



“AI based Network Resource Management (1)”

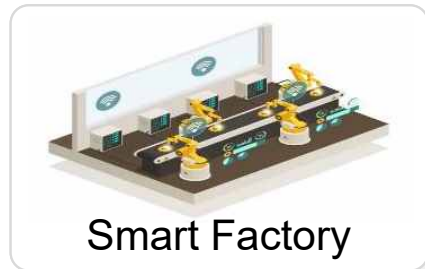
Choong Seon Hong

Department of Computer Science and Engineering
Kyung Hee University, Republic of Korea.

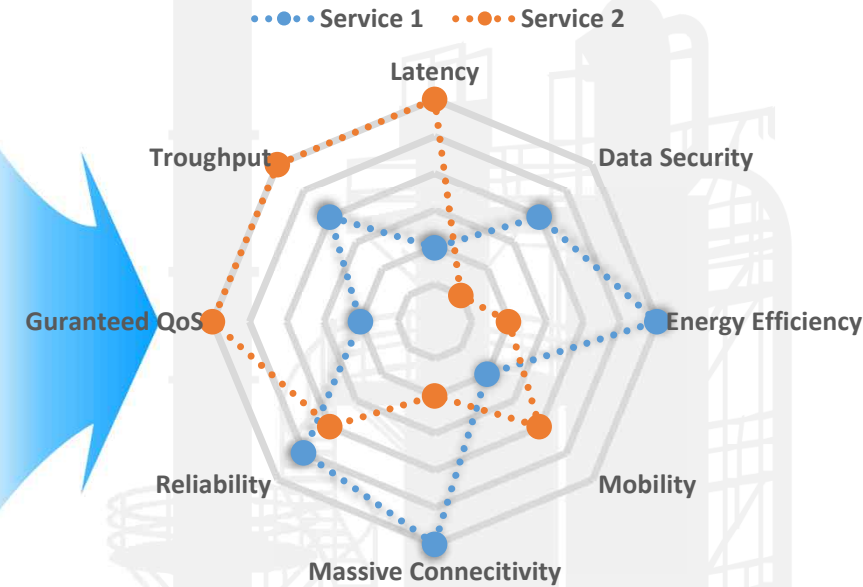
- Introduction
- Network Slicing: The Concept
- Use Case 1: Virtual Reality
- Use Case 2: Chunk-Based Resource Allocation
- Use Case 3: Energy Efficient Communication and Computation Resource Slicing for eMBB and URLLC
- Use case 4: Joint Communication, Computation, and Control for Computational Task Offloading in Vehicle-Assisted Multi-Access Edge Computing
- Use case 5: Collaboration in the Sky: A Distributed Framework for Task Offloading and Resource Allocation in Multi-Access Edge Computing
- Concluding Remarks

Introduction

- Diverse Requirements
- Requirements in
 - Manufacturing Industry
 - Transportation Industry
 - Health Sector
- Evolution of Cellular Systems
- Challenges to realize 5G Networks



Service-Specific Requirements



✓ The network resources management becomes more complex because of the very diverse requirements.

- The manufacturing industry requires
 - high-quality,
 - time-sensitive,
 - automated,
 - intelligent and
 - flexible industrial control.
- So that materials, products and processes can be monitored, optimized and controlled in real time.



Example applications

- **Low-latency, high-reliability and high-availability connectivity, mobility and precise positioning** of all the devices (e.g. sensors and actuators) for **real-time monitoring and control of processes**, and **end-to-end** logistics and asset tracking
- Connectivity and local processing for **real-time video capture** and video-based applications
- **Connectivity of massive numbers of sensors**, and platforms for collation and **processing of large amounts of data**
- **Augmented reality to optimize** and improve maintenance tasks.

- The road transportation industry expects to provide efficient, safe, environmentally-friendly and comfortable transportation, especially by exploiting the potential of artificial intelligence, to achieve connected and automatic driving through perception, decision and control.



Example applications

- **Transmission of high quality video** or images of **road condition and roadside facilities** to help navigation, remote and automatic driving, as well as identification of blind zones and other vulnerabilities for vehicles
- **Real-time communication among vehicles and road infrastructures**, coupled with precise vehicle positioning and local (edge) computing capabilities, enabling identification of potential dangers, which can help decision making including route planning and updating, emergency braking, and intelligent car collision avoidance

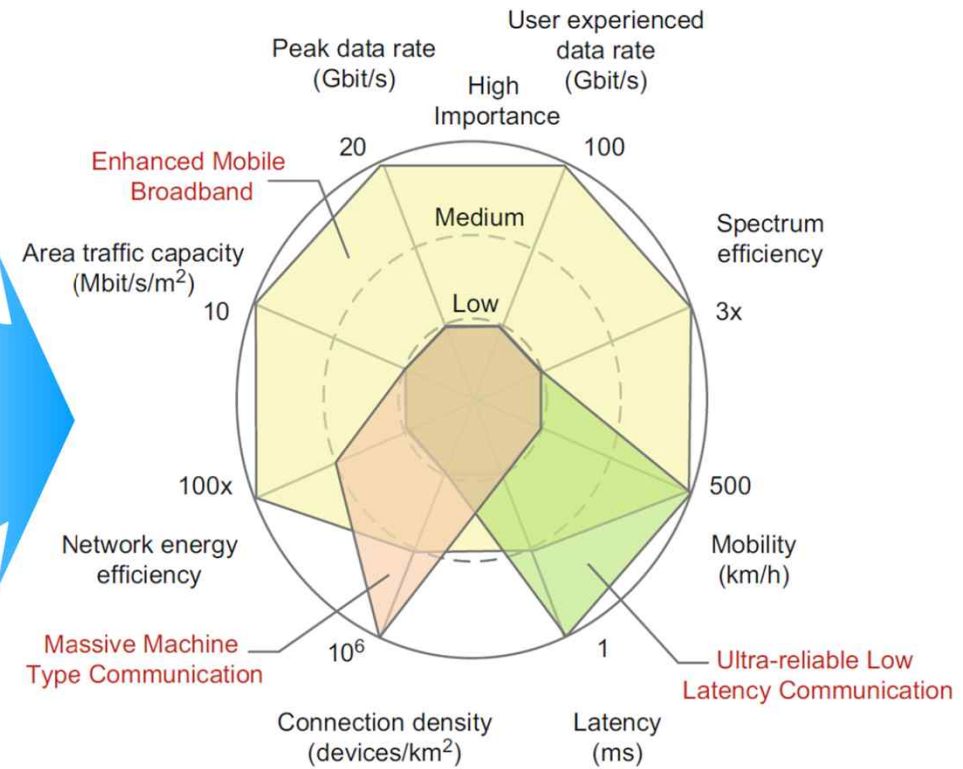
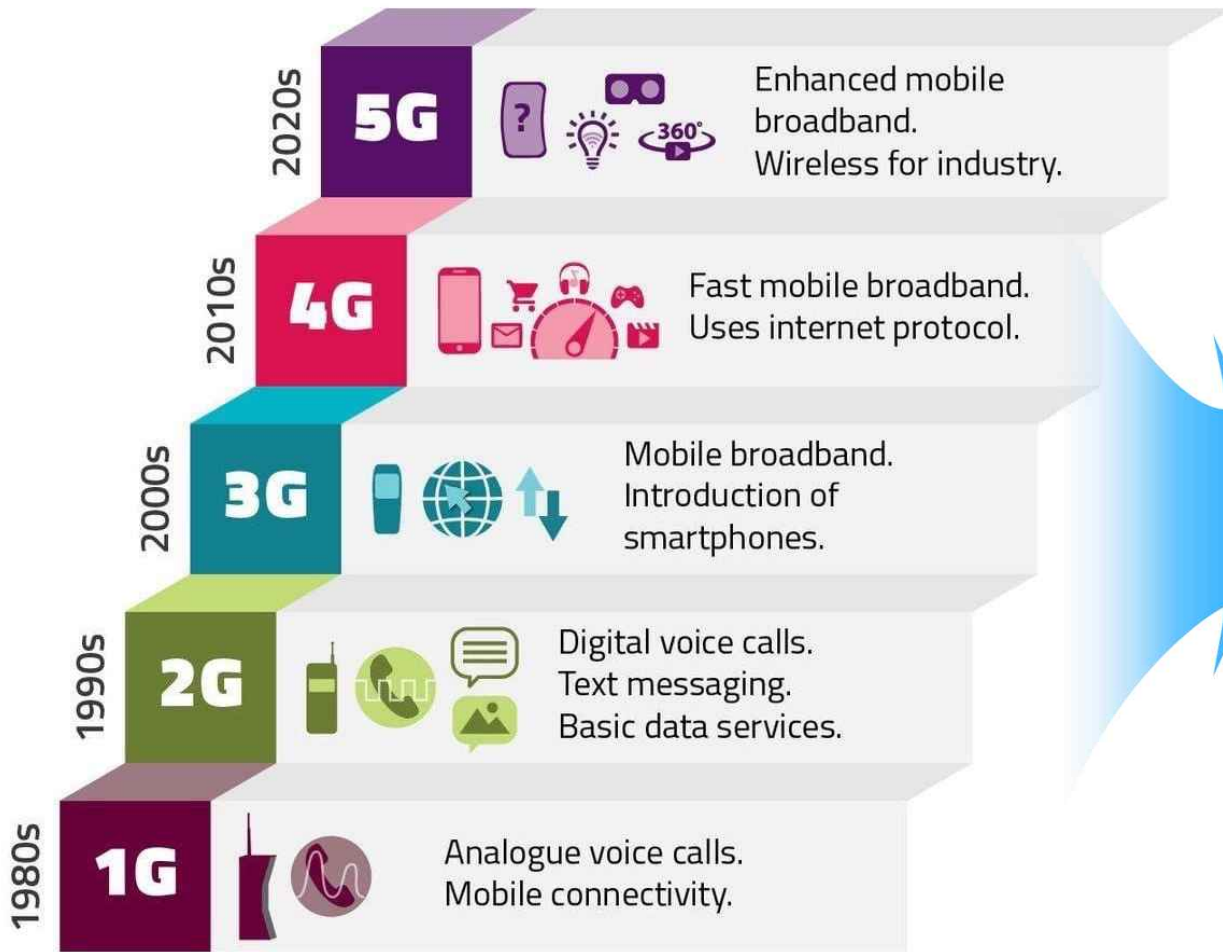
- The health industry requires a balanced allocation of medical resources, portable and intelligent medical equipment, improved medical vehicle treatment capabilities, and the transformation of surgical operations from the operating room to multiple regions.



Example applications

- Wide-area continuous coverage for ambulances, including **sending live high quality video and patient vital signs in real time** to the command center in the hospital
- **Sensors collecting vital signs from the wearable devices** of patients or the elderly, wherever they are, helping remote medical staff make timely treatment decisions and administer medication remotely

Introduction : Evolution of Cellular Systems



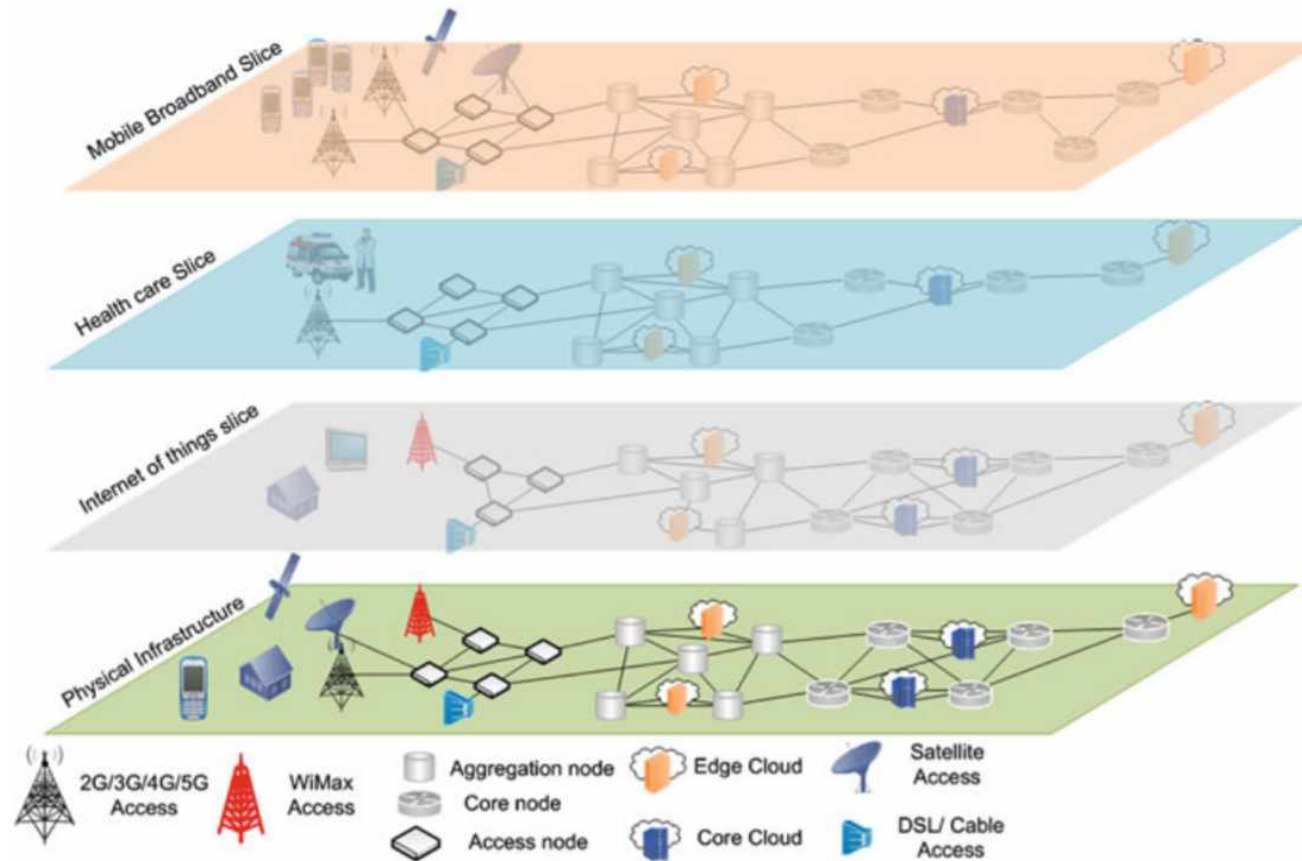
- Generally, to realize diverse 5G use cases there is need to resolve the given challenges:
 - Scalability and Reliability
 - Interoperability
 - Sustainability
 - **Network Slicing**
 - Security
 - Integration of AI

Network Slicing: The Concept

- Network Slicing
- Key Enablers
- Network Slicing: Industrial Efforts

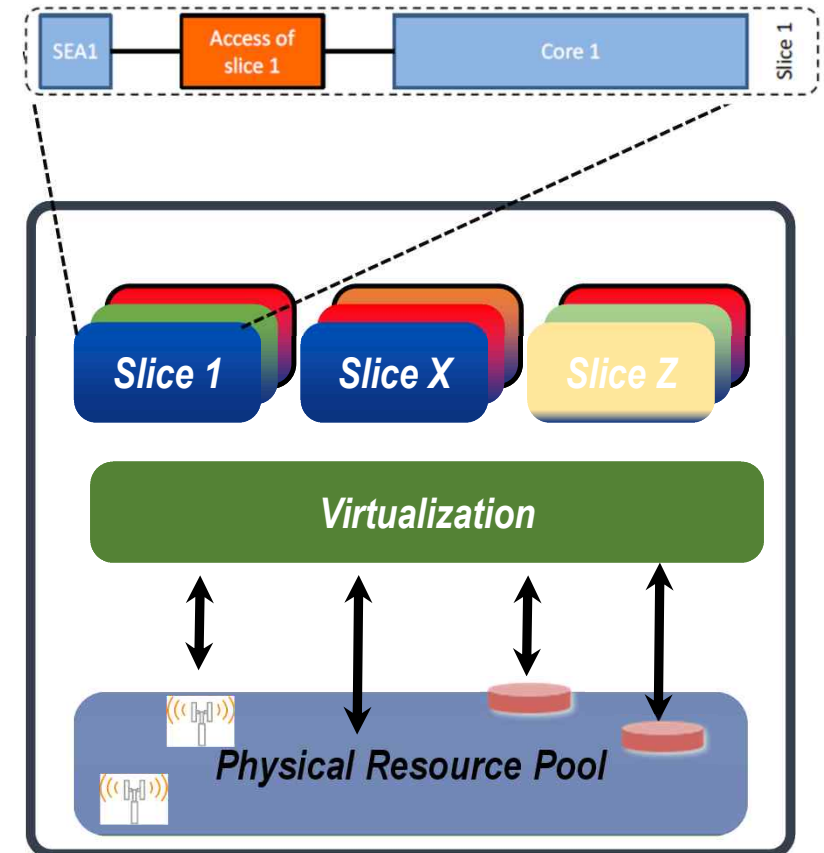
Network Slicing (1)

Network slicing is a new network approach that can provide highly tailored services to specific customer groups and even individual customers.

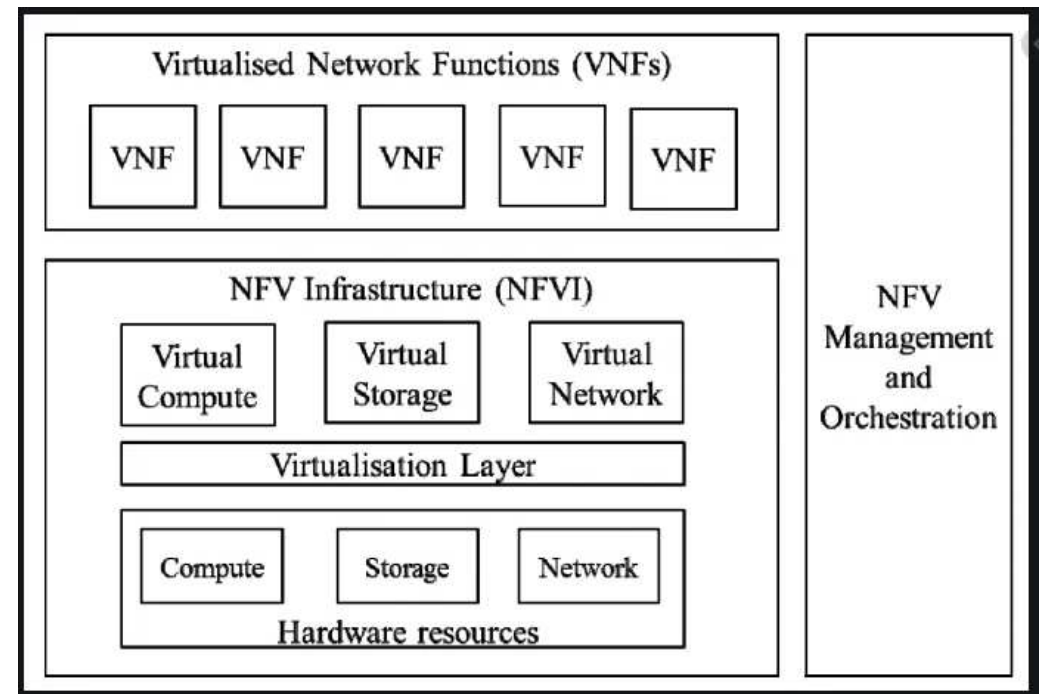
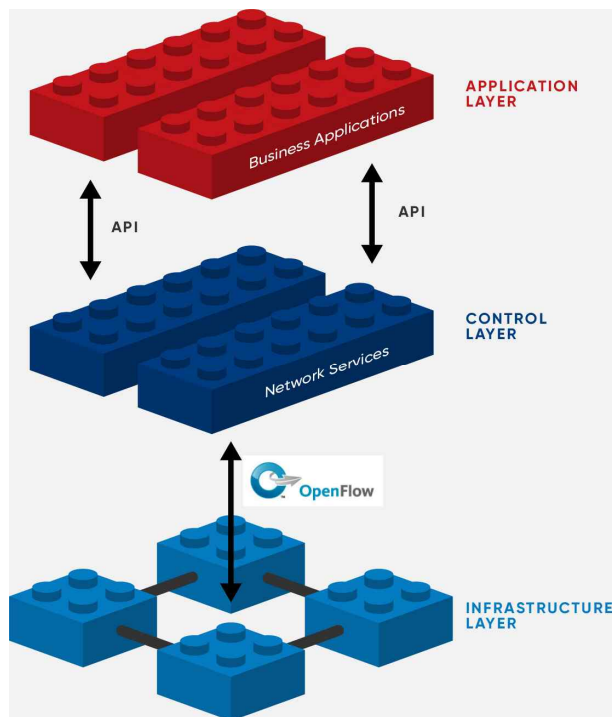


