# Looking Beyond 400G

## A System Vendor Perspective

### Beyond 400 Gb/s Ethernet Study Group

Rakesh Chopra
Cisco Fellow
February 8, 2020

rakchopr@cisco.com
www.linkedin.com/in/rakesh-chopra/
@Rakesh_Chopra1

... Many thanks to Cisco Engineers and Insightful Customers ...

# System Architectures

## Fixed*

Standalone Switch — Retimer — Port / Port

* – Can be interconnected to create a disaggregated chassis

## Centralized

Standalone Switch — Mux / Mux — Port / Port

Optional Redundancy

## Distributed

FE Switch — Retimer — LC Switch — Retimer — Port / Port

LR | VSR

# Relentless Advancement – Switch Silicon Bandwidth

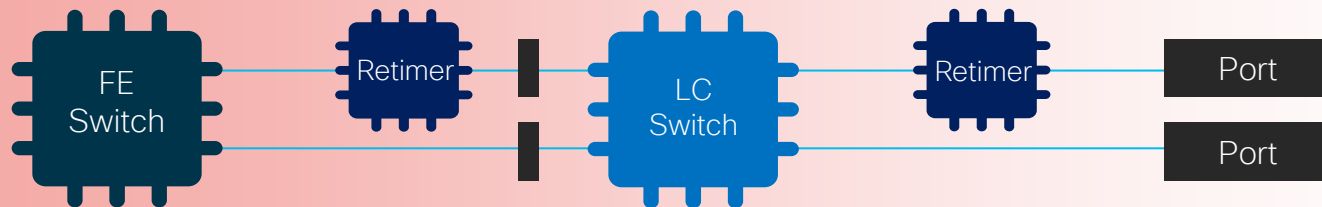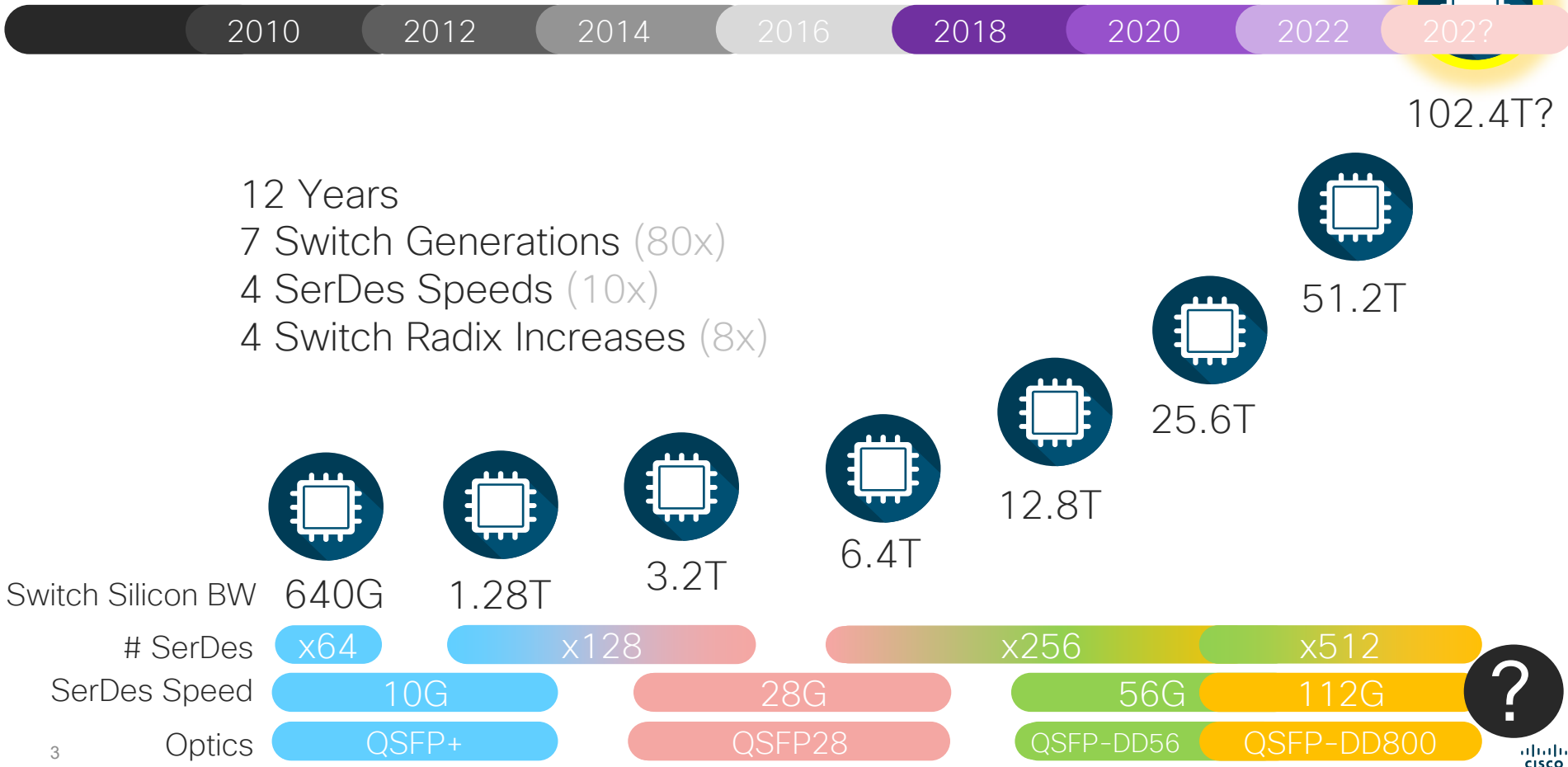Represents a combination of multiple chip families and architectures to provide historical context and future projections

| 2010 | 2012 | 2014 | 2016 | 2018 | 2020 | 2022 | 202? |

102.4T?

12 Years
7 Switch Generations (80x)
4 SerDes Speeds (10x)
4 Switch Radix Increases (8x)

51.2T

25.6T

12.8T

6.4T

640G          1.28T          3.2T

| Switch Silicon BW | 640G | 1.28T | 3.2T | 6.4T | 12.8T | 25.6T | 51.2T |
| --- | --- | --- | --- | --- | --- | --- | --- |
| # SerDes | x64 | | x128 | | x256 | | x512 |
| SerDes Speed | 10G | | 28G | | 56G | | 112G |
| Optics | QSFP+ | | QSFP28 | | QSFP-DD56 | | QSFP-DD800 |

