

# AI-Driven Infrastructure

Innovating Network for Scale

Manoj Wadekar

AI Systems Technologist



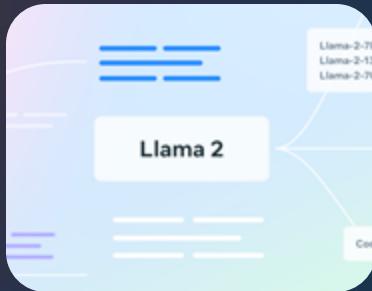


“We've come to this view that, in order to build the products that we want to build, we need to build for **general intelligence**.”



- Mark Zuckerberg, CEO Meta Inc.

# Meta AI is used for diverse cases



Large language models  
(LLMs)



Text-to-image  
generation



AI-enabled  
creation tools



Conversation topic  
growth on Instagram

# GenAI runs on Large Languages Models

Llama-2 65B  
Circa: 2023



Total Compute (PF/s)	400
Memory Capacity (TB)	10
Training Scale (GPUs)	4k

# and we are not done yet...

Llama-2  
Circa: 2023

Llama-3  
Circa: 2024

# and we are not done yet...

Llama-2  
Circa:  
2023



Llama-3  
Circa: 2024



# ...towards Multi-Modality

Llama-2  
Circa:  
2023

Llama-3  
Circa: 2024

Llama-Next  
Circa: 202x

Text  
1x Tokens

Text

7-8x Tokens

Videos

Images

Audio

# AI Silicon Evolution

Number of connected accelerators

1x

2x

8x

~512x

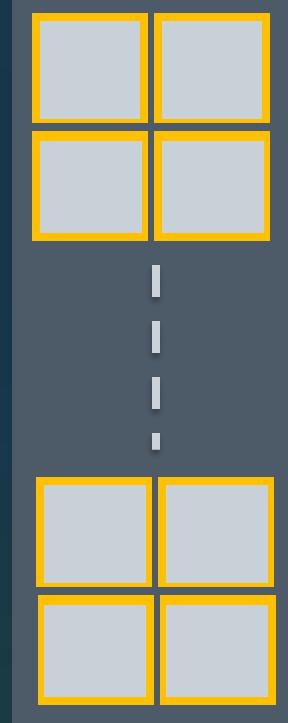
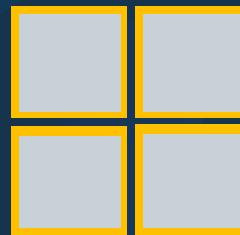
Reticle

2.5D

2.5D+3D

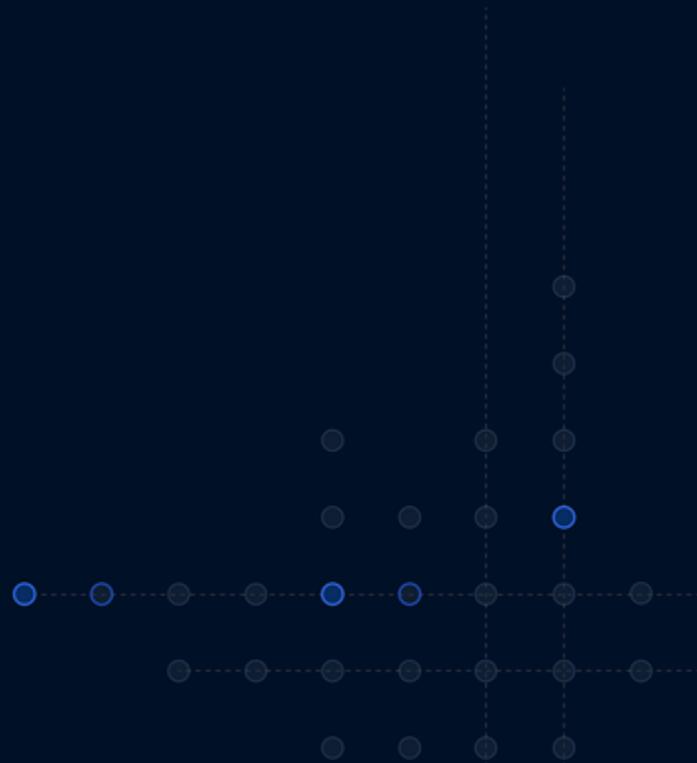
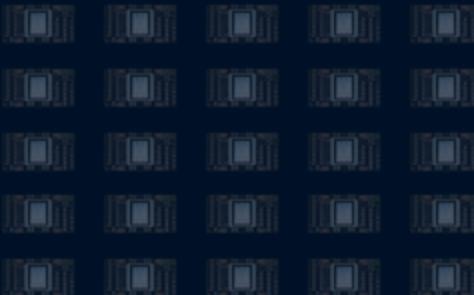
Scaleup

