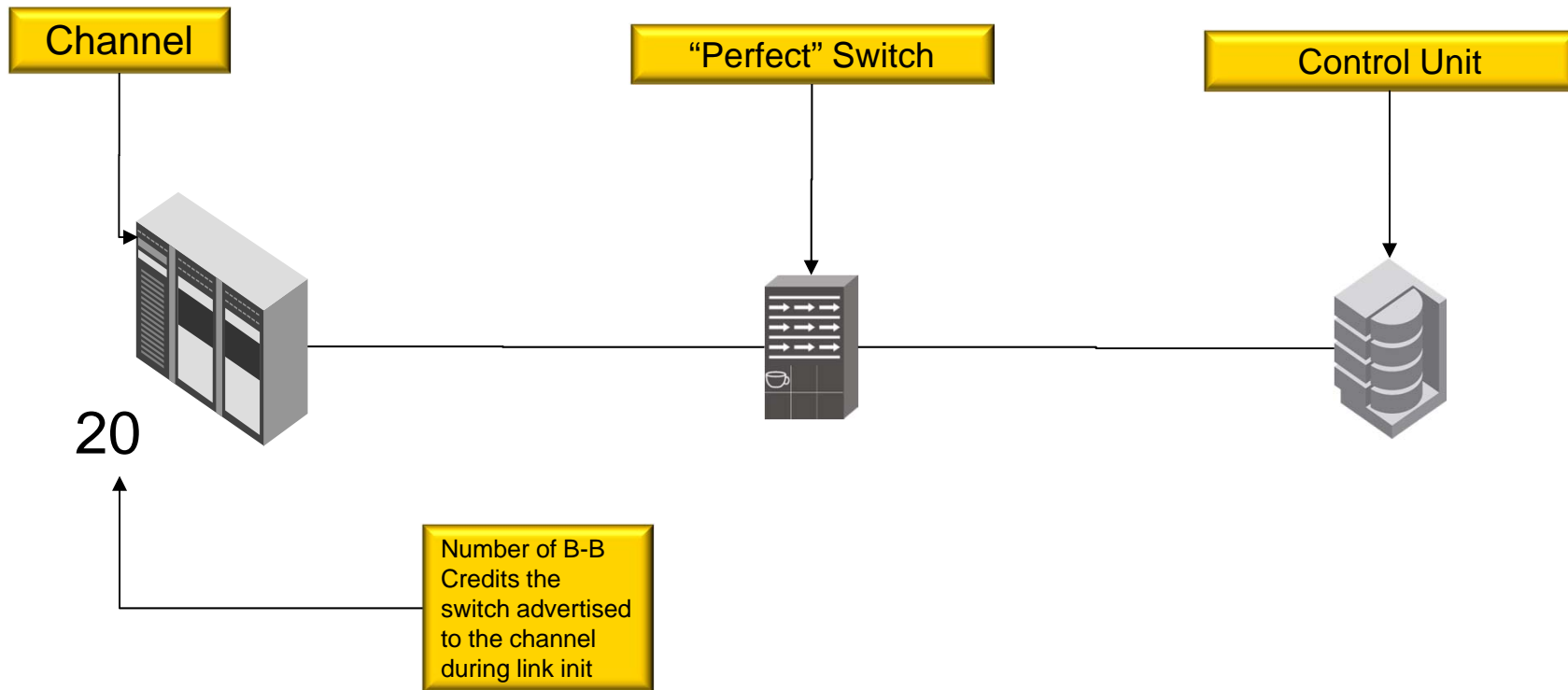
How Long is a Fibre Channel Frame?



Example: A full pipe

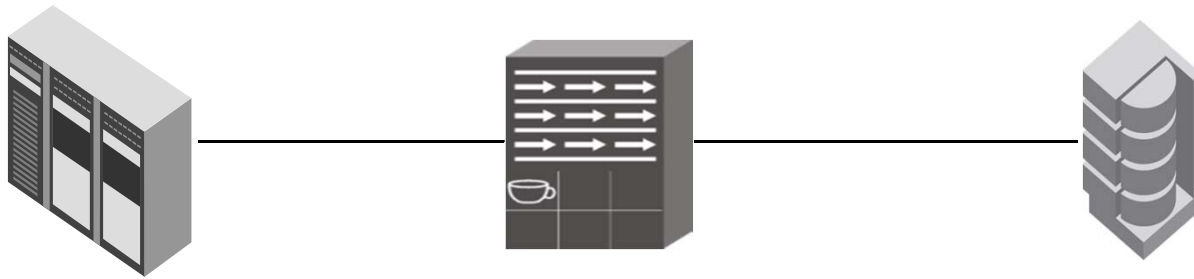
BUFFER CREDITS

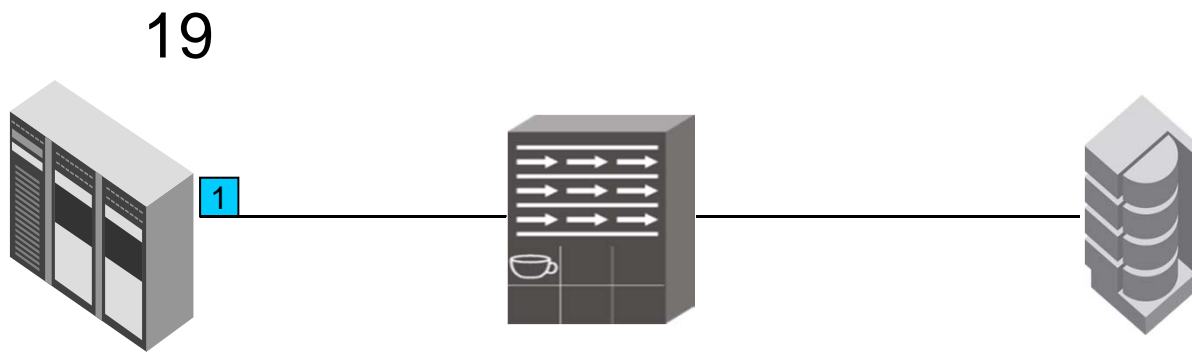
Initial Conditions



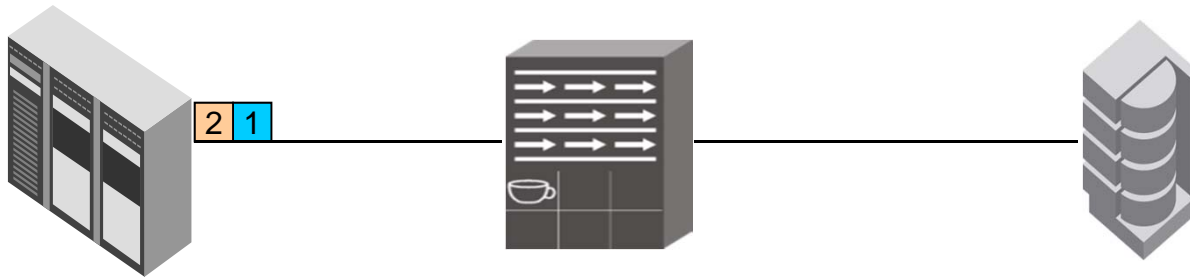
NOTE: In these animations, both the frames and R_RDY's are numbered. This is for illustrative purposes only. In reality, neither the frames nor the R_RDY's are numbered. The arrival of an R_RDY only informs the receiver that **A** frame has been forwarded, not **WHICH** frame has been forwarded.

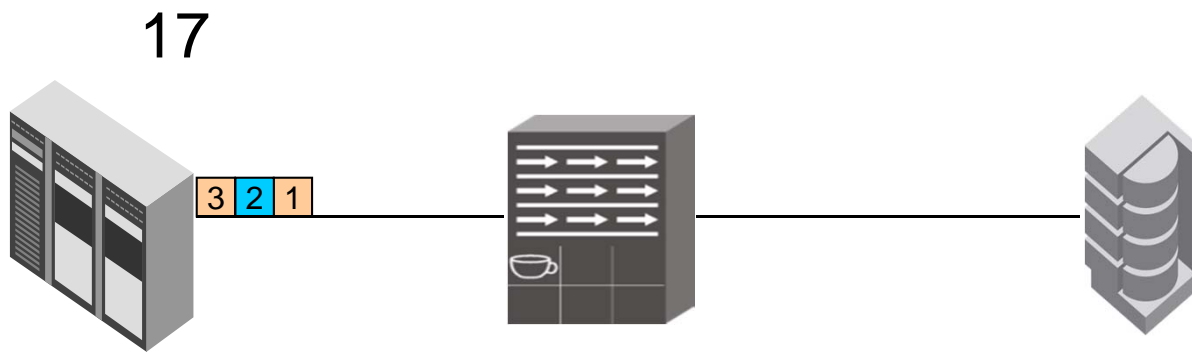
20

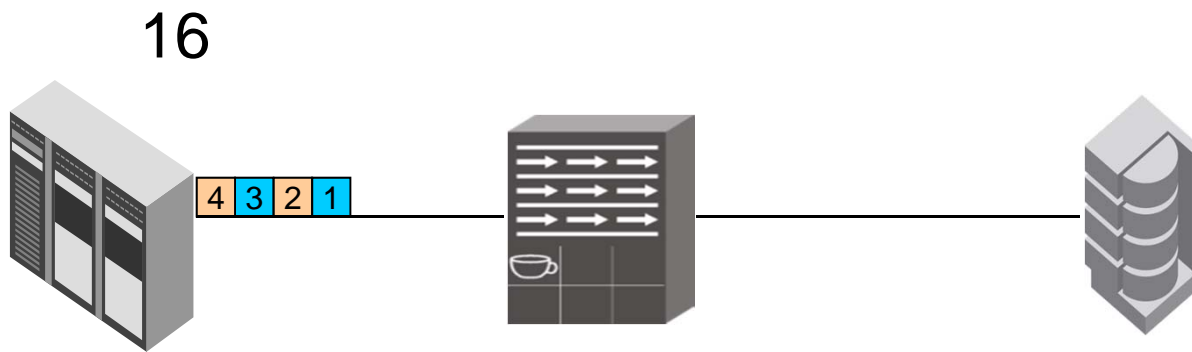


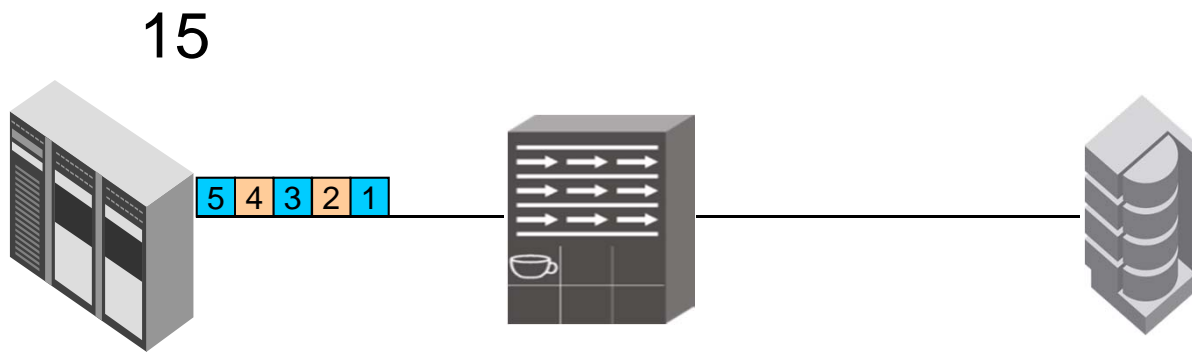


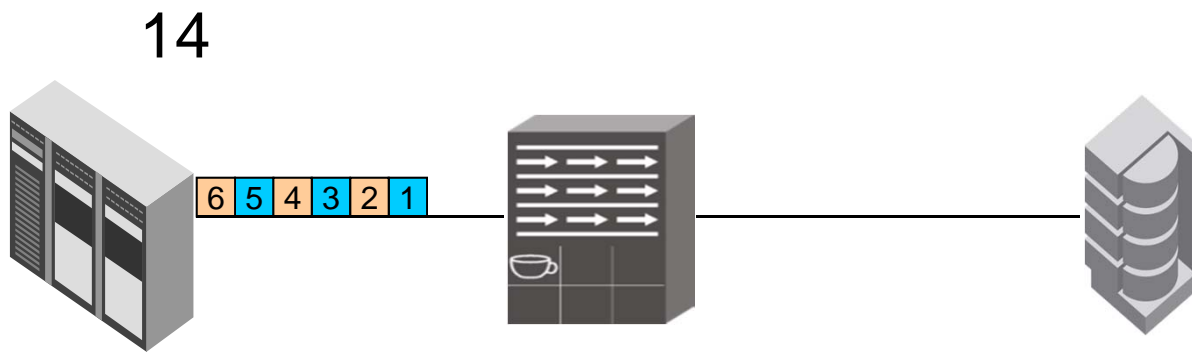
18

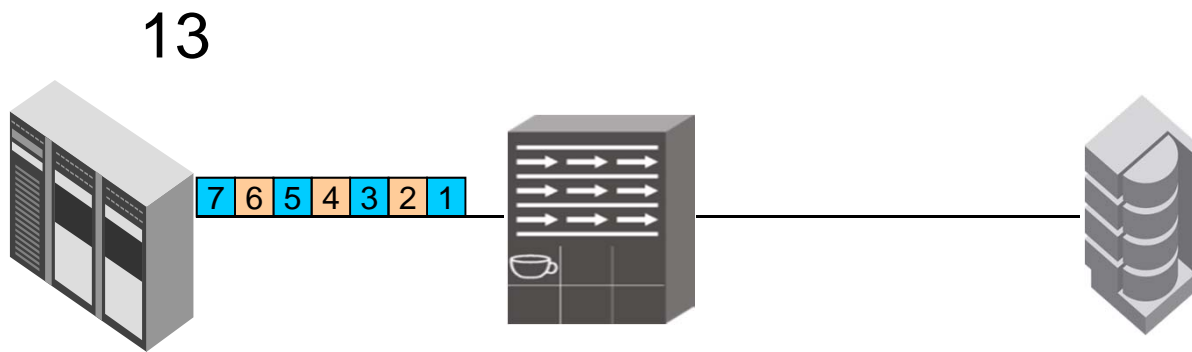


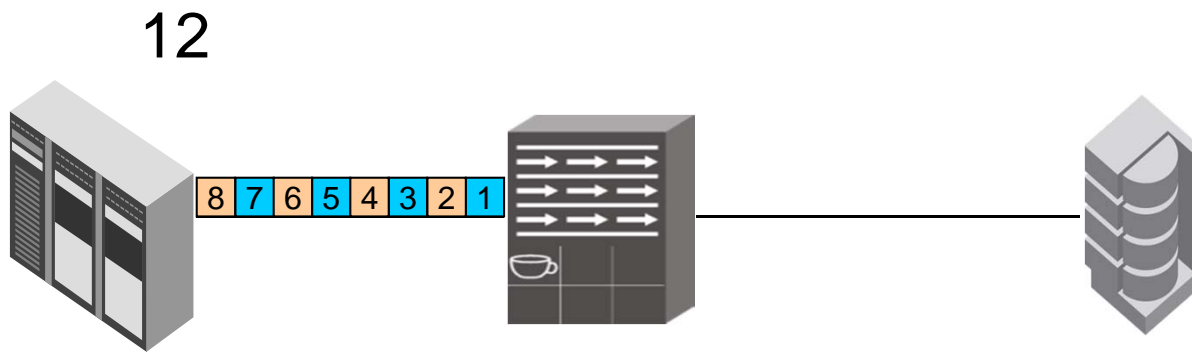


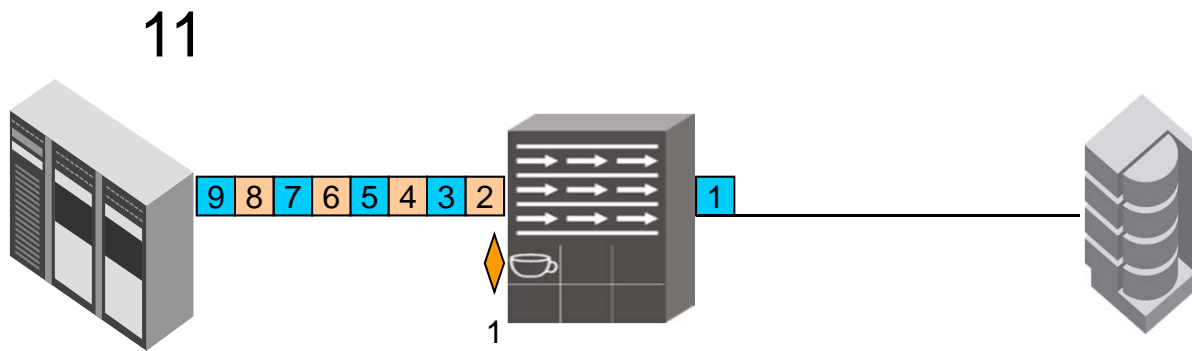


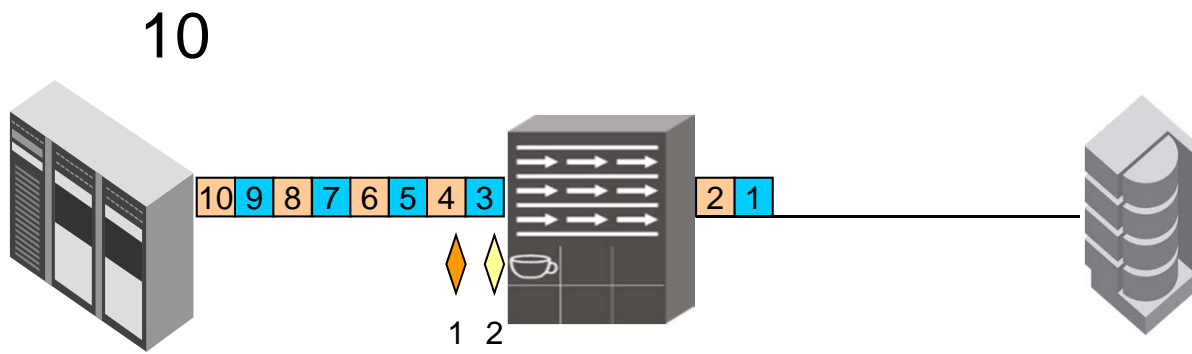


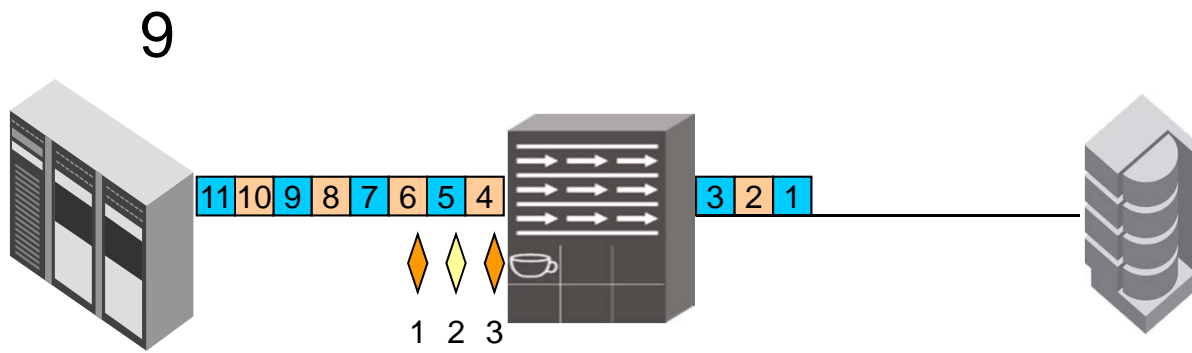


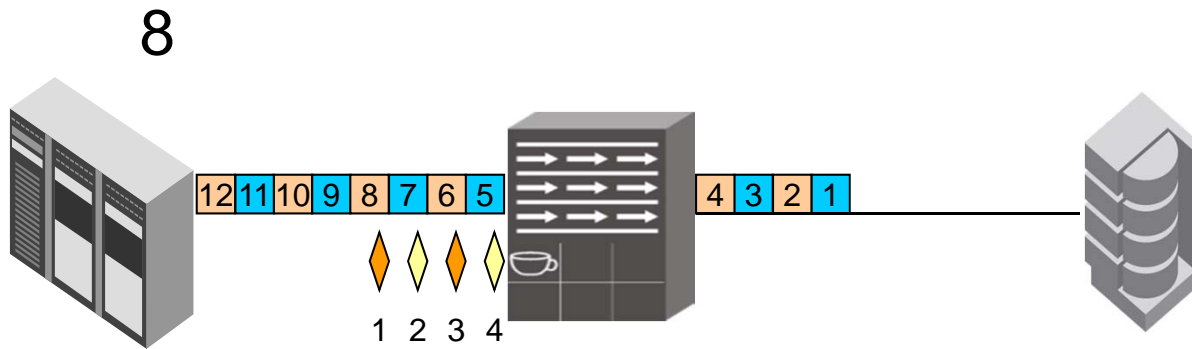


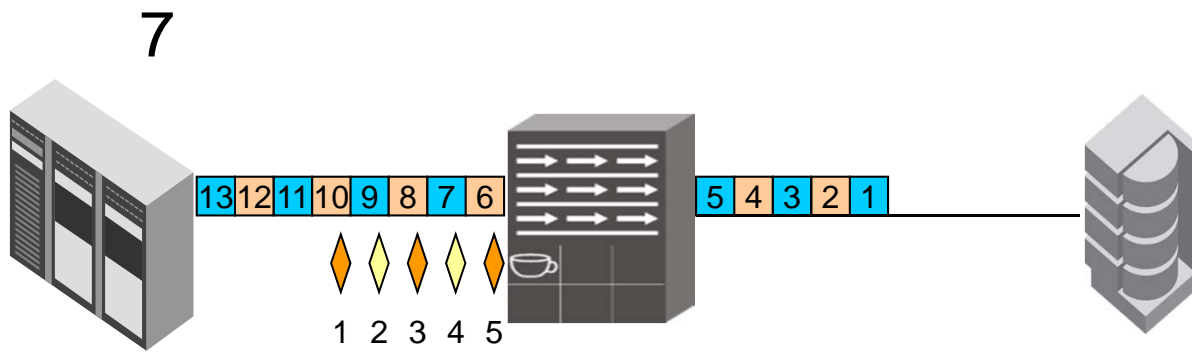


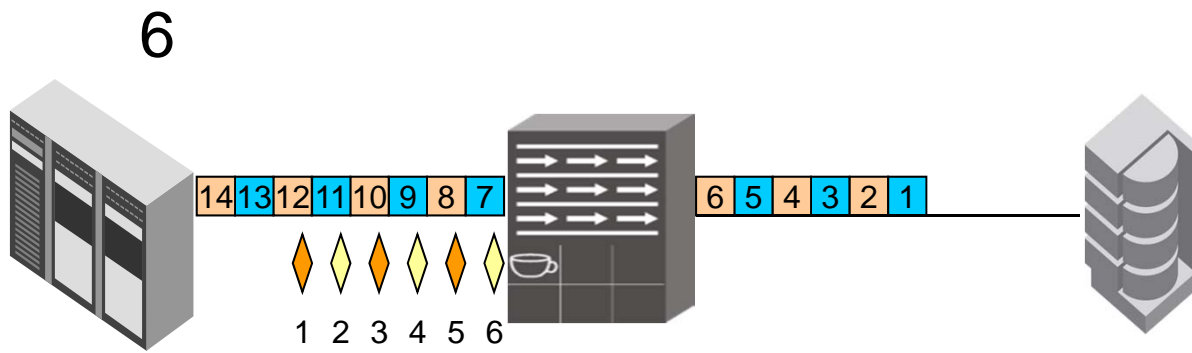


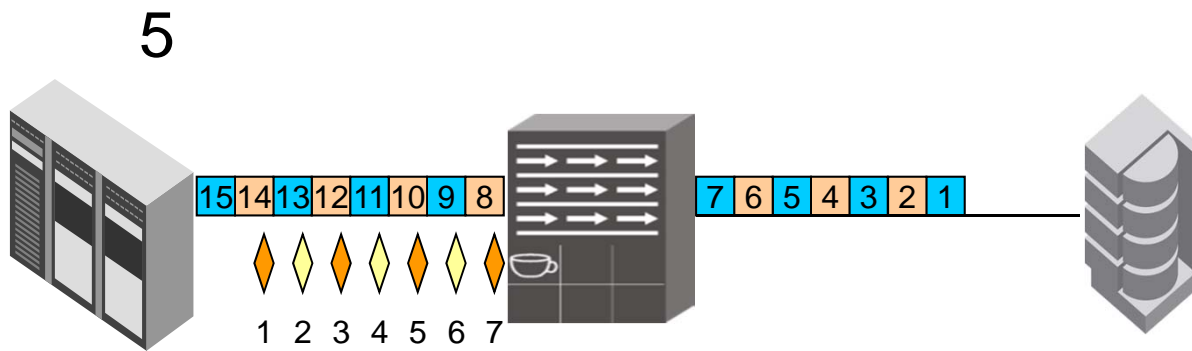


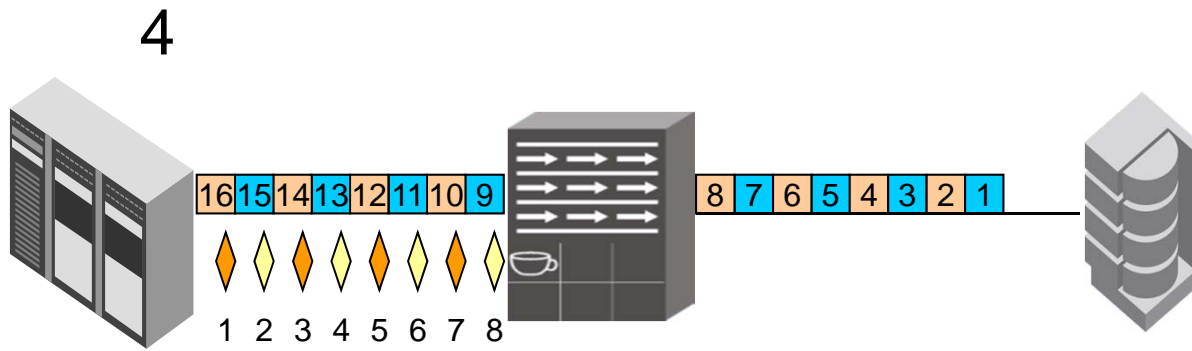


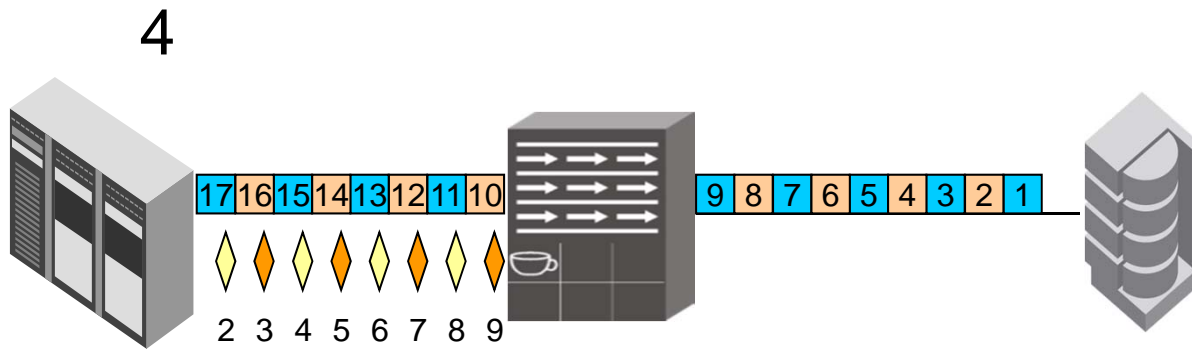


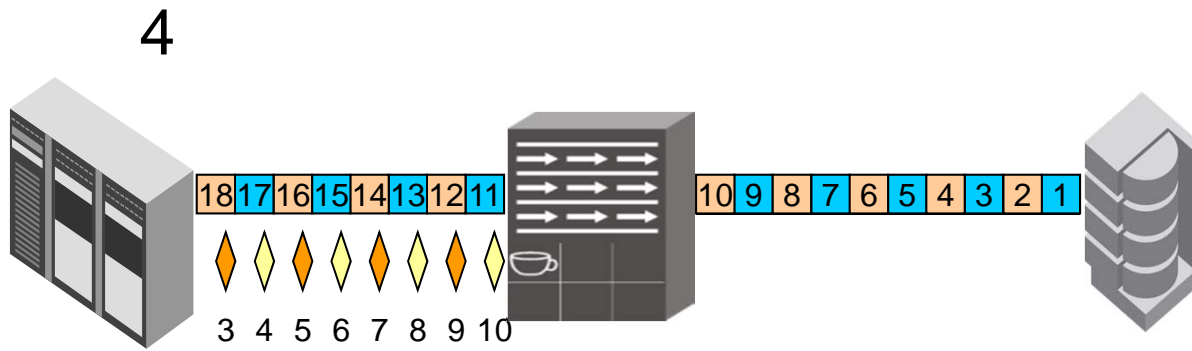


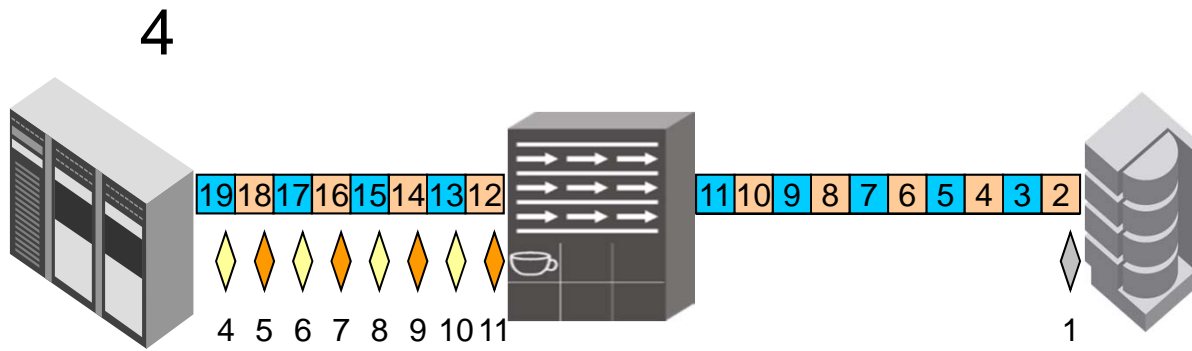










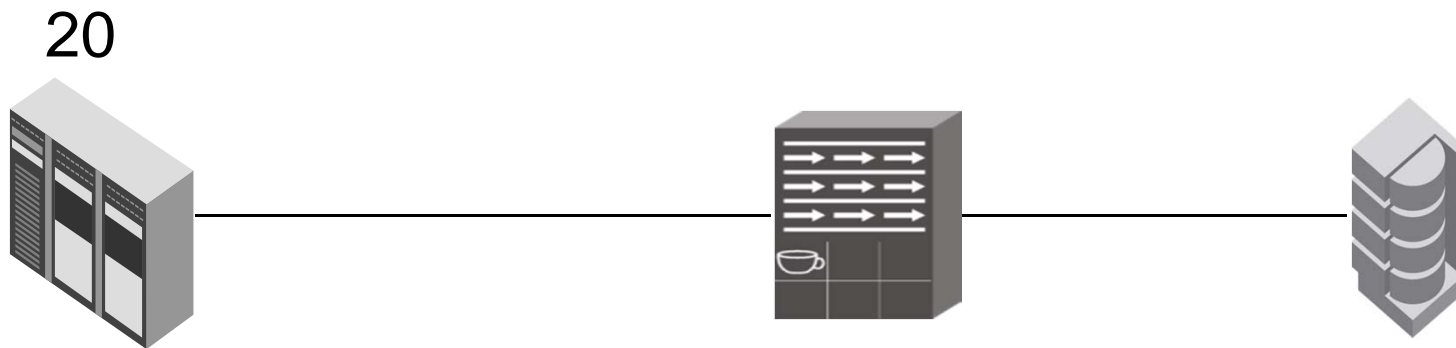


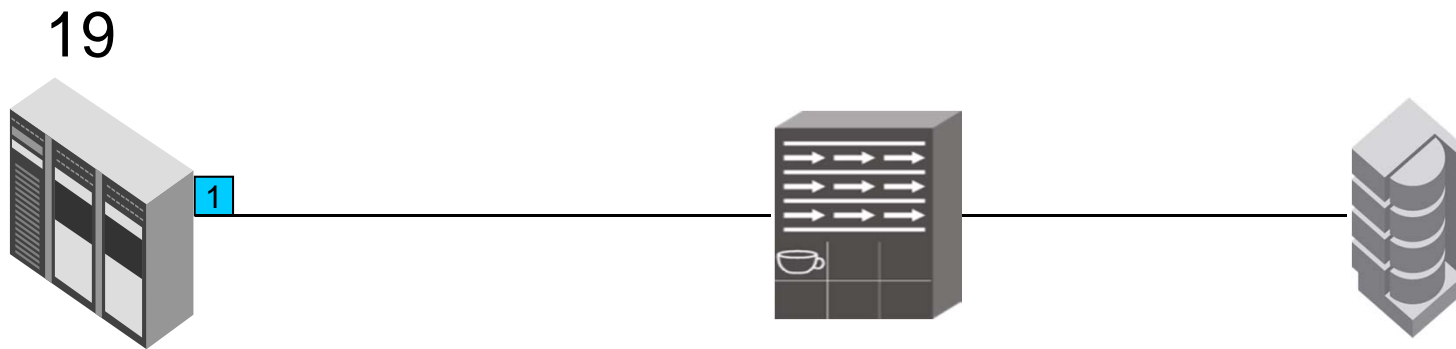
THIS PAGE INTENTIONALLY
LEFT BLANK

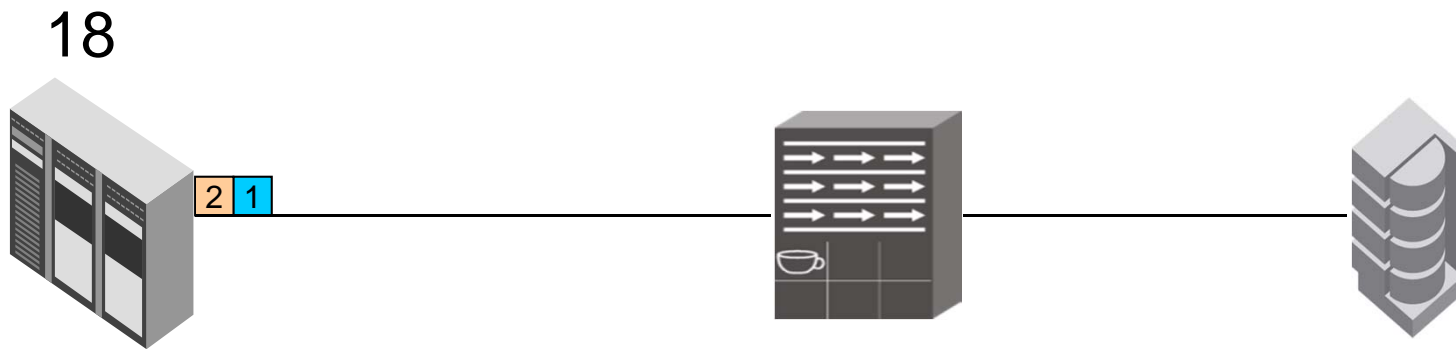
Example: A not so full pipe

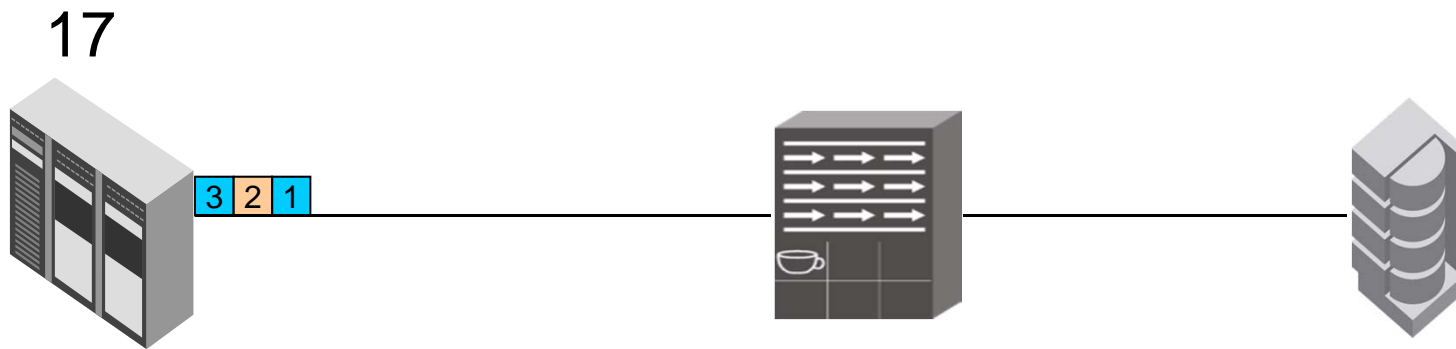
BUFFER CREDITS

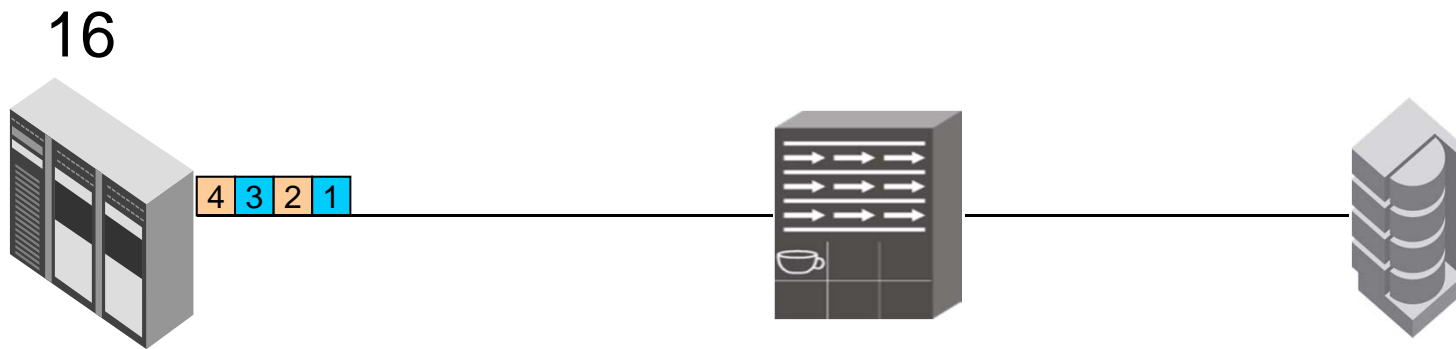
Suppose the switch is too far away from the channel for the B-B credit it advertised to the channel

