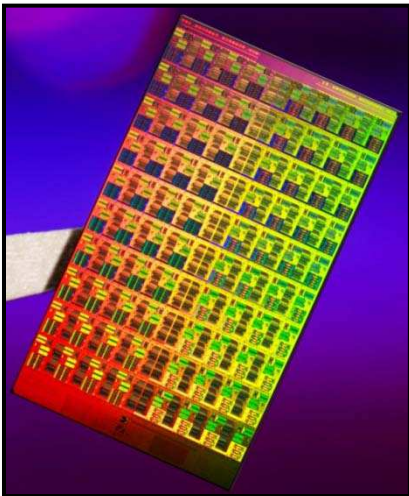# Hybrid On-chip Data Networks
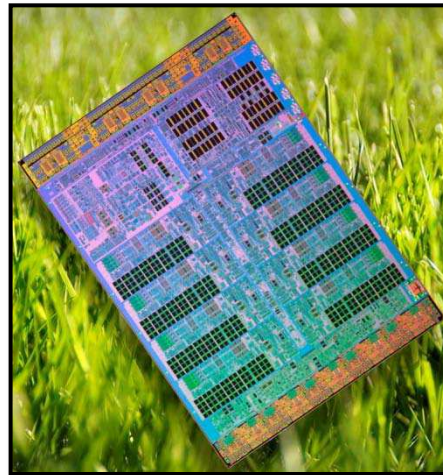
**Gilbert Hendry**

Keren Bergman
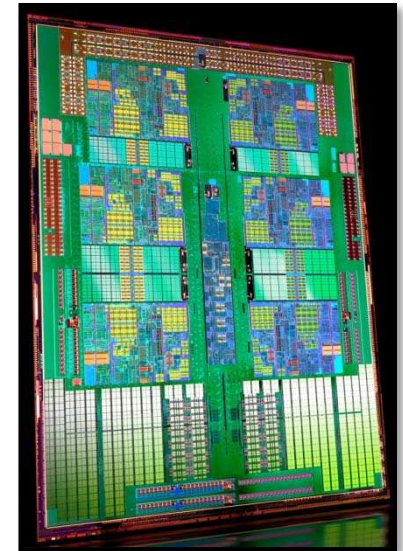
*Lightwave Research Lab*
Columbia University

- **Chip multi-processors create need for high performance interconnects**

- **Performance bottleneck of on-chip networks and I/O**

- **Power dissipation constraints of the chip package**

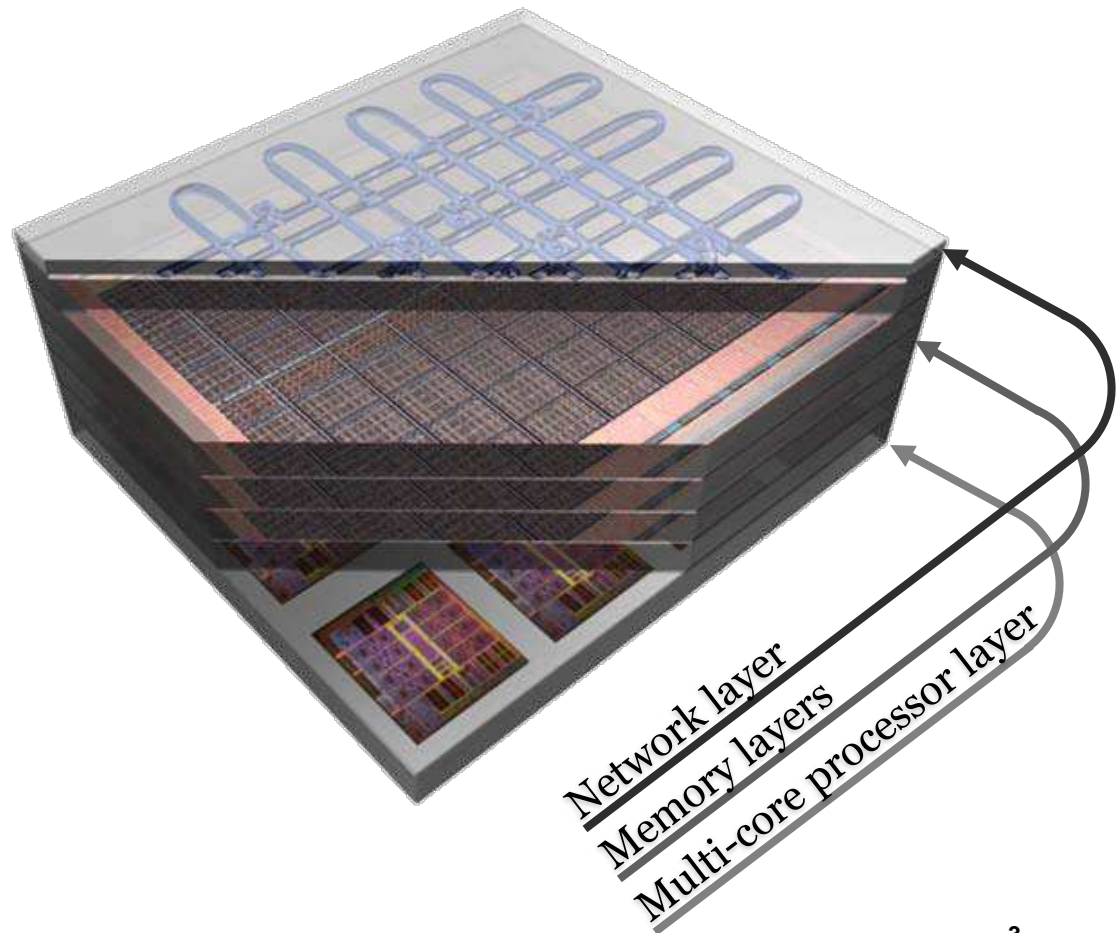  - **> 50% of total power comes from interconnects***


Intel Polaris


IBM Cell


AMD Opteron

* N. Magen *et al.*, "Interconnect-power dissipation in a microprocessor," SLIP 2004.
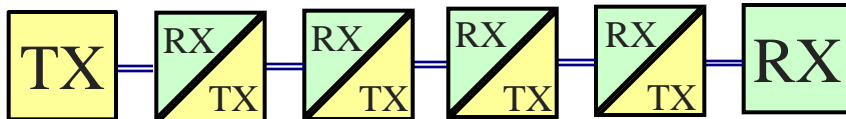
- **CMPs of the future = 3D stacking**
- **Lots of data on chip**
- **Photonics offers**
  **key advantages**



Network layer
Memory layers
Multi-core processor layer

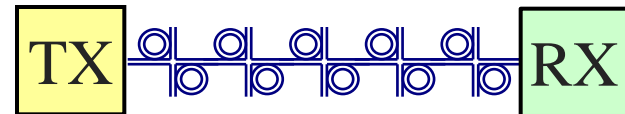**Photonics changes the rules for Bandwidth, Energy, and Distance.**

## ELECTRONICS:

- Buffer, receive and re-transmit at every router.

- Each bus lane routed independently. ($P \propto N_{LANES}$)
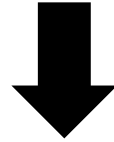
- Off-chip BW is pin-limited and power hungry.
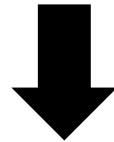
## OPTICS:

- Modulate/receive high bandwidth data stream once per communication event.

- Broadband switch routes entire multi-wavelength stream.

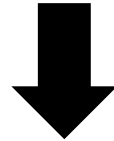- Off-chip BW = On-chip BW for nearly same power.

