



Resource Allocation for Multi-Tenant Edge Computing

Assist. Prof. **Andrea Araldo**

High level view

Network Operator

Service Providers

Entertainment

NETFLIX

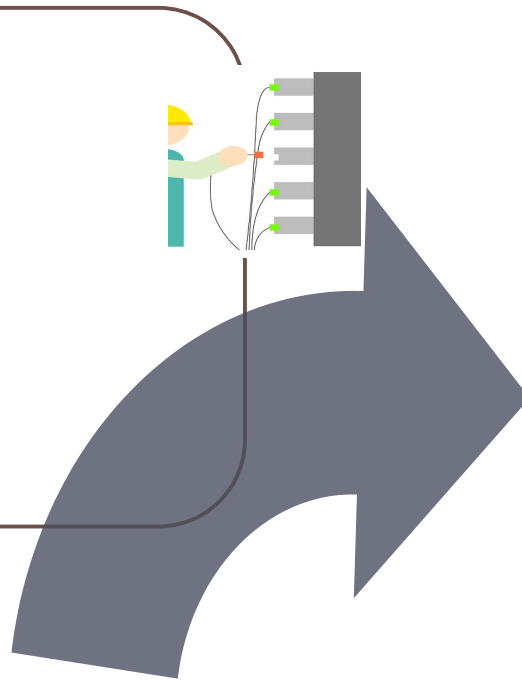


RENAULT

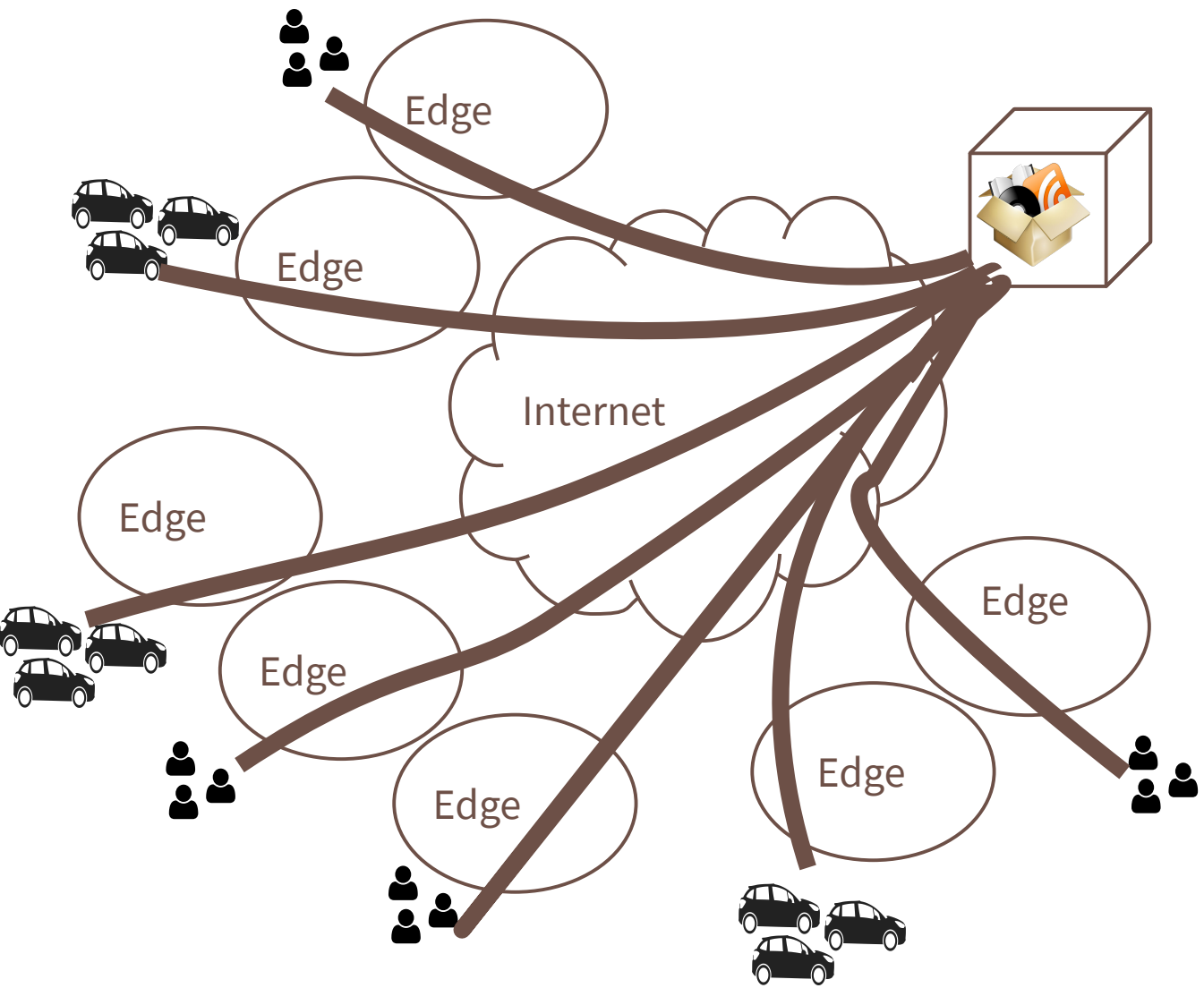
Monitor and predict hazards [1]



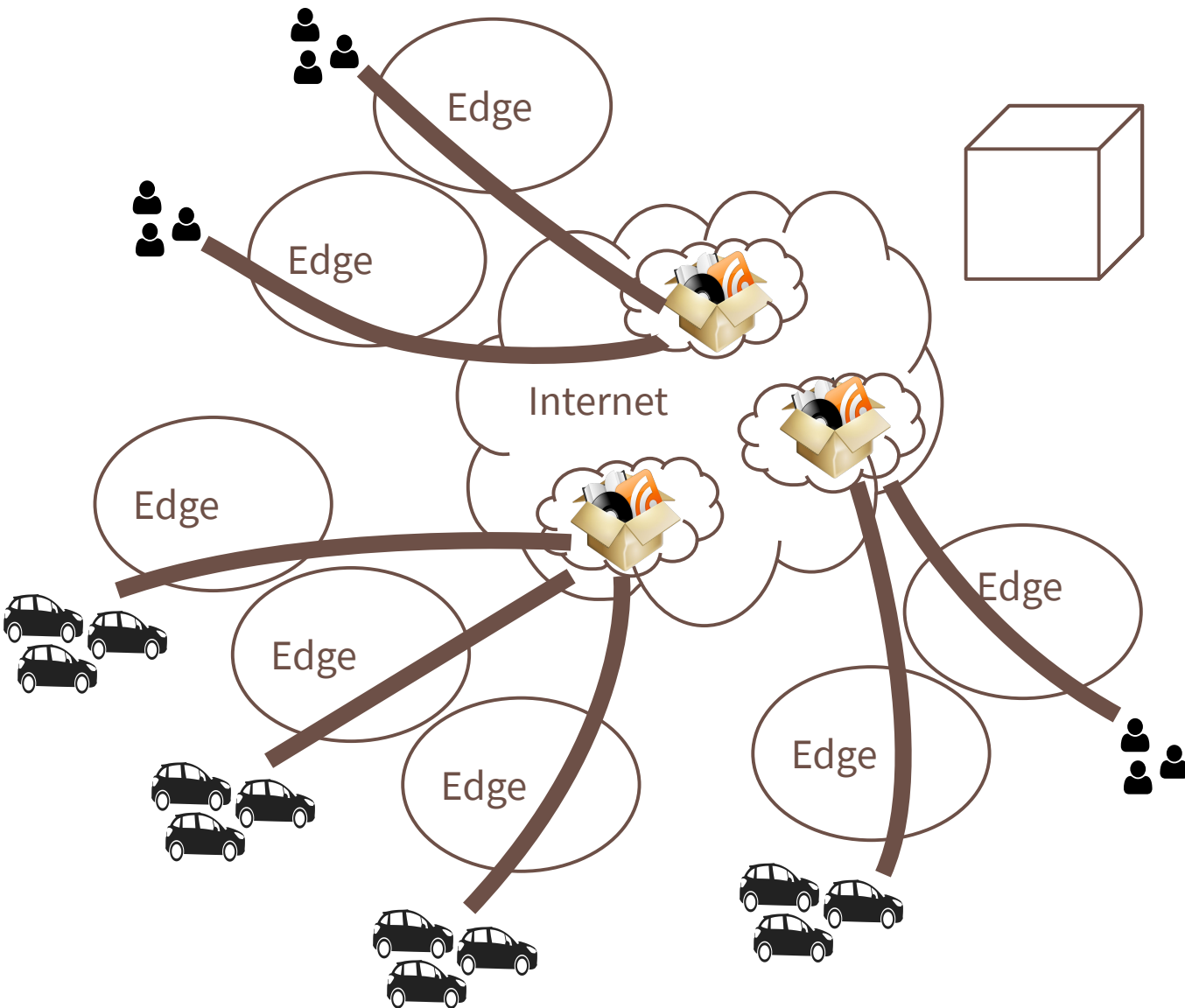
Users



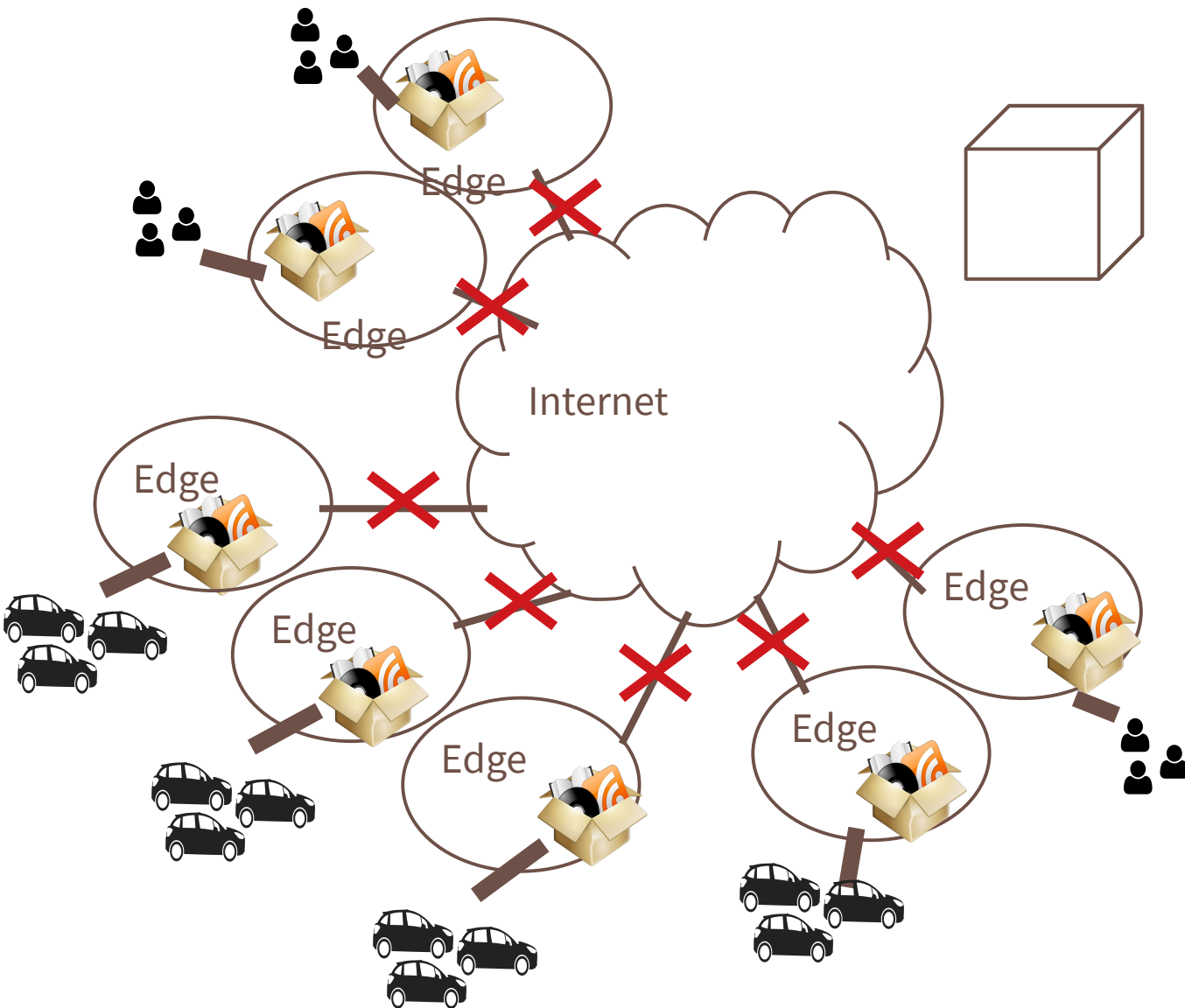
[1] SAFESPOT Final Report. EU Project, 2010



