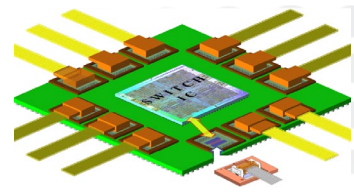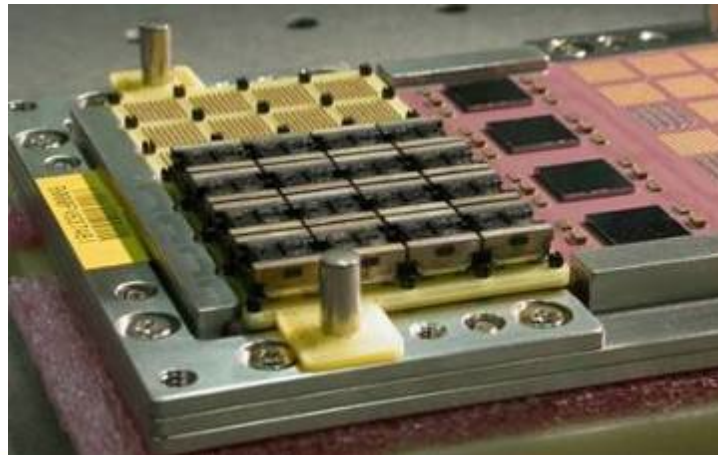**Multi-wavelength Optical Transceivers Integrated On Node**
**PI: Dan Kuchta, IBM  and II-VI (Finisar)**

# Outline
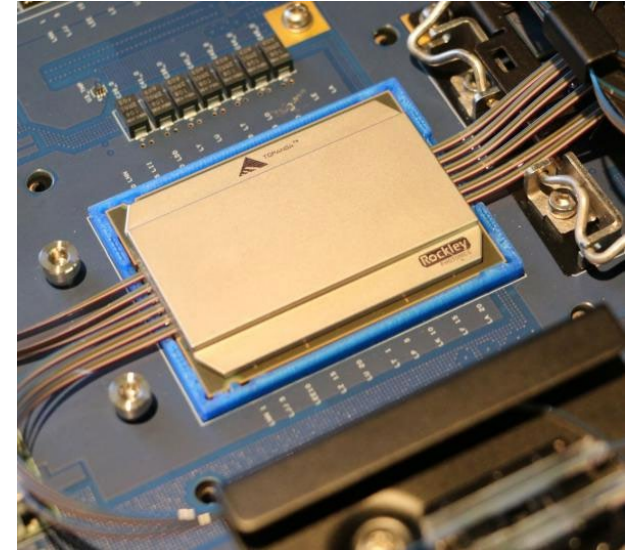
▶ IBM P775: 1st Commercial system to use Co-Packaged Optics

▶ Motivations for Co-Packaging

▶ Overview of MOTION Project and Technology

▶ SAFE (Simplified Analog Front End) ICs

▶ Flip-chip VCSELs and Photodiodes

▶ Network Modelling & Simulation results

▶ MOTION Phase 2

**Rockley Photonics c. 2018**

1.2+1.2 Tb/s
96 fibers at 25 Gbps

**IBM c. 2010**

3.4+3.4 Tb/s
672 Fibers at 10 Gbps

**Altera & Avago c. 2012**

120+120 Gbps
24 fibers at 10 Gbps

# IBM Power 775: 1st system with co-packaged Optics c.2010

IBM



1m W x 1.8m D x 10cm H

Water Connection

360VDC Input Power Supplies 16-18kW typ.

IBM's HPCS Program partially supported by

DARPA

Memory DIMM's (64x)

P7 QCM (8x)

Memory DIMM's (64x)

Hub Module (8x)

Each Drawer:
8TF/s, 4TBytes DRAM
24+24 Tb/s co-packaged optical IO

Hub Assembly

MLC Module

Avago microPOD™

PCIe Interconnect

PCIe Interconnect

L-Link Optical Interface
Connects 4 Nodes to form Super Node

PCIe Interconnect

D-Link Optical Interface
Connects to other Super Nodes

D-Link Optical Interface
Connects to other Super Nodes

- Ceramic with Switch ASIC, **μLGA**, & Optical Modules

•Not Field Replaceable*
•Fail in place Strategy
•~$1Gb/s in 500K volumes in 2010

**μLGA** Interposer

Optical Modules

Frame

Switch ASIC

Ceramic Substrate

7

4

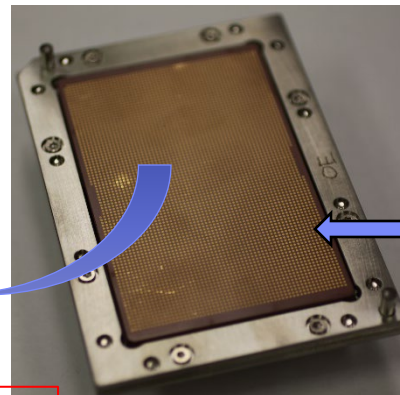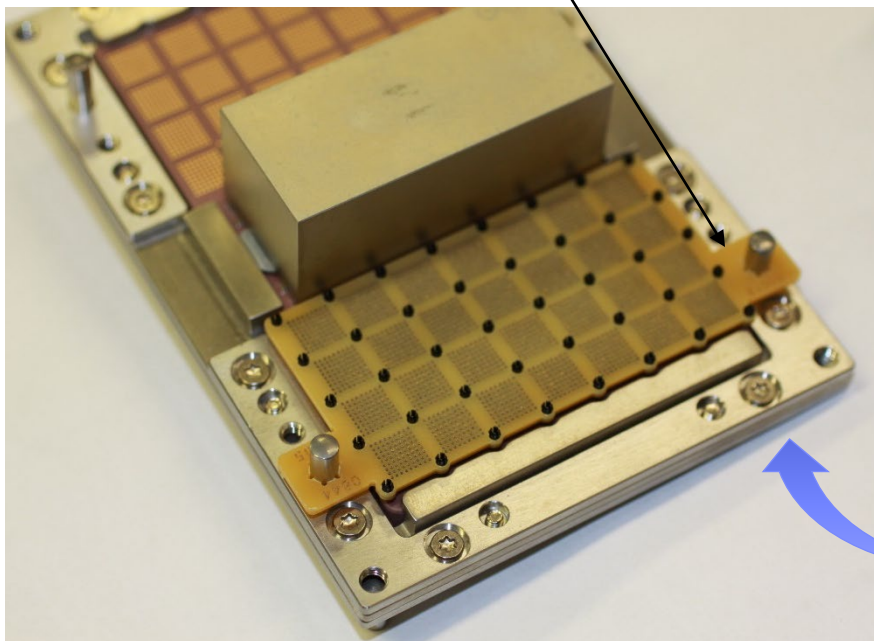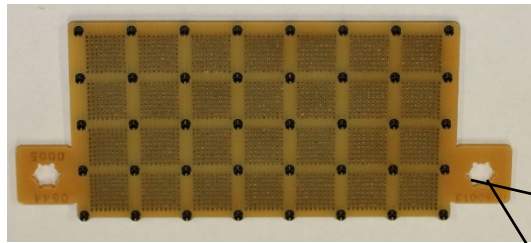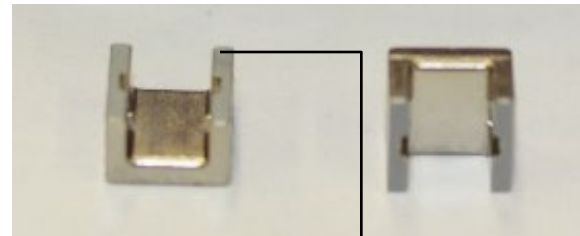- **..and after Fibers & Heat Spreaders are added**