- Network slicing can fulfil the diverse requirements of these novel network services
- Network slicing enables one physical network into multiple, virtual, end-to-end (E2E) networks, each logically isolated including device, access, transport and core network
- A slice is dedicated for different types of service with different characteristics and requirements given to a Service End-point Agent (SEA)
- Enforce strong isolation between slices, i.e., actions in one slice do not affect another

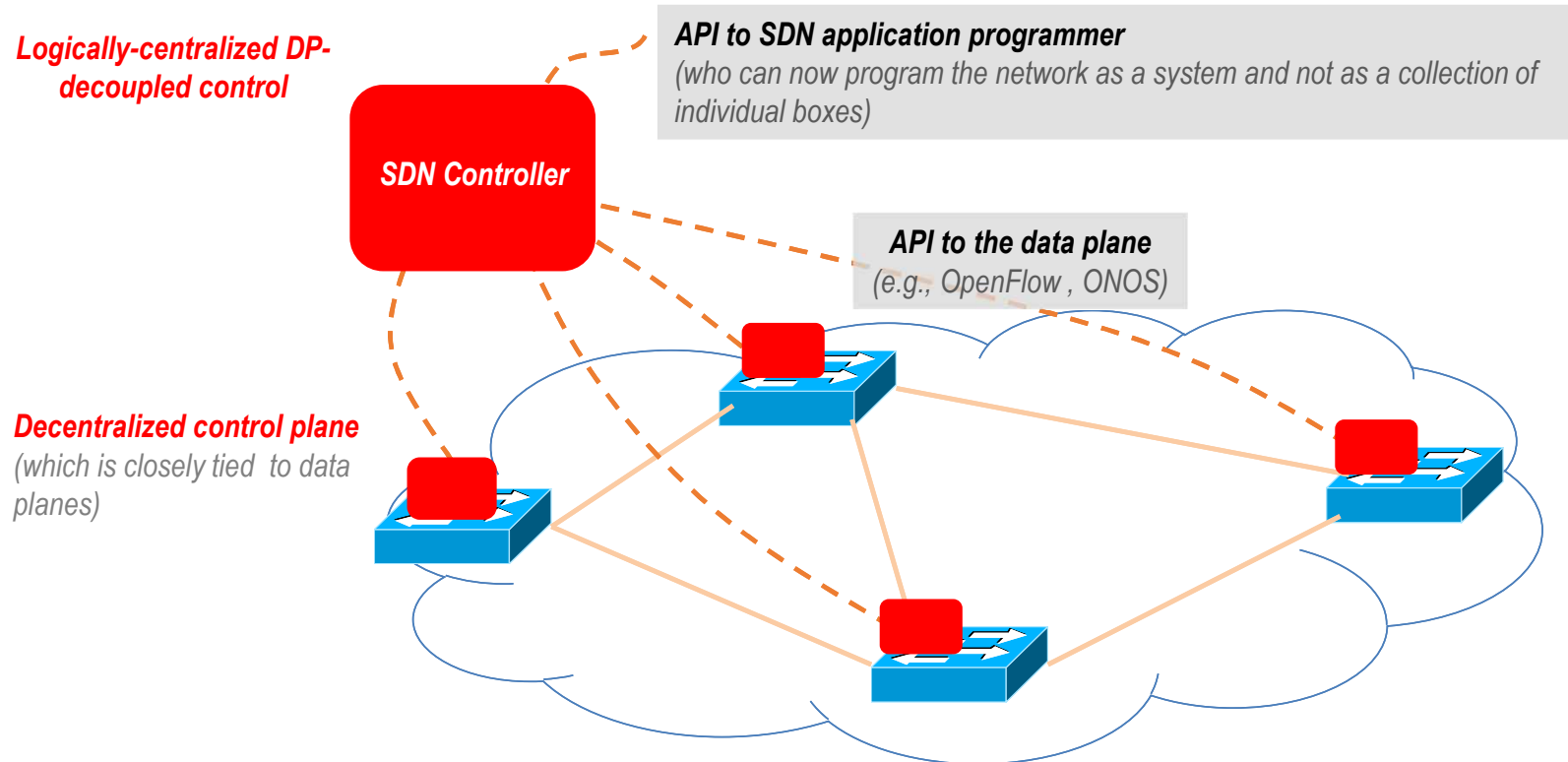
Network Slicing



- Network Slicing enablers: How to do it ?
 - Software-defined networking (SDN)
 - Network Functions Virtualization (NFV)



High-level NFV framework. Source: [ETSI](https://www.etsi.org/)



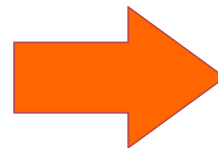
At the highest level, the SDN movement is an effort to build networks you can program at a higher level of abstraction— just as you can program a computer.



SDN enables programmability which is important for network slicing

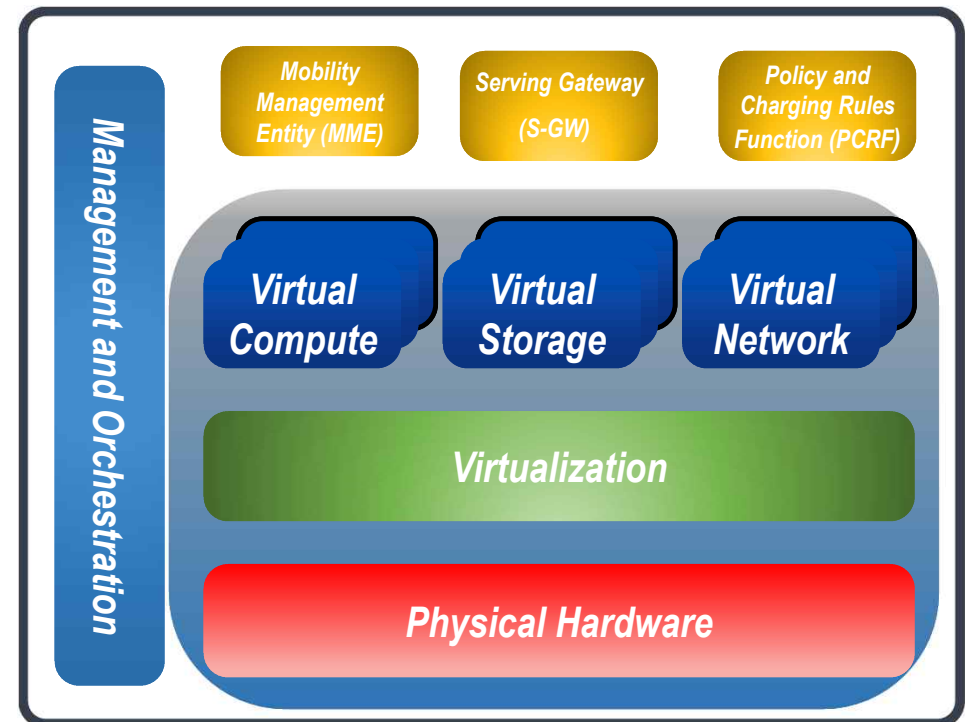


*Vertically integrated
Closed, proprietary
Slow innovation*



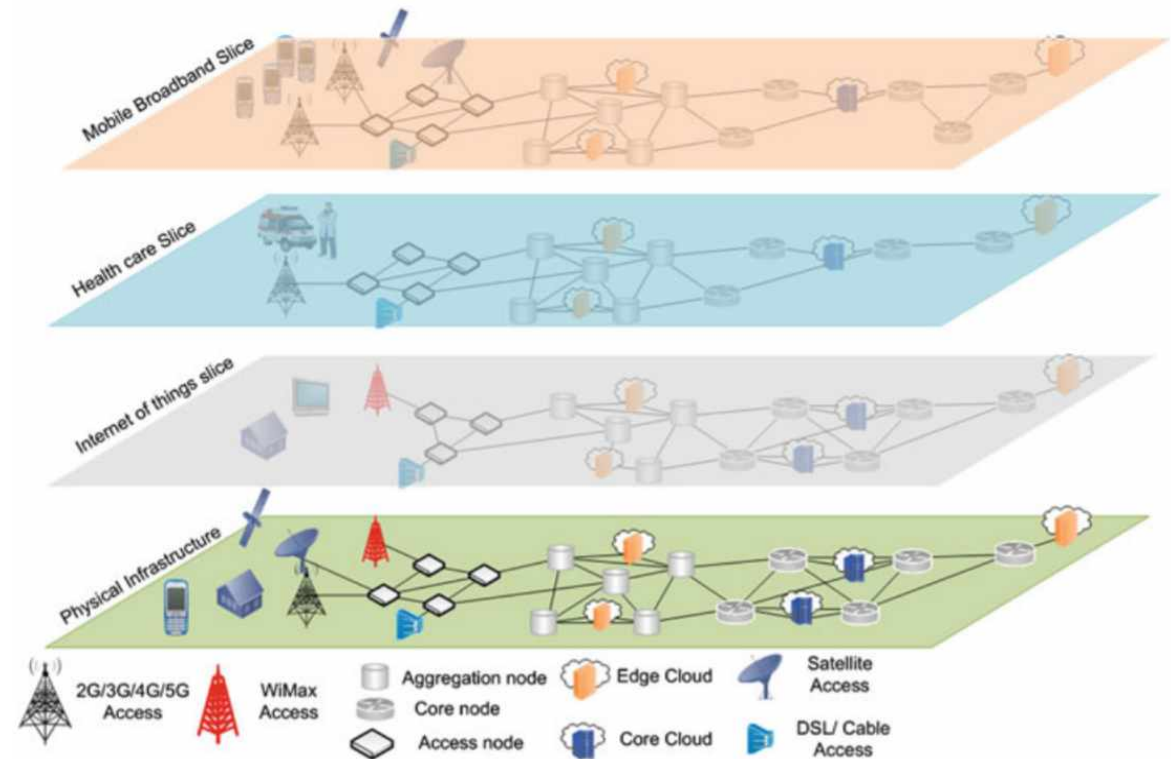
*Horizontal
Open interfaces
Rapid innovation*

- A network architecture concept that uses the technologies of *IT virtualization* to virtualize entire classes of network node functions that may connect, or chain together, to create communication services
- NFV is envisioned to play a crucial role in network slicing as it will be responsible to build isolated slices based on user service requirements



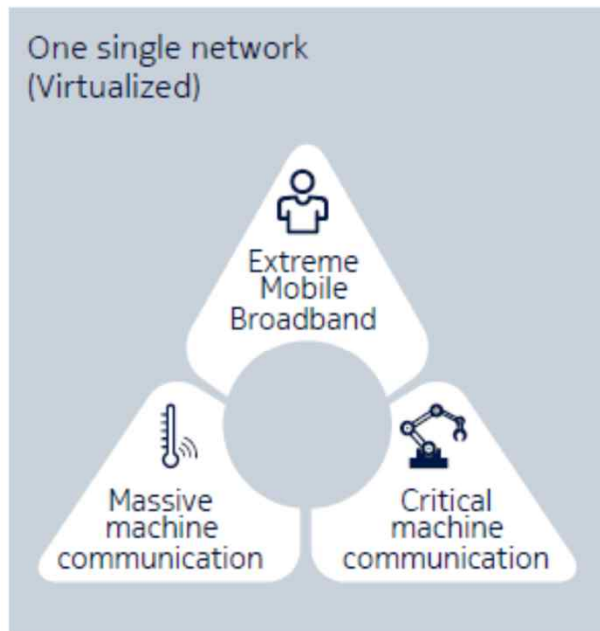
- Due to massive success of NFV and SDN in wired domain, a number of studies are being conducted to adopt them both in the core and radio access networks (RANs) for future cellular networks such as:
 - CORD (Central Office Re-architected as a Datacenter) [1]
 - Radisys M-CORD [2]
- Wireless network virtualization (WNV) is a novel concept for *virtualizing the RANs* of future cellular networks
- WNV has a very broad scope ranging from spectrum sharing, infrastructure virtualization, to air interface virtualization

- Network Slicing Principles
 - Slice Isolation
 - Elasticity
 - End-to-End Customization



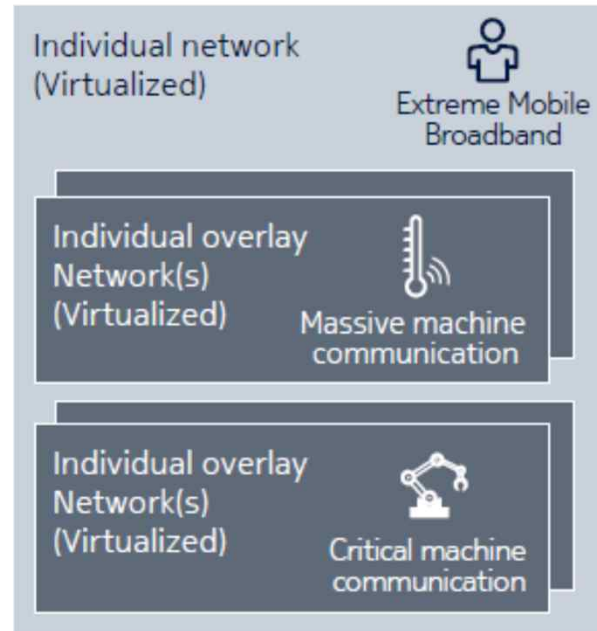
Three network scenarios provide comparative costs

Single network



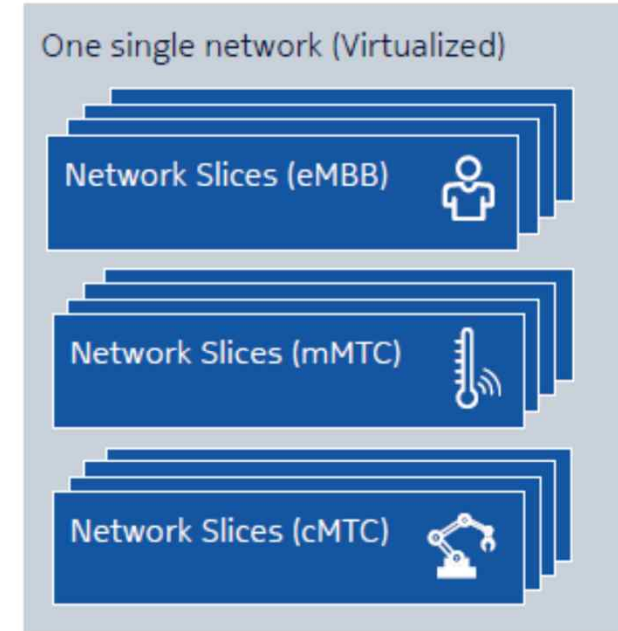
Best-effort support of eMBB, mMTC and cMTC

Dedicated networks



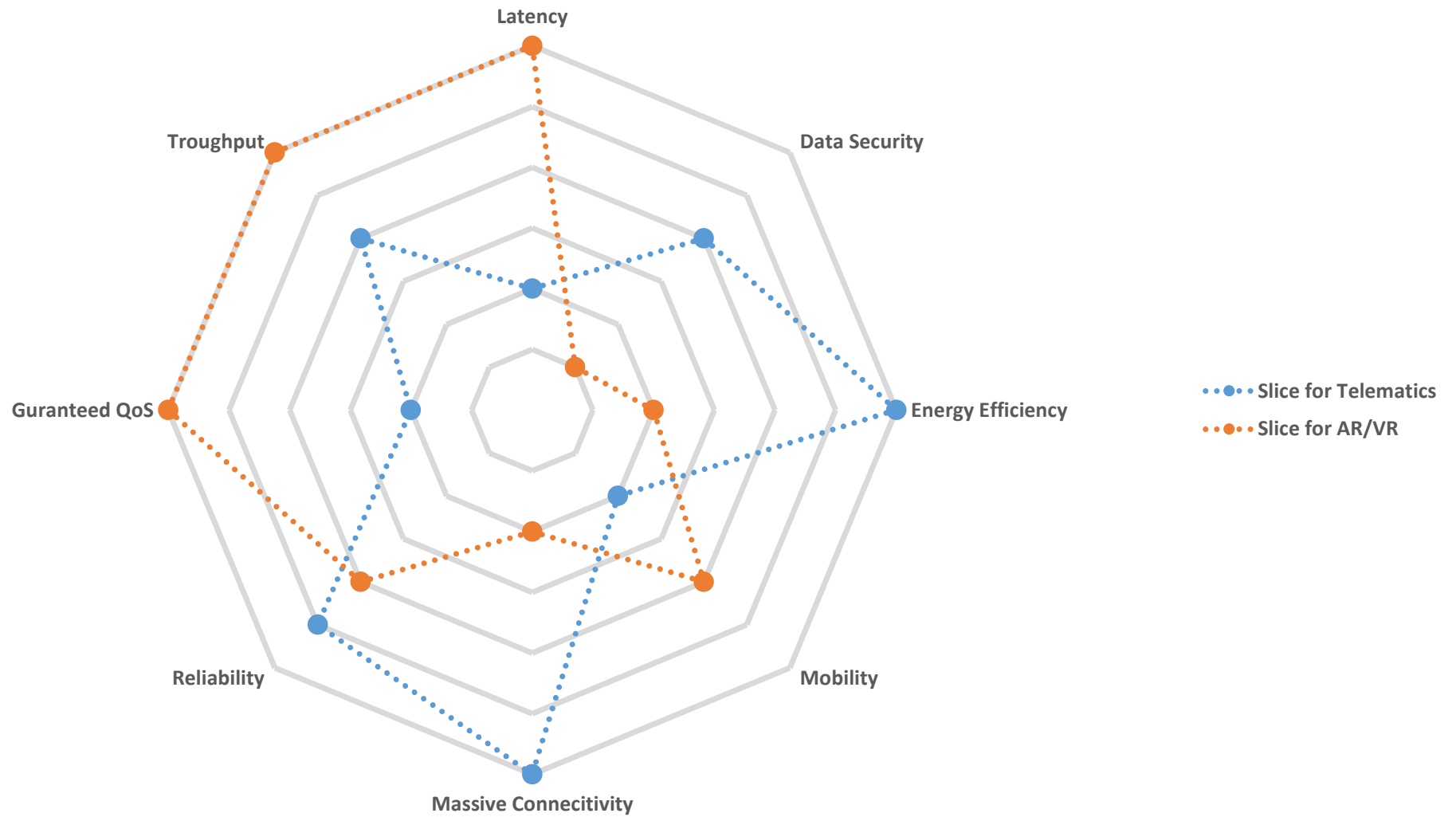
Customized dedicated networks for mMTC and cMTC