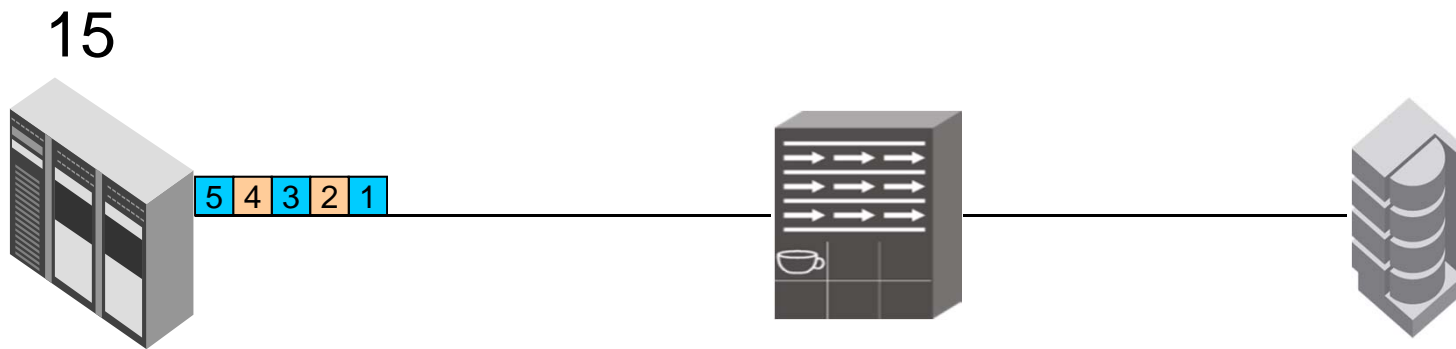


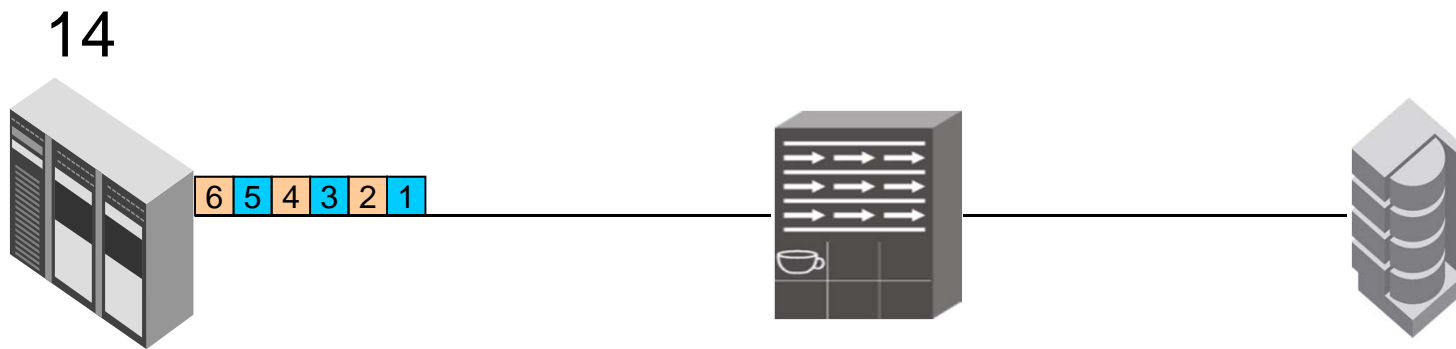


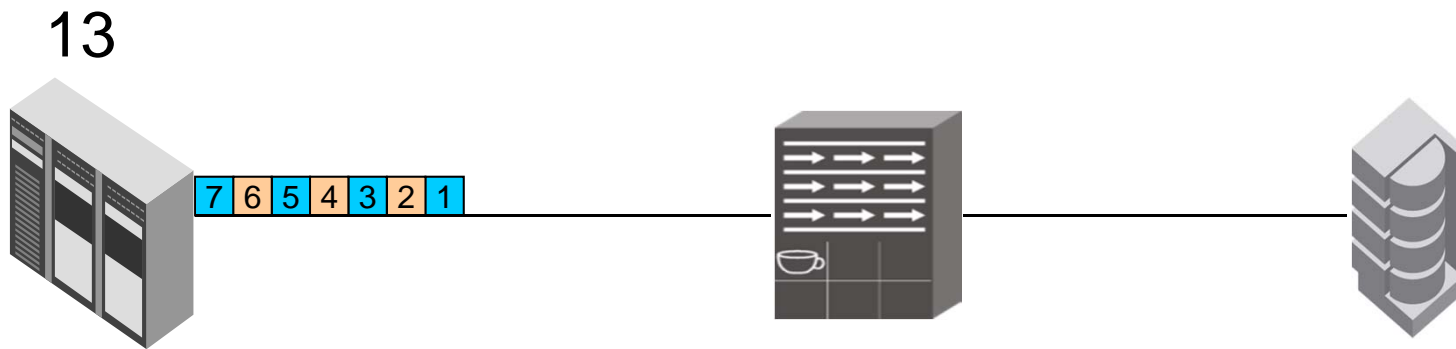


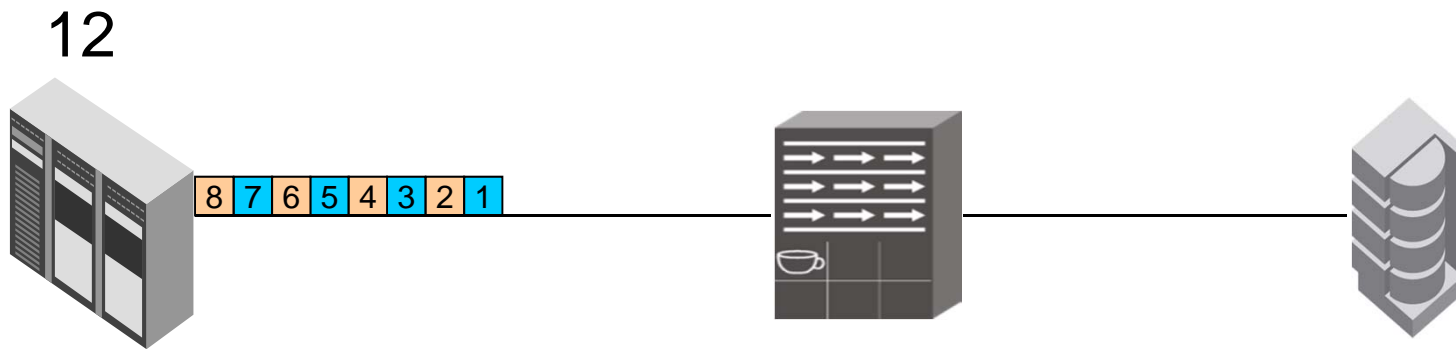


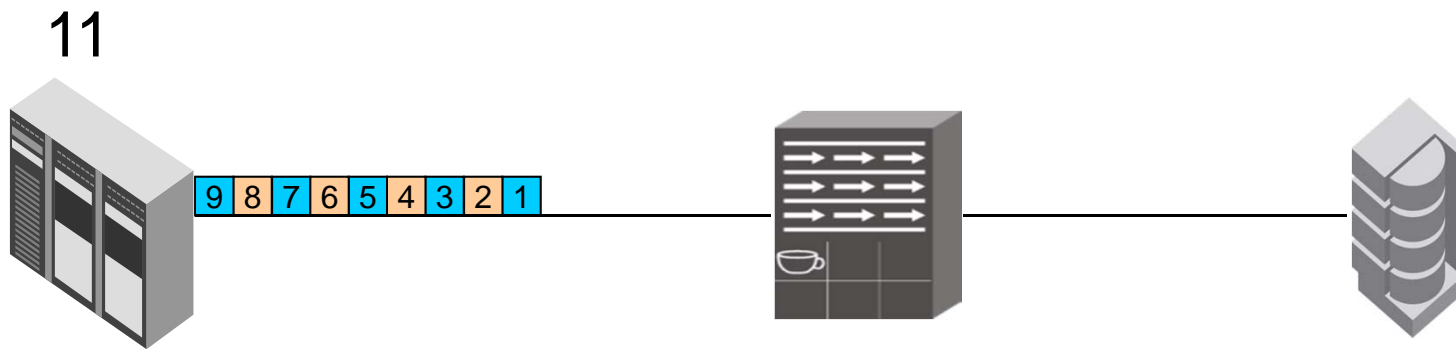


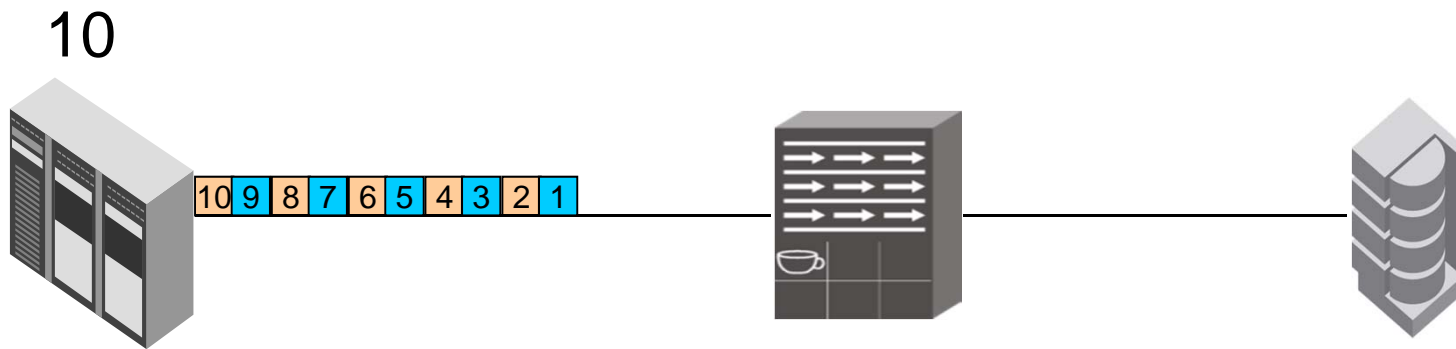


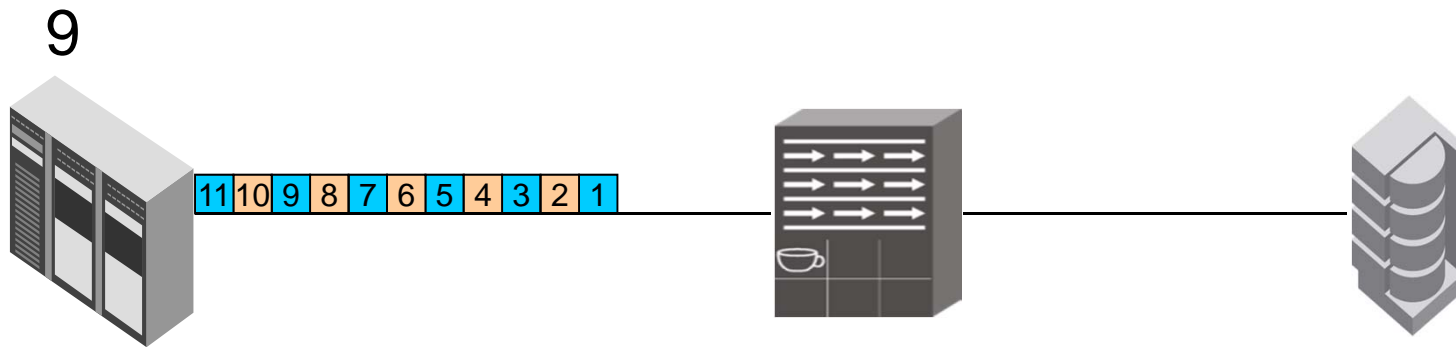


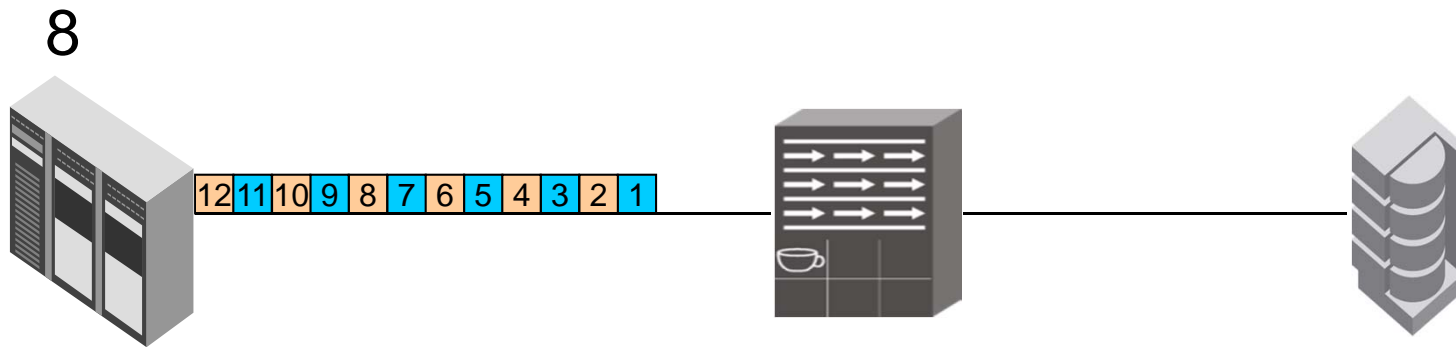


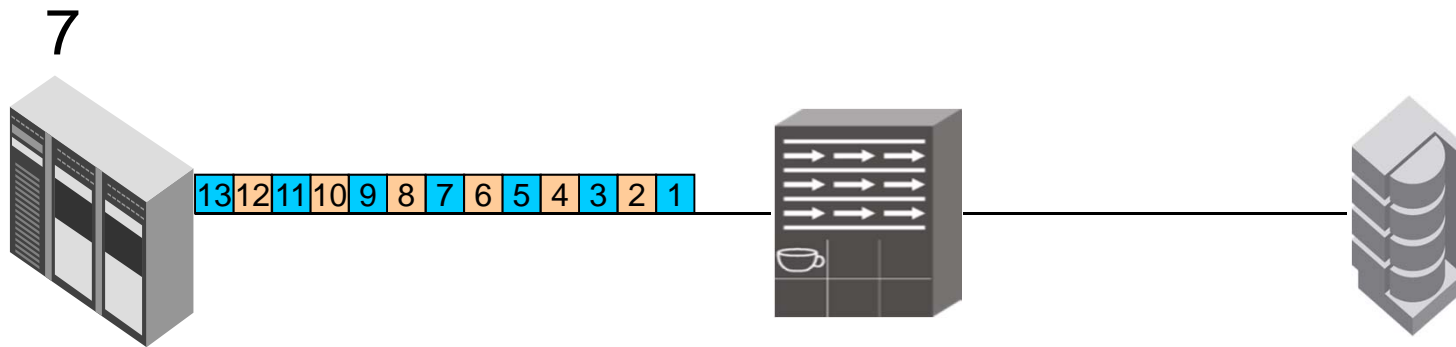


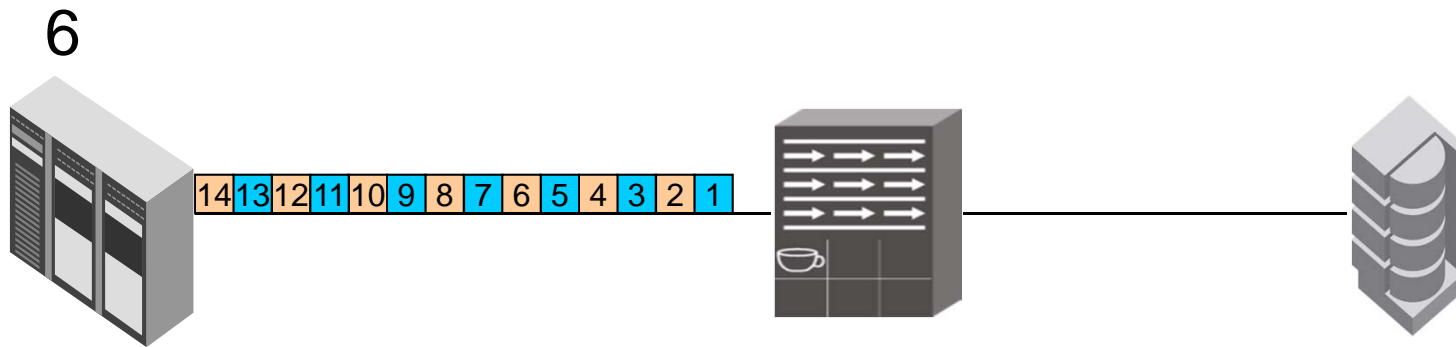


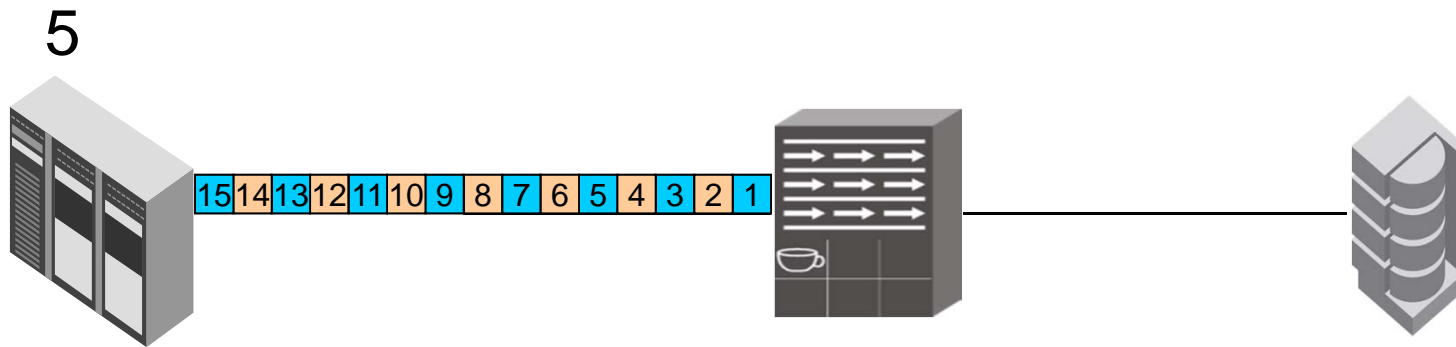


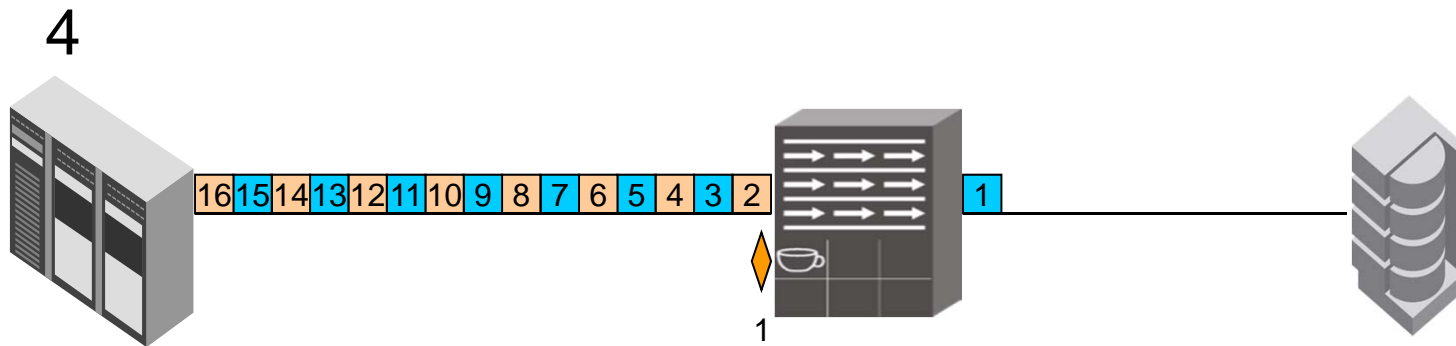


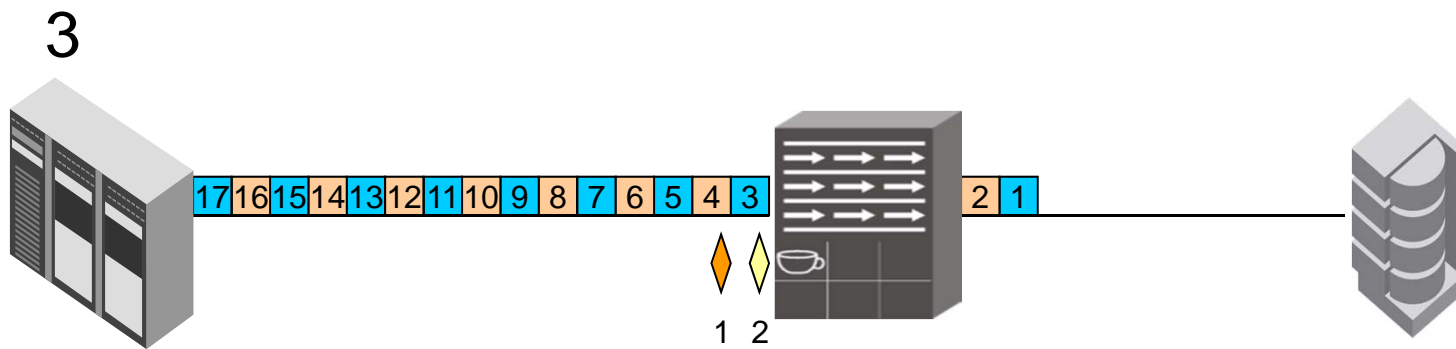


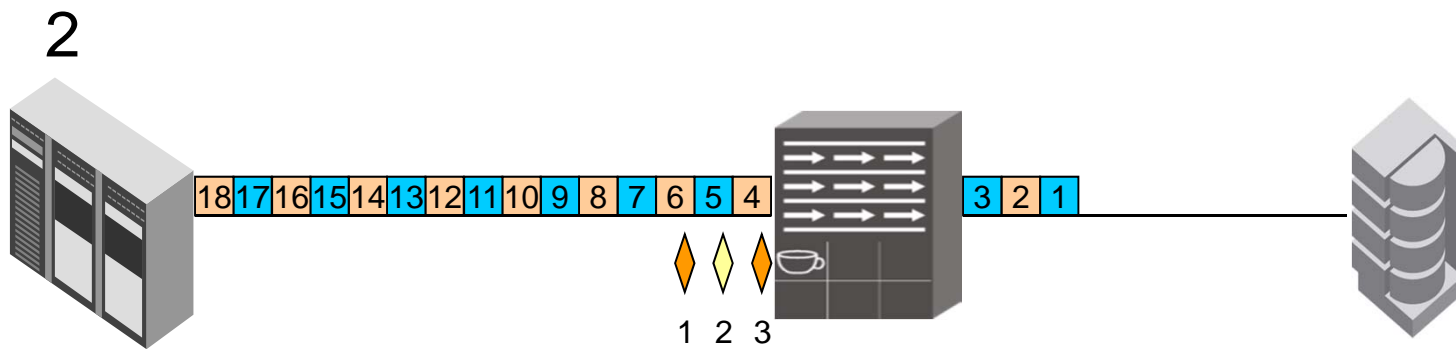


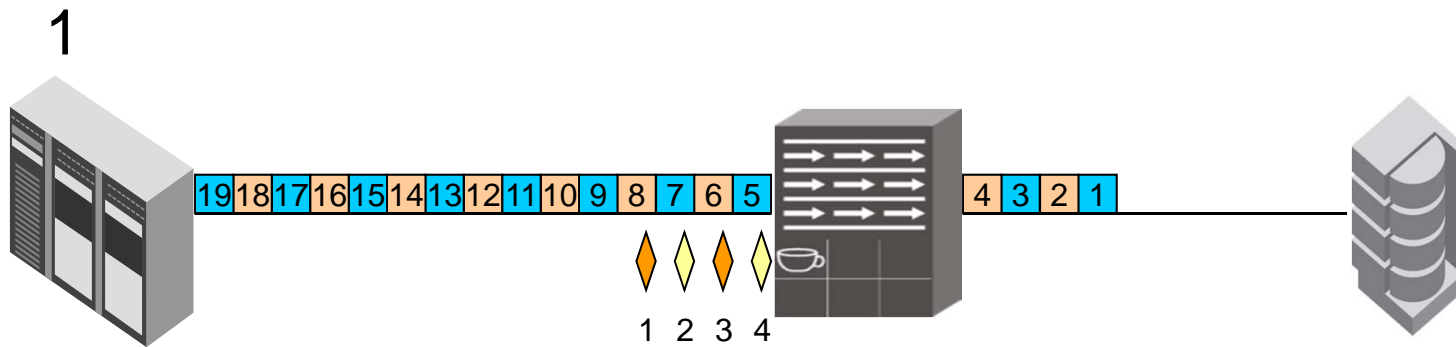


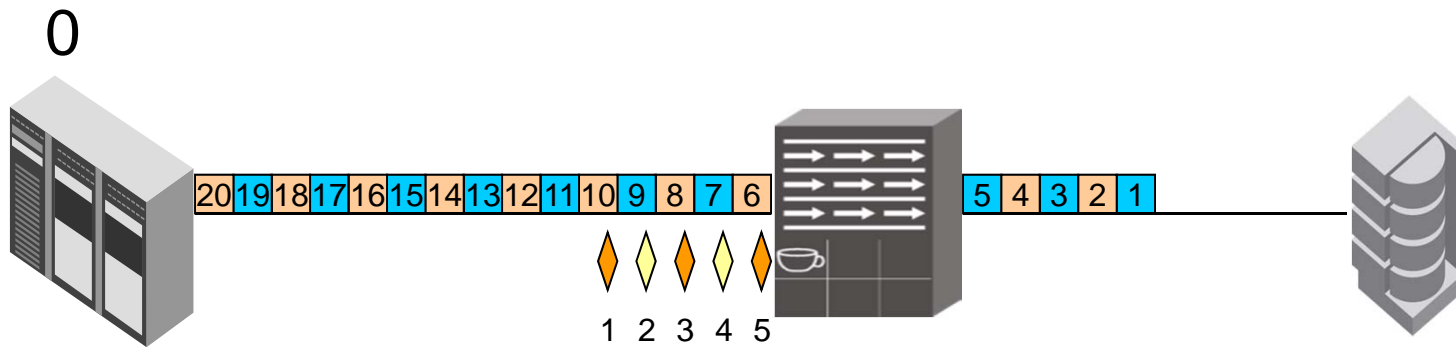


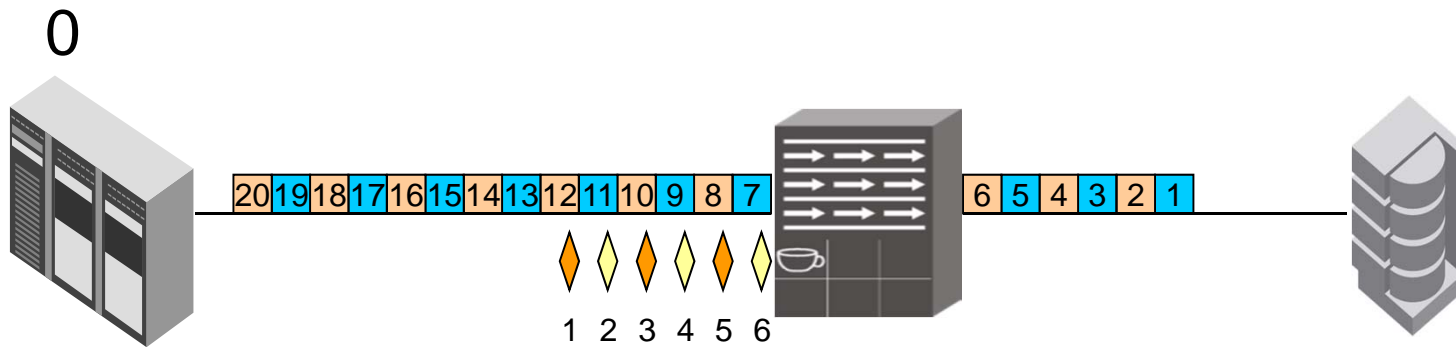


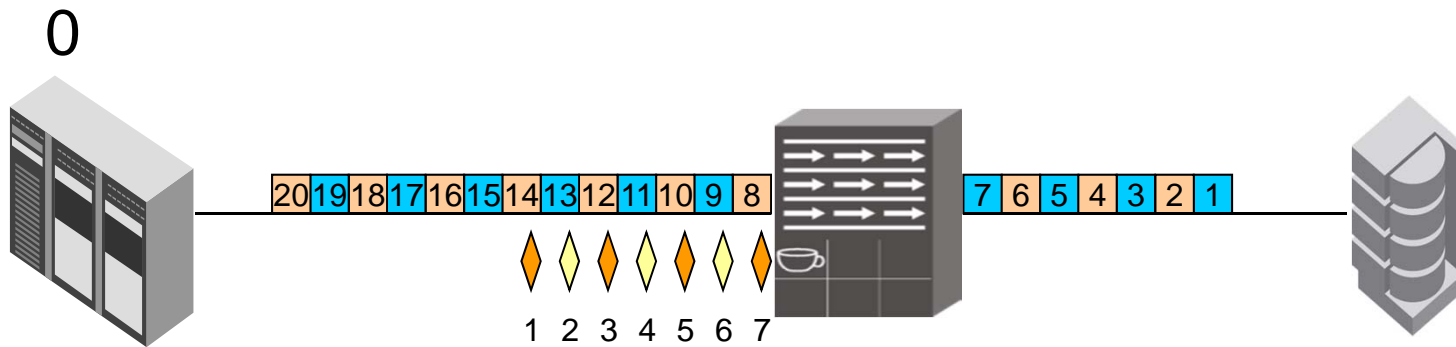


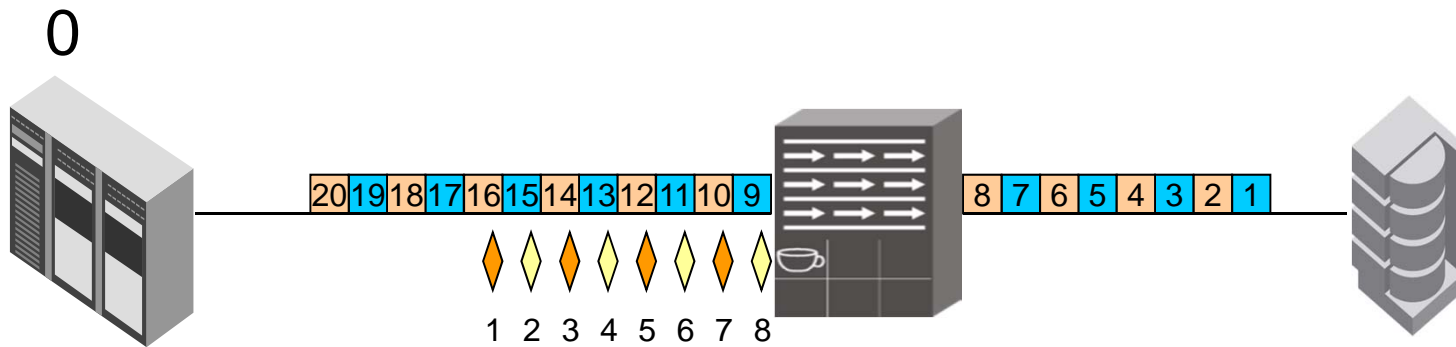


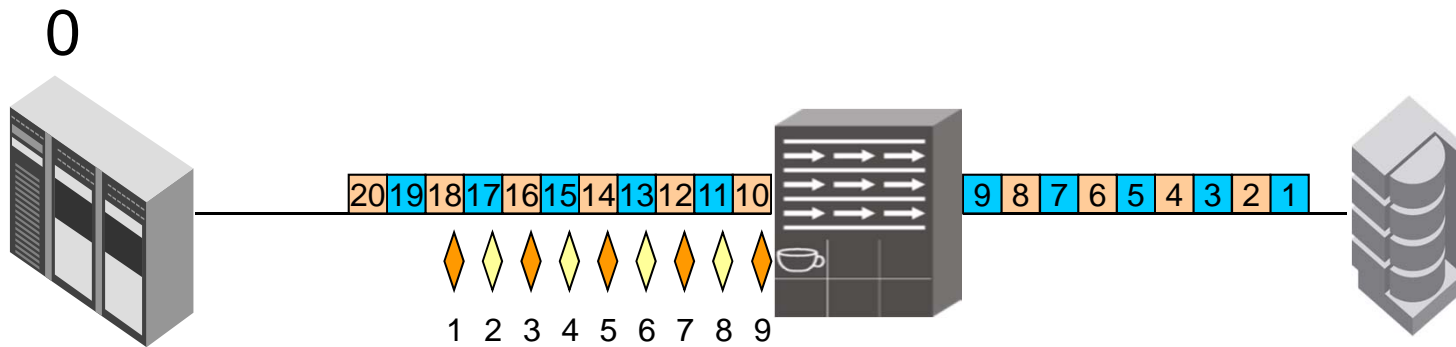


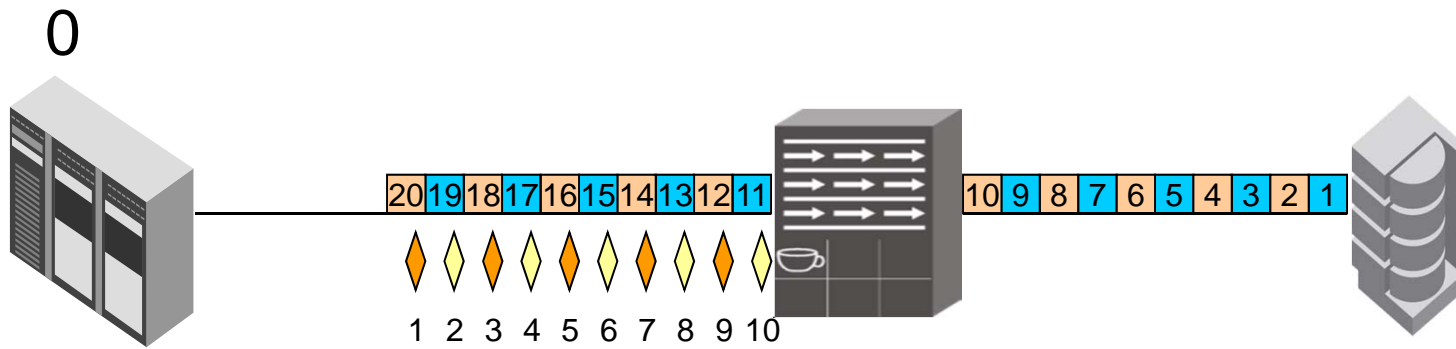


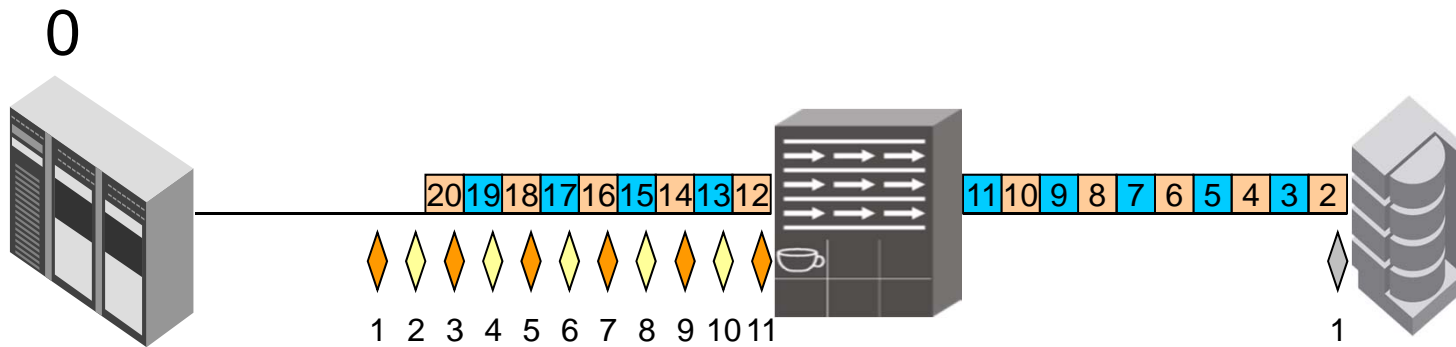


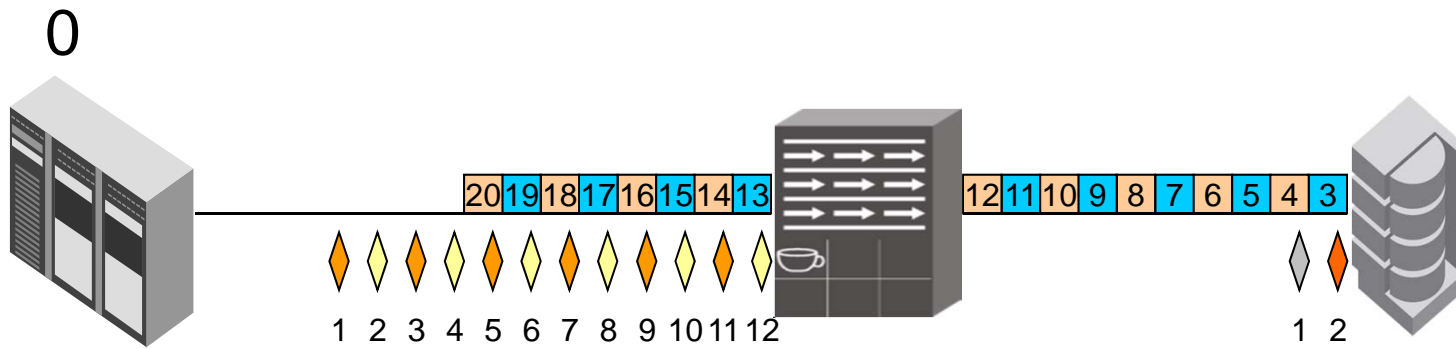


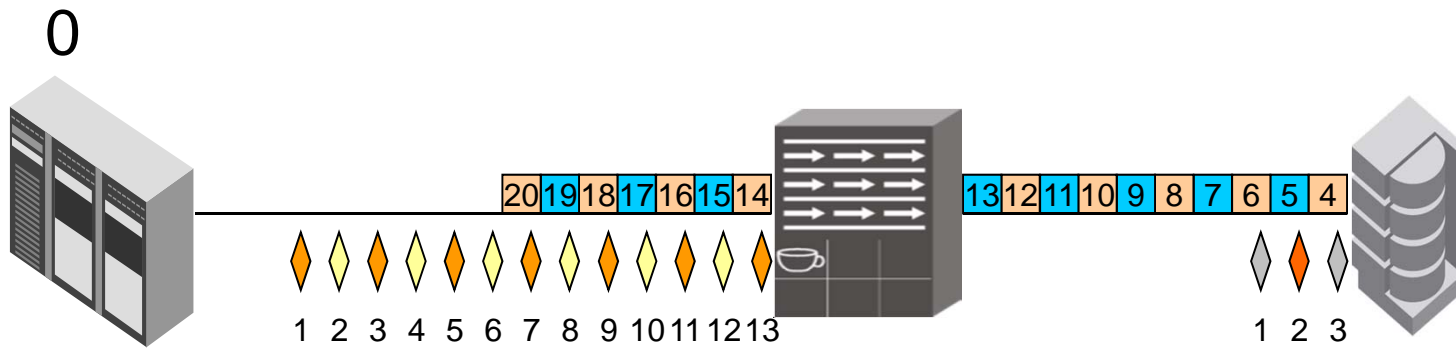


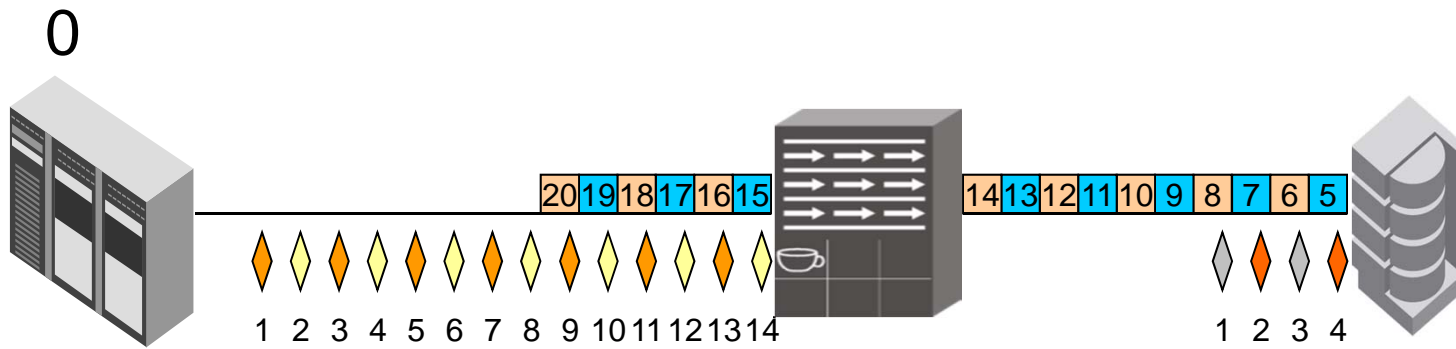


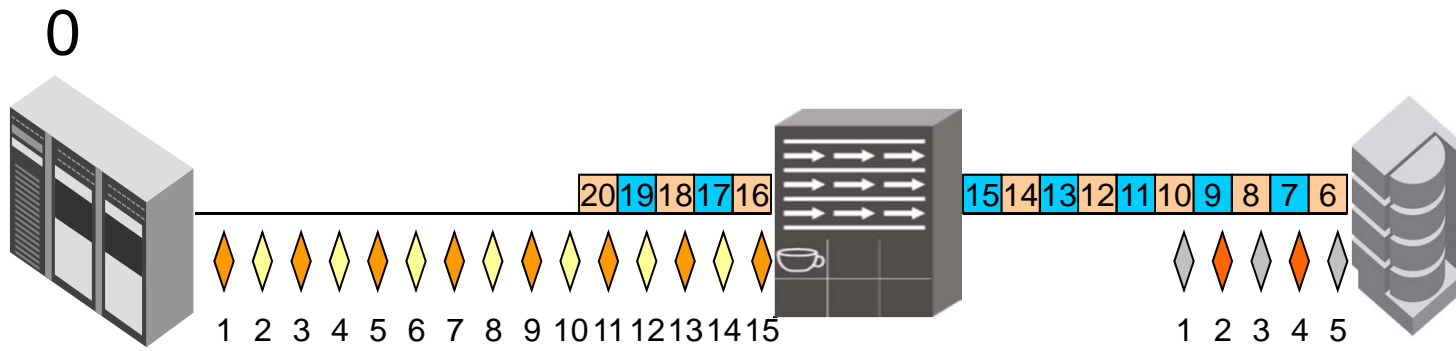


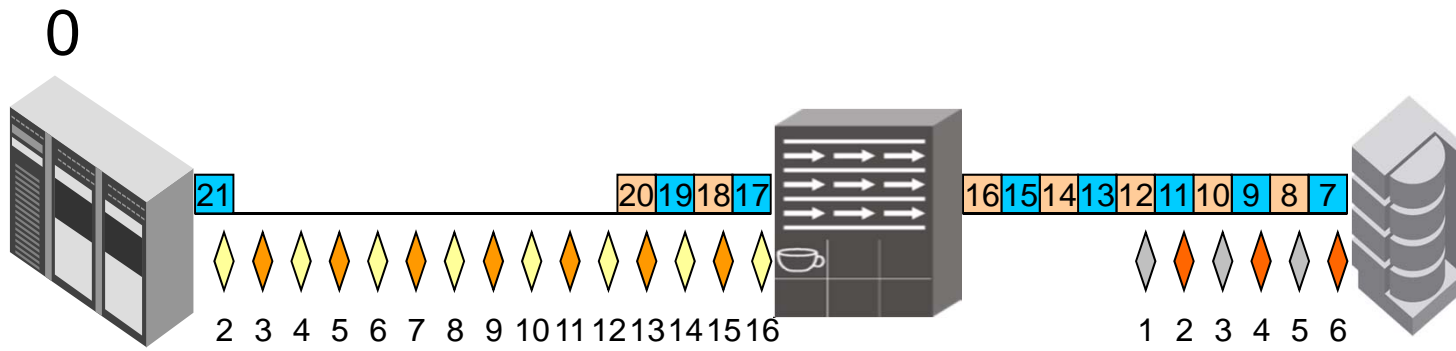


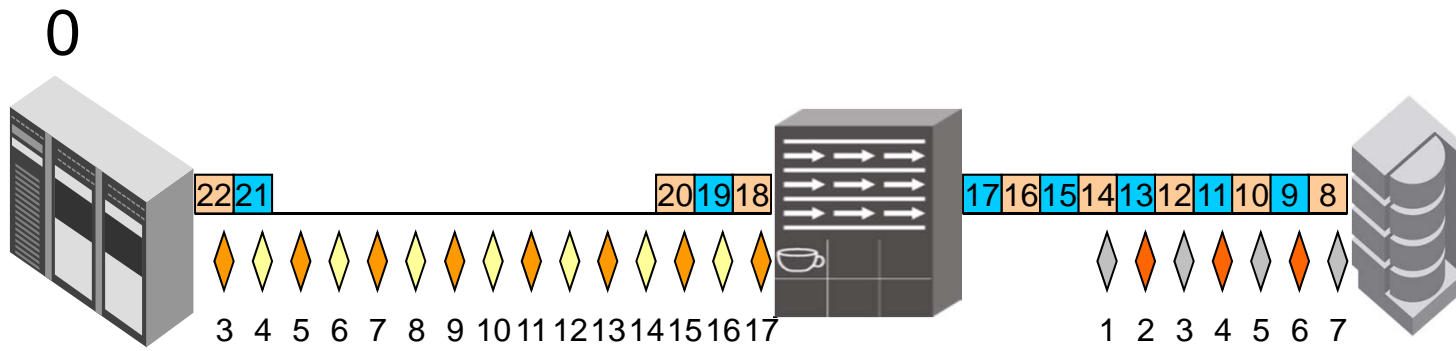


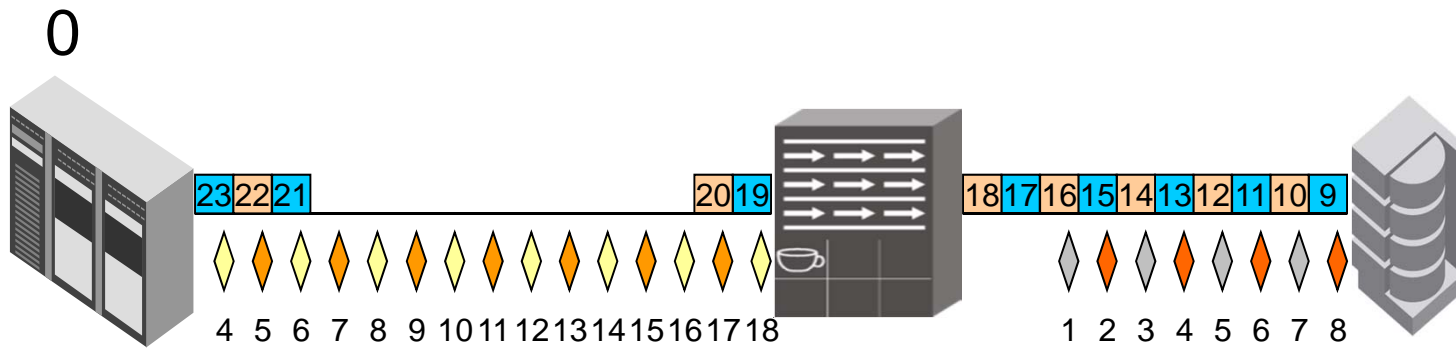


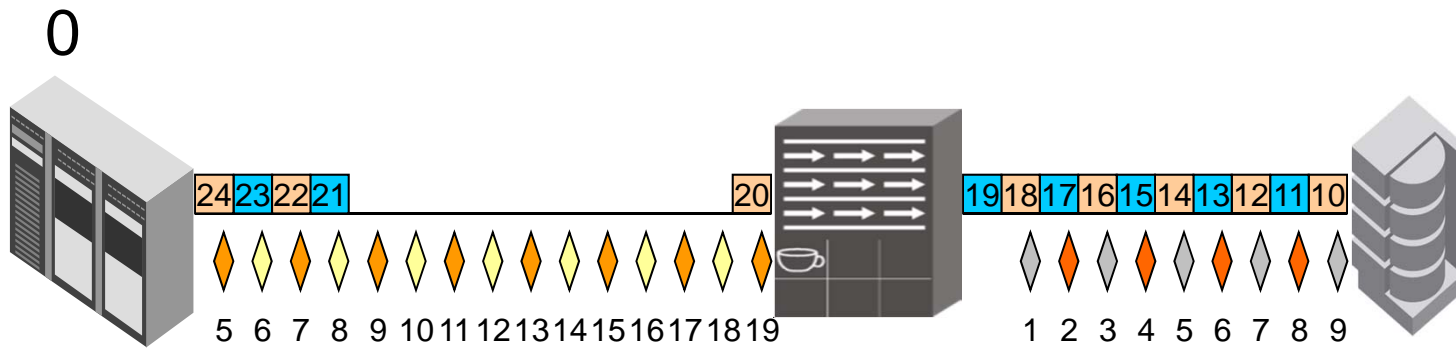


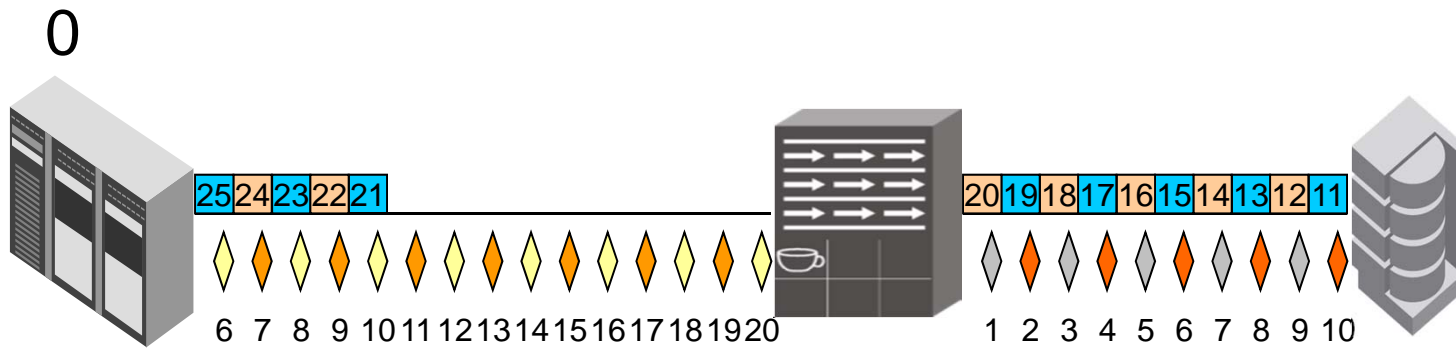


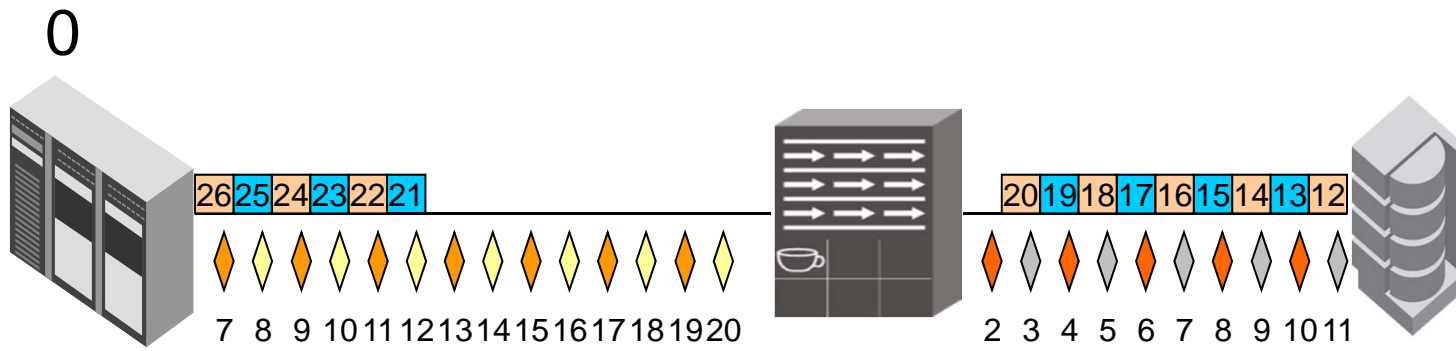


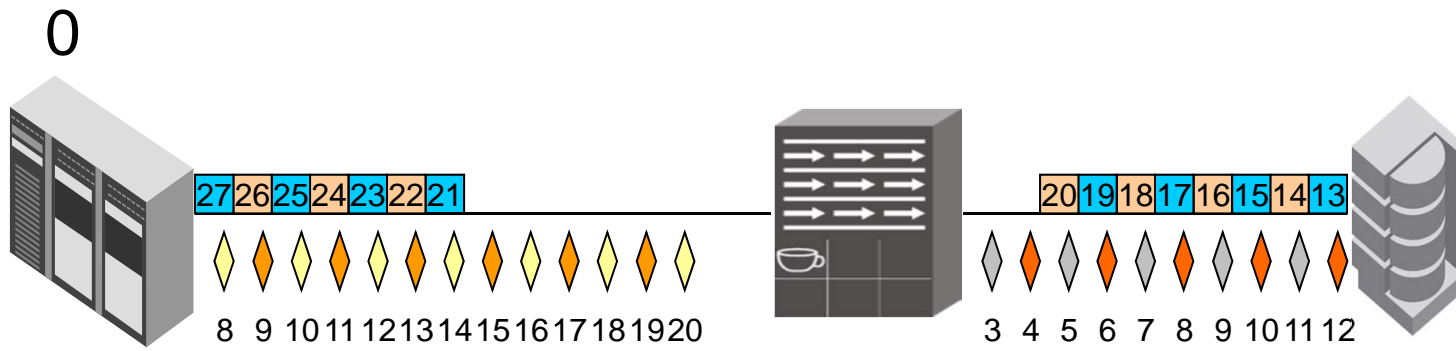


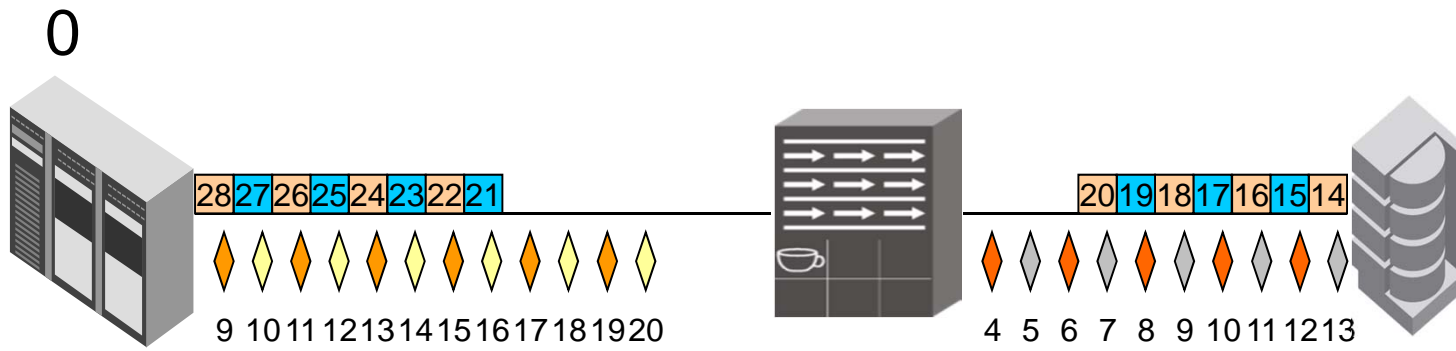


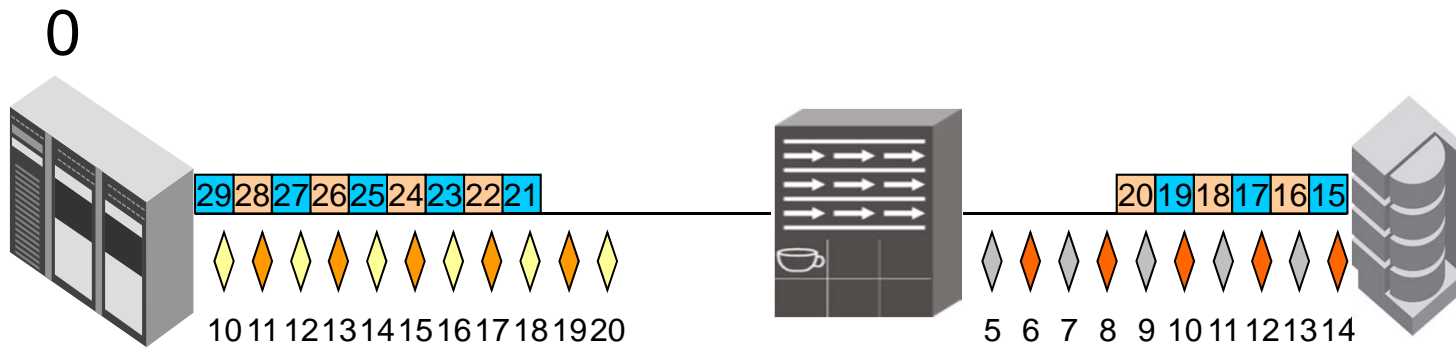


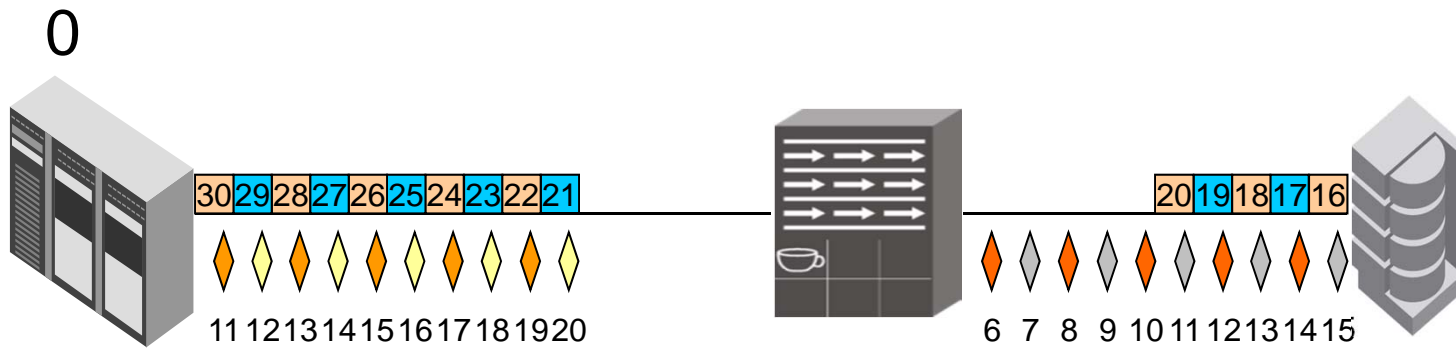


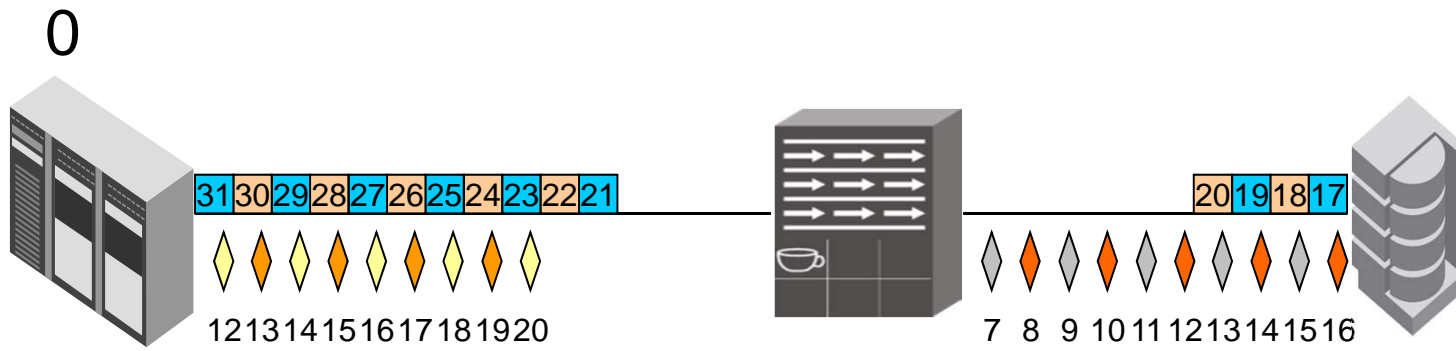


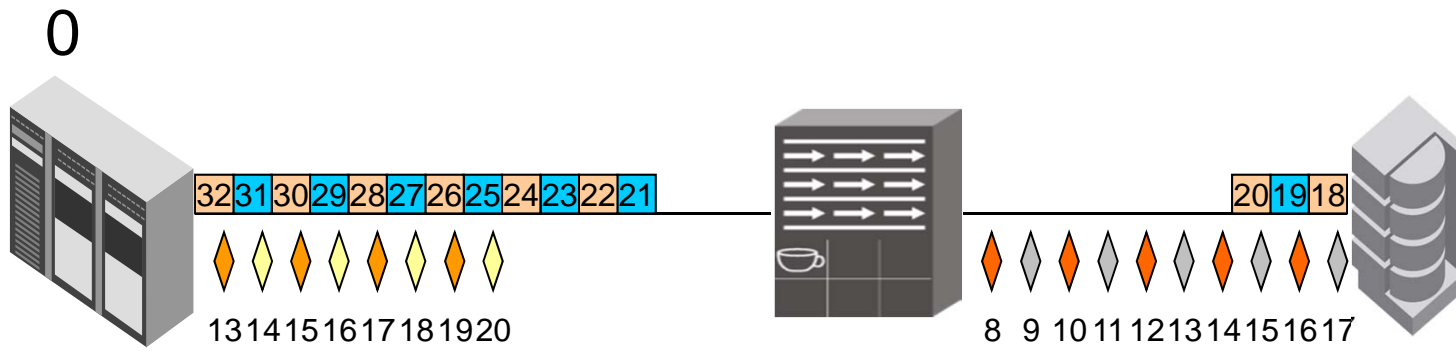


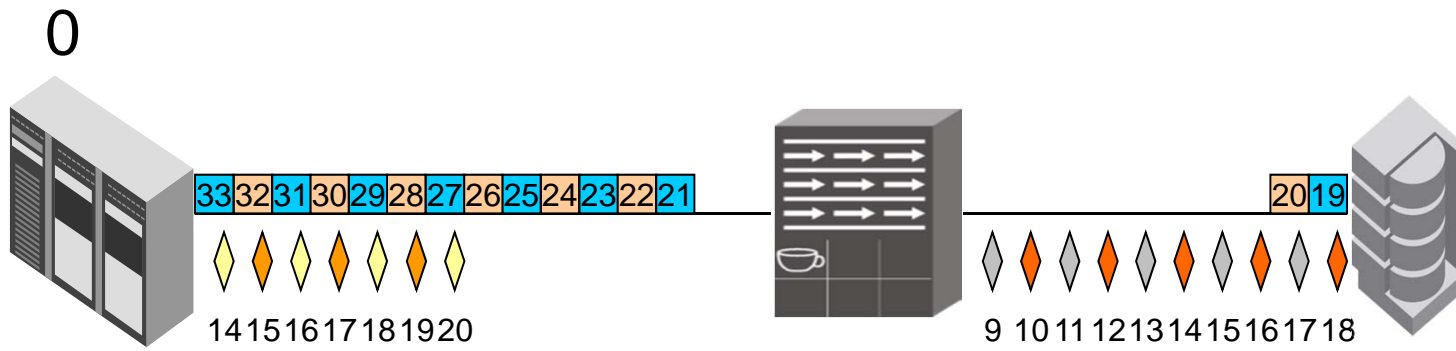


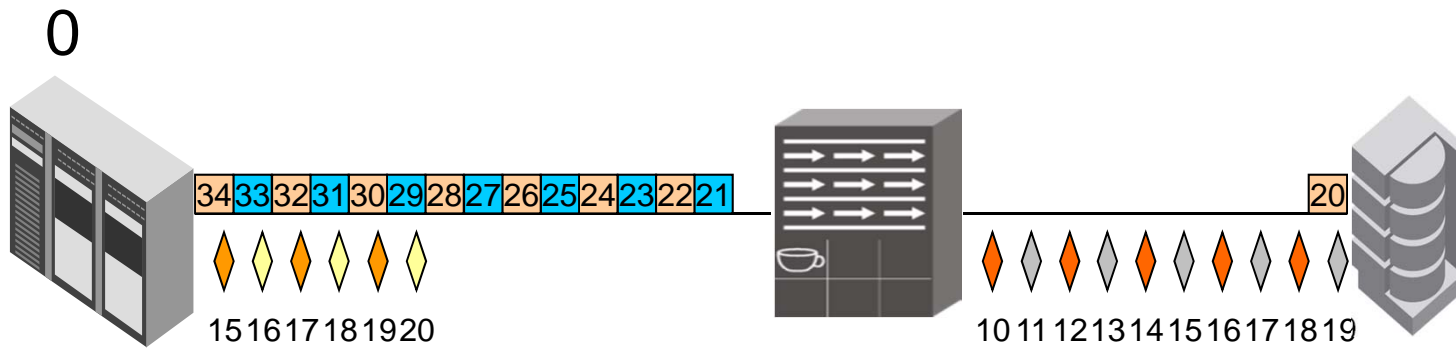


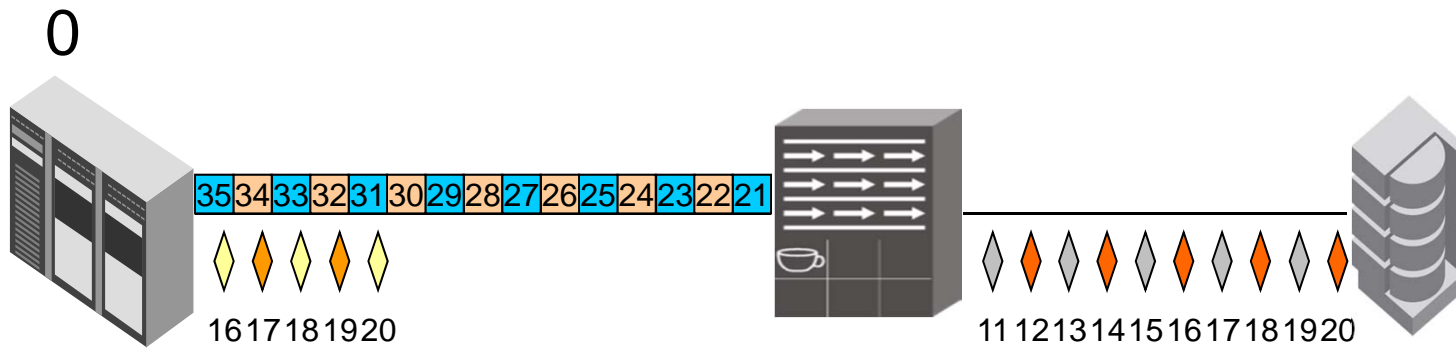


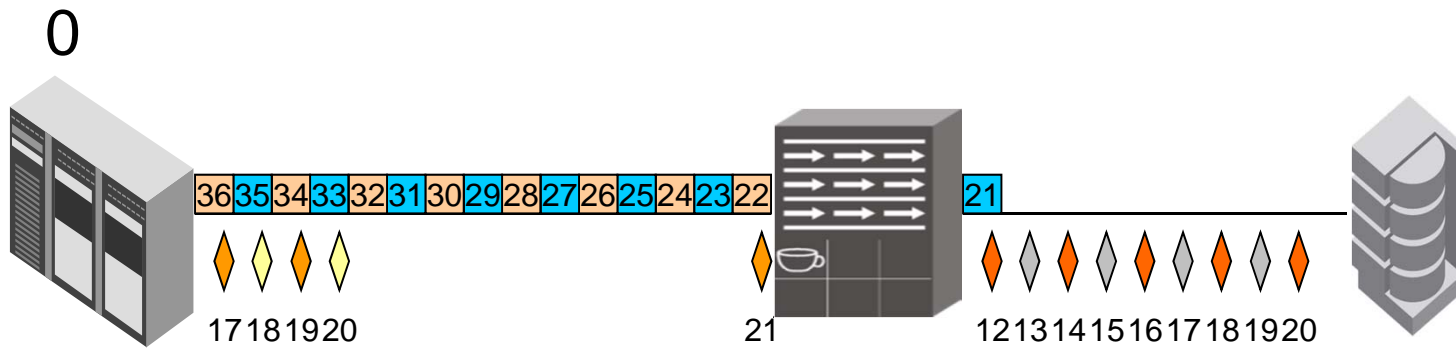


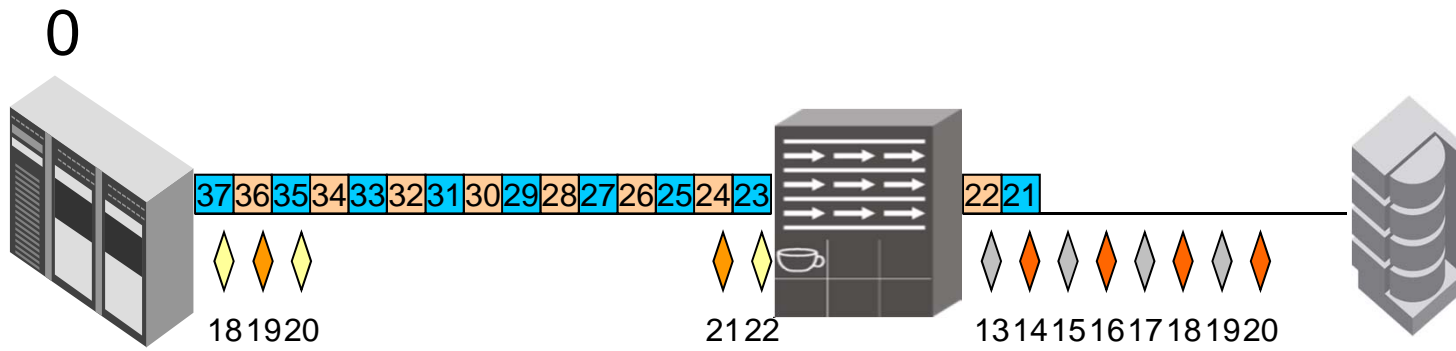


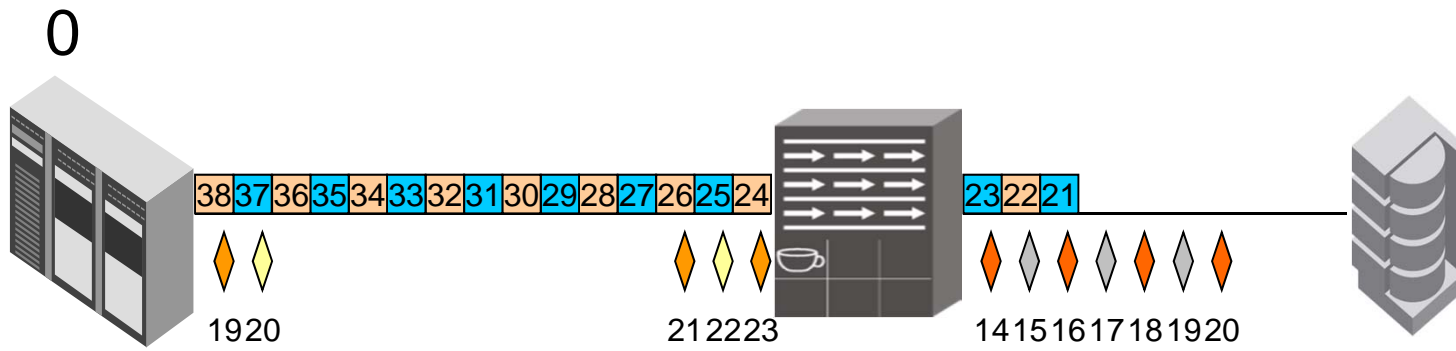




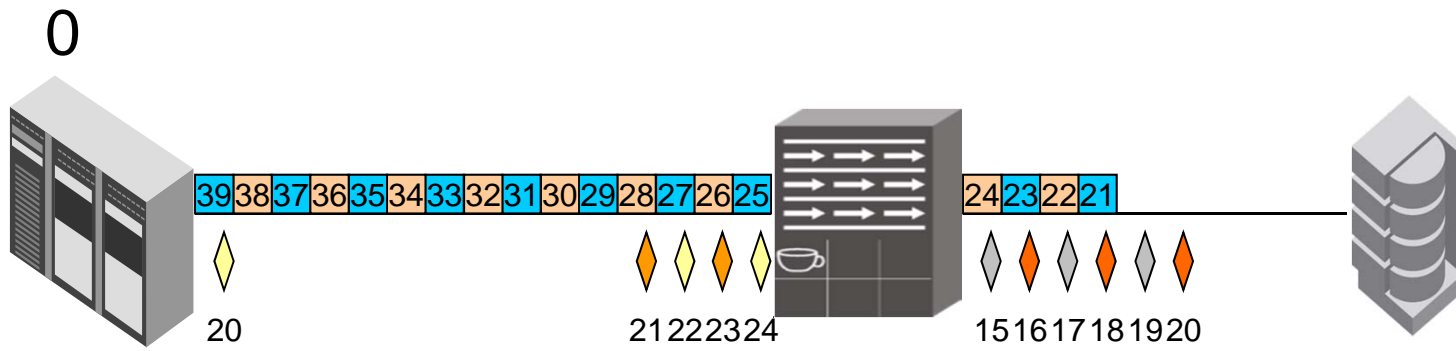


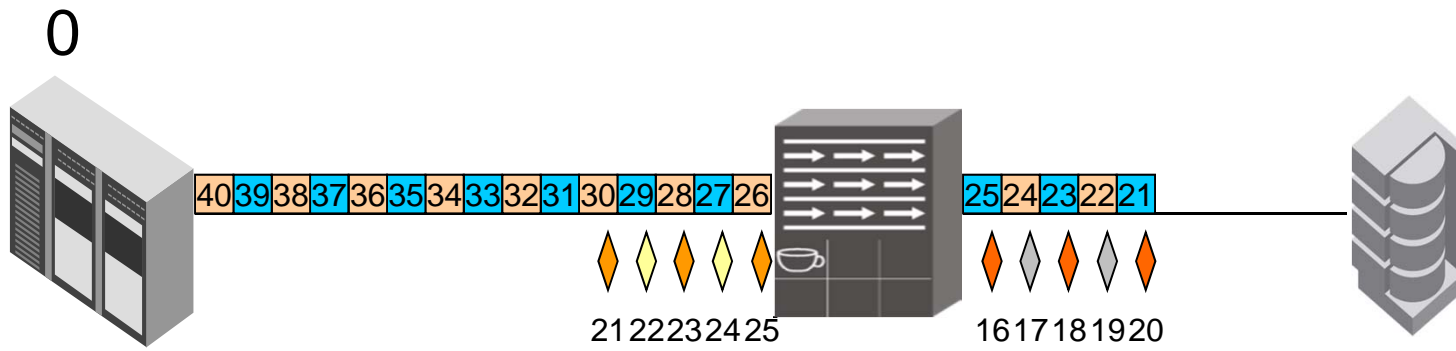


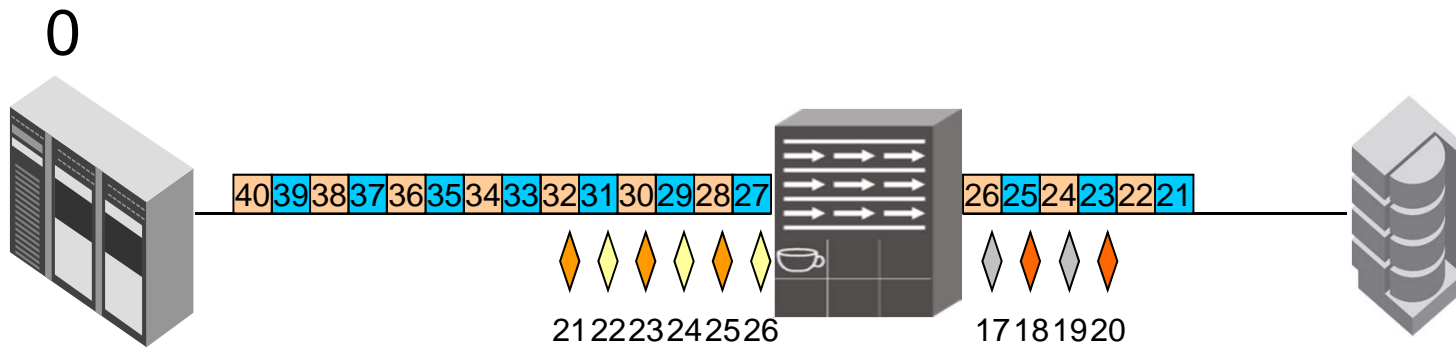


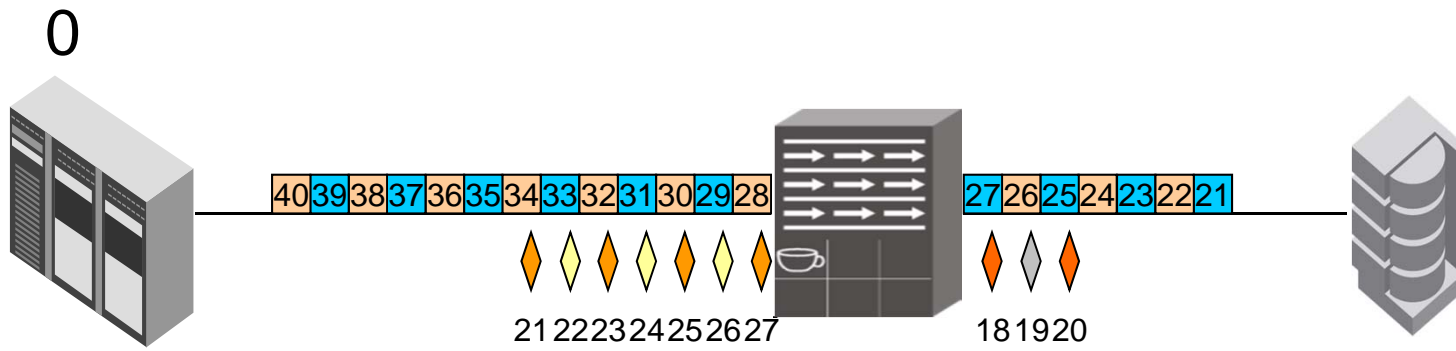


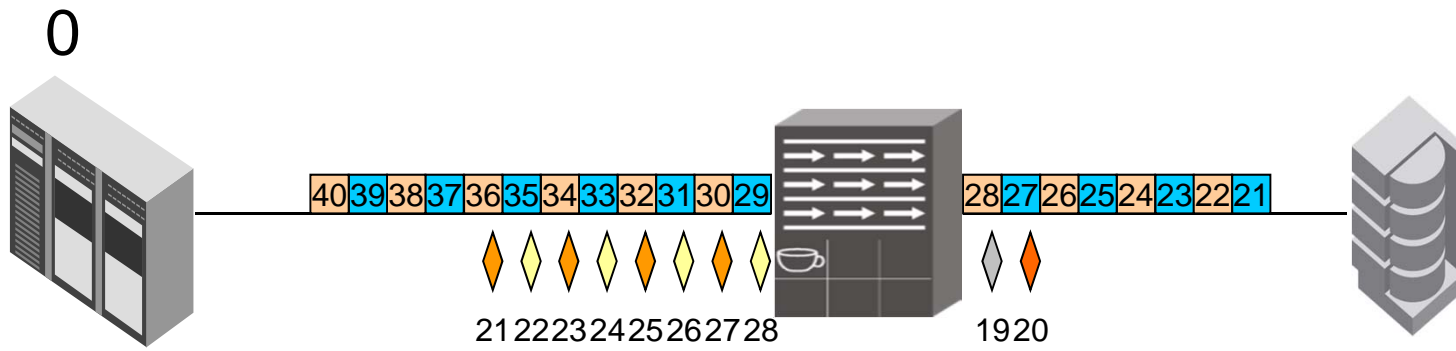
Complete your sessions evaluation online at SHARE.org/BostonEval

