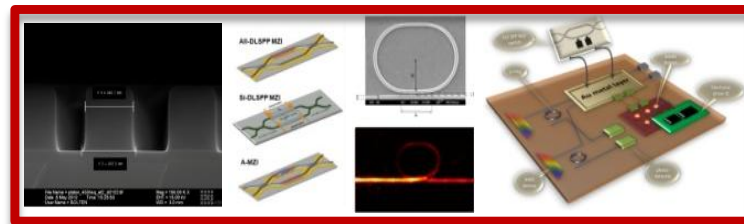




# *High-Speed Routing and Switching in Optical Communications & High-Performance Computing Systems*



*Dr. Nikos Pleros*

*Photonics Systems & Networks (PhosNET) Research Group  
Dpt. of Informatics, Aristotle University of Thessaloniki*

# Outline

---

- ☑ *HPCS systems today: status and challenges*
- ☑ *Routing in HPC systems*
- ☑ *Optics for Routing in HPC*
- ☑ *Tb/s Si-Plasmonic Routers*
- ☑ *Optical RAM*

# HPC examples..and metrics

## No.1: Jaguar (USA)



## No.2: Nebulae (China)



Has light anything to do with computing ?

❏ **1.75 Petaflop/s**

( 1 PF =  $10^{15}$  calculations per sec)

❏ **410 m<sup>2</sup> floor space**

❏ **7 MW power consumption !!**

❏ **1.271 Petaflop/s**

❏ **Total 120640 cores**

❏ **2.25 MW power cons.**

❏ **relies on BladeSystem**

# ...a look inside

## IBM's Roadrunner architecture

- ☞ *18x Connected Units*
- ☞ *270x Racks*



- ✓ *actually a small-range network*
- ✓ *...with 1.04 Pflop/sec and 384 Gb/s intra-CU traffic*
- ✗ *...and 2.5 MW power consumption !*

**...and here comes light in**



***use optical fiber for the interconnection***

✓ *...and enable Tb/s transmission speeds*

***...is there any other problem ?***

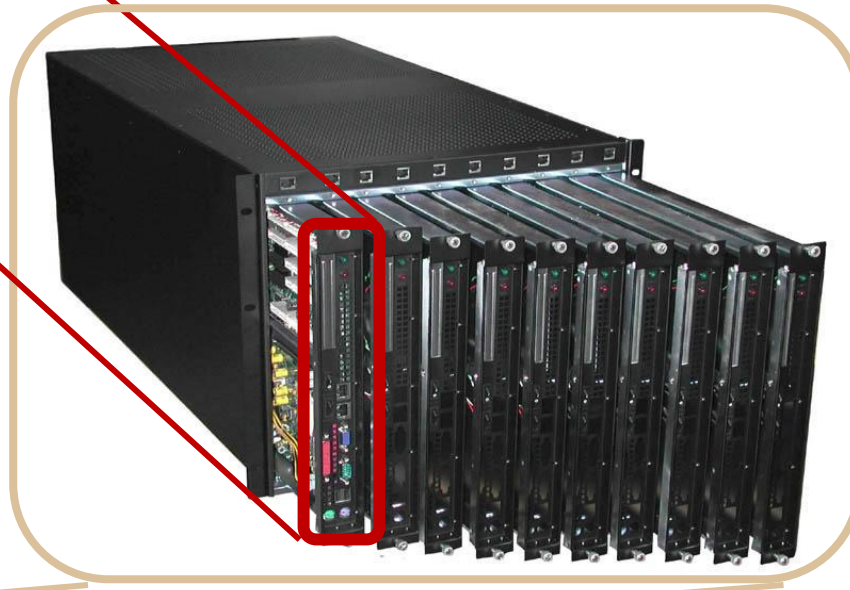
- ❖ *size and cable length*
- ❖ *ultra-small latency required*
  - *for fast and low-complexity parallelization*
- ❖ *power consumption...in MWs !!*
  - *consumes what a small plant can produce !!*



# BladeCenters: a solution ?

*HPC architecture supported by IBM*

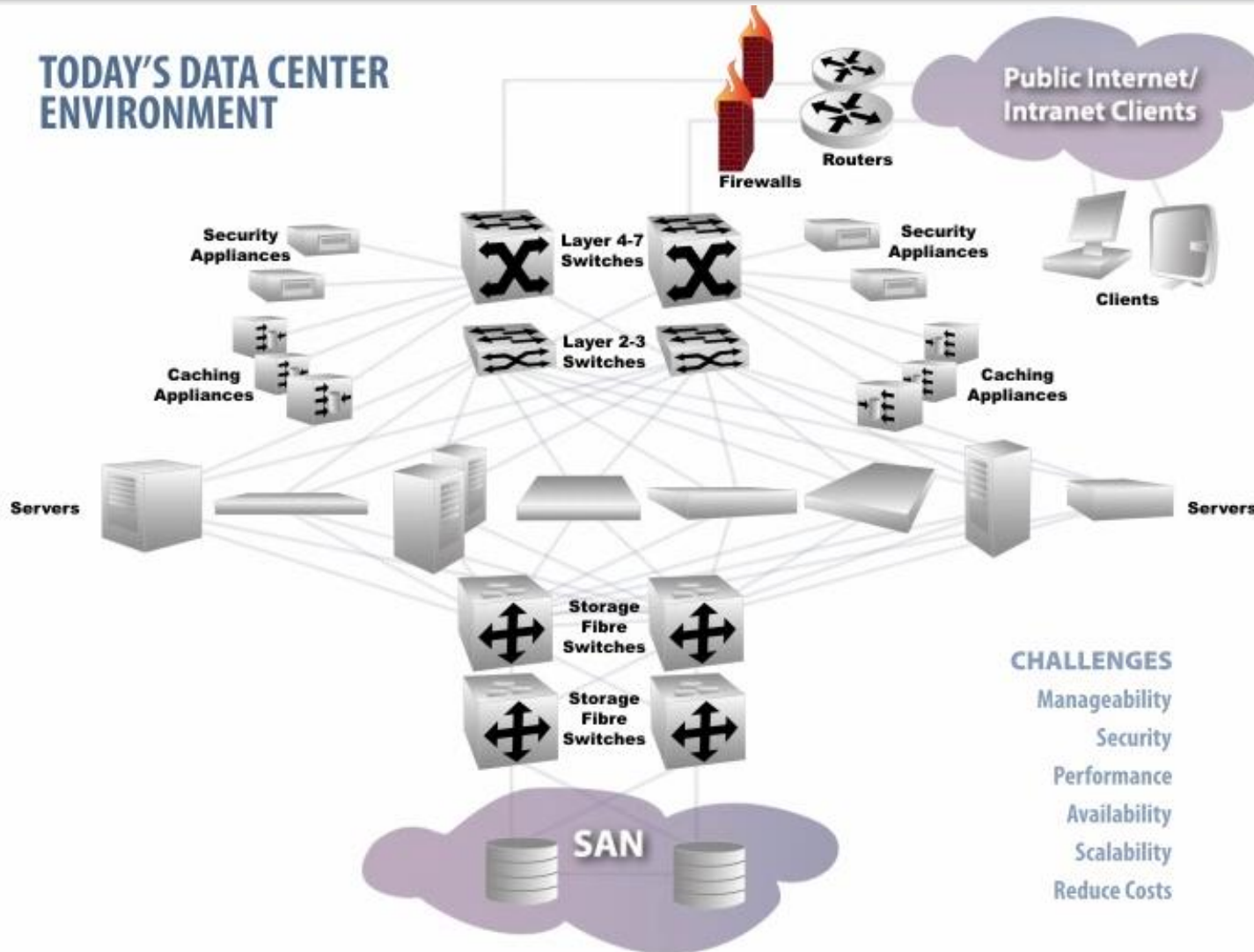
**Blade server** is a stripped down server computer, for minimizing physical space and energy requirements