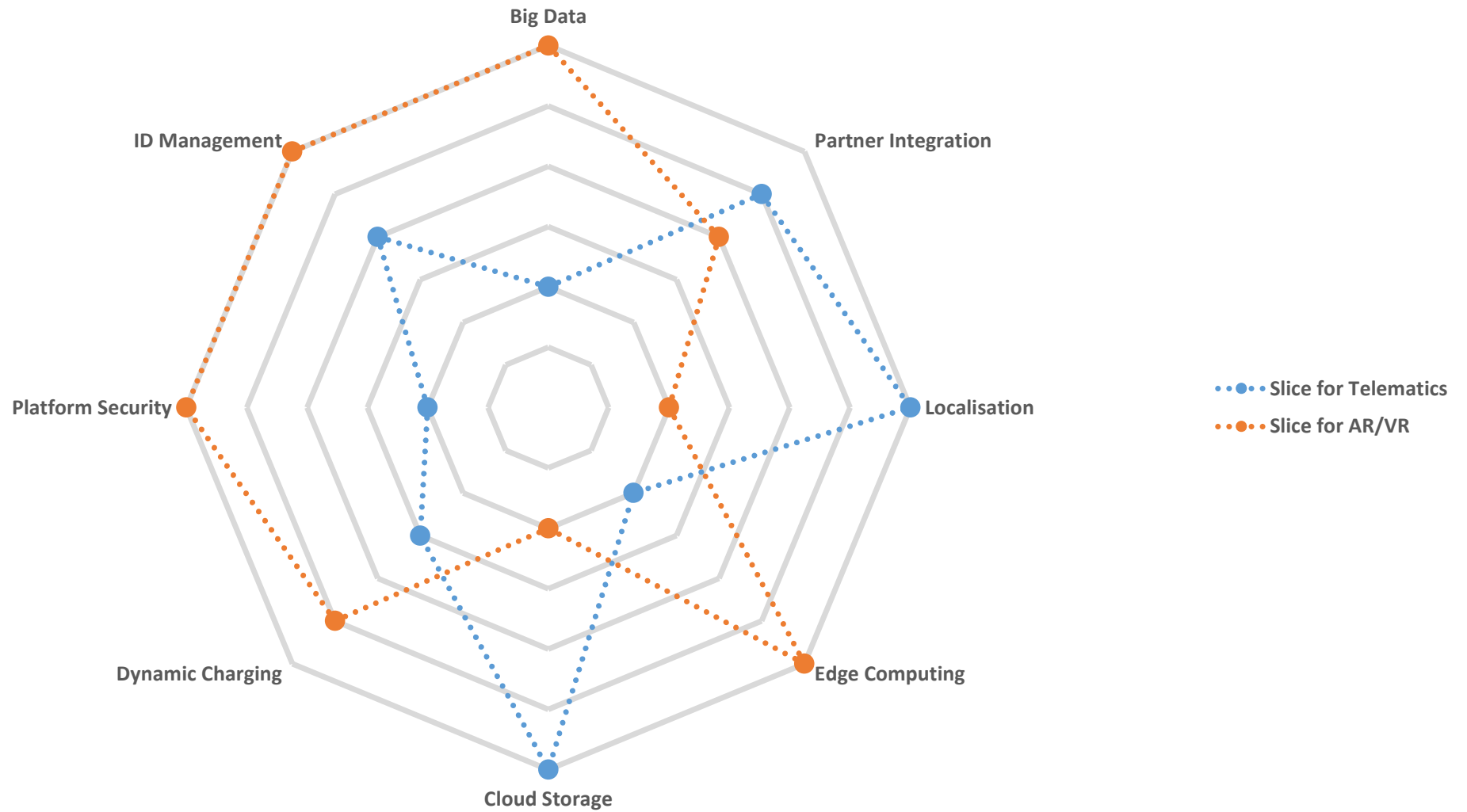
E2E Network Slicing



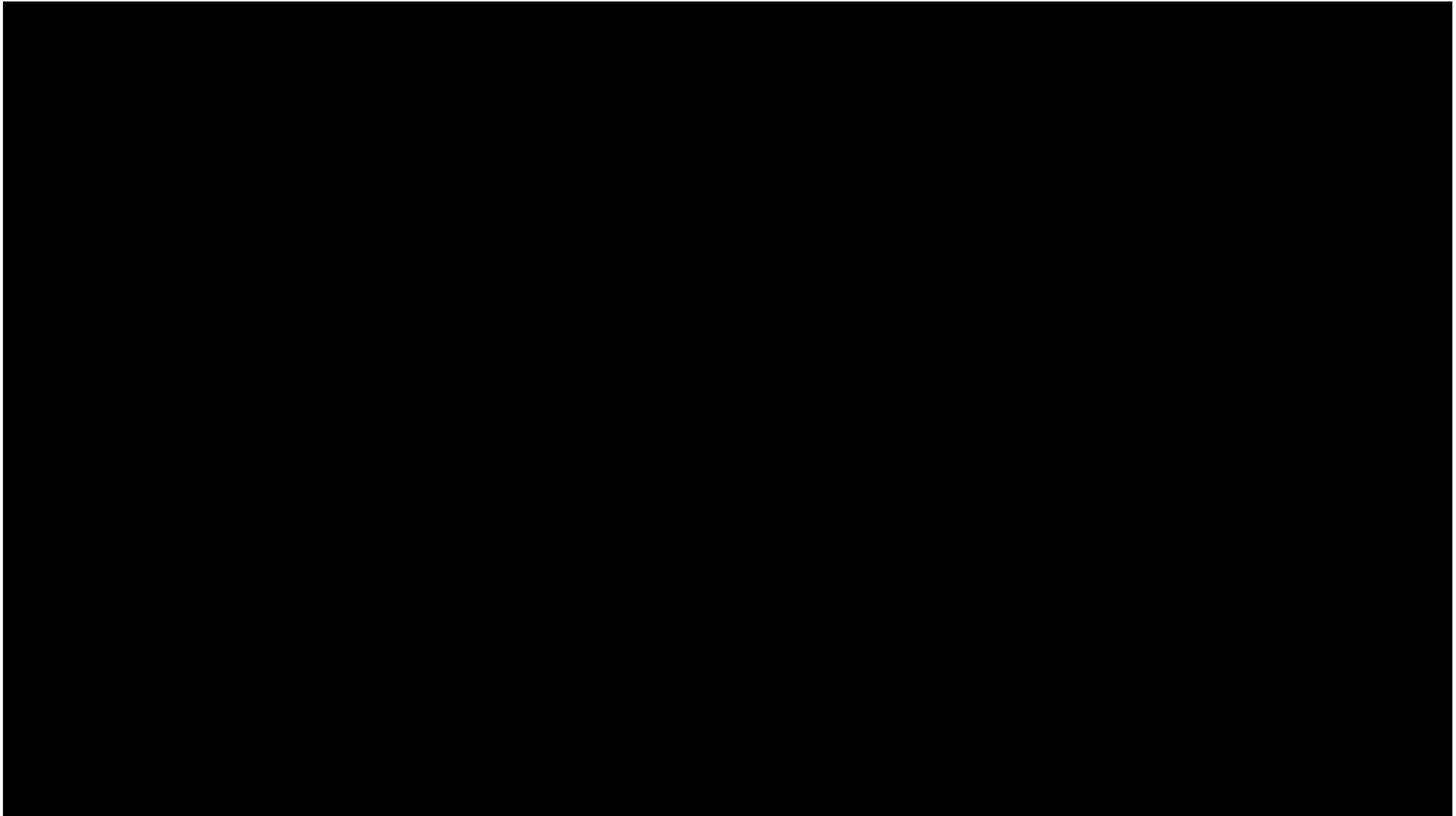
Customized network slices for eMBB, mMTC and cMTC

Network Slicing based on Network Capability

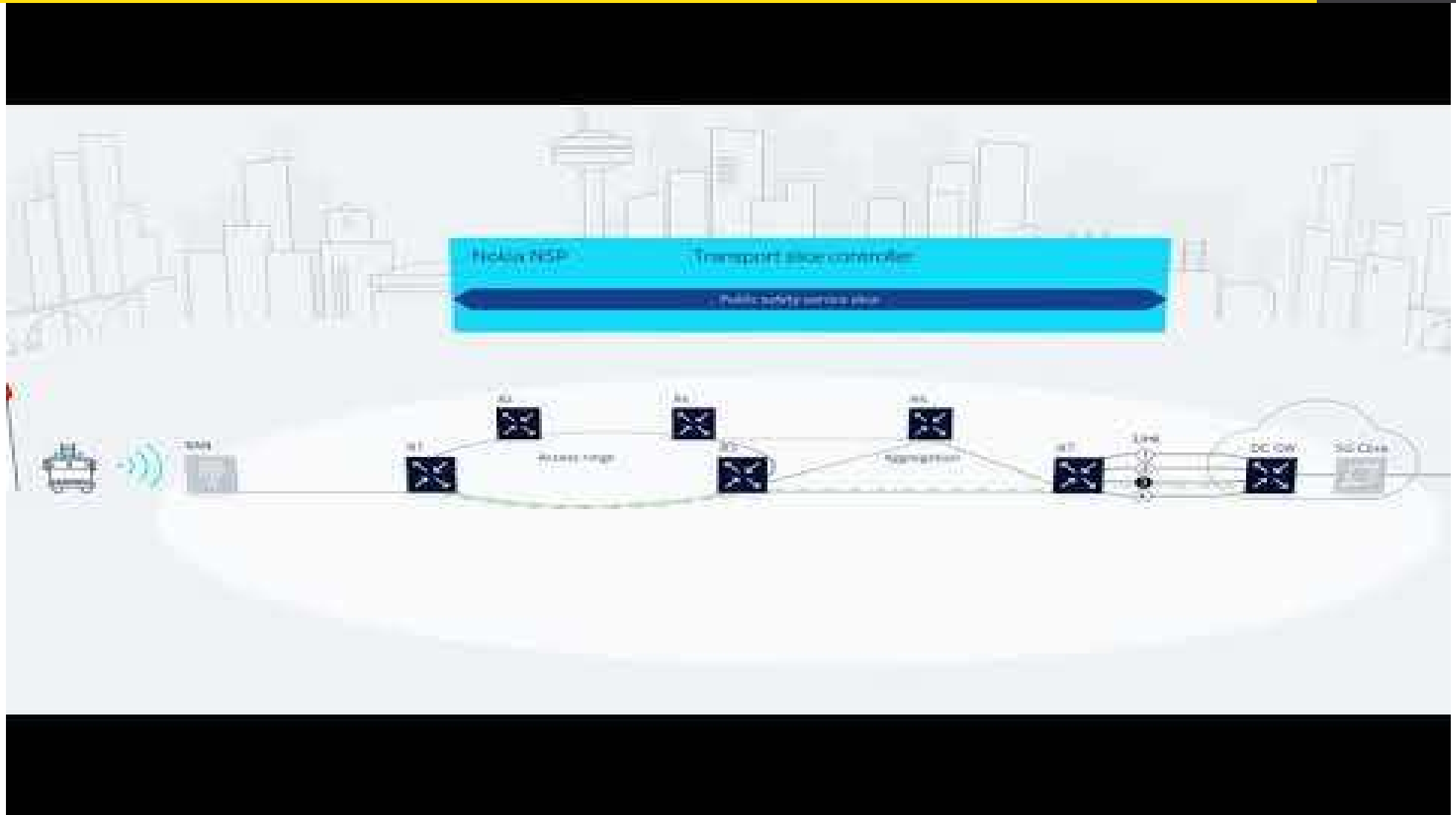


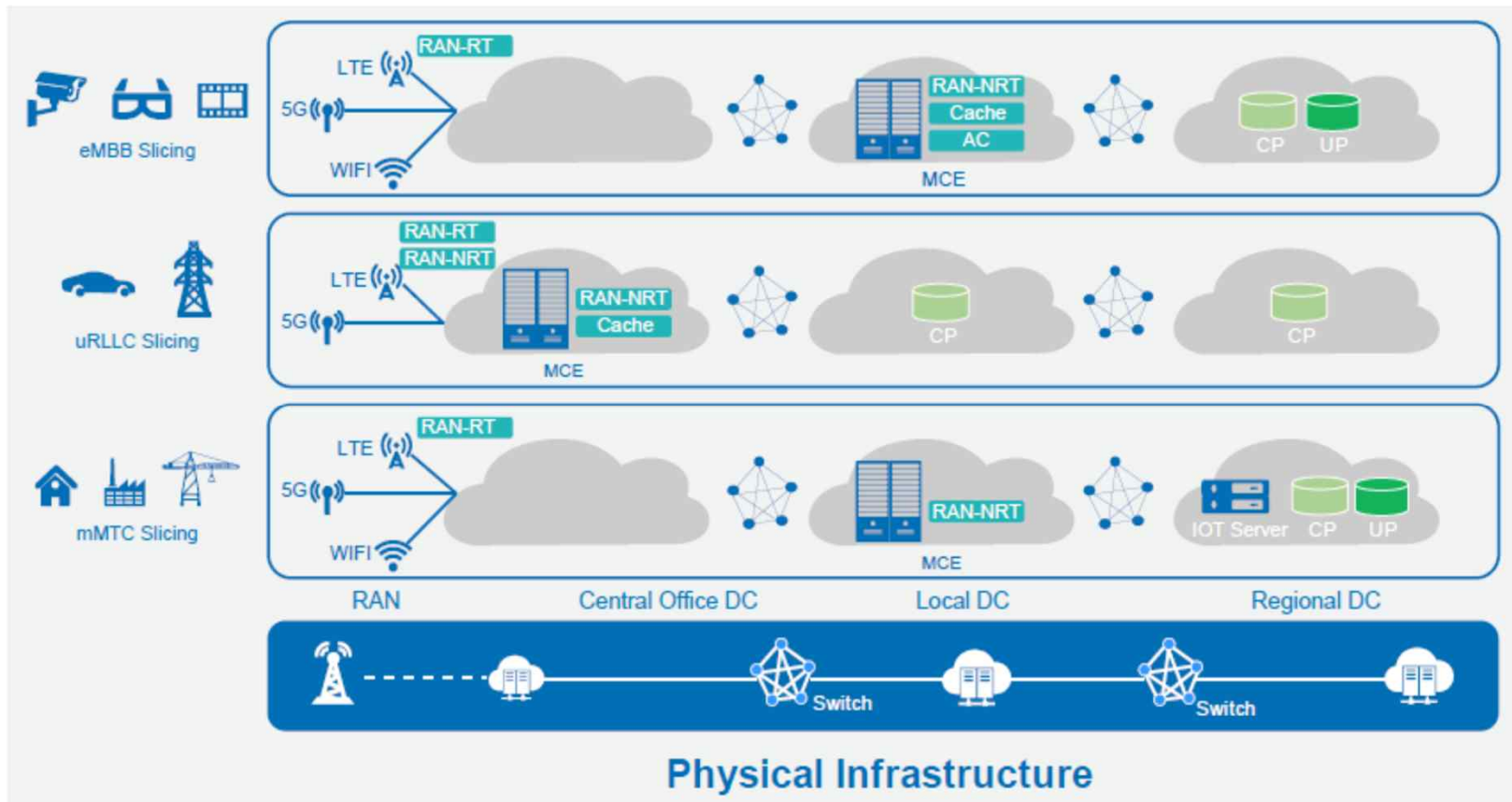


- **Network Slicing** certainly is one of the most discussed technologies these days. Network operators like KT, SK Telecom, China Mobile, DT, KDDI and NTT, and also vendors like Ericsson, Nokia and Huawei are all recognizing it as an ideal network architecture for the coming 5G era.
- **Ericsson** has been working on network slicing with **NTT DOCOMO** since 2014. In 2016 the two announced a successful **proof of concept of dynamic network slicing technology** for 5G core networks.
 - They created a slice management function and network slices based on requirements such as latency, security or capacity.
- **Samsung** and **KDDI** Complete 5G **End-to-End Network Slicing** Demonstration in **September 2020**



Network Slicing: Industrial Efforts (3)





RAN-RT: Radio Access Network-Real Time

RAN-NRT: Radio Access Network-non Real Time

AC: Access Cloud

CP: Control Plane

UP: User Plane

MCE: Mobile Cloud Engine

DC: Data Center

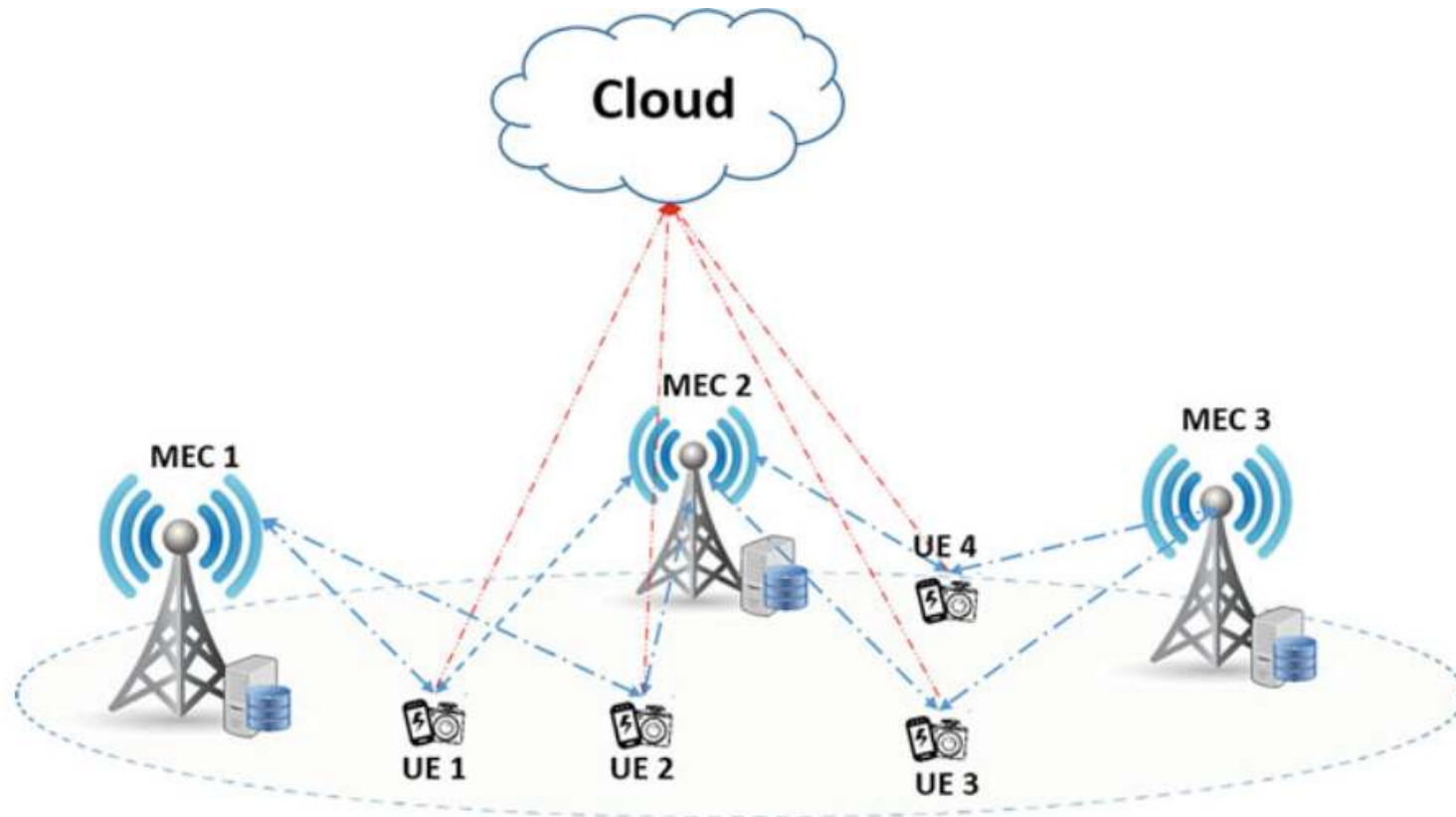
- Radio resources
 - Access network resources
 - Core network resources
- Caching
 - On-device caching
 - Edge caching
 - Core network caching
- Edge Computing Servers
 - Cloudlets
 - Fog servers
 - Multi-access edge computing servers

Use Case 1: Virtual Reality

- Introduction
- System Model
- Problem Formulation
- Solution Approach
- Simulation Results

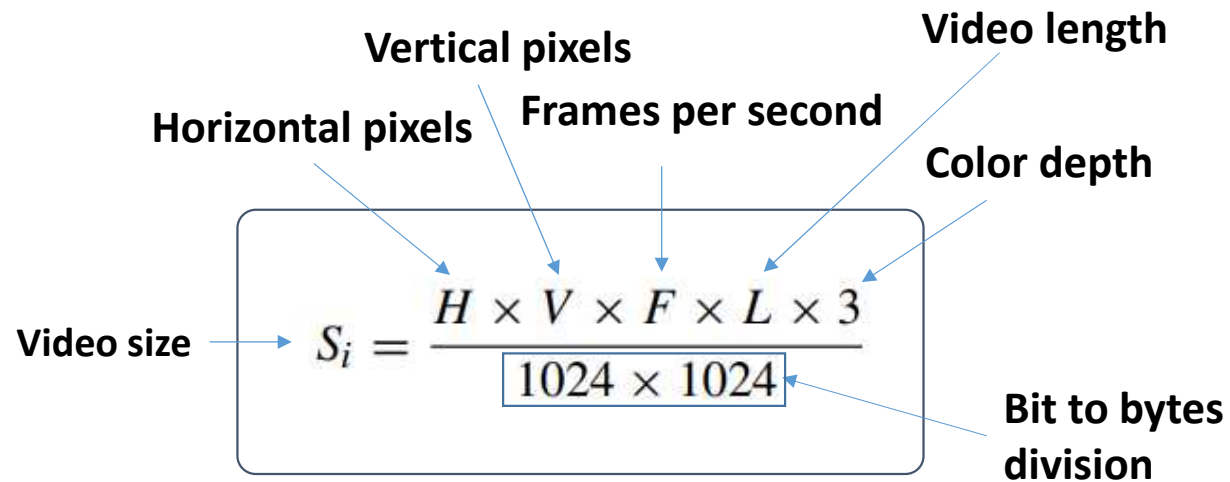
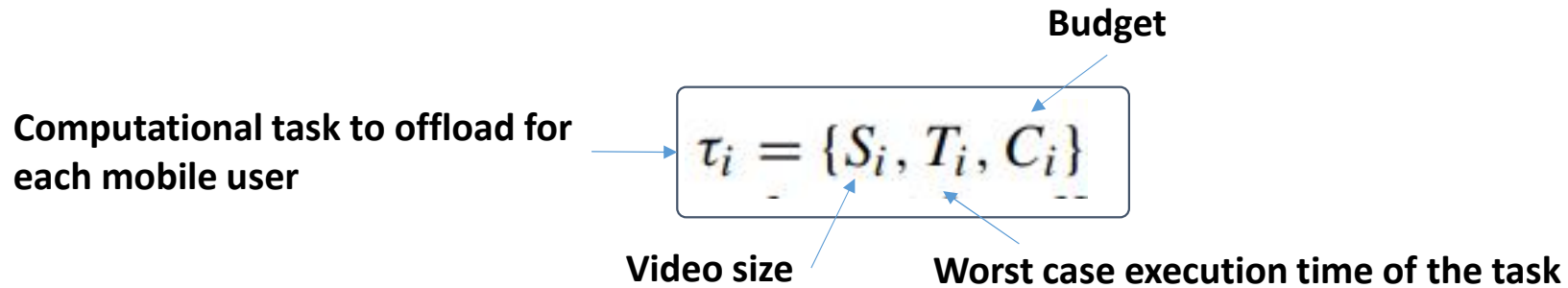
- Virtual Reality Applications
 - Smart Health-Care
 - Smart Industries
 - Smart Gaming

