

Evolution of the System z Channel

Howard Johnson – hjohnson@Brocade.com

Patty Driever – pgd@us.ibm.com

8 August 2011 (4:30pm – 5:30pm)

Session Number 9934

Room Europe 7

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

- This session examines the evolution of the FICON channel from the birth of the industry to the looming converged infrastructure. The speakers will discuss the designs and modifications that allowed System z solutions to move from ESCON to FICON and beyond. Come see where the FICON channel has been and where it's going.

Special acknowledgement goes out to Dan Casper – “Mr. Channel” in the halls of System z Development. Dan was instrumental in the design of ESCON, FICON, and zHPF. Dan retired from IBM on July 31st and he and his contributions will be sorely missed.

Agenda

- Yesterday
 - Review of the origins of the channel and its early evolution
 - Parallel channels
 - ESCON channel
 - FICON bridge
- Today
 - Examine the current channel technology
 - Native FICON
 - High Performance FICON
- Looking Ahead
 - Explore the possible future of channel technology
 - 16G / 32G FICON
 - Fibre Channel over Ethernet

Yesterday

The origins of the channel and its early evolution



The Early Years

S/360 & S/370 I/O ARCHITECTURE

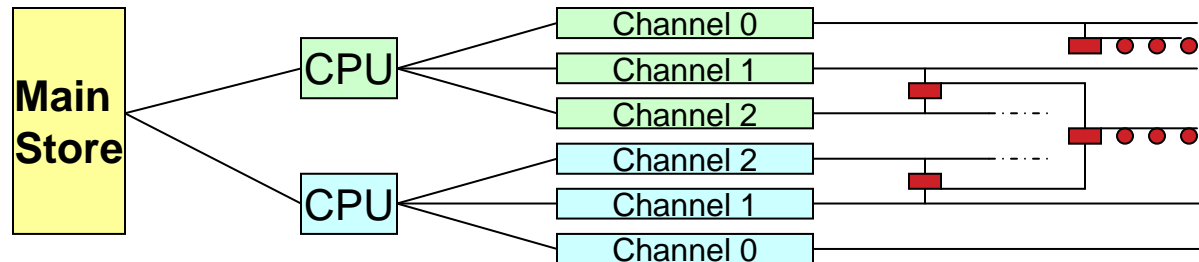
S/360

- First real computer **architecture**
 - Provided for ability to deliver a range of models with increasing price/performance characteristics
 - Fastest models used hard-wired logic, while microcode technology used to deliver wide range of performance within the S/360 family
 - Provide foundation that enabled application programs to migrate forward across models/system
 - Application level compatibility maintained through today on System z processors
 - Channels were specialized processors driven by a special instruction set that optimized the transfer of data between peripheral devices and system main memory
 - Channel architecture defined mechanism to transfer data, but was independent of device architecture
 - S/360 systems had one **byte** channel (channel 0) and one to six **selector** channels
 - New family of I/O devices developed for S/360 that used new standardized 'bus & tag' interfaces



S/360 & S/370 I/O Architecture

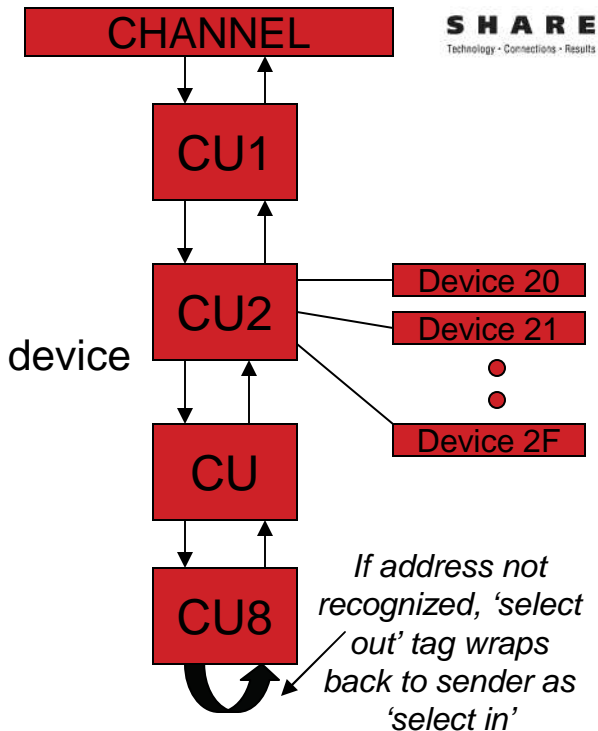
- Designed to provide value in key mainframe focus areas: security, resiliency & performance
 - Security/Integrity
 - Host-based configuration definition methodology introduced in S/370
 - *Security controls based on host definitions (IOCDs)*
 - *4 digit device number = Channel + Unit Address*
 - CU enforced atomicity
 - *Reserve/release*
 - *Extent checking*
 - Resiliency
 - Channel set switching (S/370)
 - *Allowed OS to connect the set of channels of the failing CPU to another CPU*
 - Performance
 - Bi-directional data transfers (reads/writes) and transfer of non-contiguous portions of the disk in a single I/O operation
 - *Command chaining*
 - *Data chaining*
 - *Indirect Data Address Words (S/370)*
 - Asynchronous notification for unsolicited events
 - *Busy/no longer busy status*



Parallel Channels



- Introduced with S/360 in 1964
 - Circuit connected (multi-drop CUs & devices)
 - 1 device at a time (max 256 devices per channel)
 - No dynamic switching
 - Initially support ~200 ft distances between channel and device
- Typical “bus & tag” channel/CU communications:
 - **‘Initial selection’** sequence to establish a connection with a control unit
 - **‘Data transfer’** sequence (CU always controlled data transfer)
 - Each byte of data sent requires a response
 - *Read: CU says “I’m sending a byte”, and channel says “I received a byte”*
 - *Write: CU says “I’m ready for a byte” and channel says “Here’s a byte”*
 - **‘Ending sequence’** (when one side recognizes that all of the available or required bytes have been transferred)
- Parallel channels are distance/speed limited, because of skew between the transmission lines in the channel
 - When a tag is received on one side, all of the bits of the data on the bus must be valid at that moment



Parallel Channels...initial types

- Byte Multiplex
 - Could have several devices performing a data transfer at a time because each device could disconnect between bytes of data transferred
- Selector
 - These were used with devices that have high data transfer rates, requiring the channel to remain connected to the device until the entire chain of CCWs is executed
 - Tape devices are typical examples

Parallel Channels...evolving types



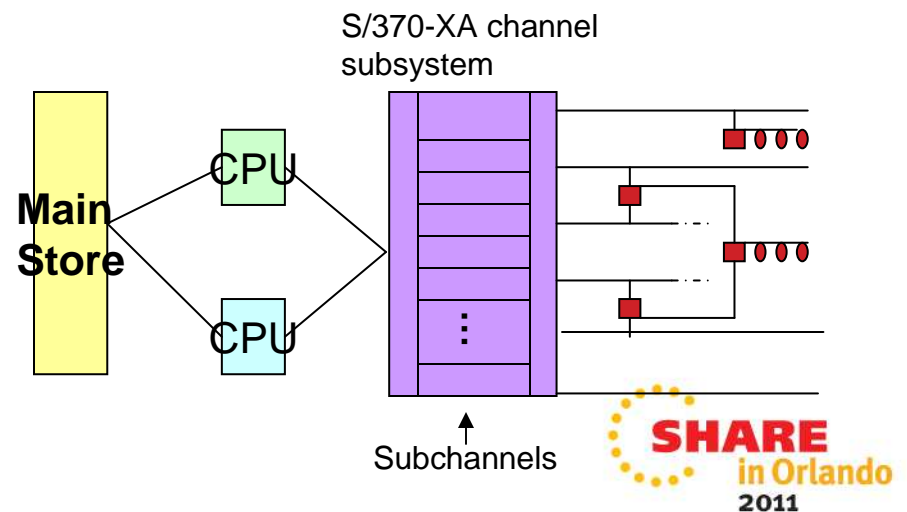
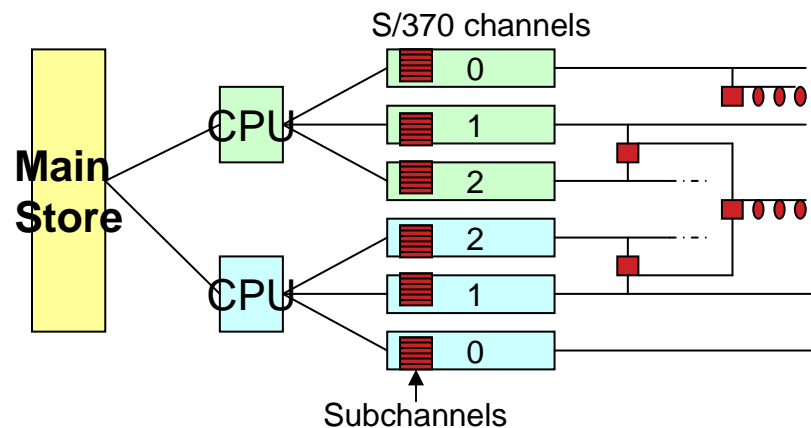
- Block Multiplex
 - Devices could disconnect only after the entire block of data for the command is transferred
 - When device is ready it presents 'request in' to request to be serviced again
- Data streaming
 - Added a new tag that enabled removal of the interlock between the channel and control unit before the next byte of data could be transferred
 - Each byte was still acknowledged, but CU kept count of responses to know how much was transferred
 - Enabled 400 ft. distances

The winds of change

TRANSITION TO ESCON

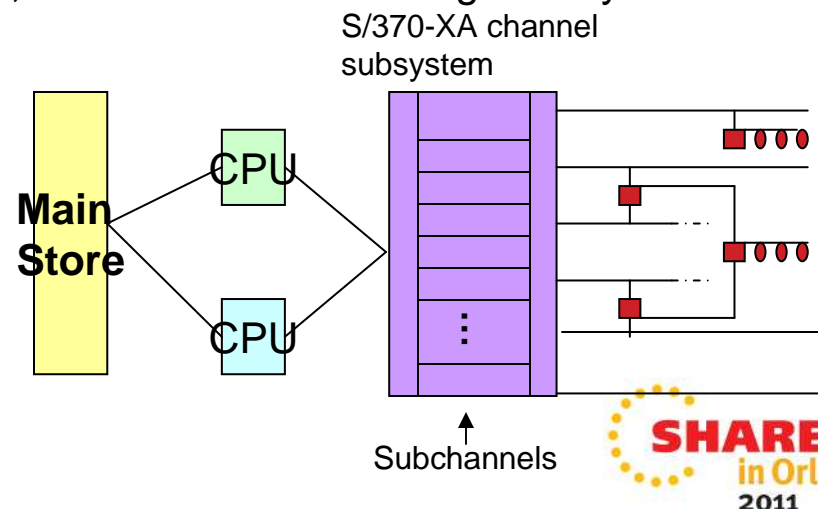
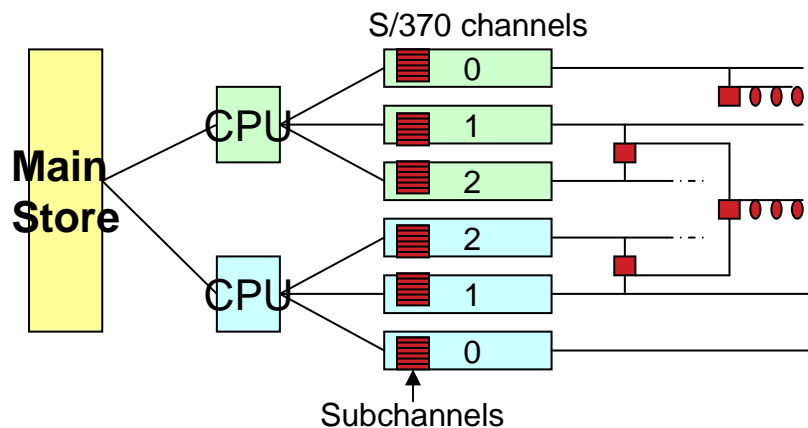
S/370-XA (Extended Architecture)

S/370	S/370-XA
Operating system selected the channel, and then the device on that channel	Operating system selected a device number (subchannel), and the channel subsystem knew all the addressing to get to that device and chose the best path
Channel can be addressed only by the single CPU to which it is connected and channel can interrupt only that CPU to which it is connected	Any CPU can initiate an I/O function with any device and can accept an I/O interruption from any device
Once a chain of operations is initiated with a device, same path must be used to complete transfer of all data, commands and status	Device can reconnect on any available path



S/370-XA (Extended Architecture)

- Enabled improved performance
 - All I/O instructions are executed asynchronously by the CPU with respect to the channel subsystem
 - All I/O busy conditions and path selection are handled by the channel subsystem rather than the CPU
 - Reduced CPU overhead in processing no-longer-busy conditions and encountering (possibly recurring) busy conditions on path selection
 - Eliminated program differences required to manage channels by type (e.g. selector vs. multiplexor)
 - Reduced the number of conditions that interrupted the CPU
 - Channel-available, control-unit-end, and device-end no-longer-busy eliminated