4 Optical Fiber Ribbons
2 incoming, 2 outgoing
12 fibers each

Strain Relief for Optical Fiber Ribbons

Heat Spreaders

(Tx+Rx) Pair
(12+12) optical channels,
10 Gb/s per fiber
Implements either
1 D-link or 2 LR-links

IBM



Cold Plate mates up here

← TIM 2

OE Heat Spreader

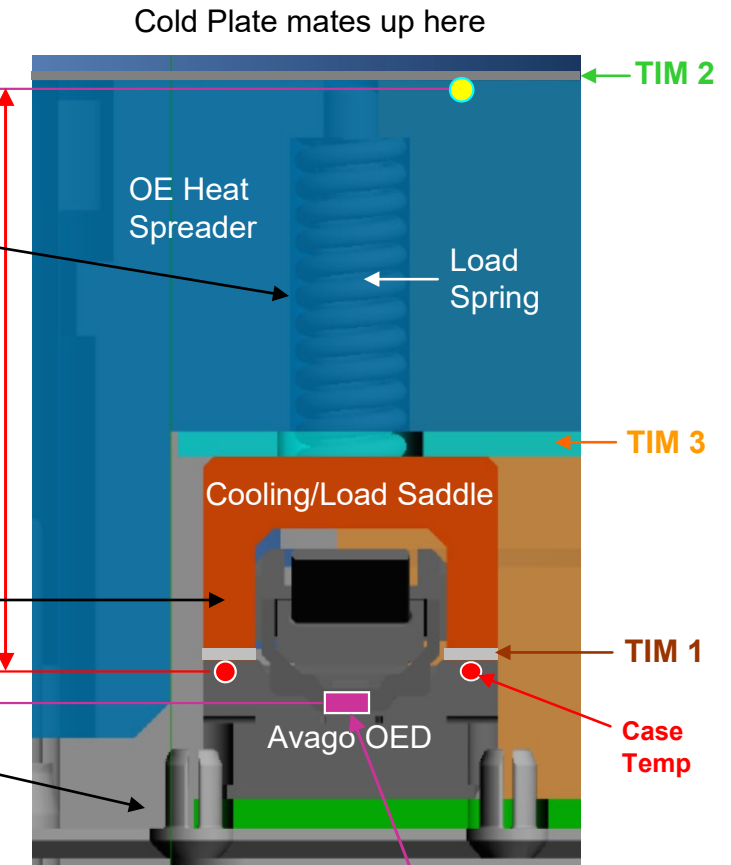Load Spring

TIM 3

Cooling/Load Saddle

TIM 1

Avago OED
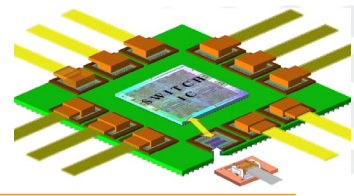
Case Temp

SHOE Thermistor

µLGA Interposer for co-packaged optics

**Total Optics Stack Thermal Resistance shown in the charts**

Actual thermal resistance owned by IBM 1st level packaging team

LGA Interposer for hub module

5

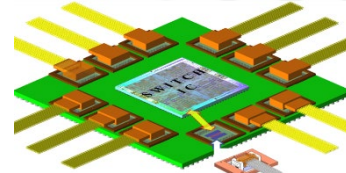*LGAs make Thermal Path and Mounting more challenging….*

© 2018 IBM Corporation

# Why are we *(still)* interested in Co-Packaging?

▶ Primarily for increasing BW from ASICs
  – Large ASICs are package pin constrained
  – Co-Packaging permits I/O from both sides of the package

▶ Reduction in Power Consumption
  – Juxtaposed die do not require high power SERDES
  – Lower power I/O cells also use less Si Area

▶ Reduction in Cost
  – Stripped down optical packages and reduced function ICs *should* cost less
  – Reduced ASIC area will have higher yield