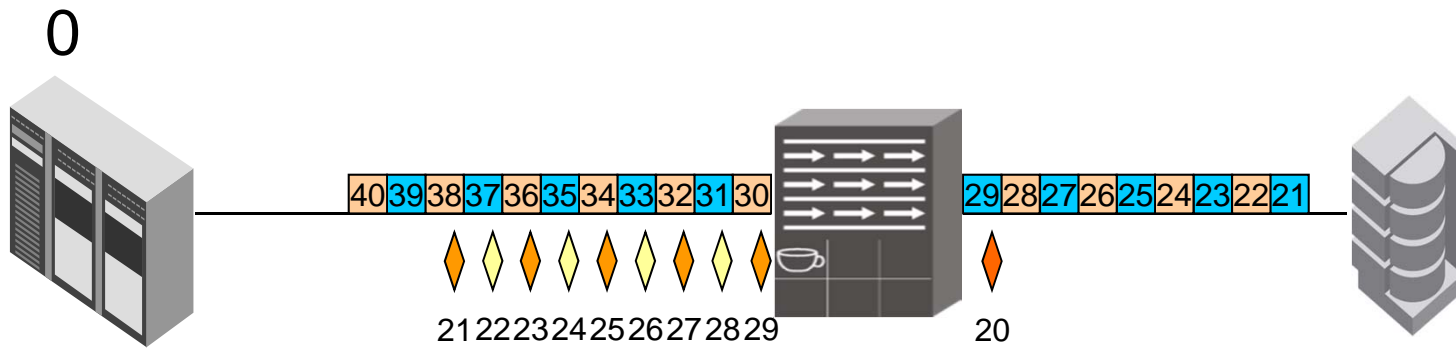


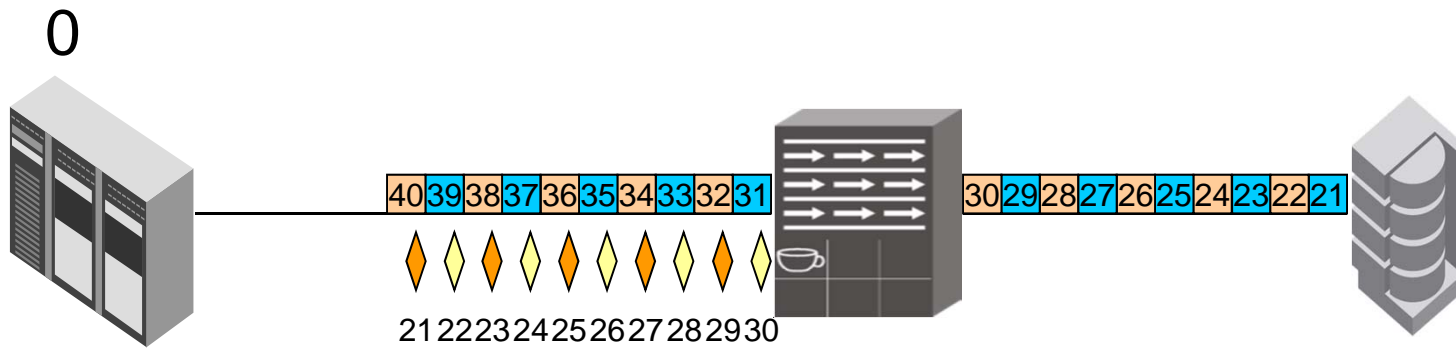


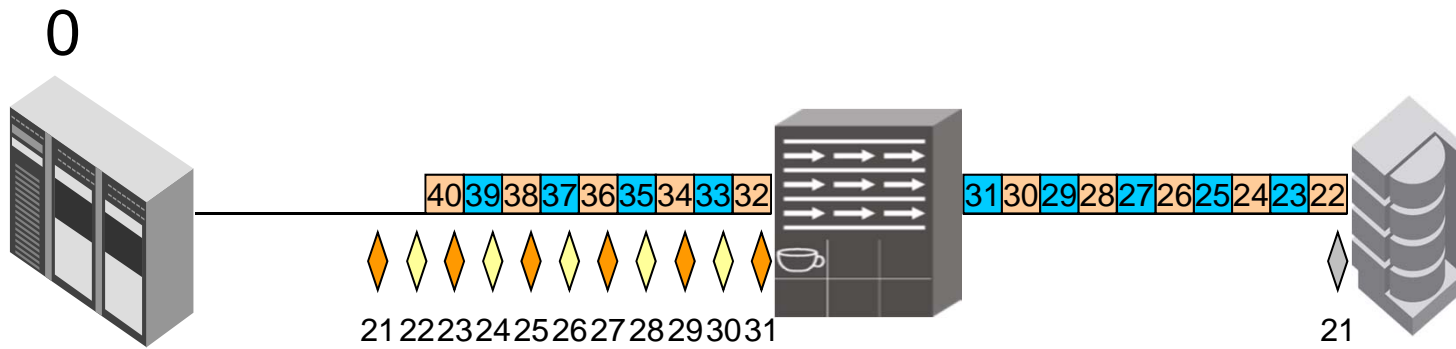


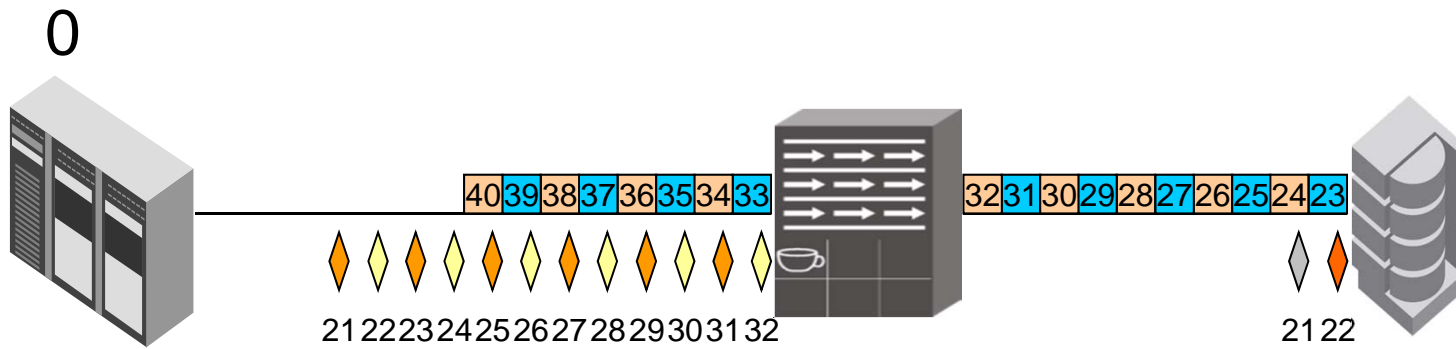










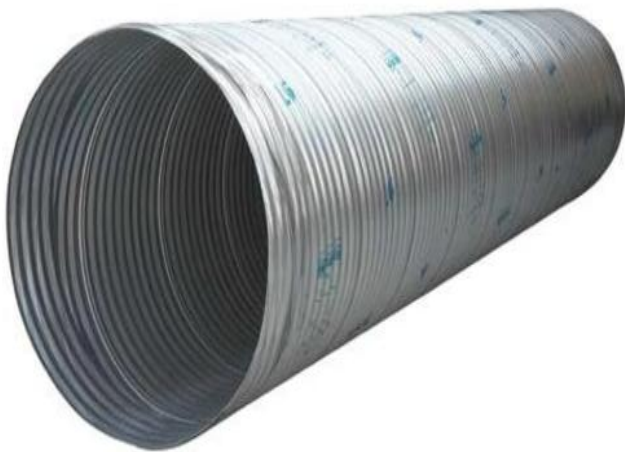


THIS PAGE INTENTIONALLY
LEFT BLANK

Example: Different sized pipes

BUFFER CREDITS

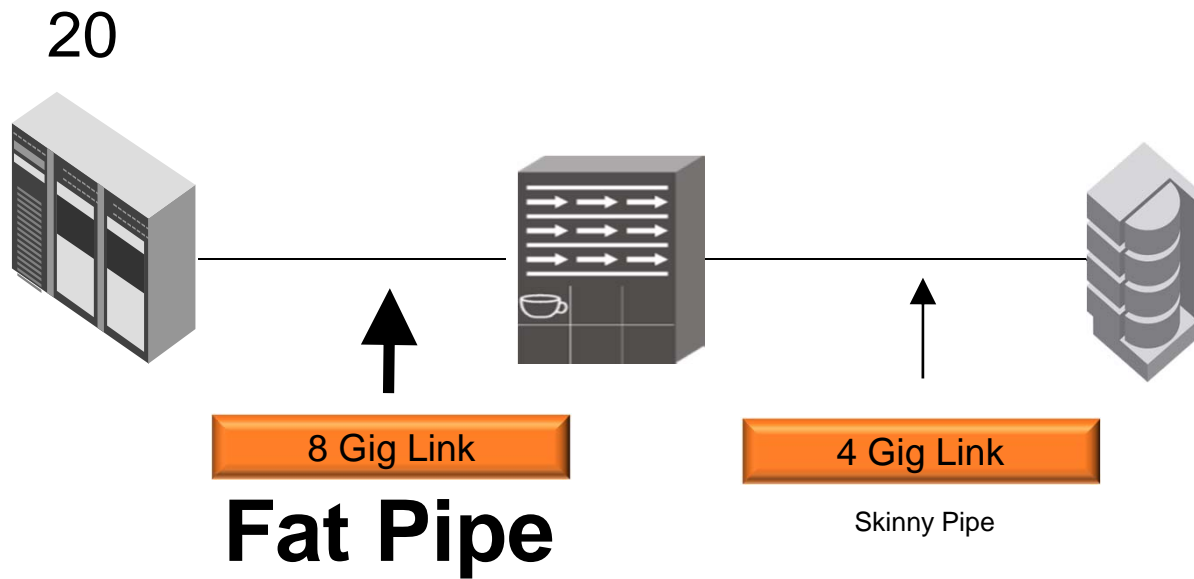
Fat Pipe



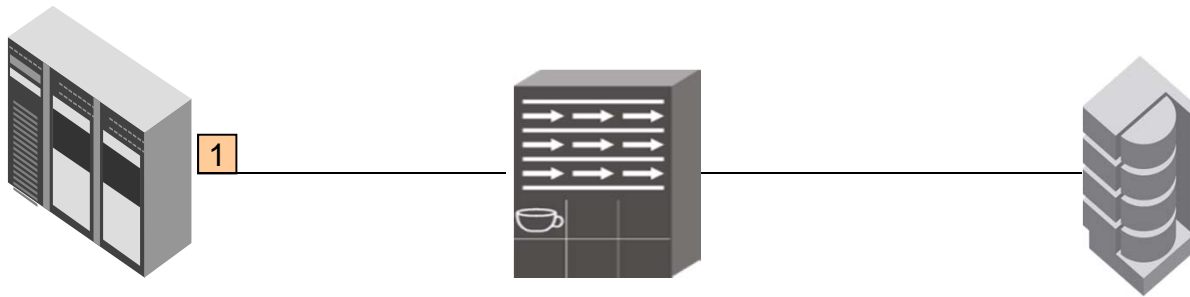
Skinny Pipe



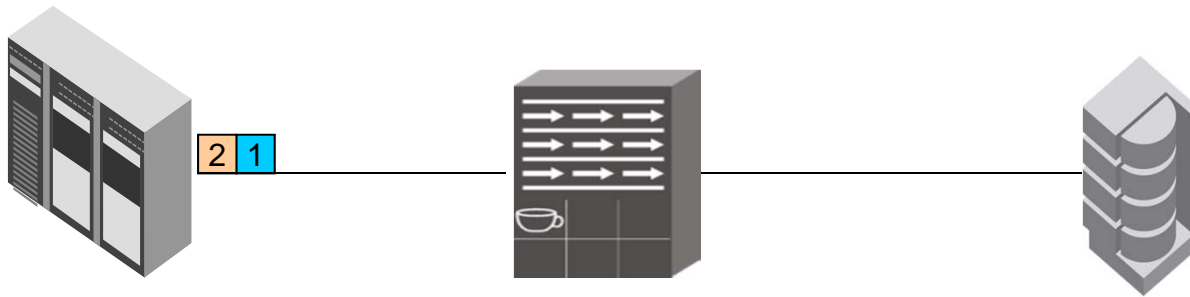
Fat Pipe / Skinny Pipe



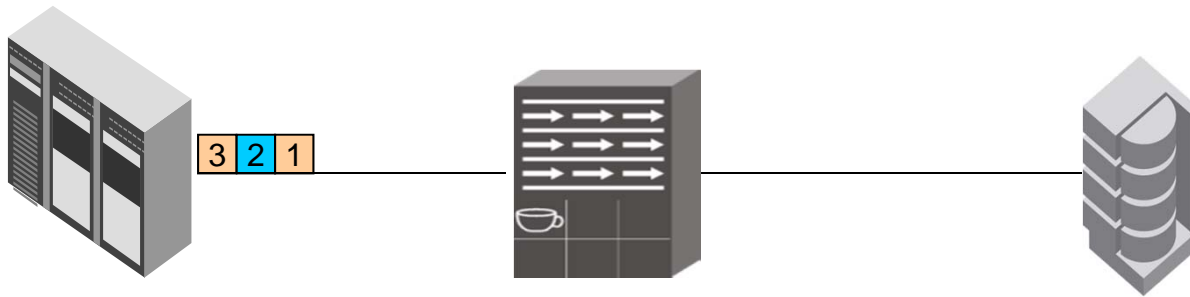
19

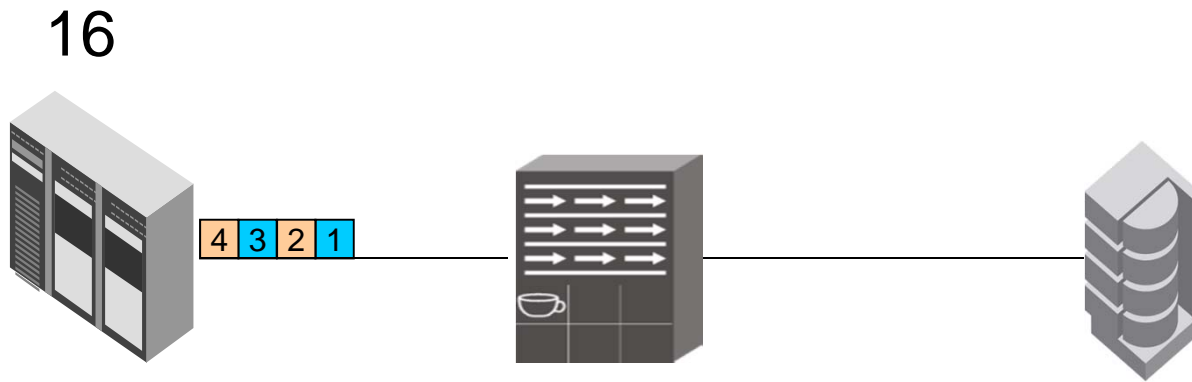


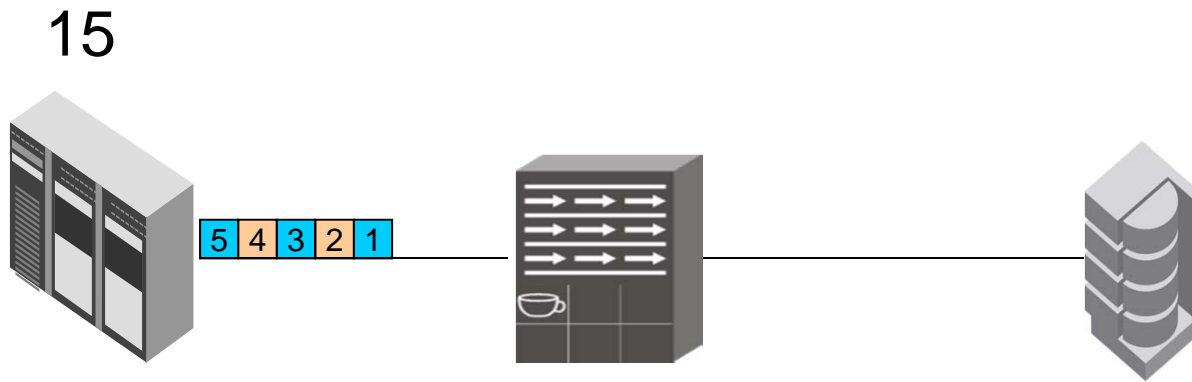
18

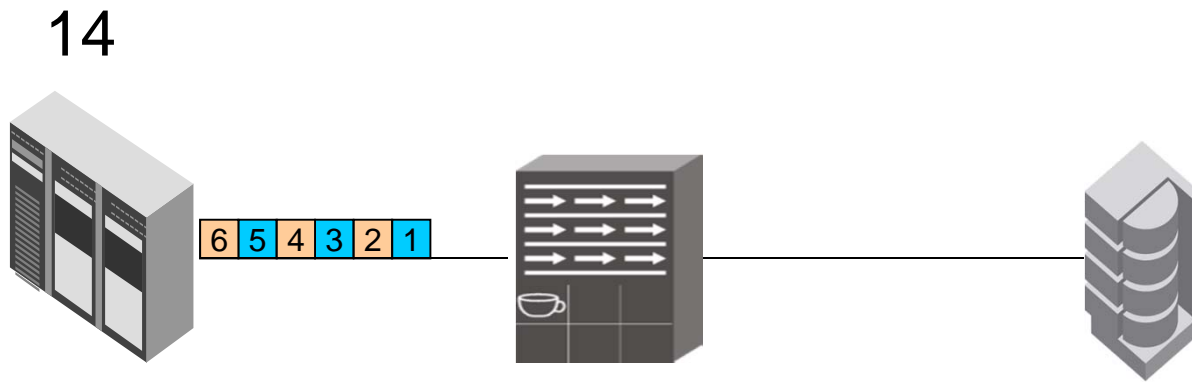


17

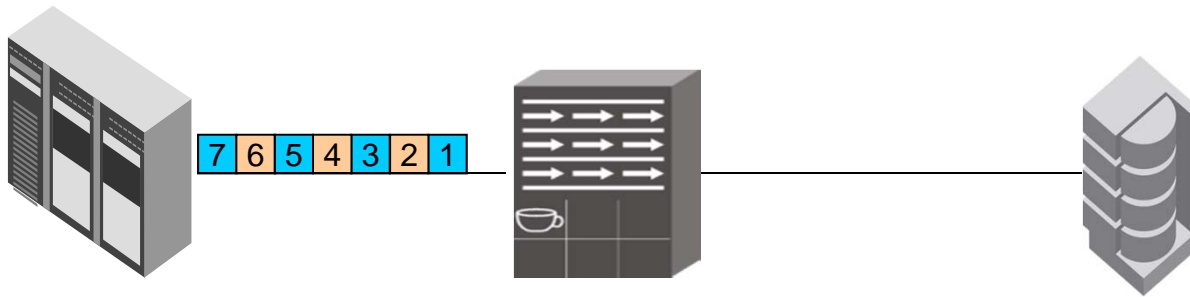


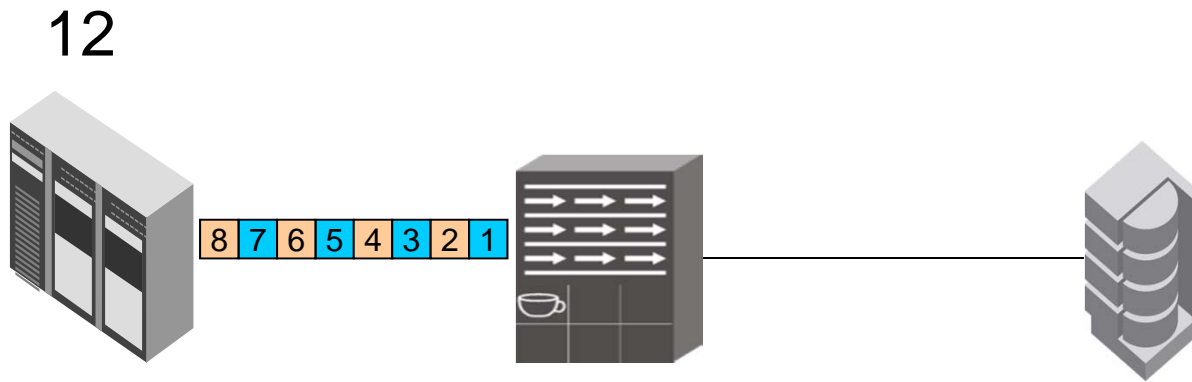


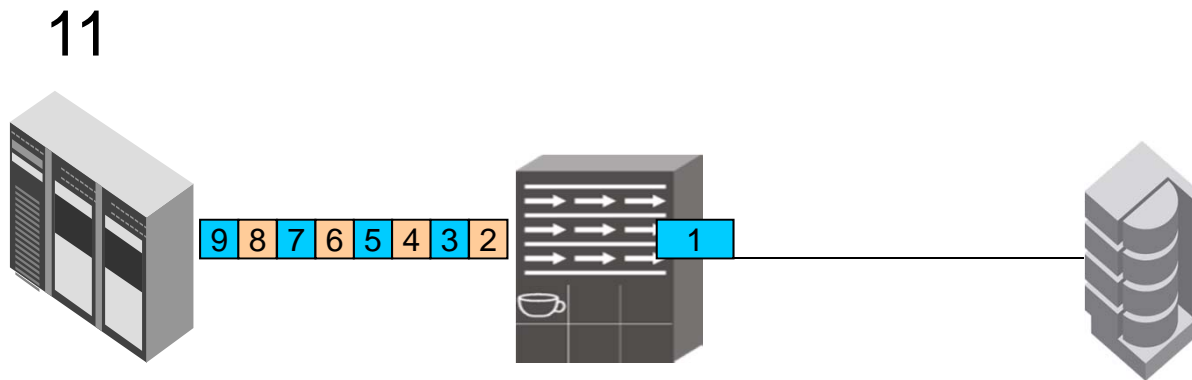


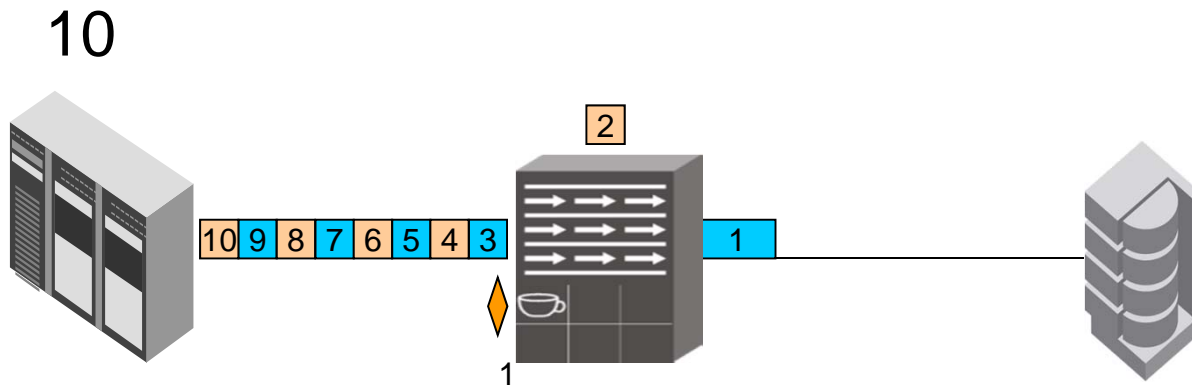


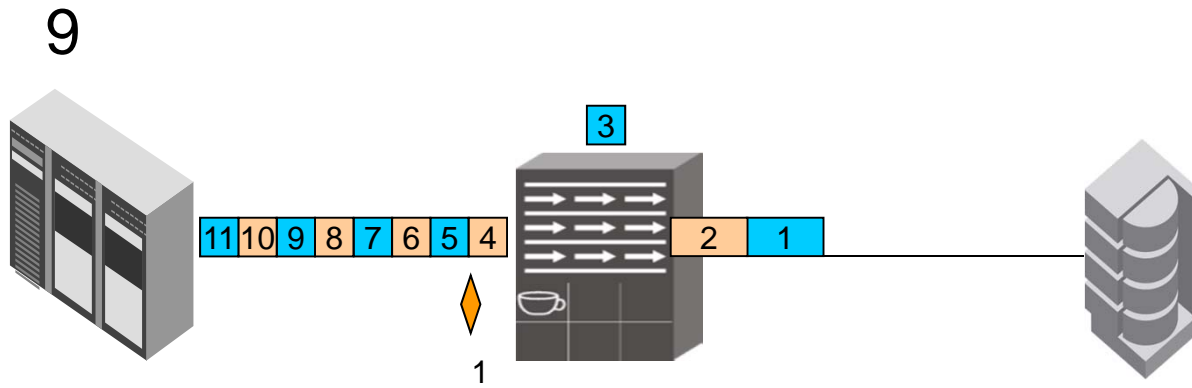
13

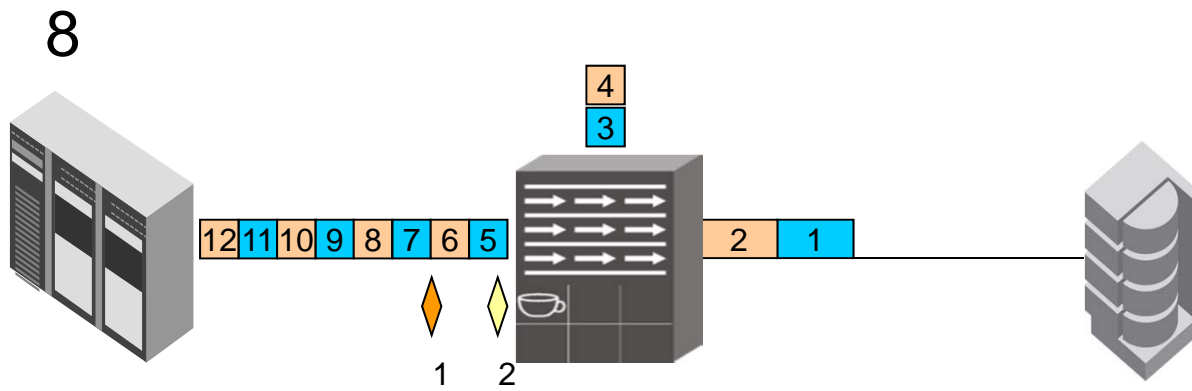


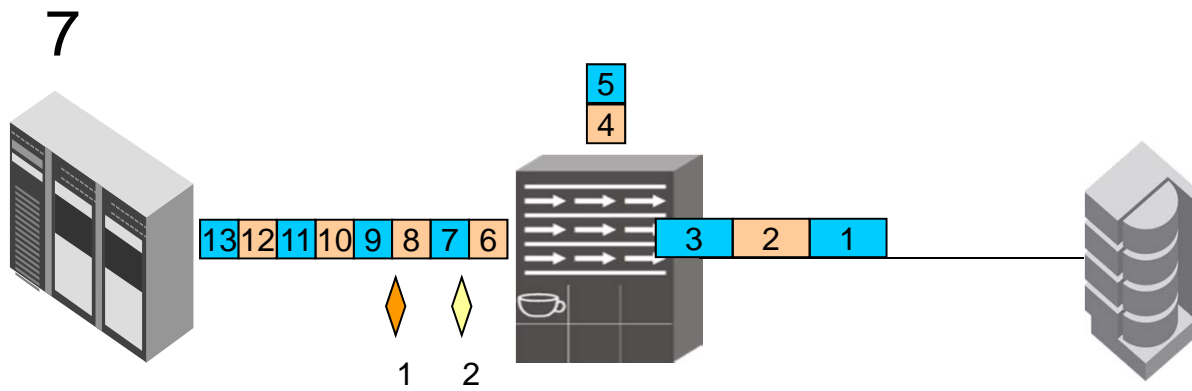


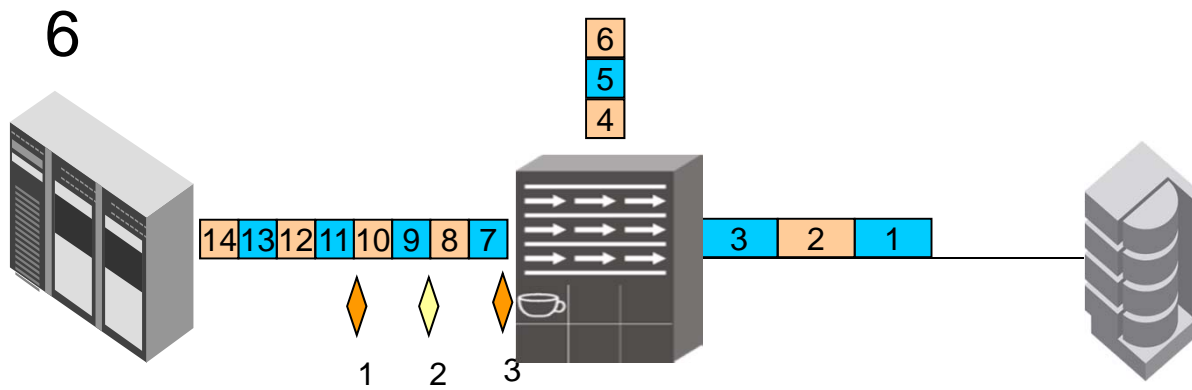


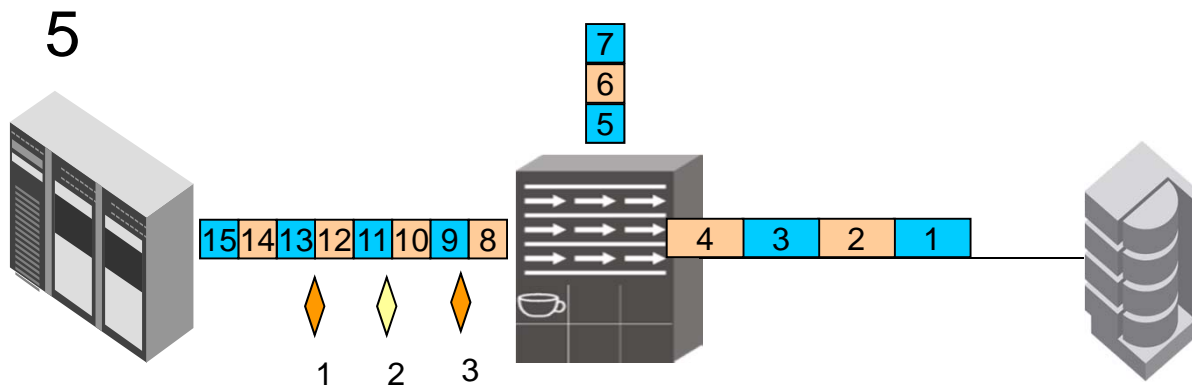


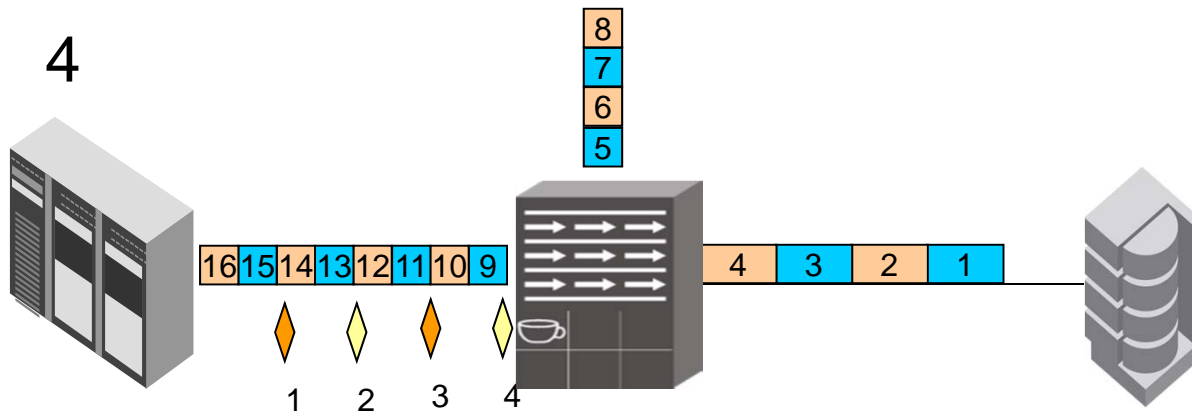


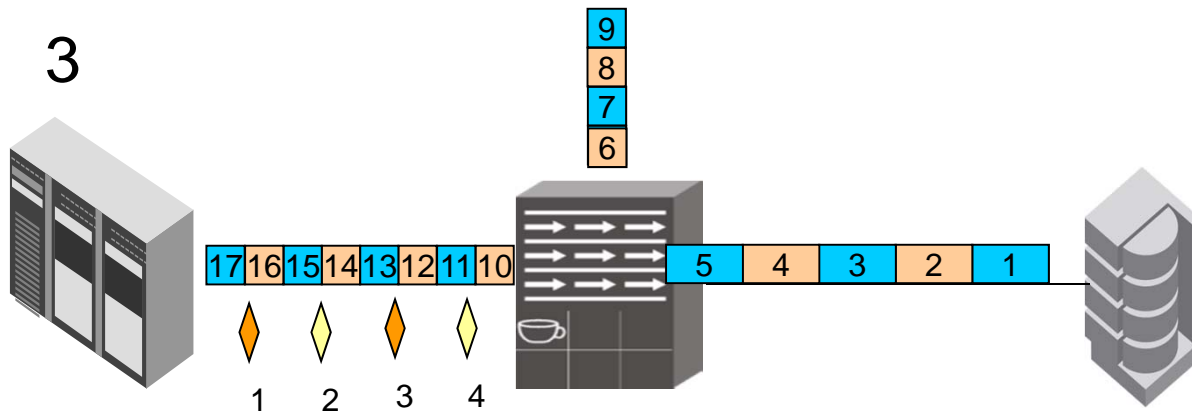


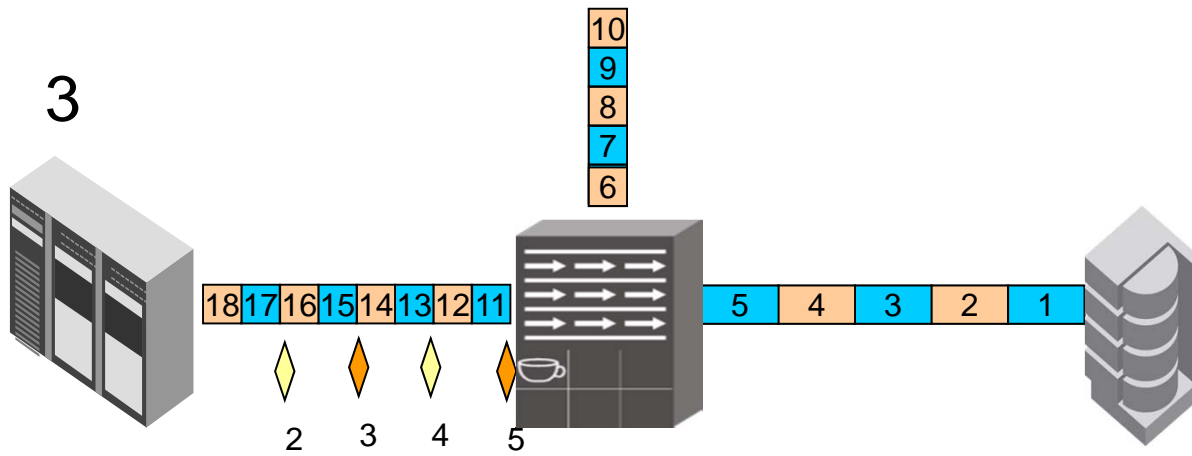


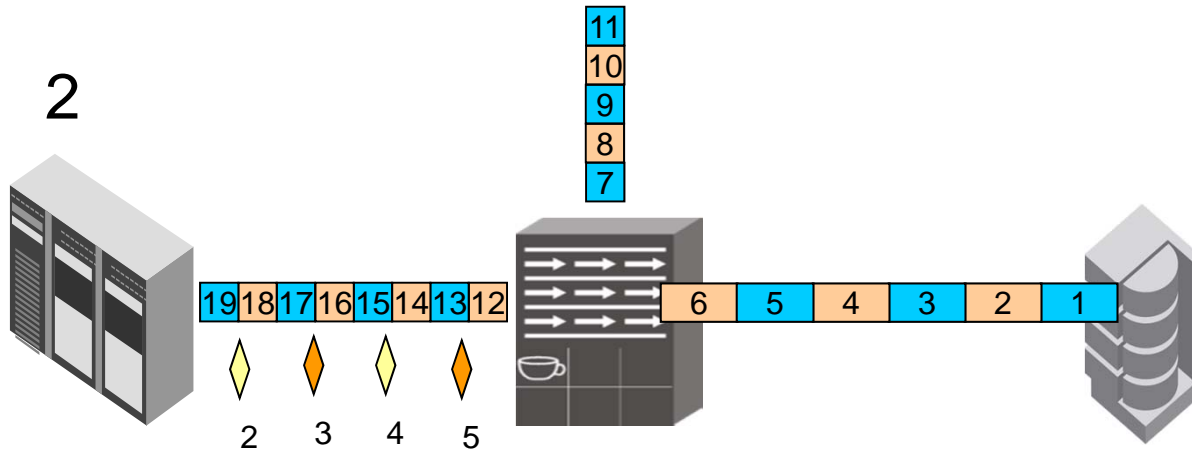


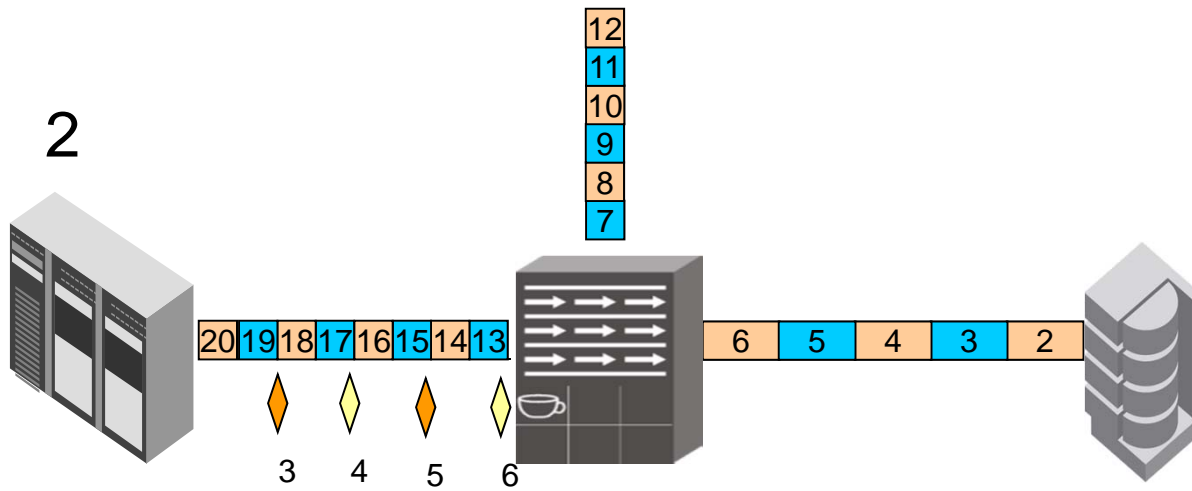


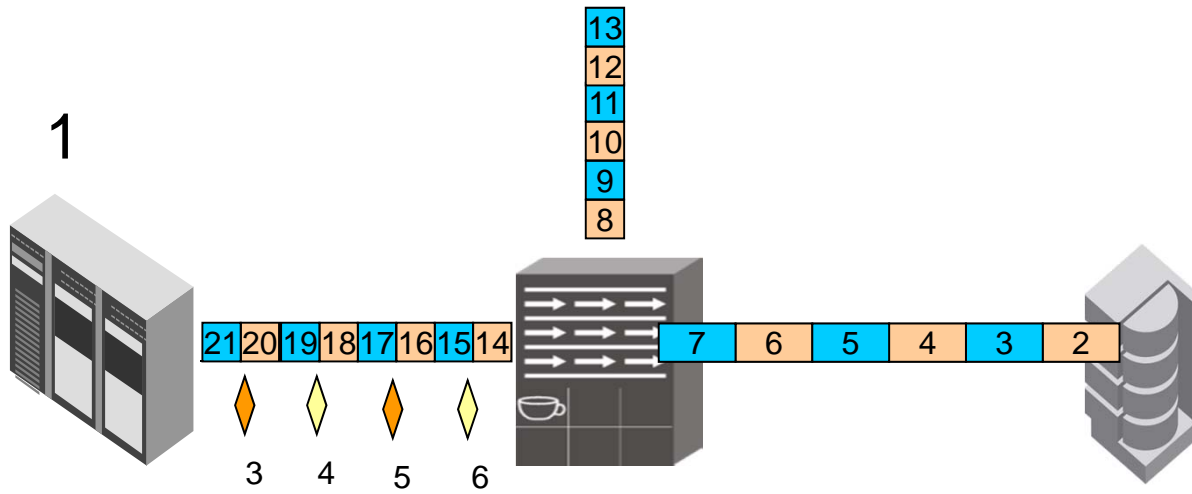


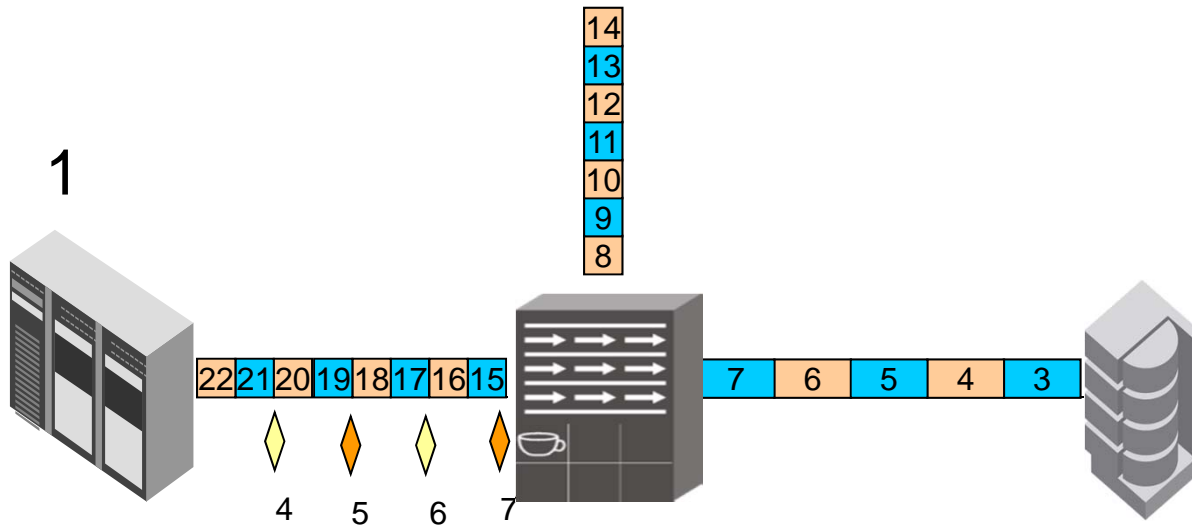


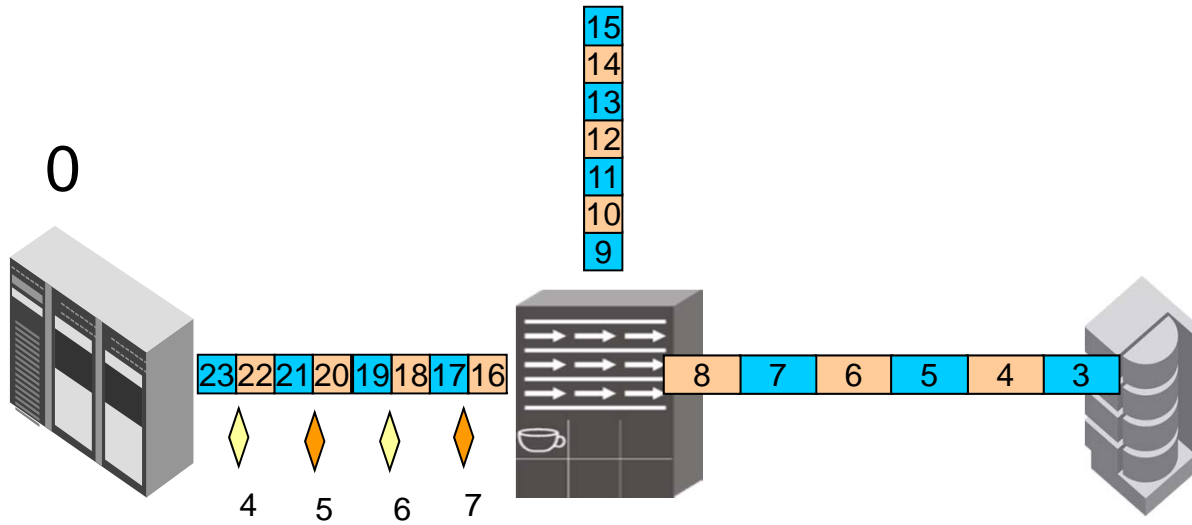


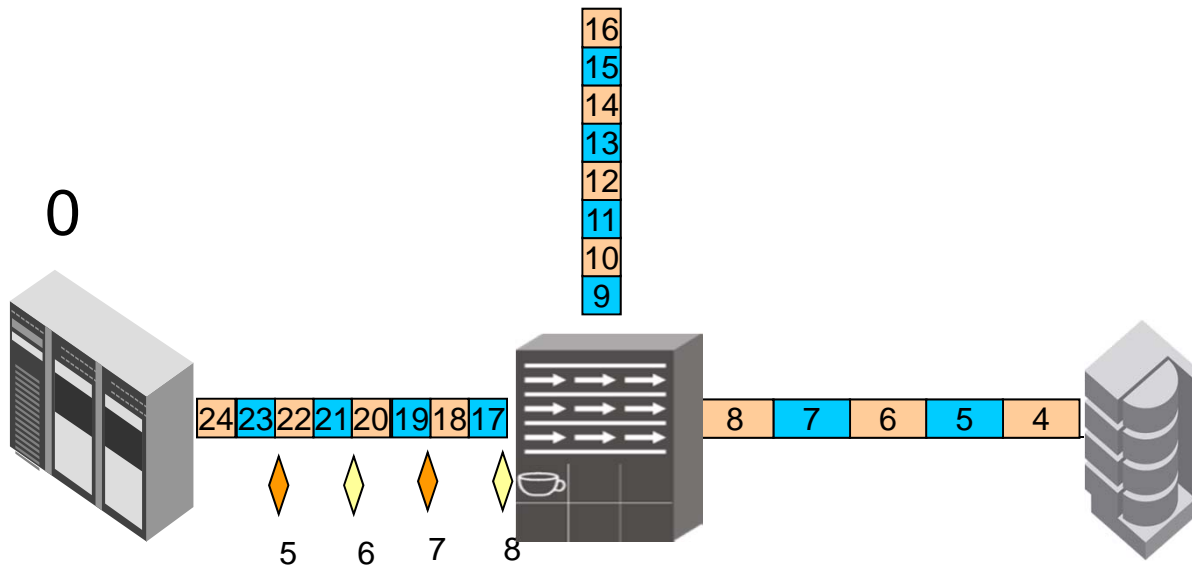


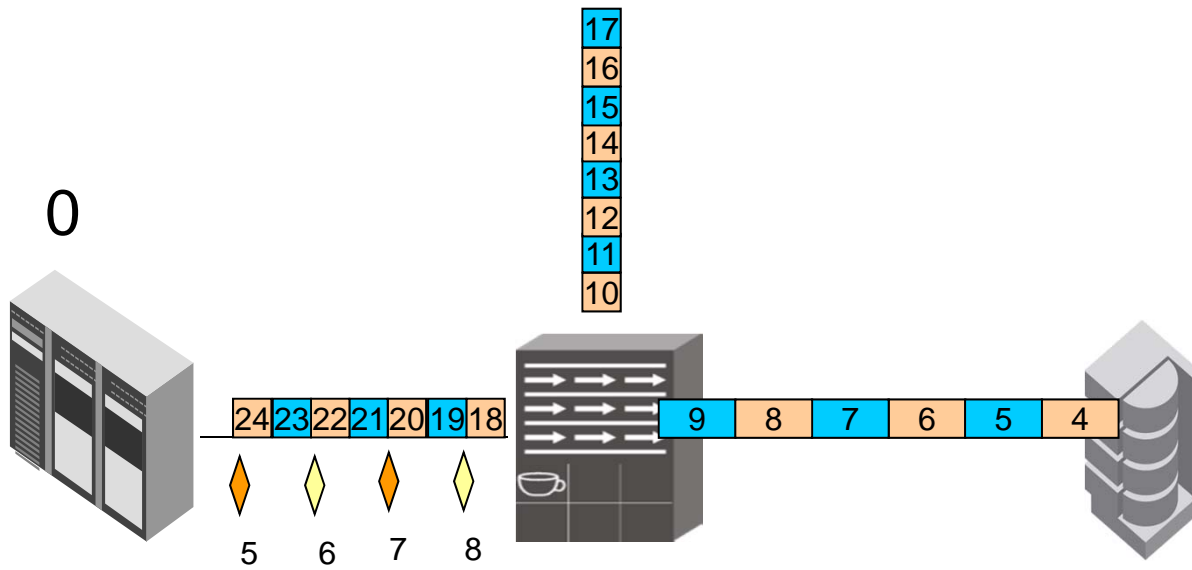


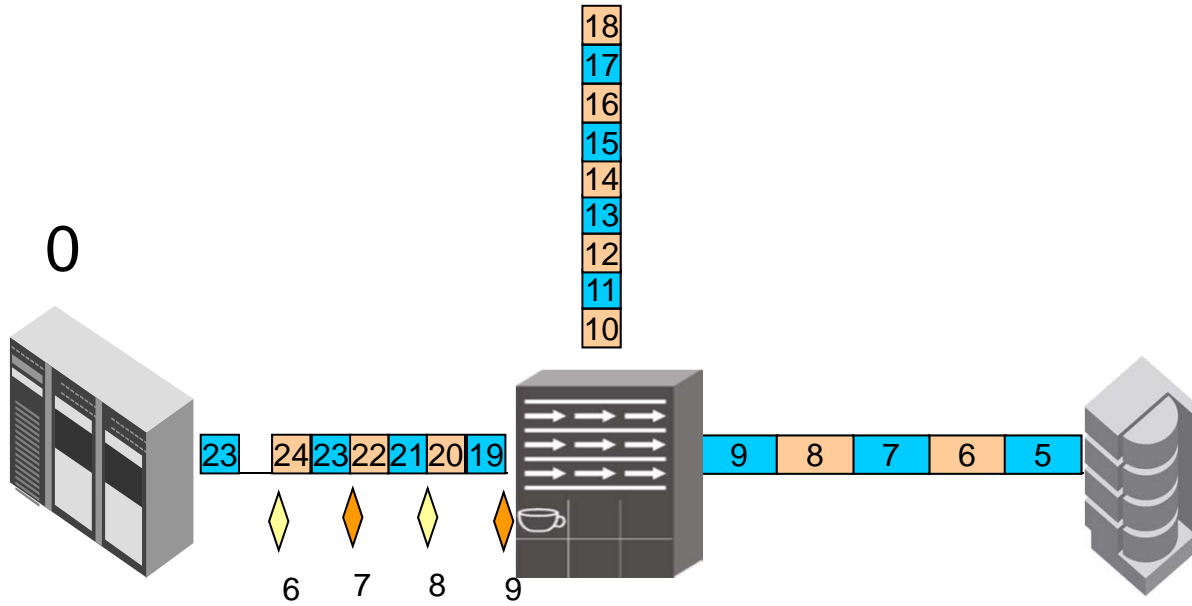


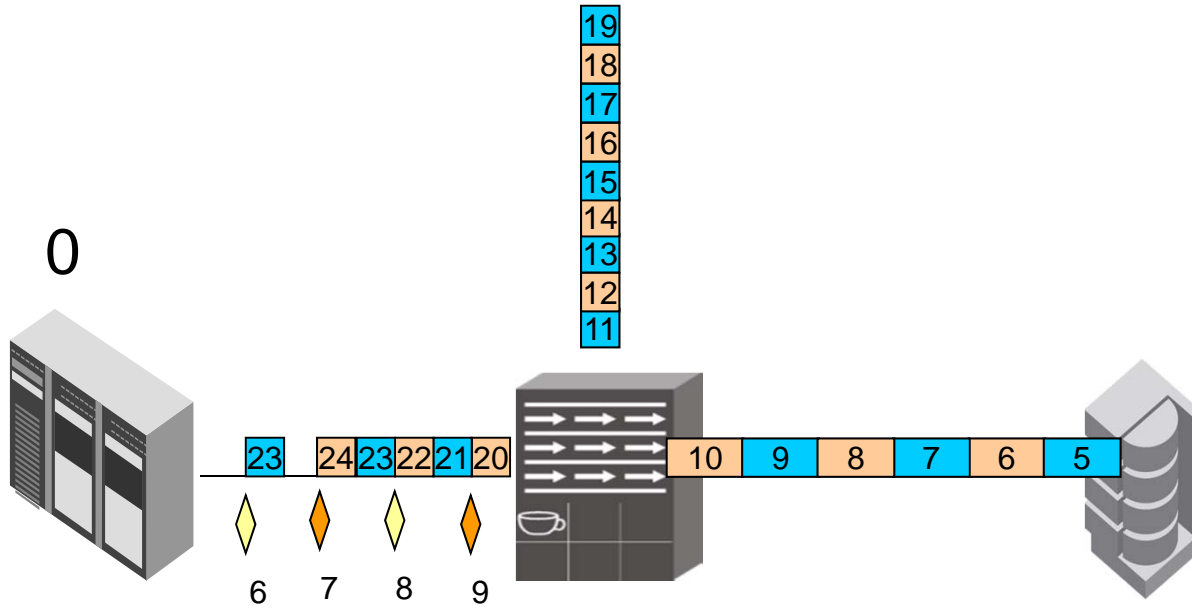


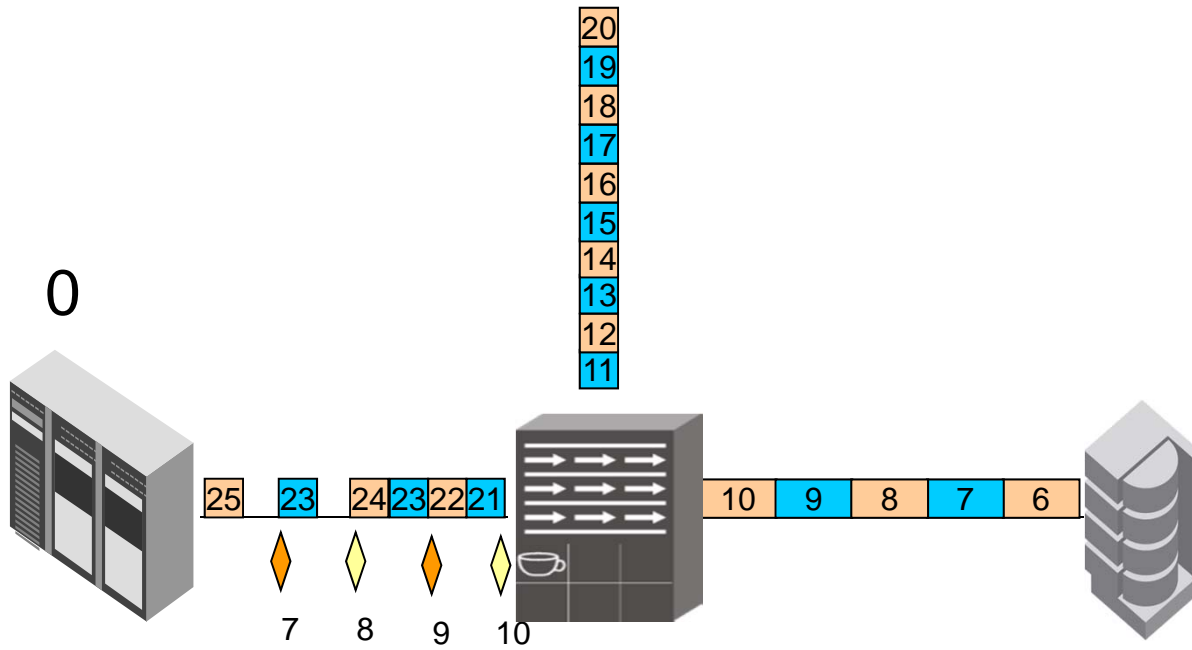


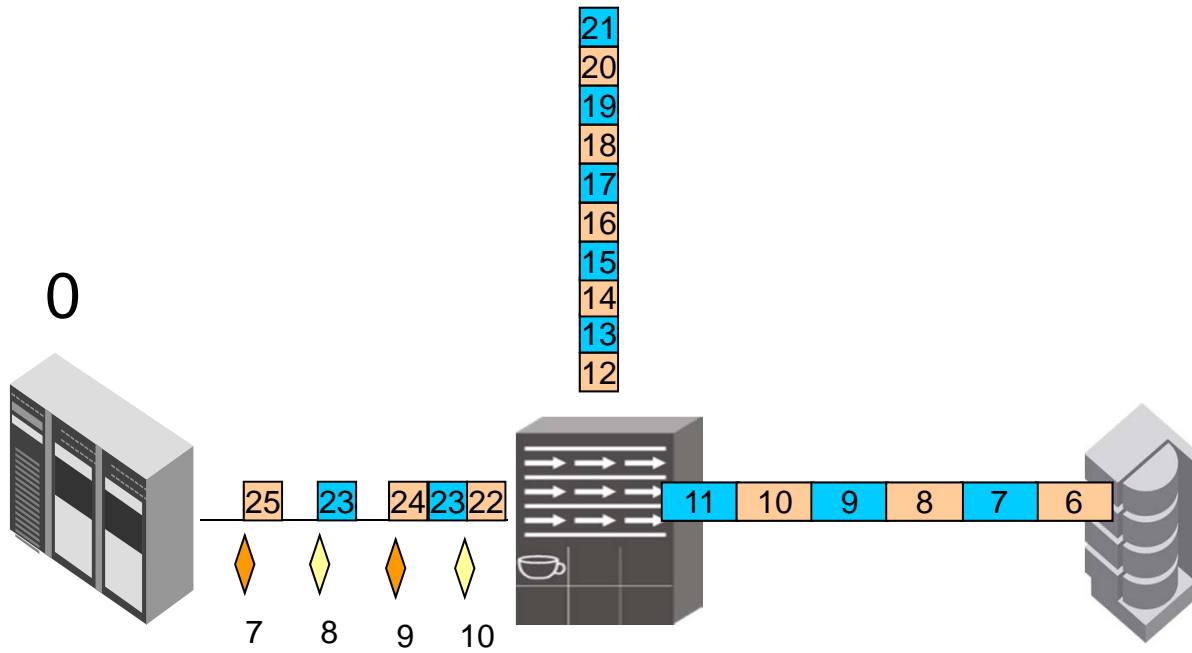


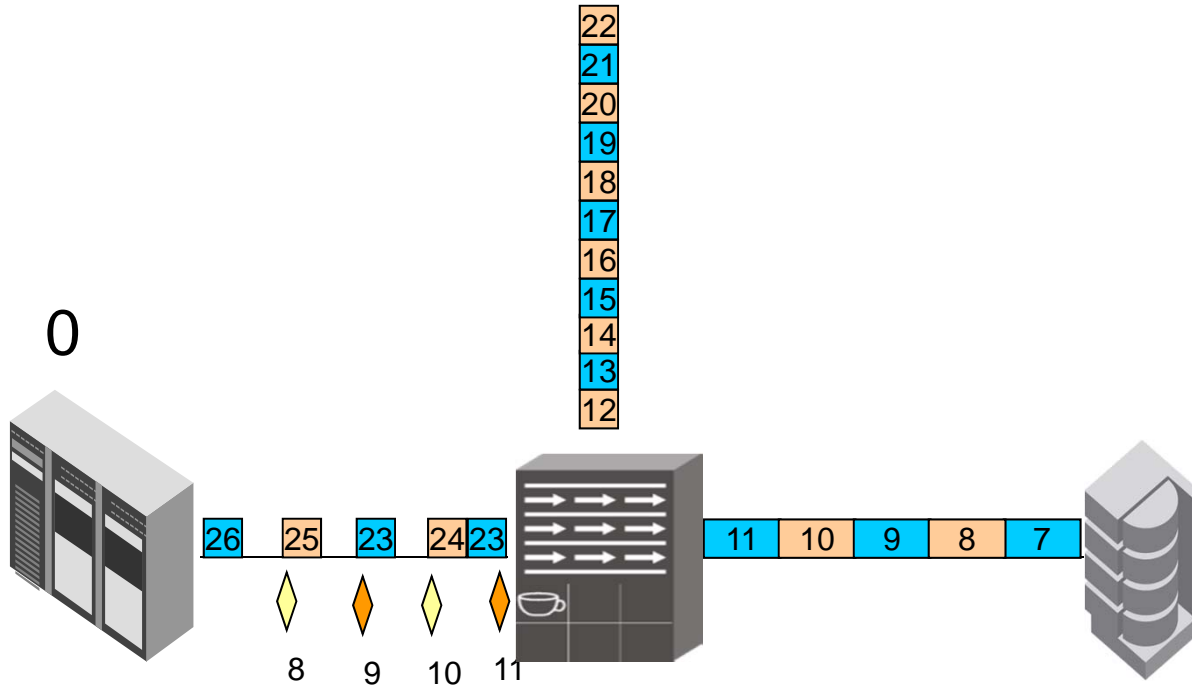


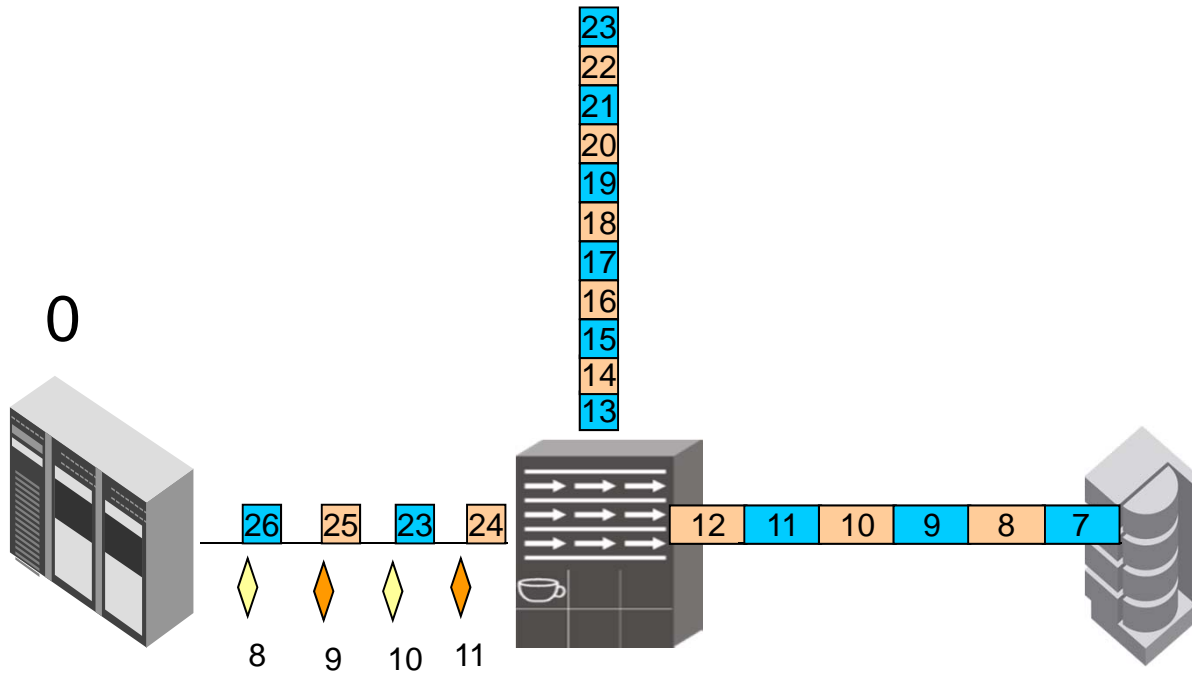


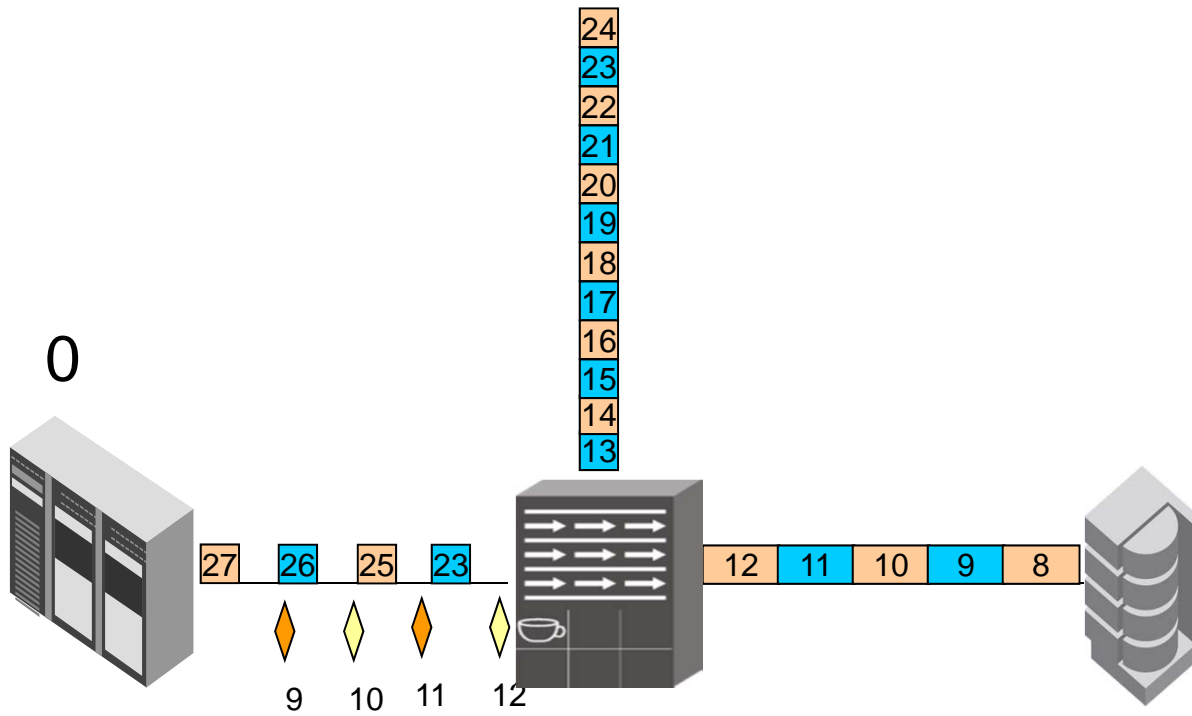


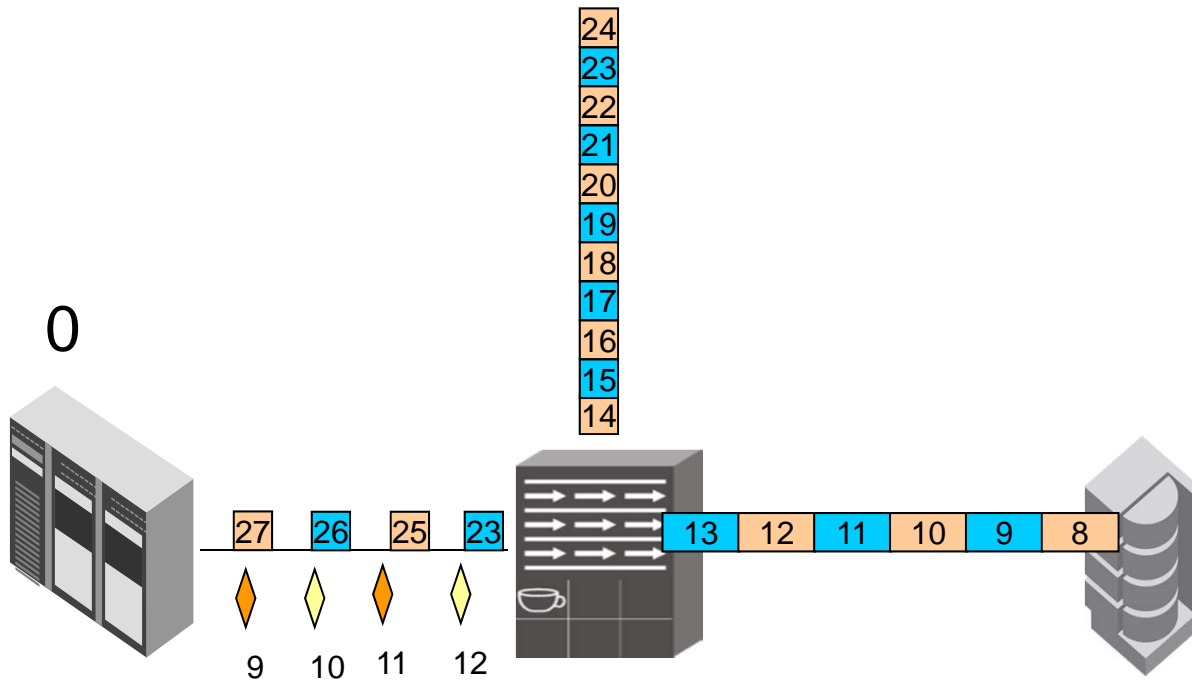


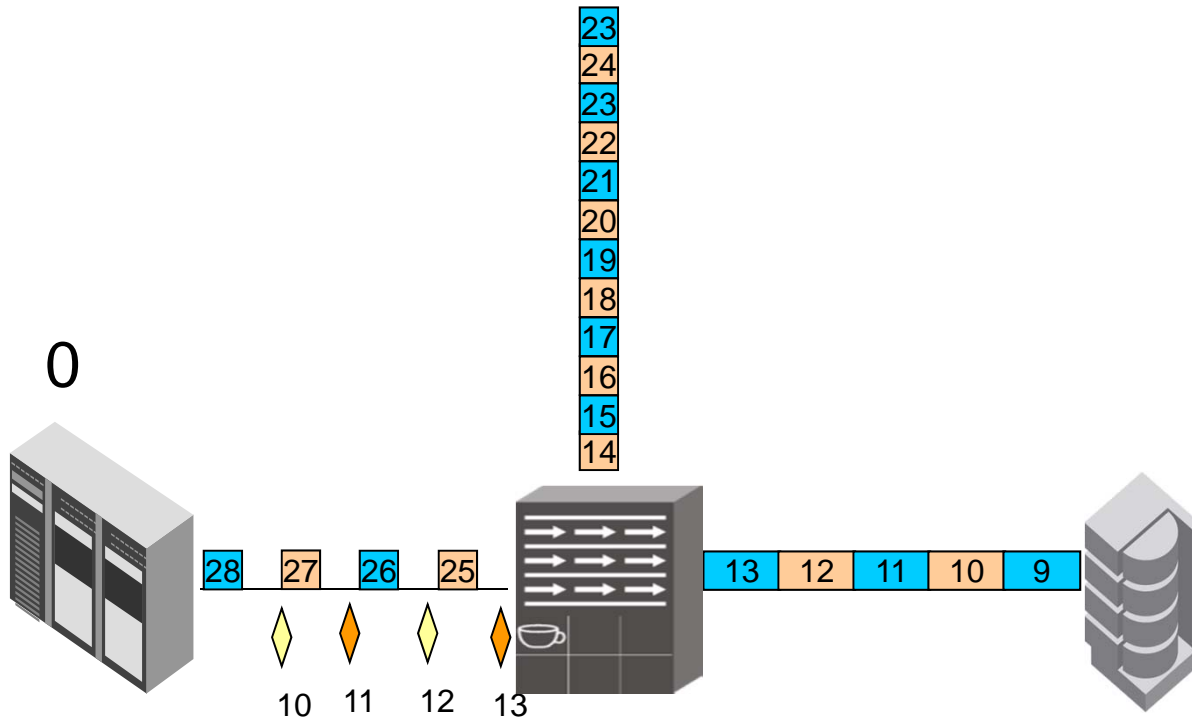


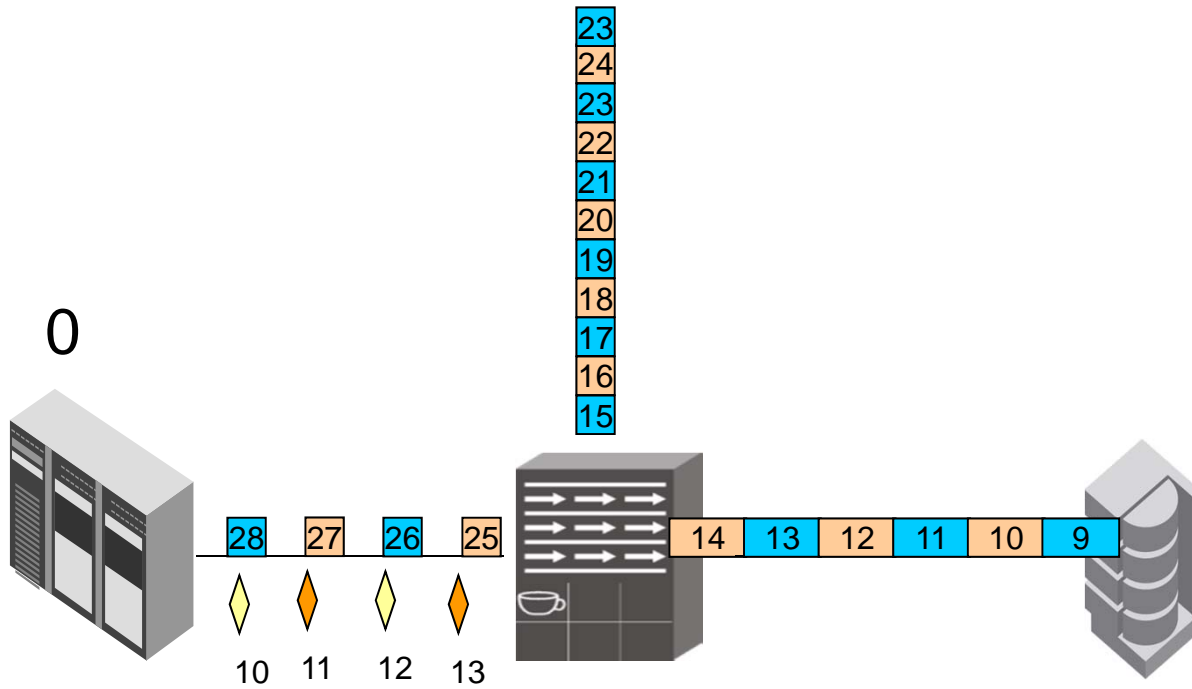


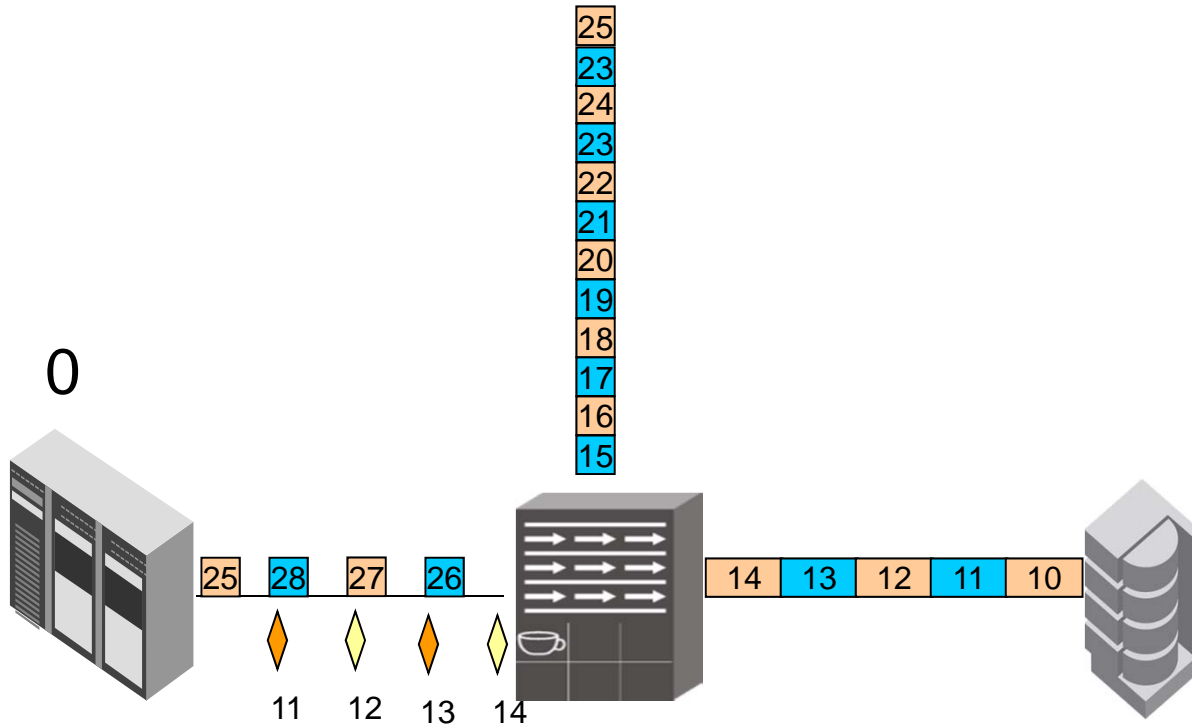










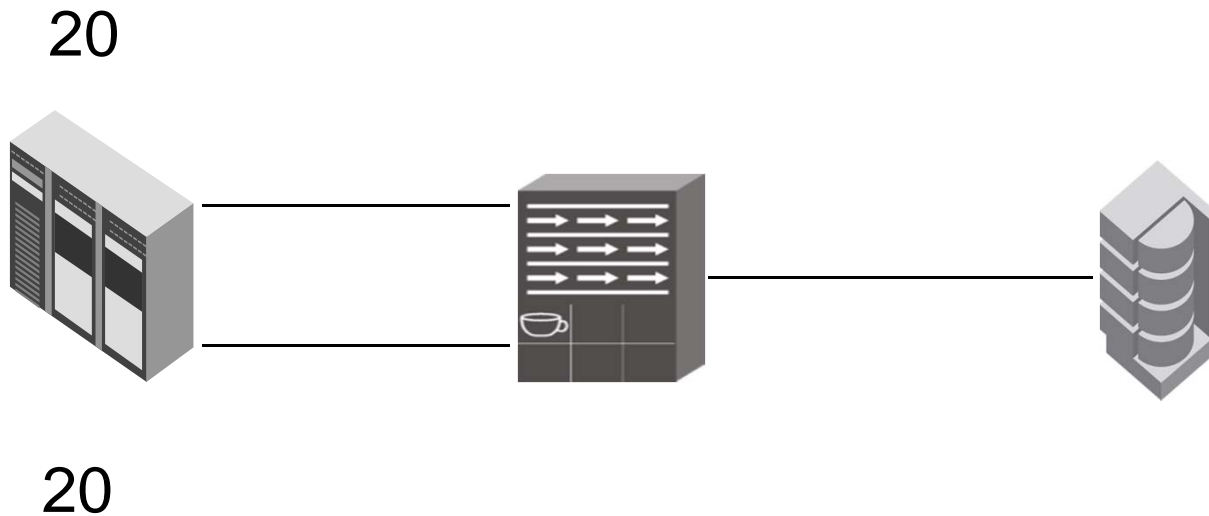


**THIS PAGE INTENTIONALLY
LEFT BLANK**

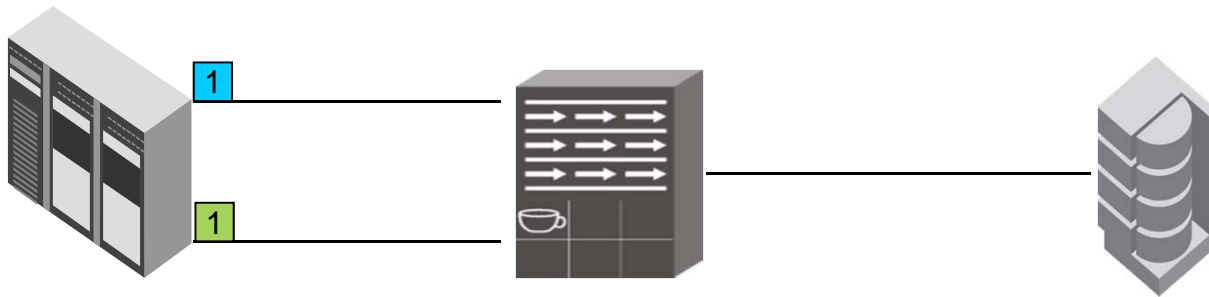
Example: Real Life?

BUFFER CREDITS

More real life example: Two senders sending at 30% - 50% link rate to one receiver