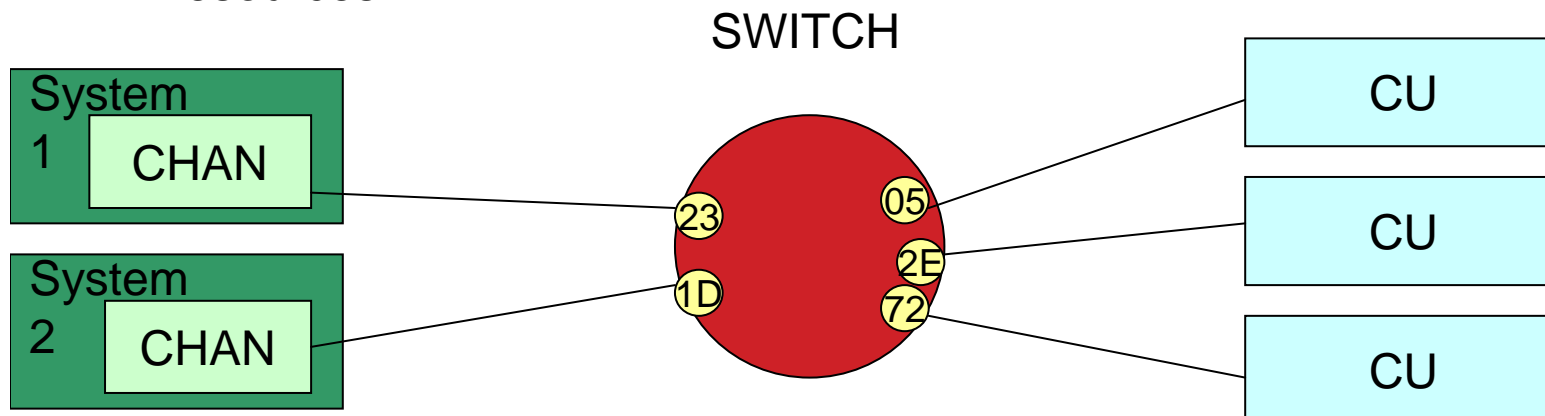




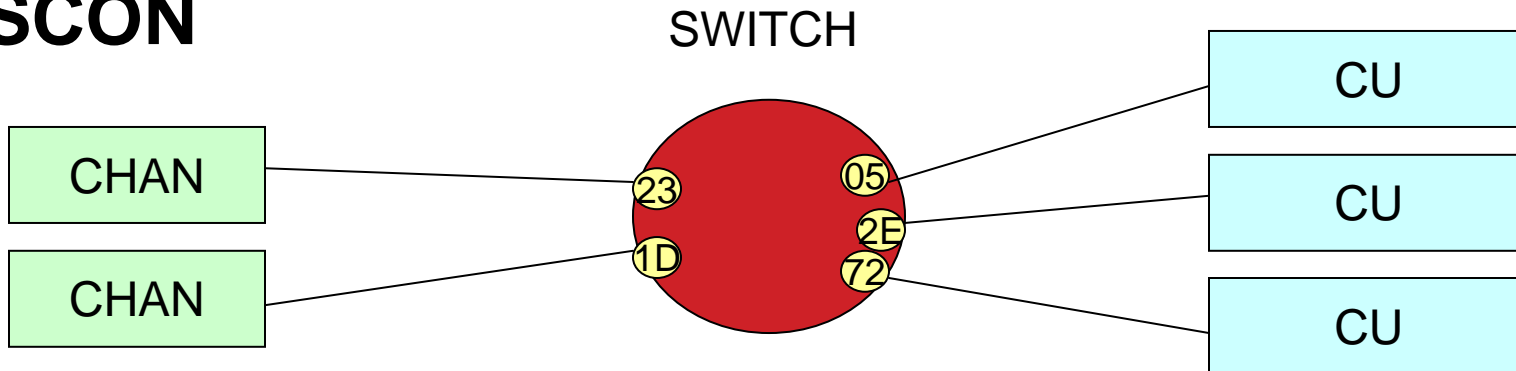
“A New Generation”
ESCON

ESCON

- Introduced in 1990 on the 3090 system
 - Used fiber optic (serial) connectivity instead of copper
- Link technology supported 20MB/sec, but achieved ~18MB/sec
- Supported introduction of the first modern SAN with dynamic switching capabilities
 - But still circuit-connect (communication with one device at a time)
- Distance limitations improved from 400ft to 9km (3km per link)
- Supported 1000 devices per channel
- I/O Configuration integrity checking with reset event architecture and self-describing devices (RNID)
- Fixed link addressing used in IOCDS controlled host access to CU resources



ESCON



Resiliency

- State Change Notification
 - Sent to each link level facility (channel or CU) which is potentially affected by state changes in switch ports or connected channels or control unit

Security

- Read Node Identifier
 - Enables host program to determine the specific physical device attached at the end of the link and revalidate attachments after link down conditions

Performance

- Streamlined command chaining
 - In parallel channels, when commands were chained together, after status was sent in for one command the channel would go through the selection sequence again
 - With ESCON, the channel responded to status directly with the new command
- Streamlined data transfer
 - Buffering in channels and control units removed some of the interlocks involved in data transfer

The winds of change

TRANSITION TO FICON

Late 1990s - What Was Going on Inside IBM

- Evolution from System/360 to System/390 saw a significant increase in MIPS, main storage, and I/O capacity
 - I/O capacity had largely been improved through the continued addition of channel paths, thus adding cost, complexity, and overall bandwidth without significant improvement in the capacity of a single channel
 - Original 7 channels provided by S/360 evolved to the 256 channels provided by S/390 at that time
 - 256 was the architectural, programming, and machine limit
 - Projected processor MIPS growth along with improved controller technology and increased I/O densities also drove the need to significantly improve the I/O throughput over a single channel path

What was going on in the Industry

- ANSI Fibre Channel had its beginnings about 1988-9
 - Initial motivation was for higher bandwidth I/O channels that operated efficiently at fiber optic distances (10s of kms)
- FC proponents argued early on that the requirements and technologies of LAN and channel applications were converging, and that a layered architecture could deliver excellent performance, cost, and migration characteristics for many applications
 - By late 1990s FC found a toehold as a storage interface

FICON – IBM Goals

- Significantly improve the I/O throughput over a single channel path
- Decrease the execution time of a single channel program
- Significantly improve the data transfer rates of a single channel path
- Provide for more addressable units (devices) attached to a single channel path
- Accomplish the above in a fashion that preserves the S/390 investment in existing programming
- Provide connections that support ESCON equivalent distances with negligible performance penalties at the increased distances
- Provide a **migration path** to support existing controllers currently deployed in the field
- Provide the above using existing industry standards where applicable, and develop new industry standards where needed

Comparing Characteristics

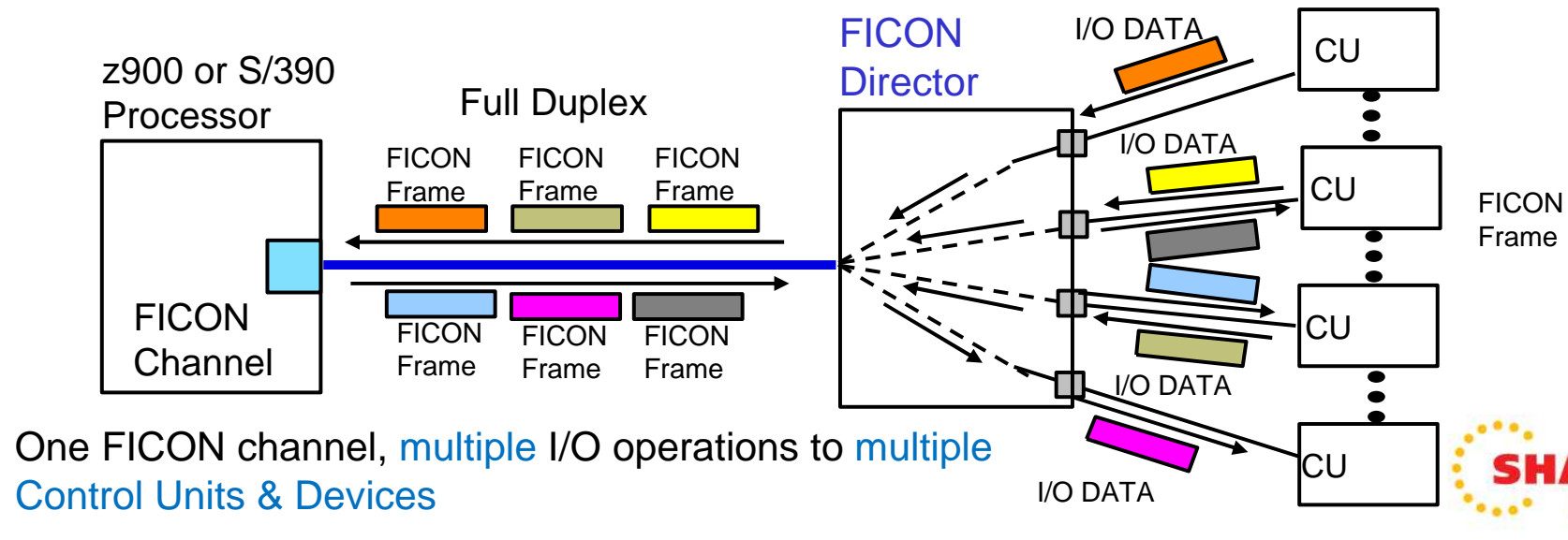
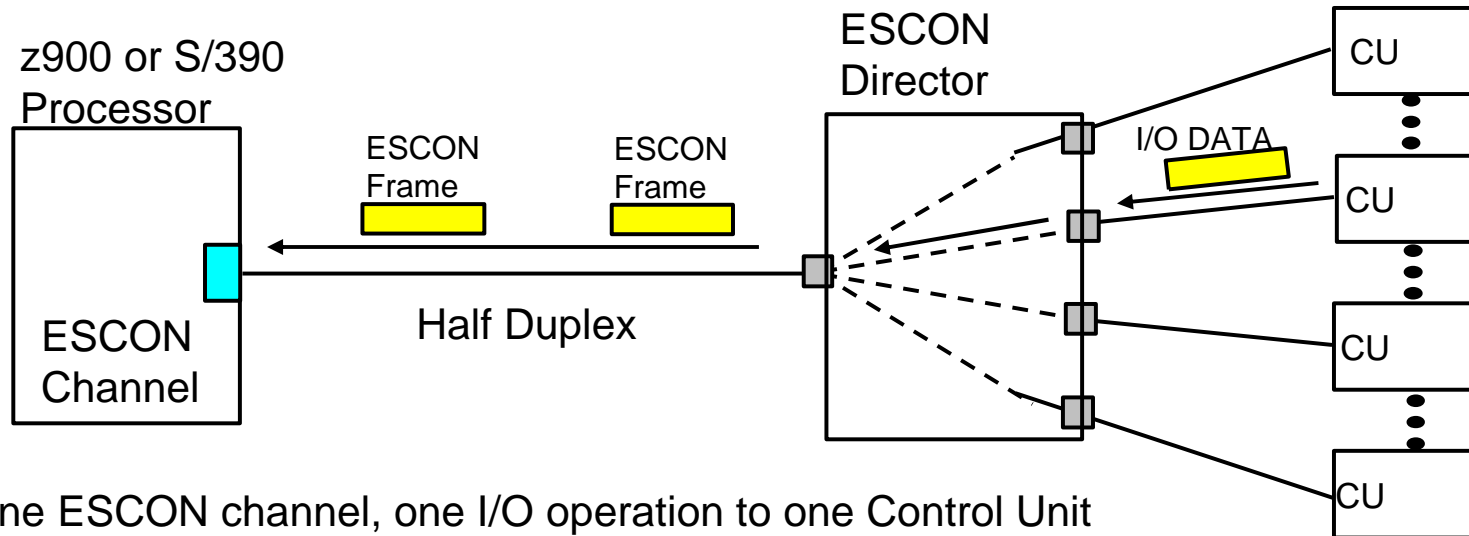
ESCON Channels

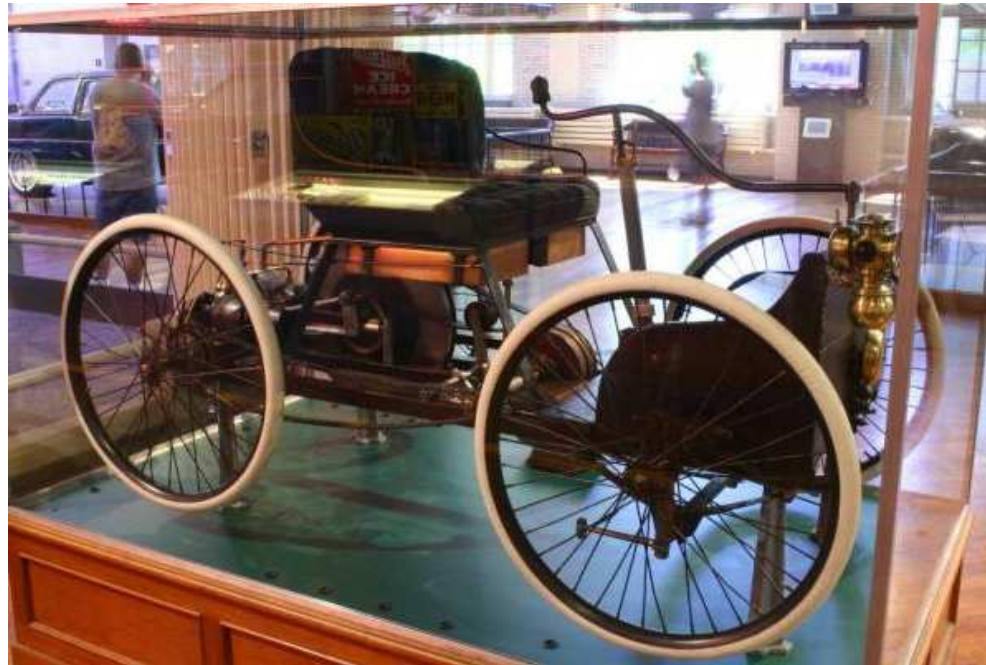
- Connection oriented
- Circuit Switching
- Read or Write
 - Half-duplex data transfers
- Dedicated path pre-established
- When packet is sent, the **connection** is locked
- Synchronous data transfer
- One operation at a time

FICON Channels

- Connectionless
- Packet Switching
- Simultaneous read/write
 - Full-duplex data transfers
- Packets individually routed
- When packet is sent, **connection** is released
- Asynchronous data transfer
- Pipelined and multiplexed operations

