![Broadcom logo]

# An AI Compute ASIC with Optical Attach to Enable Next Generation Scale-Up Architectures

Manish Mehta
VP Marketing and Operations – Optical Systems Division

August 26th, 2024

# Broadcom Co-Packaged Optics (CPO) Background

## Our objective

- Build an optical interconnect that provides substantial improvements over current optical modules on cost, power consumption reliability and latency

- Build a high-density optical interconnect that enables up to 1 Tb/s/mm duplex connectivity to support current gen and next gen scale-up and scale-out optical BW density
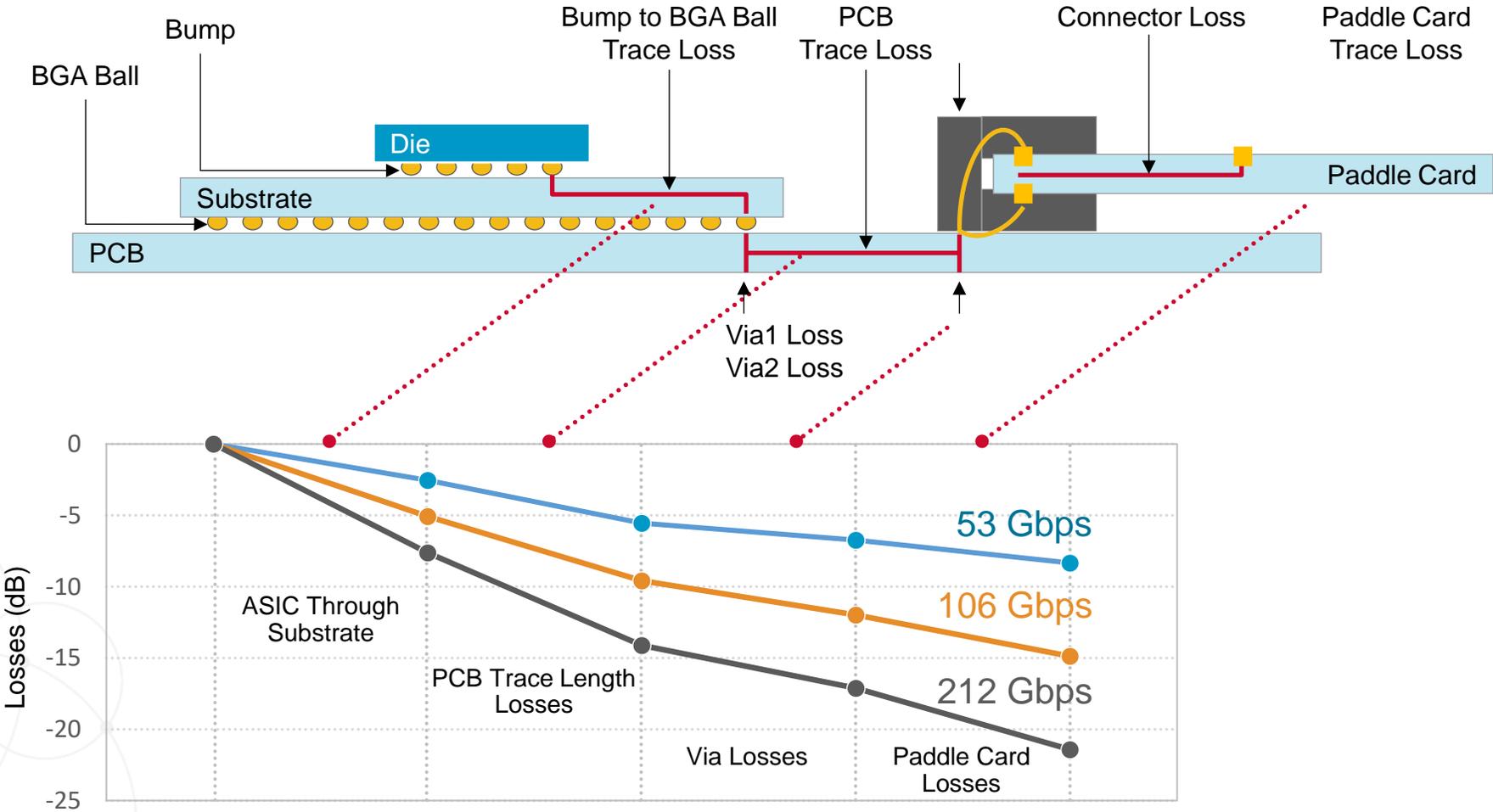
## Our Progress

- Shipped first generation proof-of-concept 25.6T CPO switch product

- Currently qualifying and ramping second generation 51.2T CPO switch product

- Demonstrated AI ASIC + CPO functionality: CPO + 2.5D Packaging

## What I will share today

- Our journey to build the first Co-Packaged Optics (CPO) deployed in datacenter Front-End and Back-End networks using high density Silicon Photonics (SiPh)

- Implication of CPO Platform on both scale-out and scale-up connectivity

BROADCOM®

# Serdes Migration to 200Gbps Limits Electrical I/O Reach



**Mandates Development of Optical Interconnects Co-Packaged with ASIC**

# Journey to an AI ASIC with CPO: Discrete III-V to SiPh



**Conventional Module Design**

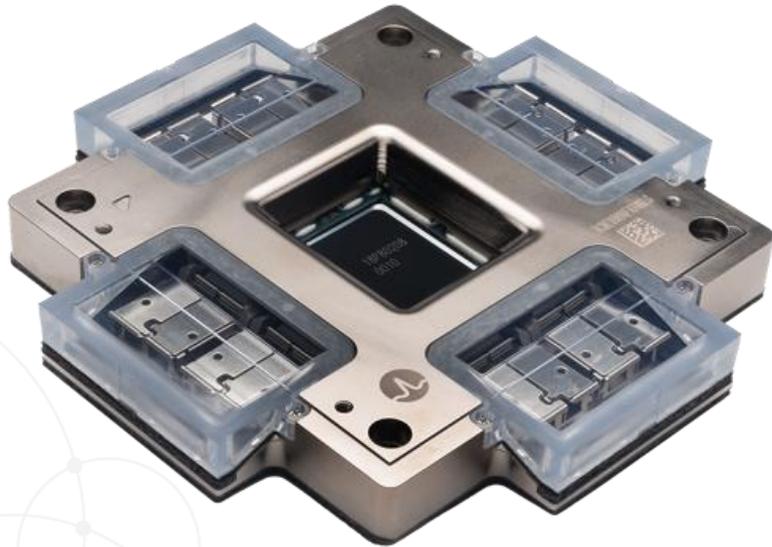**Engineering and manufacturing limits to scale**

