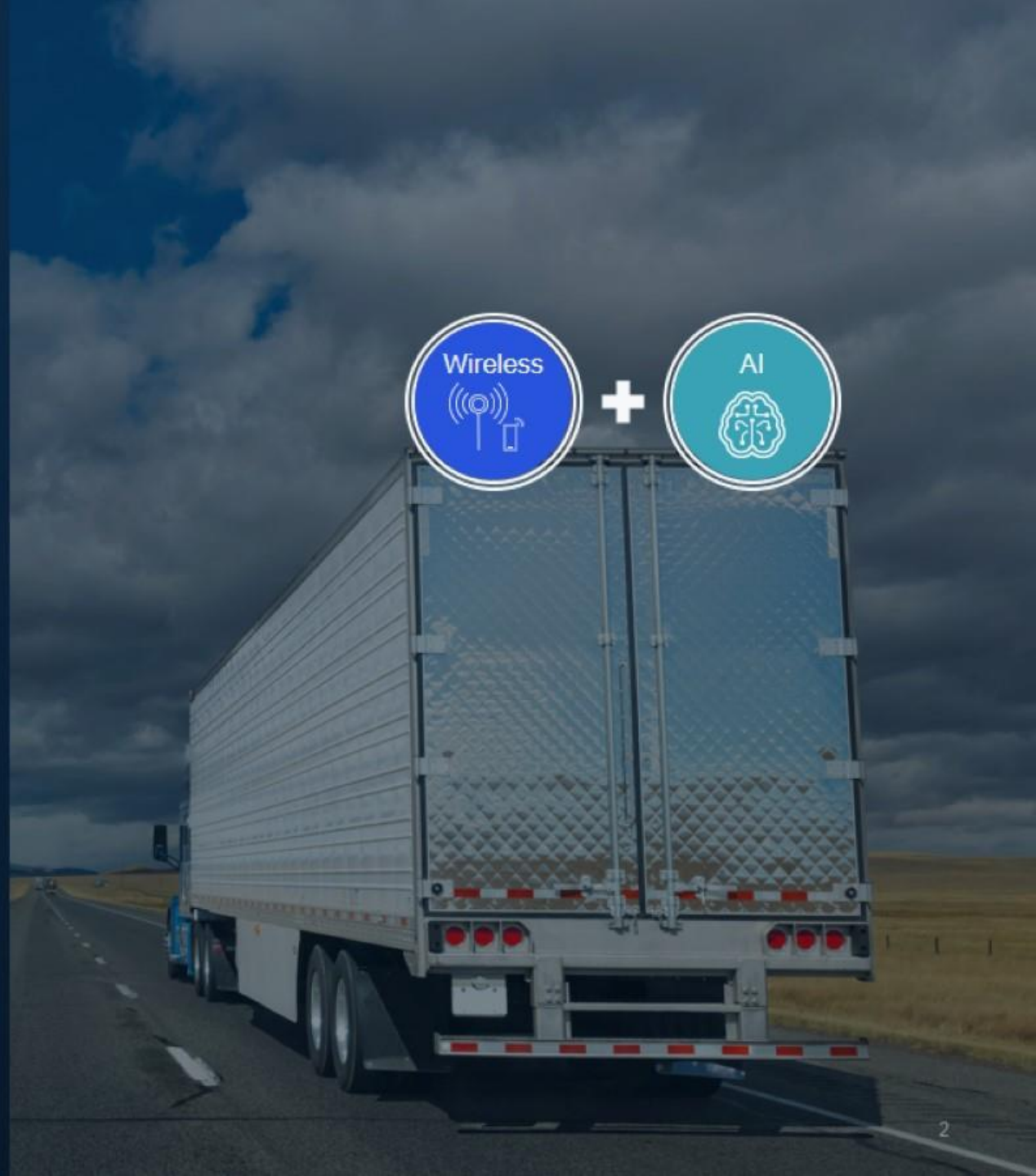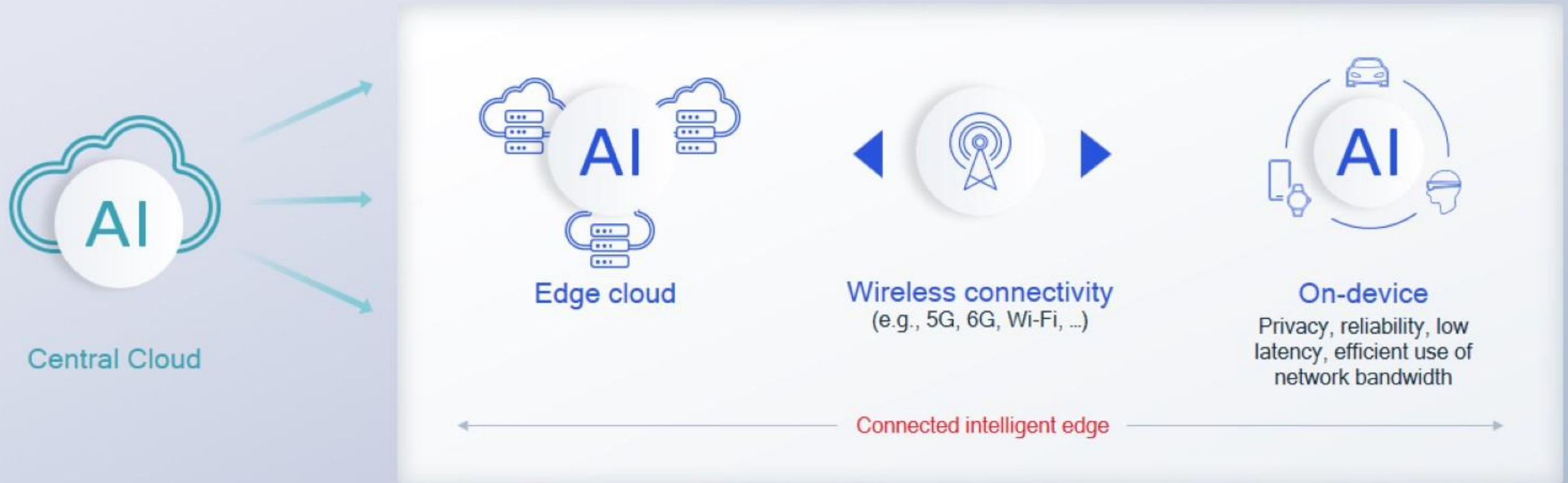# Today's agenda