ESCON vs FICON Frame Processing



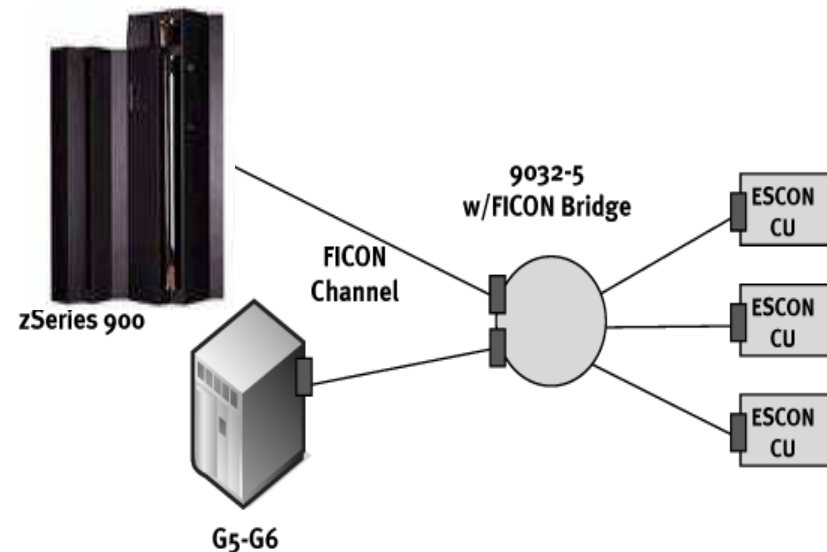


Almost a bridge too far

THE FICON BRIDGE

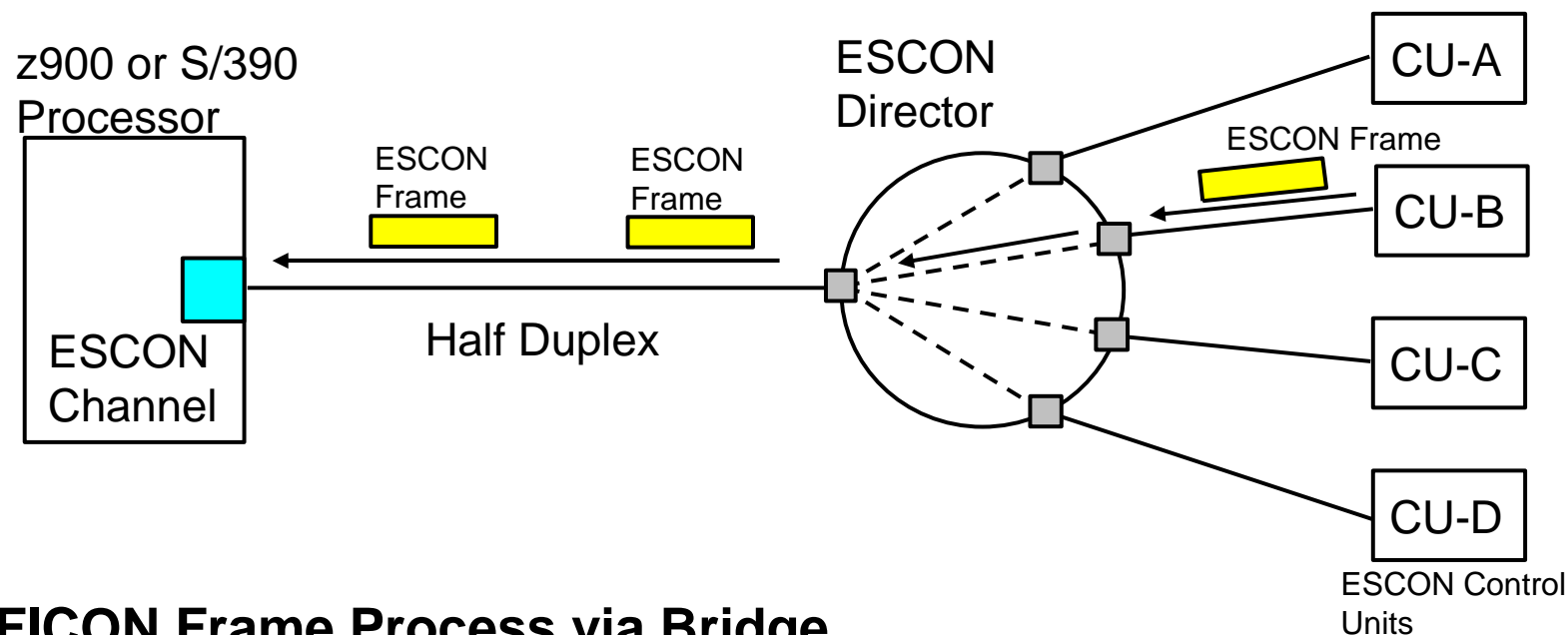
The “FICON Converter” (FCV)

- An ESCON to FICON bridge card was needed to support existing ESCON devices.
- The bridge card was a feature of the 9032-5 ESCON Directors.
- Bridged FICON gains some of the benefits of FICON:
 - 8 simultaneous transactions.
 - Increased I/O (3,200 I/O/sec).
- The bandwidth is about half that of Native FICON and is limited to 1 Gbps links.
- Bridge cards were used as a first step toward a FICON native infrastructure.

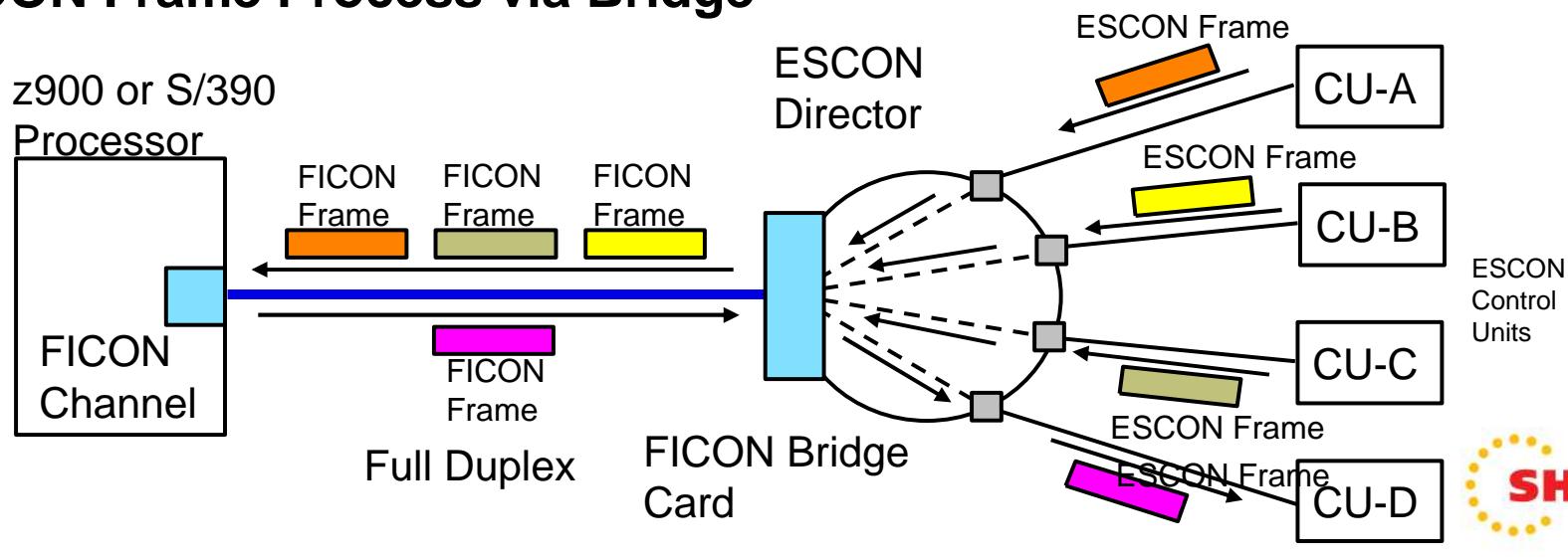


- The bridge card has 8-ESCON connections.
- The 1-external FICON port and 8-internal ESCON ports.
- Replacing an ESCON card with a FICON Bridge card swapped 8-ESCON connections for 1-FICON connection which attached to 8-ESCON devices.

ESCON vs FICON Bridge Frame Processing

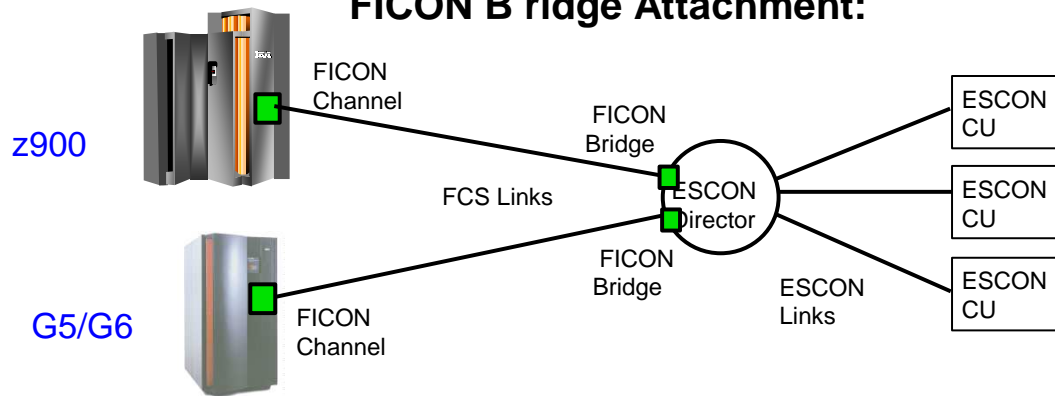


FICON Frame Process via Bridge



FICON Operating Modes

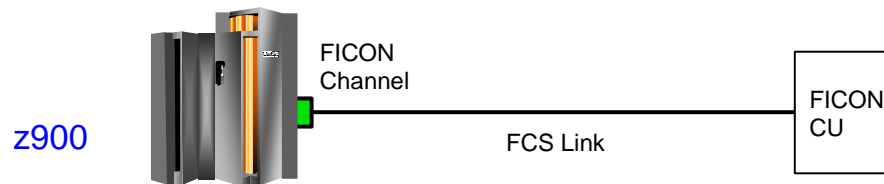
FICON B ridge Attachment:



- ★ *Exploit FICON Channel with Existing ESCON Control Units*

Type=FCV

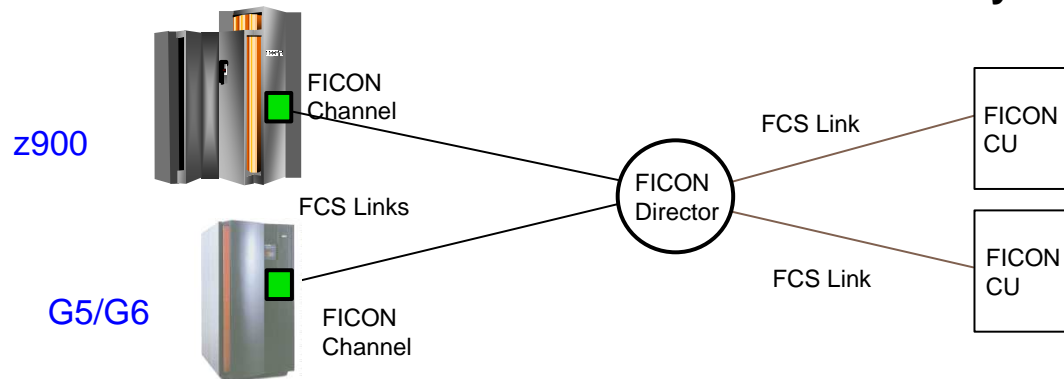
Native FICON Direct Attachment:



- ★ *Native FICON Control Units*

Type=FC

Native FICON Switched Connectivity:



- ★ *Full Dynamic Switching of FICON Control Units*

Type=FC

ESCON & FICON by the Numbers

A comparison

Capabilities	ESCON	FICON Bridge	FICON	FICON Express 4
I/O operations at a time	Only one	Any eight	32 open exchanges	64 open exchanges
Logically daisy-chained Control Units to a single channel	Any one I/O, take turns	Any 8 I/Os concurrently	Any number of I/Os concurrently	Any number of I/Os concurrently
Average I/Os per channel (4k blocks)	2,000-2,500	2,500	6,000	13,000
Unit addresses per channel	1,024	16,384+	16,384 +	16,384 +
Unit addresses per control unit	1,024	1,024	16,384 +	16,384 +
Bandwidth degradation	Beyond 9 km	Beyond 100 km	Beyond 100 km	Beyond 100 km

Today

Current channel technology



A Fibre Channel Standard

FICON



FICON – Security, Resiliency, Performance

- Security & Resiliency
 - Frame Delivery
 - High Integrity Fabrics
 - FC fabrics will validate the WWPN of a re-established ISL before allowing any data to flow on it
 - FICON **requires** high integrity fabrics
 - In Order Delivery
 - Sequence numbers, IU numbers
 - CCW numbers – enables pipelining – used to associate response with the particular command sent
 - Bit Stream Integrity
 - FC2 CRC created by adapter ensures that data as it flows on the link is not corrupted
 - FICON adds separate CRC to insure integrity of data at FC4 layer (end-to-end between hosts)
 - Carried over from ESCON:
 - Link Incident Reporting (LIRR), State Change Notification, RNID