**Modules with Silicon Photonics**

SiPh Chiplets in Module

Phy IC

Fiber Jumper

**Module Integration = First step to improved scale**
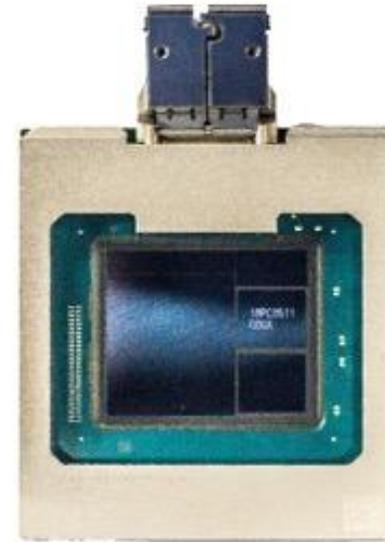
BROADCOM®

# Journey to an AI ASIC with CPO: Switch CPO to GPU CPO
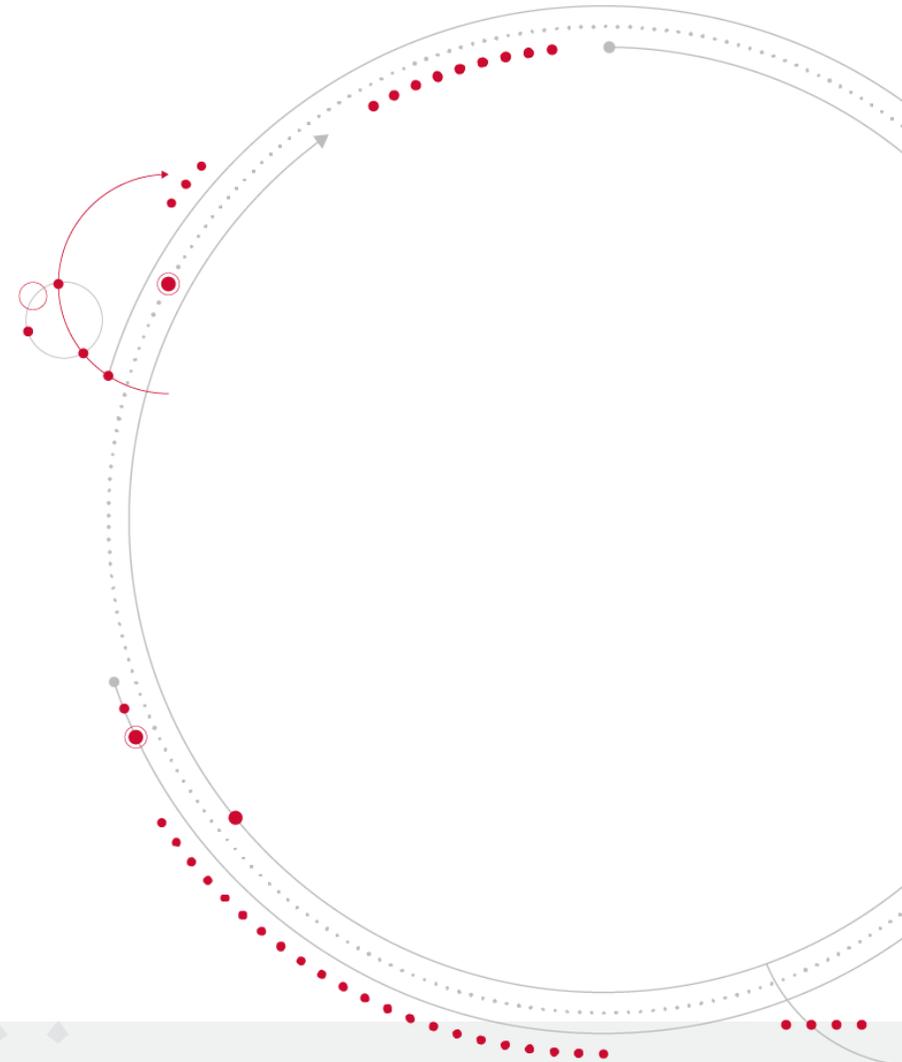
**CPO for Scale-Out Networking**



**Greater than 50Tbps of Optics Attached to Switch ASIC**
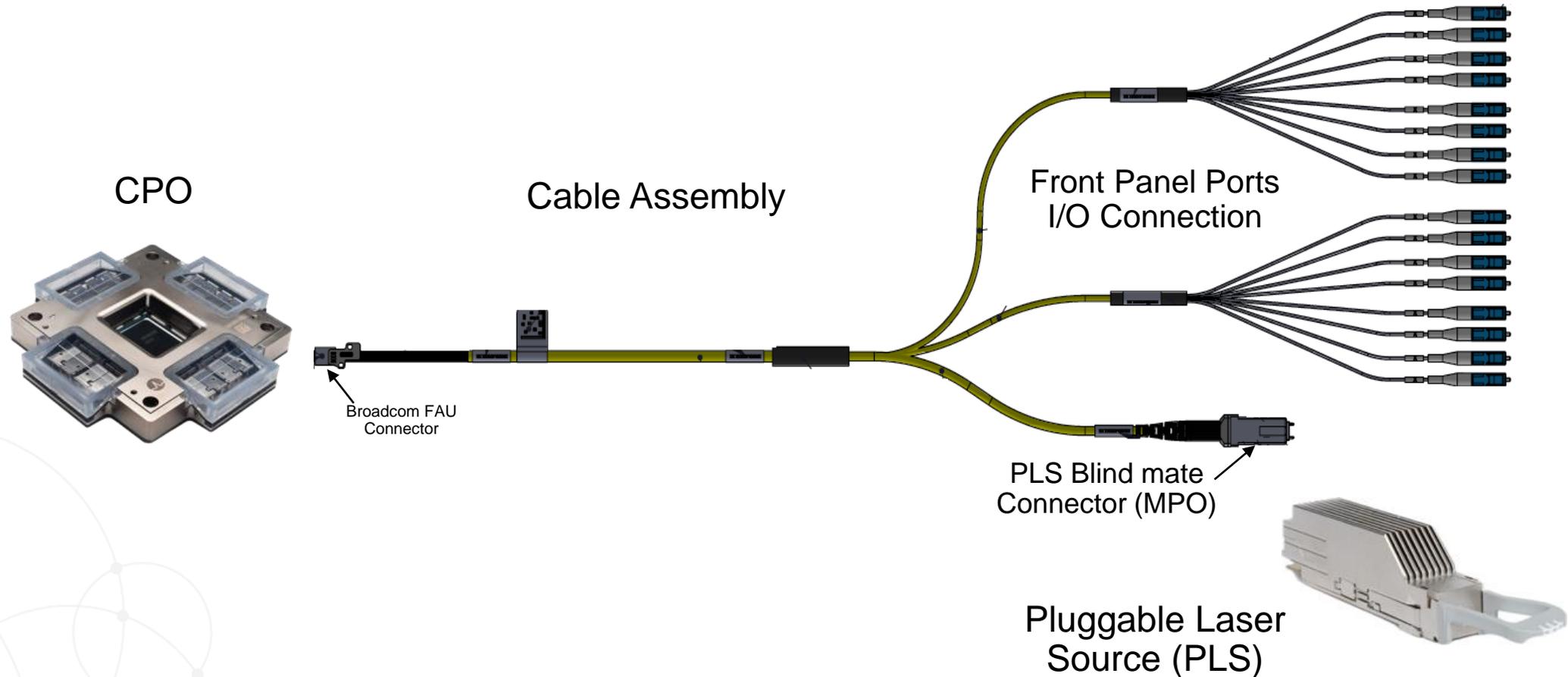
**CPO for Scale-Up Compute**



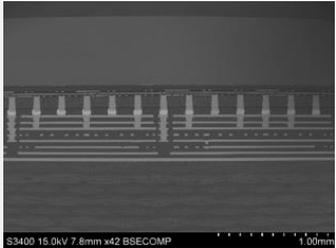**Greater than 6.4Tbps of Optics Attached to GPU**

**BROADCOM®**

# What is our CPO Platform

# CPO Schematic



CPO

Cable Assembly

Front Panel Ports I/O Connection

Broadcom FAU Connector

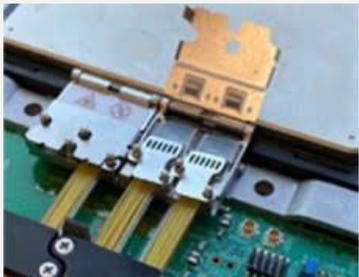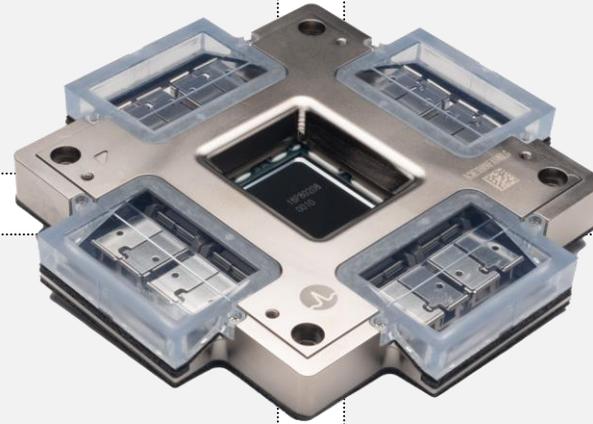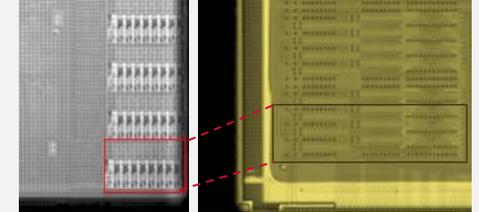PLS Blind mate Connector (MPO)

Pluggable Laser Source (PLS)