1. AI enablement of 5G wireless today

2. AI-based air interface for 5G-Advanced

3. AI for 6G

# To scale efficiently, AI processing is expanding towards the edge

**AI** (Central Cloud)

**AI** — Edge cloud

**Wireless connectivity** (e.g., 5G, 6G, Wi-Fi, …)

**AI** — On-device
Privacy, reliability, low latency, efficient use of network bandwidth

Connected intelligent edge

**Qualcomm is leading the realization of the connected intelligent edge**

**CONVERGENCE OF:**

Wireless connectivity        Distributed AI

Efficient computing         Unleashing massive amount of data to fuel our digital future

# Snapdragon
## X80 5G Modem-RF

**Qualcomm®
5G AI Suite**
Gen 3

**AI-enhanced 5G
Advanced
user experience**

Multi-antenna management to improve user experience

Contextually-aware QoS and latency improvements

60%* faster CPE service acquisition (mmWave)

10%* lower power in connected mode (mmWave)

Location accuracy improvement by 30%*

Best-cell selection time reduced by 20%*

30%* faster link acquisition

**Snapdragon X Elite**

**PERFORMANCE REBORN**

## Best-in-class performance vs. x86

Up to
**2X**
faster CPU

Qualcomm Oryon™ CPU
12 high-performance cores
Dual-Core Boost

Up to
**2X**
faster GPU

Qualcomm® Adreno™ GPU
4.6 teraflops
Triple UHD monitor support

**4nm**
process node

**136GB/s**
memory bandwidth
LPDDR5x

## Generative AI powerhouse

Snapdragon smart

Capable to run
**13B+**
parameters on device

Generates
**30**
tokens/sec for 7B LLMs

**Built for AI**

**75** TOPS
Qualcomm® AI Engine

**45** TOPS
Qualcomm® Hexagon™ NPU

**<1 seconds**
Stable Diffusion

**4.5X**
faster AI NPU processing power than competitors

1st PC processor with integrated Always-Sensing ISP

Integrated Micro NPU on Qualcomm Sensing Hub

## Smart user experiences

Lightning-fast 5G | Wi-Fi 7
Immersive lossless audio
Advanced camera ISP
Snapdragon Seamless
Chip-to-cloud security

Snapdragon connect
Snapdragon sound
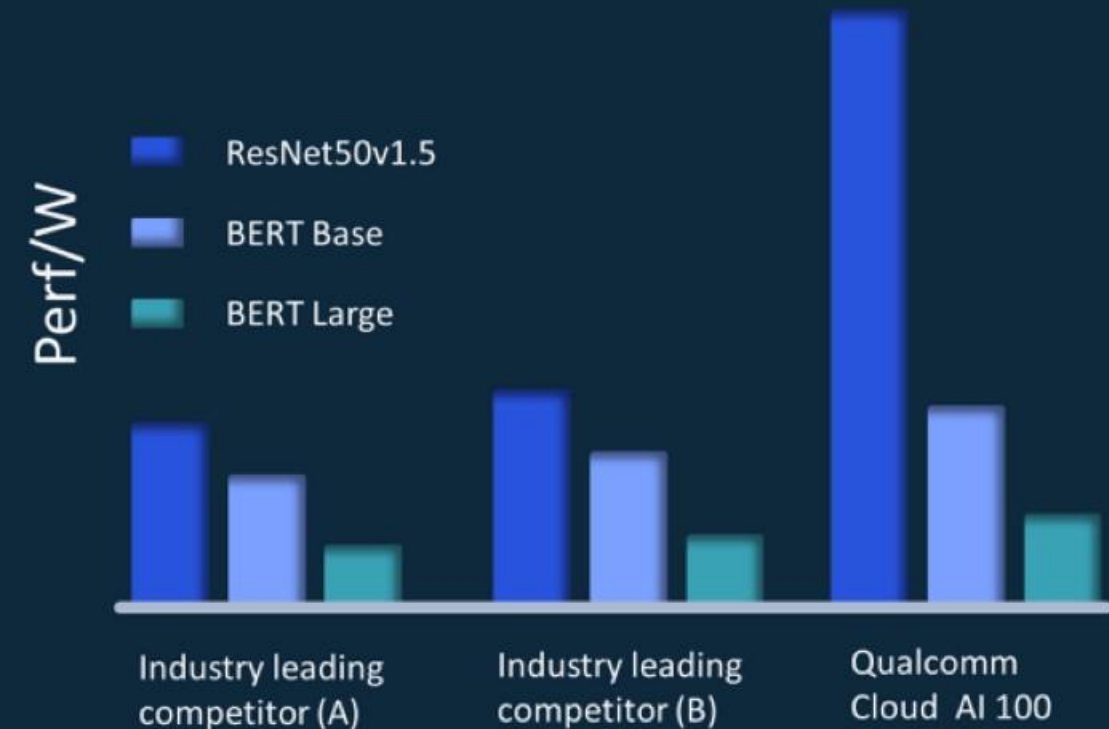Snapdragon secure
Snapdragon sight

## Leading PC performance per watt

Matches peak PC performance at
**68%**
less power vs. competition

## Scalable across a range of thermal designs and form factors

# Building high performance, power efficient, AI inference accelerator
# Qualcomm® Cloud AI 100

Perf/W

- ResNet50v1.5
- BERT Base
- BERT Large

Industry leading competitor (A)

Industry leading competitor (B)

Qualcomm Cloud AI 100

$Perf/TCO\$ = Inf/sec / (CapEx + OpEx)$

e.g., 4x Perf at same OpEx delivers up to 5x Perf/TCO$ benefit

## Qualcomm solution's 2-4x Performance advantage at similar power leads to 2-5x Perf/TCO$ advantage

# Qualcomm 5G RAN platforms



Virtual core

Qualcomm Edgewise™ Suite

Virtual core

Interoperable

Private 5G network

Private 5G network

Public 5G network

Qualcomm® 5G RAN Platforms

Qualcomm 5G RAN platform QRU100

Qualcomm 5G RAN platform QDU100

Qualcomm 5G RAN platform FSM200

Qualcomm 5G RAN platform FSM100

Qualcomm 5G RAN platform compact macro

Qualcomm® X100 5G RAN Accelerator Card

# 5G NR AI-based air interface design in Rel-18+

# Wireless AI

## 3 projects in Release 19

### Study on AI/ML for Next-Gen Radio Access Network[1]

New use cases including network slicing and coverage and capacity optimization (CCO)

Continued studies on mobility optimization for NR-DC, split architecture support, enhanced energy saving, continuous MDT, and multi-hop device trajectory

### Study on AI/ML to enhance 5G NR mobility[2]

Focusing on L3 device mobility, including RRM measurement & event prediction, device assistance information for network-side model, enhanced LCM, evaluation on testability, interoperability, impacts on RRM requirements and performance

1 RAN 3 led; 2 RAN 2 led; 3 RAN 1 led; 4 Continued study with corresponding checkpoints in RAN#105 (Sept '24)
Source: RP-234039 (AI/ML for NR air interface); RP-234054 (Study on AI/ML for NG-RAN); RP-234055 (Study on AI/ML for mobility in NR)

