# KEYSIGHT

# 6G Readiness: Technology Insights for Tomorrow
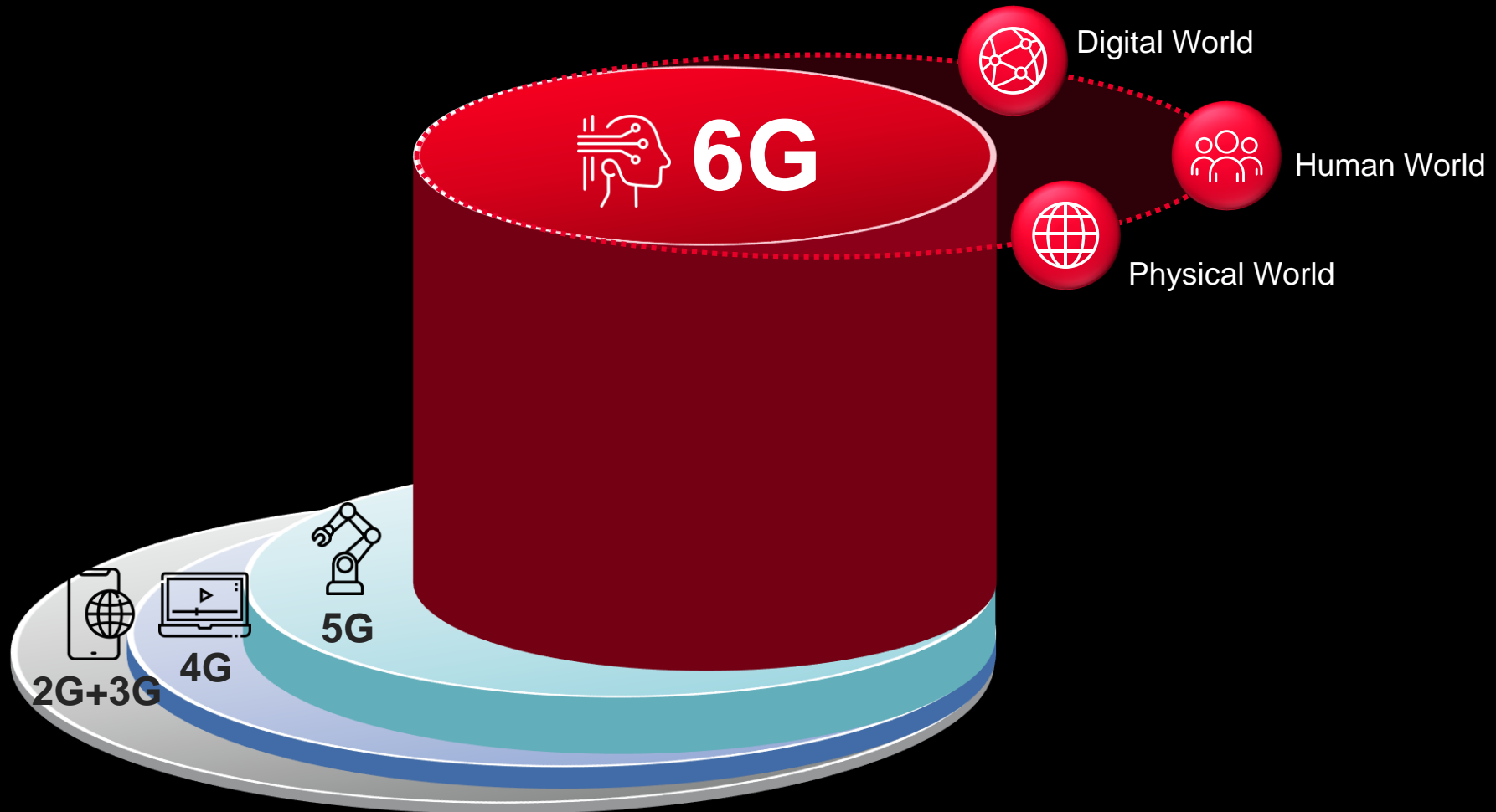
**Mombasawala Mohmedsaeed**
**May 14, 2025**
**Bharat 6G 2025**
**New Delhi**

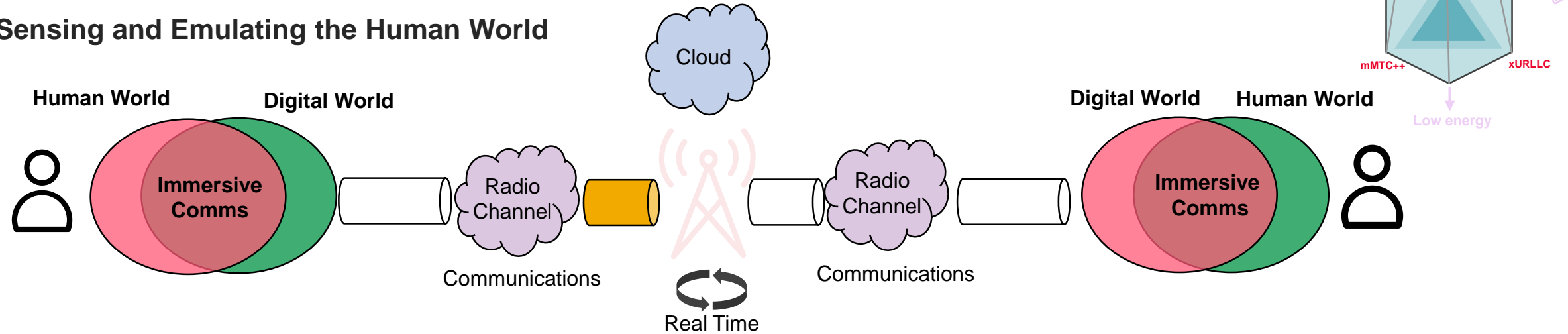# 6G Will Connect the Physical, Digital, and Human Worlds



