

# 6G as Cellular Network 2.0

## - A Networked Computing Perspective -

Mobile Korea 2023 – 6G Global 2023

2023/11/01

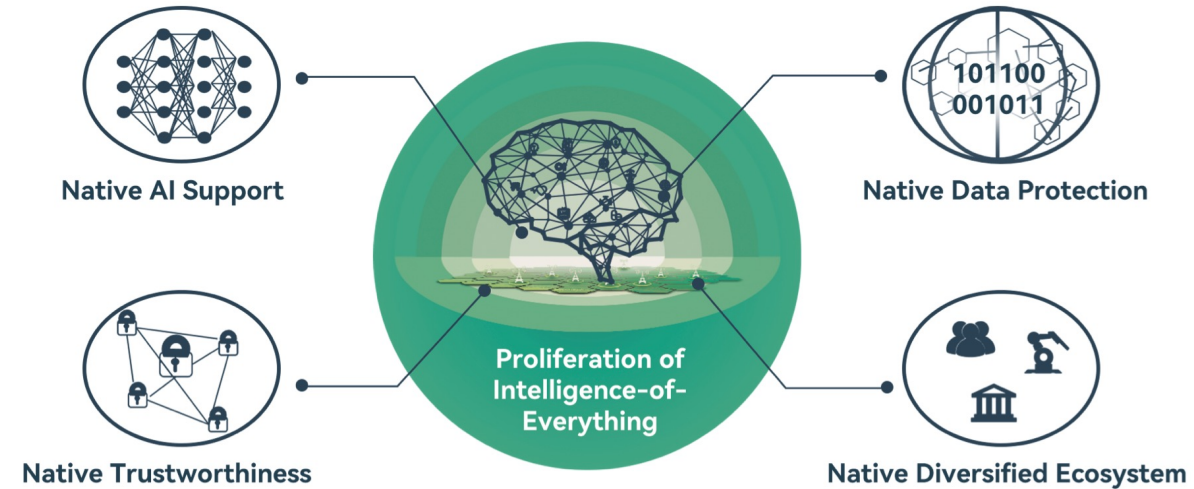
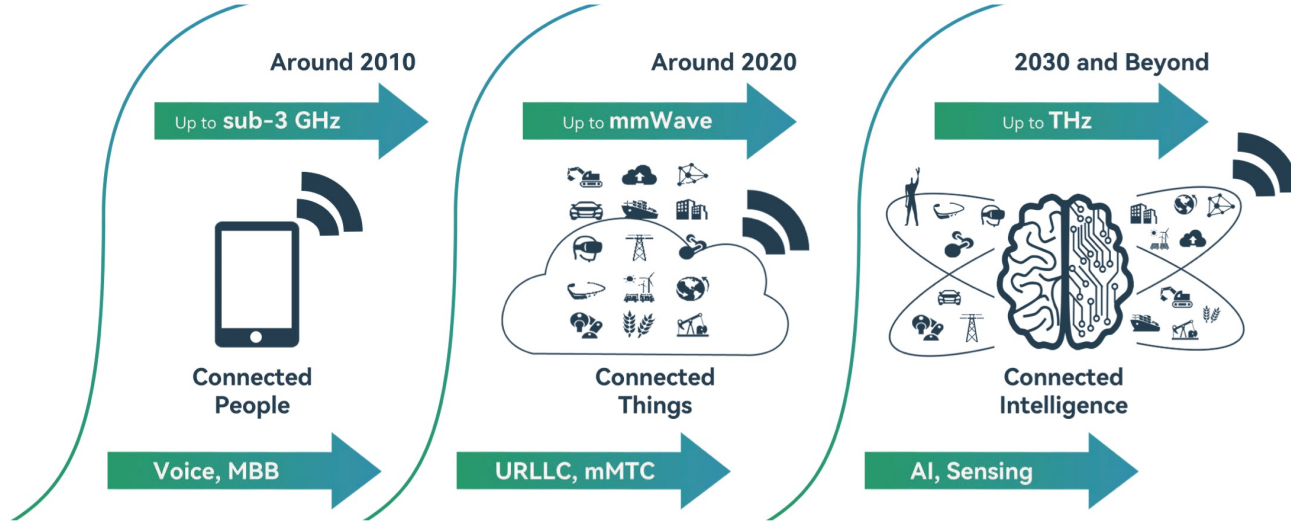
Kyunghan Lee

[kyunghanlee@snu.ac.kr](mailto:kyunghanlee@snu.ac.kr)

Networked Computing Lab.  
Electrical and Computer Engineering  
Seoul National University



# 6G? An Ordinary Perspective

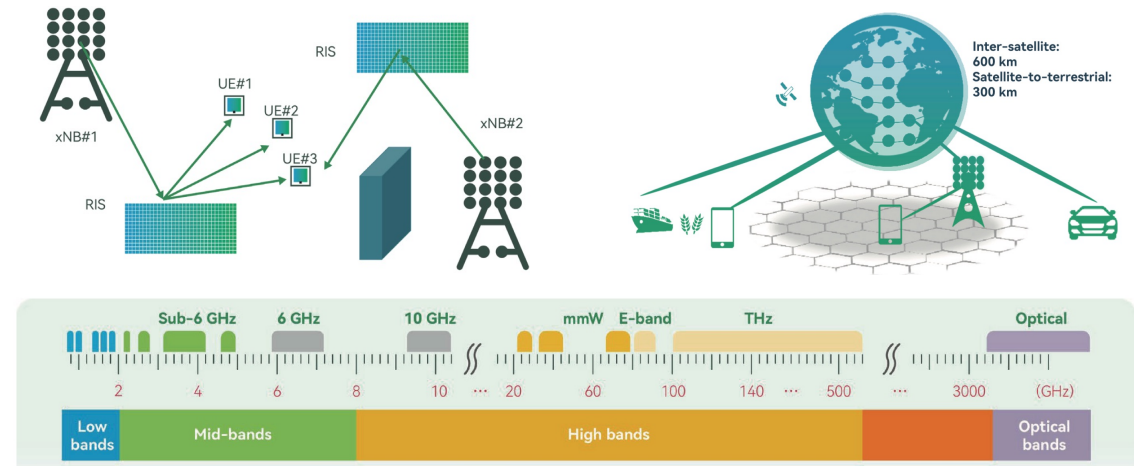


## 3/4/5G

Service	Connectivity Only
Networking	Public with Extended Private
Security	Encryption-based Security
Algorithm	Analytic Only
O&M	Automated OA&M
Business	Networking Infrastructure
Coverage	Terrestrial Only

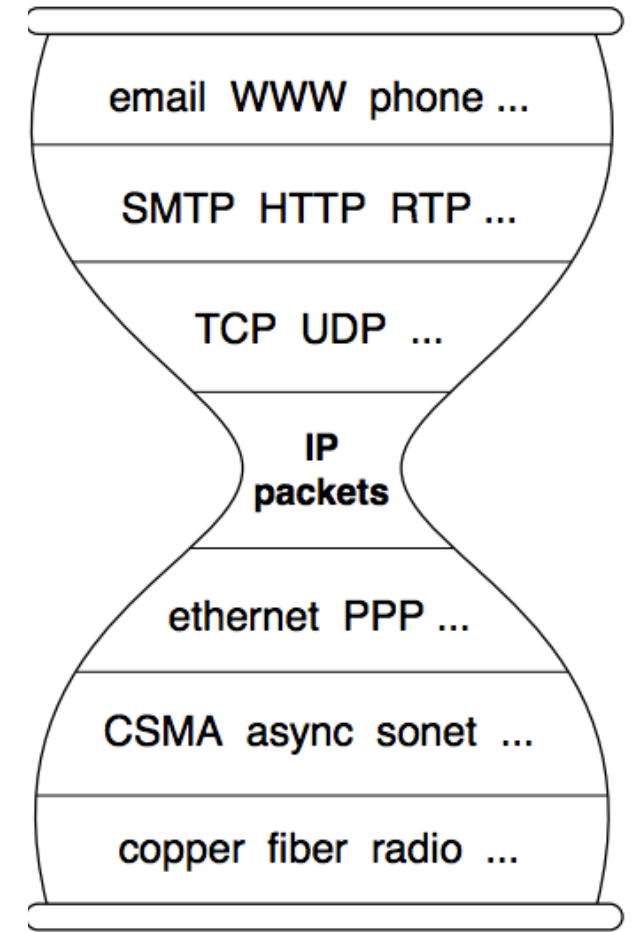
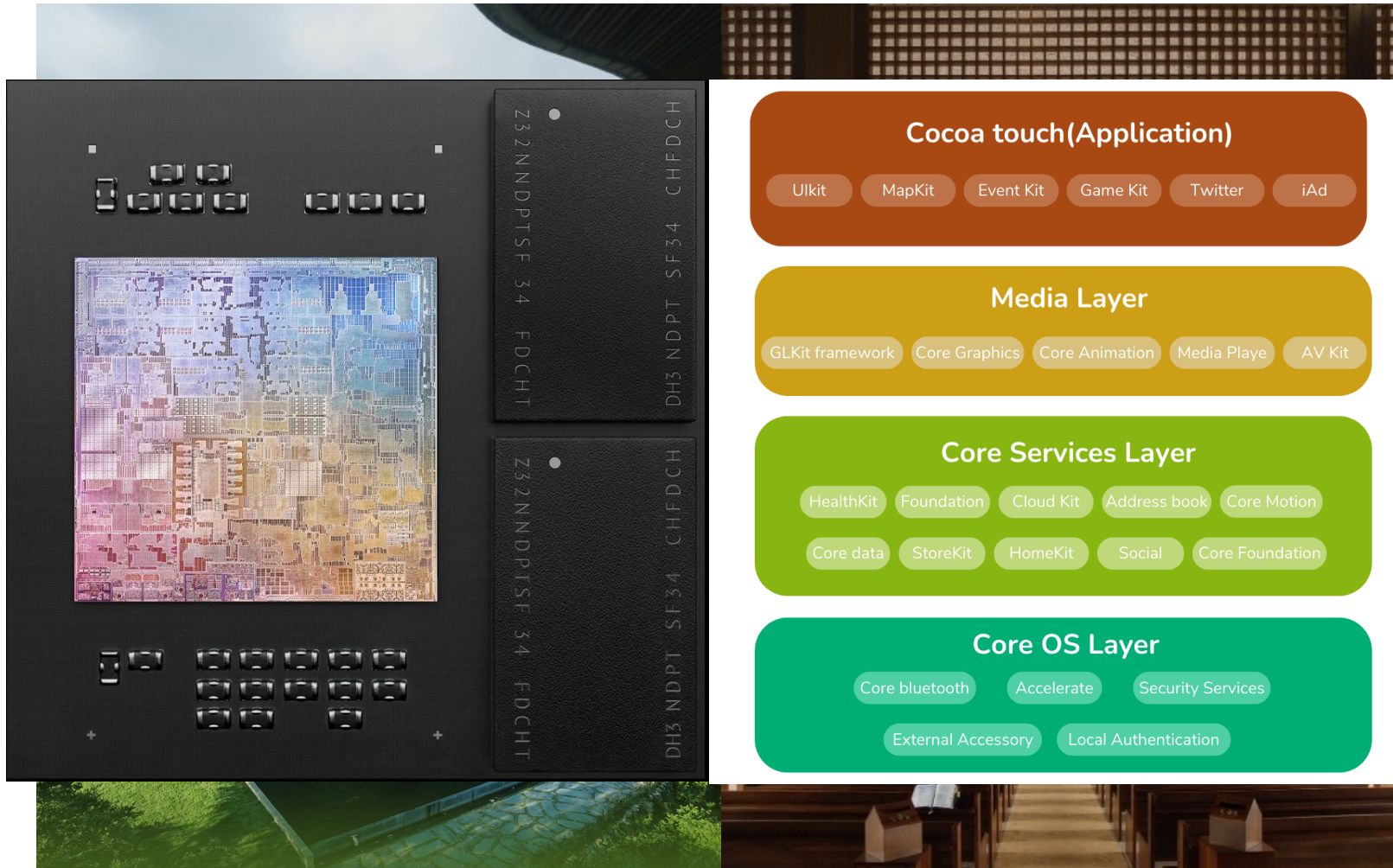
## 6G