- 51.2Tbps TH5 Switch CPO with 8x6.4T optical engines
- Pluggable Laser Modules x 16 (field serviceable)
- Fiber Cable Assembly
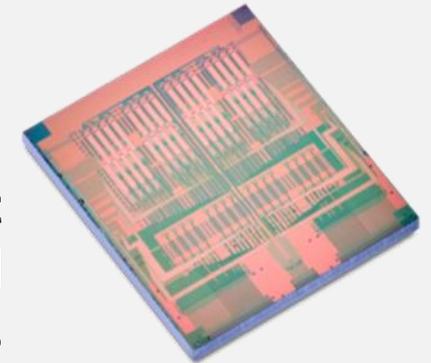
BROADCOM®

# Key Components of CPO
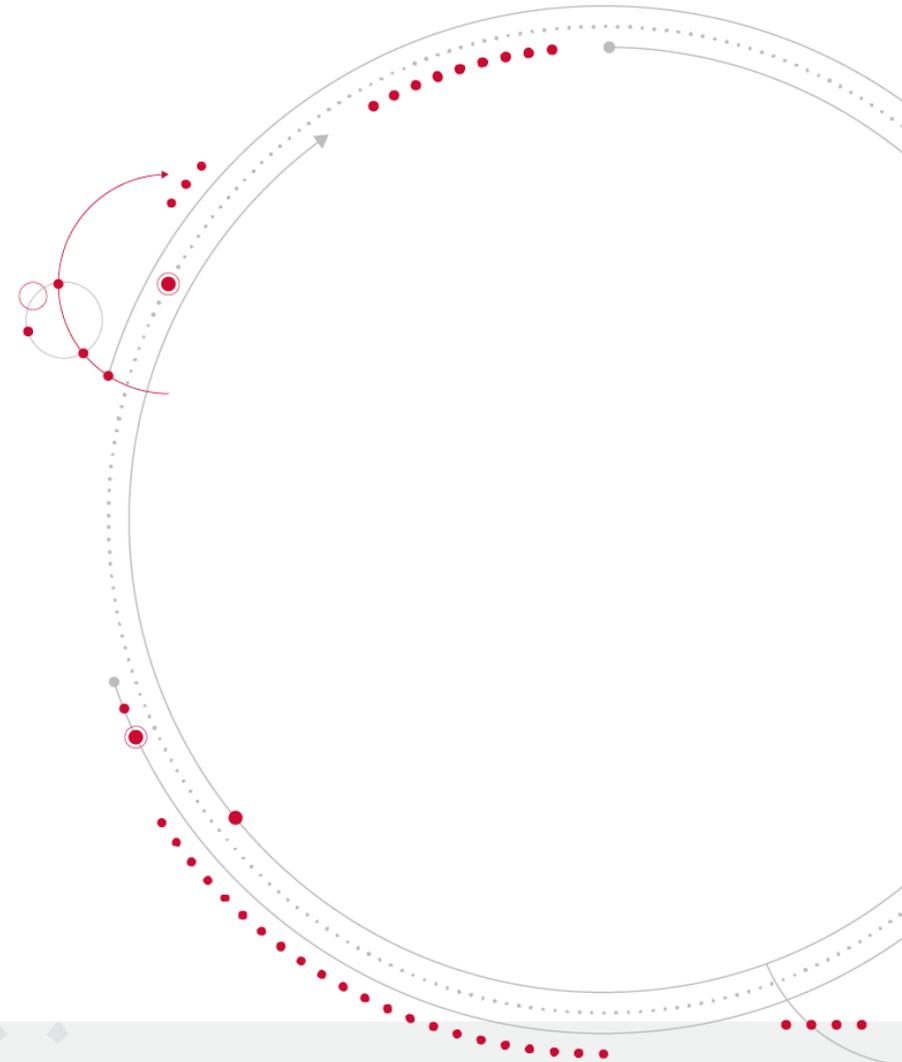
**Advanced Packaging**

**Electrical Integrated Circuit (EIC) with driver and TIA**

**High density fiber connector**

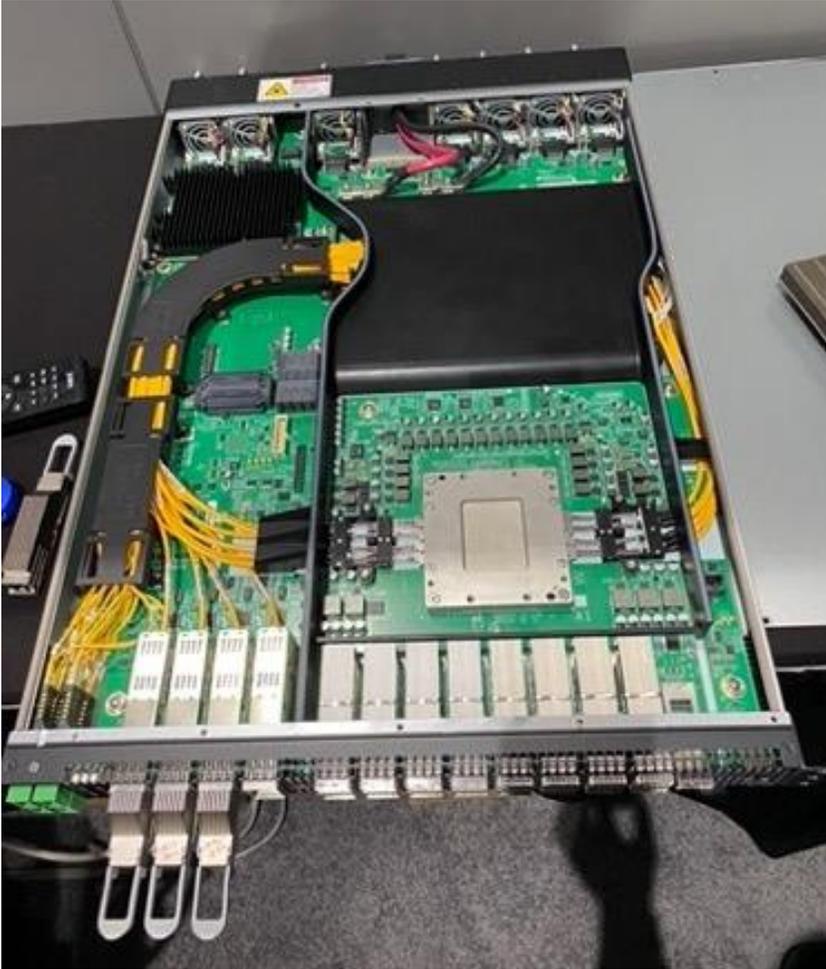**Photonic Integrated Circuit (PIC) with modulator and photo diodes**

**BROADCOM**®

# Scale-Out Networking using CPO

# TH4-Humboldt: First Generation System



**Product Features:**

- 25.6T Ethernet Switch
- Half CPO, Half Electrical connectivity
- Four 3.2T optical engines (32x100Gbps DR connectivity)
- Optical engine is a PIC bonded to a SiGe EIC
- Each optical engine has ~ 250 optical components

