# Opportunities and Challenges for Photonics in Next Generation Data Centers

## Clint Schow

Department of Electrical and Computer Engineering

The University of California Santa Barbara

schow@ece.ucsb.edu

# Outline

- Background: Interconnects and Technologies

- Optics In Today's Data Centers
  - Point-to-point interconnects
  - Multiple technologies for multiple purposes

- Opportunities
  - Photonic I/O
  - Photonic routing and switching

- Path Forward: Large-Scale Electronic/Photonic Integration

- Closing Thoughts

# Historically Two Fiber Optics Camps:
# Datacom and Telecom

**Telecom** (10's – 1000's of km)

- Expensive to install fiber over long distances
  - Single-mode fiber (SMF)
  - Wavelength Division Multiplexing (WDM)
- Cost of transceivers a secondary concern
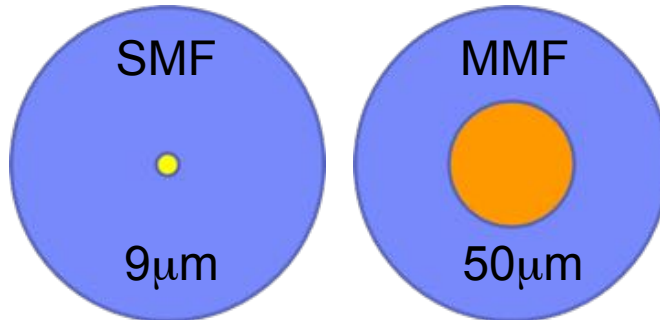- Performance is the primary objective

## Data Centers

**Datacom** (100's of meters)

- Cost of everything (transceivers, fibers, connectors) is the biggest factor
  - Multi-mode fiber (MMF)
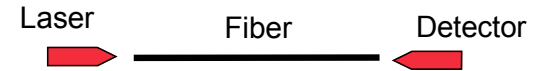  - Transceivers are commodities
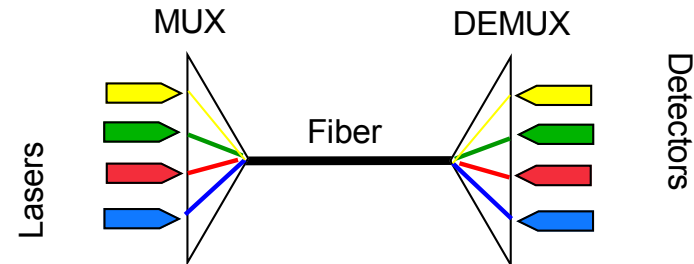
Relative core size:

MMF >30X SMF

**SMF** 9μm

**MMF** 50μm

**TDM = Time Division Multiplexing**
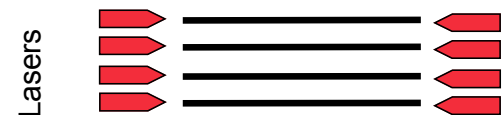Single optical channel, Electronic Mux/Demuxing

Laser — Fiber — Detector

**WDM = Wavelength Division Multiplexing**
Single optical channel data carried on separate λ's

MUX        DEMUX

Lasers — Fiber — Detectors

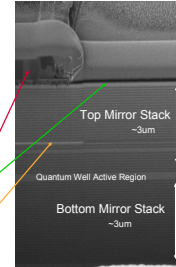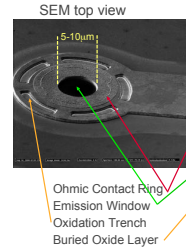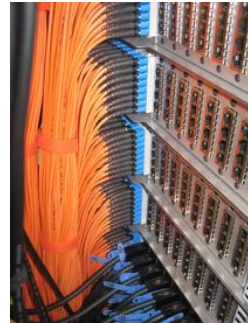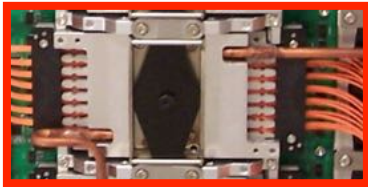**SDM = Space Division Multiplexing**
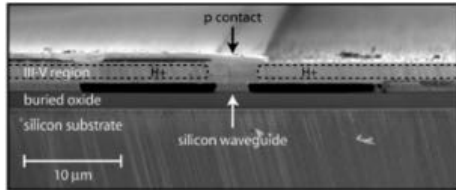Parallel fiber channels, No Mux/Demuxing

Lasers

**VCSELs** (Vertical Cavity Surface Emitting Lasers)
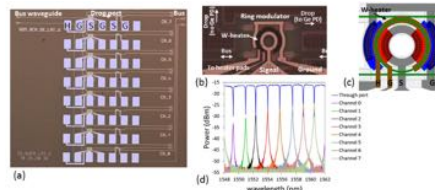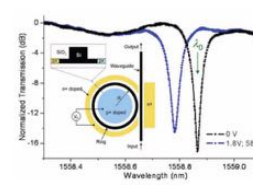Multi-mode fiber (MMF), Polymer Waveguides, Multicore fiber

**Si Photonics**
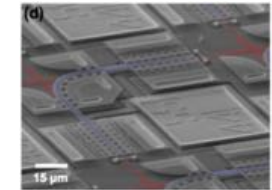Single-mode fiber (SMF), Wavelength Division Multiplexing (WDM)

**UCSB**  **IMEC/Ghent**  **Cornell/Columbia**  **UC Berkeley**

**Luxtera**  **Intel**  **Aurrion**  **IBM**  **Oracle**  **HP**

# The Old Days: 2012

Maintaining the HPC Performance Trend — IBM

PERFORMANCE DEVELOPMENT / PROJECTED

Total: #1 to #500
#1 machine
#500 machine

10x/4yrs = ~90% CAGR

10PF 2012    1EF 2020

Overly Optimistic

TOP500 NOVEMBER 2011
www.top500.org
Vendors System Share

- IBM 44.6%
- HP 28.2%
- Cray Inc.
- SGI
- Bull
- Appro
- Dell

**Performance enabled by increased parallelism:**

- Processor speed no longer primary driver
- Aggregation of massive numbers of multicore processors
- Challenging interconnect BW demands across system hierarchy
  - Intra-chip, inter-chip, on-board, intra-rack, between racks
  - Communication bottlenecks moving closer to processors
  - Optics displacing copper at ever shorter distance scales

Approved for Public Release, Distribution Unlimited.

IBM, presented at Optical Interconnects Conference, 2012
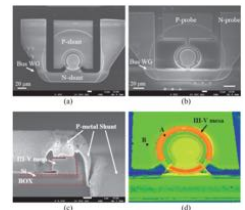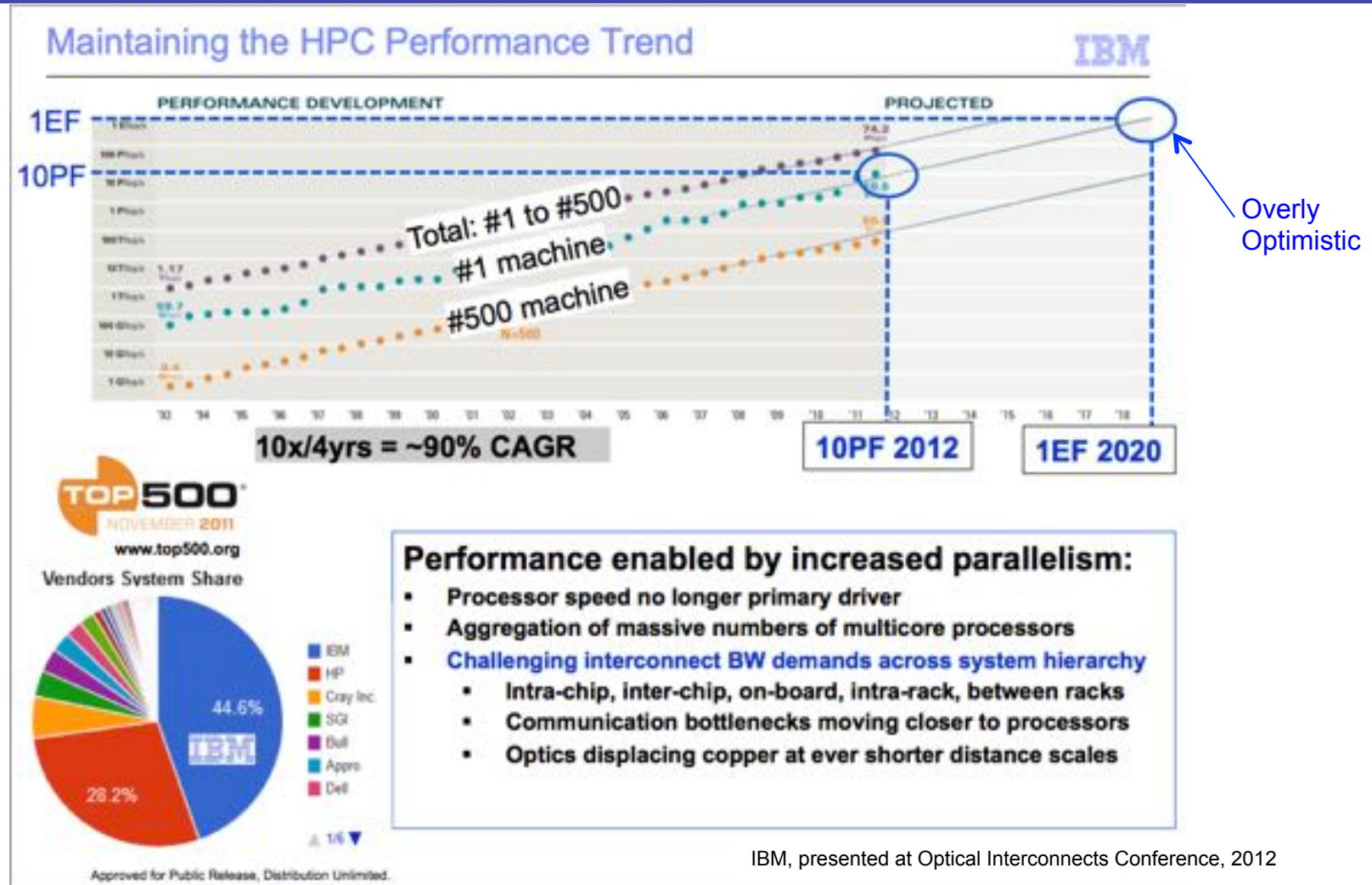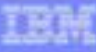
Trickle-Down: HPC drives development of highest performance components that are later picked up by commercial servers