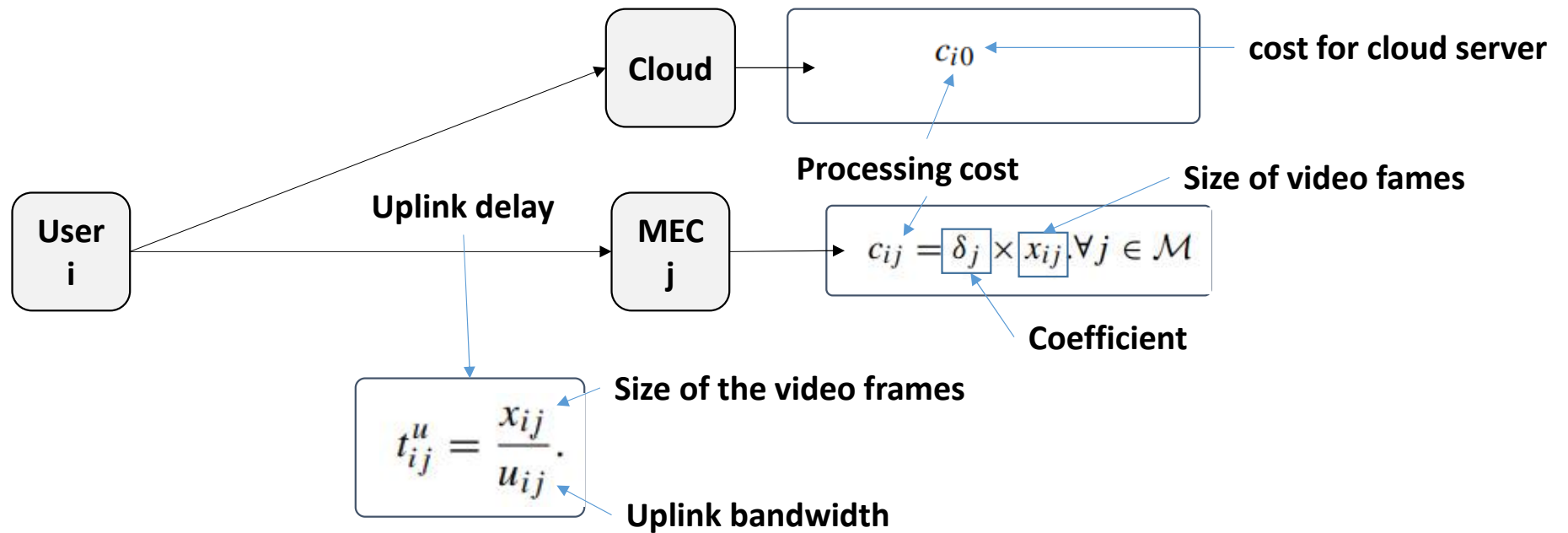
- Virtual Reality Challenges
 - High Computational Power for Processing complex Algorithms
 - Strict-Latency Constraints

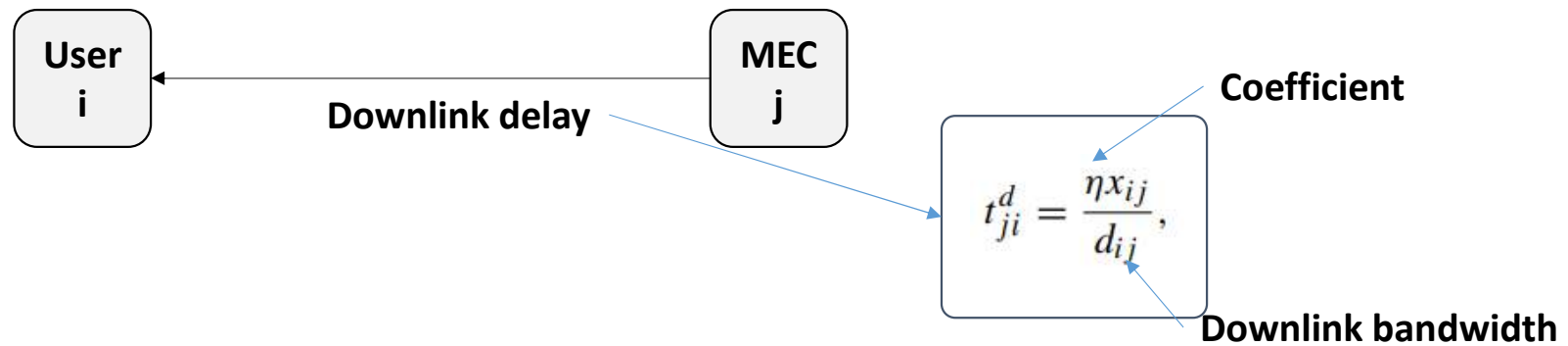


$\mathcal{N} = \{1, \dots, N\}$ be the set of N mobile users with VR capability

$\mathcal{M} = \{1, \dots, M\}$ be the set of M MEC servers.







Total transmission and processing cost minimization problem.

Uplink delay ← $\left(\frac{\sum_{j \in M} t_{ij}^u}{t_{i0}^u} + \frac{\sum_{j \in M} t_{ji}^d}{t_{i0}^d} \right)$ → Downlink delay

Processing cost ← $(1 - \alpha) \frac{\sum_{j \in M} c_{ij}}{c_{i0}}$

$$\underset{x}{\text{minimize}}: \sum_{i \in N} \alpha \left(\frac{\sum_{j \in M} t_{ij}^u}{t_{i0}^u} + \frac{\sum_{j \in M} t_{ji}^d}{t_{i0}^d} \right) + (1 - \alpha) \frac{\sum_{j \in M} c_{ij}}{c_{i0}} \quad (3.10)$$

$$\text{subject to: } \sum_{j \in M} x_{ij} = \frac{H_i \times V_i \times F_i \times L_i \times 3}{1024 \times 1024}, \forall i \in N, \quad (3.11) \quad \text{Guarantees that all users are served by network.}$$

$$\sum_{i \in N} x_{ij} \leq \Gamma_j, \forall j \in M, \quad (3.12) \quad \text{Maximum MEC server processing capacity}$$

$$\sum_{j \in M} t_i \leq T_i, \forall i \in N, \quad (3.13) \quad \text{Latency limit constraint}$$

$$\sum_{j \in M} c_i \leq C_i, \forall i \in N, \quad (3.14) \quad \text{Budget cost must not be greater than processing cost}$$

$$t'_i < 1, \forall i \in N, \quad (3.15) \quad \text{Total time constraint variable}$$

$$c'_i < 1, \forall i \in N, \quad (3.16) \quad \text{Total cost constraint variable}$$

$$x_{ij} \geq 0, \forall i \in N, \forall j \in M. \quad (3.17) \quad \text{Video frames size variable}$$

Re-write the objective function.

$$(3.10) = \alpha \left(\sum_{j \in M} \left(\frac{x_{ij}}{u_{ij}t_{i0}^u} + \frac{\eta x_{ij}}{d_{ij}t_{ji}^d} \right) \right) + (1 - \alpha) \left(\frac{\sum_{j \in M} \delta_j x_{ij}}{c_{i0}} \right) \quad (3.18)$$

$$= \alpha \left(\sum_{j \in M} \left(\frac{1}{u_{ij}t_{i0}^u} + \frac{\eta}{d_{ij}t_{ji}^d} \right) \right) x_{ij} + (1 - \alpha) \left(\frac{\sum_{j \in M} \delta_j}{c_{i0}} \right) x_{ij} \quad (3.19)$$

$$= \left(\alpha \sum_{j \in M} \left(\left(\frac{1}{u_{ij}t_{i0}^u} + \frac{\eta}{d_{ij}t_{ji}^d} \right) + (1 - \alpha) \left(\frac{\sum_{j \in M} \delta_j}{c_{i0}} \right) \right) \right) x_{ij} \quad (3.20)$$

$$= f_i(\mathbf{x}_i) \quad (3.21)$$

ADMM: Alternating Direction Method of Multipliers

Modified Problem

$$\underset{x}{\text{minimize}} : \sum_{i \in N} f_i(\mathbf{x}_i) \quad (3.22)$$

$$\text{subject to} : \mathbf{1}^T \mathbf{x}_i = S_i, \forall i \in N \quad (3.23)$$

$$\mathbf{1}^T \mathbf{x}_j \leq \Gamma_j, \forall j \in M \quad (3.24)$$

$$\sum_{j \in M} t_j \leq T_i, \forall i \in N \quad (3.25)$$

$$\sum_{j \in M} c_j \leq C_i, \forall i \in N \quad (3.26)$$

$$t'_i < 1, \forall i \in N \quad (3.27)$$

$$c'_i < 1, \forall i \in N \quad (3.28)$$

$$x_{ij} \geq 0, \forall i \in N, \forall j \in M \quad (3.29)$$

Guarantees that all users are served by network.

Maximum MEC server capacity constraint

Latency limit constraint

Budget cost must not be greater than processing cost

Total time constraint variable

Total cost constraint variable

Video frames size variable

For ADMM-based solution new variable z is introduced.

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} : \sum_{i \in N} f_i(\mathbf{x}_i) + h(z) \\
 & \text{subject to} : \mathbf{x}_i = z \\
 & \quad \quad \quad \mathbf{x}_i \in \mathcal{X}, \forall i \in N
 \end{aligned} \tag{3.31}$$

where $h(z) = 0$ when $\mathbf{x}_i \in \mathcal{X}$.

$$h(z) = I_{\mathcal{X}}(z) = \begin{cases} 0, & \mathbf{x}_i \in \mathcal{X} \\ \infty, & \text{otherwise} \end{cases} \tag{3.32}$$

ADMM-Based Solution (3)

For ADMM-based solution new variable z is introduced.

Then, the augmented Lagrangian function of (3.31) is as follows:

$$\mathcal{L}(\mathbf{x}, z, \lambda) = \sum_{i \in N} \left(f_i(\mathbf{x}_i) + \lambda_i^T (x_i - z) + \frac{\rho}{2} \|x_i - z\|_2^2 \right) \quad (3.33)$$

Lagrangian Lagrangian penalty term Augmentation

Lagrange multiplier penalty parameter

Based on the solution from [43], the resulting ADMM variables update are the following:

$$x_i^{k+1} = \arg \min \left(f_i(x_i) + \lambda_i^{kT} (x_i - z^k) + \frac{\rho}{2} \|x_i - z^k\|_2^2 \right) \quad (3.34)$$

$$z^{k+1} = \arg \min \left(h(z) + \sum_{i=1}^N \left(-\lambda_i^{kT} z + \frac{\rho}{2} \|x_i^{k+1} - z\|_2^2 \right) \right) \quad (3.35)$$

Lagrange multiplier \rightarrow $\lambda_i^{k+1} = \lambda_i^k + \rho(x_i^{k+1} - z^{k+1}) \quad (3.36)$

ADMM-Based Task Offloading Algorithm.

Algorithm 1 ADMM-based task offloading

- 1: **input:** Initialization for $\mathcal{N}, \mathcal{M}, \mathbf{D}, \mathbf{U}$
 - 2: **Output:** Minimal offloading cost
 - 3: Initialization
 - 4: $max_iteration = 1000, \rho = 0.5, \alpha = 0.5, \mathbf{x}_i^0 \geq 0, \lambda_i^0 \geq 0, z \geq 0, t_{i0}^u, t_{i0}^d, c_{i0}, \forall i \in \mathcal{N}$
 - 5: **for** $k \in max_iteration$ **do**
 - 6: Each user $i \in \mathcal{N}$ update its offloading decision by Eqs. (3.34), (3.35), (3.32) respectively, parallelly
 - 7: After getting updated values from all users each MEC will update λ using (3.36), parallelly
 - 8: After all variable updated, update the objective function (3.10)
 - 9: **end for**
- return** Optimal value of objective (3.10)
-

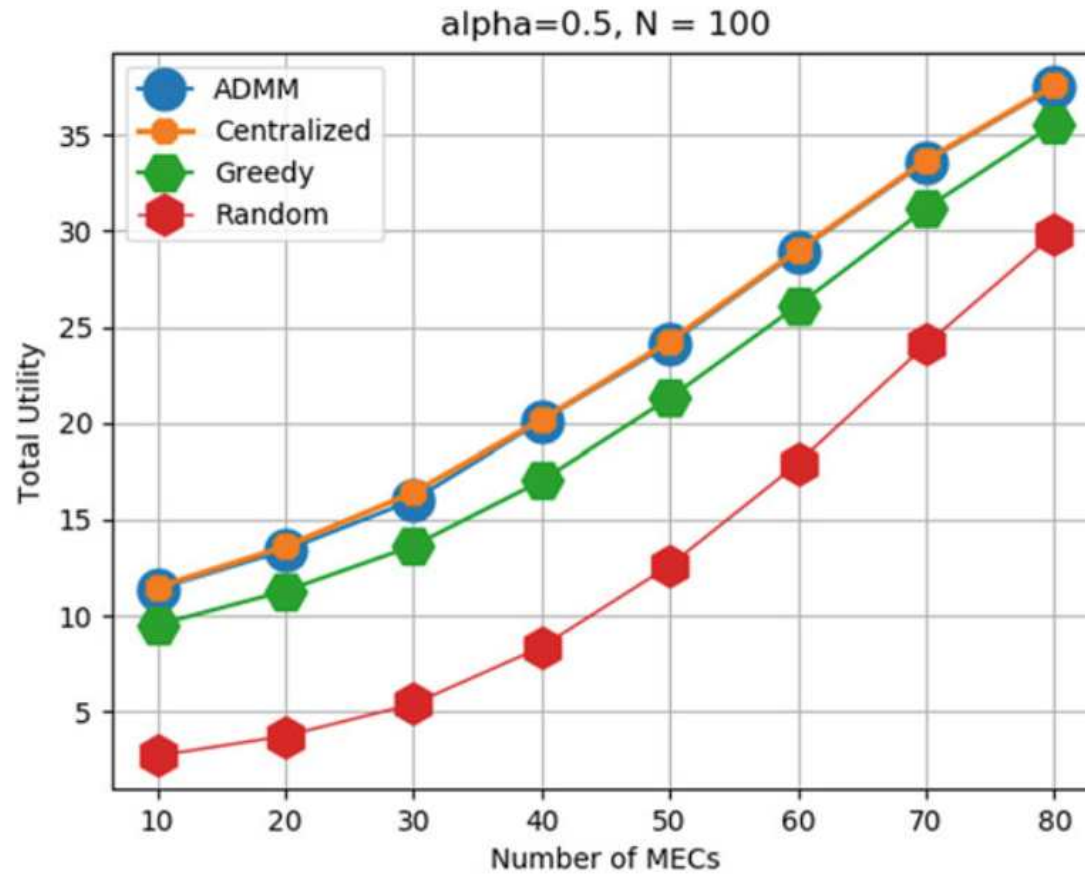
$$x_i^{k+1} = \arg \min \left(f_i(x_i) + \lambda_i^{kT} (x_i - z^k) + \frac{\rho}{2} \|x_i - z^k\|_2^2 \right) \quad (3.34)$$

$$z^{k+1} = \arg \min \left(h(z) + \sum_{i=1}^N \left(-\lambda_i^{kT} z + \frac{\rho}{2} \|x_i^{k+1} - z\|_2^2 \right) \right) \quad (3.35)$$

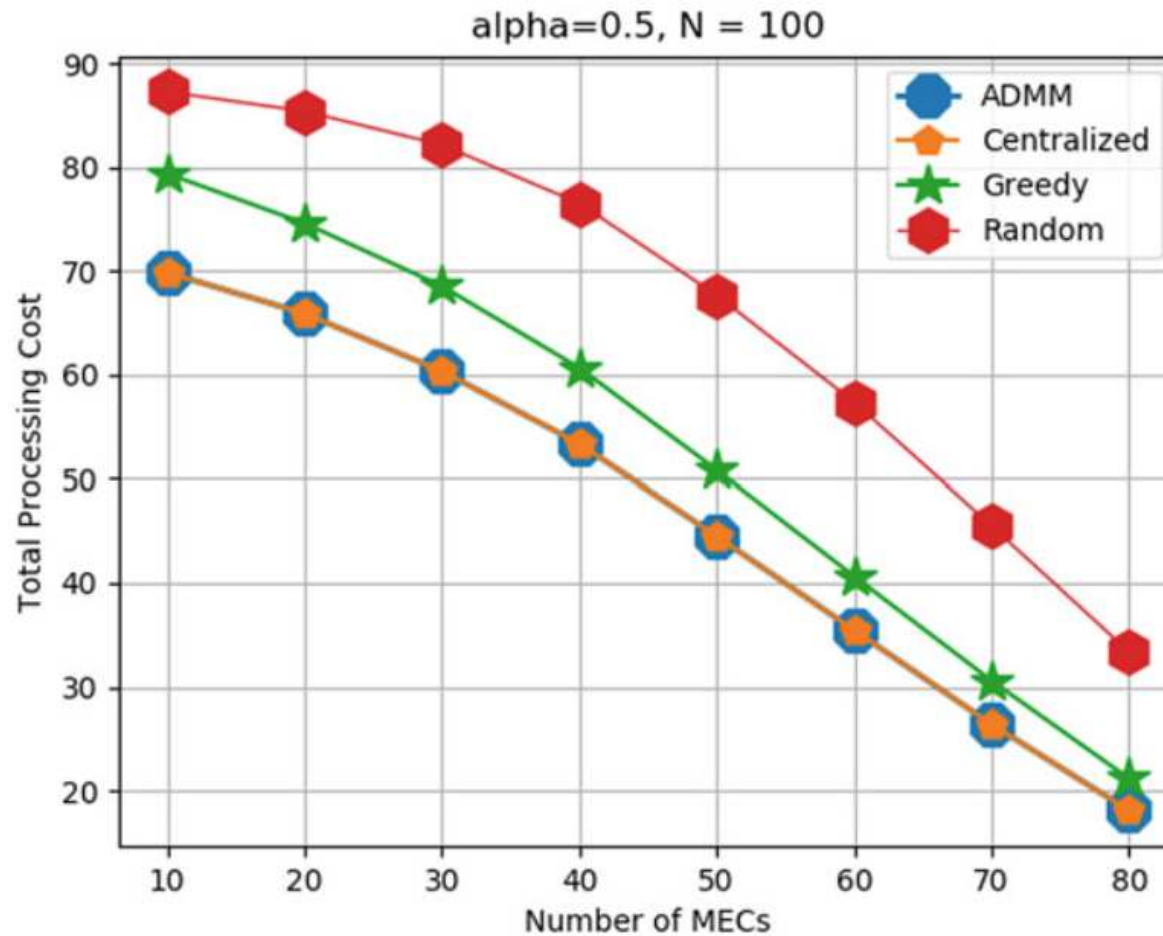
$$h(z) = I_{\mathcal{X}}(z) = \begin{cases} 0, & \mathbf{x}_i \in \mathcal{X} \\ \infty, & otherwise \end{cases} \quad (3.32)$$

$$\lambda_i^{k+1} = \lambda_i^k + \rho(x_i^{k+1} - z^{k+1}) \quad (3.36)$$

$$\underset{x}{\text{minimize}} : \sum_{i \in \mathcal{N}} \alpha \left(\frac{\sum_{j \in \mathcal{M}} t_{ij}^u}{t_{i0}^u} + \frac{\sum_{j \in \mathcal{M}} t_{ji}^d}{t_{i0}^d} \right) + (1 - \alpha) \frac{\sum_{j \in \mathcal{M}} c_{ij}}{c_{i0}} \quad (3.10)$$



Total utility vs. number of MEC servers



Total Processing cost vs. number of MEC servers

- An overview of resource management for network slicing has been presented in this chapter.
- Numerous key resources for network slicing are discussed.
- A use case of virtual reality is along with its ADMM-based solution is presented.