10x



# 6,000

2022



# 16,000

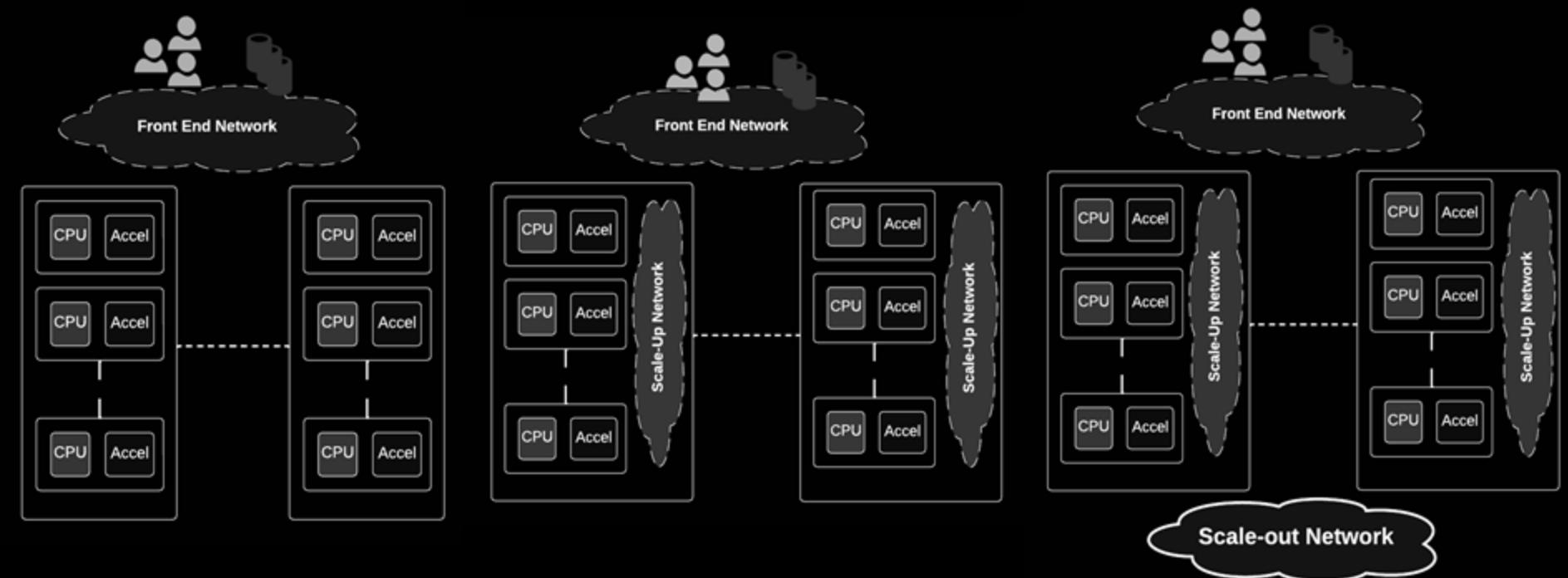
2023



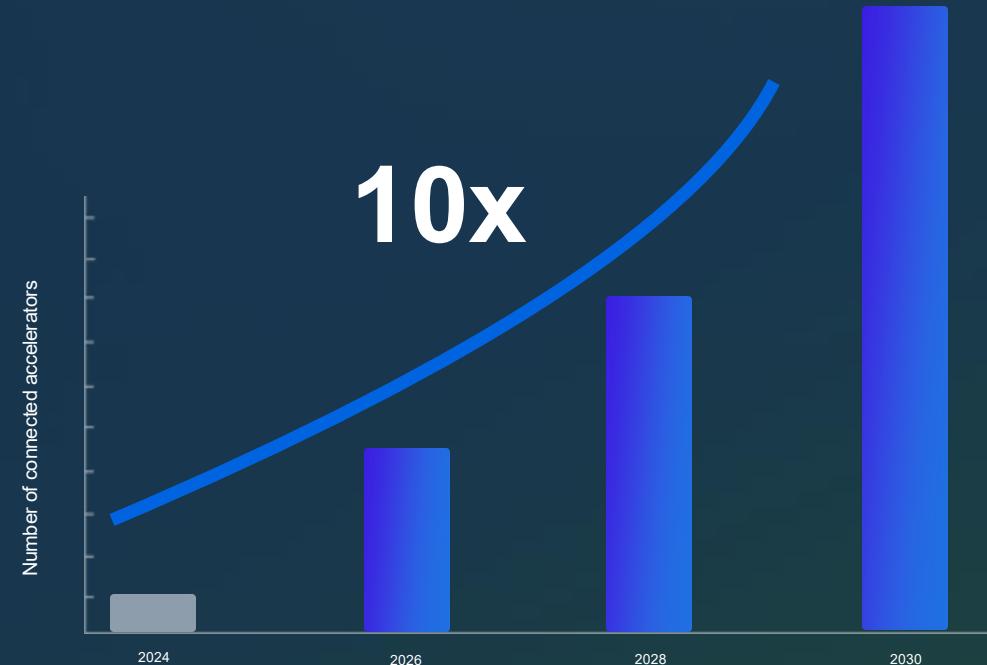
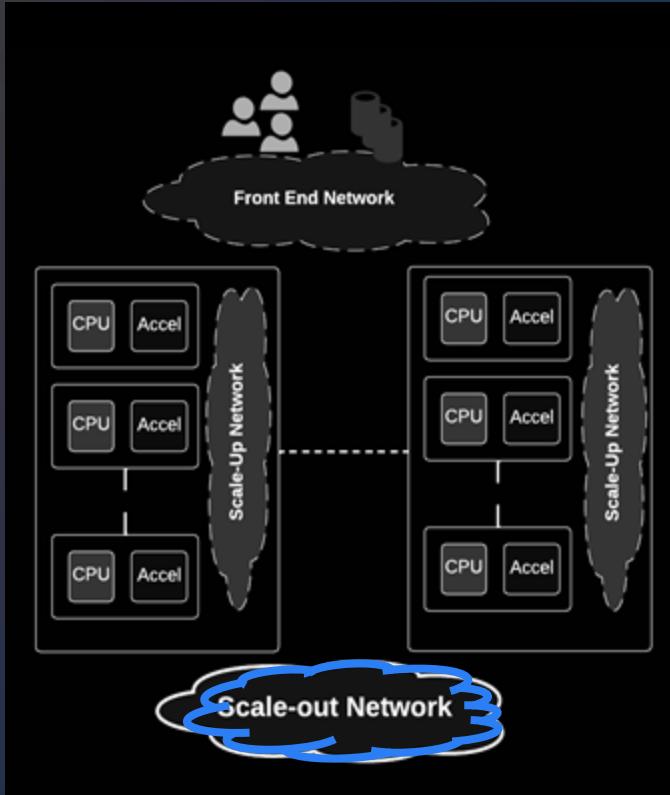
600,000