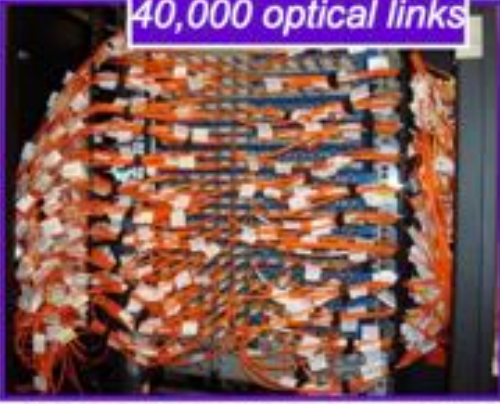# The Promise of A Machine with 2,000,000 VCSELs



Computercom copper displacement

**Current System Implementations** — **IBM Blue Waters (2011)**
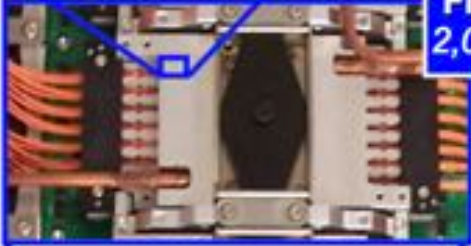
**IBM Roadrunner (2008)**

**Fiber to the Rack**
40,000 optical links

Active optical cables plugged into back of switch rack
5 Gb/s

2008→2011: 100X increase

microPOD™ parallel optical TX/RX    10 Gb/s
[M. Fields, Avago, OFC 2010, OTuP1]

**Fiber to the Module**
2,000,000 optical links

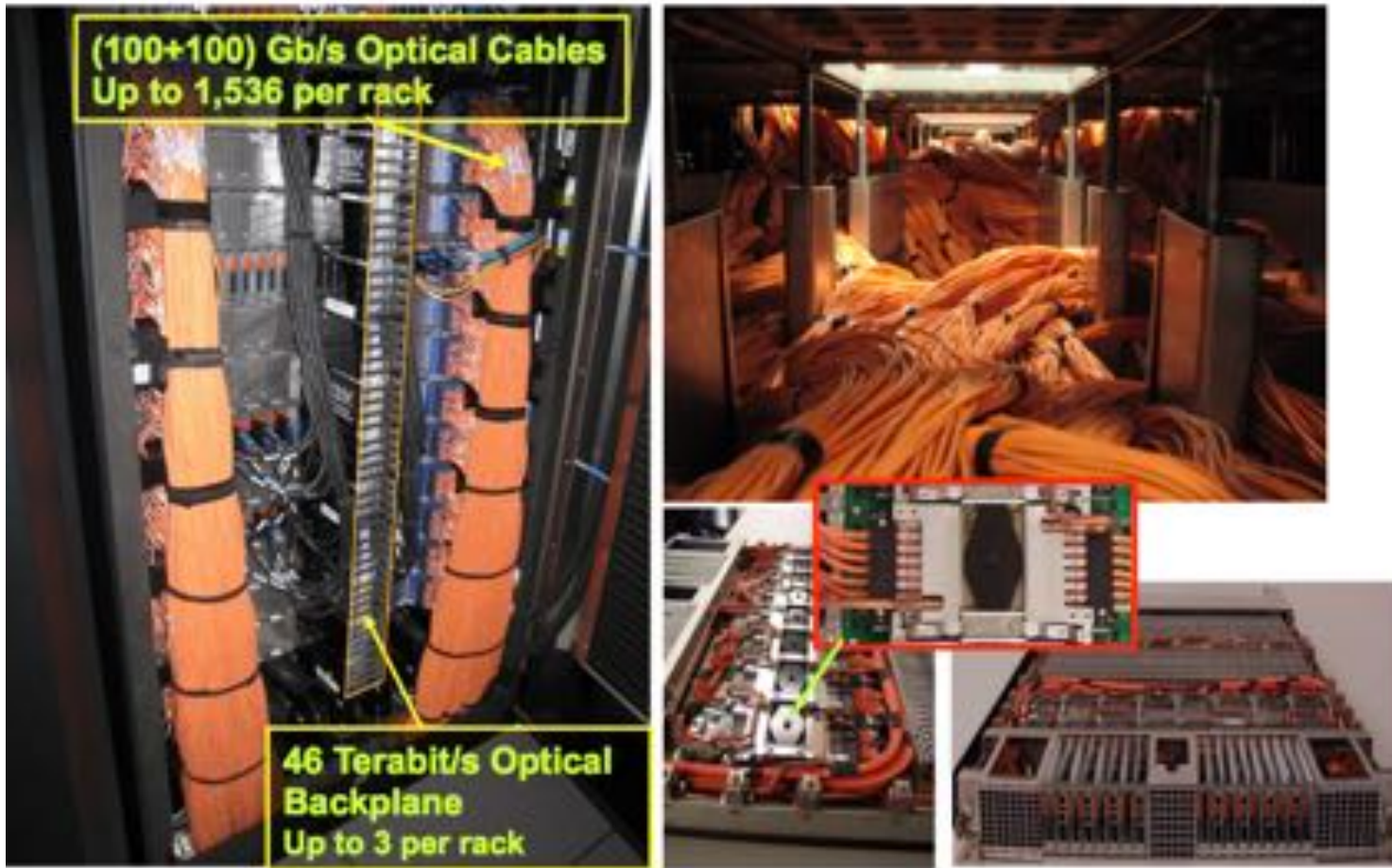Hub/switch module, with IC and 56 microPODs

Node Drawer

[A. Benner, IBM, OFC 2010, OTuH1]

© 2010 IBM Corporation

IBM, public presentation 2010

- A. Benner, "Optical Interconnect Opportunities in Supercomputers and High End Computing," OFC Tutorial 2012.

IBM Power 775: Pushing the Limits — IBM

(100+100) Gb/s Optical Cables Up to 1,536 per rack

46 Terabit/s Optical Backplane Up to 3 per rack

IBM, public presentation 2010

**Need more BW/fiber: WDM, multicore fiber**

- A. Benner, "Optical Interconnect Opportunities in Supercomputers and High End Computing," OFC Tutorial 2012.

# Optics in HPC: IBM Sequoia

**96 IBM Blue Gene/Q Racks**
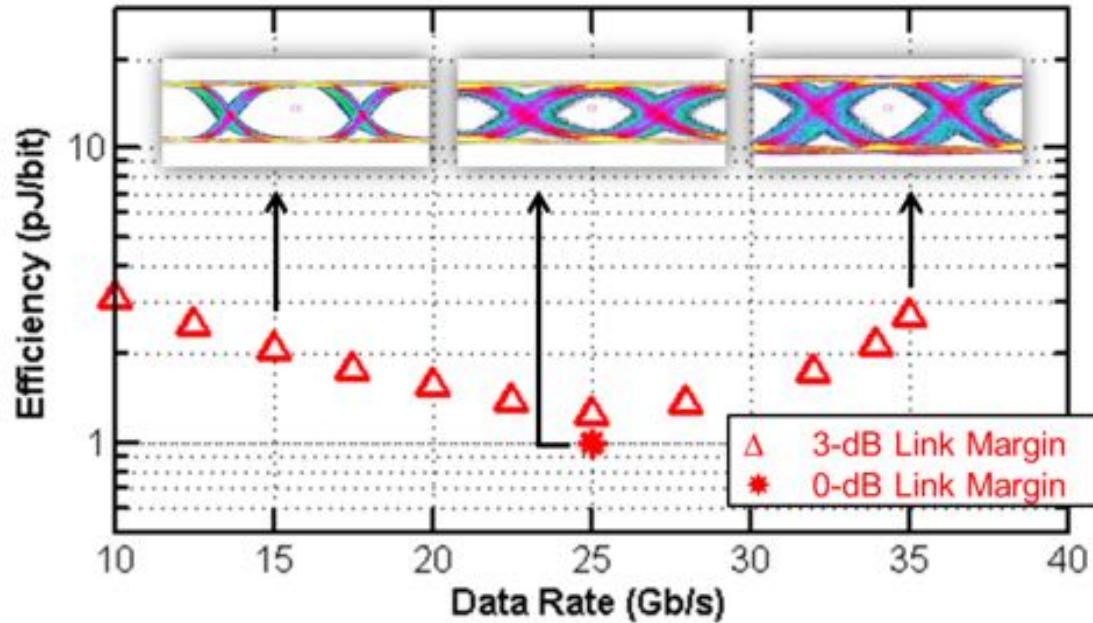**20.013 Pflops Peak … 1.572M Compute Cores … ~8MW … 2026 Mflops/Watt**

620,000 VCSEL links

IBM, presented at IPC, 2013

- HPC requires technologies optimized for short reach ~50m
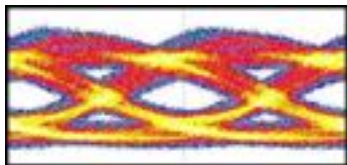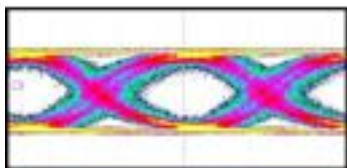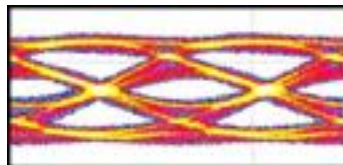
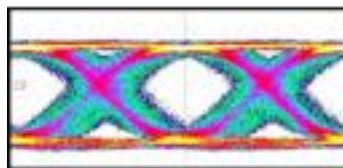# VCSELs: Efficient and Fast

