

# Transitioning to 400G SerDes: Key Drivers and System Design Implications for Future AI Workloads

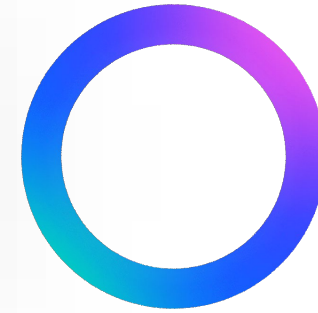
Halil Cirit, AI Architect & META

This presentation has been developed within the Ethernet Alliance, and is intended to educate and promote the exchange of information. Opinions expressed during this presentation are the views of the presenters, and should not be considered the views or positions of the Ethernet Alliance

# Keynote Speaker

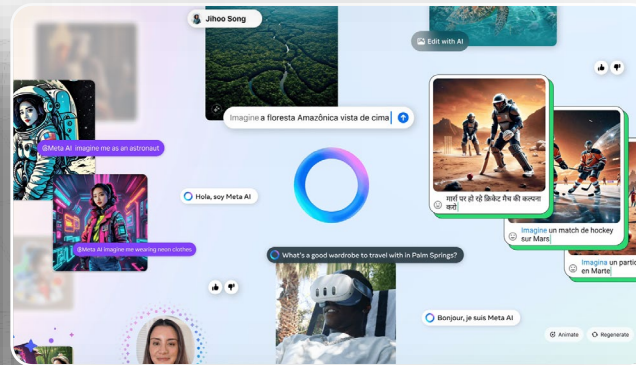
Seasoned SerDes and System expert with over 19 years of experience in the semiconductor industry. Since November 2018, has been with Meta (formerly Facebook), where he is responsible for the SerDes and Systems Interconnect, focusing on next-generation high-speed interfaces. Prior to his current role, Halil served as a Senior Manager at Inphi Corporation where he led the development of 56/112 Gb/s PAM4 DSP chips. Prior to that, he worked on 25 Gb/s Ethernet receivers at Broadcom Inc. and high-speed I/O design at NVIDIA.





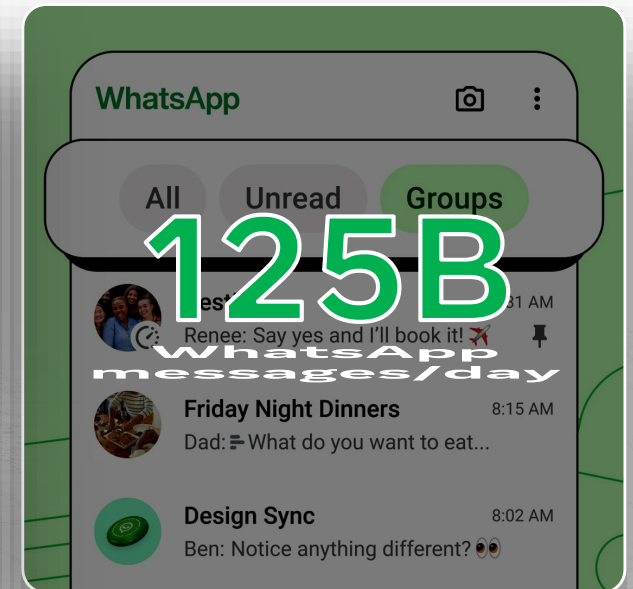
3.4B

daily active users



200B+

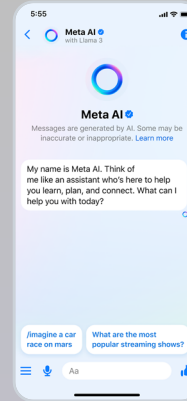
content recommendations/day



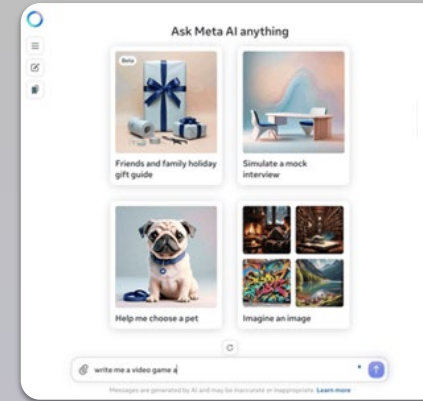


# Built with Llama

Meta AI, launched 1.5 years ago, is available across our apps and on the web. It is more creative, smarter — and used by 1B people every month.



CHAT



WEB

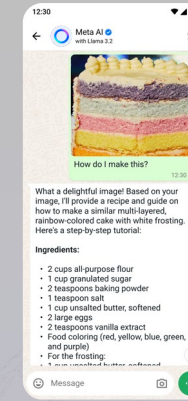
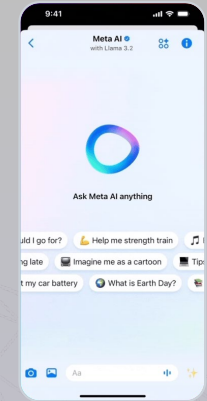