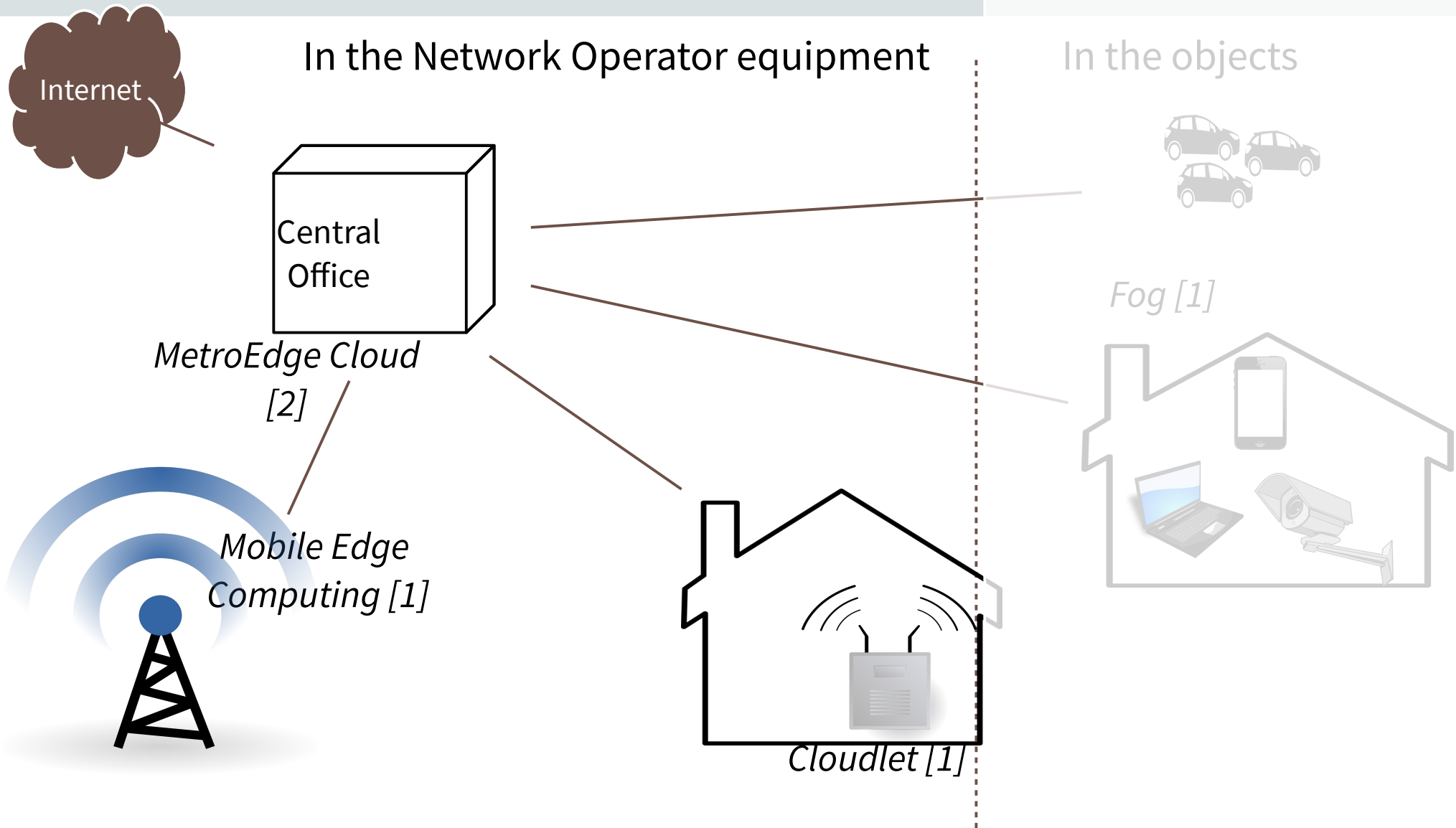
Cloud Computing



Edge Computing



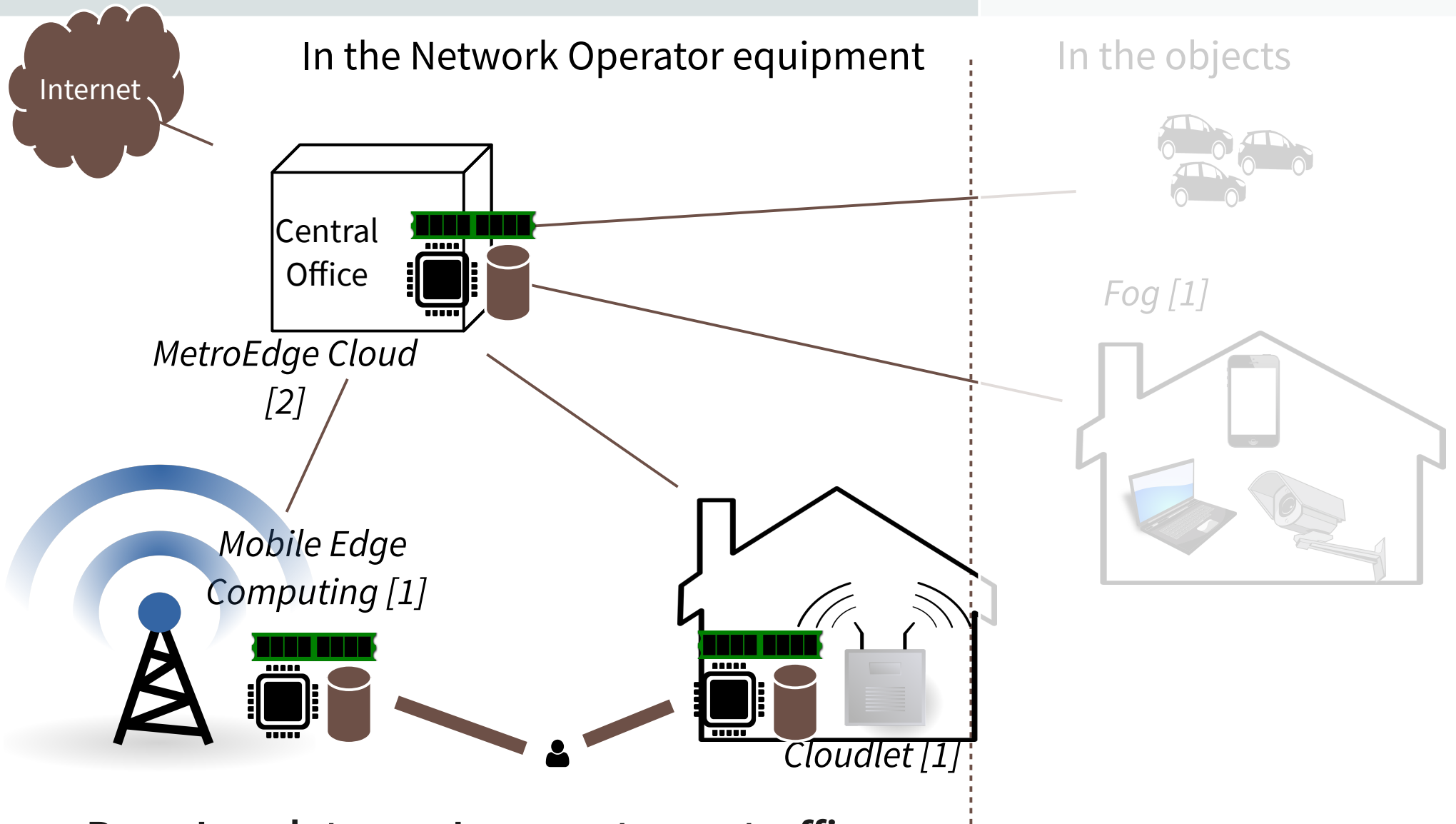
Where is the Edge? [1]



[1] Dolui, K. (2017). Comparison of Edge Computing Implementation. In IEEE GloTS

[2] Rimal et Al. (2018). Experimental Testbed for Edge Computing. IEEE ComMag

Where is the Edge? [1]



Pros: Less latency, Less upstream traffic

[1] Dolui, K. (2017). Comparison of Edge Computing Implementation. In IEEE GloTS

[2] Rimal et Al. (2018). Experimental Testbed for Edge Computing. IEEE ComMag

Example of Edge Computing

- **Netflix Open Connect Appliance**

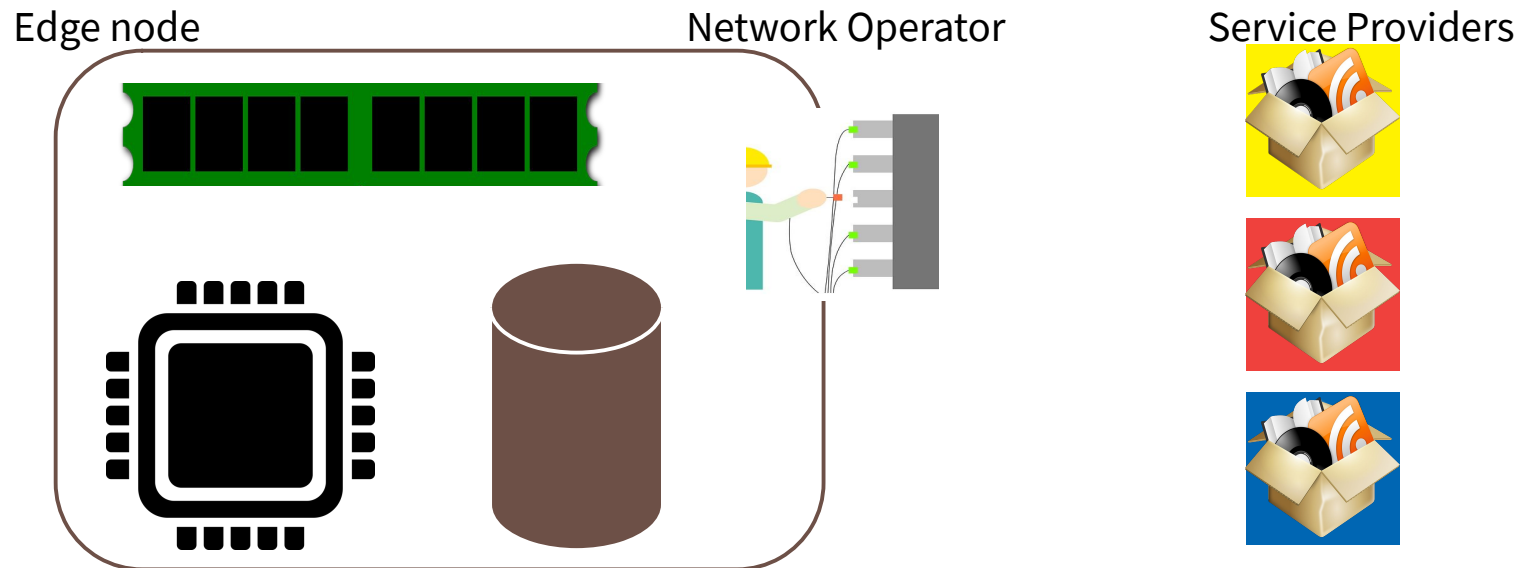


- **Limitations**

- Install physical boxes of many service providers: impossible!
- Impossible to reach extreme edge (e.g. Access Points)