Optical processing difficult and limited

⬇

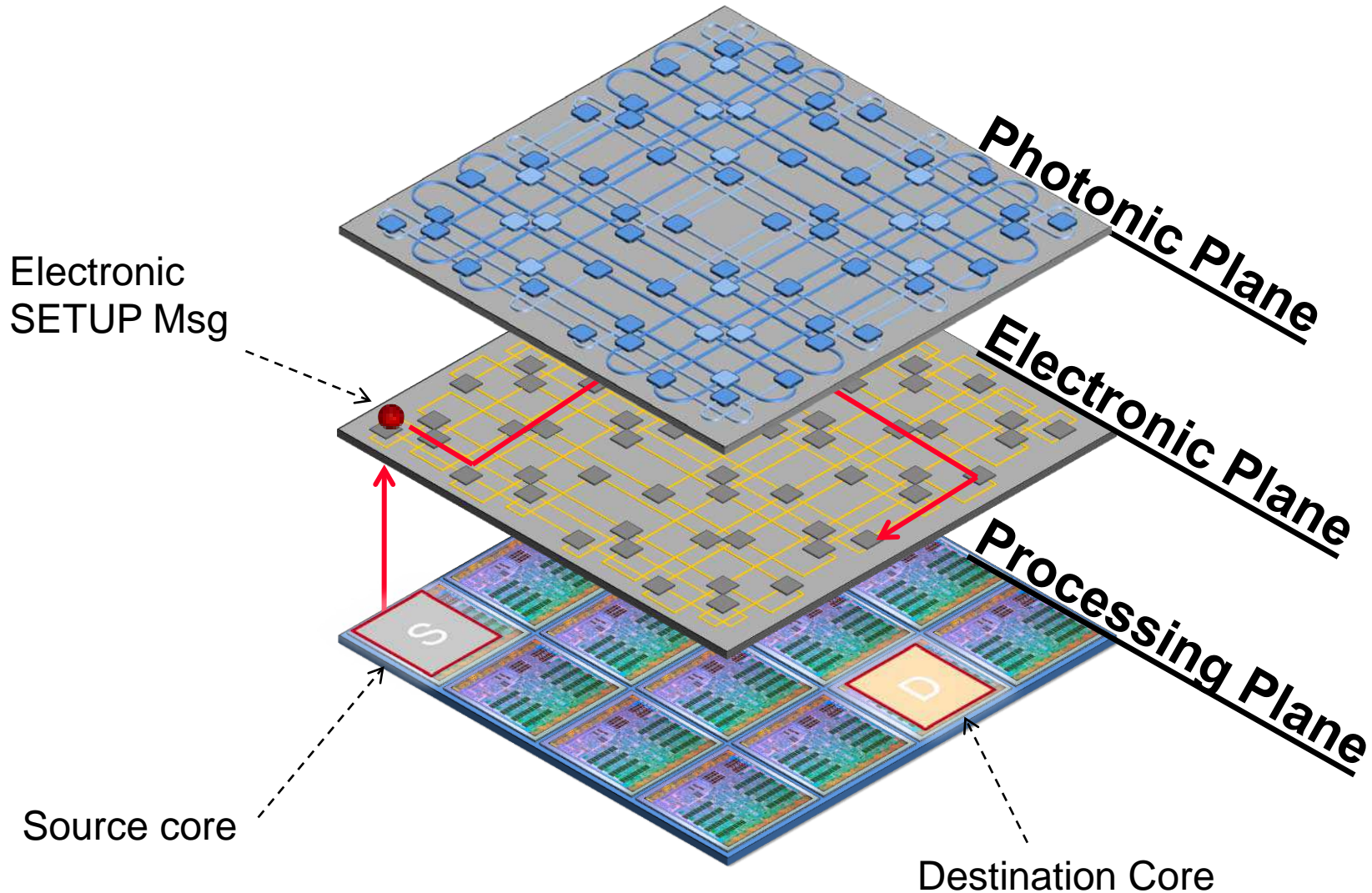Source, destination routing inefficient

⬇

Use electronics for routing,
optics for switching and transmission
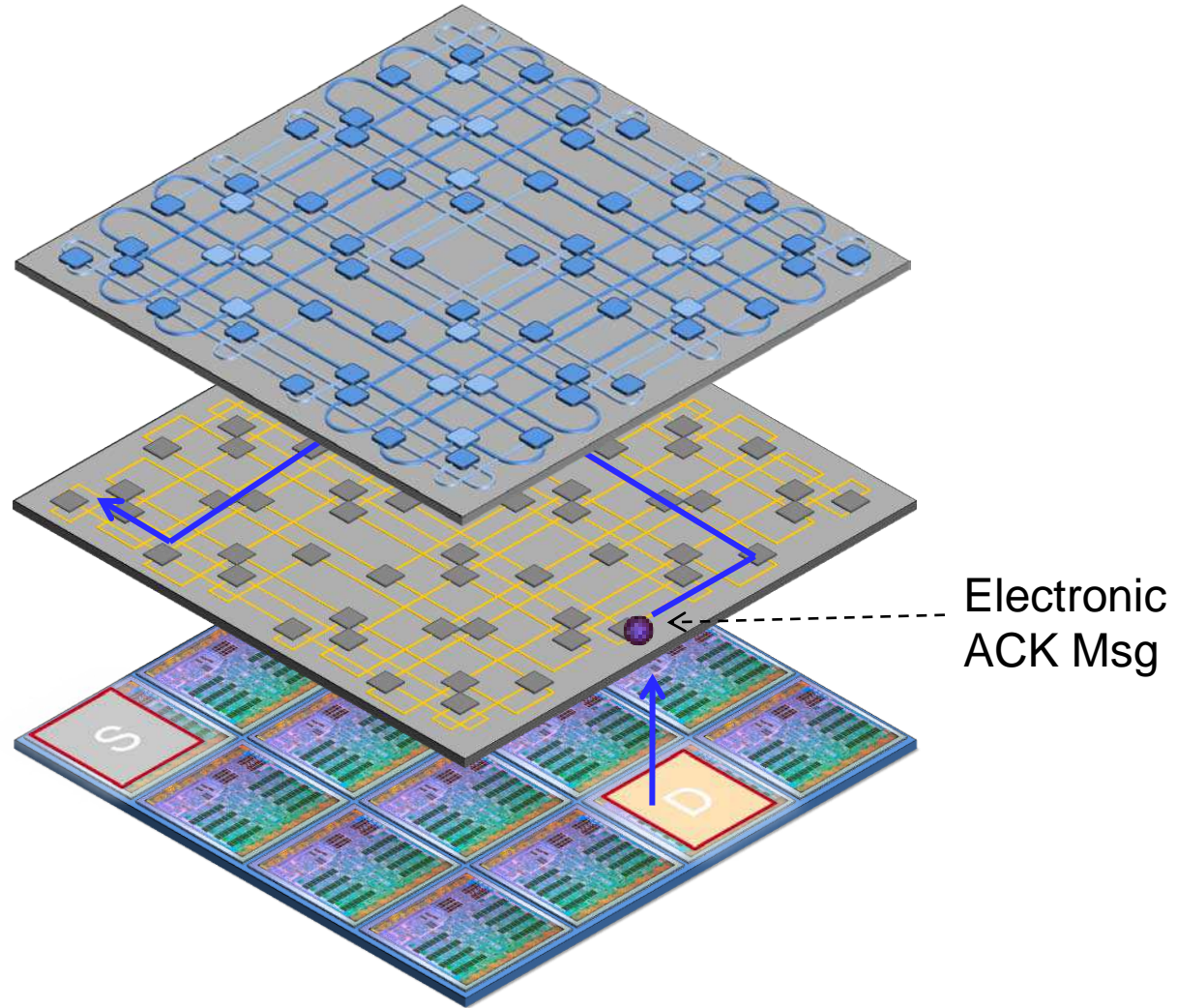
⬇

Hybrid Circuit-Switching

Step 1: Path SETUP request



Photonic Plane

Electronic
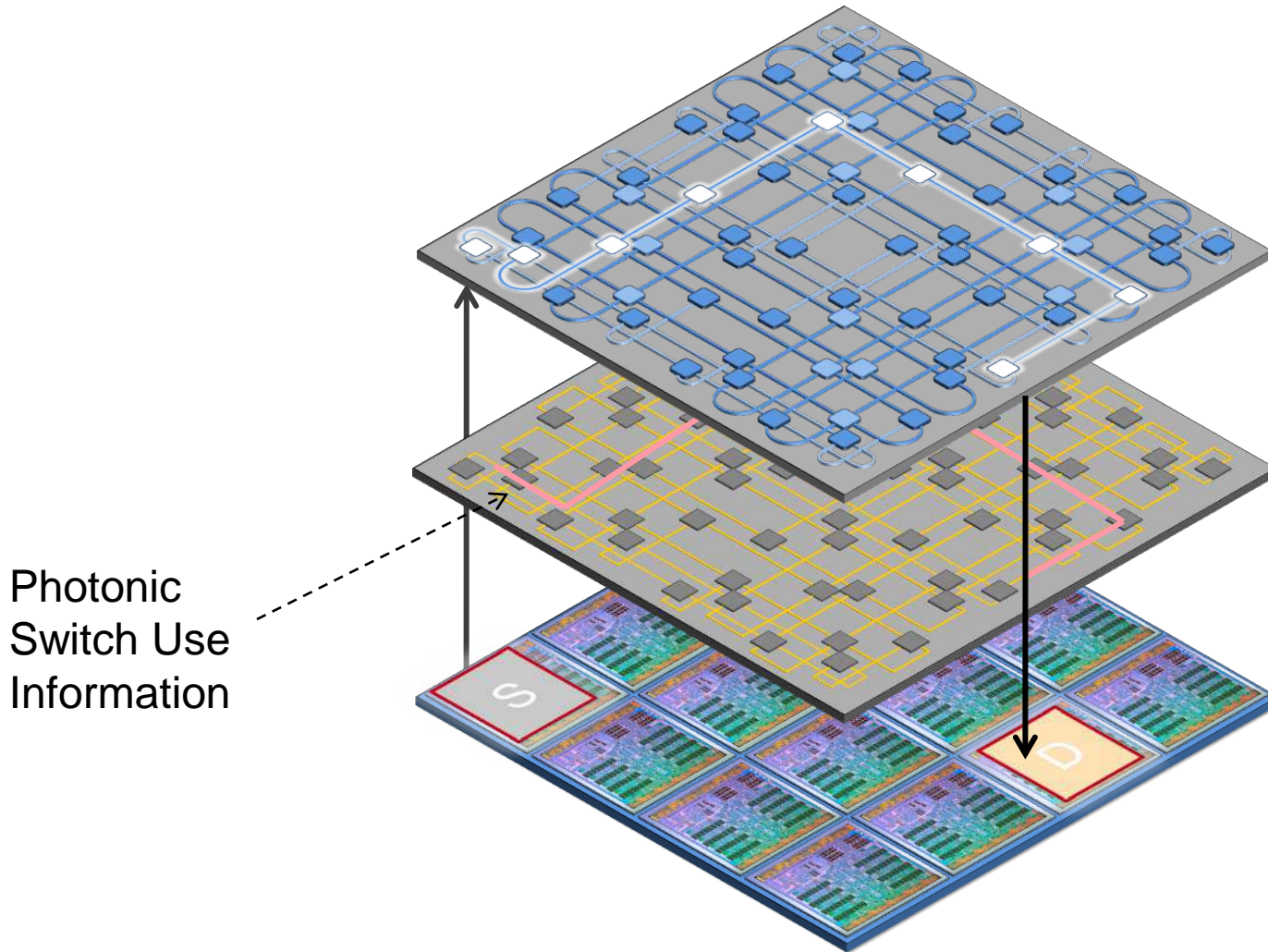SETUP Msg

Electronic Plane

Processing Plane

Source core

Destination Core

## Step 2: Path ACK



Electronic
ACK Msg

## Step 3: Transmit Data

Photonic
Switch Use
Information

Meanwhile: Path Contention



Path
BLOCKED Msg
(Backoff)

## Step 4: Path TEARDOWN

Electronic
TEARDOWN
Msg

Source core

Destination Core

# Hybrid Circuit-Switched Networks

Pros:

- **Energy-efficient end-to-end transmission**
- **High bandwidth through WDM**
- **Electronic network still available for small control messages***
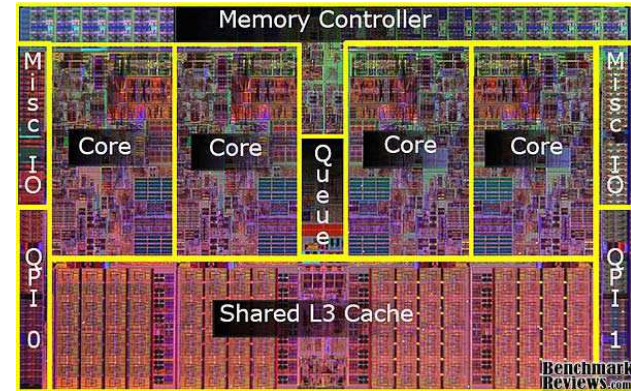- **Network-level support for secure regions**

Cons:

- **Path setup latency**
- **Path setup contention (no fairness)**

* [G. Hendry et al. *Analysis of Photonic Networks for a Chip Multiprocessor Using Scientific Applications.* In NOCS, 2009]
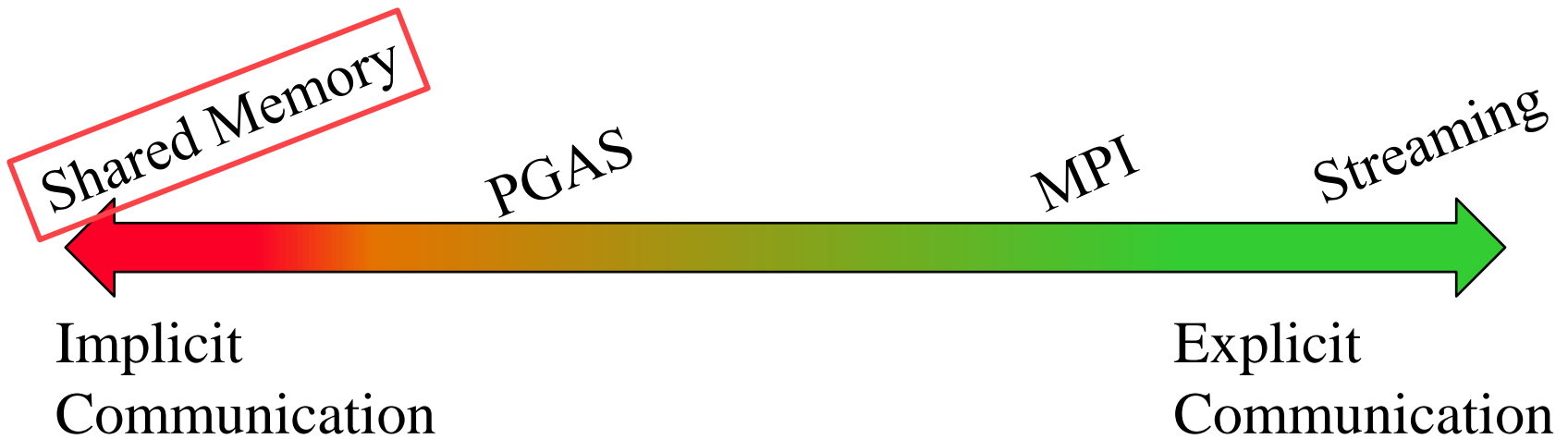
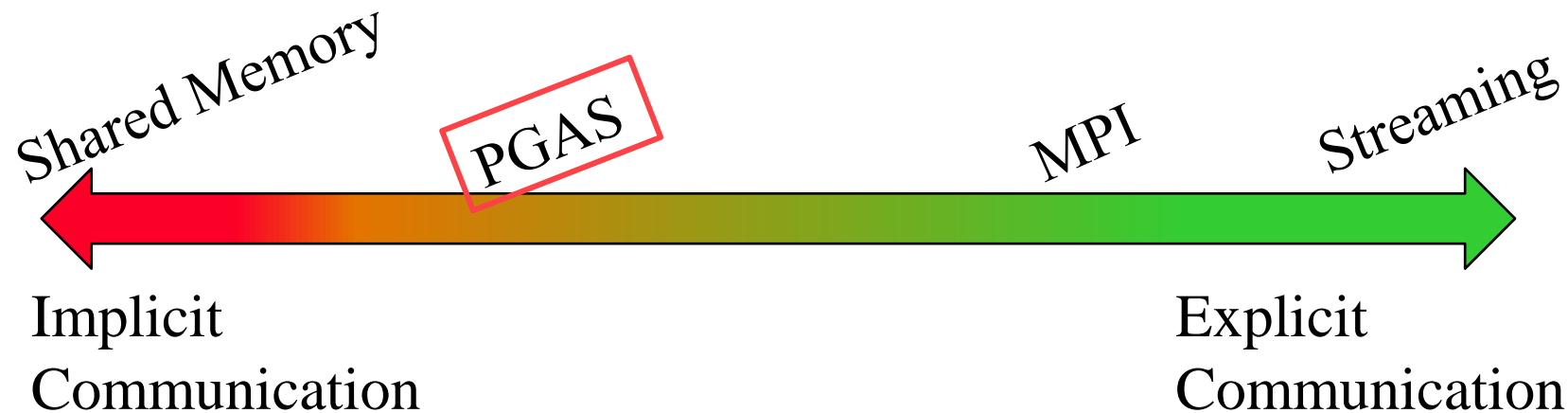# Programming and Communication