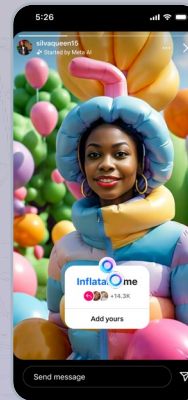
IMAGE UNDERSTANDING



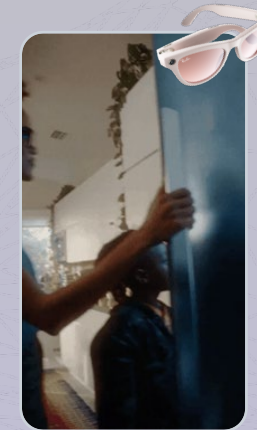
VOICE



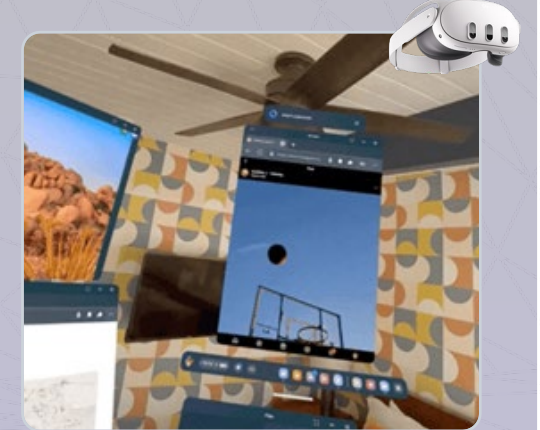
IMAGINE EDIT



IMAGINE YOURSELF

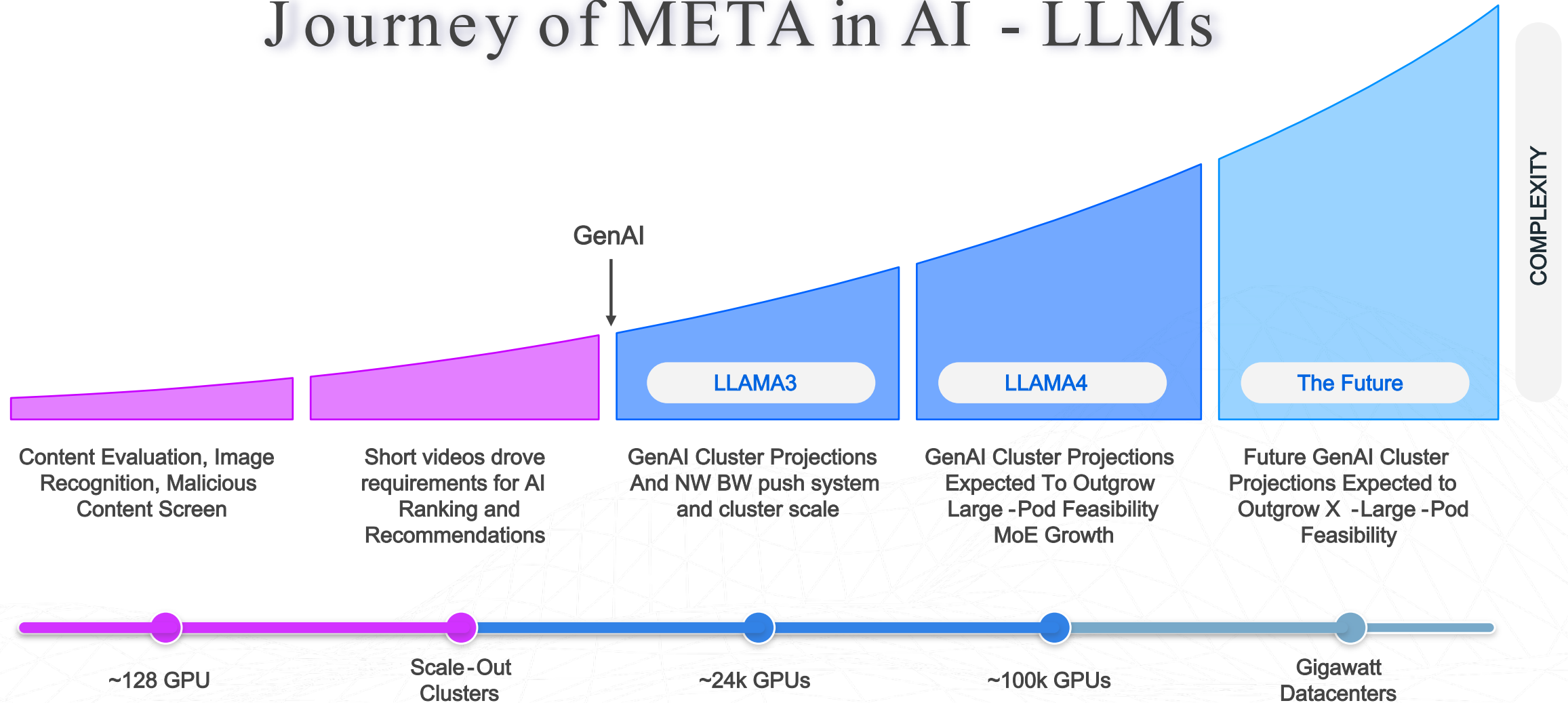


RAY-BAN META

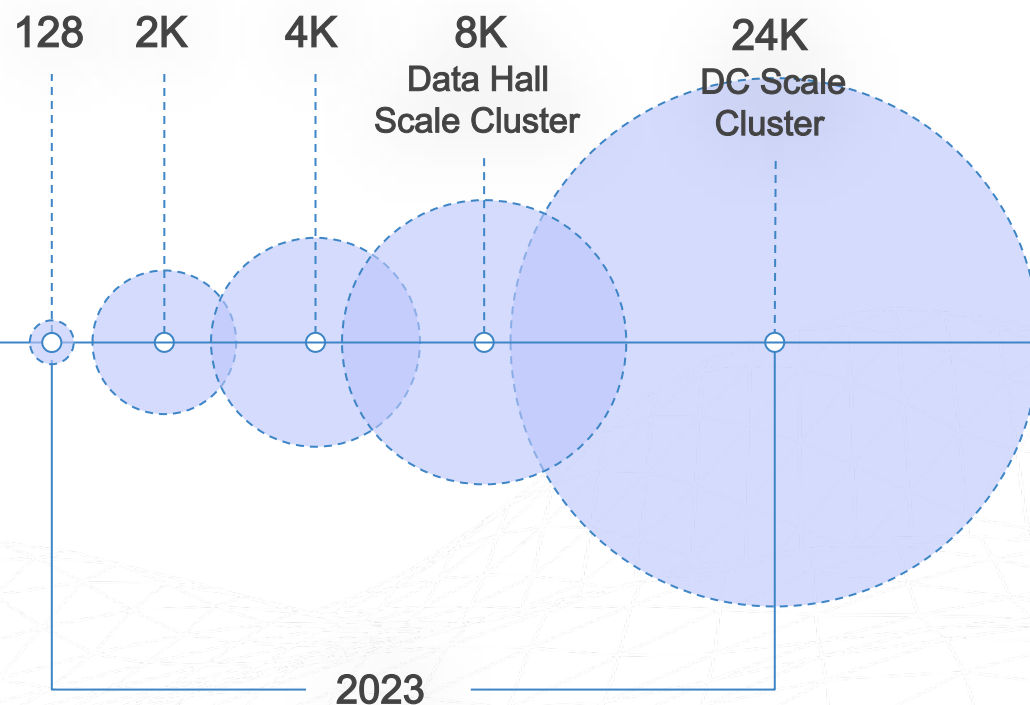


META QUEST

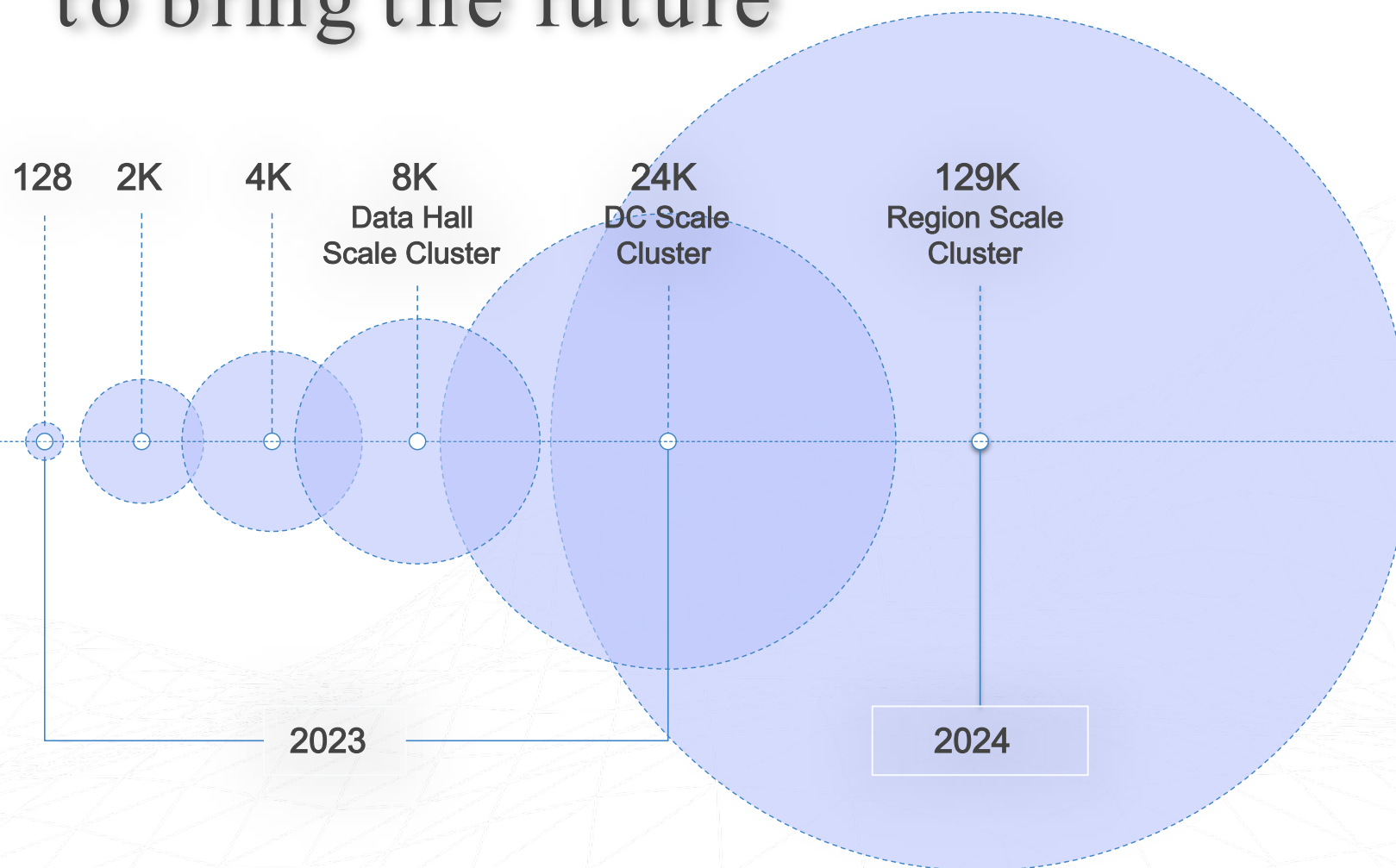
# Journey of META in AI - LLMs



# We are building large clusters to bring the future

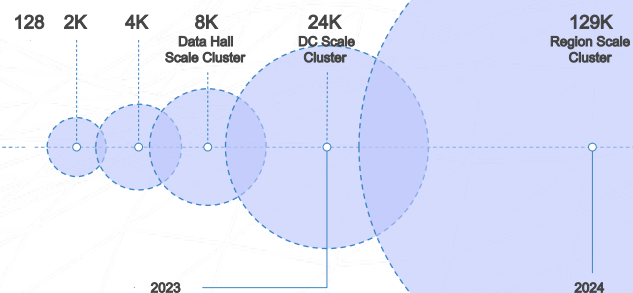


# We are building large clusters to bring the future





# We are building large clusters to bring the future



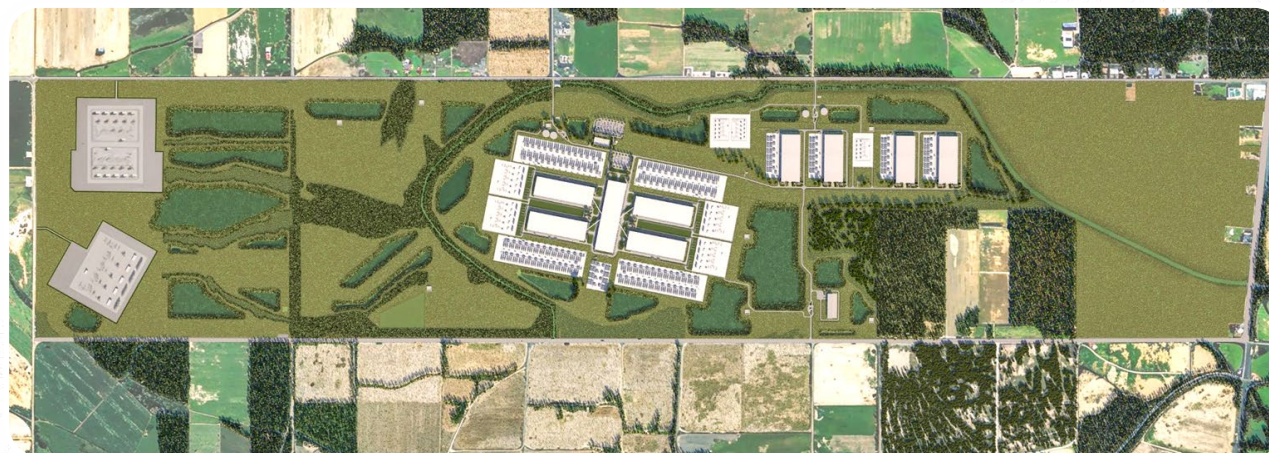
Multi - Million  
Multi - Region Scale  
Cluster

2025 and  
Beyond



## Talking about the massiveness...

- Currently Developing - Prometheus : 1GW+ cluster by 2026
- Our Big Goal - Hyperion : 5GW over the next few years





Manhattan shapes our  
dreams in the skyline;  
Hyperion powers them in  
the network. In the end,  
it's all about our dreams.



Mark Zuckerberg

July 14 at 8:01 AM · 🌐

For our superintelligence effort, I'm focused on building the most elite and talent-dense team in the industry. We're also going to invest hundreds of billions of dollars into compute to build superintelligence. We have the capital from our business to do this.

SemiAnalysis just reported that Meta is on track to be the first lab to bring a 1GW+ supercluster online. 🙌

We're actually building several multi-GW clusters. We're calling the first one Prometheus and it's coming online in '26. We're also building Hyperion, which will be able to scale up to 5GW over several years. We're building multiple more titan clusters as well. Just one of these covers a significant part of the footprint of Manhattan.

Meta Superintelligence Labs will have industry-leading levels of compute and by far the greatest compute per researcher. I'm looking forward to working with the top researchers to advance the frontier!