Connectivity and Sensing, AI as a Service
Public Native and Private Native
Technology-based Trustworthiness
Analytic + Data (AI)
Level 5 Native OA&M
Networking & Computing Infrastructure
Integrated Terrestrial and Non-Terrestrial



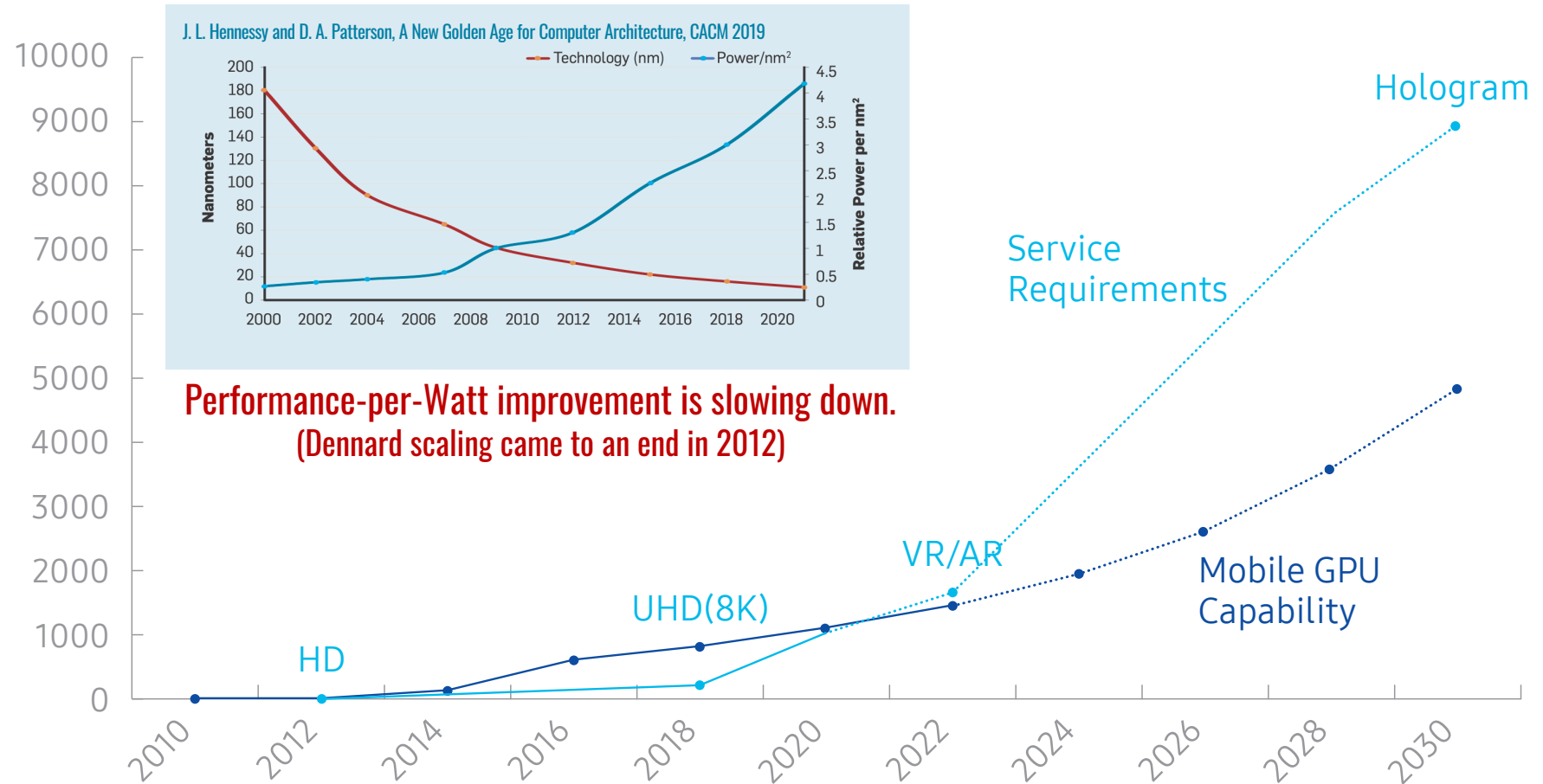
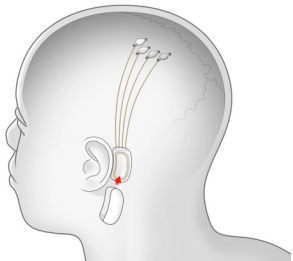
# Architecture

**Architecture** is the art and technique of designing and building, as distinguished from the skills associated with construction. It is both the process and the product of sketching, conceiving, planning, designing, and constructing buildings or structures.



# Computing over Networks (Connected Computing) is the Inevitable Future

Mobile/Wearable/Implantable Devices need **ExoComputing Capability** for Next-Generation Services





# Realtime Connected Computing is the Key to Human Augmentation

How Can we Enable **Realtime ExoComputing (Connected Computing)** over Wireless Networks?