cisco

# Relentless Advancement – 80x BW over 12 Years
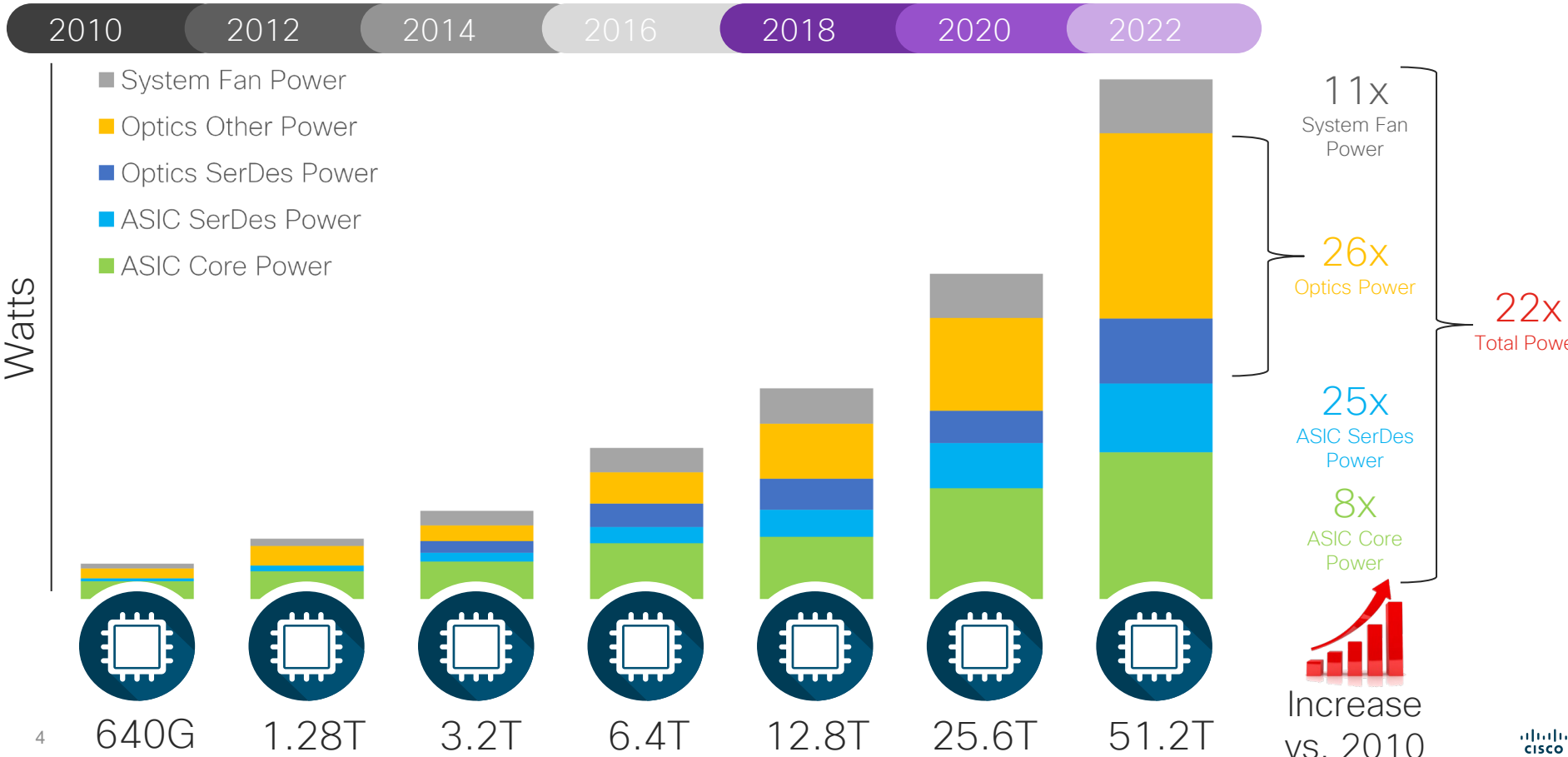
Represents a combination of multiple chip families and architectures to provide historical context and future projections
Fixed Box Power Breakdown
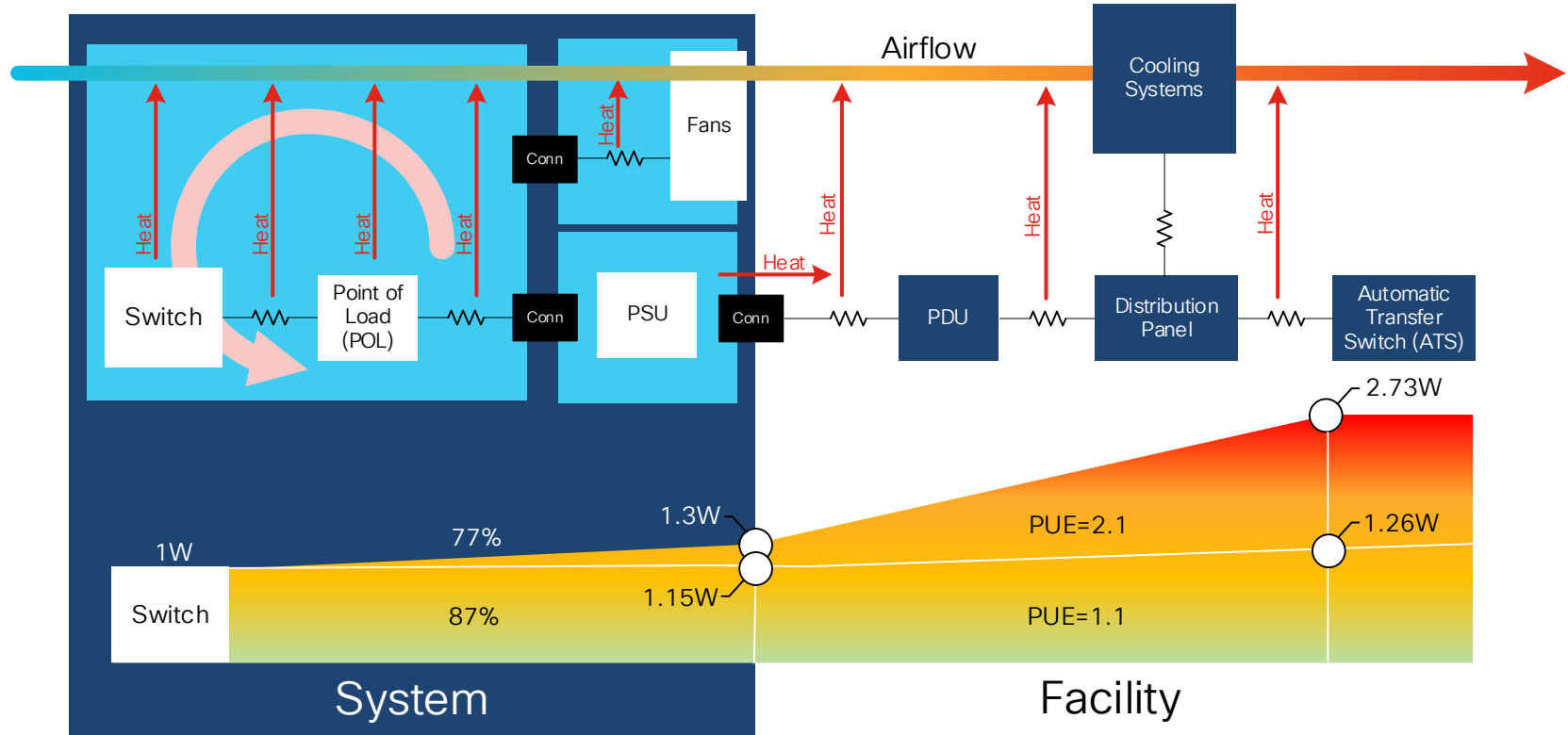Retimer Power and other system components not included