👍❤️ 211K

73.4K comments 18K shares



Like



Comment

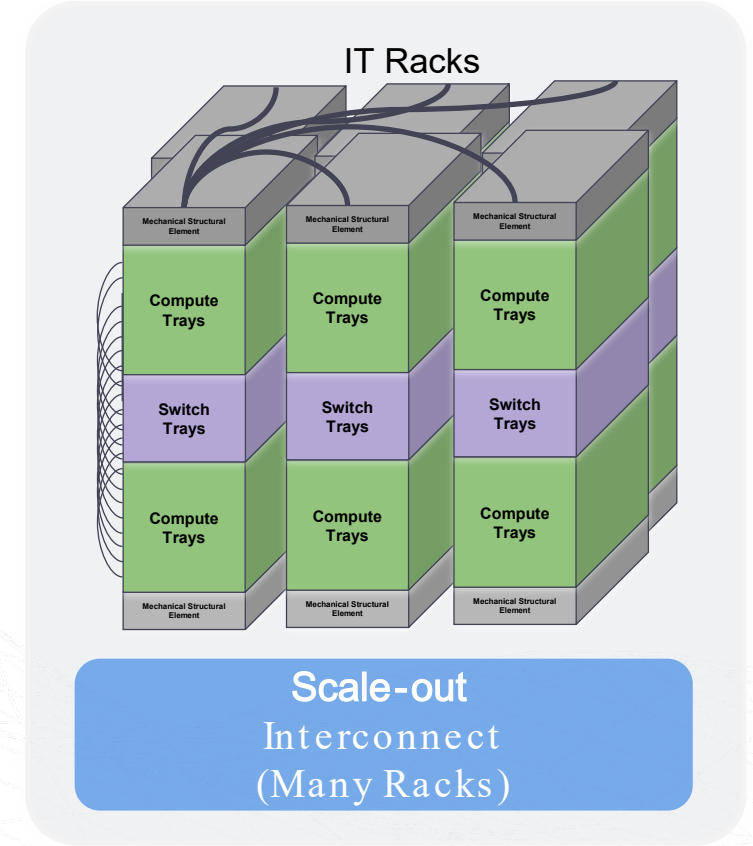
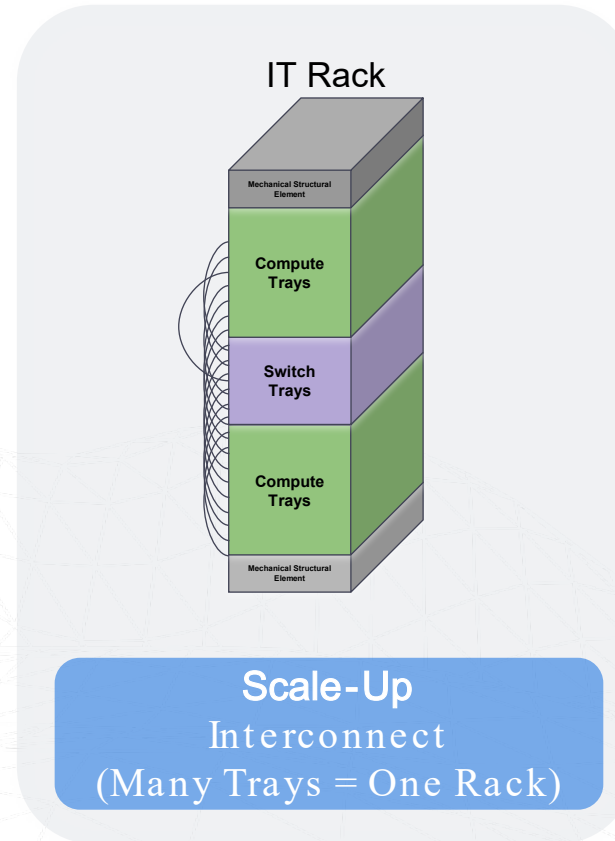
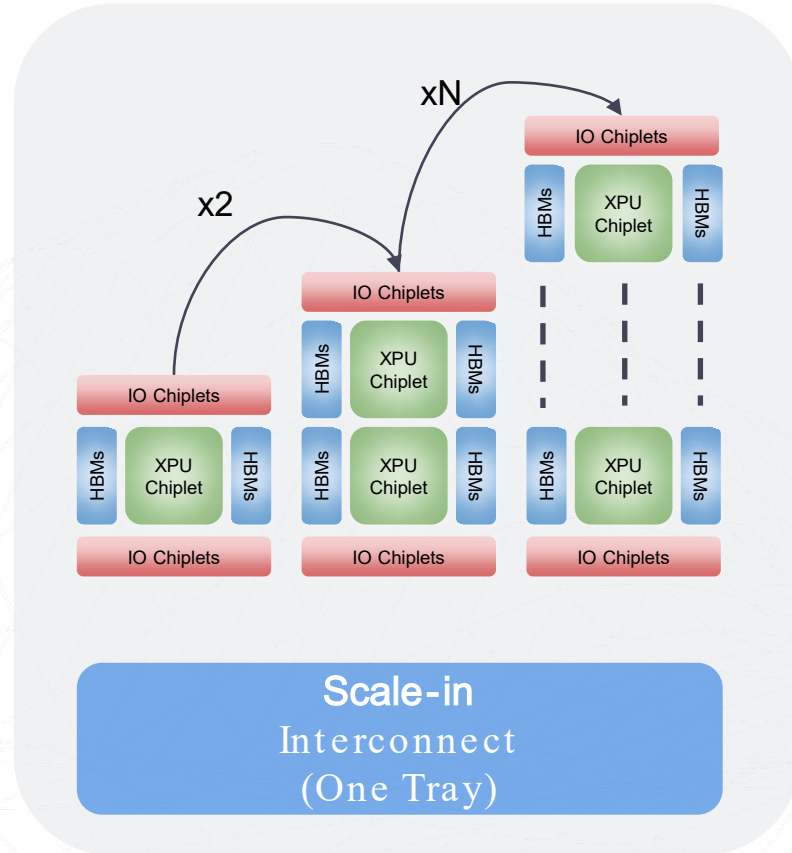