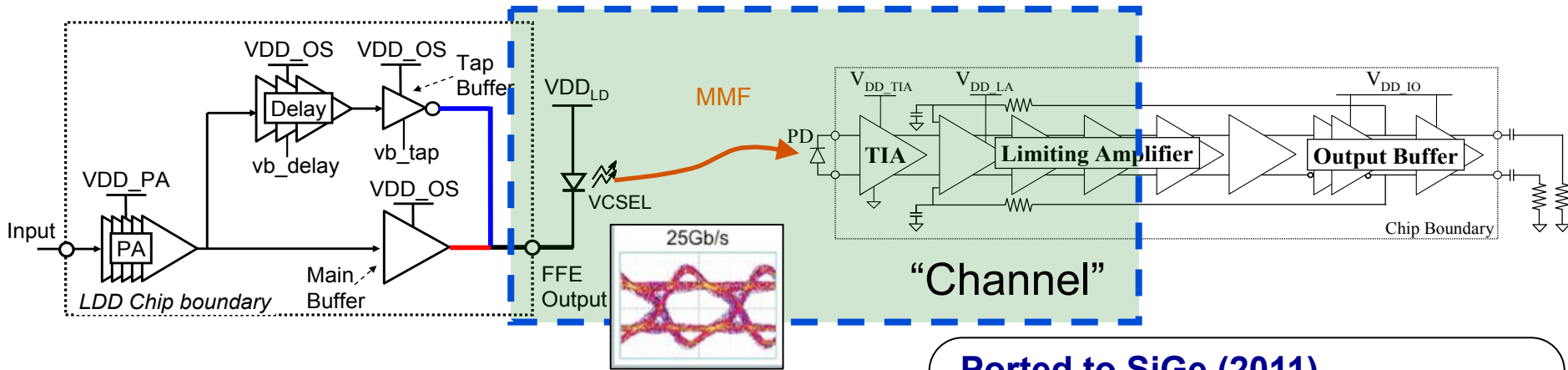## 32-nm CMOS-Driven Link



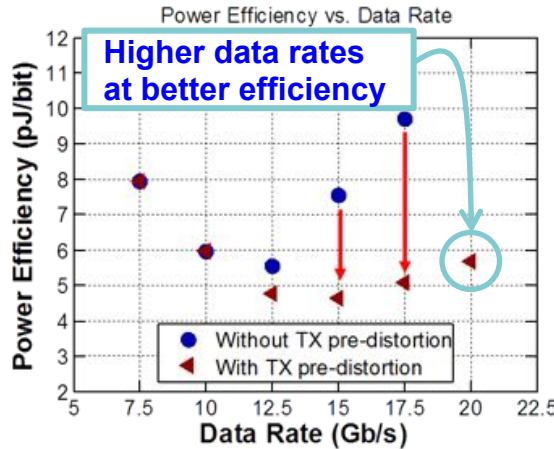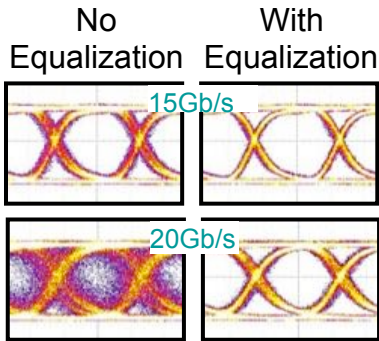25 Gb/s, 24mW          35 Gb/s, 95mW

TX out

RX out

**Wall-plug efficiency:
1pJ/bit at 25 Gb/s**

- J. E. Proesel *et al.*, "35-Gb/s VCSEL-based optical link using 32-nm SOI CMOS circuits," *OFC 2013*, Paper OM2H2, Mar. 2013.

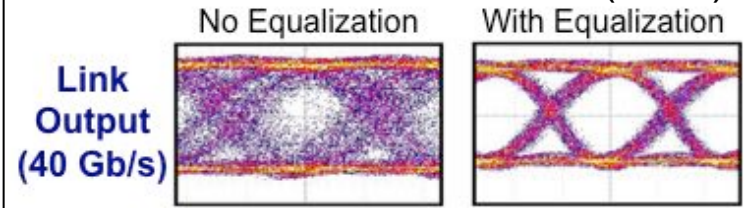# Rethinking Equalization: Optimizing the Performance of Complete Links



25Gb/s

"Channel"

## First Implemented in CMOS (2011)

No Equalization    With Equalization

15Gb/s

20Gb/s

Higher data rates at better efficiency

Power Efficiency vs. Data Rate

- Without TX pre-distortion
- With TX pre-distortion

## Ported to SiGe (2011)

First 40 Gb/s VCSEL Links (2012)

No Equalization    With Equalization

Link Output (40 Gb/s)

71 Gb/s (2015)

Link Output (71 Gb/s)

## Next opportunity: Si photonic WDM links

- A.V. Rylyakov *et al.*, "Transmitter Pre-Distortion for Simultaneous Improvements in Bit-Rate, Sensitivity, Jitter, and Power Efficiency in 20 Gb/s CMOS-driven VCSEL Links, *JLT* 2012.
- A. V. Rylyakov *et al.*, "A 40-Gb/s, 850-nm VCSEL-based full optical link," *OFC* 2012.
- D. Kuchta *et al.*, "A 71-Gb/s NRZ Modulated 850-nm VCSEL-Based Optical Link," *PTL,* March 2015.

# And Then The Cloud Rolled In
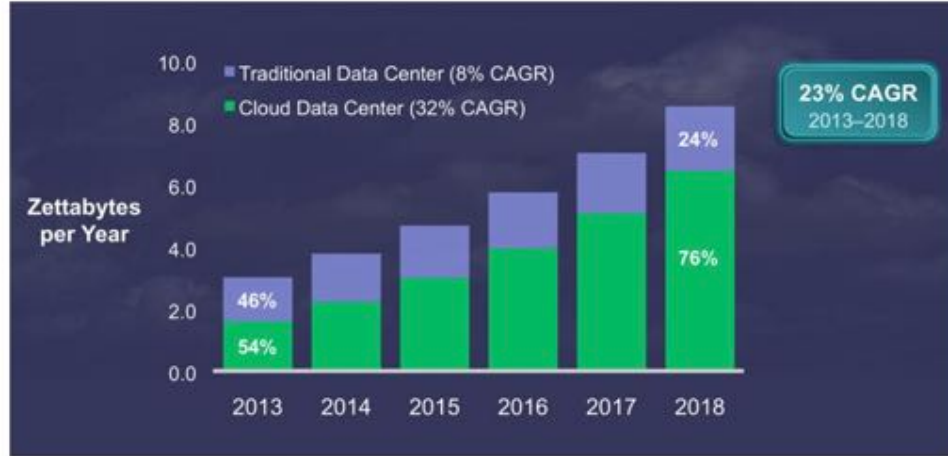


View of Goleta Beach from UCSB, source www.theinertia.com
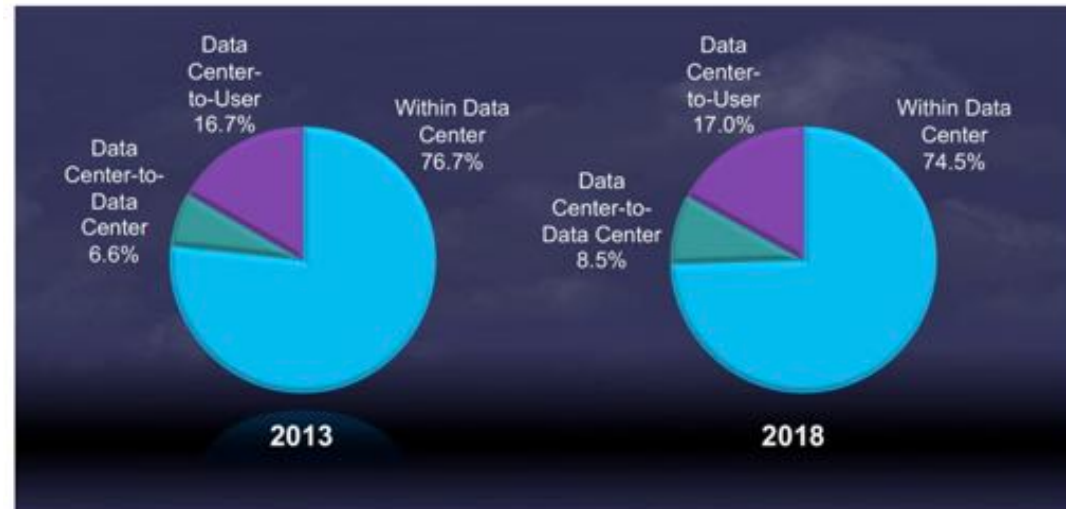
# Growth in Cloud Data Centers

**Figure 3.** Total Data Center Traffic Growth



Source: Cisco Global Cloud Index, 2013–2018
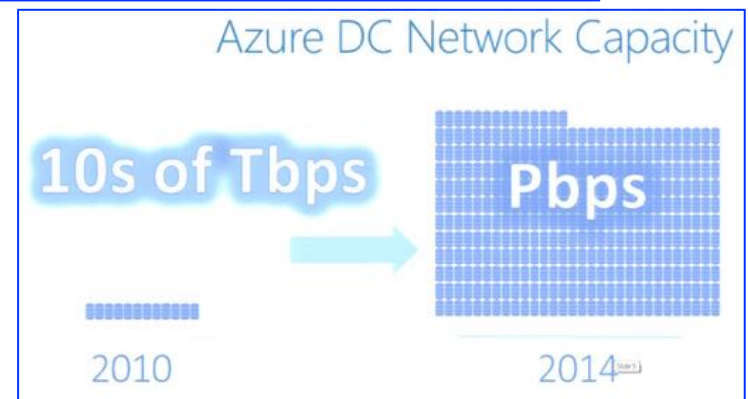
Cloud is the growth market