# The Multiplication Effect of a Watt

# Power is THE Problem to Solve

Apollo 13 – Universal Pictures

- ❌ Limits what we can build

- ❌ Limits what can be deployed

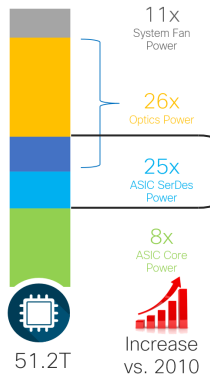- ❌ Limits what our planet can sustain

"Power is Everything"*

John Aaron– Apollo 13 Flight Controller

Adopt a power first design and deployment methodology

* – Thanks to Kraig Owen for the reference

# Co-packaged Optics Is Inevitable
## Power savings drives requirement



11x
System Fan Power

26x
Optics Power

25x
ASIC SerDes Power

8x
ASIC Core Power

51.2T

Increase vs. 2010

Must minimize SerDes power

SerDes power increases with distance

*Trends plot from premier CMOS wireline 2018 conference*

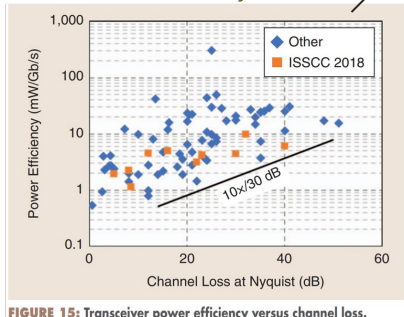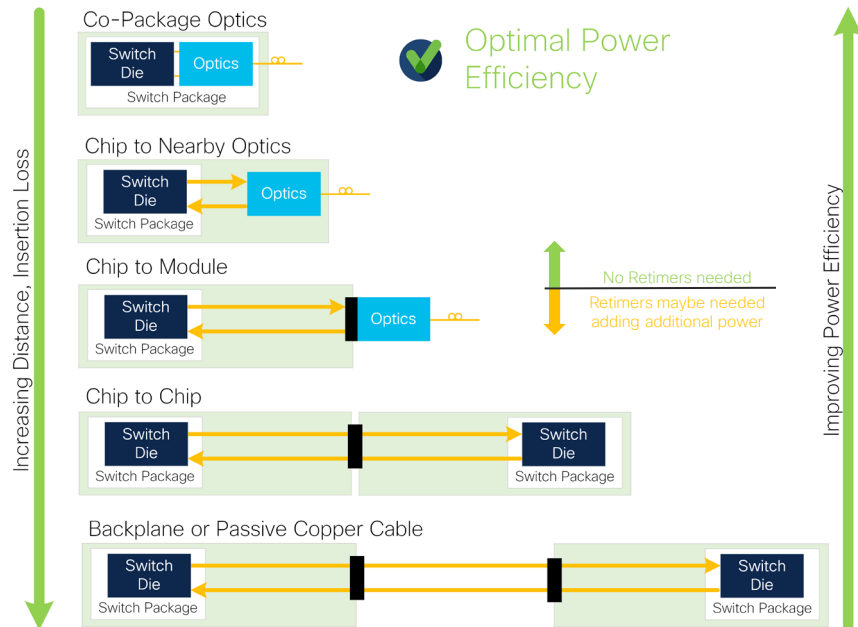**FIGURE 15:** Transceiver power efficiency versus channel loss.

Power Efficiency (mW/Gbi/s)

Channel Loss at Nyquist (dB)

Other
ISSCC 2018

10x/30 dB

Daly, Denis C., Laura C. Fujino, and Kenneth C. Smith. "Through the Looking Glass-The 2018 Edition: Trends in Solid-State Circuits from the 65th ISSCC." *IEEE Solid-State Circuits Magazine* 10.1 (2018): 30-46.

