

# Lecture 9

## Advanced Optical Technologies in DCNs

Dr Shuangyi Yan & Dr George T. Kanellos

# Optical Technologies in DCNs

## Bandwidth

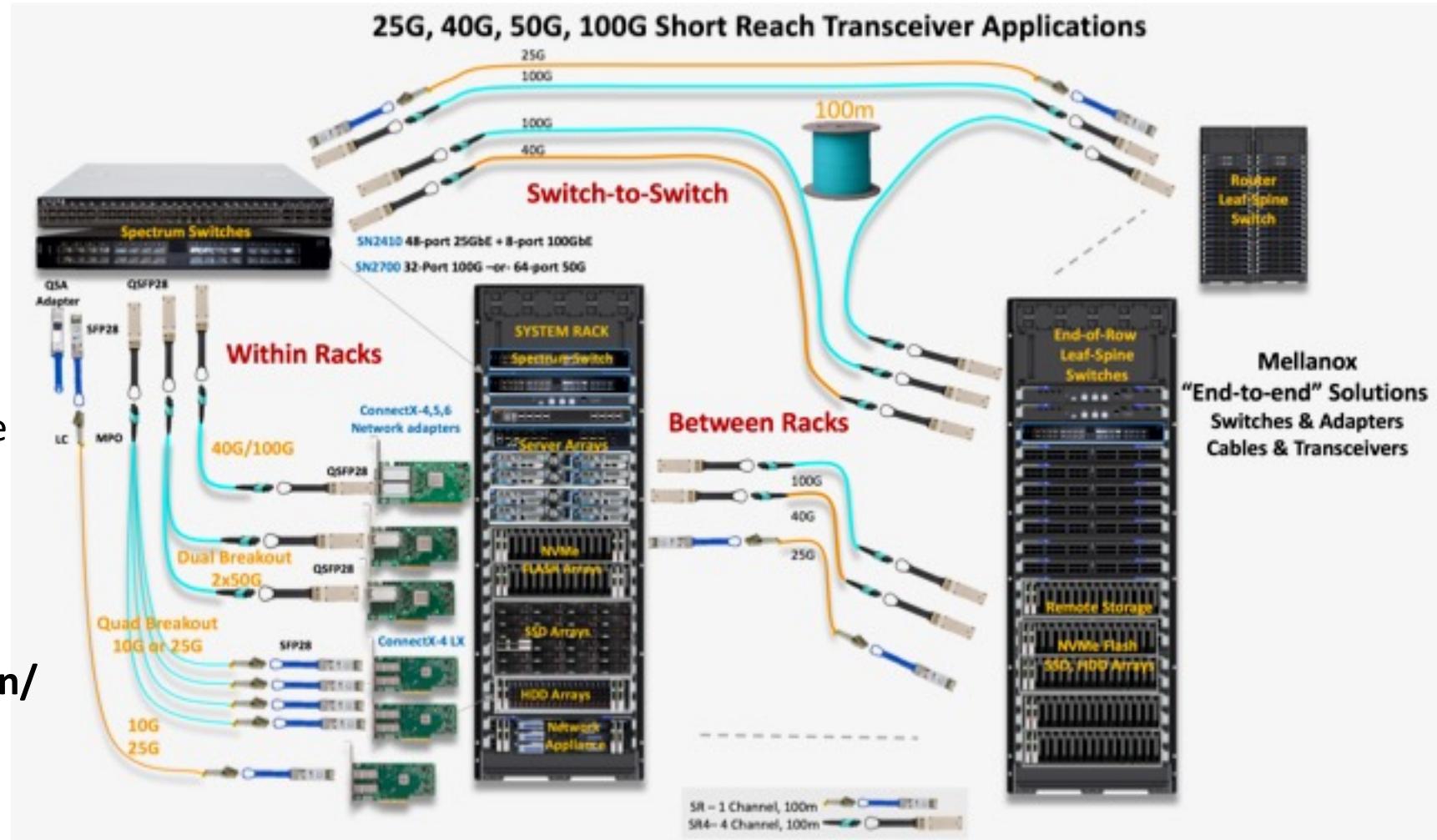
10Gbs  
25Gbs  
40Gbs  
100Gbs

## Reach

Multi-Mode  
Single-Mode  
WDM

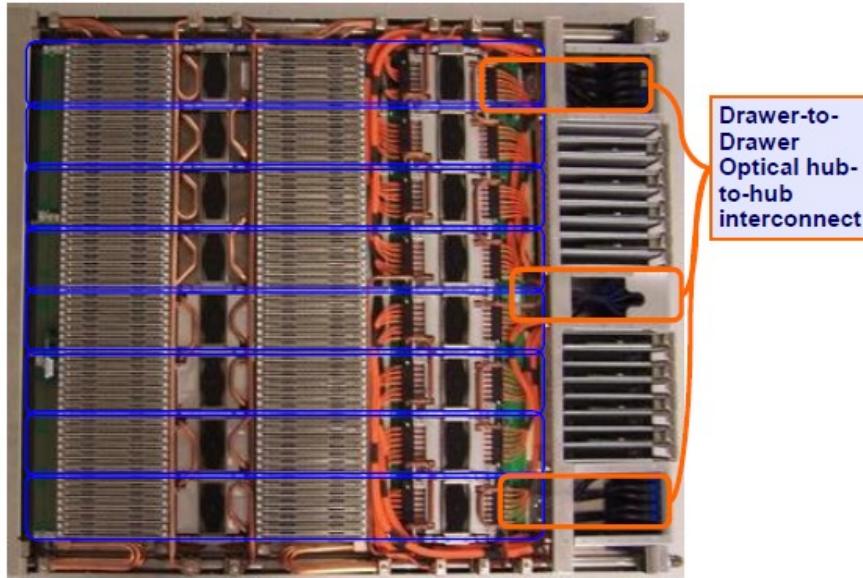
## Power consumption/ Size

SFP+  
QSFP  
CFP



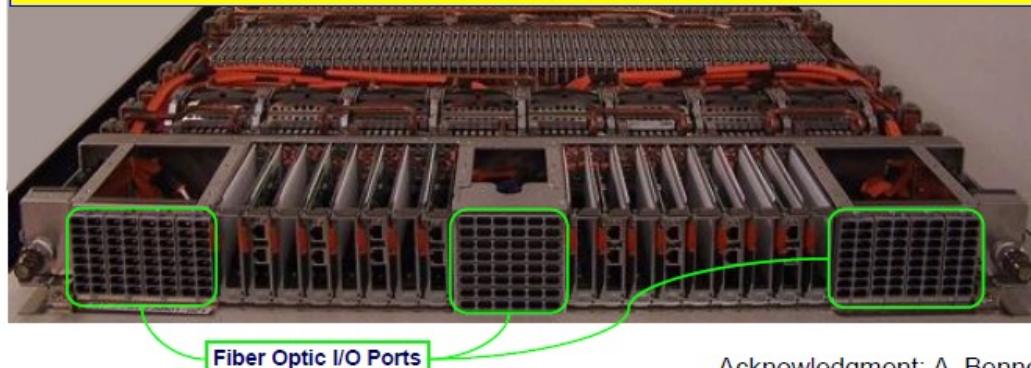
- 
- Transmission medium
  - Transmission technologies
  - Optical switching technologies (Potential technologies)

# Intra-Rack Parallel Optics



256-core Node Drawer

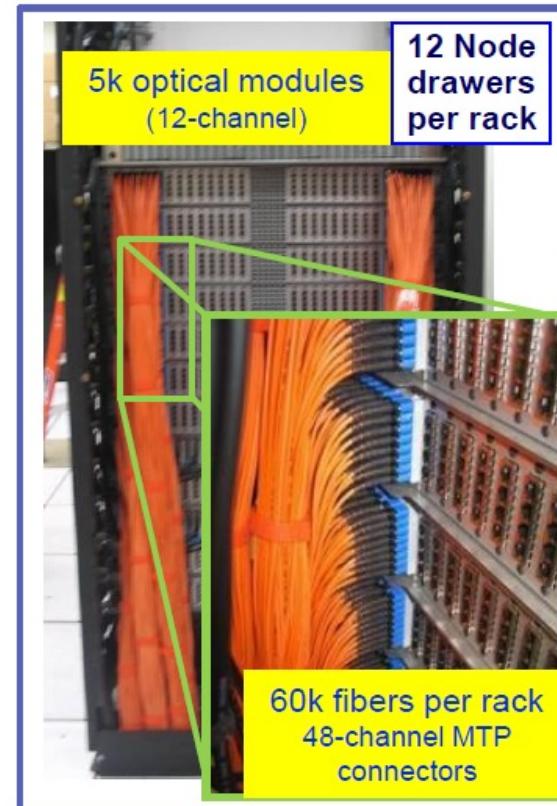
- Optical transceivers tightly integrated, mounted within drawer
- 8 Hub/switch modules (8 x 56 optical modules)



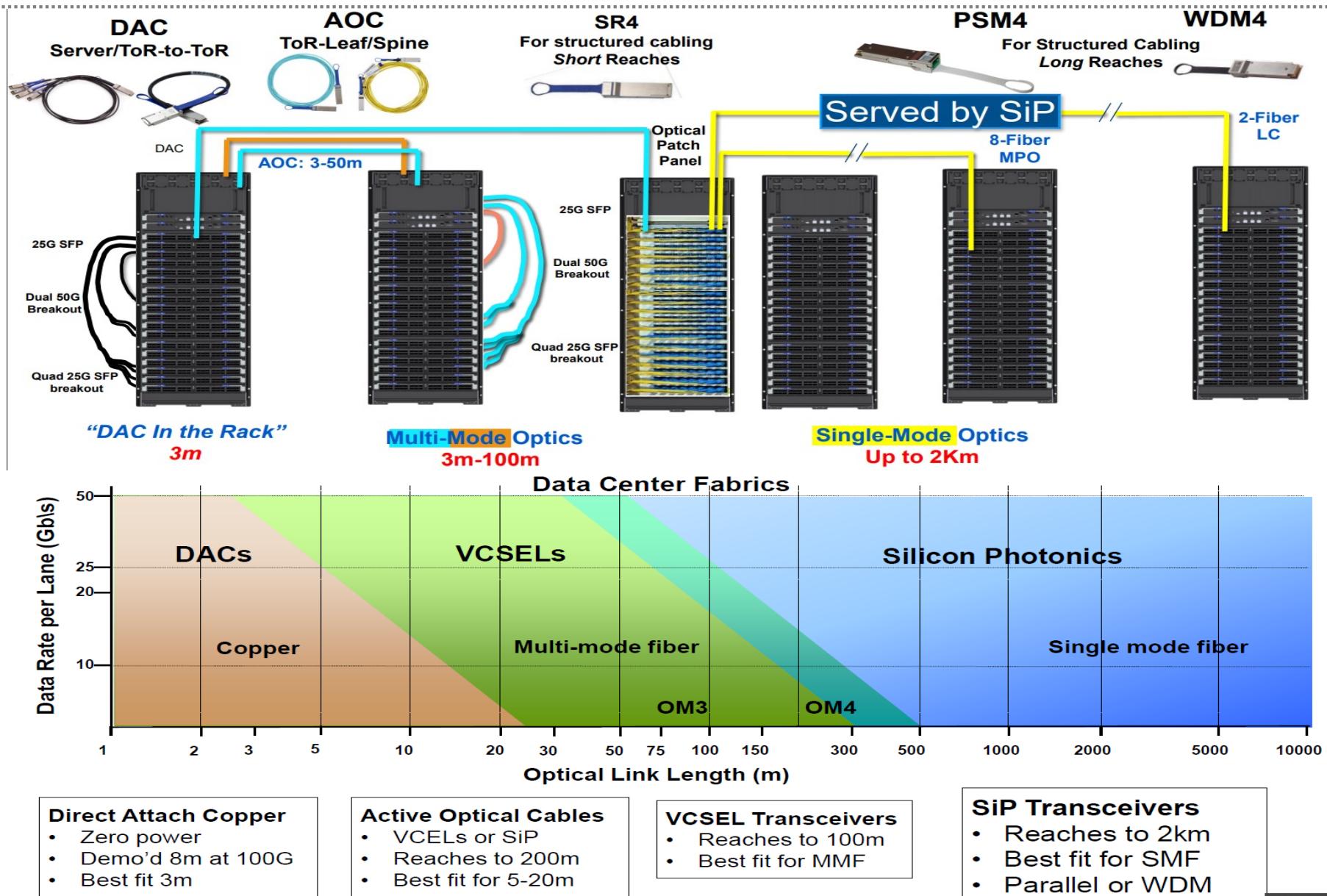
Acknowledgment: A. Benner

## P775 Drawer

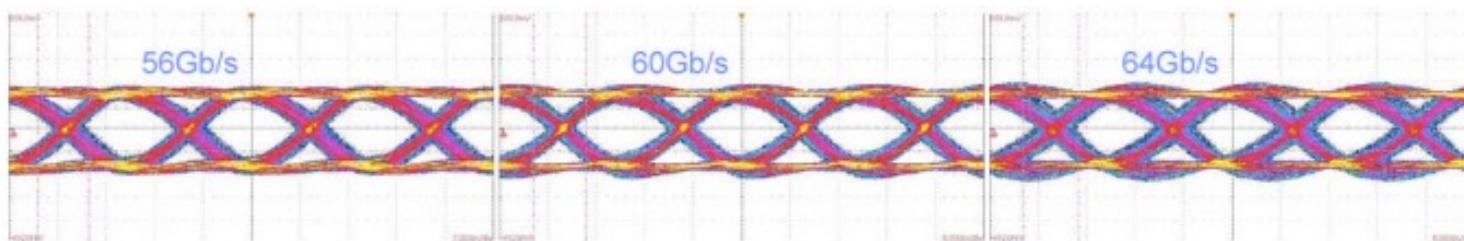
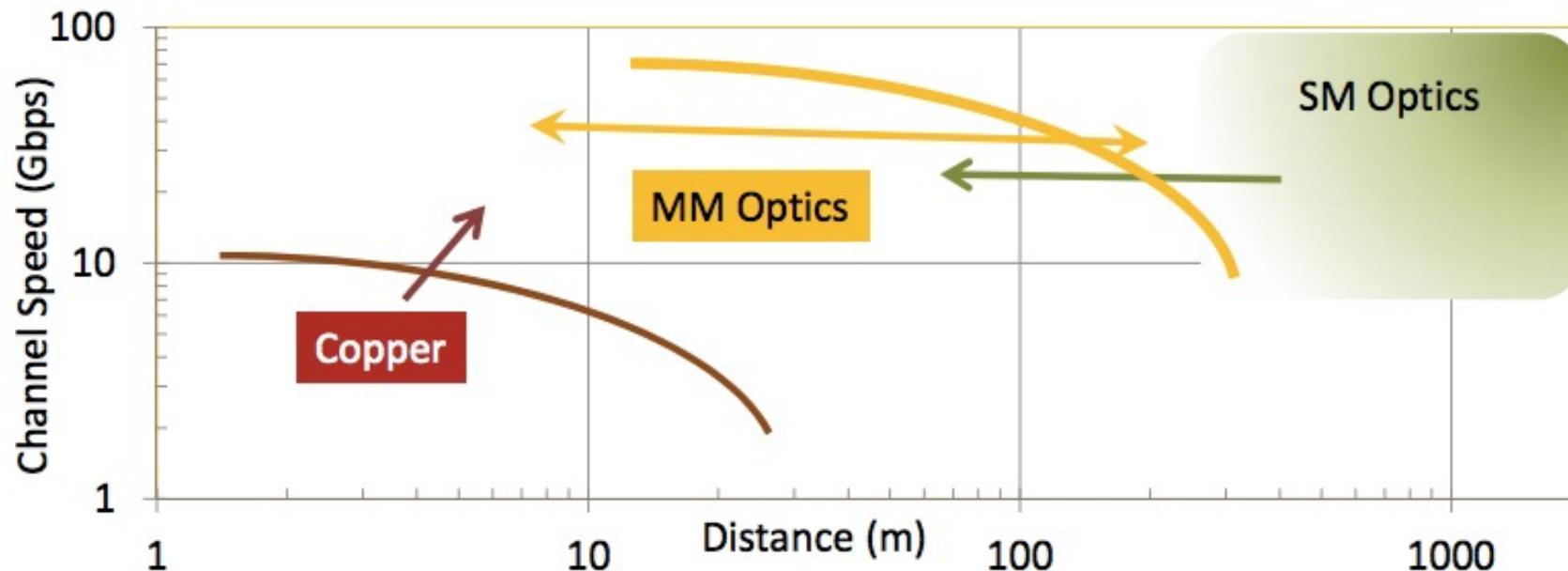
- 8 32-way SMP nodes
- Per SMP node:
  - 1 TF
  - 128 GB DRAM
  - >512 MB/s memory BW
  - >190 GB/s network BW



# Transceivers vs Reach



# Competing transmission medium



*OFC 2014 Th3C.2: D. Kuchta et al "64Gb/s Transmission over 57m MMF using an NRZ Modulated 850nm VCSEL" (also >250m at 40G and > 100m at 60G)*

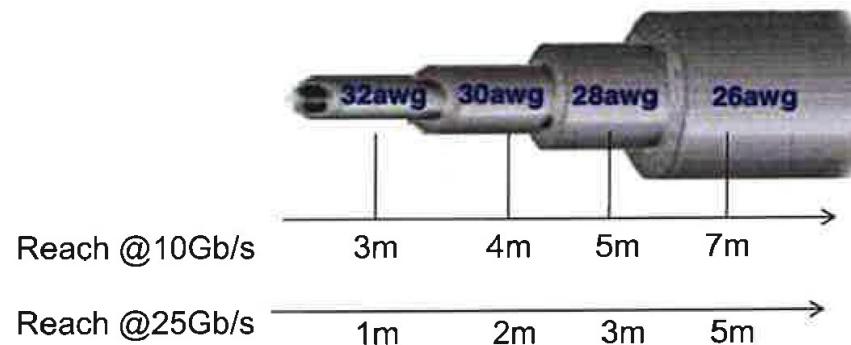
MM: Multi-Mode

SM: Single-Mode

# System packaging: Copper Cables

- 

	32AWG	30AWG	28AWG	26AWG
<b>Wire diameter</b>	0.20	0.25	0.32	0.41
<b>8x Cable outer diameter (mm)</b>	4.4	5.5	6.5	7.1
<b>Attenuation @ 5GHz (dB/m)</b>	-4.0	-3.3	-2.4	-2.0
<b>Attenuation @ 12.5GHz (dB/m)</b>	-6.4	-5.2	-3.8	-3.2



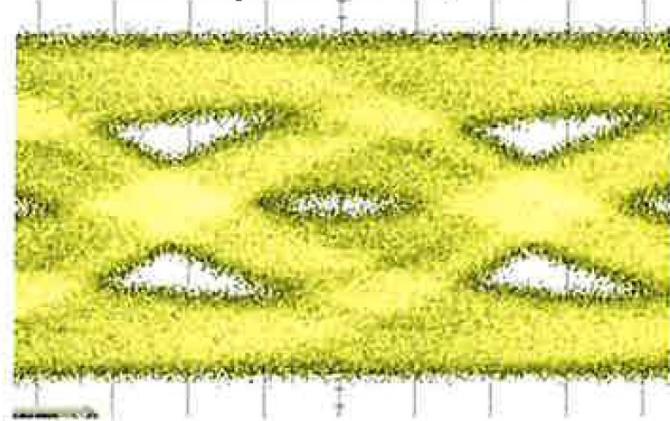
- A single copper cable can support a single channel.
- The attenuation is dependent on the signal bandwidth.
- Copper cables experience 3dB per meter attenuated @ 12.5 GHz.
- They can offer up to 5 meters reach for 25 Gb/s data channels.