2024

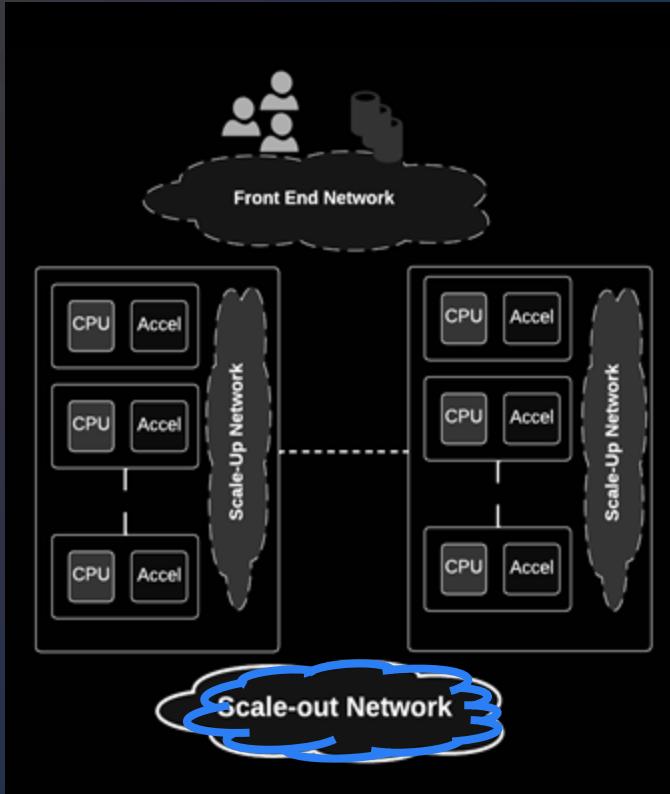
# AI Network



# AI Scale-out Network (cluster size)

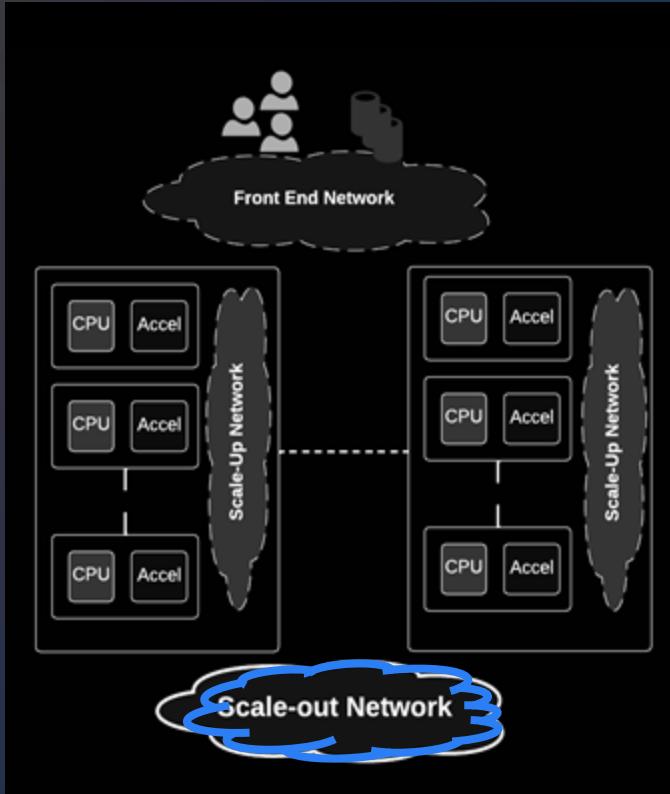


# AI Scale-out Network



- High Bandwidth, low latency
  - RDMA (RoCEv2, IB)
- Large scale (multiple 000's of GPUs)
- Zones, DCs.. and beyond
  - Longer links, higher latencies
- Larger clusters, bursty workloads
  - Need for congestion management and avoidance
  - Newer enhancements (vendor specific)
- Silo'ed infrastructure
  - Homogeneous/Heterogeneous network
- Fabric enhancements, NIC enhancements

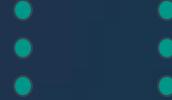
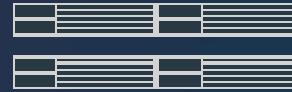
# AI Scale-out Network - Innovation Area



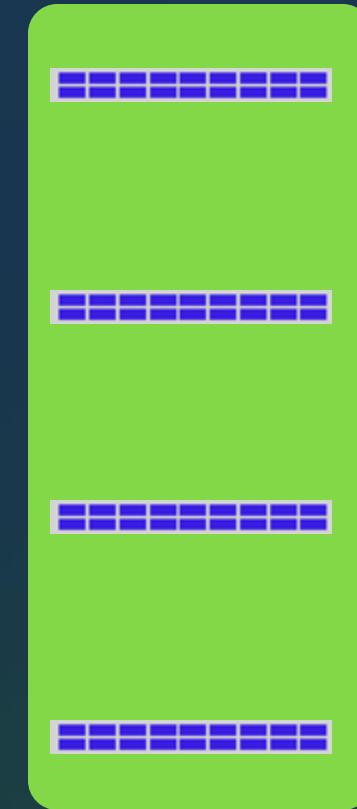
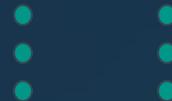
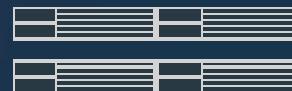
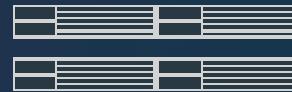
- Scale for cluster size
  - 000's → Millions
- Scale for bursty-ness
  - Job completion
  - Effective throughput
- Scale for distance
  - Driven by Cluster size and physical constraints
- Scale for Heterogeneity
  - GPU and Network
  - Practical considerations