## Architectural Approach to Power Optimization

Increasing Distance, Insertion Loss

Improving Power Efficiency

Co-Package Optics

Switch Die | Optics

Switch Package

Optimal Power Efficiency

Chip to Nearby Optics

Switch Die | Optics

Switch Package

Chip to Module

Switch Die | Optics

Switch Package

No Retimers needed

Retimers maybe needed adding additional power

Chip to Chip

Switch Die

Switch Package

Switch Die

Switch Package

Backplane or Passive Copper Cable

Switch Die

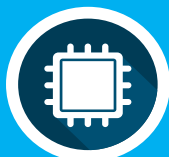Switch Package

Switch Die

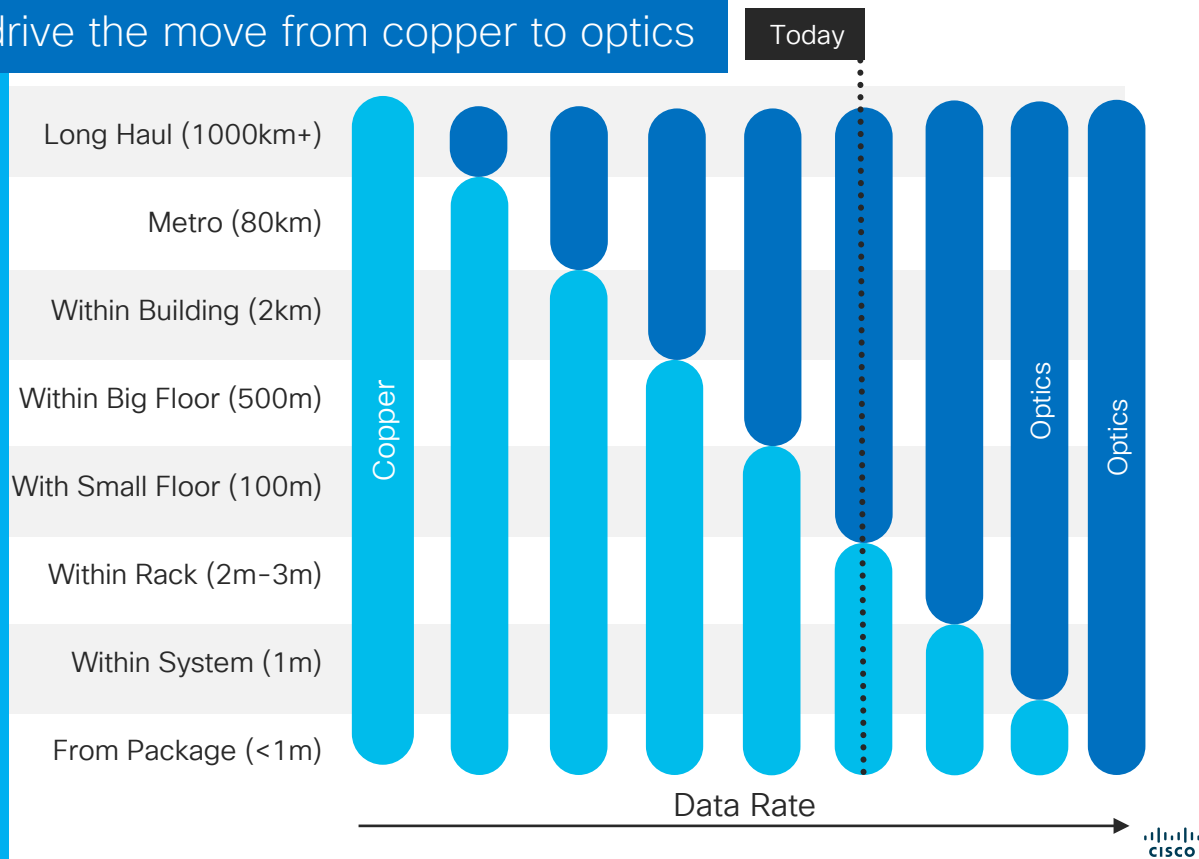Switch Package

# Co-packaged Optics Is Inevitable
and viable in the 51.2T generation

Higher data rates and distance drive the move from copper to optics

Today

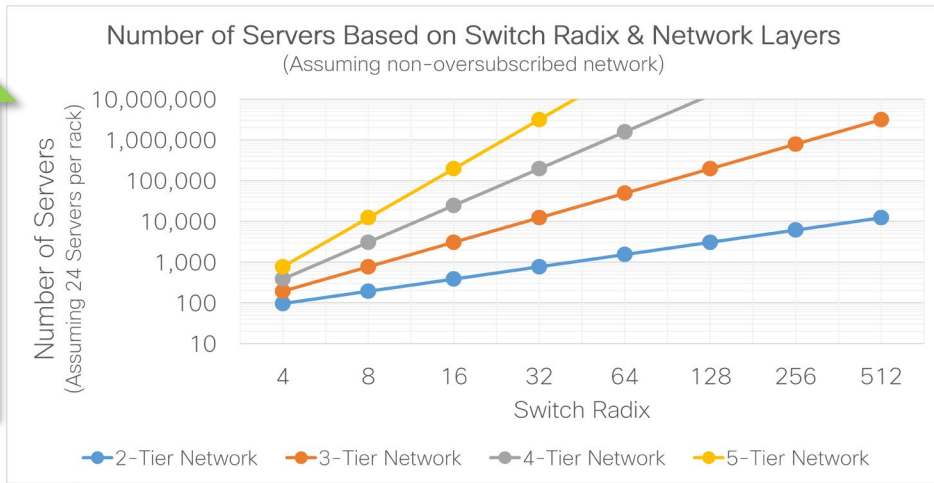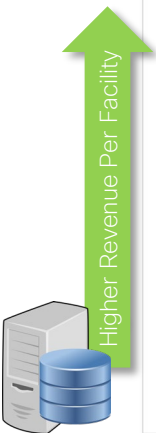Future innovations will only be possible with silicon and optical integration

**Si** Silicon **+** **Op** Optics

51.2T

| | Copper | | | | | | Optics | Optics |
|---|---|---|---|---|---|---|---|---|
| Long Haul (1000km+) | | | | | | | | |
| Metro (80km) | | | | | | | | |
| Within Building (2km) | | | | | | | | |
| Within Big Floor (500m) | | | | | | | | |
| With Small Floor (100m) | | | | | | | | |
| Within Rack (2m–3m) | | | | | | | | |
| Within System (1m) | | | | | | | | |
| From Package (<1m) | | | | | | | | |

Data Rate

# Building Your Data Center
## Impact of Switch Radix



Doubling Radix adds 2x-16x more servers

Server Multiplier Effect with Doubling Switch Radix

16x — 5-Tier Network
8x — 4-Tier Network
4x — 3-Tier Network
2x — 2-Tier Network

Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

Number of Servers (Assuming 24 Servers per rack)

Switch Radix: 4, 8, 16, 32, 64, 128, 256, 512

- 2-Tier Network
- 3-Tier Network
- 4-Tier Network
- 5-Tier Network

Higher Revenue Per Facility

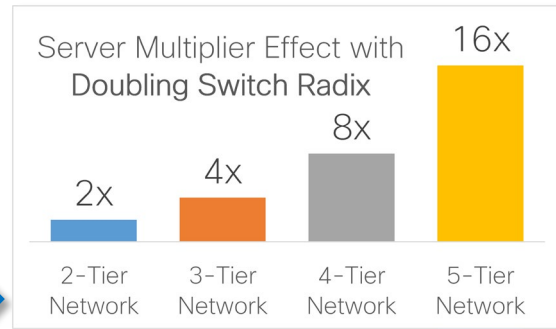Scale Out — Wider Radix

Graph concept leveraged from R. Nagarajan, Ilya Lyubomirsky, "Next-Gen Data Center Interconnects: The Race to 800G" Adjusted to hold servers per rack constant
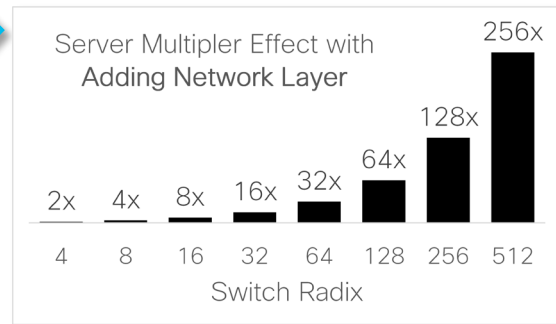
Adding a layer adds 2x-256x more servers

Scale Up — More Layers

Server Multipler Effect with Adding Network Layer

2x (4), 4x (8), 8x (16), 16x (32), 32x (64), 64x (128), 128x (256), 256x (512)

Switch Radix
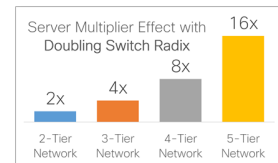
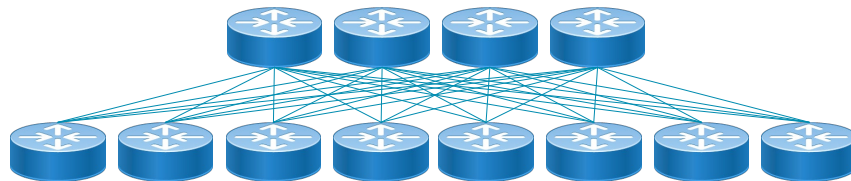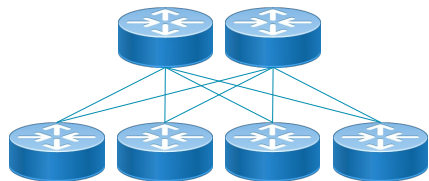# Impact of Switch Radix
## Case against increasing Switch Radix