**Figure 2.** Global Data Center Traffic by Destination



Source: Cisco Global Cloud Index, 2013–2018

75% of traffic is within the data center

# Huge Growth in the Cloud (Microsoft)

UCSB

**Windows Azure**

Global datacenters ●

Over **25** datacenters    Over **1** billion customers,    **20 million** businesses    **76** markets worldwide

Compute Instances

**100K** → **Millions**

2010                2014

Azure DC Network Capacity

**10s of Tbps** → **Pbps**

2010                2014

Source: D. Maltz, Microsoft, OFC 2014

# Traditional Data Center Network Hierarchy

Traditional hierarchical networks grew out of campus/WAN installations

Good for North-South traffic

Core (Spine)

Aggregation (Leaf)

Edge (TOR=Top-of-Rack)

Server

**Bad for East-West traffic
(75% of data center traffic)**

See L. A. Barroso and U. Hölzle, "*The Datacenter as a Computer—An Introduction to the Design of Warehouse-Scale Machines,*"

# Data Center Hardware
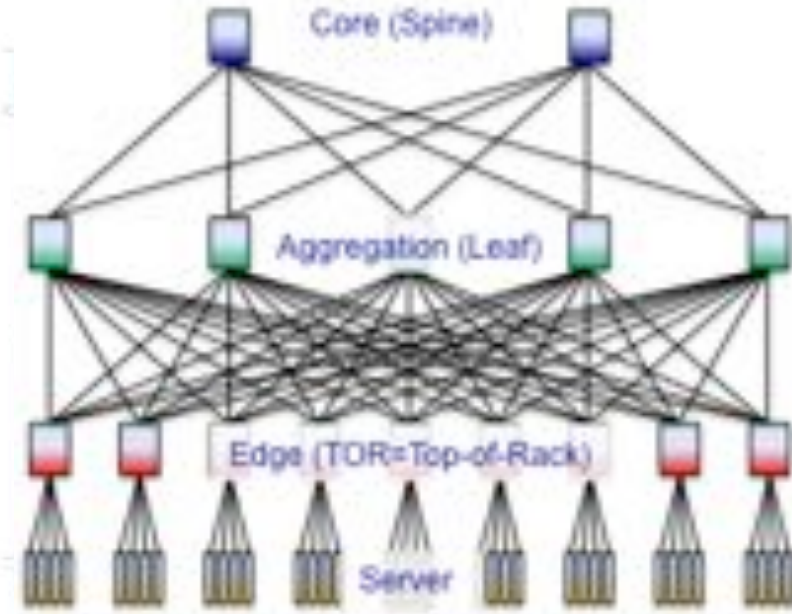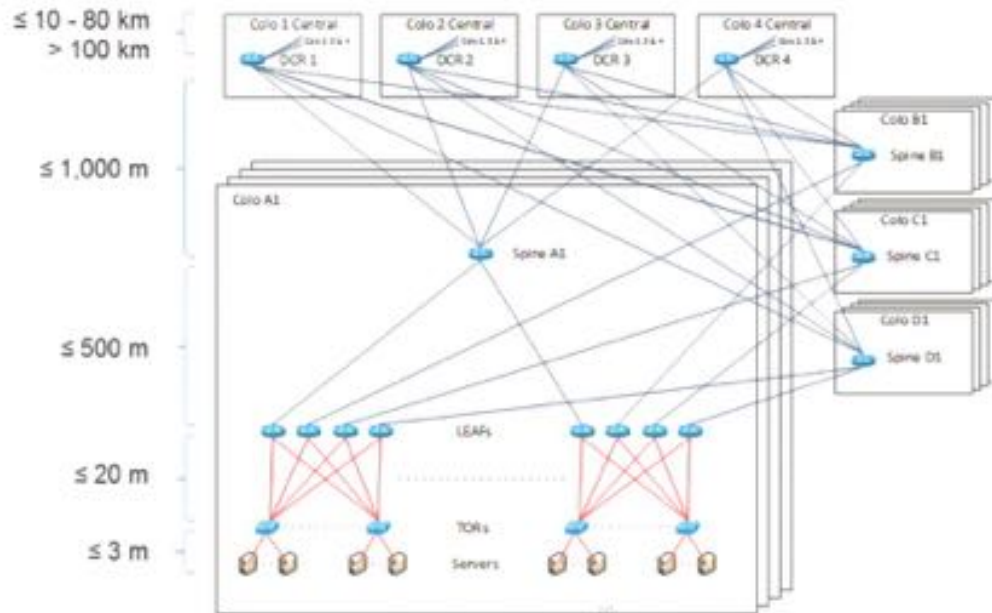
### Racks



Top of Rack
(ToR) switch

Servers

### Lots of Racks



Photos of Facebook data centers found on Google images

# Cloud Data Centers (Microsoft)

Cloud Data Center & Campus Interconnections



| A End | Z End | Link Quanity | Link Length | Type of interconnection |
|---|---|---|---|---|
| Server | TOR | 10,000s | .5-3m | TwinAx |
| TOR | LEAF | 1,000s | 1-20m | AOC |
| LEAF | SPINE - local | 100s | 20-300m | SM fiber |
| LEAF | SPINE - inter building | 1,000s | 100-400m | SM fiber |
| SPINE | DCR | 100s | 100-1,000m | SM fiber |
| INTRA METRO | | 100s | 1,000m+ | SM fiber |

WDM identified as path to lower cost if transceivers are cheap

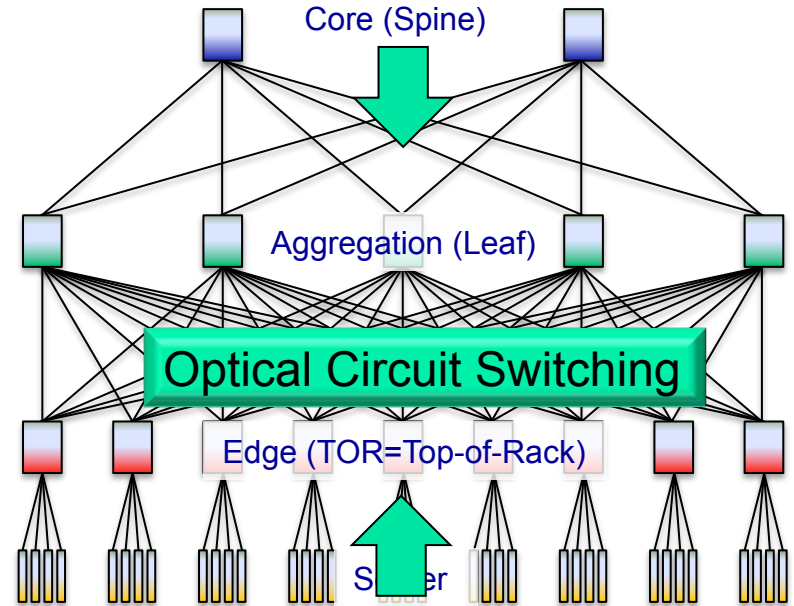Source: D. Maltz, Microsoft, OFC 2014

# Largest Data Centers, 2015

**1. Range International Information Group** (Langfang, China)
Area: 6,300,000 Sq. Ft.

**2. Switch SuperNAP** (Nevada, USA)
Area: 3,500,000 million Sq. Ft.

**3. DuPont Fabros Technology** (Virginia, USA)
Area: 1,600,000 million Sq. Ft.

**4. Utah Data Centre** (Utah, USA)
Area: 1,500,000 million Sq. Ft.

**5. Microsoft Data Centre** (Iowa, USA)
Area: 1,200,000 Sq. Ft.

**6. Lakeside Technology Centre** (Chicago, USA)
Area: 1,100,000 Sq. Ft.

**7. Tulip Data Centre** (Bangalore, India)
Area: 1,000,000 Sq. Ft.

**8. QTS Metro Data Centre** (Atlanta, USA)
Area: 990,000 Sq. Ft.

**9. Next Generation Data Europe** (Wales, UK)
Area: 750,000 Sq. Ft.

**10. NAP of the Americas** (Miami, USA)
Area: 750,000 Sq. Ft.

Switch SuperNAP

Utah Data Centre

Huge facilities need lots of longer-distance links:
2km becoming the magic number for data centers

| Number of ports per switch | #Nodes Connected Two Level | #Nodes Connected Three Level | #Nodes Connected Four Level | #Nodes Connected Five Level |
|---|---|---|---|---|
| 8 | 32 | 128 | 512 | 2048 |
| 16 | 128 | 1024 | 8192 | 65536 |
| 32 | 512 | 8192 | 131072 | 2.10E+06 |
| 64 | 2048 | 65536 | 2.10E+06 | 6.71E+07 |
| 128 | 8192 | 524288 | 3.36E+07 | 2.15E+09 |
| 256 | 32768 | 4.19E+06 | 5.37E+08 | 6.87E+10 |
| 512 | 131072 | 3.36E+07 | 8.59E+09 | 2.20E+12 |
| 1024 | 524288 | 2.68E+08 | 1.37E+11 | 7.04E+13 |



Core (Spine)

Aggregation (Leaf)

Optical Circuit Switching

Edge (TOR=Top-of-Rack)

Server

## 1) Flatten the Network:

Higher radix (larger port count) switches

Electrical or Optical Cores

## 2) Change the Network:

Photonic switching

# Photonic I/O:
# Optics to the Chip

# Shared Visions:
# Photonically Connected Chips

**UCSB**

## Today: Electrical Chip Packaging

IC packages with course BGA/LGA electrical connectors
- Poor scalability, signal integrity
- Reduced system performance
- Reduced system efficiency

## Future: Photonic Packages

Interposer
(if needed)

High-density, ultra-low
power photonics

IC's: one or more
can be 3-D stacks

Multi-chip package
(BGA, LGA)

**Photonic integration must provide more I/O bandwidth at better efficiency**

# Limitations of Electrical Switches

## Mellanox

| Ordering Part Number | Description | Typical Power |
|---|---|---|
| MT52236A0-FDCR-E | Switch-IB, 36 Port EDR InfiniBand Switch IC | 83W |

| Ordering Part Number | Description | Typical Power |
|---|---|---|
| MT52132A0-FCCR-C | Spectrum, 32 Port Ethernet 100GbE Switch IC (RoHS R6) | 135W |

Source: mellanox.com

## Broadcom

### Tomahawk
128 x 25Gb/s SerDes = 32 100GbE ports
7 Billion transistors

Source: broadcom.com

**Two primary limitations**
- **Power (mostly electrical I/O)**
- **Density**

## Challenge: Limited chip radix

- Do we have enough package area for all the SERDES needed?
- 2015 – 150 SERDES (3000mm²)
- 2017 – 200/250 SERDES (5000mm²)
- 2020 – end of the road?