▶ Expansion of ASIC performance
  – Instead of reducing ASIC Area and Power, other choice is to maintain size and add more functionality to the chip up to the original Power constraint
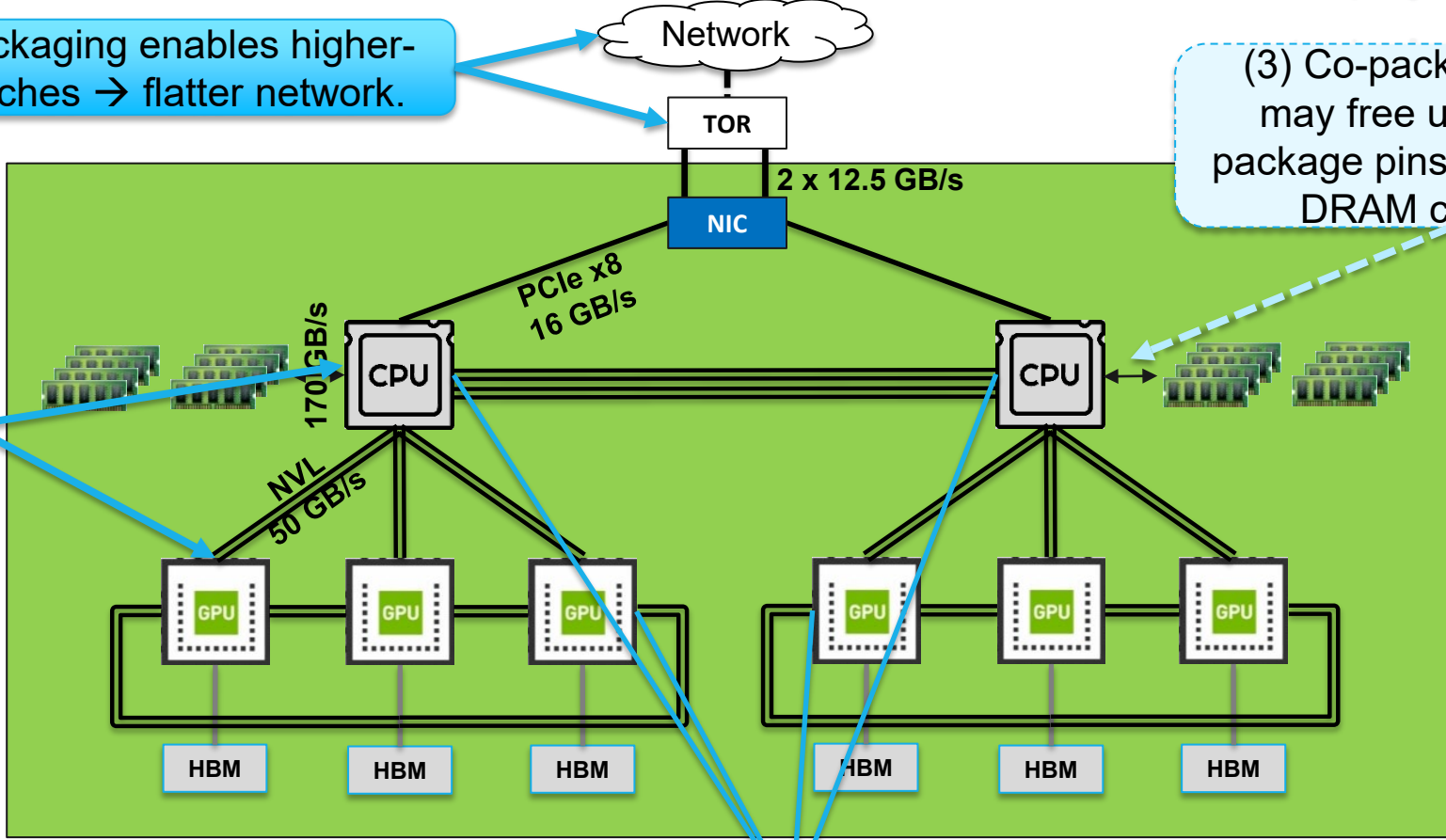  – i.e. Additional switch ports or Additional Memory Hubs

CHANGING WHAT'S POSSIBLE

# Where can Co-Packaging alleviate bandwidth bottlenecks?

Network

TOR

**(1) Co-packaging enables higher-radix switches → flatter network.**

**(3) Co-packaged optics may free up electrical package pins for additional DRAM channels.**

**2 x 12.5 GB/s**

NIC

PCIe x8
16 GB/s

**(2) Optics on CPU & GPU modules enables higher on-node BW, or moving them on different boards.**

CPU

170 GB/s

CPU

NVL
50 GB/s

GPU   GPU   GPU        GPU   GPU   GPU

HBM   HBM   HBM        HBM   HBM   HBM

**(4) Disaggregated network enables:**
I.   On-demand allocation of Accelerators
II.  NVMe / storage disaggregation
III. Memory "stealing"
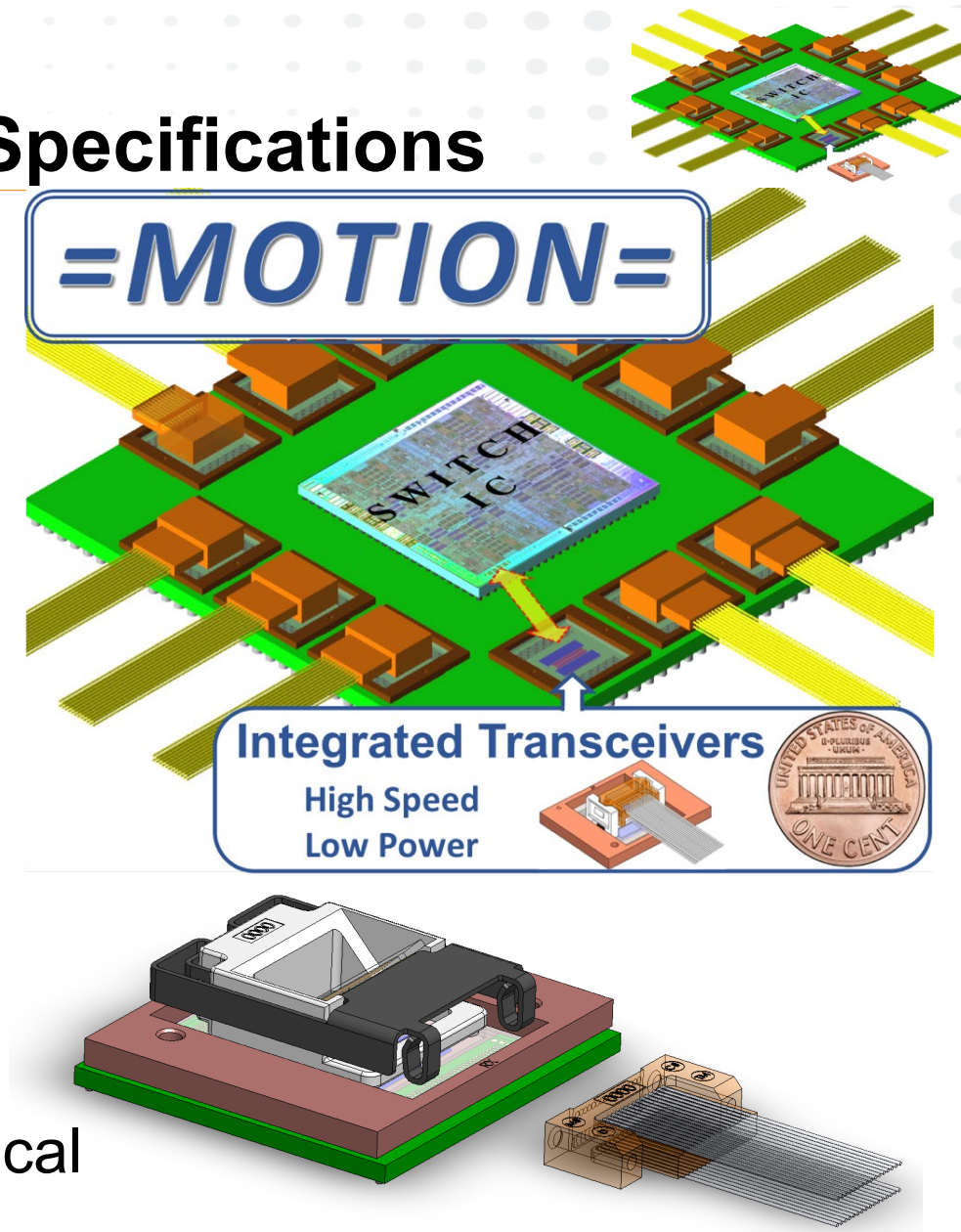
Disaggregated Network