# Work on AI/ML Air Interface[3]

## General Wireless AI Framework

Support Life Cycle Management (LCM) of one-sided (i.e., device or network) AI/ML models

### Channel feedback[4]

Further study 2-sided CSI compression, 1-sided CSI prediction, model transfer/deliver, …

Improve user downlink throughput and reduce uplink overhead

### Beam management

Support device/network-sided beam prediction model in time/spatial domain

Reduce overhead, latency, and improve beam selection accuracy
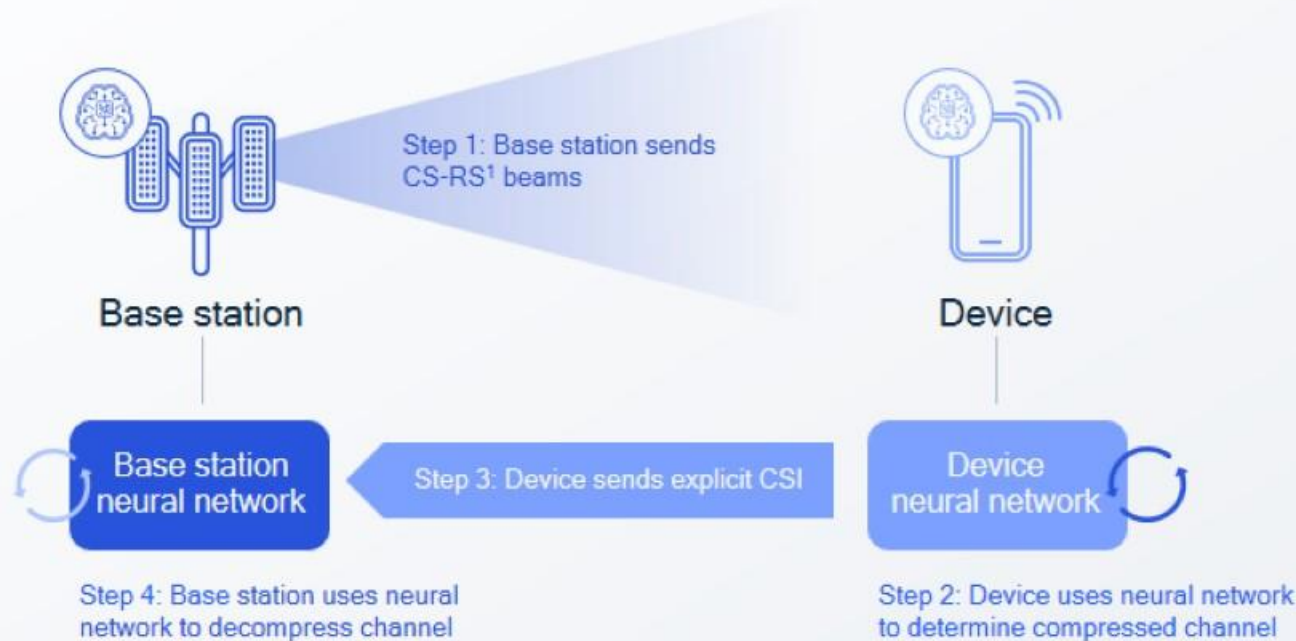
### Precise positioning

Support single-sided model for both AI-direct and AI-assisted positioning

Improve positioning accuracy for different indoor/outdoor scenarios

10

# Cross-node machine learning based channel state information

Explicit channel feedback framework for CSI compression and prediction utilizing domain knowledge and neural networks



Reconstructed DL channel estimates — CSI decoder — Data or control channel
Decoder at the base station

Data or control channel — CSI encoder — Downlink channel estimates
Encoder at the device

Step 1: Base station sends CS-RS[1] beams

Base station

Device

Base station neural network — Step 3: Device sends explicit CSI — Device neural network

Step 4: Base station uses neural network to decompress channel

Step 2: Device uses neural network to determine compressed channel

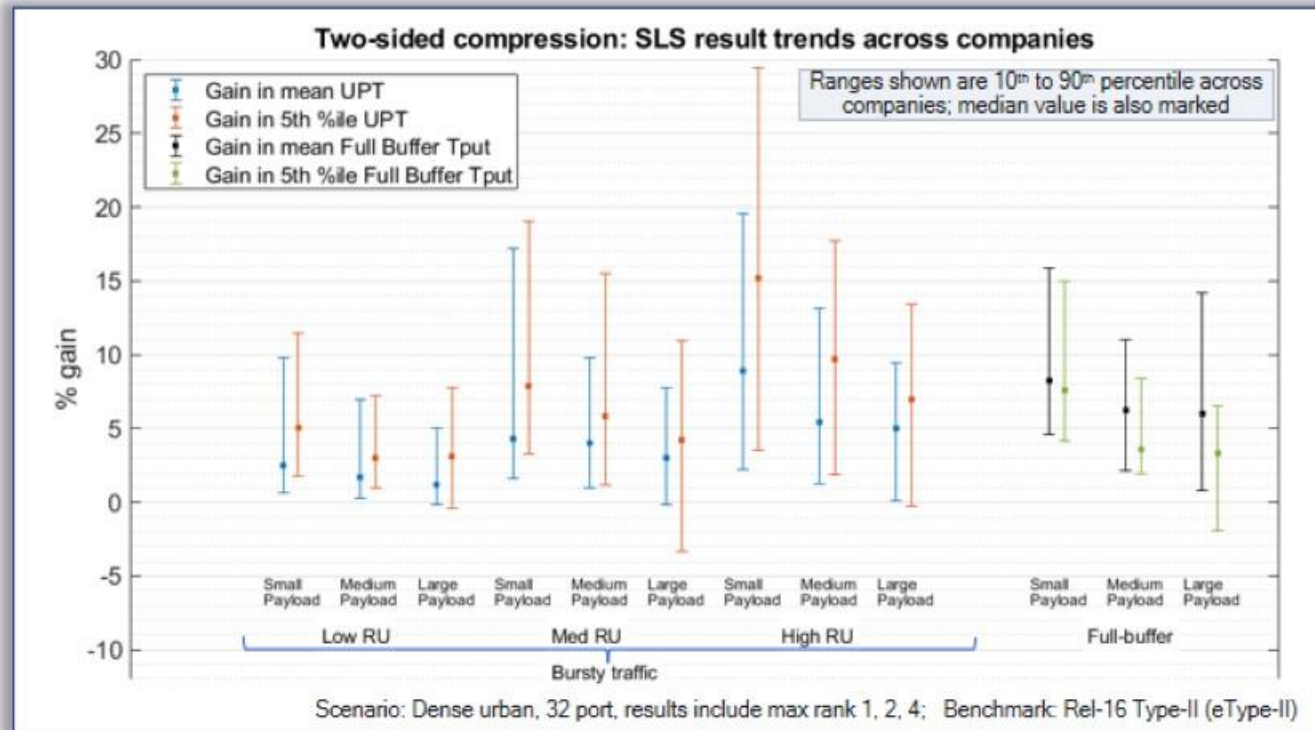Improve system efficiency with neural network framework for CSI on non-linear encoding and decoding