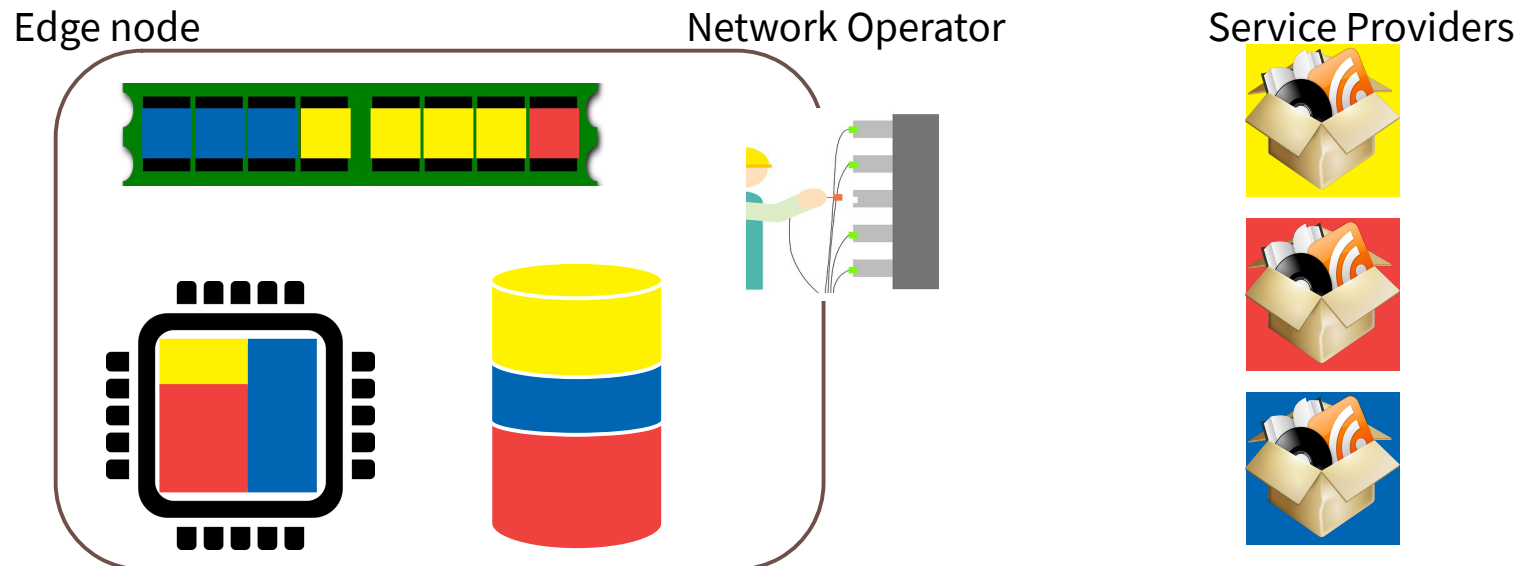
Resource allocation at the Edge

- The Network Operator should own physical resources



Resource allocation at the Edge

- The Network Operator should own physical resources
- Resources are virtualized
- ... and made available to Service Providers



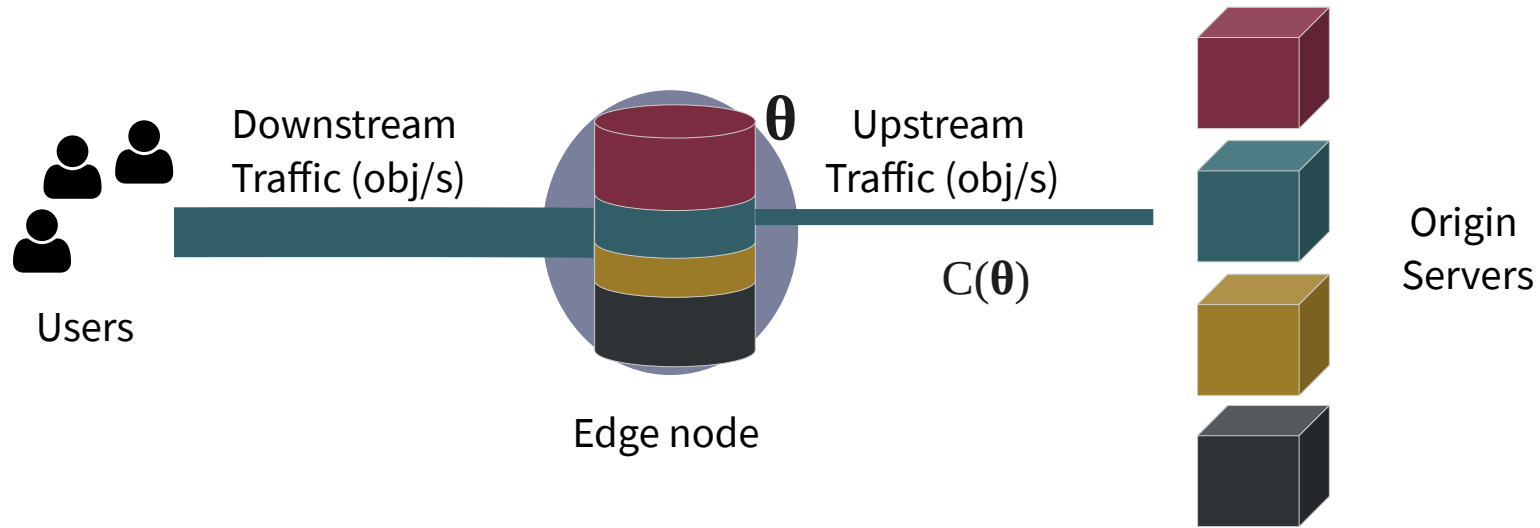
- The Network Operator is not simply a pipe, but a Micro-Cloud service provider
- Business opportunity for the Network Operator!
- 3rd party Service Providers can offer low-latency applications

Resource allocation strategies

- **Data-driven**
- **Multiple Option Resource Allocation (MORA)**

Data-driven optimization for cache partitioning

25



- Allocation of cache slots: $\theta = (\theta_1, \dots, \theta_p)$
- Traffic is encrypted
- Black box optimization
 - The function $C(\theta)$ is unknown
- Data-driven cache allocation
 - Just based on measured traffic amounts

Data-driven allocation

- **Stochastic Perturbation**

- Araldo et Al.,
“*Caching Encrypted Content*”,
IEEE Transactions on Networking 2018

IEEE/IFIP ComSoc Best Paper Award



Strong theoretical
convergence guarantees



Continuously perturbs
the allocation

- **Reinforcement Learning**

T. Bouganim, A. Araldo et Al.,
“*The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation*”,
ITC PhD Workshop, 2020

Data-driven allocation

- **Stochastic Perturbation**

- Araldo et Al.,
“*Caching Encrypted Content*”,
IEEE Transactions on Networking 2018

IEEE/IFIP ComSoc Best Paper Award



Strong theoretical
convergence guarantees



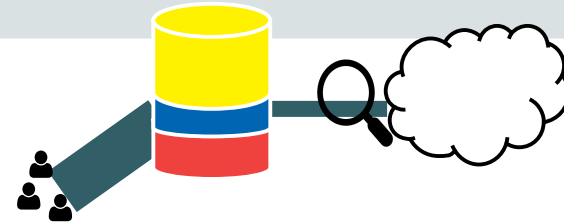
Continuously perturbs
the allocation

- **Reinforcement Learning**

T. Bouganim, A. Araldo et Al.,
“*The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation*”,
ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $a=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

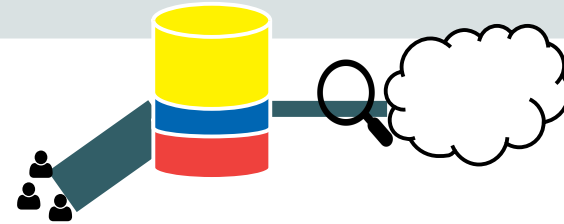


	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $a=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)



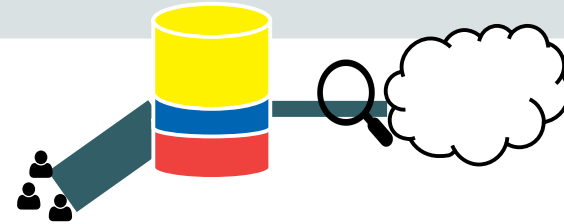
	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta = (\theta_1, \dots, \theta_p)$
- Action: perturbation, e.g. $\mathbf{a} = \Delta \cdot (1, 0, -1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)



	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

The table shows a Q-table with actions and allocations. A red circle highlights the cell for allocation 2 and action 2, labeled C_{22} . A red arrow points from this cell to the label "best" above it.

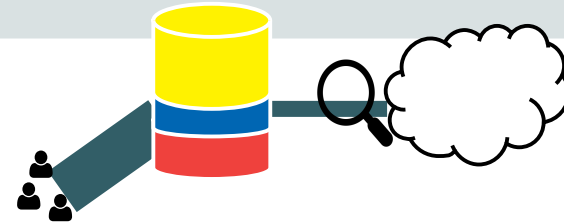


$$\mathbf{a} = \Delta \cdot (1, 0, -1)$$

[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta = (\theta_1, \dots, \theta_p)$
- Action: perturbation, e.g. $\mathbf{a} = \Delta \cdot (1, 0, -1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)



	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

The table shows a Q-table with rows for different allocations and columns for different actions. The cell containing C_{22} is circled in red, and a red arrow points to it from the word "best" written in red. The cell containing "allocation 2" is also circled in red.



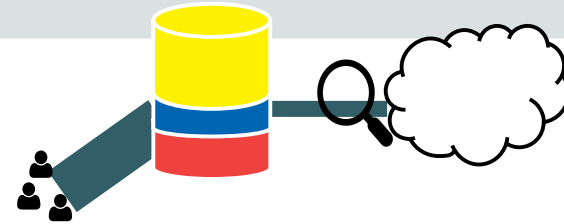
$$\mathbf{a} = \Delta \cdot (1, 0, -1)$$



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

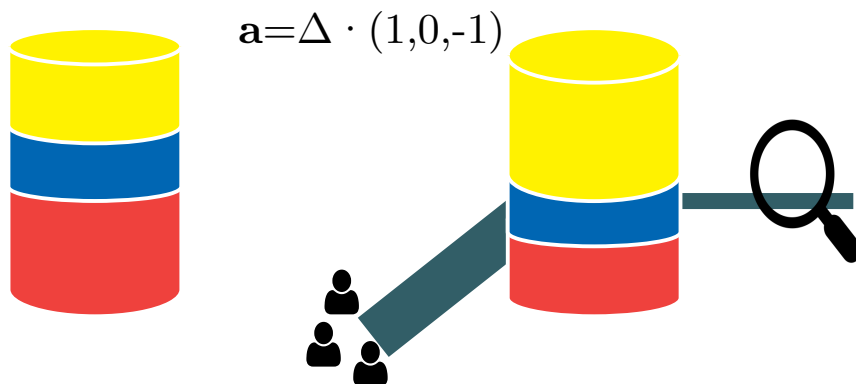
Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta = (\theta_1, \dots, \theta_p)$
- Action: perturbation, e.g. $a = \Delta \cdot (1, 0, -1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)



	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

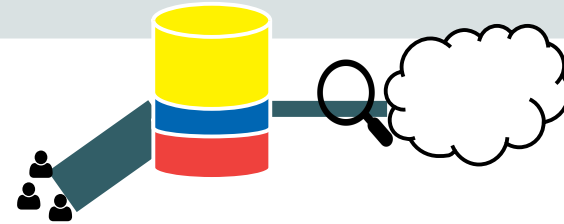
The table shows a Q-table with rows for different allocations and columns for different actions. The cell containing C_{22} is circled in red, and a red arrow points to it from the word "best" written in red. The cell containing "allocation 2" is also circled in red.



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta = (\theta_1, \dots, \theta_p)$
- Action: perturbation, e.g. $a = \Delta \cdot (1, 0, -1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

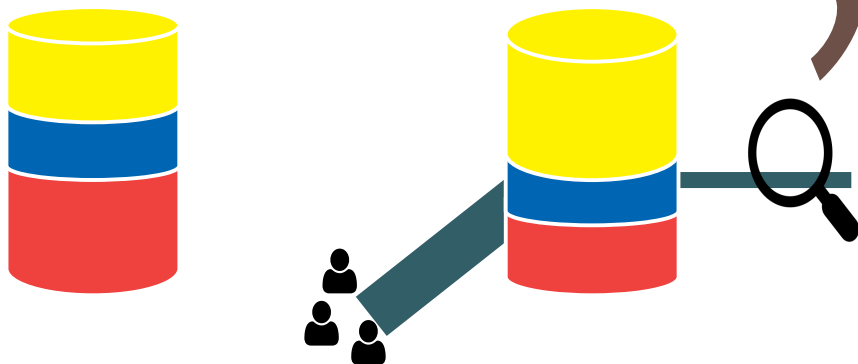


	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

Q-learning algorithm

$$Q^{new}(s, a) = Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s + a, a') - Q(s, a)]$$