- Protocol matters: e.g. PCIe, CXL, ...
- Functions can be on different boards

7

IBM

# MOTION Phase 1
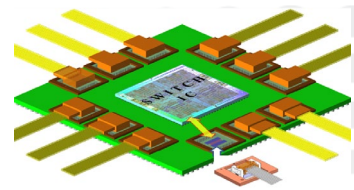# Co-packaging for CPU/GPU High-level Specifications

- ARPA-E Sponsored Project on co-packaging

- IBM and Finisar collaboration

- 56GBd NRZ; BER tested to <1E-12 pre-FEC

- $0^o$C to $70^o$C Case

- 6dB (electrical) link budget  (XSR-like)

- 2 dB optical link margin (30m w/connectors)

- Solderable onto ASIC 1st level substrate

- < 4 pJ/bit   (3.2W, 16 channels)

- W:13mm x D:13mm x H:4mm

- 25¢/Gb/s

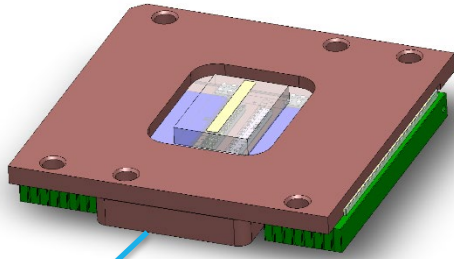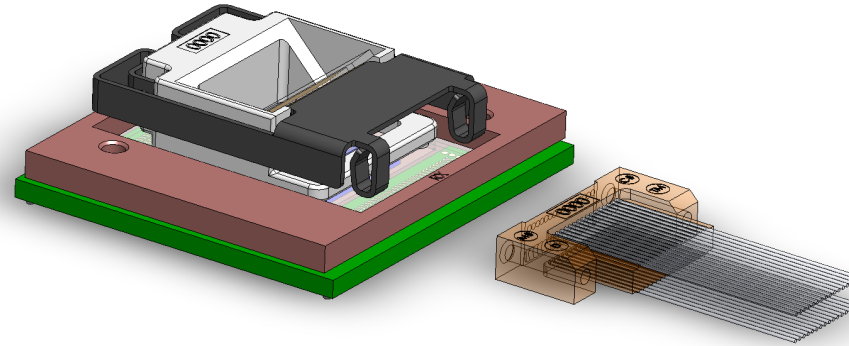*=MOTION=:* Multi-wavelength Optical
Transceivers Integrated on Node

=MOTION=

Integrated Transceivers
High Speed
Low Power

arpa·e
CHANGING WHAT'S POSSIBLE
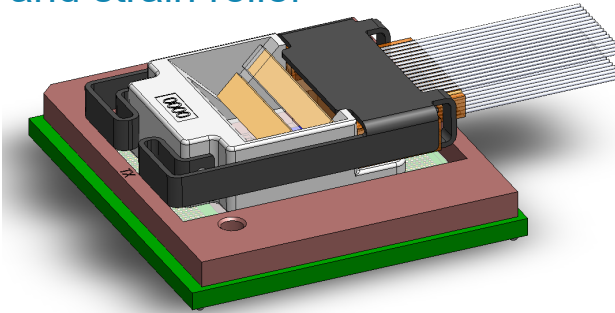
IBM

# MOTION Transceiver Package Overview



**Chip-Scale Optical Package (CSOP)**

**Final Assembly with lens and clip attached**
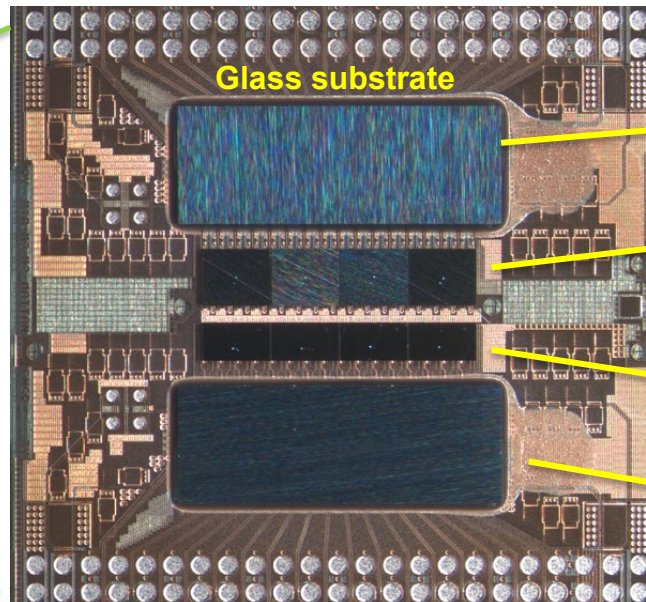
**Fully Assembled with fiber cable and strain relief**