# Realistic Metaverse is an Example of Realtime ExoComputing



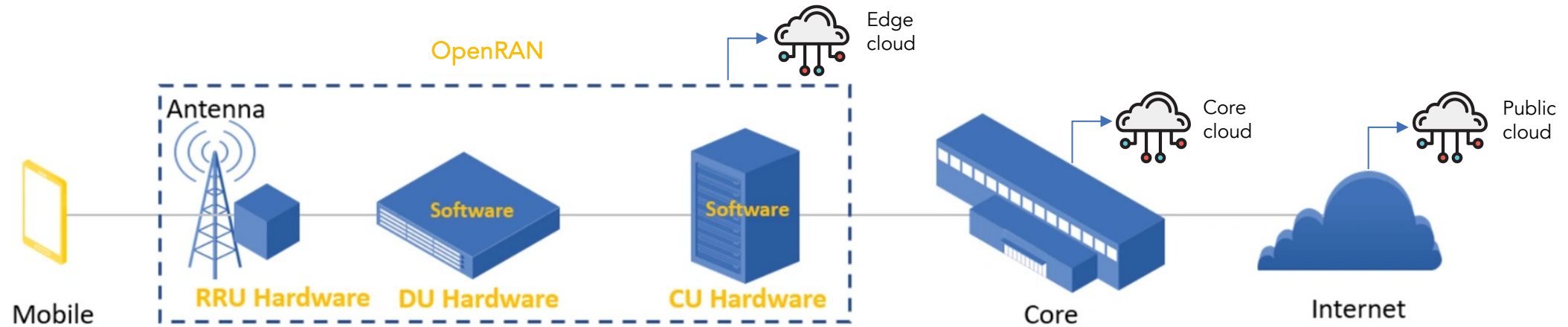


# Offloaded AI Analytics is another Example of Realtime ExoComputing

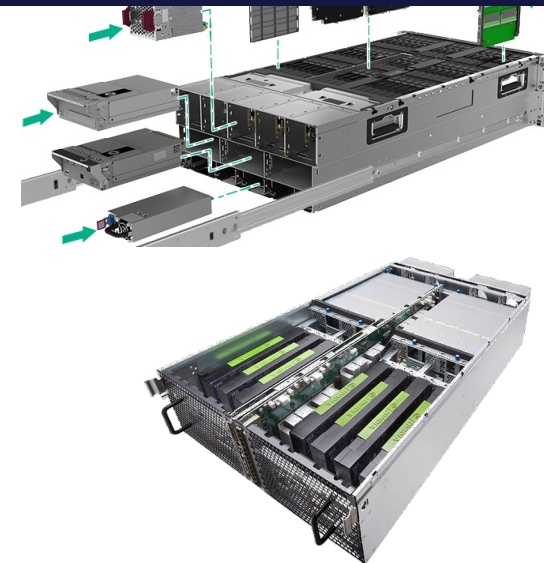
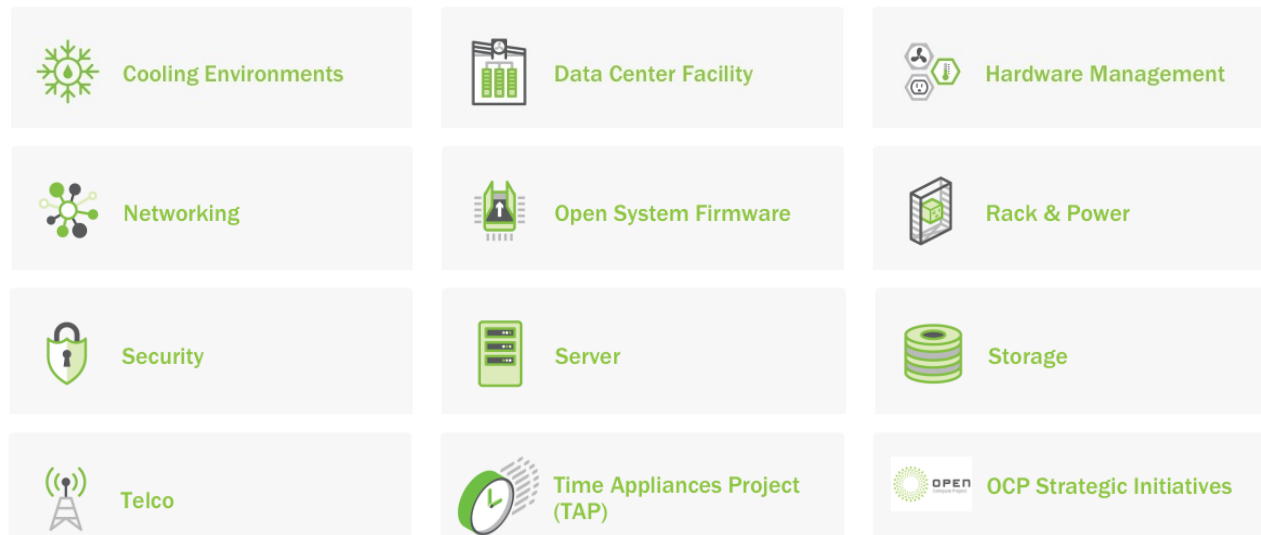
ETP (event to photon) Latency < 20ms



# Can Cloud Computing over 5G Networks Enable Realtime ExoComputing?

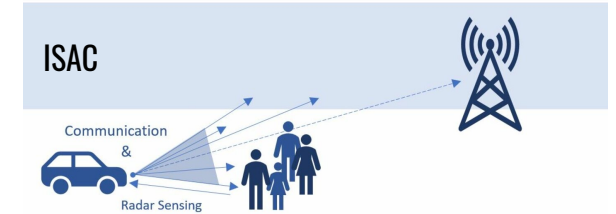
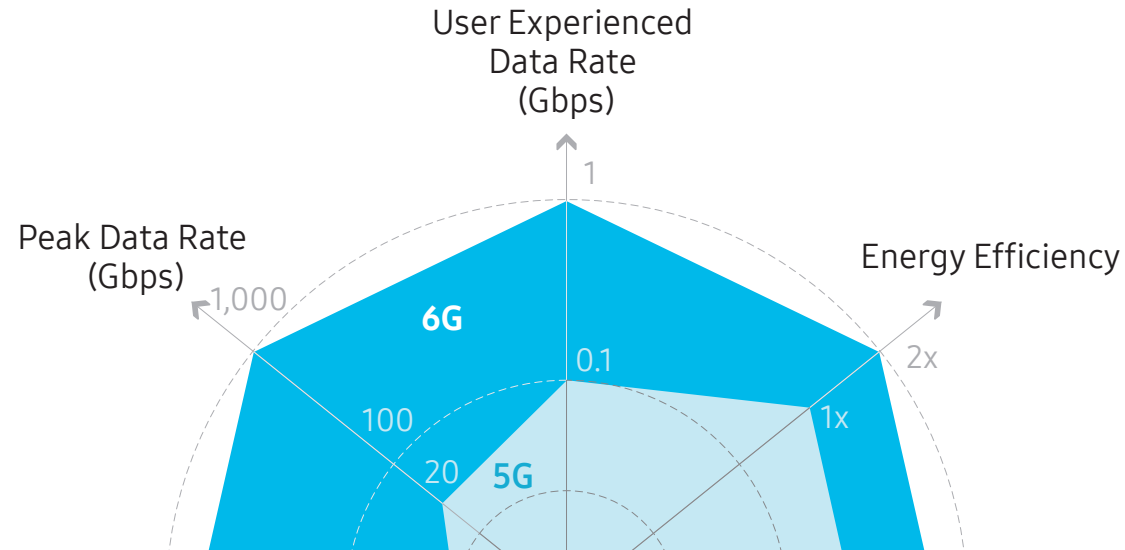


**No. 5G Network is Too Unstable and Unpredictable.**

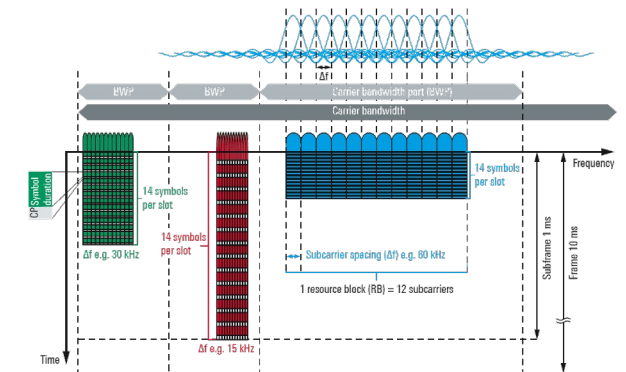
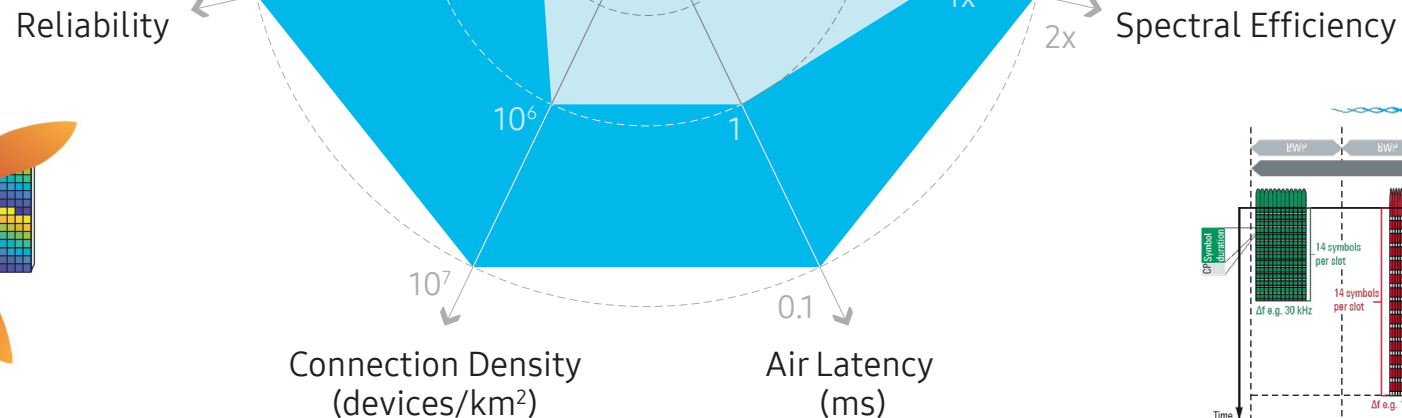
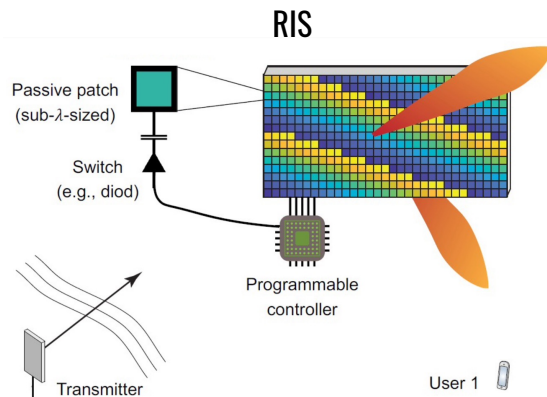