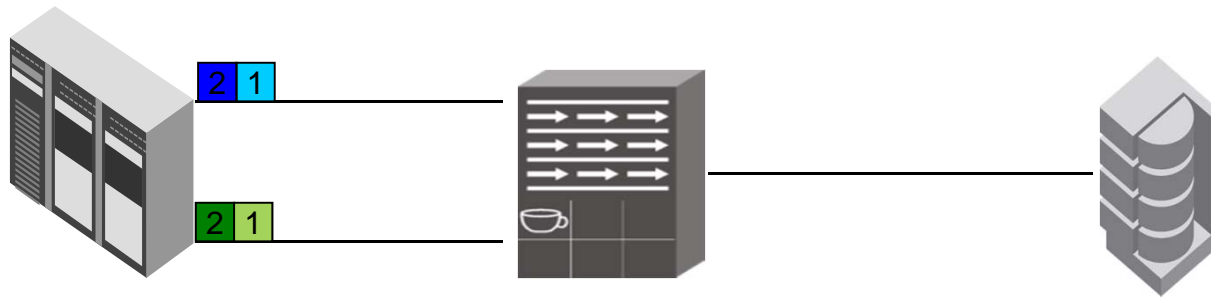


19



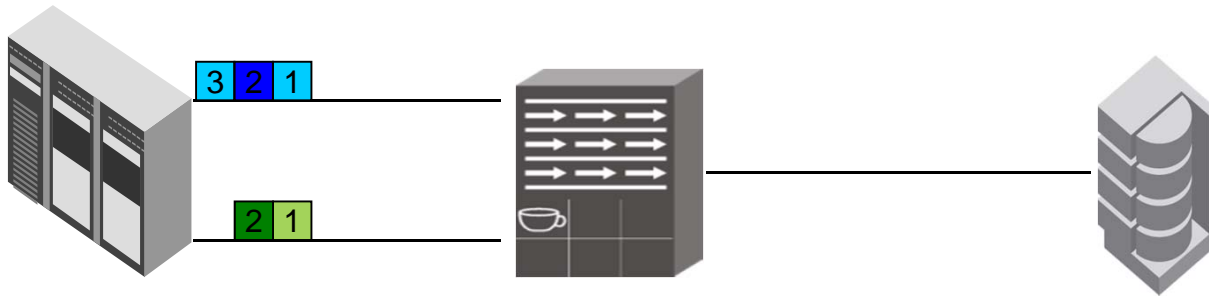
19

18

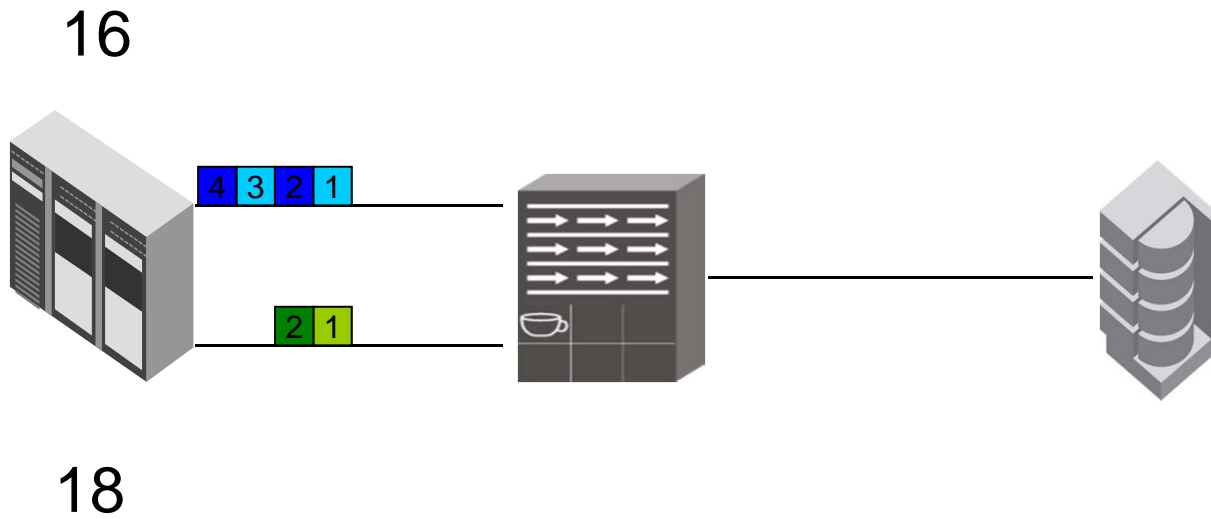


18

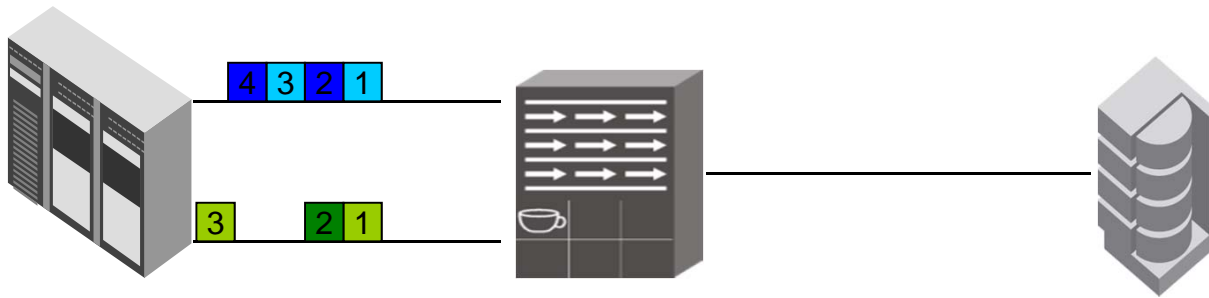
17



18

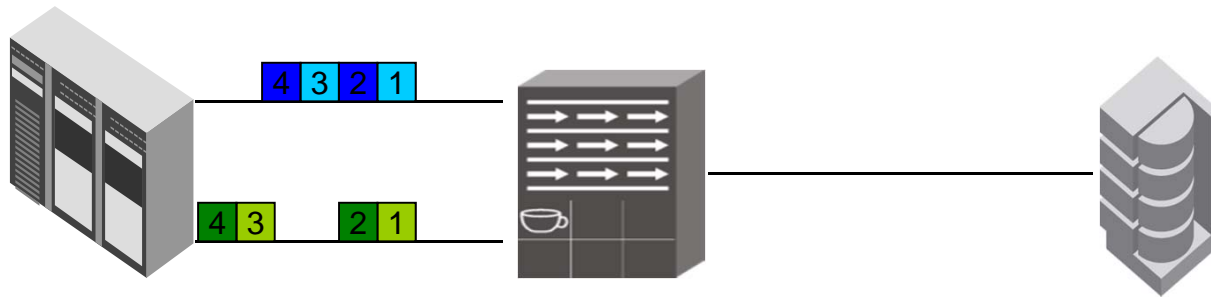


16

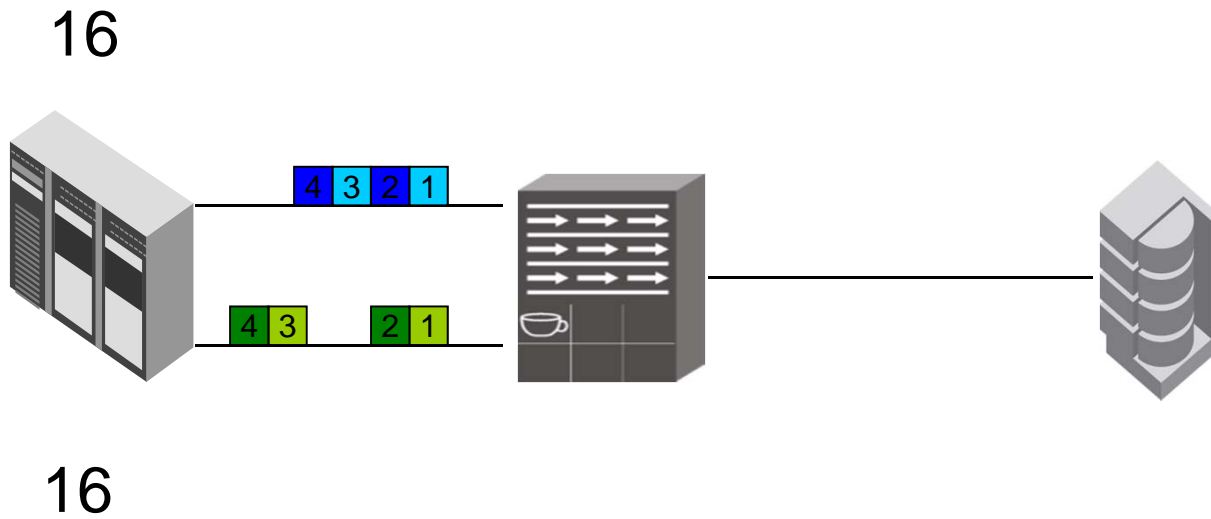


17

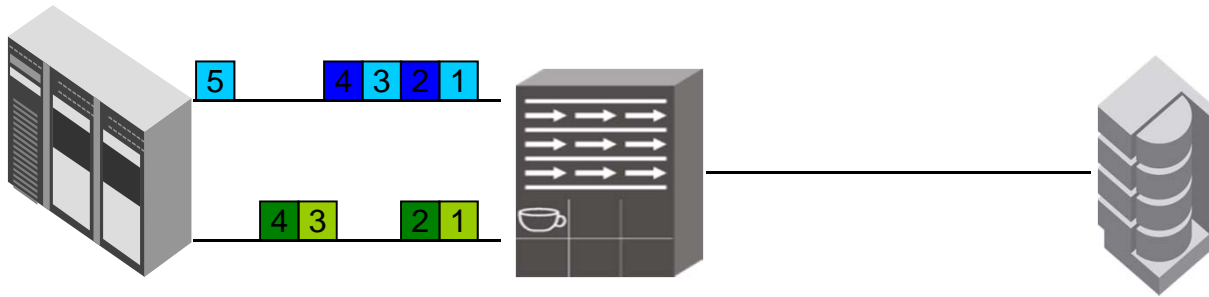
16



16

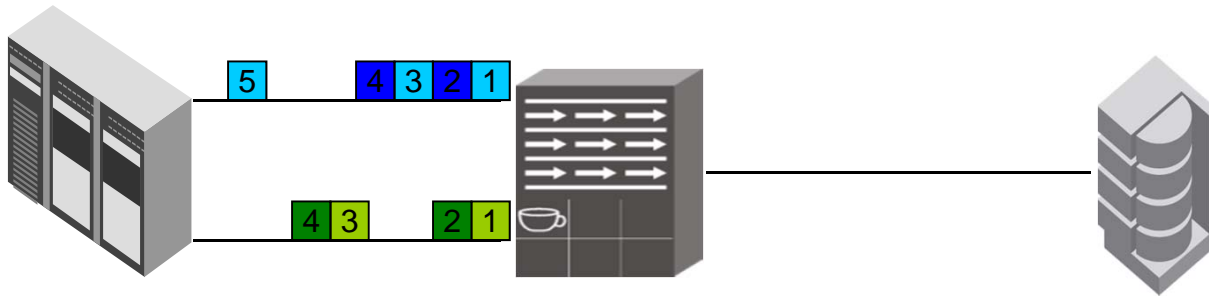


15

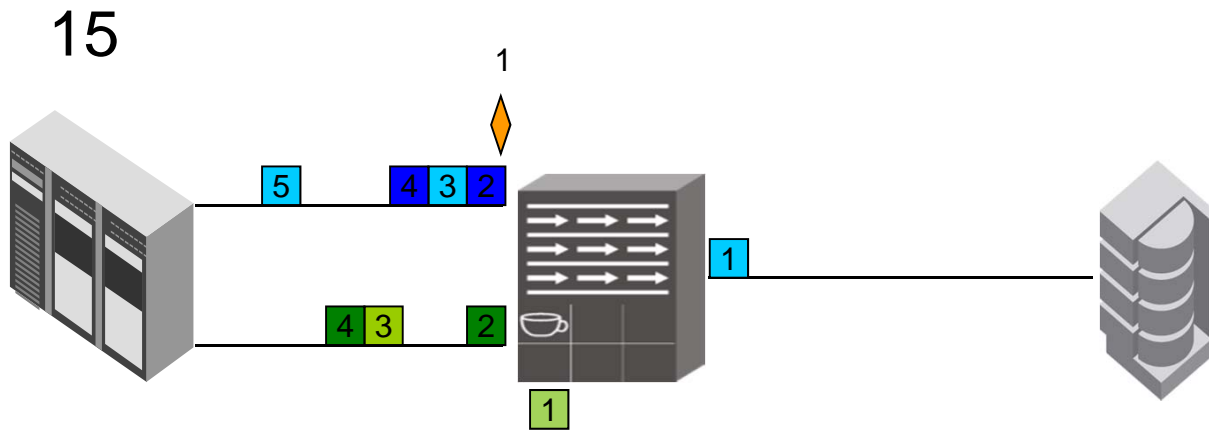


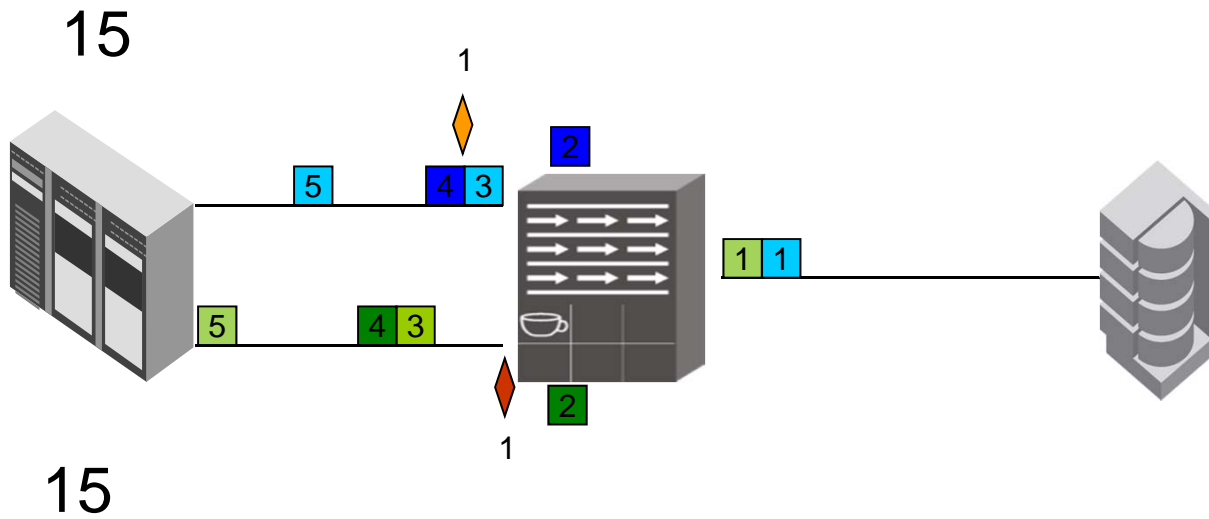
16

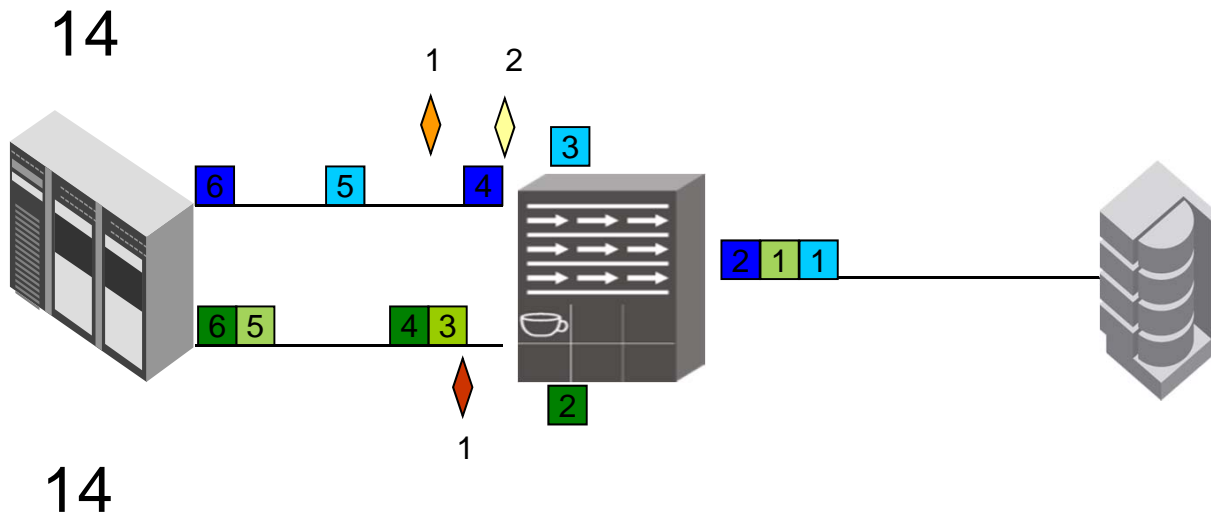
15

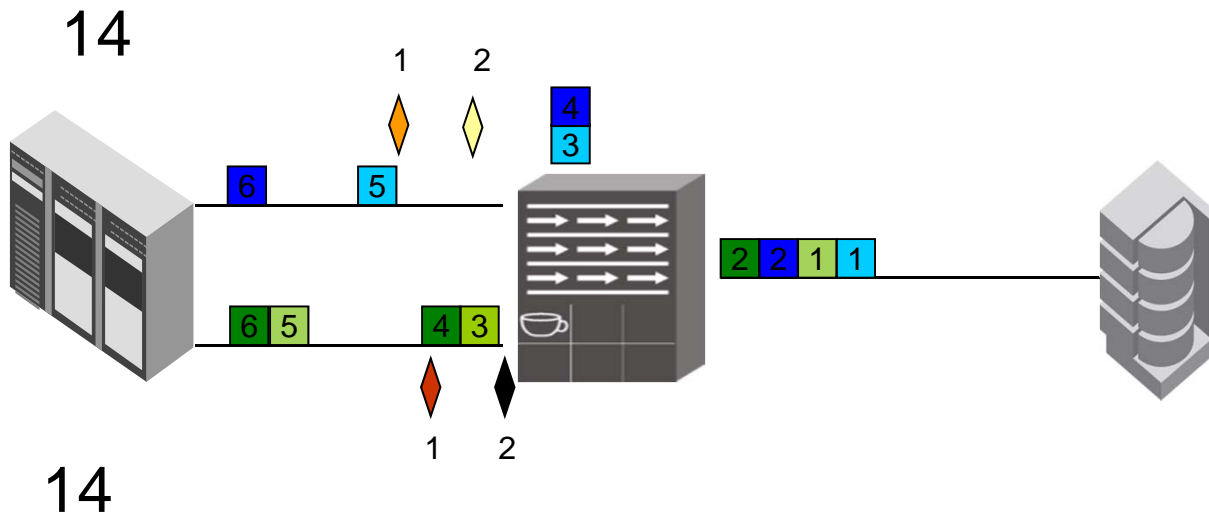


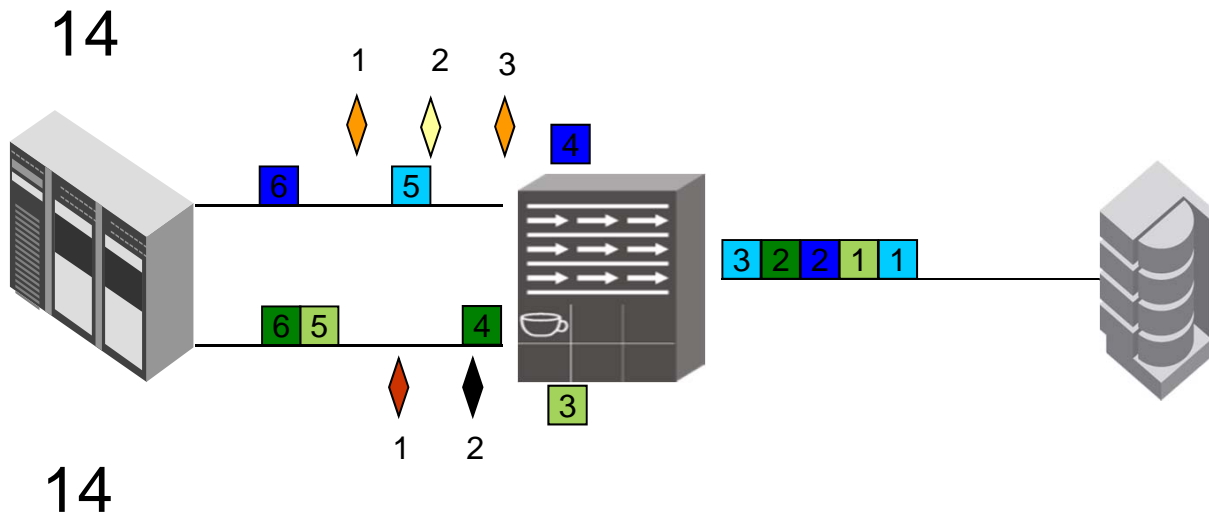
16

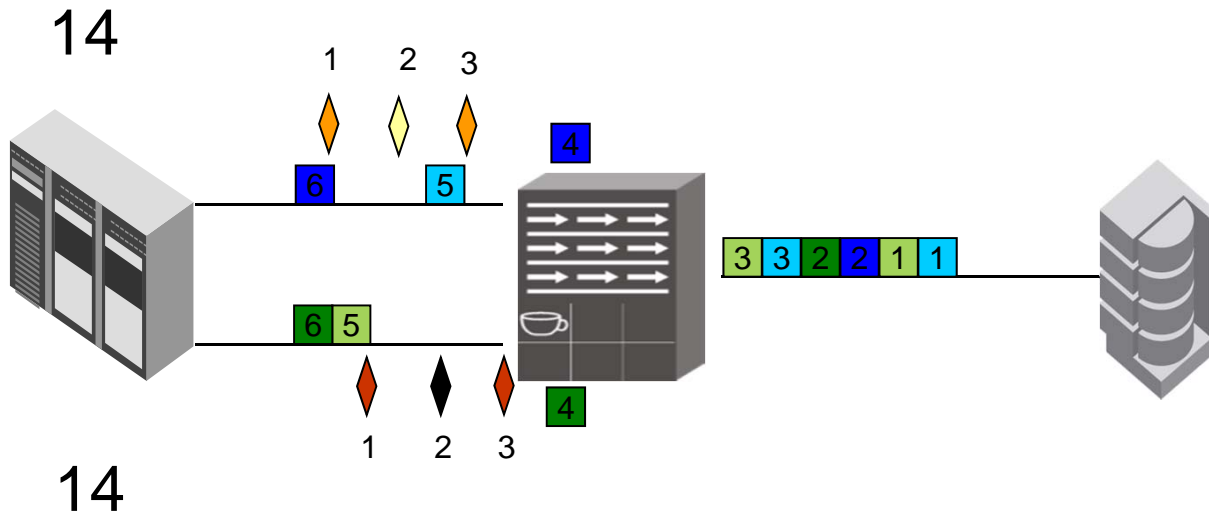


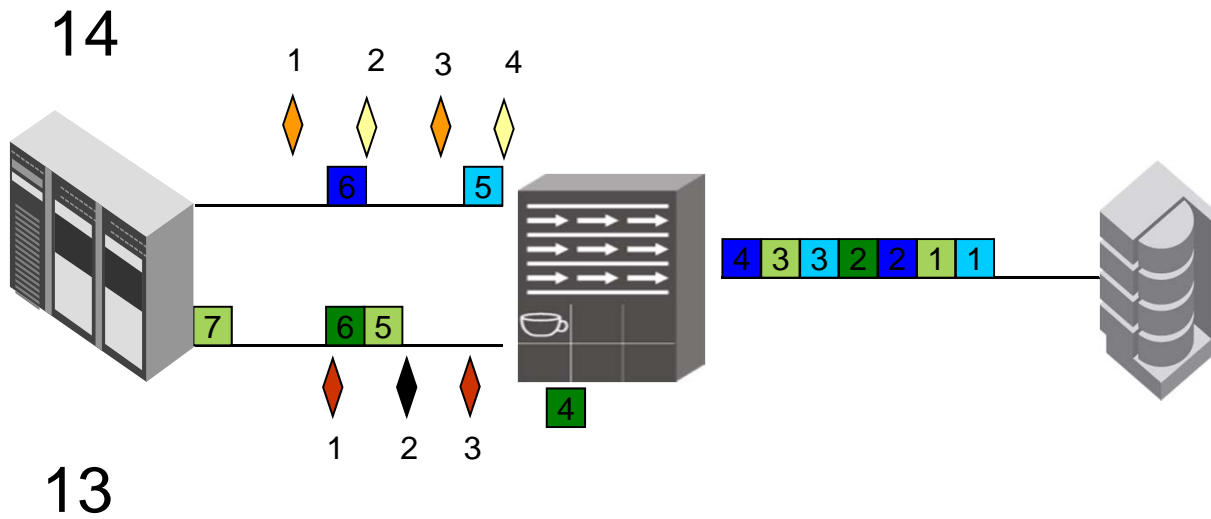


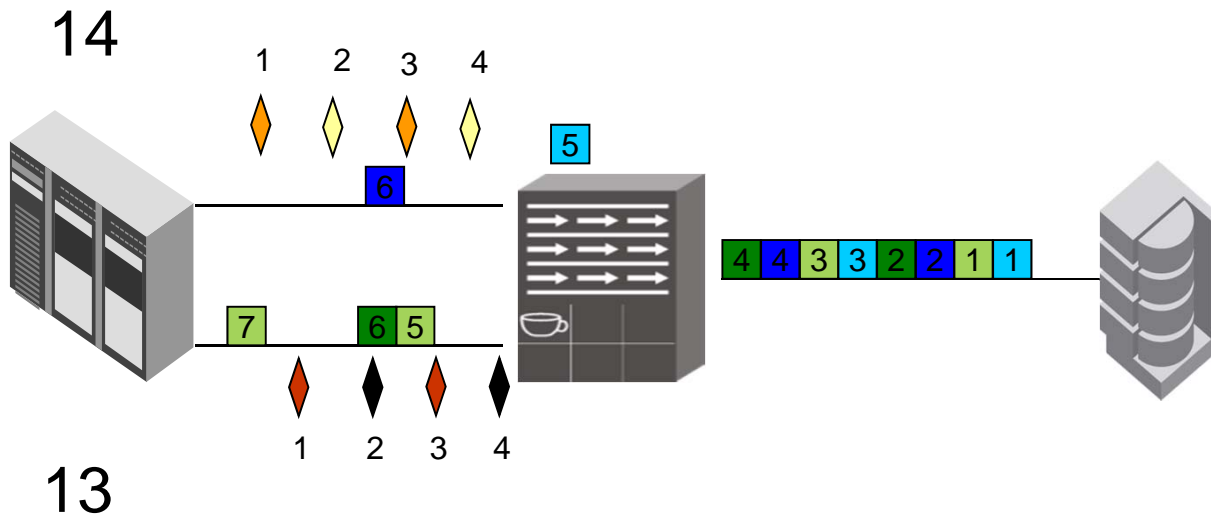


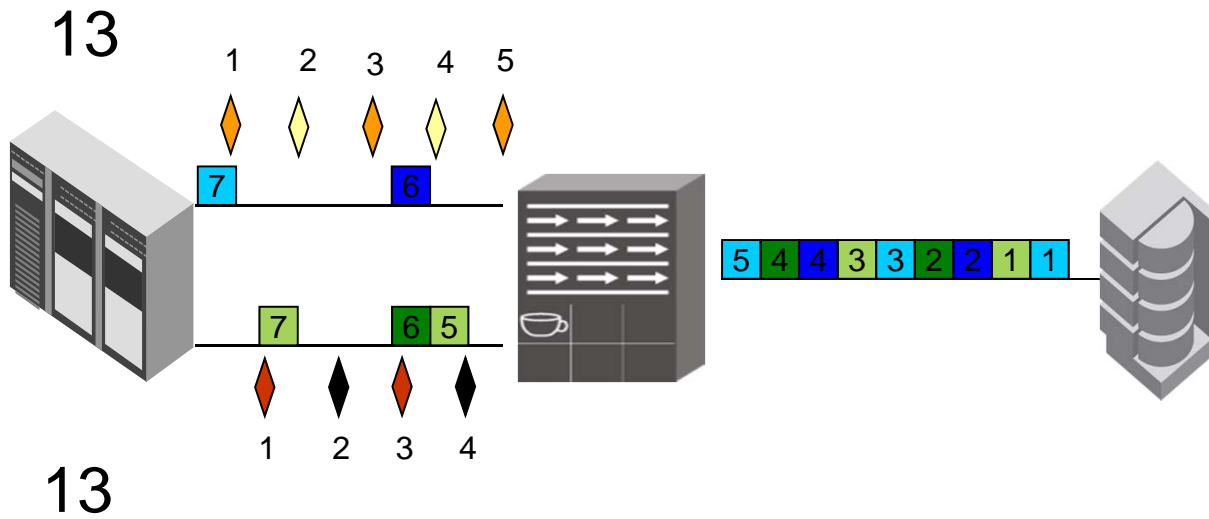


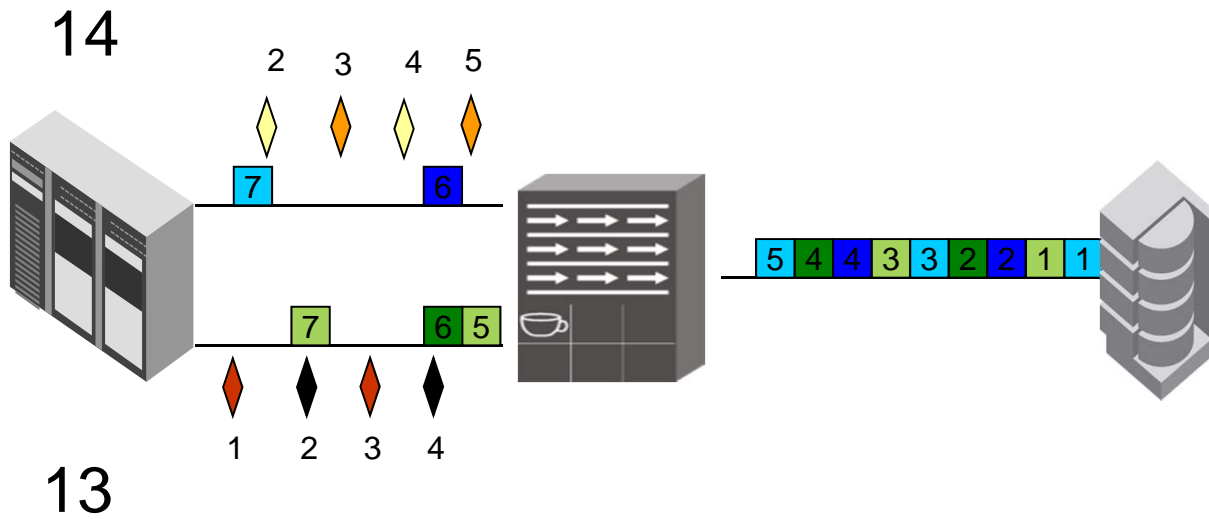


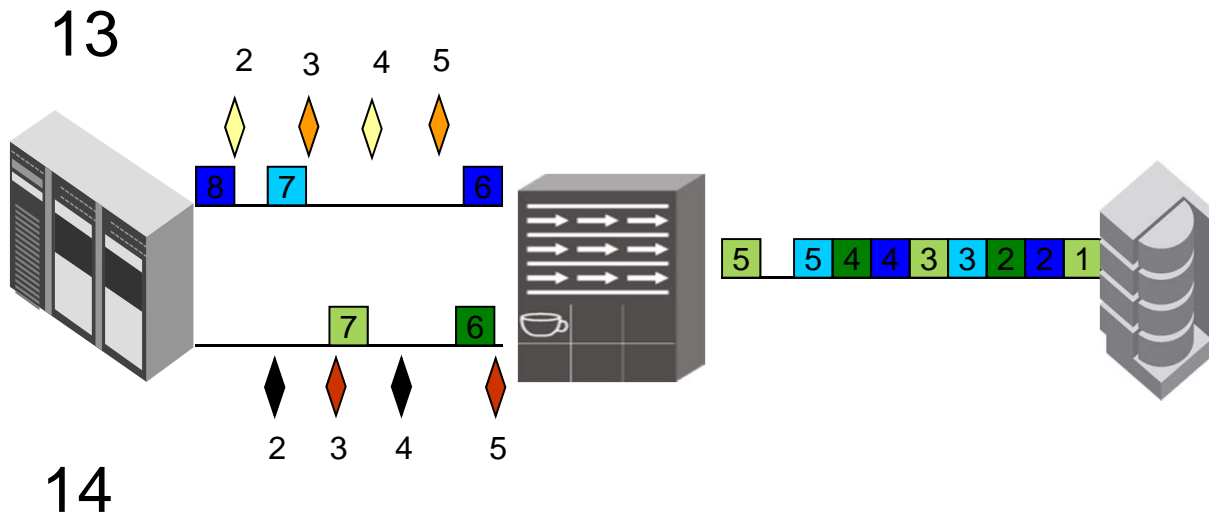


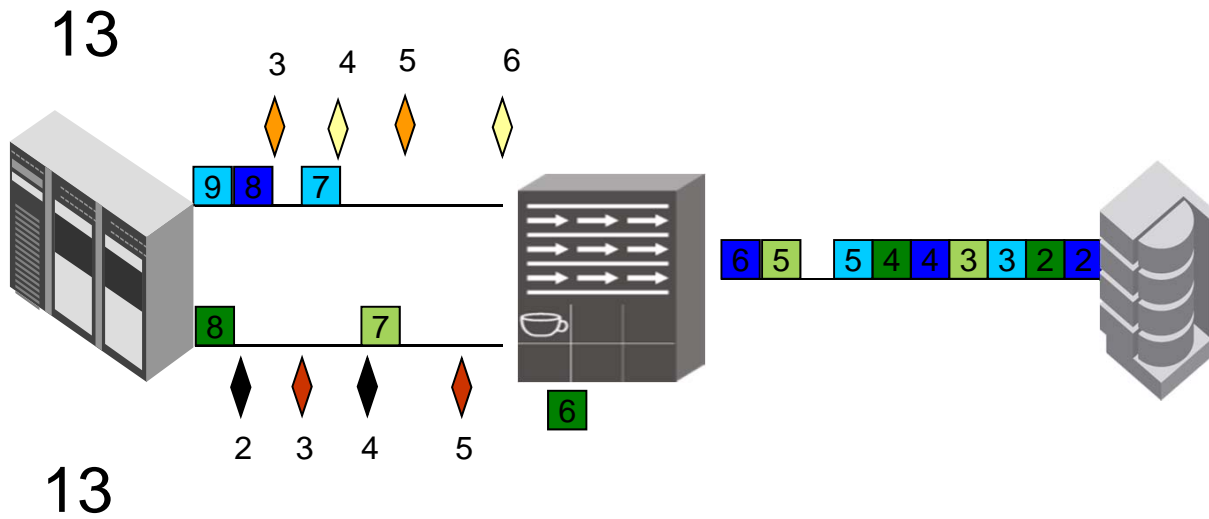


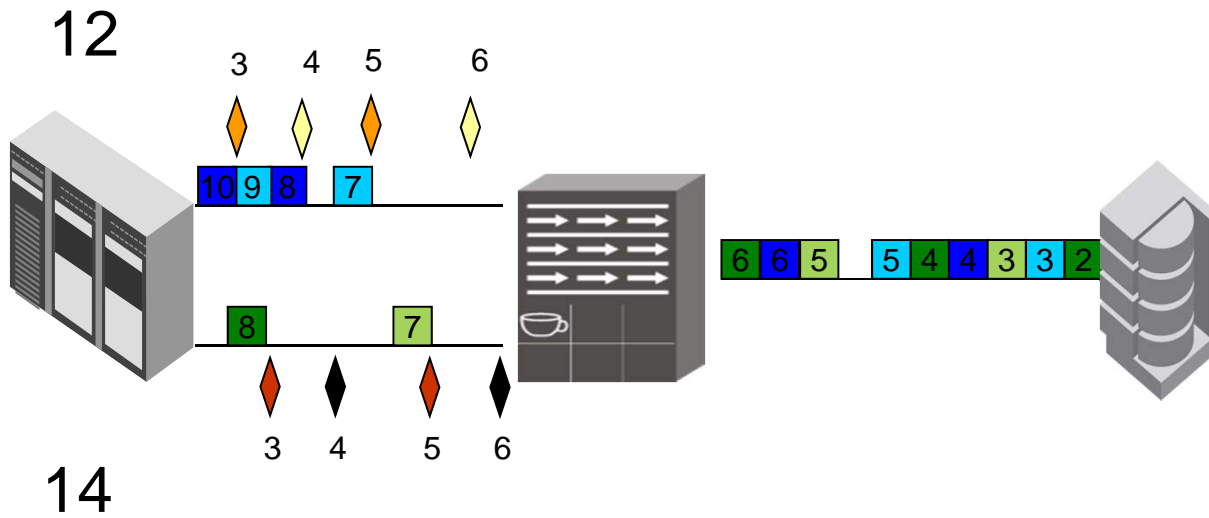


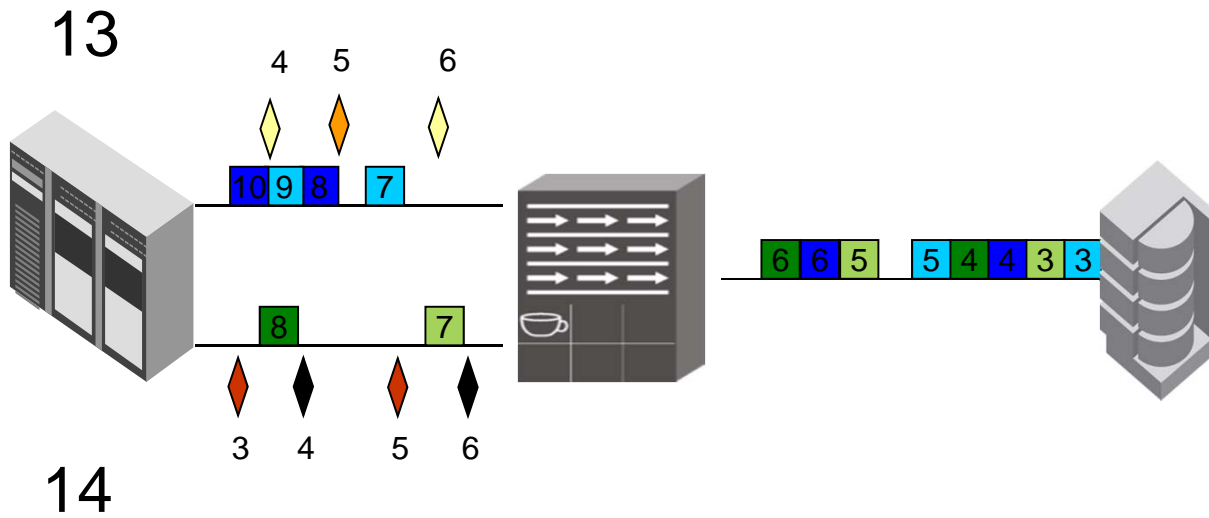


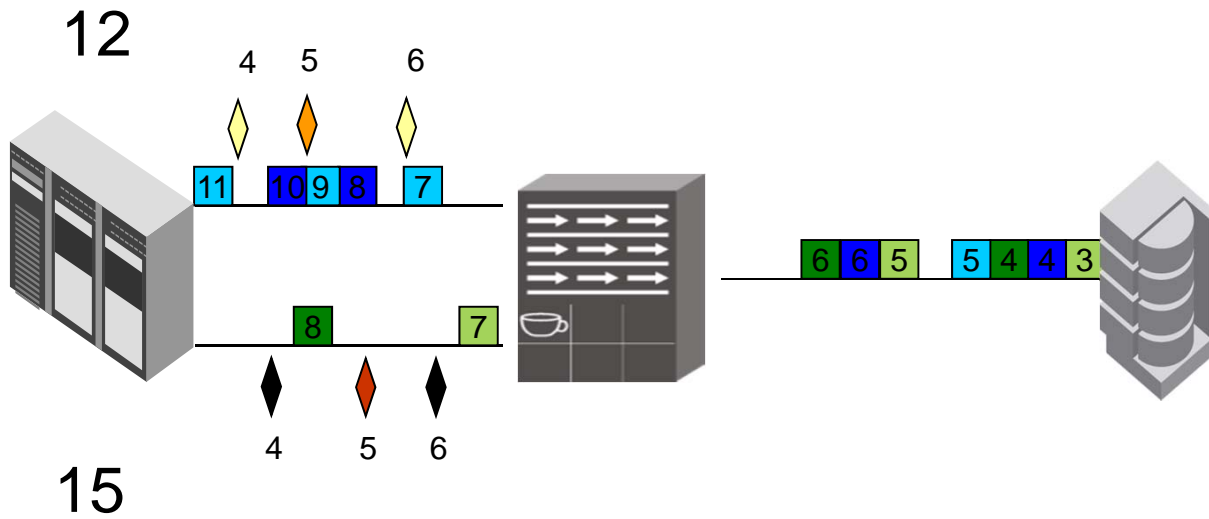


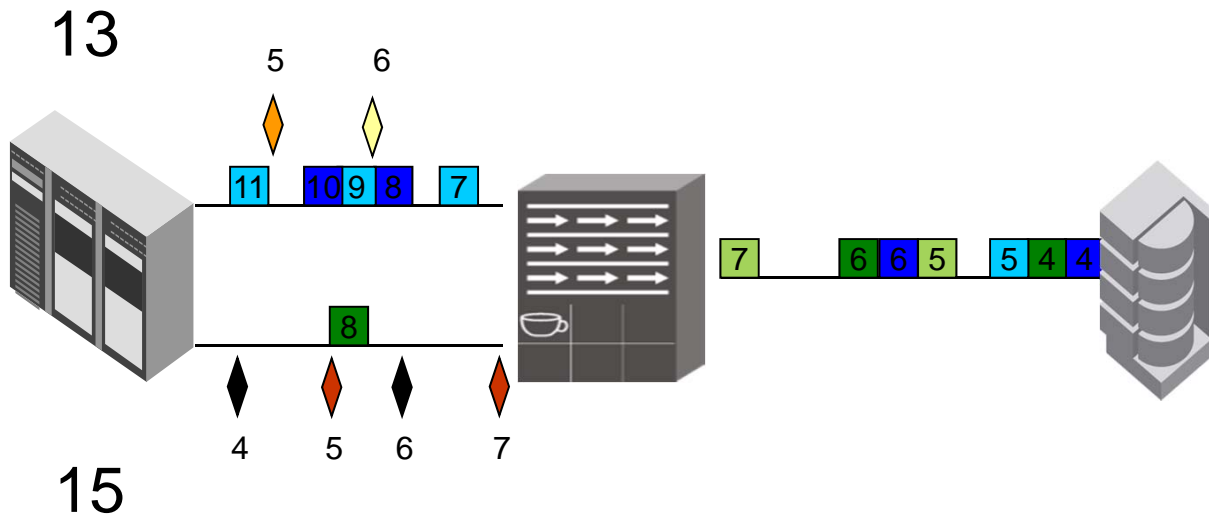


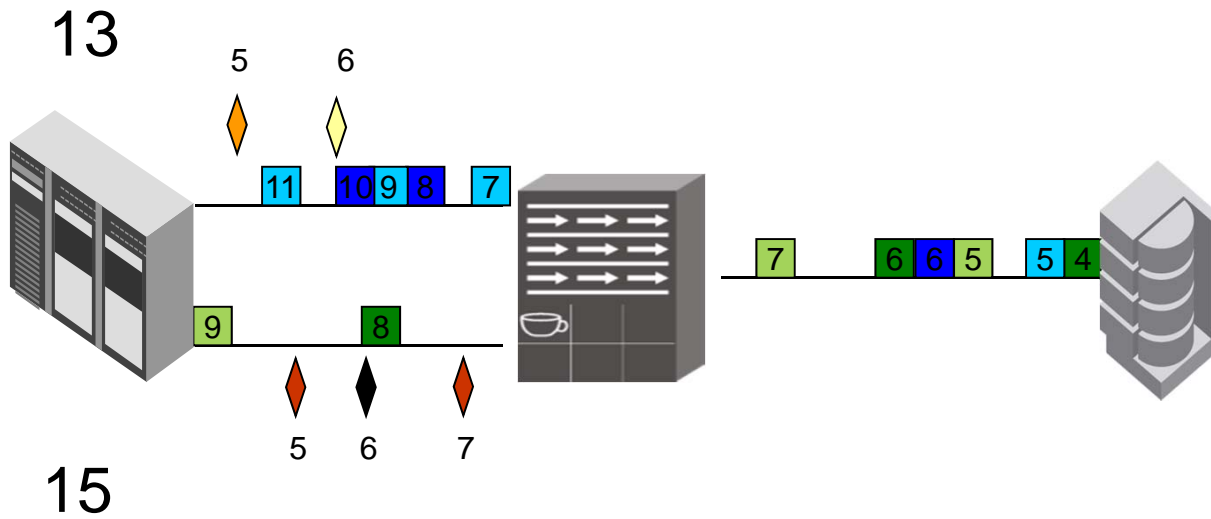


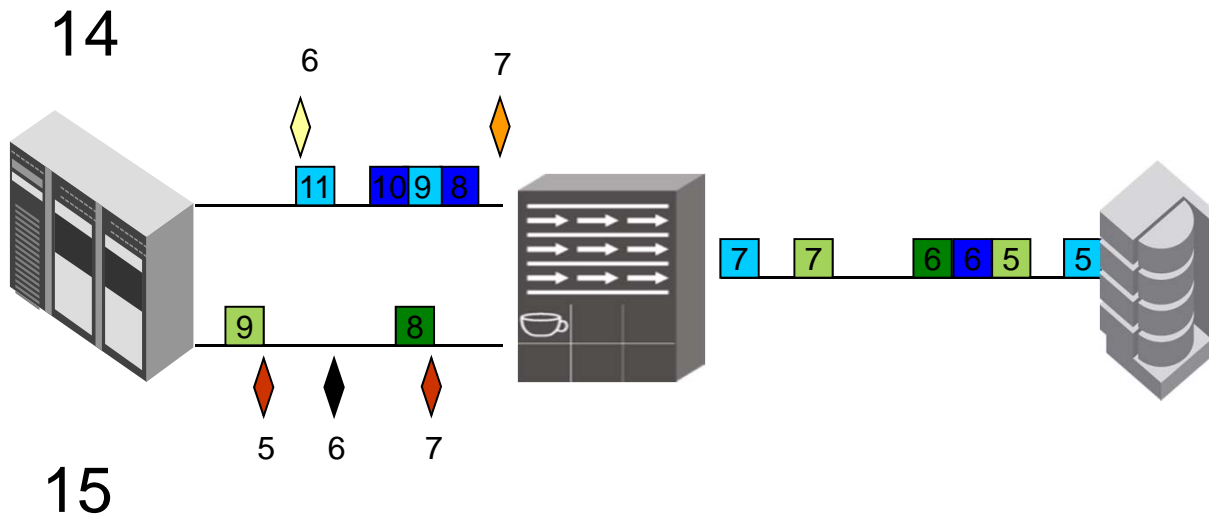


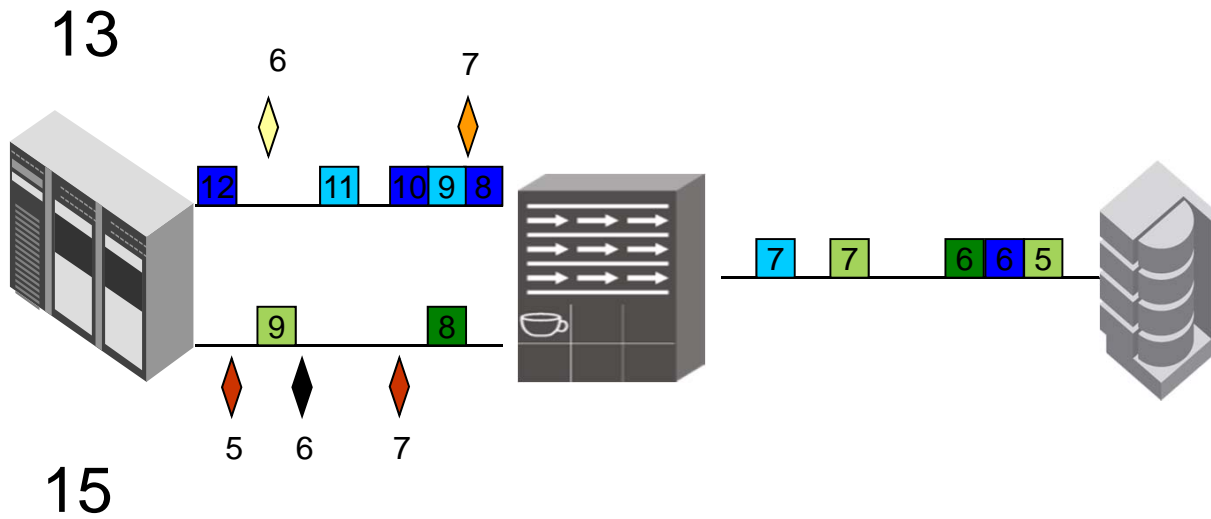


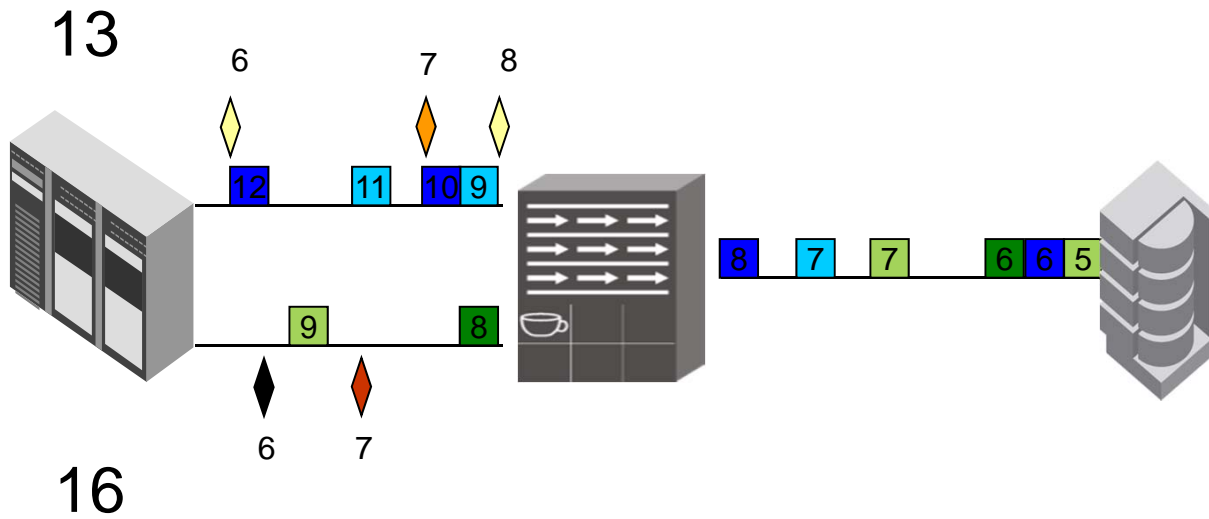


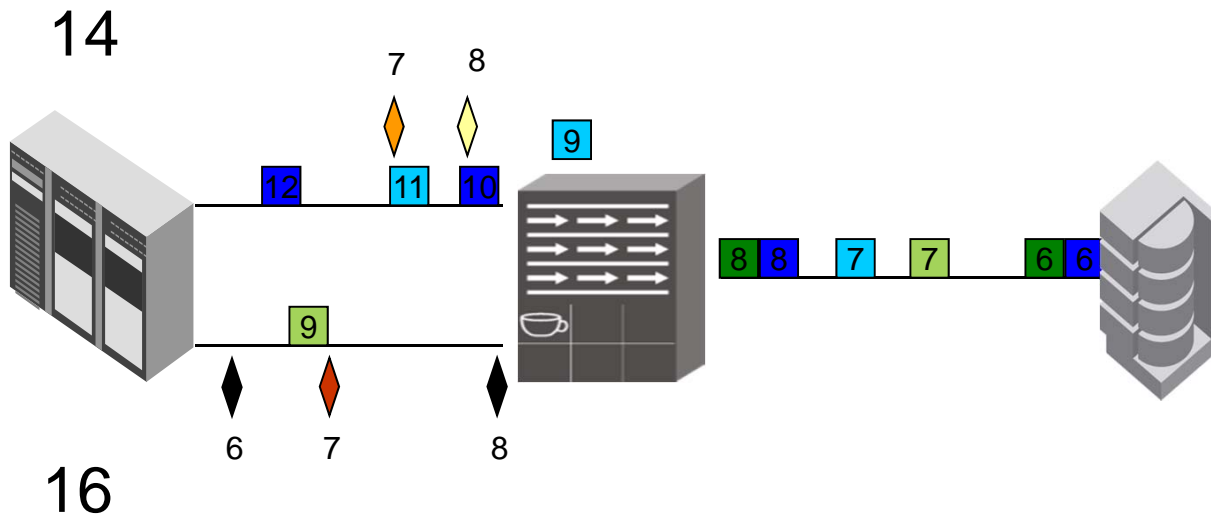






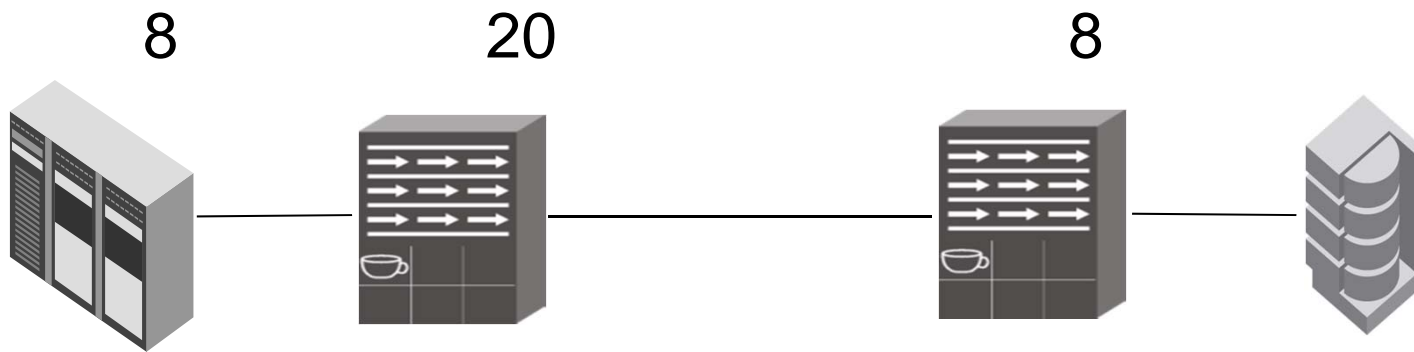


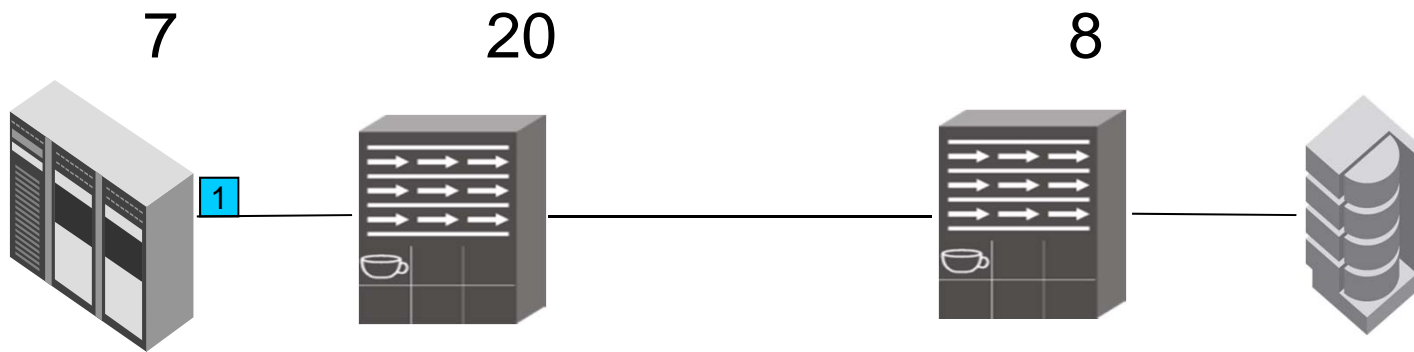


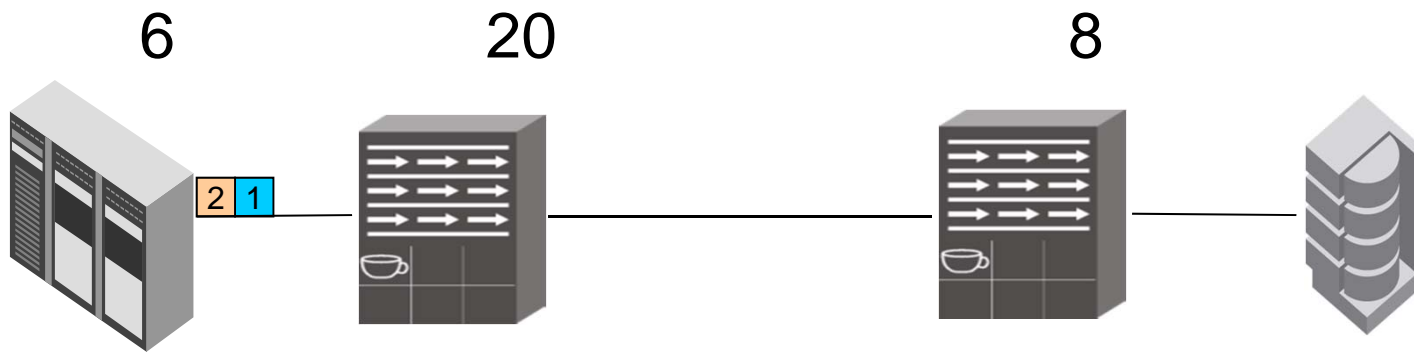


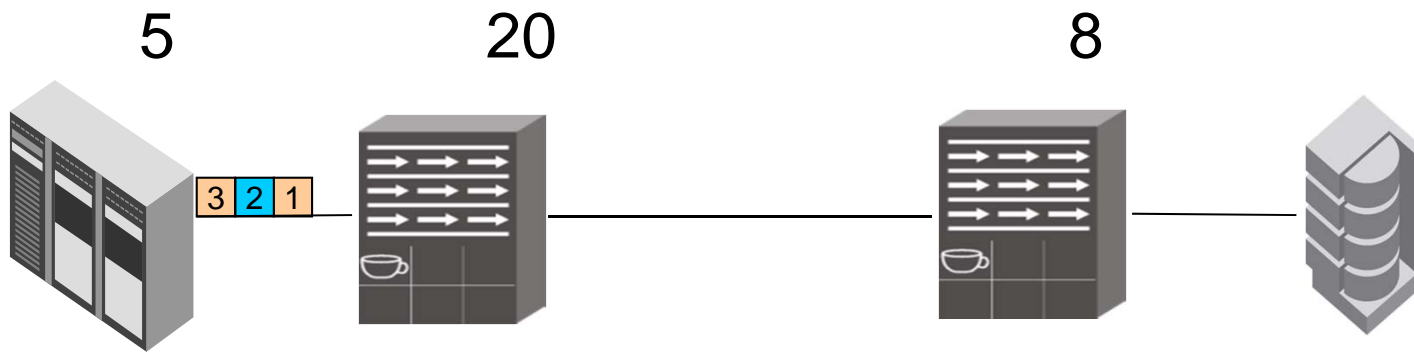
Example: Cascaded Directors