# Shared Memory



scaling



"… [OpenMP on large systems] often performs worse than message passing due to a combination of false sharing, coherence traffic, contention, and system issues that arise from the difference in scheduling and network interface moderation"
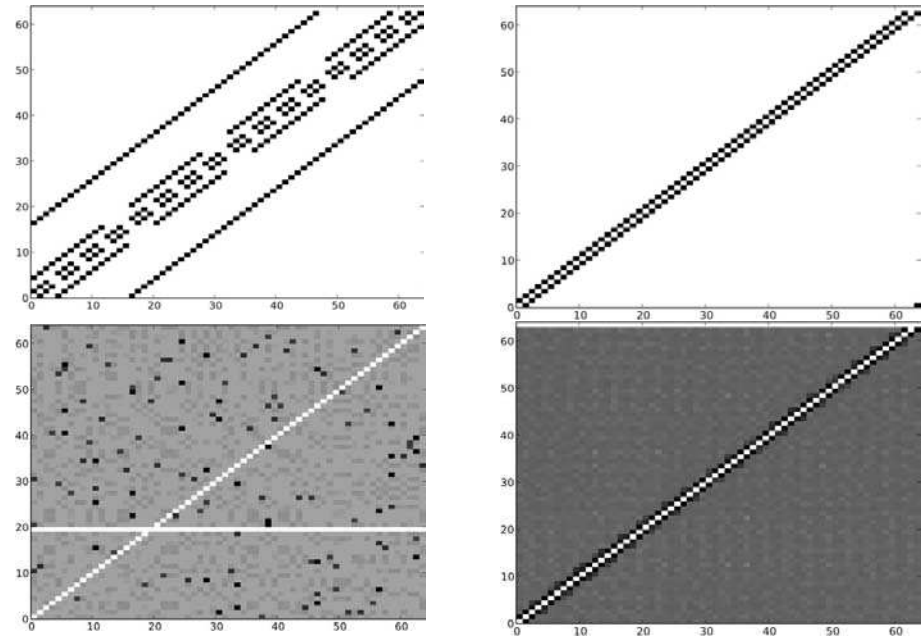
~ Exascale Report

Shared Memory    PGAS    MPI    Streaming

Implicit
Communication

Explicit
Communication

# Partitioned Global Address Space

| Access | Method |
|--------|--------|
| Local Read | Optical Receive |
| Local Write | Optical send |
| Remote Read | Electronic request, optical receive |
| Remote Write | Optical send |
| Shared R/W | ? |

Shared Memory      PGAS      MPI      Streaming

Implicit
Communication

Explicit
Communication

[G. Hendry et al. *Circuit-Switched Memory Access in Photonic Interconnection Networks for HPEC*. In Supercomputing, Nov. 2010]

# Message Passing

- Complex, dynamic access patterns
- Relatively larger blocks of data
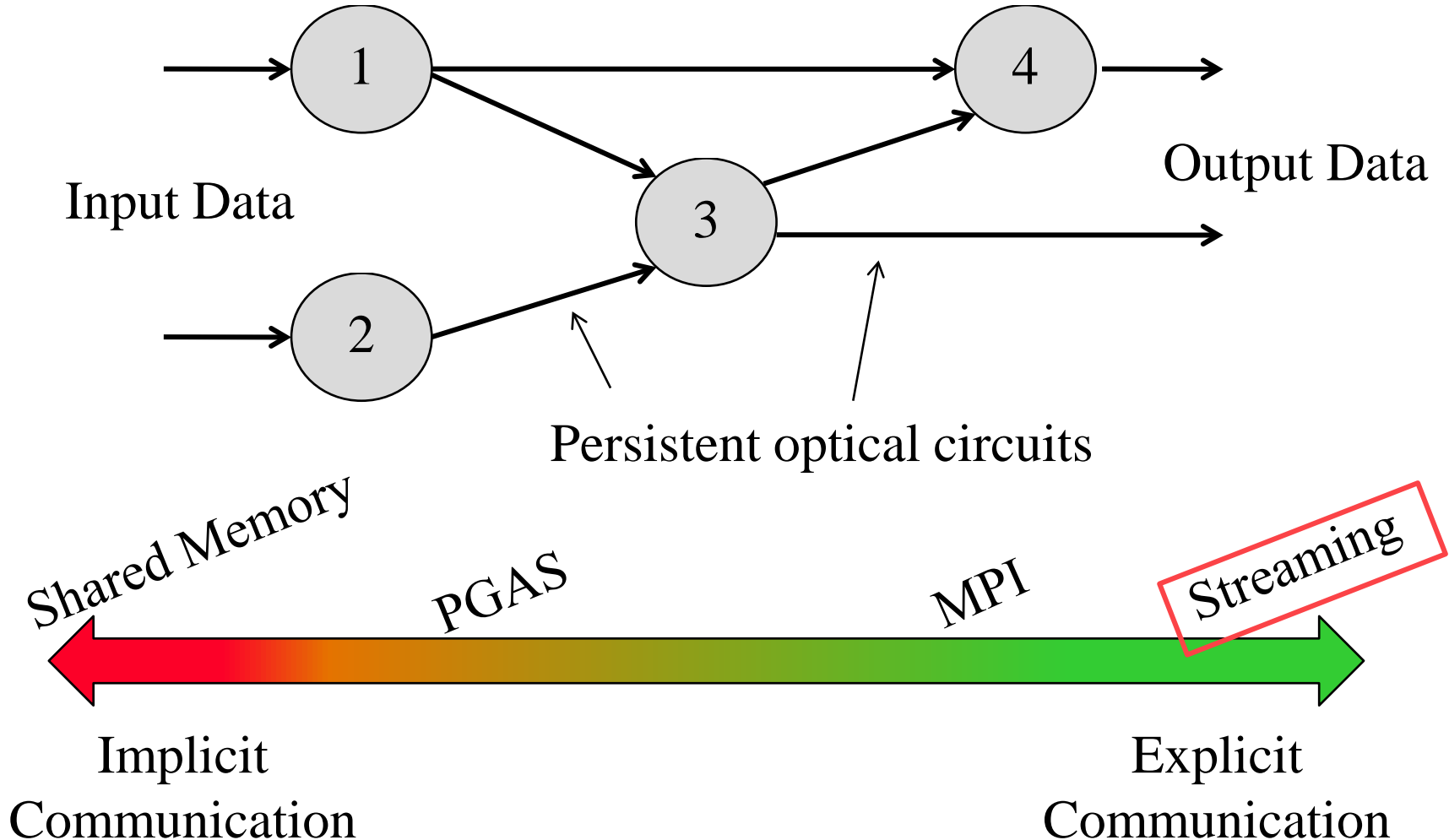- Scientific computing →



Shared Memory    PGAS    MPI    Streaming
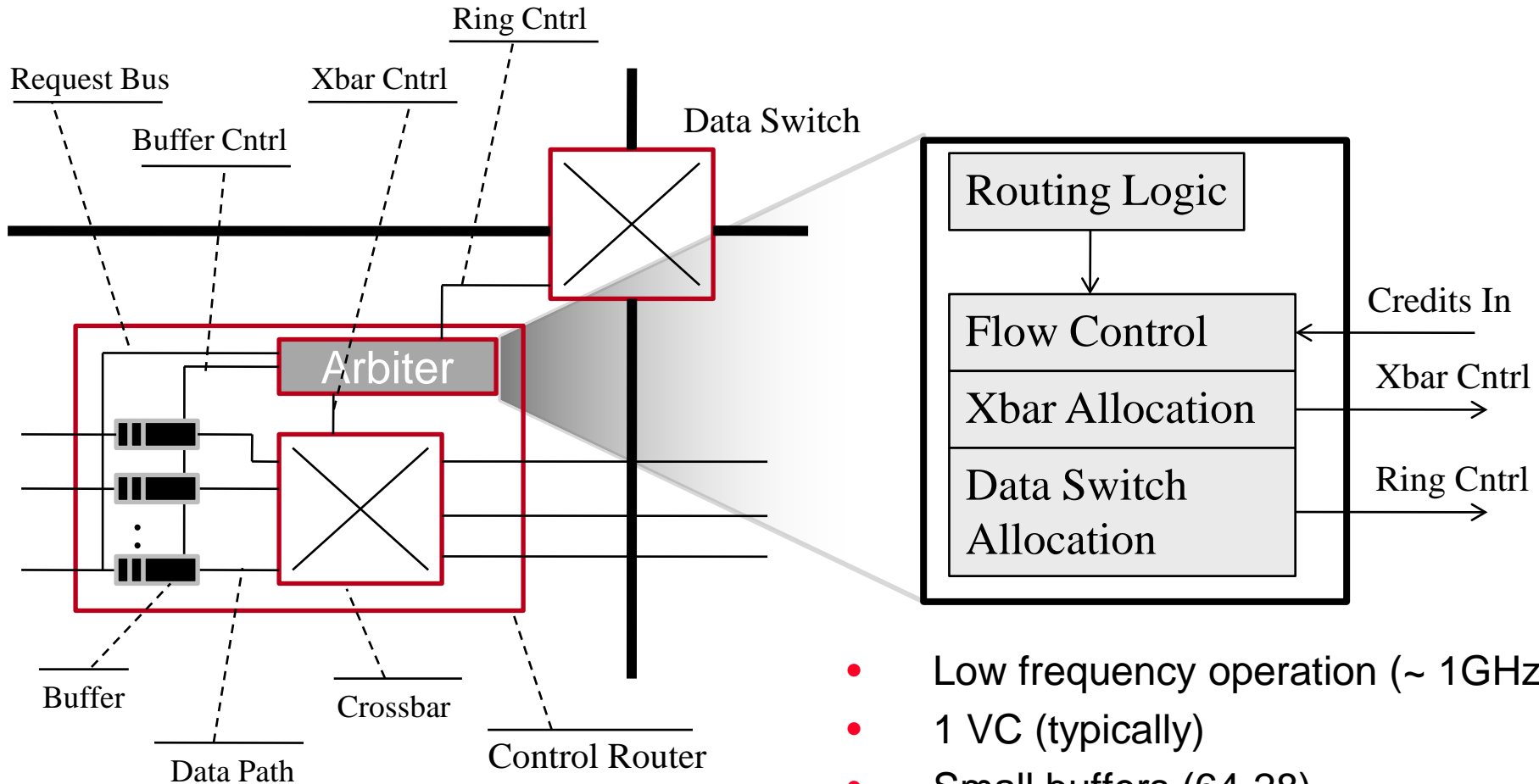
Implicit
Communication

Explicit
Communication

* [G. Hendry et al. *Analysis of Photonic Networks for a Chip Multiprocessor Using Scientific Applications.* In NOCS, 2009]