6G

Digital World

Human World

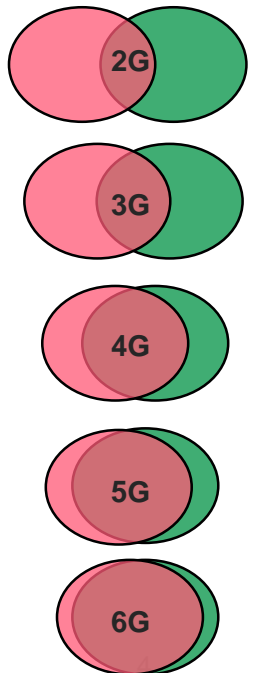Physical World

2G+3G

4G

5G

# 6G Monetizable Use Cases

# Immersive Communication – 2G to 6G Evolution

## Sensing and Emulating the Human World



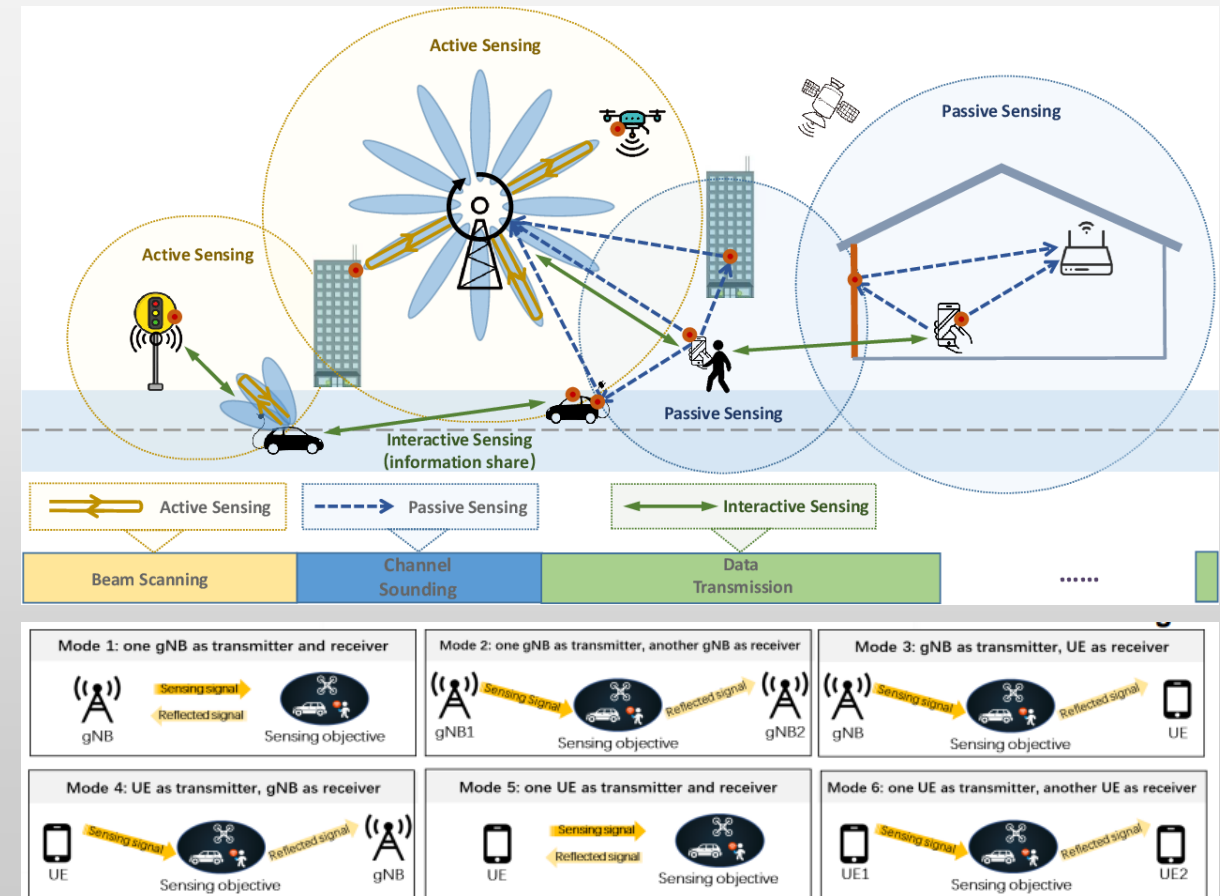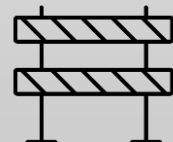| "G" | Sensor | Frequency Req | Tput Req | Latency Req | Comments / Technology Requirements |
|-----|--------|---------------|----------|-------------|-------------------------------------|
| 2G | Audio | 900MHz | 8-32kbps | <100ms | |
| 3G | Video | 900MHz / 1.8GHz | 2Mbps | <100ms | |
| 4G | Location | 900MHz / 1.8GHz | 100's bps | <1s | |
| 5G | XR | 3.5GHz | 100Mbps+ | <20ms | |
| 5G | High accuracy positioning | 3.5GHz | 1kbps | <1ms | Need wide bandwidth signal to cross correlate for delays and location |
| 6G | Immersive cloud XR | 3.5GHz | 10Gbps+ | <10ms | High resolution immersive XR<br>Watch live football game from the referee's viewpoint<br>FR3 development |
| 6G | Haptic information | n/a | | 0.1ms | Touch, motion, vibration<br>Low latency communication |
| 6G | Integrated Sensing and Communication | Mid band? THz? | ? | ? | Mid band / THz Technology development |
| 6G | Holographic display | THz? | >1Tbps | <1ms | |

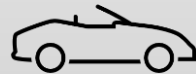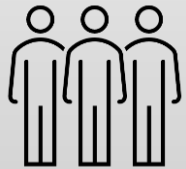# 6G- New Market Segments

## Integrated sensing and communications (ISAC)

## ISAC

R19 SI Focus: Define channel modelling aspects to support object detection & tracking

# Four Key Technology Areas Driving 6G

**New Spectrum Technologies**
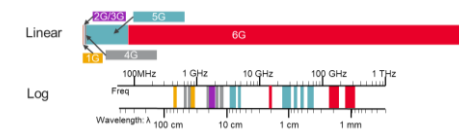
**Artificial Intelligence and Machine Learning**

**Digital Twins**

**New Network Topologies**

# 6G Candidate Spectrum: Specifics

| | 6G Research Topics | | | | | | | | | | Mobile Regulatory Situation | Technical Challenges |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NOMA | Waveforms | Channel Coding | Unlicensed/WiFi | Advacned MIMO | Satellite | Mobility/Coverage | Radar/ISAC | PA & LNA | Antennae | | |
| <7GHz | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | • Moderate changes ongoing e.g. 3.4 GHz and 6-7 Ghz. Most allocations/auctions complete. | • Coverage<br>• Spectral Efficiency |
| 7-16GHz | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | • Entire band has co-primary use<br>• Heavy Federal/DoD Allocation<br>• Most EU states ambivalent at best<br>• Passive (EES) Satellite & Radio Astronomy co-existence<br>• ITU Decisions WRC-27 or later | • Co-existence/Sharing<br>• Coverage and Link Budget vs. Cell Density |
| 16-24GHz | | | | | | ✓ | | ✓ | | | | • "FR2-like" (more challenging than <16GHz) |
| 24-52 GHz | | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | • 24-52 Allocated allocated to Mobile IMT use | • Coverage<br>• Energy Efficiency<br>• Mobility |
| 52-71GHz | | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | • 57-71 Unlicensed | |
| 71-110GHz | | ✓ | ✓ | | ✓ | | | ✓ | ✓ | ✓ | • Point-To-Point (71-76/81-86) & Automotive Radar<br>• Inadequate contiguous sub-bands.<br>• Heavy constraints 90-110 | • Coverage<br>• Energy Efficiency<br>• Noise BW<br>• Mobility |
| 110-170GHz | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | • Lightly regulated<br>• ITU RR-5.340 Constraints: Radio Astronomy/EES<br>• ITU decisions WRC-31 or later | • Coverage<br>• Energy Efficiency<br>• Link Budget<br>• Noise BW<br>• Mobility |
| >170GHz | | ✓ | | | | | ✓ | ✓ | ✓ | ✓ | • Lightly regulated so far<br>• ITU RR-5.340 Constraints: Radio Astronomy/EES<br>• ITU Decisions WRC-31 or later | |

KEYSIGHT

# 6G New Spectrum

## "FR3" and sub-THz under Evaluation

Precision positioning

Hi Res 3D Imaging

Mass spectrometry

Best coverage, lower capacity (1Gbps)
Congested bands
More allocations needed to track traffic growth

High capacity (10Gbps)
Wide bandwidth
cm level positioning

Ultra high capacity (1Tbps)
Ultra wide bandwidth
Sensing applications

**WRC-23**

| Re-farming, harmonization, explore new bands | Harmonization, explore new bands | Create new Sub-THz bands |

**6G**

| Lower Mid | Upper Mid | mmWave | sub THz |
| Low | Mid | High | |

**5G 3GPP**

| FR1 | cm Wave / FR3 – not official term | FR2-1 | FR2-2 | |

. . .   . . .