4mm total height

**Cu Heat Spreader**

**Glass Carrier**

**SAFE ICs, VCSELs, PDs**

**Keel**

Glass substrate

SAFE Rx

4xPDs (1x4) w/Monitor

4xVCSELs (2x4) Primary+Spare
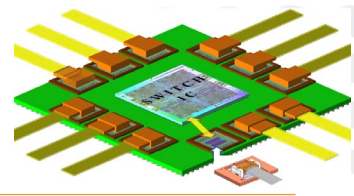
SAFE Tx

**Chip join process**
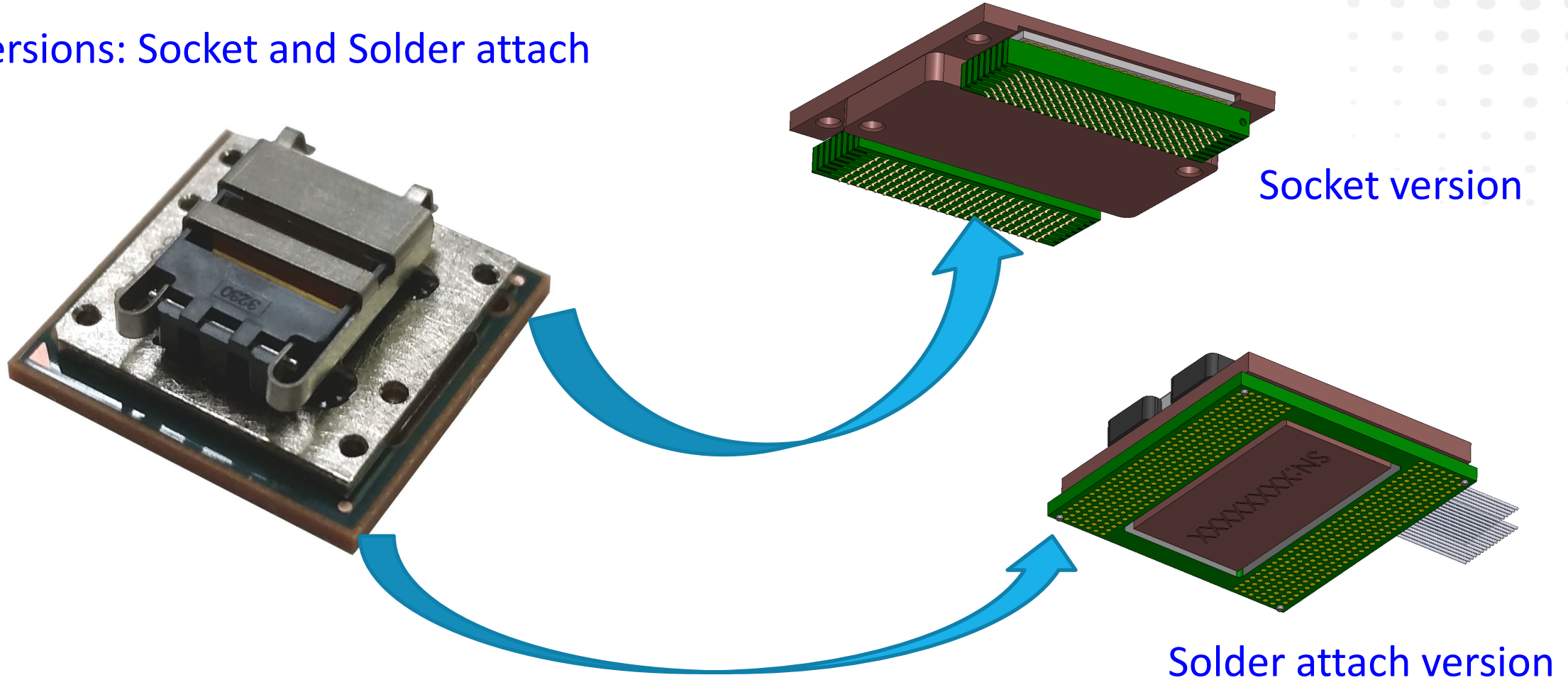- SAFE chip – Cu pillar w/SnAg cap
- PD/VCSEL – AuSn

**Underfill process**
- Structural UF for electrical chips
- Epoxy based Optical UF for OEs

# A 13mm x 13mm package with socket insert or interposer

Two Versions: Socket and Solder attach

Socket version

Solder attach version

Sockets may be helpful but incur cost as well as signal & area loss!

# Optical Transceiver Architecture

- Simplified Analog Front End (SAFE)

- SAFE2: 55 nm SiGe BiCMOS

- 16 channels × 56 Gb/s NRZ ~900 Gb/s/IC

- No retiming / CDR

- 4 pJ/bit power consumption

- Fully DC-coupled: passes 64b/66b & PRBS31

- CMOS-compatible electrical signal levels

- Built-in pattern generators and error detectors

- Only 2 Power supply voltages: 1.8v and 3.3v

# SAFE2 IC Summary

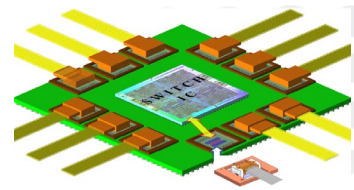## Tx Power Consumption

| Power supply | Data mode | PRBS mode |
|---|---|---|
| 1.8 V | 820 mA | 1.6 A |
| 3.3 V | 256 mA | 256 mA |
| Total power | 2.3 W | 3.7 W |
| Energy efficiency | 2.5 pJ/bit | 4.1 pJ/bit |

## Rx Power Consumption

| Power supply | Data mode | PRBS mode |
|---|---|---|
| 1.8 V | 480 mA | 1.2 A |
| 3.3 V | 120 mA | 130 mA |
| Total power | 1.3 W | 2.7 W |
| Energy efficiency | 1.5 pJ/bit | 3 pJ/bit |

**TX IC Photo**

4.63mm

1.64mm

## 50Gbps NRZ data mode

NRZ Data

51.28G clock

## 50Gbps NRZ PRBS mode

# 56GBd VCSELs and Photodiodes

- Designed and fabricated flip chip 56GBd 940nm VCSELs on a production epi and wafer fab process
    - Extended to 112G PAM4 for Phase 2
- Implemented a dual aperture structure to realize 1:1 laser sparing
- Designed and Fabricated two different 56GBd Photodiodes structures: GaAs and InP substrate based.