**Blade enclosure**, hosts multiple blade servers, provides power, cooling, networking, interconnects & management

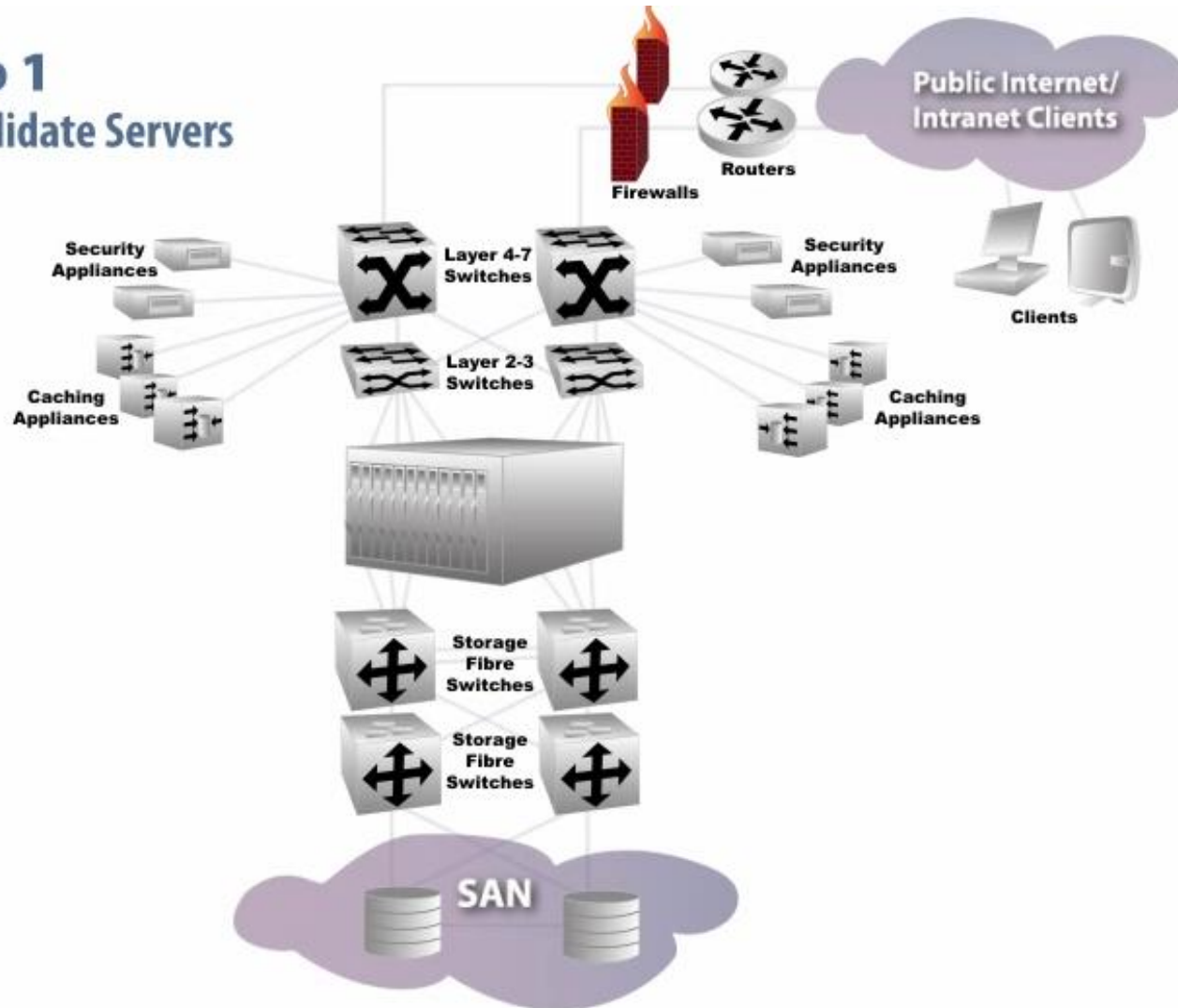
# BladeCenters: The vision

## TODAY'S DATA CENTER ENVIRONMENT



# BladeCenters: The vision

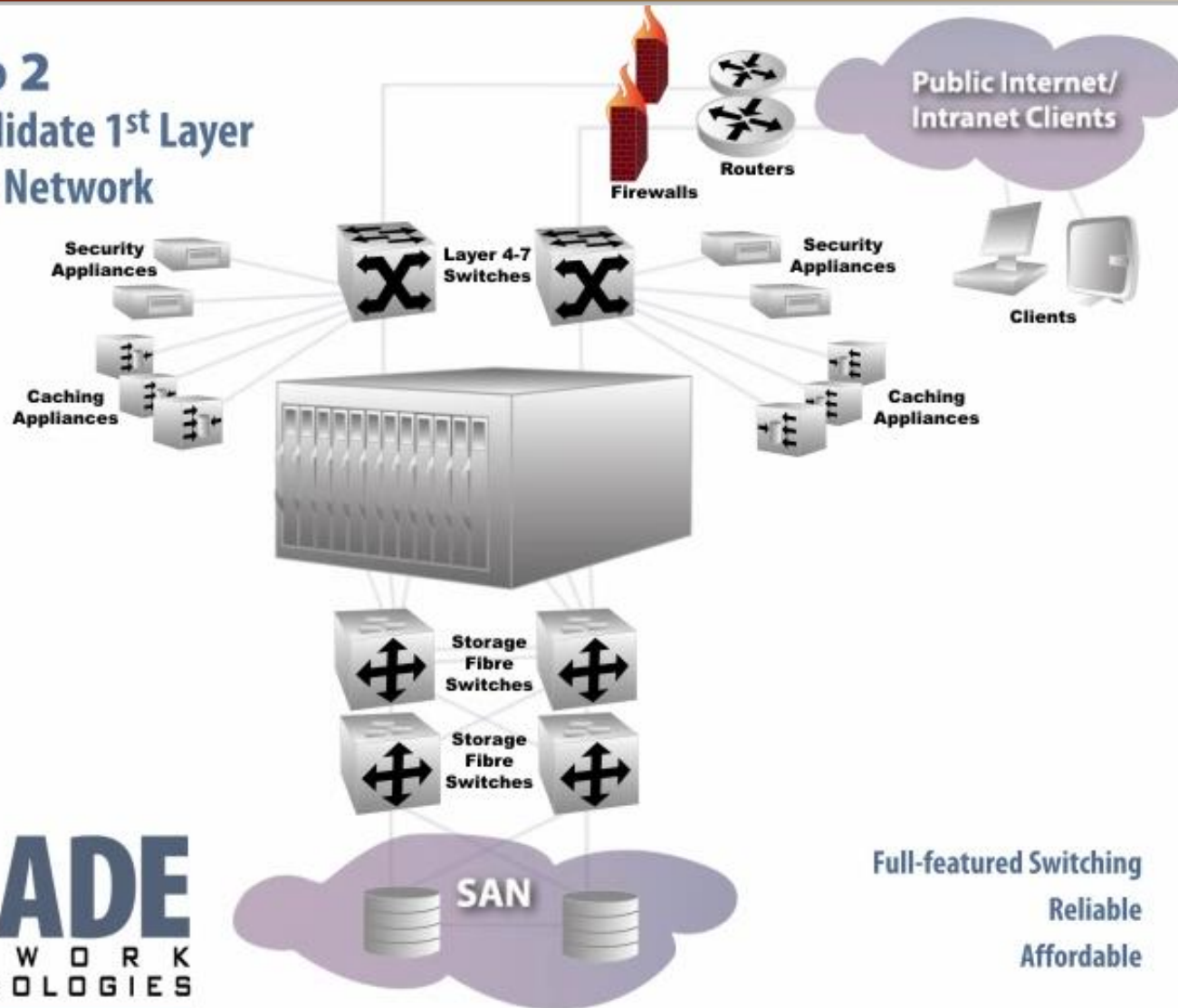
## Step 1 Consolidate Servers





# BladeCenters: The vision

## Step 2 Consolidate 1<sup>st</sup> Layer of the Network

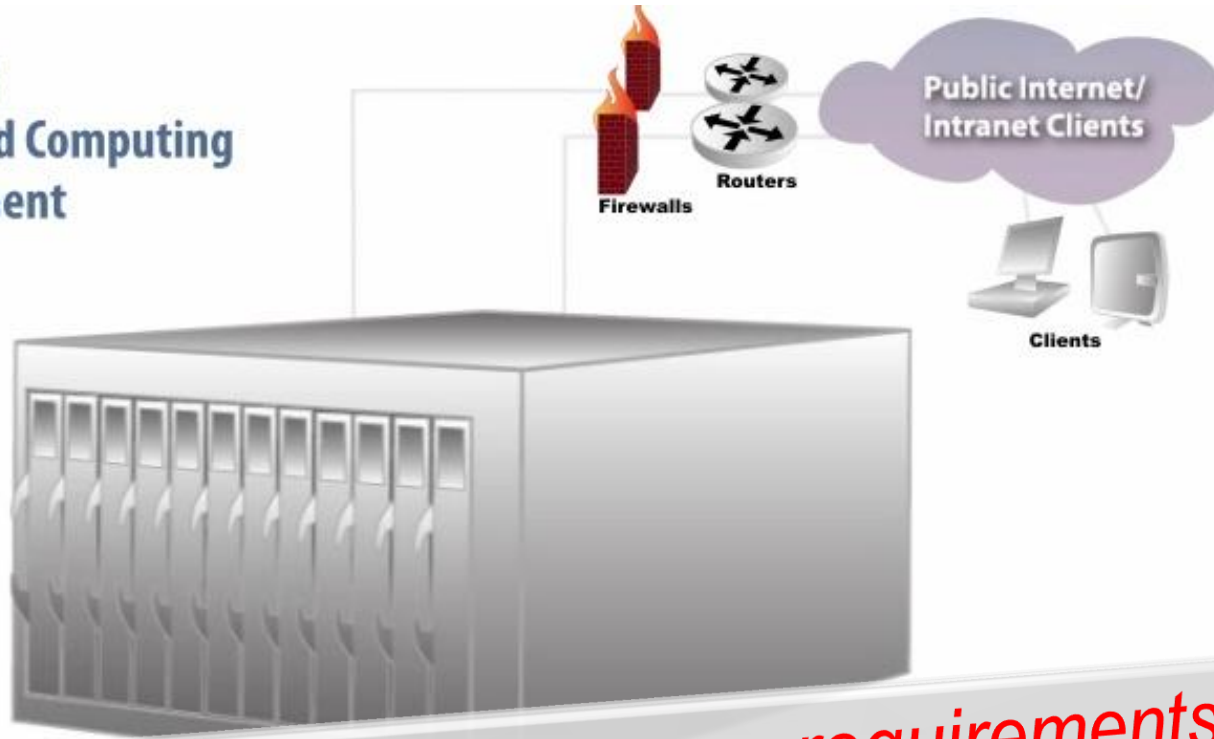


**BLADE**  
NETWORK  
TECHNOLOGIES

Full-featured Switching  
Reliable  
Affordable

# BladeCenters: The vision

**Result**  
Simplified Computing  
Environment



*reduces power and space requirements  
through sharing of resources within the blade*

NETWORK  
TECHNOLOGIES

# BladeCenters: a solution ?

---

- # *creates large aggregate traffic...*
- # *100's of Gb/s in miniature networks !*

- ***The Question: How to route this ?***
  - *...in consolidated network environment*
  - *...at inter-, intra-blade, backplane level*
  - *without consuming most of the blade power*