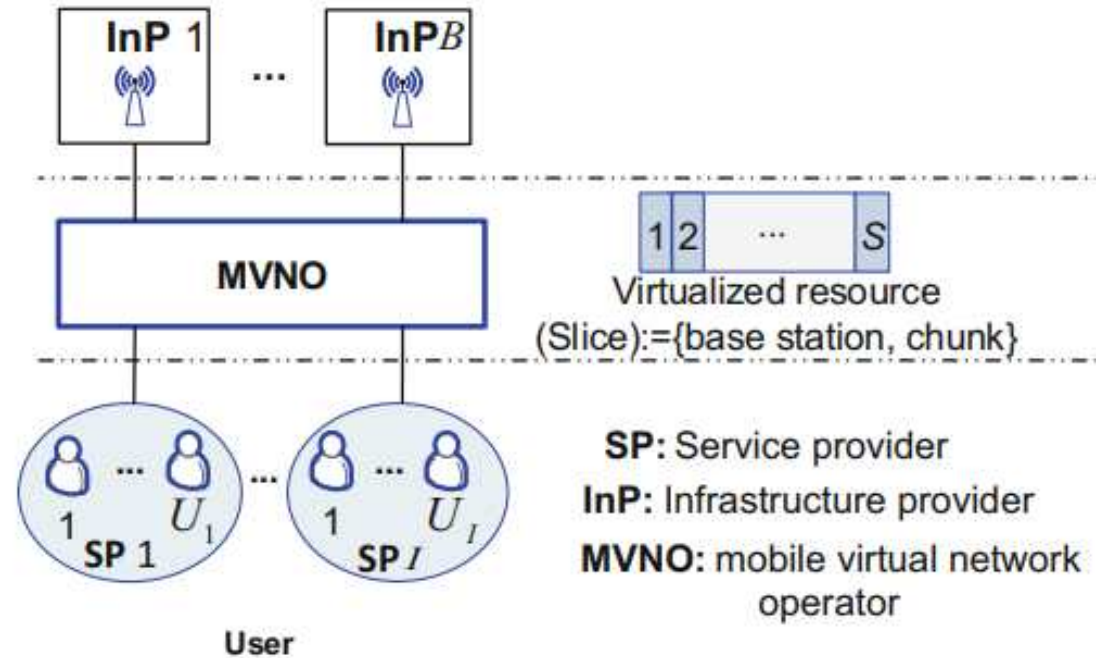
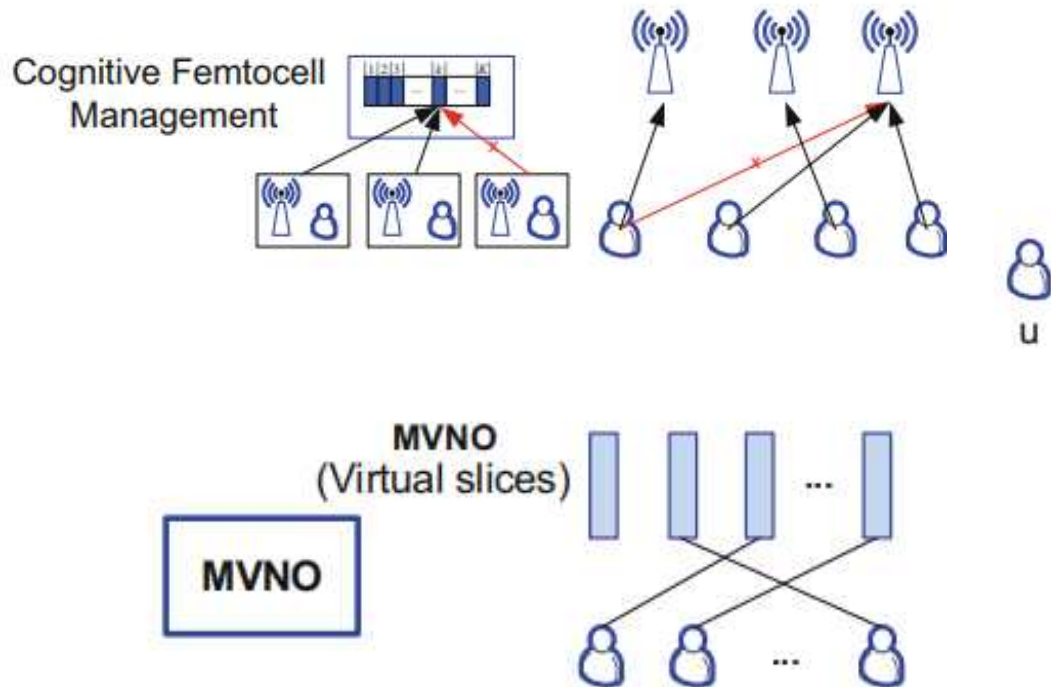
Use Case 2: Chunk-Based Resource Allocation

- Introduction
- System Model
- Problem Formulation
- Solution Approach
- Simulation Results

- Traditionally in cellular networks radio resources were only considered as a performance bottleneck
- A number of solutions were devised to only cater the radio resource allocation challenge
- The proliferation of end users and novel applications have also imposed limitations on other network resources such as backhaul and cache spaces

System Model

Set of Infrastructure provider $\mathcal{B} = \{1, 2, \dots, B\}$
 Set of Service provider $\mathcal{I} = \{1, 2, \dots, I\}$
 Set of user $U^i = \{1, 2, \dots, U_i\} \quad i \in \mathcal{I}$

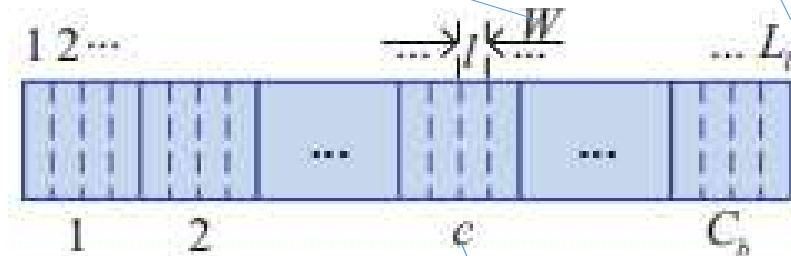


SP: Service provider
InP: Infrastructure provider
MVNO: mobile virtual network operator

Virtualized cellular network hierarchical model

The bandwidth of narrowband orthogonal subcarriers is W

Subcarriers set $\mathcal{L}_b = \{1, 2, \dots, L_b\}$,



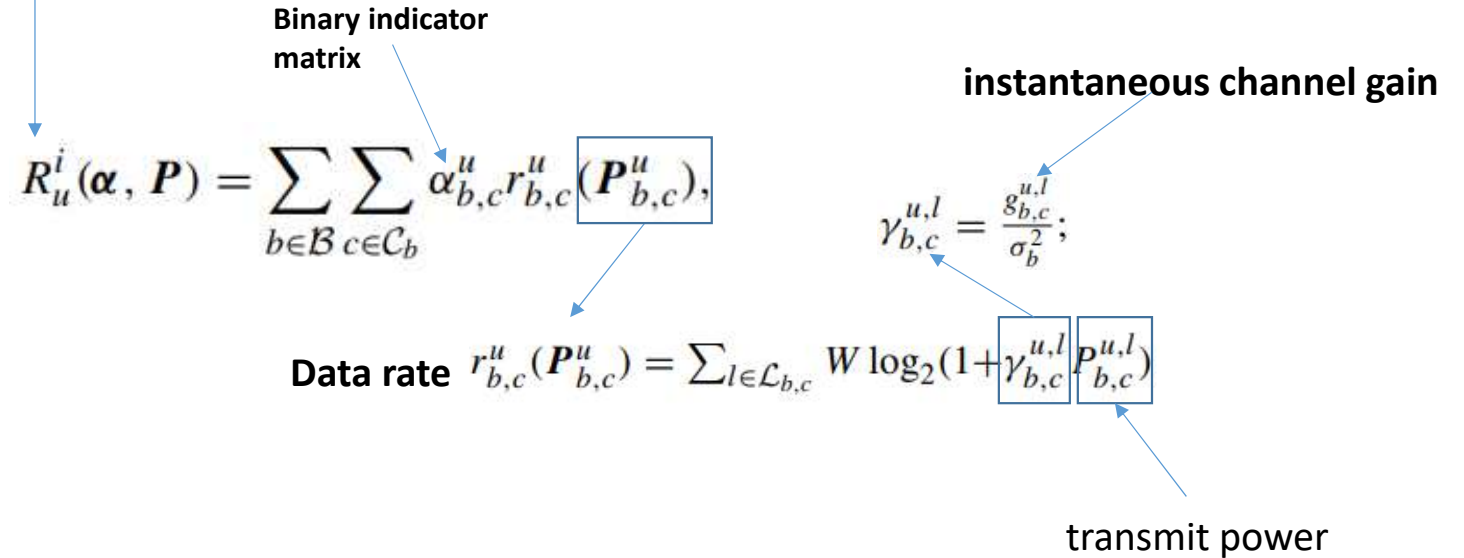
$\mathcal{C}_b = \{1, 2, \dots, C_b\}$ chunks

Chunk-based radio resources of $\ln P_b$

Every chunk is formed by aggregating $\mathcal{L}_{b,c} = \{1, 2, \dots, L_b/C_b\}$ subcarriers

Data rate of a user u associated with service provider i

\mathcal{B}	Set of InPs
\mathcal{I}	Set of virtual resources given by MVNO to SPs
\mathcal{U}_m	Set of customers connected with SPs
\mathcal{L}_b	Set of radio resources subcarriers owned by InPs
\mathcal{C}_b	Set of subcarriers chunks
α	Slice allocation matrix
$\alpha_{b,c}^\mu$	Binary indicator matrix
R_μ^i	Data rate
σ_b^2	Background noise
W	Bandwidth
$P_{b,c}^\mu$	Transmission power vector
$Z_{b,bh}$	Predefined Backhaul capacity
ϕ_i^{SP}	The payment (in units/Mbps) of each SP i to the MVNO
ϕ_b^{bh}	Unit price (in units/Mbps) of the Backhaul set by InP
$L_{\alpha,P,\lambda,\beta}$	Lagrangian function
λ	Lagrangian nonnegative multiplier
β	Lagrangian nonnegative multiplier
μ	Lagrangian nonnegative multiplier
$\omega_{b,c}^\mu$	Used in Lagrangian dual function
ϕ_u^k	Utility function



$$U_{MVNO}(\alpha, P) = U^{\text{rev}}(\alpha, P) - U^{\text{cost}}(\alpha, P), \quad (6.4)$$

$$\max_{(\alpha, \mathbf{P})} U_{\text{MVNO}}(\alpha, \mathbf{P}) \quad (6.5)$$

$$\text{s.t. } R_u^i(\alpha, \mathbf{P}) \geq R_{u,\min}^i, \quad \forall u \in \mathcal{U}, \forall i \in \mathcal{I}. \quad (6.2) \text{ (Guaranteeing the required minimum rate)}$$

$$\sum_{i \in \mathcal{I}} \sum_{u \in \mathcal{U}; \alpha_{i,u}^u = 1} R_u^i(\alpha, \mathbf{P}) \leq Z_{b,\text{bh}}, \quad \forall b \in \mathcal{B}, \forall c \in \mathcal{C}_b, \quad (6.3) \text{ (Aggregated data rate of users)}$$

$$\sum_{b \in \mathcal{B}} \sum_{c \in \mathcal{C}_b} \alpha_{b,c}^u \sum_{l \in \mathcal{L}_c} P_{b,c}^{u,l} \leq \bar{P}_u, \quad \forall u \in \mathcal{U}, \quad (6.6) \text{ (Total transmit power)}$$

$$P_{b,c}^{u,l} \geq 0, \quad \forall b \in \mathcal{B}, \forall c \in \mathcal{C}_b, \forall u \in \mathcal{U}, \quad (6.7)$$

$$\alpha_{b,c}^u \in \Pi_\alpha, \quad \forall b \in \mathcal{B}, \forall c \in \mathcal{C}_b, \forall u \in \mathcal{U}, \quad (6.8)$$

$$\sum_{u \in \mathcal{U}} \alpha_{b,c}^u \leq 1, \quad \forall c \in \mathcal{C}_b, \forall b \in \mathcal{B}, \quad (6.9) \text{ (Restricts the allocation of a slice to at most user)}$$

$$\sum_{b \in \mathcal{B}} \sum_{c \in \mathcal{C}_b} \alpha_{b,c}^u \leq 1, \quad \forall u \in \mathcal{U}, \quad (6.10) \text{ (Isolation of the slices)}$$

$$\sum_{u \in \mathcal{U}} \sum_{c \in \mathcal{C}_b} \alpha_{b,c}^u \leq 1, \quad \forall b \in \mathcal{B}, \quad (6.11) \text{ (Isolation of the slices)}$$

$$\sum_{u \in \mathcal{U}} \sum_{b \in \mathcal{B}} \alpha_{b,c}^u \leq 1, \quad \forall c \in \mathcal{C}_b, \quad (6.12) \text{ (Isolation of the slices)}$$

$$\alpha_{b,c}^u \in [0, 1], \quad \forall u \in \mathcal{U}, \forall b \in \mathcal{B}, \forall c \in \mathcal{C}_b. \quad (6.13)$$



The Lagrangian function of (6.5) is as follows

$$\begin{aligned}
 L(P, \lambda, \mu, \beta) = & U_{MVN}(\alpha, P) + \\
 & \lambda_i (R_u^i(\lambda, \beta) - R_{u,m}^i) + \\
 & \beta_i \left(Z_{b,bh} - \sum_{i \in I} \sum_{u \in U} R_u^i(\alpha, \beta) \right) + \\
 & \mu_i (\bar{P}_u - \sum_{b \in B} \sum_{c \in C_b} \alpha_{b,c}^u \sum_{l \in L_c} P_{b,c}^{u,l})
 \end{aligned}$$

Algorithm 1 JSPA-HSA: JSPA with Hungarian-based slice allocation

1: Initialization: $\mathcal{I}, \mathcal{B}, \mathcal{C}_b, \mathcal{U}_i, \mathbf{P}^{(0)}, \boldsymbol{\lambda}^{(0)}, \boldsymbol{\mu}^{(0)}$, and $\boldsymbol{\beta}^{(0)}$.

2: Repeat:

3: **Power allocation phase:**

4: *At the subscribed user u :

5: Update λ_u as: (Transmission rate)

$$\lambda_u(t+1) = [\lambda_u(t) - s_1(t)(R_u^i(\boldsymbol{\alpha}, \mathbf{P}) - R_u^{\min})]^+; \quad (6.18)$$

6: Update μ_u as: (Transmit power)

$$\mu_u(t+1) = \left[\mu_u(t) - s_2(t) \left(\sum_{b \in \mathcal{B}} \sum_{c \in \mathcal{C}_b} \alpha_{b,c}^u \sum_{l \in \mathcal{L}_c} P_{b,c}^{u,l} - \bar{P}_u \right) \right]^+; \quad (6.19)$$

7: Update transmit power $P_{b,c}^{u,l}(t+1)$ by $P_{b,c}^{u,l*} = \left[\frac{\varphi_i^{\text{sp}} - \varphi_b^{\text{bh}} + \lambda_u - \beta_b}{(\ln 2/W)\mu_u} - \frac{1}{\gamma_{b,c}^{u,l}} \right]^+$, (6.16)

8: *At the SBS b :

9: Update congested backhaul link price $\beta_b(t+1)$:

(Backhaul data rate) $\beta_b(t+1) = \left[\beta_b(t) + s_3(t) \left(\sum_{i \in \mathcal{I}} \sum_{u \in \mathcal{U}_i} R_u^i(\boldsymbol{\alpha}, \mathbf{P}) - Z_{b,\text{bh}} \right) \right]^+; \quad (6.20)$

10: **Slice allocation phase:**

11: *At the MVNO:

12: Update $\alpha_{b,c}^u(t+1)$ using the Hungarian algorithm to maximize

$$\max_{(\boldsymbol{\alpha}, \mathbf{P})} \sum_{i \in \mathcal{I}} \sum_{u \in \mathcal{U}_i} \sum_{b \in \mathcal{B}} \sum_{c \in \mathcal{C}_b} \alpha_{b,c}^u \left[\Omega_{b,c}^u(\mathbf{P}_{b,c}^{u*}) - \varphi_{b,c}^{\text{slice}} \right] \quad (6.17)$$

13: Until $|\lambda_u(t+1) - \lambda_u(t)| \leq \epsilon_1$, $|\mu_u(t+1) - \mu_u(t)| \leq \epsilon_2$, and $|\beta_b(t+1) - \beta_b(t)| \leq \epsilon_3$ are simultaneously satisfied.

Lagrangian multiplier

(Backhaul data rate)

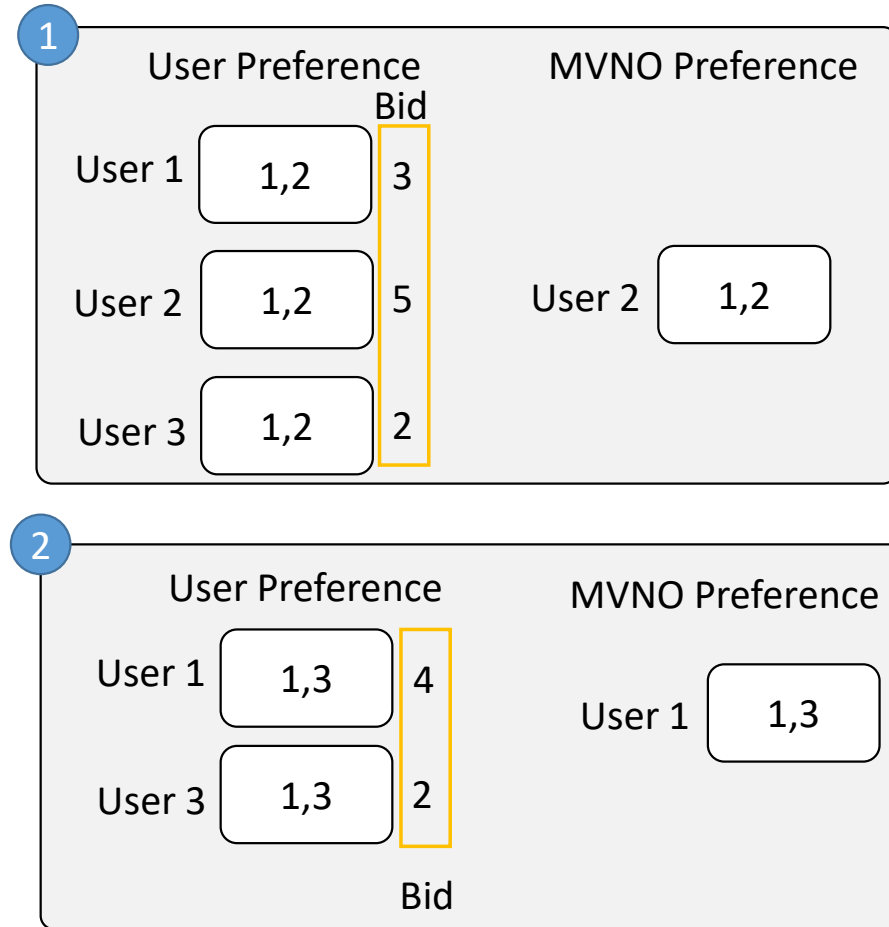
Solution Approach-2

Matching-Based Low-Complexity Algorithm

Algorithm 2 MSA: matching-based slice allocation

- 1: **while** $\sum_{u, \{b,c\}} b_{u \rightarrow \{b,c\}} \neq 0$ or convergence not achieved **do**
- 2: *At the subscribed users:*
- 3: Send a bid for the slice $\{b, c\}^* = \arg \max_{\{b,c\} \in \succ_u} \phi_u(\{b, c\})$.
- 4: *At the MVNO:* Users side
- 5: Construct $\succ_{\{b,c\}}$ based on (6.21).
- 6: Update $\{b, c\}^* = \mu(\{b, c\}) | u^* = \arg \max_{u \in \succ_{\{b,c\}}} \phi_{\{b,c\}}(u)$.
- 7: Update the rejected user lists on the slices and the preference \succ_u .
- 8: **end while** MVNO side

$$\phi_u(k) = \Omega_{b,c}^u(\mathbf{P}_{b,c}^u) - \varphi_{b,c}^{\text{slice}}. \tag{6.21}$$