0   2   4   6   8   10   12   14   16   20   40   60   80   100   200   300   400



**KEYSIGHT**

8

# Satellite Communications – 3GPP NTN Architectures
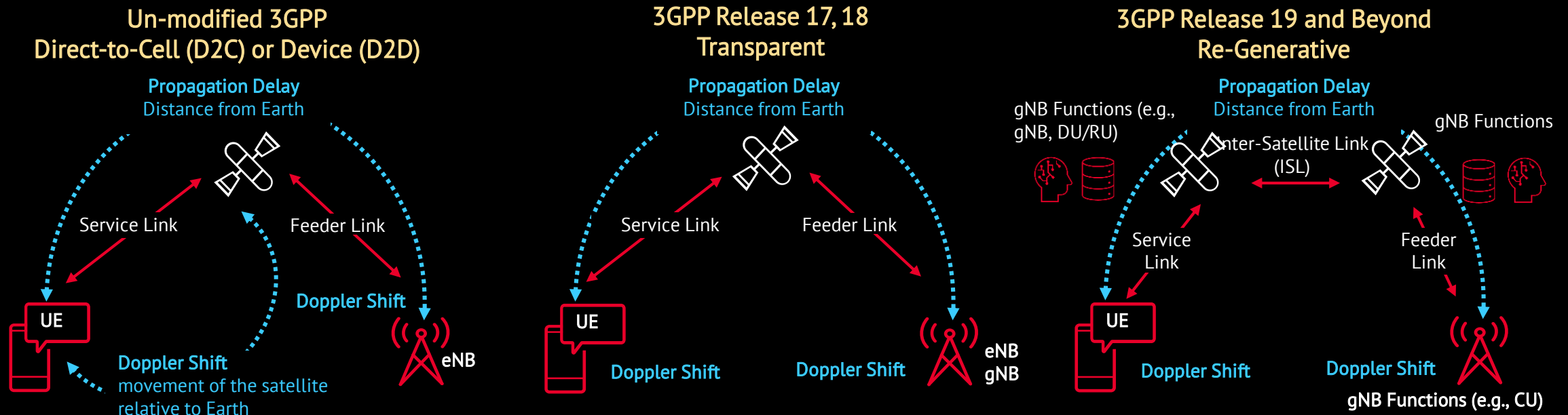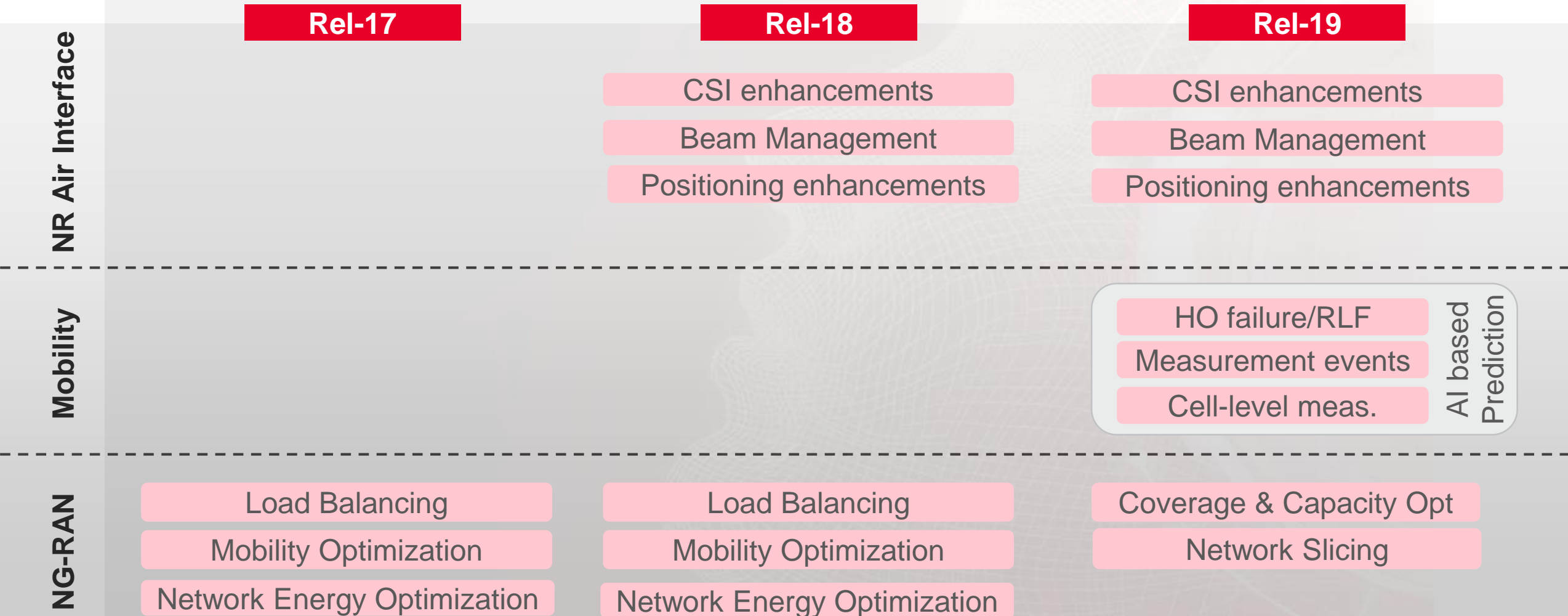
## Space to Earth – Satellite Options

**Common Challenges**: round-trip delays and frequency shifts due to the movement of the satellite relative to Earth (doppler shift)



Un-modified 3GPP
Direct-to-Cell (D2C) or Device (D2D)

Propagation Delay
Distance from Earth

Service Link

Feeder Link

Doppler Shift

UE

Doppler Shift
movement of the satellite
relative to Earth

eNB

3GPP Release 17, 18
Transparent

Propagation Delay
Distance from Earth

Service Link

Feeder Link

UE

Doppler Shift

Doppler Shift

eNB
gNB

3GPP Release 19 and Beyond
Re-Generative

gNB Functions (e.g.,
gNB, DU/RU)

Propagation Delay
Distance from Earth

Inter-Satellite Link
(ISL)

gNB Functions

Service
Link

Feeder
Link

UE

Doppler Shift

Doppler Shift

gNB Functions (e.g., CU)

**KEYSIGHT**

# 5G Advanced leading to 6G– Smarter with AI/ML

|  | **Rel-17** | **Rel-18** | **Rel-19** |
|---|---|---|---|
| **NR Air Interface** |  | CSI enhancements | CSI enhancements |
|  |  | Beam Management | Beam Management |
|  |  | Positioning enhancements | Positioning enhancements |
| **Mobility** |  |  | HO failure/RLF · Measurement events · Cell-level meas. — AI based Prediction |
| **NG-RAN** | Load Balancing | Load Balancing | Coverage & Capacity Opt |
|  | Mobility Optimization | Mobility Optimization | Network Slicing |
|  | Network Energy Optimization | Network Energy Optimization |  |

KEYSIGHT

# AI Infrastructure

## Adapting Hyperscale DC to Edge AI

- Training Clusters: 100k+ GPUs in 2024 and **path to 600k**

- 800G/1.6T links, 112/224G lanes and path to 448G

- Power need **100+MW, 160% increase by 2030**

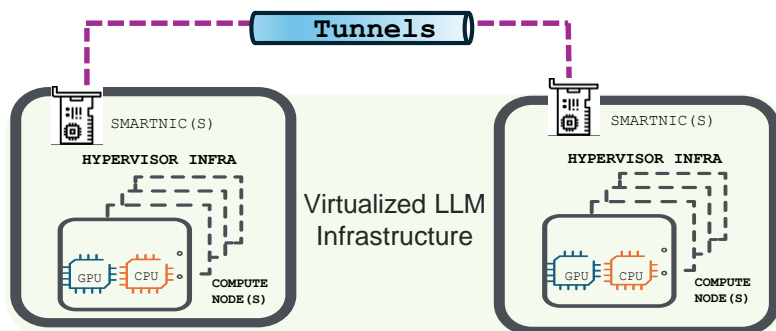- New protocols for transport and congestion management