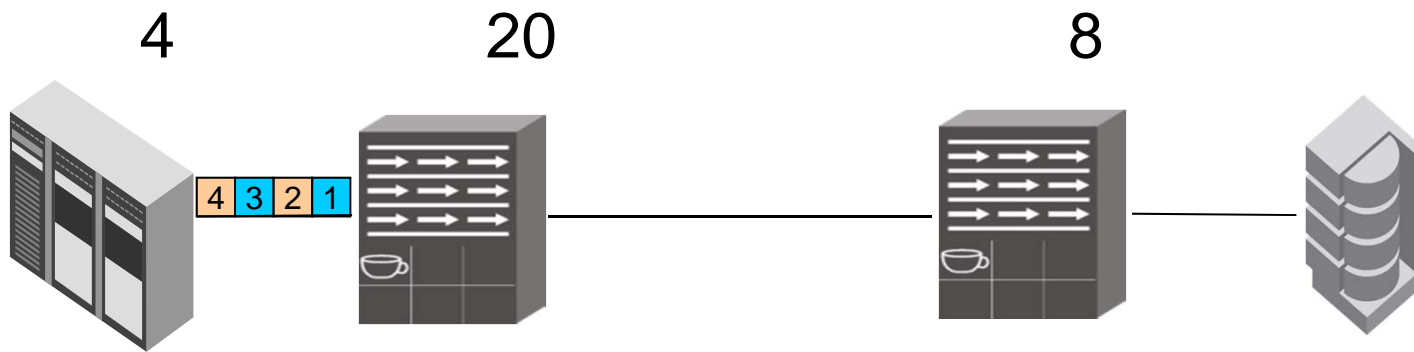
BUFFER CREDITS

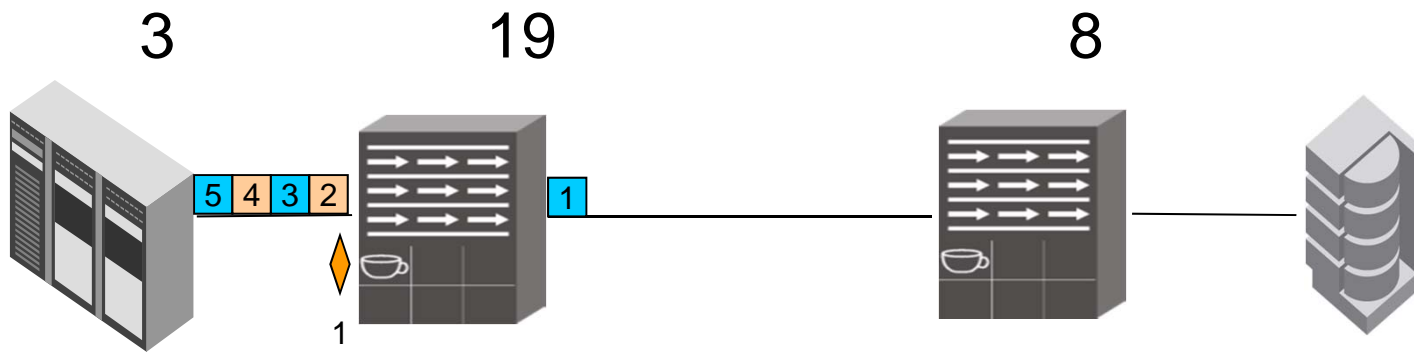


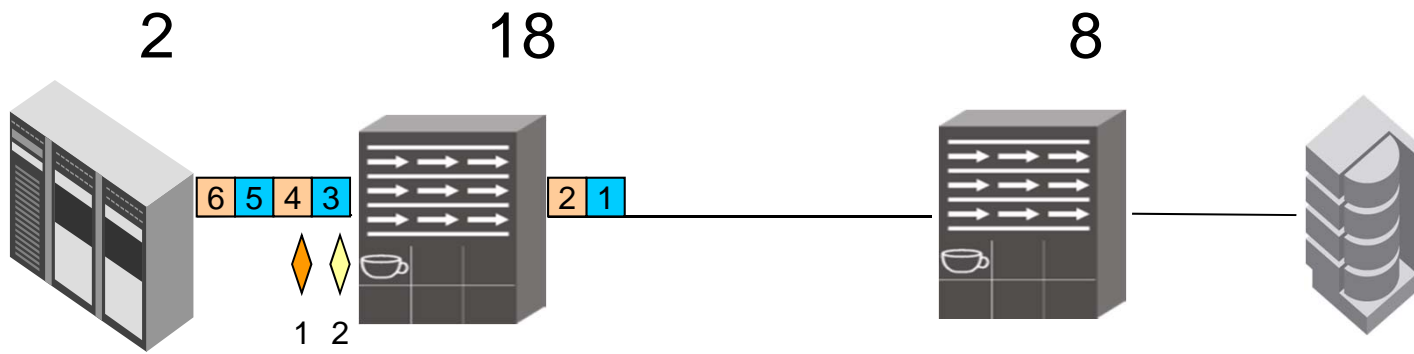


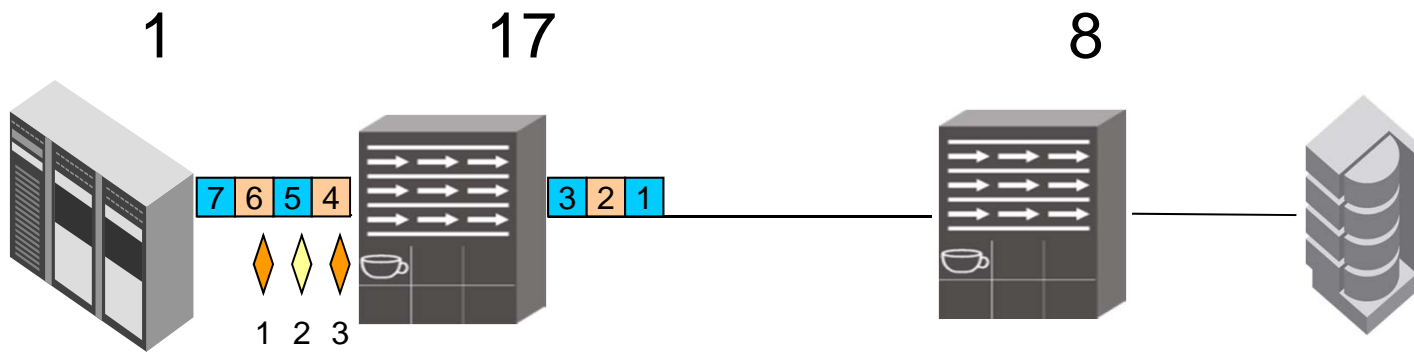


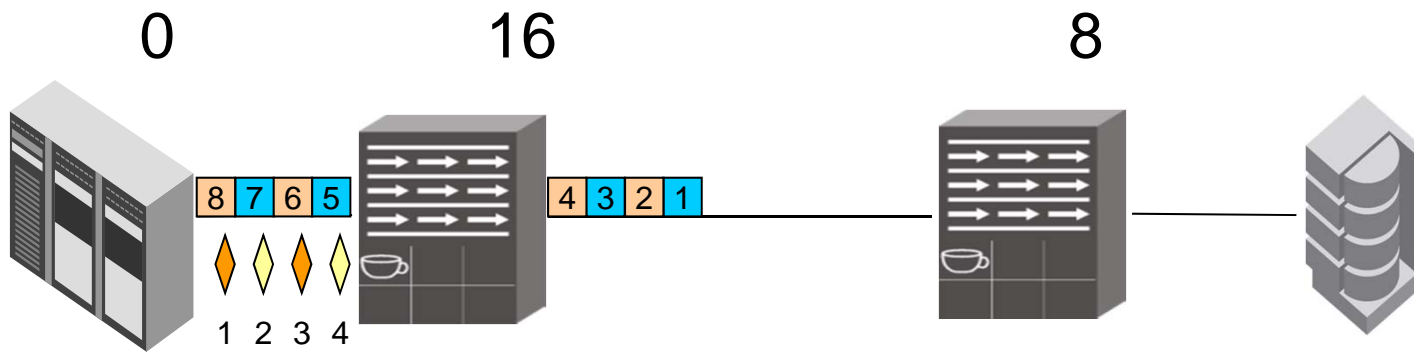


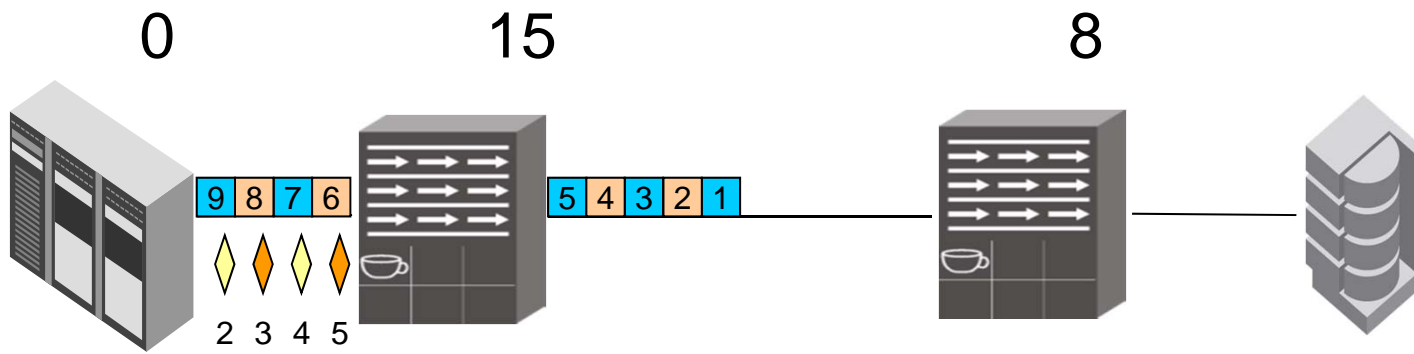


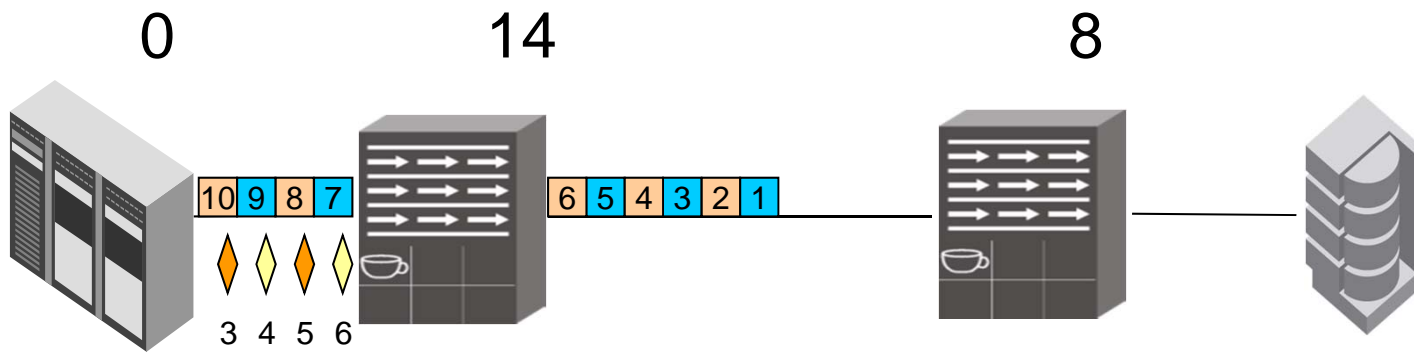


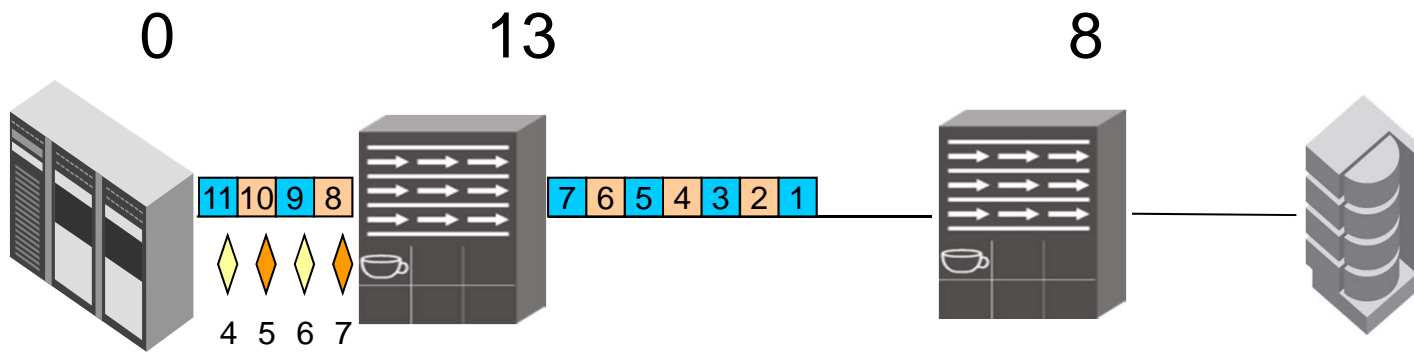


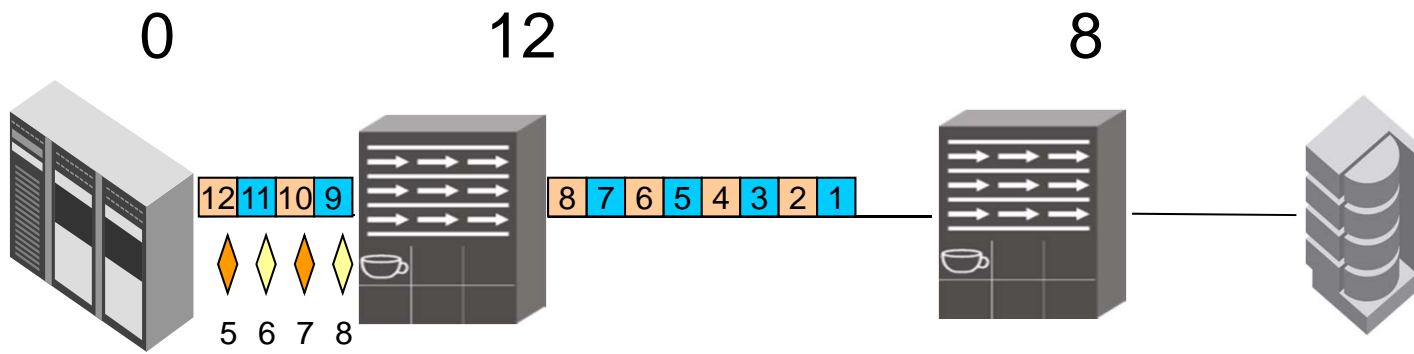


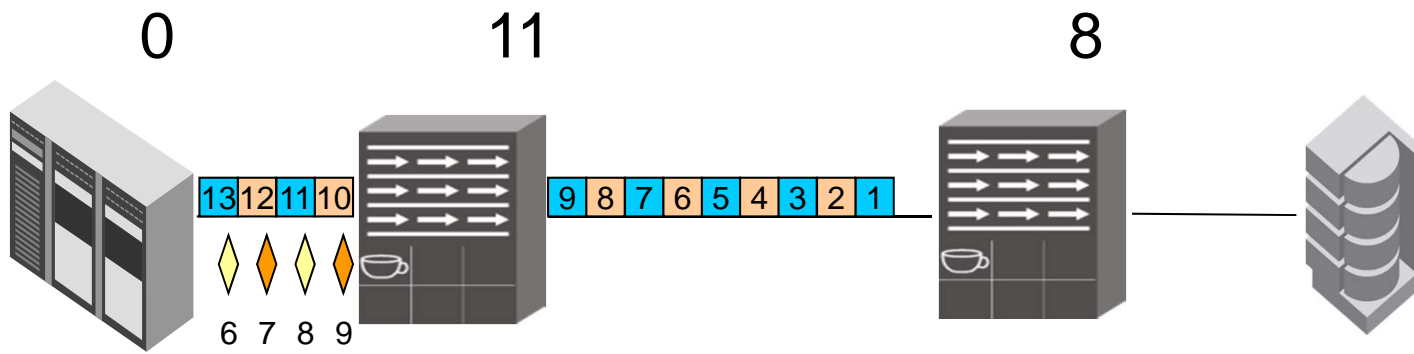


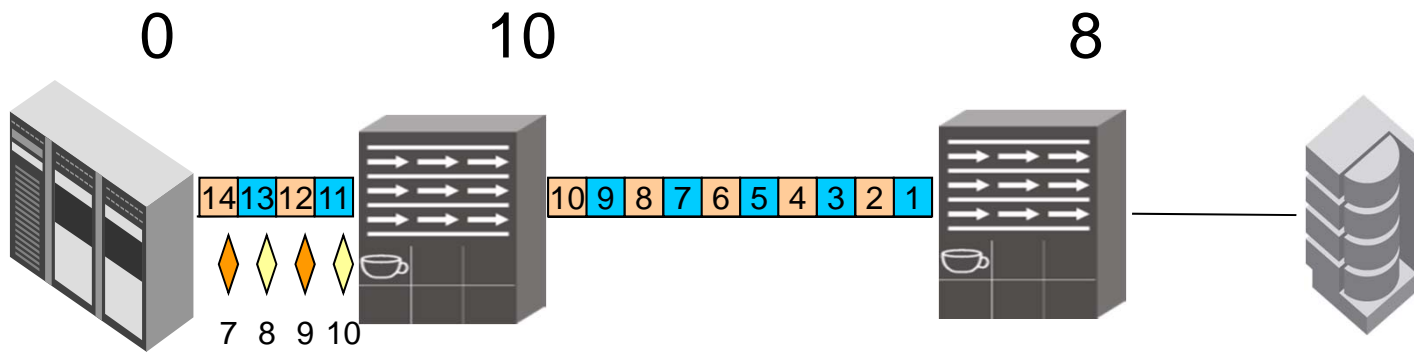


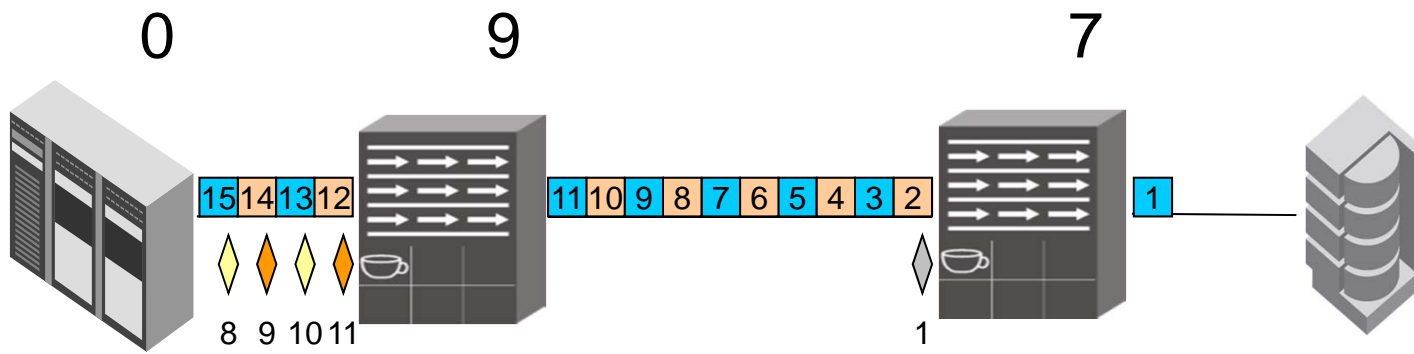


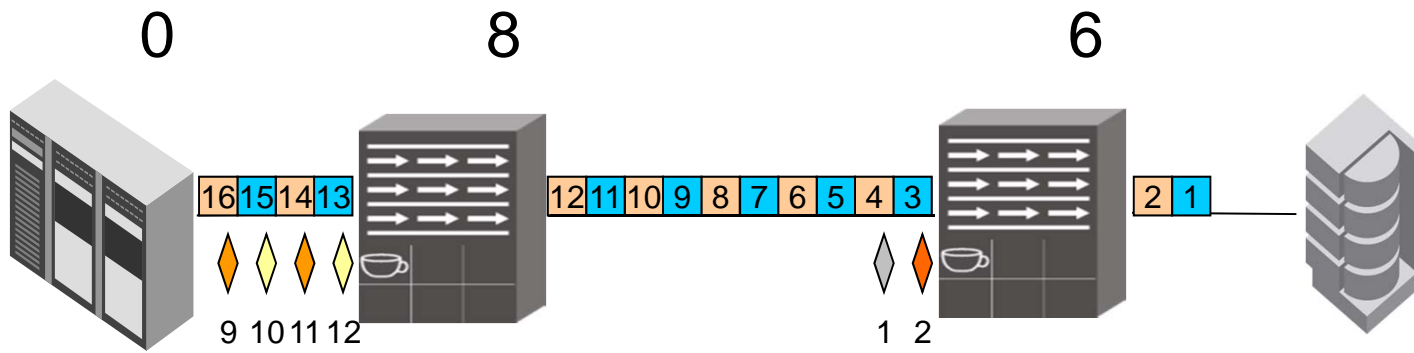


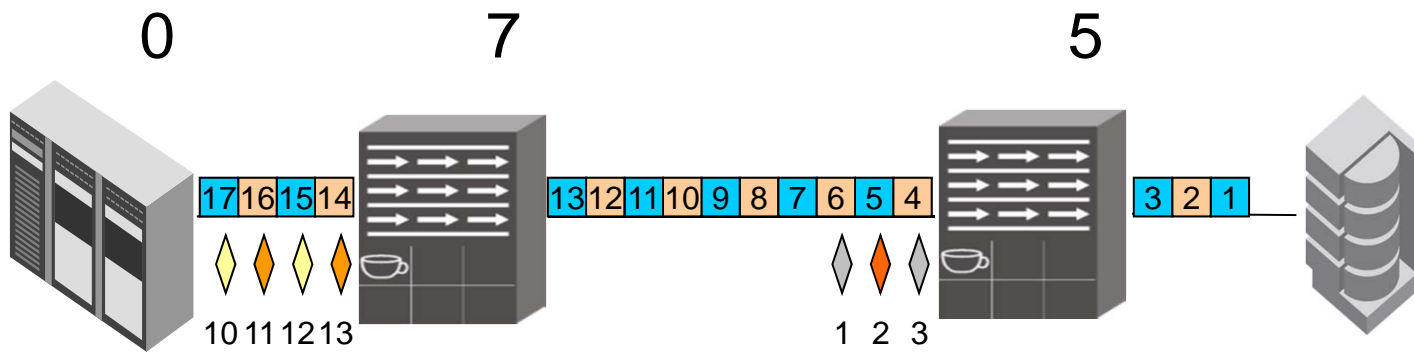


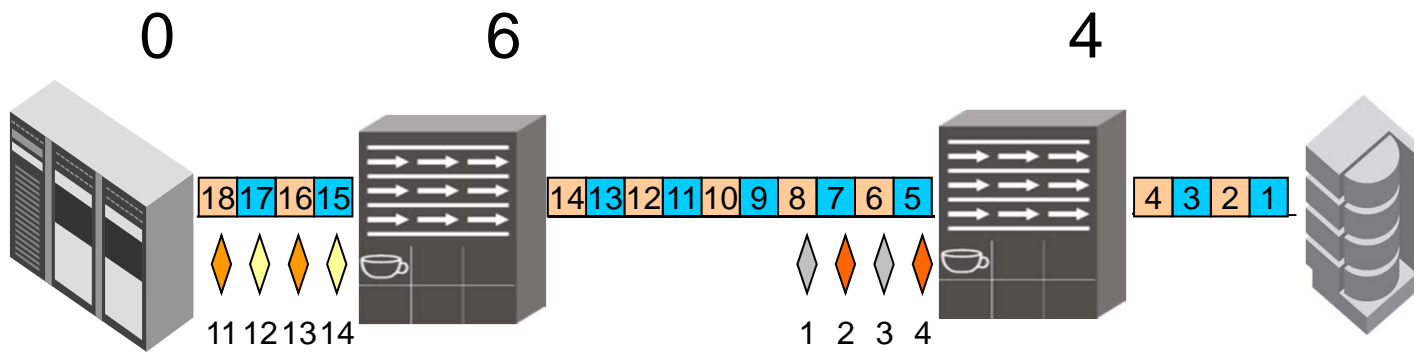


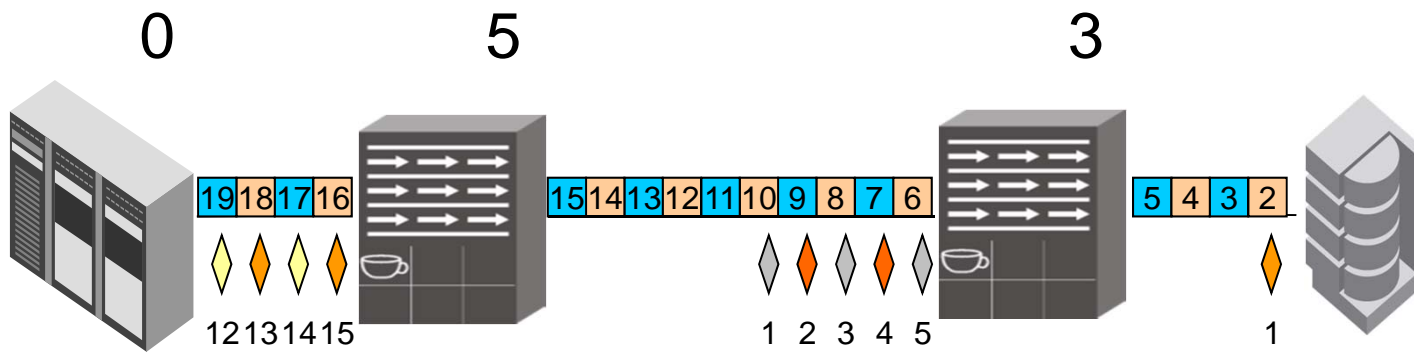


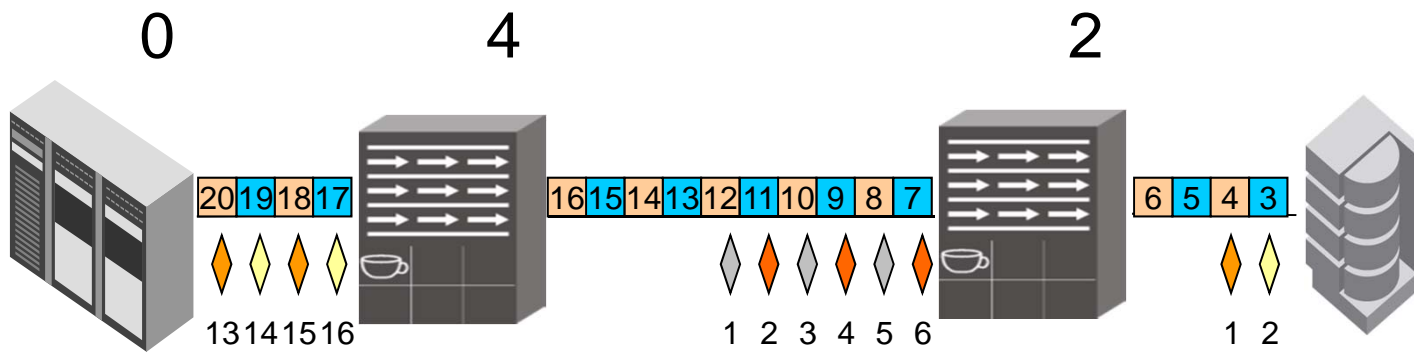


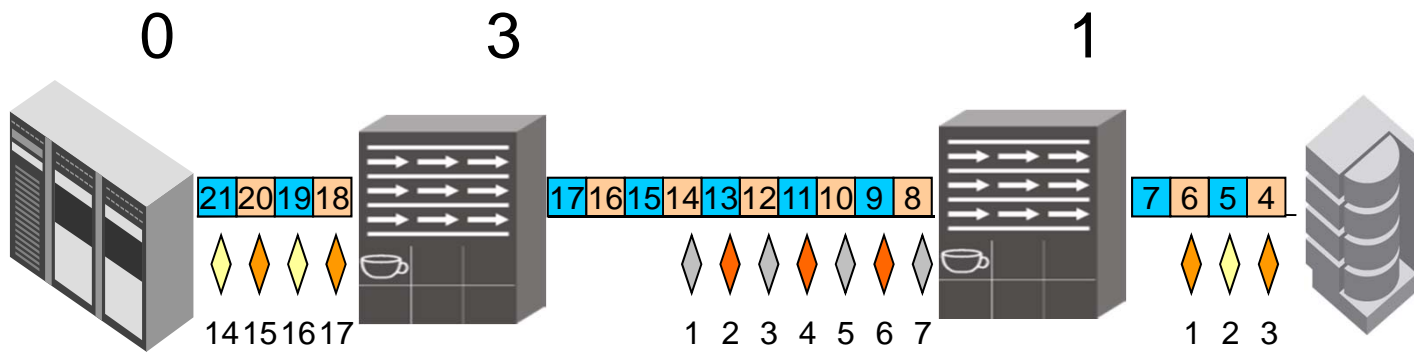


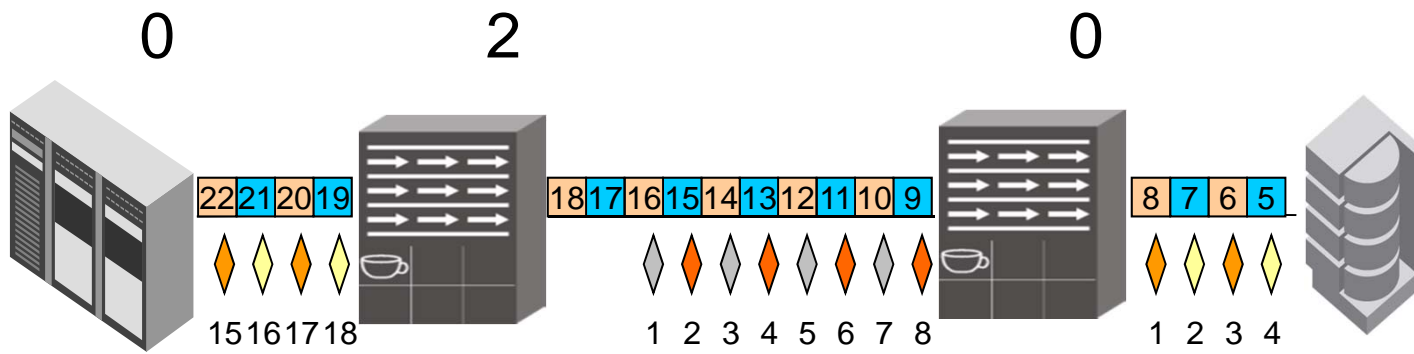


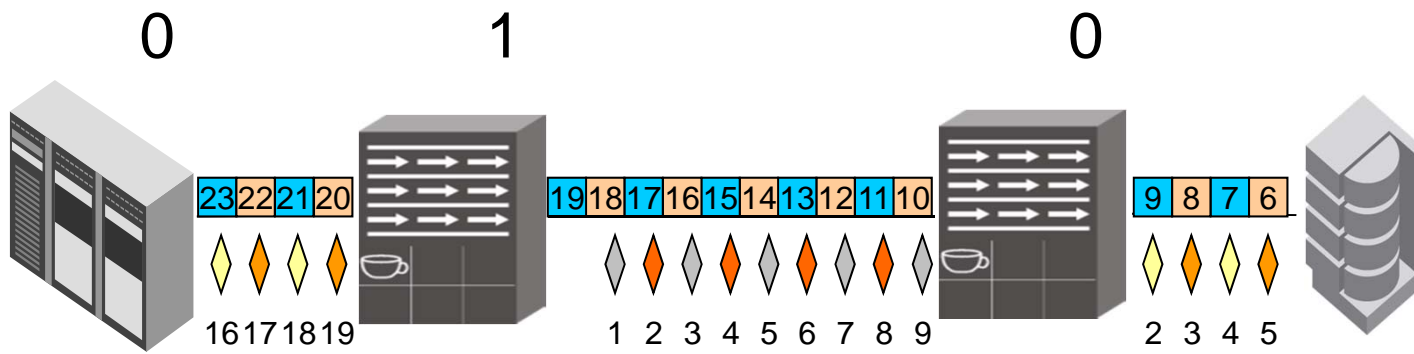


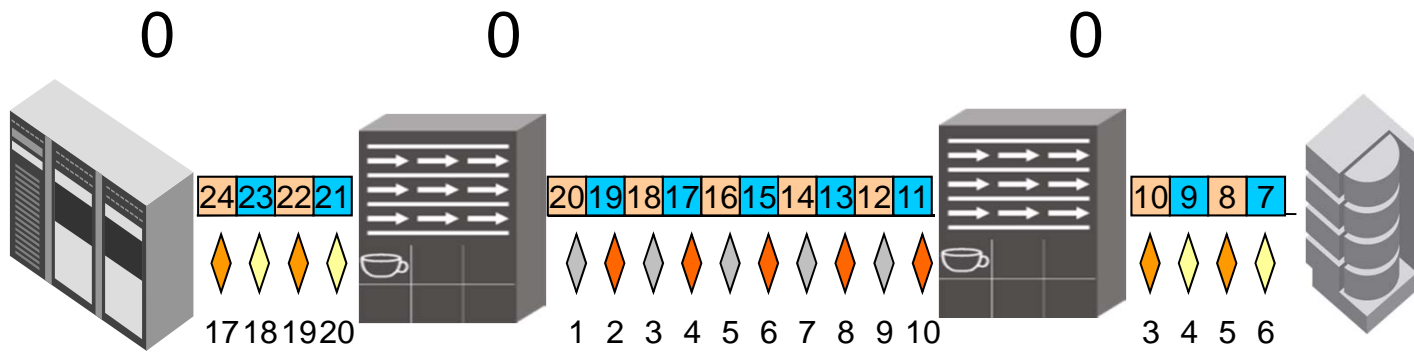


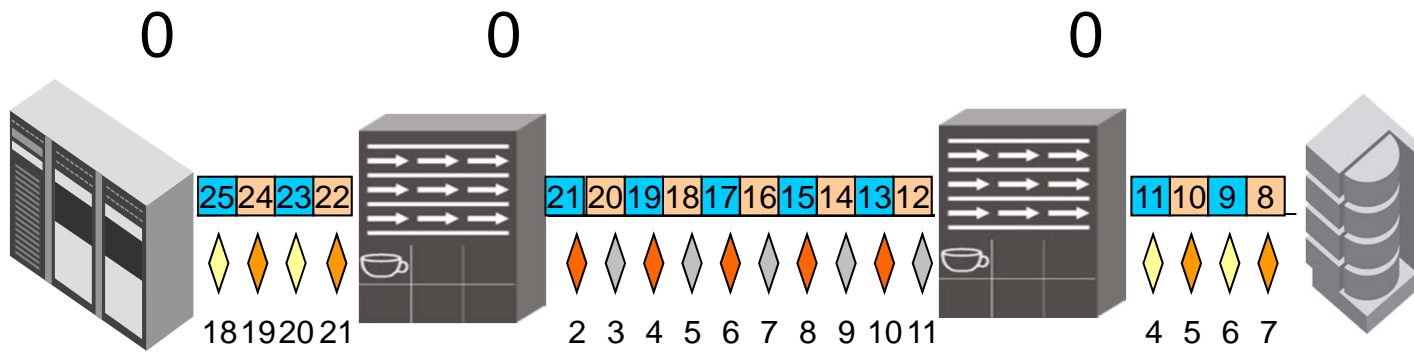










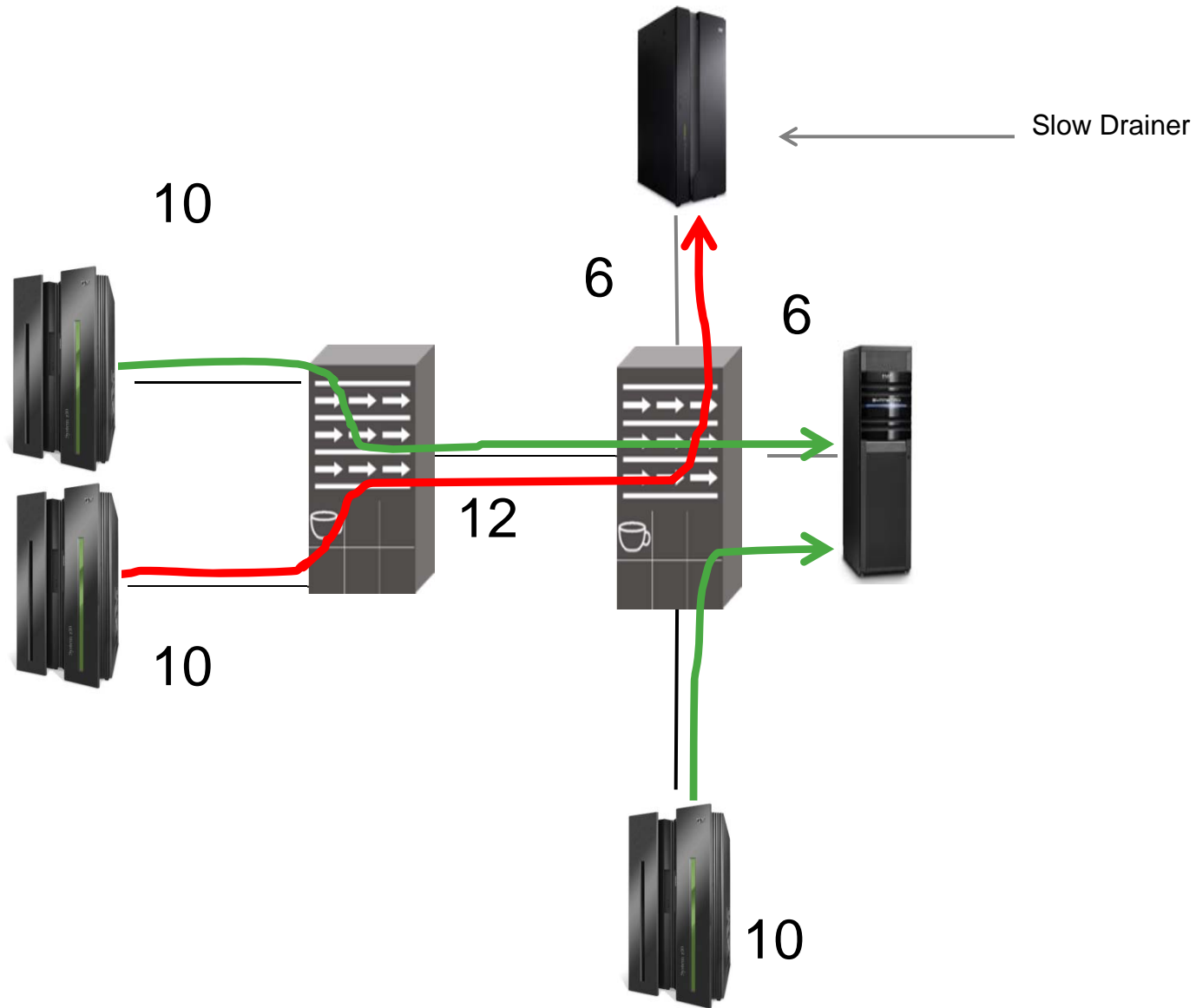


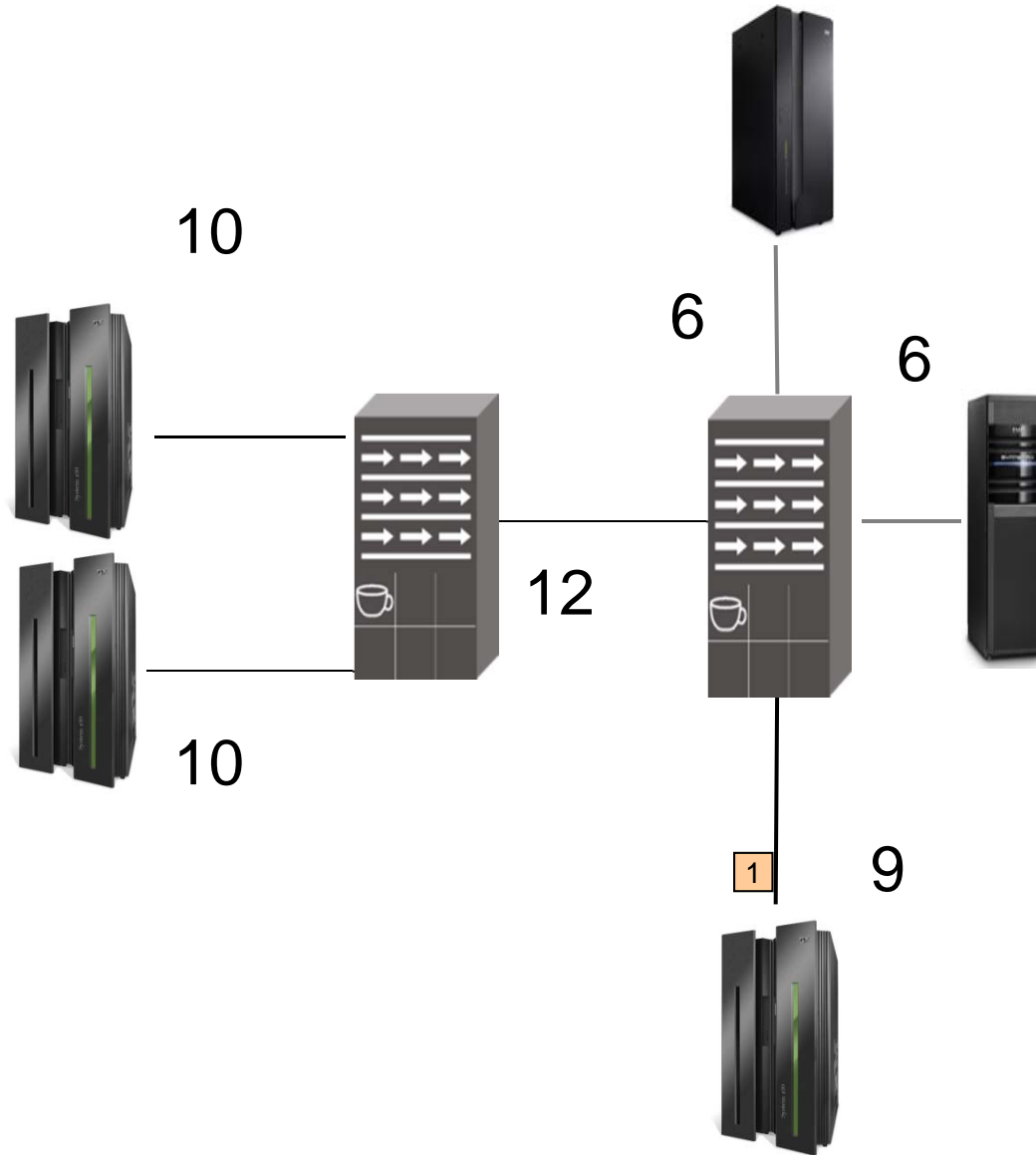
THIS PAGE INTENTIONALLY
LEFT BLANK

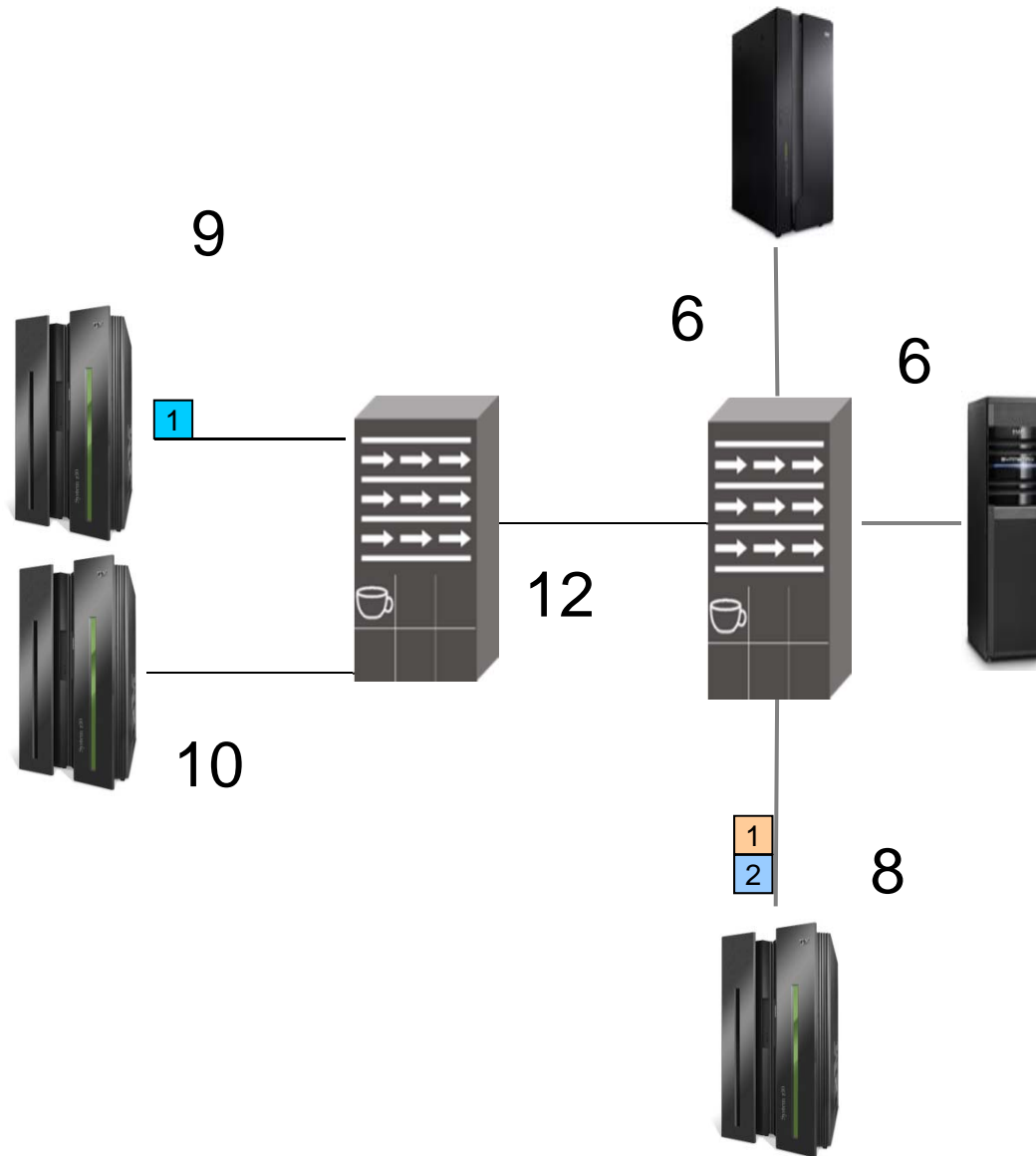
Slow Drainers

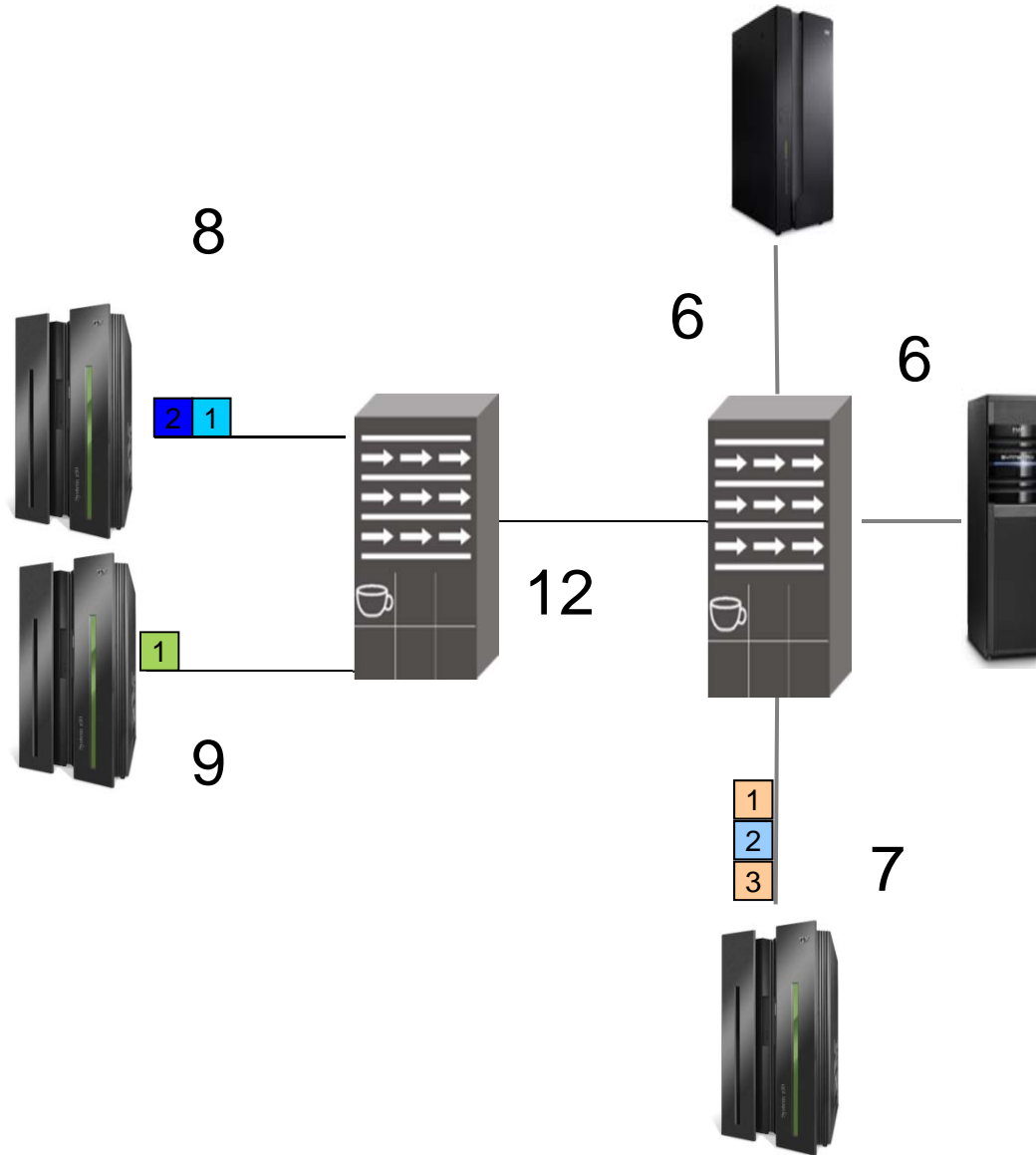


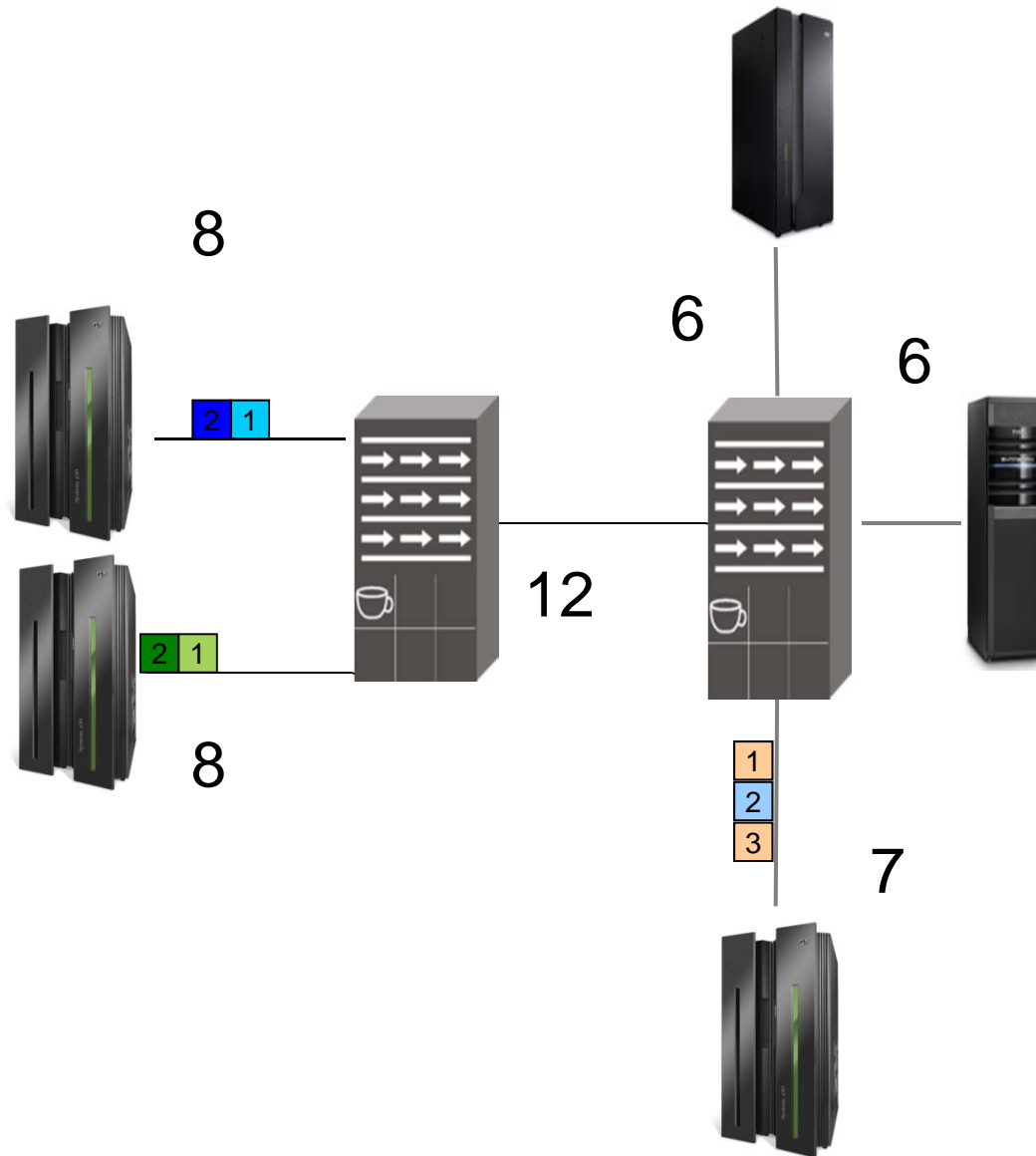
- A device that delays the return of BB credit significantly slower than the desired load
- Moderate delays can cause performance impacts even to other uninvolved devices
- Severe delays can cause errors (IFCCs), even to other uninvolved devices