# AI Scale-Up Connectivity...

Accelerator Bank

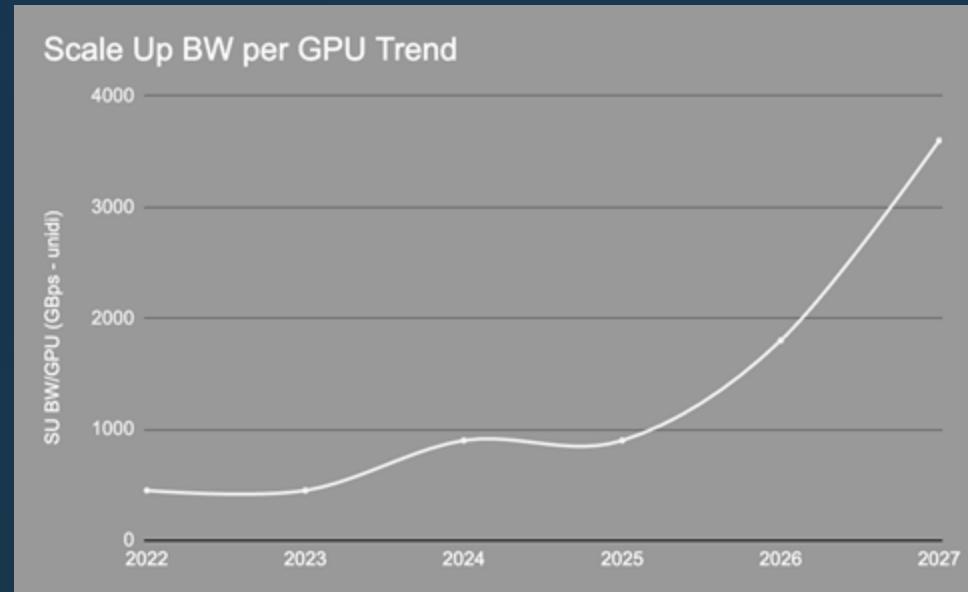
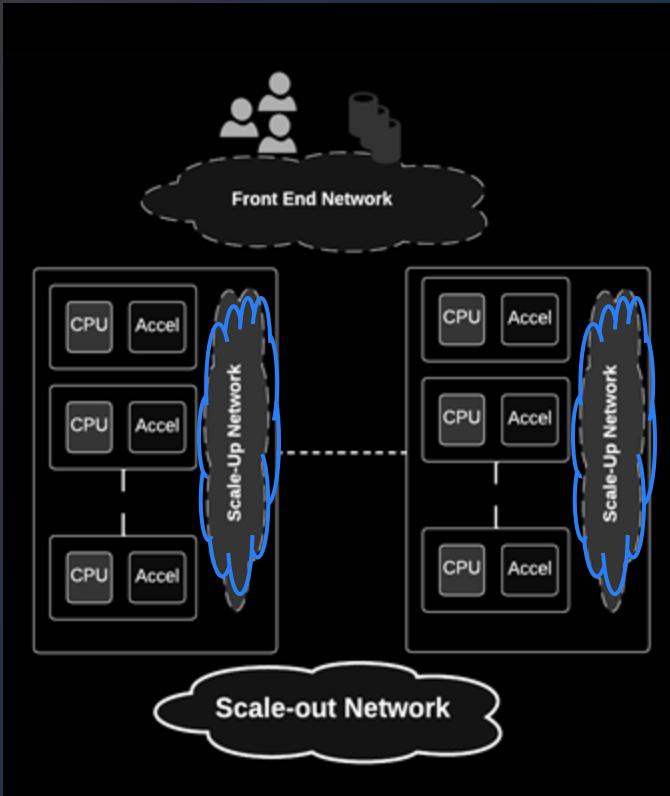