56Gbps (error free)




112G PAM4 (1E-8)

# MOTION Approach to Reliability: 2-to-1 Sparing

- Lack of field replacement drives stringent reliability requirements
- Laser wearout dominates: Sparing is desired

- MOTION has 2:1 laser redundancy on every channel
- Simulation shows ~1000x improvement in reliability at the end of 10 years of service →

Spare

Primary

Ch n

Ch n+3

### Reliability Simulation for 4 CSOPs on MCC

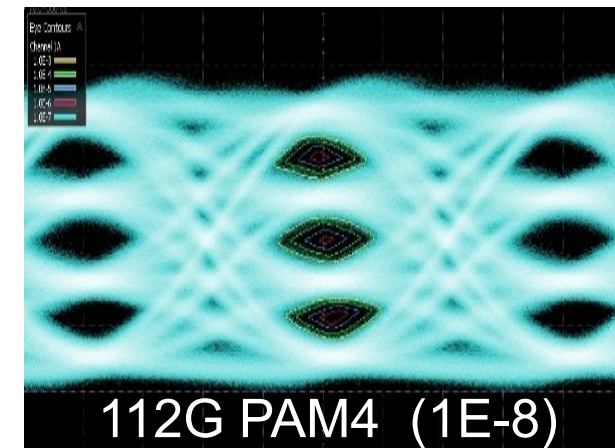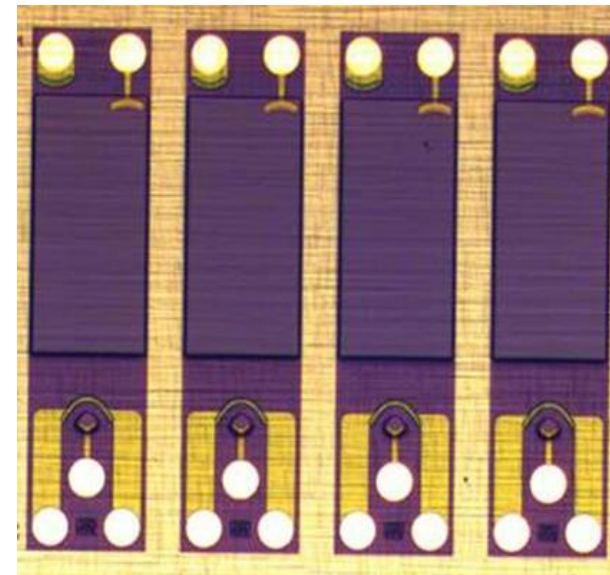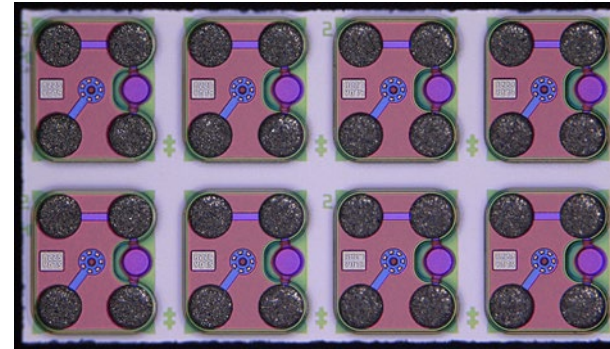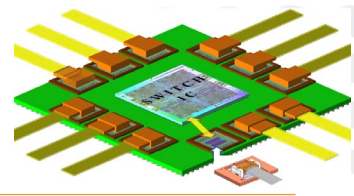Reliability: System with 64 VCSELs from Finisar at Tcase=66C,
VCSEL FIT=20, VCSEL Median Life =75 yrs, Fixed FIT=100

~77K ppm

No Sparing

~110 ppm

w/Hot Sparing

Fixed FIT only

Cumulative Failures (ppm)

$10^4$

$10^3$

$10^2$

$10^1$

$10^0$

$10^{-1}$

$10^{-2}$

Fix FIT dominates in this period | Laser wearout dominates in this period

time (yrs)

1  2  3  4  5  6  7  8  9  10

Fixed FIT: package components that can not be spared = Sum of everything with a non-zero FIT

Assumed Parameters:
VCSEL MTTF = 75 yrs, VCSEL FIT = 20
Fixed FIT = 50 FIT per module
Ibias = 9 mA

arpa·e
CHANGING WHAT'S POSSIBLE

IBM

# Fast Laser Sparing

Spare

Primary

Ch N           Ch N+3

- Switching from Primary VCSEL to Spare is observed to take < 100ns
  - Orders of magnitude improvement over pluggables!
- Waveform shows the sum of both outputs and a brief 2ns period when the primary and spare VCSELs are both on.

# MOTION Phase 2 Hardware Changes

| Parameter | Phase 1 | Phase 2 |
|---|---|---|
| Electrical Interface | 16 channels @ 56G NRZ | 32 channels @ 112G PAM4 |
| IC Technology | SiGe | CMOS |
| Optical Interface | 16 channels @ 56 G NRZ | 32 channels @ 112G PAM4 |
| # of Wavelengths | 1 | 2 |
| # of Fibers | 16 Tx + 16 Rx | 16 Tx + 16 Rx |
| Fiber Type | 50/125 MMF | 50/125 MMF |
| Package I/O Pitch | 400um | 300um |
| Glass Carrier Size | 13x13mm | 13x13mm |
| Energy consumption | 4 pJ/bit | 2 pJ/bit |
| Projected cost | 25¢/Gig | TBD but <25¢/Gig |
| Laminate interface | Soldered or LGA | Soldered or LGA |

IBM

# Phase 2: IBM SYSTEMS TECHNOLOGY EVALUATION



- Goal: <u>To assess the technology readiness of co-packaged optics</u>
  - Optical transceivers soldered directly on the top surface of a laminate package will be built
- The IBM Systems group will perform this technology evaluation focusing on the thermal & mechanical robustness of packaging:
  - Four (4) optical transceivers and one (1) test site die on the top of a single FC-PLGA laminate, assembled with a thermal lid
- Evaluations:
  - Preconditioning
  - Deep Thermal Cycle
  - Accelerated Thermal Cycling
  - Thermal Aging
  - Temperature & Humidity
  - Temperature, Humidity & Bias
  - Power Cycling