**SiGe dissipates additional 3 pJ/bit power consumption compared to CMOS solutions**

BROADCOM®

# TH4-Humboldt: SiPh PIC + SiGe EIC + TSV



Laser

**Photonic and Electronic Circuits**

Signal Processor ASIC

PIC

EIC

Substrate

BROADCOM®
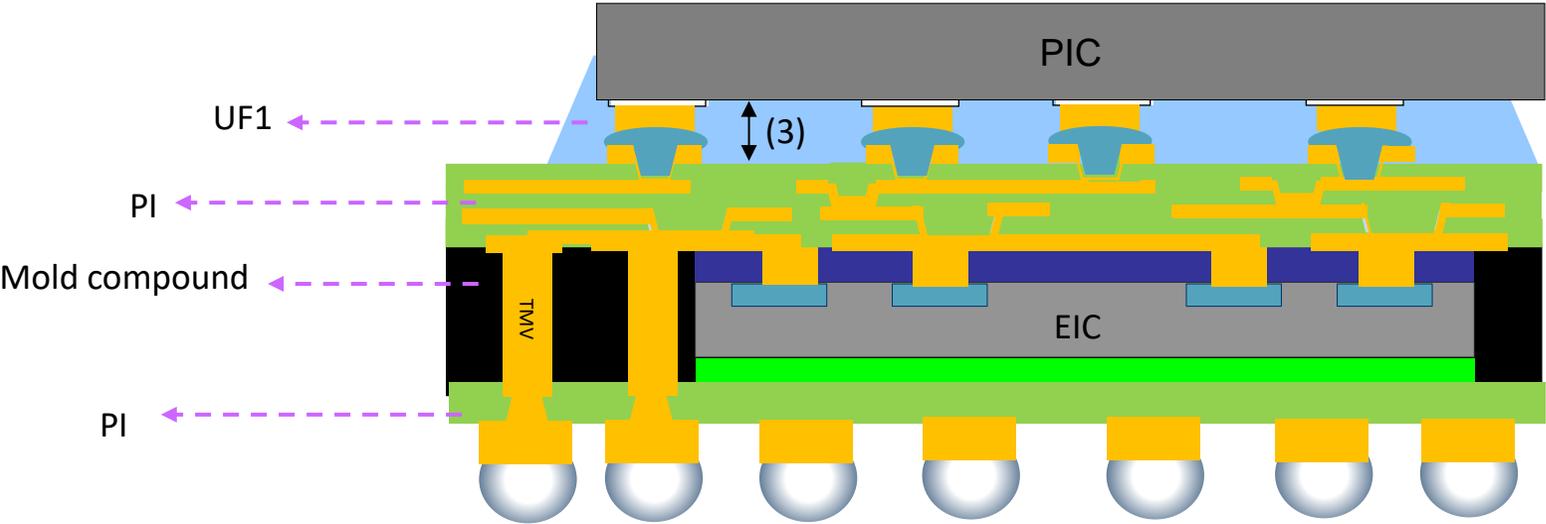
# TH5-Bailly: Second Generation System



**Product Features:**

- 51.2T Ethernet Switch
- All Optical CPO connectivity
- Eight 6.4T optical engines (64x100Gbps FR4 connectivity)
- Optical engine is a PIC bonded to a CMOS EIC
- Each optical engine has ~ 1000 optical components

**BROADCOM**®

# TH5-Bailly: SiPh PIC + 7nm CMOS EIC + FOWLP



PIC

UF1

(3)

PI

Mold compound

TMV

EIC

PI

*FOWLP: Fan-out Wafer-level Packaging*



RDL

PIC

EIC

## FOWLP Improved Scalability of PIC to EIC bonding

**BROADCOM**®

# FOWLP Innovation: Dual-Side Attach for Co-Packaged Optics

# Optical Engine Cross Sectional Images



□ **Cross-section on U2**

**Cross-sections taken AFTER eight (8) optical engines bonded to substrate**

# Network Placement and Module Substitution

## Scale-Out Network

## Modules vs CPO



Spine

128 ports 400G

Leaf

GPU GPU GPU GPU • • • GPU GPU

BROADCOM®

# 51.2T TH5-Bailly Demonstration



Fully Functional 51.2T TH5-Bailly inside 4RU MP3 Chassis

| 800G Interconnect (Watts) |
| :---: |
| 4.8 |

**TDECQ: 1.07**

BROADCOM®

# Error Free Operation on a Fully Integrated System



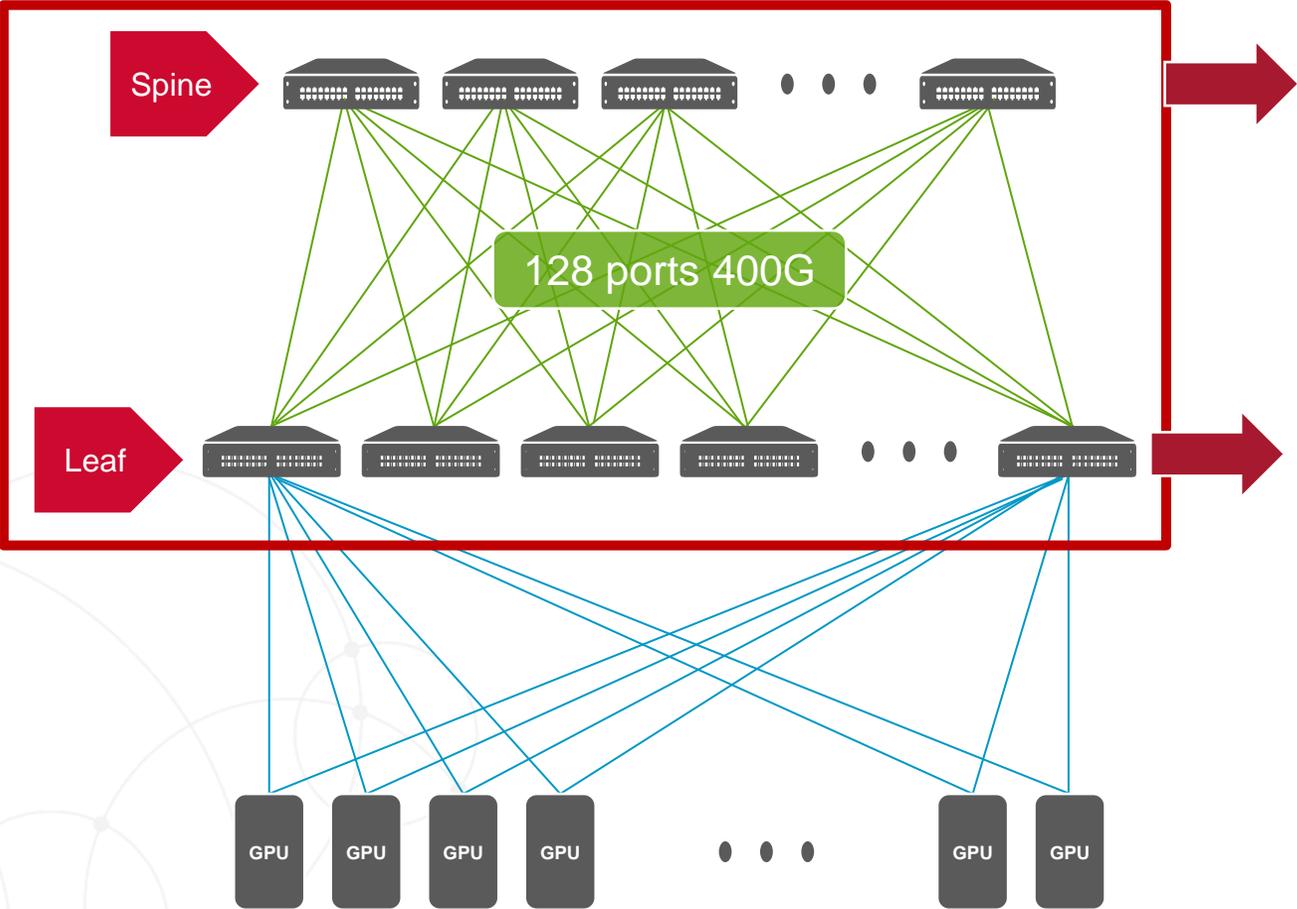- FEC tail distributions collected on all 128 optical ports of a fully integrated 51.2T switch
- Production chassi assembled by ODM in a manufacturing environment
- Image to the left shows FEC tail decay curves for 72 ports

BROADCOM®

# Further Optimization at Lab Level

- Data on a single 64 channel optical engine (6.4T)
- Recent optimizations (still on-going) demonstrate even more rapid FEC tail decay achievable

