# What's beyond 400G ?

## Mark Filer

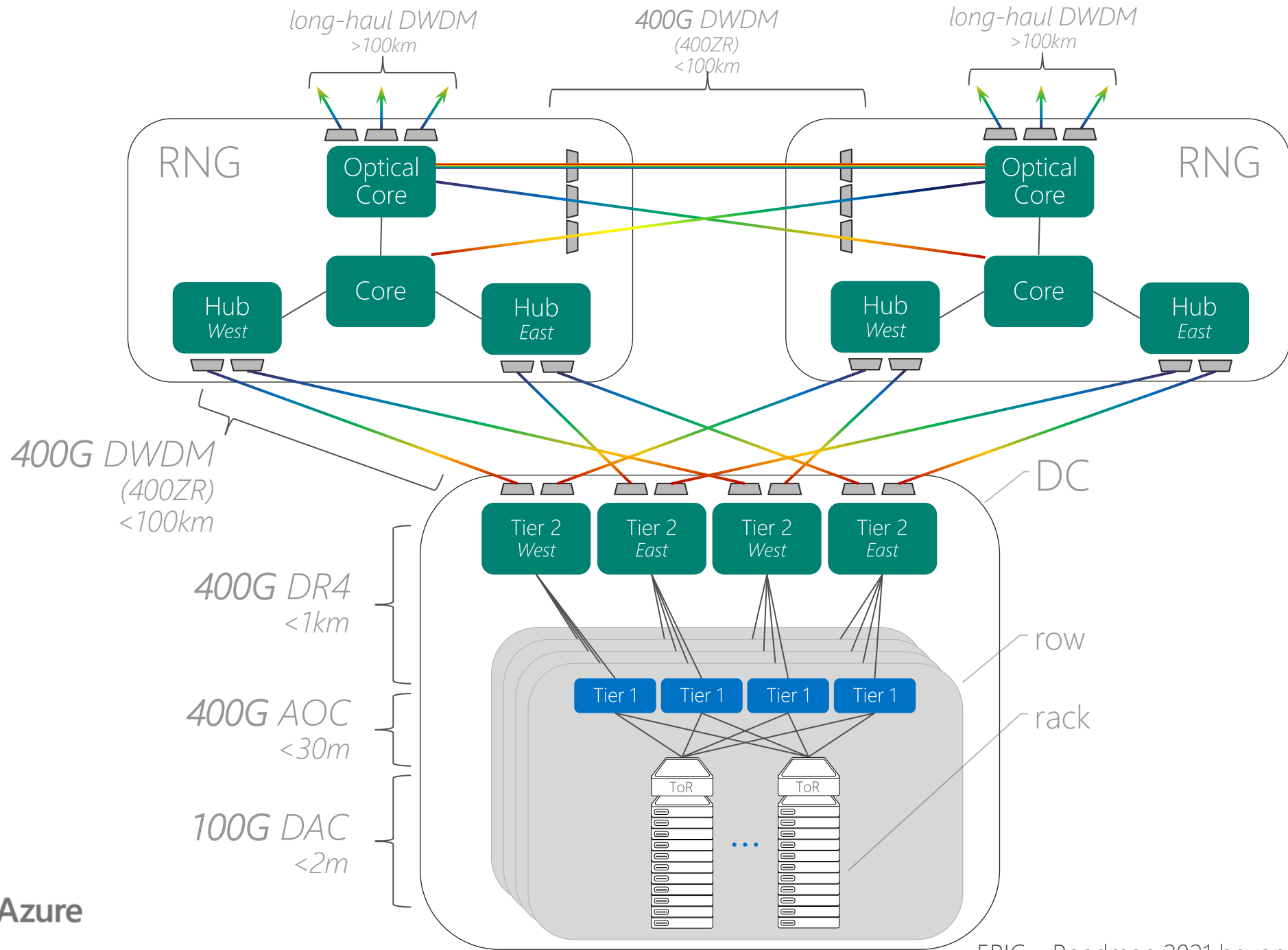Principal Engineer, Azure Hardware Architecture (AHA)