# A new framework for photonics



**Wide Area Networks**

End 80's – early 90's



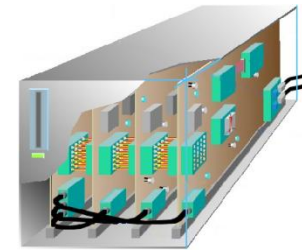
**LANs**

...early 2000

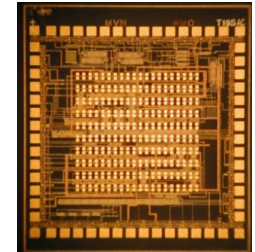


**rack-to-rack**

...now



**Backplane & chip-to-chip**



**On-chip**

**Network dimensions**

1000 km

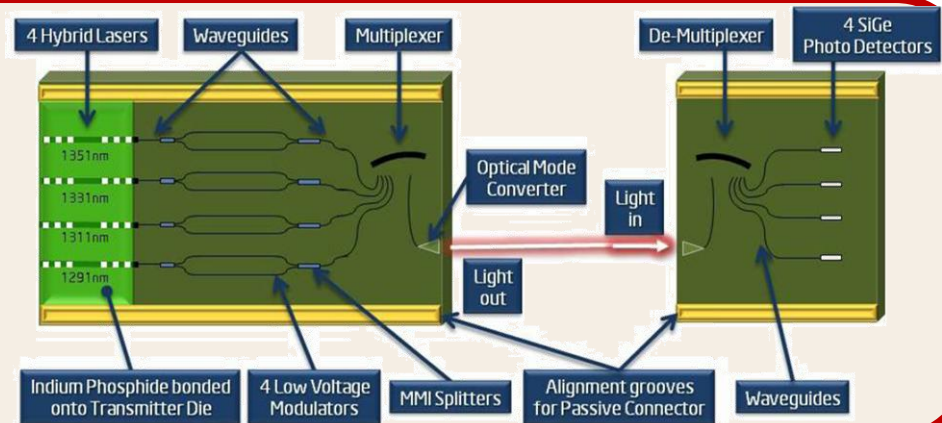
1 m

1 cm

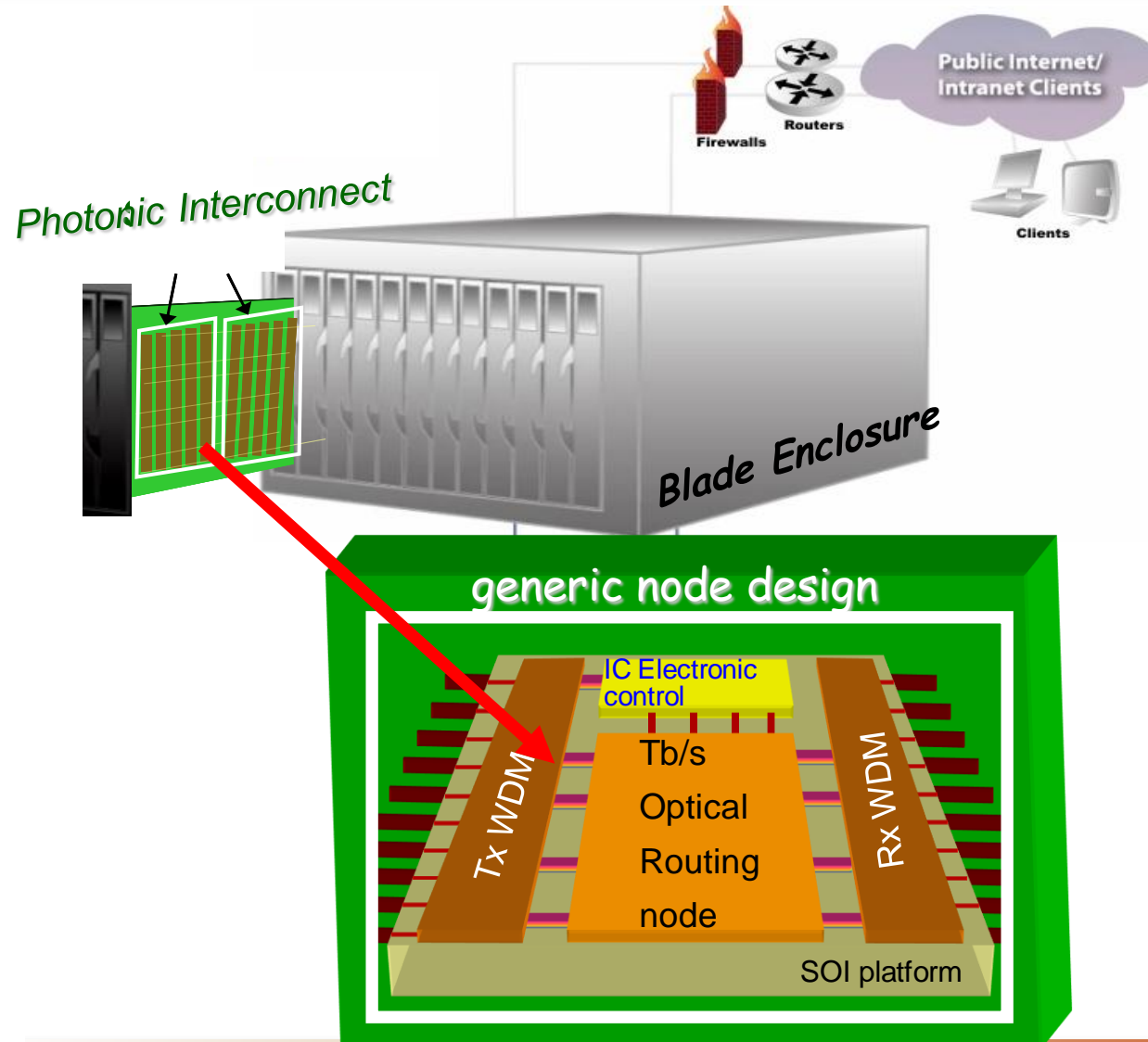
1 mm

## ...and a new roadmap

- **Silicon Photonics integration platform**
- **Recent example: 50Gb/s optical bus (Intel USA, 2010)**



# Need for chip-scale routers

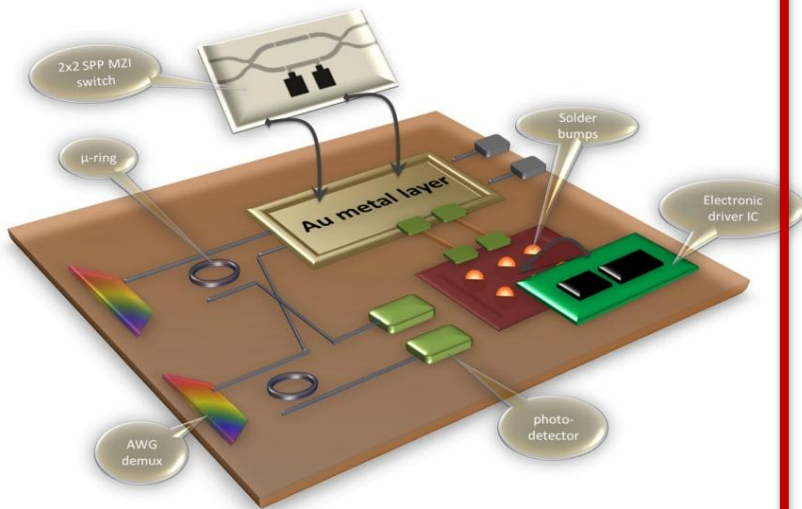


# Tb/s optical routers on-chip

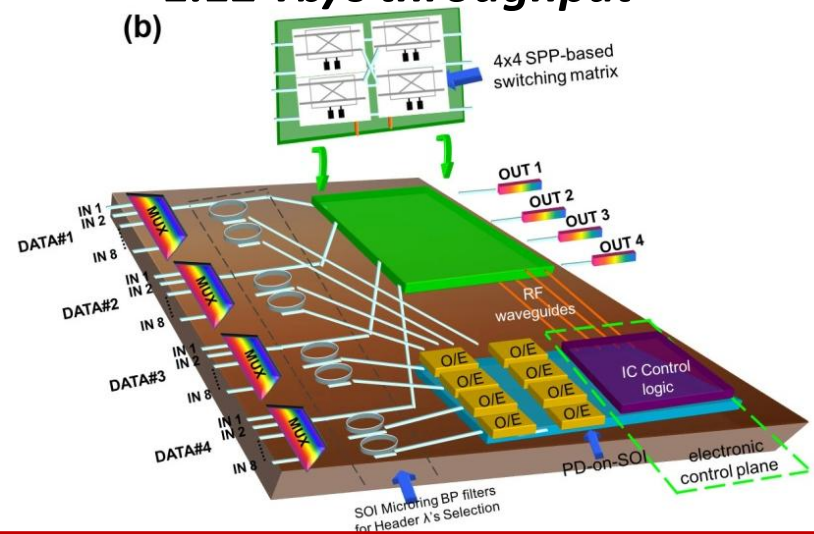
- integrate plasmonics and silicon photonics platforms
- demonstrate integrated Tb/s routers:
  - ✓  $\text{mm}^2$  footprint
  - ✓ a few Watts power consumption



**2x2 Router**  
**560 Gb/s throughput**



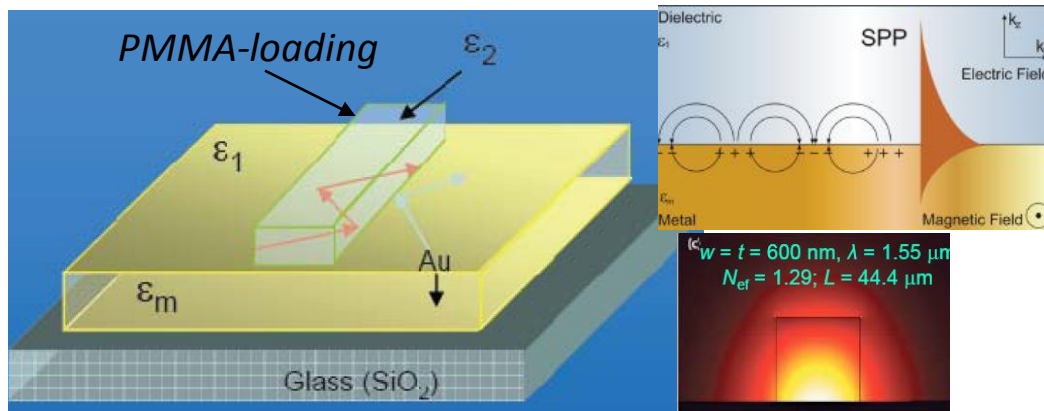
**4x4 Router**  
**1.12 Tb/s throughput**



# Plasmonics for switching

## Dielectric-Loaded Surface Plasmon Polaritons

*polymer strip (PMMA) on top of Au film*



- ☑ EM waves guided at the dielectric-gold interface
- ☑ small footprint (500x600nm waveguide dimensions)

- ✓ appropriate for interfacing photonics and electronics
- ✓ allows for thermo-optic-induced switching phenomena
- ✓ low switching power consumption (few mWs)
- ✗ ...but high propagation losses  $L_{\text{prop}} \sim 45 \mu\text{m}$  (while  $L_{\pi} \sim 90 \mu\text{m}$ )