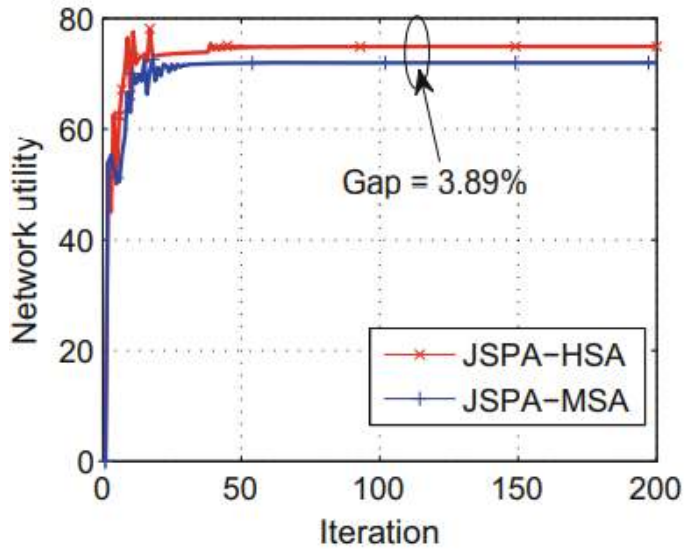


Fig. 6.3 Evaluation results of Network utility with $Z_{b,bh} = 10$ Mbps

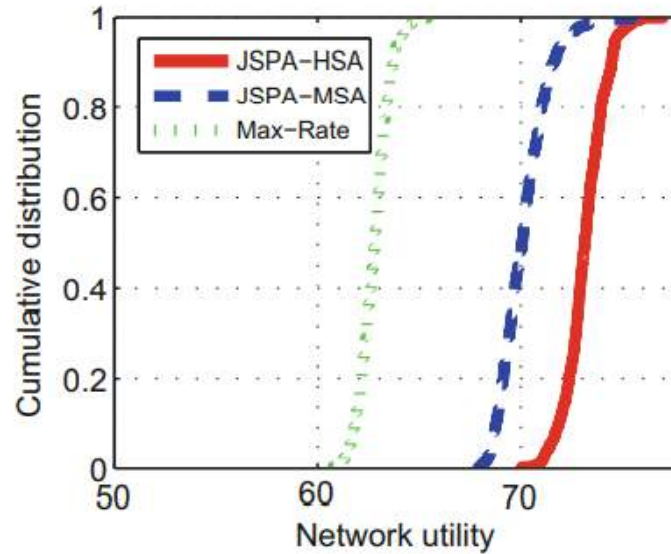


Fig. 6.4 CDF of network utilities

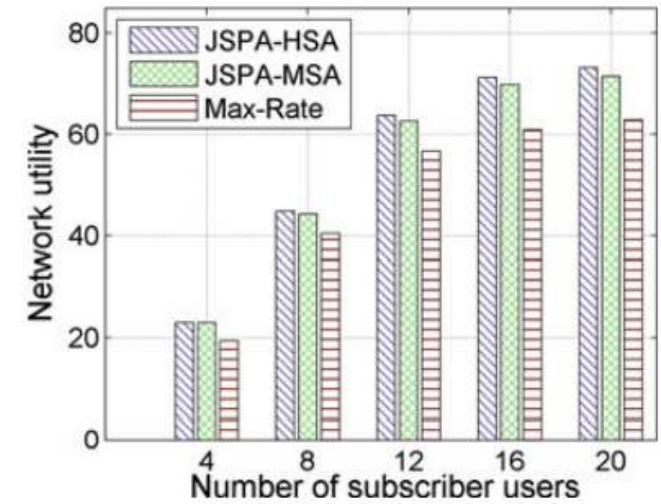
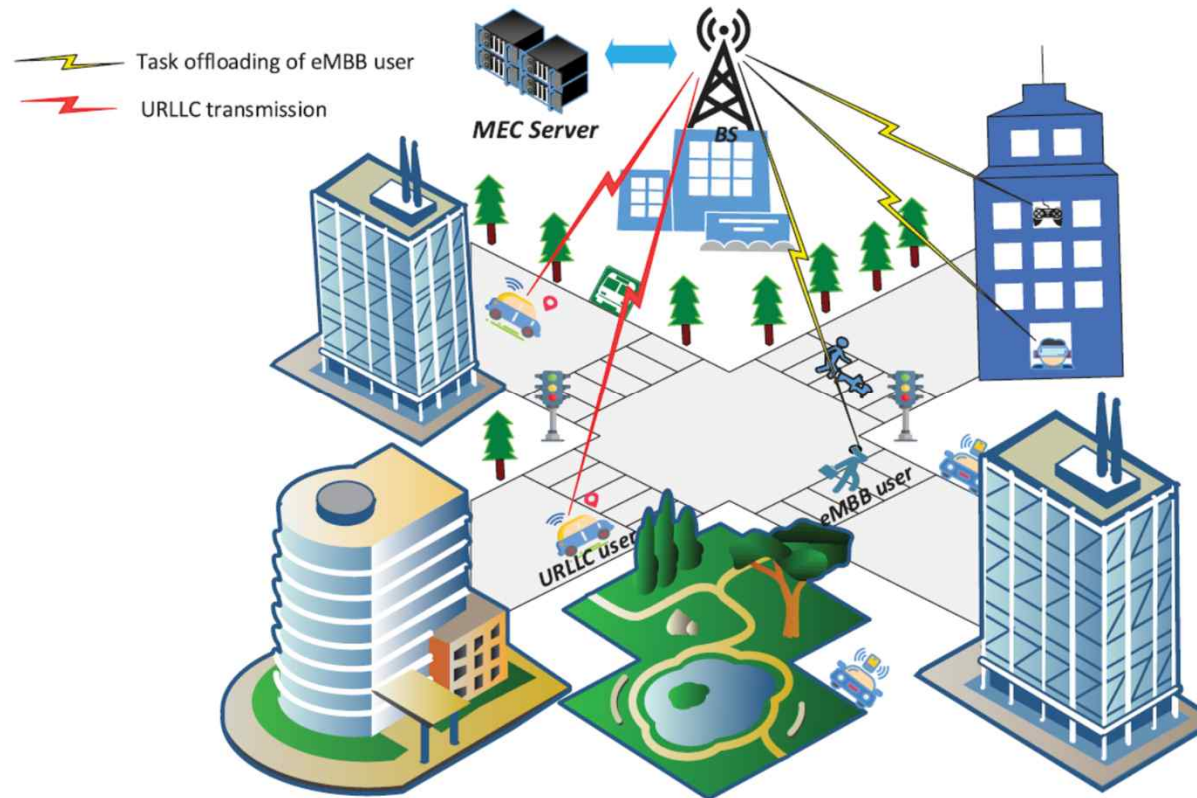


Fig. 6.5 Network utility versus number of subscriber users

Use Case 3: Energy Efficient Communication and Computation Resource Slicing for eMBB and URLLC Coexistence in 5G and Beyond

- System Model
- Problem Formulation
- Solution Approach
- Simulation Results



Yan Kyaw Tun, Do Hyon Kim, Madyan Alsenwi, Nguyen H. Tran, Zhu Han, Choong Seon Hong, "Energy Efficient Communication and Computation Resource Slicing for eMBB and URLLC Coexistence in 5G and Beyond," IEEE Access, Vol.8, pp.136024-136035, Jul. 2020

Minimize the overall energy consumption

Energy Consumption of eMBB users

Data rate of eMBB user

weight

$$\min_{y,l,w} \left(\sum_{u=1}^U E_u^{\text{Off}} + \sum_{u=1}^U E_u^L \right) - \phi \sum_{u=1}^U \sum_{b=1}^B y_u^b R_{u,b} \quad (18)$$

$$\text{s.t. C1 : } \frac{l_u}{R_u} + \frac{c_u l_u}{f_u^C} \leq T_u, \quad u \in \mathcal{U}, \quad (19)$$

execution latency constraint of the eMBB users

$$\text{C2 : } \frac{c_u(d_u - l_u)}{f_u^l} \leq T_u, \quad u \in \mathcal{U}, \quad (20)$$

execution latency constraint of the eMBB users

$$\text{C3 : } l_u \leq d_u, \quad \forall u \in \mathcal{U}, \quad (21)$$

offloading data size of user u has to be less than the total input data size

$$\text{C4 : } 0 \leq w_u \leq 1, \quad \forall u \in \mathcal{U}, \quad (22)$$

weight parameter for eMBB user u that can be punctured by traffic of URLLC users

$$\text{C5 : } \text{CVaR}_\beta(R) \leq \alpha, \quad (23)$$

reliability constraints of URLLC users

$$\text{C6 : } \Pr[R_{urllc} \leq L] \leq \epsilon, \quad (24)$$

reliability constraints of eMBB users

$$\text{C7 : } \sum_{u=1}^U y_u^b \leq 1, \quad \forall b \in \mathcal{B}, \quad (25)$$

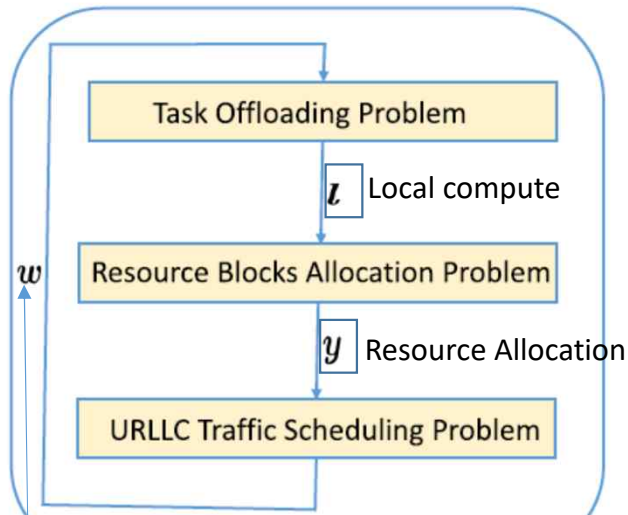
one resource block can be allocated to only one eMBB user

$$\text{C8 : } y_u^b \in \{0, 1\}, \quad \forall u \in \mathcal{U}, \forall b \in \mathcal{B}, \quad (26)$$

Notation	Definition
\mathcal{U}	Set of eMBB users, $ \mathcal{U} = U$
F	Total system bandwidth
F_b	Fraction of system bandwidth allocated to eMBB users
F_s	Fraction of system bandwidth allocated to URLLC users
\mathcal{B}	Set of resource blocks, $ \mathcal{B} = B$
d_u	Total input data size of eMBB user u
c_u	Required CPU cycles to accomplish one bit of the input data of eMBB user u
T_u	The execution deadline of the task of eMBB user u
l_u	The offloaded data size of the task of eMBB user u
t_u^L	The local computation execution time of eMBB user u
E_u^L	The local computation energy of eMBB user u
y_u^b	Resource block assignment variable
M	Number of minislots divided in each resource block
L_m	Traffic of URLLC users at minislot m
L_{\max}	Maximum traffic of URLLC users that can be served at a time slot
w_u	Weight of puncturing eMBB user u
g_u^b	Achievable channel gain of eMBB user u
P_u^b	Transmit power of eMBB user u
$R_{u,b}$	Achievable data rate of eMBB user u on resource block b
t_u^{up}	The uplink transmission delay experienced by eMBB user u
f^C	The total CPU capacity of the MEC server
f_u^C	The CPU capacity of the MEC server that is allocated to eMBB user u
E_u^{Off}	The energy consumption of eMBB user u for offloading data
R_s	The achievable data rate of URLLC user s
P_s	The transmit power of URLLC user s
g_s	The achievable channel gain of URLLC user s
R_{urllc}	The total achievable data rate of URLLC users

Yan Kyaw Tun, Do Hyon Kim, Madyan Alsenwi, Nguyen H. Tran, Zhu Han, Choong Seon Hong, "Energy Efficient Communication and Computation Resource Slicing for eMBB and URLLC Coexistence in 5G and Beyond," IEEE Access, Vol.8, pp.136024-136035, Jul. 2020

BLOCK COORDINATE DESCENT BASED SOLUTION



weight parameter for eMBB user u that can be punctured by traffic of URLLC users

Optimization framework

$$(P1) : \min_l \sum_{u=1}^U E_u^{\text{Off}} + \sum_{u=1}^U E_u^L \quad (30)$$

s.t. (C1)-(C3), (31)

$$(P2) : \min_y \sum_{u=1}^U E_u^{\text{Off}} - \phi \sum_{u=1}^U \sum_{b=1}^B y_u^b R_{u,b} \quad (34)$$

s.t. (C1), (C7), and (C8), (35)

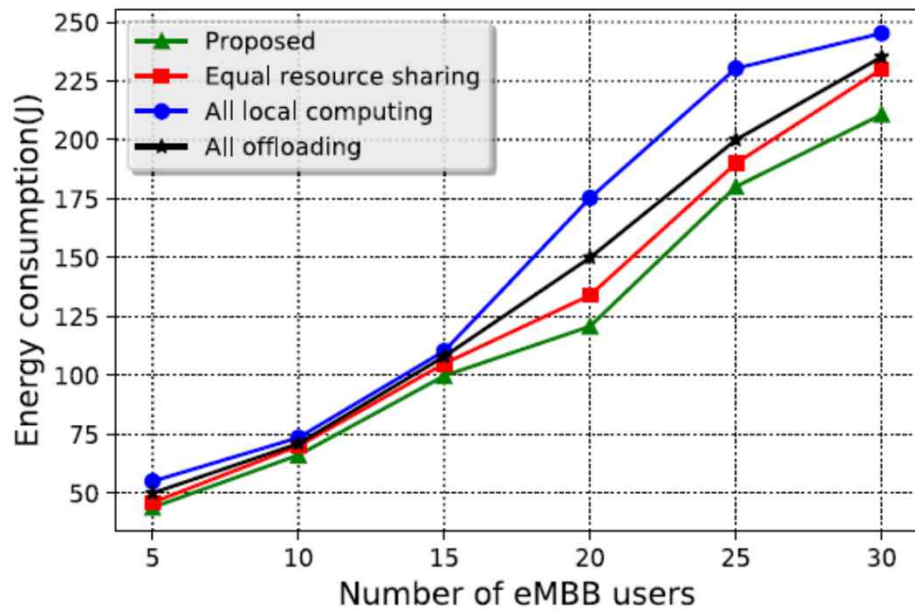
$$(P3) : \min_w \sum_{u=1}^U E_u^{\text{Off}} - \phi \sum_{u=1}^U \sum_{b=1}^B y_u^b R_{u,b} \quad (41)$$

s.t. (C4), (C5), and (C6), (42)

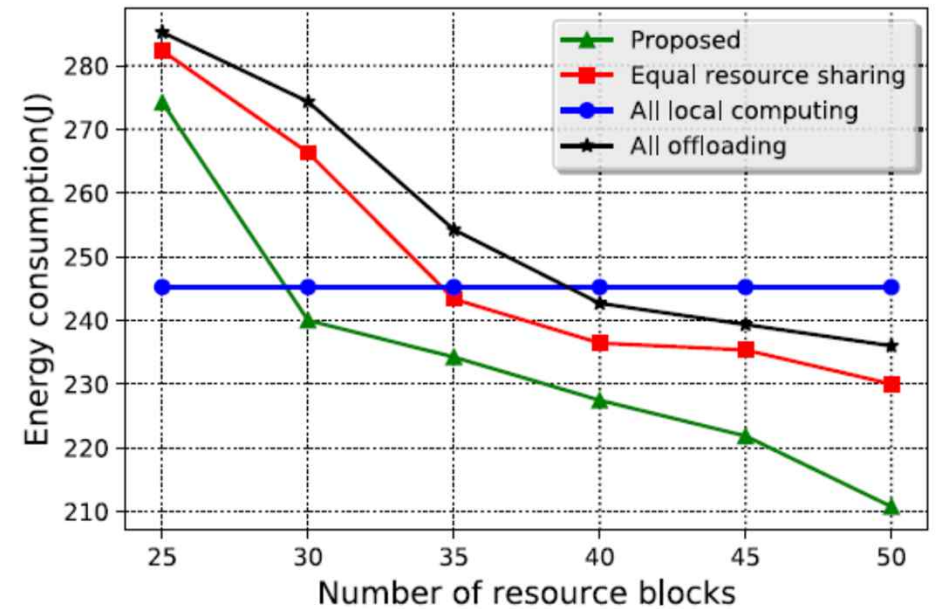
Algorithm 1 Iterative Algorithm for the Relaxed Problem

- 1: **Initialization:** Set $k = 0$, $\epsilon_1, \epsilon_2, \epsilon_3 > 0$, and initial solutions $(l^{(0)}, y^{(0)}, w^{(0)})$;
- 2: **repeat**
- 3: Compute $l^{(k+1)}$ from (P1) at given y^k , and w^k ;
- 4: Compute $y^{(k+1)}$ from (P2) at given $l^{(k+1)}$, and w^k ;
- 5: Compute $w^{(k+1)}$ from (P3) at given $l^{(k+1)}$, and $y^{(k+1)}$;
- 6: $k = k + 1$;
- 7: **until** $\|l^{(k+1)} - l^{(k)}\| \leq \epsilon_1$, $\|y^{(k+1)} - y^{(k)}\| \leq \epsilon_2$, and $\|w^{(k+1)} - w^{(k)}\| \leq \epsilon_3$;
- 8: Then, set $(l^{(k+1)}, y^{(k+1)}, w^{(k+1)})$ as the desired solution.

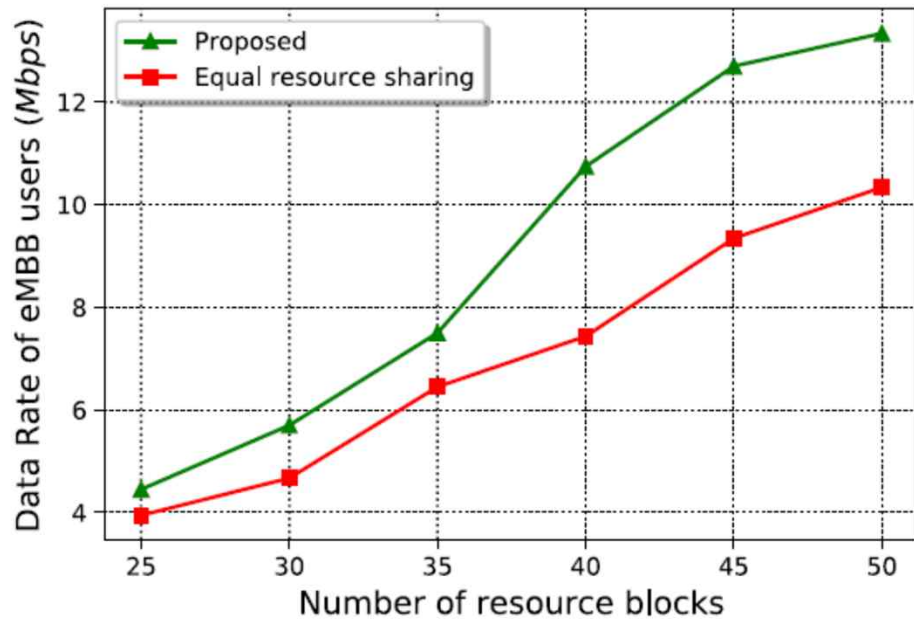
Yan Kyaw Tun, Do Hyon Kim, Madyan Alsenwi, Nguyen H. Tran, Zhu Han, Choong Seon Hong, "Energy Efficient Communication and Computation Resource Slicing for eMBB and URLLC Coexistence in 5G and Beyond," IEEE Access, Vol.8, pp.136024-136035, Jul. 2020



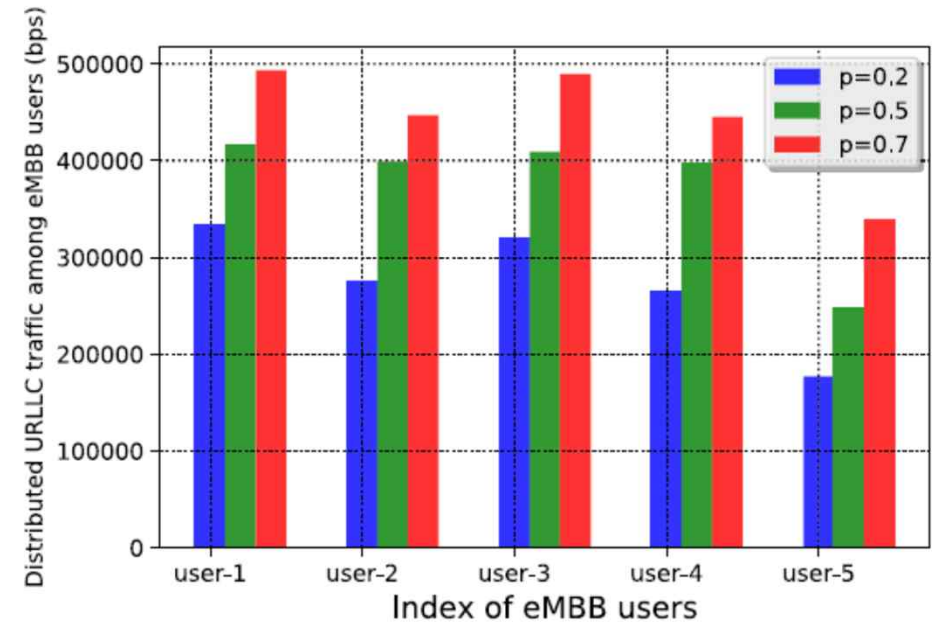
Energy consumption under different number of users.