More effective multi-user multiplexing minimizing interference

Customized, lower overhead feedback based on individual device

# Performance gains from 3gpp spatio-frequency CSI compression study

System-level results: Mean and 5th %ile Throughput

Assumptions: 3gpp Dense Urban scenarios, 4GHz, 32 CSI-RS ports, 4 Rx, 20MHz comparing ML CSF and eType 2



**Ongoing and future directions**

- CSI prediction
- Spatial-frequency-temporal compression
- Joint CSI prediction and compression
- Hyperlocal models
- Joint CSI-RS optimization, feedback, and precoder
- Joint source channel coding on CSI
- Utilizing DL/UL reciprocity for CSI

- CSI feedback overhead reduction: 30-70% reduction in feedback overhead

- Throughput gains: up to 20-30% gain in median and tail user perceived throughput
  - Gains are more pronounced for smaller payload, higher resource utilization, and larger cells.

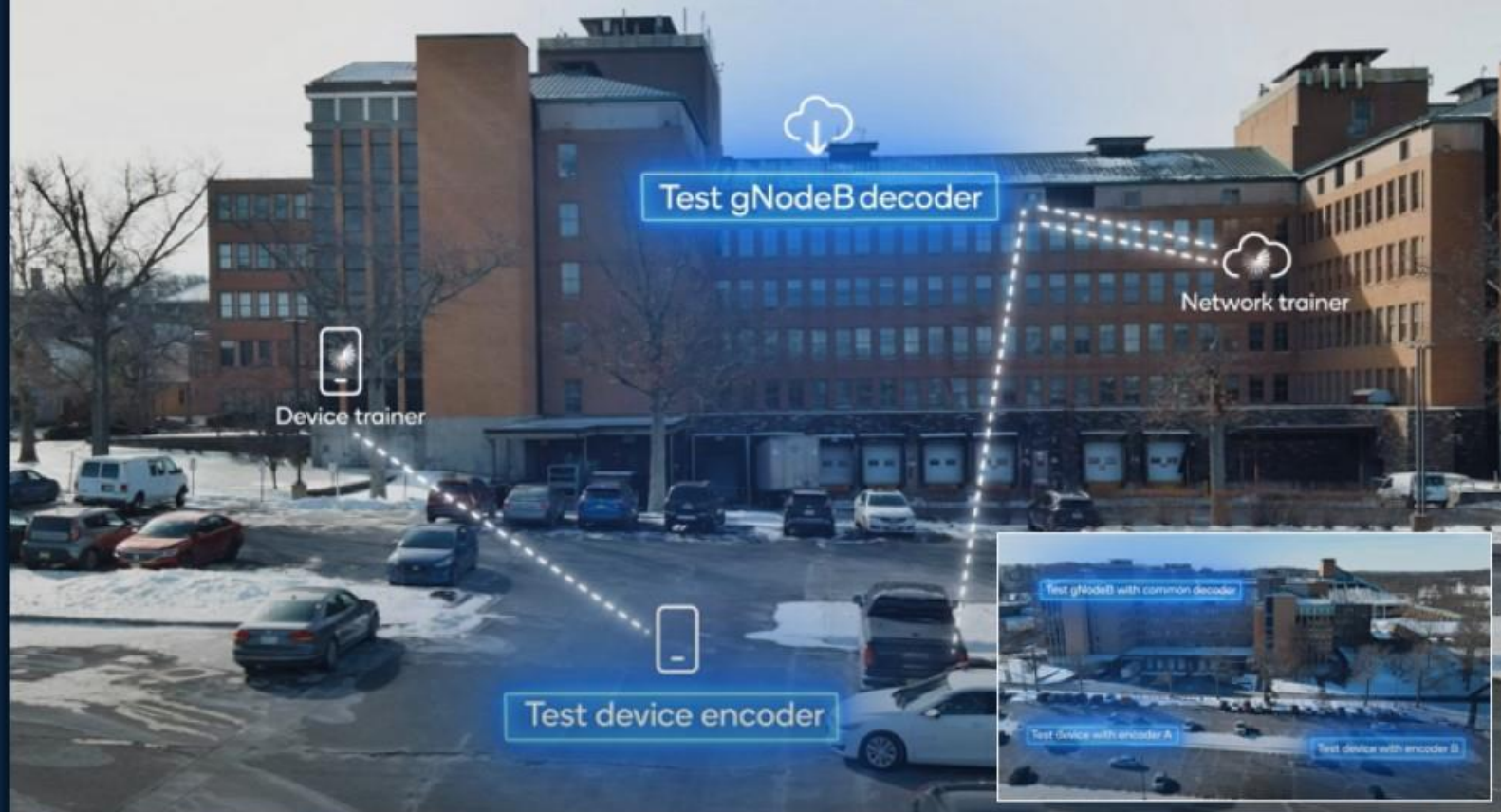# Wireless AI Interoperability
## For multi-vendor system

Close collaboration with Nokia Bell Labs on an over-the-air prototype of two-sided channel state information (CSI) feedback

Test network in Murray Hill, NJ, with Nokia infra and Snapdragon 5G Modem-RF

Sequential training enables data sharing but not neural network structures (models)

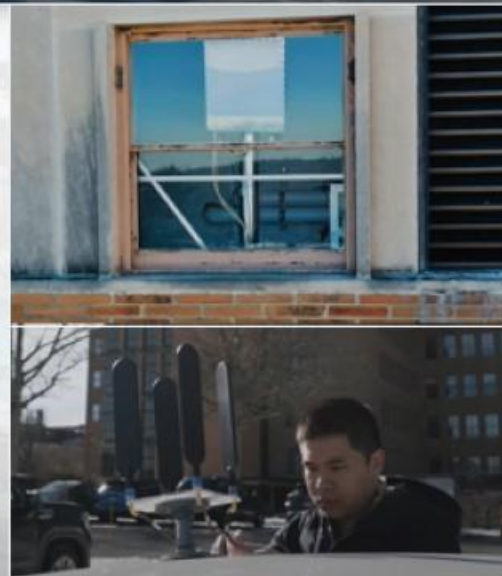3GPP global standards compliant, potentially a part of 5G Advanced Rel-19+