FICON – Security, Resiliency, Performance

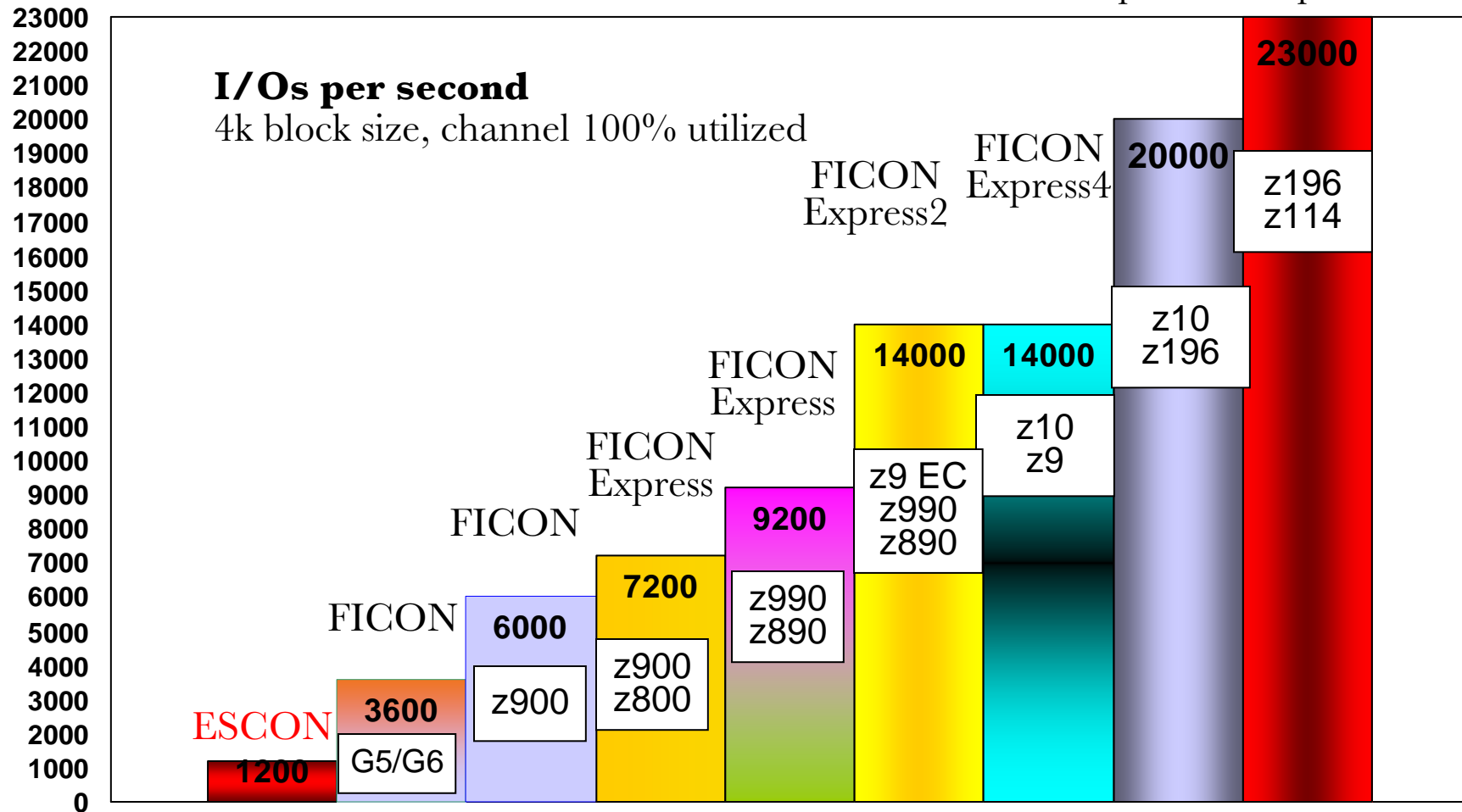
- Performance Management and Droop
 - Fibre Channel standard provides for link flow control through buffer-to-buffer credit scheme
 - FC2 function that controls stream of frames between end points on a link (i.e. between nearest neighbors)
 - Determines the distance two nodes can be apart and still maintain full link frame rate
 - IU Pacing - FICON architectural provision for end-to-end flow control
 - Prevents flooding of target N-Port
 - *With command pipelining needed mechanism to prevent over-running the control unit*
 - ‘Extended distance’ support enabled CU to dynamically change pacing count on each Command Response

FICON – Security, Resiliency, Performance

- Performance Management and Droop
 - I/O Priority
 - Separate priority mechanisms in I/O Subsystem, Channel and Control Unit
 - Modified Indirect Data Addressing Words (MIDAWs)
 - A method of gathering and scattering data into & from non-contiguous System z storage locations during an I/O operation
 - Removed IDAW restriction that appended data must be on 2K storage boundary
 - Improved performance of certain applications (e.g. DB2 sequential workloads) that process small records with Extended Format data sets
 - Measurements – granular to service class
 - Allows algorithms for WLM based I/O priority, DCM & intelligent data placement
 - Dynamic CHPID Management
 - allows adding/removing bandwidth to a control unit as workload needs dictate

FICON performance – Start I/Os

Historical Actuals





Rev'ing the Engine

ZHPF

High Performance FICON

The host communicates directly with the control unit

- The channel is acting as a conduit
- No individual commands or state tracking
- The CCW program is sent to the control unit in one descriptor
- Uses the Fibre Channel FCP link protocol. The channel provides both the new and old protocols
- HPF provides increased performance for small block transfers & enables greater bandwidth exploitation
- Complex channel programs that are not easily converted to the new protocol still execute with the existing FICON protocol
- Devices accessible using both old and new protocols

Link Protocols for 4K Read

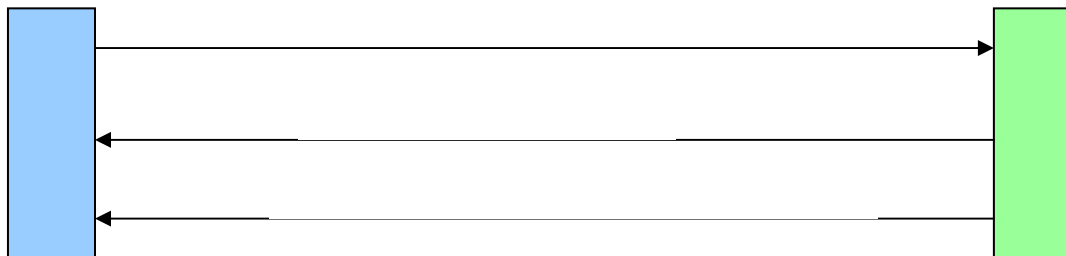
FICON:

Two Exchanges opened and closed
Six Sequences opened and closed



zHPF:

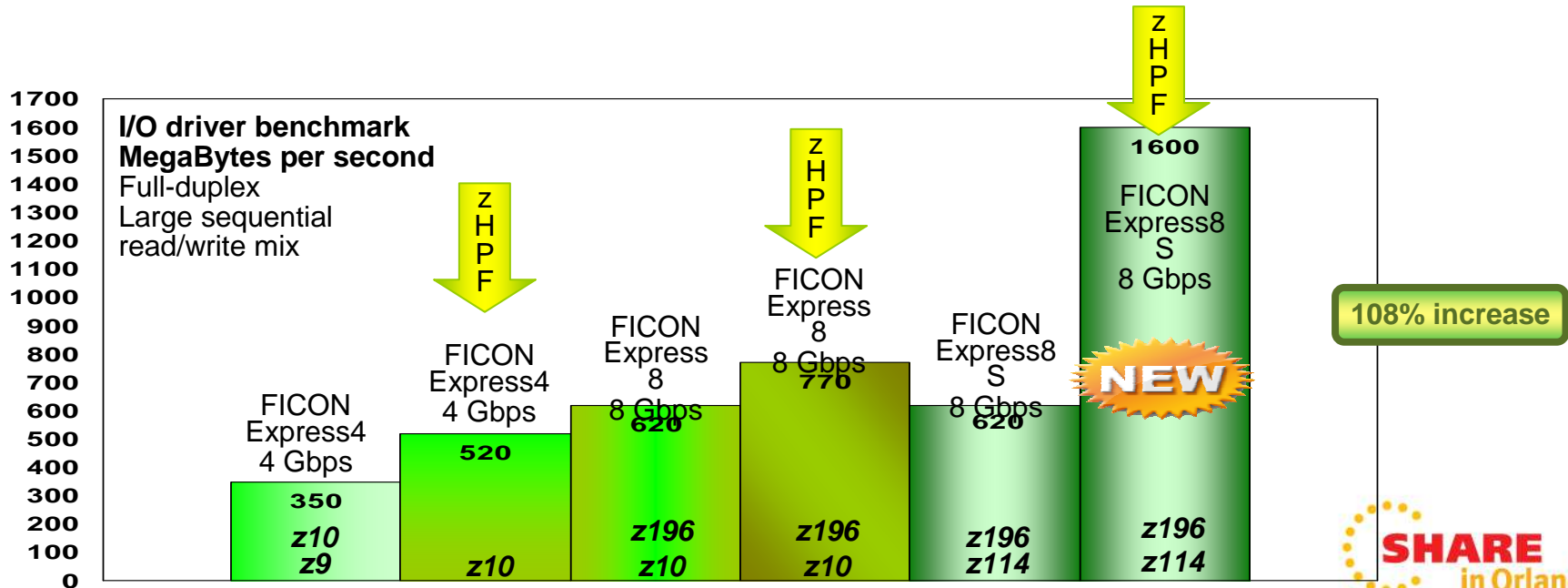
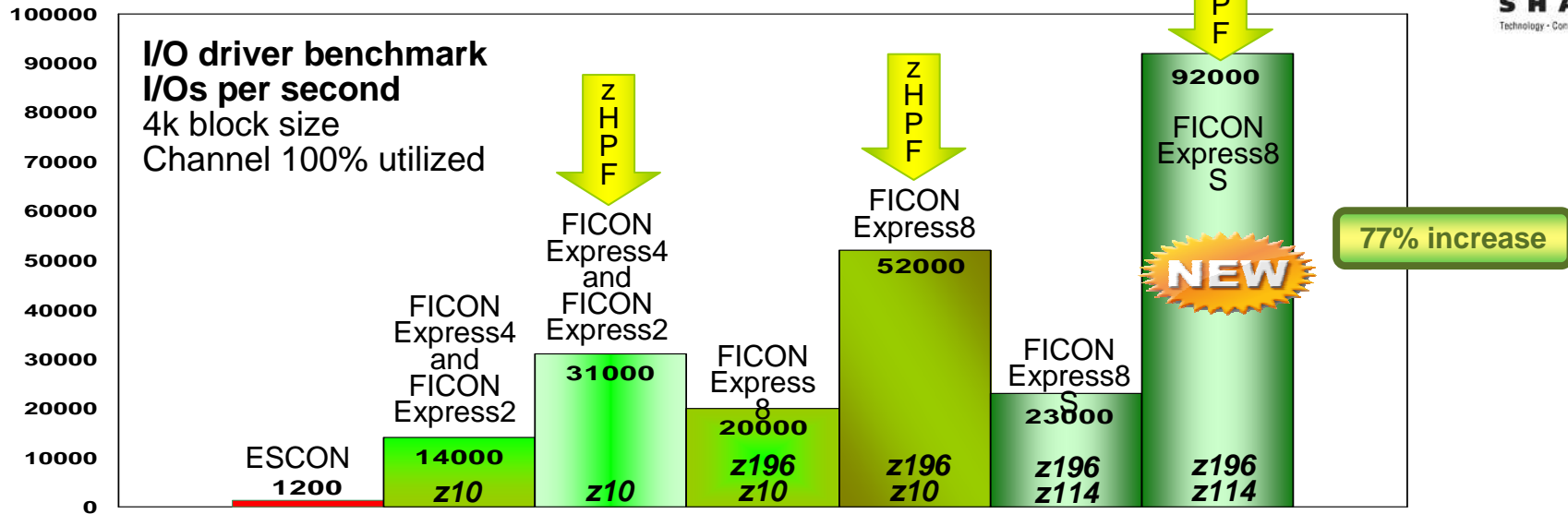
One exchange opened and closed
Three Sequences opened and closed



CHANNEL

CONTROL UNIT

FICON performance on System z



Looking Ahead

But no reading of the tea leaves



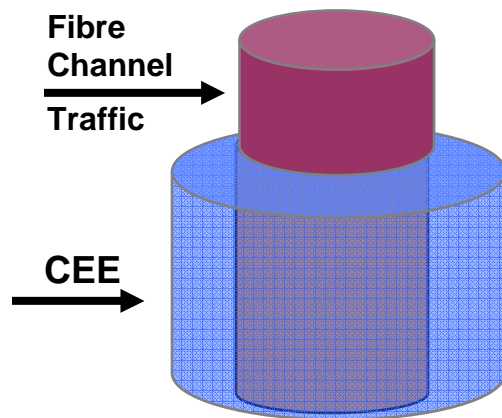
**How fast can this go? And
over what roads?**

What's Rumbling About in the Industry?

- 16Gb Fibre Channel
 - Physical interface approved as an ANSI INCITS T11.2 standard in Sept. 2010
 - Auto-negotiated backward compatibility to 4Gb
 - Aimed at improved price-performance (as market matures)
 - Twice the bandwidth of 8Gb
 - Leapfrogs fiber channel storage SAN fabric over 10GbE
- 32Gb Fibre Channel
 - Next logical step for Fibre Channel after 16Gb
 - Work on-going in ANSI with proposals put forward
 - Future will depend on market demand
- Fibre Channel over Ethernet (FCoE)
 - See next charts for description
 - Future will depend on market demand

What is FCoE?