# 4x4 Si-Plasmonic Router

## Technology & Architecture

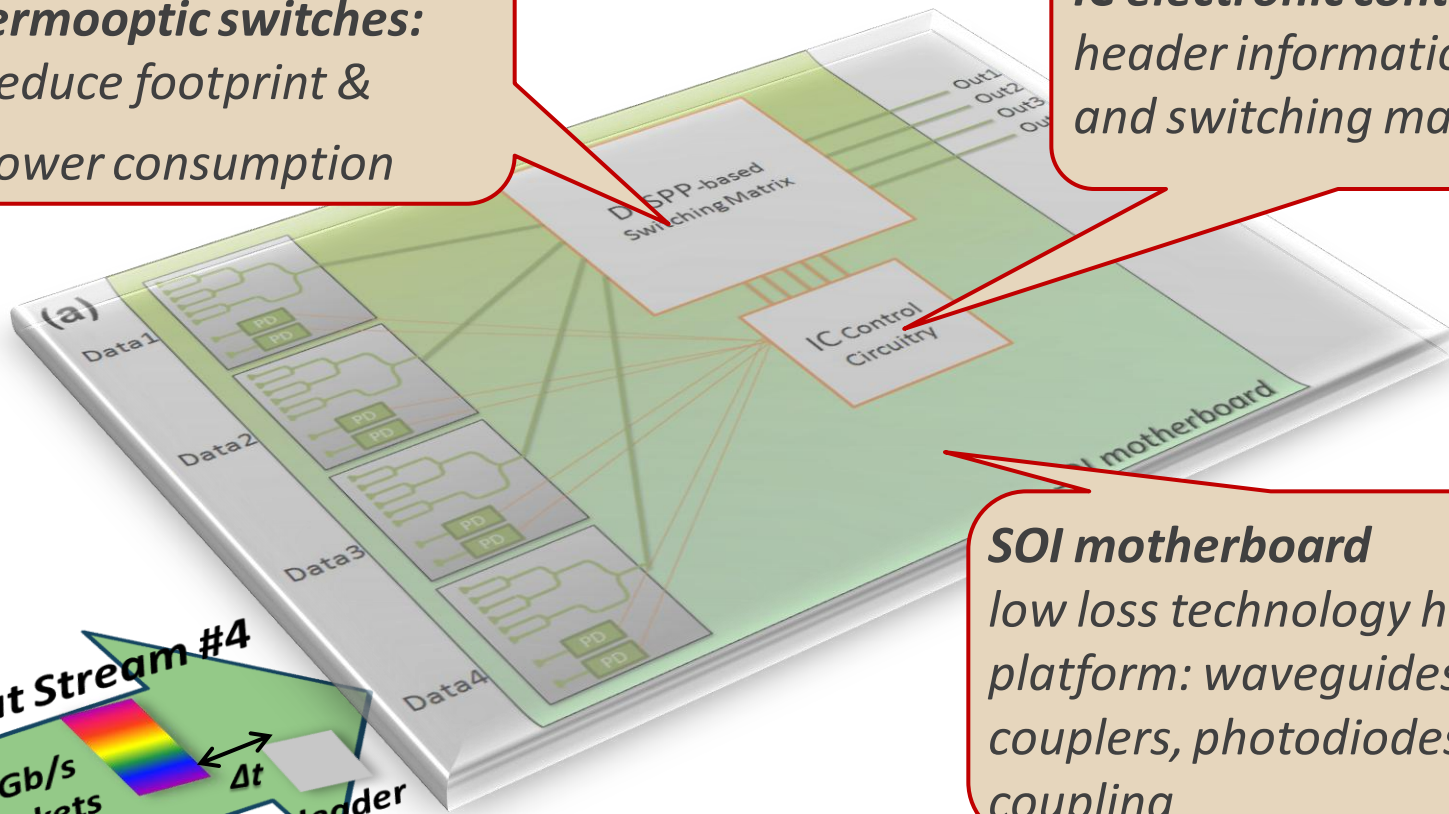
*2x2 / 4x4 plasmonic thermo-optic switches:*

- ✓ reduce footprint &
- ✓ power consumption

*IC electronic control circuit header information processing and switching matrix control*

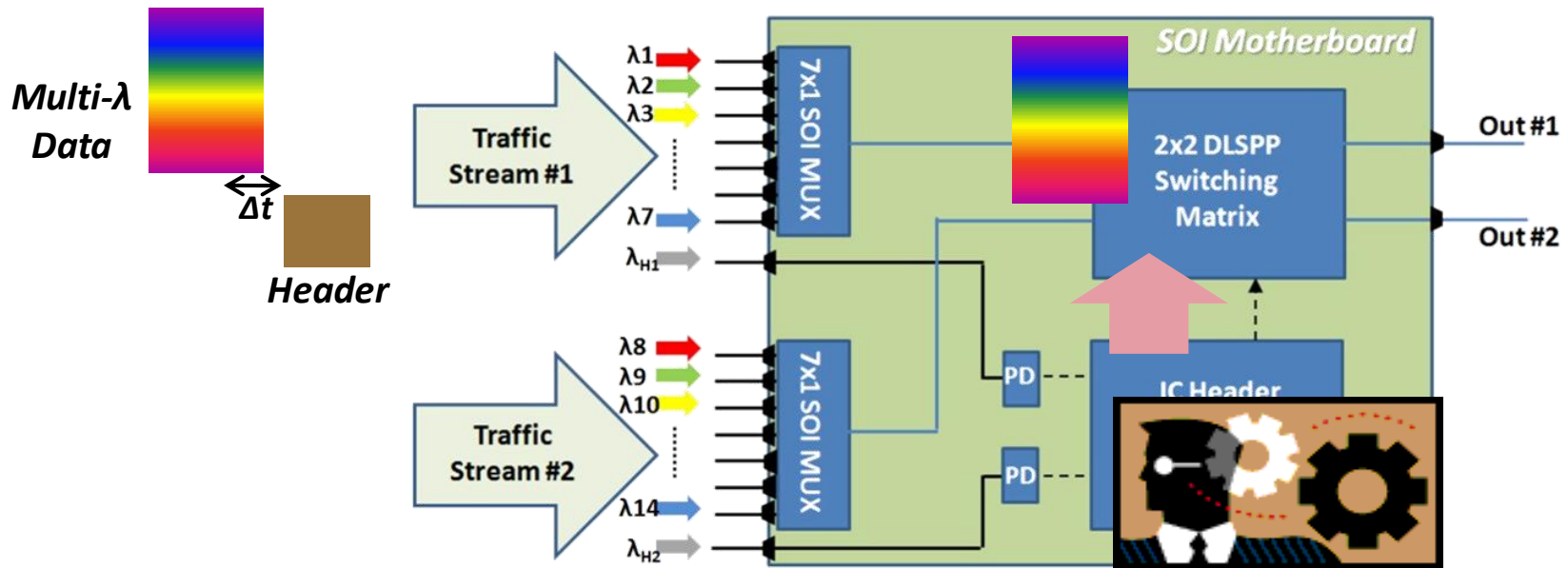
*SOI motherboard low loss technology hosting platform: waveguides, MUX, couplers, photodiodes, fiber coupling*

*Input Stream #4*  
*7- $\lambda$  40Gb/s data packets*  
*Header*  
 *$\Delta t$*