  - Low Temperature Storage,
  - Shock & Vibe



*Demonstrating a viable path to system integration using established IBM Server Group processes*

# IBM Processor module with MOTION Devices

Phase 2 hardware

~80mm

~70mm

Processor Laminate with 4 MOTION devices, w/o heat spreader

Processor Laminate with 4 MOTION devices, w/ heat spreader

- Large IC and MOTION devices are assembled with a single solder reflow step.
- Two different orientations for fiber escape: Overlapping and Non-overlapping

# Ribbon Fiber Orientation experiment

- **Ribbons parallel to air flow**
  - ➢ Stacked ribbons
  - ➢ Single strain relief



Air Flow

- **Ribbons perpendicular to air flow**
- ➢ **Individual Strain Relief**



Data Center Switches do not have this option

# Thermomechanical package simulations

Lid

TIM

TIM 1.5

sealband

OEDs

Testsite Chip

JohnsonOED Chip Carrier

Components in Finite Element Analysis simulation

Challenge: TIM to cover a large gap with high variability

XSEC of Lid

59μm   400μm   78μm   460μm

198μm

CTE mismatches between ASIC, laminate, and OED result in considerable warpage prior to the lid attach but the range is still acceptable

# Summary of Reliability Testing

Thermal readout at different interval of stressing, the results are reasonably stable

| Stress | Samples | T/S -40/60°C | Time 0 | Readout 1 | Readout 2 | Readout 3 | Readout 4 | Readout 5 | Readout 6 |
|--------|---------|--------------|--------|-----------|-----------|-----------|-----------|-----------|-----------|
| DTC -40/125°C | 5 | 5X | T, R | 250 c | 500 c | 750 c | 1000 c | 1250 c | 1500 c |
| ATC 0/100°C | 5 | 5X | T, R | 500 c | 1000 c | 1500 c | 2000 c | 2500 c | 3000 c |
| HTS 125°C | 5 | 5X | T | 494 hrs. | 989 hrs. | 1486 hrs. | 2009 hrs. | | |
| T&H 85°C/85% RH | 5 | 5X | T | 279 hrs. | 438 hrs. | 721 hrs. | 1002 hrs. | | |

Laminate feature fail rate vs DTC cycles.

| Laminate features | | | | | | | |
|---|---|---|---|---|---|---|---|
| C4 nets, LW/LS, RFPs, Die edge crack sensors, FC&BC side resin cracking sensor, Power and Ground | | | | | | | |
| | 250 | 500 | 750 | 1000 | 1250 | 1500 | 1750 |
| DTC | 0/12 | 0/12 | 0/12 | 0/12 | 3/10 | 3/10 | |
| DTC on card | 0/5 | 0/5 | 0/5 | 0/5 | 0/5 | 0/5 | 3/5 |

The first fail occurs after 1000 cycles of DTC which is not a reliability concern (DTC is a highly accelerated stress as opposed to field application).

# The benefits of higher-radix switches enabled by MOTION



**SUMMIT-like: 1620 36x36 SWITCH MODULES**

CORE SWITCH  18x  CORE SWITCH
↓648x
↑18x
↓18x
TOR  TOR  ...  648x  TOR  TOR
2 links / node
NIC  ...  NIC
P  P  324x  P  P
A A A A A A    A A A A A A
18x

SPINE  SPINE  18x  SPINE
↓36x                    (1x)
↑18x
LEAF  LEAF  LEAF  ... 36x  LEAF
↓18x

2-level fat tree in a box built from 54 36-port switches

36 leaf switches x ↓18 ports → 648 ports

**MOTION: 1280 128x128 SWITCH MODULES**

CORE SWITCH  64x  CORE SWITCH
↓512x
↑64x
↓64x
TOR  TOR  512x  TOR  TOR
6 links / node
NIC  ...  NIC
P  324x  P    P  P
A A A A A A    A A A A A A
18x

SPINE  SPINE  4x  SPINE
↓128x                    (16x)
↑64x
LEAF  LEAF  LEAF  8x  LEAF
↓64x

2-level fat tree in a box built from 12 128-port switches

8 leaf switches x ↓64 ports → 512 ports

- 3x more network end points for 21% fewer switch modules
- 2.8x higher bisection BW for 100 Gb/s per port (11.2x for 400 Gb/s)
- MOTION opens the way to direct-network-attached accelerators

arpa·e
CHANGING WHAT'S POSSIBLE

22

# Network performance analysis: MOTION vs. Summit

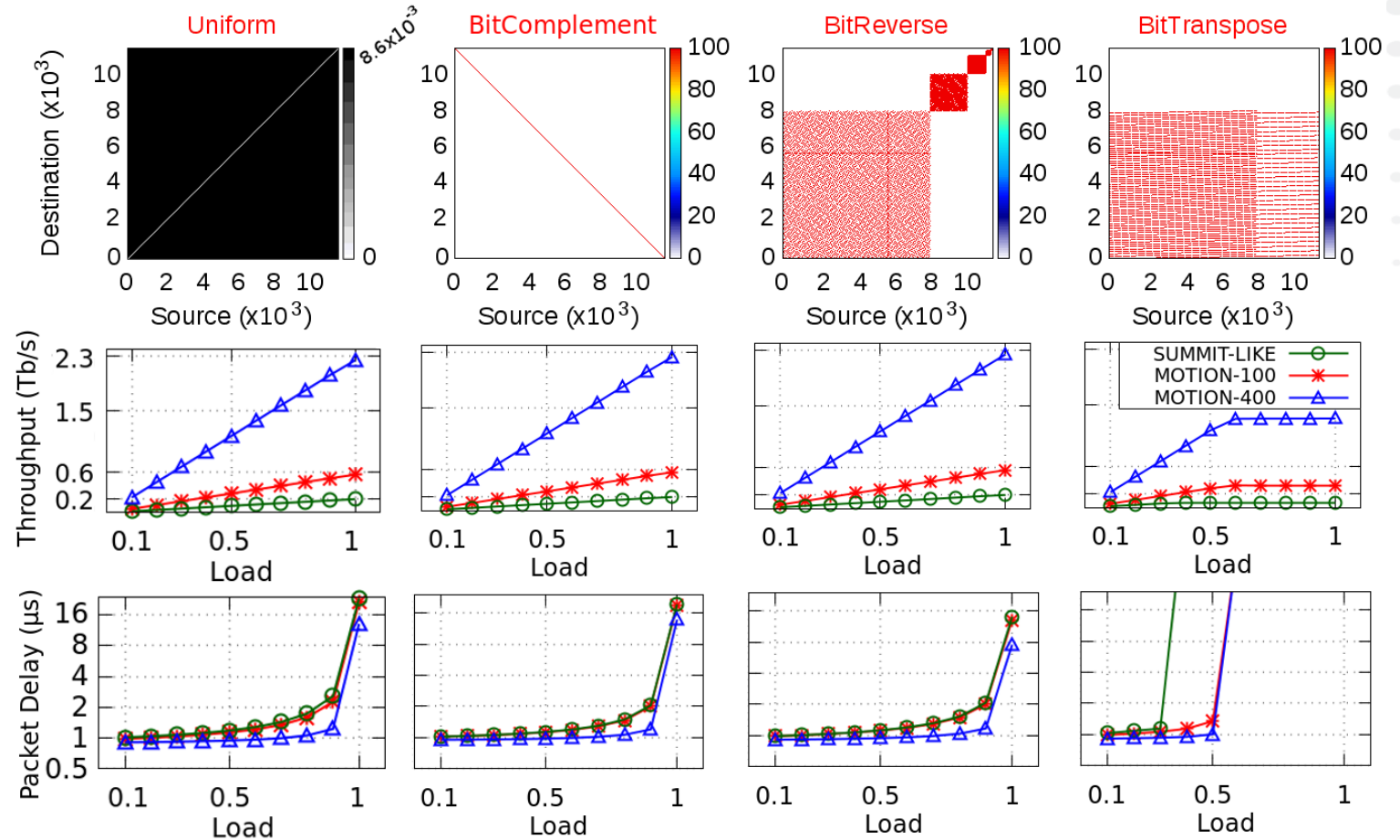## Venus discrete event network simulator