# Would Deploying 6G Enable Realtime ExoComputing?

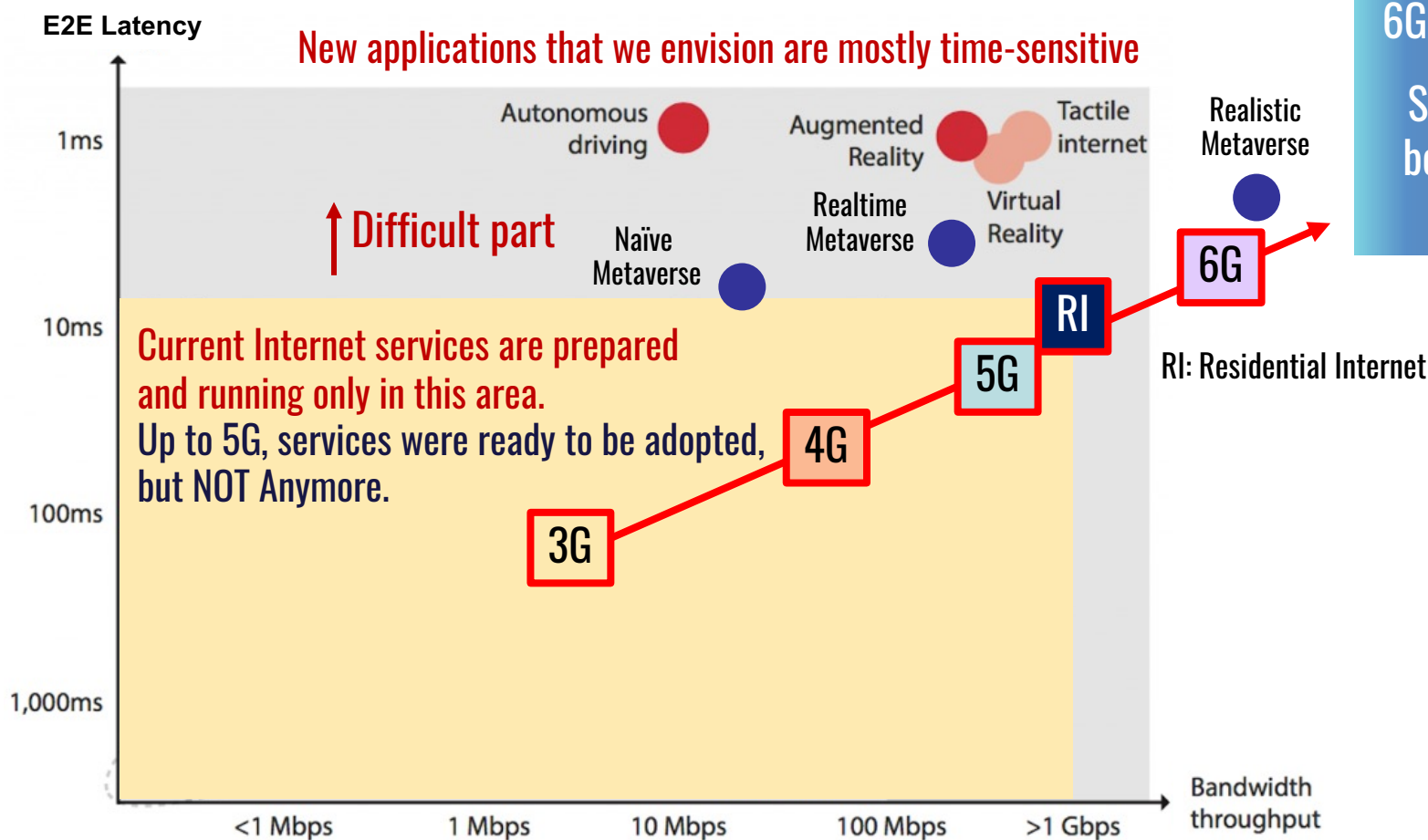
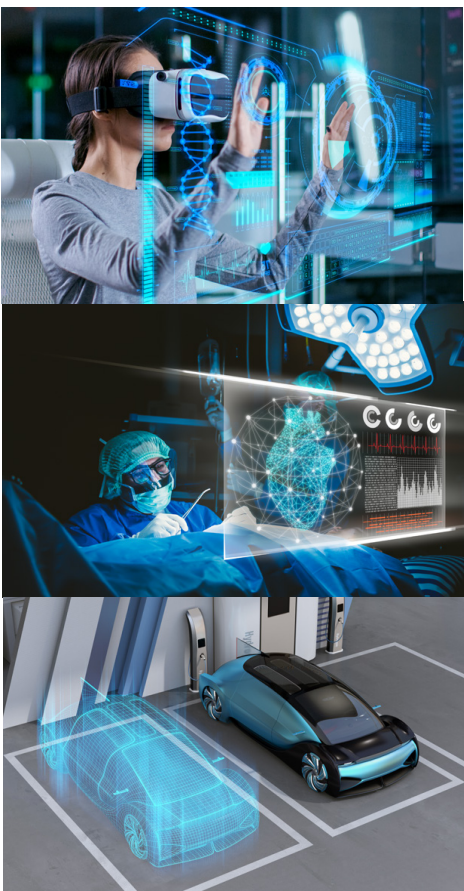


Unfortunately, Still No. Why?



Samsung Research, "6G The Next Hyper Connected Experience for All," white paper, 2020

# What is Wrong with 6G Cellular Network System?

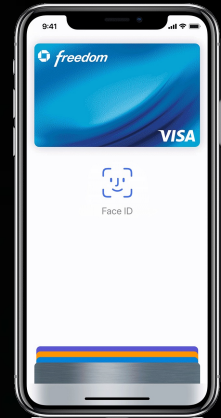


6G with MEC will surpass RI.

So, NEW Services will not be ready for 6G unless we prepare our own plan.

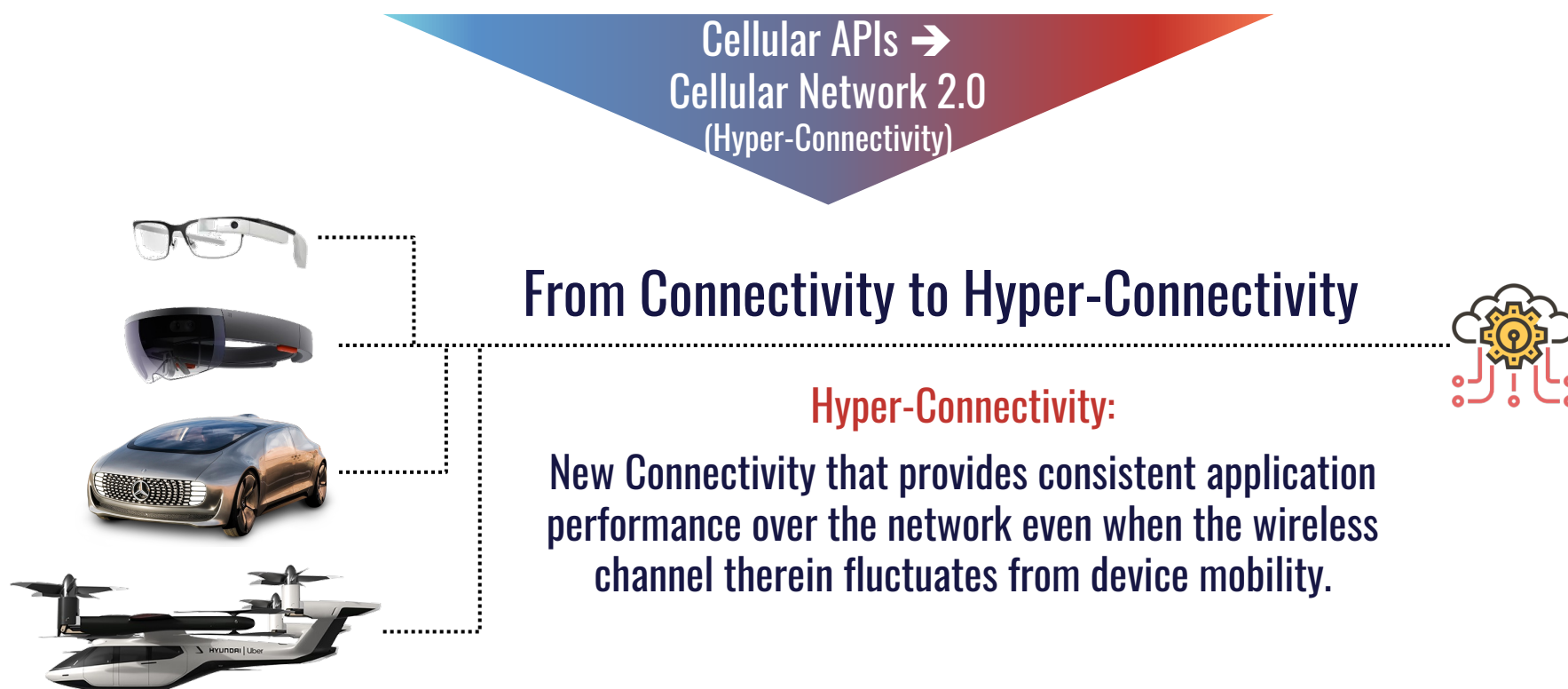
# Low-level Evolution $\neq$ Perceivable Performance Improvement

**Full potential** of low-level evolutions can only be delivered to users via **sophisticated software**.  
It's time to think about developing **Cellular APIs** for Realtime ExoComputing Services like Realistic Metaverse.



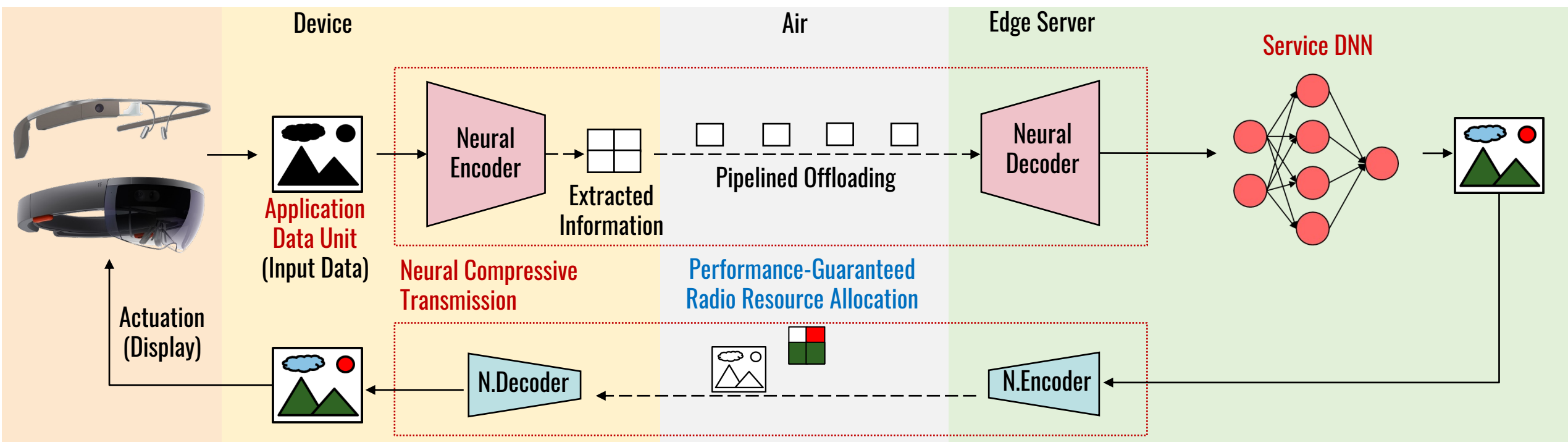
# Realtime ExoComputing needs Cellular Network 2.0, not 6G

Cellular APIs can make the **Network Application Performance Indistinguishable** from the Performance of **Applications Running on Device**.