Energy consumption under different number of resource blocks



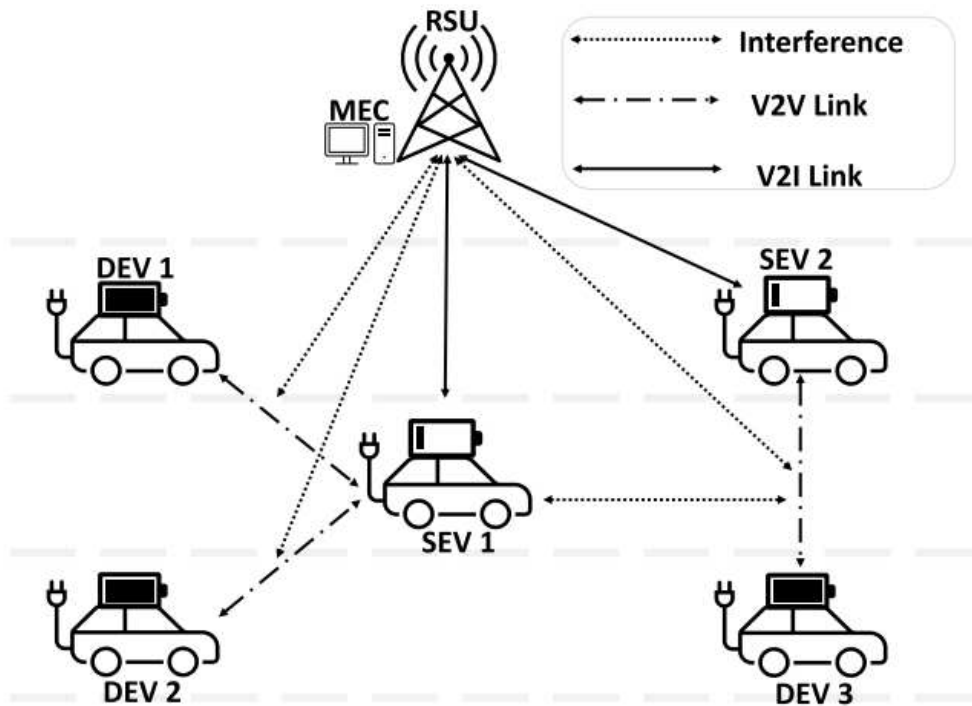
Data rate of eMBB users under different number of resource blocks.



Data rate of URLLC users on the puncturing resource of eMBB users.

Use case 4: Joint Communication, Computation, and Control for Computational Task Offloading in Vehicle-Assisted Multi-Access Edge Computing

- System model
- Problem formulation
- Solution approach
- Simulation results



- Decision making for task-offloading in multiple location: EVs, MEC, etc.
- Communication resource allocation for V2V, V2I links
- Interference Management for V2V, V2I links w.r.t. protection constraint, *i.e.*, maximum interference for each link

Goal: Minimize the trade-off between Energy Consumption and Latency in EV-assisted MEC

$$\mathcal{F}(\zeta, p, h, x, y, z) = \sum_{i \in \mathcal{O}} \left[\psi \Phi \left(\underbrace{\kappa_i h_i^2 x_i}_{\text{local energy}} + \underbrace{\sum_{j \in \mathcal{D}} \left(\kappa_j h_{i,j}^2 + p_{i,j} \frac{S_i}{\sum_{b \in \beta} R_{i,j}^b \zeta_{i,j}^b} \right) y_{i,j}}_{\text{V2V energy}} + \underbrace{\left(\kappa_0 h_{i,0}^2 + p_{i,0} \frac{S_i}{\sum_{b \in \beta} R_{i,0}^b \zeta_{i,0}^b} \right) z_{i,0}}_{\text{V2I energy}} \right) \right. \\
 \left. + (1 - \psi) \left(\underbrace{\frac{S_i}{h_i} x_i}_{\text{local latency}} + \underbrace{\sum_{j \in \mathcal{D}} \left(\frac{S_i}{h_{i,j}} + \frac{S_i}{\sum_{b \in \beta} R_{i,j}^b \zeta_{i,j}^b} \right) y_{i,j}}_{\text{V2V latency}} + \underbrace{\left(\frac{S_i}{h_{i,0}} + \frac{S_i}{\sum_{b \in \beta} R_{i,0}^b \zeta_{i,0}^b} \right) z_{i,0}}_{\text{V2I latency}} \right) \right].$$

Trade-off coefficient: ψ

Normalization factor: Φ

Energy consumption: $\kappa_i h_i^2 x_i$, $\kappa_j h_{i,j}^2 + p_{i,j} \frac{S_i}{\sum_{b \in \beta} R_{i,j}^b \zeta_{i,j}^b}$, $\kappa_0 h_{i,0}^2 + p_{i,0} \frac{S_i}{\sum_{b \in \beta} R_{i,0}^b \zeta_{i,0}^b}$

Latency: $\frac{S_i}{h_i} x_i$, $\frac{S_i}{h_{i,j}} + \frac{S_i}{\sum_{b \in \beta} R_{i,j}^b \zeta_{i,j}^b}$, $\frac{S_i}{h_{i,0}} + \frac{S_i}{\sum_{b \in \beta} R_{i,0}^b \zeta_{i,0}^b}$

Note that: Energy and latency are not the same unit.
That's why we need to normalize them in the same scale

Use Cases 4: Problem formulation

$$\min_{\zeta, p, h, x, y, z} \mathcal{F}(\zeta, p, h, x, y, z) \quad (29a)$$

SEV/DEV/MEC decision variables

s.t.:

Resource Block allocation variable

Transmit power variable

Achievable data rate

Channel gains

Interference threshold

Computing resource capacity of SEV

Computing resource capacity of MEC

Computing resource capacity of DEV

Energy Capacity of SEV

Energy Capacity of DEV

Latency requires of SEV

Mobility constraint between SEV and DEV

$$\sum_{b=1}^B \zeta_{i,b} \leq 1, \forall i \in \mathcal{O}, \quad (29b)$$

$$0 \leq p_i \leq p_i^{\max}, \forall i \in \mathcal{O}, \quad (29c)$$

$$R_i \geq R_{\min}, \forall i \in \mathcal{O}, \quad (29d)$$

$$\sum_{i \in \mathcal{O}} p_i^b g_{i,0}^b \zeta_{i,b} \leq I_b^{\max}, \forall b \in \beta, \quad (29e)$$

$$0 \leq h_i x_i \leq H_i^{\max}, \forall i \in \mathcal{O}, \quad (29f)$$

$$\sum_{i \in \mathcal{O}} h_{i,0} z_{i,0} \leq H_0^{\max}, \forall i \in \mathcal{O}, \quad (29g)$$

$$\sum_{i \in \mathcal{O}} h_{i,j} y_{i,j} \leq H_j^{\max}, \forall j \in \mathcal{D}, \quad (29h)$$

$$E_i \leq E_i^{\max}, \forall i \in \mathcal{O}, \quad (29i)$$

$$E_j \leq E_j^{\max}, \forall j \in \mathcal{D}, \quad (29j)$$

$$L_i \leq L_i^{\max}, \forall i \in \mathcal{O}, \quad (29k)$$

$$L_{i,j} \leq \tau_{i,j}, \forall i \in \mathcal{O}, \forall j \in \mathcal{D}, \quad (29l)$$

$$x_i + \sum_{j \in \mathcal{D}} y_{i,j} + z_{i,0} = 1, \quad (29m)$$

Communication resource allocation constraints

Interference management constraint

Computation resource allocation constraints

Energy constraints

Latency constraint

QoS constraint which is coupling among variables that makes the problem intractable!

$$\text{RBA : } \max_{\zeta} \mathcal{F}(\zeta) = \sum_{i \in \Omega_S} \sum_{b \in B} R_{i,b} \zeta_{i,b} \quad (30a)$$

s.t. Set of EV pairs

$$\sum_{b=1}^B \zeta_{i,b} \leq 1, \forall i \in \Omega_S, \quad (30b)$$

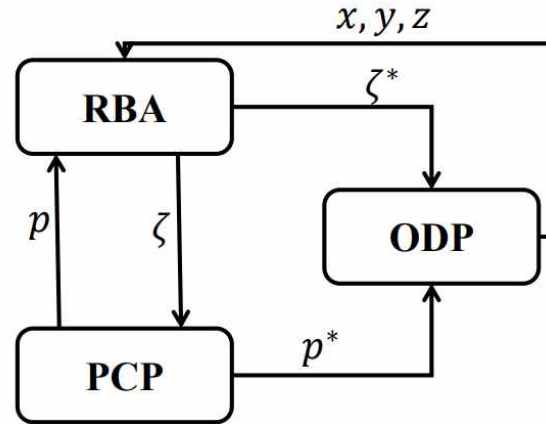
$$\sum_{i \in \Omega_S} p_{i,0,b} g_{i,0,b} \zeta_{i,b} \leq I_b^{\max}, \forall b \in \beta, \quad (30c)$$

$$\zeta_{i,b} \in \{0, 1\}, \forall b \in \beta. \quad (30d)$$

RBA: Resource Block Allocation

1. Iteratively solve each subproblem to achieve an sub-optimal solution

2. Update the solution based on last iteration



ODP: Offloading Decision Problem

$$\text{ODP : } \min_{x,y,z} \mathcal{F}(x, y, z) \quad (39a)$$

s.t. $0 \leq h_i \leq H_i^{\max}, \quad (39b)$

$$\sum_{i \in \mathcal{O}} h_{i,0} \leq H_0^{\max}, \quad (39c)$$

$$\sum_{i \in \mathcal{O}} h_{i,j} y_{i,j} \leq H_j^{\max}, \forall j \in \mathcal{D}, \quad (39d)$$

$$E_i \leq E_i^{\max}, \forall i \in \mathcal{O}, \quad (39e)$$

$$E_j \leq E_j^{\max}, \forall j \in \mathcal{D}, \quad (39f)$$

$$L_i \leq L_i^{\max}, \quad (39g)$$

$$L_{i,j} \leq \tau_{i,j}, \forall i \in \mathcal{O}, \forall j \in \mathcal{D}, \quad (39h)$$

$$x_i + \sum_{j \in \mathcal{D}} y_{i,j} + z_{i,0} = 1, \quad (39i)$$

$$h_{i,j} \geq 0, \quad (39j)$$

$$h_{i,0} \geq 0, \quad (39k)$$

$$x_i \in \{0, 1\}, \quad (39l)$$

$$y_{i,j} \in \{0, 1\}, \quad (39m)$$

$$z_{i,0} \in \{0, 1\}, \quad (39n)$$

3. Always guarantee a stationary solution by the BCD technique - global optimal solution

BCD: Block-Coordination Descent

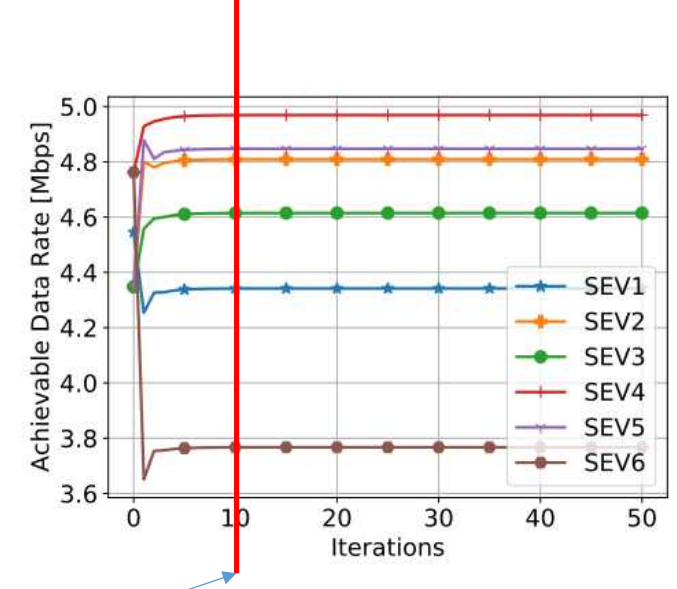
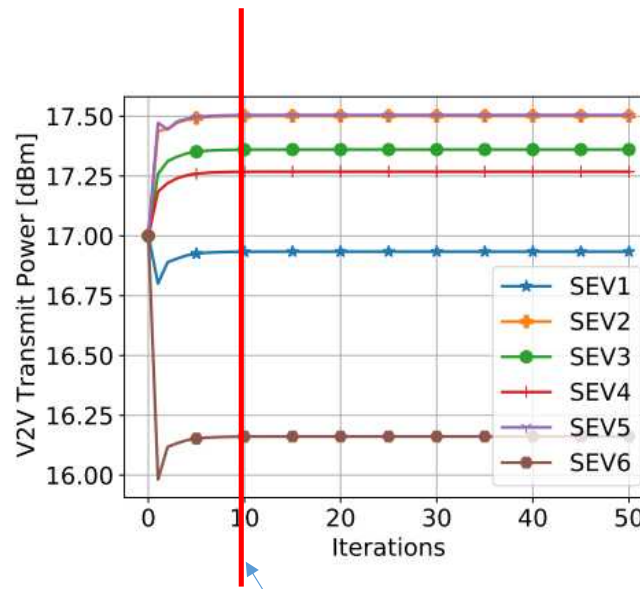
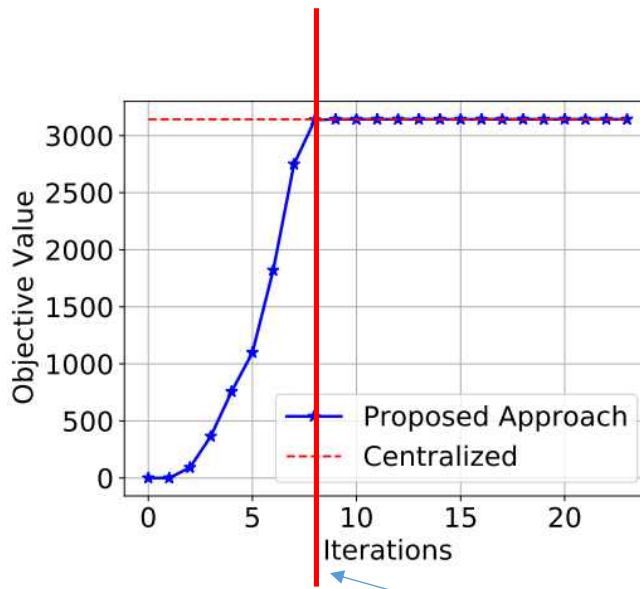
$$\text{PCP : } \max_p \mathcal{F}(p) = \sum_{i \in \mathcal{O}} \sum_{b \in \beta} F(p_{i,b}) = R_i^b \quad (34a)$$

s.t. $0 \leq p_i^b \leq p_i^{\max}, \quad (34b)$

$$R_i^b \geq R_{\min}, \quad (34c)$$

$$\sum_{i \in \mathcal{O}} p_i^b g_{i,0} \zeta_{i,b} \leq I_b^{\max}. \quad (34d)$$

PCP: Power Control Problem

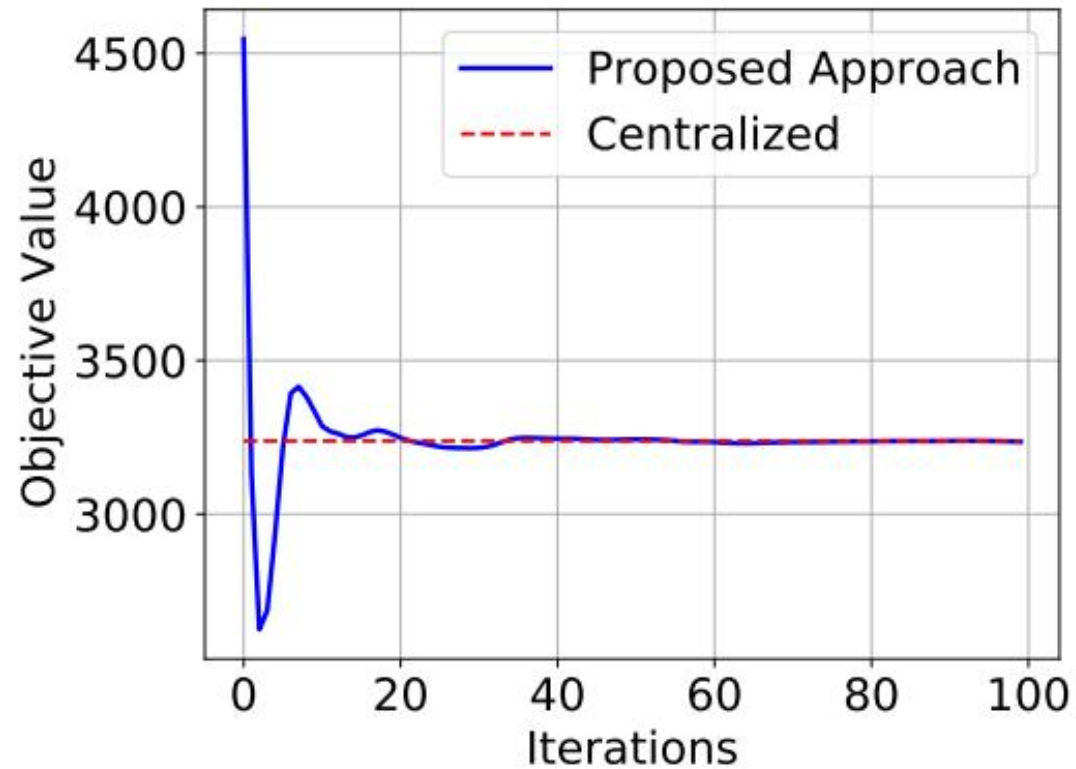


Convergence of ODP

Convergence of PCP

Convergence of RBA

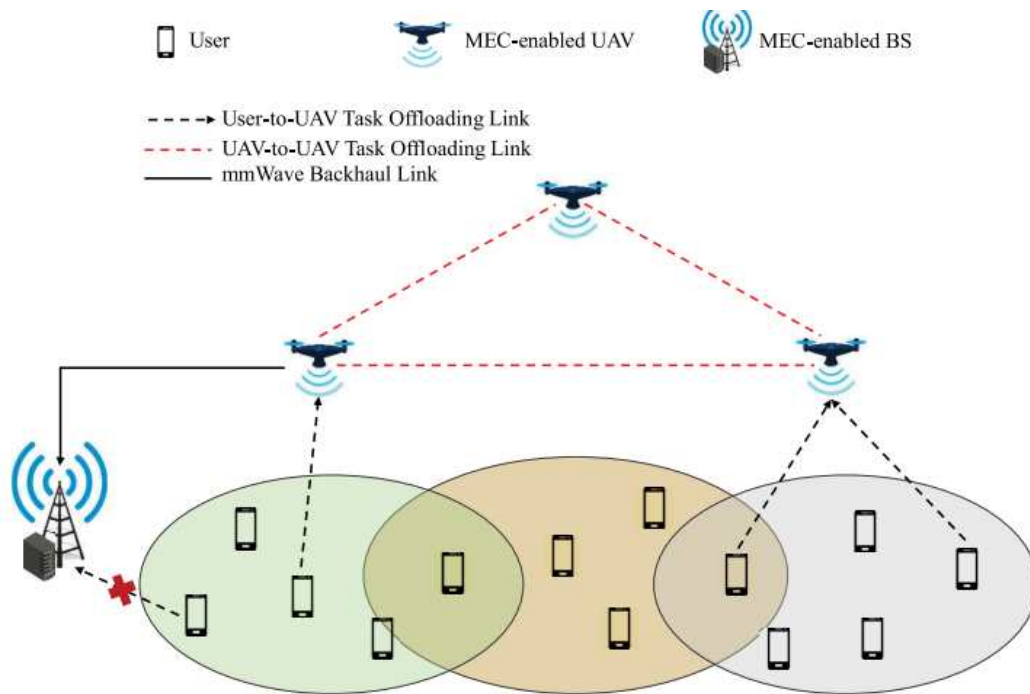
Fast convergence rate within 10 iteration