One SerDes

Courtesy of M. Laor, Compass Networks

# Higher BW Density:
# Demands Integration at (First-Level) Chip Package

## Cross-sectional view of chip module on board

Chip
C4s

Chip
Carrier

LGA
or BGA
connector

PCB

Courtesy of M. Ritter, IBM

### Supported bandwidth
(50x50mm module, 20Gb/s bitrate)

C4s **160Tb/s**

Chip carrier **114Tb/s**

Connector **12.6Tb/s**

PCB **70 Tb/s**

- Electrical Packaging is mature
  - # Signal Pins & BW/pin not increasing at rate of silicon
  - All-electrical packages will not have enough pins or per-pin BW
- Chokepoint is at the module-to-circuit board connection

• D. Kam et al., "Is 25 Gb/s On-Board Signaling Viable?, *IEEE Adv. Packag.*, 2009

## Electrical Link Example: Backplane



**INPUT**

Packaged SerDes
Backplane trace
Line card trace
Edge connector
Via stub

**OUTPUT**

| Pkg | | Line card trace | Edge connector | Backplane 16" trace | Edge connector | Line card trace | | Pkg |
|---|---|---|---|---|---|---|---|---|
| Tx IC | | | | The Channel | | | | Rx IC |

> 20dB loss @ 5GHz

*Tyco 16" Backplane Channel Response*

- Backplane: even higher loss than most chip-to chip links (extra connectors, longer distance). Very hard to scale to higher data rates.

- Why keep pushing it then?

- Ease of use: plug in a new card into the existing legacy backplane and get a speed boost.

- Could be a very cost-effective solution, but power dissipation is an issue, especially at higher data rates.

Courtesy A. Rylyakov, OFC Short Course #357

# Integration to Maximize BW and Efficiency

short trace on carrier (mm)

connector

long trace on card (many cm)

IC

Chip Carrier

TX/RX

Host Printed Circuit Board (PCB)

**Move from bulky optics located far away**

**To highly-integrated optics close to the logic chips**

*Integrating photonics into the most expensive, constrained, and challenging environment in the system*

- **Cost**
- **Reliability**
- **Thermal**
- **Power delivery**

short trace on carrier

IC

TX/RX

Chip Carrier

Host PCB

**UCSB**



# High-Speed Optical Interconnect Evolution II

| CONTEMPORARY – Today | EMERGING – 2014/15 | STRATEGIC DIRECTION – 2018+ |
| --- | --- | --- |
| • Traditional **MSA compliant pluggable modules and AOCs** on card edge<br>• Considerable **SI issues** (electrical connectors, long traces on host PCBA) require re-timers.<br>• Front panel interconnect **density limited by module size** (physical implementation + module power dissipation) | • **Embedded optical transceivers** located closer around ASIC<br>• Shorter traces on PCB **alleviate SI issues**<br>• Optical fibers bring IOs to optical connectors on front panel<br>• Front panel interconnect **density limited by size optical connectors**<br>• **Very high reliability/quality required** | • Optical **transceivers co-packaged w/ ASIC**<br>• Minimized electrical interconnect **eliminates SI issues**<br>• Optical fibers bring IOs to optical connectors on front panel<br>• **Lowest system power dissipation**<br>• Highest front panel density and smallest potential system form factor<br>• **Very high reliability required** |

Luxtera Proprietary     **LUXTERA**     8/11/2015 Page 10

Slide courtesy of P. De Dobbelaere, Luxtera

Reducing system power dissipation by integration

Slide courtesy of P. De Dobbelaere, Luxtera

**A macrochip with WDM links**

- Silicon lattice engineered with waveguides carries CPUs & DRAMs
  - Bridges convert electrons to photons & couple to waveguides in lattice
  - A high-bandwidth fully connected point-to-point optical network connects the sites

DRAM stack

CPU

ORACLE

Approved for public release. Distribution Unlimited

Courtesy of A. Krishnamoorthy, Oracle

Courtesy of A. Krishnamoorthy, Oracle

# 32 nm SOI CMOS-Driven Link (Aurrion/IBM)

*Current practice:*
*Low integration level,*
*inefficient 50 $\Omega$ interfaces*



*Maximum efficiency: directly drive the EAM*



**Aurrion's heterogeneous integration platform for photonics**

InP

SOI



- Integration of lasers, modulators, & detectors on the same wafer
- Adds III-V functionality to Si Photonics

**30 Gb/s**

TX out



RX out



- **3 pJ/bit at 30 Gb/s (not including laser power)**

- **No measured penalty for 10km transmission at 25 Gb/s**

- N. Dupuis *et al.*, "30Gbps Optical Link Utilizing Heterogeneously Integrated III-V/Si Photonics and CMOS Circuits," *OFC 2014 (*post deadline).
- N. Dupuis *et al.*, "30Gbps Optical Link Combining Heterogeneously Integrated III-V/Si Photonics with 32nm CMOS Circuits," *JLT*, 2015.

# Circuit/Photonics/Packaging Co-Design

## Electronic and photonic chip integration

CMOS IC
(IBM) →

*EAM bias
decoupling
capacitors*

*Photonic IC
(Aurrion)* →

## Demonstrated in hardware

*CMOS IC
(IBM)* →

*Photonic IC
(Aurrion)* →

*Transmitter output*

### 32 Gb/s

CH 1
9.5 dB

CH 2
7.6 dB

CH 3
9.4 dB

CH 4
8.2 dB

→ 31 ps ←

- B. G. Lee *et al.*, "A WDM-Compatible 4 × 32-Gb/s CMOS-Driven Electro-Absorption Modulator Array," *OFC* 2015.

# Wall-Plug WDM links (Aurrion/IBM)



| Parameter | CH 1 | CH 2 | CH 3 | CH 4 |
|---|---|---|---|---|
| EA Bias (V) | 4.0 | 4.6 | 3.9 | 5.1 |
| Extinction Ratio (dB) | 5.9 | 5.5 | 6.6 | 7.0 |
| EAM Power (mW) | 7.22 | 4.39 | 3.28 | 3.64 |
| Laser Power (mW) | 225 | 245.5 | 287.0 | 246.3 |
| Total Laser Power (mW) | 1004 | | | |
| Total EAM Power (mW) | 18.5 | | | |
| CMOS Driver Power (mW) | 97.6 | | | |

Low modulation power → directly driving EAM, no 50Ω interfaces

- A. Ramaswamy *et al.*, "A WDM 4x28Gbps Integrated Silicon Photonic Transmitter driven by 32nm CMOS driver ICs," *OFC 2015* (post deadline).

# Photonic Switching

# What do we mean by *"fast"* optical switching?

## ms-scale



**Mice flows over packet switch, elephant flows over OCS**
**Promise:**
- Lower cost/power, fewer cables, software control (SDN)

**Challenges:**
- Scalability of software scheduler
- Slow reconfiguration time may limit applications

## μs-scale



Examples: Mordia/
REACTOR (UCSD)

OCS at first switch level, hybrid but much faster then 3D-MEMS
**Promise:**
- Reconfigure network at flow-level, based on workloads
- Hardware control (FPGA).

**Challenges:**
- Custom NICs
- Scalability: switch radix

## ns-scale



All-optical switching, electronic buffering at end points
**Promise:**
- Switching times ~ packet durations
- More power-efficient than electrical networks of equal BW

**Challenges:**
- Switch hardware and fast synchronizing links
- Scalability: switch radix: losses, fast control plane, flow control

C. Schow, UCSB

Adapted from L. Schares, IBM

Courtesy of Prof. G. Papen and G. Porter, UCSD

# Calient: 3-D MEMS OCS



Fig. 3. Switch configuration using 3-D MEMS.



X. Zheng et al., "Three-dimensional MEMS photonic cross-connect switch design and performance," JSTQE, 2003.



**Figure 4: Hybrid Packet-OCS Datacenter Network Architecture**

"The Software Defined Hybrid Packet Optical Datacenter Network" whitepaper available at www.calient.net

## High port count (320), low insertion loss, low crosstalk, <50ms reconfiguration

# Photonic Switches in Data Centers

## UCSD Hybrid Networking Research



Courtesy of Prof. G. Papen and G. Porter, UCSD

**ToR photonic switches:**
- Fast reconfiguration (dynamic traffic)
- High-radix
- Low cost

- N. Farrington *et al.,* "Helios: a hybrid electrical/optical switch architecture for modular data centers." *SIGCOMM*, 2011.
- R. Aguinaldo *et al.*, "Energy-efficient, digitally-driven "fat pipe" silicon photonic circuit switch in the UCSD MORDIA data-center network." *CLEO* 2104.
- H. Liu *et al.*, "REACToR: A reconfigurable packet and circuit ToR switch," *Photonics Society Summer Topicals*, 2013.