Send



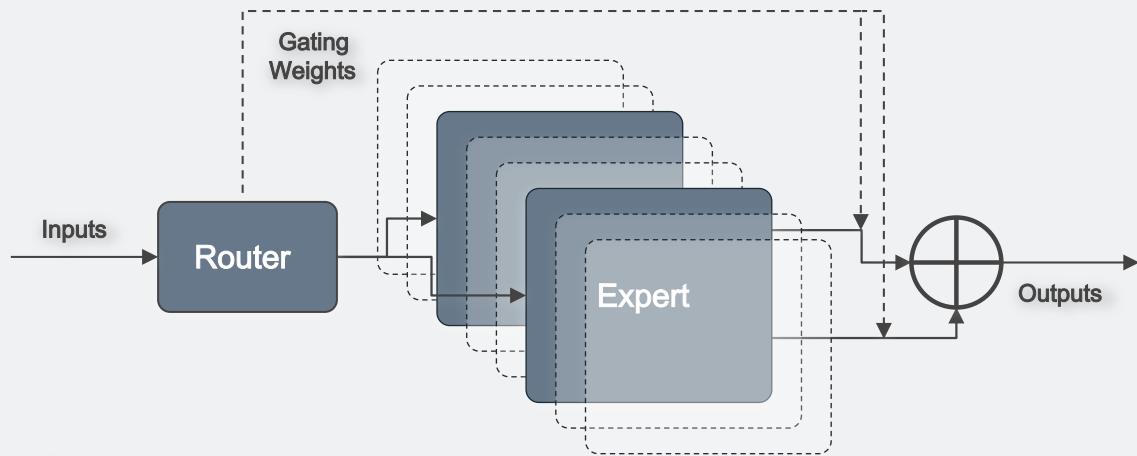
Share

# Scaling to Reach Connected Networks

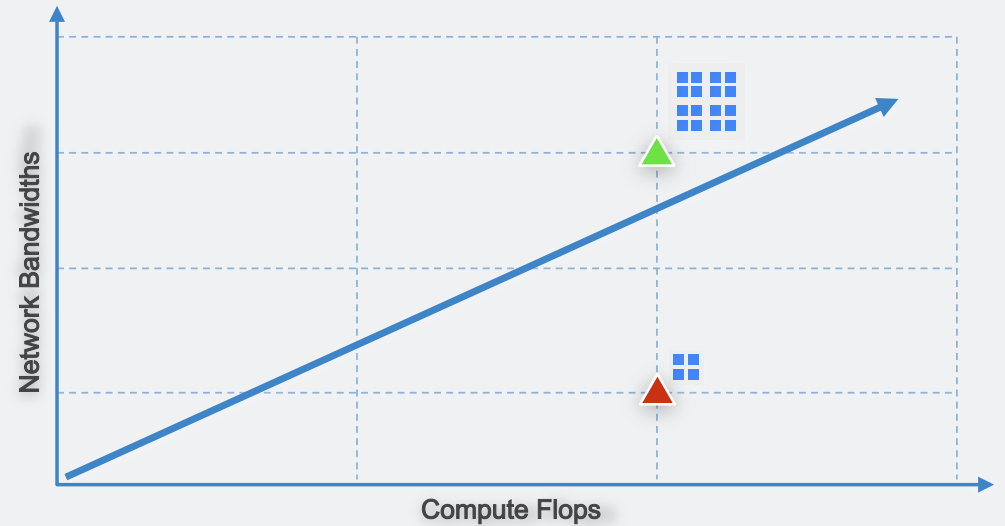




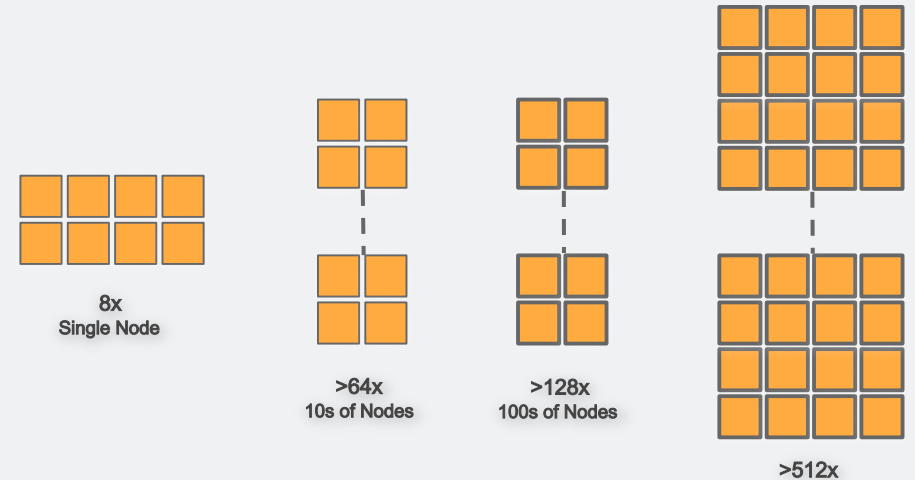
# Scale Up domains are important for next generation data centers



Scaling number of experts in LLMs



Balanced designs prevent stranded resources



Silicon itself is growing in size to have more computing power

Higher number of  
accelerators is  
required for  
performance

Basic Idea: More  
nodes in racks?

ORv3 HPR  
MP



$\leq 72$  Accelerators

Cabled Backplane  
48 VDC /  $\pm 400$  VDC  
Air / Liquid Cooled  
IT/Power Rack Single -Wide

ORW  
2026 -Q3



$\leq 144$  Accelerators

Cabled Backplane  
48 VDC /  $\pm 400$  VDC  
Air / Liquid Cooled  
IT Rack Double -Wide

What is next?  
2027 -Q3



$\geq 256$  Accelerators

Greater than 900kW  
 $\pm 400$  VDC  
Primarily Liquid Cooled  
IT Rack Size TBD