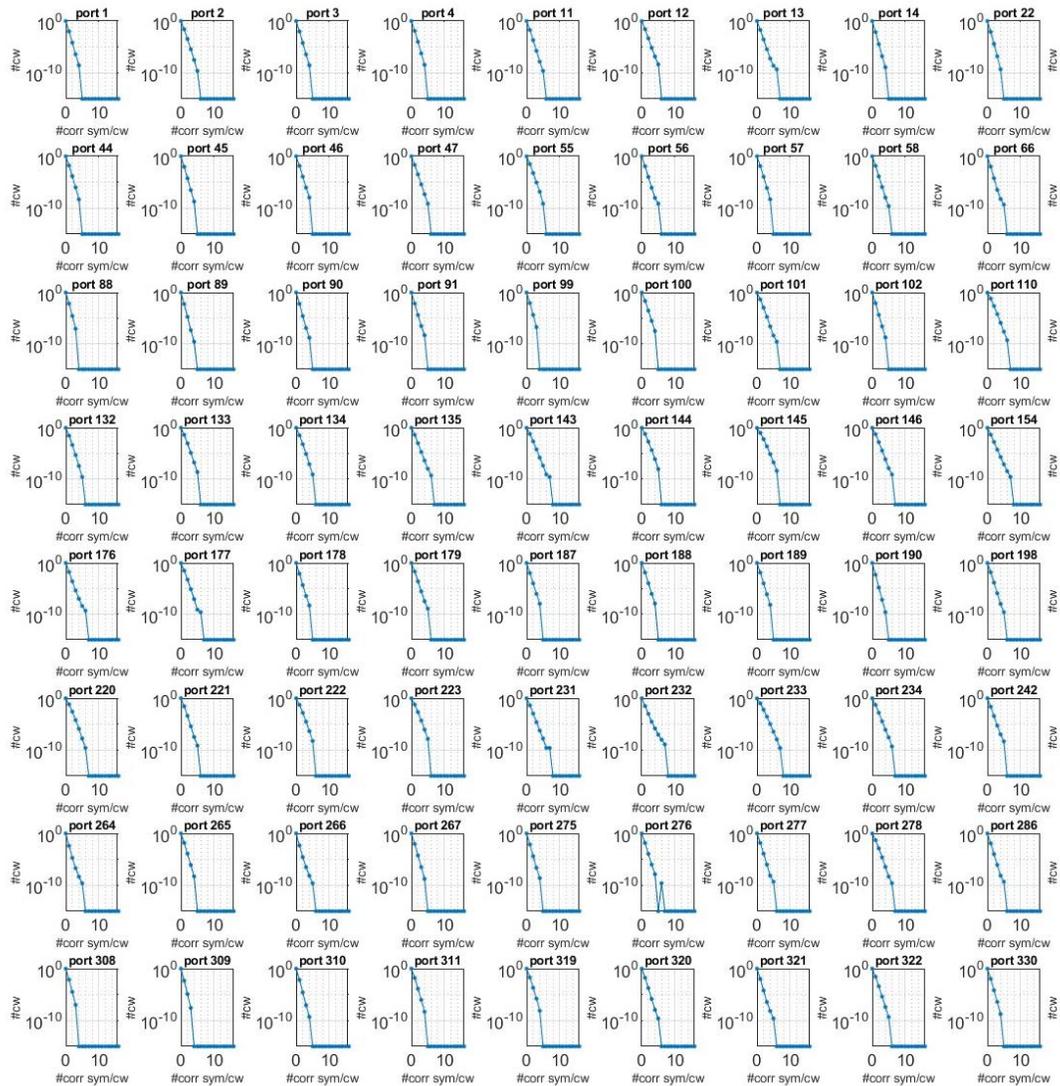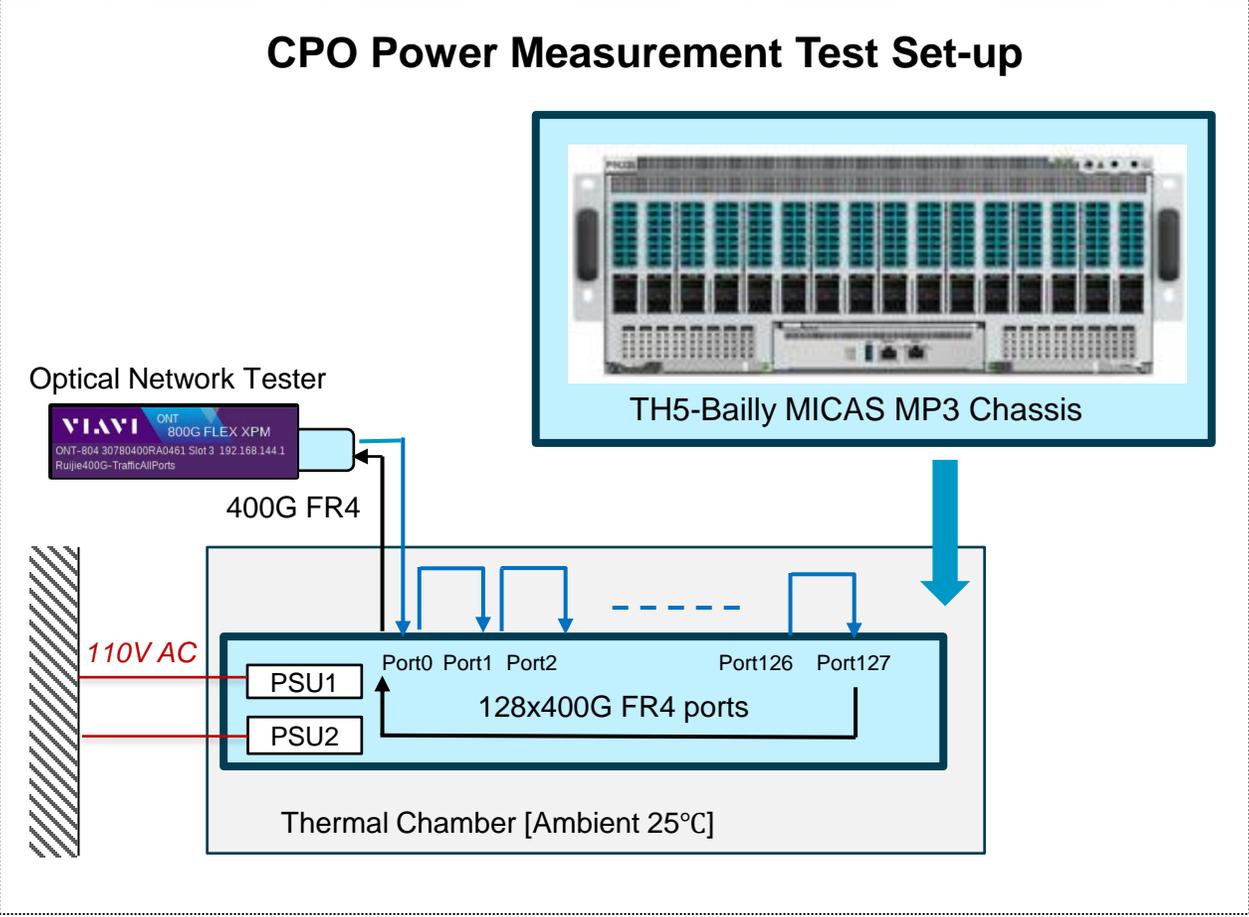| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SymbolErrCt_acc | 3290186 | 7671256 | 7932118 | 24649884 | 1.25E+08 | 63221638 | 1.32E+08 | 10589737 | 20564834 | 50187617 | 23595621 | 37519305 | 84686979 | 25715431 | 15865029 | 20820220 |
| CWErrS0_acc | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.53727E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.54E+11 | 7.53764E+11 | 7.54E+11 | 7.54E+11 |
| CWErrS1_acc | 3283282 | 7665603 | 7923532 | 24618841 | 1.25E+08 | 63173934 | 1.32E+08 | 10562832 | 20557065 | 50125524 | 23593049 | 37513916 | 84549673 | 25698273 | 15847603 | 20814721 |
| CWErrS2_acc | 3452 | 2825 | 4293 | 15499 | 55816 | 23702 | 64263 | 13451 | 3880 | 31030 | 1286 | 2693 | 68569 | 8503 | 8713 | 2748 |
| CWErrS3_acc | 0 | 1 | 0 | 15 | 24 | 92 | 16 | 1 | 3 | 11 | 0 | 1 | 56 | 39 | 0 | 1 |
| CWErrS4_acc | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 |
| CWErrS5_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |
| CWErrS6_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS7_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS8_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS9_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS10_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS11_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS12_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS13_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS14_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS15_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| CWErrS16_acc | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# 51T Switch Chassis Power Consumption with CPO

## CPO Power Measurement Test Set-up



TH5-Bailly MICAS MP3 Chassis

Optical Network Tester

VIAVI ONT 800G FLEX XPM
ONT-804 30780400RA0461 Slot 3 192.168.144.1
Ruijie400G-TrafficAllPorts