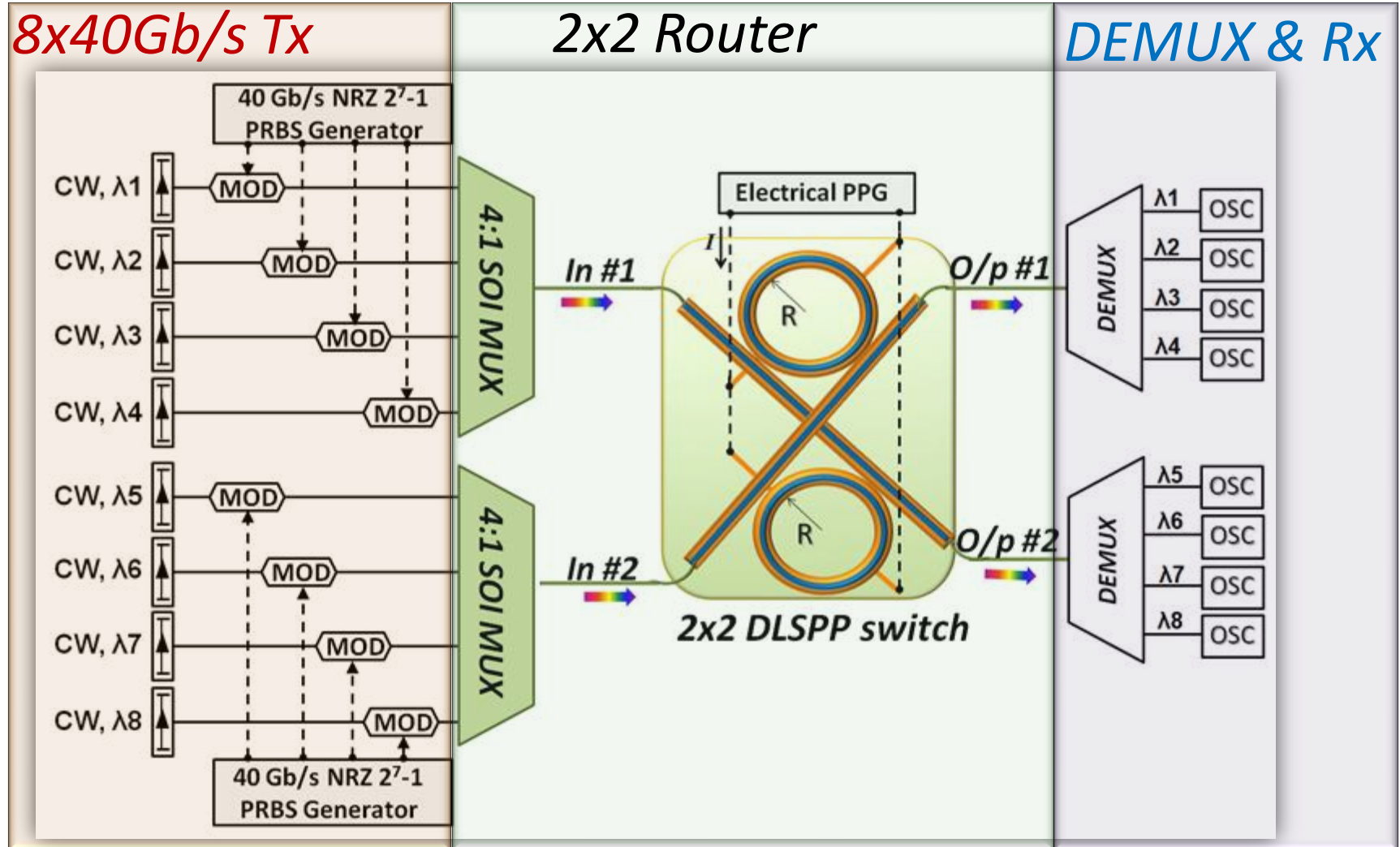


# Si-Plasmonic Router



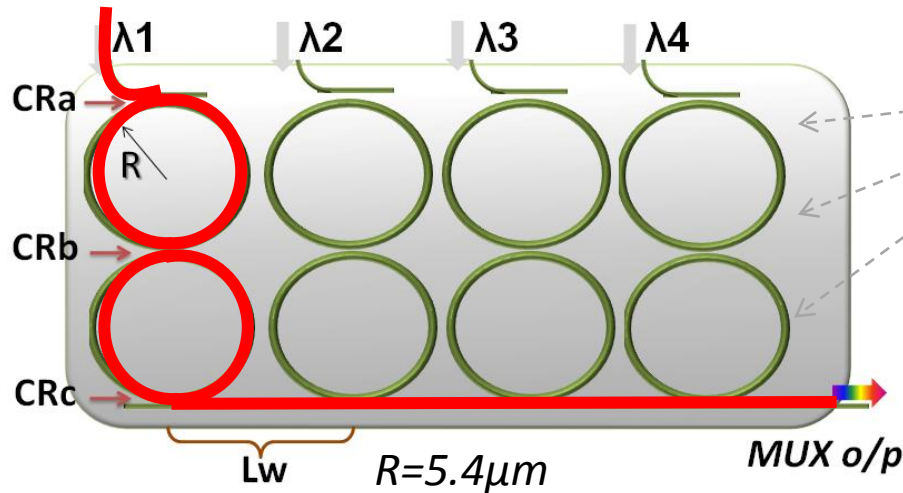
- 7x $\lambda$  data packets at 40Gb/s : 280 Gb/s per input port
- 1 extra wavelength for header (MHz data pulses)
- Time-offset between Header and Payload* information for ensuring header processing in the IC (burst-mode network concept)

# A 320Gb/s 2x2 architecture



# 40Gb/s NRZ 4:1 SOI MUX

✓ 4 cascaded 2<sup>nd</sup> order silicon rings

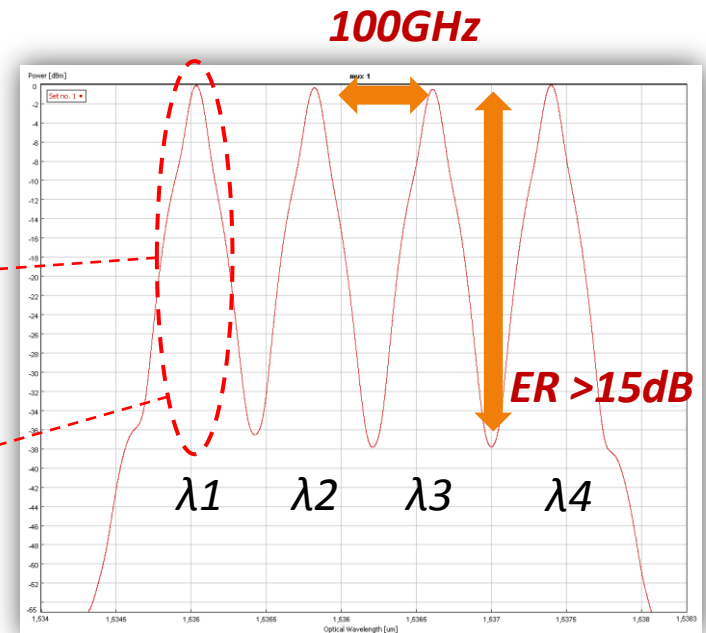
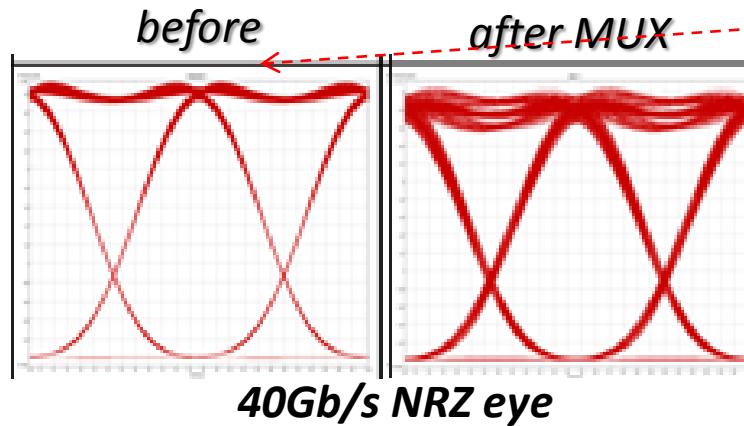


Gaps:

$g_1=200nm$  (power coupling of 0.06)

$g_2=460nm$  (power coupling of 0.0007)

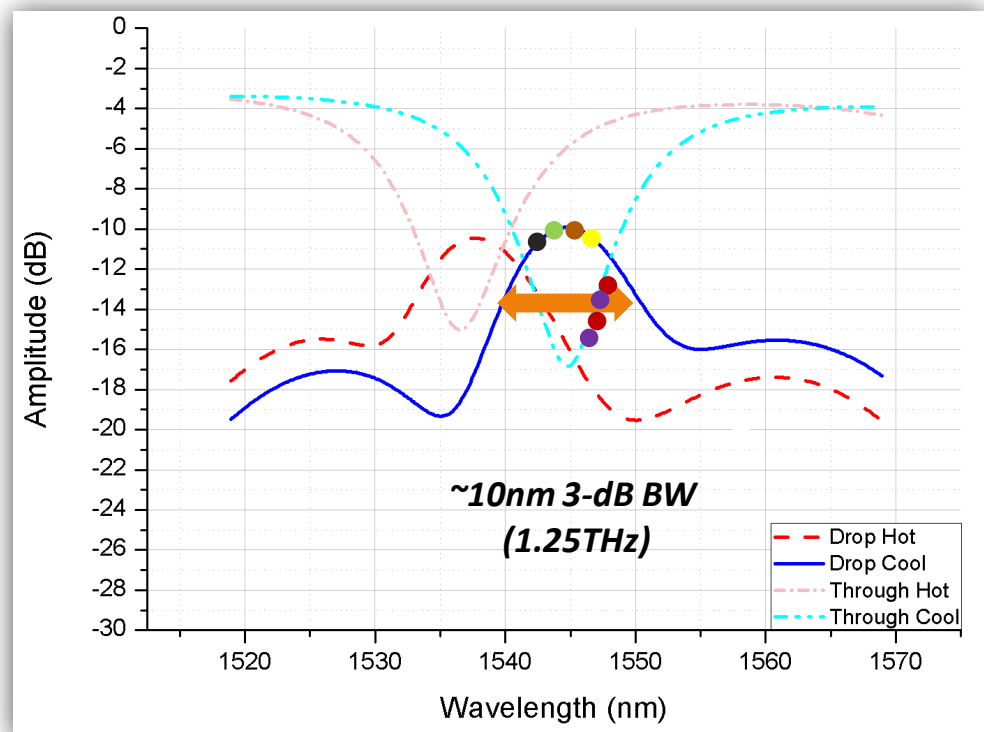
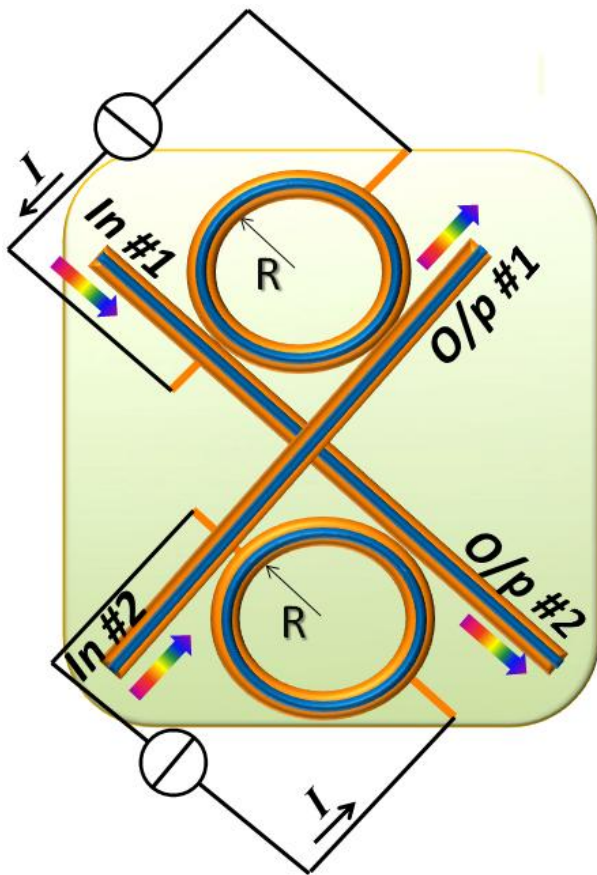
$g_3=200nm$  (power coupling of 0.06)