Convergence of RBA + PCP + ODP

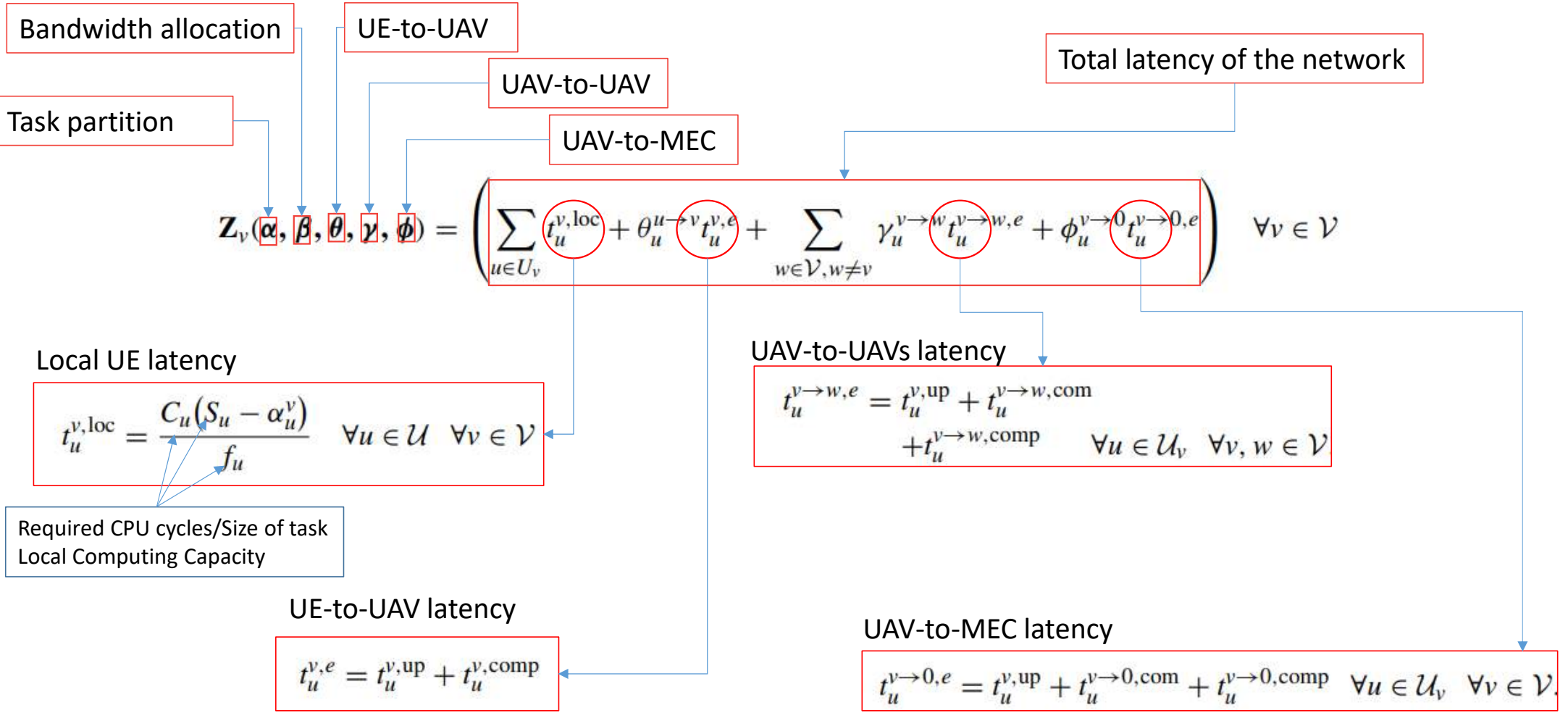
Use case 5: Collaboration in the Sky: A Distributed Framework for Task Offloading and Resource Allocation in Multi-Access Edge Computing

- System model
- Problem formulation
- Solution approach
- Simulation results



Goal: Minimize total latency of the network

- UAVs-assisted multi-access edge computing:
 - ✓ UE is able to offload task to UAVs
 - ✓ UAV sharing resource with each other to speedup the task processing
 - ✓ Increase the capability of MEC
- We propose a joint optimization problem:
 - ✓ Bandwidth allocation
 - ✓ Task offload decision
 - ✓ Collaboration among UAVs



$$\mathbf{P}: \underset{\alpha, \beta, \theta, \gamma, \phi}{\text{minimize}} \quad \mathbf{Z}(\alpha, \beta, \theta, \gamma, \phi) \quad (40a)$$

$$\text{subject to} \quad E_u^{\text{loc}} + E_u^{v, \text{up}} \leq E_u^{\text{max}} \quad \forall u \in \mathcal{U}_v \quad (40b)$$

Energy for offloading and local processing less than or equal to capacity of local user

Base load energy + energy to help other UAVs + energy for offloading to MEC + hovering energy must less than maximum energy of UAV

Bandwidth allocation variable

Whenever to offload, the task must be 100% processed!

Task partitioning: percentage of task to offload

Decision of UAV v on user u's task

Decision of UAV v on request of UAV w

Decision of MEC for request of user u

$$\sum_{u \in \mathcal{U}_v} \theta_u^{u \rightarrow v} E_u^{v, \text{comp}} + \sum_{w \in \mathcal{V}, w \neq v} E^{v \rightarrow w} \quad (40c)$$

$$+ E^{v \rightarrow 0} + E^{v, \text{hov}} \leq E_v^{\text{max}} \quad \forall v \in \mathcal{V} \quad (40c)$$

$$\sum_{u \in \mathcal{U}_v} E_u^{v \rightarrow w, \text{comp}} + E^{w, \text{hov}} \leq E_w^{\text{max}} \quad (40d)$$

$$\forall w \in \mathcal{V}, w \neq v \quad (40d)$$

$$\sum_{u \in \mathcal{U}_v} \beta_u^v \leq 1 \quad \forall v \in \mathcal{V} \quad (40e)$$

$$\beta_u^v \in [0, 1], \quad \forall u \in \mathcal{U}_v, \forall v \in \mathcal{V} \quad (40f)$$

$$\theta_u^{u \rightarrow v} + \sum_{w \in \mathcal{V}, w \neq v} \gamma_u^{v \rightarrow w} + \phi_u^{v \rightarrow 0} = 1 \quad \forall u \in \mathcal{U}_v \quad (40g)$$

$$0 \leq \alpha_u^v \leq S_u \quad \forall u \in \mathcal{U} \quad \forall v \in \mathcal{V}, \quad (40h)$$

$$\theta_u^{u \rightarrow v} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v \in \mathcal{V} \quad (40i)$$

$$\gamma_u^{v \rightarrow w} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v, w \in \mathcal{V} \quad (40j)$$

$$\phi_u^{v \rightarrow 0} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v \in \mathcal{V} \quad (40k)$$

Energy constraints

Bandwidth constraints

QoS constraint

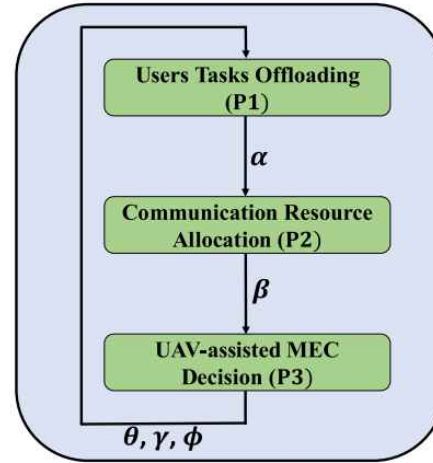
Decision variable constraints

P1: minimize $Z(\alpha)$ (41a)

subject to $E_u^{loc} + E_u^{v,up} \leq E_u^{max} \quad \forall u \in \mathcal{U}_v$ (41b)

$0 \leq \alpha_u^v \leq S_u \quad \forall u \in \mathcal{U} \quad \forall v \in \mathcal{V}$. (41c)

1. Iteratively obtain solution for P1 by fixing P2, and P3
2. Fixed P3, updated P2 based on P1
3. Update P3 based solution of P1 and P2 on last iteration



P2: minimize $Z(\beta)$ (42a)

subject to $E_u^{v,up} \leq E_u^{max} \quad \forall u \in \mathcal{U}_v$ (42b)

$\sum_{u \in \mathcal{U}_v} \beta_u^v \leq 1 \quad \forall v \in \mathcal{V}$ (42c)

$\beta_u^v \in [0, 1] \quad \forall u \in \mathcal{U}_v \quad \forall v \in \mathcal{V}$. (42d)

P3: minimize $Z(\theta, \gamma, \phi)$ (50a)

subject to $\sum_{u \in \mathcal{U}_v} \theta_u^{u \rightarrow v} E_u^{v,comp} + \sum_{w \in \mathcal{V}, w \neq v} E^{v \rightarrow w} + E^{v \rightarrow 0} + E^{v,hov} \leq E_v^{max} \quad \forall v \in \mathcal{V}$ (50b)

$\sum_{u \in \mathcal{U}_v} E_u^{v \rightarrow w,comp} + E^{w,hov} \leq E_w^{max} \quad \forall w \in \mathcal{V} \quad w \neq v$ (50c)

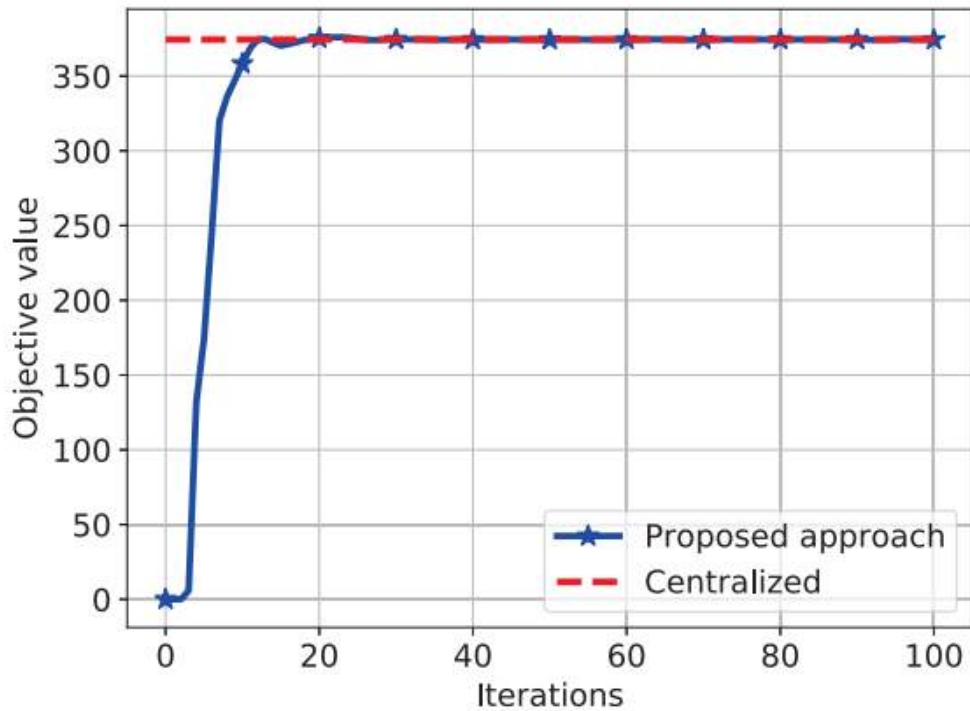
$\theta_u^{u \rightarrow v} + \sum_{\substack{w \in \mathcal{V}, \\ w \neq v}} \gamma_u^{v \rightarrow w} + \phi_u^{v \rightarrow 0} = 1 \quad \forall u \in \mathcal{U}_v$ (50d)

$\theta_u^{u \rightarrow v} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v \in \mathcal{V}$ (50e)

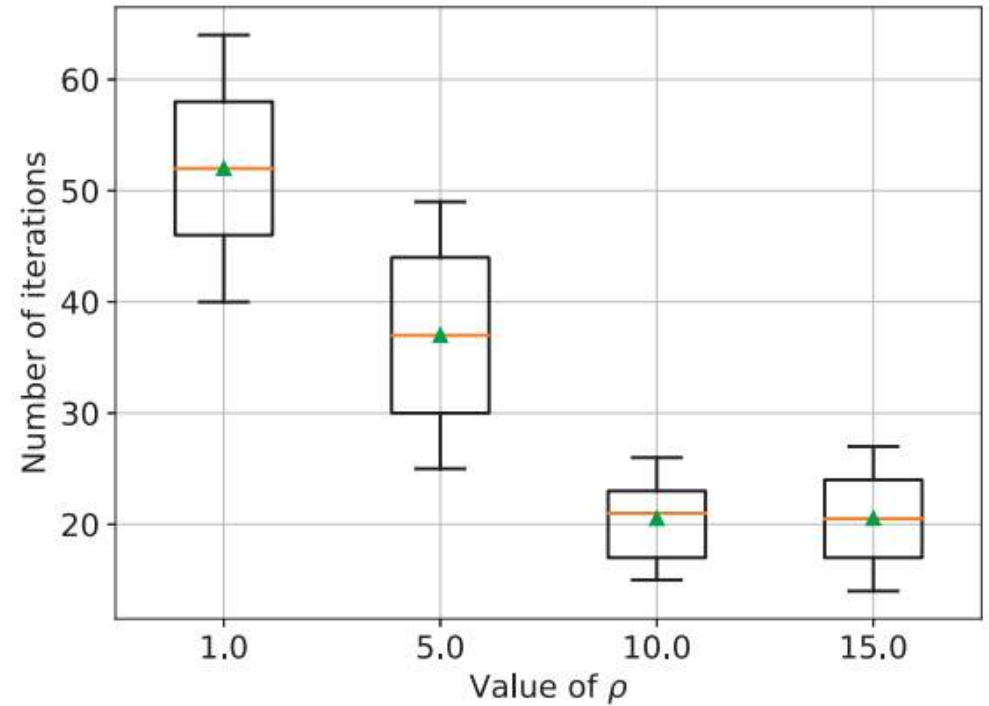
$\gamma_u^{v \rightarrow w} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v, w \in \mathcal{V}$ (50f)

$\phi_u^{v \rightarrow 0} \in \{0, 1\} \quad \forall u \in \mathcal{U}_u \quad \forall v \in \mathcal{V}$. (50g)

This iterative algorithm always guarantee an optimal solution with in $O(\frac{1}{\epsilon^2})$, where ϵ is convergence rates which is positive and strictly small, i.e., 1e-4

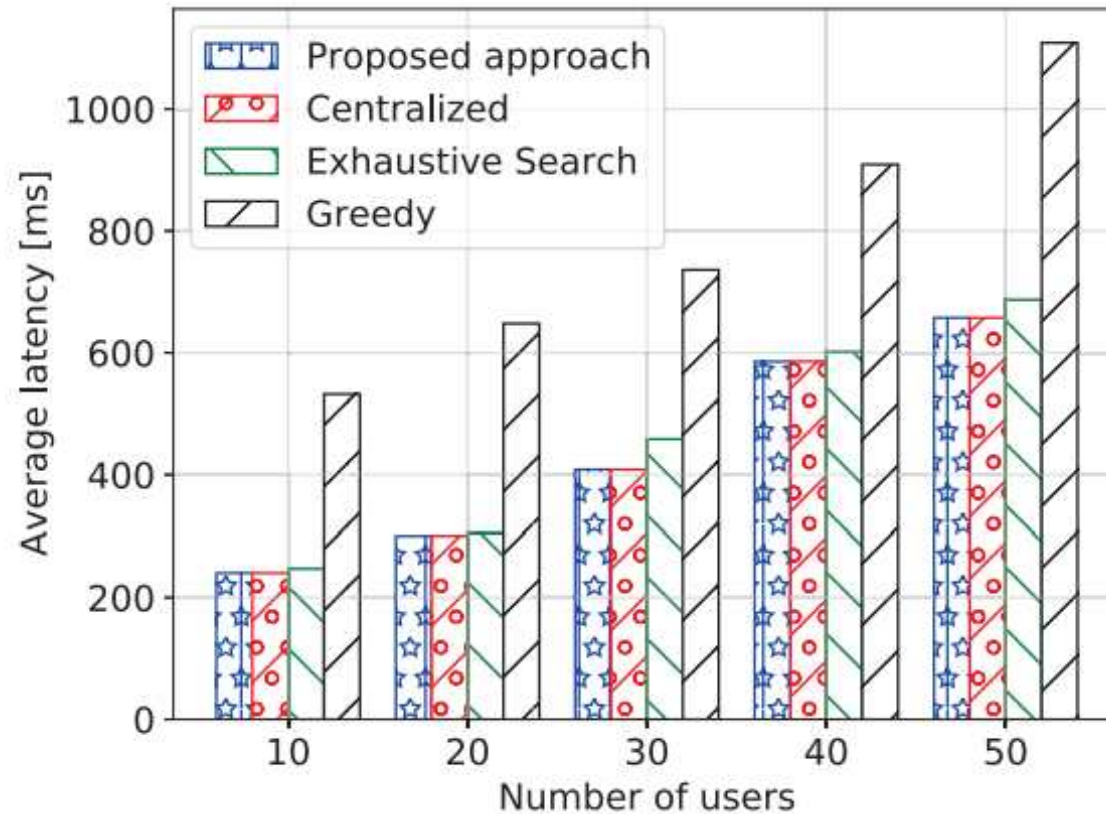


Convergence performance of proposed algorithm



Convergence rate versus penalty parameters

*Rho is penalty parameter for controlling the augmented Lagrangian in P3.



Average latency of proposed approach versus based-lines

Centralized Algorithm: the BS serves as a central coordinator

Greedy Algorithm: Considering availability of the neighboring UAVs on the bandwidth between a UAV and it's neighbors and it's neighbors computing resources

Concluding Remarks

- Open Issues
- Conclusion

- **Dynamic Slice Allocation**
 - A practical system would have users arriving and leaving a system with different demands at different time slots.

- **Mobility Aware Network slicing**
 - The current approaches for network slicing are not designed to handle mobility in the network.
 - Handling and orchestrating the radio access and core network will be very challenging in case of mobility.
 - Require migration of services from one point to other points in the network.

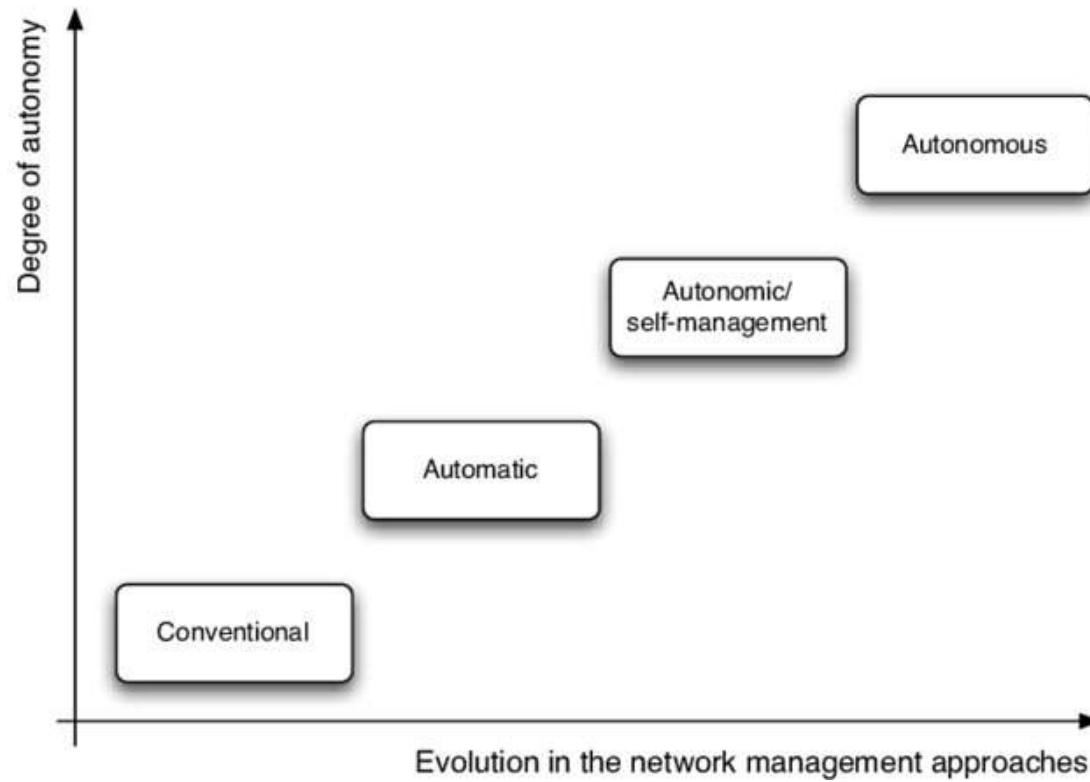
- This lecture is mainly focus to understand a full view of the resource management problem in 5G networks.
- We learned
 - The requirements and enabling technologies of 5G networks.
 - A detailed overview of network slicing that can be adapted to fulfill the 5G deliverables.
 - The recent research works' motivation, issues, challenges, and solutions.
 - Some open issues for future research and their potentials.

- This lecture is mainly focus to understand a full view of the resource management problem in 5G networks.
- We learned
 - The requirements and enabling technologies of 5G networks.

What is *missing*?

deliverables.

- The recent research works' motivation, issues, challenges, and solutions.
- Some open issues for future research and their potentials.



- The use of artificial intelligence will play a vital role for enabling a variety of applications in 5G and beyond wireless networks.
- AI definitively provides precious opportunities to analyze trends and recognize patterns. However, it is difficult to perfectly predict the desired results by using traditional simple models such as shallow ANNs
- Deep Neural Networks are envisioned to fill this gap and serve as key predicting enabler to support the 5G networks
- Network slicing coupled with AI will be defining the future of wireless networks

Thanks !!!

Q&A