400G FR4

110V AC

PSU1

PSU2

Port0  Port1  Port2          Port126  Port127

128x400G FR4 ports

Thermal Chamber [Ambient 25℃]

## Power Measurement Results

### 51T Switch Box Total Power

| | Bailly CPO | Pluggable LPO | Pluggable w/DSP |
|---|---|---|---|
| 8x OE Chiplets | 241 | 630 | 1024 |
| 16x PLS | 118 | | |
| ASIC, CPU, Other | 975 | 975 | 975 |
| Total (Watts) | 1,334 | 1,605 | 1,999 |

- Optical Interconnect using Bailly CPO is 70% lower power
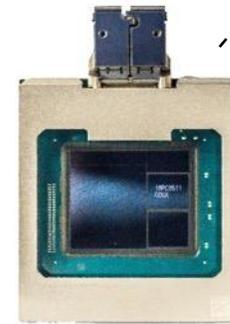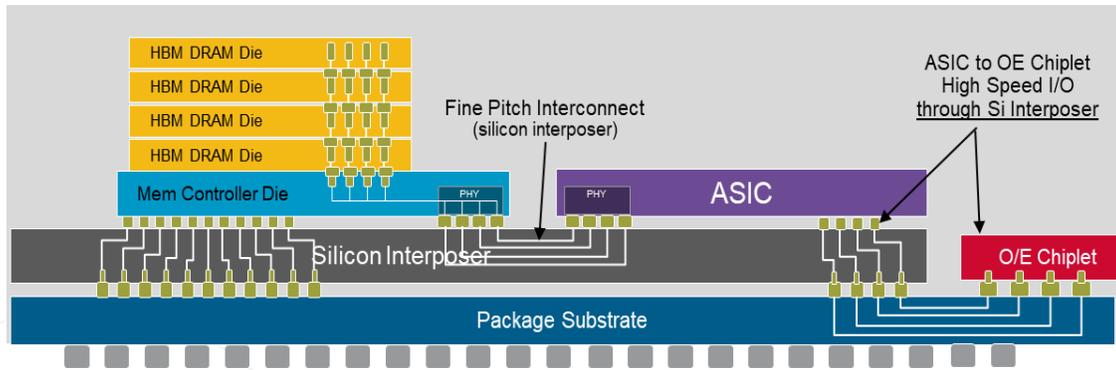- Total Switch Box is ~30% lower power using Bailly CPO

## For a 32k GPU cluster, achieve between > 1MW power consumption savings

BROADCOM®

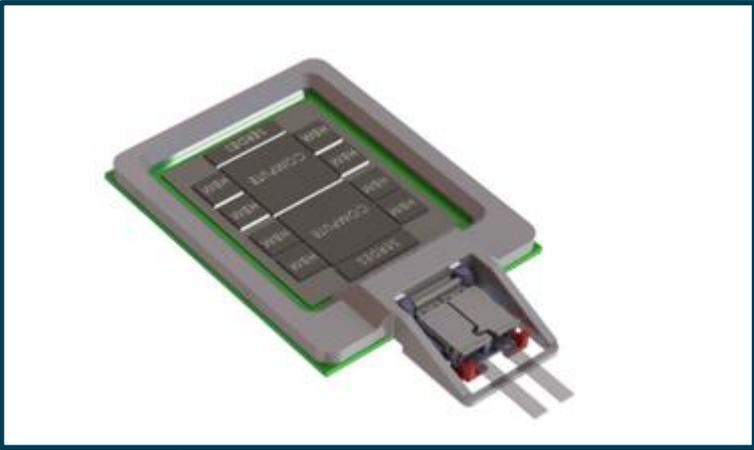# Scale-Up Compute using CPO

# Stage 3: Compute ASICs with CPO
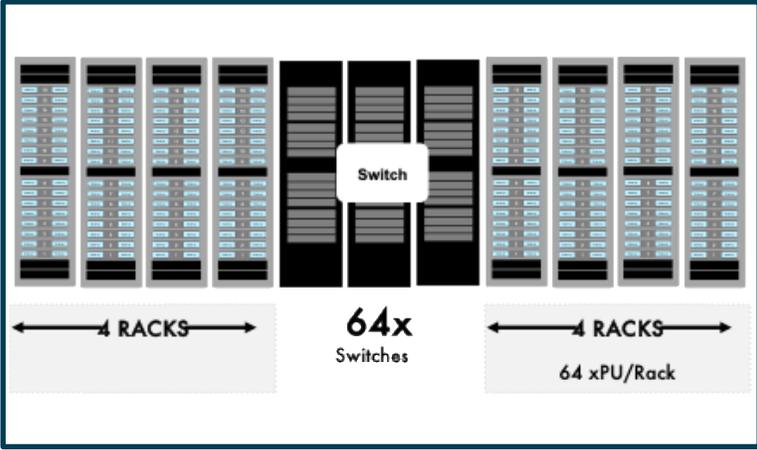


**AI scaling with 2.5D multi-die packaging**

**CoWoS Package with Si Interposer, O/E chiplets and HBM**

**CPO with 6.4Tbps I/O BW per optical Engine**

BROADCOM®

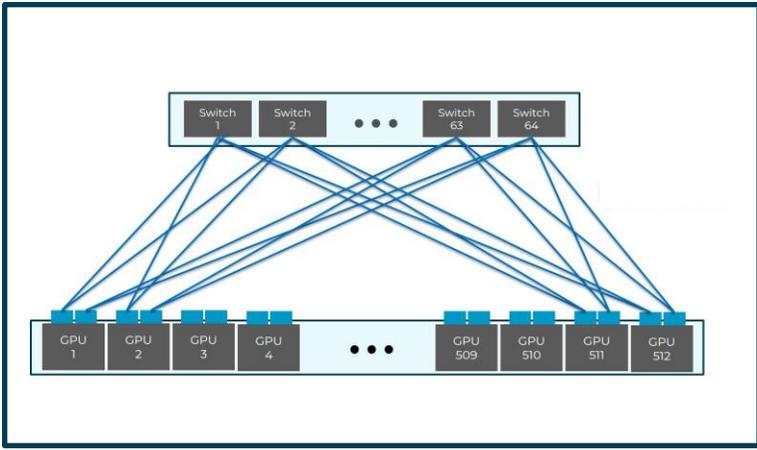# Single-Stage Scale-up Fabric with All-to-All 512 GPU Connectivity

## CPO GPU Attach

## Physical Diagram

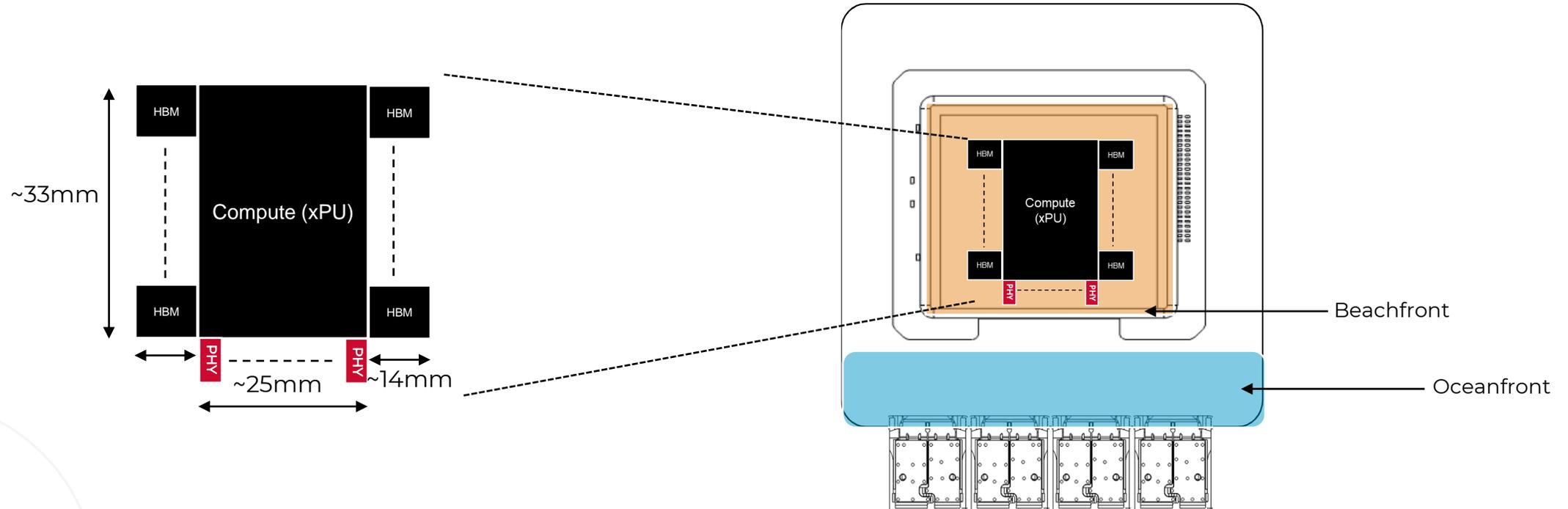- 512 GPU in single row of connectivity
- Optical links 5m-30m (single layer)