# Broadband 2x2 Plasmonic Switch

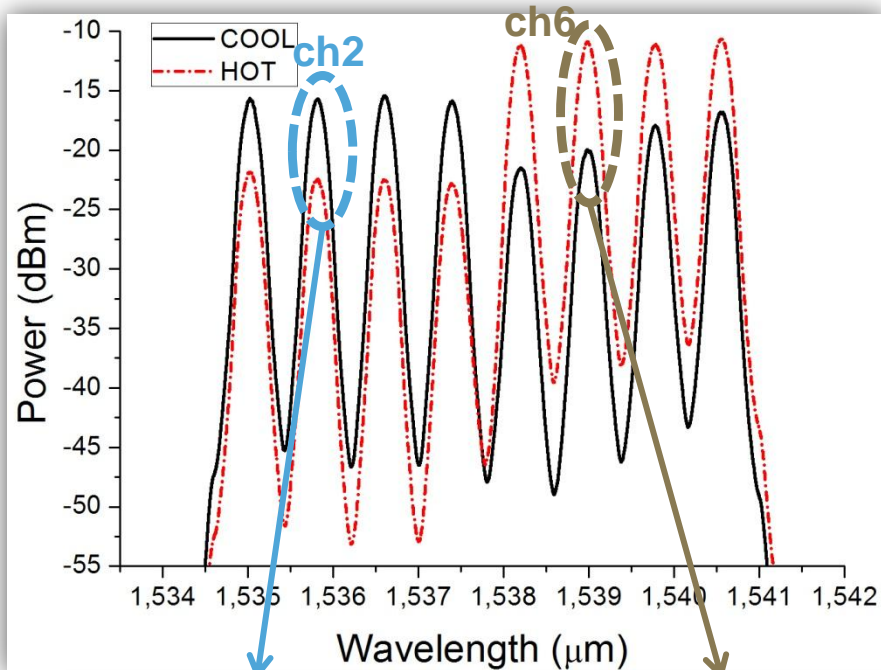
✓ *dual plasmonic ring resonator*

-  $R=5\mu\text{m}$

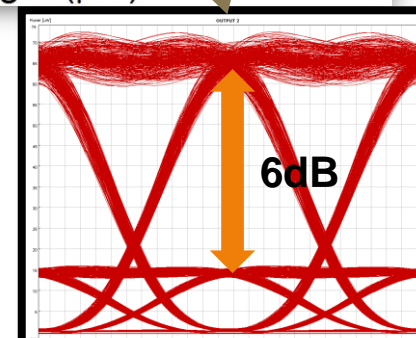
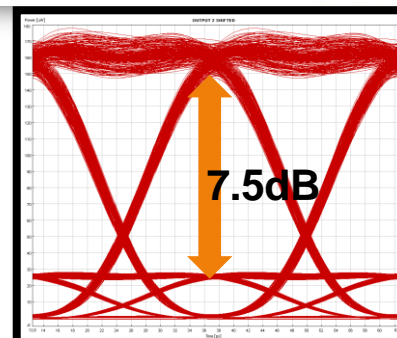
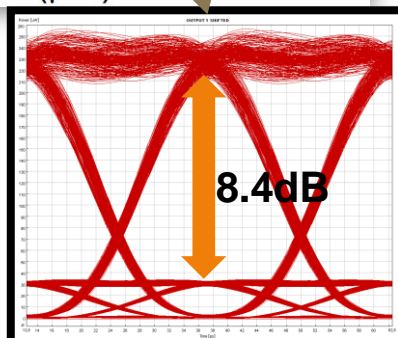
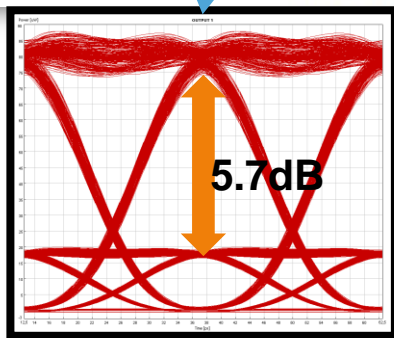
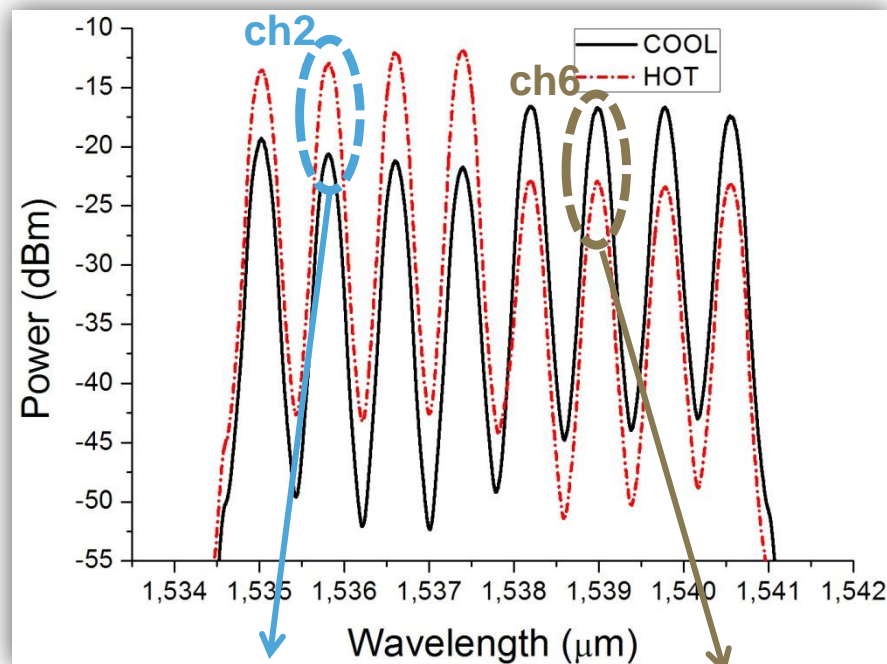


# 320Gb/s throughput routing

## Output 1

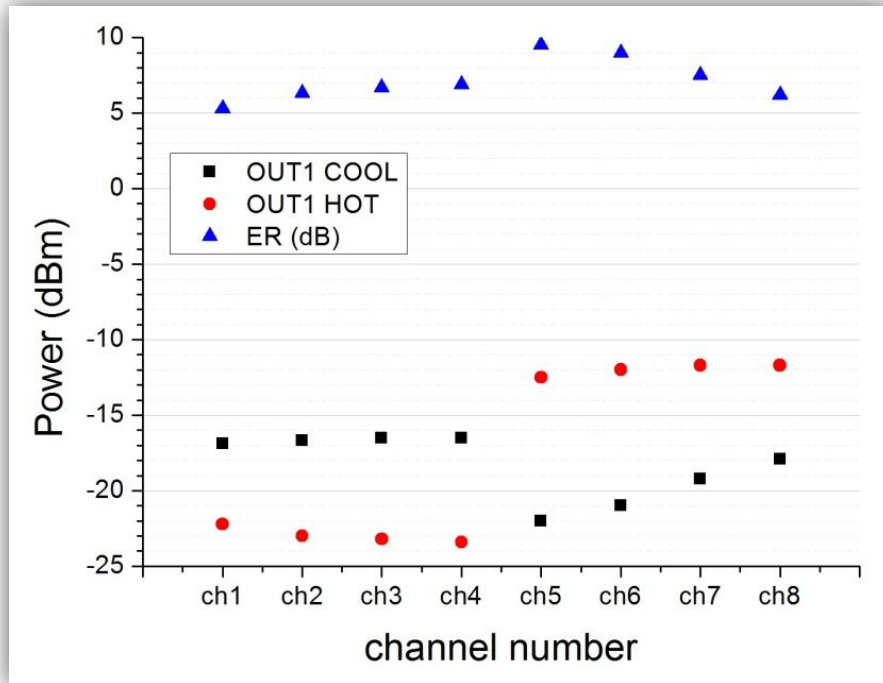


## Output 2

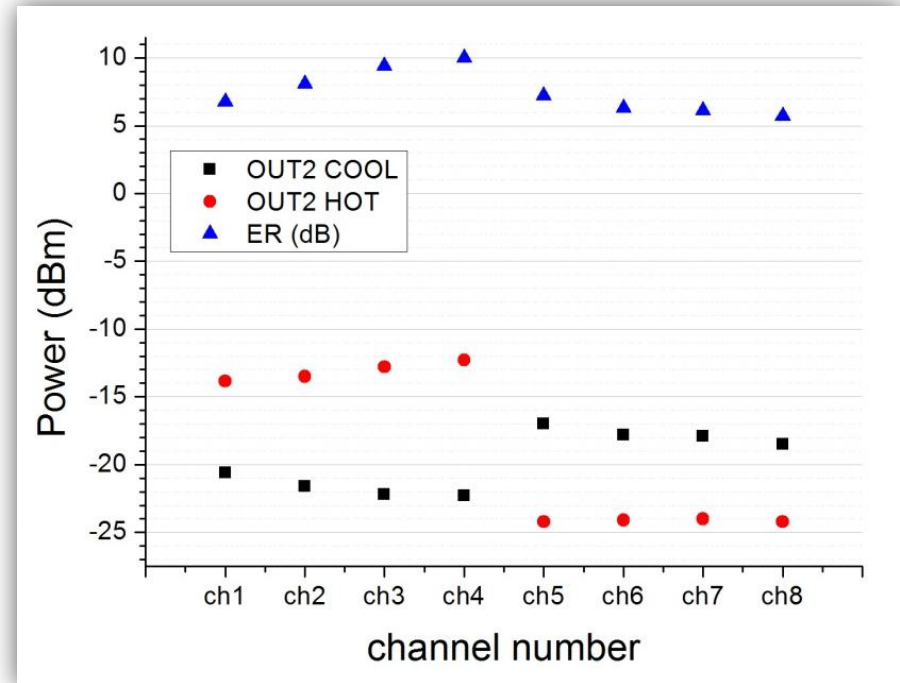


# 320Gb/s throughput routing

## Output 1



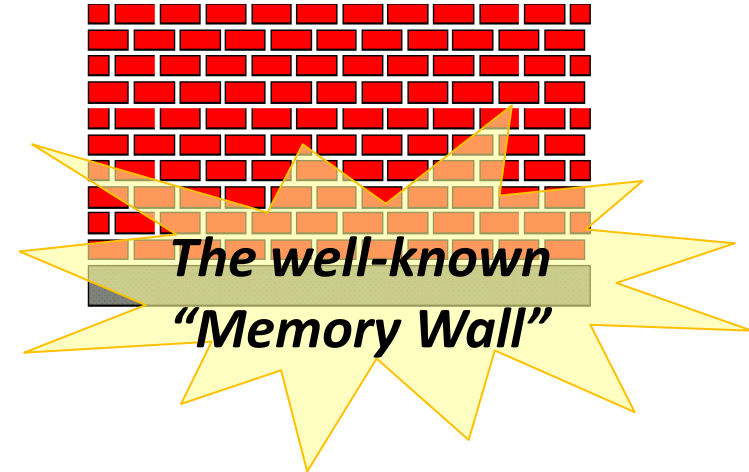
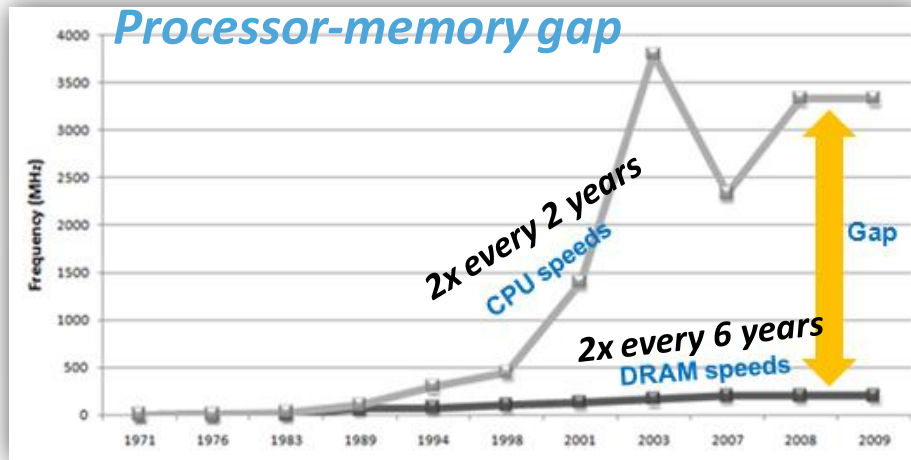
## Output 2