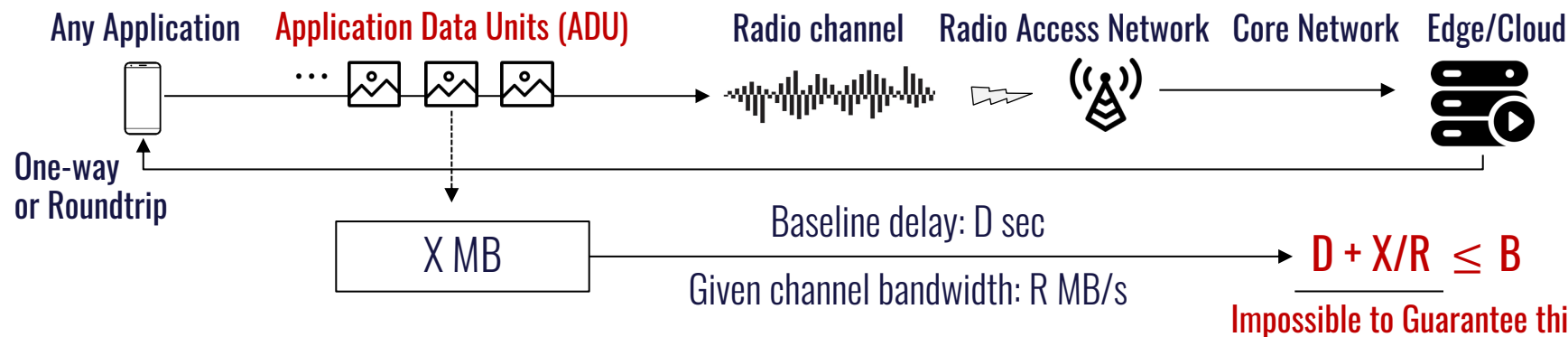
# A Reference ExoComputing Service Model with Cellular Network 2.0



$$ADU.CT = T_{wait} = t_{en/decoding} + t_{tx\ up} + t_{computing} + t_{tx\ down} + t_{en/decoding} \leq T_{budget}$$

Networked Computing can Pipeline and Hide most of these Delay Components!

# Enabling Cellular API (Hyper-Connectivity) $\approx$ ADU Completion Guarantee



## Conventional Approach

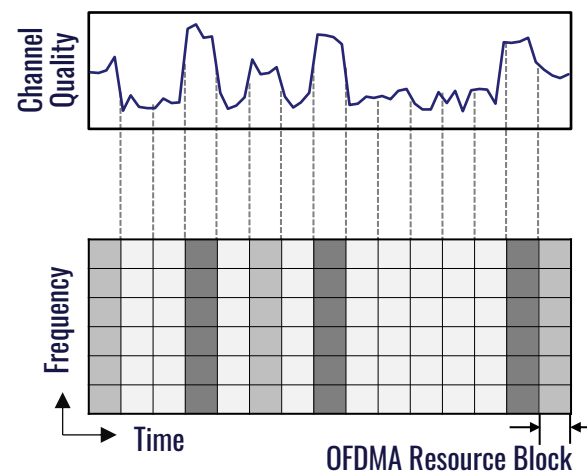
### → Network Slicing (Radio Resource Reservation)

- No guarantee
- Resource waste

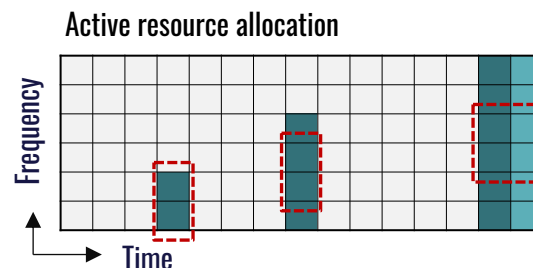
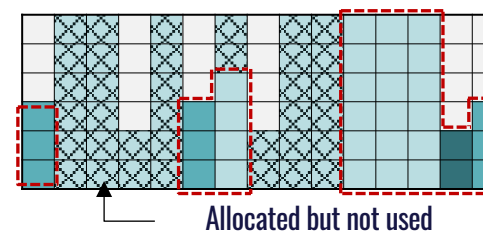
## Cellular API Approach for Hyper-Connectivity

### → Active Networking

- Persistent guarantee
- New pricing needed (to service providers)



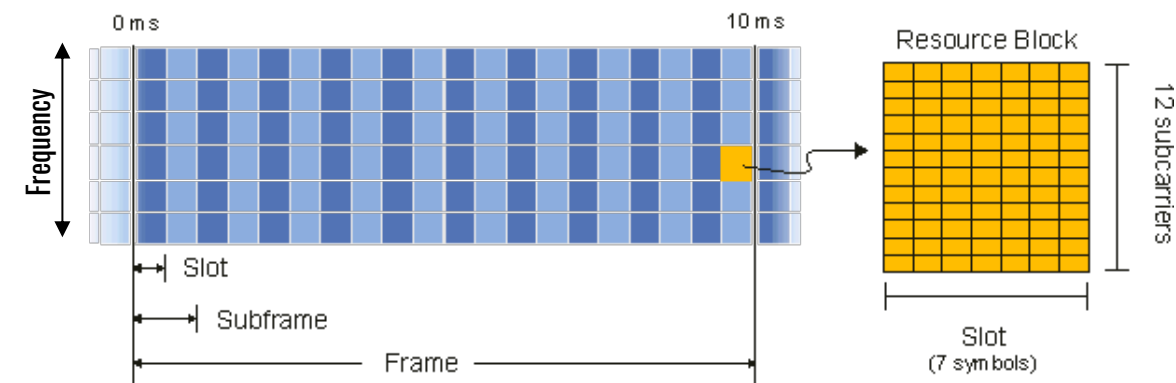
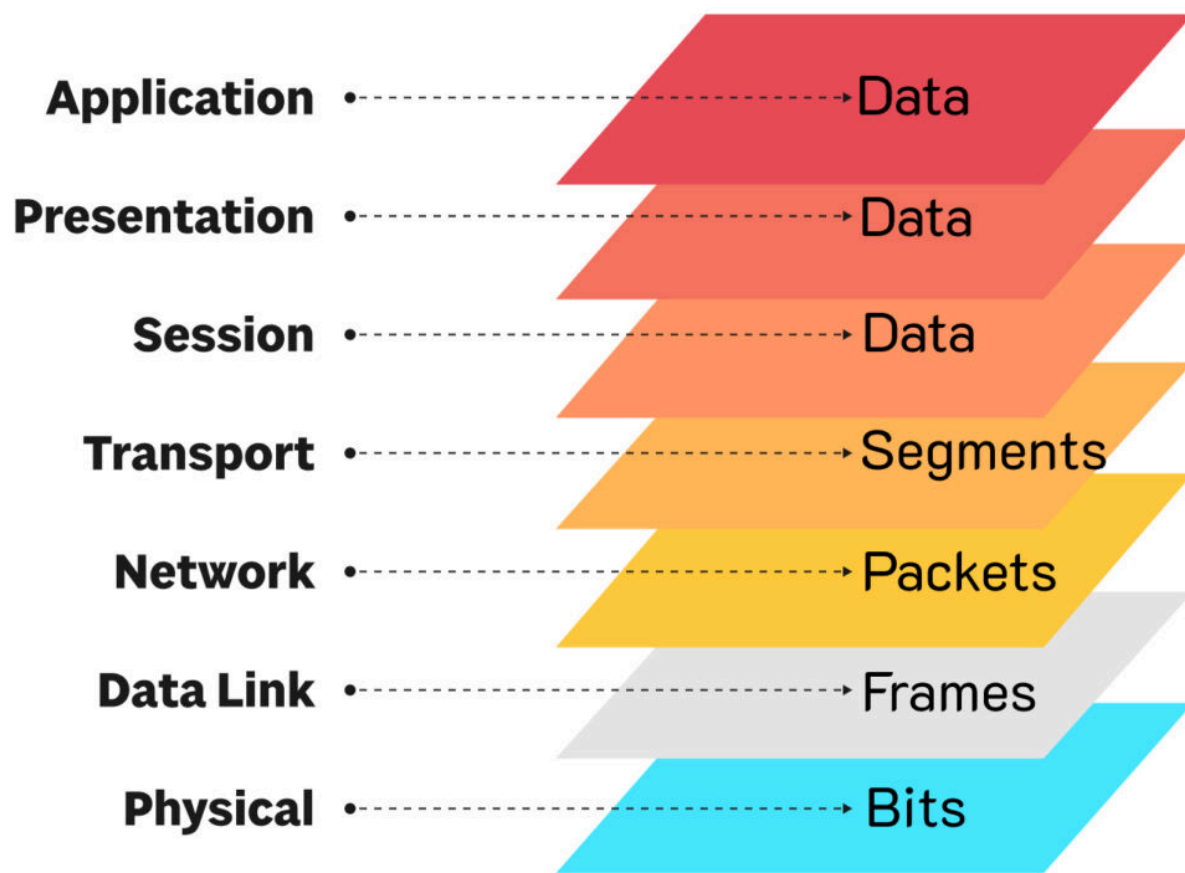
### Worst-case driven resource reservation



Resource efficiency:  
18/18

Resource efficiency:  
33/69

# From the Conventional Bottom-Up Philosophy to a New Top-Down Philosophy



## Shannon (Information) Capacity

Maximum achievable data rate (in bits/sec)

$$C = B \cdot \log_2 \left( 1 + \frac{S}{N} \right)$$

Radio Channel Bandwidth (in Hz)

Signal Power (in Watts)