# Bigger racks brings more challenges



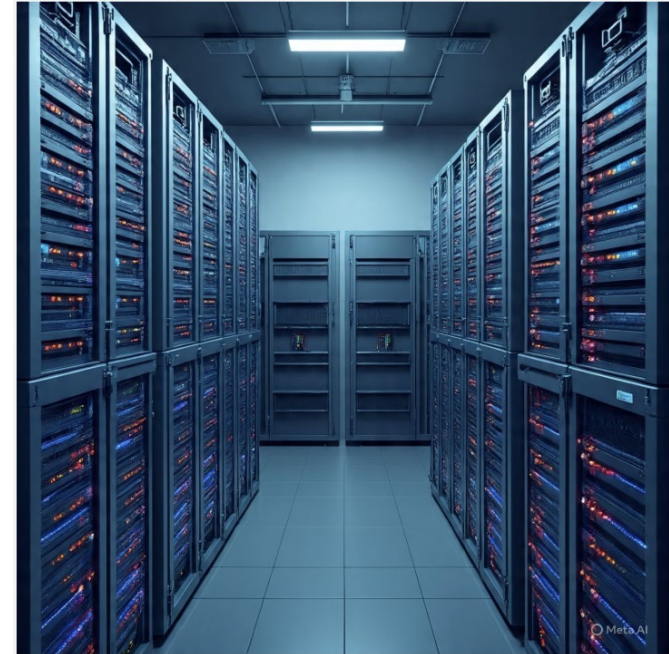
Hard to design  
and manufacture



Expensive  
to transport



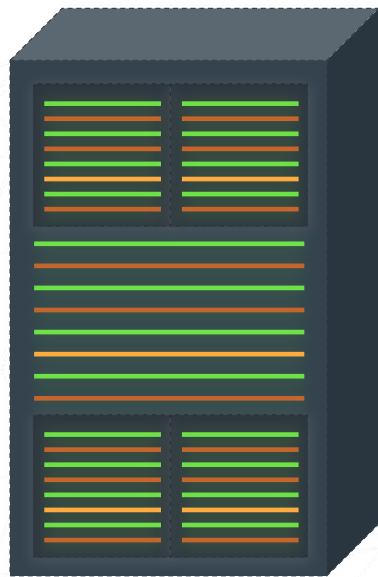
Operationally  
challenging



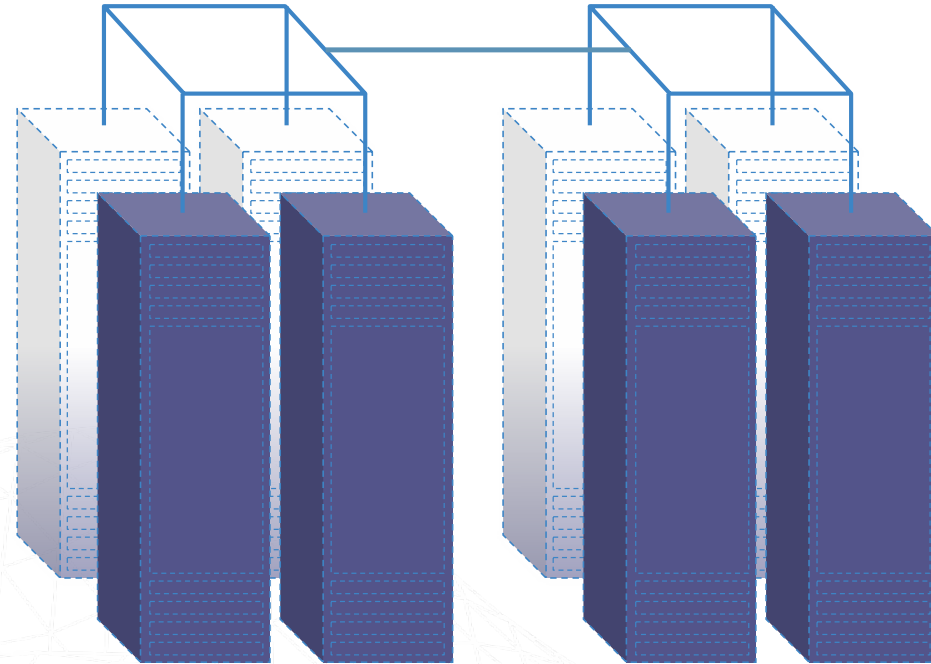
Requiring DC  
infra upgrades

## Next Idea: Improving rack scale systems

# The future is disaggregation



BFR  
(Backplane View)



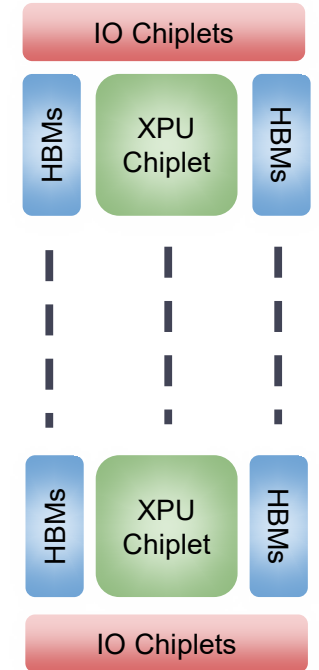
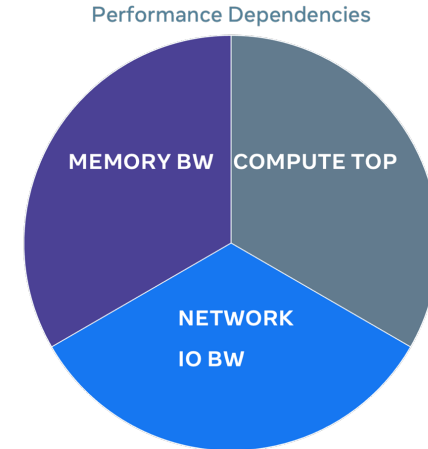
ORV3 ?  
Optically interconnected  
lower density racks

## Next Idea: Improving rack scale systems

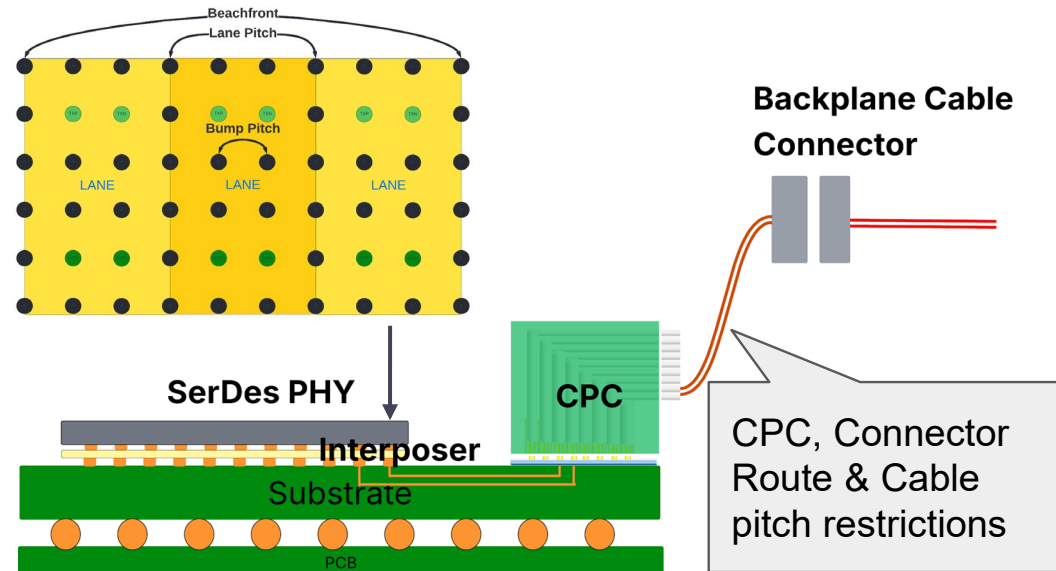


# Beachfront Bottleneck of Scale in Systems

- AI system performance is dependent on 3 major pillars which need to improve all together
- Memory and Compute takes full benefit (Linear scaling) of advanced nodes while SerDes design can't
- Beachfront is the major limitation on scaling the Network IO BW