# System packaging: Copper Cables

- Electrical dispersion compensation continues to drive copper cable to higher data rate and longer reach
- However, it is increasingly difficult and costly in terms of:
  - cost, power and density

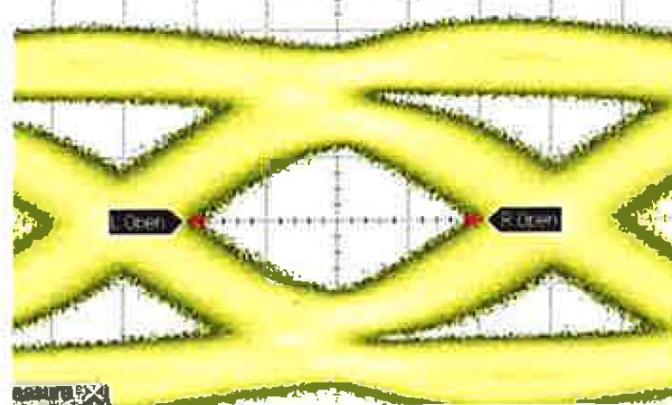
10Gbps Electrical Link in 90nm CMOS over 1m, 5m of 26 Gauge Cable

**No Equalization, 1m**



Cable + PCB loss: -10dB at 5GHz

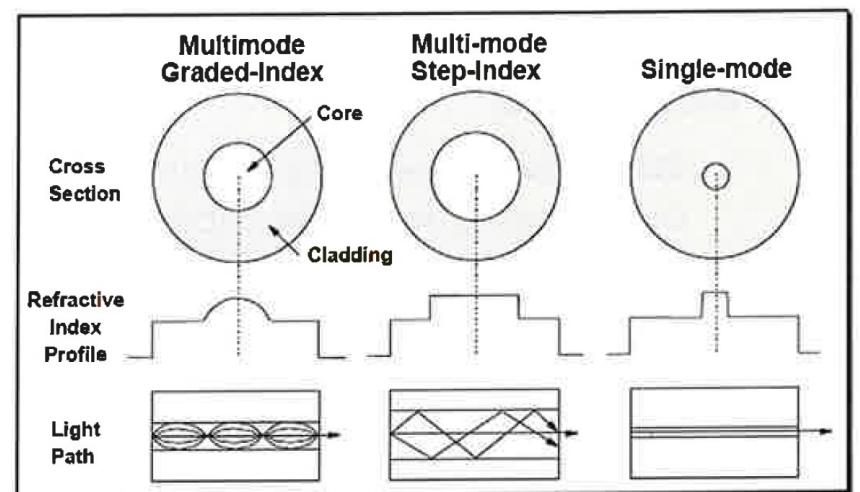
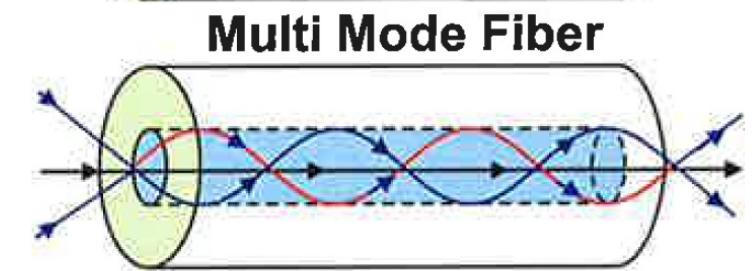
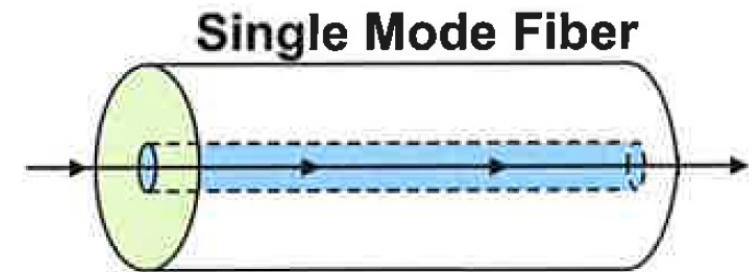
**Equalization, 5m**



Cable + PCB loss: -20dB at 5GHz  
Transmitter has one pre, one post tap

# Fibres

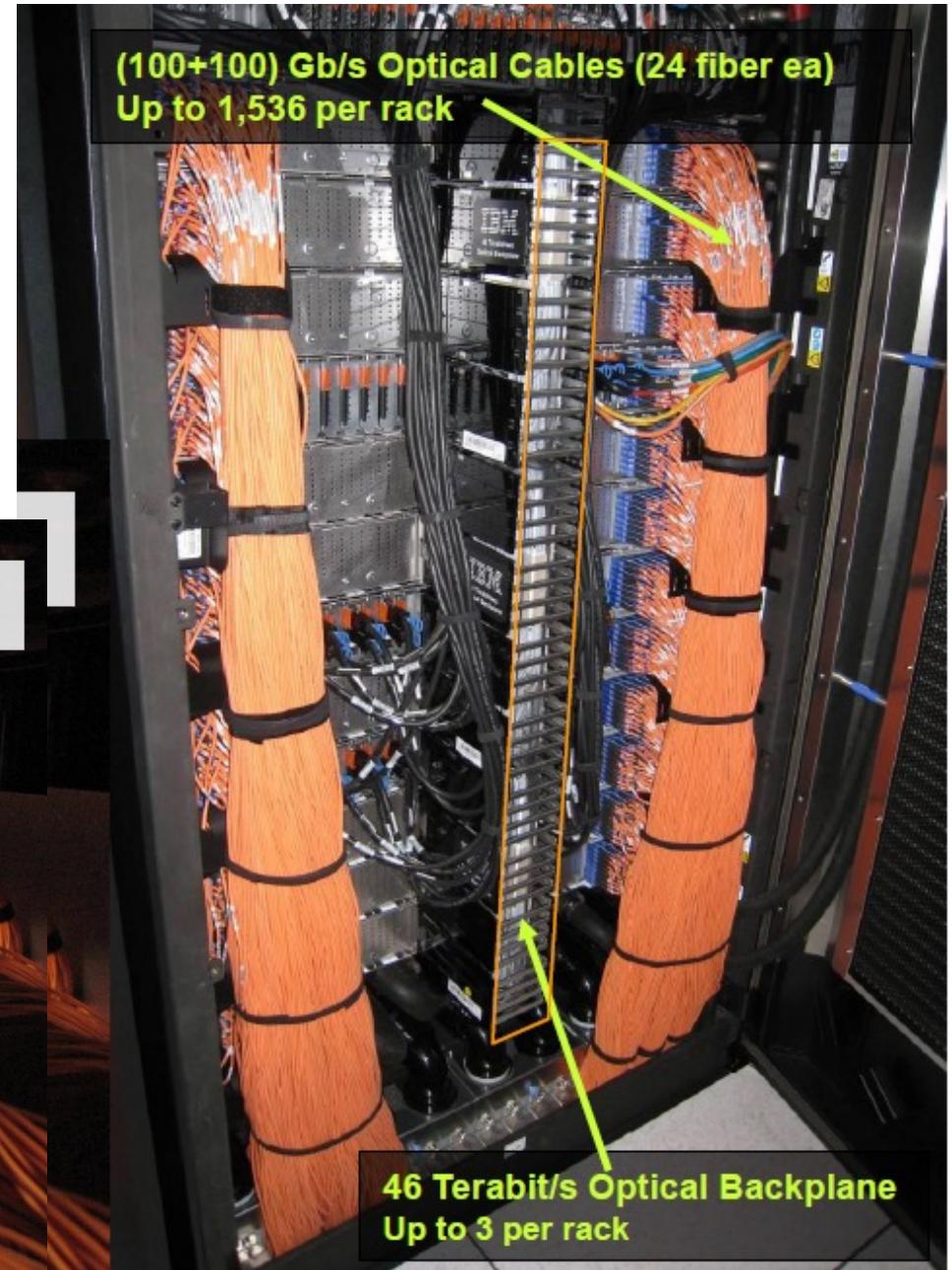
- SMF: core diameter: 9  $\mu\text{m}$  (single core)
  - Attenuation:  $\sim 0.2 \text{ dB/Km}$
- MMF: core diameter: 50  $\mu\text{m}$ , 62.5  $\mu\text{m}$
- MMF (Multi Mode Fiber) has limited bandwidth
  - Modal dispersion limited: Multiple modes propagate at different velocities causing pulse spreading
  - Graded parabolic index profile reduces modal dispersion
    - Requires careful control of doping concentration layer by layer from the center of the core
- SMF (Single Mode Fiber) has 10s of Tb/s bandwidth by use of WDM (Wavelength division multiplexing)



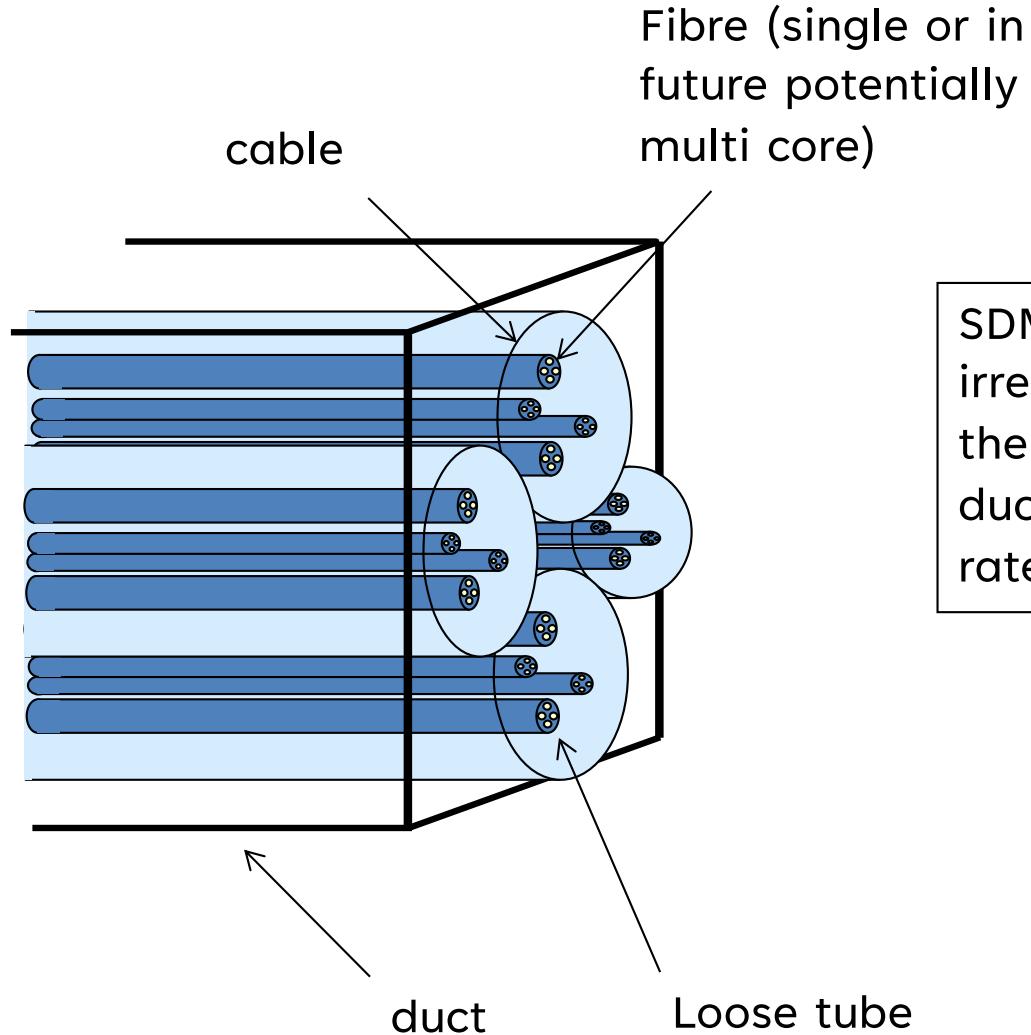
# More fibres and bandwidth

More fibres are needed for connection.

- Wavelength division multiplexing (WDM)
- Space division multiplexing (SDM) based on Multi-core fibres



# What is ‘SDM’ (Space Division Multiplexing)?



SDM is any use of multiple cores, irrespective of whether they are in the same fibre, same cable, same duct, to increase the transmission rate of a link in a network

# Homogeneous hexagonal Multi-Core Fiber (MCF)

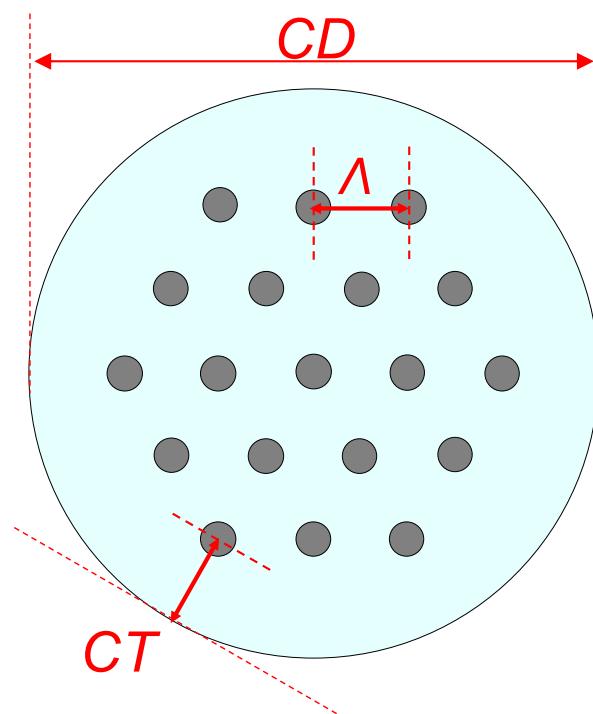
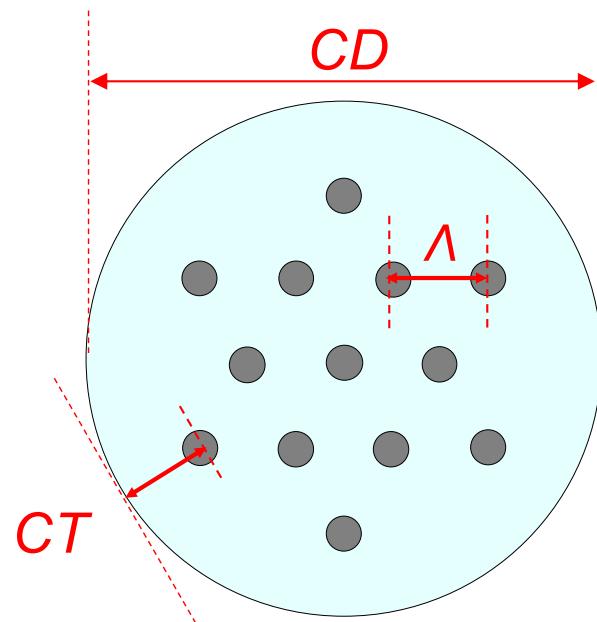
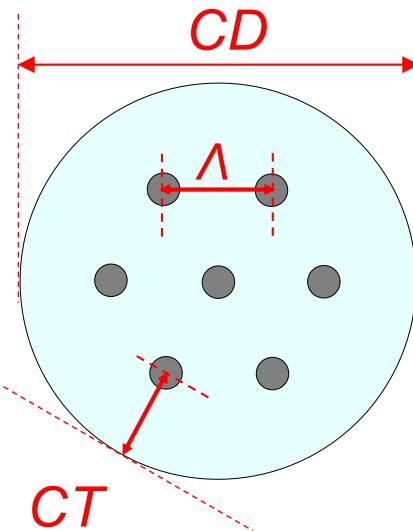
- ✓ Cladding diameter (CD) limit :  $\sim 230 \mu\text{m}$
- ✓ Outer clad thickness (CT):  $\sim 40 \mu\text{m}$

1-mm decrement of core pitch results in about 3-dB degradation of XT

19-core  
( $\Lambda < 37 \mu\text{m}$ )

7-core  
( $\Lambda < 70 \mu\text{m}$ )

13-core  
( $\Lambda < 43 \mu\text{m}$ )



- Cross-section Area
  - $A_{\text{cross}} = \pi/4 * (\text{Cladding Diameter})^2$

# Relation between number of cores and crosstalk level

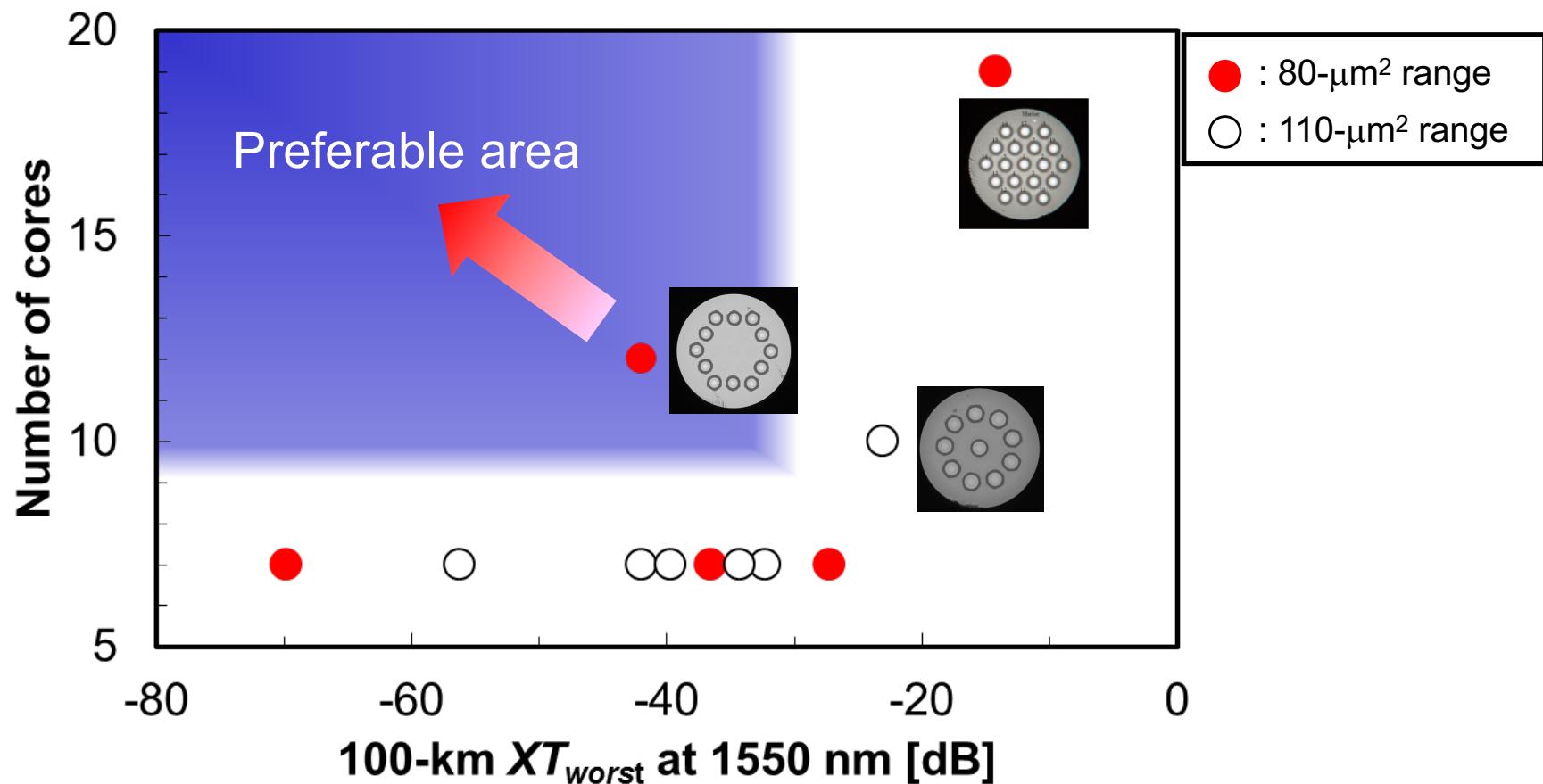


University of  
BRISTOL

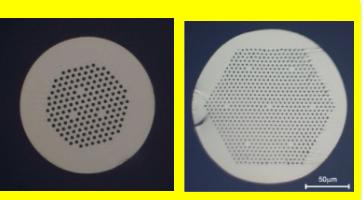
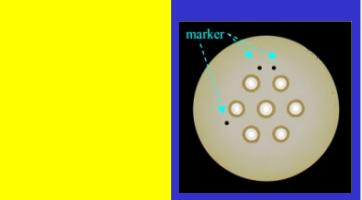
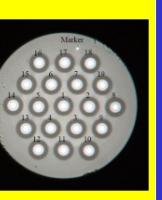
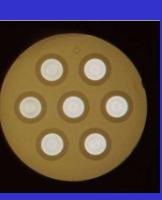
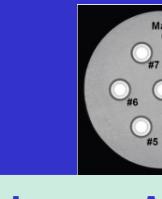
Core arrangement?  
Maximum number of cores?

depend on ...