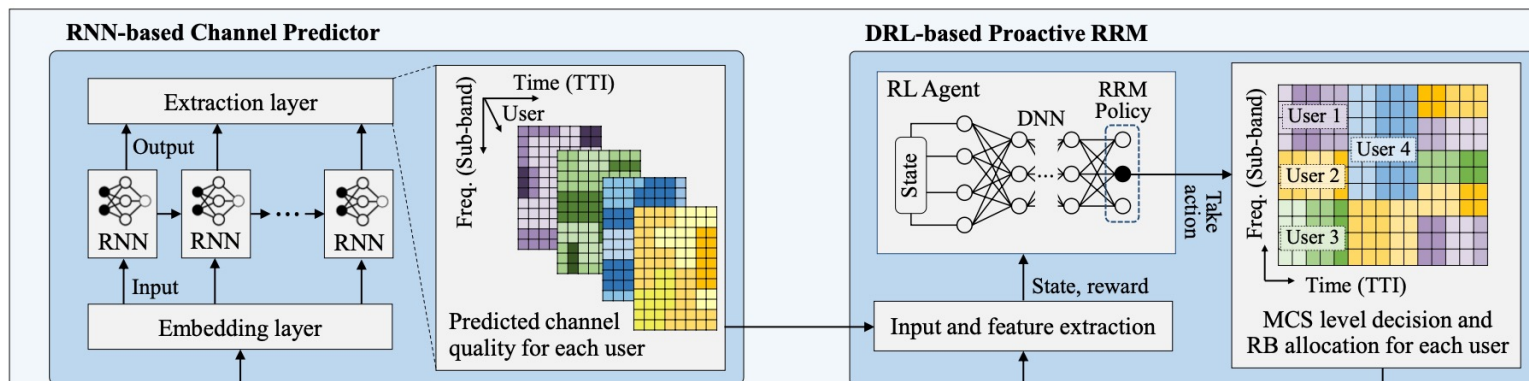
Noise Power (in Watts)

SNR (Linear Scale, not in dB)

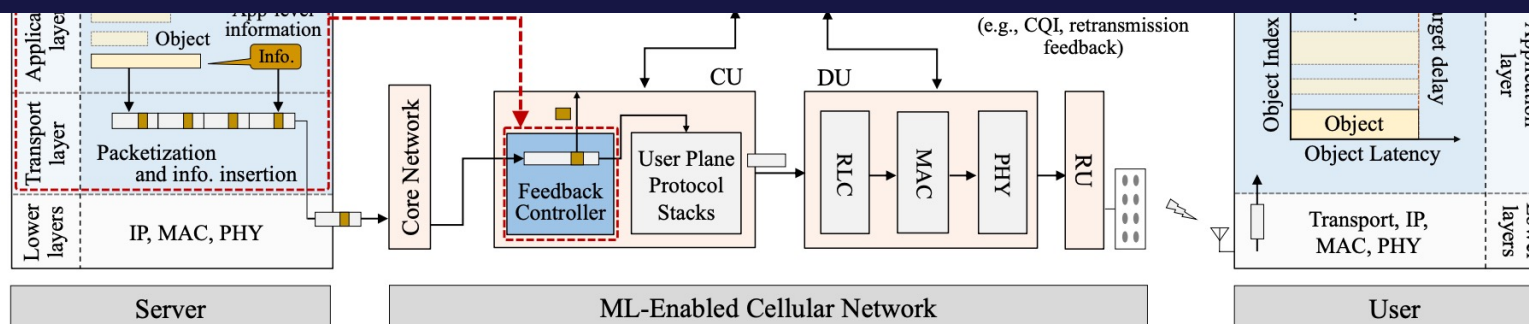
As this gets larger, C (Capacity) gets larger

# Cellular Network should be No More a Blackbox to Applications!

## Towards Enabling Performance-Guaranteed Networking in Next-Generation Cellular Networks



## Cellular Network should be able to Interact with Applications! Cellular Network 2.0: Active Networking for Cellular APIs





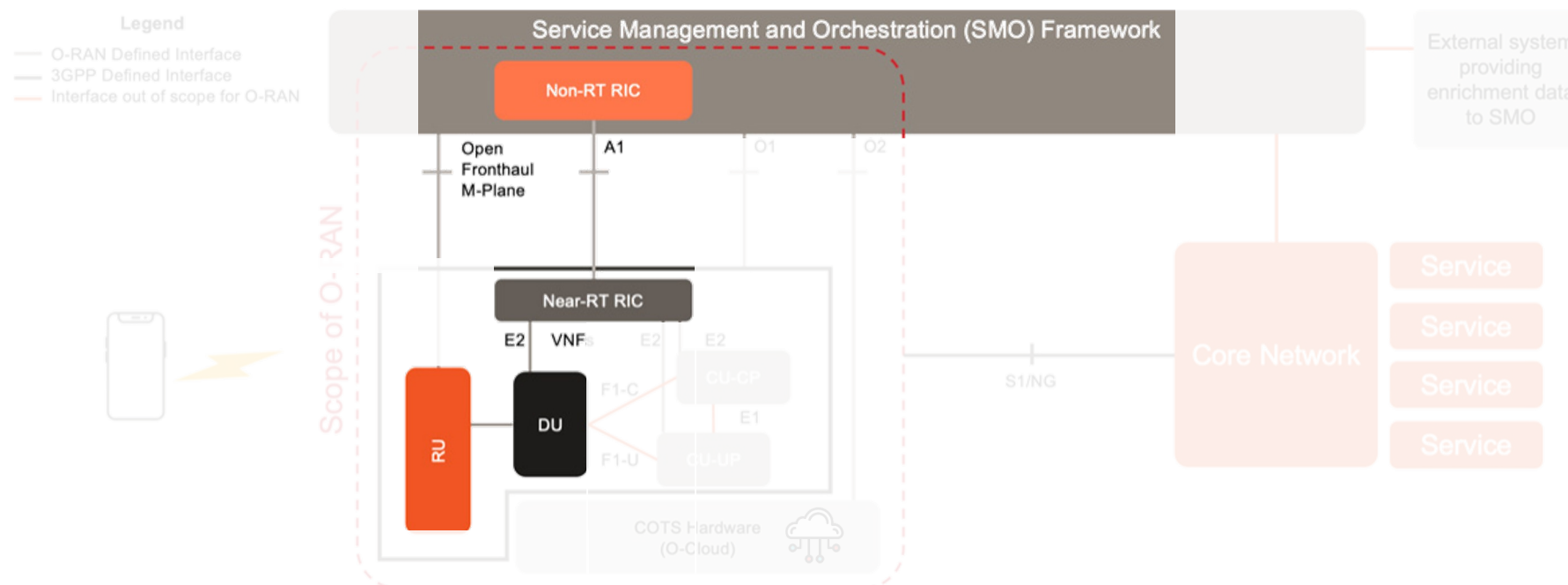
# Active Networking Can Make Cellular Network a Whitebox to Applications!



Cellular API supported by  
Softwarized Network Orchestrator



RAN time budget  
Core-network time budget  
Computing time budget



# True Benefits of Cellular Network 2.0 for Realtime ExoComputing



Hyper Connectivity to Applications from Cellular APIs



Infinite Computing Capability to Network Devices

Services running over CN2.0 becomes indistinguishable from services running on device.

Guaranteed  
ADU CT

Network application programming becomes straight-forward.  
(socket programming was a disaster.)

```
frame = capture( sensor );  
display( Cellular_API( frame, destination, service deadline ) );
```



Pricing/incentive model for **Cellular\_API()**  
\* Difficulty-based

CN2.0 Immediately Enables Tons of New ExoComputing Services

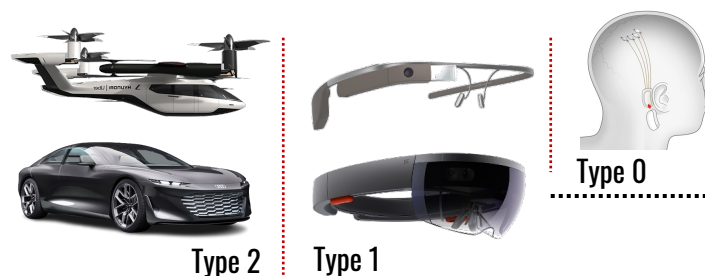


$$1) CT_{UL} \leq T_{budget}$$

$$2) CT_{DL} \leq T_{budget}$$

$$3) CT_{RoundTrip} \leq T_{budget}$$

# Is Hyper-Connectivity a Matter of Choice?



Computing over Hyper-Connectivity

vs.

On-device Computing

Near Real-time  
AI/ML Analytics/Computing  
Services (e.g., GPT-4)



No. Hyper-Connectivity is Indispensable! However, Is Networking Really Power-Efficient?

Higher Performance ← Higher Power



RTX3090  
(350W)