xN Computing  
xN Memory BW  
x1 Network IO BW

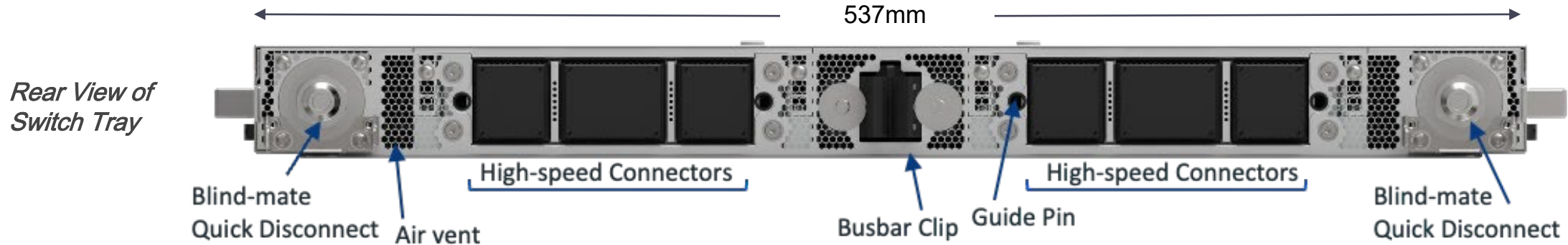


CPC = Co-Packaged Copper

# Bandwidth Bottleneck in Rack-Scale Systems

## Busy Backplane Beachfront

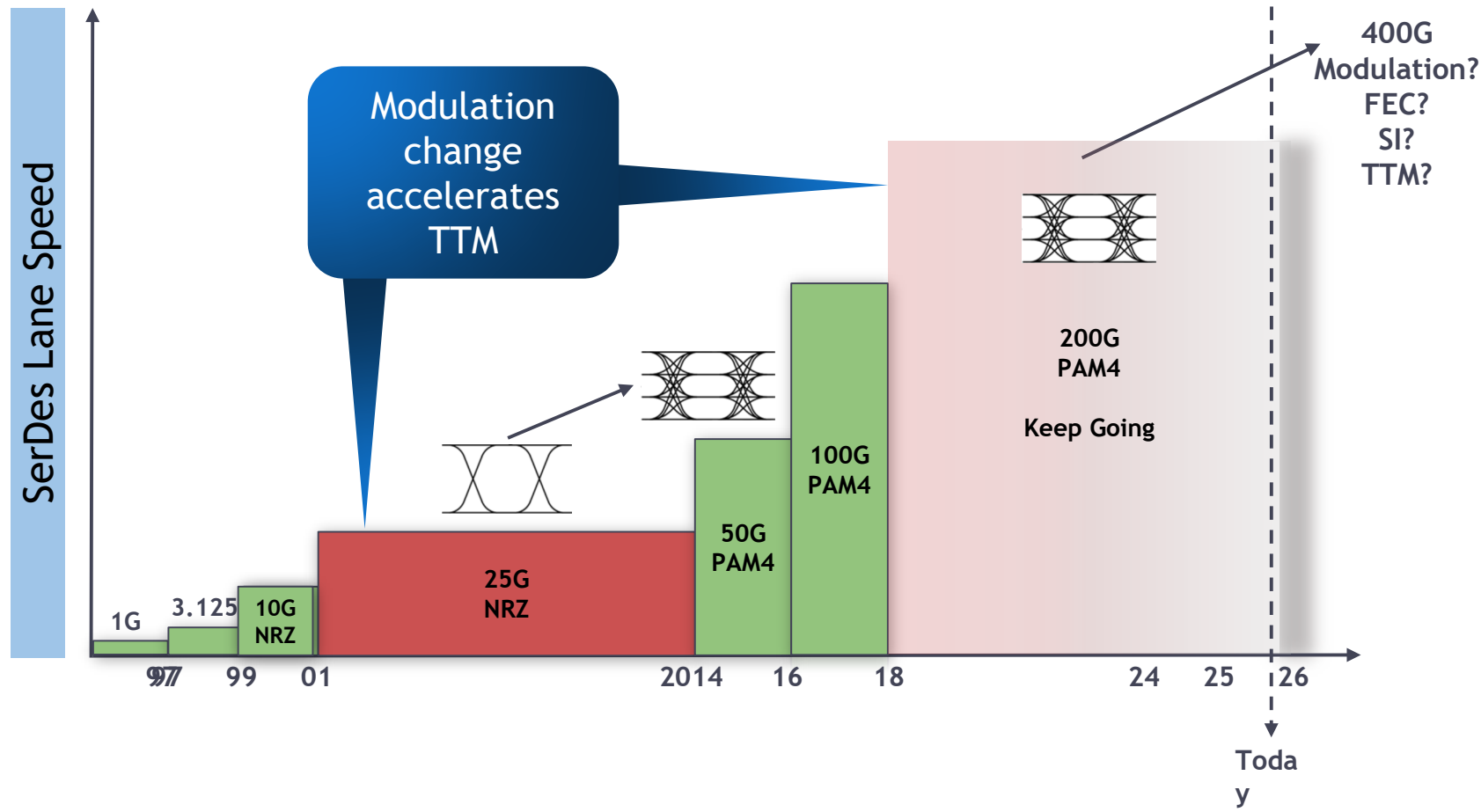
- Crowded backplane beachfront space usage



- 1,024~1,152 differential pairs per OU (ORv3 / 600mm) cross-sectional area for cable backplanes
  - Enables 102.4-115.2T of bandwidth per OU in an ORv3 / 600mm tray with 200G SerDes
  - Typically a bottleneck for NW switch trays
- Space required for different needs including connector alignment

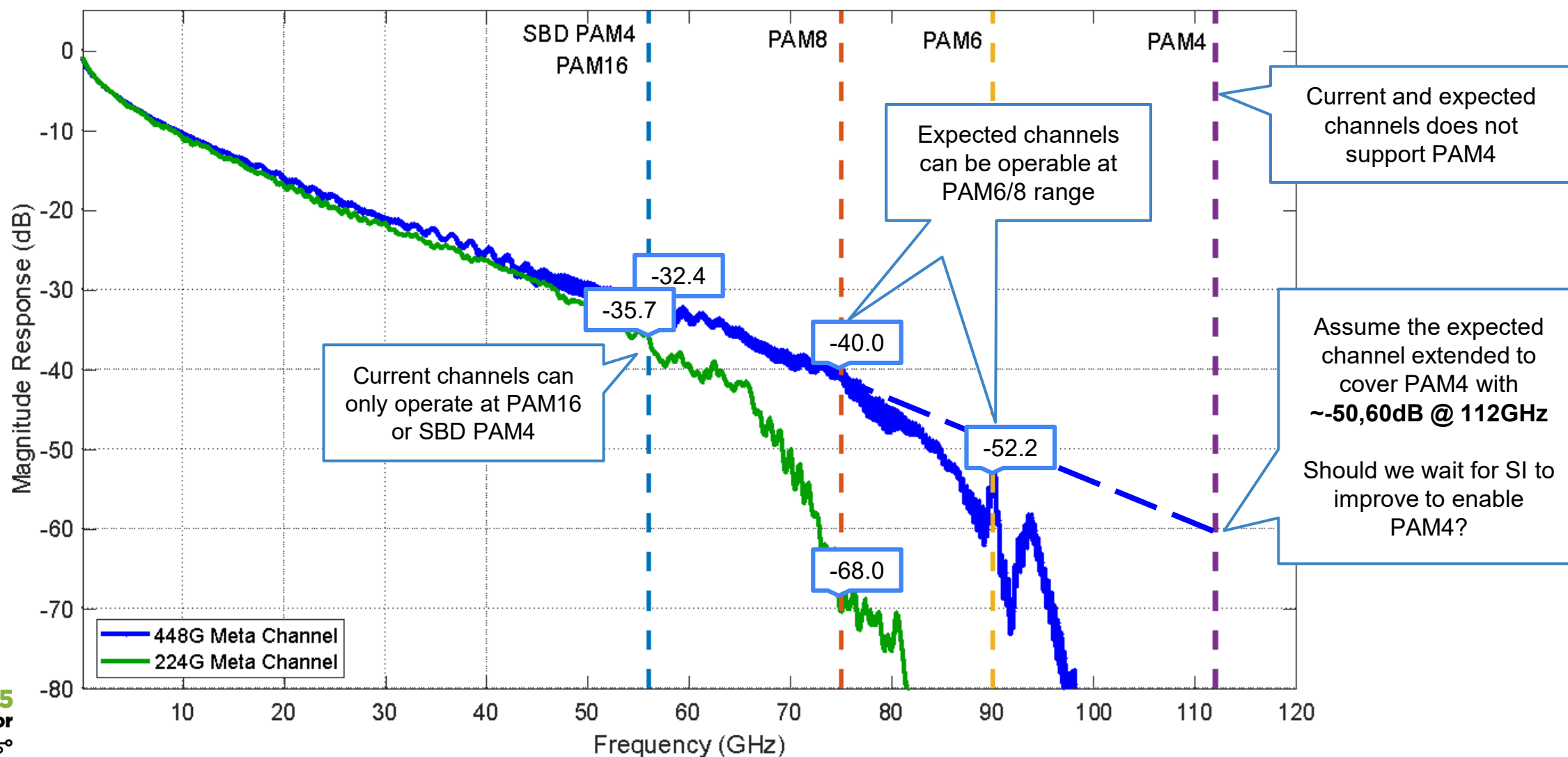
Higher data rate per lane required for aggregate bandwidth scaling