Source: Marvell



Edge Network

**AI Edge / IoT**

>100 km - 60 mi

10 km - 6 mi

Frontend Network

Spine

500m - 0.31 mi

Leaf

50m – 164 ft

ToR

**Pod**

3 m - 10 ft

Backend Network

**AI Servers**

CPU  NIC

PCIe switch

**RDMA NICs**

GPU  GPU

**GPU Fabric**

**xPU + HBM**

**SAN**

Die-to-Die

**Chiplets**

1 mm -.04 in

# The Operation of AI ML Network Infrastructure

**Backend Data Center for AI Models Training**
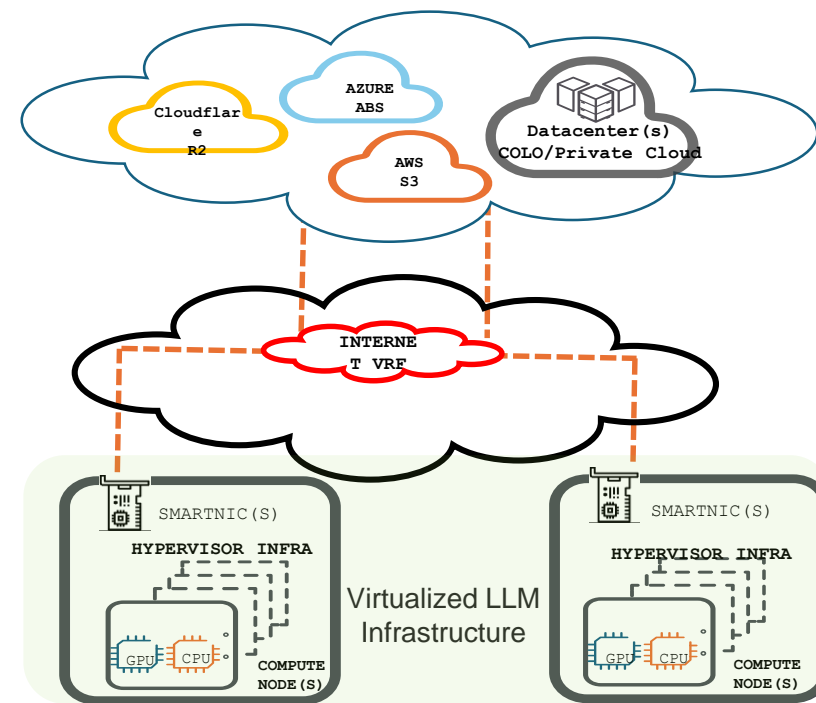
East-West Traffic Test Demands -

o   Distributed GPU/CPU architectures
o   Collective communications & parallel processing among GPU nodes
o   Hyper-virtualized infrastructures for multi tenancy
o   Immense performance needs for lossless connectivity and minimum tail-end latency

**Front-end Data Center for Inference Workloads**

North-South Network Traffic Test Demands -

o   GPUs need high-speed access to block/remote storages
o   Provisions to secure data in motion
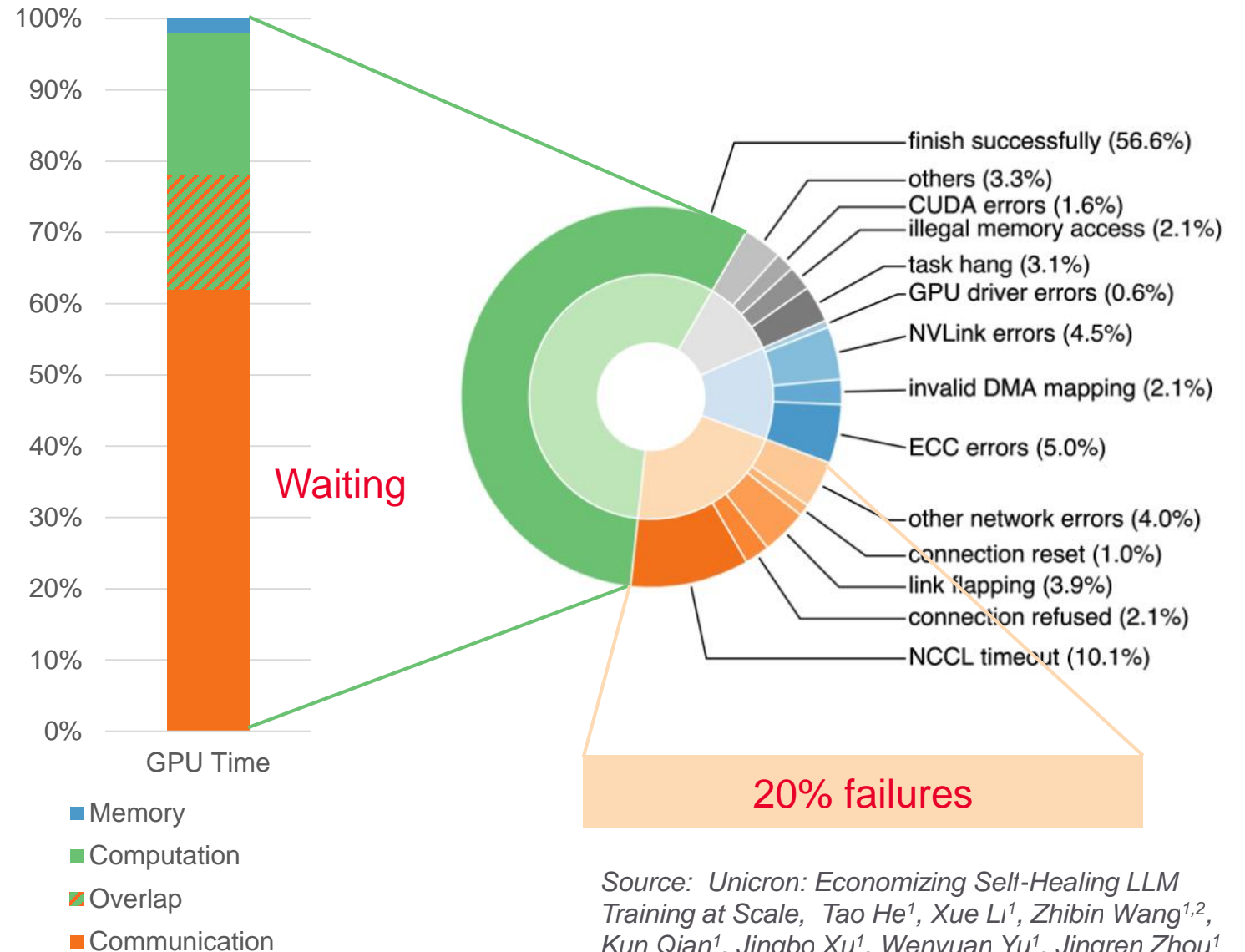o   Ultra-low latency demands

# Why the Network & Components Matters in an AI Cluster

AI is Compute, Network & Data Intensive and requires validation at System Scale

Network failures
**>20%**

GPUs waiting on data
**>50%**



Waiting

**20% failures**

Legend:
- Memory
- Computation
- Overlap
- Communication

GPU Time

finish successfully (56.6%)
others (3.3%)
CUDA errors (1.6%)
illegal memory access (2.1%)
task hang (3.1%)
GPU driver errors (0.6%)
NVLink errors (4.5%)
invalid DMA mapping (2.1%)
ECC errors (5.0%)
other network errors (4.0%)
connection reset (1.0%)
link flapping (3.9%)
connection refused (2.1%)
NCCL timeout (10.1%)

*Source: Unicron: Economizing Self-Healing LLM Training at Scale, Tao He[1], Xue Li[1], Zhibin Wang[1,2], Kun Qian[1], Jingbo Xu[1], Wenyuan Yu[1], Jingren Zhou[1] [1]Alibaba Group, [2]Nanjing University*

*Vision transformer example. Source: https://github.com/facebookresearch/HolisticTraceAnalysi*