FCoE is a technology that straddles two worlds:
Fibre Channel and Ethernet

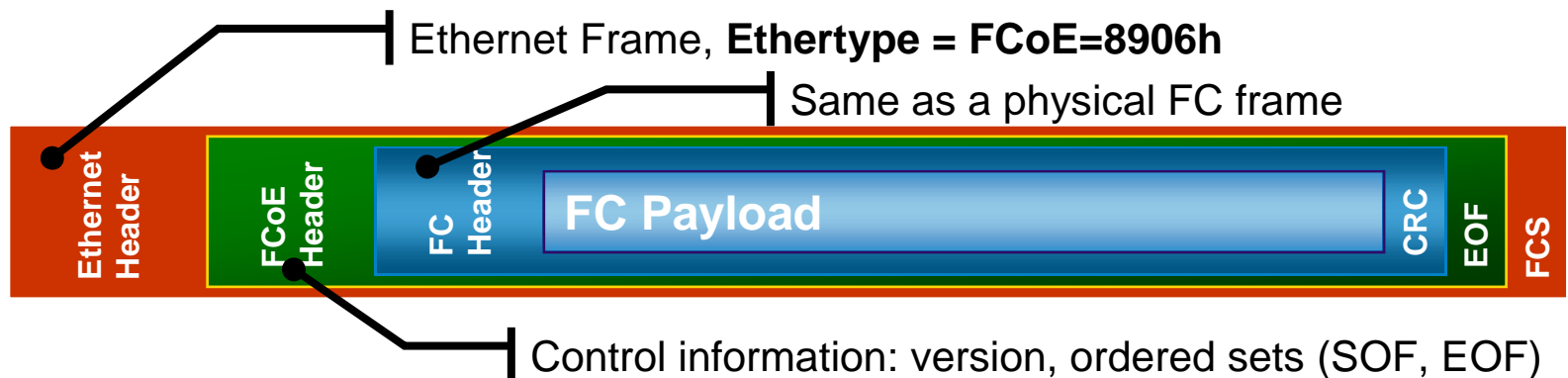
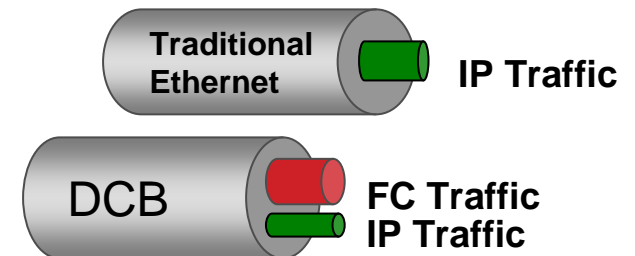


- **FC View:** FC gets a new transport in the form of lossless Ethernet (CEE)
- **Ethernet view:** A new upper-layer protocol, or storage application, that runs over a new lossless Ethernet

What does FCoE look like?

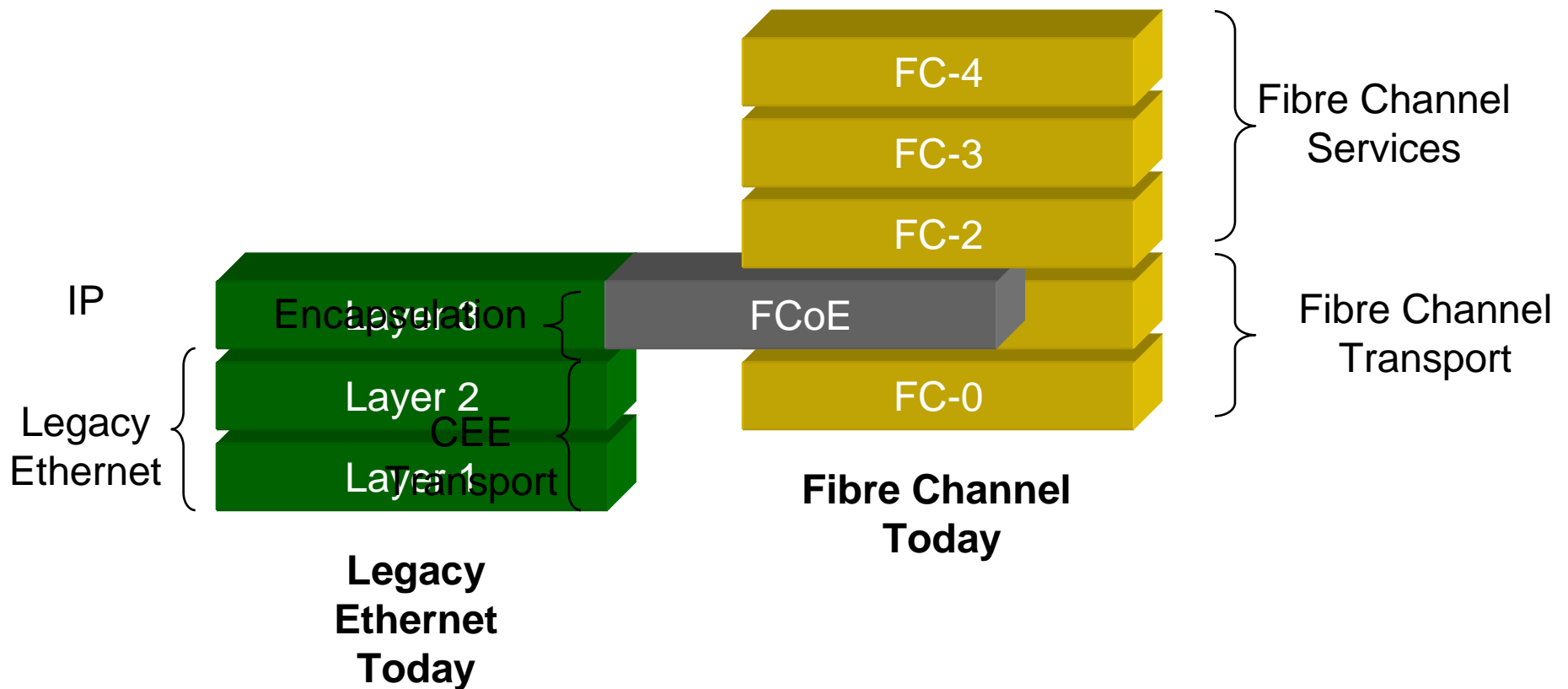
It's an encapsulation protocol

Encapsulation protocol for transporting FC over Ethernet (Lossless Ethernet: DCB)



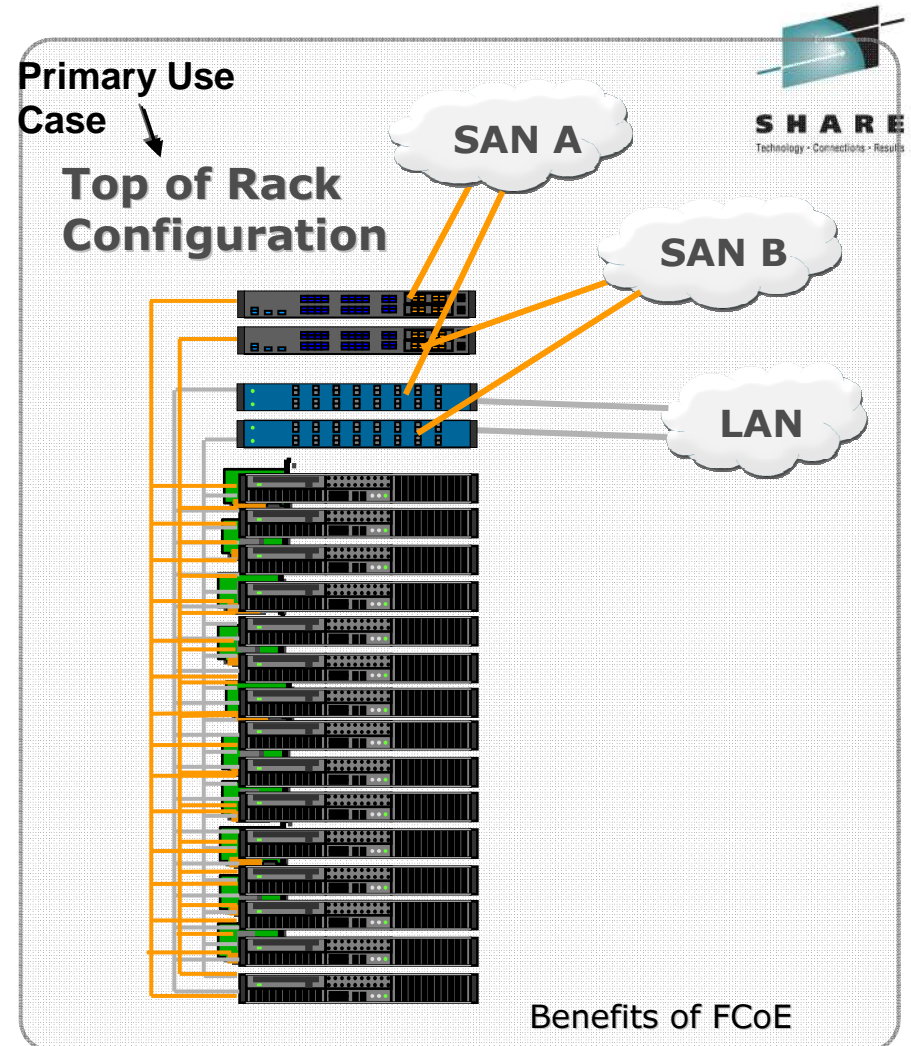
- FC frame remains intact: FC does not change
- Ethernet needs a larger frame: Larger than 1.5 KB
- Ethernet must become lossless to carry storage data with integrity

Where does FCoE fit in the stack?



FCoE & CEE Benefits

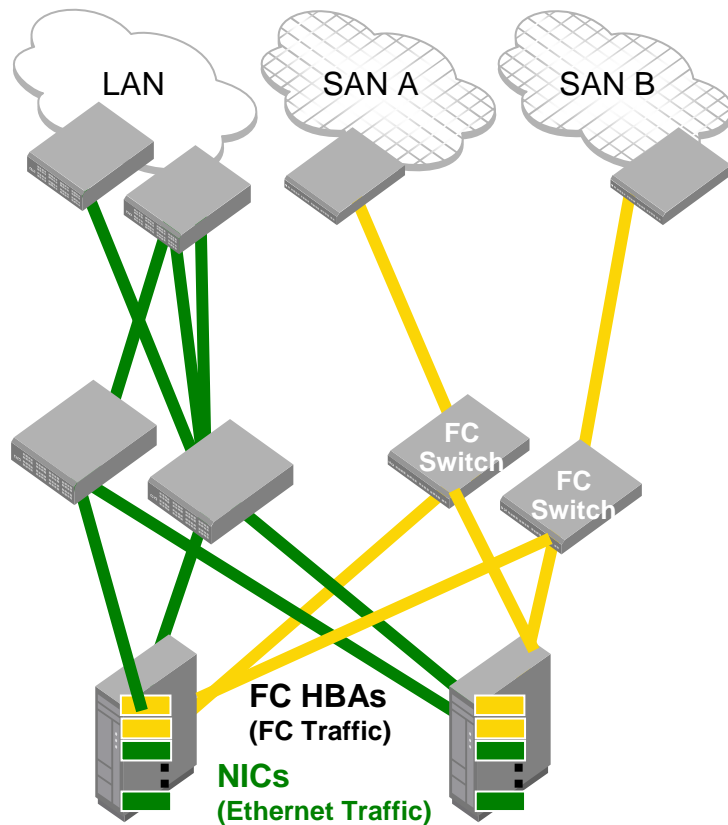
- Reduce number of server ports
- Reduce number of switches ports
- Reduce cabling
- Reduce power consumption
- Increase speed of links
- Increase utilization of links