Doubling Radix adds 2x-16x more servers depending on the layers in the network


Server Multiplier Effect with **Doubling Switch Radix**
2x — 2-Tier Network
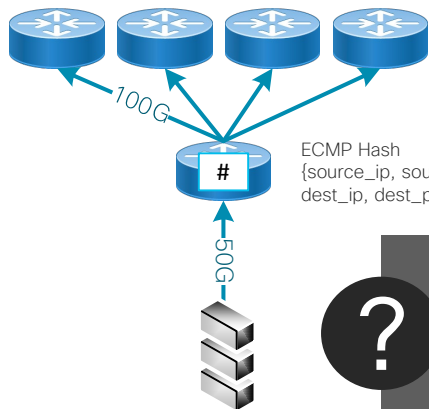4x — 3-Tier Network
8x — 4-Tier Network
16x — 5-Tier Network

Increasing radix adds cabling complexity, cost and weight



❌ Complicated Cabling

Increasing radix decreases link (mac) speed for the same switch bandwidth
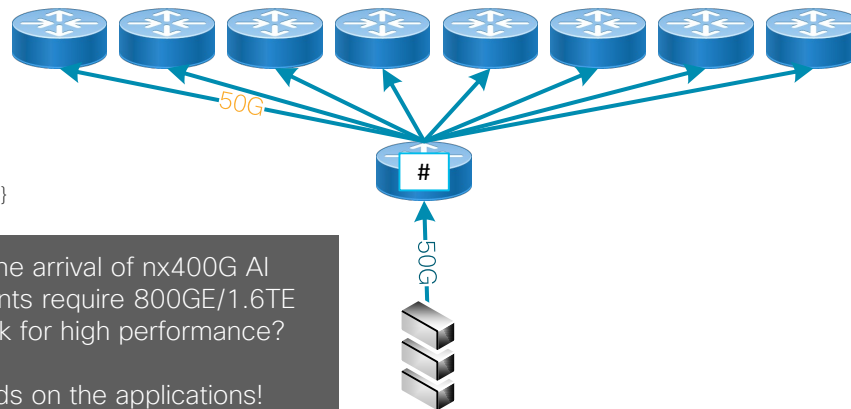As the flow speed approaches the link speed link utilization decreases



100G

ECMP Hash
{source_ip, source_port,
dest_ip, dest_port, protocol}

50G

50G

50G

? Does the arrival of nx400G AI endpoints require 800GE/1.6TE network for high performance?

Depends on the applications!

| 12.8T | |
|-------|-------|
| x32 | 400GE |
| x64 | 200GE |
| x128 | 100GE |
| x256 | 50GE |

❌ Poor link utilization

# LAG vs. ECMP
## The Basic Topology

**Service Provider**

**Data Center**

Lag Bundle

ECMP

Create a "fat pipe" between two boxes

Create multipe "equal" pipes between many boxes

Both use the same hash function and expose link utilization issues

Note : Service Providers use ECMP as well but not in an equivalent fundamental way

# LAG vs. ECMP

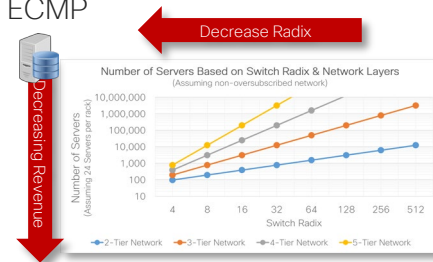The advantages of higher speed MACs aren't as clear as they used to be

## Service Provider

Increase MAC speed

Create a "fat pipe" between two boxes with no hash inefficiencies

No downside to replacing a LAG bundle with a higher speed Ethernet MAC

## Data Center

Increase MAC speed

Shrink the number of boxes we can connect to

Downside for higher speed MAC with ECMP

For same speed silicon:
- As you increase MAC speed
- Decrease your radix
- Decreases your switches per DC
- Lower revenue potential

Decrease Radix

Decreasing Revenue

Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

Number of Servers
(Assuming 24 Servers per rack)

10,000,000
1,000,000
100,000
10,000
1,000
100
10

Switch Radix
4    8    16    32    64    128    256    512

2-Tier Network    3-Tier Network    4-Tier Network    5-Tier Network

# Impact of Switch Radix
## Case against adding Network Layers

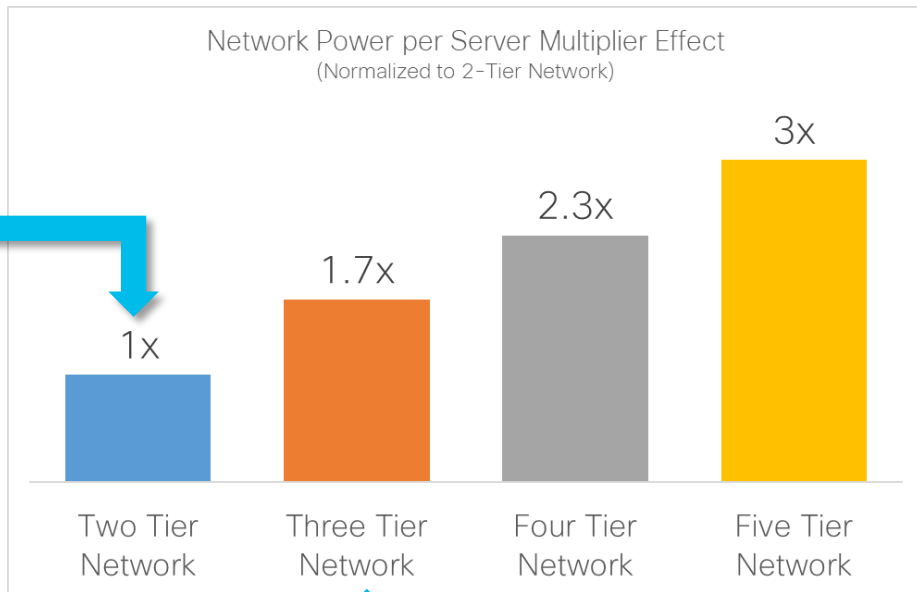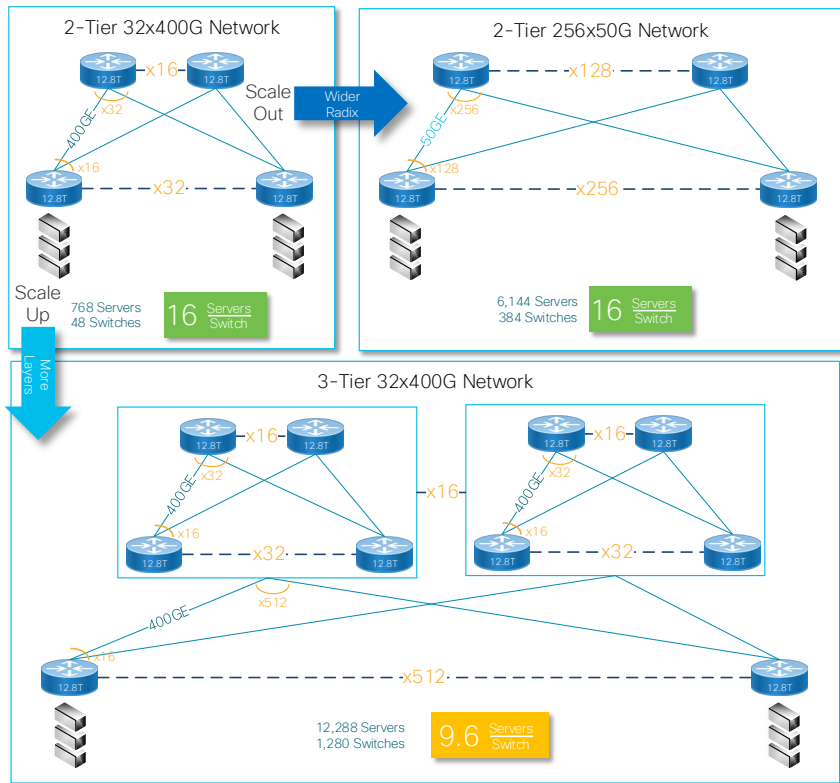Adding a layer adds 2x-256x more servers depending on the switch radix