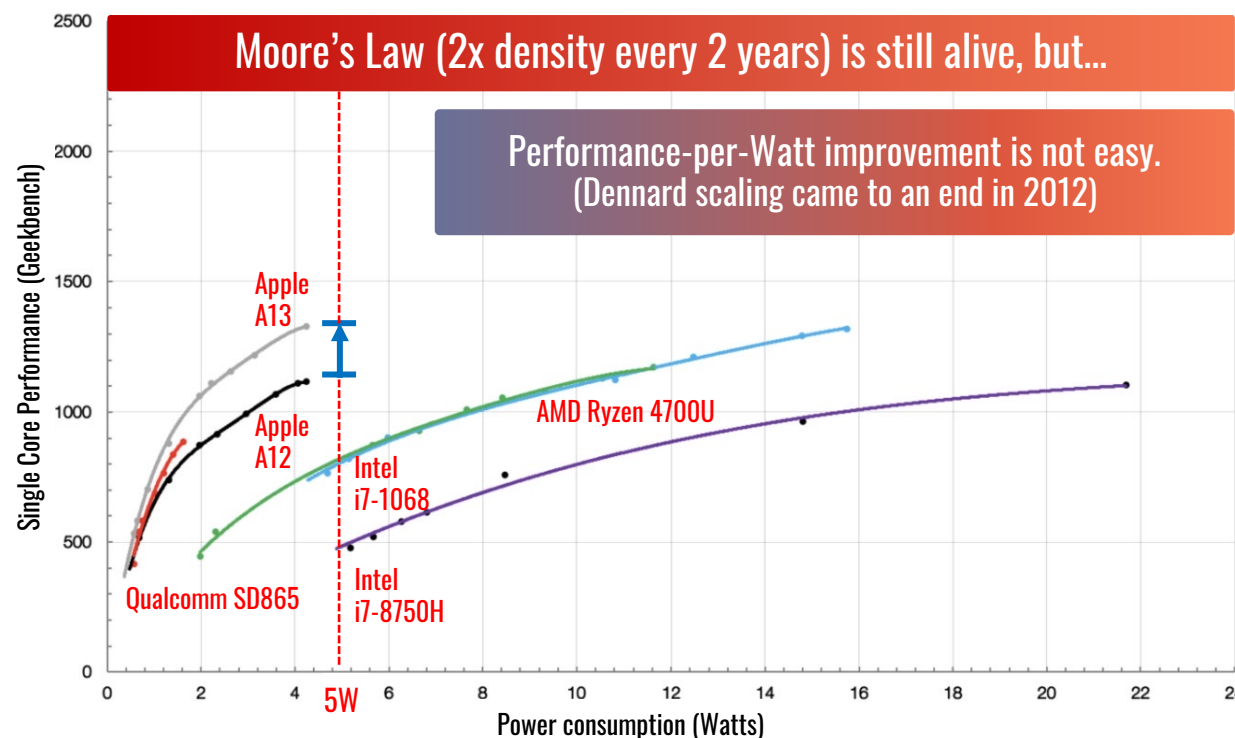


RTX2080  
(200W)

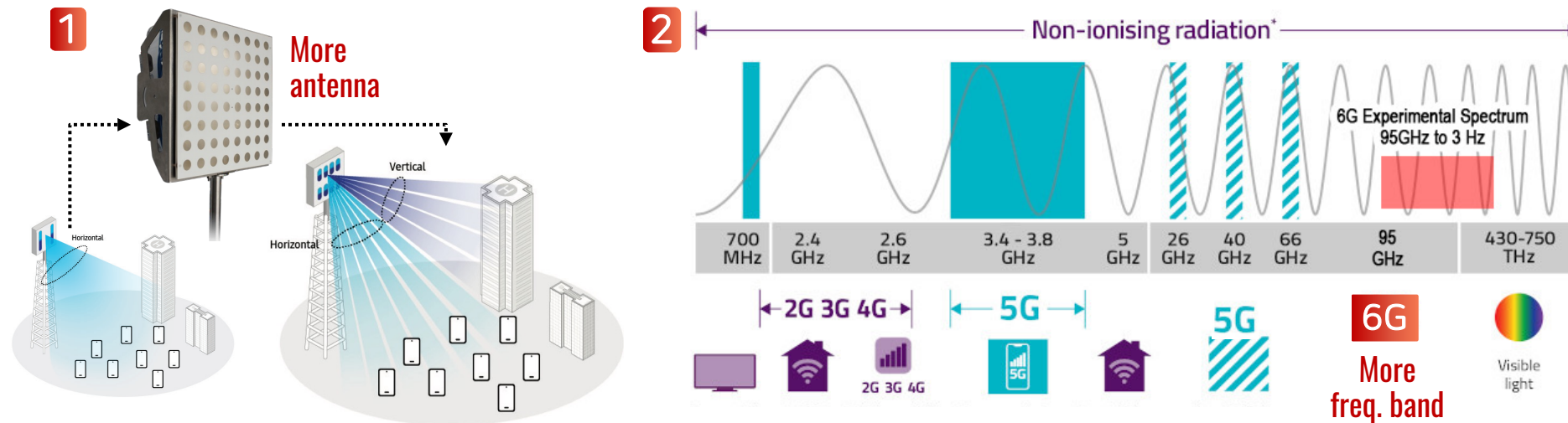
Thermal & Battery Capacity Limited (Type 1 & 0)



SD888  
(5W)

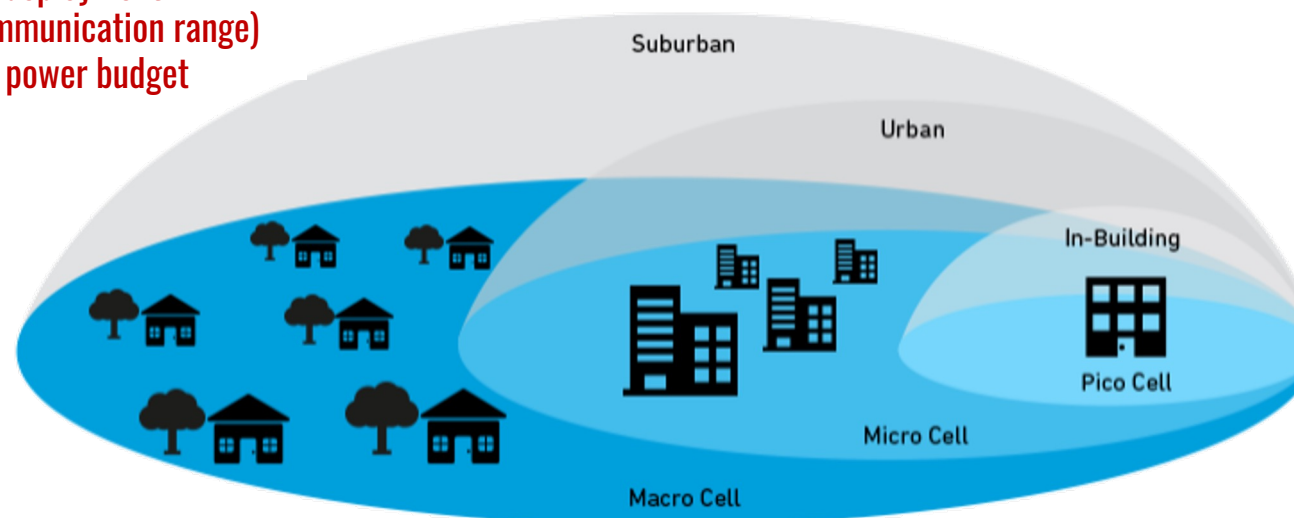


# Is Hyper-Connectivity Achievable with Limited Power Consumption?



Unlike On-device Computing, We have Multiple Options to Limit the Communication Power!

- 3** Ultra dense deployment  
(shorter communication range)  
to meet the power budget



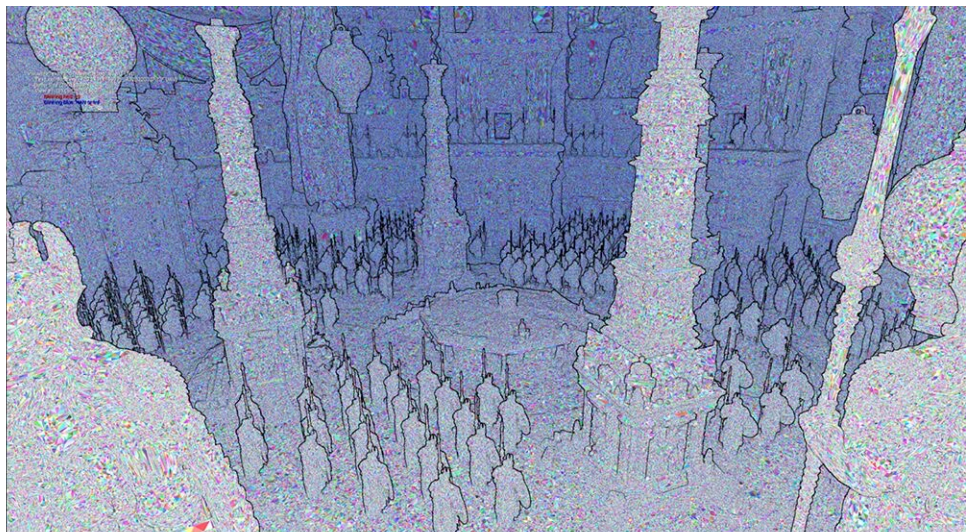


# ExoComputing with Hyper-Connectivity vs. On-device Computing



Not a pre-rendered video. PS5 Unreal Engine 5 (at 1440p, 30fps, 10TFlops from 150W)

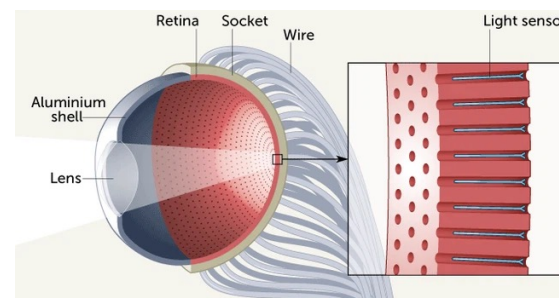
# ExoComputing with Hyper-Connectivity vs. On-device Computing



## 1440p Rendering on Device (3.5 MP)

30 fps - 33 ms (budget that a frame is rendered)

**150 W** (10 TFlops, about 16 billion triangles)



Human Eyes Spec  
~ 2 x 63 MP

RD: x(6x6x6)

SR: x(6x6)

**32 KW**



## 1440p Streaming from Edge (3.5 MP)

14 MB (BMP, 2560x1440x32 bits)

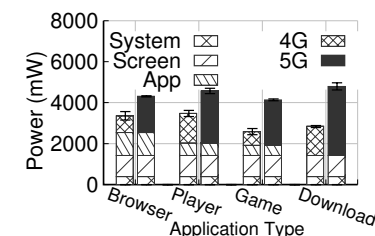
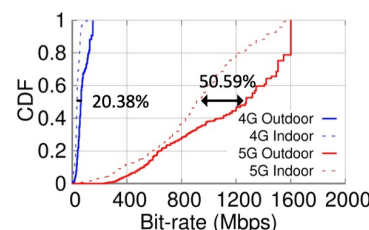
<1.4 MB (JPEG, BPG, Web-P)

9.3 ms (1.4 MB / 1.2 Gbps)

**2.5 W** (1.2 Gbps)

**5 W @ 25ms**

(50 MB / 16 Gbps)