Microsoft Azure

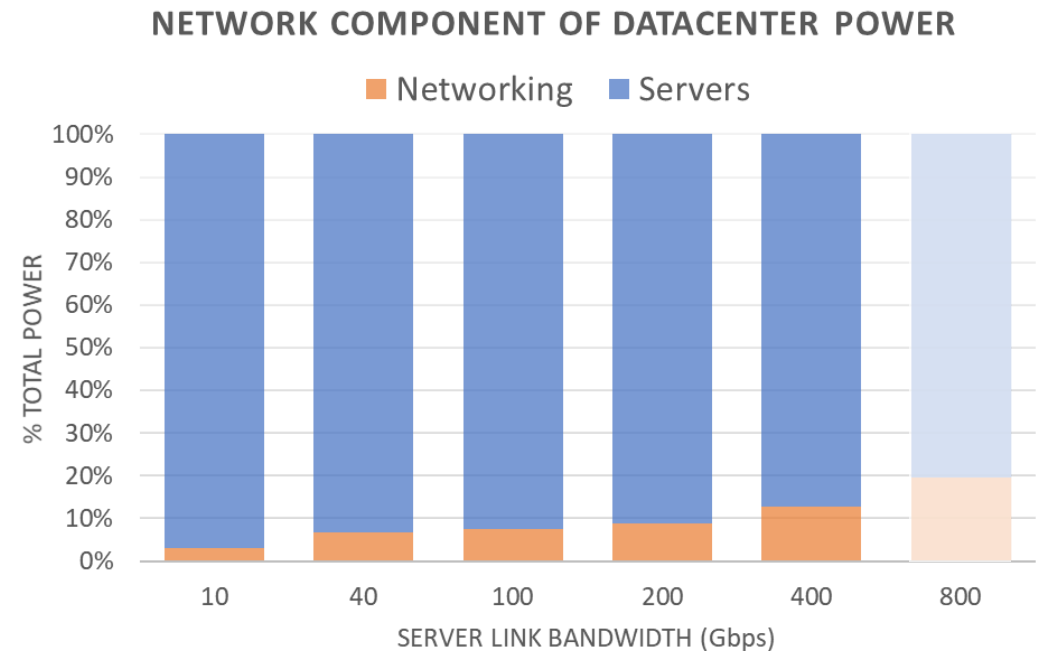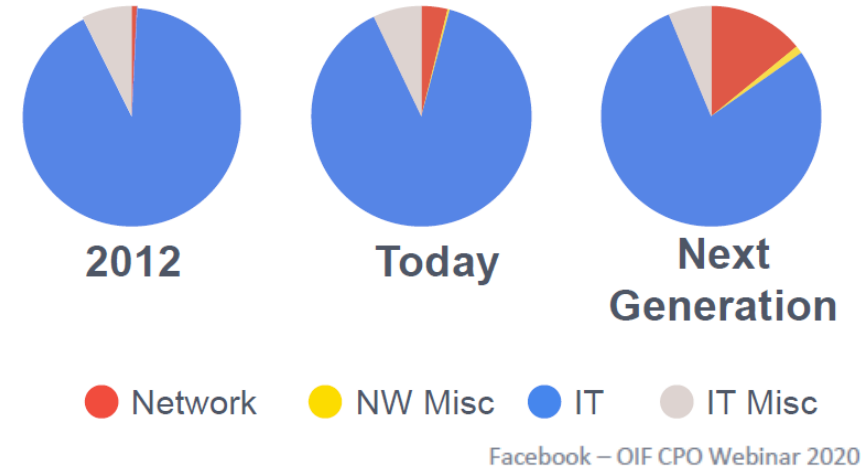# Disclaimer – Statements of Future State