**15**

- Embedded / specialized systems (Graphics, Image + Signal Proc.)
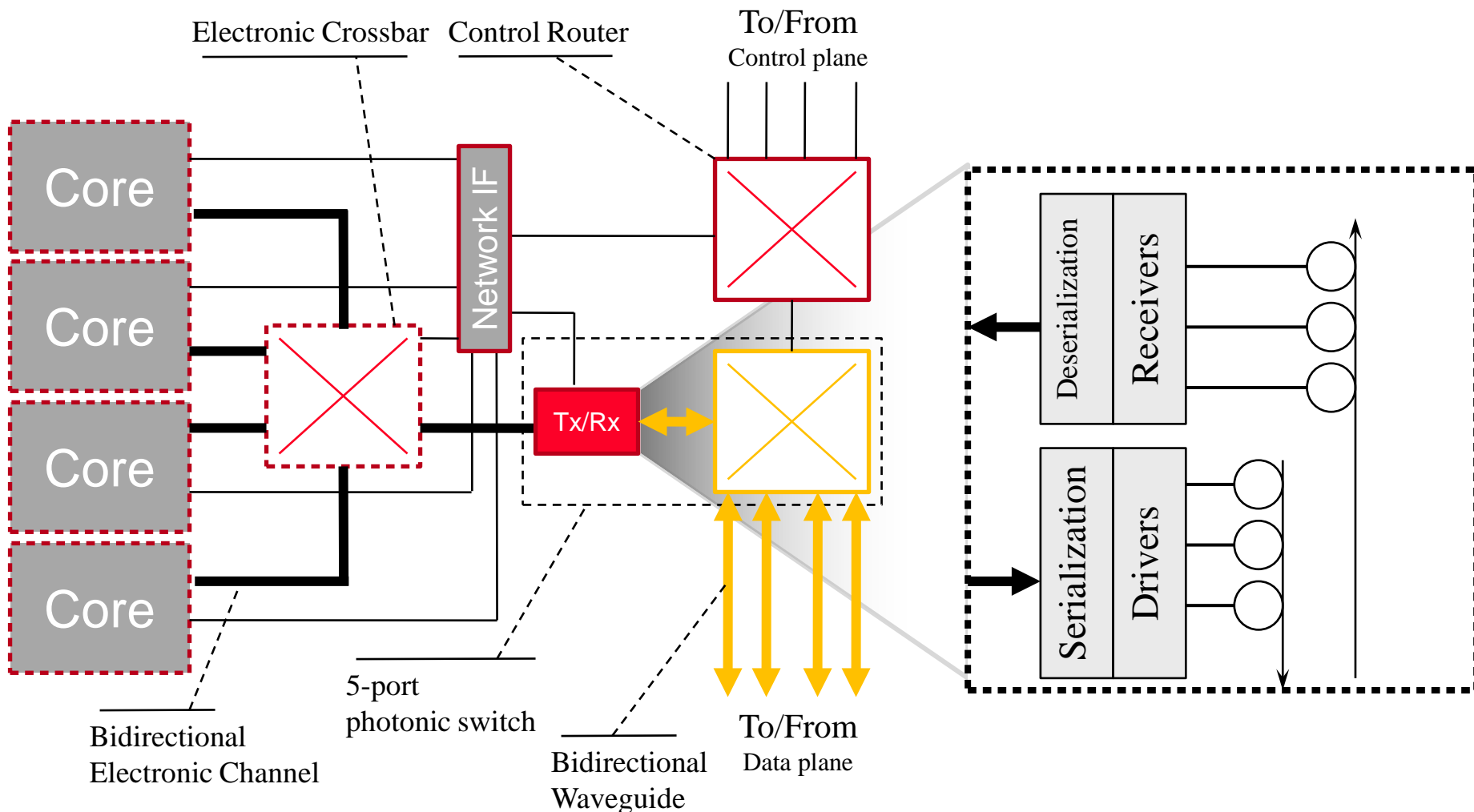- Execution mode of general-purpose systems (Cell Processor)



Input Data

Output Data

Persistent optical circuits

Shared Memory

PGAS

MPI

Streaming

Implicit
Communication

Explicit
Communication

# Electronic Plane

Ring Cntrl

Request Bus

Xbar Cntrl

Buffer Cntrl

Data Switch

Arbiter

Routing Logic

Credits In

Flow Control

Xbar Allocation

Xbar Cntrl

Data Switch Allocation

Ring Cntrl

Buffer

Crossbar

Data Path

Control Router

- Low frequency operation (~ 1GHz)
- 1 VC (typically)
- Small buffers (64-28)
- Narrow Channels (8-32)

Electronic Crossbar  Control Router  To/From Control plane

Network IF

Tx/Rx

Core

Core

Core

Core

Deserialization  Receivers

Serialization  Drivers

5-port photonic switch

Bidirectional Electronic Channel

Bidirectional Waveguide

To/From Data plane

# External Concentration

[P. Kumar et al. *Exploring concentration and channel slicing in on-chip network router*. In NOCS, 2009]

# The Photonic Plane

waveguide

λ

# Silicon Photonic Waveguide Technology

**1.28 Tb/s Data Transmission Experiment**
(*occupies small slice of available WG BW*)



before injection into waveguide

after 5-cm waveguide and EDFA

**C23** (1559 nm)  **C28** (1555 nm)  **C46** (1541 nm)  **C51** (1537 nm)

[B. G. Lee *et al.*, *Photon. Technol. Lett.* **20** (10) 767 (2008)]

[Vlasov and McNab, *Optics Express* **12** (8) 1622 (2004)]

Silicon photonic waveguides provide low-power optical interconnects in CMOS-compatible platform.

Low-loss (1.7 dB/cm), high-bandwidth (> 200 nm) silicon photonic waveguides can be fabricated in commercial CMOS process.

22

modulator/filter

Broadband spatial switch

LIGHTWAVE RESEARCH LABORATORY
COLUMBIA UNIVERSITY



**(CW) LASER** — **modulator** — **detector**

- 18 Gb/s demonstrated



*Ge-on-Si Detectors*:
- 40-GHz bandwidths
- 1 A/W responsivities

*Receivers (detectors w/ CMOS amplifiers)*:
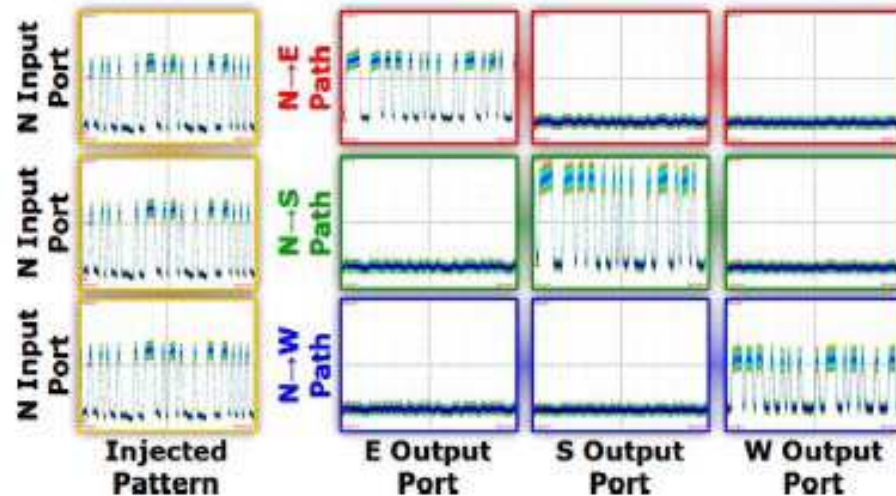- 1.1 pJ/bit demonstrated at 10 Gb/s
- Scalable to < 50 fJ/bit

[M Lipson, *Optics Express* (2007)]

- **85 fJ/bit demonstrated at 10 Gb/s**
- **Scalable to < 25 fJ/bit**



[M Watts, *Group Four Photonics* (2008)]



Ge-on-SOI photodiode

[S Koester, *J. Lightw. Technol.* (2007)]

24

Single-Channel (1546-nm) Routing Verification for Three Switch Configurations: N→E, N→S, and N→W

[A. Biberman, *IEEE Phot. Tech. Letters* (2010)]

# On-Chip Topology Exploration

- **Photonic Torus**

- **Nonblocking Photonic Torus**
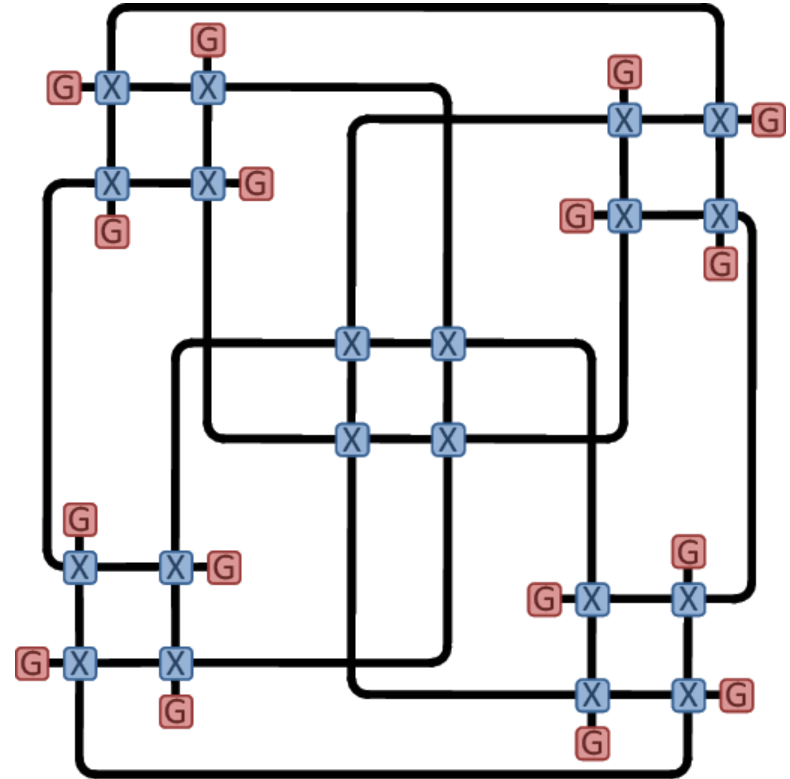


[A. Shacham *et al.,* Trans. on Comput., 2008]

[M. Petracca et al. IEEE Micro, 2008]