Accelerator Bank

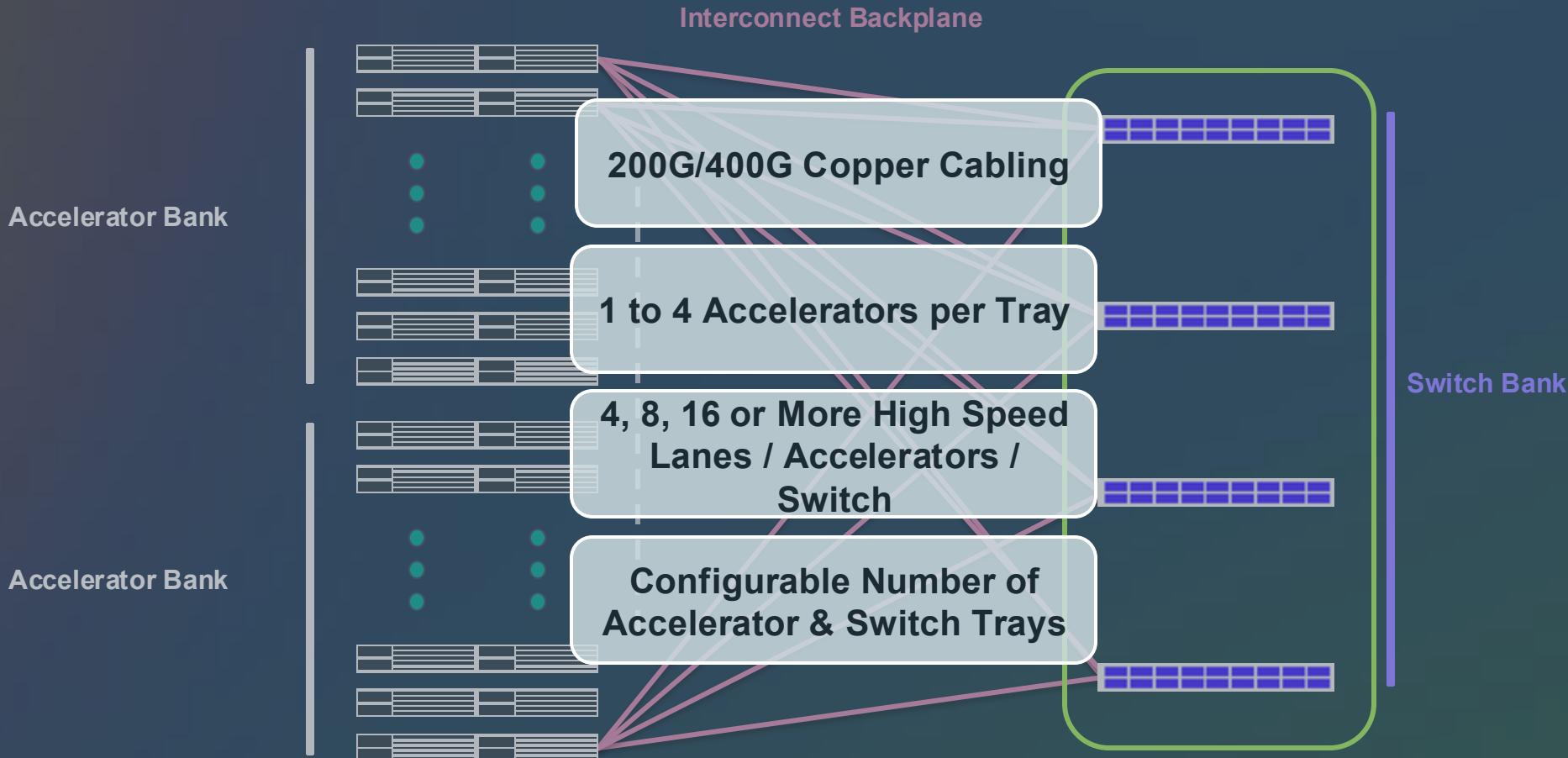


Switch Bank