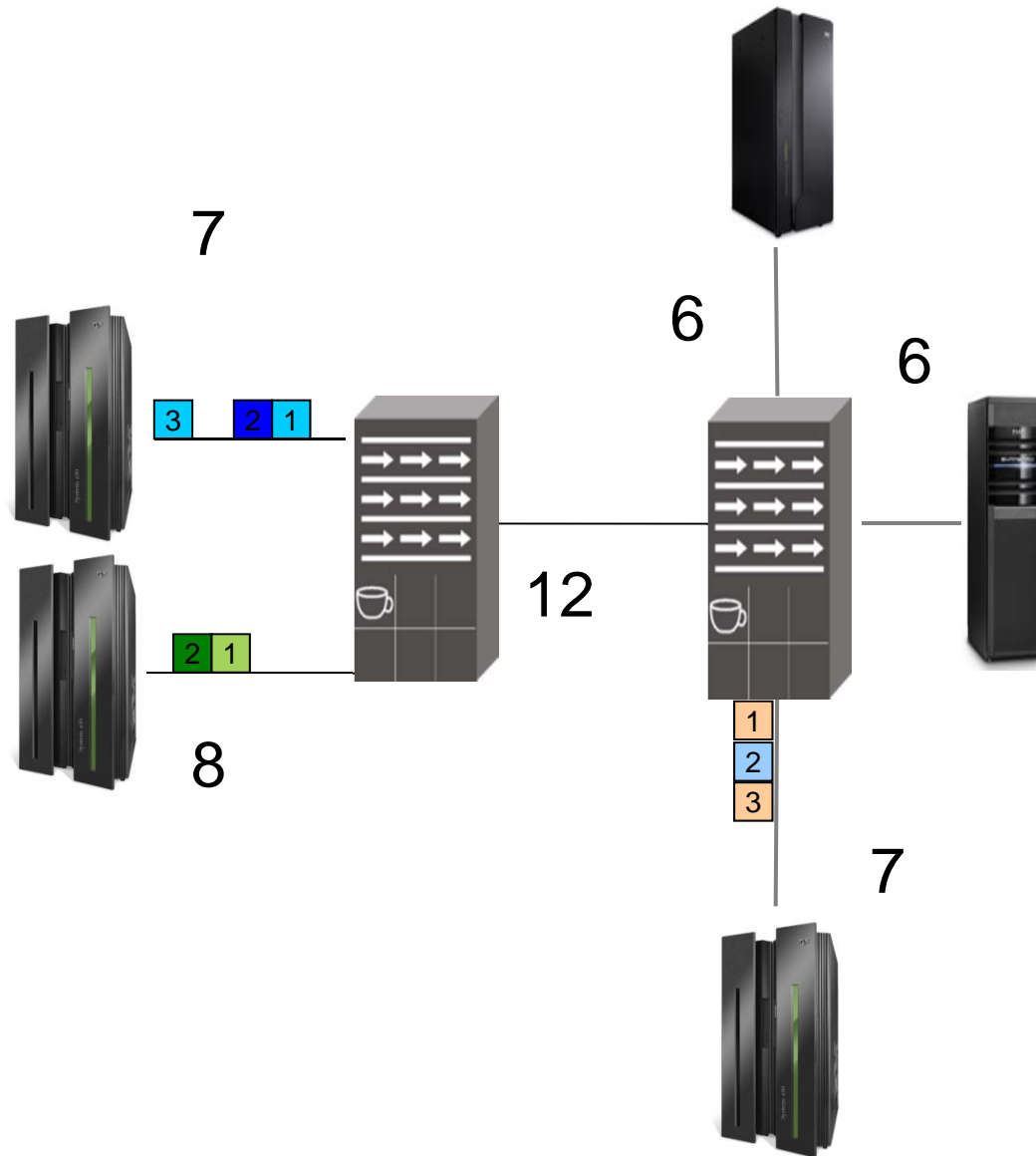


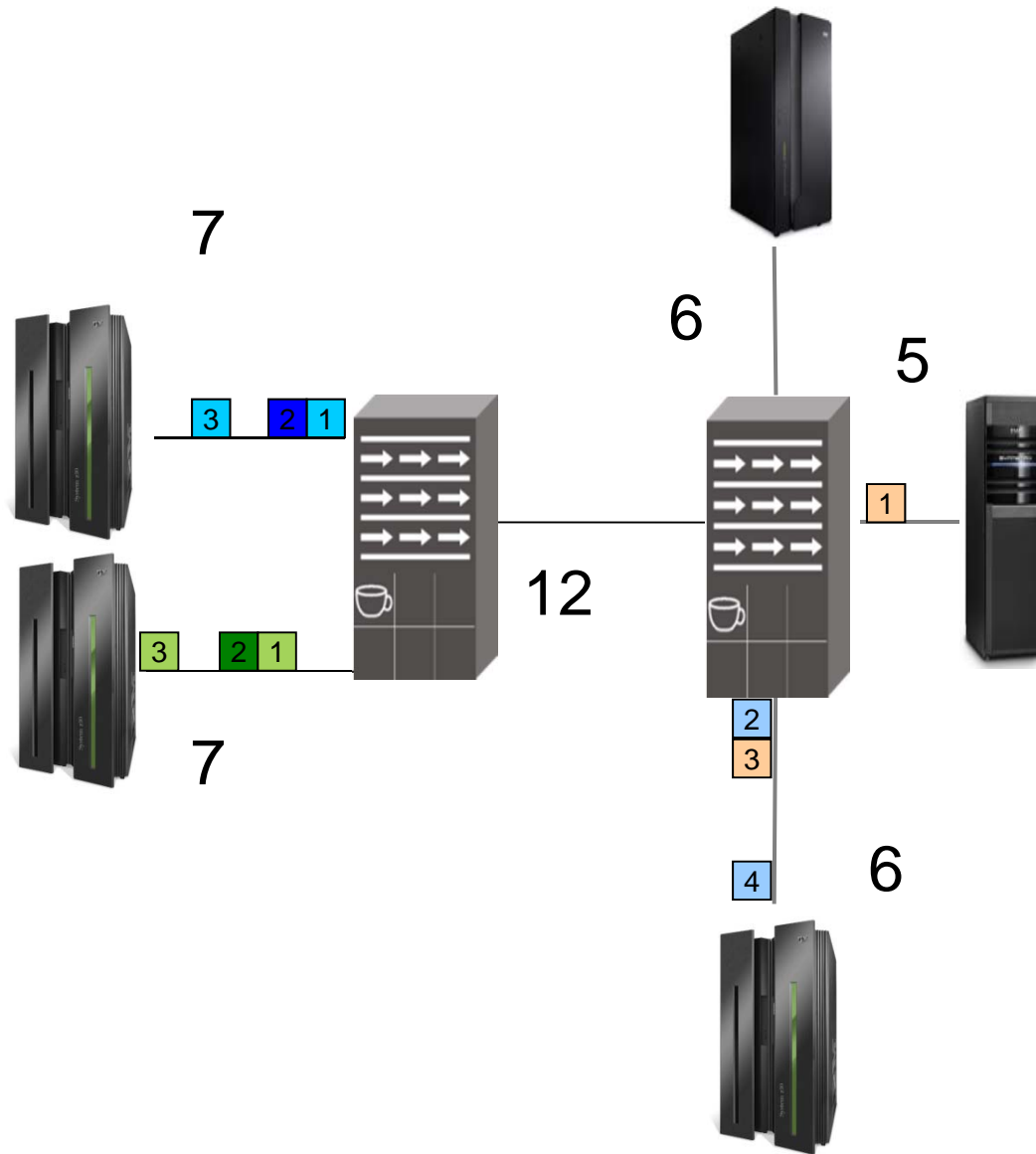


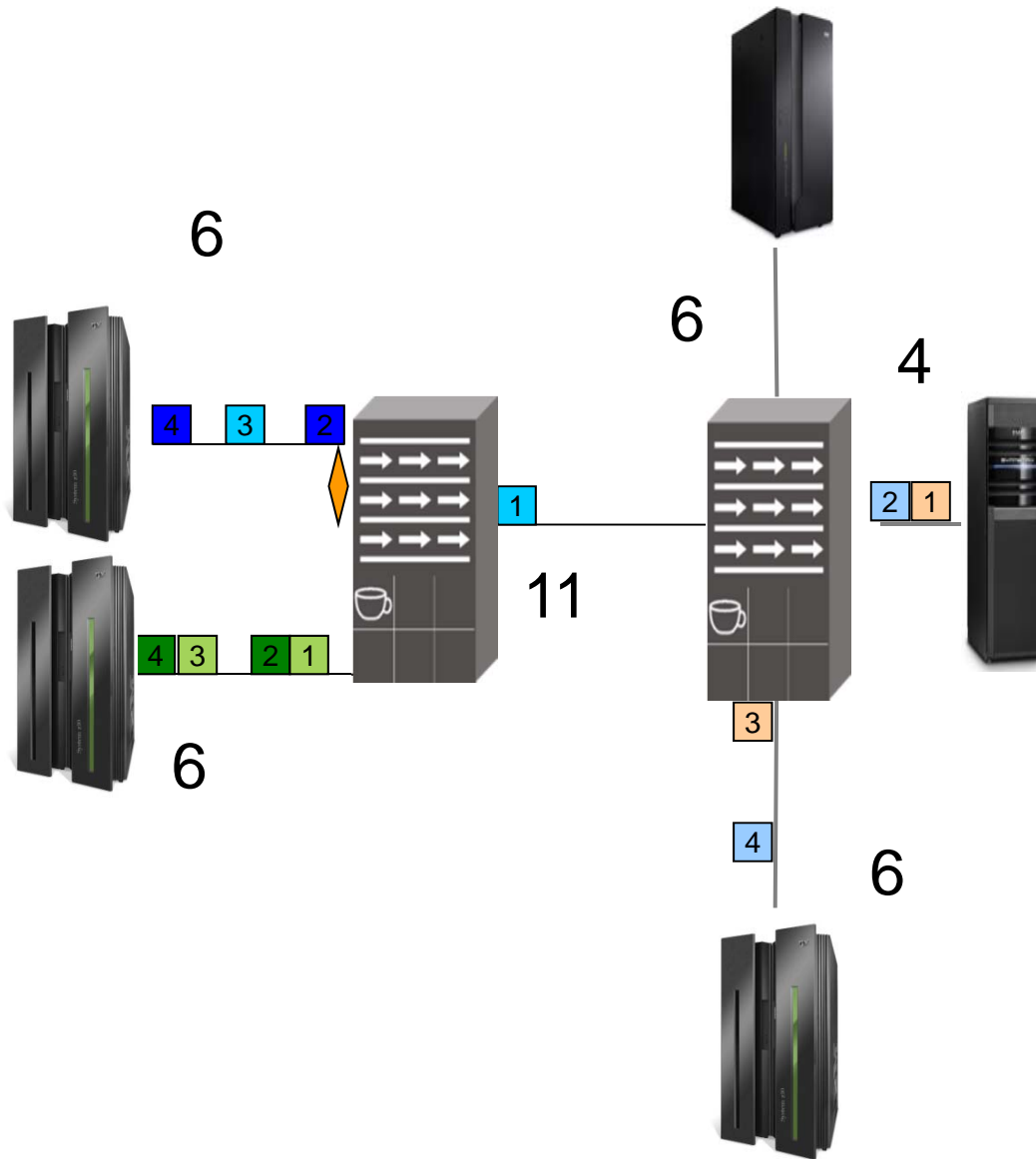


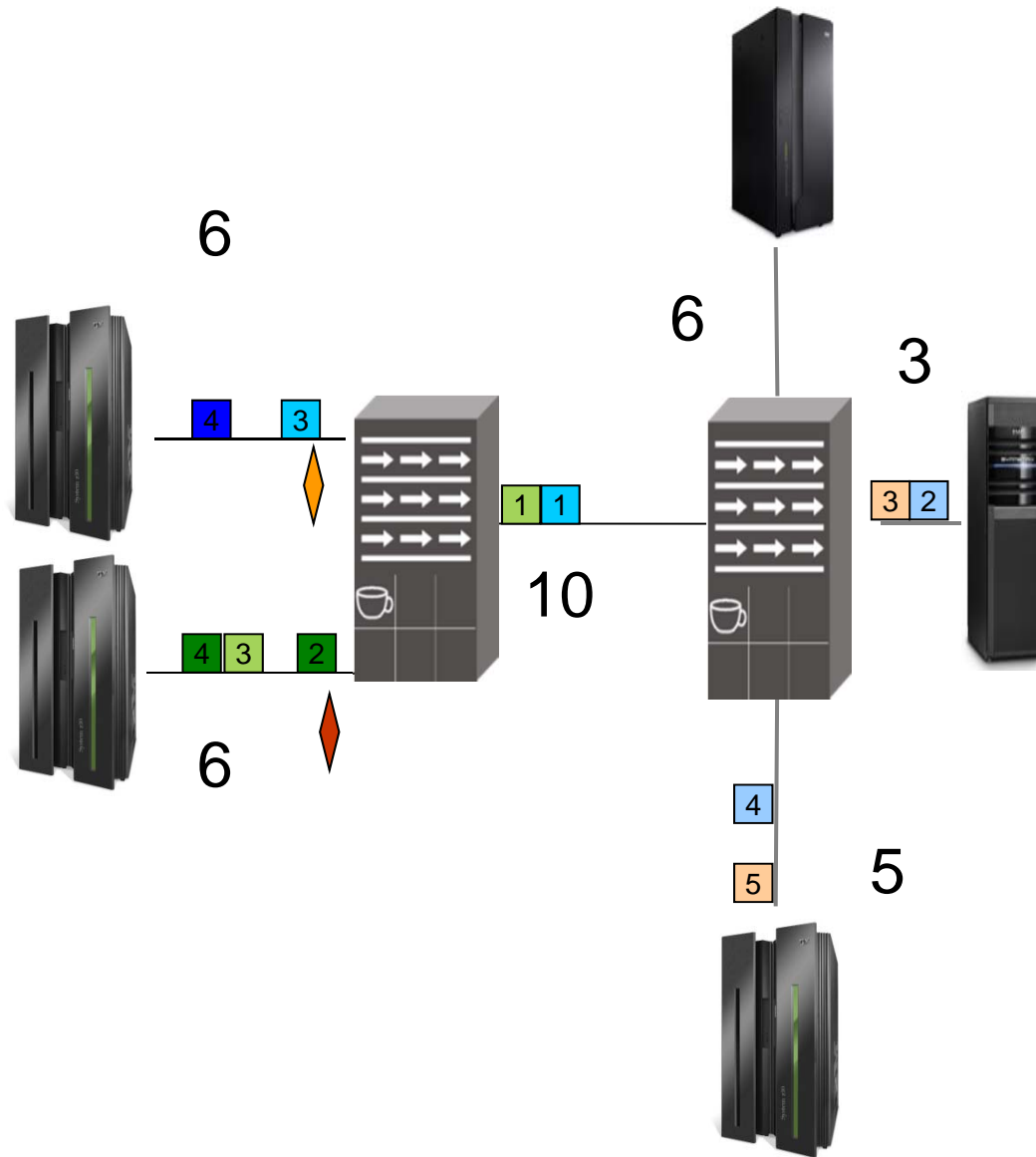


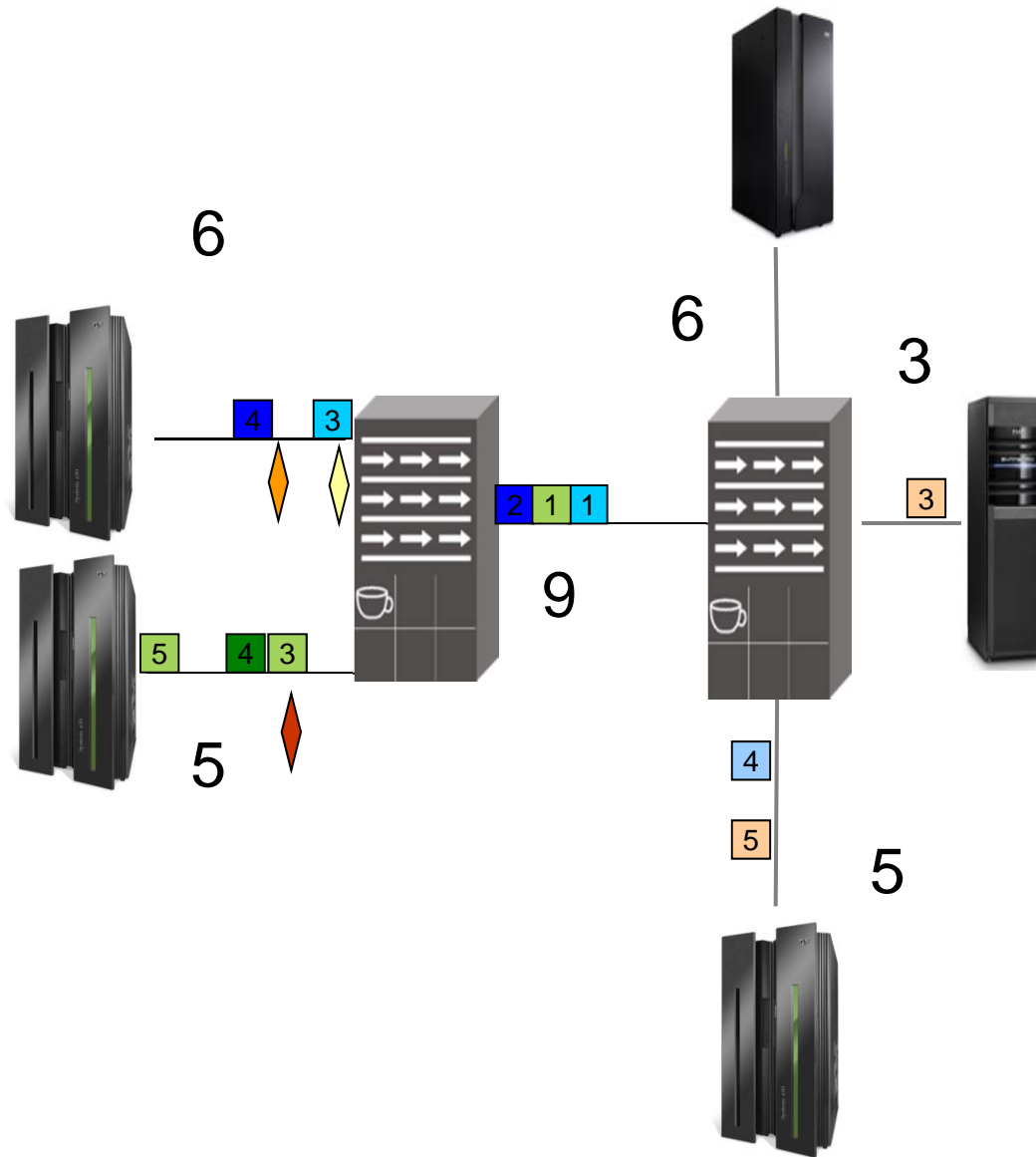


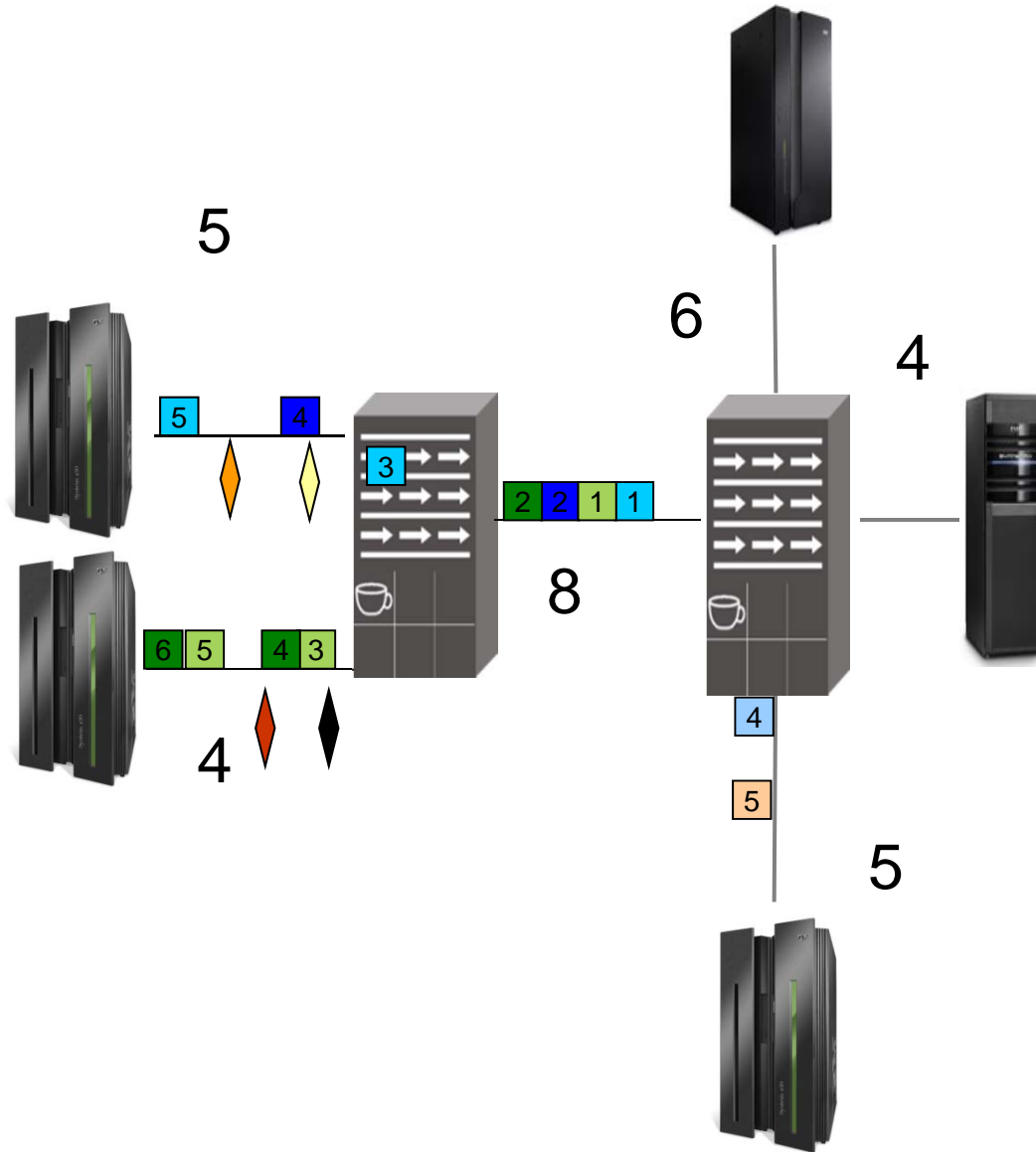


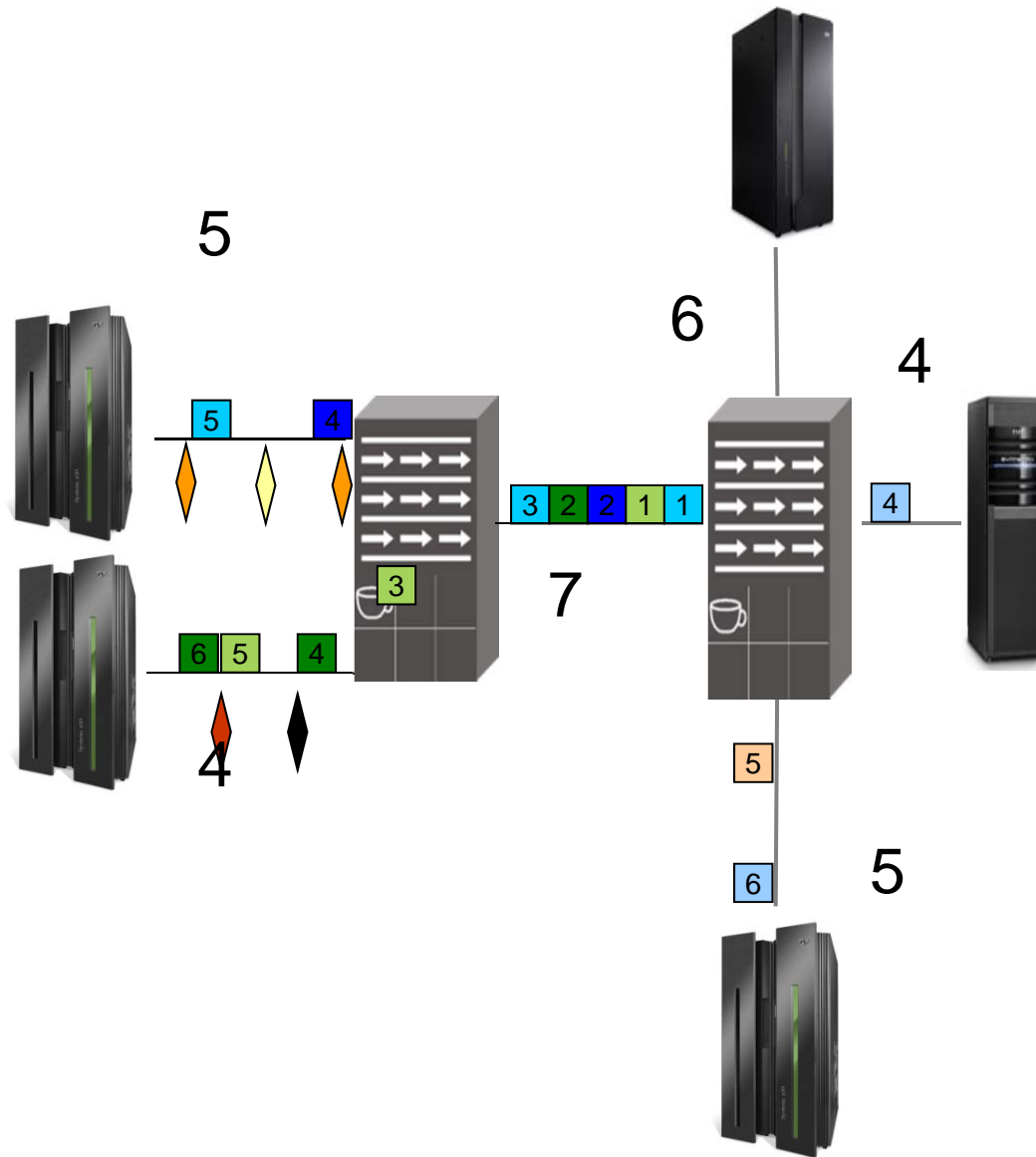


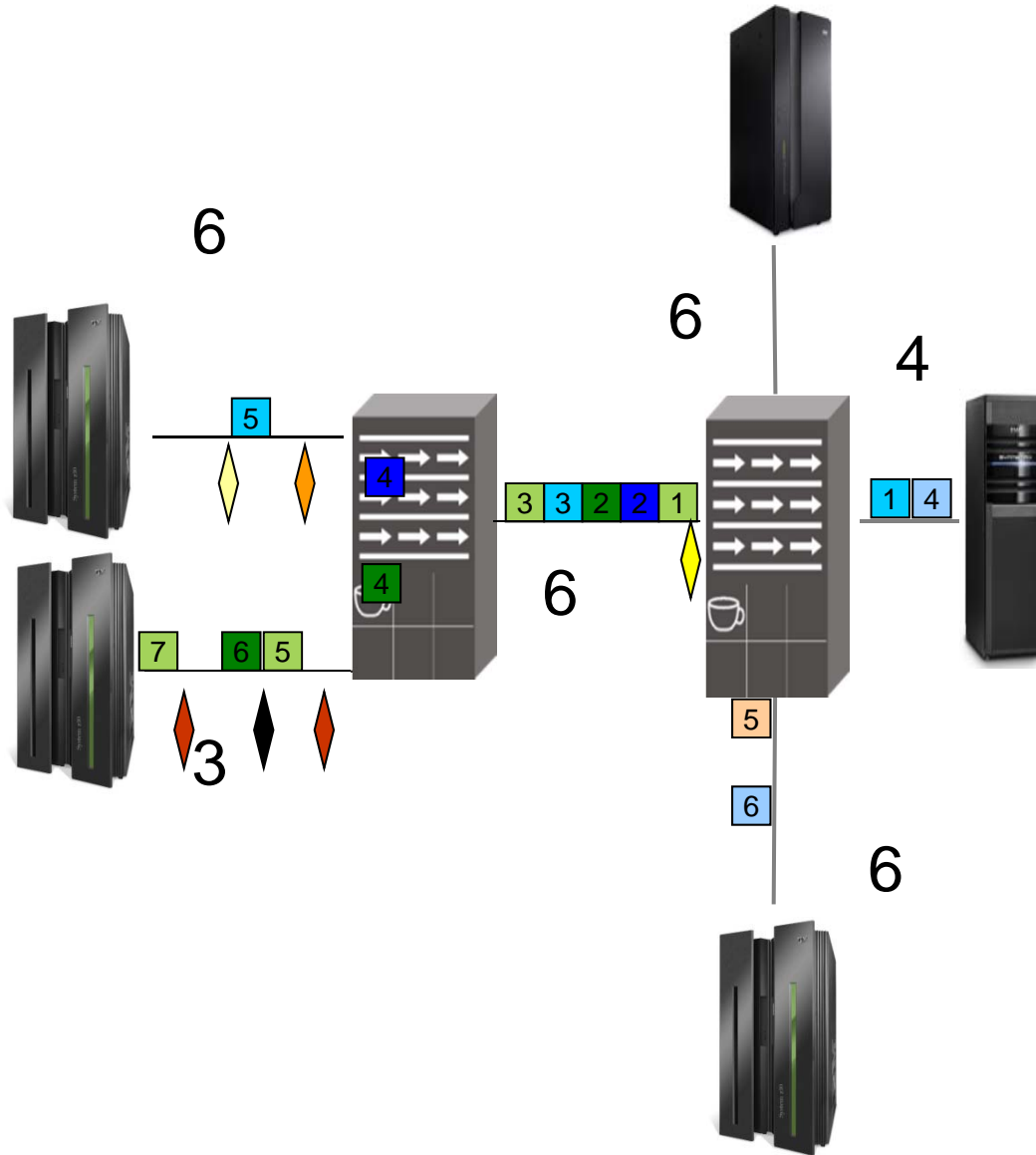


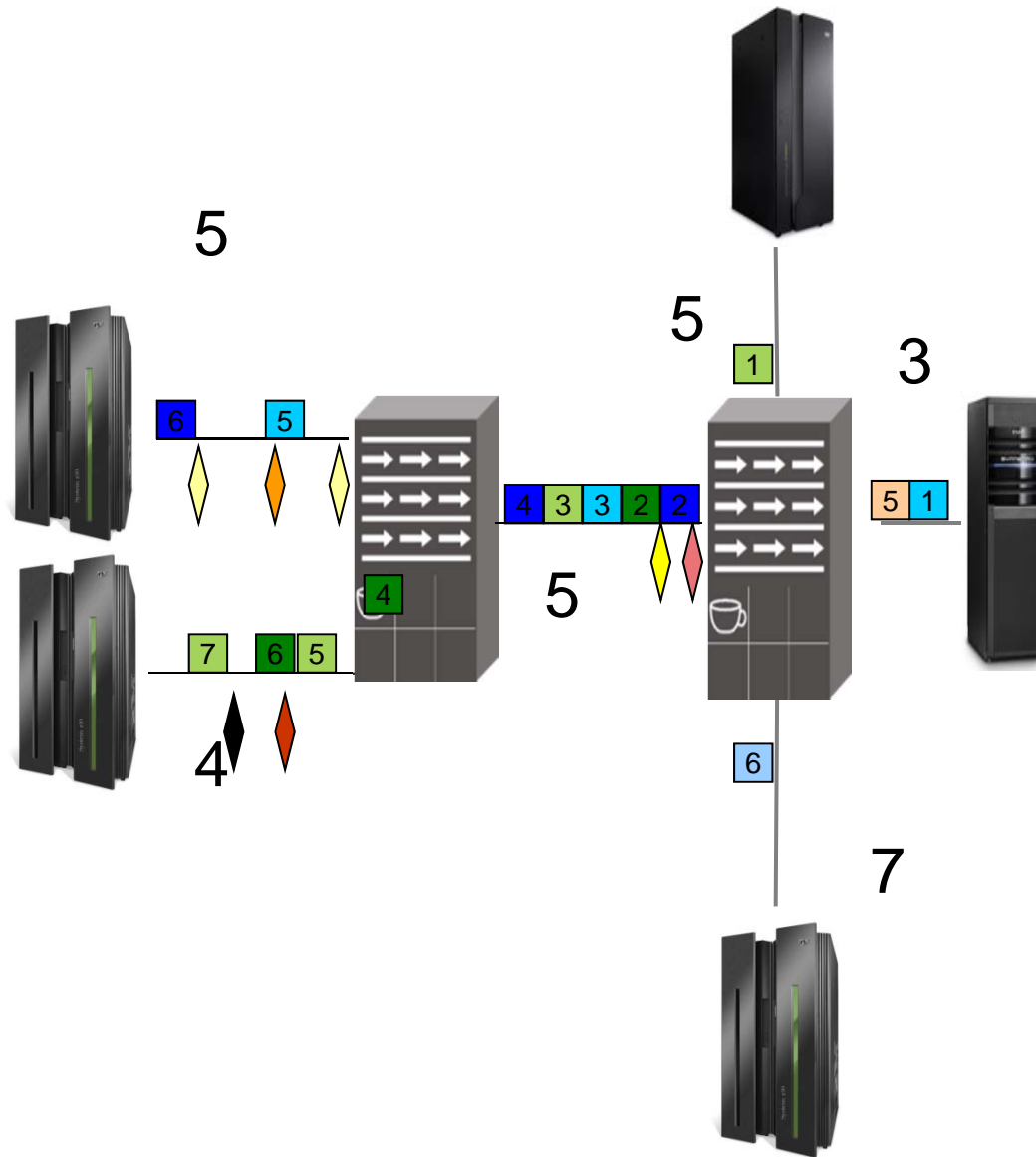


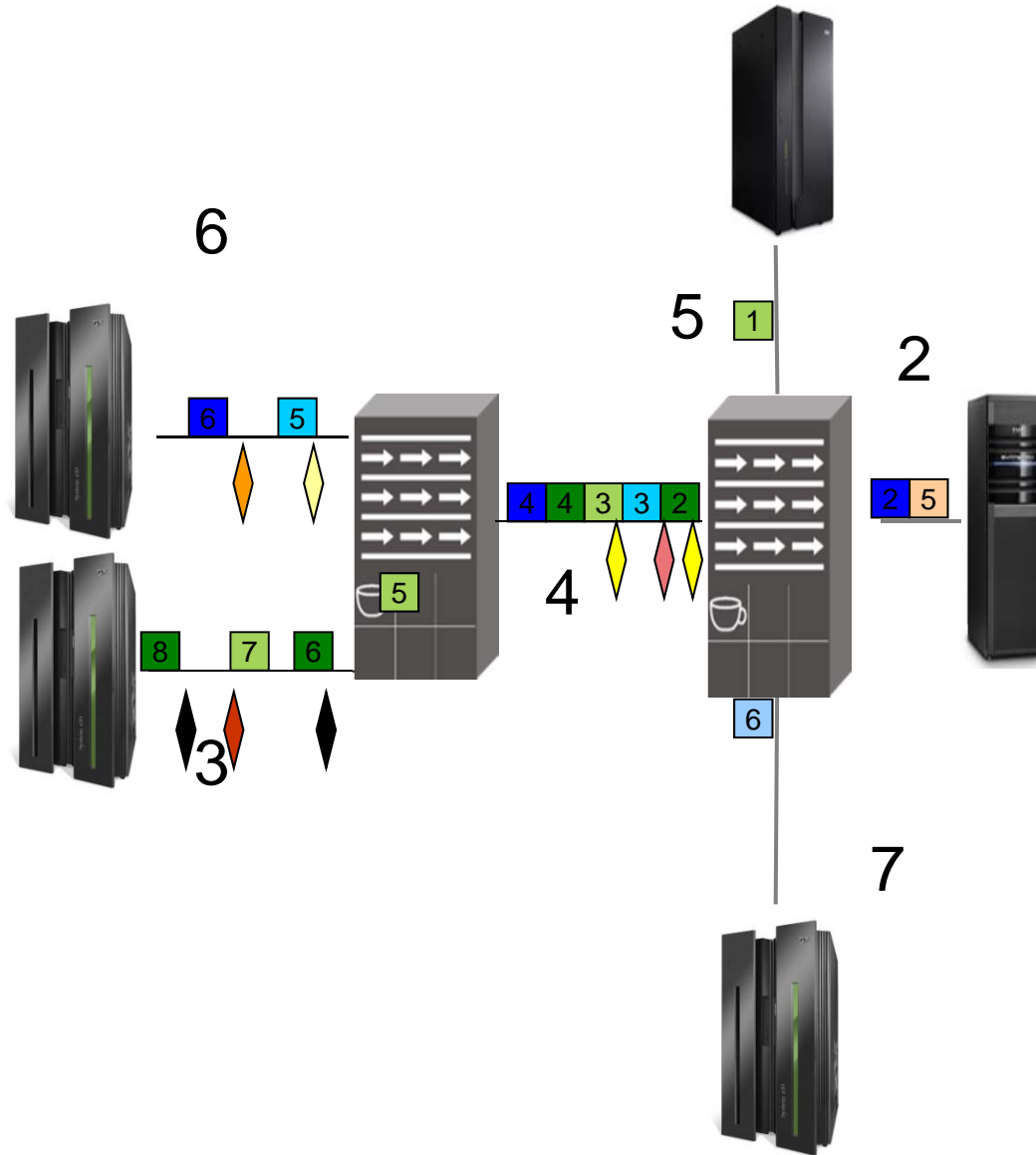


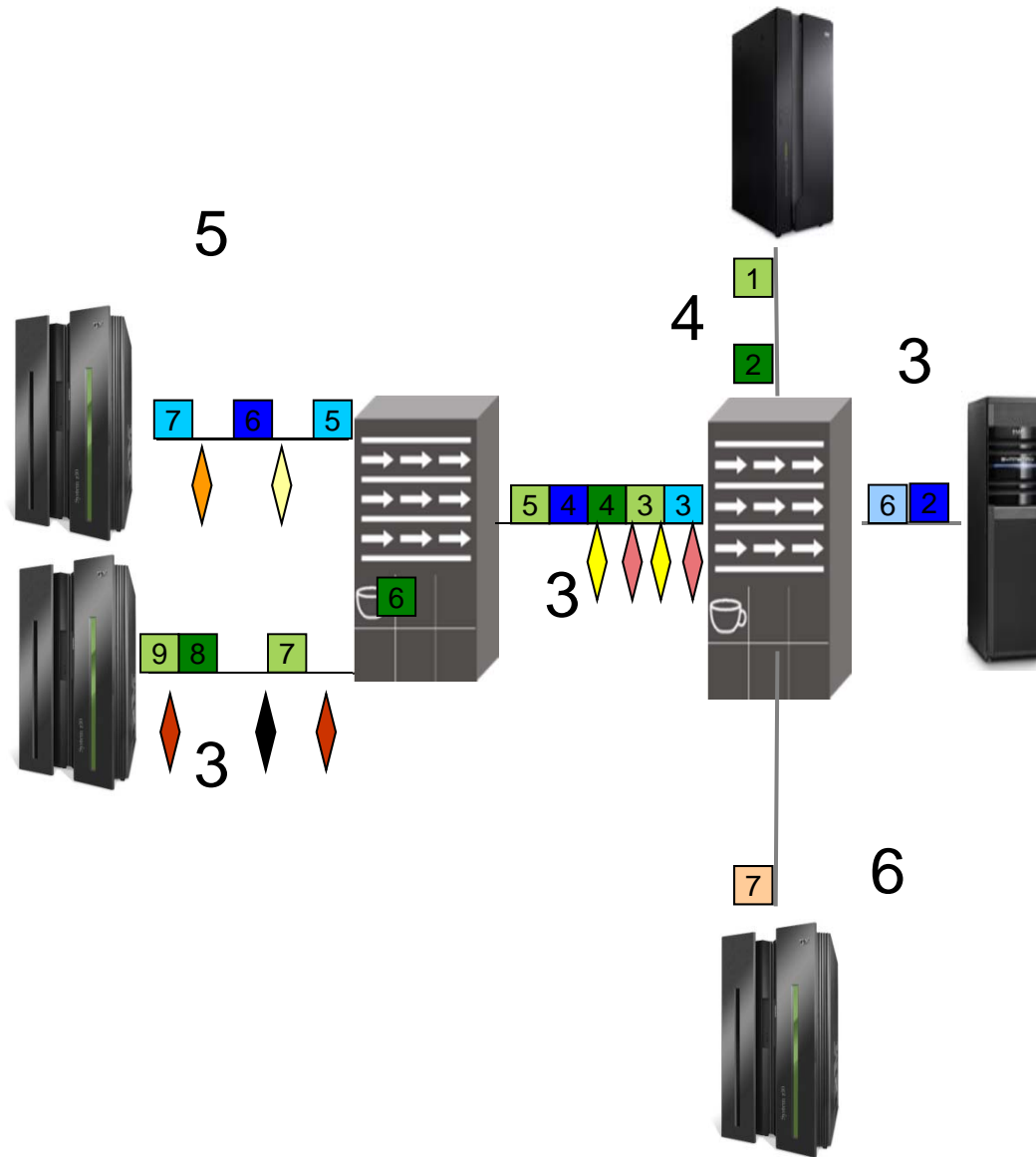


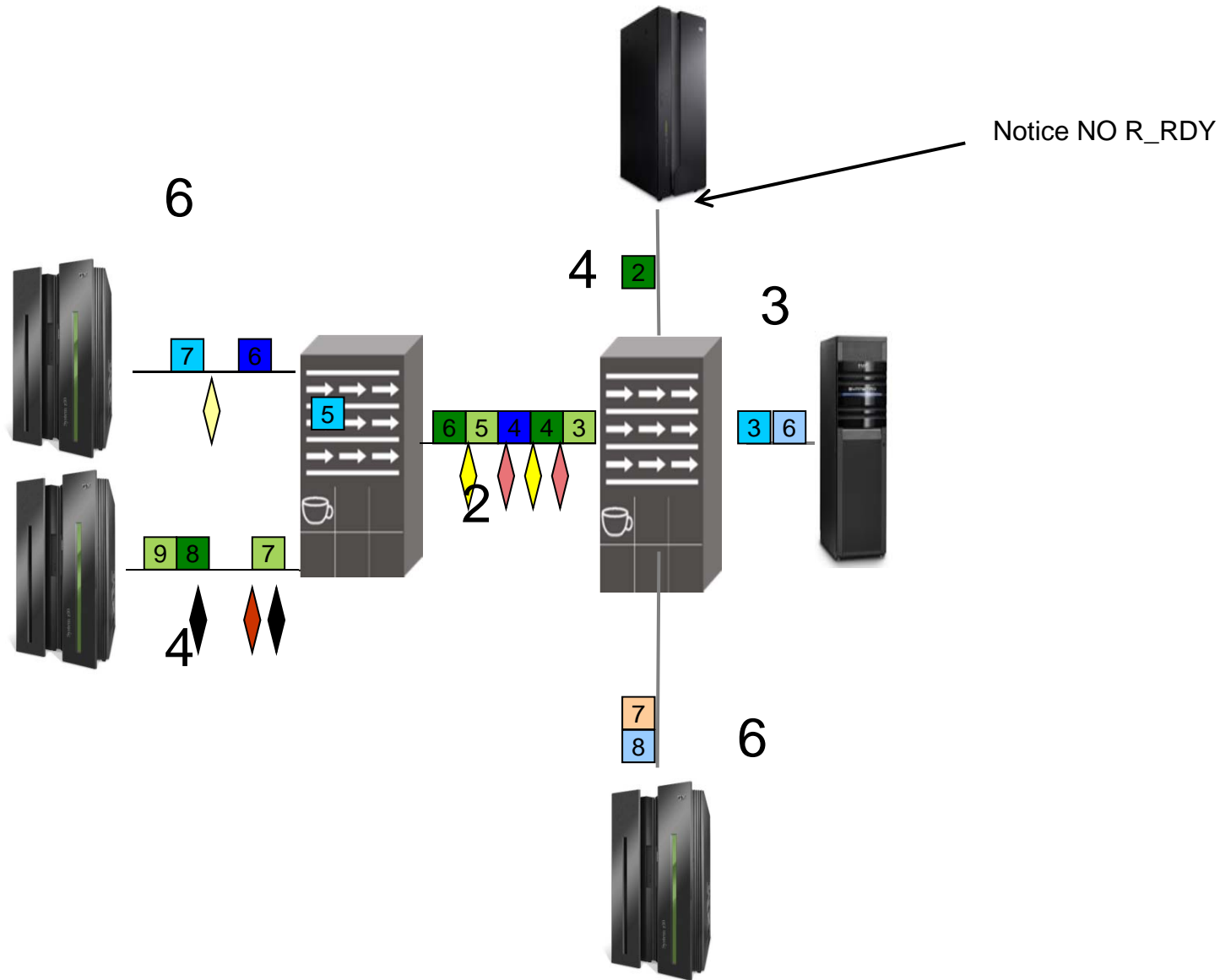


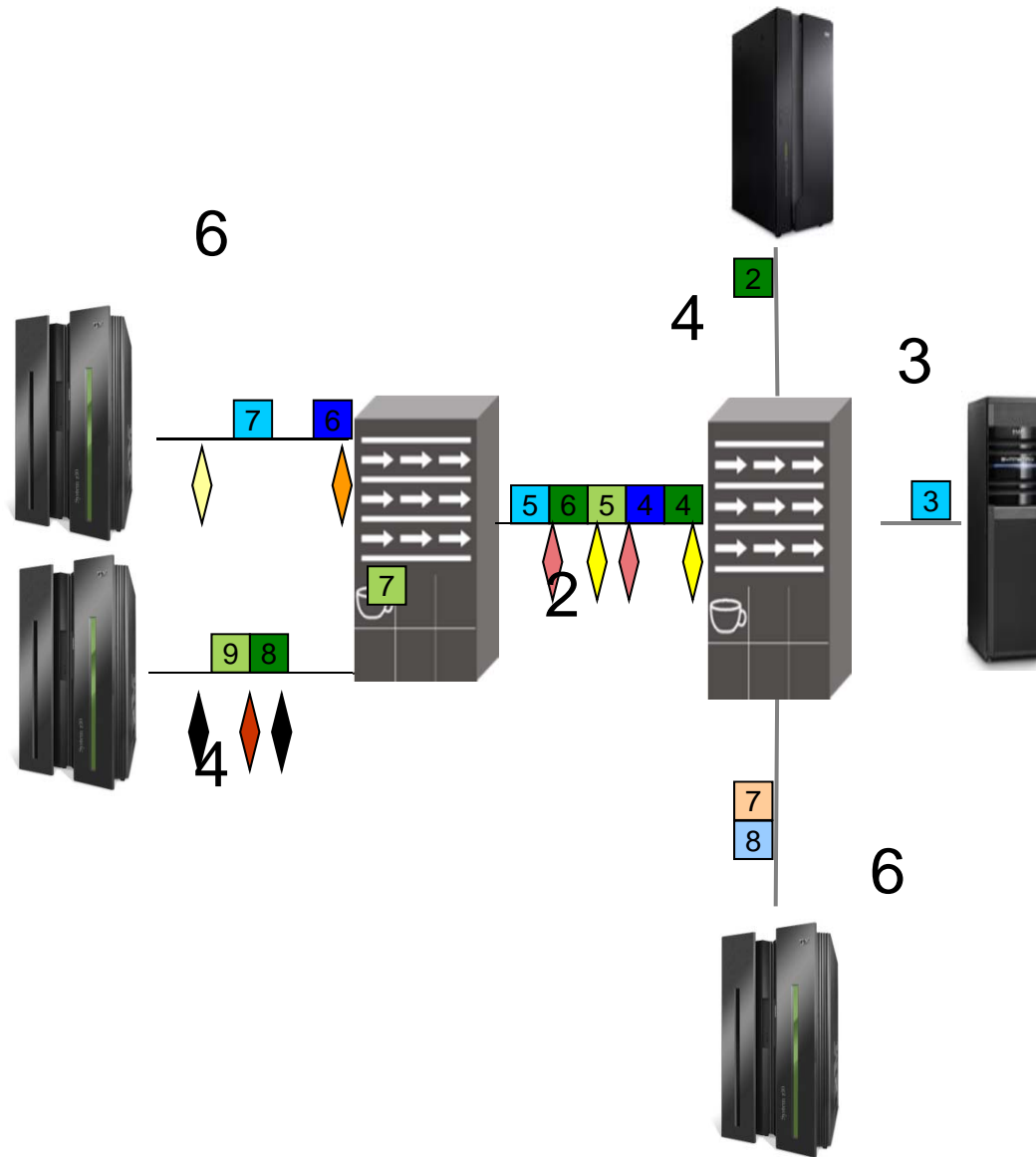


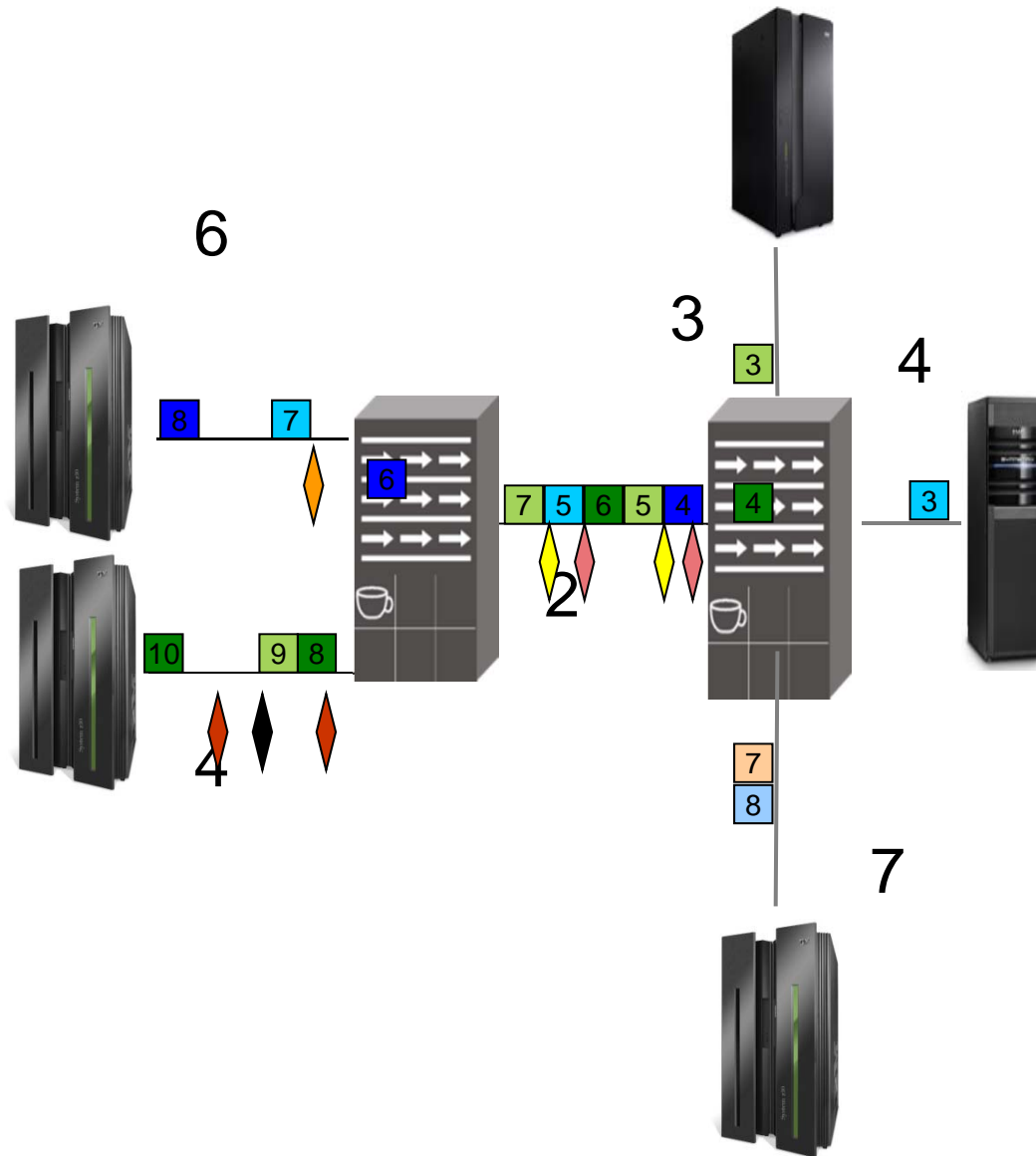


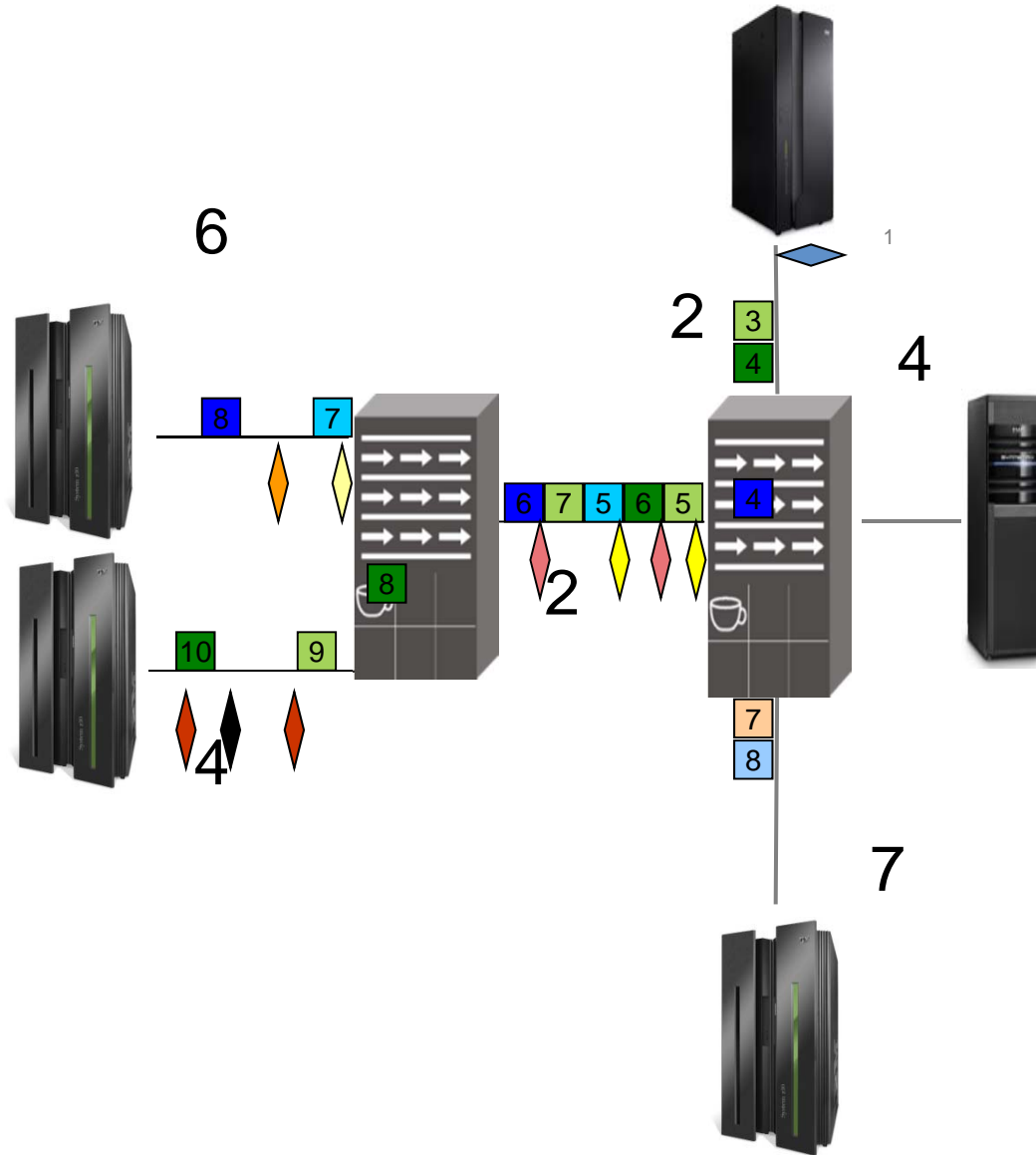


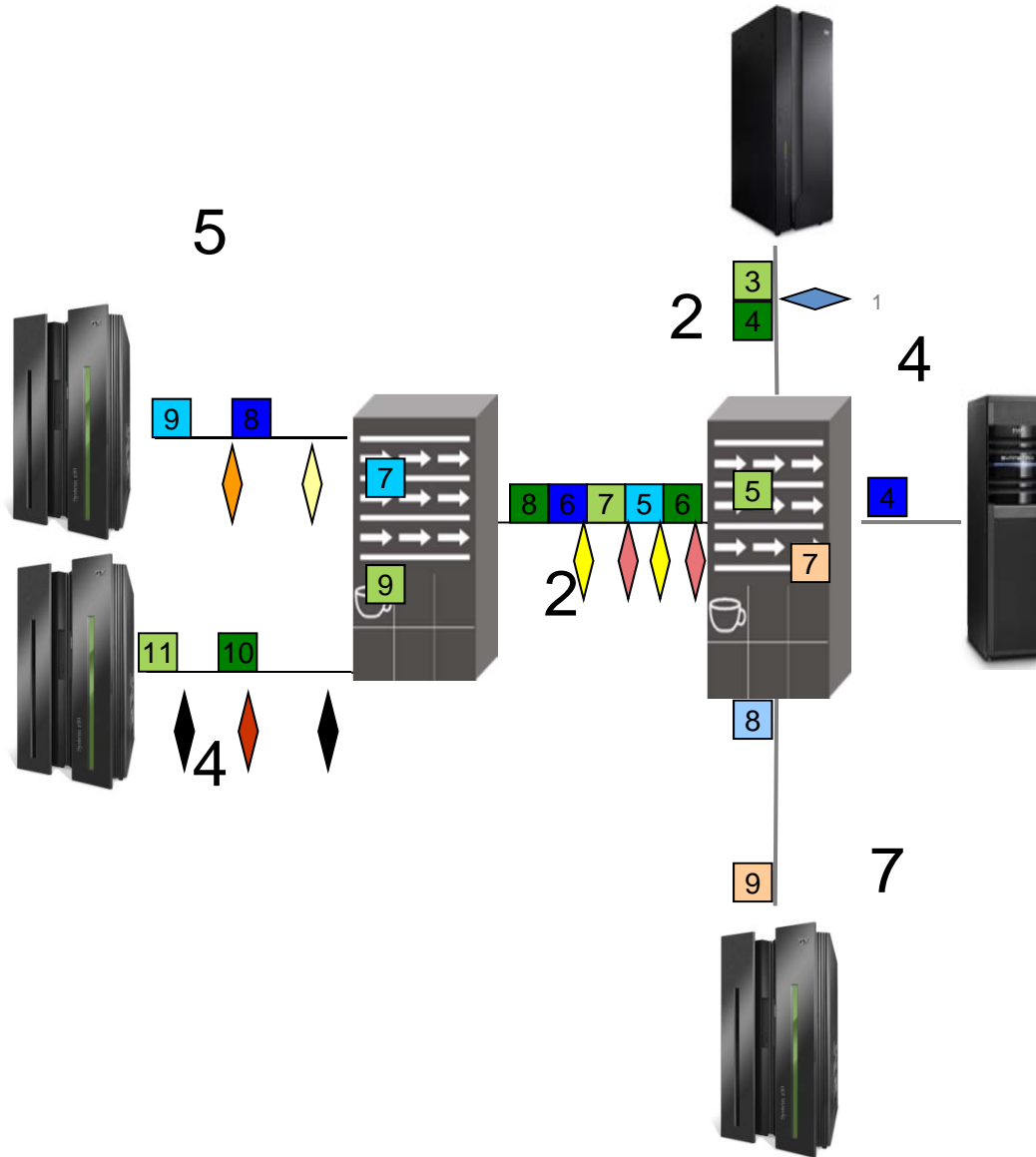


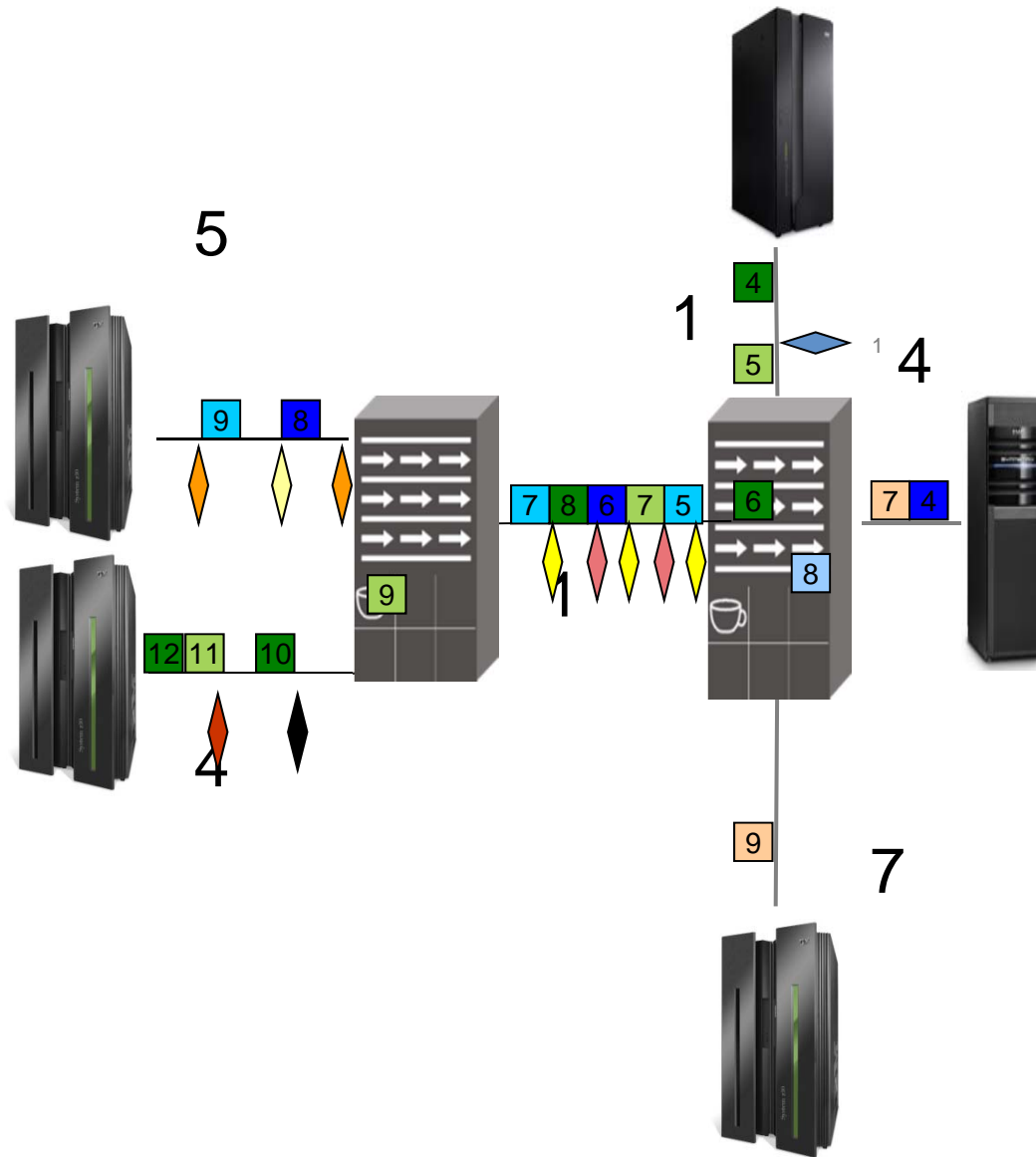


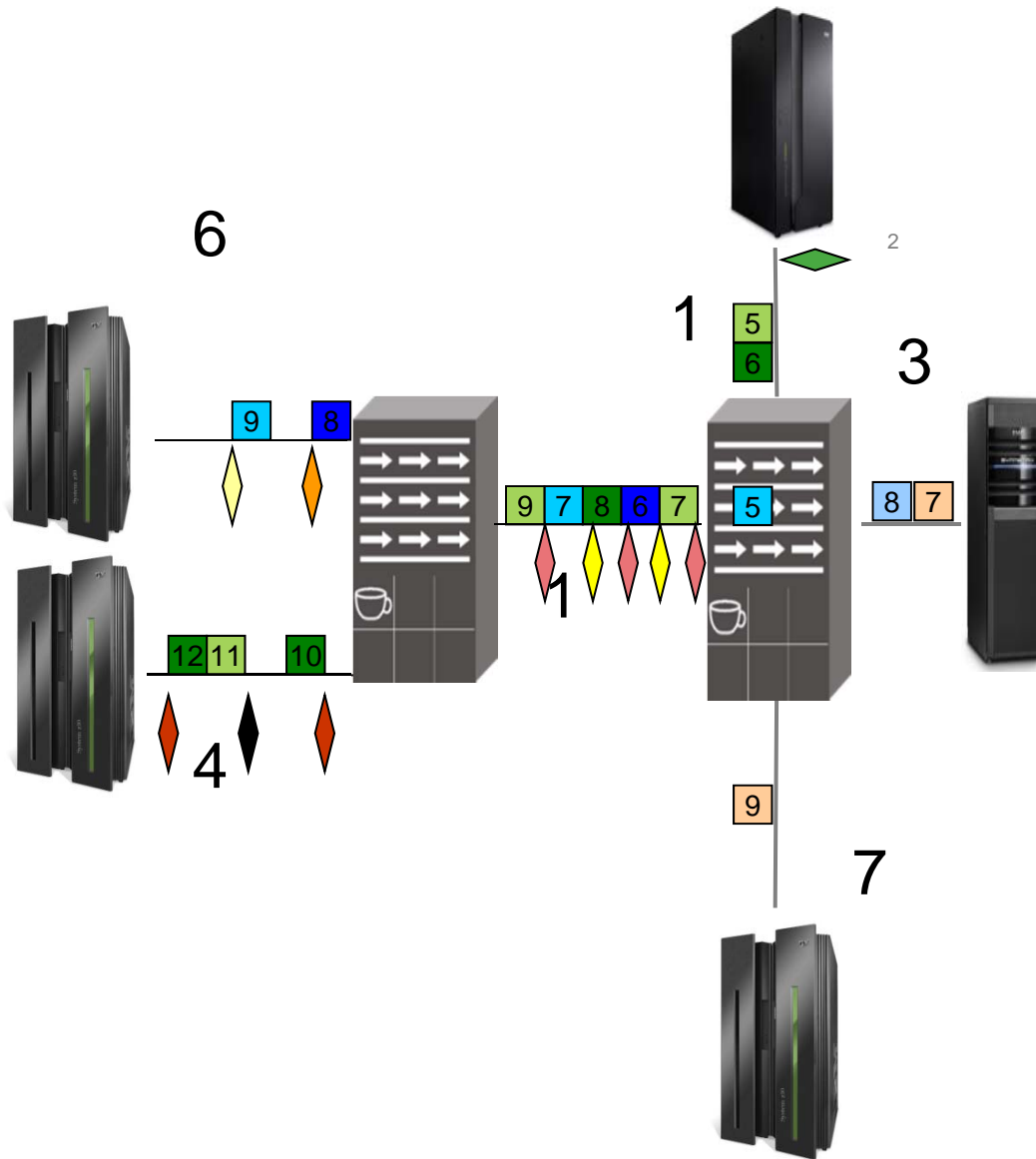


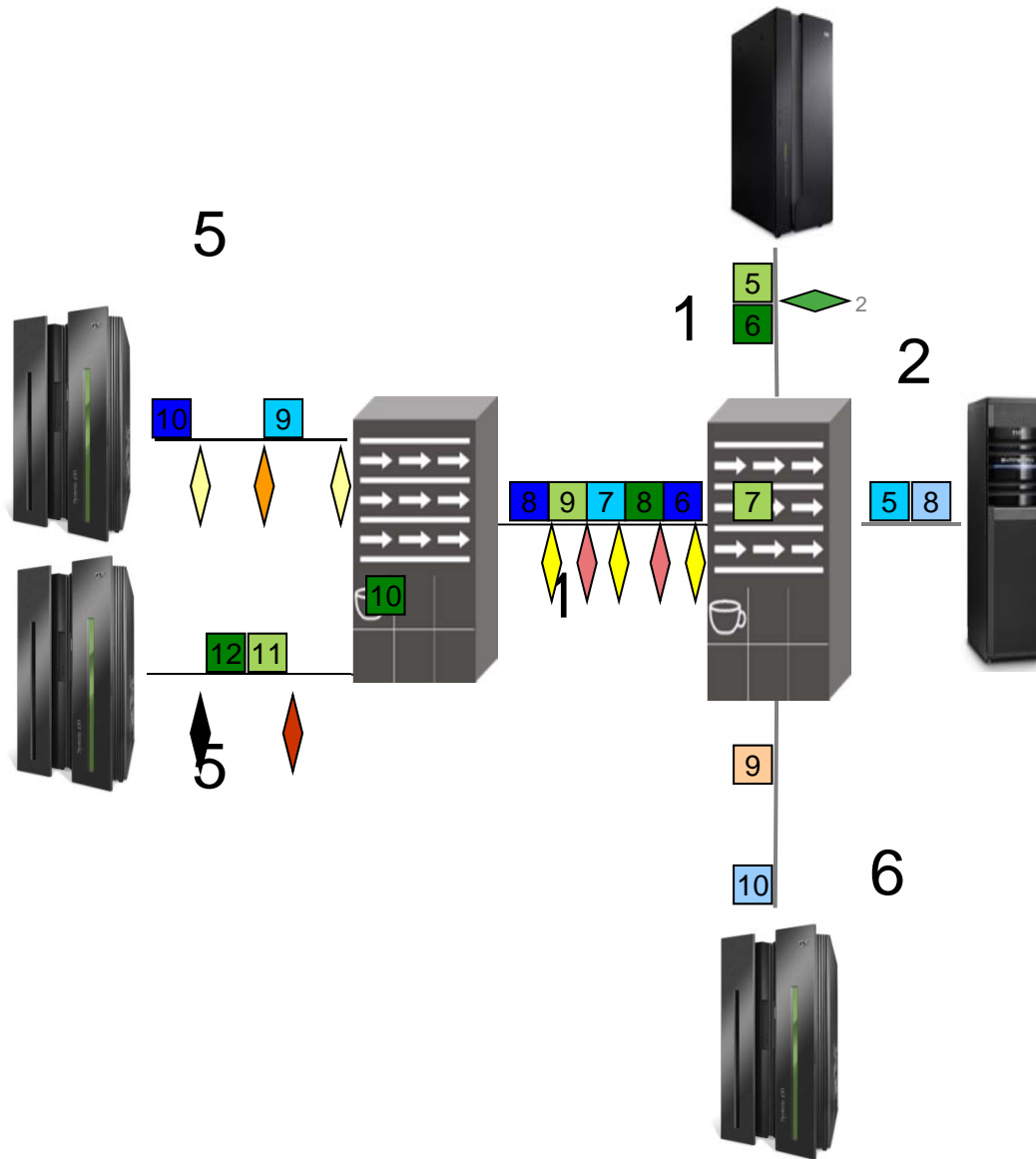


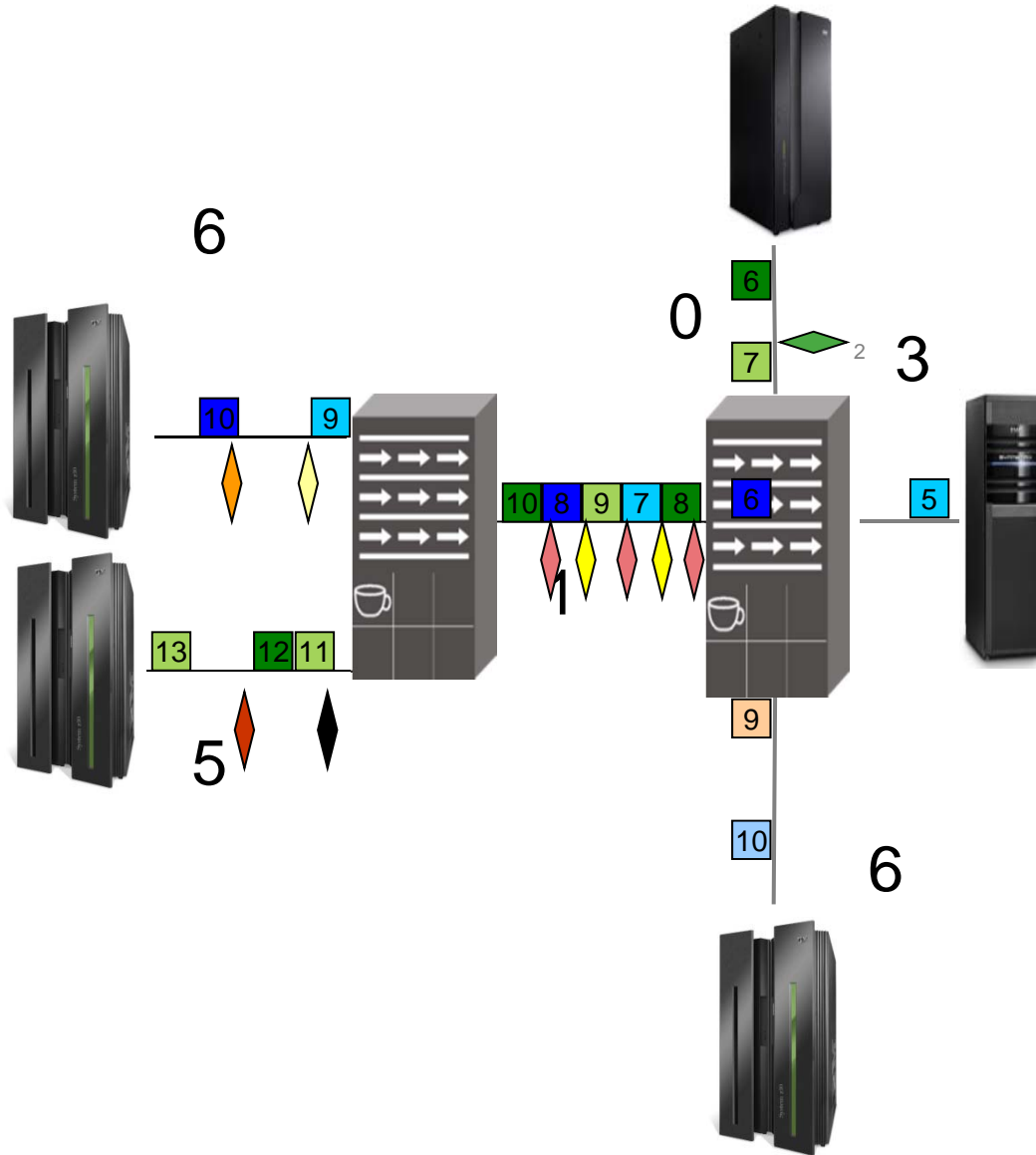


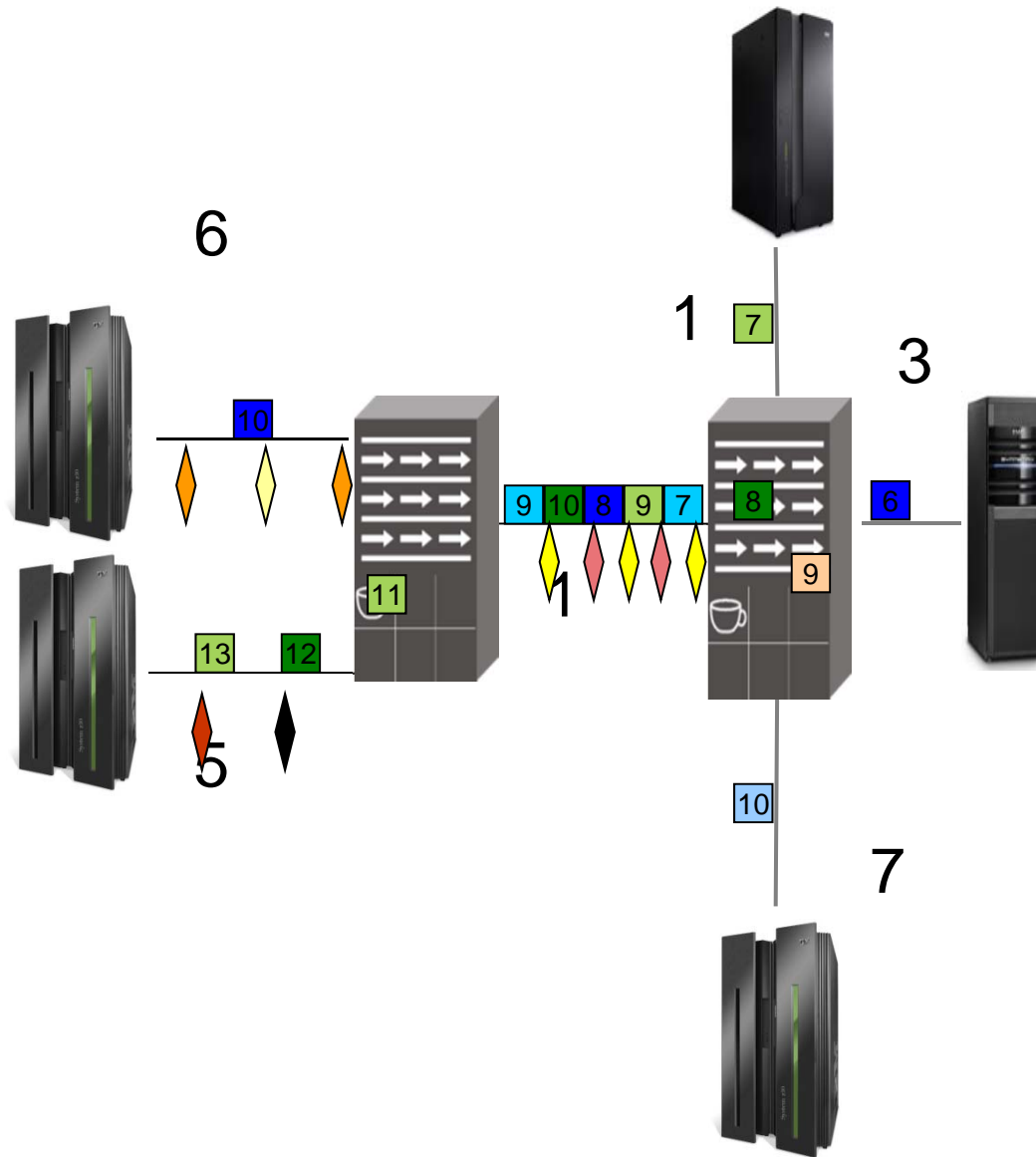


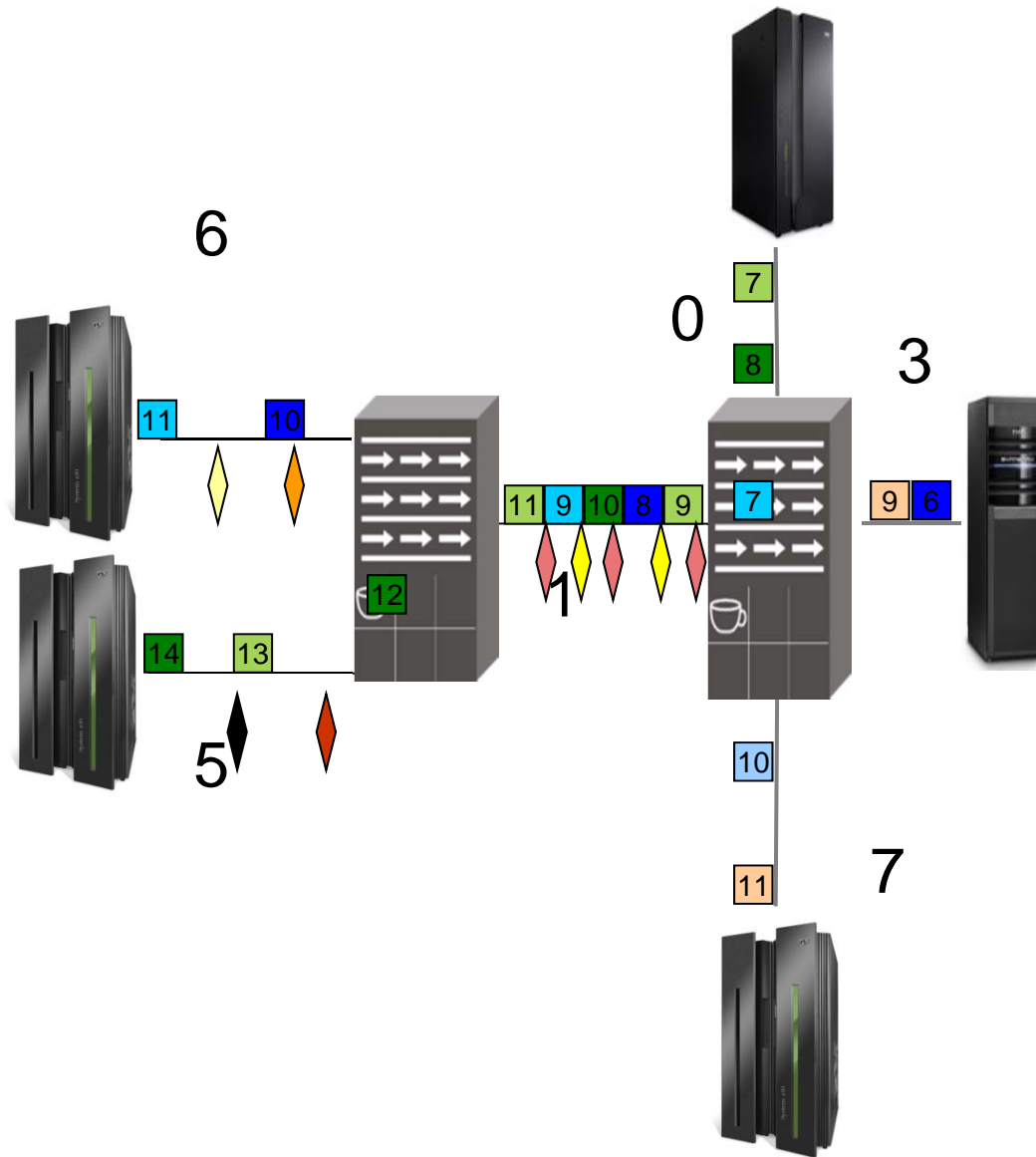


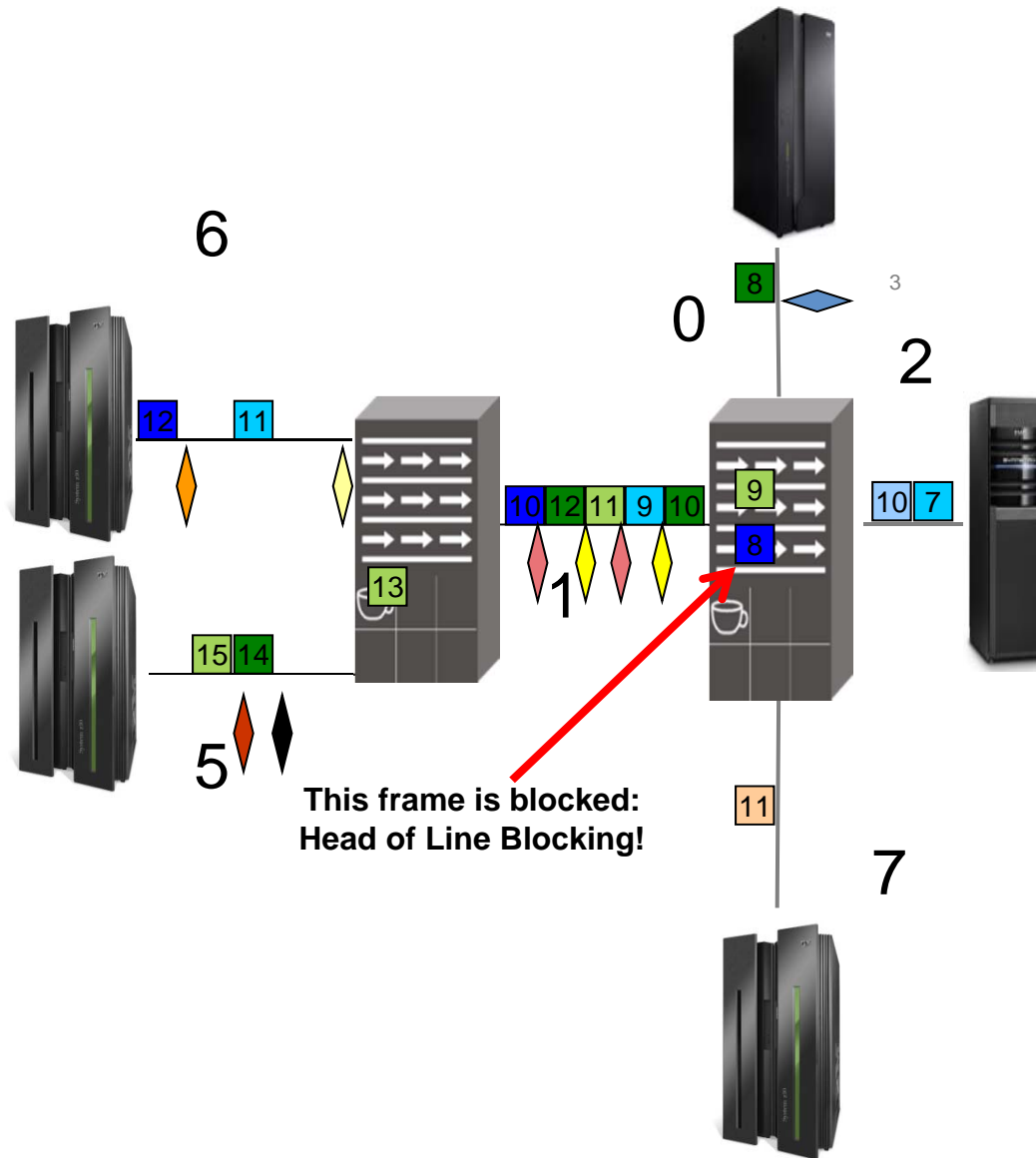


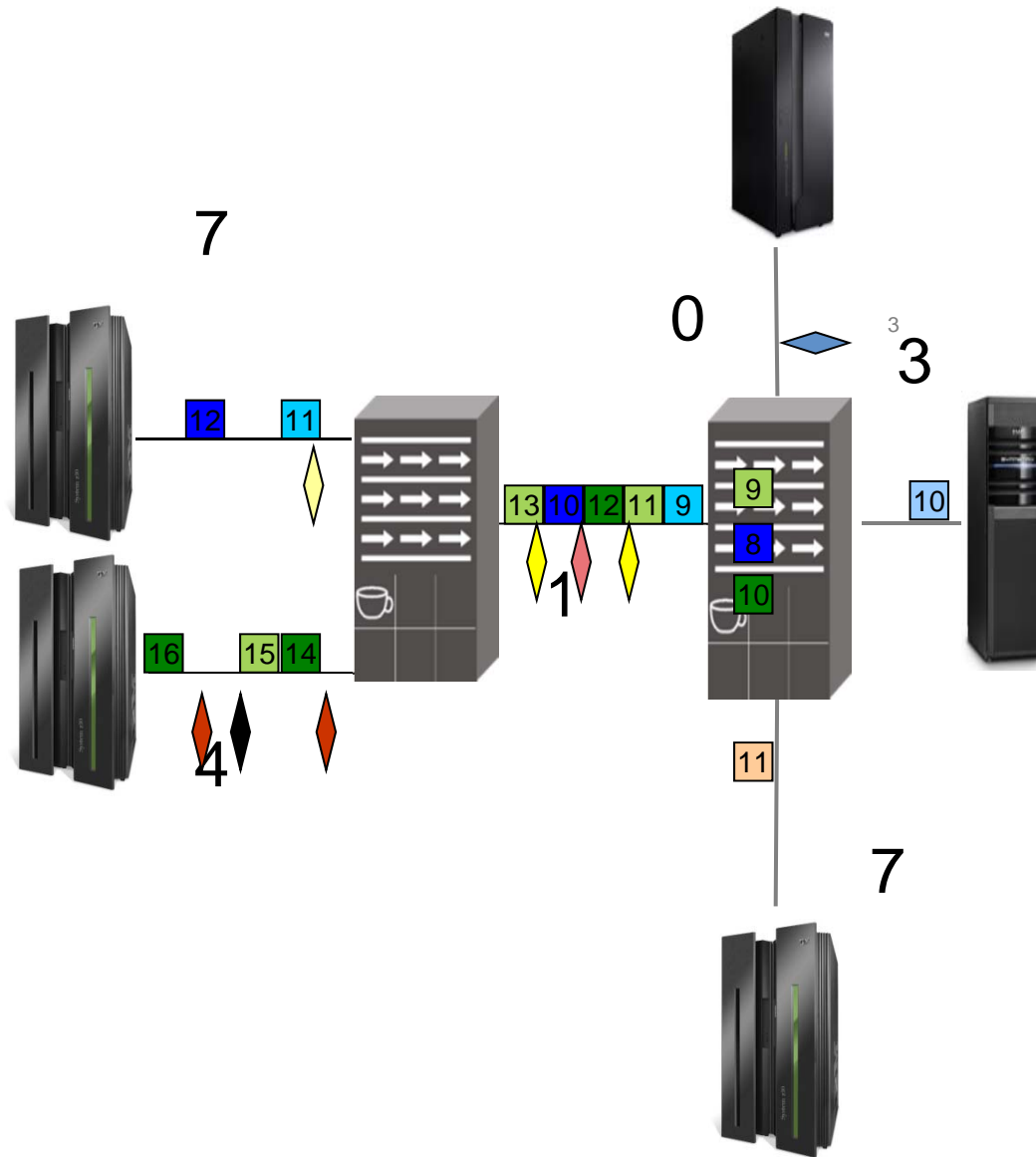


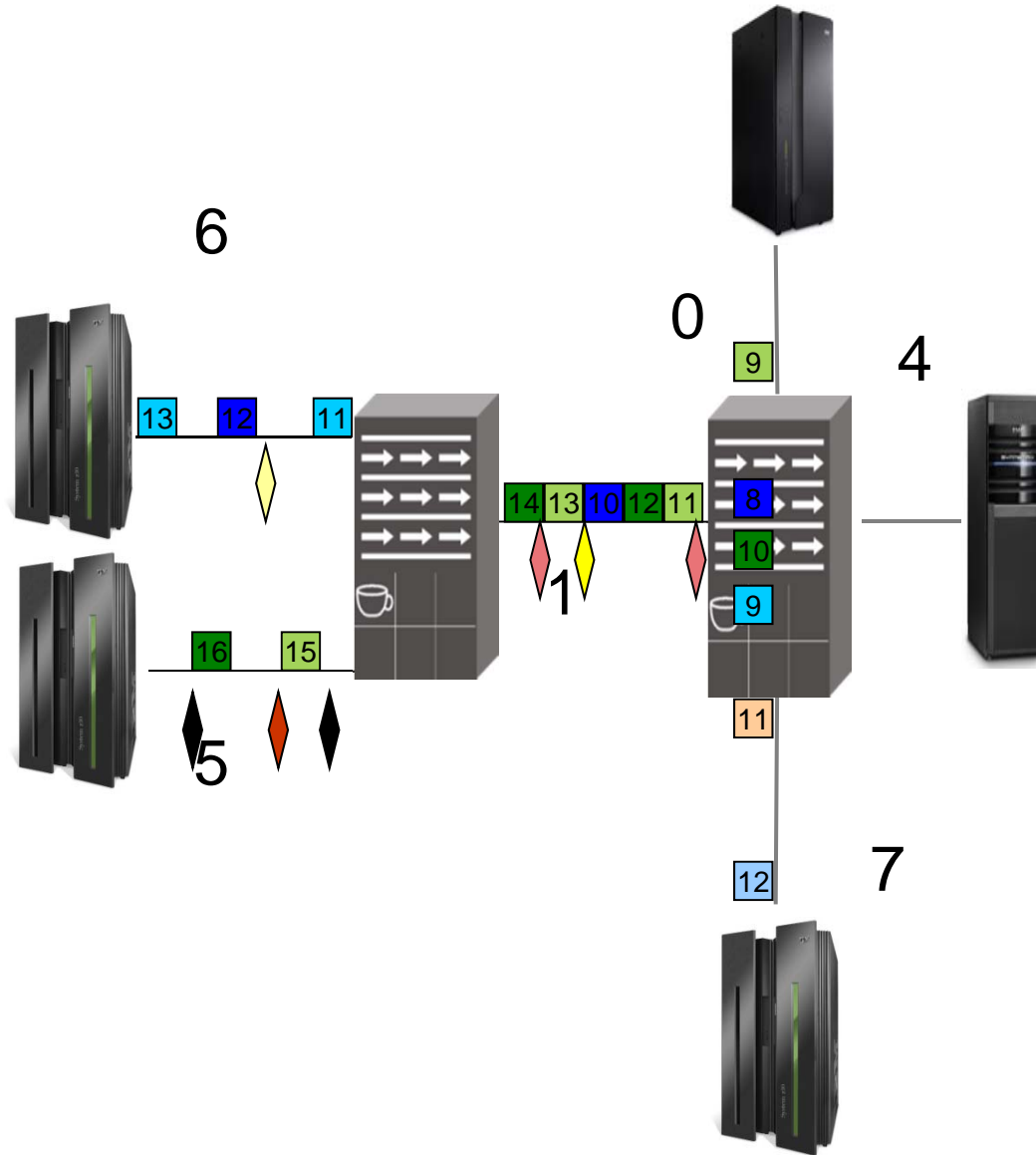


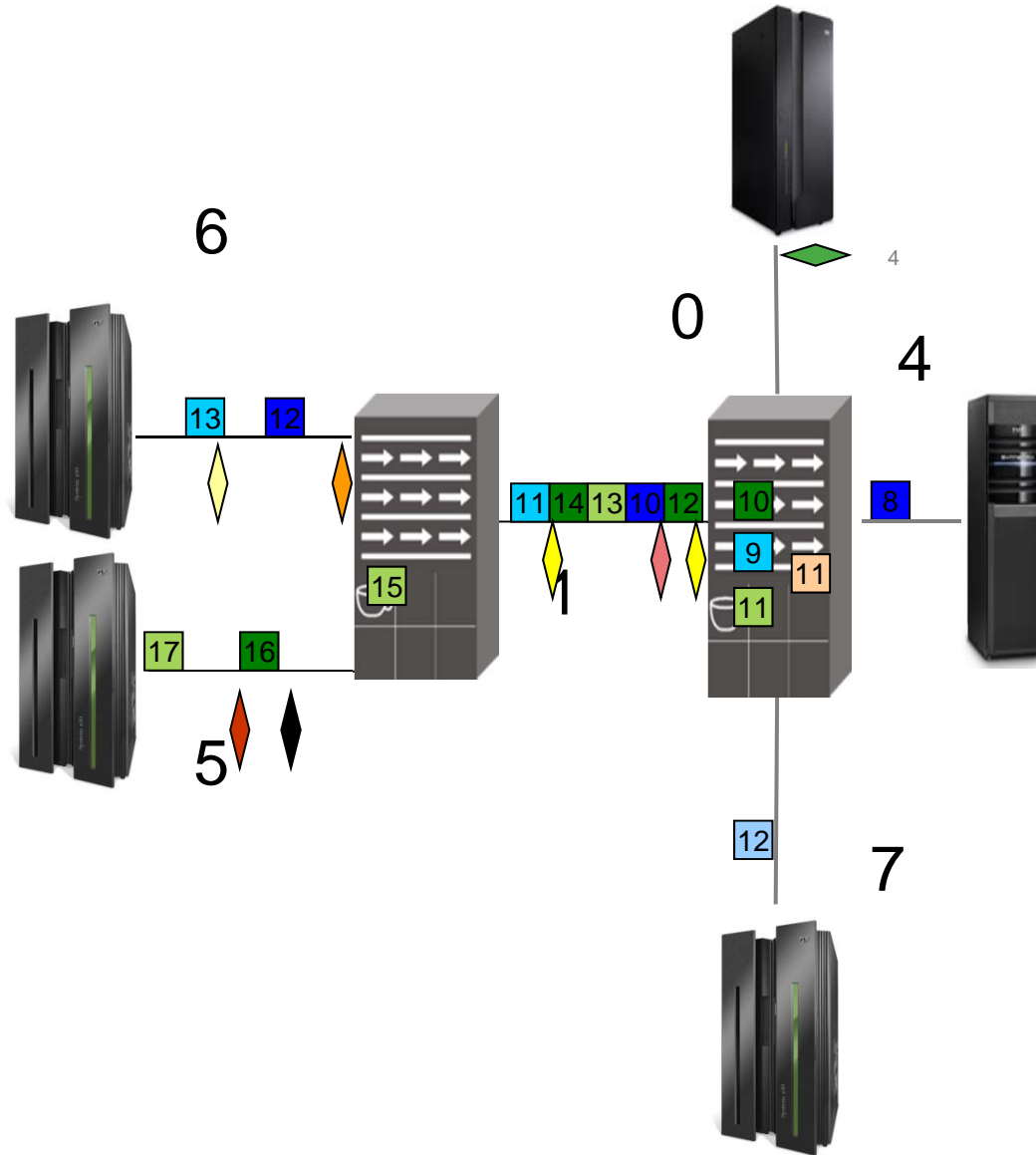


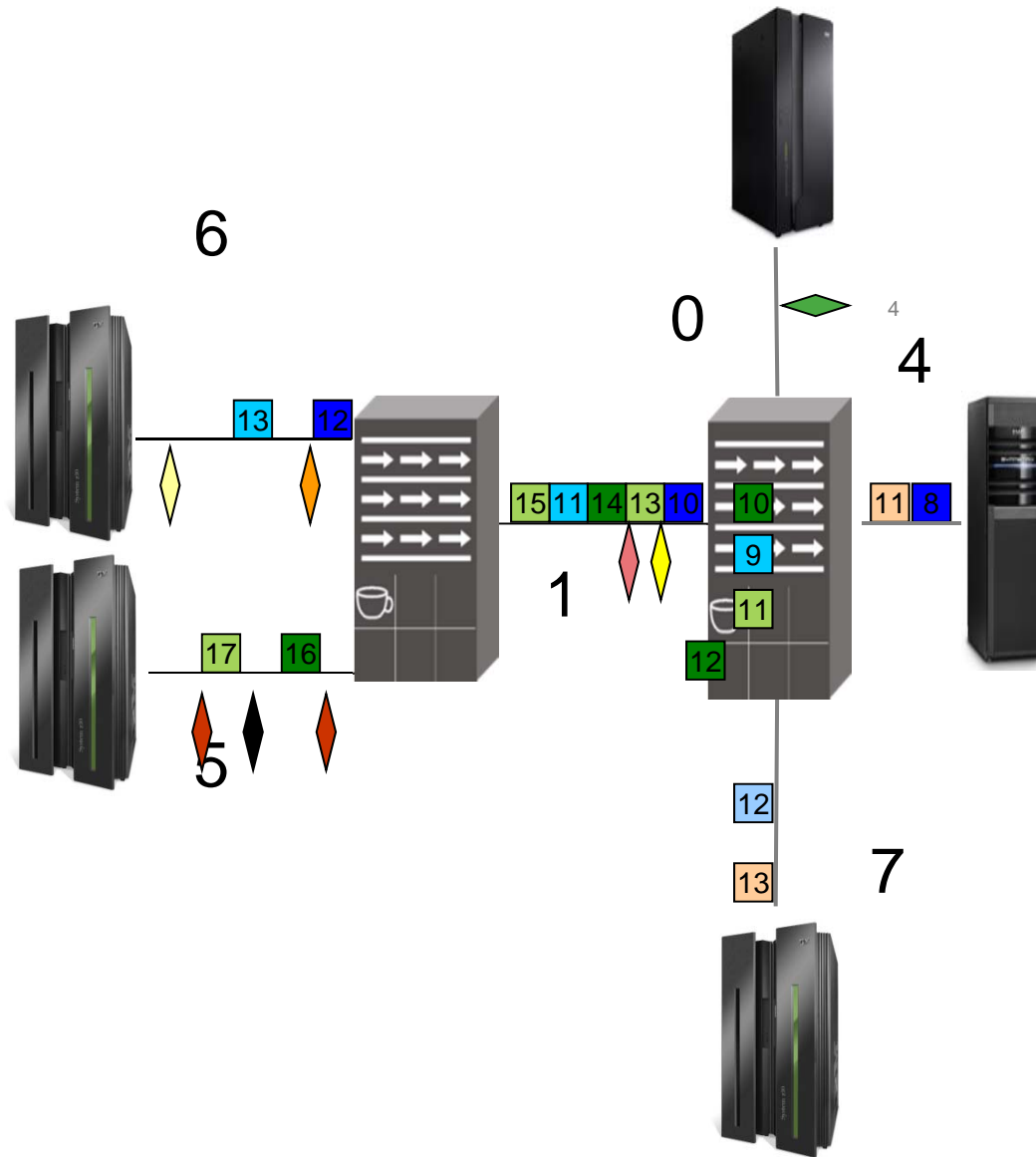


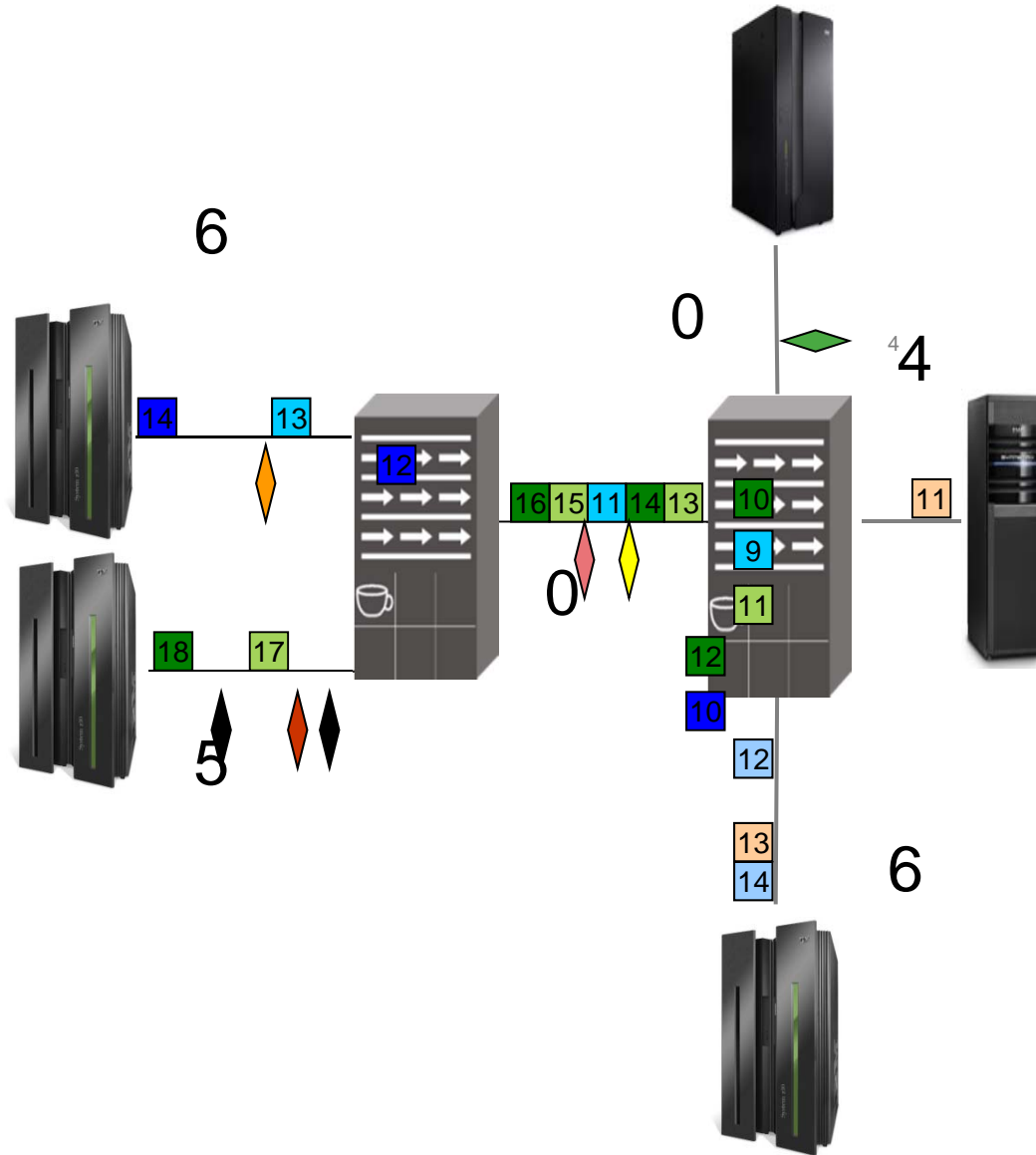


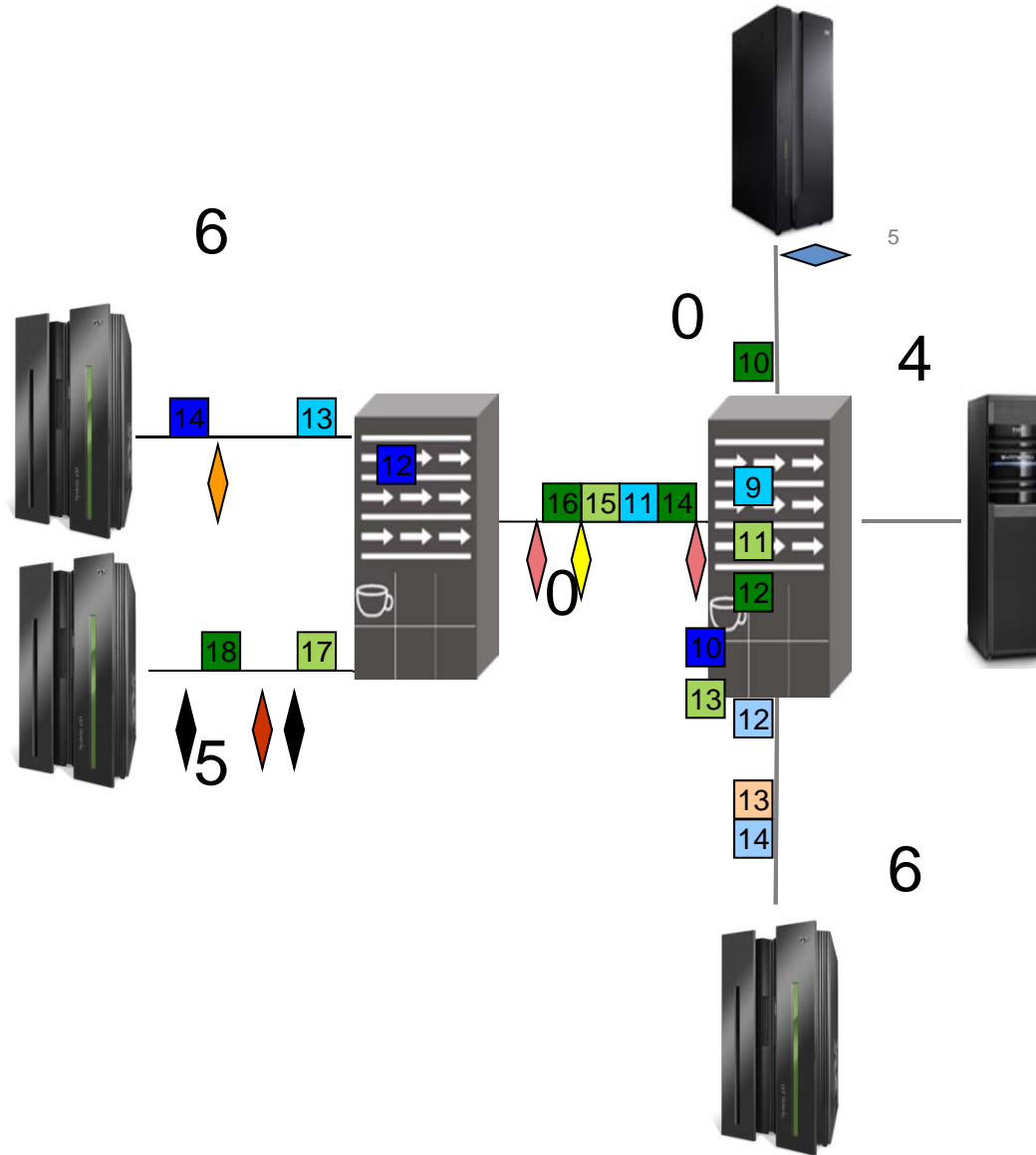


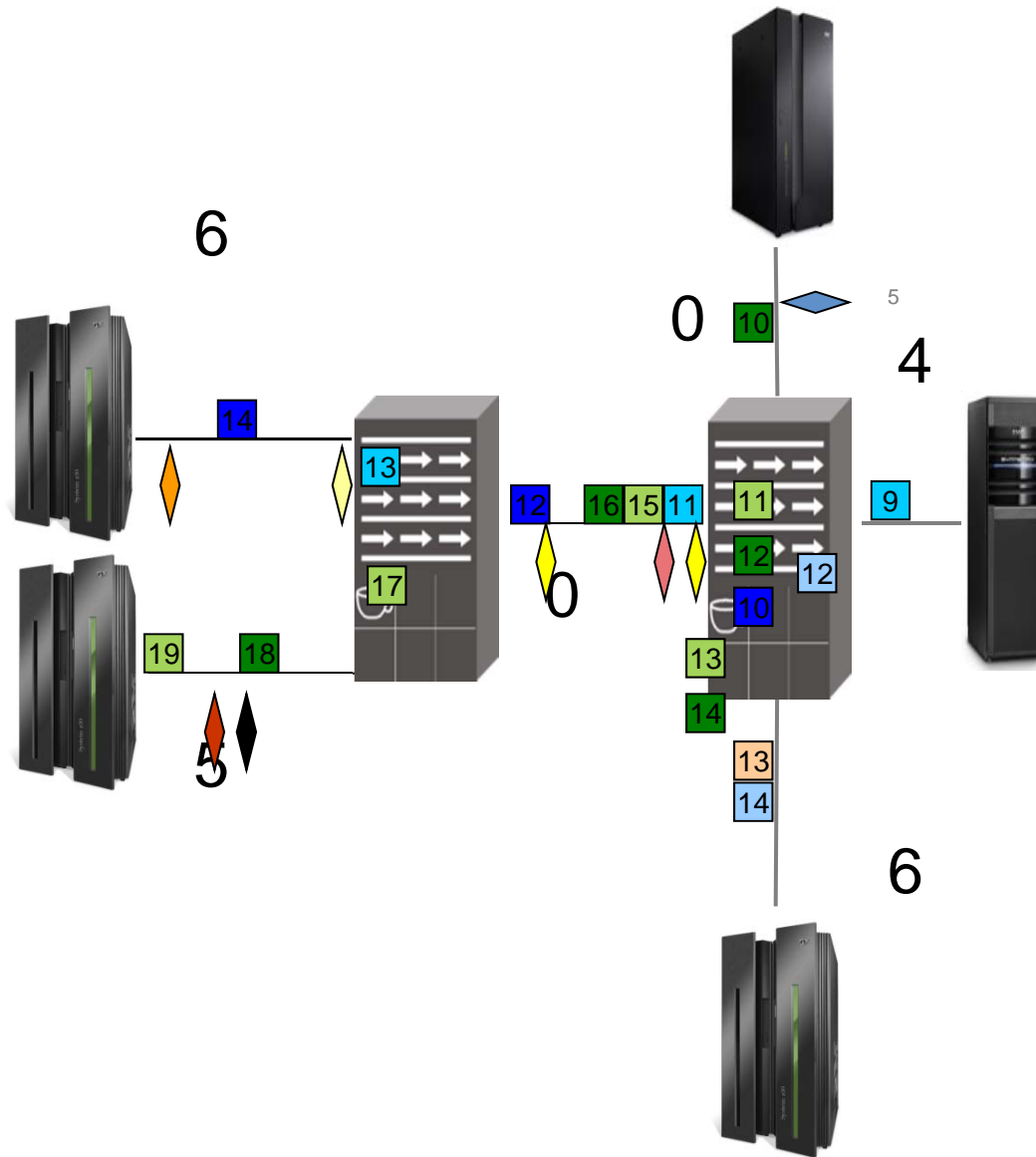


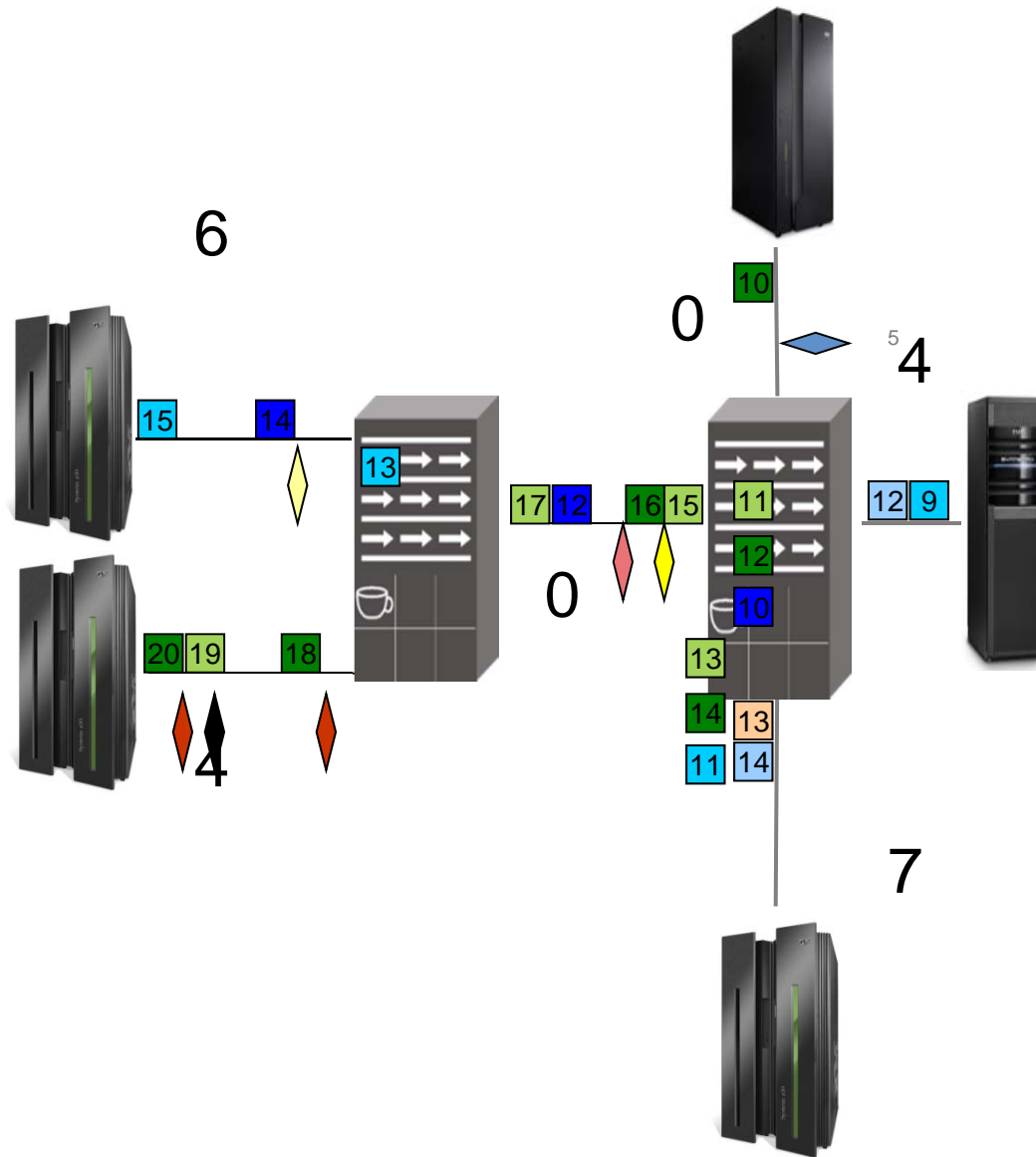


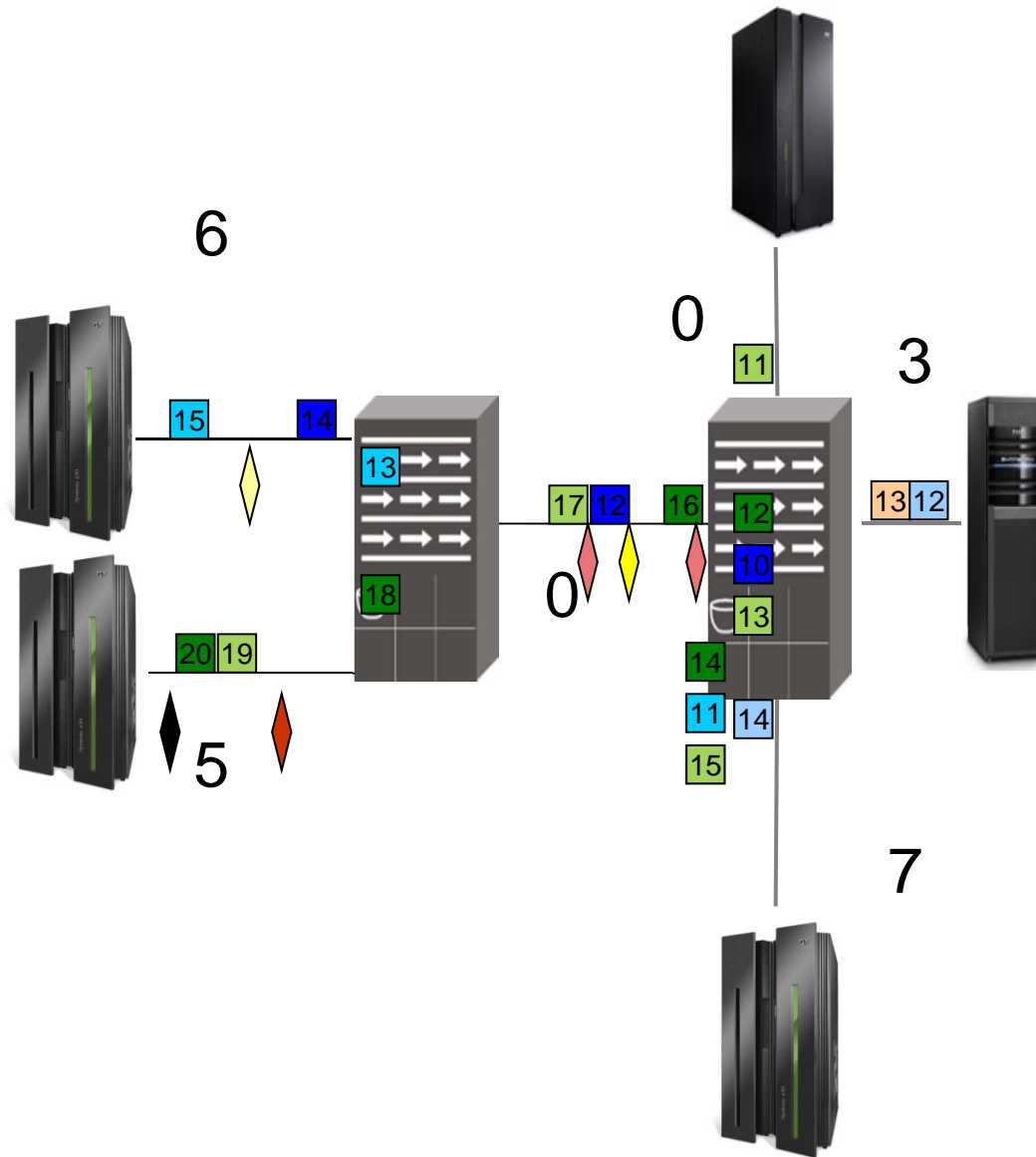


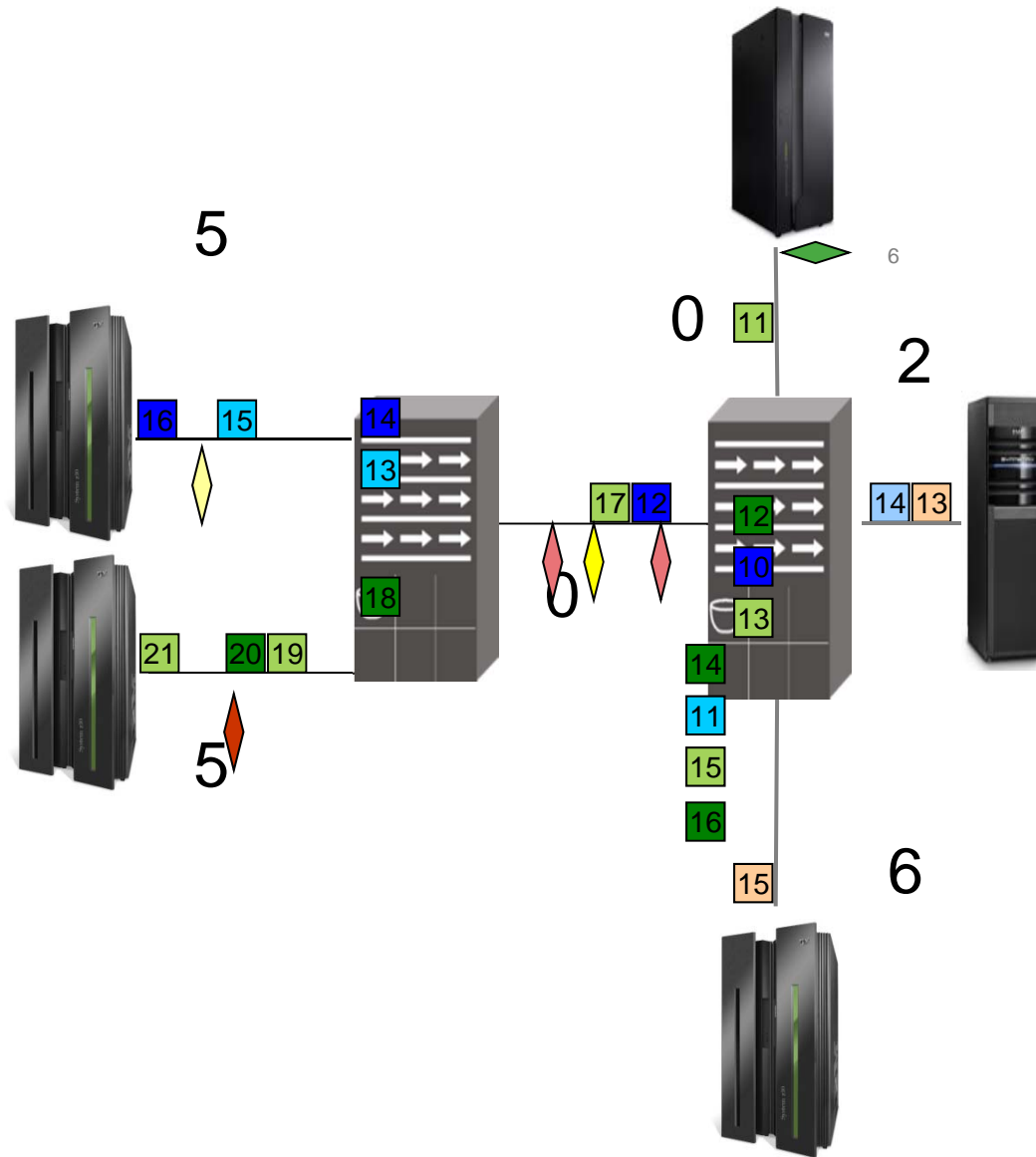


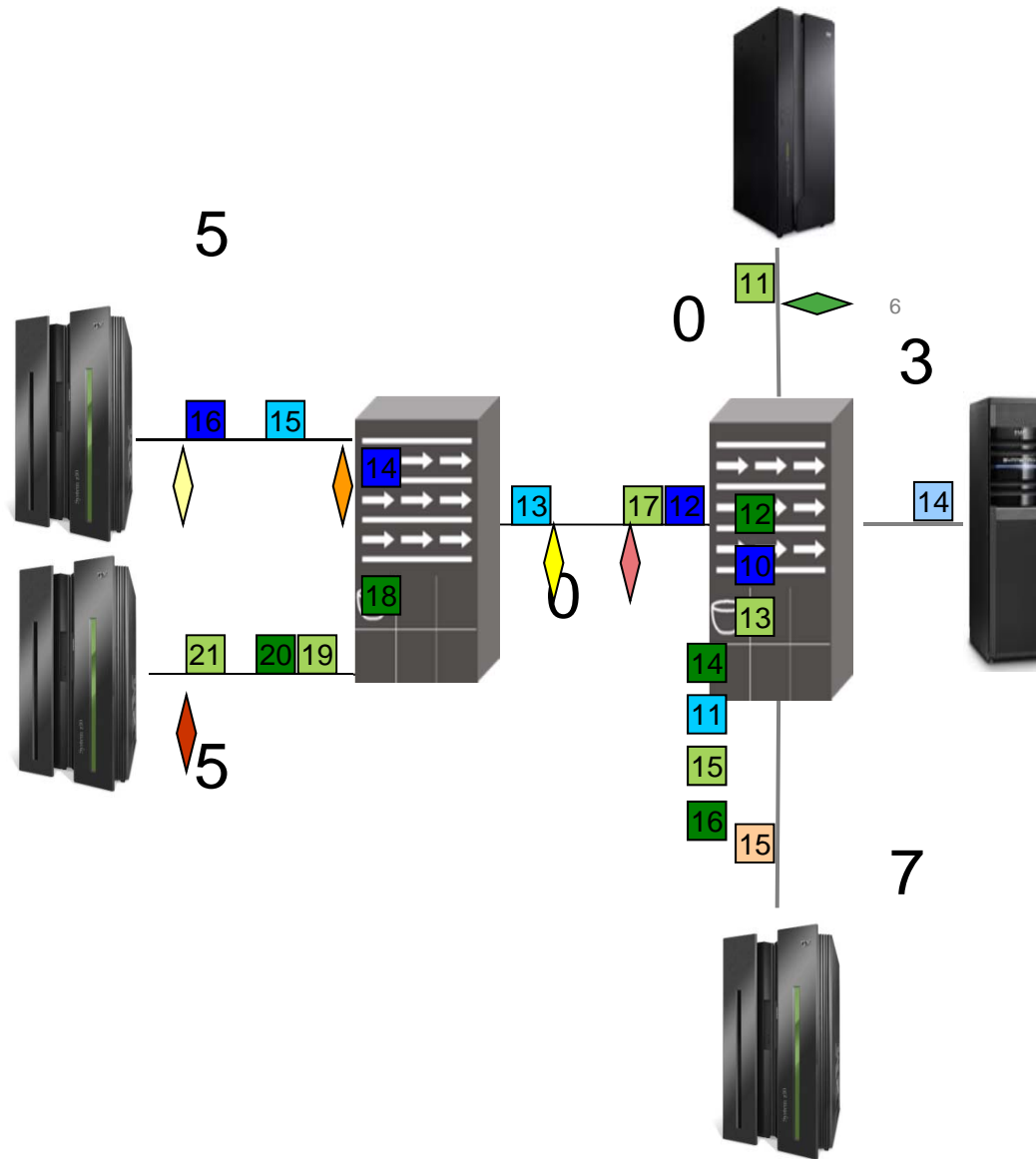


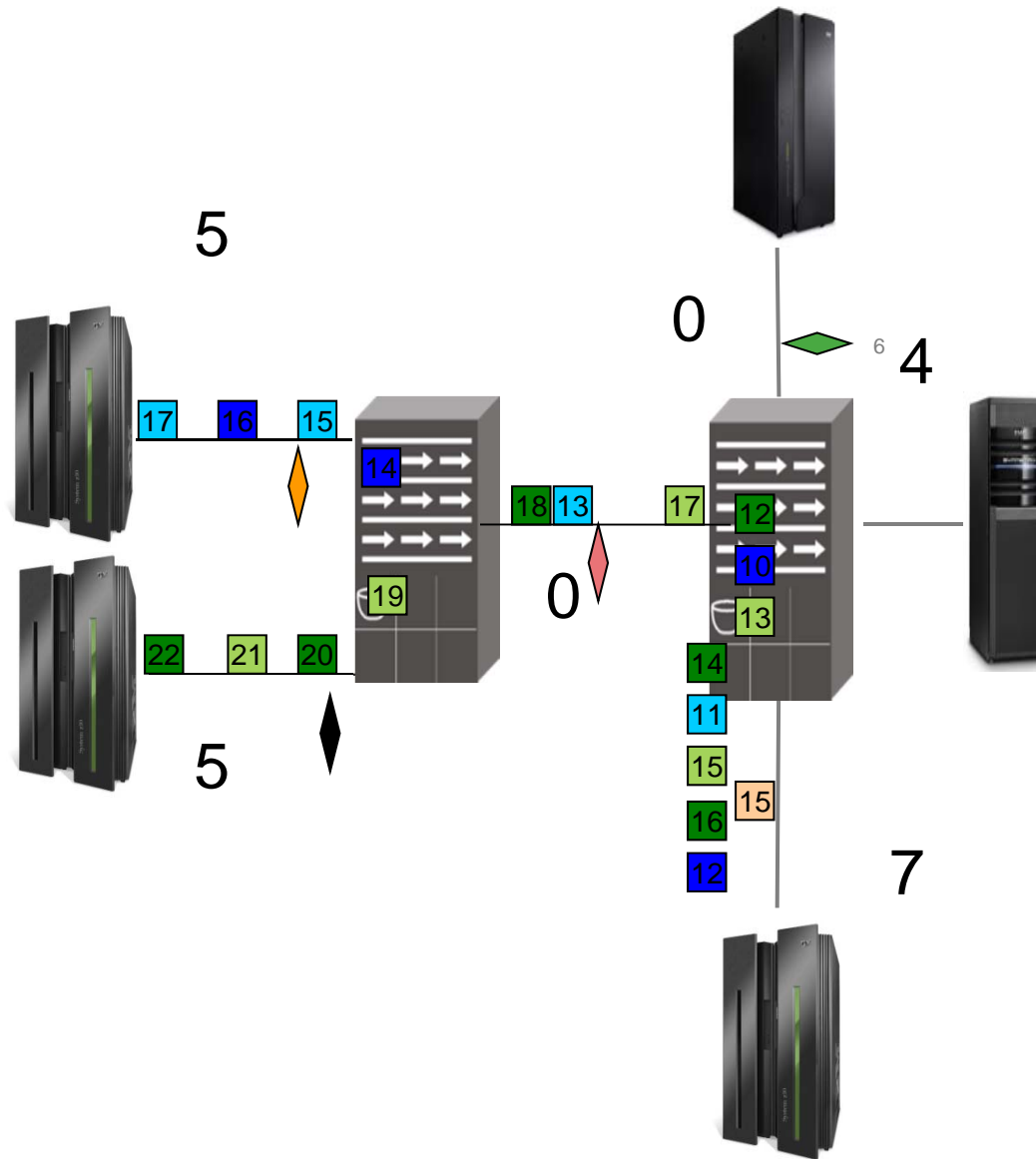


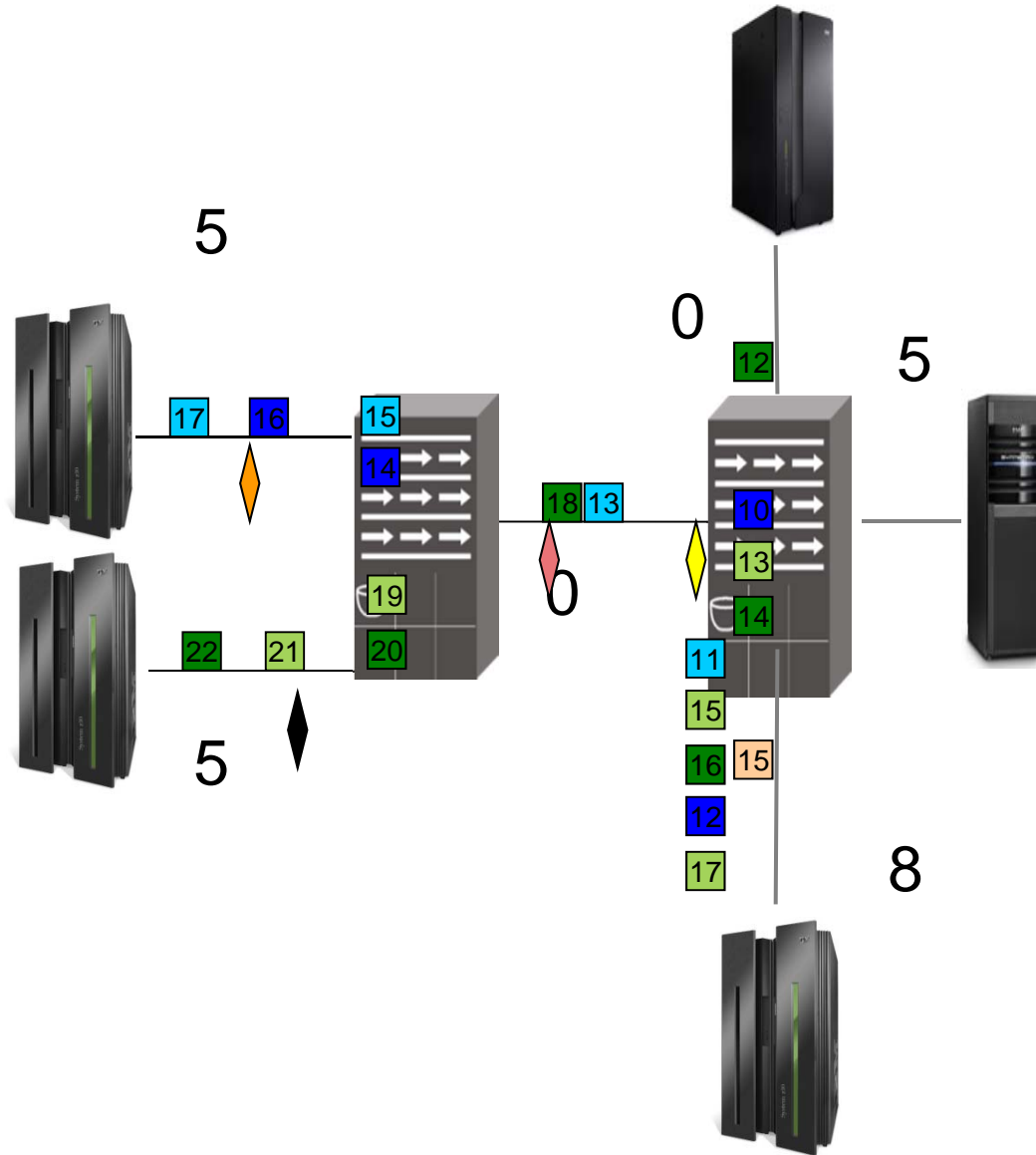


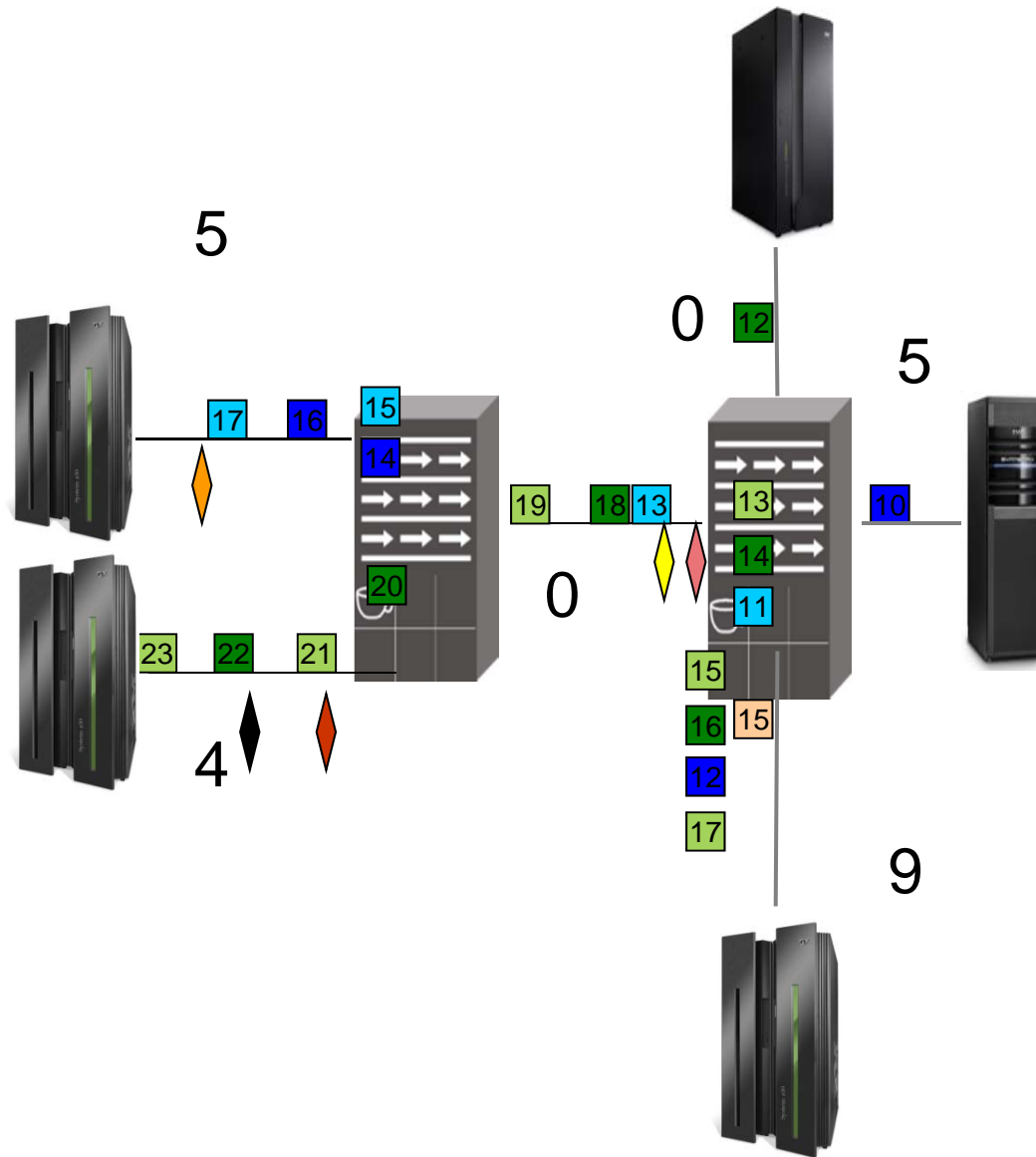


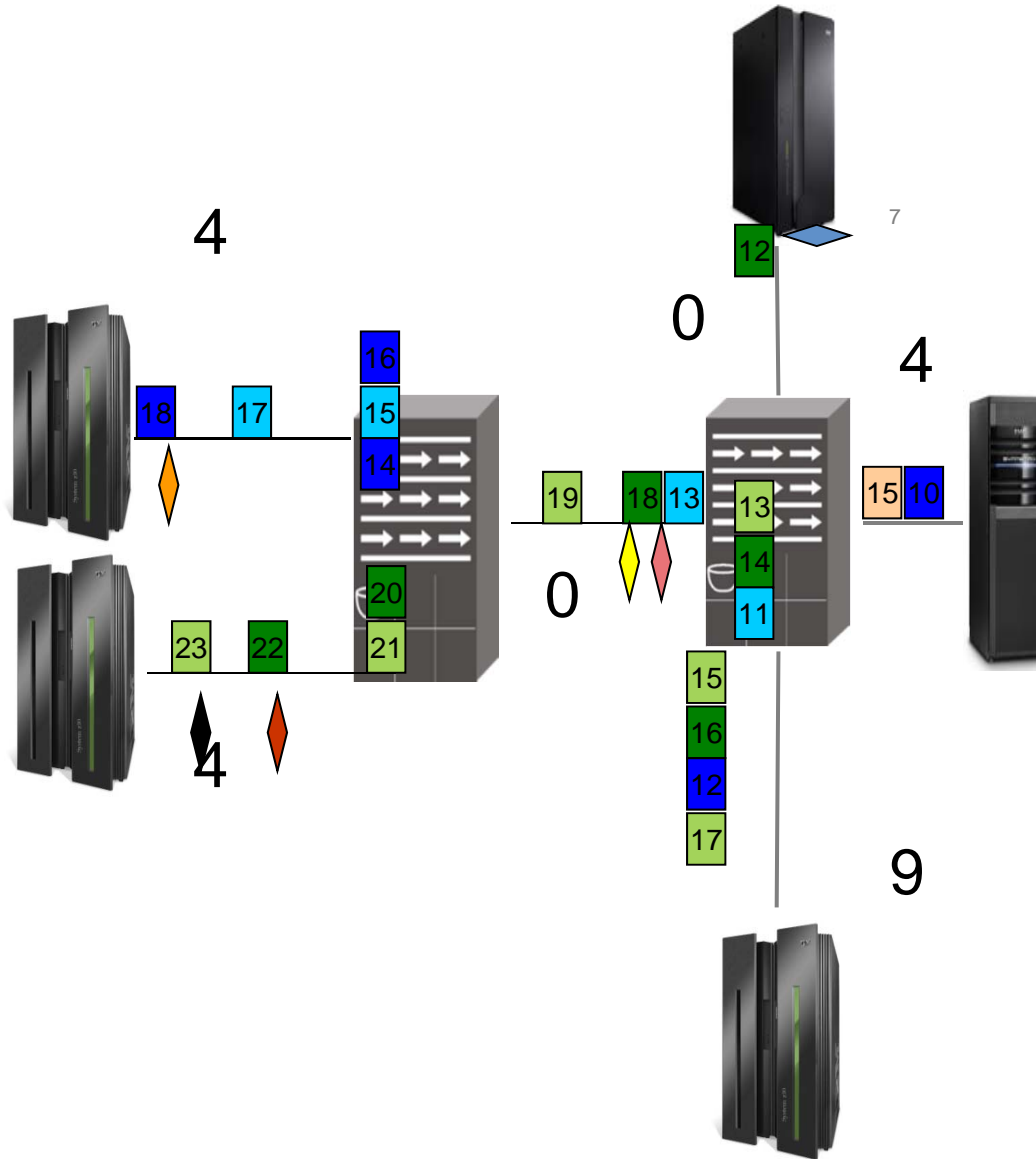


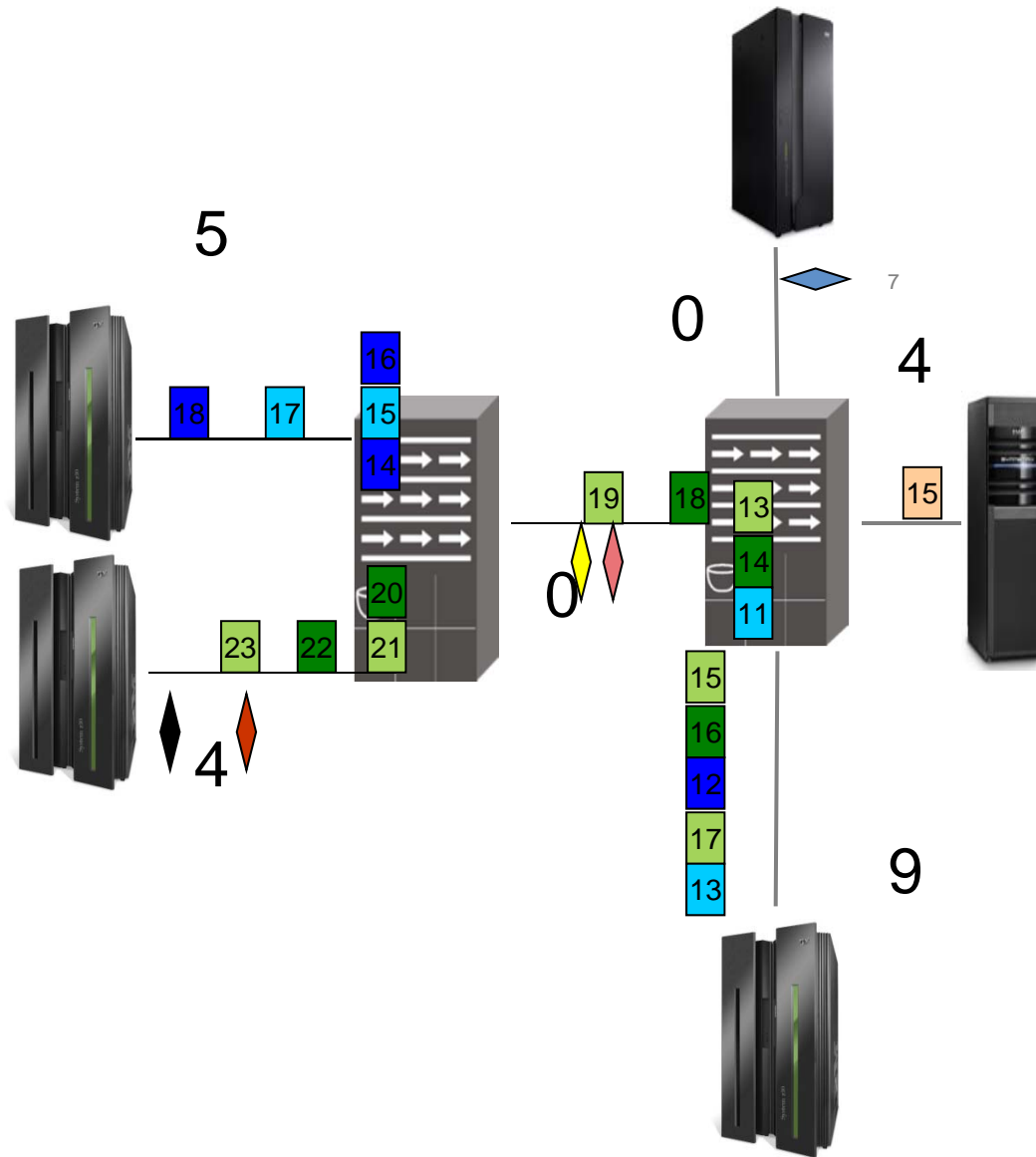


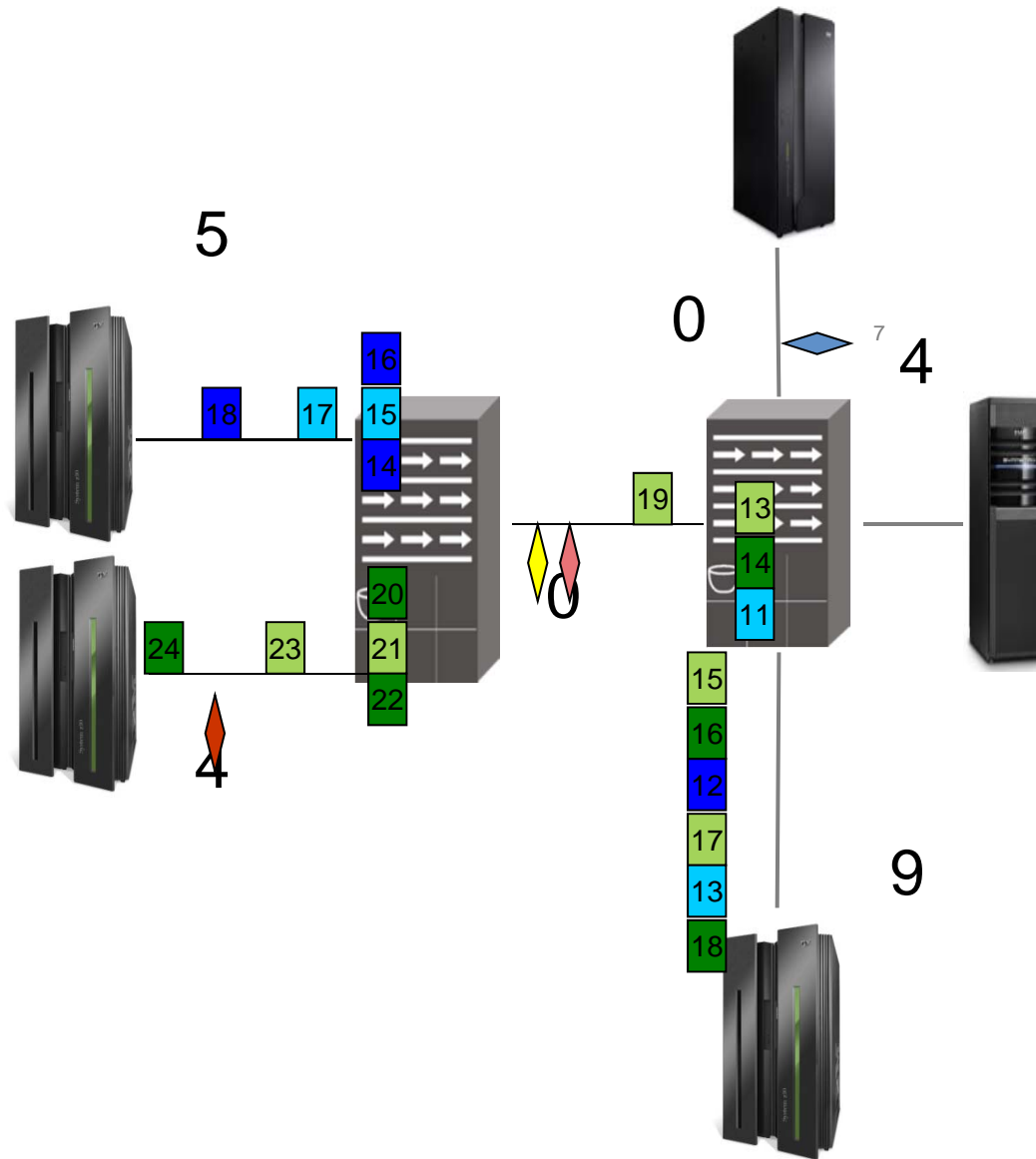












What Can Be Done About A Slow Drainer?



- Difficult to identify the culprit
 - Sometimes issues are intermittent
 - There may be more than one gremlin

- Tools are available to help!
 - “Bottleneck Detection”

See Also:

https://elabnavigator.emc.com/vault/pdf/FC_SAN.pdf

https://www.ibm.com/developerworks/mydeveloperworks/blogs/sanblog/entry/how_to_deal_with_slow_drain_devices20?lang=en

http://www.brocade.com/downloads/documents/html_product_manuals/NA_SAN_1130/wwhelp/wwhimpl/common/html/wwhelp.htm#href=Ch_Performance.31.7.html&single=true

Complete your sessions evaluation online at [SHARE.org/BostonEval](https://www.share.org/BostonEval)



What Can Be Done About A Slow Drainer?



- Creates an alert when a port(s) have zero credit for a user selectable amount of time
- Use detection mechanisms to monitor ISLs ports and F_Ports

SAN Performance Gotcha's



- Oversubscription
 - Trying to send more than can be received
- Slow draining
 - Unable to receive what is being sent in a timely manner