*Monolithically integrated switch +driver chip
(IBM 90nm photonics-enabled CMOS)*



*Integrated digital switch drivers*



**Fast reconfiguration:** 4 ns



*Transient response of 2x2 MZ switch*

**Broad spectral bandwidth:**
Routing many wavelength channels
- <-20dB crosstalk over 60 nm BW
- 32 channels at 200GHz spacing



*Spectral response of 2x2 MZ switch*

**Losses too high, need significant feedback and control to manage crosstalk**
- **High level of electronic/photonic integration demanded**

- N. Dupuis *et al.*, " Modeling and Characterization of a Non-Blocking 4 × 4 Mach-Zehnder Silicon Photonic Switch Fabric," *JLT* 2015.
- N. Dupuis *et al.*, "Design and Fabrication of Low-Insertion-Loss and Low-Crosstalk Broadband 2 × 2 Mach–Zehnder Silicon Photonic Switches," *JLT* 2015.
- B. G. Lee *et al.*, "Monolithic Silicon Integration of Scaled Photonic Switch Fabrics, CMOS Logic, and Device Driver Circuits,"  *JLT* 2014.
- B. G. Lee *et al.*, "Silicon Photonic Switch Fabrics in Computer Communications Systems," *JLT* 2015.

**Detailed Schematic of Vertical Adiabatic Coupler Switch**



Courtesy Professor M. Wu, UC Berkeley

**UCSB**



Courtesy of Professor S. J. Ben Yoo, UC Davis
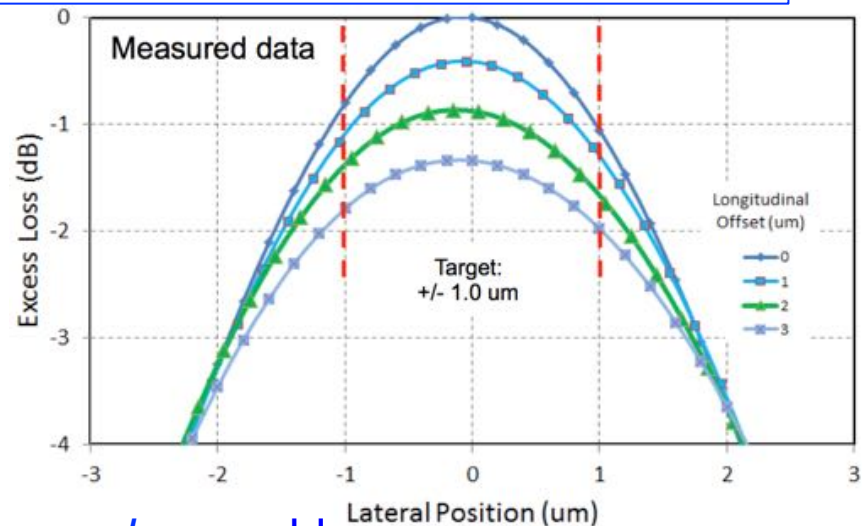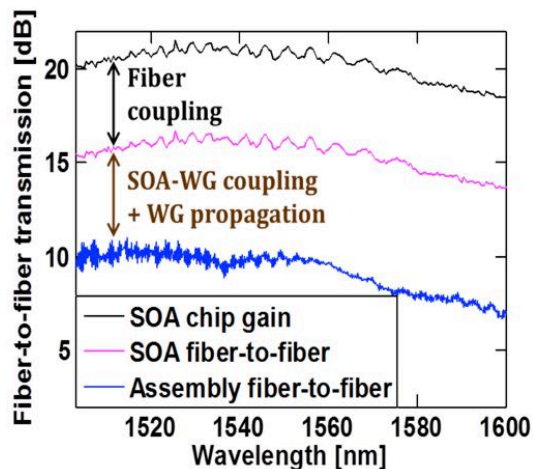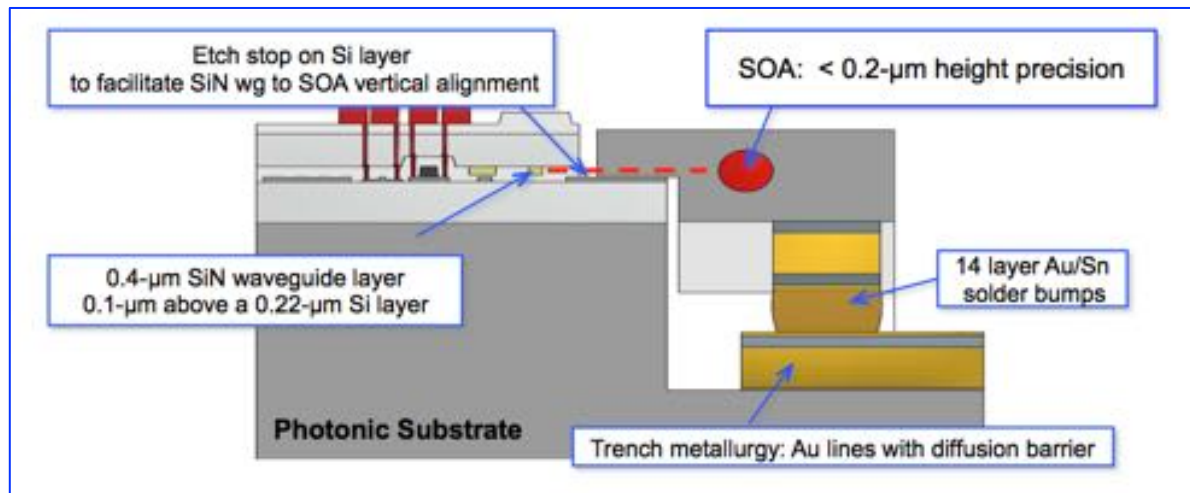
# Integrating Optical Gain for Scalability

Hybrid Approach (IBM):
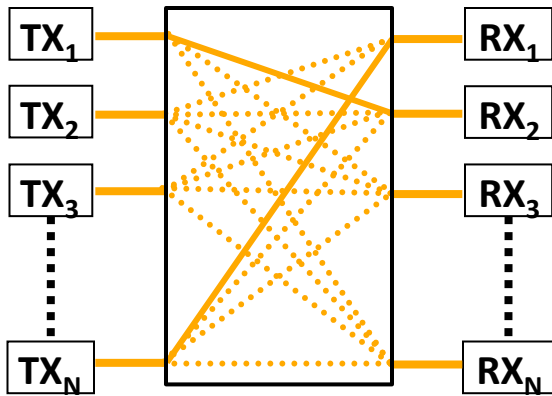


**Solder attachment of SOA arrays to photonic carriers**

- R. A. Budd et al., "Semiconductor Optical Amplifier (SOA) Packaging for Scalable and Gain-Integrated Silicon Photonic Switching Platforms," *ECTC 2015.*
- L. Schares et al., "Etched-Facet Semiconductor Optical Amplifiers for Gain-Integrated Photonic Switch Fabrics," *ECOC* 2015.
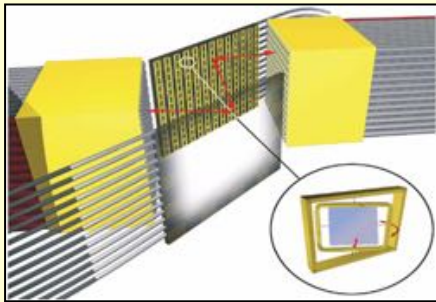
# Challenging Packaging

- **Challenging fabrication *and* assembly**
- **Precise alignment required for *each* SOA chip**

# New Electronic Capabilities Needed for Photonic Switching



- Optical circuits configured through a switch or fabric to connect $TX_i$ with $RX_j$

- Switching time defines context of usage

- But, switching time is not all hardware reconfiguration
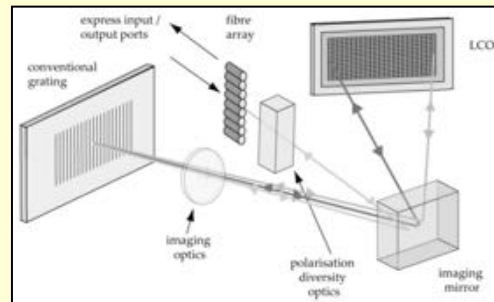
---

## Millisecond-scale



### e.g. Calient: 3D-MEMS
[J. Opt. Netw. **6** (1) 19]

- Hybrid (circuit + packet) networks
- Coarse reconfiguration
- Software control (SDN)
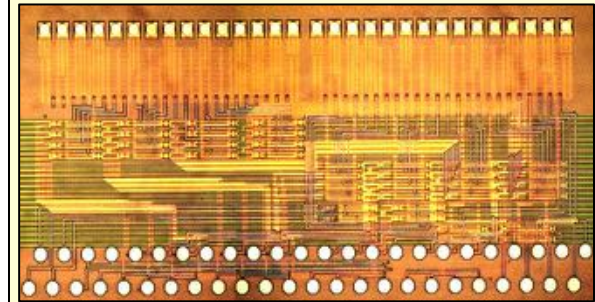- Highly scalable (> 1000 ports)

---

## Microsecond-scale



### e.g. Finisar: LCOS
[OFC 2006, OTuF2]

- Hybrid networks
- Reconfigure at flow-level
- Hardware control (FPGA)
- Scalable (10's to 100 ports)
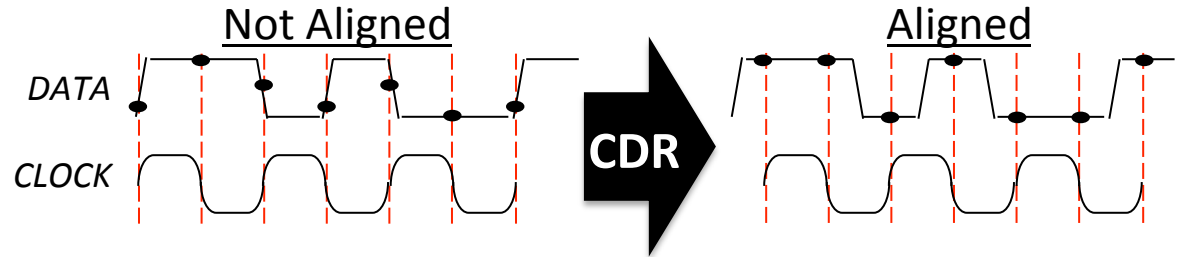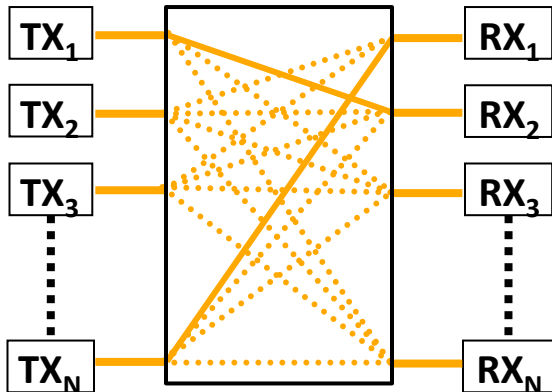
---