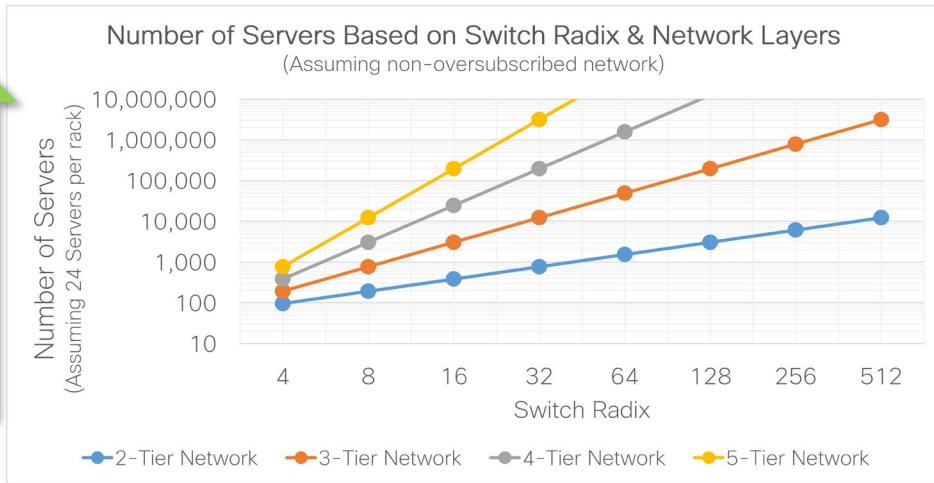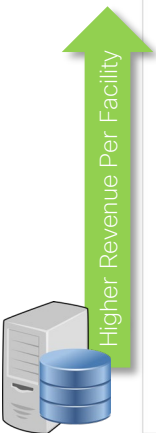Server Multipler Effect with Adding Network Layer

### Increasing layers adds network cost and power
more switches and optics per server



**2-Tier 32x400G Network**
x16
x32
400GE
x16
x32
12.8T
Scale Out
Wider Radix
Scale Up
More Layers

768 Servers
48 Switches
**16** Servers/Switch

**2-Tier 256x50G Network**
x128
256
50GE
x128
x256
12.8T

6,144 Servers
384 Switches
**16** Servers/Switch

**3-Tier 32x400G Network**
x16
x32
400GE
x16
x32
x16
x16
x32
400GE
x16
x32
x512
400GE
x16
x512
12.8T

12,288 Servers
1,280 Switches
**9.6** Servers/Switch

## Network Power per Server Multiplier Effect
(Normalized to 2-Tier Network)

3x
2.3x
1.7x
1x

Two Tier Network | Three Tier Network | Four Tier Network | Five Tier Network

Assuming no extra components needed to scale out (reverse gearboxes, etc….)
Ignoring ECMP hash efficiency impact for "goodput" of the network

cisco

# Building Your Data Center
## Impact of Switch Radix

**Doubling Radix** adds **2x-16x more servers**

Higher Revenue Per Facility

### Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

Number of Servers (Assuming 24 Servers per rack)

| Switch Radix |
|---|

10,000,000 — 1,000,000 — 100,000 — 10,000 — 1,000 — 100 — 10

Switch Radix: 4, 8, 16, 32, 64, 128, 256, 512

- ● 2-Tier Network
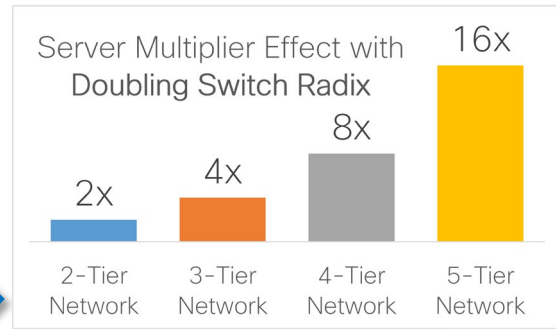- ● 3-Tier Network
- ● 4-Tier Network
- ● 5-Tier Network

Graph concept leveraged from R. Nagarajan, Ilya Lyubomirsky, "Next-Gen Data Center Interconnects: The Race to 800G"
Adjusted to hold servers per rack constant

**Scale Out**
Wider Radix

### Server Multiplier Effect with Doubling Switch Radix

16x — 8x — 4x — 2x

| 2-Tier Network | 3-Tier Network | 4-Tier Network | 5-Tier Network |
|---|---|---|---|