Qualcomm
MWCB 2024

Test gNodeB decoder

Network trainer

Device trainer

Test device encoder

Test gNodeB with common decoder

Test device with encoder A

Test device with encoder B

### Downlink throughput - device with encoders

Encoder A

Encoder B

Using a common decoder across all devices performs as well as utilizing a dedicated decoder trained individually for each device

gNodeB with common decoder
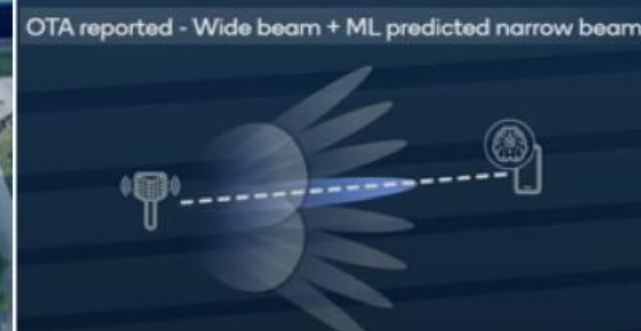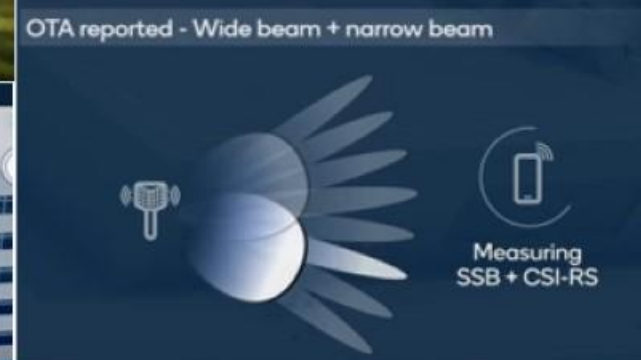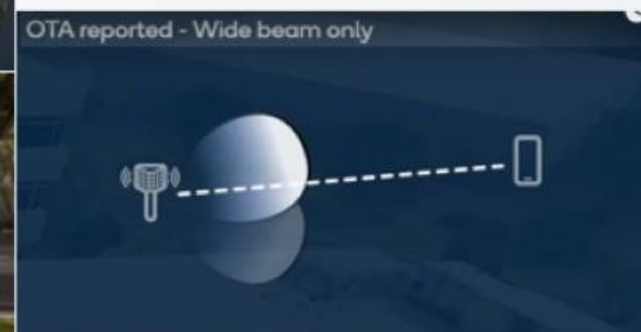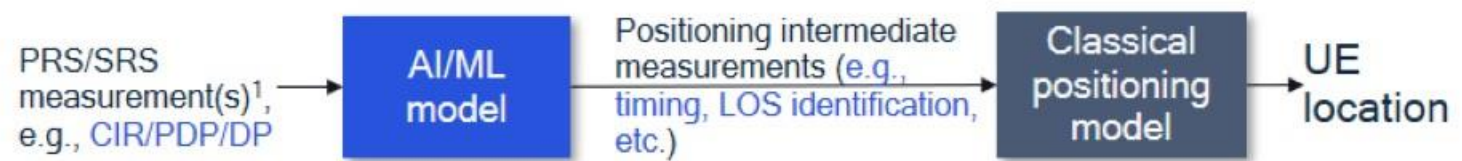gNodeB with dedicated decoder

# Performance gains from 3gpp positioning study

AI/ML learns multipath and enhances positioning in challenging NLOS scenarios, reducing positioning error from >10 meters to submeter level

## Direct AI/ML positioning

PRS/SRS measurement(s), e.g., CIR/PDP/DP [1] → AI/ML model → UE location

## AI/ML assisted positioning

PRS/SRS measurement(s)[1], e.g., CIR/PDP/DP → AI/ML model → Positioning intermediate measurements (e.g., timing, LOS identification, etc.) → Classical positioning model → UE location
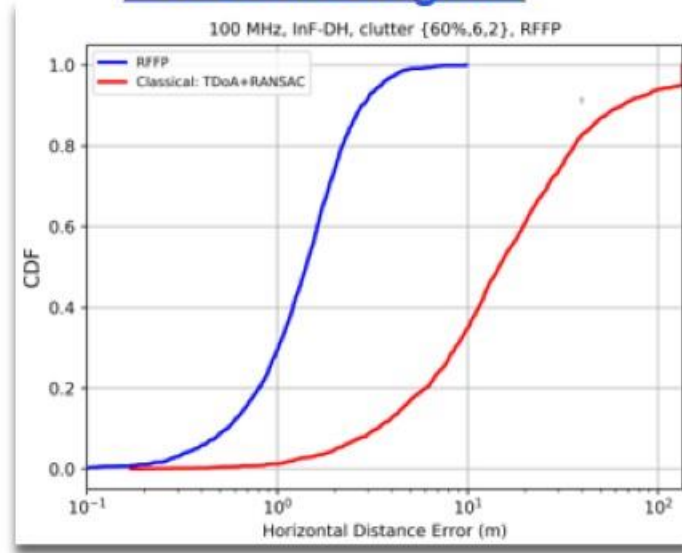
## Performance gains



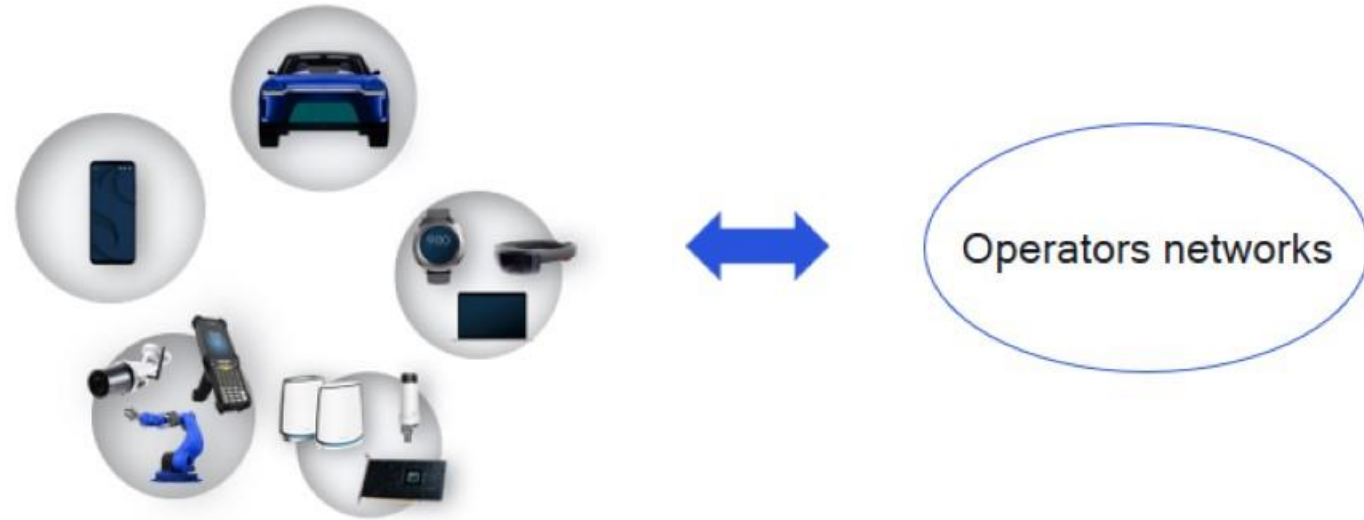RFFP: direct AI/ML positioning | Extreme NLOS condition (LOS < 1%) | 3GPP InF-DH scenario| Classical: Time difference of arrival with outlier rejection using RANSAC algorithm