## Nanosecond-scale



### e.g. IBM: Photonics
[OFC 2013, PDP5C.3]

- Reconfigure at packet granularity
- Quasi-packet switching with buffering at end points
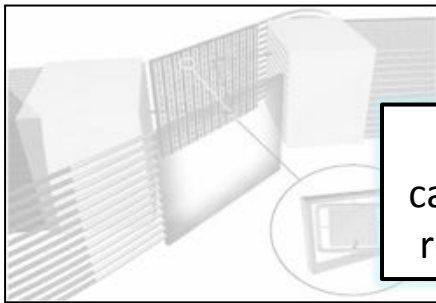- Limited scalability (10's of ports)

Courtesy of B. Lee, IBM

# Switching Time = Switch Reconfiguration Time + Link Synchronization Time



DATA

CLOCK

Not Aligned

CDR

Aligned

*Receivers must adapt to abrupt changes in incoming data amplitude and phase*

**Millisecond-scale**

**Microsecond-scale**

**Nanosecond-scale**

Receivers with burst-mode capabilities addressed through relevant standards (e.g. PON)

Receivers with burst-mode capabilities needed

**e.g. Calient: 3D-MEMS**
[J. Opt. Netw. **6** (1) 19]

- Hybrid (circuit + packet) networks
- Coarse reconfiguration
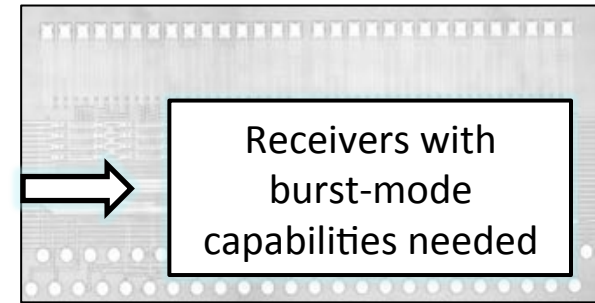- Software control (SDN)
- Highly scalable (> 1000 ports)

**e.g. Finisar: LCOS**
[OFC 2006, OTuF2]

- Hybrid networks
- Reconfigure at flow-level
- Hardware control (FPGA)
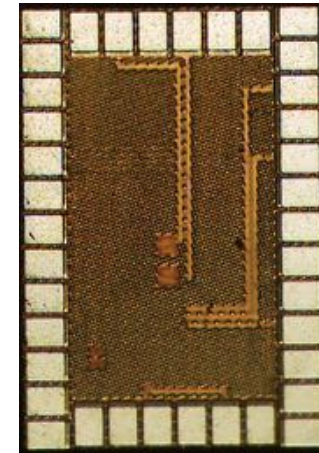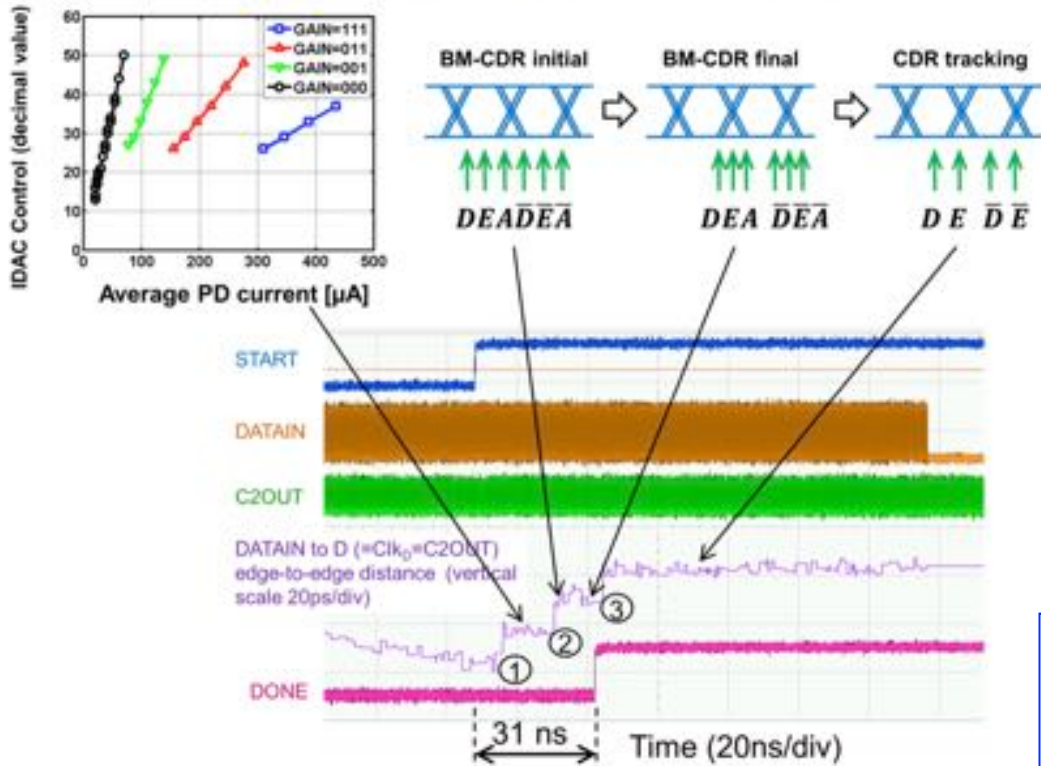- Scalable (10's to 100 ports)

**e.g. IBM: Photonics**
[OFC 2013, PDP5C.3]

- Reconfigure at packet granularity
- Quasi-packet switching with buffering at end points
- Limited scalability (10's of ports)

Courtesy of B. Lee, IBM
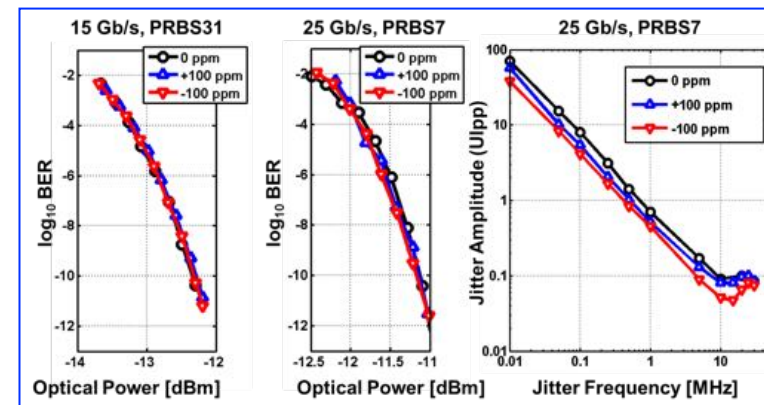
## Measured Dynamics of Burst-Mode Receiver



32nm SOI CMOS
(1.3 mm × 0.9 mm)



**Fast photonic routing and switching fabrics must have companion electronics**
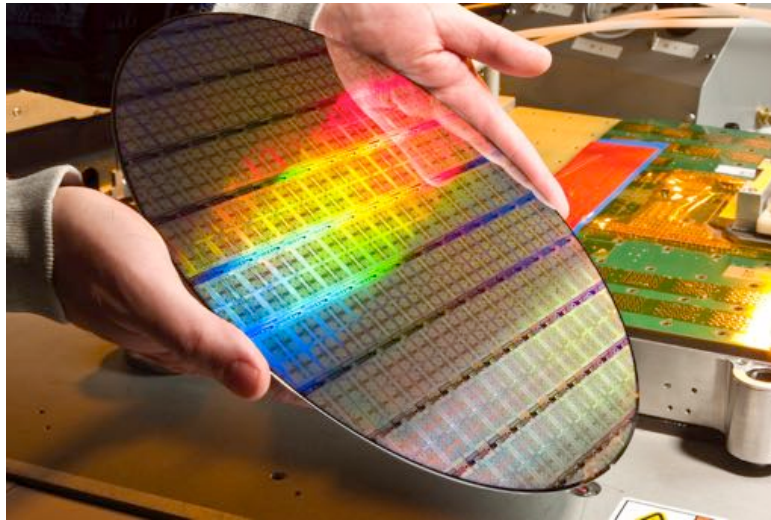- **Needs more research in the community**

**Enables fine-grained power management**

- A. R. Rylyakov *et al.*, "A 25 Gb/s Burst-Mode Receiver for Low Latency Photonic Switch Networks ," *OFC 2015.*
- A. R. Rylyakov *et al.*, "A 25 Gb/s Burst-Mode Receiver for Low Latency Photonic Switch Networks ," *JSSC 2015.*