Required crosstalk level  
Maximum cladding size  
Target cross-section area ( $A_{\text{cross}}$ )

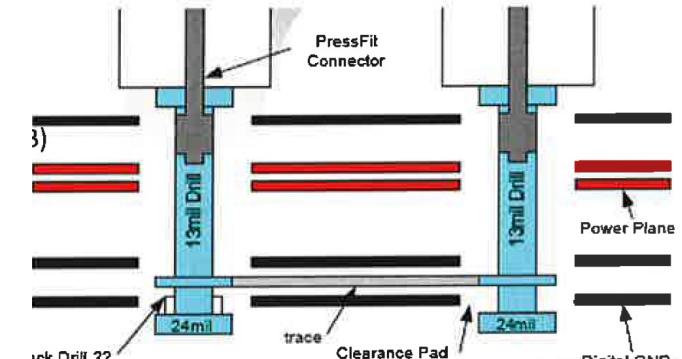


# Fabricated single-mode uncoupled MCFs

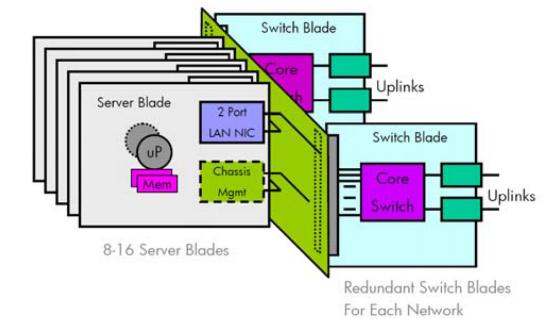
	2008	2009	2010	2011	2012	
Fujikura Ltd.			Air-Clad			
Furukawa Elec.						 
Sumitomo Elec.				 		 
OFS			All-Solid		Crosstalk suppression with index trench	 
Corning						Inter connection

# System packaging: High Speed Backplane

- IEEE P802.3ap – Ethernet Backplane
  - Support up to 1m trace on improved FR-4 PCB + 2 connectors
    - 10GBASE-KR: 10 Gb/s serial
      - 64B/66B encoding
      - Tx pre-emphasis,
      - Adaptive Rx equalization
      - + DFE, FFE , FEC (3dB)
    - 100GBASE-KR4: 25 Gb/s over 4 lanes
      - CAUI 25 Gb/s/lane
      - 64B/66B encoding
      - Tx pre-emphasis,
      - Adaptive Rx equalization
      - + DFE, FFE , FEC (5dB)
  - Key considerations
    - Improved design practices for higher speeds
      - Better board material
      - Better SI (Signal Integrity) practices: connector and pin-out, vias, back drill, etc.



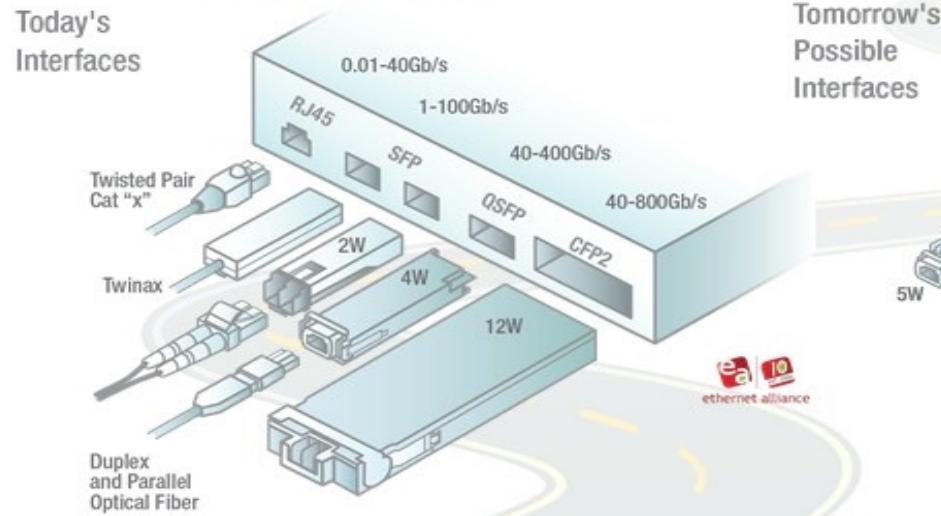
Server Blade & Switch Architecture



Typical configuration includes a daughter card on the Server Blades

# Transceivers size/ form factors

## Modules



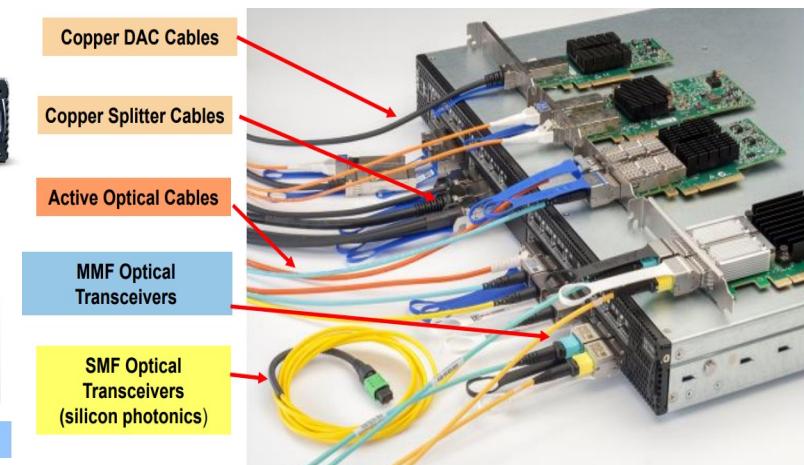
- SN2700: 32-port Non-blocking 100GbE Open Ethernet Spine Switch System (3.2Tb/s)



- SB7700 - 36-port EDR 100Gb/s InfiniBand Switch System (3.6Tb/s)

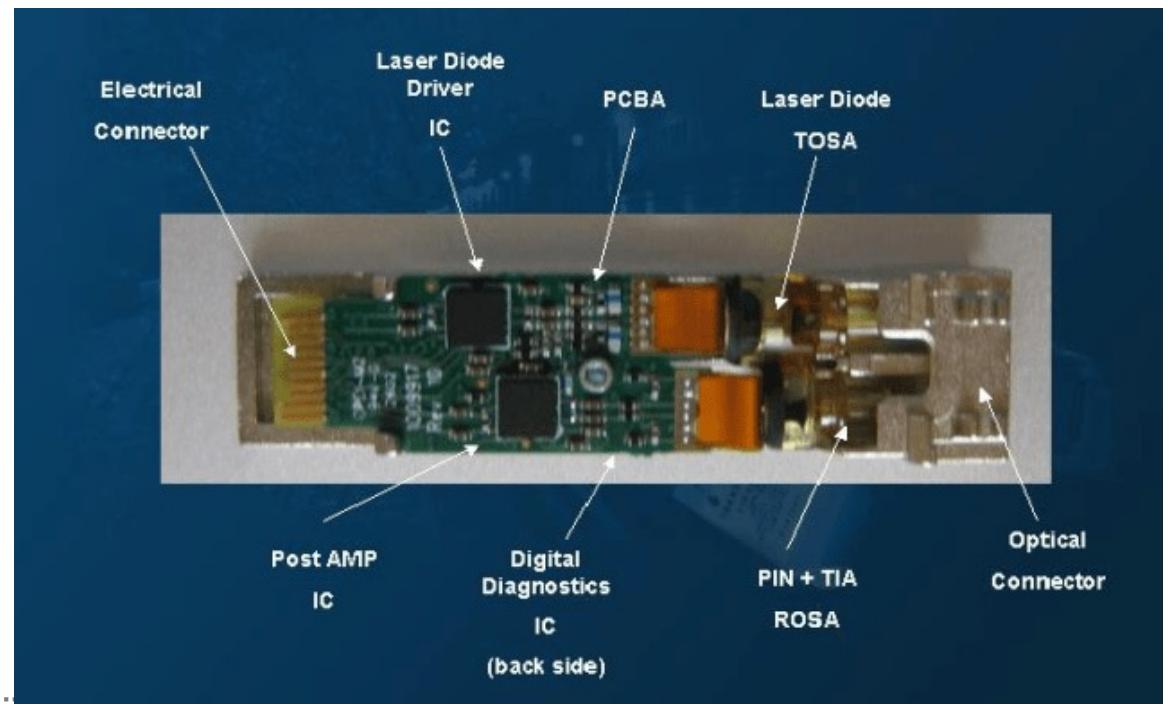
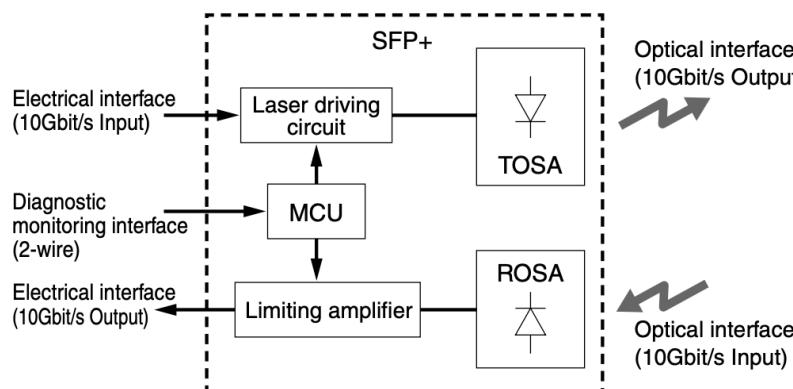


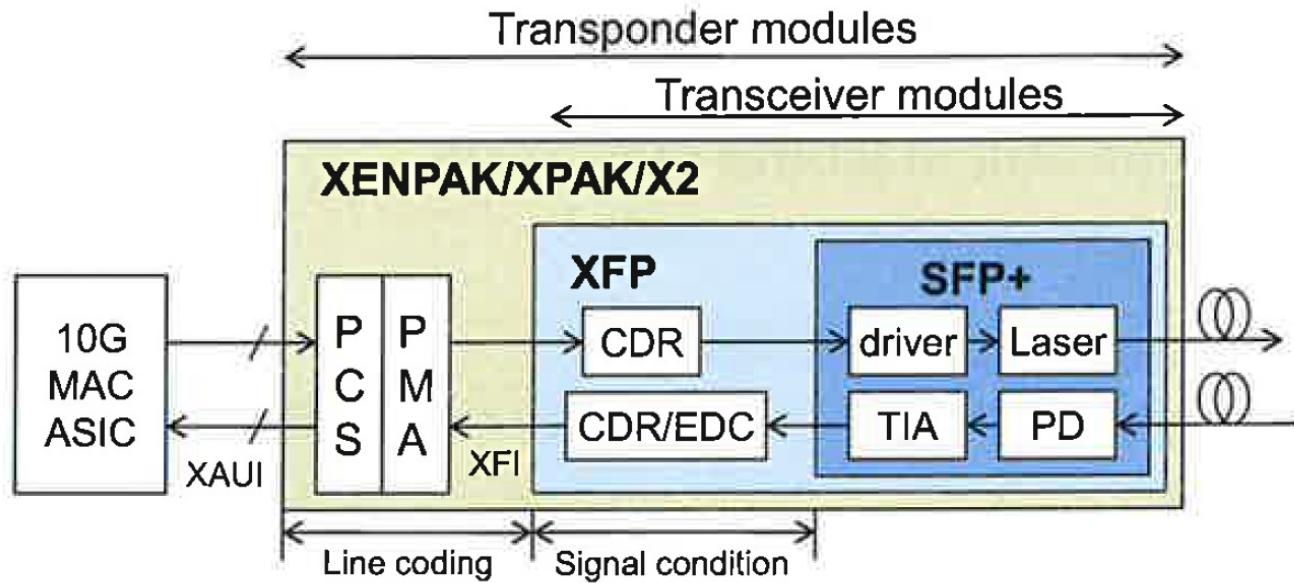
**Note:** QSFP was a great improvement over SFP; it doubled the density, but we need more



# SFP+: Small Form-factor Pluggable Transceiver

SFP+ is an enhanced version of the SFP that supports data rates up to 10 Gbit/s. SFP+ supports 8 Gbit/s Fibre Channel, 10-gigabit Ethernet and Optical Transport Network standard OTU2.

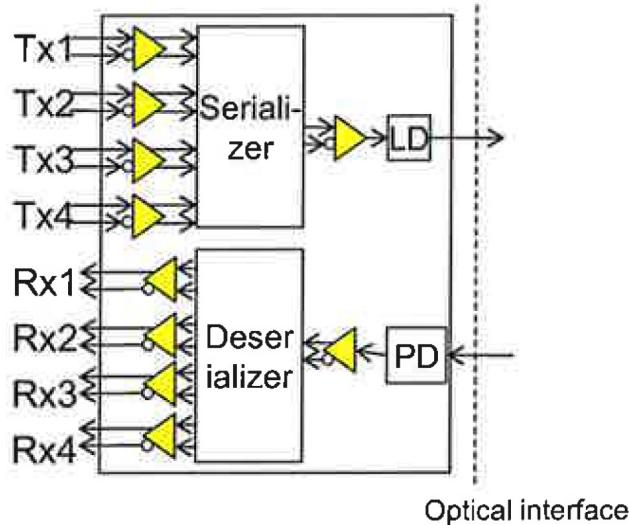




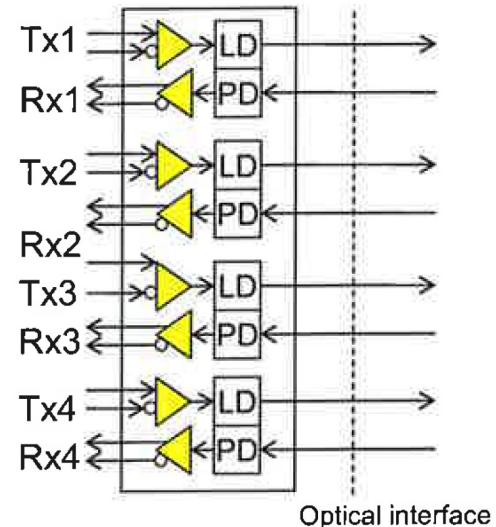
- Transponders (XENPACK, XPAK, X2) with PHY+ PCS (Physical Coding Sublayer) gave way to transceiver-based solutions
  - Footprint density, power density and price
- Coding and signal conditioning functions can be integrated with ASICs to improve density and power footprint
- Unified transceiver design can be shared with multiple applications: Ethernet, FiberChannel and Infiniband

# More Bandwidth: multiple channels

Serial Transceiver (TDM)

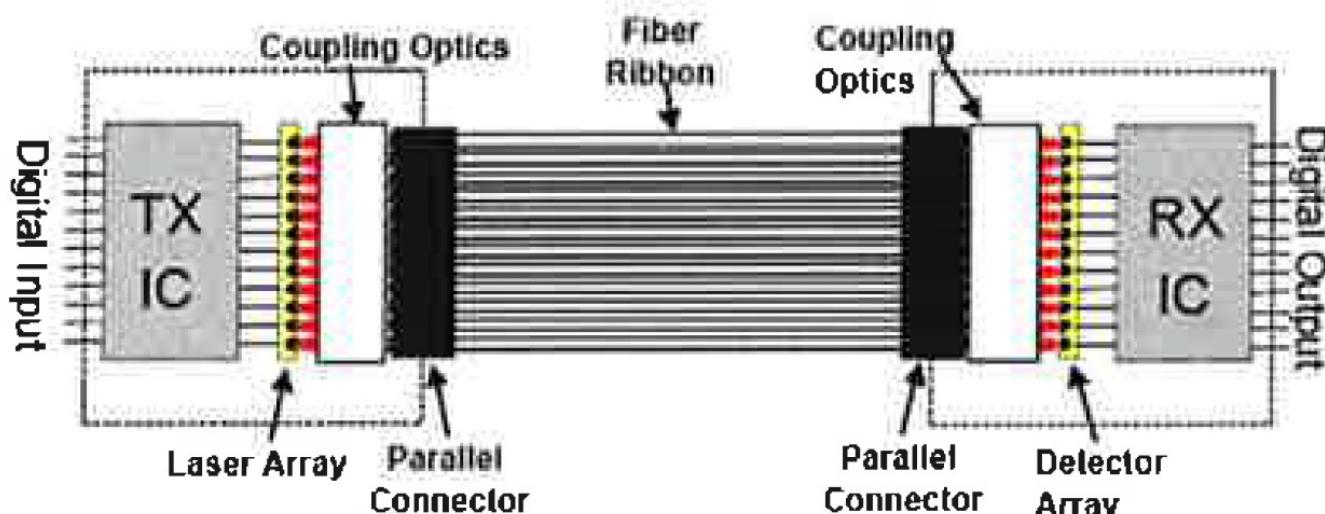


Parallel optical transceiver (SDM)



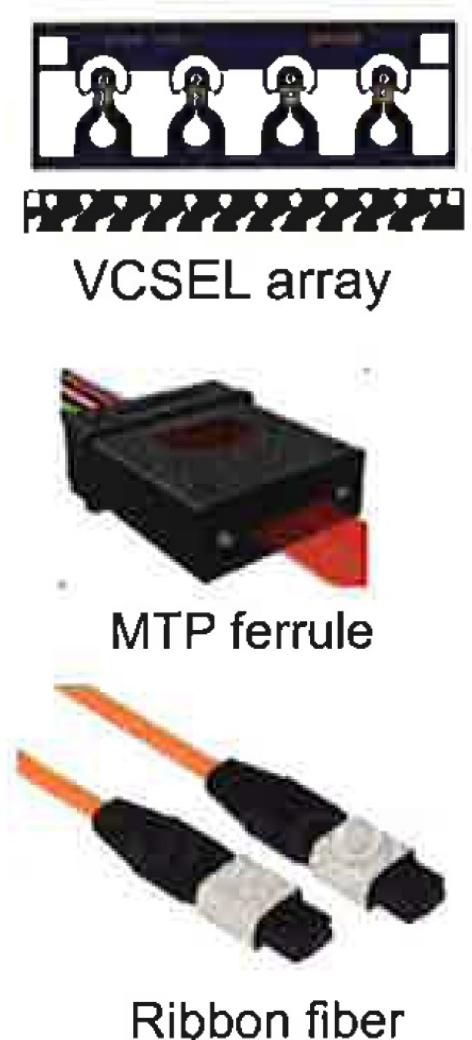
- Serial Transceiver: Beyond 20G, electronic and optical performance becomes challenging
  - Gear box consumes power and £
  - Speed mismatch between electrical and optical lanes
- Parallel transceiver takes advantage of parallel data lanes of microprocessor systems
  - Less power consumption due to the absence of SerDes or “gear box”
  - Limited reach and does not solve cabling problem