## Deployment cases



Note[1] : CIR: Channel impulse response | PDP: Power delay profile | DP: Delay profile

# Fundamental 6G motivation

Operators networks

- Increasing revenue with new use cases/services

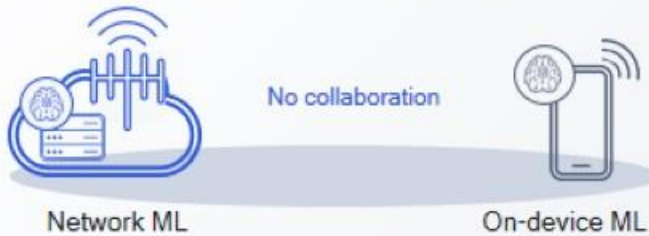- Reducing network TCO (total cost of ownership)

# Evolving towards an AI-native wireless system

## Multiple wireless AI/ML training and inference scenarios



### Overlay AI/ML
**INDEPENDENTLY AT THE DEVICE OR NETWORK**

No collaboration

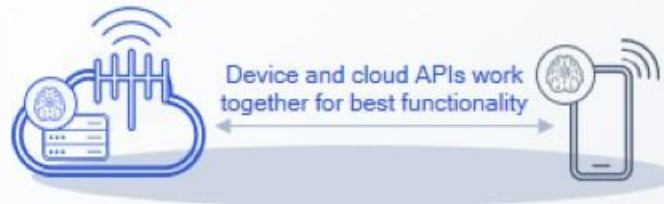Network ML                    On-device ML

ML operates independently at the device and network as an optimization of existing functions

Proprietary ML procedures including model development and management

Proprietary and standardized data collection used as input to training

### Cross-node AI/ML
**COORDINATED BETWEEN DEVICE AND NETWORK**

Device and cloud APIs work together for best functionality

ML operates in a coordinated manner between the device and network

Proprietary and standardized ML procedures including model development and management

Further data collection used as input to training as well as monitoring

### Native AI/ML
**AT ALL DEVICE AND NETWORK LAYERS**

Device and network exchange control/input across all layers

ML operates autonomously between the device and network across all protocols and layers

Integrated ML procedures across to train performance and adapt to different environments

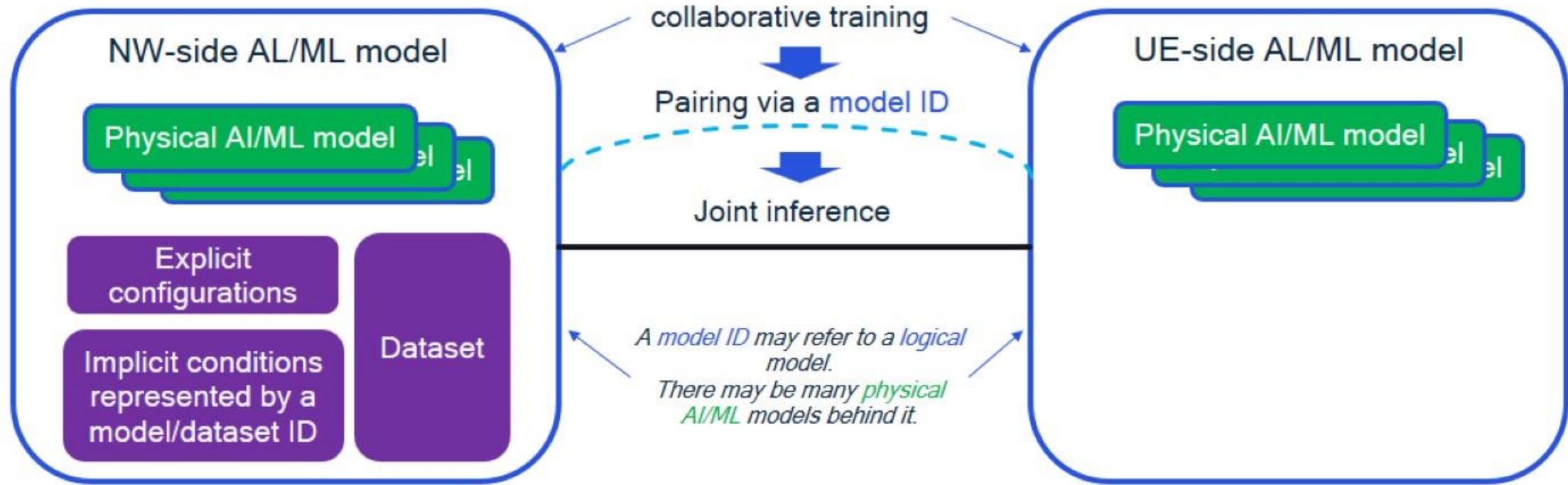Data fusion for integrated dynamic ML lifecycle management