*Material does not necessarily represent opinions of Microsoft and certainly cannot be construed as any form of commitment by Microsoft towards pursuing concepts described herein*

400G
(2021-)

long-haul DWDM
>100km

400G DWDM
(400ZR)
<100km

long-haul DWDM
>100km

RNG

RNG

Optical Core

Optical Core

Core

Core

Hub West

Hub East

Hub West

Hub East

400G DWDM
(400ZR)
<100km

DC

Tier 2 West

Tier 2 East

Tier 2 West

Tier 2 East

400G DR4
<1km

row

Tier 1

Tier 1

Tier 1

Tier 1

rack

400G AOC
<30m

100G DAC
<2m

ToR

ToR

Microsoft Azure

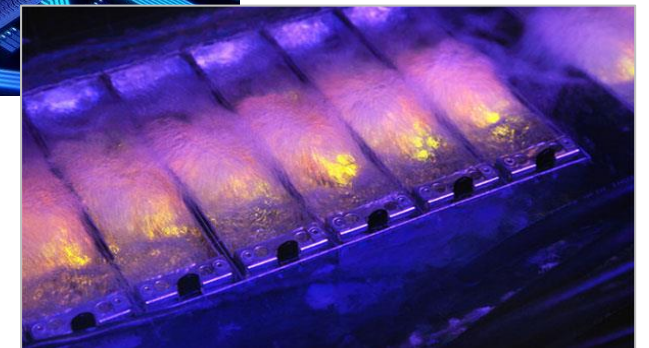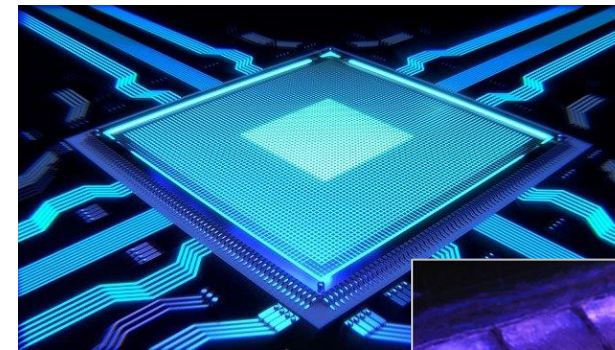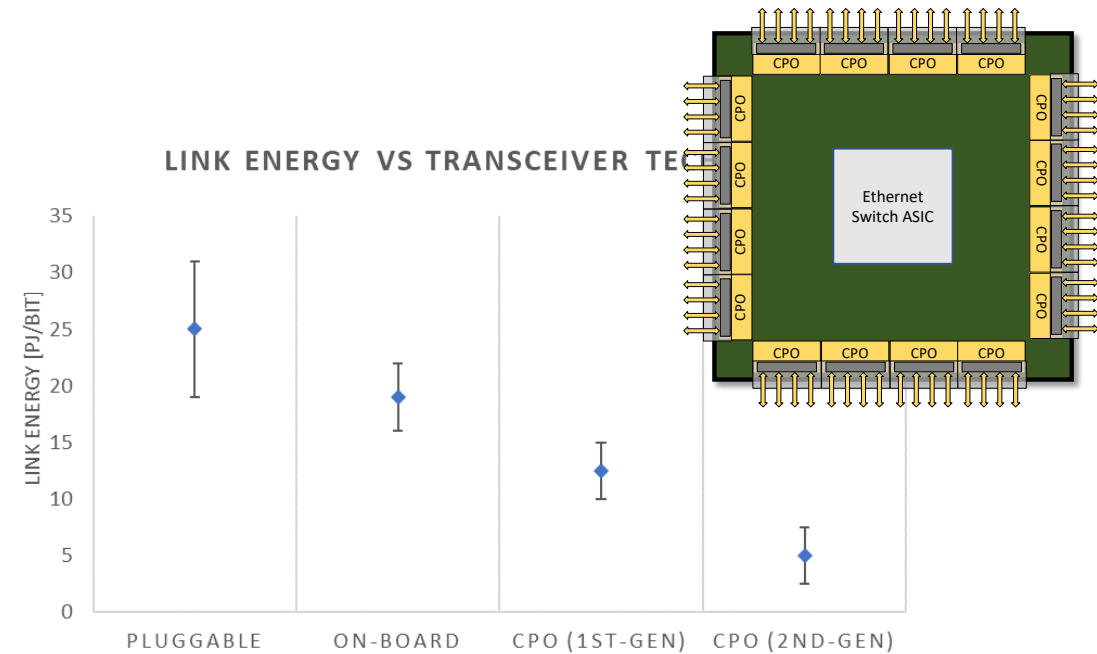EPIC – Roadmap 2021 beyond 400G | 21 Apr 2021