- **TorusNX**
- **Square Root**



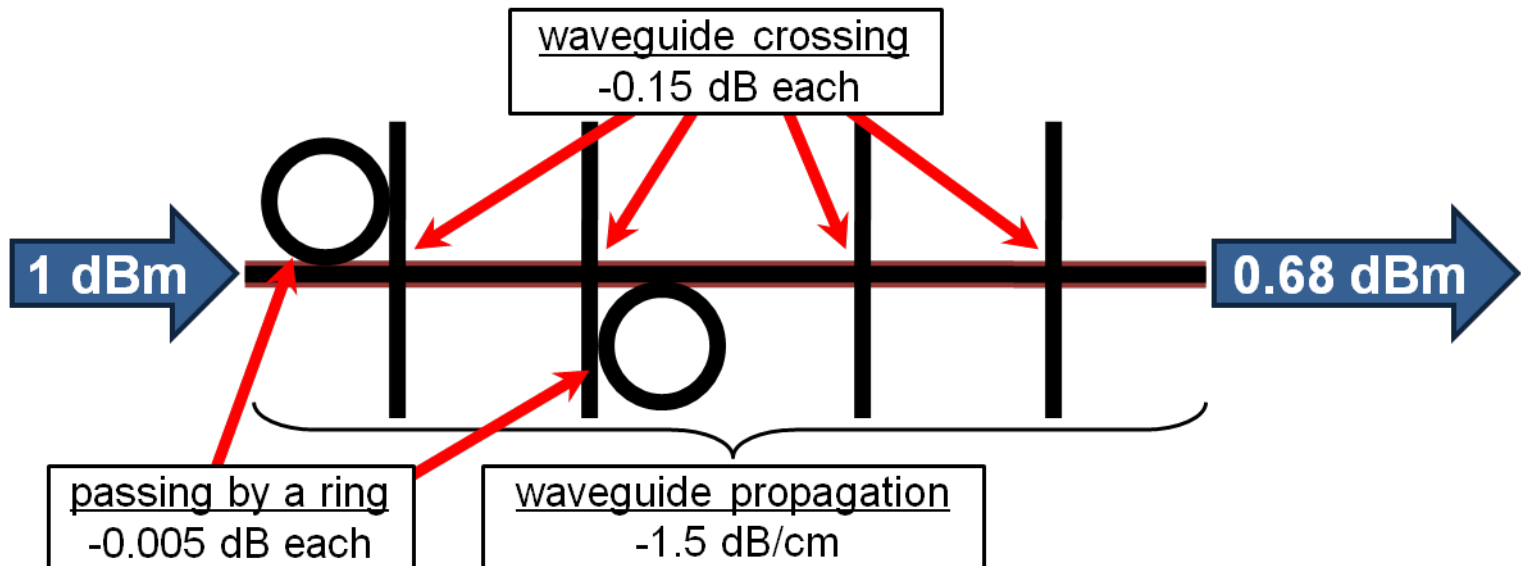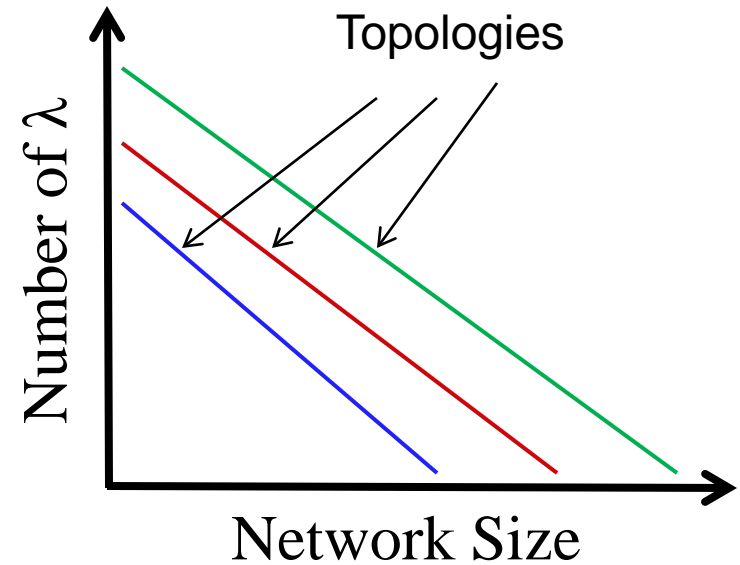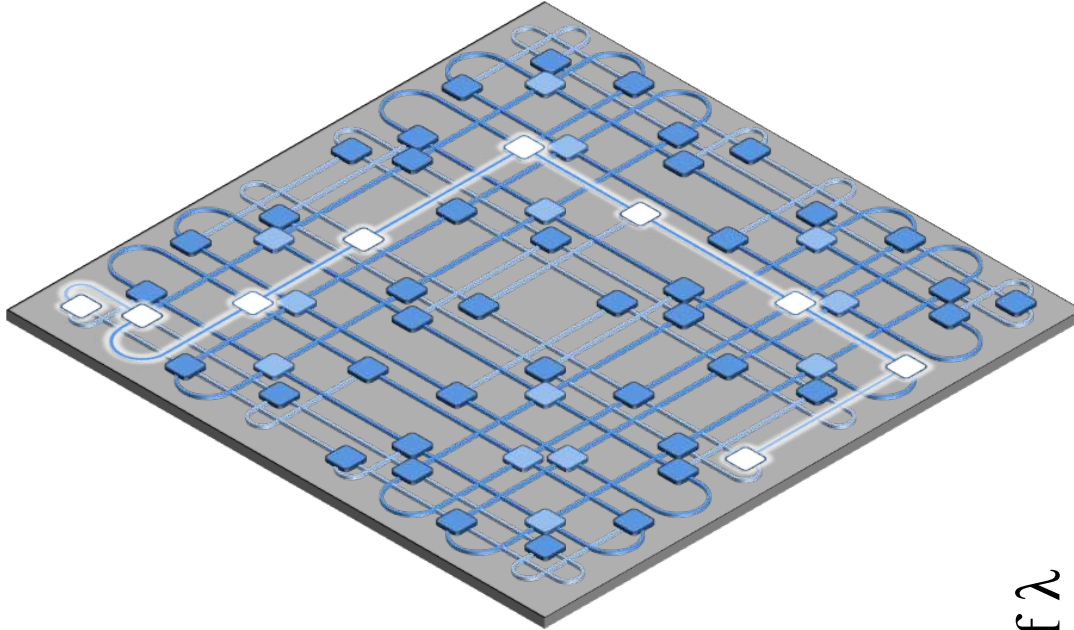[J. Chan et al.  JLT, May 2010]

# Photonic Plane Characteristics

- **Insertion Loss**
- **Noise**
- **Power**

Nonlinear Effects

Optical Power Budget

Total Power of WDM Signal

WDM Factor

Worst-case Insertion Loss

Detector Sensitivity

waveguide crossing
-0.15 dB each

1 dBm

0.68 dBm

passing by a ring
-0.005 dB each
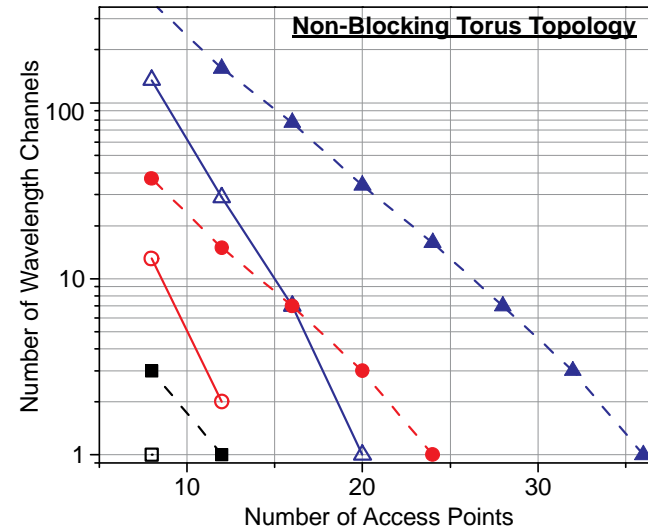
waveguide propagation
-1.5 dB/cm
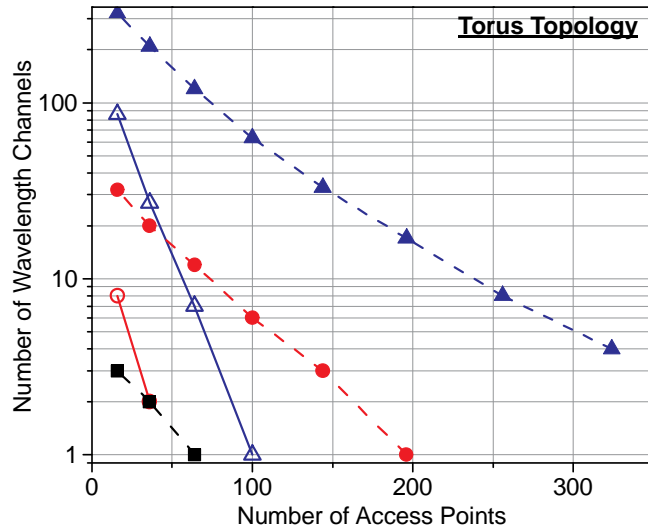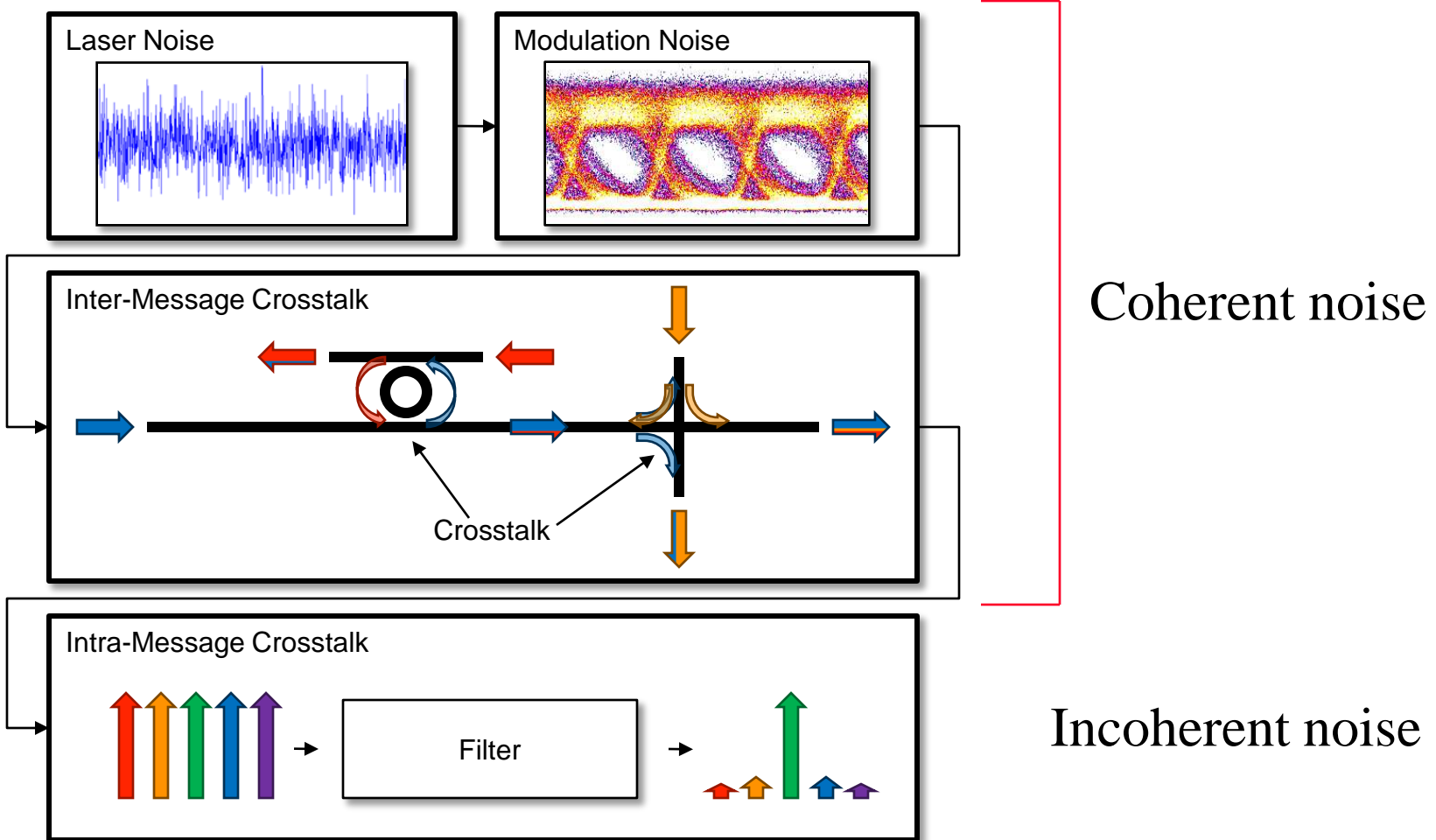
# Simulation Results