5G    5G ADVANCED    6G

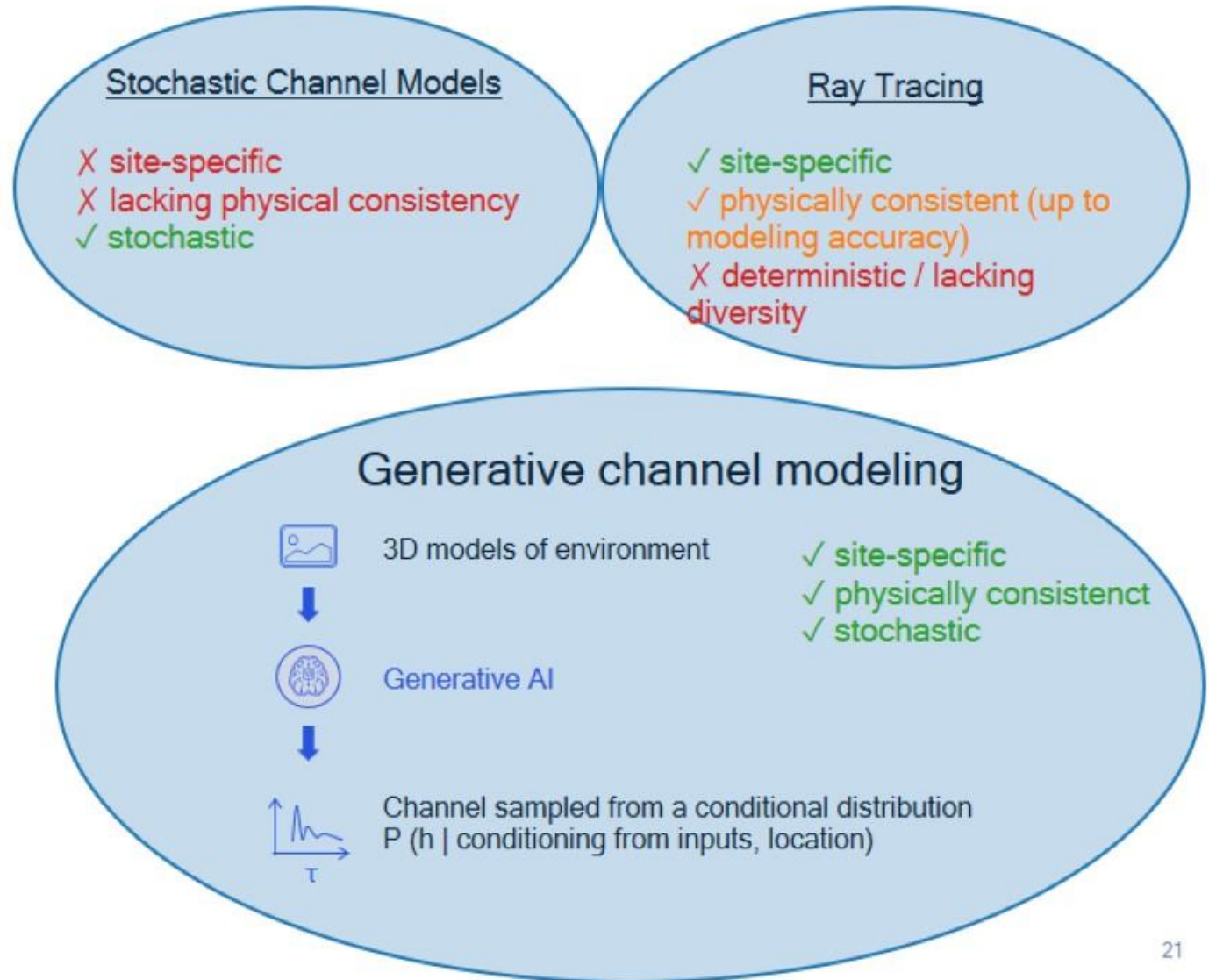# AI-native air interface open issues – Model pairing

# AI-native air interface open issues – channel models

- **Conventional communication algorithms have been designed and evaluated based on synthetic channels following stochastic modeling assumptions.**

- **AI-native air interface design necessitates new channel modeling framework.**
  - AI/ML models often achieve higher gain when optimized to a given site-specific propagation environment.
  - AI/ML models may not work well in real-world channels when trained on synthetic channels.
  - Channels from Ray Tracing is too deterministic and easily lead to AI/ML model overfitting.

- **Data-driven approach may be used toward new channel modeling.**

### Stochastic Channel Models

X site-specific
X lacking physical consistency
✓ stochastic

### Ray Tracing

✓ site-specific
✓ physically consistent (up to modeling accuracy)
X deterministic / lacking diversity

### Generative channel modeling

3D models of environment

Generative AI

Channel sampled from a conditional distribution
P (h | conditioning from inputs, location)

✓ site-specific
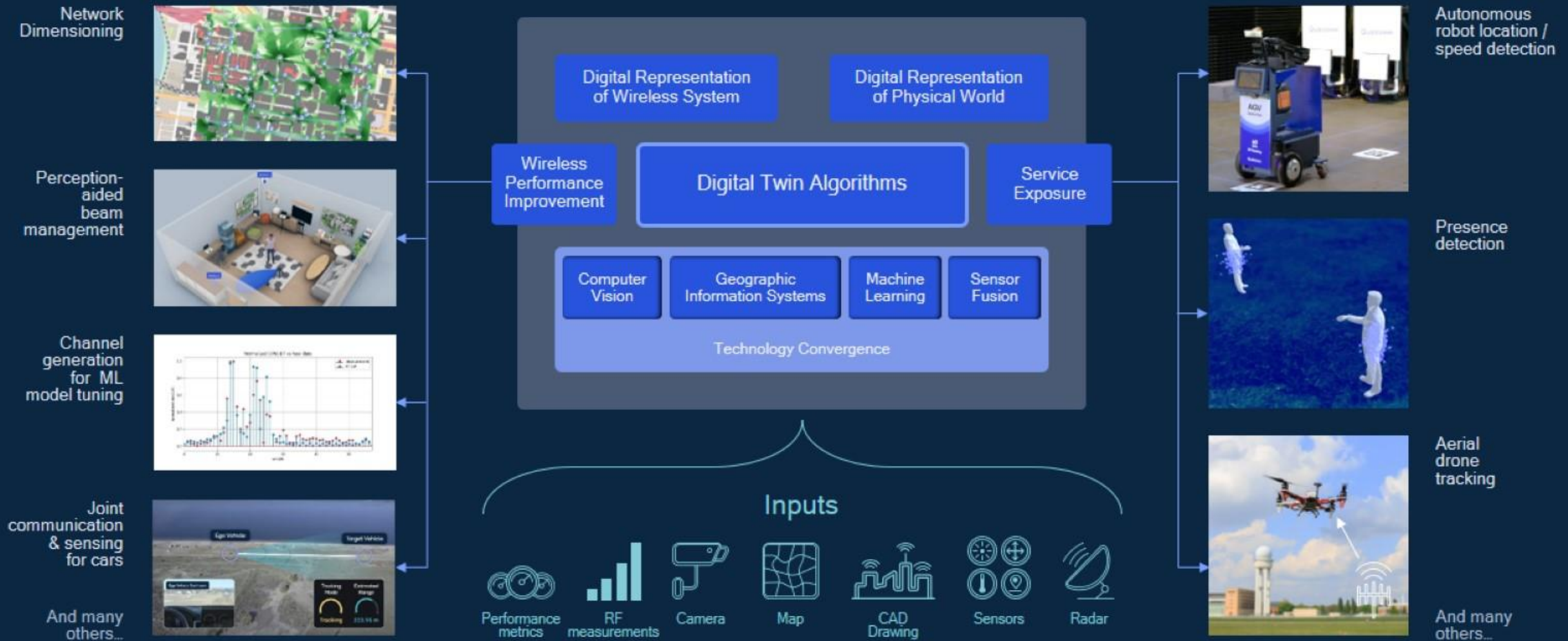✓ physically consistenct
✓ stochastic

# AI-native air interface open issues – Life cycle management



Source: TR 38.843 v18.0.0, "Study on Artificial Intelligence (AI)/Machine Learning (ML) for NR air interface" (Release 18), December 2023

- Backend ML infra for data collection and model Life Cycle Management (LCM) are important.