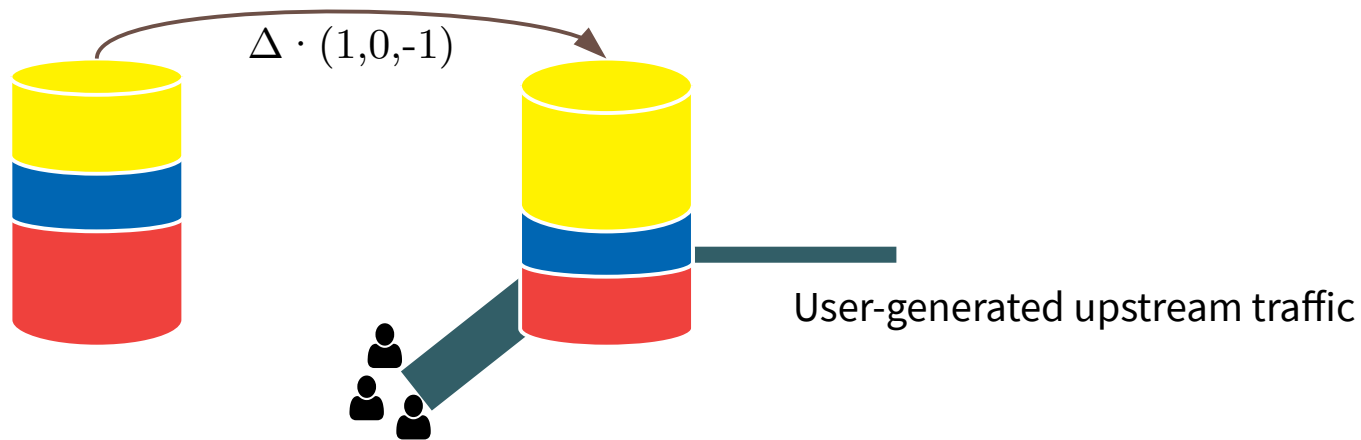
- We learn a good Q-table by perturbing-and-observing the system



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

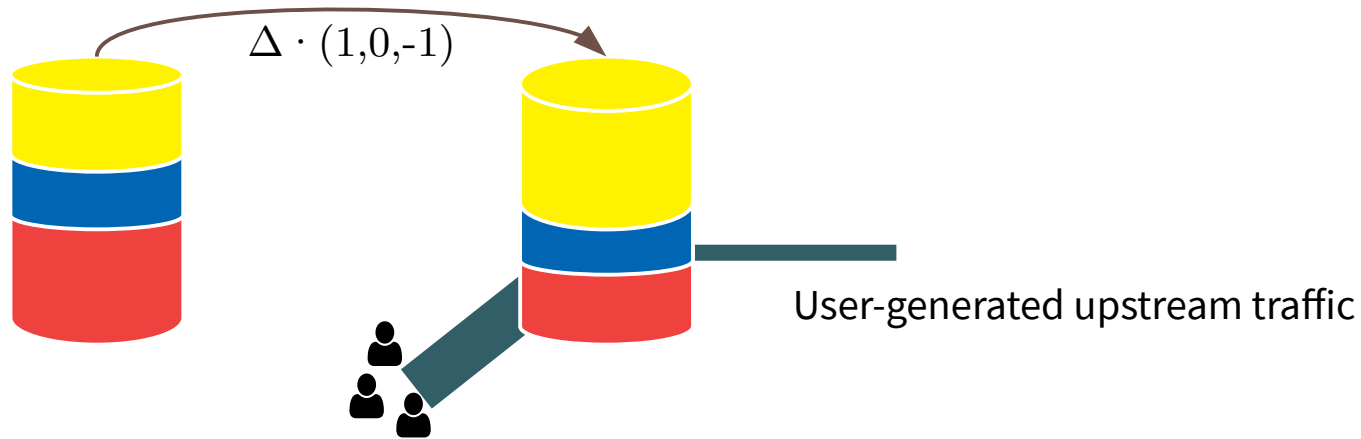
Offline vs. Online Reinforcement Learning