# More Bandwidth: Parallel Optical Tranceiver



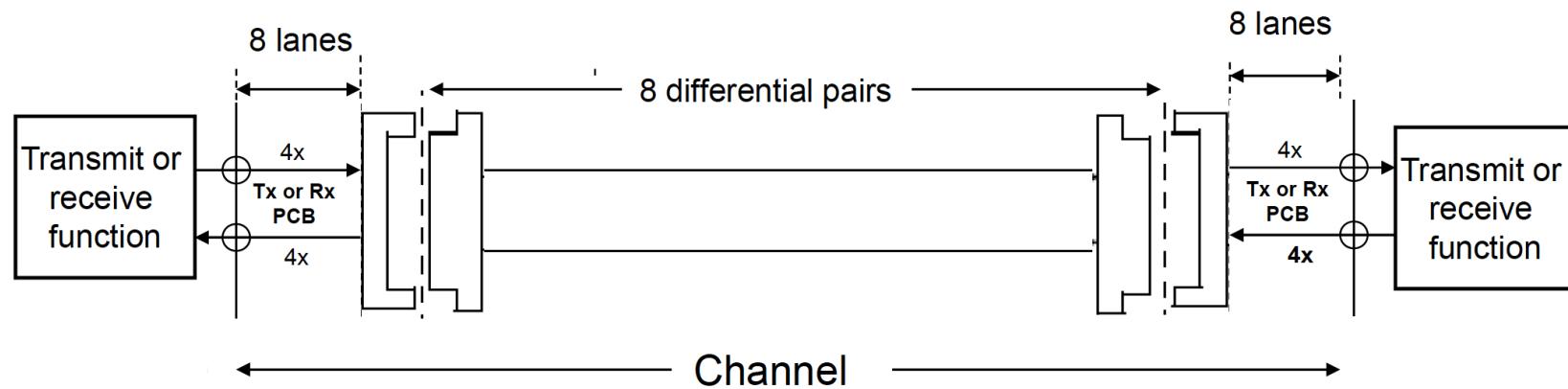
**TX Module**

**RX Module**

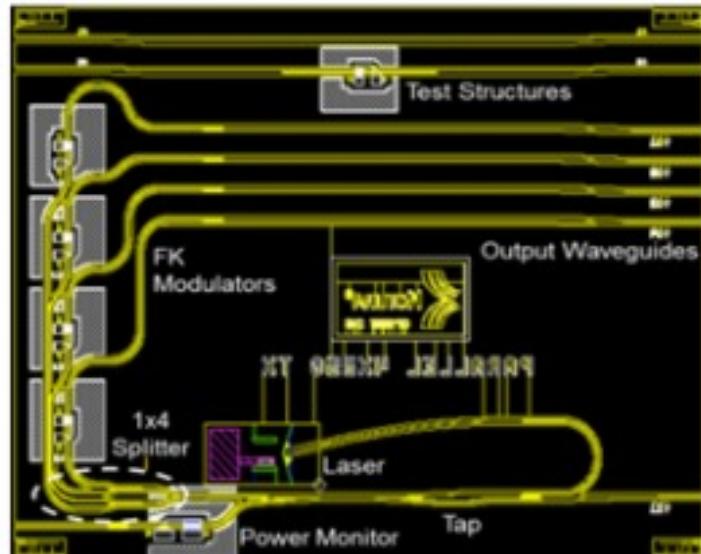


- Challenges for VCSELs as high-speed transmitter
  - Low output power due to short gain region => limited reach of 10
  - tight power budget for patch panel and MTP
  - Trade-off between power and speed

- 8 lanes/pairs versus 20 lanes/pairs
  - Smaller form factor enables higher port densities.
- Cost factors
  - Fewer cable assembly pairs reduces cost.
    - 20-pair assembly ~1.5 times cost of 8-pair assembly.
  - Reduction in number of Tx/Rx lanes simplifies host routing.
  - Reduction in cable assembly pairs enables reduction in cable diameter.
  - Fewer pairs => higher levels of integration => higher port density
- Could define plug compatible port types to enable plug-and-play for 4x10Gb/s and 4x25Gb/s.

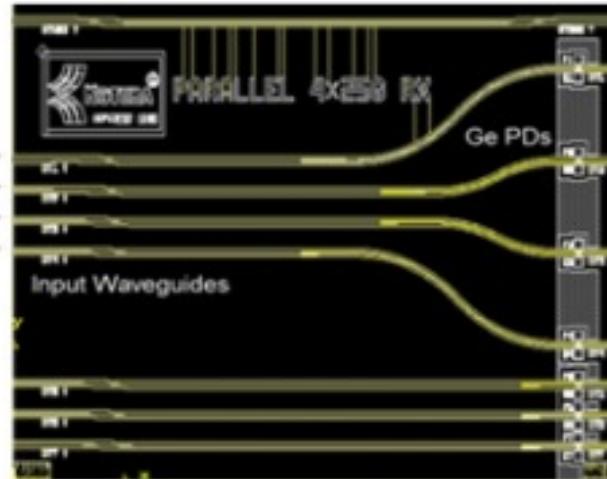
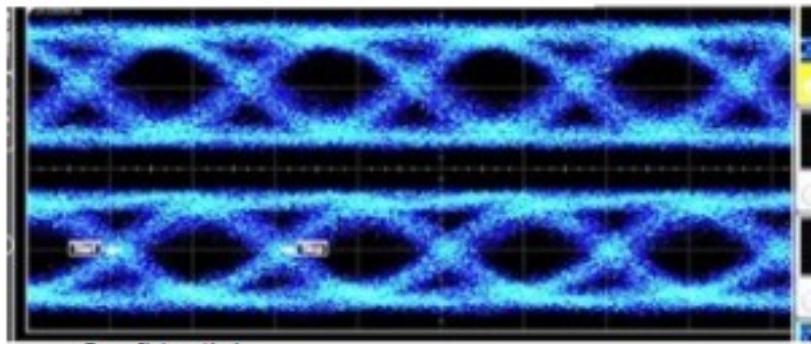