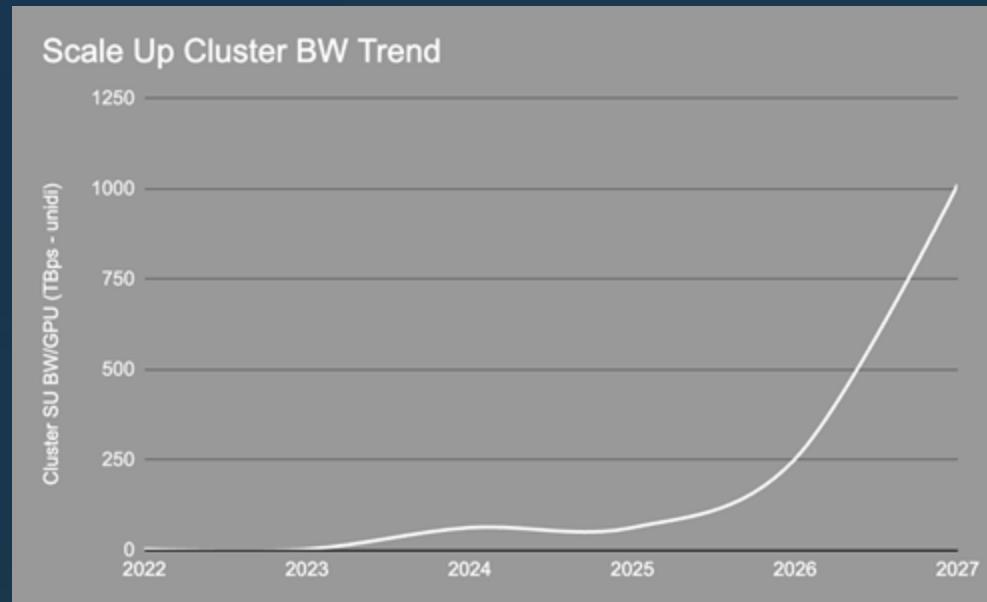
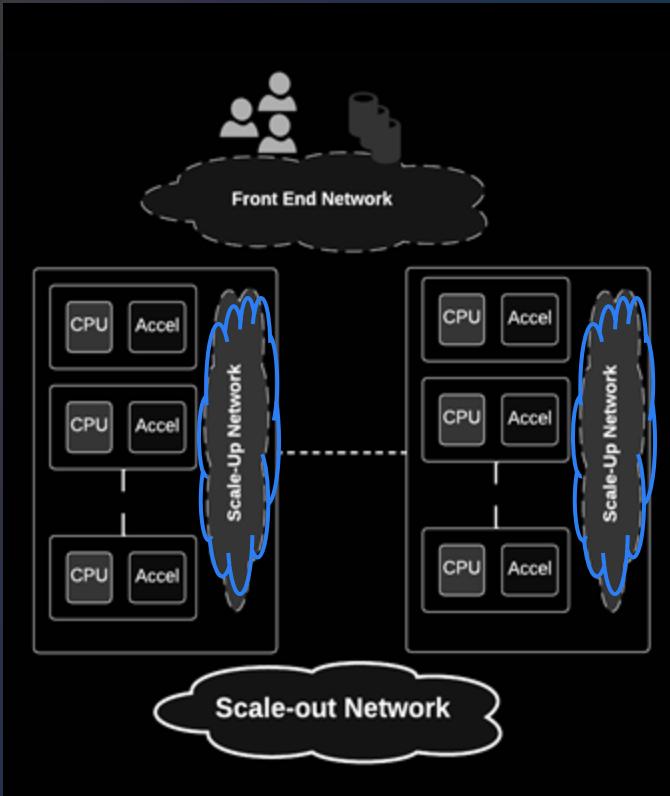
# AI Scale-Up Network