Consolidating I/O Interfaces Lowers CapEx and OpEx

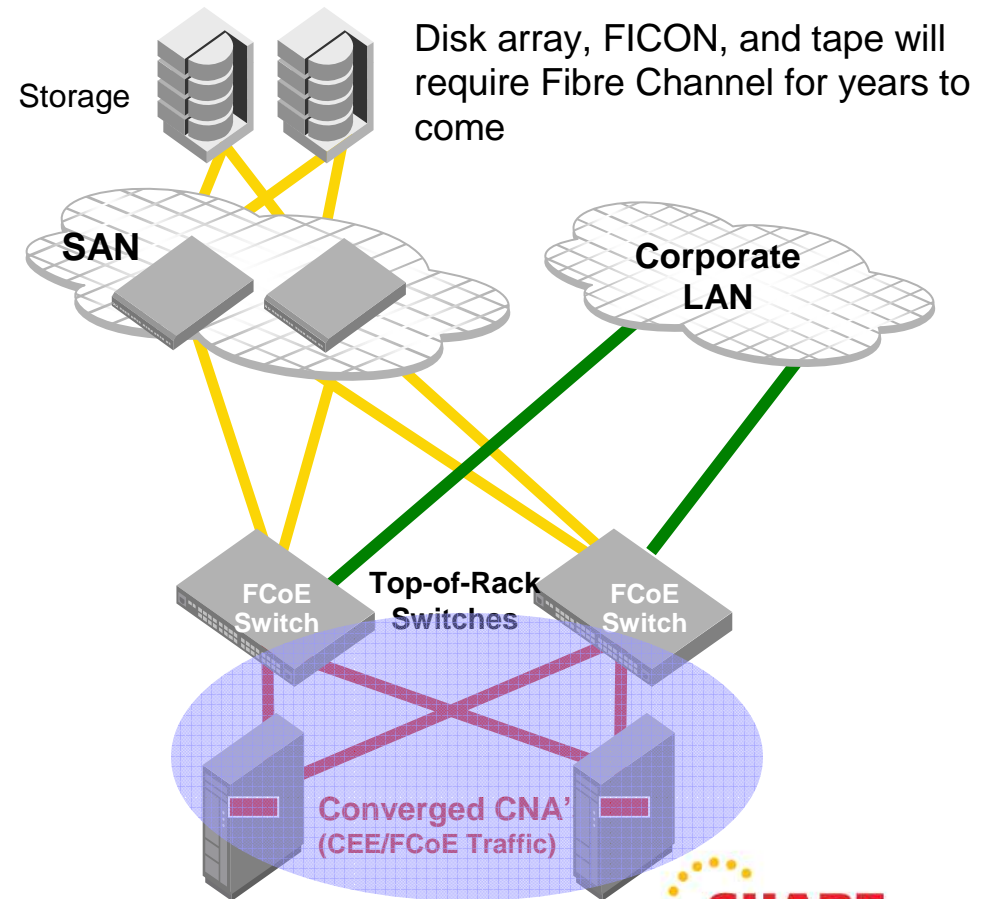
Primary FCoE Use Case

Before Unified I/O

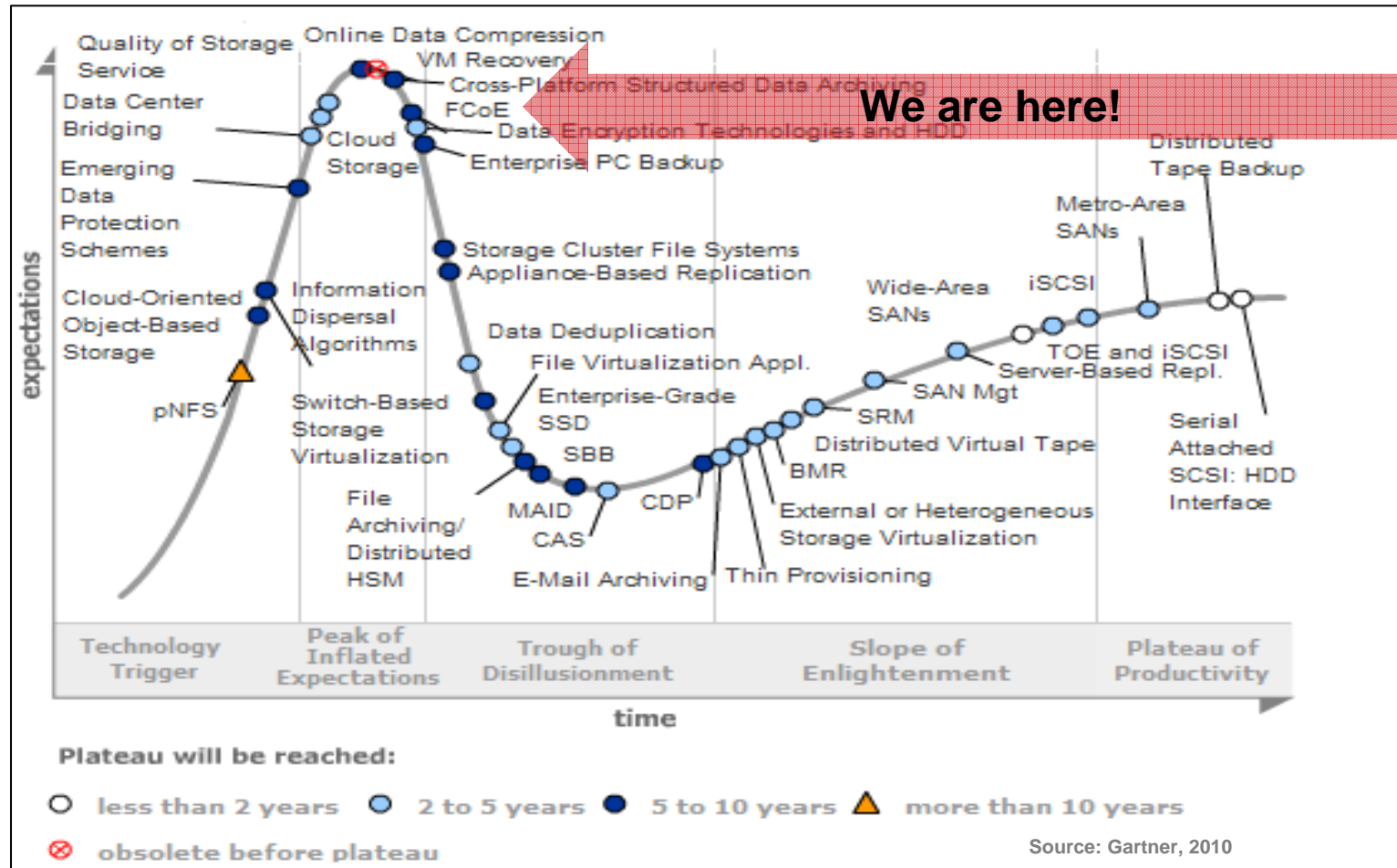


— FC — Ethernet — FCoE and CEE

After Unified I/O



Storage Technology Hype Cycle Curve



As the Evolution Continues....

- **System z continues focus on key Enterprise storage fabric attributes:**
 - **Security**
 - **Resiliency**
 - **Performance**
- through:**
- **Continued innovation in channel and control unit architectures**
 - **Working within ANSI Standards Committees to drive requirements into the standards of evolving technologies (such as FCoE)**
 - **Working with eco-system vendor partners to provide differentiating functions**

SHARE

Orlando

August 2011

THANK YOU!

Fun Links

- <http://www.vikingwaters.com/htmlpages/MFHistory.htm>
- http://en.wikipedia.org/wiki/Mainframe_computer
- http://en.wikipedia.org/wiki/Channel_program#Channel_Program
- http://en.wikipedia.org/wiki/IBM_System/360#Channels
- http://en.wikipedia.org/wiki/Channel_I/O
- <http://en.wikipedia.org/wiki/ESCON>
- <http://en.wikipedia.org/wiki/FICON>

REFERENCES

Speaker Biography

- Patty Driever
 - IBM
 - System z I/O and Networking Technologist
- Contact Information
 - pgd@us.ibm.com

Speaker Biography

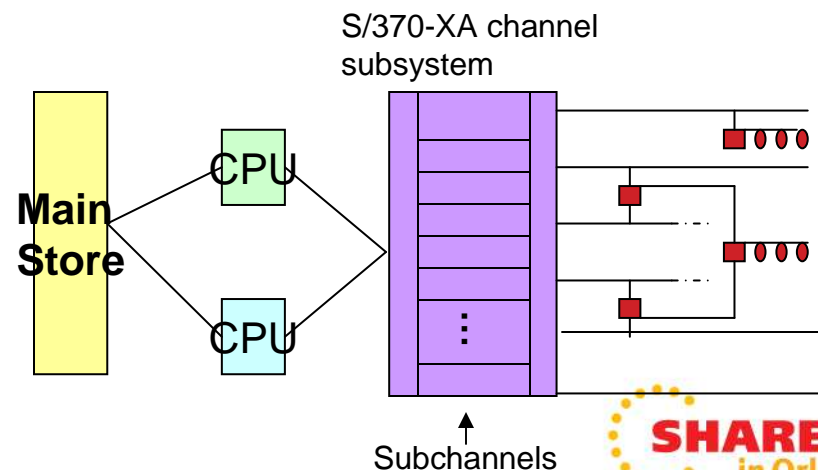
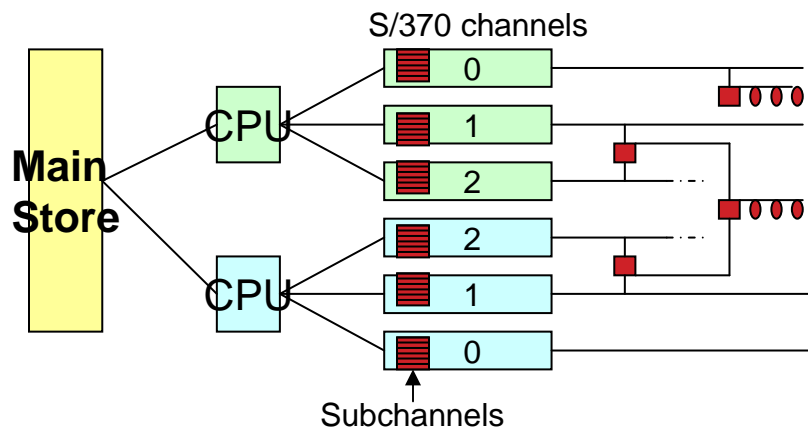
- Howard L. Johnson
 - BROCADE
 - Technology Architect, FICON
 - 27 years technical development and management
- Contact Information
 - howard.johnson@brocade.com

Stuff we didn't get to

BONUS INFORMATION

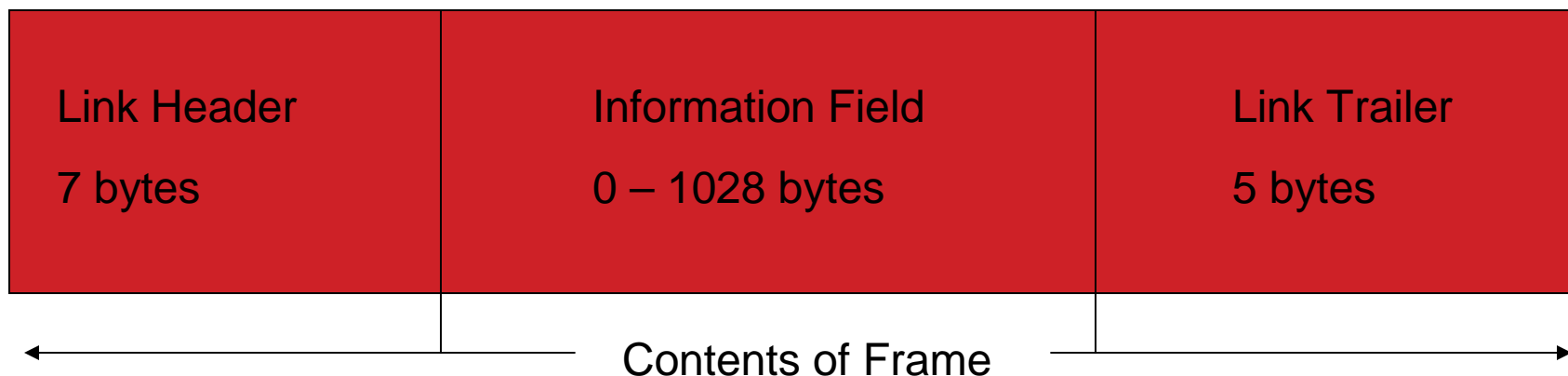
Move to S/370-XA (Extended Architecture)

- During S/370 SIO processing, the CP was hung up the entire time the channel was connected to a device (~100 usec)
 - On a 1MIP machine this was ~100 instructions
- IBM 3033 in 1978 introduced initial architectural concepts of a channel subsystem that allowed CP I/O instruction processing execution to be disconnected from the channel processing
 - Buffering of status in subchannels
 - Queuing of I/O requests in subchannels (Start-I/O-fast queuing)
 - CP stored information in the subchannel control block and went on its way
 - At the end of the operation, the CP received an interrupt indicating result of the operation
 - Suspend-resume facility



ESCON

- Basic frame structure:



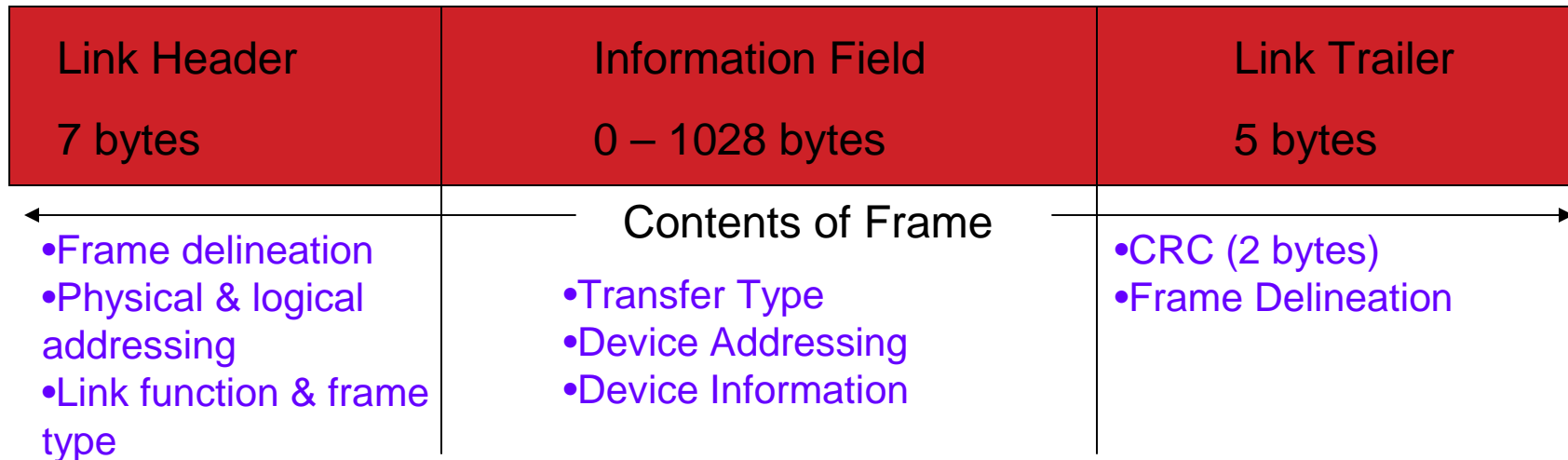
- Frame delineation
- Physical & logical addressing
- Link function & frame type

- Transfer Type
- Device Addressing
- Device Information

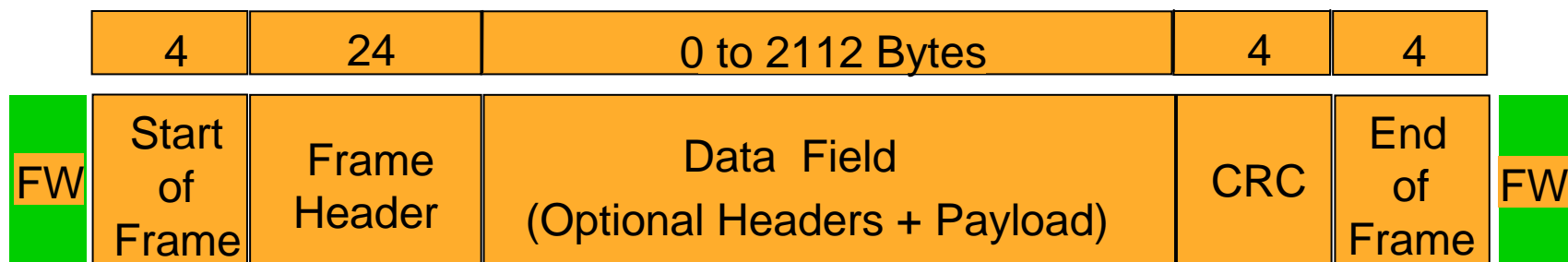
- CRC (2 bytes)
- Frame Delineation

FICON – Basic Building Block

- ESCON basic frame structure:



- FICON basic frame structure:



- Four byte **SOF** delimiter
- Twenty-four byte fixed-format **Frame Header**
- Variable-size **Data Field**:
 - 0 to a maximum of 2112 bytes (2048+64)
 - May contain Optional Headers as well as payload
- Four byte **CRC**
- Four byte **EOF** delimiter

ESCON to FICON Comparison

- ESCON
 - A somewhat proprietary protocol created by IBM
 - *Circuit switched*, only a single switch allowed
 - ANSI Standard
 - Information Technology--Single-Byte Command Code Sets CONnection (SBCON) Architecture (formerly ANSI X3.296-1997)
- FICON
 - A standard, Fibre Channel FC4 level protocol
 - *Packet switched*, multiple switches are allowed
 - ANSI Standards
 - Information Technology - Fibre Channel - Single-Byte-2
 - *(FC-SB-2) (formerly ANSI NCITS 349-2001)*
 - Fibre Channel - Single Byte Command Set – 3
 - *(FC-SB-3) (currently T11 Project 1569-D)*
- Both feature I/O operations that are address-centric, definition oriented, and host assigned.