| Traffic and Network | |
|---|---|
| **Generator** | |
| RX buffer size | Infinite |
| Message size | 1024 B |
| Generation distribution | Bernoulli |
| Data rate | 100/400 Gbps |
| Load | [0.1-1] |
| **Adapter/Switch** | |
| Type | InfiniBand |
| Data rate per link | 100/400 Gbps |
| Delay | 100 ns |
| Switch Buffer per port | 128 KB |
| Packet size | 1024 B |
| Routing algorithm | Random |



- *Uniform, BitComplement, BitReverse*: Linear throughput increase for all systems
  → 2.8x and 11.2x higher throughput for 100 and 400 Gb/s data rates
- *BitTranspose*: earlier saturation due to the significantly fewer destination nodes
  → 4.3x and 17.2x higher throughput for 100 and 400 Gb/s data rates
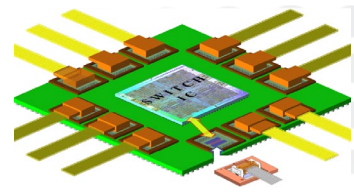- MOTION-400: best mean packet delay for all cases

# Challenges for Co-packaged Optics

- Reliability
- Field Replacement/Serviceability
  - Fail-in-place Strategy?
- Yield
  - Who is responsibility for final yield?
- Assembly
  - Who does what and when?
- Standards and/or MSAs
  - Minimum Time for Standards seems to be > 2 years
  - Proprietary solutions likely to emerge first
- Compatible Technologies
  - MMF or SMF
- Field Upgradeable Firmware!

CHANGING WHAT'S POSSIBLE

# Acknowledgements:  MOTION Phase 1 & 2 Team members & Sponsor

- IBM Research
  - C. Baks, A. Benner, R. Budd, T. Dickson, F. Doany, W. Lee, M. Meghelli, P. Pepeljugoski, J. Proesel, M. Taubenblatt, L. Schares, M. Schultz, P. Maniotis, P. Stark, H. Ainspan, Z. Toprak Deniz, S. Dhawan, T. Dickson, N. Dupuis, P. Francese, B. Sadhu, M. Kossel, T. Morf, M. Brändli, S. Rylov, M. Cochet, C. Ozdag, A. Watanabe, H. Rahmani, D. Kuchta,

- IBM Bromont
  - L-M. Achard, P. Fortier, C. Dufort, E. Tucotte, C. Bureau, M. Pion, Y. Cossette, P. Ducharme, S. Desputeau, A. Janta-Polczynski, P. Minier, G. Jutras, P. McInnes, S. Whitehead, B. Sow

- IBM Server Packaging
  B. Parikh, S. Ostrander,  S. Li, C. Setzer, H. Toy, J. Ross, K. Lange, M. Kapfhammer, B. Meiswinkel, C. Muzzy, E. Steiner, D. Smith, T. Saunders, G. Pomerantz, J. Coffin, K. Marston, K. Smith, T. Ahmed, L. Rapp, Y.Yao, T. Wassick, M. Warbrick, T. Lombardi, B. Singh, C. Walker, S. Iruvanti, D. Yannitty, P. Ramaglia, T. Weiss, M. Interrante, J. Sorbello, C. Arvin, M.  Stalter, A. Perez, P. Torbet, T. Olowofela, J. Bunt, M. Fisher, T. Childress, J. Mingo, C. Savoy, S. Ruiz, D.Kohler, R. Seifts, R. Olson, H. Polgrean, J. Rowland, C. Thomas, E. Kastberg, A. Schetter, D. Babcock, A. Greenberg, D. Lord, R. Rodriguez, C. Taylor

- IBM Server Development
  - D. Becker, R. Laning, D. Dreps, M. Hoffmeyer, J. Eagle, F. Gholami, K. O'Connell, S. Canfield, S. Chun, R. Frota

- IBM Supply Chain Engineering
  - H Bagheri, K. Akasofu, C. Grosskopf, A. Tiano, T. Sass, E. Mallery

- II-VI Finisar Corporation
  F. Flens, D. Case, P. Chen, J. Glover, C. Kocot, K. Koski, G. Light, T. Nguyen, S. Pandy, S Quadery, K. Szczerba, B. Wang, P Westbergh, S. Pandey, H. Hayashigatani  (131)

- Texas A&M University
  - N. Kim, A. Kumar, T. Liu, I. Yi, S. Palermo

CHANGING WHAT'S POSSIBLE