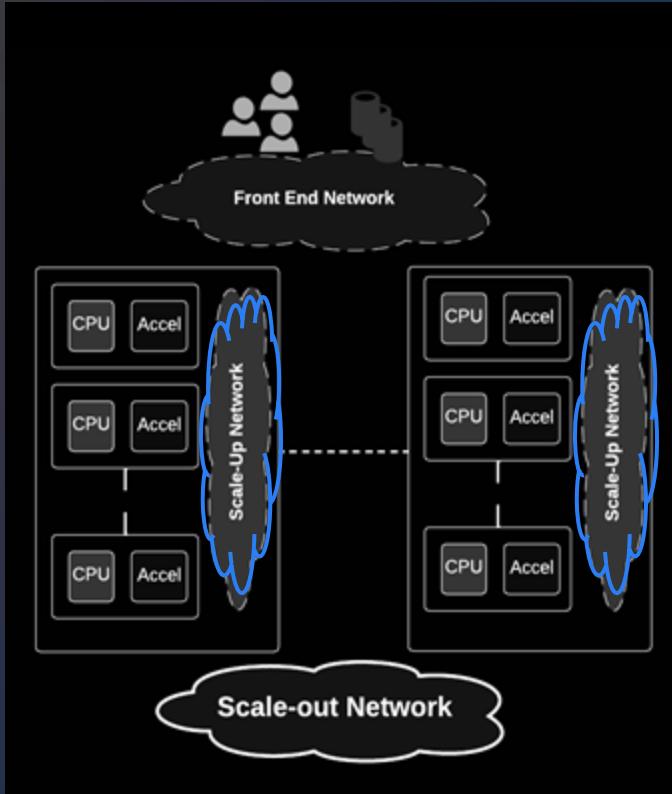
# Scale-Up Connectivity...