# 4x25G Parallel optical link

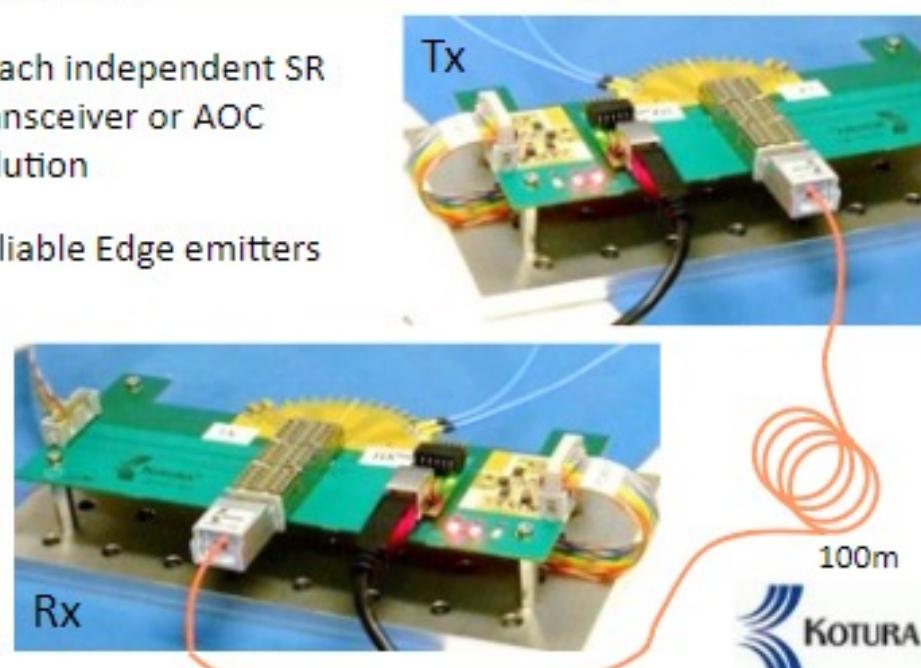


- 4x25G Parallel in QSFP

Differential Rx Eye at 25Gb

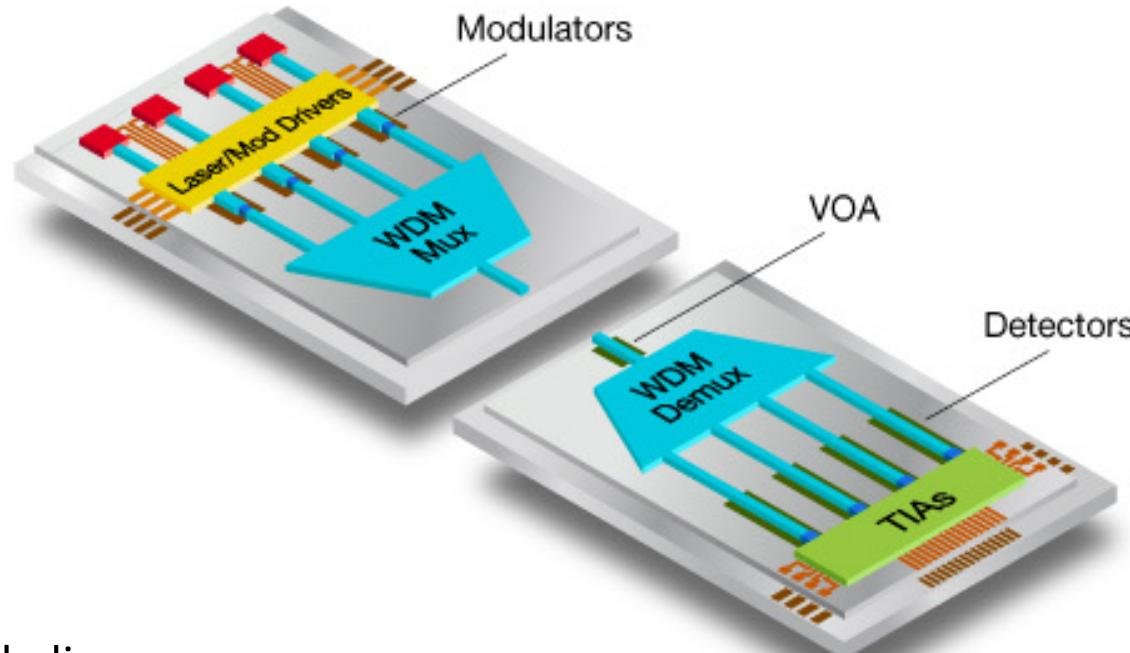


- Reach independent SR transceiver or AOC solution
- Reliable Edge emitters

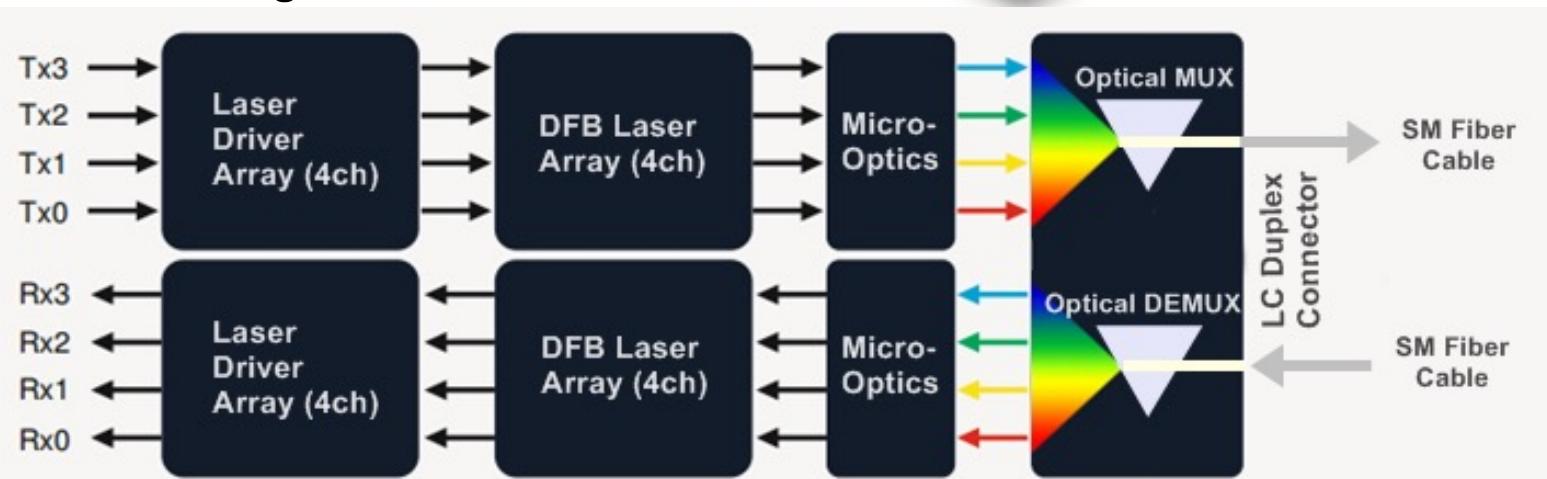


- Current:  
~1.5W
- Next gen:  
0.8W
- Assuming no  
CDR
- SMF low  
jitter

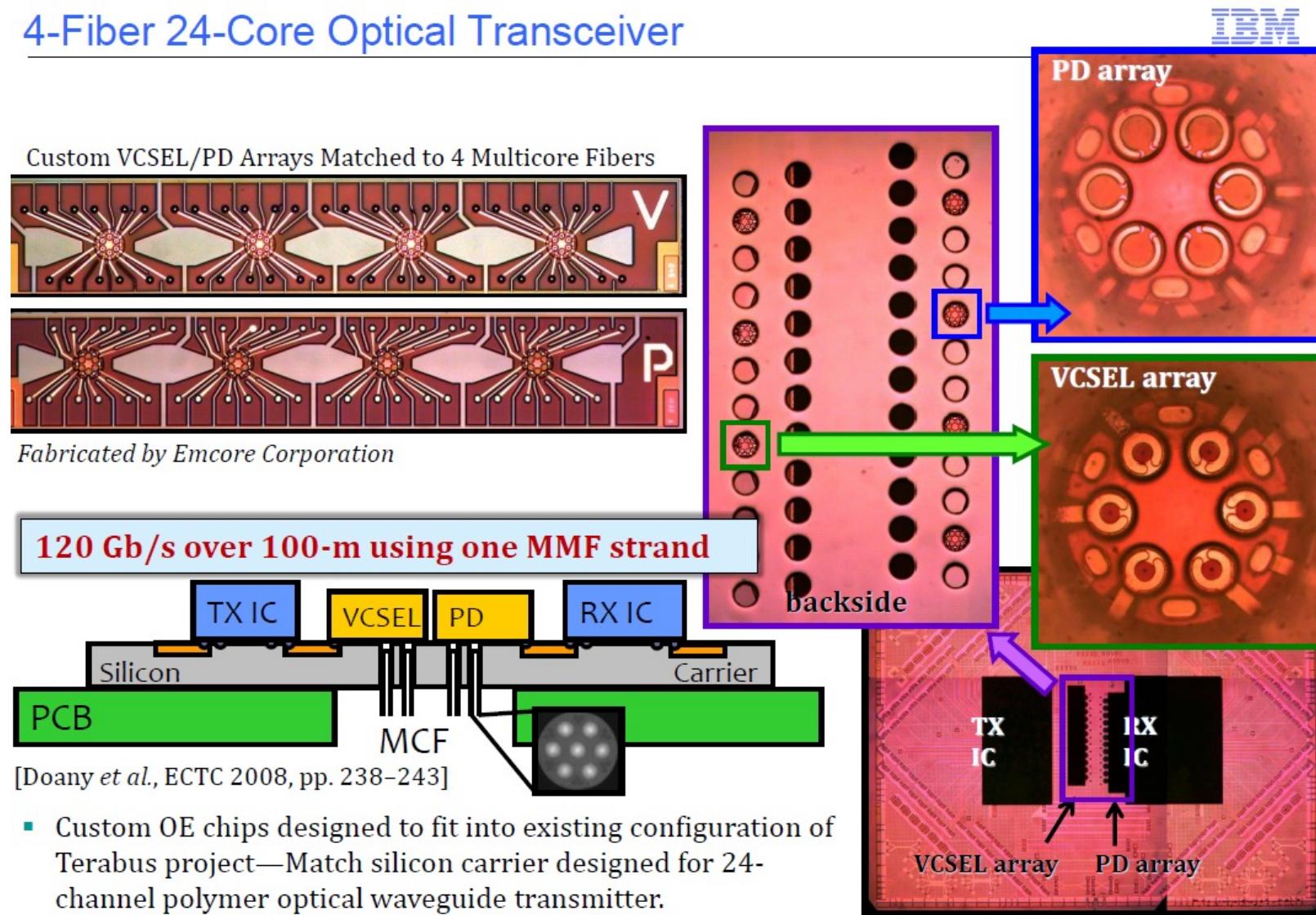
# WDM transceivers: 4 CWDM channels



QSFP block diagram



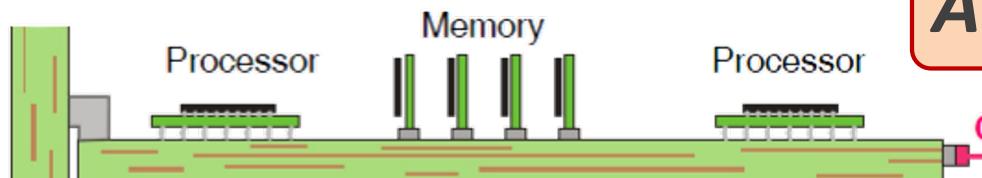
## 4-Fiber 24-Core Optical Transceiver



© 2011 IBM

# Optical interconnects hierarchy

Backplane



Electrical system, optical fibers at card edge

## Active Optical Cables

Optics



Development



Optical fibers across the boards

## On-board subassemblies

Optics

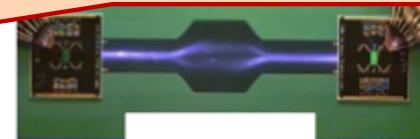


Development



Optical waveguides in/on boards

## Optical PCBs & C2C



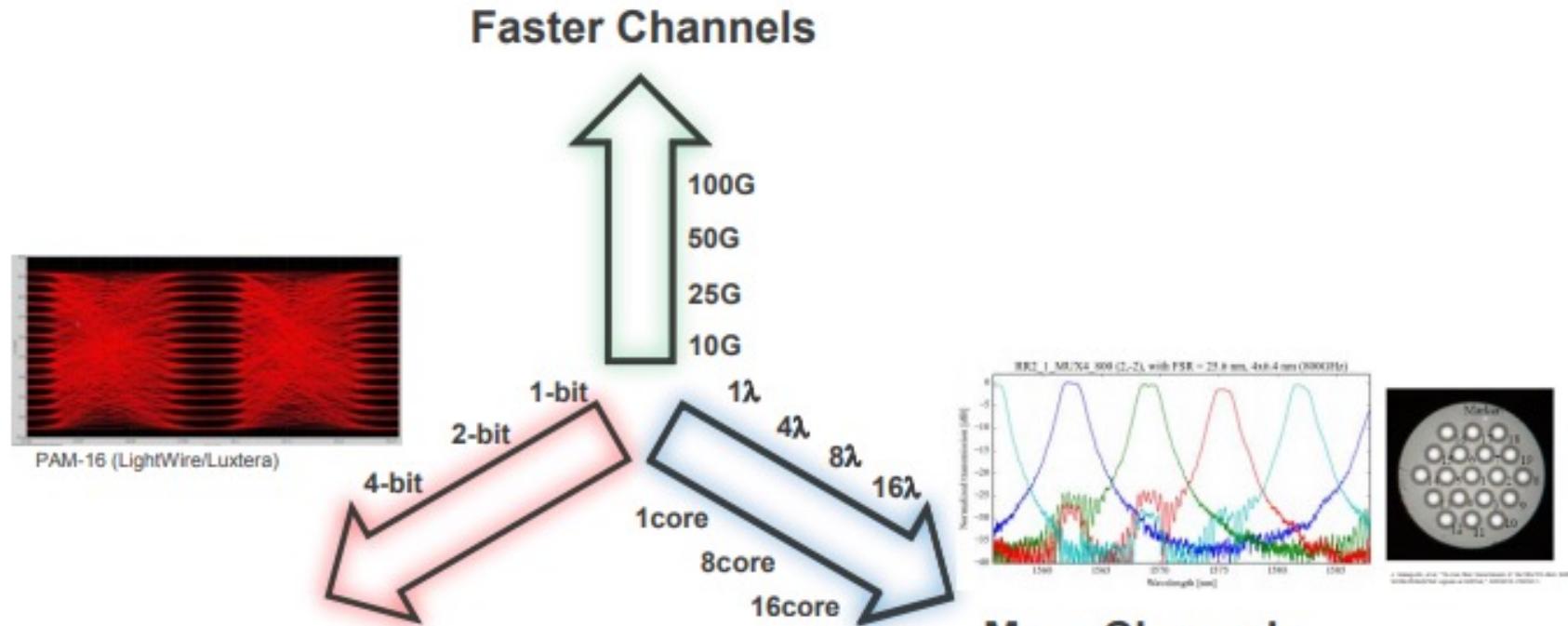
Research



Optical interconnects integrated with the processor

## Network-on-chip / Crossconnects





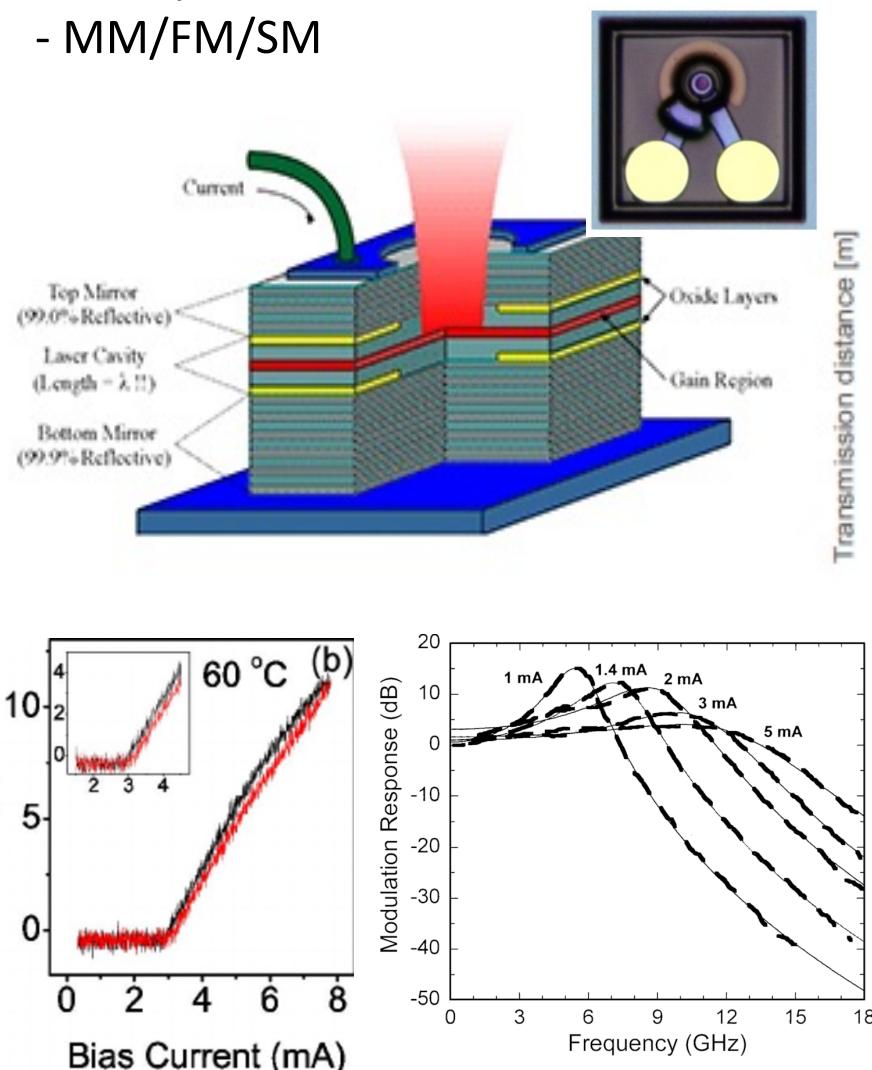
### More Bits per Symbol

- Amplitude: PAM-X
- Phase and Amplitude: DP-QPSK, QAM-X, ...

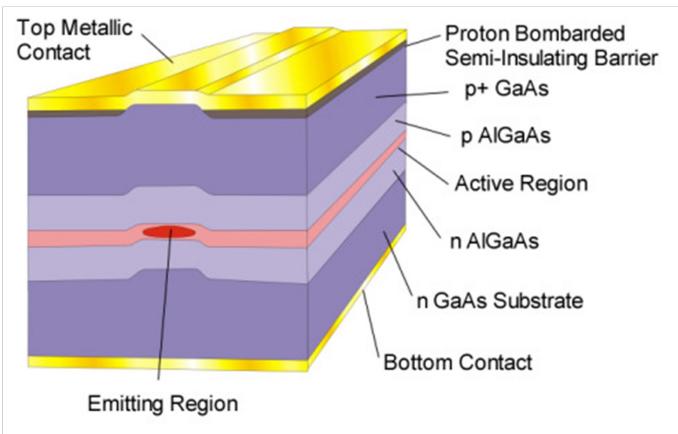
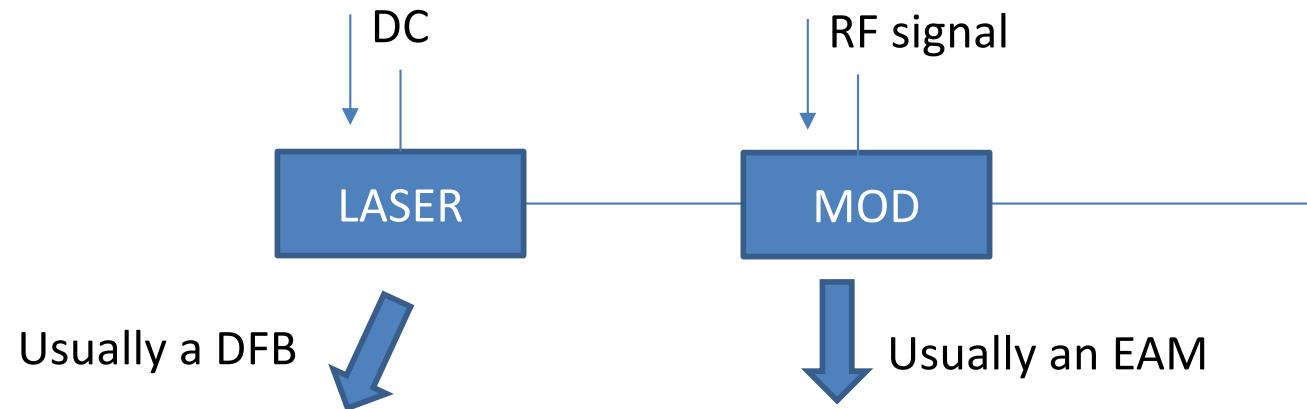
### More Channels

- Parallel (Single-Mode) Fiber [PSM]
- Multi-Core Fiber, Spatial Division Multiplexing [SDM]
- Wavelength Division Multiplexing [WDM]

## Directly modulated lasers - MM/FM/SM



# Externally modulated lasers (EML)

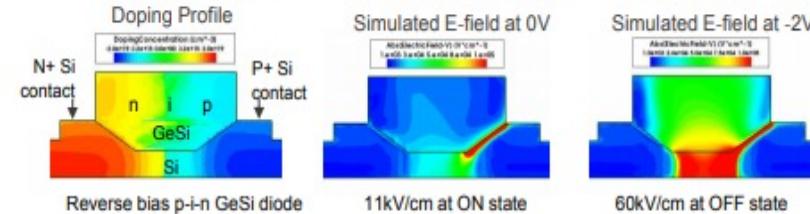


Single-mode  
High optical power output

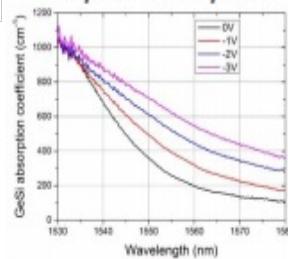
## Franz-Keldysh effect

- Absorption coefficient increases with applied field (C-band)
- Sub-picosecond effect

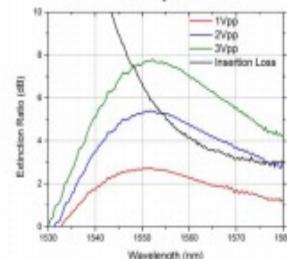
*Device cross section, perpendicular to light propagation*



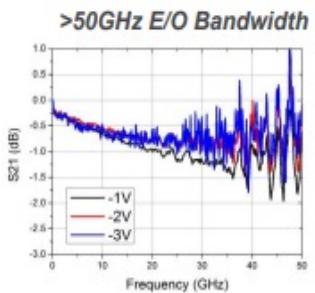
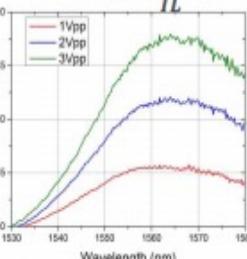
## Optical absorption



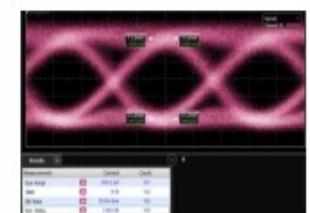
## ER, IL



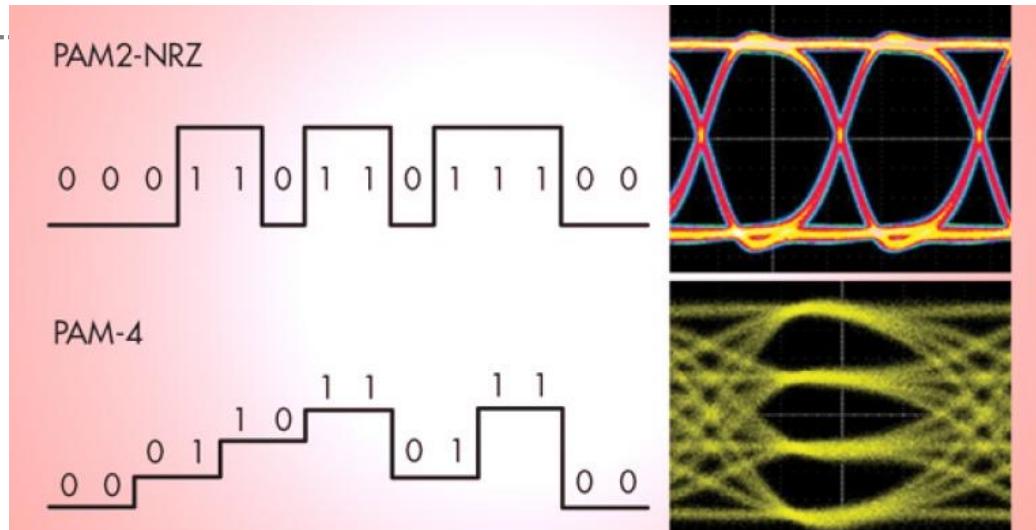
$$\text{FOM} = \frac{\text{ER}}{\text{IL}}$$



## 56Gb/s NRZ Eye

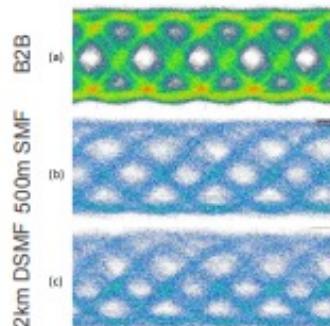
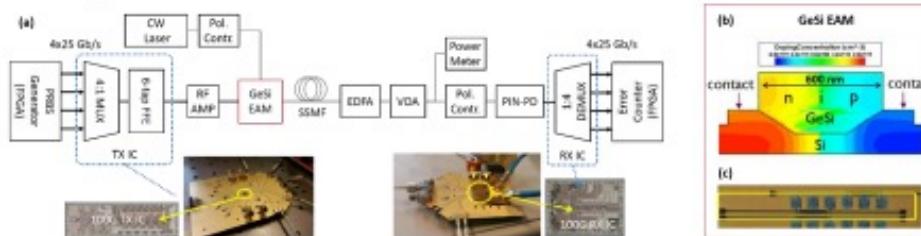


# PAM-4 format

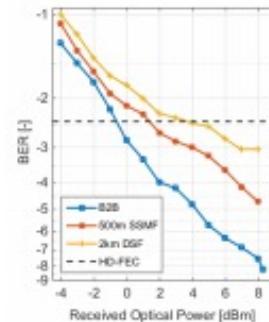


GeSi EAM-PD at 100Gb/s

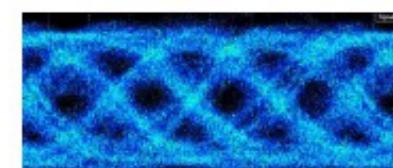
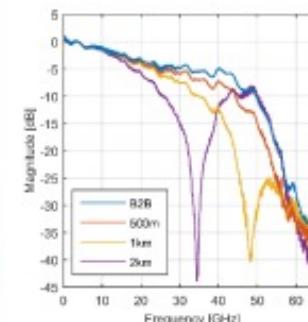
*First single-wavelength 100Gb/s NRZ-OOK modulation demo in Silicon Photonics*



100Gb/s Eye Diagrams  
(III-V discrete PD)



BER below threshold for FEC

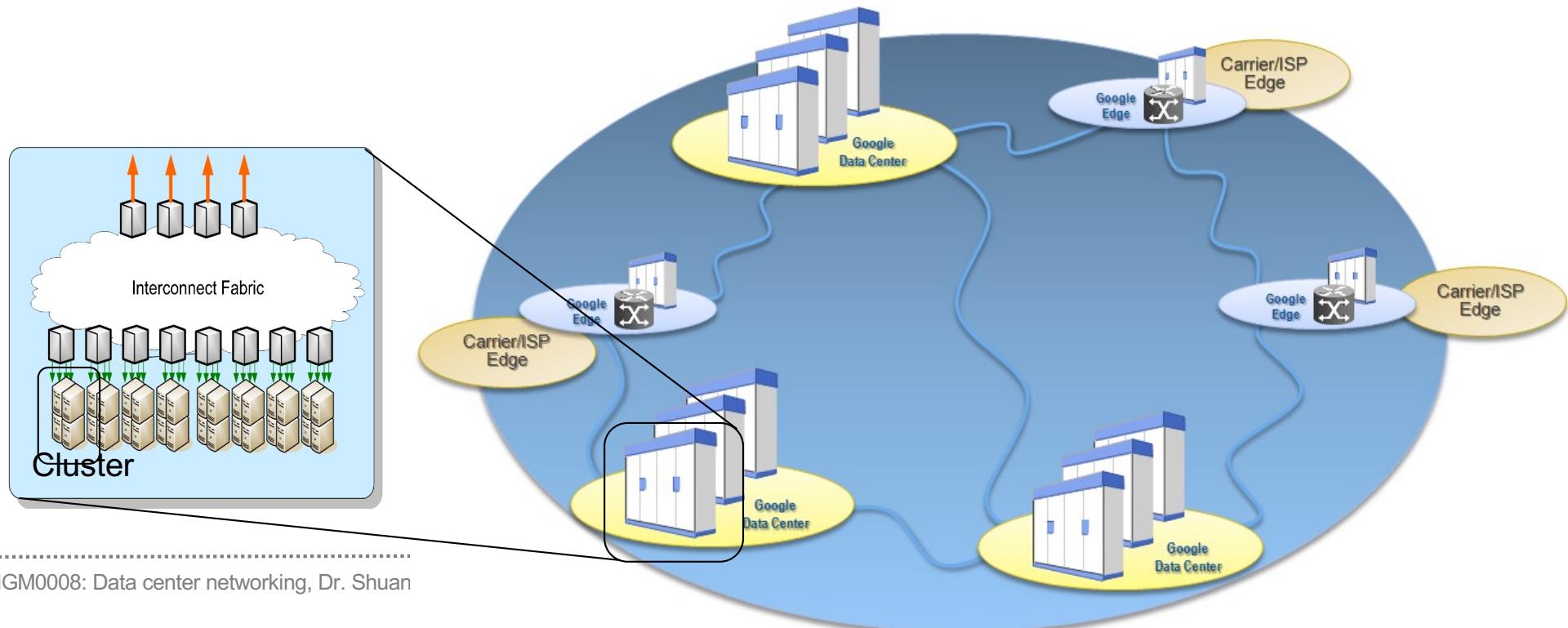


100Gb/s Eye Diagram  
for GeSi EAM-to-GeSi PD

Collaboration with **IDLab**  
INTERNET & DATA LAB

Name	Fiber	Reach	Modulation
<b>400G-ZR/ZR+</b>	Duplex SMF	10km-1000km	16-QAM
<b>400G-FR4/LR4</b>	Duplex SMF	2km/10km	100G-PAM4
<b>400G-DR4</b>	8xSMF	500m/2km	100G-PAM4
<b>400G-SR8</b>	16xMMF	50m	50G-PAM4
<b>400G-CR8</b>	copper	3m	50G-PAM4