*Original* is based on the IL results from previous slide, *Improved* is based on a hypothetical improvement in crossing loss from 0.15 dB to 0.05 dB.
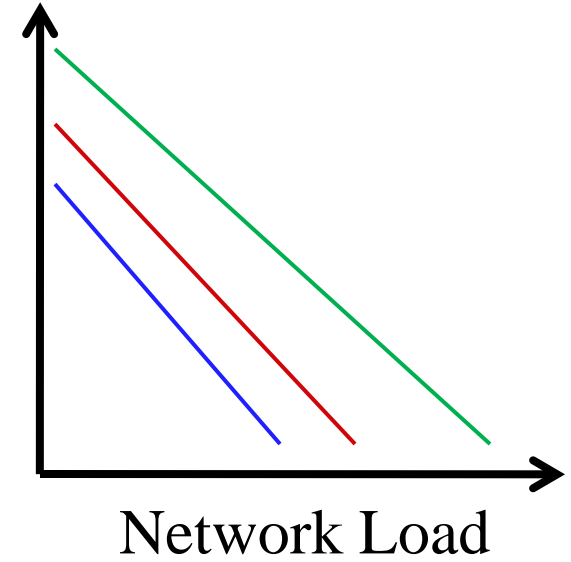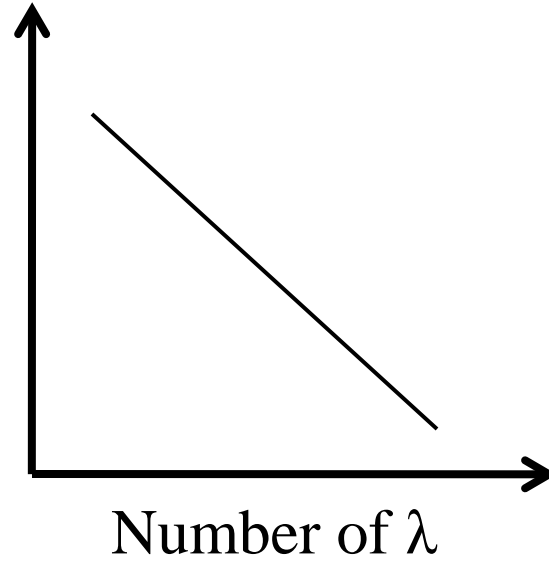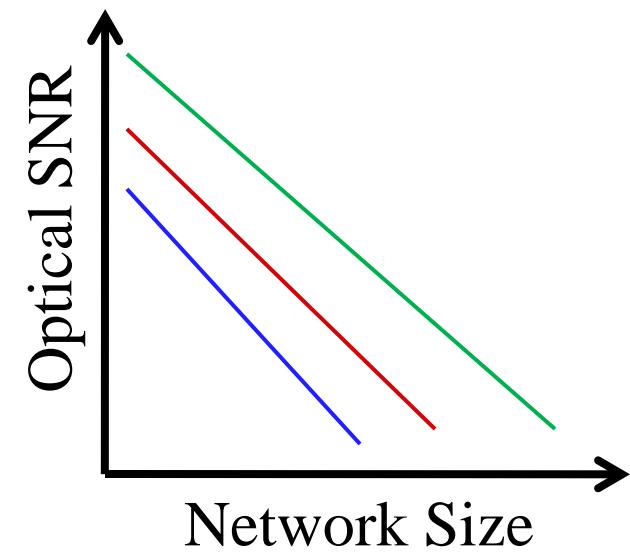
# Photonic Plane Characteristics

- **Insertion Loss**
- **Noise**
- **Power**

Laser Noise

Modulation Noise

Coherent noise

Inter-Message Crosstalk

Crosstalk
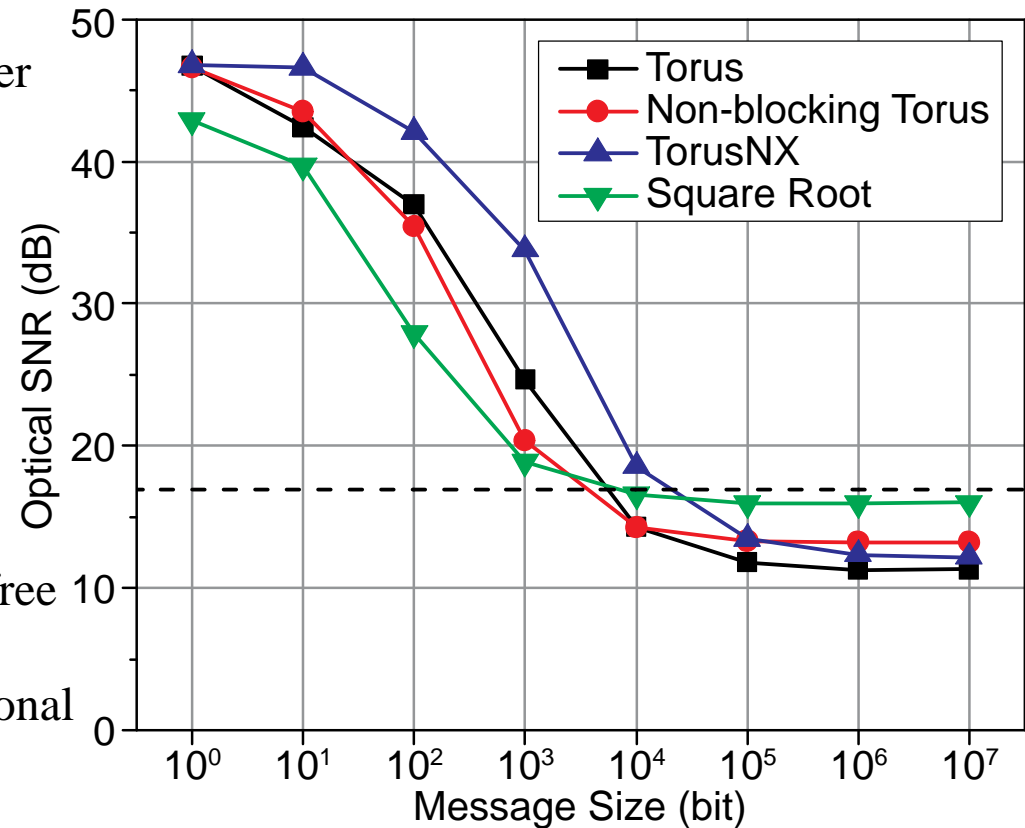
Intra-Message Crosstalk

Filter

Incoherent noise

## Results

•Results are plotted for network size of 8×8 at saturation, at the detectors.

• Maximum OSNR = ~45 dB (due to laser noise)

• Minimum OSNR < 17 dB (due to message-to-message crosstalk)

• Variations between networks due to varying likelihood of two message intersecting on network topology.

## System Performance

• SNR measures the likelihood of error-free transmission.

• Lower SNR designs will require additional retransmission, resulting in lower throughput performance.
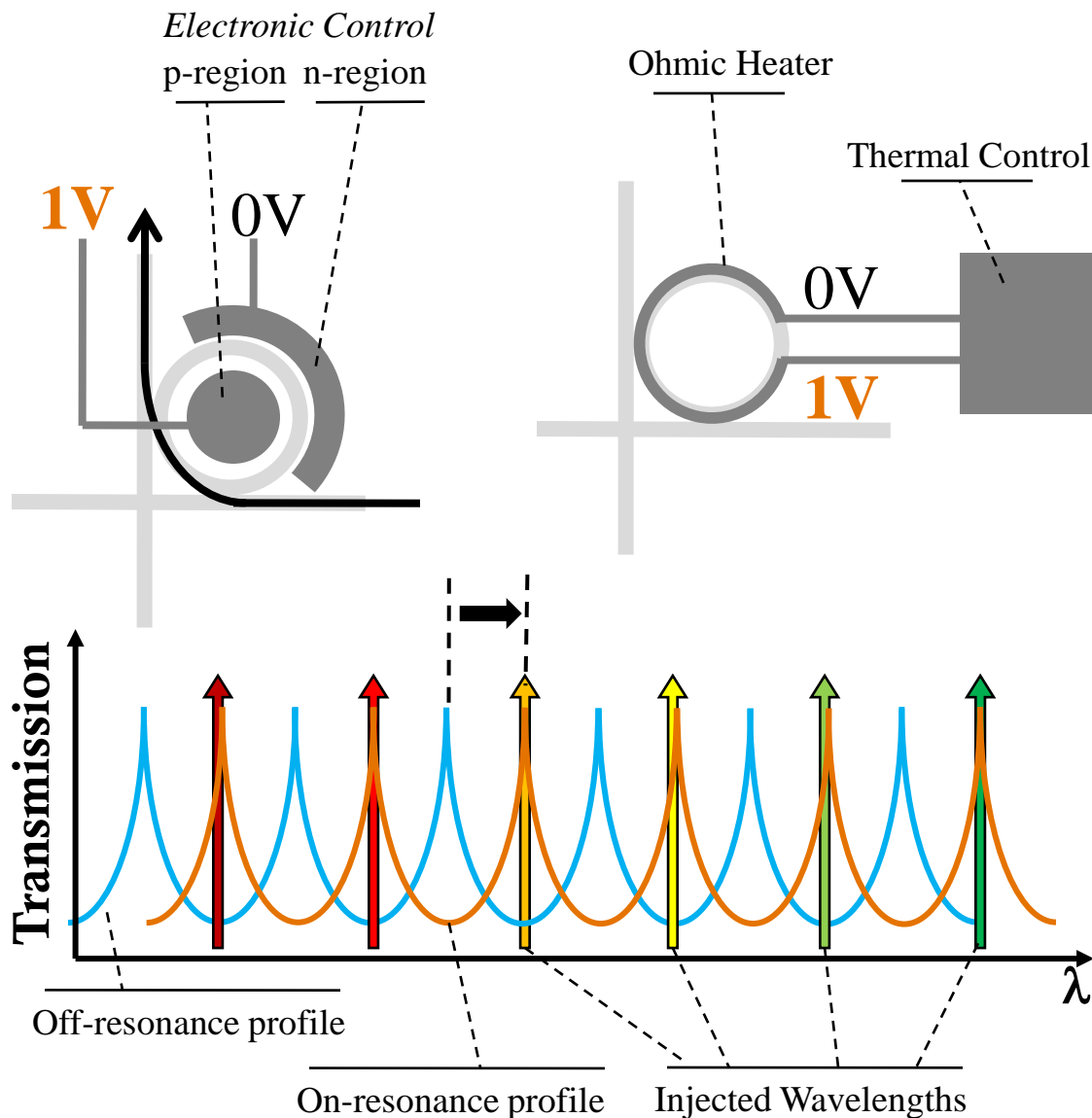


The line at OSNR=16.9 dB is where a bit-error-rate of $10^{-12}$ can be achieved, assuming an ideal binary receiver circuit and orthogonal signaling.
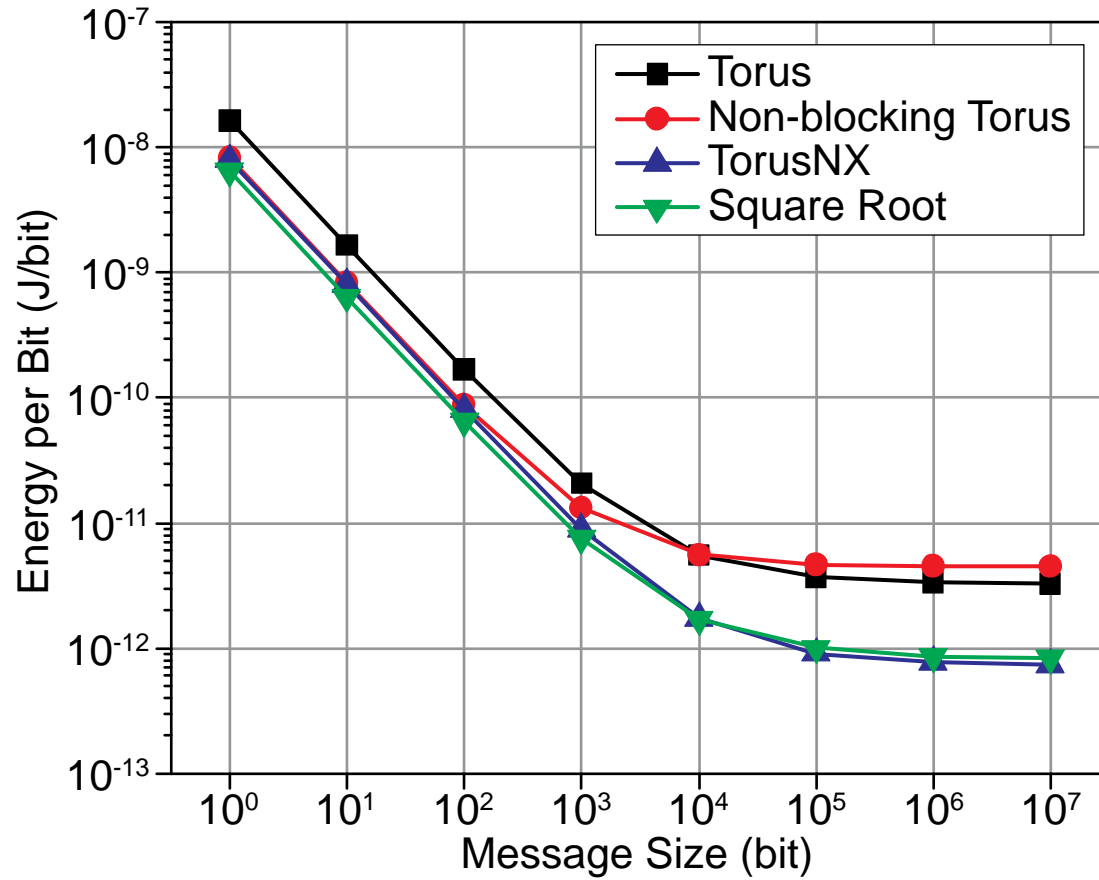
# Photonic Plane Characteristics