Offline

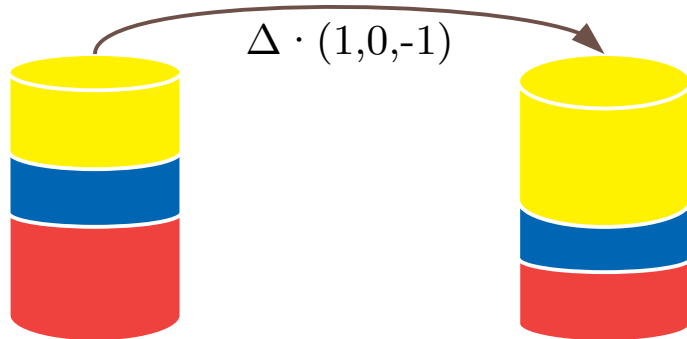


Offline vs. Online Reinforcement Learning

Offline

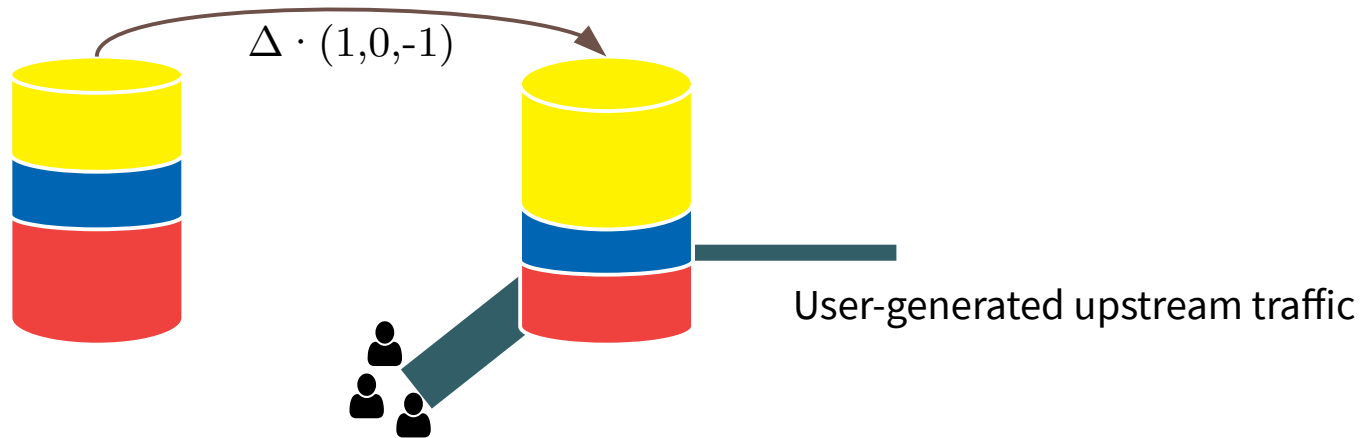


Online

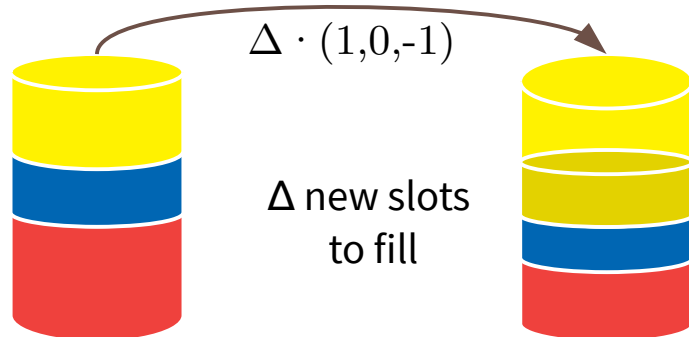


Offline vs. Online Reinforcement Learning

Offline

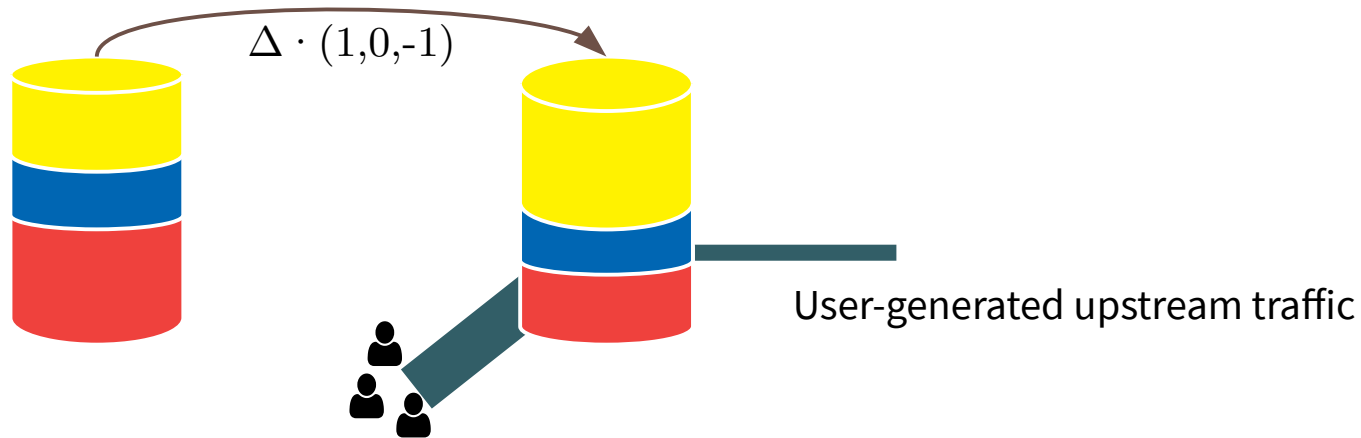


Online

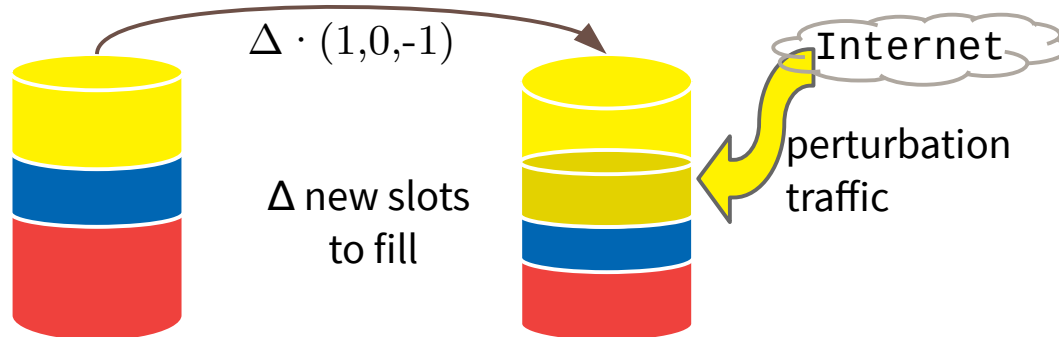


Offline vs. Online Reinforcement Learning

Offline

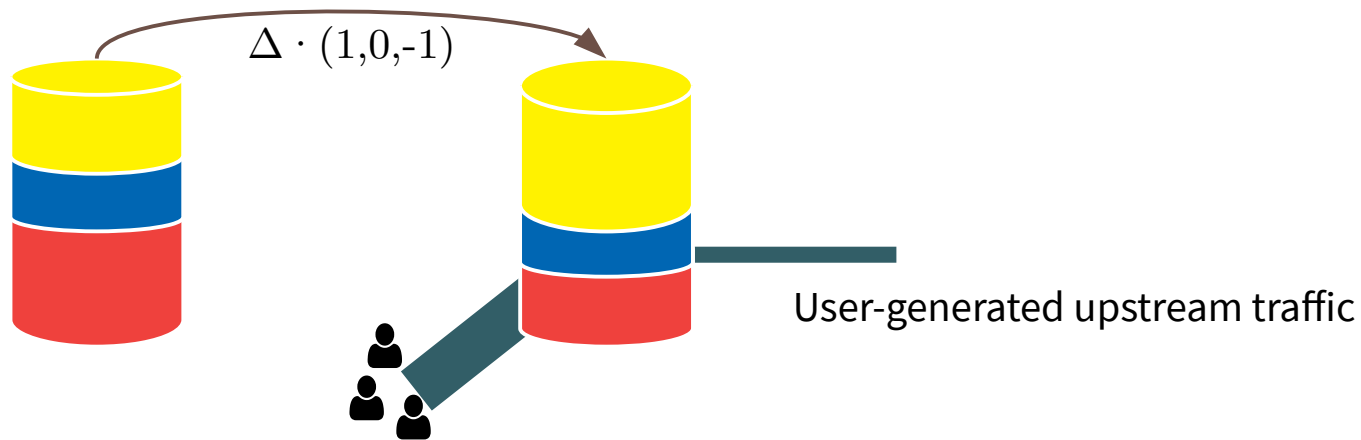


Online

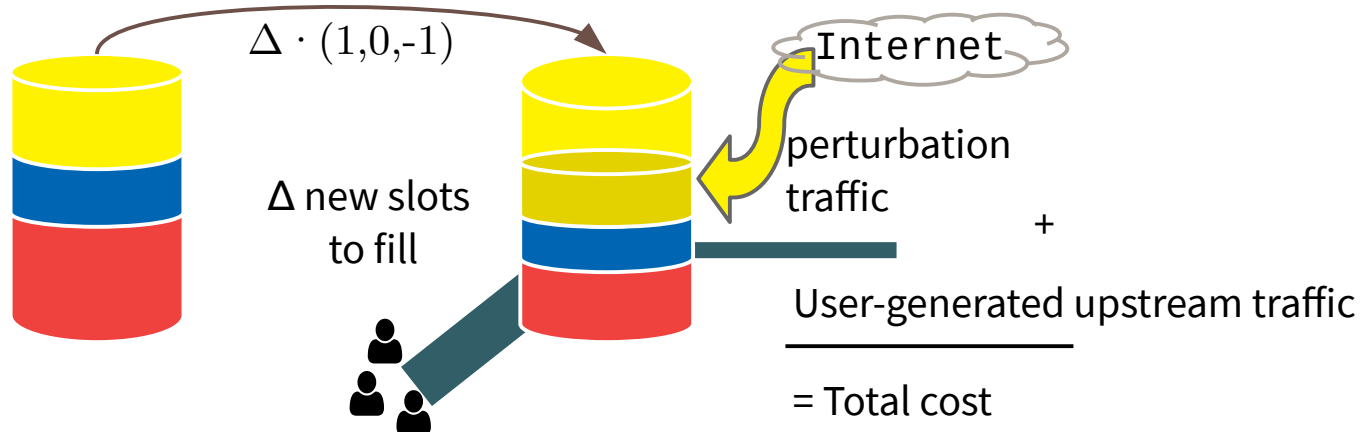


Offline vs. Online Reinforcement Learning

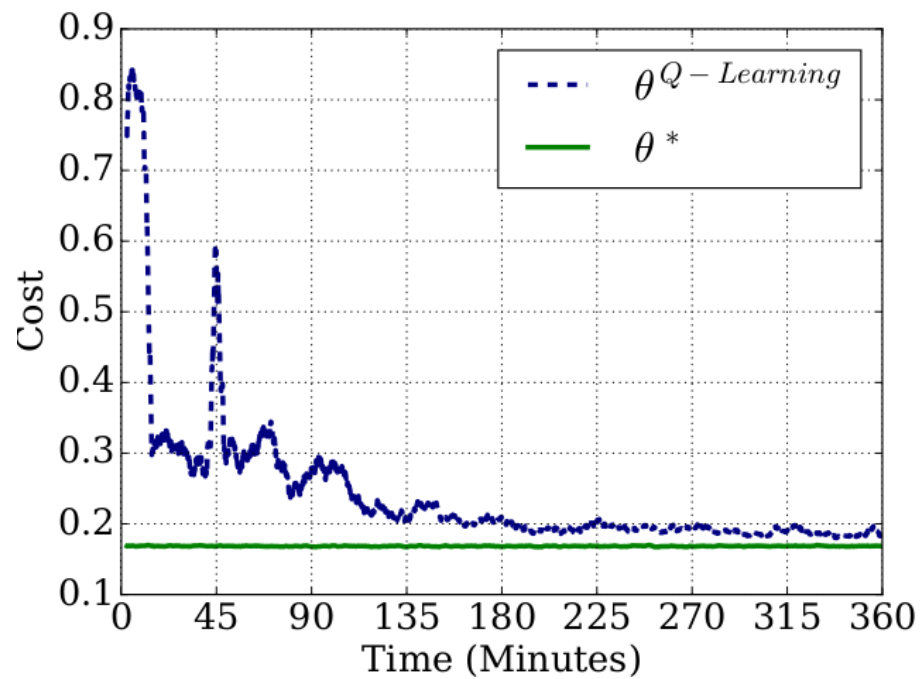
Offline



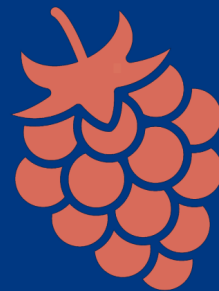
Online



Preliminary results



MORA: Multiple Option Resource Allocation

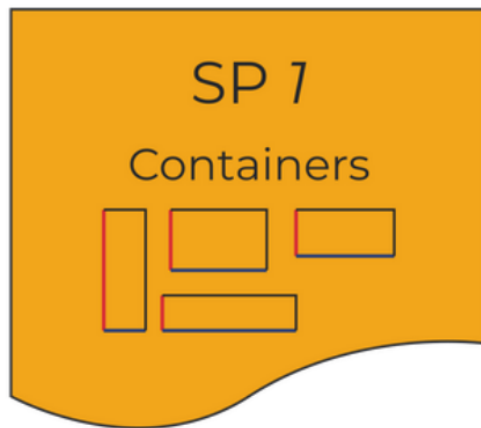


Araldo, A., Di Stefano, A., & Di Stefano, A. (2020). Resource Allocation for Edge Computing with Multiple Tenant Configurations. In **ACM/SIGAPP** Symposium On Applied Computing **SAC**.

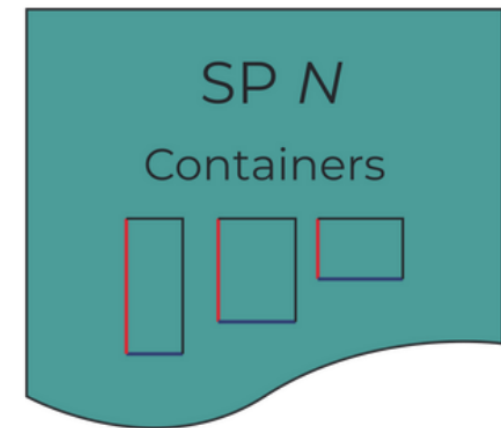


Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]



Service Providers



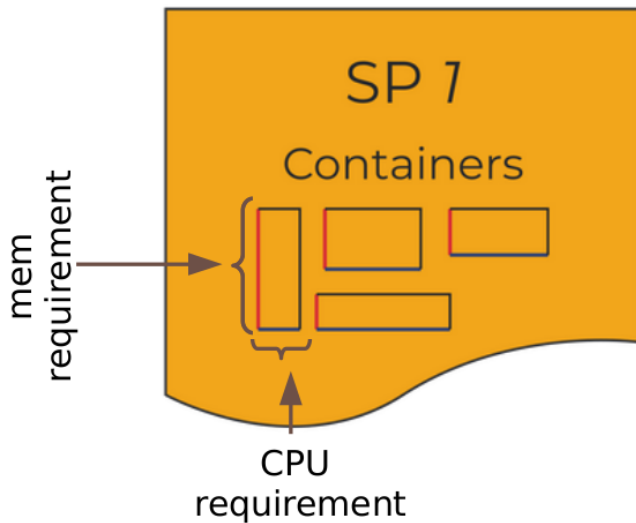
[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.

<https://github.com/Ressource-Allocation/CDN-Transcode-Sample>

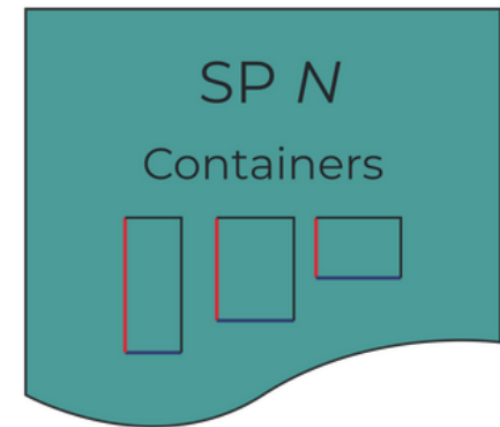


Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]