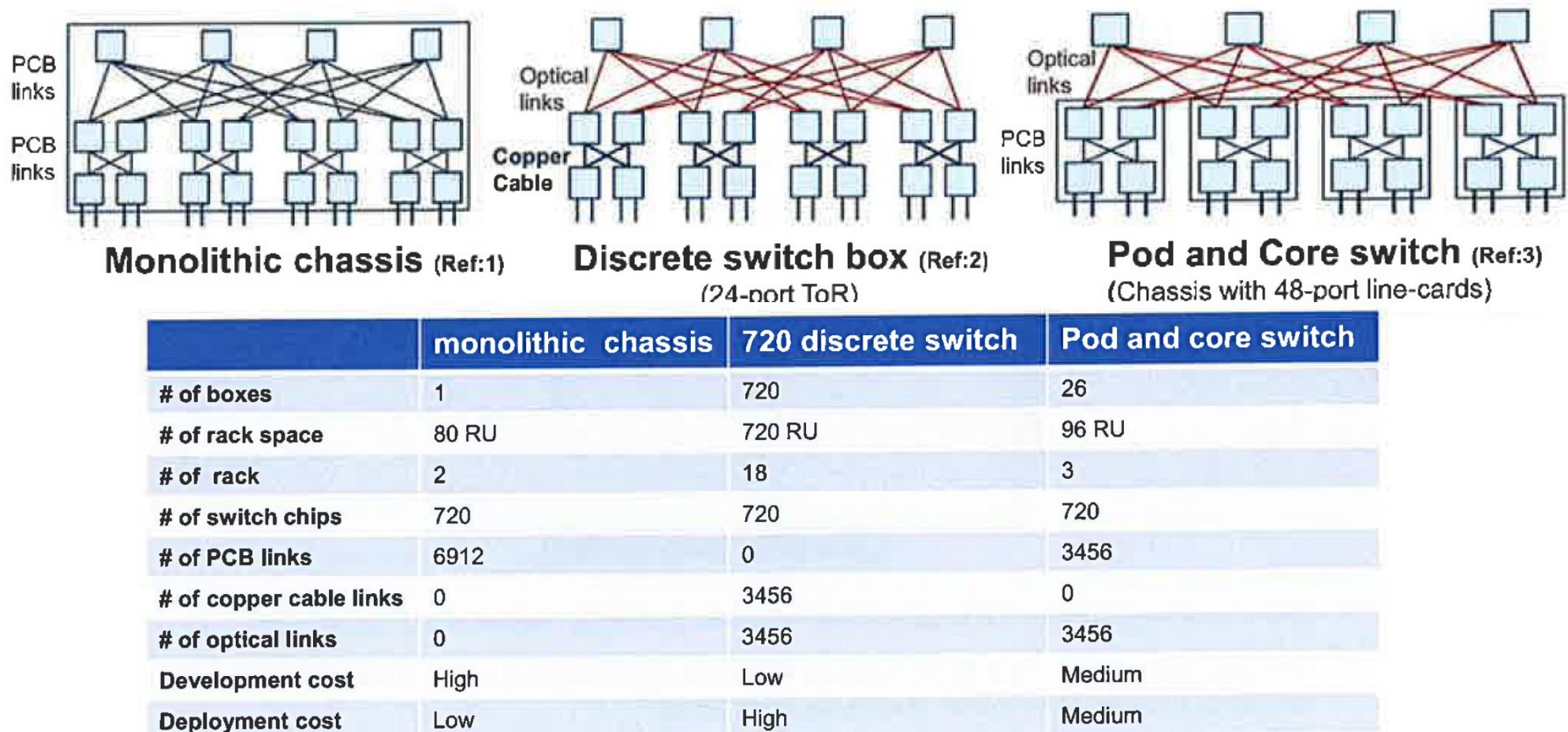
# Datacenter Optics - Overview



# Optical Switching Solutions

A possible solution for large scale switch fabrics

# System packaging for switching

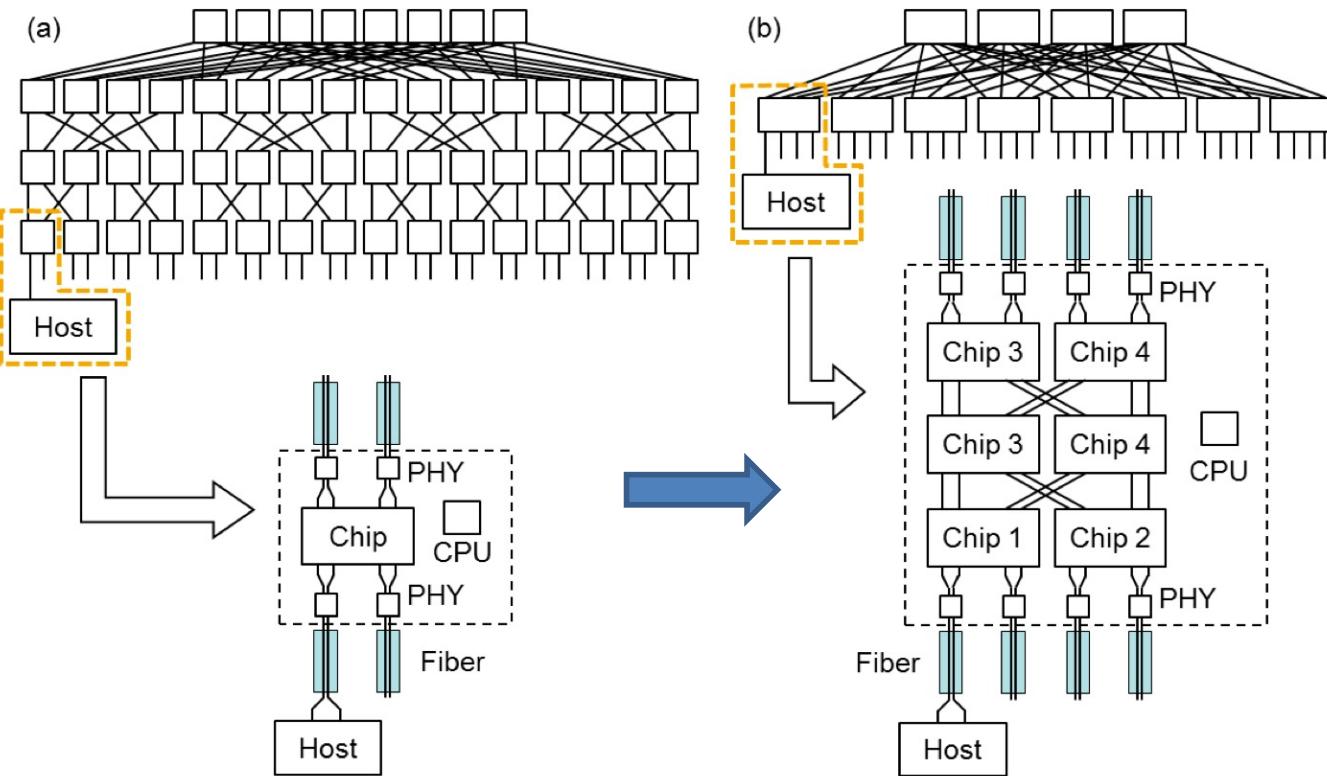


1. [http://www.sun.com/products/networking/datacenter/ds3456/ds3456\\_wp.pdf](http://www.sun.com/products/networking/datacenter/ds3456/ds3456_wp.pdf)

2. M. Al-Fares, et al, "A Scalable, Commodity Data Center Network Architecture," SIGCOMM'08, Aug. 17-22, 2008, Seattle, Washington, USA

3. N. Farrington, et al, "Data Center Switch Architecture in the Age of Merchant Silicon," IEEE HotInterconnect '09, New York, USA.

# Multi-chip



Single channel traditional

Single channel chassis-based

## Advantages:

- Multiple switch chips in a chassis
- Copper backplane connection
- Reduce fiber and optical transceiver

## Disadvantages:

- more switch chips in a chassis
- Larger worst-case number of hops (latency)
- High chassis power consumption

W. M. Mellette, A. C. Snoeren, and G. Porter, “P-FatTree: A multi-channel datacenter network topology,” 2016, pp. 78–84.

# Scale up for large scale DCNs

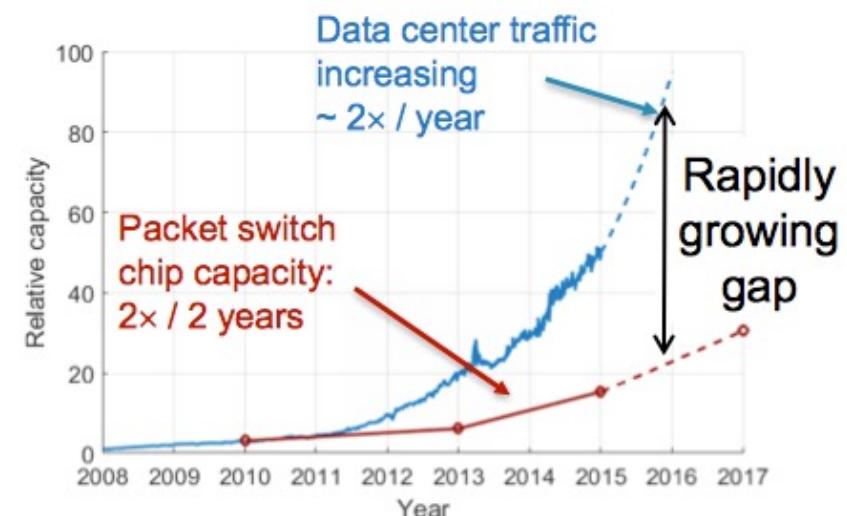
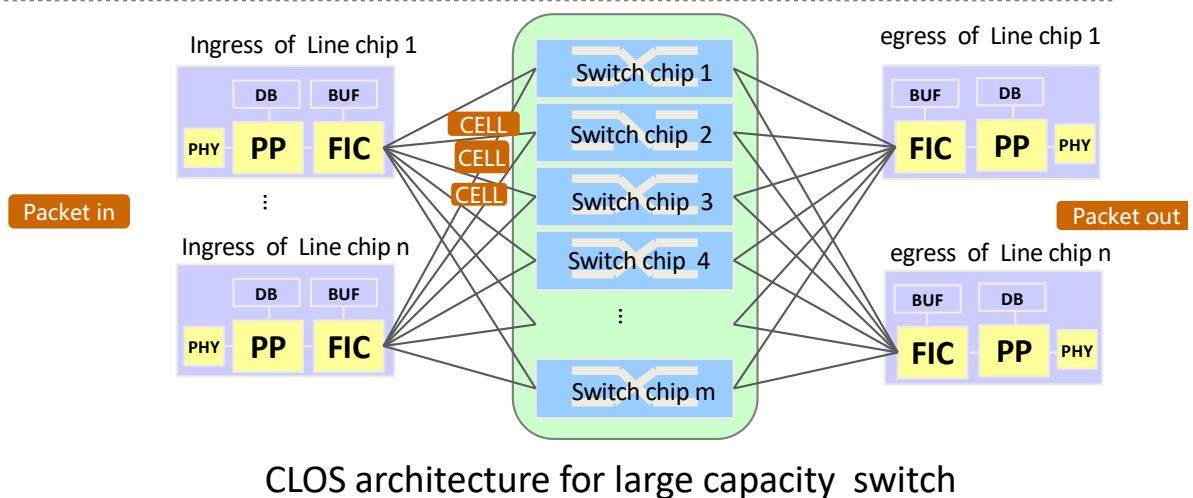
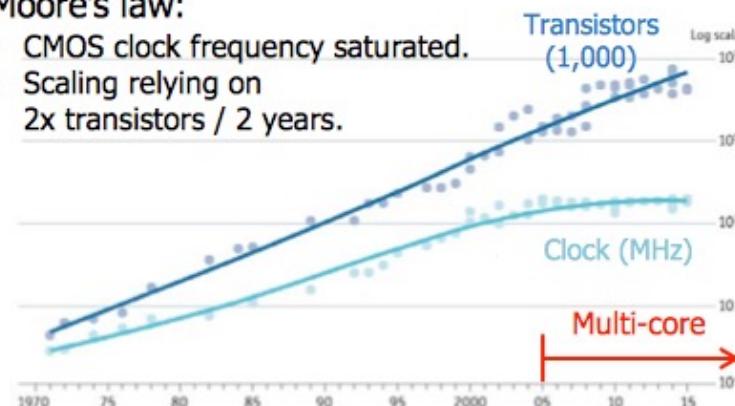
- Scale up electrical switch
  - Switch Capacity
  - Port count ( radix)

Source from:

George Rapen, OFC 2017, M3k.1

Moore's law:

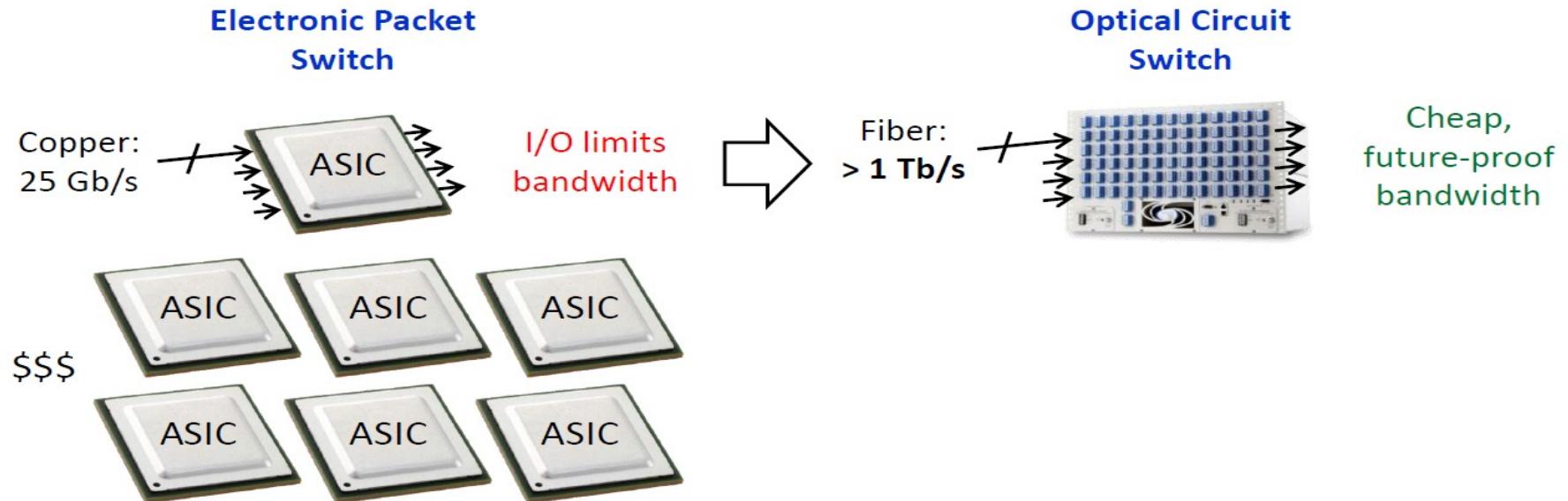
- CMOS clock frequency saturated.
- Scaling relying on 2x transistors / 2 years.



The ITRS projections for signal-pin count and per-pin bandwidth are nearly flat over the next decade. Single-chip bandwidth saturates.

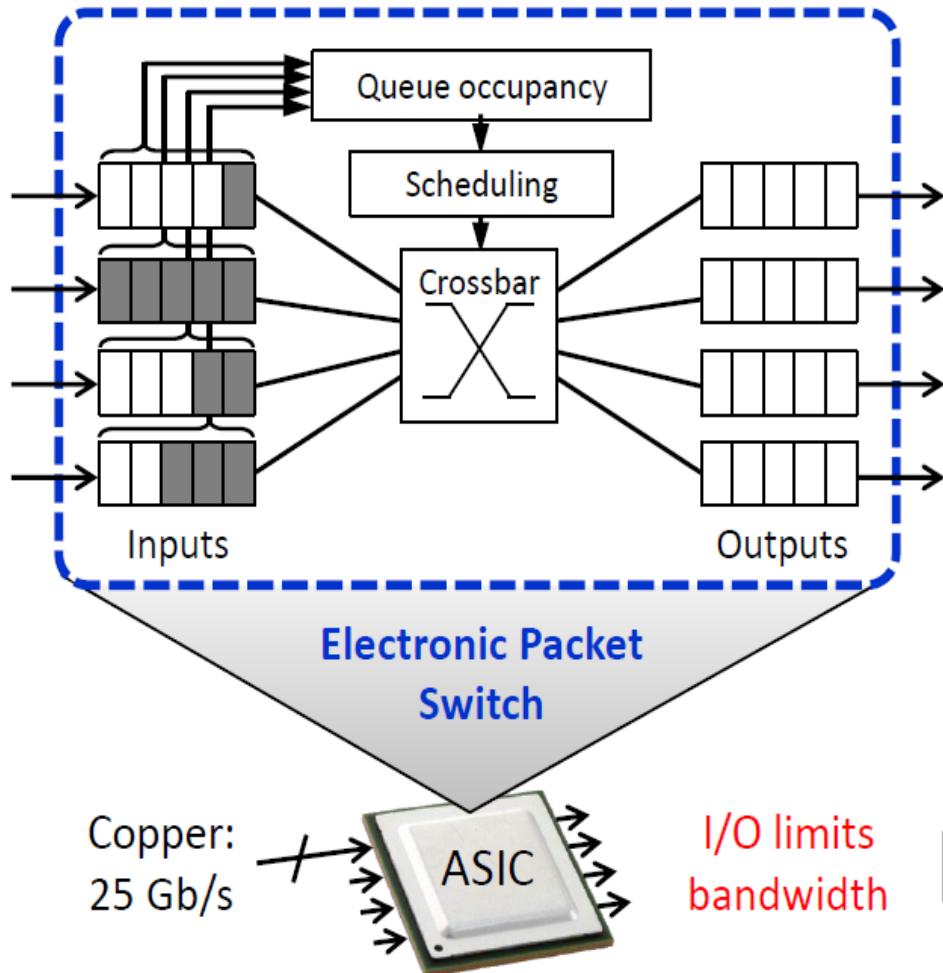
A. Singh et al., "Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network," SIGCOMM 2015.

# What about optic switches?

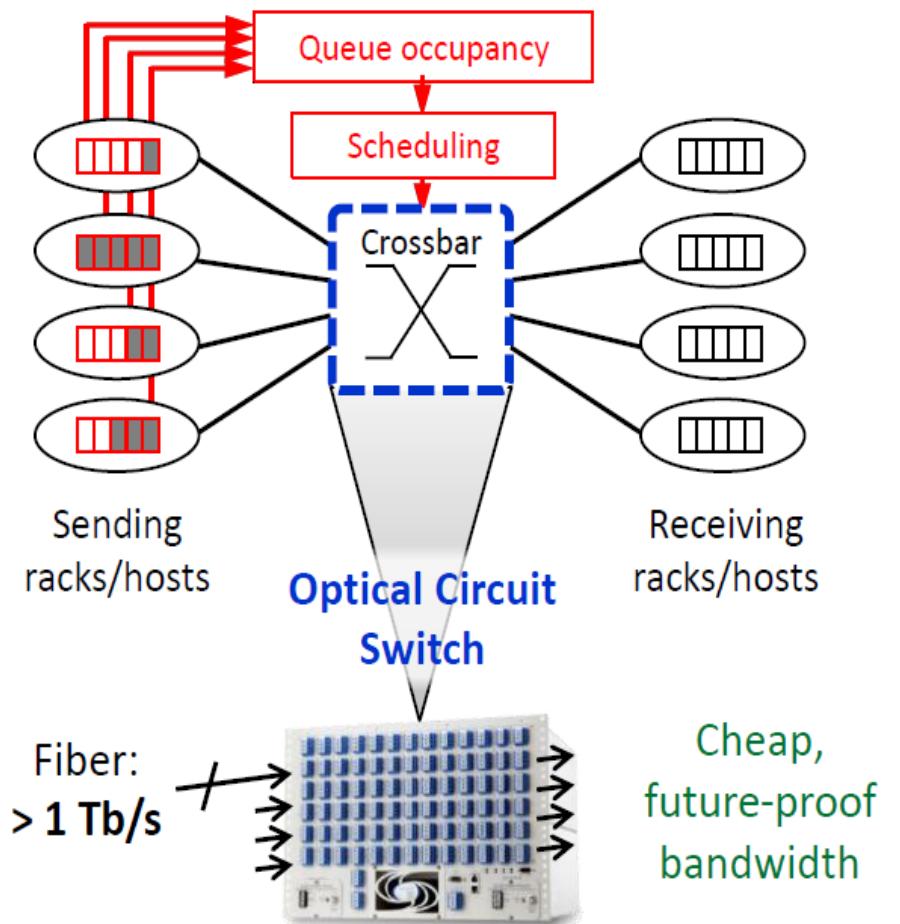


- ✓ **No O/E/O**
- ✓ **Practically unlimited Bandwidth per port** (limited by spectral efficiency)
- ✓ **Low Power:** Energy irrespective of throughput (power spent on steering pipes rather than processing/transmitting bits)
- ✓ Potentially **scales to thousands ports**

# Barrier no1: lack of buffering



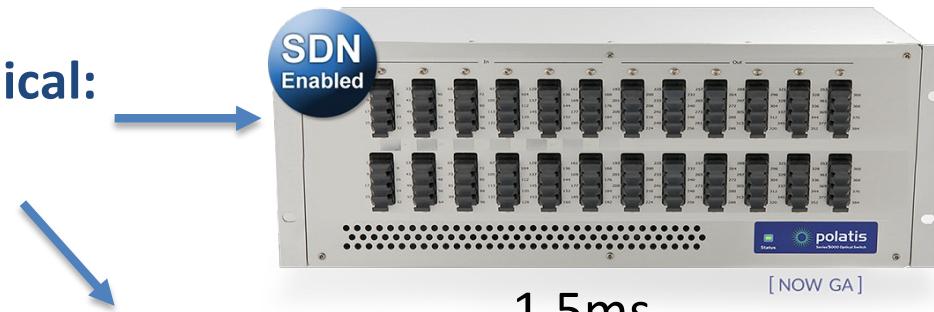
Data plane doesn't scale to entire datacenter!