✔ Power Efficiency

**Adding a layer** adds **2x-256x more servers**

**Scale Up**
More Layers

### Server Multiplier Effect with Adding Network Layer

2x, 4x, 8x, 16x, 32x, 64x, 128x, 256x

Switch Radix: 4, 8, 16, 32, 64, 128, 256, 512

✔ Link Efficiency

There is no free lunch, every engineering choice has trade-offs

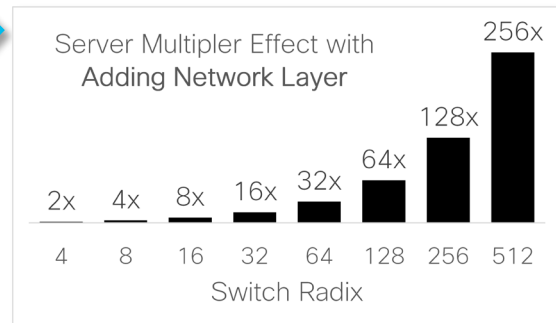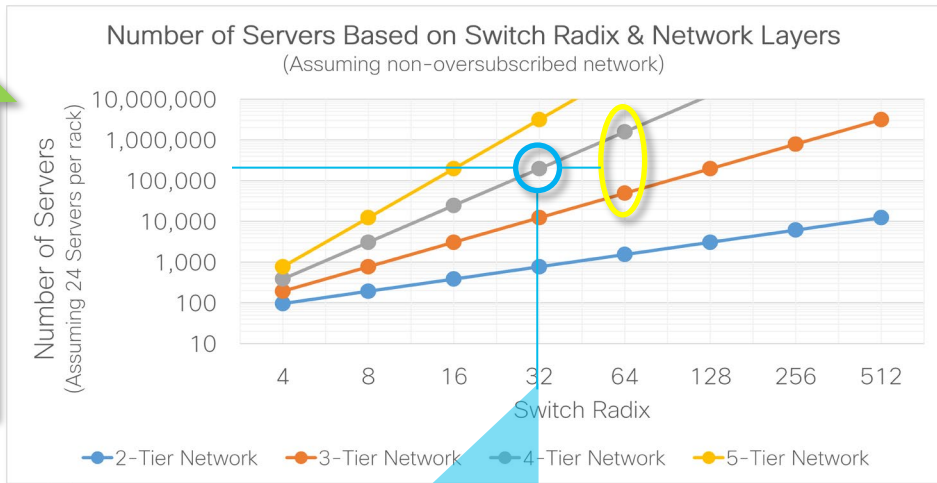Balancing act between radix, MAC speed, and layers in the network...

# Building Your Data Center
## Scale-Out vs. Scale-Up– A Balancing Act

### Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

**Higher Revenue Per Facility**

**Number of Servers** (Assuming 24 Servers per rack)

| Y-axis values |
|---|
| 10,000,000 |
| 1,000,000 |
| 100,000 |
| 10,000 |
| 1,000 |
| 100 |
| 10 |

Switch Radix: 4, 8, 16, 32, 64, 128, 256, 512

Legend:
- 2-Tier Network
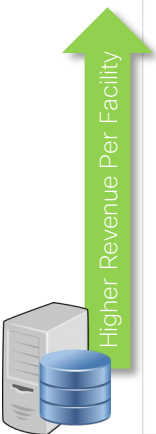- 3-Tier Network
- 4-Tier Network
- 5-Tier Network

Graph concept leveraged from R. Nagarajan, Ilya Lyubomirsky, "Next-Gen Data Center Interconnects: The Race to 800G"
Adjusted to hold servers per rack constant

| Switch BW | SerDes | Radix x32 | |
|---|---|---|---|
| 12.8T | 56G | 400GE | x8 |
| 25.6T | 112G | 800GE | x8 |
| 51.2T | 112G | 1.6TE | x16 |
| 102.4T? | 212G? | 3.2TE | x16 |

- x32 and x128 radix are prominent today
  - Ethernet rates are lagging for x32 radix
  - Will x32 networks migrate to x64?

# Building Your Data Center
## Scale-Out vs. Scale-Up– A Balancing Act

**Higher Revenue Per Facility**

### Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

Number of Servers (Assuming 24 Servers per rack)

Switch Radix

● 2-Tier Network ● 3-Tier Network ● 4-Tier Network ● 5-Tier Network

Graph concept leveraged from R. Nagarajan, Ilya Lyubomirsky, "Next-Gen Data Center Interconnects: The Race to 800G"
Adjusted to hold servers per rack constant

| Switch BW | SerDes | Radix x32 | | Radix x64 | |
|-----------|--------|-----------|-----|-----------|-----|
| 12.8T | 56G | 400GE | x8 | 200GE | x4 |
| 25.6T | 112G | 800GE | x8 | 400GE | x4 |
| 51.2T | 112G | 1.6TE | x16 | 800GE | x8 |
| 102.4T? | 212G? | 3.2TE | x16 | 1.6TE | x8 |

Wider Radix – Scale Out
More Layers – Scale Up

Improved Power Efficiency
Improved Link Utilization

- x32 and x128 radix are prominent today
  - Ethernet rates are lagging for x32 radix
  - Will x32 networks migrate to x64?

**Radix 64**

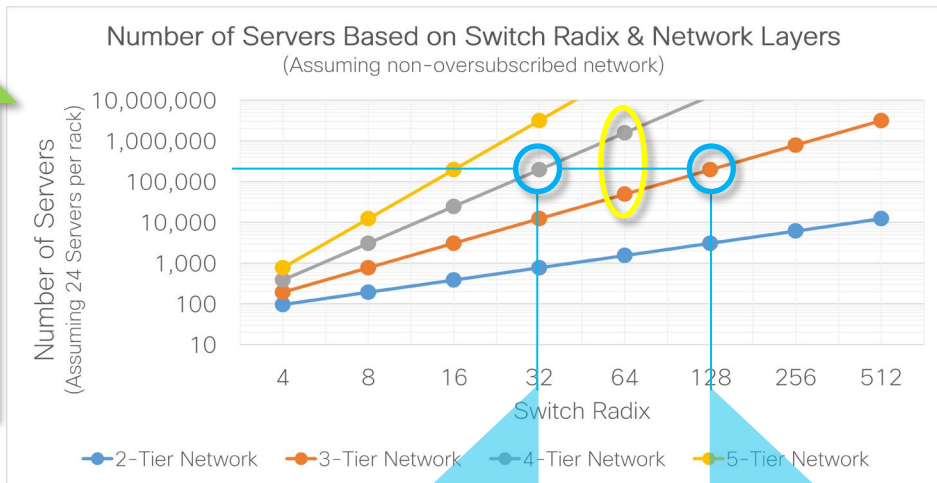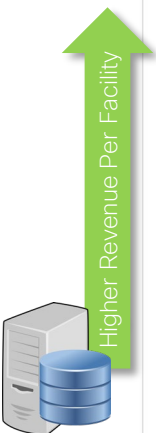- **Potential** need for 800GE with 8x112G Lanes
  - 51.2T
  - 64 x QSFP-DD800 (carrying 1x800GE) – 2RU
- **Potential** need for 1.6TE with 8x224G Lanes
  - 102.4T
  - 64 x QSFP-DD1600 (Carrying 1x1.6TE) – 2RU

cisco

# Building Your Data Center
## Scale-Out vs. Scale-Up– A Balancing Act

### Number of Servers Based on Switch Radix & Network Layers
(Assuming non-oversubscribed network)

**Higher Revenue Per Facility** ↑

Number of Servers (Assuming 24 Servers per rack)

Y-axis: 10 / 100 / 1,000 / 10,000 / 100,000 / 1,000,000 / 10,000,000

X-axis (Switch Radix): 4 / 8 / 16 / 32 / 64 / 128 / 256 / 512

Legend: ●— 2-Tier Network   ●— 3-Tier Network   ●— 4-Tier Network   ●— 5-Tier Network

Graph concept leveraged from R. Nagarajan, Ilya Lyubomirsky, "Next-Gen Data Center Interconnects: The Race to 800G"
Adjusted to hold servers per rack constant

| Switch BW | SerDes | Radix x32 | | Radix x64 | | Radix x128 | |
|---|---|---|---|---|---|---|---|
| 12.8T | 56G | 400GE | x8 | 200GE | x4 | 100GE | x2 |
| 25.6T | 112G | 800GE | x8 | 400GE | x4 | 200GE | x2 |
| 51.2T | 112G | 1.6TE | x16 | 800GE | x8 | 400GE | x4 |
| 102.4T? | 212G? | 3.2TE | x16 | 1.6TE | x8 | 800GE | x4 |

Wider Radix - Scale Out
More Layers – Scale Up

→ Improved Power Efficiency →
← Improved Link Utilization ←

- x32 and x128 radix are prominent today
  - Ethernet rates are lagging for x32 radix
  - Will x32 networks migrate to x64?