# Power limits future DC scaling

- Equipment power consumption at 400G is already problematic!
  - Switches projected @ 3x power of 100G
  - Optics projected @ 3-4x power of 100G
- Challenges power envelopes of facilities
- Uses power that could be generating revenue (lost server capacity)
- Costs $$$ and not green
- Trajectory makes transition to >400G appear all but impossible



2012    Today    Next Generation

● Network    ● NW Misc    ● IT    ● IT Misc

Facebook – OIF CPO Webinar 2020



NETWORK COMPONENT OF DATACENTER POWER

■ Networking    ■ Servers

% TOTAL POWER

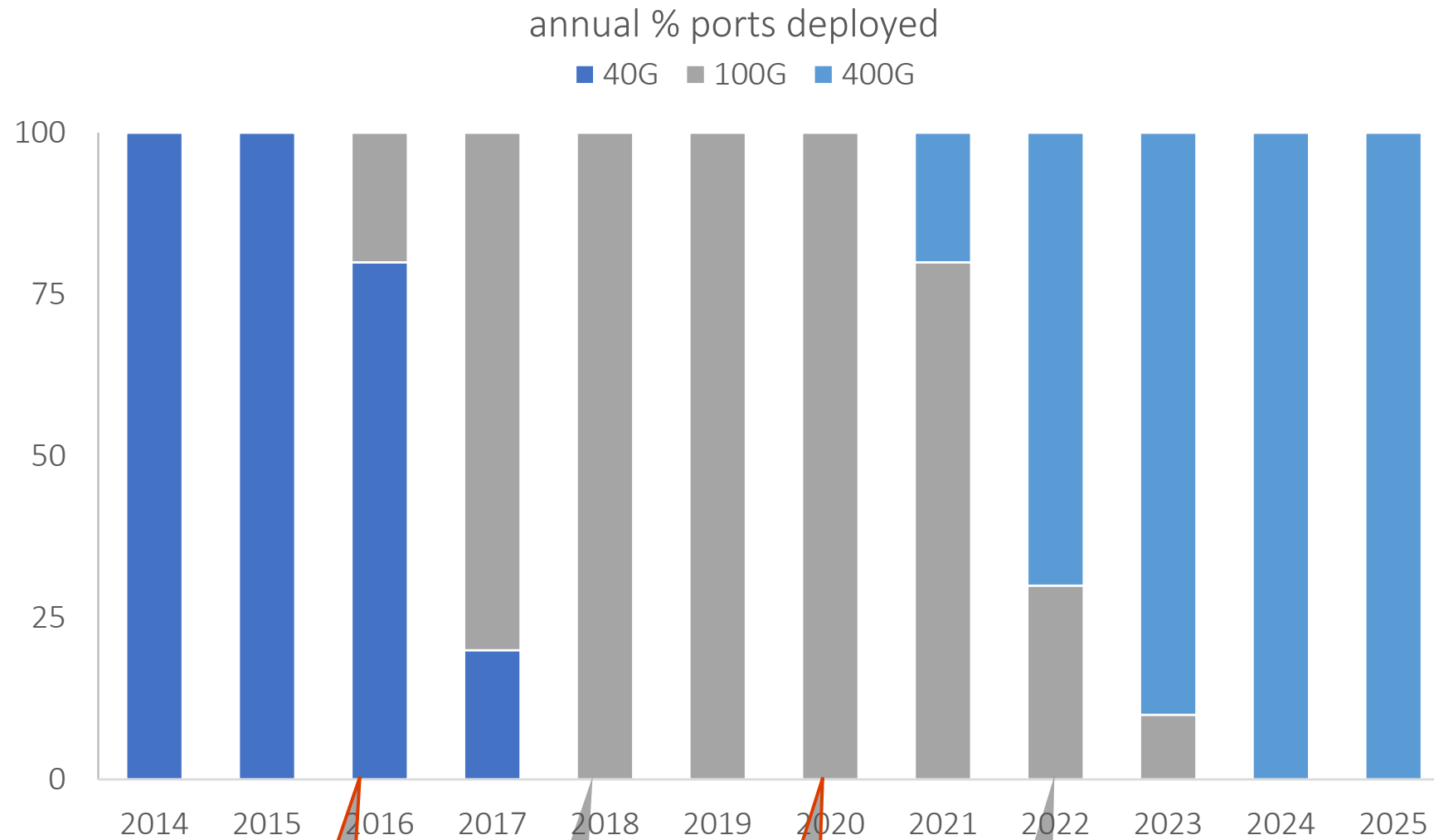SERVER LINK BANDWIDTH (Gbps)