## Barrier no2: speed

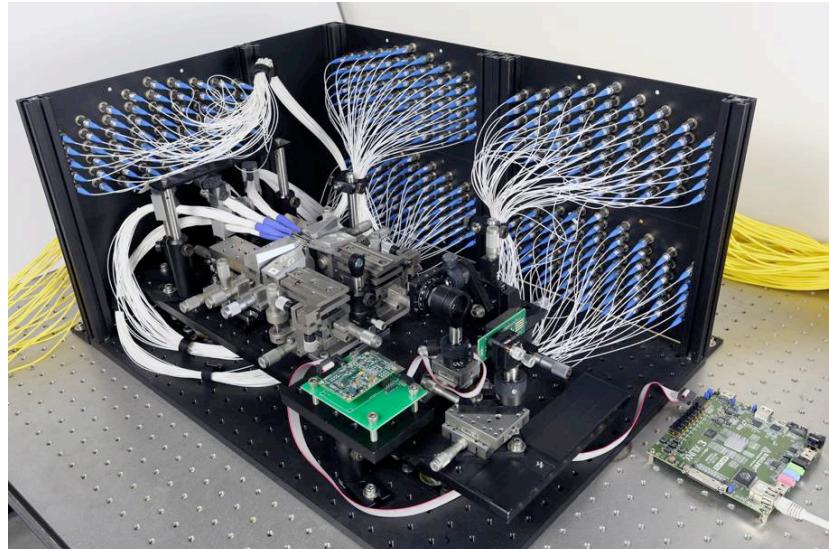
**Electromechanical:**

- High port
- low speed



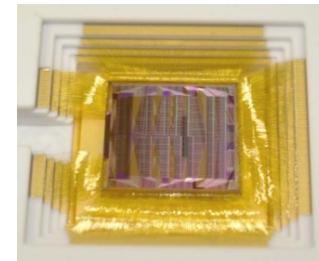
1.5ms

**UCSD/BERKELEY 2017**  
2D MEMS (3.2 ms)



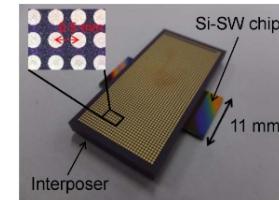
**Thermo-Optic**  
• **Huawei**

- 32x32, TO, 23 dB, 750  $\mu$ s



• **AIST**

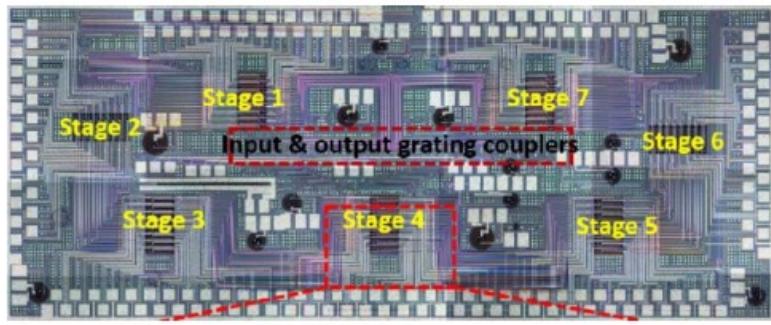
- 32x32, TO, 14.5 dB, 30  $\mu$ s



# Photonic integration: Fast Optical switches

2016

**16×16 EO Mach-Zehnder Switch**

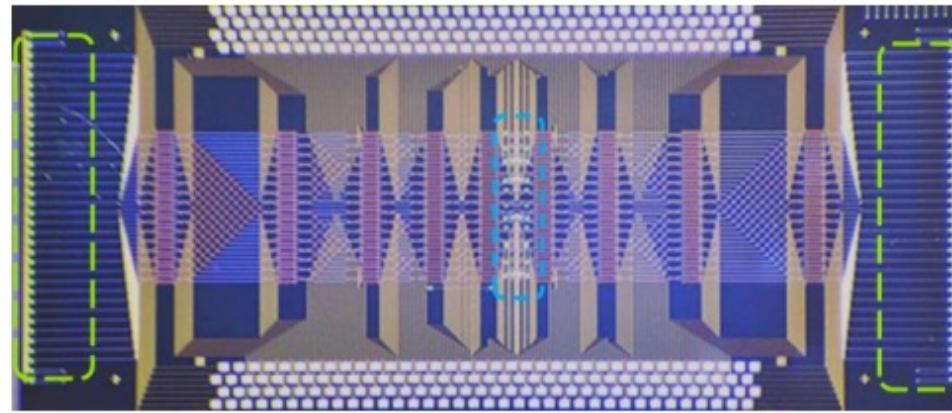


[Lu, Opt. Exp. 2016, 24 (9) 9295]

3.2/2.5ns

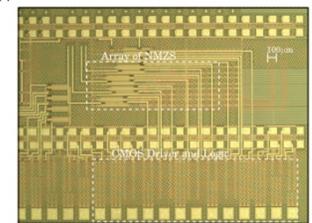
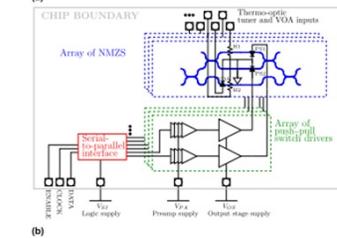
2017

**32×32 EO Mach-Zehnder Switch**



[Qiao, Nature Sci. Rep. 2017, 7 42306]

IBM 2016 nested MZ with electronic



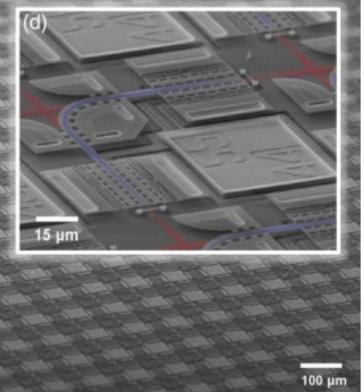
4ns

Benj. Lee

Berkeley 2016

**64×64 2D MEMS Switch**

0.91 μs



[Seok, Optica, vol. 3, no. 1, pp. 64]

Tae Joon Seok, Niels Quack, Sangyoon Han, Richard S. Muller, and Ming C. Wu

- Port to Port switch
  - 3D- MEMS switch
    - Coarse-grained switching technology (Fiber to fiber, port to port)
    - Up-to 320×320 switch fabric
    - Several tens ms switch time

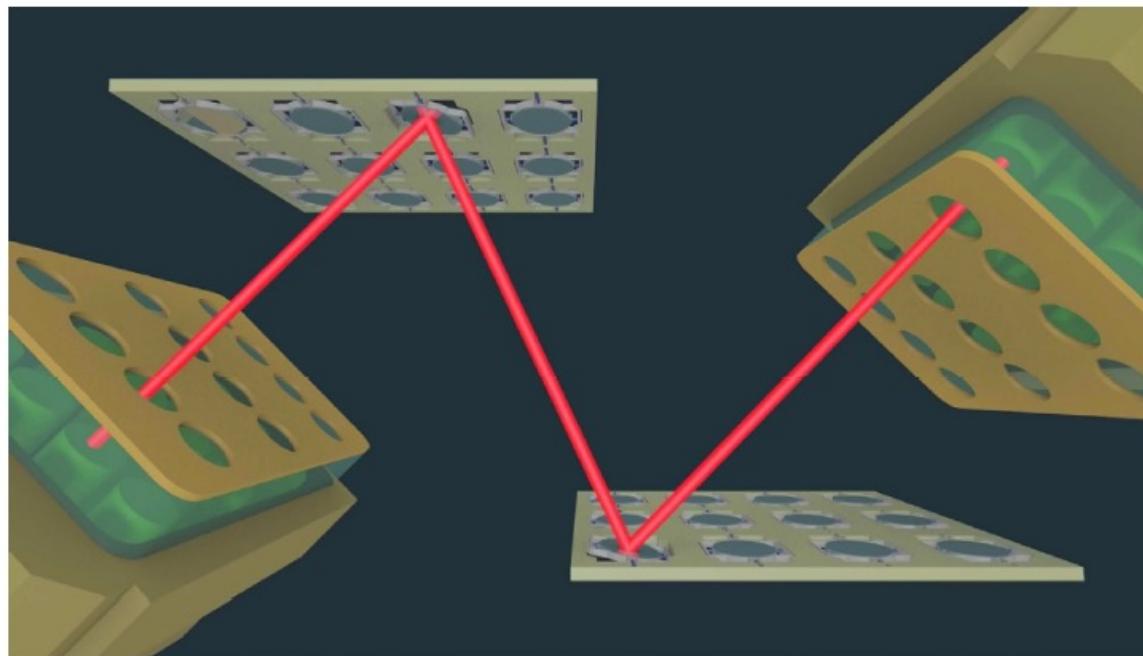
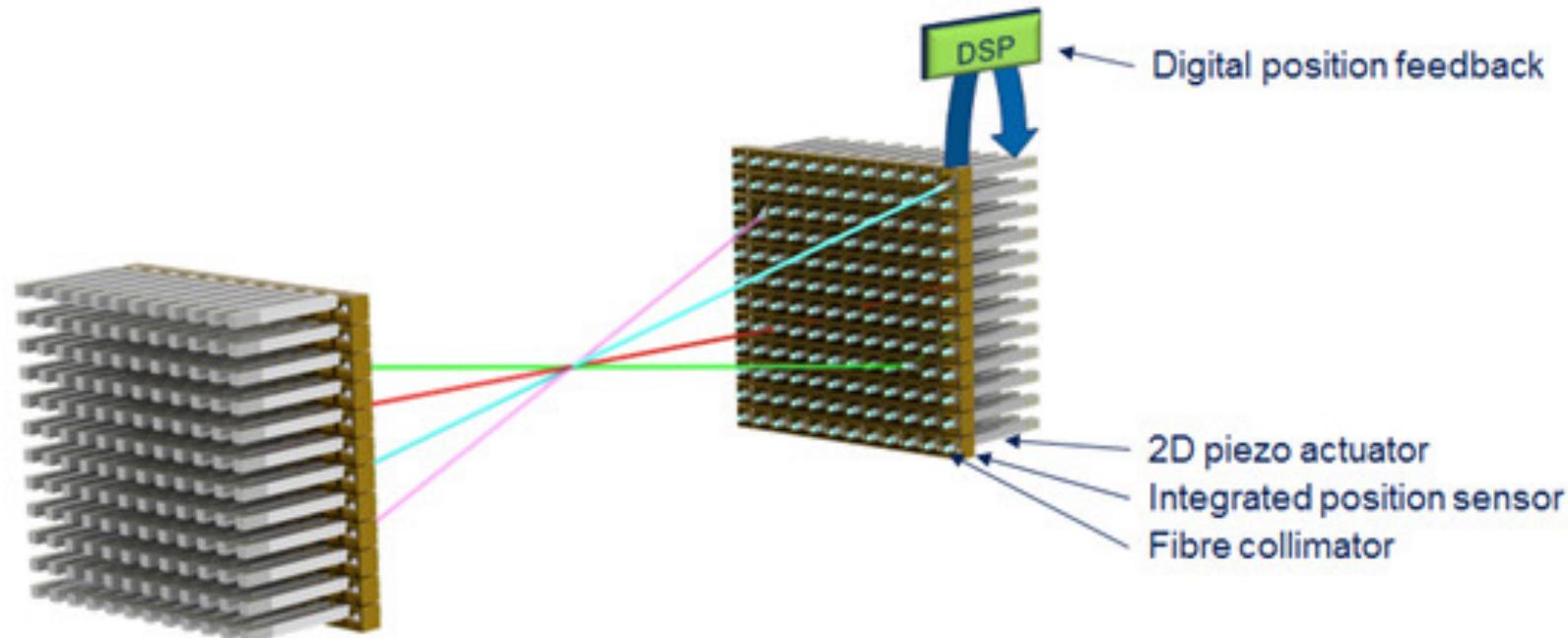


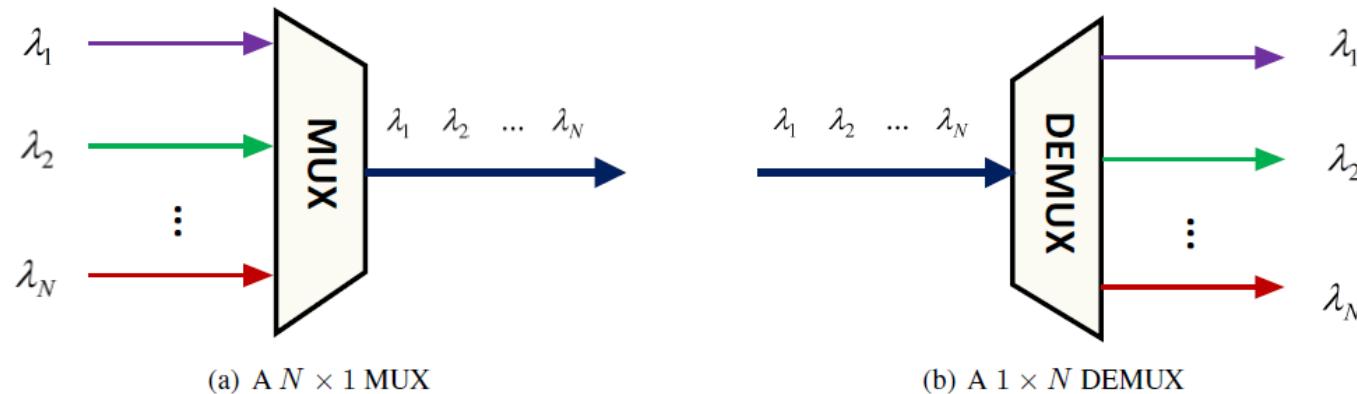
Figure 1.11: A three-dimensional optical MEMS switch.

- Port to Port switch
  - Beam steering technologies (Polatis)
    - Coarse-grained switching technology (Fiber to fiber, port to port)
    - Up-to  $320 \times 320$  switch fabric
    - Several tens ms switch time

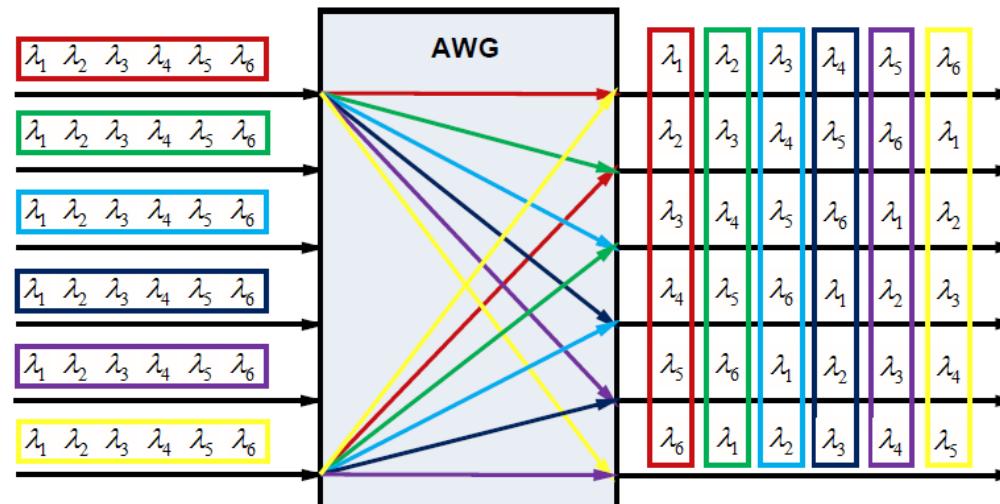


# Optical Switching technologies

- Array waveguide grating (AWG) based optical switch
  - MUX and DEMUX

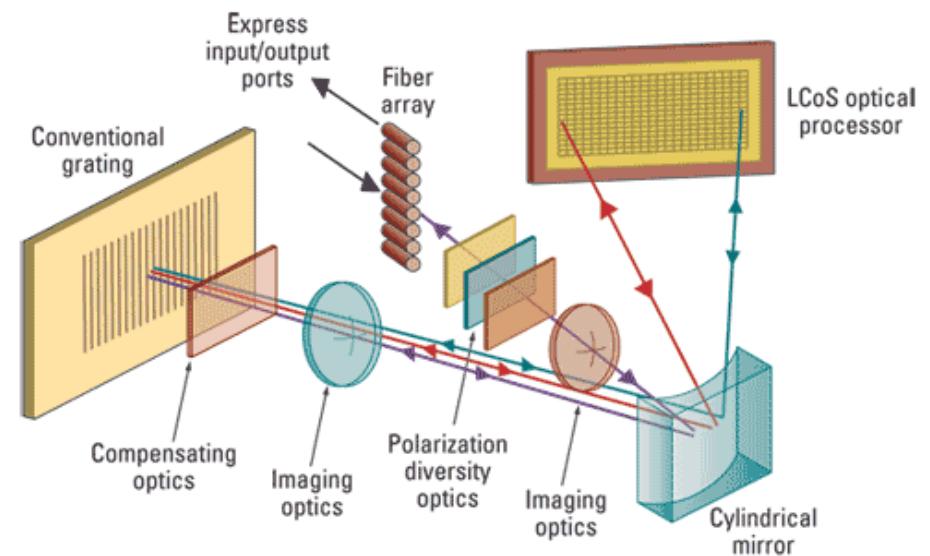
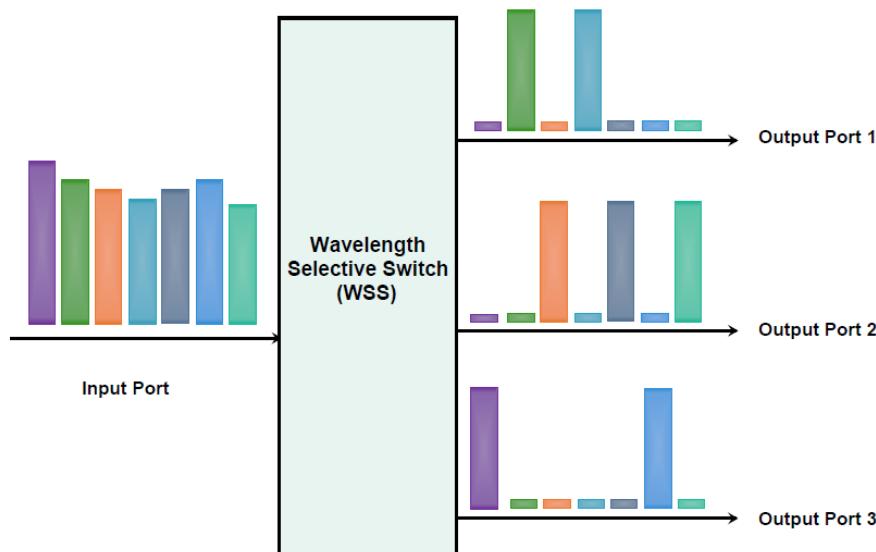