**KEYSIGHT**

13

# AI Model Training
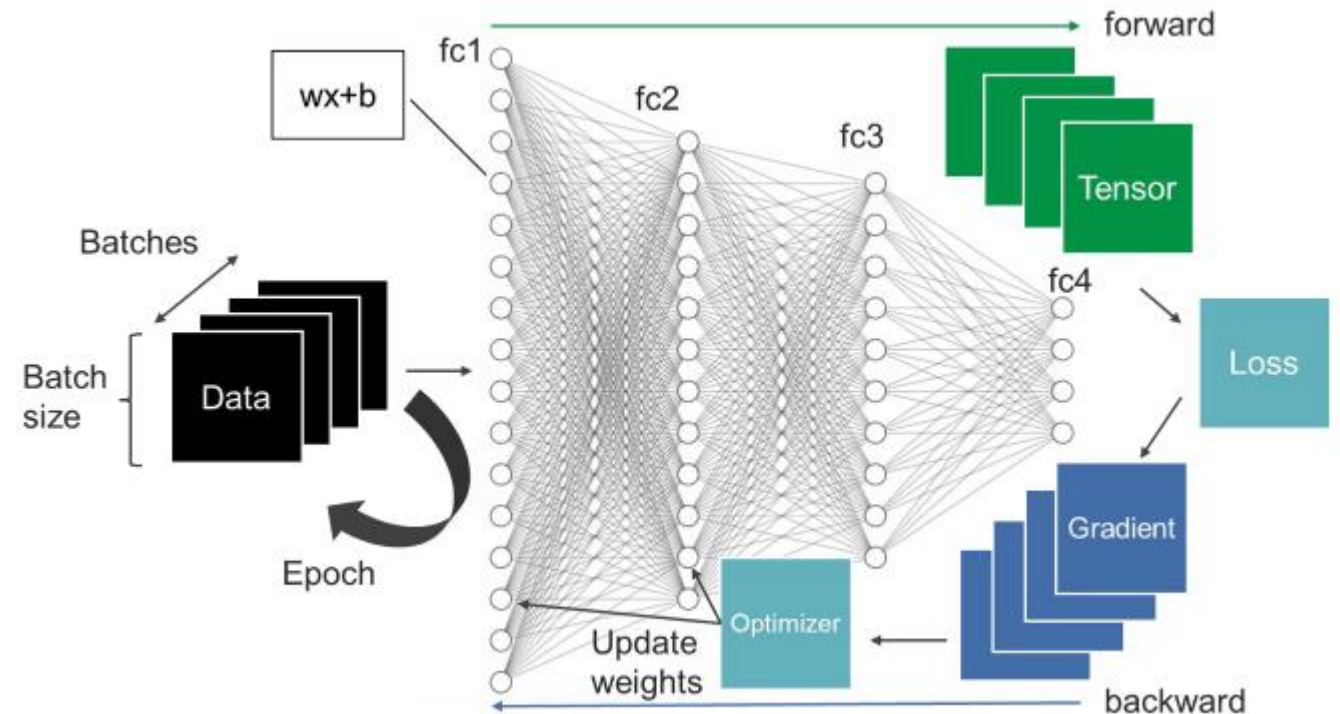
## 3 Step Process

### Step 1: Data preparation
- Collect and preprocess large datasets (for example, text files, images, and audio).
- Tokenize and normalize data to ensure consistency and efficiency.
- Split data into training, validation, and testing sets.

### Step 2: Model definition
- Define the architecture of the AI model (for example, neural network and decision tree).
- Specify hyperparameters (for example, learning rate, batch size, and number of layers).
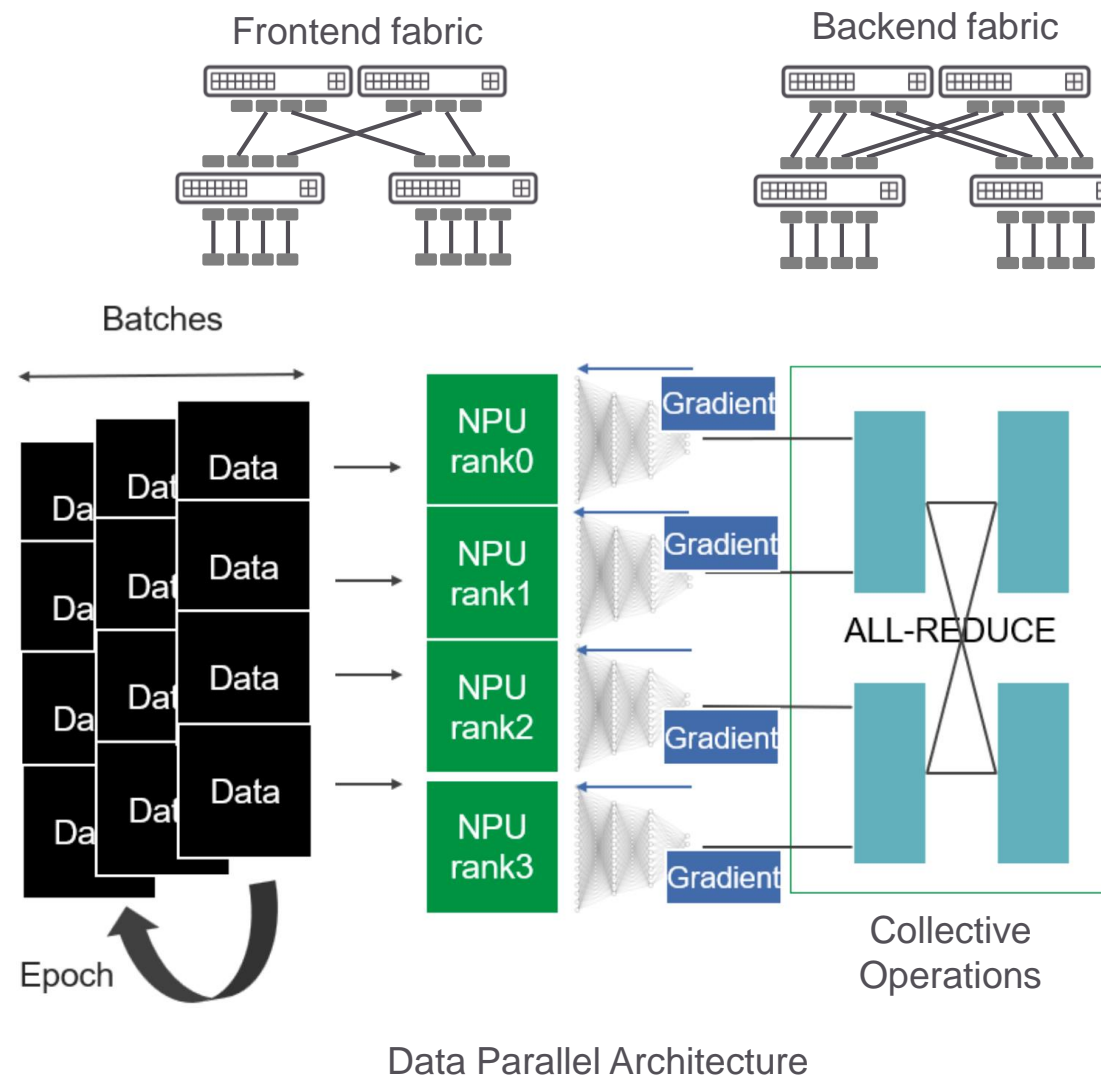
### Step 3: Model training
- Initialize the model's weights and biases.
- Feedforward pass: Compute outputs for each sample in the training set.
- Backpropagation: Calculate gradients and update model parameters by using an optimization algorithm (for example, Stochastic Gradient Descent and Adam).
- Repeat the preceding steps until convergence or a stopping criterion is reached.
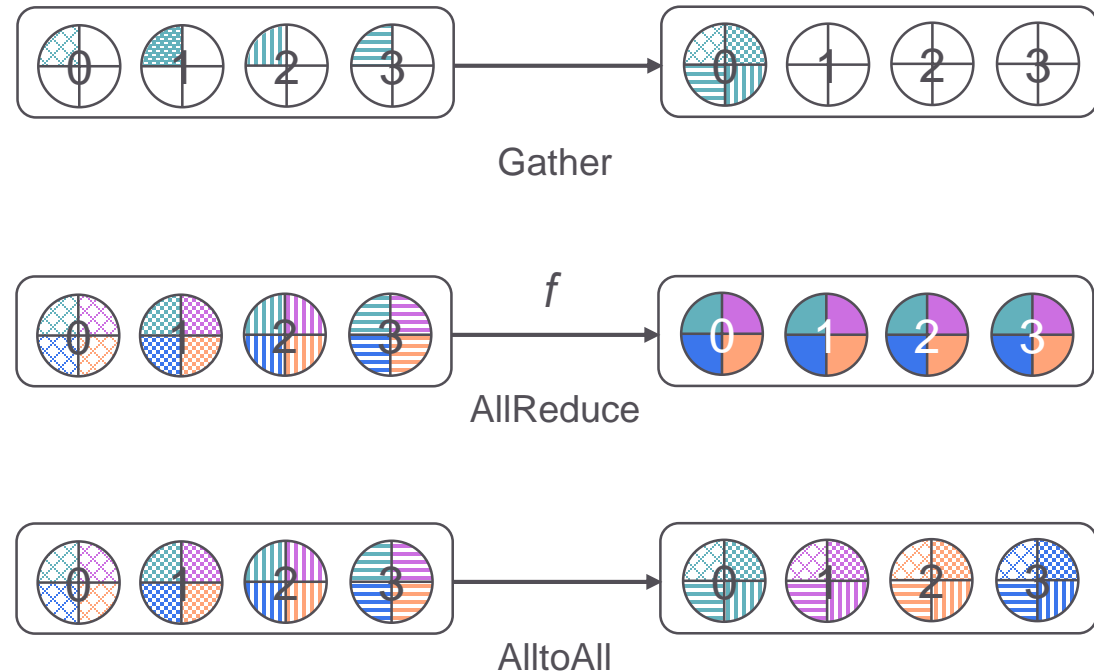
# Network role in AI clusters