Microsoft Azure

# Possible Solutions

- Photonics
  - Co-Packaged Optics (CPO)
  - Novel optical approaches

- Network architecture + HW changes
  - Collapsed tiers with multi-homed NICs (fanning out horizontally)
  - Simplified forwarding requirements → cooler ASICs
  - Additional integration, e.g. encryption on switch ASIC
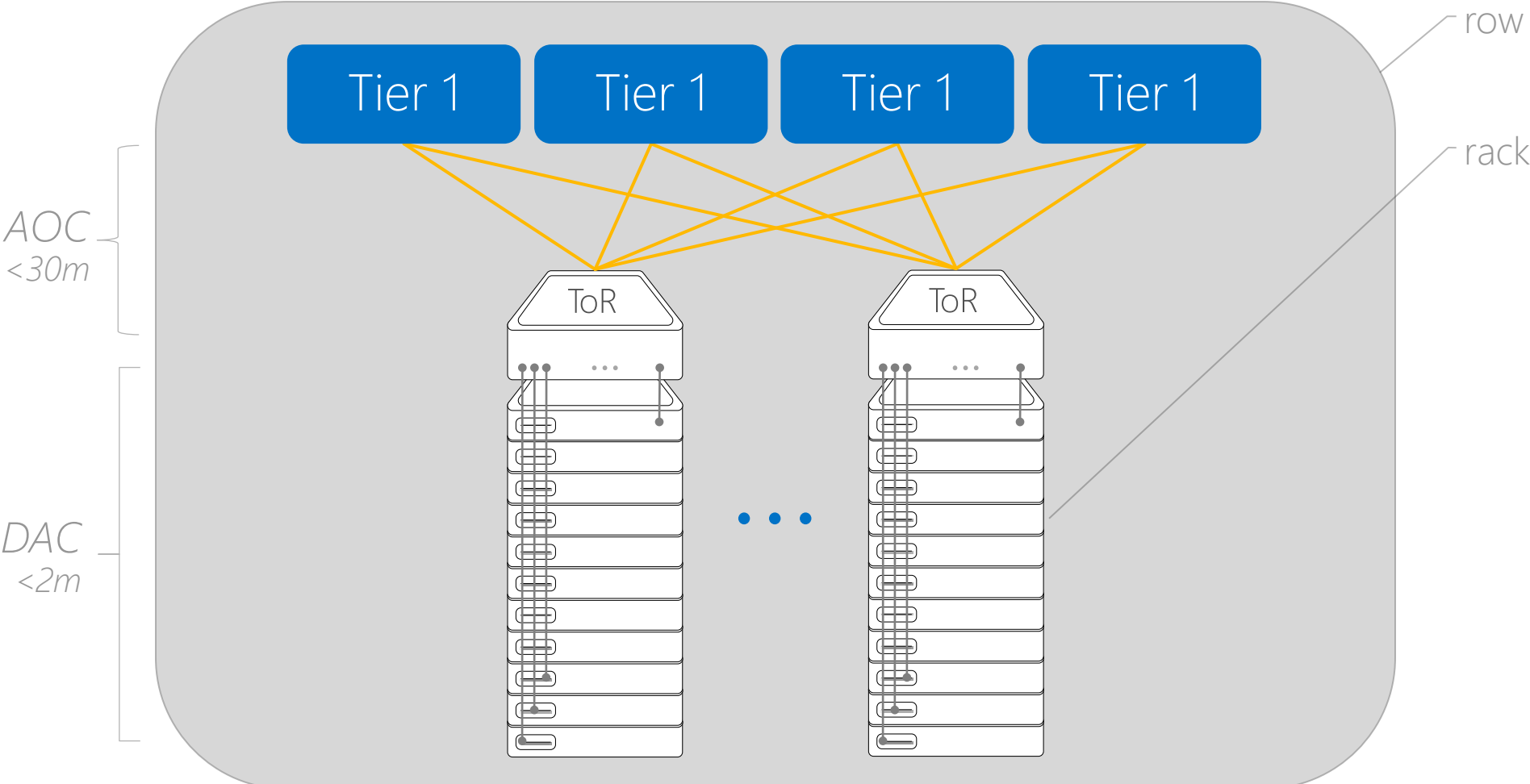  - Liquid cooling