- AWGR (Arrayed waveguide grating router)



# Optical Switching technologies

- Wavelength selective switching ( $1 \times 3$  WSS)



# SOA-based time switch

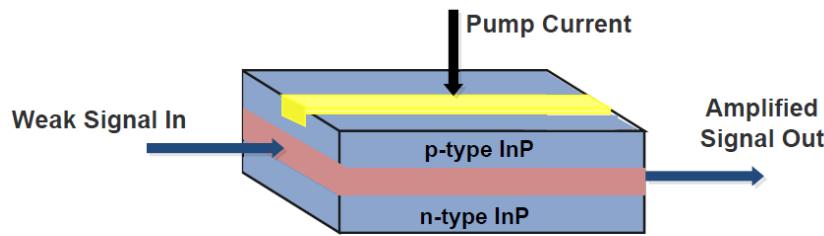
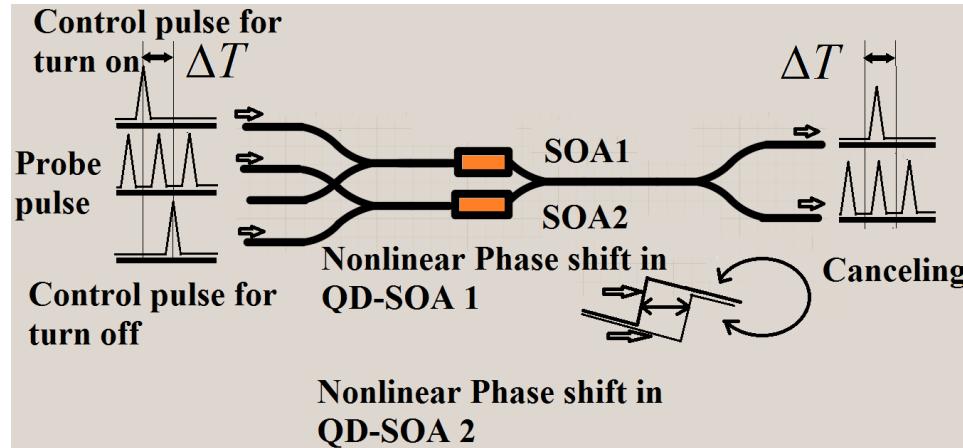
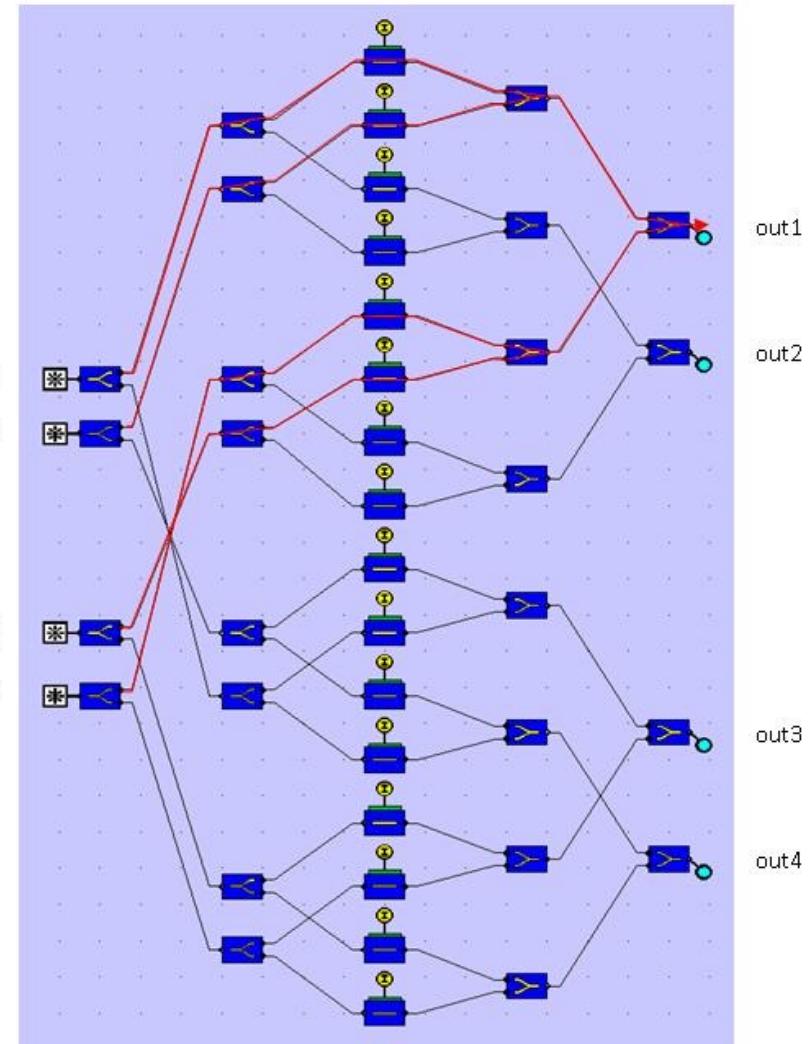


Figure 1.15: A Semiconductor Optical Amplifier (SOA).



All-optical switching

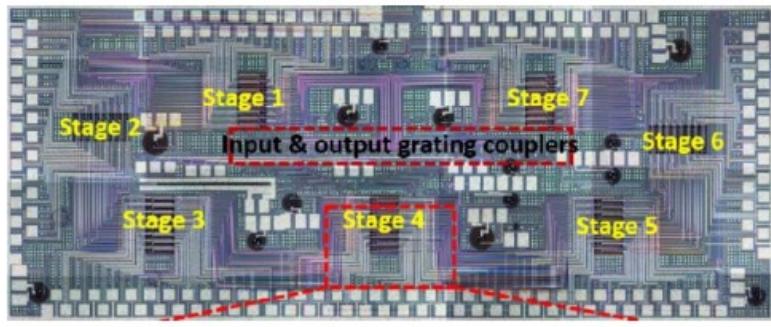


Optical Switch Circuit

# Photonic integration: Fast Optical switches

2016

**16×16 EO Mach-Zehnder Switch**

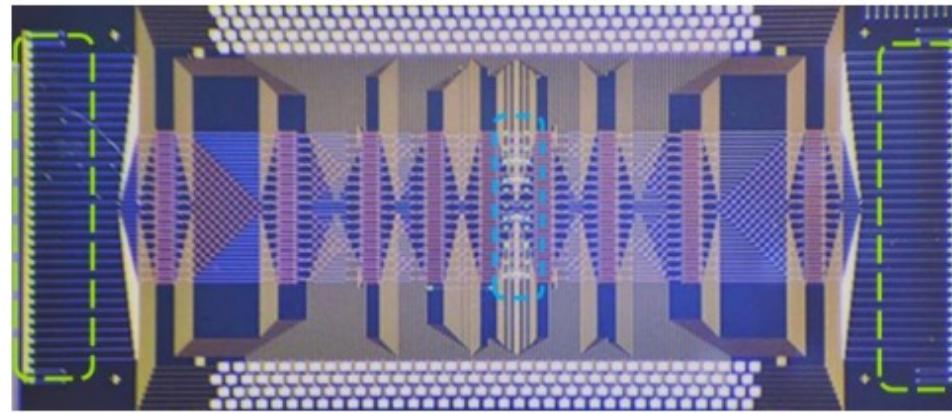


[Lu, Opt. Exp. 2016, 24 (9) 9295]

3.2/2.5ns

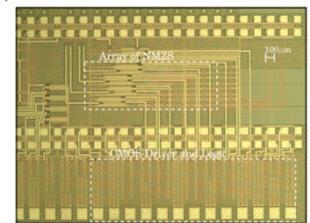
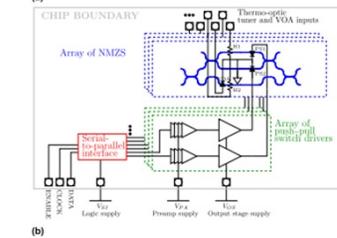
2017

**32×32 EO Mach-Zehnder Switch**



[Qiao, Nature Sci. Rep. 2017, 7 42306]

**IBM 2016 nested MZ with electronic**



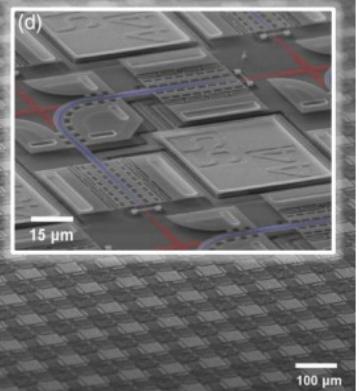
4ns

Benj. Lee

Berkeley 2016

**64×64 2D MEMS Switch**

0.91 μs

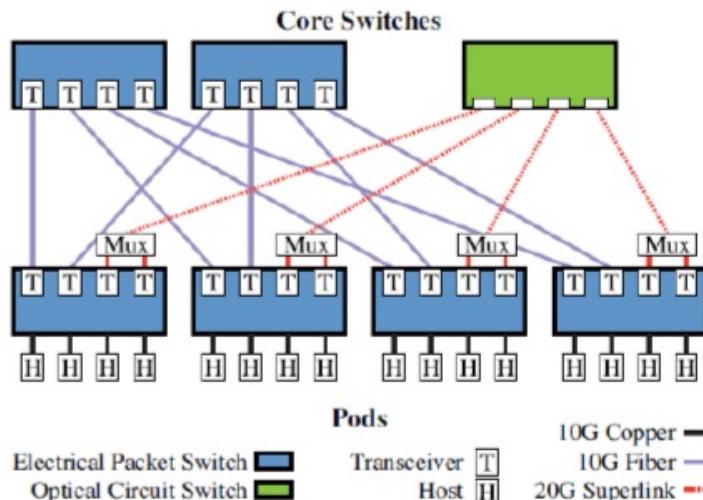


[Seok, Optica, vol. 3, no. 1, pp. 64]

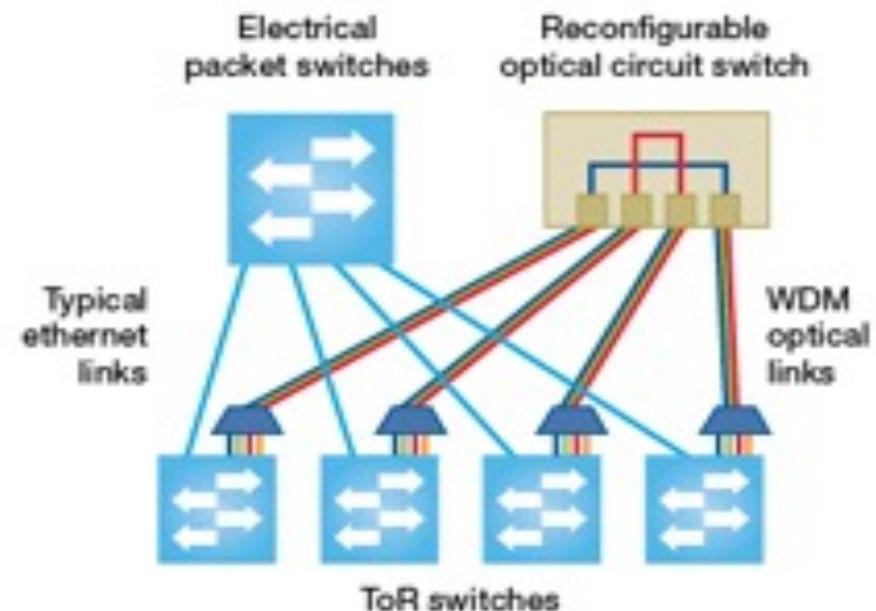
Tae Joon Seok, Niels Quack, Sangyoon Han, Richard S. Muller, and Ming C. Wu

- **Hybrid architectures (Fat tree architecture is enhanced using Optical Circuit Switching)**
  - Easily implemented (commodity switches)
  - Slow switching time (MEMs). Good only for bulky traffic that lasts long
  - Not scalable (constraint by Optical switch ports)

**Optical switches handle large elephant flows**

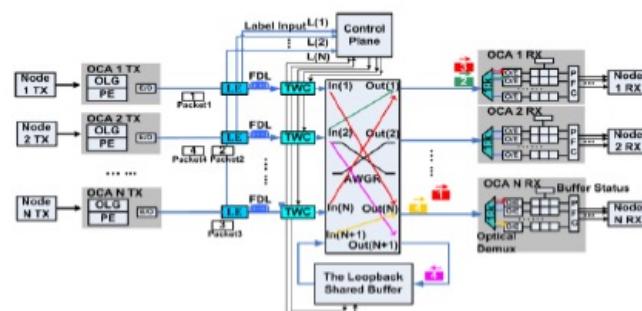


Farrington (SIGCOMM 2011)

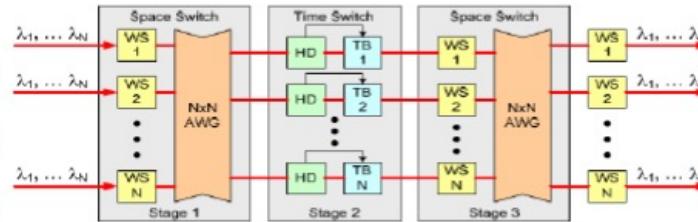


# High radix optical architectures

Optical switch architectures with high radices in order to lead to more flat DC architectures (less tiers) by replacing electrical switches in the upper tiers of the fat-tree

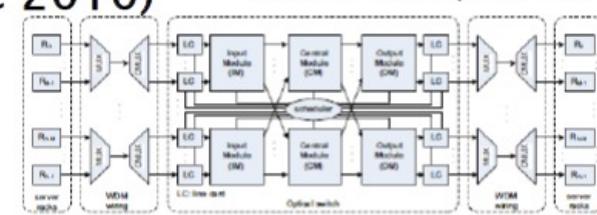


Ye et al (SANCS 2010)



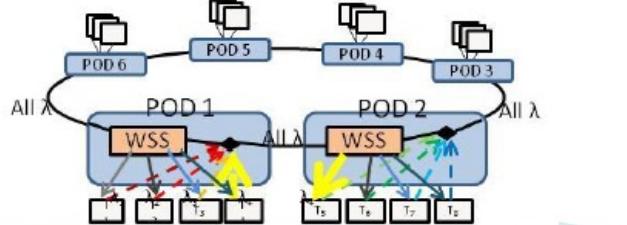
Gripp et al (OFCC 2010)

Xia et al (TR 2010)

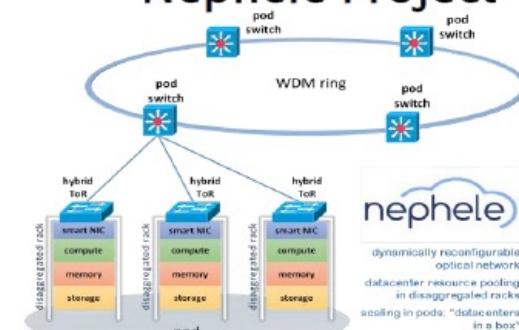


## ► Alternatives to fat trees

Farrington et al (OFCC 2013)

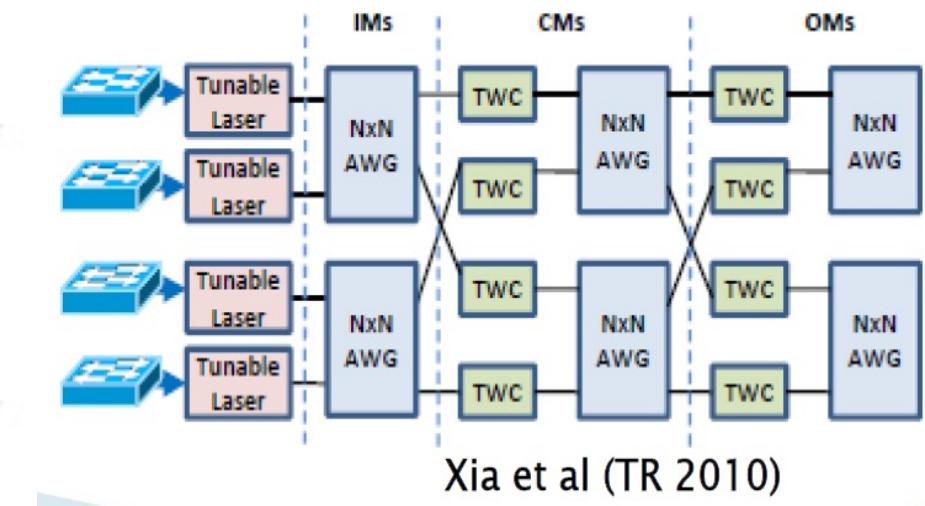
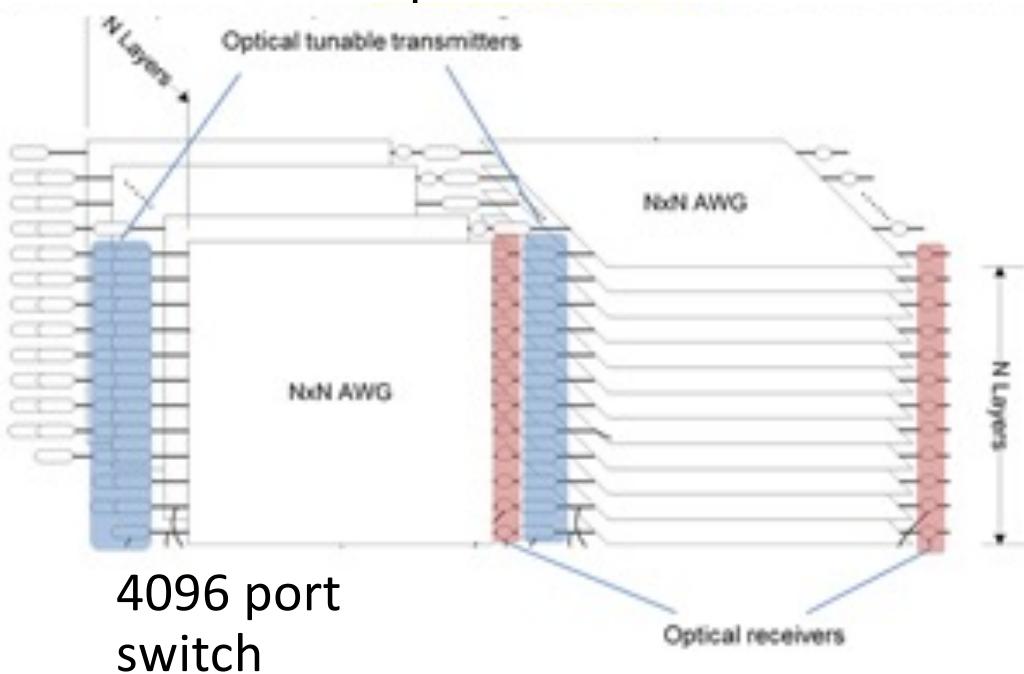


Nephele Project



# Petabit switch fabric

- Petabit switch fabric: three-stage Closnetwork and each stage consists of an array of AGWRs that are used for the passive routing of packets.
- In the first stage, the tunable lasers are used to route the packets through the AWGRs, while in the second and in the third stage TWC are used to convert the wavelength and route accordingly the packets to destination port.



# Summary

---

- Optical fibres as transmission medium
- Transmitter design for high-speed low-cost interconnection
- Possible opportunities for optical switching in DCN