**All channels having ER between 5.5 and 10 dB**

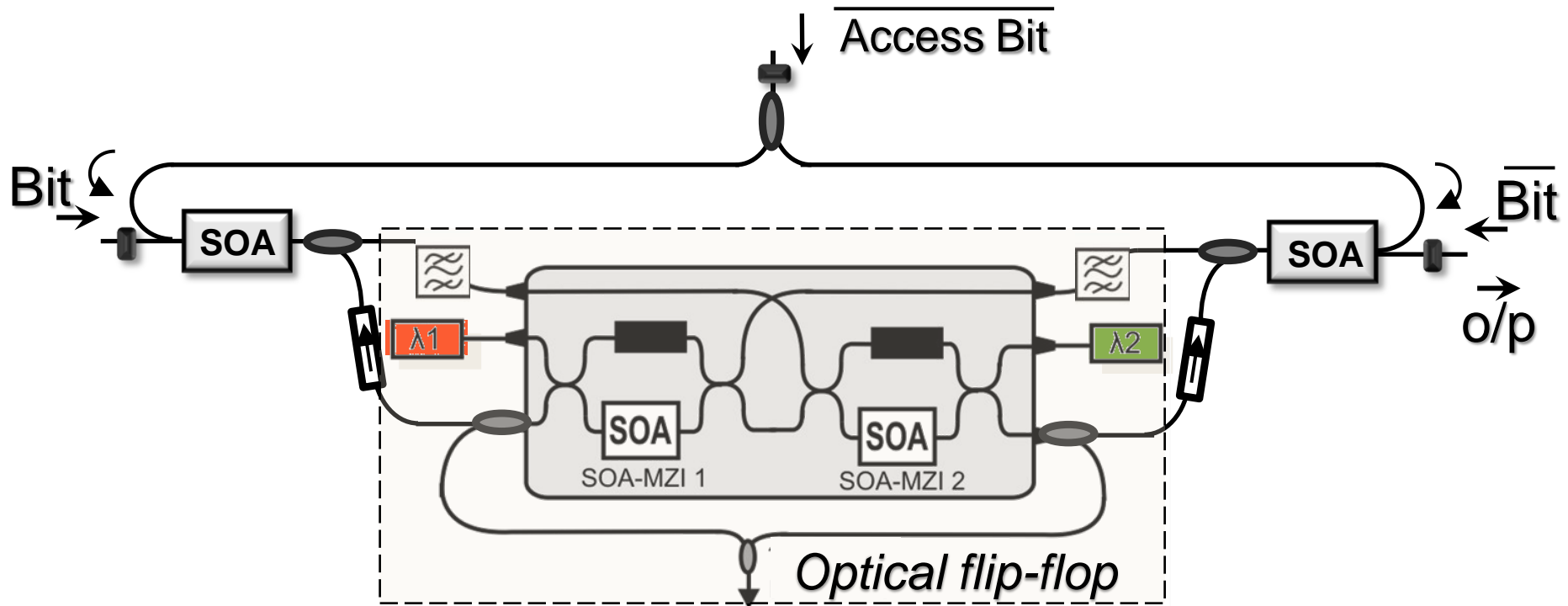
# What about buffering in HPC?

- ⌘ *Latency of the entire HPC is limited by the nsec access time of electronic RAM*



- ⌘ *...but electronic RAM is the only available solution for the HPC Storage Area*

# Optical RAM



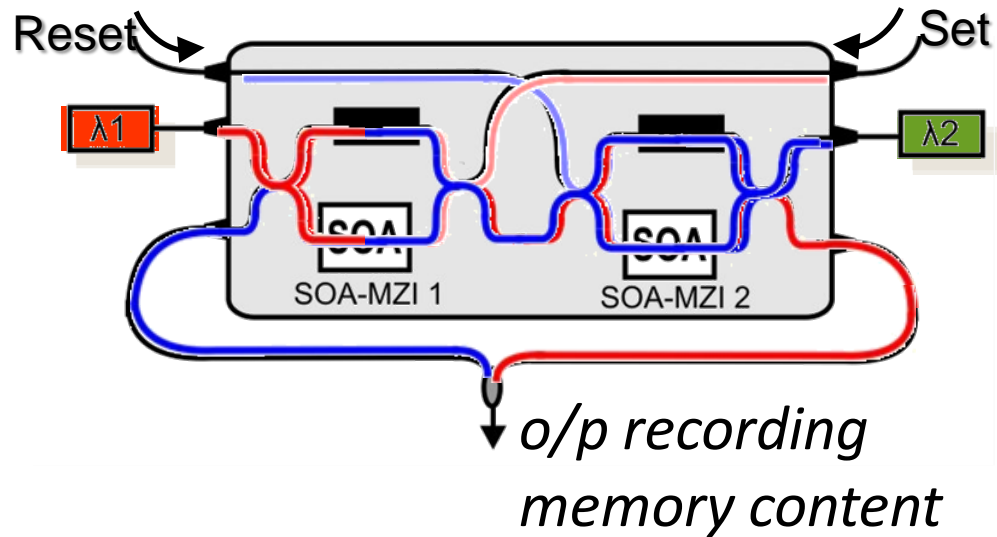
comprises:

- *integrated optical flip-flop as memory unit*
- *2 'ON-OFF' SOA switches controlled by  $\overline{\text{Access Bit}}$*



# Optical RAM

**Memory unit:**

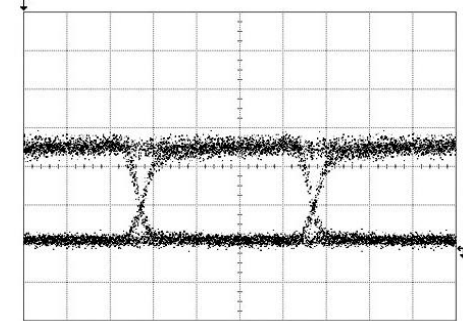
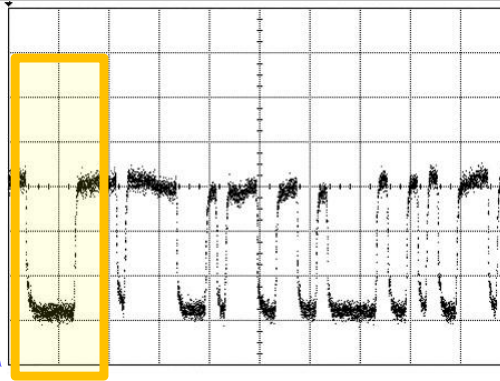


*Optical flip-flop using 2 coupled optical switches*

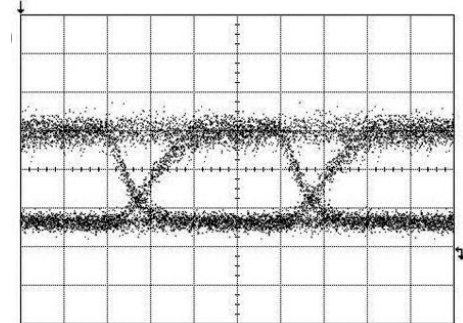
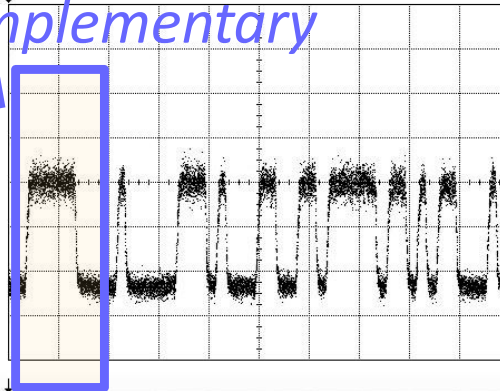
- Memory content = *logical '1' when  $\lambda 1$  dominant*
- Memory content = *logical '0' when  $\lambda 2$  dominant*

# 5GHz Optical Random Access Read

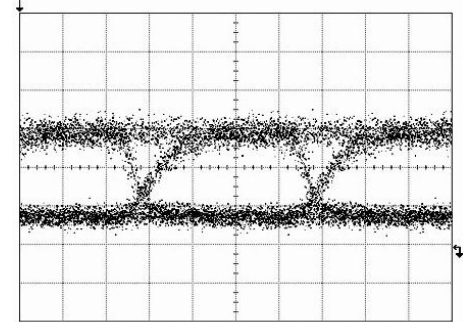
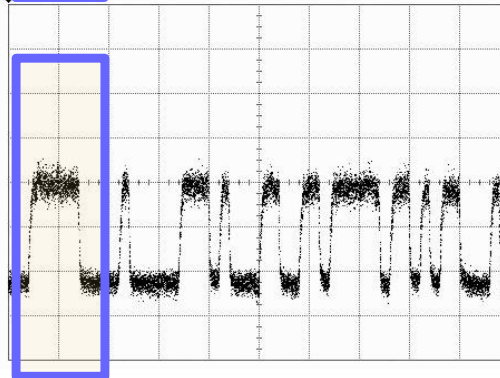
*Inverted  
Access Bit*