FICON Operating Modes

- There are two FICON operating modes: FCV and FC
 - For System zSeries and 9672 G5 and G6 servers, there were two modes supported:
 - FICON Bridge Mode – referred to as FCV which is a FICON channel mode designed to enable access to ESCON interfaces using the FICON Bridge Adapter in the ESCON Director (9032-5).
 - FICON Native Mode – referred to as FC which enables access to the FICON channel mode from native FICON control units in a point-to-point mode through a FICON Director. (No cascaded FICON.)
 - For zSeries and System z, support was added for:
 - FICON Native (FC) but not FICON Bridge.
 - Cascaded FICON.
 - Fibre Channel Protocol (FCP) which provided support for open systems via industry standard SCSI devices.

What does FCoE do?

Lower CapEx and OpEx



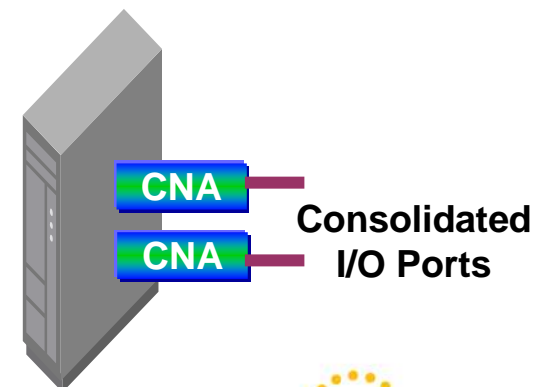
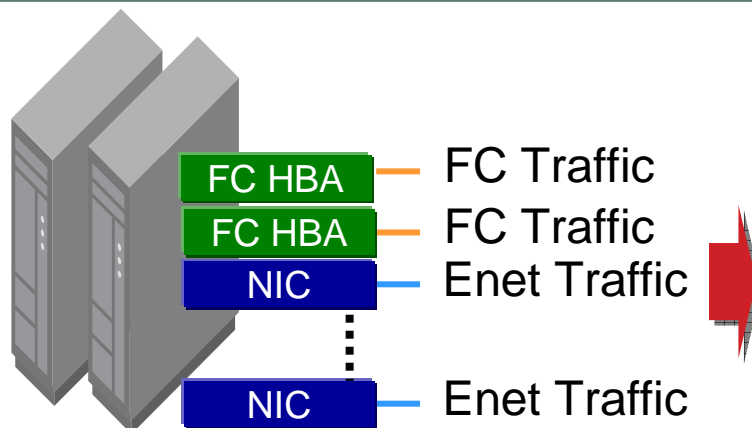
Consolidate and simplify server connectivity to LAN and SAN



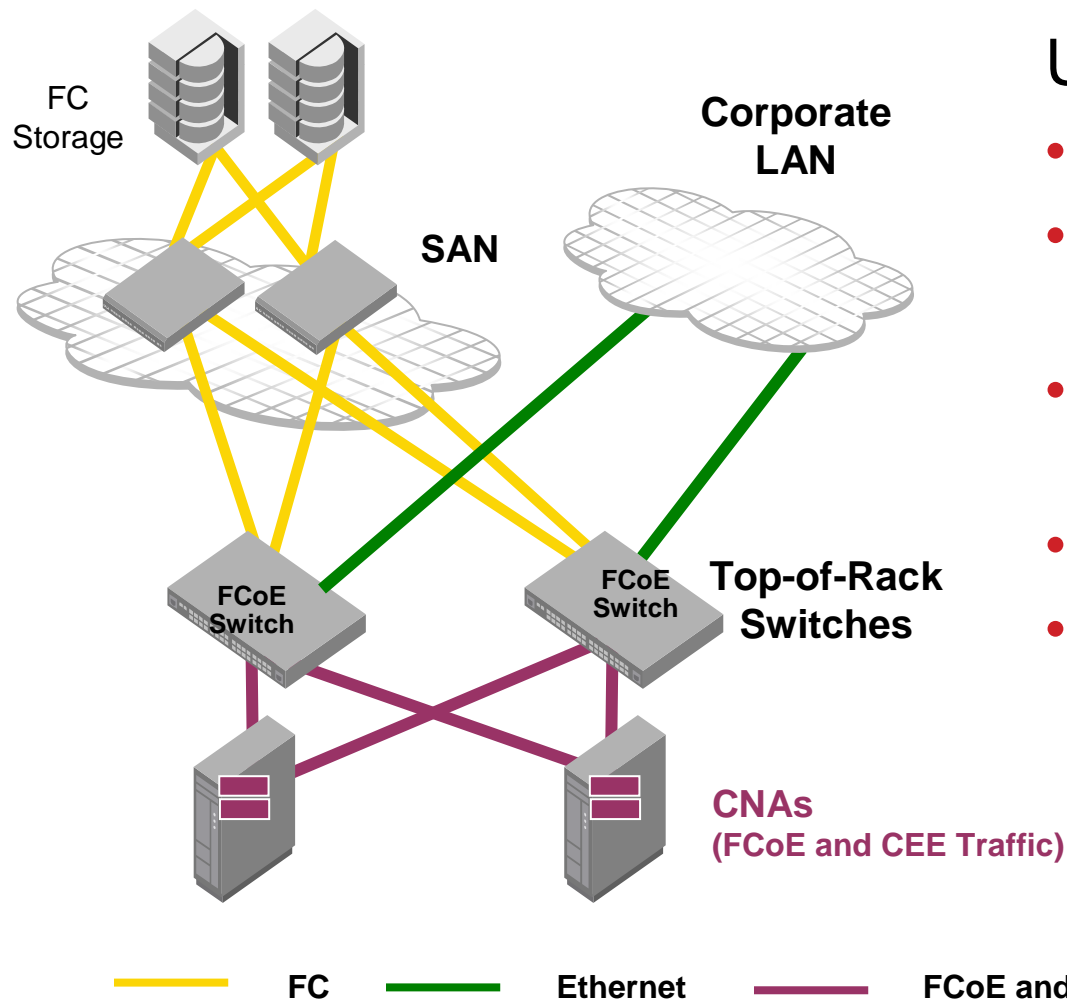
Reduce server adapters, cabling, switch ports, and power usage



Lower CapEx and OpEx in Enterprise Data Centers



Unified I/O Usage Case

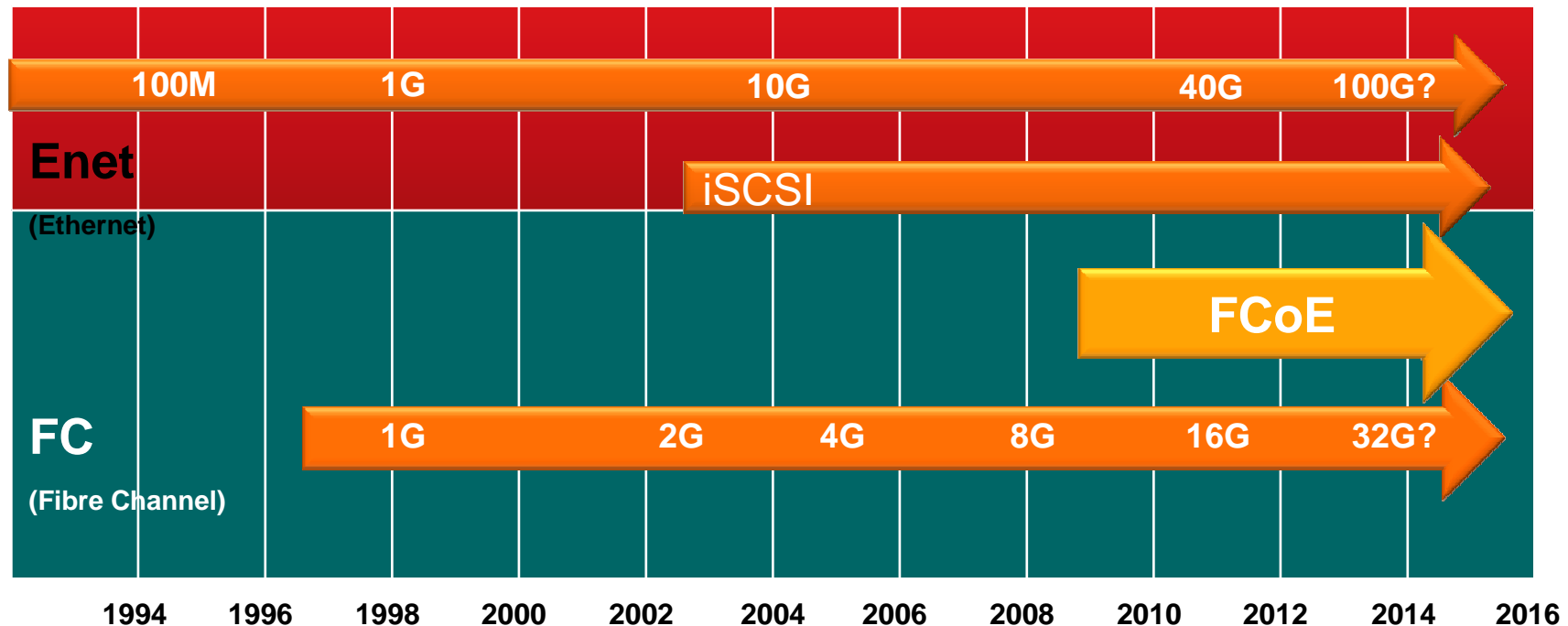


Unified I/O:

- One fabric infrastructure
- FC and CEE L2/L3 multipathing end to end
- Fewer components to deploy
- Lower TCO (future)
- Disk array, FICON, and tape will require Fibre Channel for years to come

Ethernet and FC Roadmaps

Parallel Evolution & Potential for Convergence



- FC and Ethernet evolved in parallel paths with FC dominating storage SANs and Ethernet supporting IP networking
- Lossless Ethernet & FCoE open the door for server I/O consolidation

FCoE & DCB Industry Standards Activities

Making 10GbE Lossless



- T11 FCoE & FIP standards
 - FCIA approved standards
 - INCITS approved and expected to publish soon
- IEEE draft DCB components
 - 802.1Qbb priority-based flow control
 - 802.1Qaz enhanced transmission selection
 - DCBX capability exchange protocol
 - 802.1Qau Congestion Notification
 - Completion & approvals expected H2 of 2010
- IETF status
 - Transparent Interconnection of Lots of Links (TRILL)
 - Expected completion H2 2010