## Logical Diagram

- 512 GPUs, 64 high radix switches
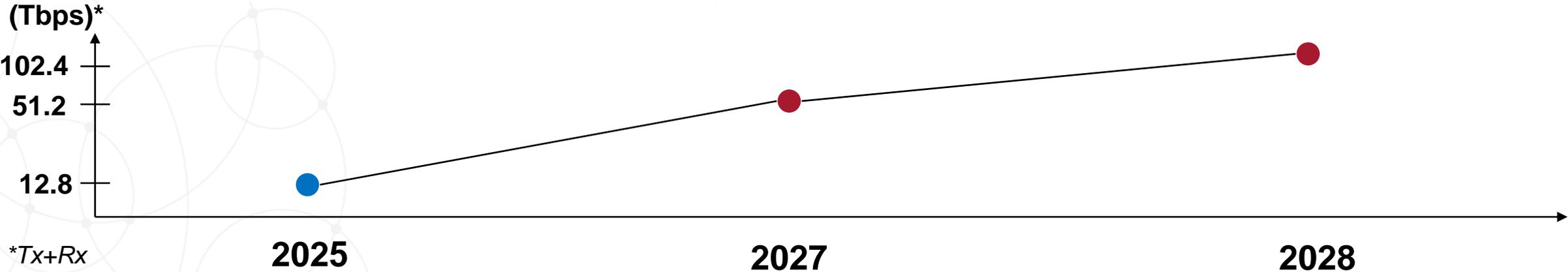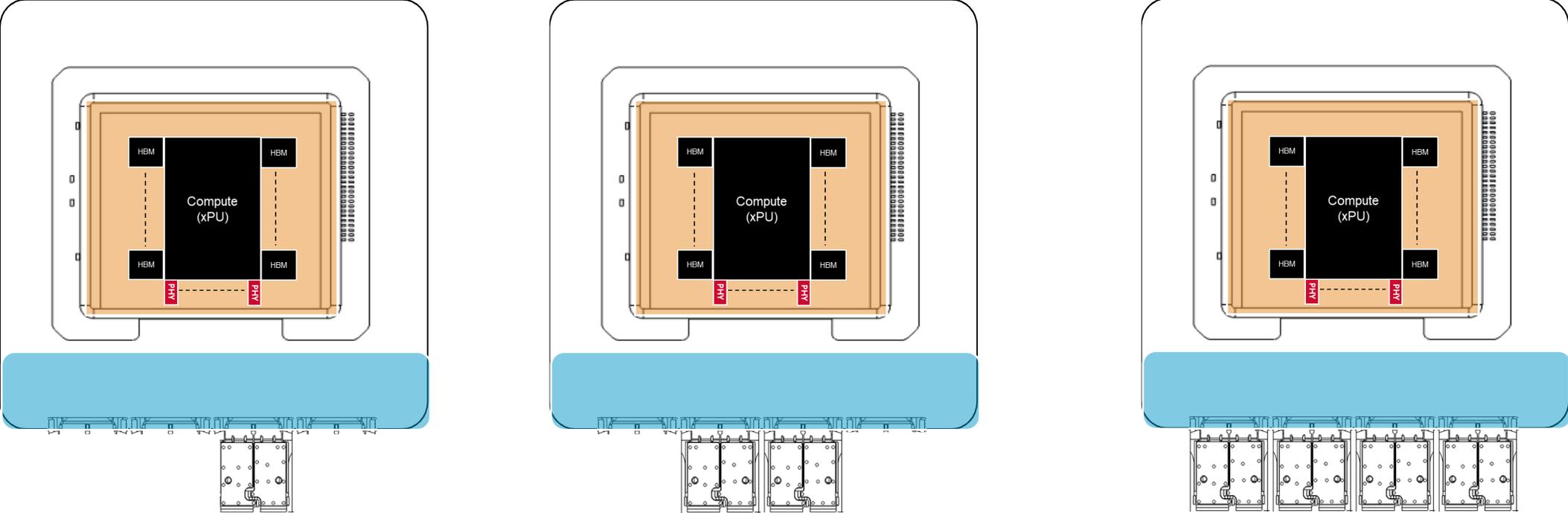- Each GPU connects to all 64 switches via CPO optics

**CPO can enable even larger scale-up domains**

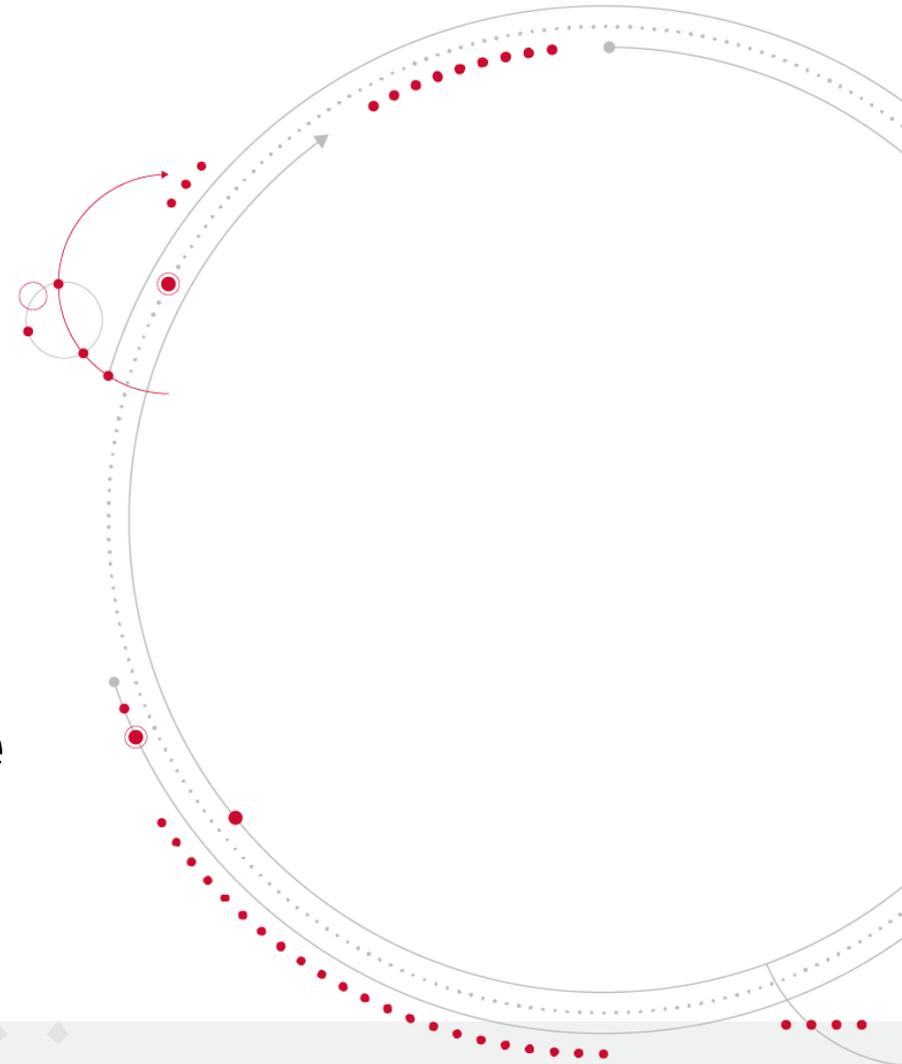BROADCOM®

# Beachfront vs Oceanfront: Utilizing Fan-Out



- **Can escape four optical engines along single oceanfront**

- **Much more reliable and cost effective**
  - Optics is farther away from high power dissipation GPU
  - Known Good Optical Engines can be attached last to the package → high manufacturing yield