- **Insertion Loss**

- **Noise**

- **Power**

- **Laser Power**
- **Active Power**
  - **Modulating**
  - **Detecting**
  - **Broadband**
- **Static Power**
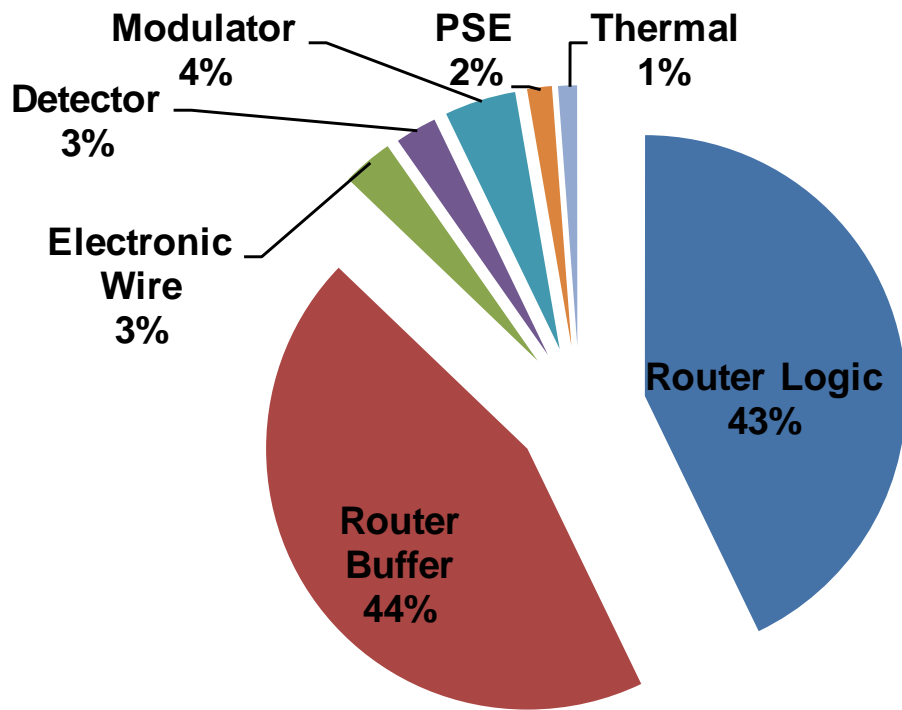  - **Thermal tuning**
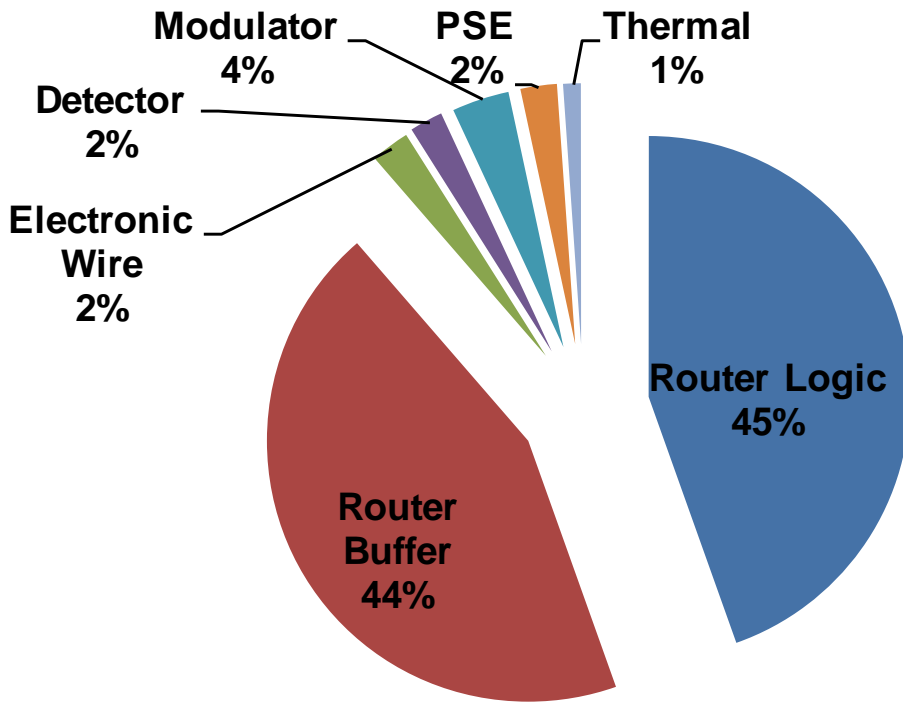- **Tx\Rx Power**
  - **Drivers**
  - **TIAs**



*Electronic Control*
p-region  n-region

**1V**  0V

Ohmic Heater

Thermal Control

0V

**1V**

Transmission

λ

Off-resonance profile

On-resonance profile  Injected Wavelengths

# Energy Per Bit

Torus Topology

Modulator
4%

PSE
2%

Thermal
1%

Detector
3%

Electronic
Wire
3%

Router Logic
43%

Router
Buffer
44%

- 12 wavelengths @ 10 Gbps/each
- Power Dissipation = 4.31 W

Nonblocking Torus Topology

Modulator
4%

PSE
2%

Thermal
1%

Detector
2%

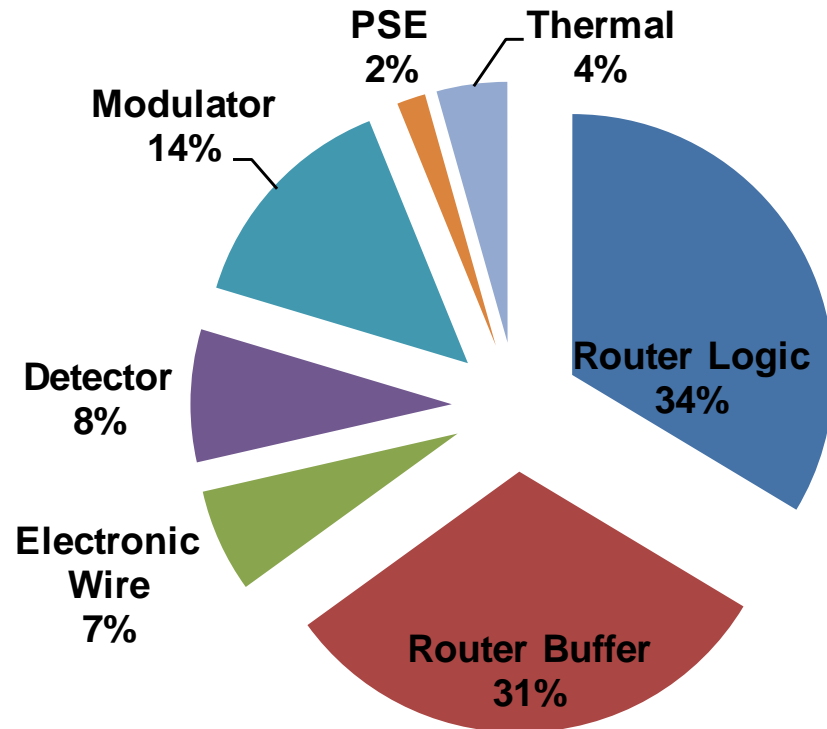Electronic
Wire
2%

Router Logic
45%

Router
Buffer
44%

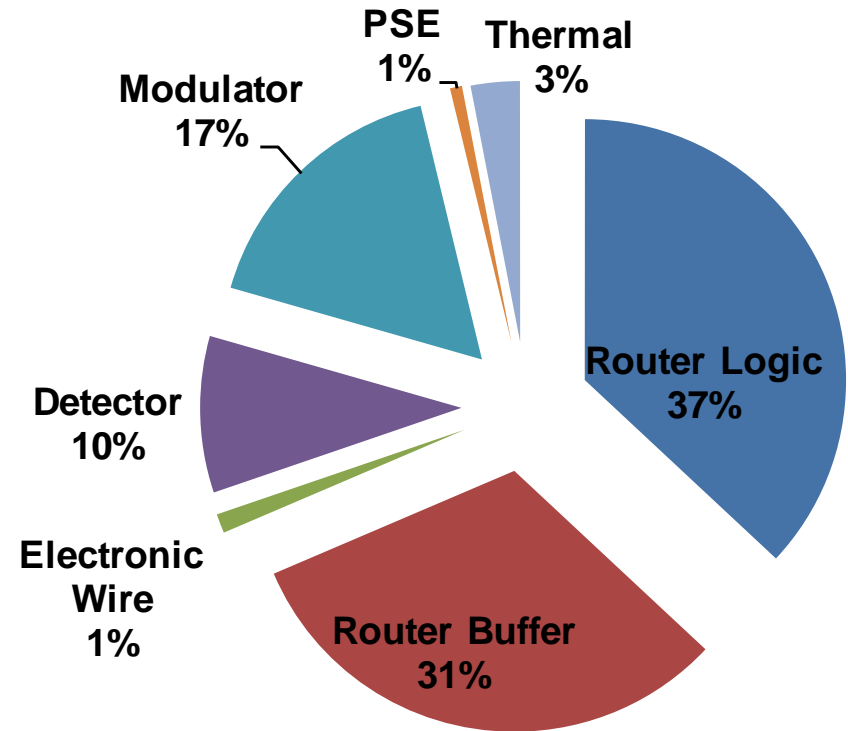- 7 wavelengths @ 10 Gbps/each
- Power Dissipation = 1.59 W