Culprits

- Oversubscription
 - Big pipe to small pipe
- Signal degradation
 - Failing SFP
 - Poor connections
 - Insufficient cable hygiene
- Unexpected slow down
 - Big pipe to small pipe (16G to 2G)
 - Narrow bridge (reduced ISL bandwidth)
- Odd stuff
 - Long links
 - Zoning conflicts
 - Invalid Attachment / Decommissioned Channel/CU Port
 - Queue depth

Stuff to Look For

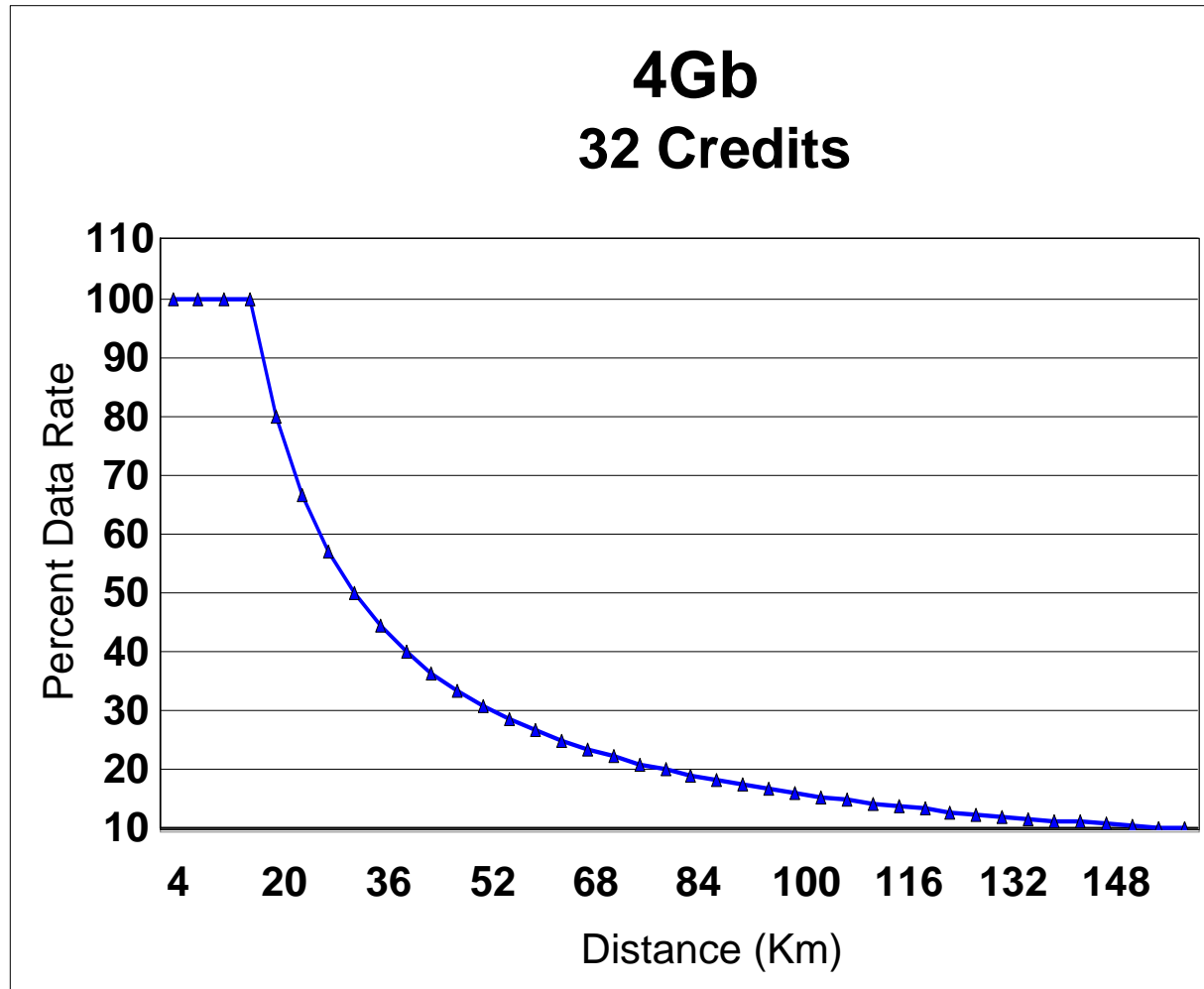
- Physical Errors
 - Detected (Obvious stuff)
 - Logins & Logouts
 - Link Failures
 - Undetected (Not Obvious)
 - CRC Errors
 - *Encoding errors (inside of a frame)*
 - Code Violation Errors
 - *Loss of synchronization*
 - *Loss of signal*
 - Loss of Sync
 - *Encoding errors (outside of a frame)*
- Logical Errors
 - Credit Starvation
 - Time at zero credit
 - Link Resets
 - Class 3 Discards
 - Response Times
 - High CMR
 - Exchange Completion Times
 - SCSI Reservation Conflicts

Other Neat Stuff

BUFFER CREDITS

Complete your sessions evaluation online at SHARE.org/BostonEval

How Does Credit Affect Throughput?



How Much Credit Do I Need For A Given Distance?

- **Absolute Minimum Formula:**

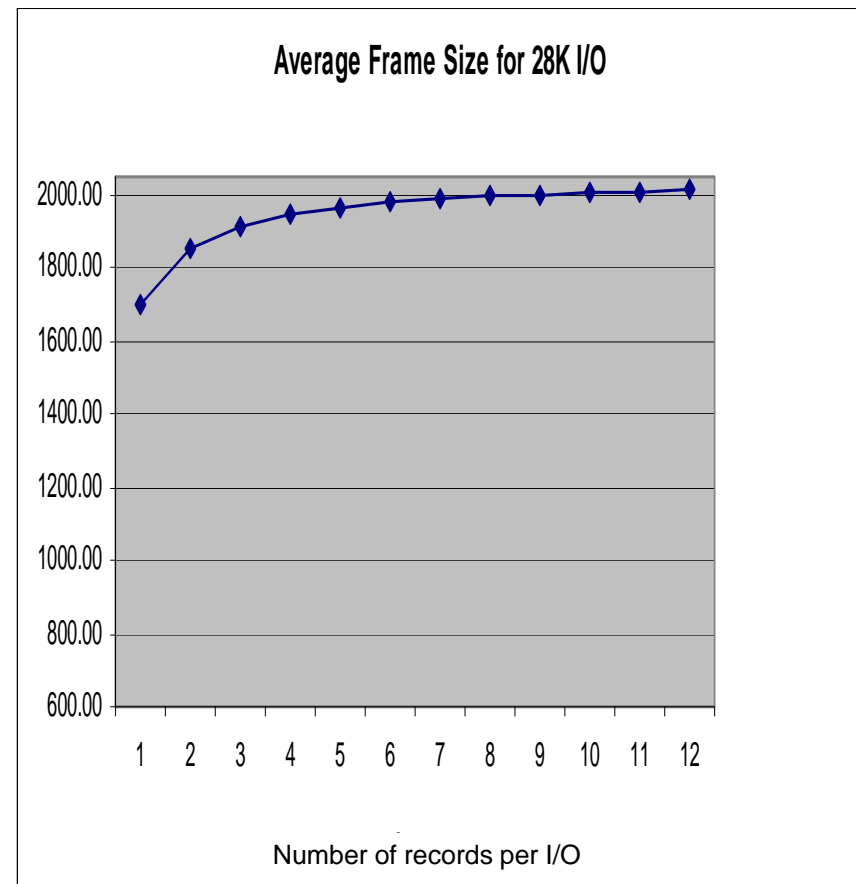
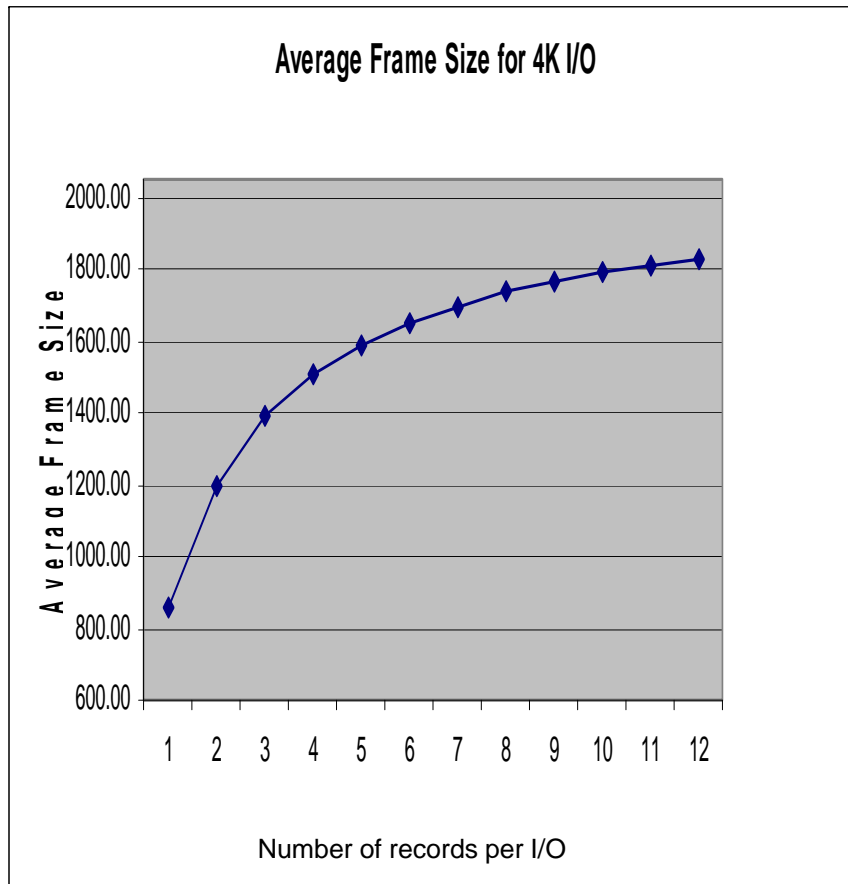
Number of credits needed = $1 + \frac{\text{Link speed in Gb/s} * \text{Distance in Km}}{\text{Average Frame Size in KB}}$

↑
This "1" is to account for the processing delay in a real device

↖
This accounts for round trip time on the fiber

↖
This is the **HARD** part

Average Ficon Frame Size vs Block Size



How Should I Configure My BB Credit?



- Local links under 1Km in length:
 - Even at 8G and small frame size, < 16 credits required

$$1 + ((8*1) / .8) = 11$$

- Provision ISLs with extra headroom
 - Use formula, and round UP

$$1 + ((8*10) / 1.5) = 54$$

EXCHANGES

Complete your sessions evaluation online at SHARE.org/BostonEval

Urban Legend: FICON uses fewer Exchanges Than FCP

- In Ficon, each concurrent I/O operation uses two Exchanges
 - One unidirectional Exchange for IUs from the Channel to the CU
 - A different unidirectional Exchange for IUs from the CU to the Channel
- The PAIR is commonly know as a “Ficon Exchange”

How many Exchanges do I need?

- Little's Law states:
 - *The number of “things” in a system can be determined by multiplying the average arrival rate of those “things” by the average time each “thing” stays in the system.*
- Applied to Ficon:
 - The average number of Exchanges active at any given time = Average I/O rate * Average response time
 - Example: 5000 Ficon I/Os / Second on a given channel with .4ms service time¹ needs 2 Active Exchanges (pairs) at any given time

¹ The amount of time the I/O is active in the channel