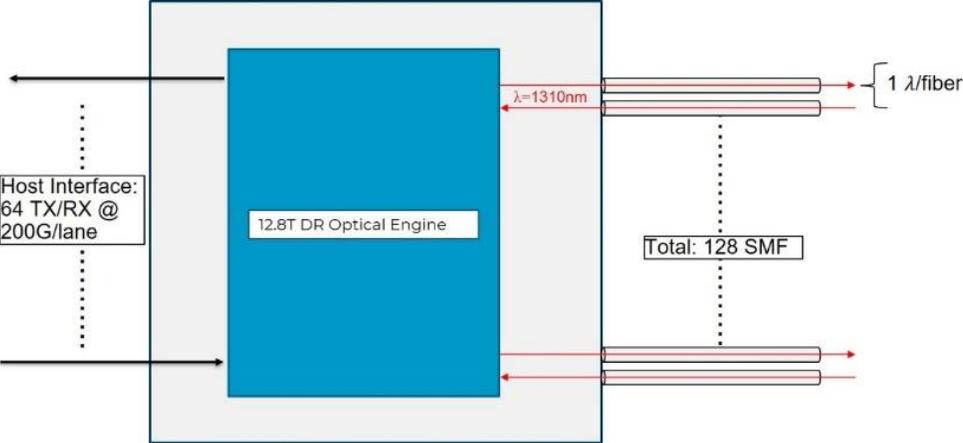# Scale-Up Optical Oceanfront Density Roadmap

# Using Integration and CPO to Minimize Network Cabling Cost in High Radix Clusters

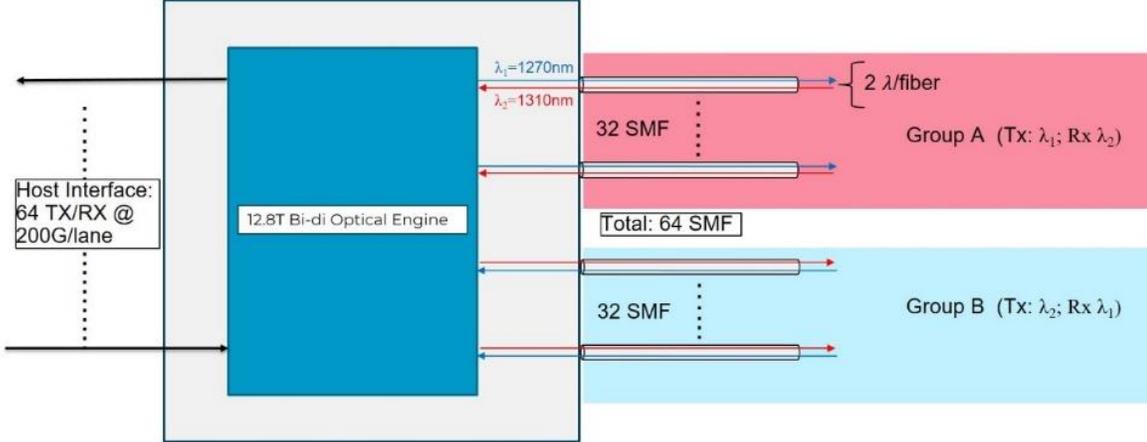**BROADCOM**®

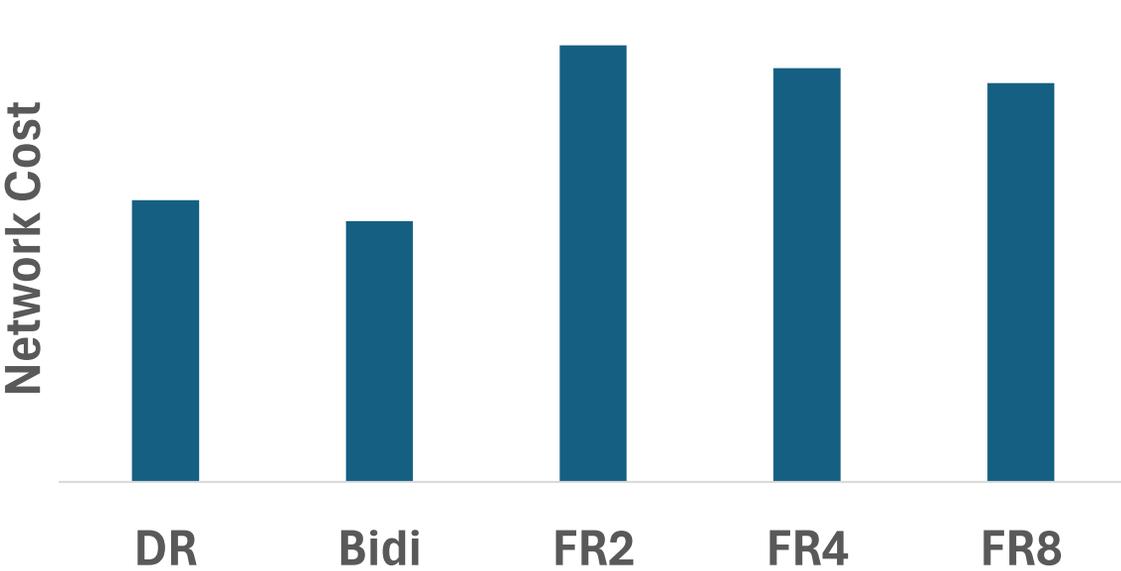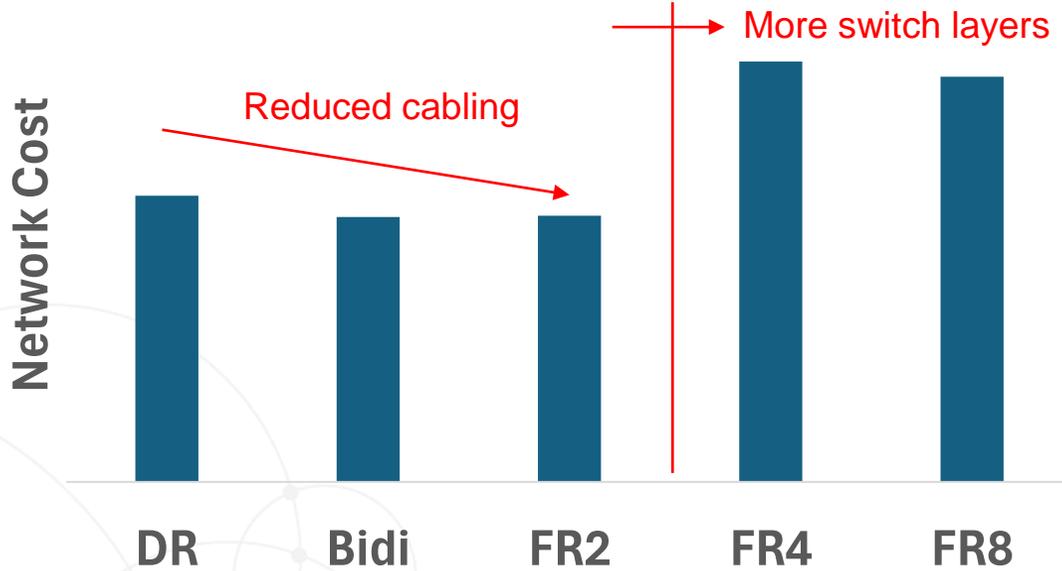# The Lowest Cost High Radix Solution: Bidi Optics



**Bidi widely deployed in FTTx for 20 years, now an attractive solution for AI**

# Implication of Radix and Cabling Selection on Network Cost



**32k GPU cluster, 512 radix capable switch**

Network Cost

Reduced cabling

More switch layers

DR    Bidi    FR2    FR4    FR8

**64k GPU cluster, 512 radix capable switch**

Network Cost

DR    Bidi    FR2    FR4    FR8

## High Radix + Efficient Cabling = Lowest Latency and Lowest Cost

BROADCOM®

# Thank You