**Scaling up systems, scaling out clusters**

- Accelerate model training with **Data Parallelism**

- Split large models across GPUs with **Tensor** and **Pipeline Parallelism**

- Subdivide complex problems among several models with **Mixture of Experts**



Data Parallel Architecture
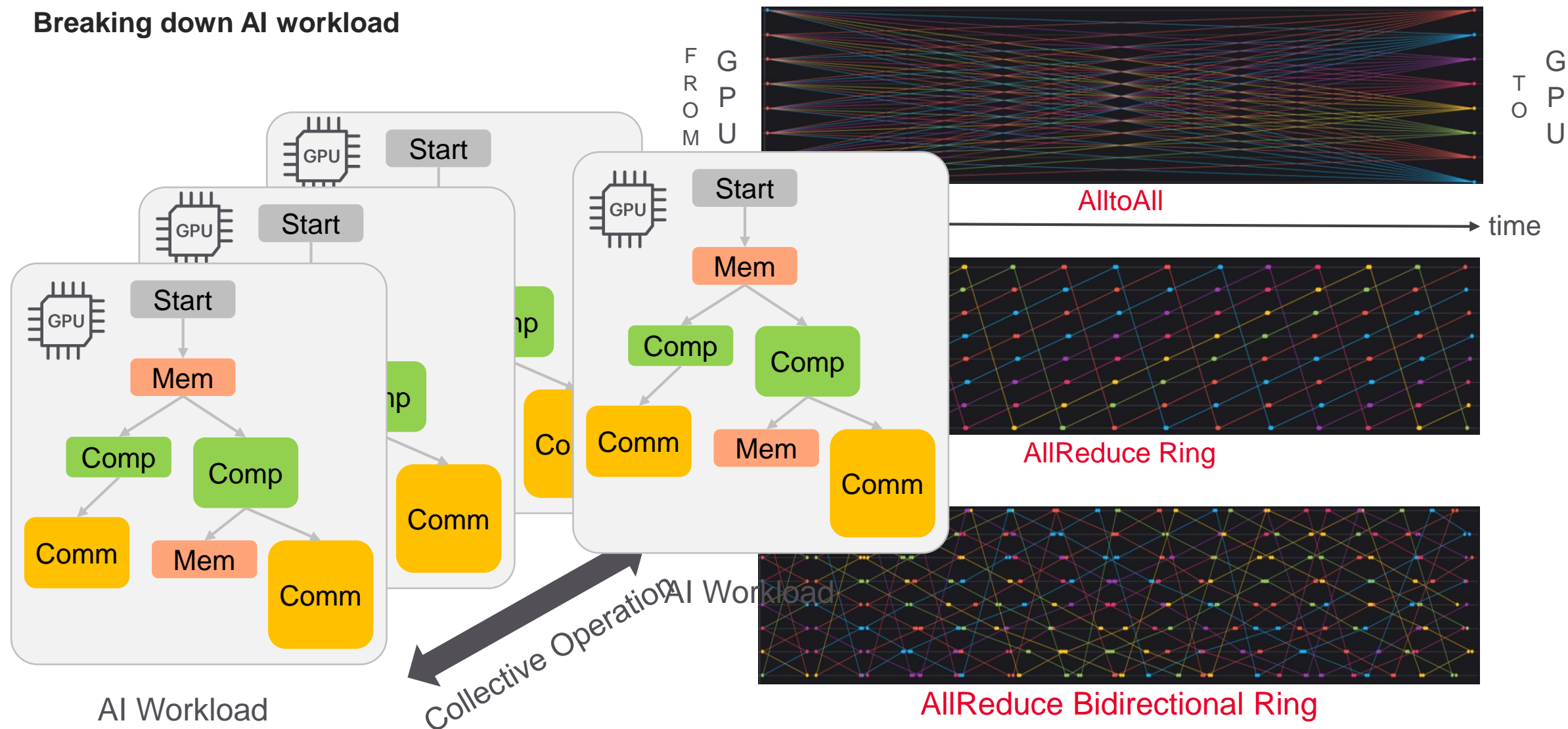
# Types of Collective Operations

- Common types for AI workloads:
  - Broadcast
  - Gather
  - AllReduce
  - AlltoAll
  - ReduceScatter
  - AllGather

- Reduce implies math with data ($f$)