*Read o/p @  
 $\lambda_{FF\#1}$  1556nm*



*Read o/p @  
 $\lambda_{FF\#2}$  1559nm*



*complementary*

# 5GHz Optical Random Access Write

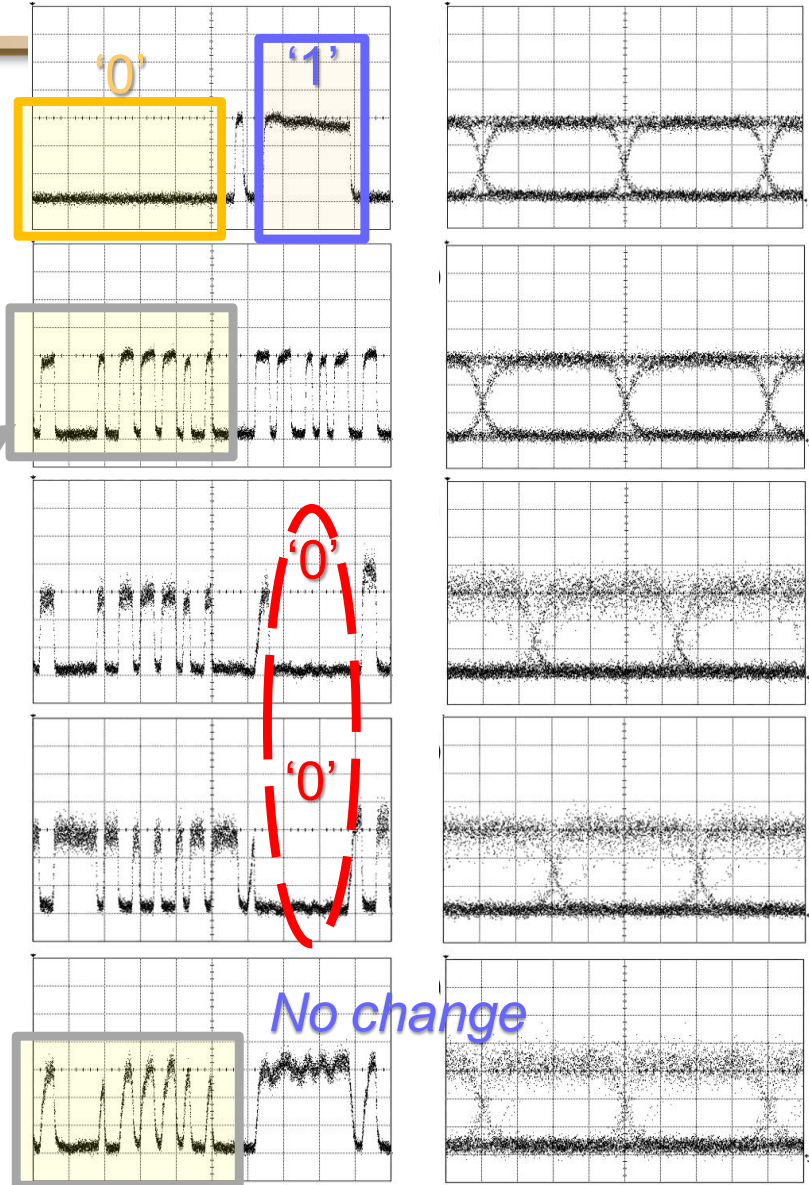
*Inverted Access Bit*

*Incoming Bit signal*

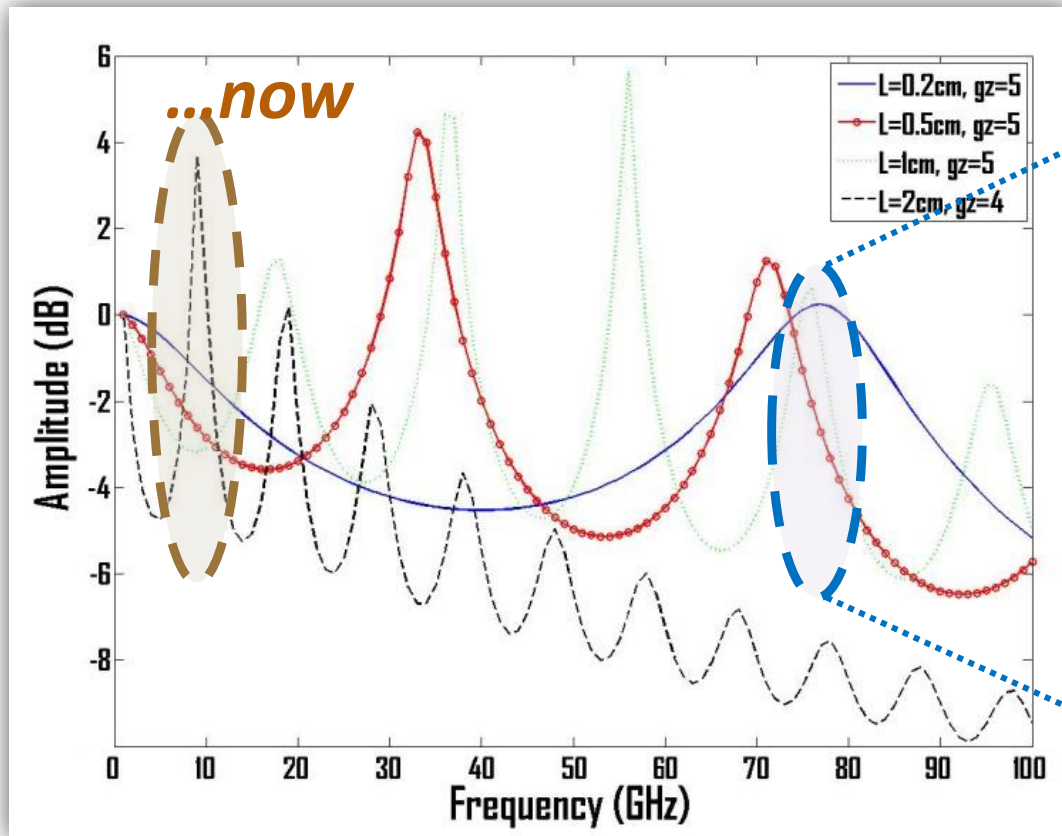
*'Reset' signal*

*'Set' signal*

*Memory content*



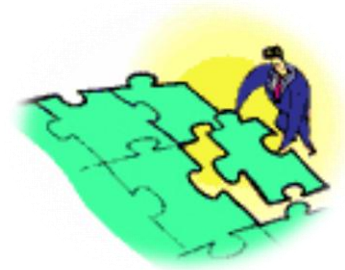
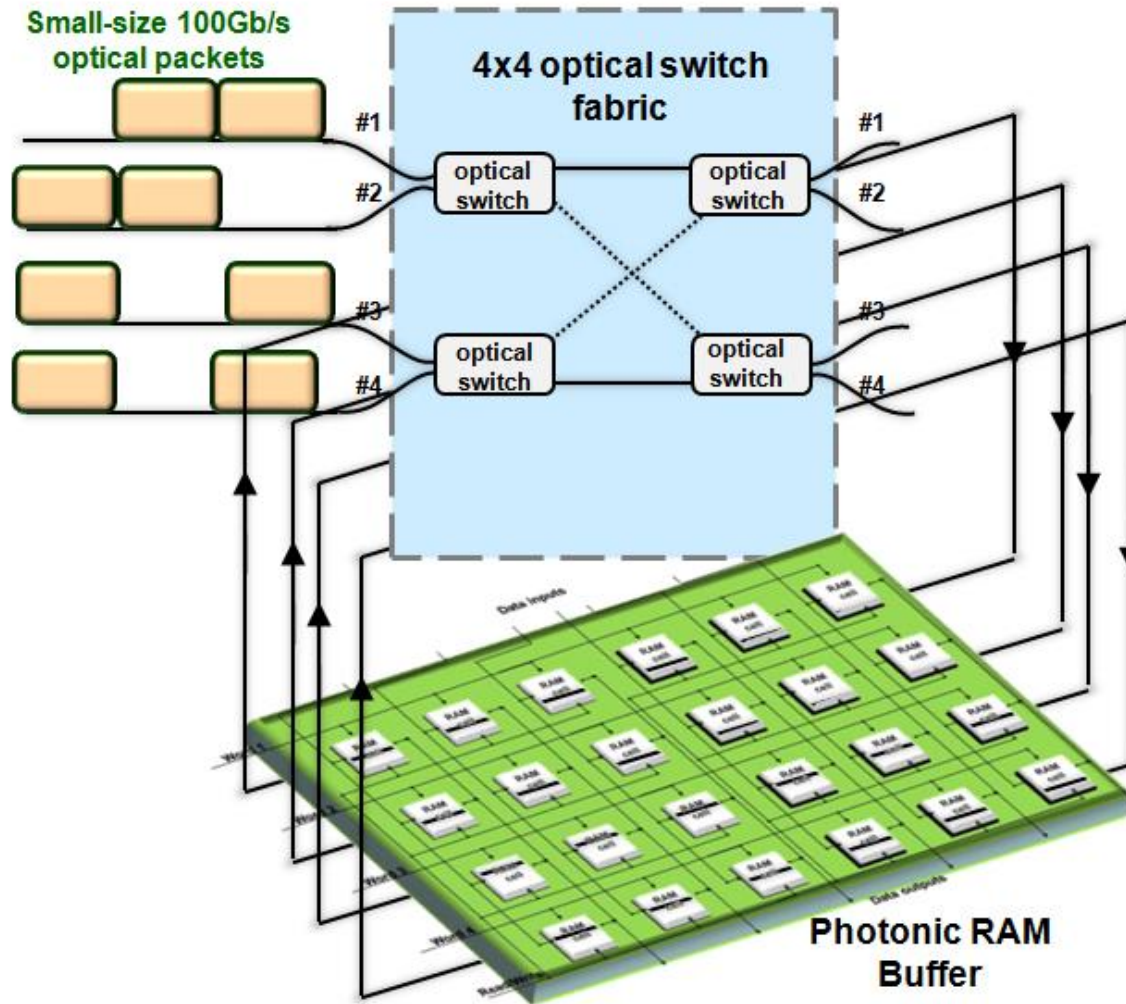
# Towards 100GHz Optical RAM



*Optimized circuit design and silicon-integration can lead to 100GHz Read/Write*

*RAM Speed  $\sim c/nL$*

# Towards true all-optical routers



---

# THANK YOU !

## The PhosNET team...

■ *T. Alexoudi, D. Fitsios, G. Kalfas, G.T. Kanellos, A. Miliou, S. Papaioannou, D. Tsiokos , K. Vyrsokinos*