# AI Scale-Up Network

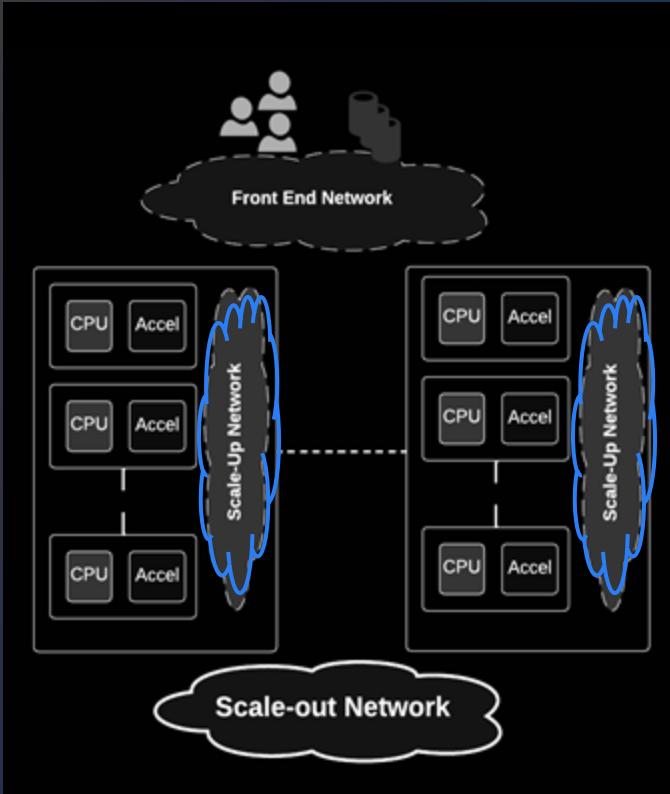


# AI Scale-Up Network



- Very High Bandwidth, low latency
  - Memory semantic (load/store), channel semantic (send/receive)
- Limited scale (8, 72 .. Growing number of GPUs)
  - Shoreline, link lengths, power (pJ/bit), reliability
- Homogeneous network
  - E.g. NVLink™, UALink™, Ethernet (RoCE) etc.

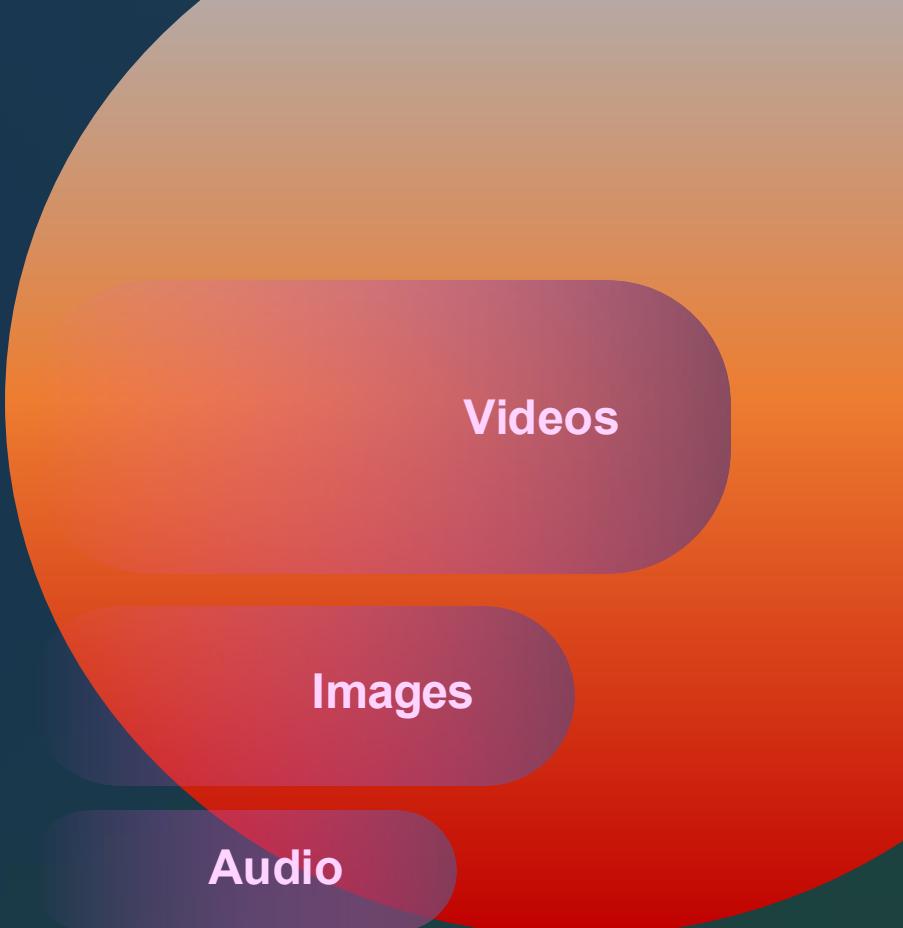
# AI Scale-Up Network - Innovation Opportunities



- Very High Bandwidth, low latency
  - Collective libraries to take advantage of new semantic support
- Growing scale-up BW and Size
  - How to optimally trade off scale-up and scale-out
  - Model innovation
- New Network Topologies and Protocols
  - What would SW like to see here?
  - **We are at crossroads of new technology here!**

# Summary

- Network critical part of AI cluster growth
- Scale Up and Scale Out networks growing exponentially
- New protocols are being defined
- How can software take advantage; what can we do better?



Videos

Images

Audio

# An Interactive Supercomputer



Thank You