Takeaway: we can't just keep scaling link bandwidths... "next gen" systems will require all of the above
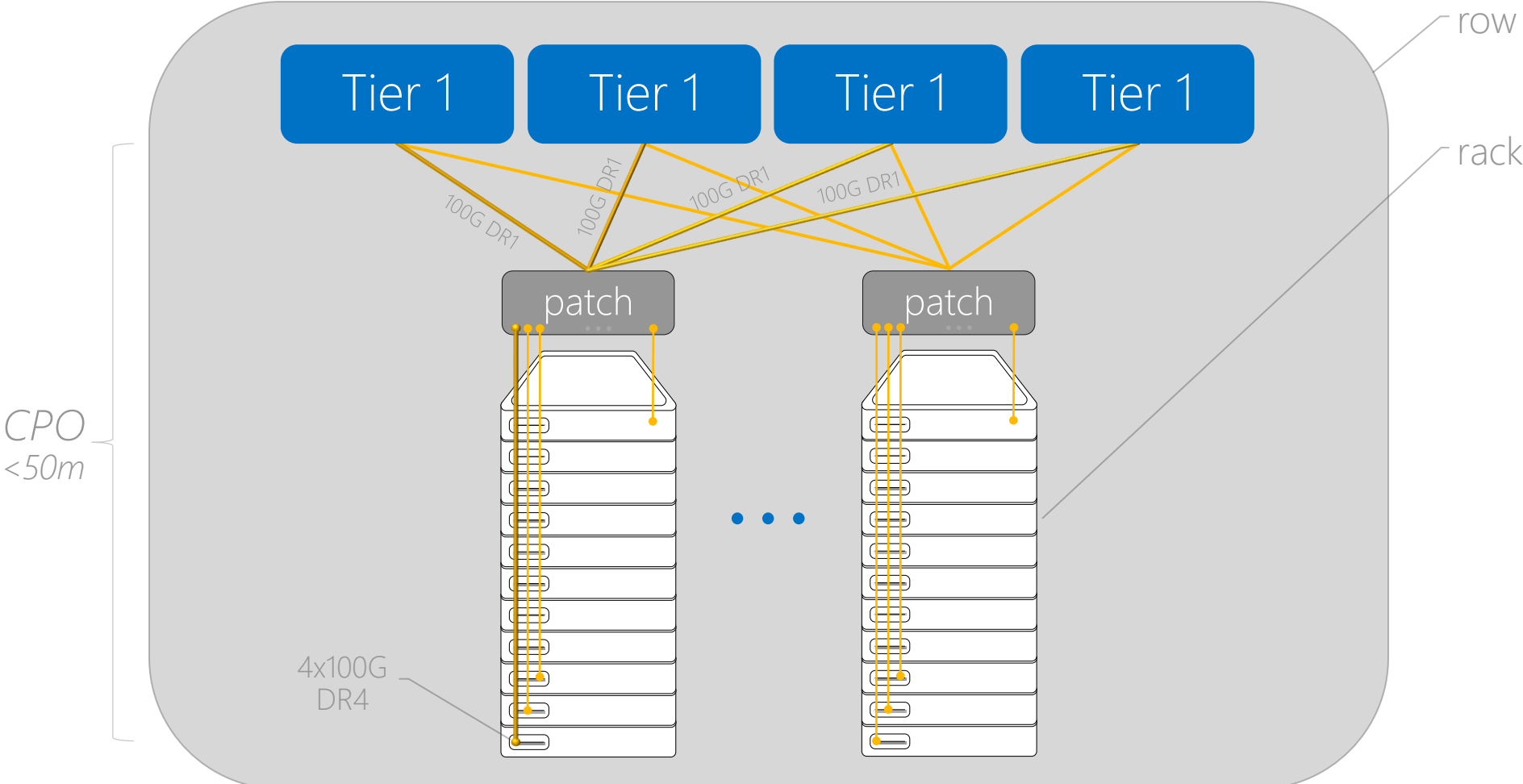