- *All* or *Scatter* – symmetry



Gather

AllReduce

AlltoAll

# GPU Communications

**Breaking down AI workload**

**Examples of Collective Operations**



AlltoAll



AllReduce Ring



AllReduce Bidirectional Ring

AI Workload

Collective Operation

# Network is the bottleneck in AI model training

Job Completion Time Factors
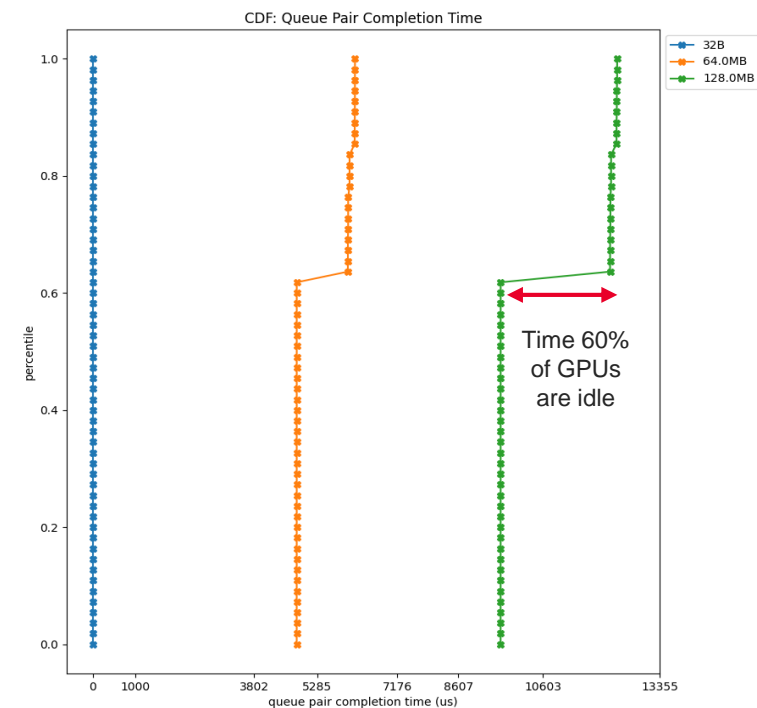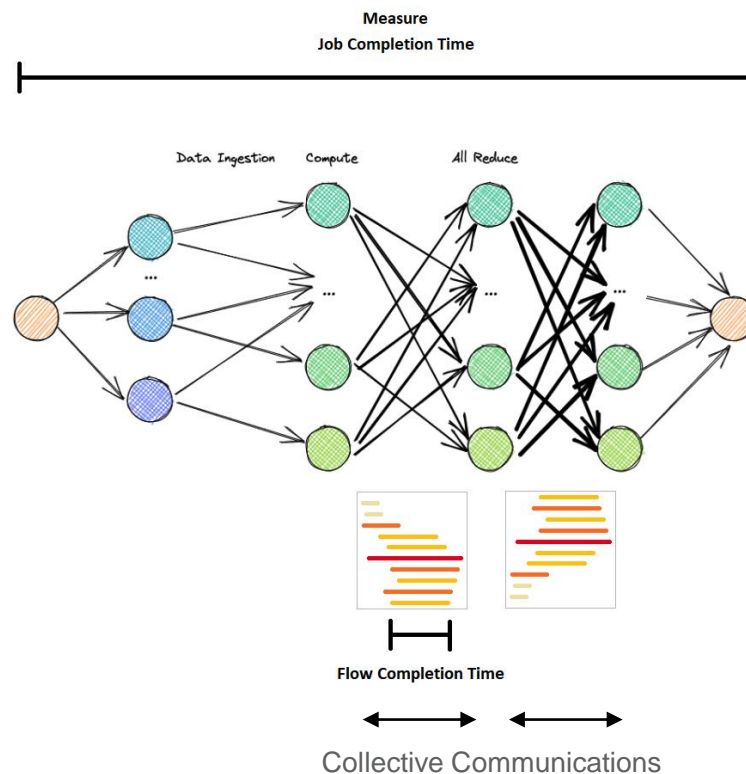- Data Ingestion
- Computation
- Collective Communications

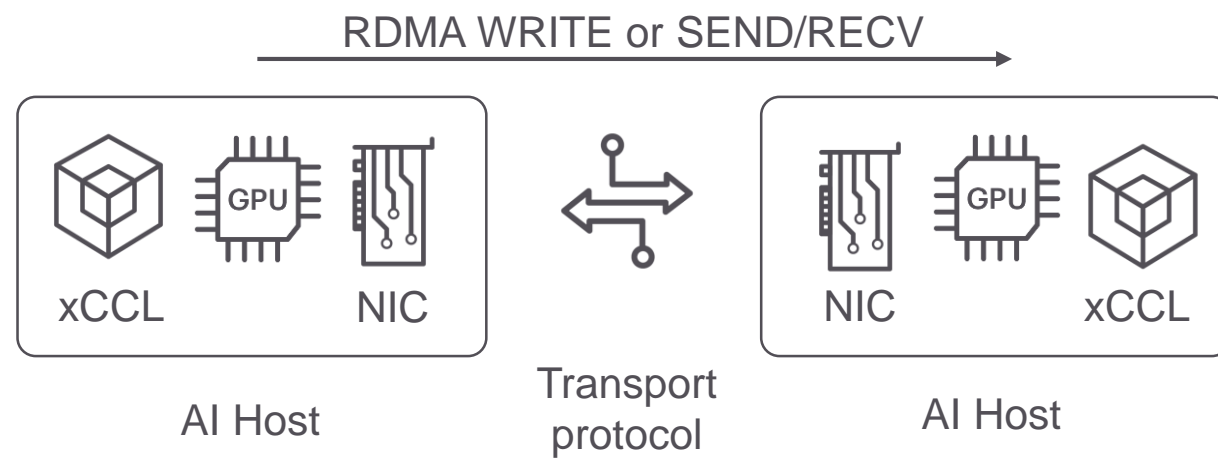Network tail latency
- Defines wasted GPU time

Contributors
- Data exchange algorithm
- Software stack
- System I/O
- DPU (NIC)
- Network fabric



**Measure**
**Job Completion Time**

Data Ingestion    Compute    All Reduce

**Flow Completion Time**

Collective Communications



CDF: Queue Pair Completion Time

— 32B
— 64.0MB
— 128.0MB

percentile

Time 60% of GPUs are idle

queue pair completion time (us)

# RDMA and Transport Protocols

**Hardware accelerated Remote Direct Memory Access**

RDMA WRITE or SEND/RECV →

xCCL   GPU   NIC

AI Host

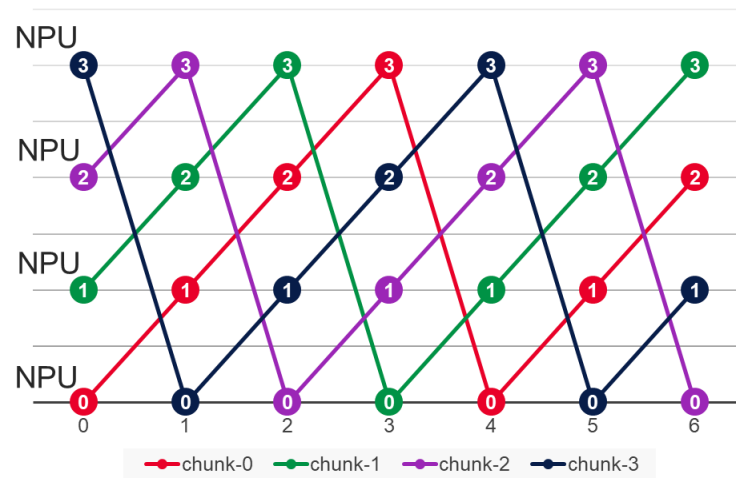Transport protocol

NIC   GPU   xCCL

AI Host

## Ethernet transport options

- RoCEv2

- Falcon

- Custom / Proprietary

- Ultra Ethernet (future)