# Overcoming Beachfront Bottleneck



- Signaling rate doubled every two years, except 25G and 200G where modulation and FEC are at their limits
- NRZ→PAM4 gave two more generations
- 400G now needs a new breakthrough
- Expected stronger FEC & higher order modulation
- Hyperscalers involvement is needed to accelerate the 400G TTM

# BW Problem of Current Connectors & Cables





# Accelerating Progress with Collaboration: Vendors Side

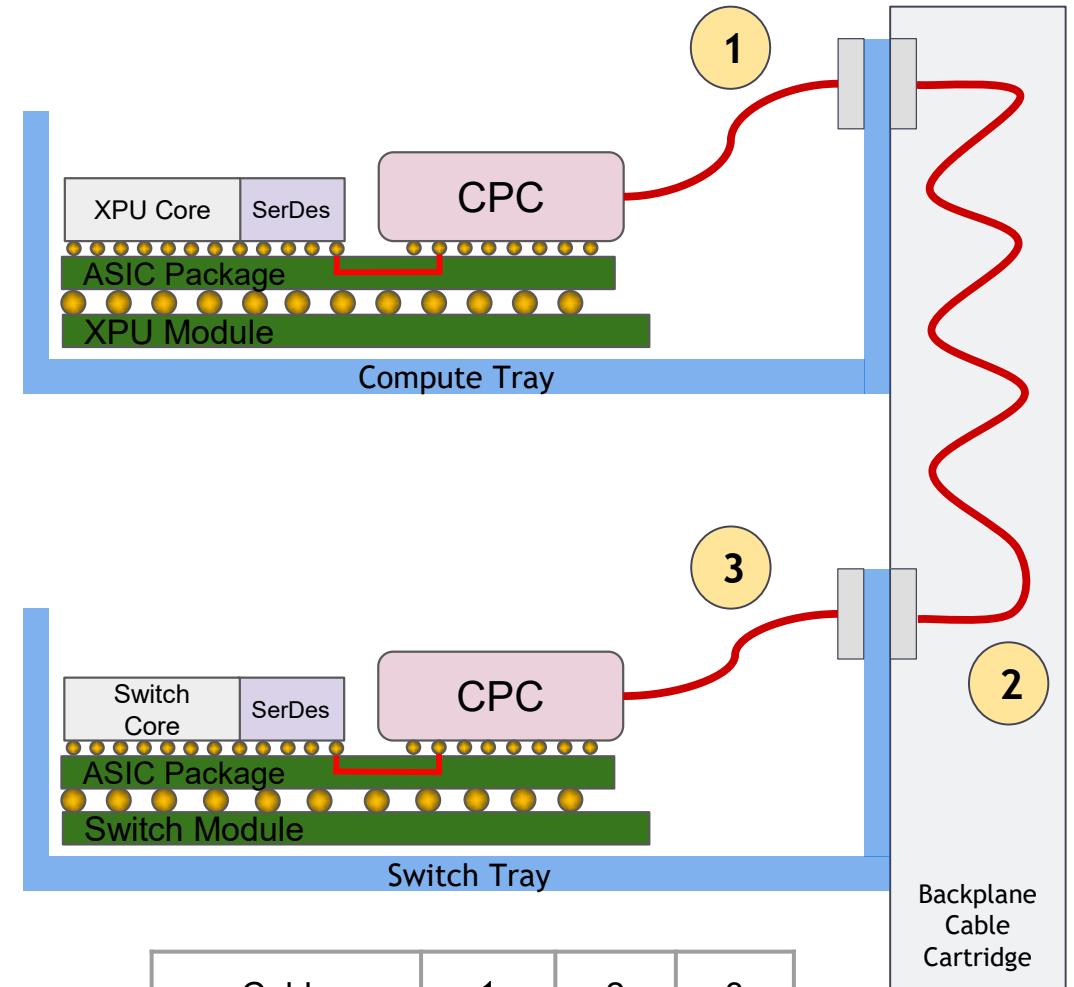
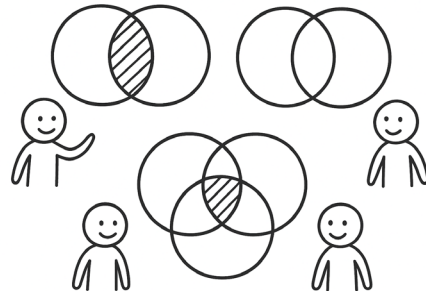
Closing the link WITHOUT a retimer

- Low power consumption
- Manageable latency

Fostering close active collaboration among

- SerDes vendors
- Package manufacturers
- Connector manufacturers
- Cable suppliers
- Rack architects

is highly needed



Cable	1	2	3
≈Length (cm)	15	140	70

# Openness Wins: Standards Side Works

- Standards-based framework
- Open discussions
- Data sharing environment



Openness over Lock -in



Compatibility First



Enable the Stack



**ESUN**



**UltraEthernet**  
Consortium

**OIF**

# Conclusion

## AI-Driven Data Center Growth is Happening FAST

- AI integration is accelerating across daily life
- Hyperscalers are expanding data centers for greater scale and efficiency
- Rising AI demand pressures interconnects to scale:
  - In, Up, Out, Accross

## Need for Next-Gen Speeds

- Transition to 400G essential:
  - Electrical communication
  - Optical communication
- Not only for enabling higher performance, but for better efficiency, and improved cost-effectiveness



# QUESTIONS?