**Radix 64**
- **Potential** need for 800GE with 8x112G Lanes
  - 51.2T
  - 64 x QSFP-DD800 (carrying 1x800GE) – 2RU
- **Potential** need for 1.6TE with 8x224G Lanes
  - 102.4T
  - 64 x QSFP-DD1600 (Carrying 1x1.6TE) – 2RU

**Radix 128**
- **Clear** need for 800GE with 4x224G Lanes
  - 102.4T with 128-Radix
  - 128 x QSFP-800 (carrying 1x800GE) – 4RU
    or
  - 64 x QSFP-DD1600 (carrying 2x800GE)–2RU

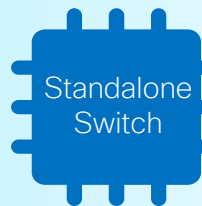# 212G Generation Traditional System Architectures

## Viable with Traditional System Designs

VSR – Optimize for Optics
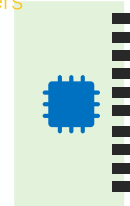112G last major passive copper generation ➜ Active Copper
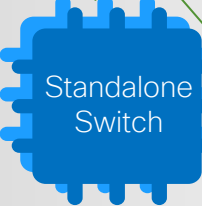
8x212G VSR – No Re-timers

**Fixed**

Standalone Switch

Retimer

2x800GE
1x1.6TE

2x800GE
1x1.6TE

AEC

**Centralized**

212G MR-LR Required

Standalone Switch

Mux

Mux

2x800GE
1x1.6TE

2x800GE
1x1.6TE

Optional Redundancy

212G LR Required

**Distributed**

FE Switch

Retimer

LC Switch

Retimer

2x800GE
1x1.6TE

2x800GE
1x1.6TE

# 212G Generation CPO System Architectures

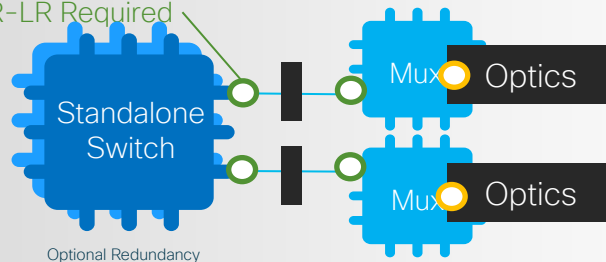Power Optimized ; Introduced first on Client-Side Optics



**Fixed**

212G/2x112G

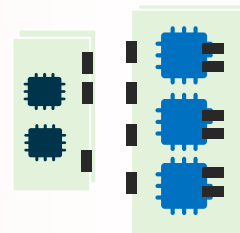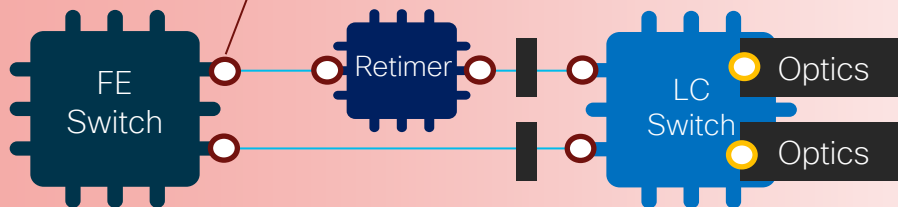Standalone Switch

Optics

Optics

**Centralized**

212G MR-LR Required

Standalone Switch

Mux — Optics

Mux — Optics

Optional Redundancy

**Distributed**

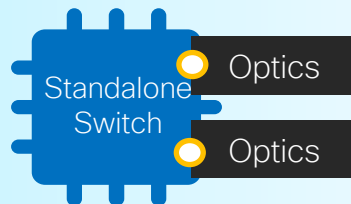212G LR Required

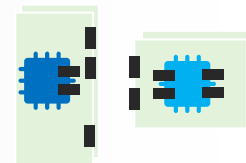FE Switch — Retimer — LC Switch

Optics

Optics

cisco

# Future CPO Architectures
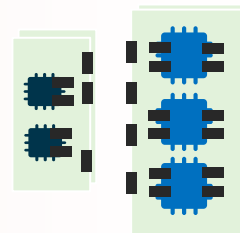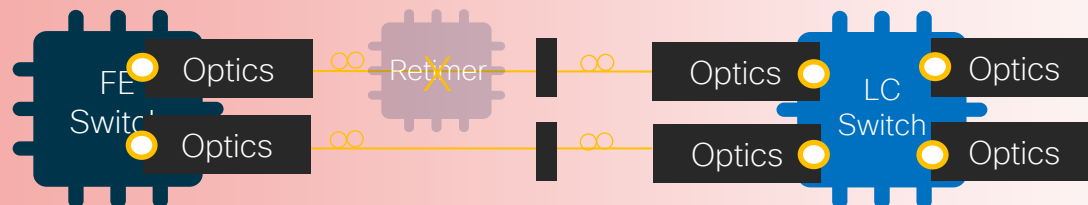## Eventually Optics replace high speed data interconnect



**Fixed**

Standalone Switch — Optics / Optics

**Centralized**

Standalone Switch — Optics / Optics — Optics / Optics — Mux / Mux

Optional Redundancy

**Distributed**

FE Switch — Optics / Optics — Retimer — Optics / Optics — LC Switch — Optics / Optics

# Call to Action
## Power Driven Architecture

✔ ## 3 Main System Architectures
Fixed, Centralized, Modular

✔ ## BW Doubling every 2 Years
Not Slowing Down, Power Too High

✔ ## Co-package Optics are Coming
51.2T Generation



Limits what our Planet can
Sustain
Moral Imperative

Business Imperative

Technology Imperative

Build
Limits what we can

Deploy
Limits what We can

Perfect Catalyst to Innovate

# Next Steps

**Increasing Priority**

**1** Define 212G Electrical

- XSR, VSR as first priority to optimize power efficiency
  - Define VSR standard to ensure retimer-less designs
- Define MR, LR as second priority
- Focus on 212G* instead of 224G to optimize for Ethernet rates

**2** Define 800GE MAC

- Over 212G to enable 102.4T with radix 128 (128x800GE)
- Over 112G to enable 51.2T with radix 64 (64x800GE)

**3** Define 1.6TE MAC

- Only if there is a cost effective PMD solution
- Over 212G to enable 102.4T with radix 64 (64x1.6TE)

* – Final rate depends on future work

Thank You!