Path Forward: Large-Scale
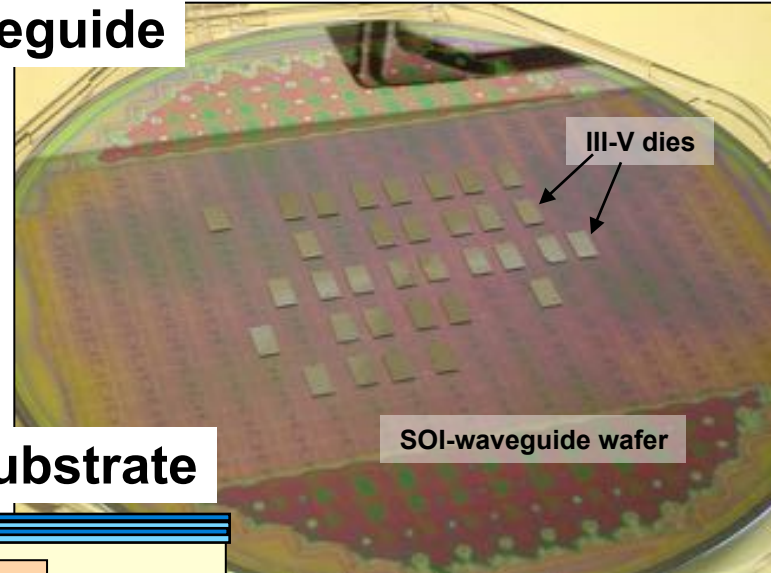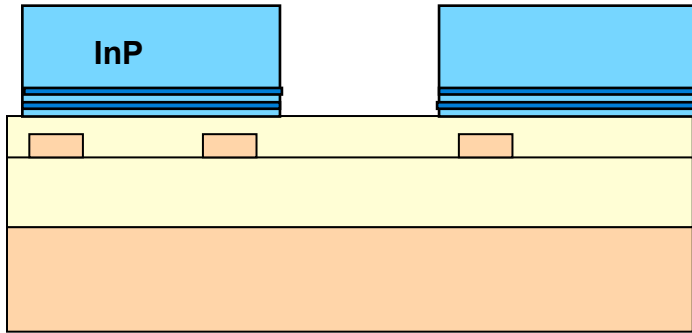Electronic/Photonic Integration

**UCSB**



**Si Photonics:** integrating photonic devices into Si platforms to leverage the huge Si electronics manufacturing infrastructure

- Very large scale integration
- Tight process control
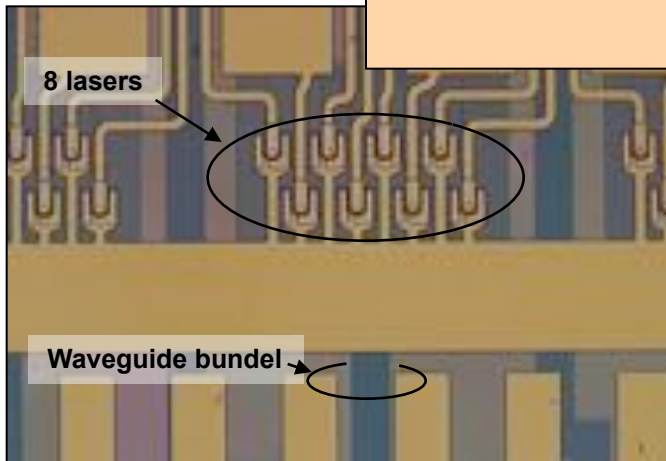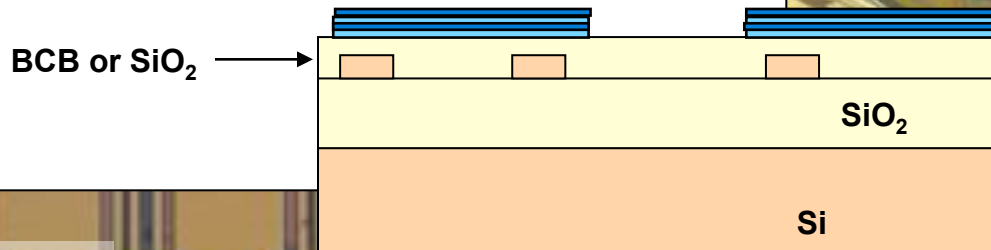- Wafer-level testing
- High-yield and low cost

**Advantages for the data center:**

- Single-mode → multi-kilometer links
- Wavelength-division multiplexing (WDM) → high BW/fiber
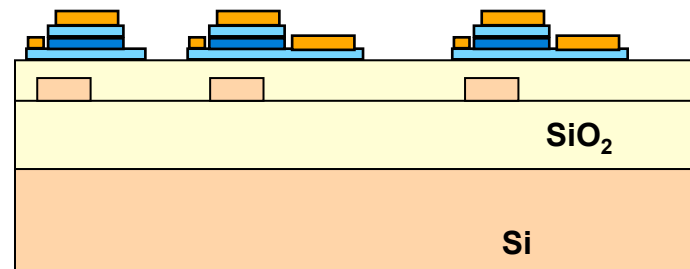- High integration level → many devices needed for switching and high-speed interconnects

# Heterogeneous integration

UCSB, Intel, HP, Ghent, IMEC,TIT, Caltech, NTT, Aurrion

## Step 1: Bond InP-dies on SOI waveguide

InP

III-V dies

SOI-waveguide wafer

## Step 2: Remove substrate

BCB or SiO$_2$

SiO$_2$

Si

8 lasers

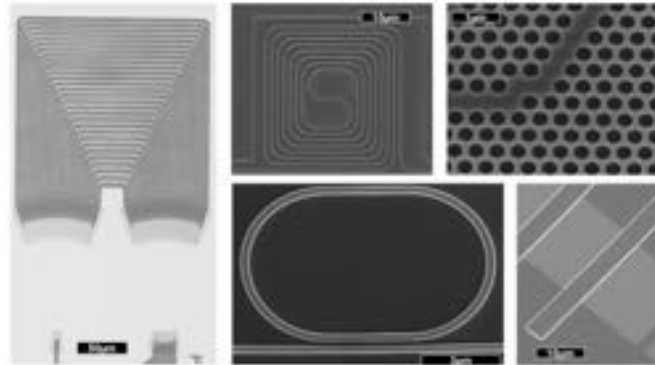Waveguide bundel

## Step 3: Process lasers at wafer scale

SiO$_2$

Si

Figures: IMEC/Ghent; Slide courtesy of Prof. J. Bowers, UCSB

GaAs

Silicon
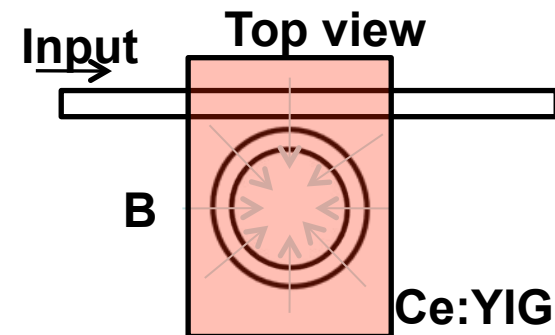
LiNbO$_3$

InP

SiN/SiON/SiO2

Ce:YIG Isolator

Input    **Top view**

**B**

**Ce:YIG**

(b)

Heck et al. JSTQE 2013

# Large Scale Integration for Switching

## Large-scale photonic integration

- Switching
- Amplification
- Control



### Hybrid III-V on Si SOAs Demonstrated in Multiple Platforms
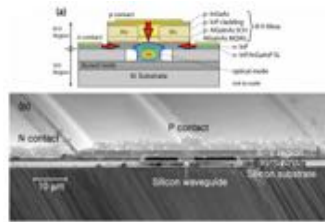
UCSB[1], 2006      Ghent University-IMEC[2], 2012      III-V Lab[3], 2014



1. H. Park *et al.*, "A Hybrid AlGaInAs-Silicon Evanescent Amplifier, *PTL* 2007.
2. S. Keyvaninia et al., "A Highly Efficient Electrically Pumped Optical Amplifier Integrated on a SOI Waveguide Circuit," *Group IV Photonics* 2012.
3. G.-H. Duan *et al.*, "New Advances on Heterogeneous Integration of III-V on Silicon," *JLT* 2015.

# Large Scale Integration for Photonic I/O

**Photonics must deliver system advantages**
- More I/O bandwidth at less power for processors
- Larger port switches to enable flatter networks
- Higher Integration levels: processor, memory, network



**Hundreds of photonic interfaces**: all off-module high-speed I/O
- High BW/interface → WDM
- Low cost → compatibility with high-volume electronics manufacturing
- Low loss → maximize link power efficiency by minimizing required laser power

**Thousands of electrical interfaces:** connecting electronics to photonics
- High density, high-speed, low-power
- New functions: rapid synchronization for low latency switching, power management
- Specialized chip I/O co-designed with and only for photonics, no general purpose cores

**Reliability**: components either don't fail or can be spared (also for yield)

- **Large-scale integration is fundamentally required**
- **Holistic design of photonics, electronics, packaging and assembly**
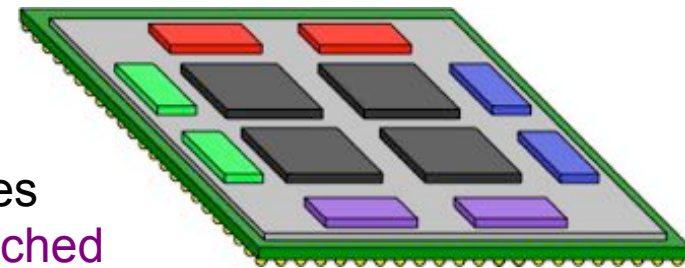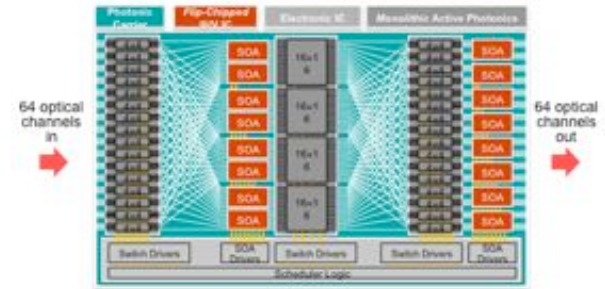
# Closing Thoughts

## Potential: More highly connected systems

### Photonic Switching
- Routing 10's of Tb/s at sub pJ/bit power

### Photonic I/O
- Higher radix electrical switches, flatter networks
- More processor/memory bandwidth

- Multiple photonic technologies for multiple purposes
  - VCSEL, PSM, WDM point-to-point, WDM switched

## Challenge: Integrating large-scale electronics with large-scale photonics

### Re-thinking manufacturability and supply chain
- Moving from electronic to photonic-centric packaging
- Electronics/Photonics/Package co-design required
- What are the highest-value components and what assembly flow makes sense
- Who does what?

# Thank You!

## schow@ece.ucsb.edu