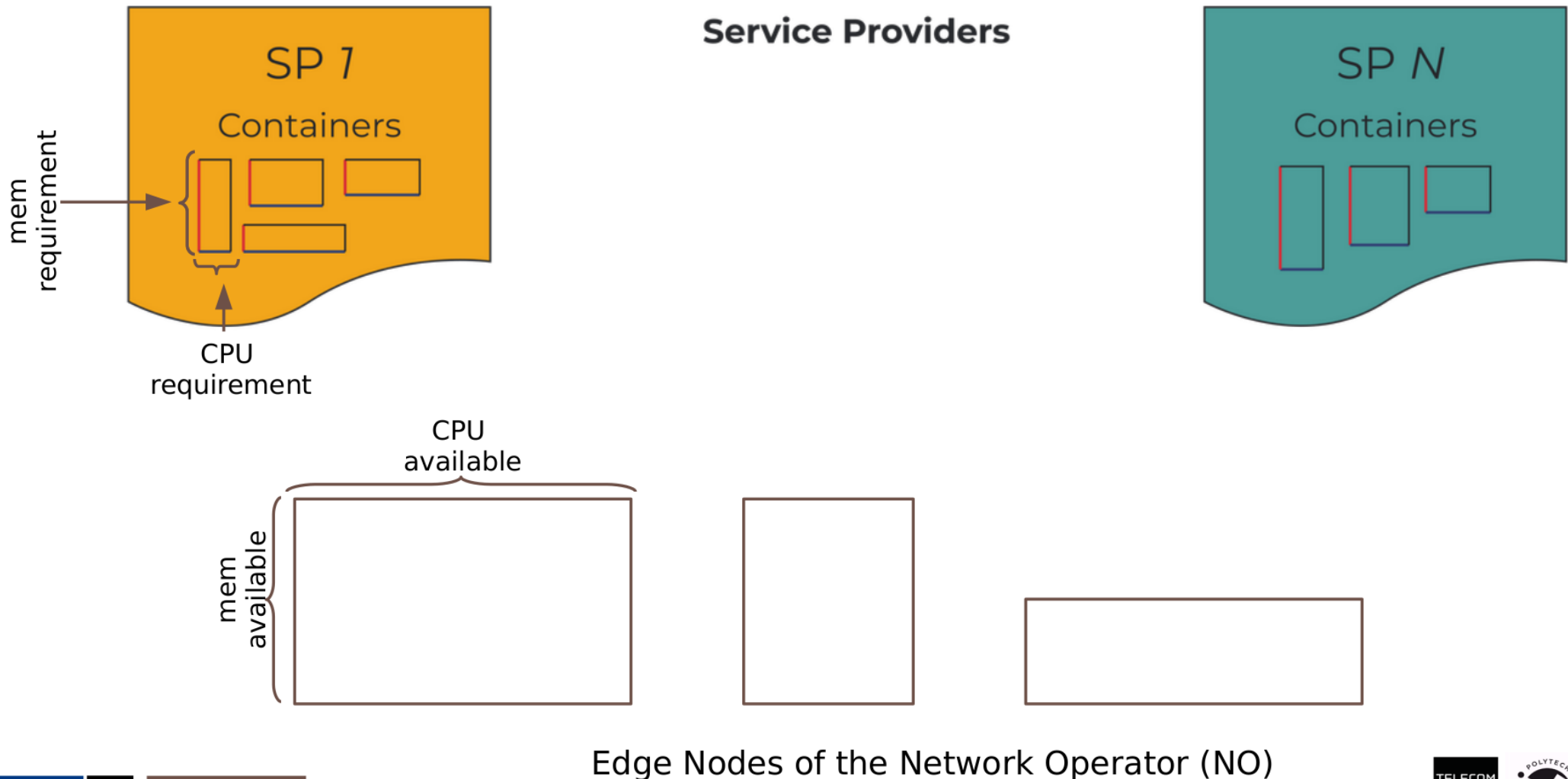
Service Providers





Microservice architecture

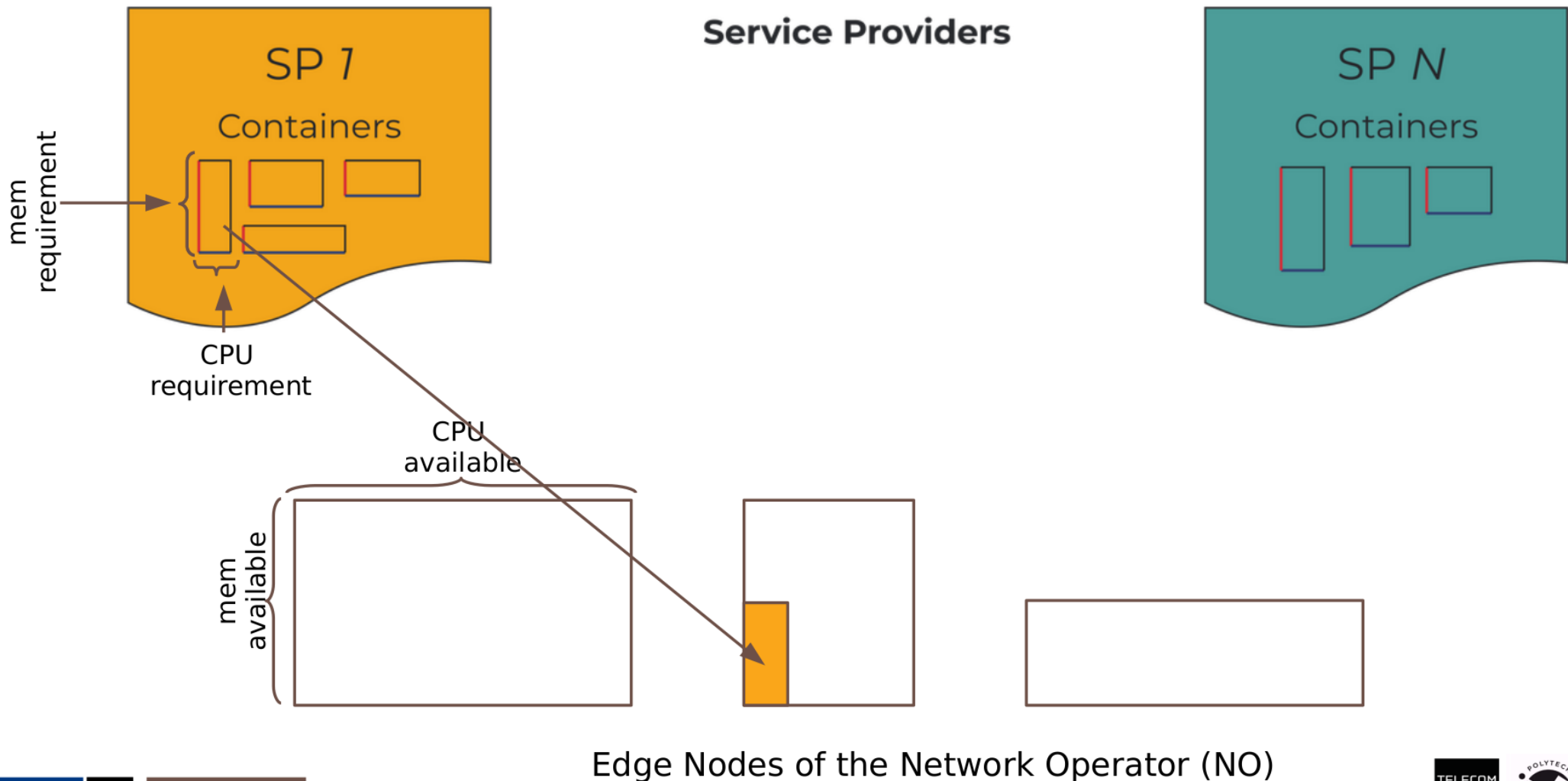
Netflix launches hundreds of thousands of containers every day [1]





Microservice architecture

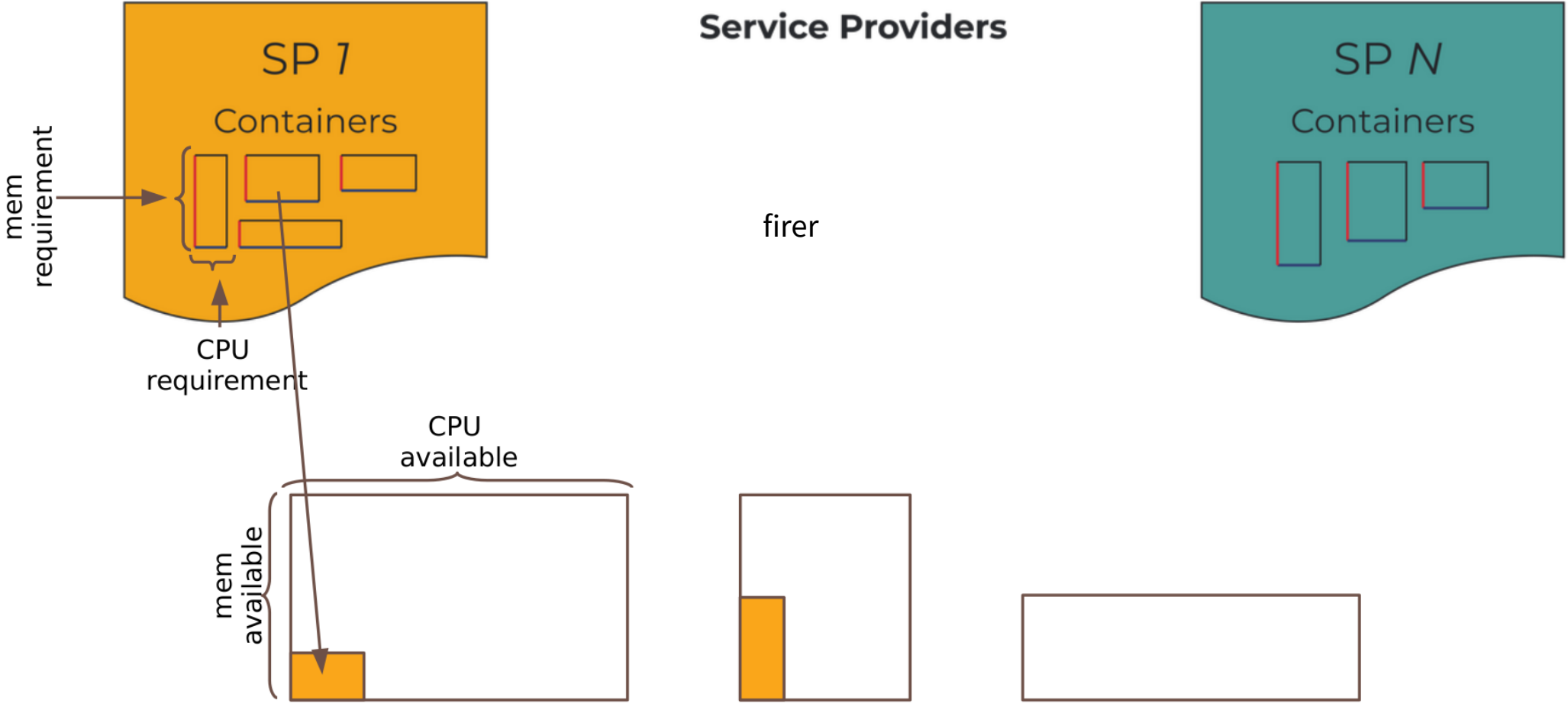
Netflix launches hundreds of thousands of containers every day [1]





Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]



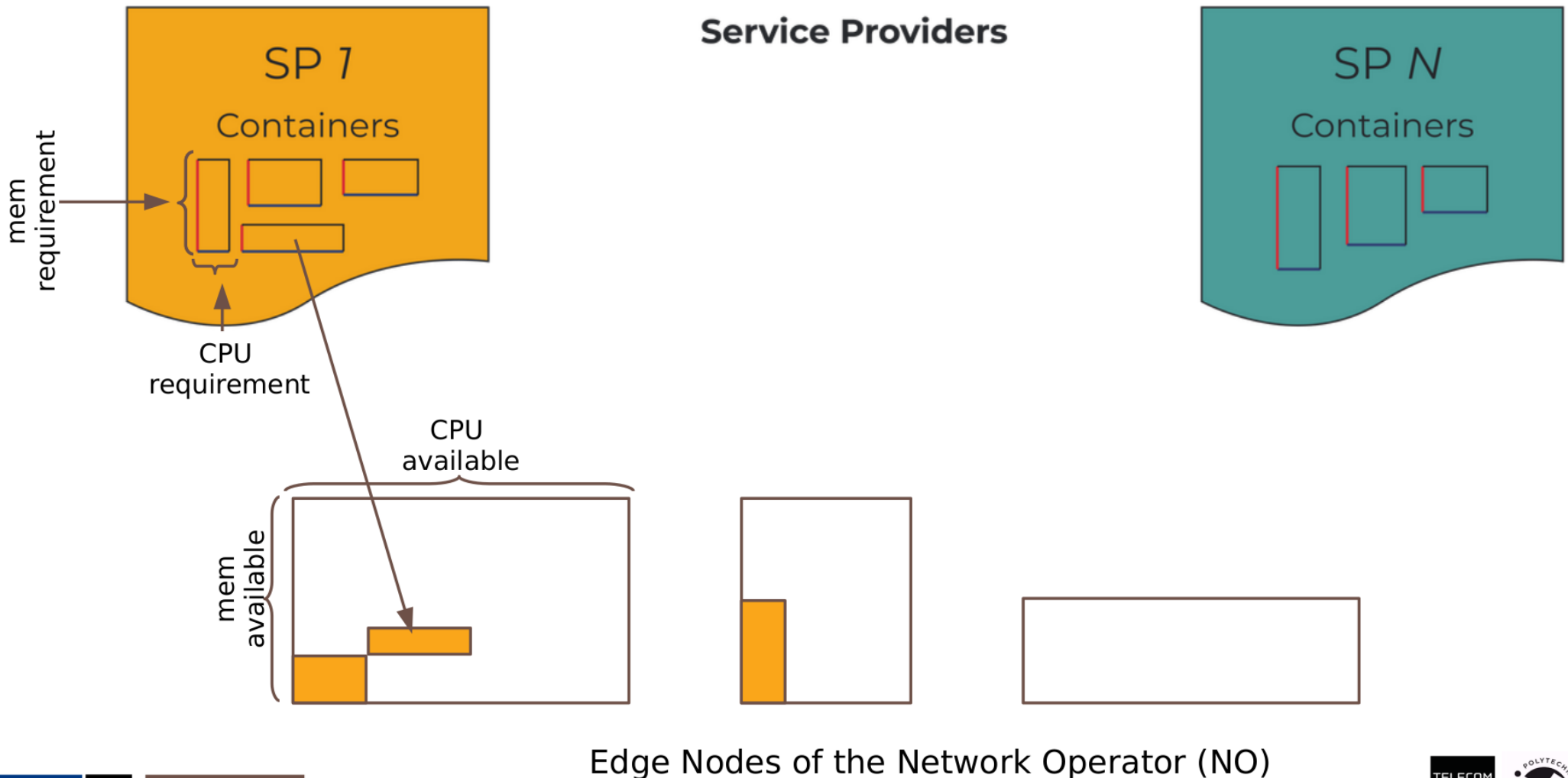
Edge Nodes of the Network Operator (NO)

[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.



Microservice architecture

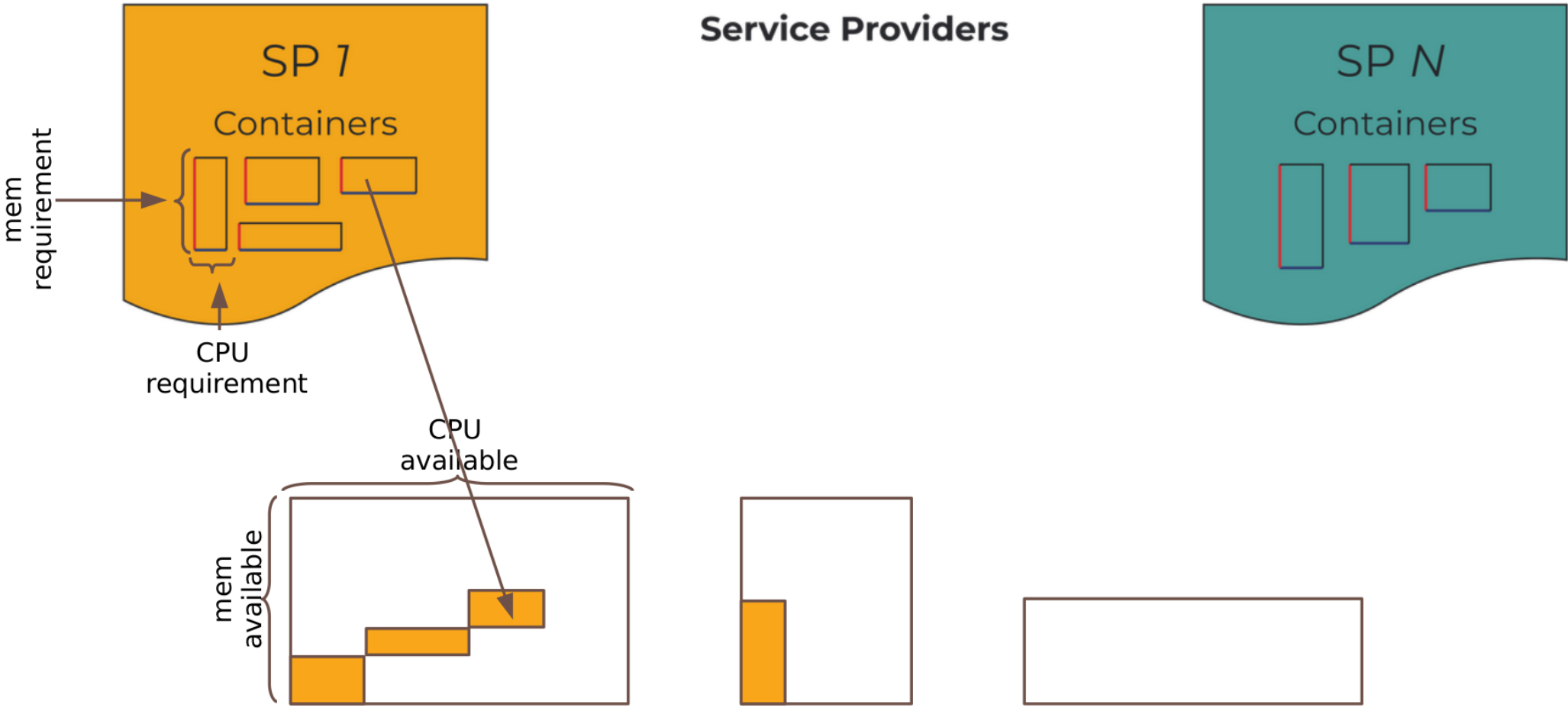
Netflix launches hundreds of thousands of containers every day [1]





Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]



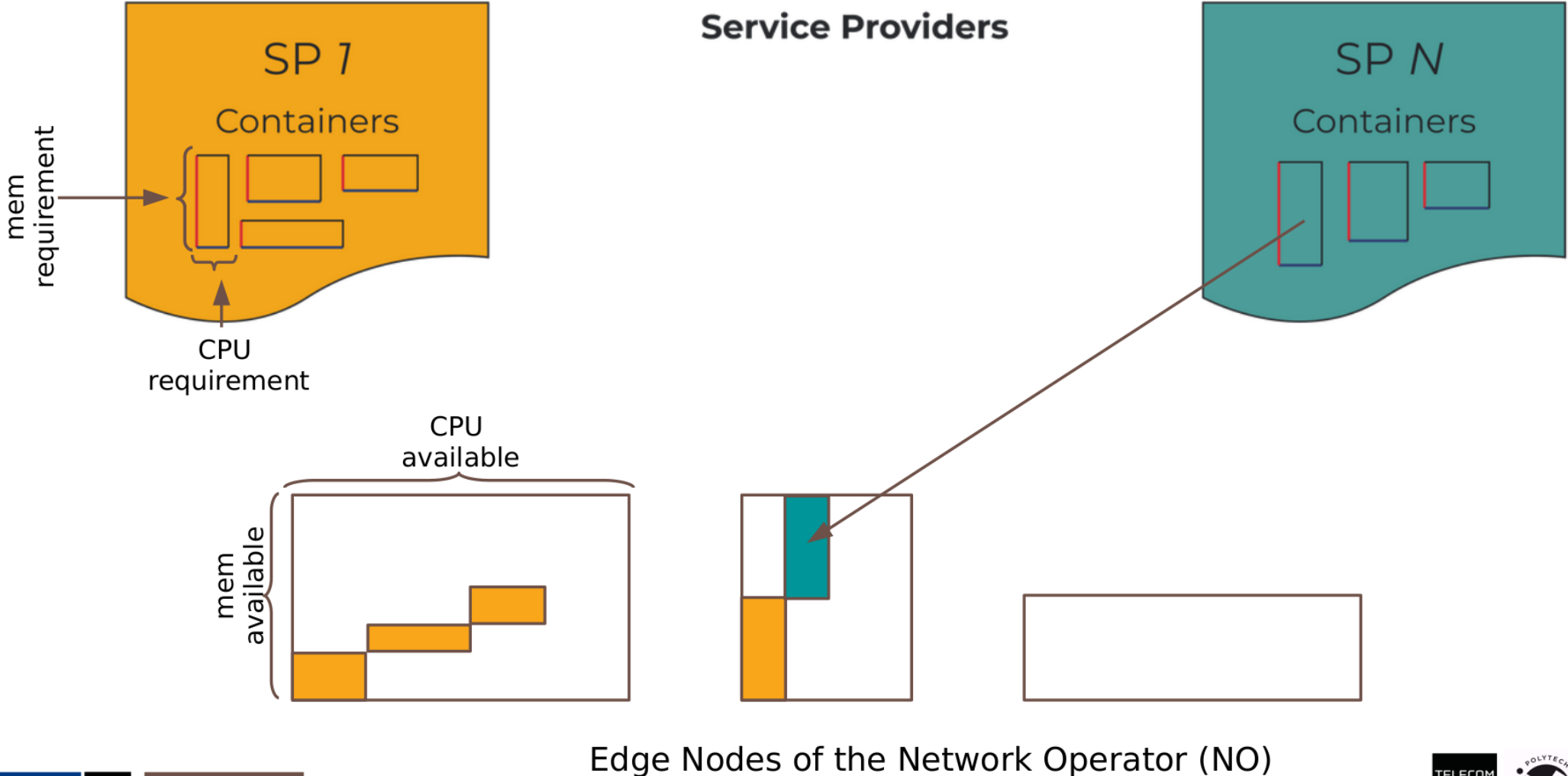
Edge Nodes of the Network Operator (NO)

[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.



Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]

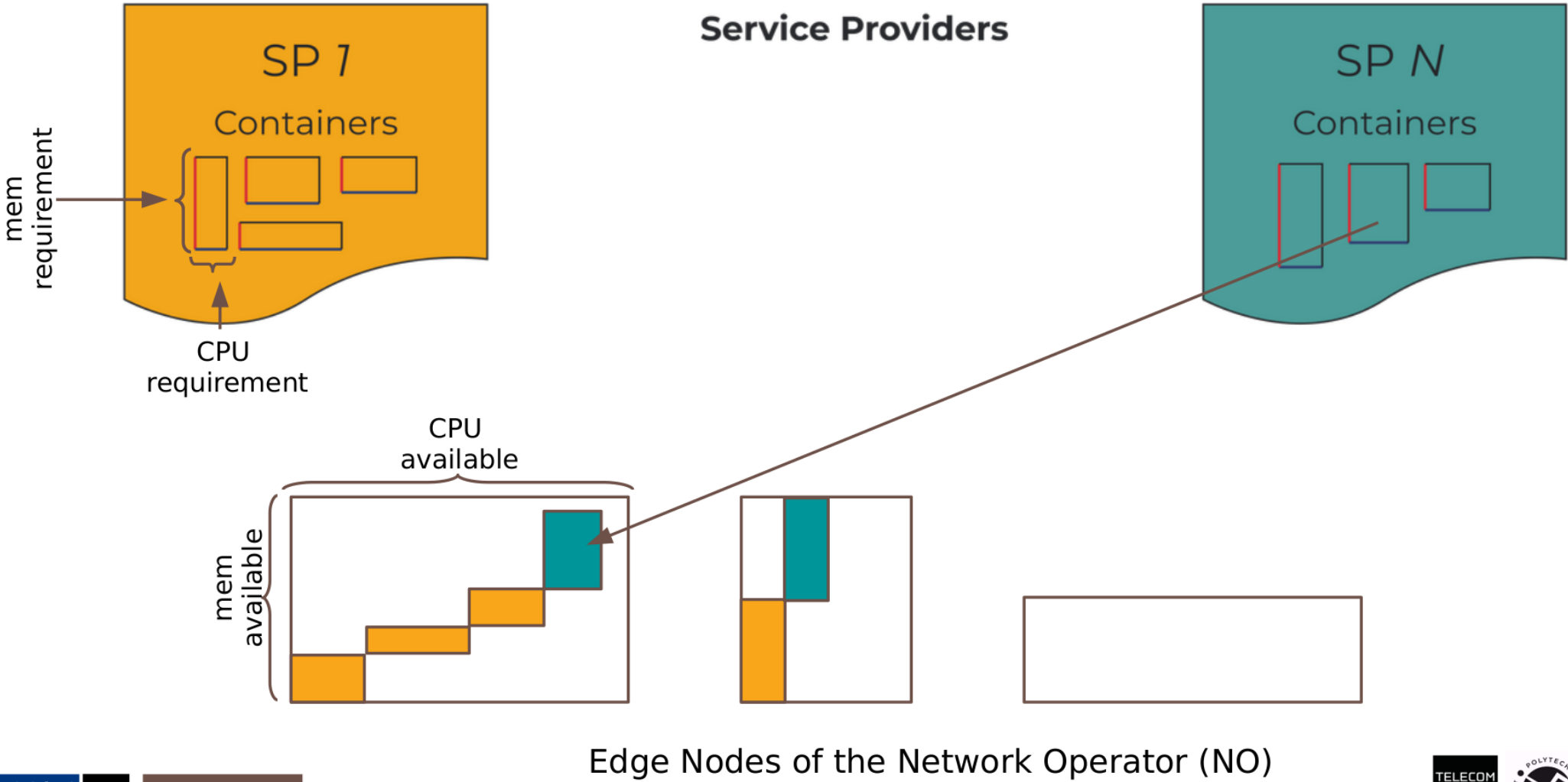


[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.



Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]

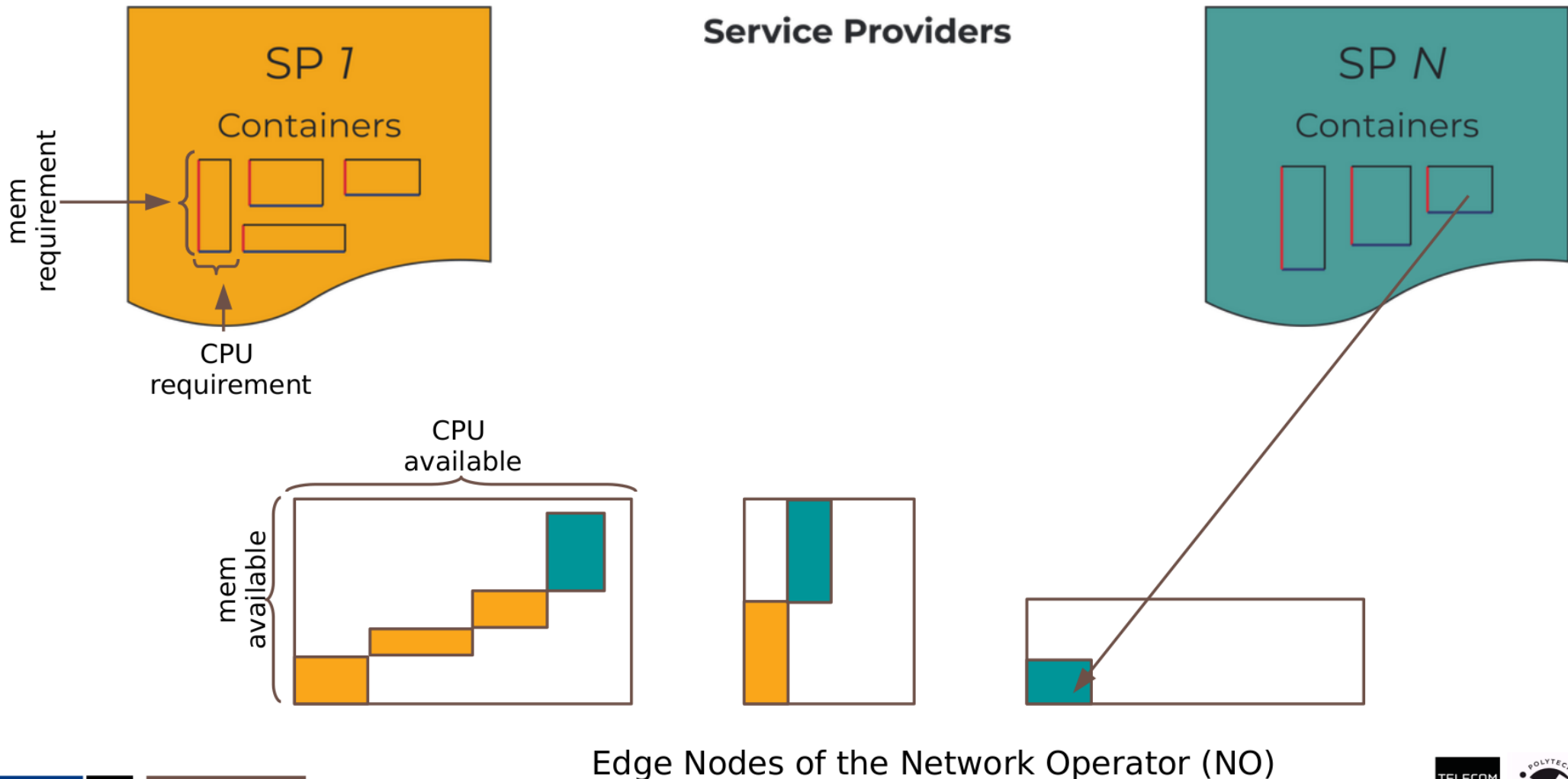


[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.