- Results based on randomly generated traffic with message sizes of 100 kbit, with network in saturation.
- Data was collected on 64 nodes topologies constrained to a total surface area of 2 cm × 2 cm.
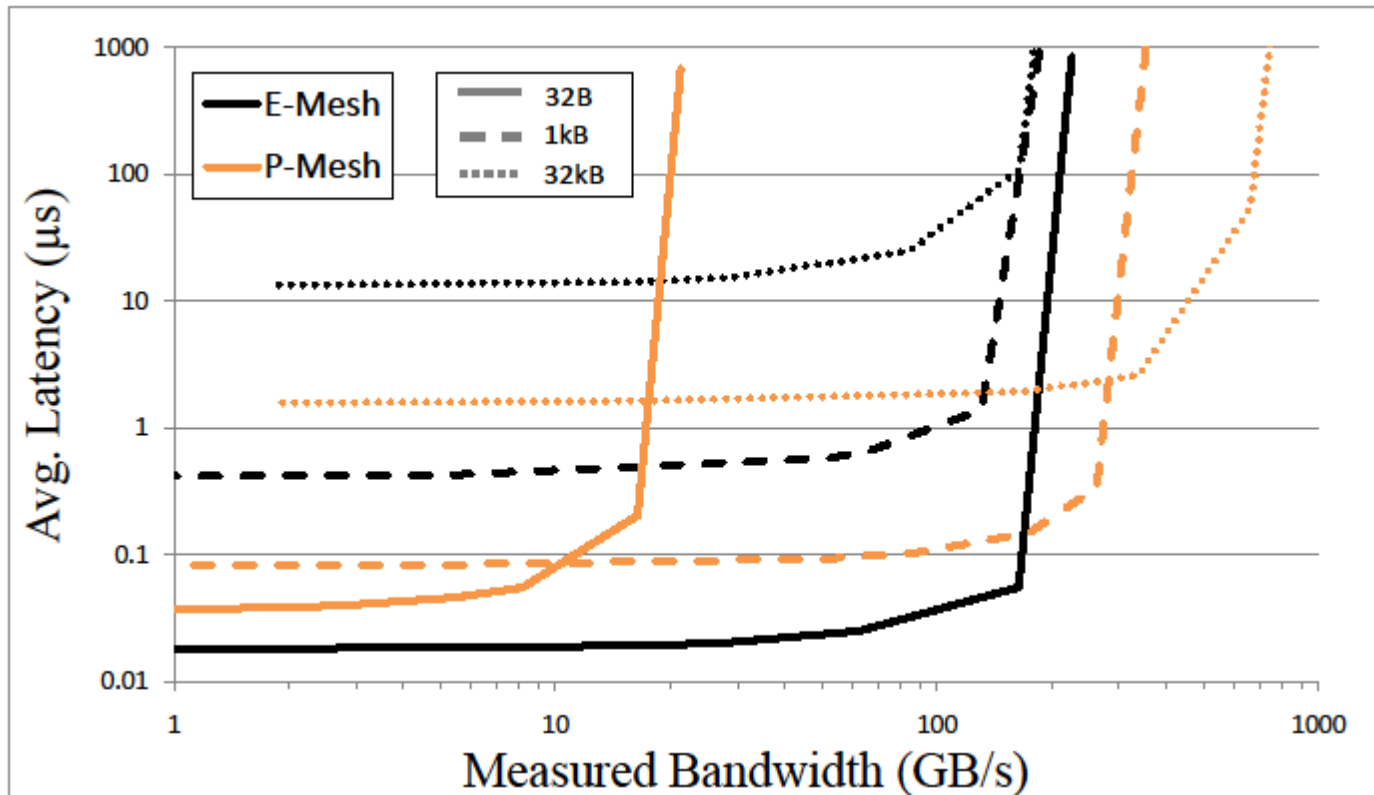
# Power Breakdown



**Square Root Topology**

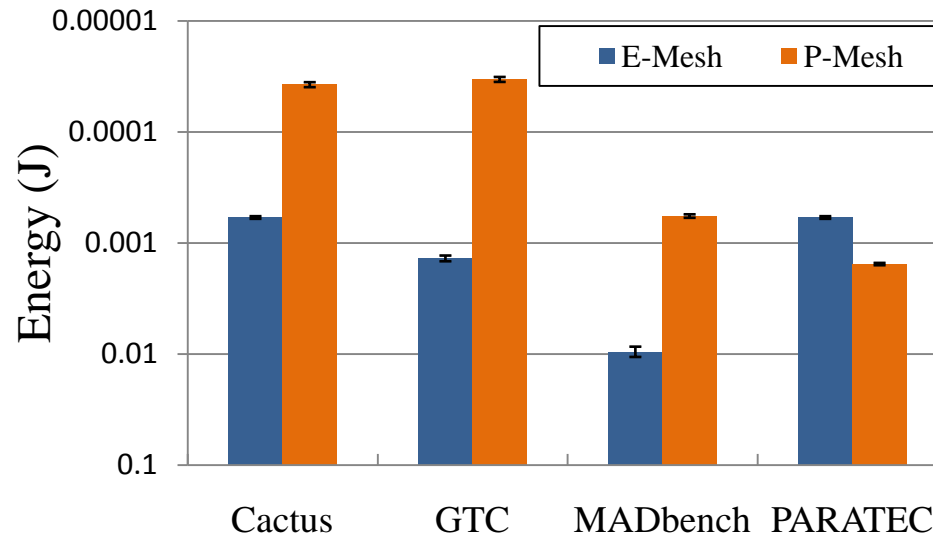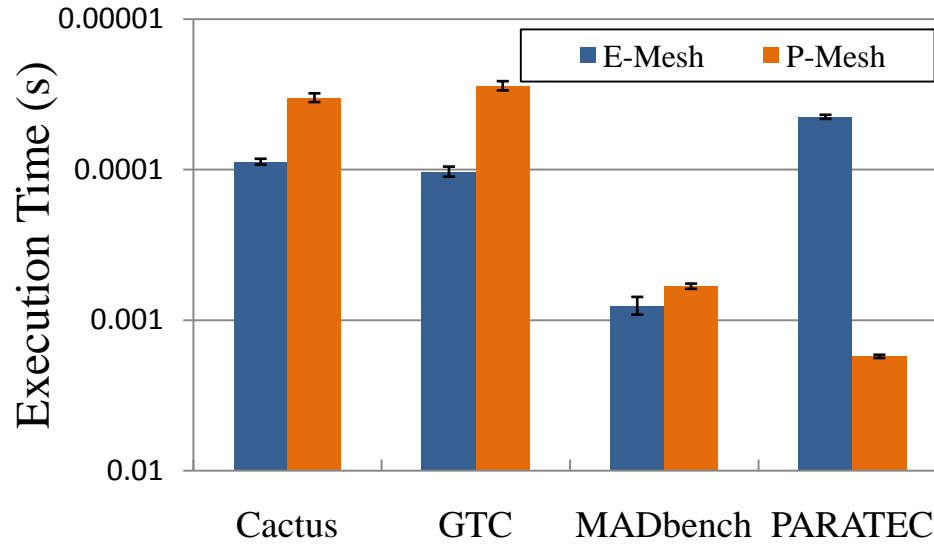- 27 wavelengths @ 10 Gbps/each
- Power Dissipation = 1.89 W

Router Logic 34%
Router Buffer 31%
Electronic Wire 7%
Detector 8%
Modulator 14%
PSE 2%
Thermal 4%

**TorusNX Topology**

- 38 wavelengths @ 10 Gbps/each
- Power Dissipation = 3.22 W

Router Logic 37%
Router Buffer 31%
Electronic Wire 1%
Detector 10%
Modulator 17%
PSE 1%
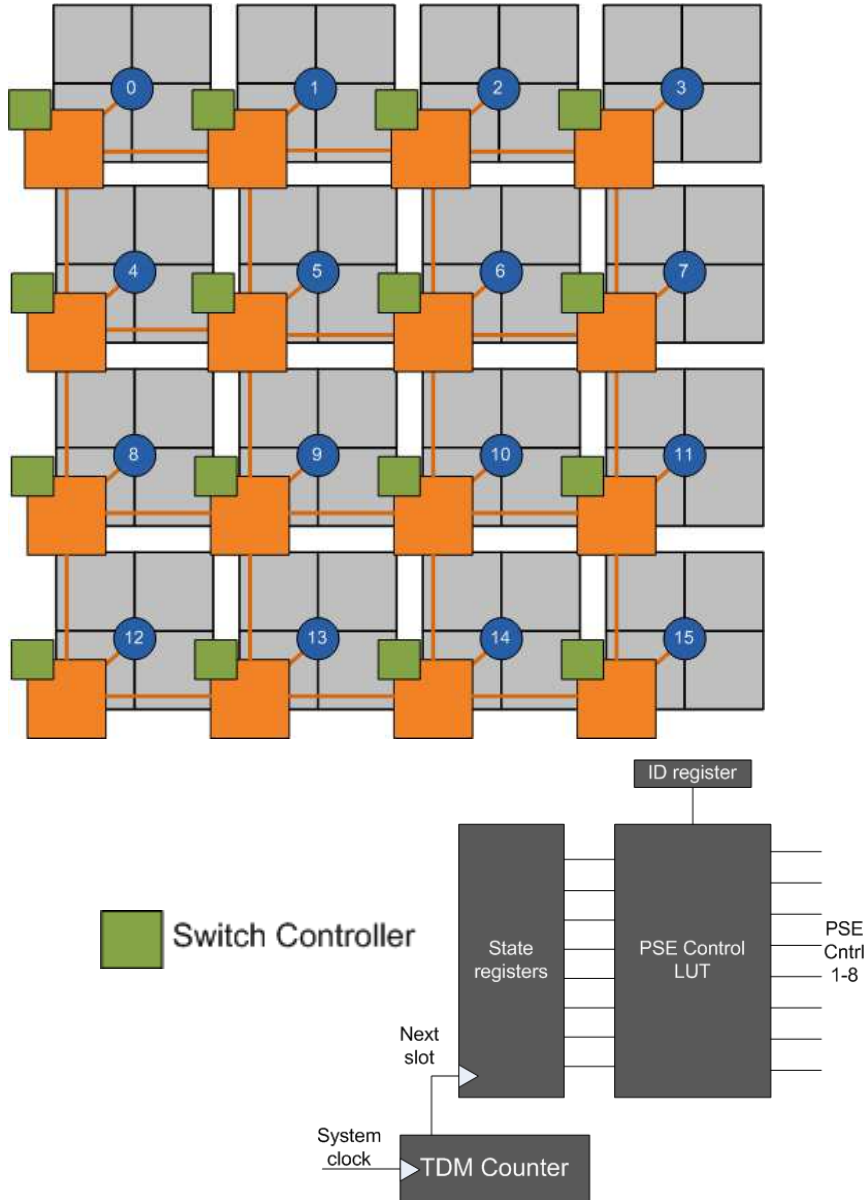Thermal 3%

# Performance

- Uniform random traffic
- 256 cores, 64-node network

# Other Interesting Issues

# Memory Access



Processor Core

Network Router

Memory Access Point

[G. Hendry et al. *Circuit-Switched Memory Access in Photonic Interconnection Networks for HPEC*. In Supercomputing, Nov. 2010]

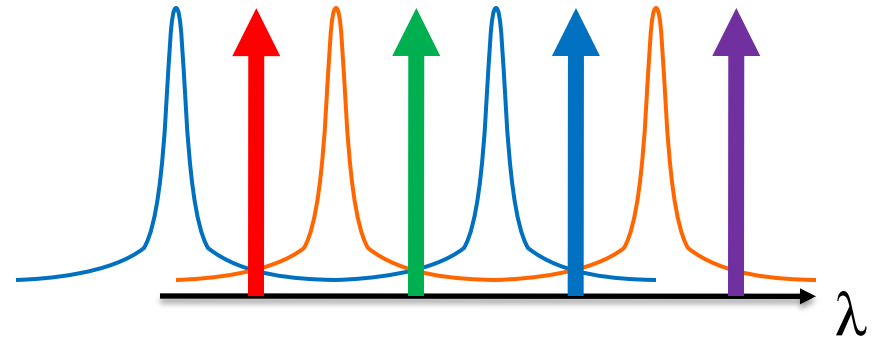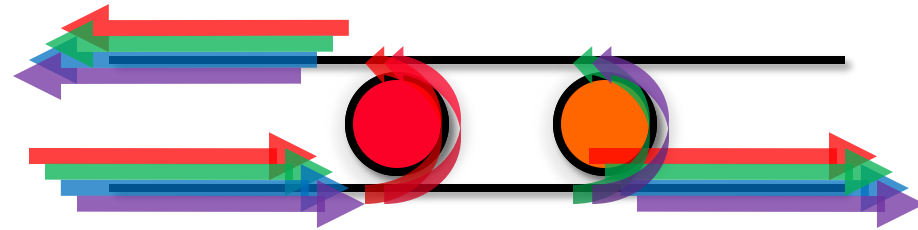[G. Hendry et al. *Silicon Nanophotonic Network-On-Chip Using TDM Arbitration*. In HOTI, Aug. 2010]

- **Original**

- Re-design



- Scalable number of WDM channels

# Conclusion

- **Some applications / programming models definitely well-suited to a circuit-switched photonic network**

- **Interesting tradeoffs and design space**
  - **Photonic physical layout / design**
  - **System-level benefits from device improvement**
  - **Network-level improvements**