Urban Legend: Ficon is More Sensitive to Errors than FCP

- Is a Ficon frame more likely to get lost, damaged or corrupted than FCP?
 - No, the probability is the same
- When a Ficon frame gets lost, damaged or corrupted, is the recovery action different from FCP?
 - Not really. They are both retried, FCP by the Device Driver, Ficon by IOS/ERP

So What are the Differences?

- z Operating Systems tend to provide more detailed messages
- Ficon does provide additional debug data and actions
 - RNID
 - Link Error Status Blocks
 - Extensive State Change Processing

Table 89 - Link Error Status Block format for RLS command

| Word | Bits | 31 | .. | 00 |
|------|------|-----------------------------------|----|----|
| 0 | | Link Failure Count | | |
| 1 | | Loss-of-Synchronization Count | | |
| 2 | | Loss-of-Signal Count | | |
| 3 | | Primitive Sequence Protocol Error | | |
| 4 | | Invalid Transmission Word | | |
| 5 | | Invalid CRC Count | | |

Source: FC-FS-3 INCITS/T11 Draft Standard v0.92

See www.t11.org

Complete your sessions evaluation online at SHARE.org/Bo

Summary

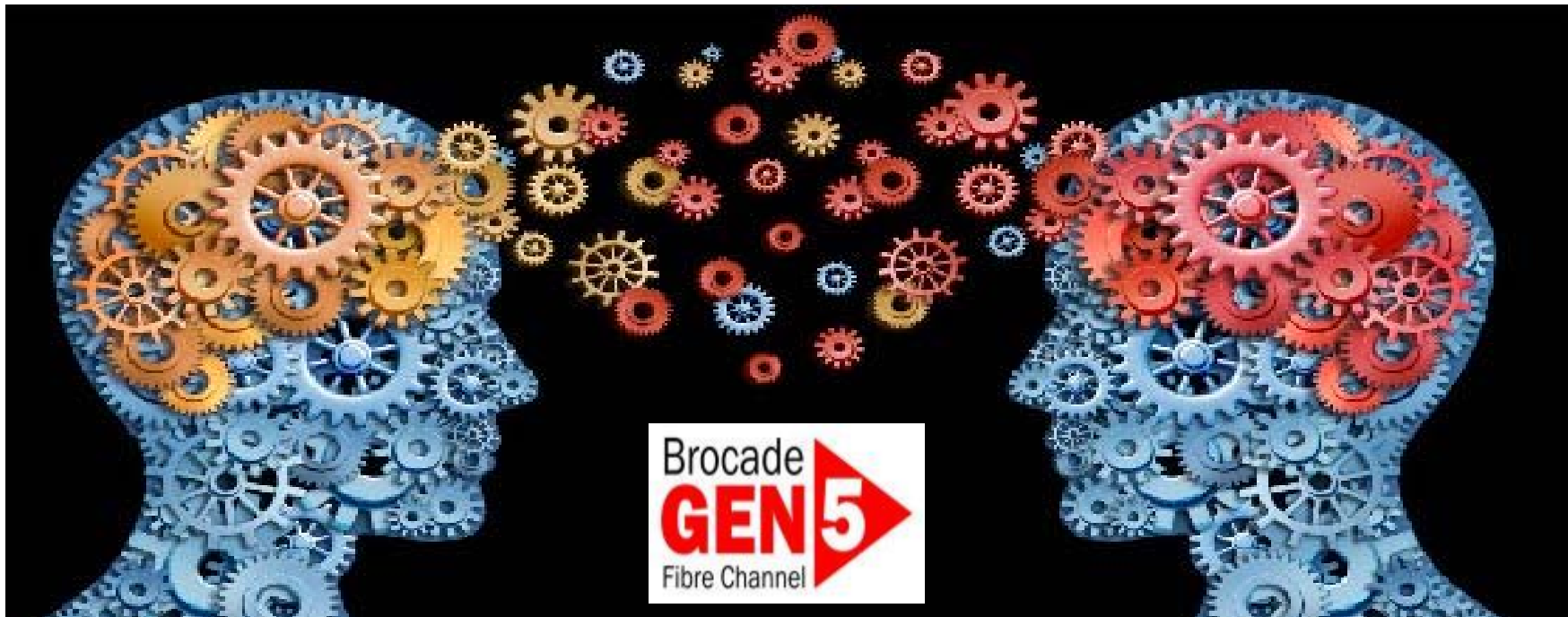
- Buffer Credits
 - Distance
 - Flow Control
- Exchanges
 - Unidirectional
 - Bidirectional

.....Newly Scheduled Presentation.....



Brocade SAN and FICON Update

**Please consider attending to discover the innovation
of Brocade's Gen 5 Fibre Channel Architecture**



Wednesday August 14, 2013 – 11:00pm to 12:00pm -- Session 14482



SHARE, Boston, August 2013

Buffer-to-Buffer Credits, Exchanges, and Urban Legends
Session 14281

THANK YOU!

QR Code



REFERENCES

Speaker Biography

- Lou Ricci
 - EMC
 - 1 Year
 - IBM
 - 34-years
 - 24-years in channel development
 - An inventor of FICON
 - FICON Firmware Team Leader
- Contact Information
 - Louis.ricci@emc.com

Speaker Biography

- Howard L. Johnson
 - BROCADE
 - Technology Architect, FICON
 - 29 years technical development and management
- Contact Information
 - howard.johnson@brocade.com

BONUS SLIDES

Complete your sessions evaluation online at SHARE.org/BostonEval

End to End Credit

- Device to Device Flow Control
 - Between source and destination
 - Not the links
 - Similar to buffer-to-buffer flow control
 - At N_Port Login
 - Report available receive buffers (EE_Credit)
 - Transmitter counts buffers transmitted (EE_Credit_CNT)
 - Receiver acknowledges frame (ACK)
 - *ACK 1 (a single data frame in a sequence) – most common*
 - *ACK n (several (N) consecutive data frames in a sequence)*
 - *ACK 0 (all data frames in a sequence) – not used*

Virtual Channels

- Technology to allocate BB_Credits to particular data flows
 - Class F traffic has one data flow
 - Assigned with Zoning by using special Zone names