**KEYSIGHT**
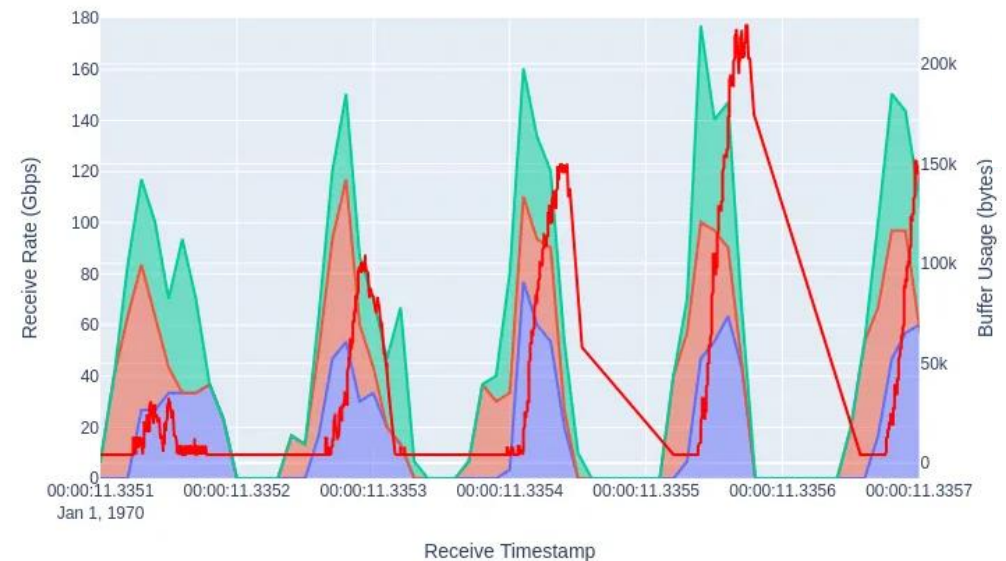
# Practical challenges

- Synchronized start – from 0 to line rate on all ports

- Flow dependencies – latencies accumulate

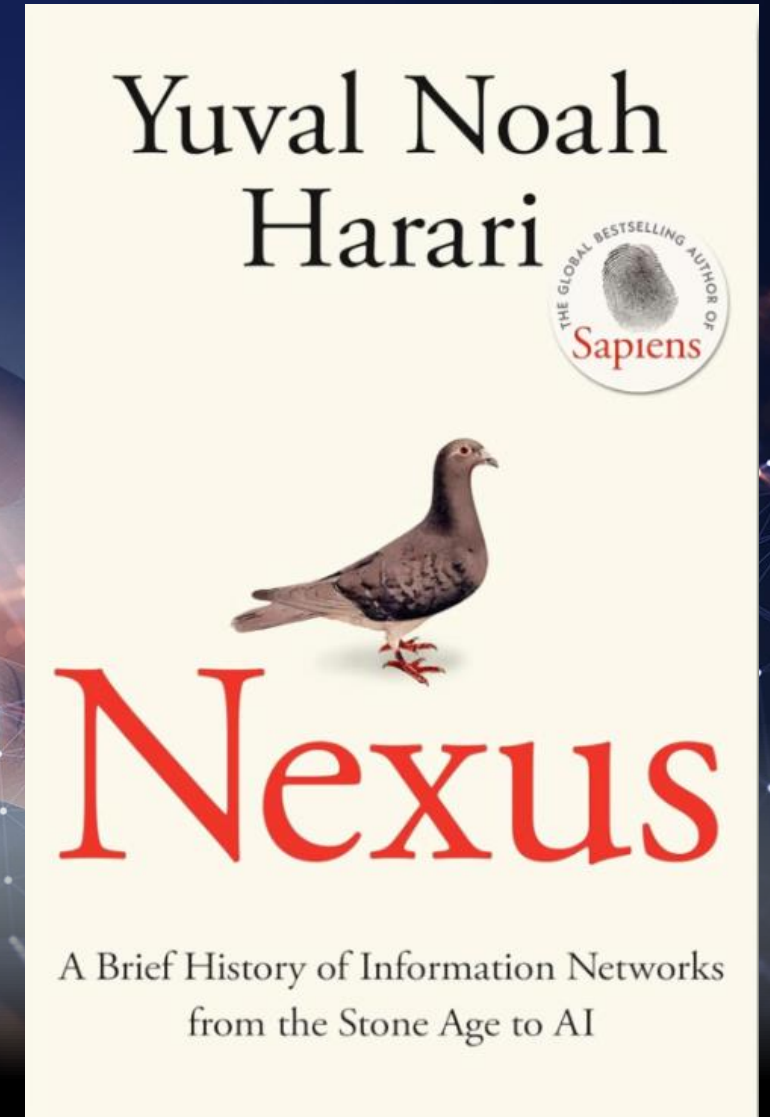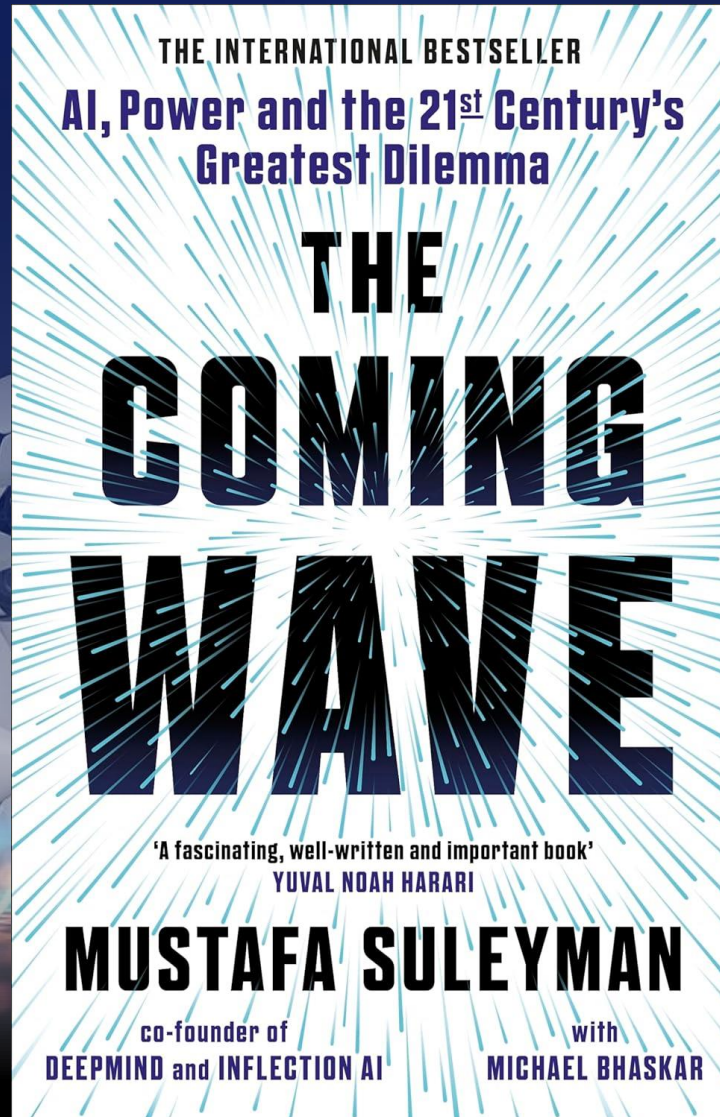- Low entropy – hard to load balance

- RDMA message bursts – incast



Unequal Load Balancing



Latencies



Incast

# Suggested Reading………



**The Age of AI** — A WALL STREET JOURNAL BESTSELLER
Henry Kissinger × Eric Schmidt × Daniel Huttenlocher
'Should be read by anyone trying to make sense of geopolitics today' *Financial Times*

**THE INTERNATIONAL BESTSELLER**
AI, Power and the 21st Century's Greatest Dilemma
**THE COMING WAVE**
'A fascinating, well-written and important book' YUVAL NOAH HARARI
**MUSTAFA SULEYMAN**
co-founder of DEEPMIND and INFLECTION AI — with MICHAEL BHASKAR

**Yuval Noah Harari** — THE GLOBAL BESTSELLING AUTHOR OF Sapiens
**Nexus** — A Brief History of Information Networks from the Stone Age to AI

# Finally, Some of my favorite quotes



IT ALWAYS SEEMS **IMPOSSIBLE** UNTIL IT'S DONE.
-NELSON MANDELA

"Insanity is doing the same thing over & over again & expecting different results."
*Albert Einstein*

WE CAN DO NO GREAT THINGS ONLY SMALL THINGS WITH GREAT LOVE
*Mother Teresa*
celebquote.com

To err is human, to persist in error is diabolical.
Georges Canguilhem

One Of The Greatest Diseases Is To Be Nobody To Anybody.
~ MOTHER TERESA ~
Statustown.com

I have decided to stick with love. Hate is too great a burden to bear.
Martin Luther King, Jr.

Do not lower your goals to the level of your abilities. Instead, raise your abilities to the height of your goals.
KEYSIGHT

"If you buy things you do not need, soon you will have to sell things you need."
- Warren Buffett

"The good thing about science is that it's true whether you believe in it or not"
Neil deGrasse Tyson

22

KEYSIGHT