Understanding Operational 5G: A First Measurement Study on Its Coverage, Performance and Energy Consumption, ACM SIGCOMM 2020



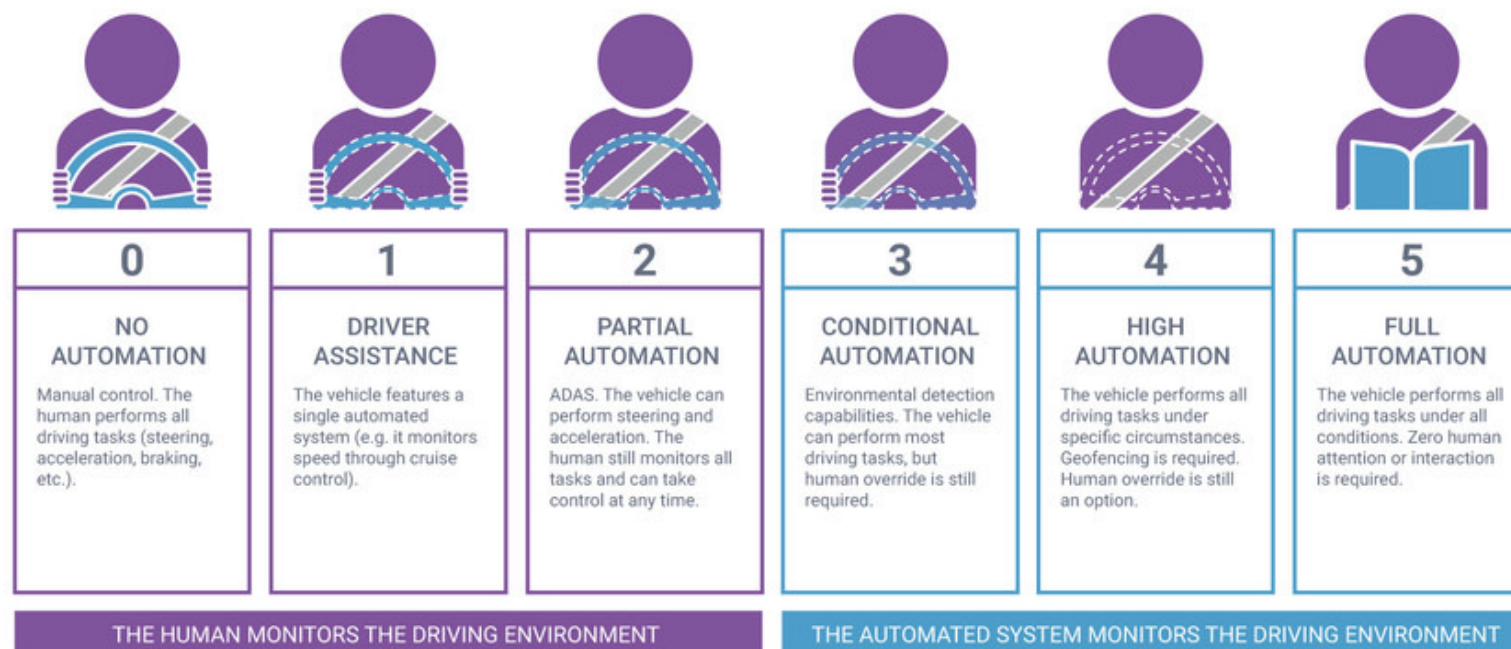
# ExoComputing with Hyper-Connectivity vs. On-device Computing



한번.. 누르면?

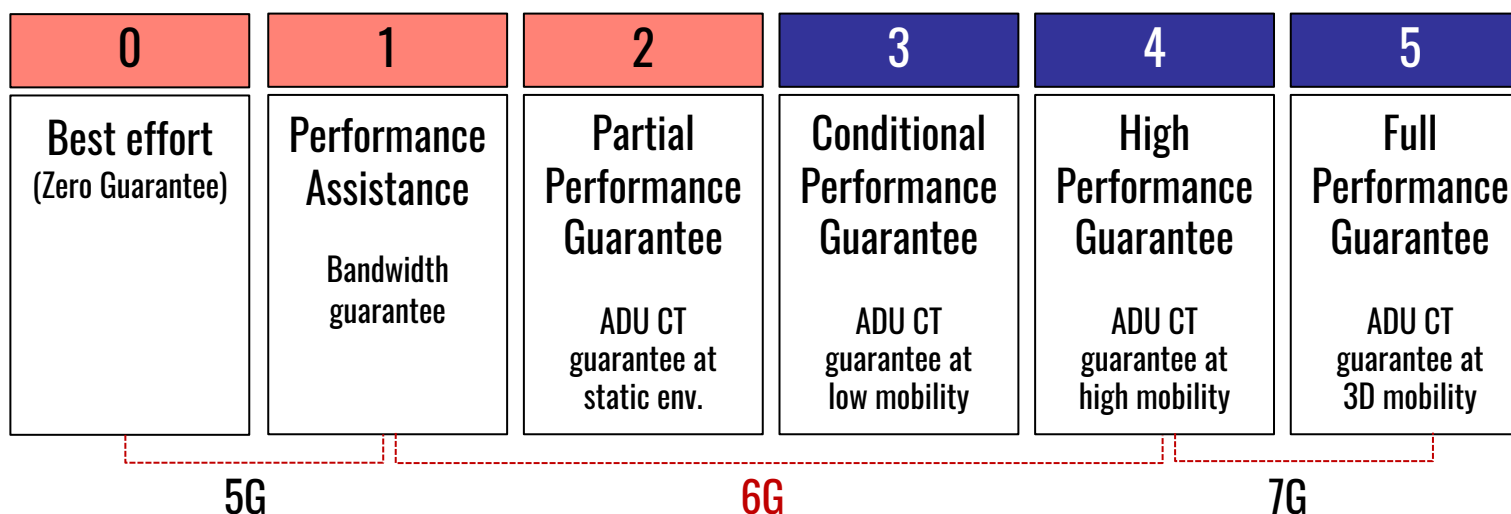
# A Roadmap of Enabling Hyper-Connectivity (6G & Beyond)

## Levels of Driving Automation

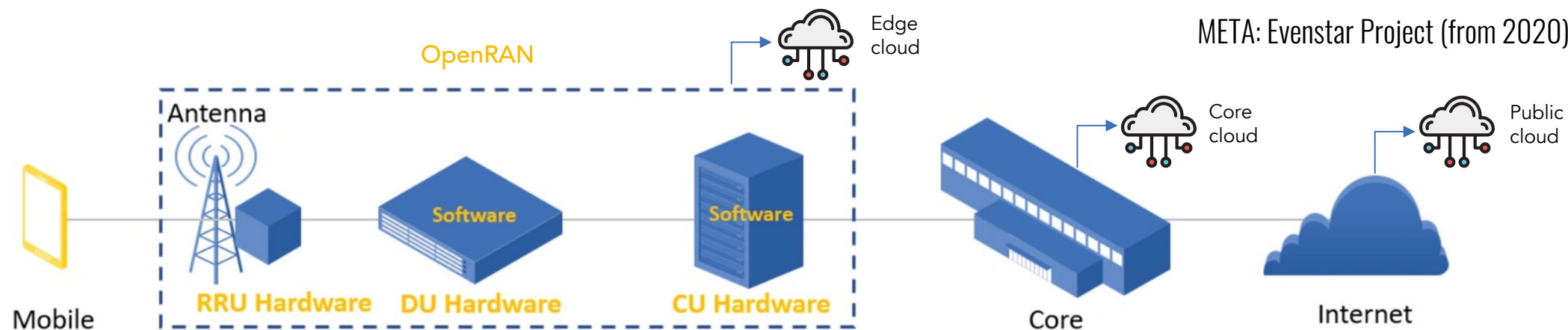


## Levels of Performance Guarantee

(per-Region Certification)

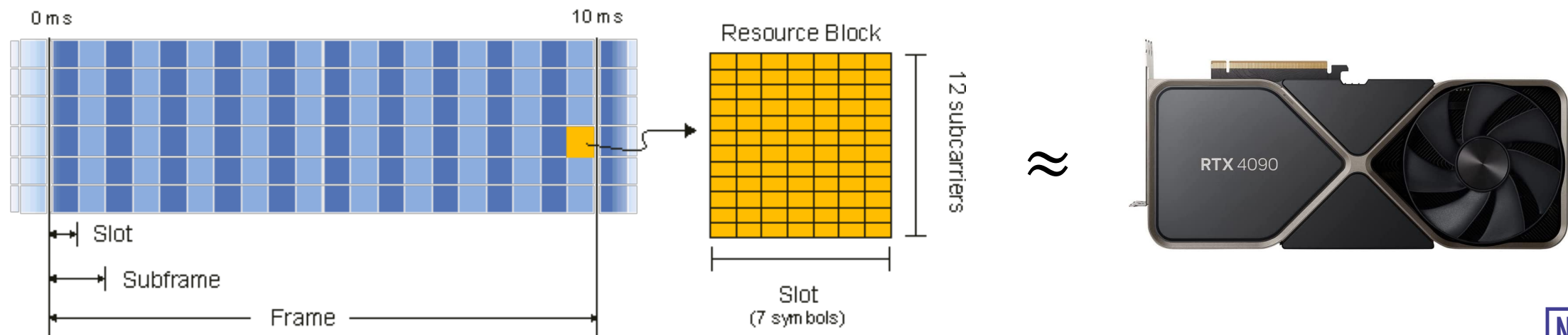


# Emergence of a New Networked Computing Service Operator



$$JobCT = T_{wait} = t_{en/decoding} + t_{tx up} + t_{computing} + t_{tx down} + t_{en/decoding} \leq T_{budget}$$

New Business Question: What is the Value of One Radio Resource Block in Terms of GPU Cores?





# Thank You

[kyunghanlee@snu.ac.kr](mailto:kyunghanlee@snu.ac.kr)

[nxc.snu.ac.kr](http://nxc.snu.ac.kr)