Mainframe Channel Cards

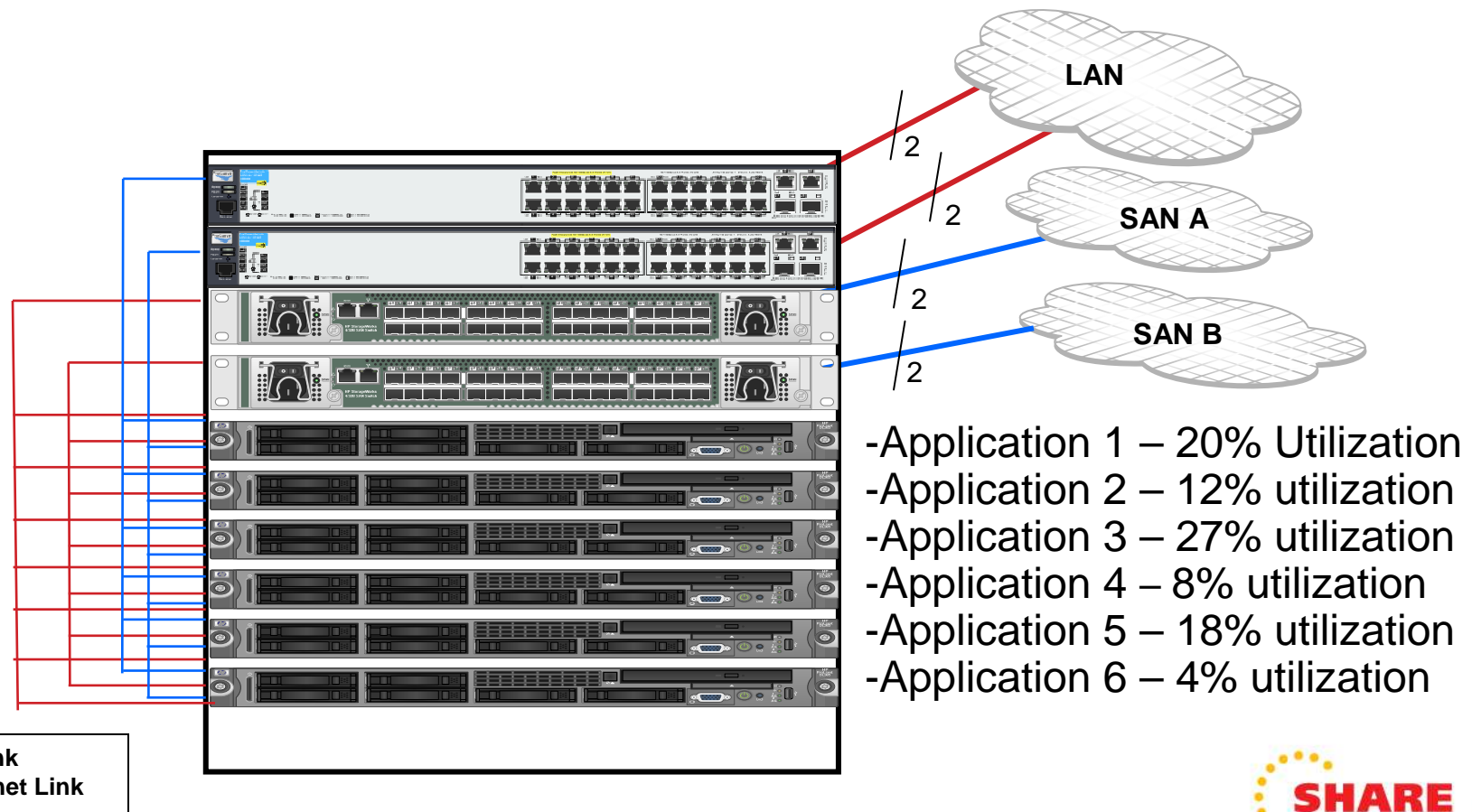
- ESCON Channel Card
 - Architecture of 20MBps; Sustainable Data Rate of about 12MBps
- 74 MBps FICON
 - First FICON Channel Card did not meet customer expectations
 - SC connector with 2 Ports
- 100 MBps FICON Express (certified for 120MBps FD)
 - The replacement for the poor performing original FICON card
 - LC connector with 2 Ports
- 200 MBps FICON Express (certified for 270MBps FD)
 - LC connector with 2 ports
- 200 MBps FICON Express2 (certified for 270MBps FD)
 - LC connector with 4 ports
 - I/O Performance Improvements
- 400 MBps FICON Express4 (certified for 350MBps FD)
 - LC connector with 4 ports
 - Probable I/O Performance Improvements

ESCON and FICON Standards

- ESCON
 - *Circuit switched*, only a single switch allowed
 - ANSI Standard
 - Information Technology--Single-Byte Command Code Sets CONnection (SBCON) Architecture (formerly ANSI X3.296-1997)
- FICON
 - *Packet switched*, multiple switches are allowed
 - ANSI Standards
 - Information Technology - Fibre Channel - Single-Byte-2 (FC-SB-2) (formerly ANSI NCITS 349-2001)
 - Fibre Channel - Single Byte Command Set - 3 (FC-SB-3) (currently T11 Project 1569-D)

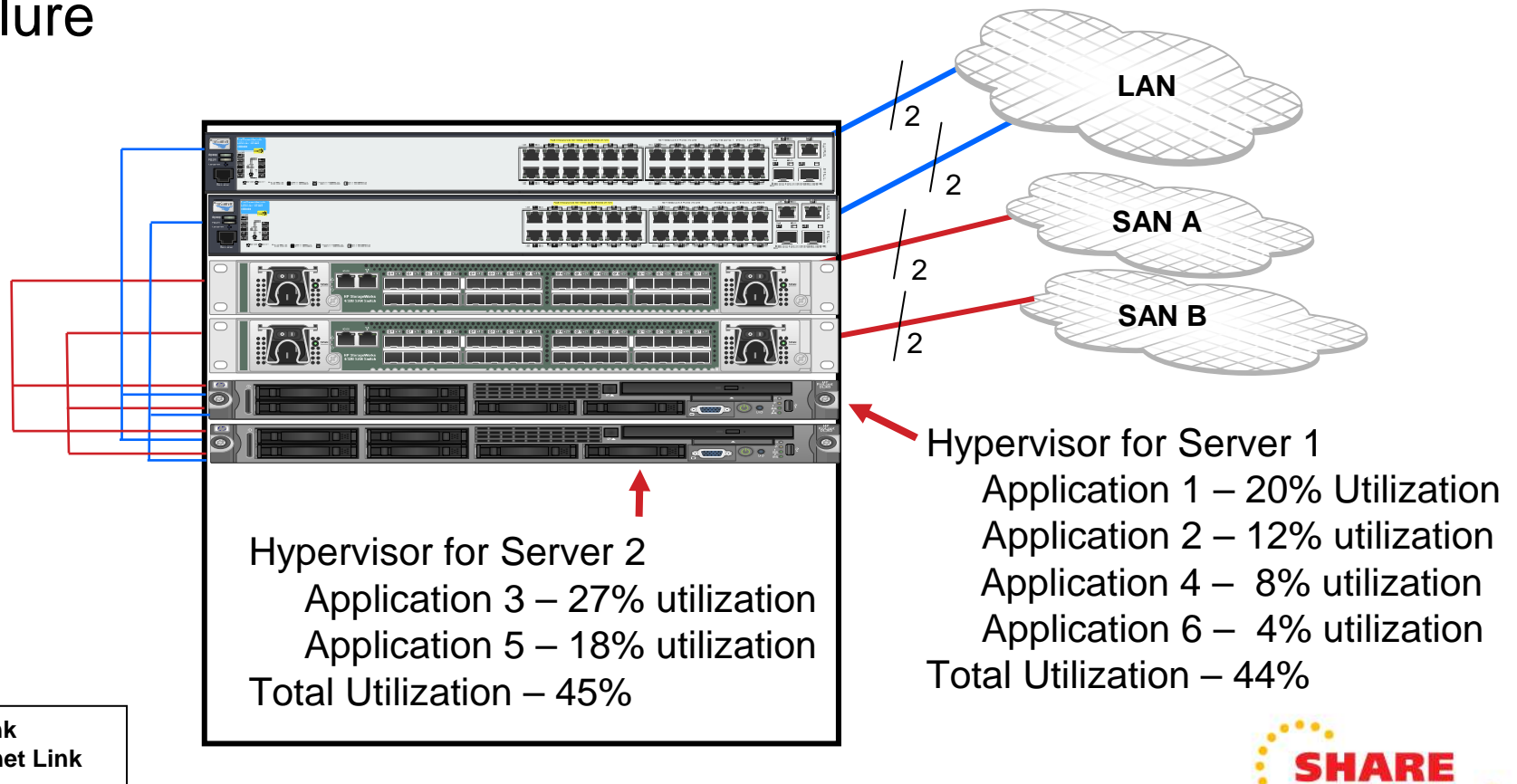
Virtualization of Servers (1 of 2)

- Applications usually underutilize network and server



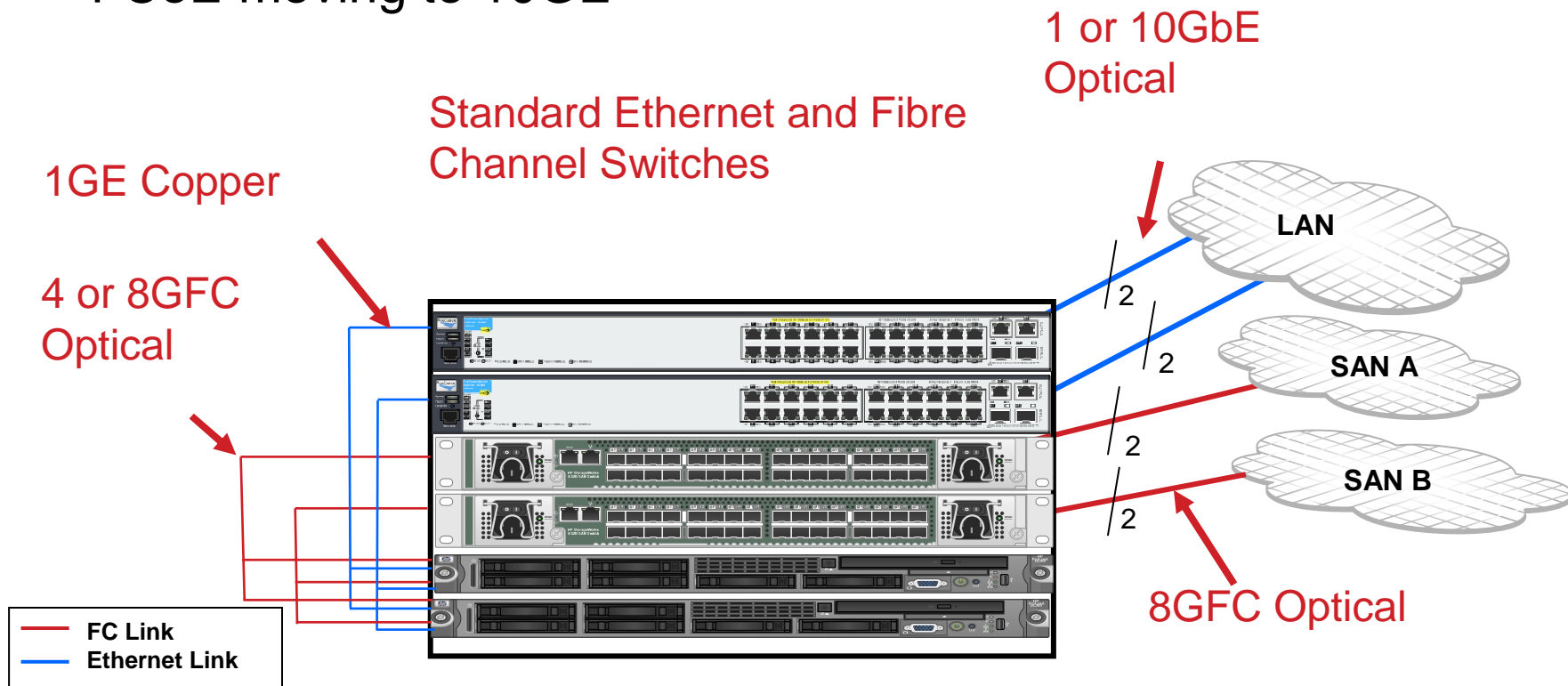
Virtualization of Servers (2 of 2)

- Consolidating servers through application virtualization
- Applications migrate between server 1 and 2 in the event of a failure



Higher Speeds

- Most server links still classic 1 Gigabit Ethernet (1GE)
- Ethernet uplinks moving to 10GE
- Fibre Channel links moving to 8GFC
- FCoE moving to 10GE



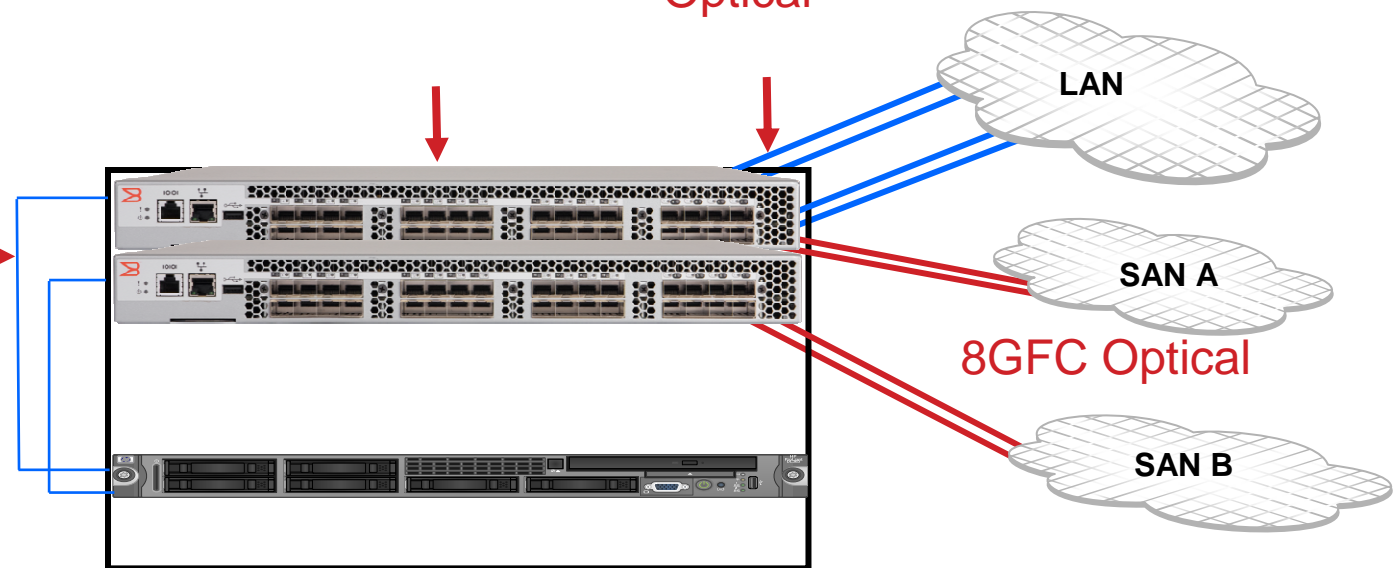
FCoE Moving to 10G CEE at Server

- Reduce the server links to two 10G CEE links
- Uplinks are 10GE links with Link Aggregation
- Fewer ports and switches save power and money

10G CEE with FCoE
SFP interface supports Twinax to 1, 3 or 5 meters or Optical to 300 meters

FCoE Switches

nx10GE Optical



**THIS SLIDE INTENTIONALLY
LEFT BLANK**