LINK ENERGY VS TRANSCEIVER TE

# Microsoft DC ecosystem technology life cycles

annual % ports deployed

■ 40G   ■ 100G   ■ 400G

100G GA   200G GA   400G GA   800G GA

Microsoft Azure

EPIC – Roadmap 2021 beyond 400G | 21 Apr 2021

# Server-ToR-Tier1 topology

# ToR bypass – multi-homed NIC



row

rack

Tier 1   Tier 1   Tier 1   Tier 1

100G DR1   100G DR1   100G DR1   100G DR1

patch   patch

CPO
<50m

4x100G
DR4

Microsoft Azure

# ToR bypass efficiencies (100G lane speeds)



| | Tier1-ToR-server | ToR-bypass |
|---|---|---|
| failure domain | ToR is SPOF for rack | no SPOF – multi-homed NIC |
| switch ASIC count | 4X-8X | 1X |
| switch space + power | baseline | reduced space and ~**1/3 power** |
| switch radix | can't leverage higher radix chips (stranded capacity at ToR) | leverages full switch radix |
| oversubscription | 3:1 typical | fully non-blocking in row |
| reach limits | DAC < 3m; AOC < 30m | 1m-2km+ |

Microsoft Azure

# Summary

- Power is the main limiter for "beyond 400G" data centers

- We can't continue to simply scale link bandwidths while building networks exactly as we do today

- Historical ecosystem life cycles would indicate we won't be ready for "800G" when the industry is (32x100G CPO will suit our needs better)

- 100G electrical lanes will be a foundational building block for power-efficient data center designs for the foreseeable future

- Future data center networks will require a combination of photonic innovation (e.g., CPO), optimized network architectures, and advanced hardware implementations

Microsoft Azure

Thank you.

mark.filer@microsoft.com