- What and how much to specify?
  - Proprietary vs. standardized data collection
  - Proprietary vs. standardized training procedure
  - Proprietary vs. standardized model performance monitoring

- Proprietary innovation vs. standardized inter-operability

# Digital Twin - virtual representation of physical system

Monitoring and performance optimization of its real-world counterpart

# Vision for gen AI-augmented and autonomous networks

NETWORK DIGITAL TWIN

OPERATOR NETWORK

Has there been any documentation of cell malfunctions today or in previous days?

To my knowledge, there are no reports of malfunctions for the designated cell.

Have there been any severe weather conditions in the vicinity or any tower construction near the cell?

There have been no reports of extreme weather and no recent construction activity near the cell.

Please suggest actions that can be done to rectify situation

Based on similar scenarios I suggest an azimuth tilt of 3-8 degrees to the right and 1-3 up.

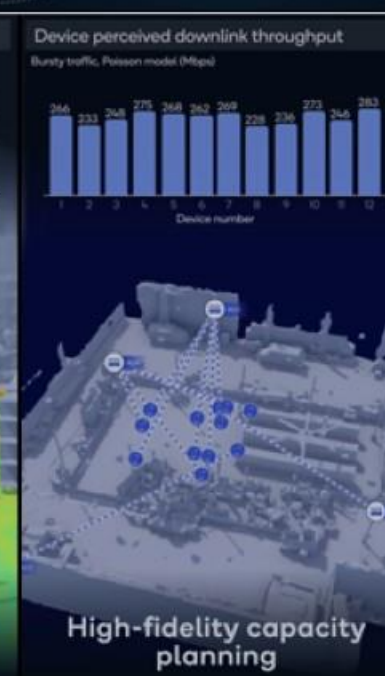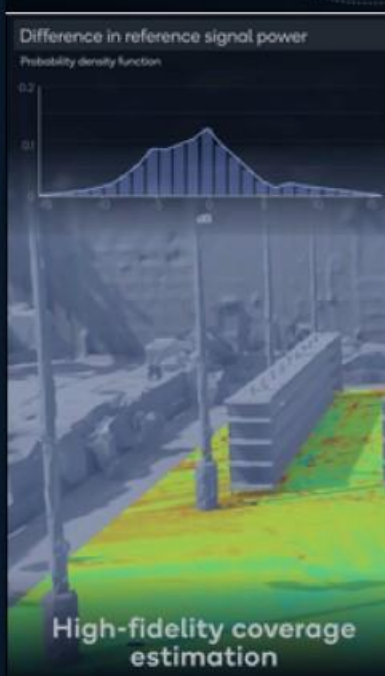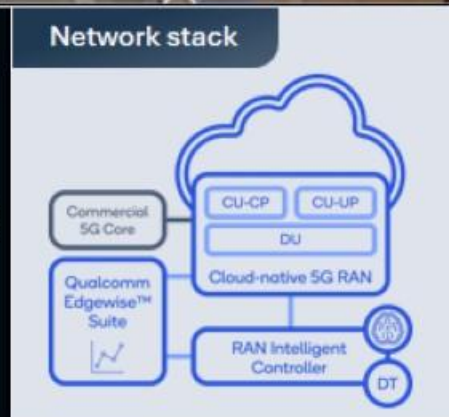| Intelligent monitoring and management | On-the-fly modeling | Proactive alerts | Programming AI-assistants | 'Level-3' autonomous networks |

# Digital Twin Network (DTN)

**Converging expertise in wireless system modeling, computer vision and AI to create the high-fidelity DTN**

**Generating synthetic data to address the data collection challenge from real world deployments**

**Sophisticated dynamic modeling of 5G RAN infrastructure**

**Over-the-air testbed operating in the 3.35 GHz band with cloud-native 5G RAN, RAN Intelligent Controller (RIC) and the Qualcomm Edgewise™ Suite**

# Applying generative modeling to improve wireless communications system design

## Wide applicability for Generative Modeling

Real-time use cases for air interface

Propagation channel
Beam management
Interference prediction

Scheduler optimization
Traffic source
Mobility enhancement

Link / system simulation

Deployment optimization

Positioning and sensing

Network and device optimization

Others...

## Application examples

### Channel rendering

Text description of image or semantic map

↓

Diffusion model
(To generate channel information)

↓

Channel sampled from a conditional distribution $P(h \mid$ conditioning from inputs, location)

### Network / device prediction

Context in text, e.g., history of device reports and base station responses

↓

Large language model
To learn link, beam, protocol languages

↓

Next action for base station and/or device, sampled from a conditional distribution $P$ (next action | conditioning from inputs)

## Our on-going wireless generative AI research areas

3D mapping and material learning

Foundation models (e.g., link and protocol level use cases, beam prediction, and others)

Neural channel rendering (e.g., map-based, ray tracer augmented, site-specific, and others)

Customized ML-based stochastic channel

Neural surrogate for base station scheduler and applications traffic

And others...

# University Collaborations

Qualcomm Innovators Development Kit AI/ML and Compute

Snapdragon Spaces™ XR Developer Platform

Robotics



- Develop AI/ML and compute applications
- Snapdragon® 8 Gen 2 SoC with AI HW/SW

- Develop immersive AR experiences
- Lenovo Think Reality A3 Glass and Motorola edge+ smartphone kit

- Develop power-efficient robots
- Qualcomm® RB5 Development Kit with Qualcomm® QRB5165 processor (15 TOPS AI, 7 cameras, VSP)
- 5G mezzanine card accessory

Expanding to other platforms in the future

University Courses

Hackathons

Research projects

Private-Public projects

Platform Training

Onboarding Assistance

27

# Thank you

**Qualcomm**

Follow us on: in 𝕏 ⊙ ▶ 𝗳

For more information, visit us at:

qualcomm.com & qualcomm.com/blog

References in this presentation to "Qualcomm" may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes our licensing business, QTL, and the vast majority of our patent portfolio. Qualcomm Technologies, Inc., a subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of our engineering, research and development functions, and substantially all of our products and services businesses, including our QCT semiconductor business.

Snapdragon and Qualcomm branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries. Qualcomm patented technologies are licensed by Qualcomm Incorporated.