Microservice architecture

Netflix launches hundreds of thousands of containers every day [1]

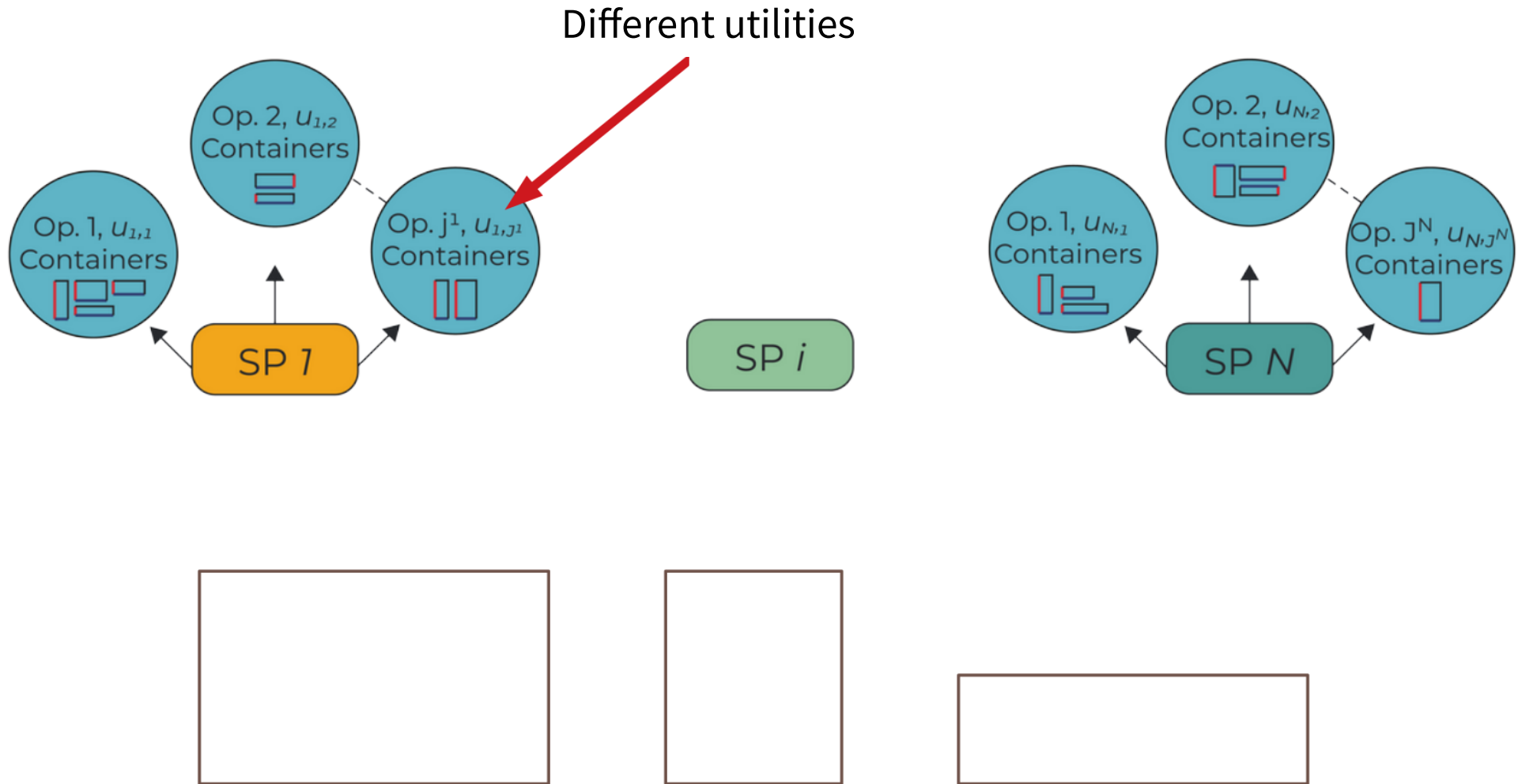


[1] Netflix. Titus. <https://netflix.github.io/titus/>, 2018.



Multiple Options

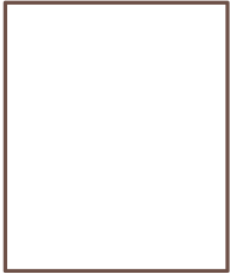
- From Resource Elasticity to **Service Elasticity**





Multiple Options

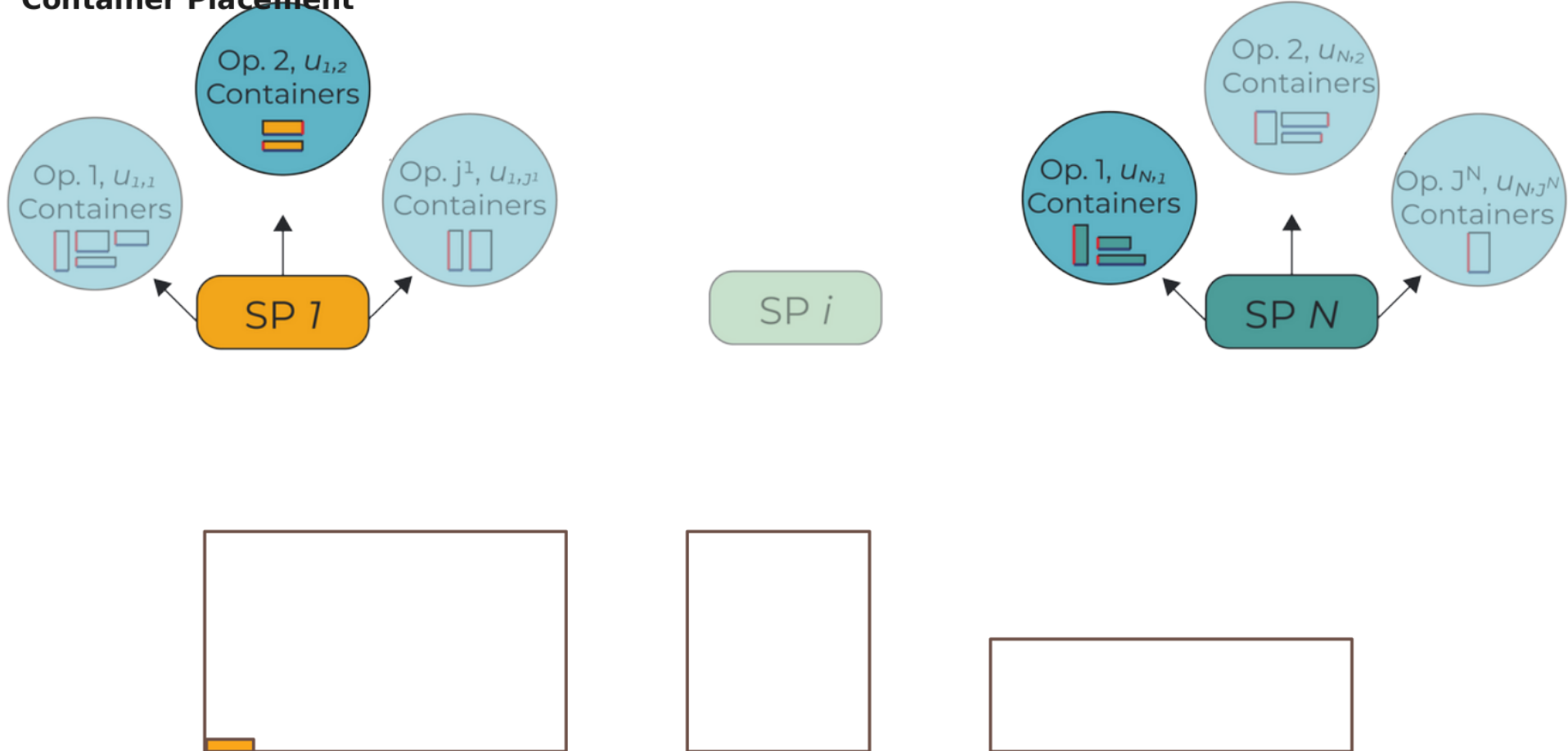
- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**





Multiple Options

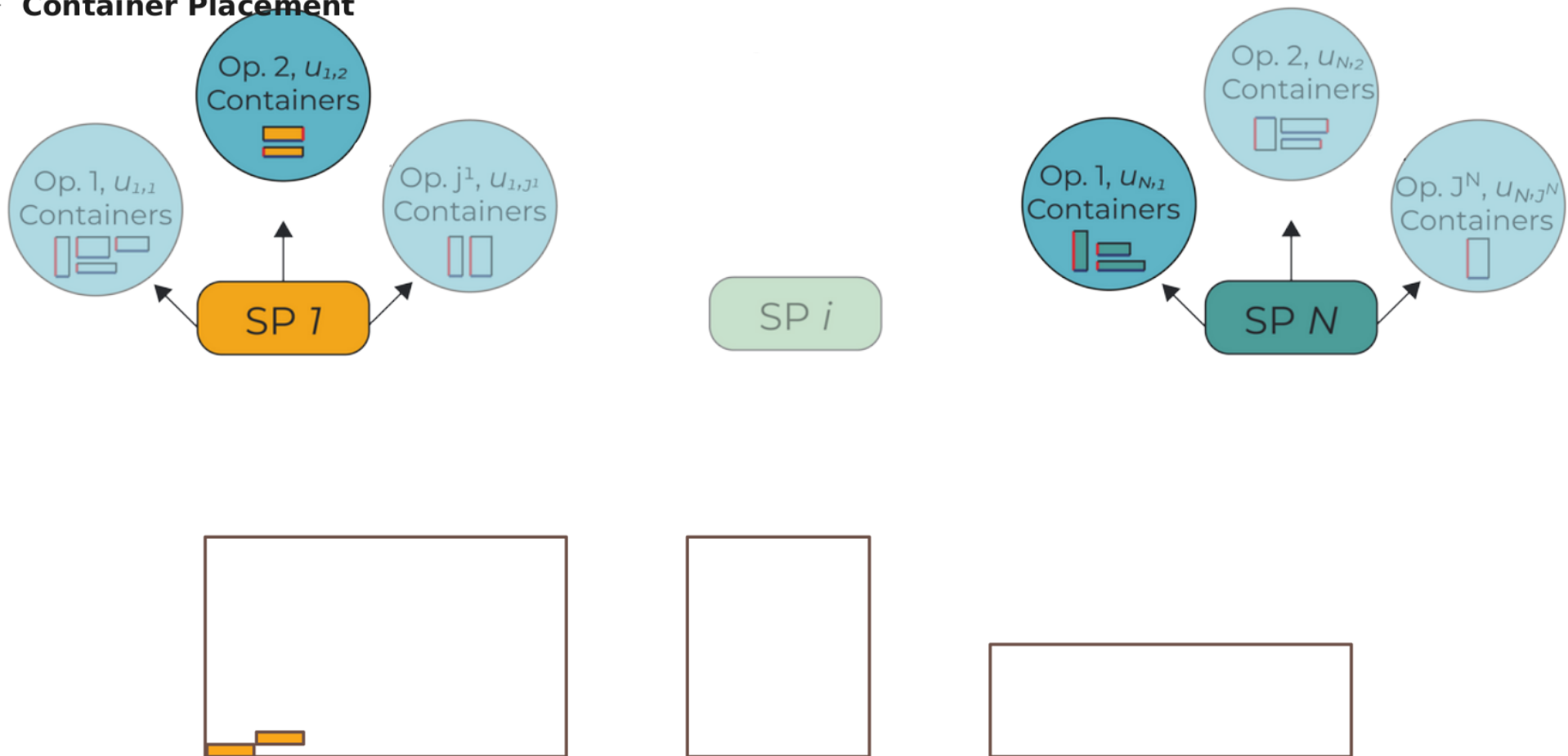
- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**





Multiple Options

- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**





Multiple Options

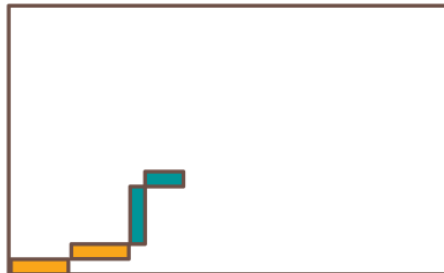
- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**





Multiple Options

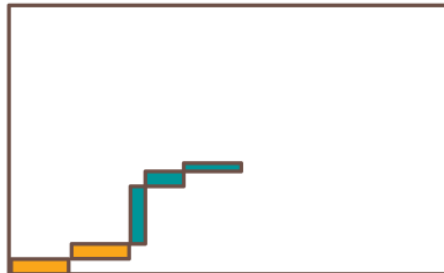
- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**





Multiple Options

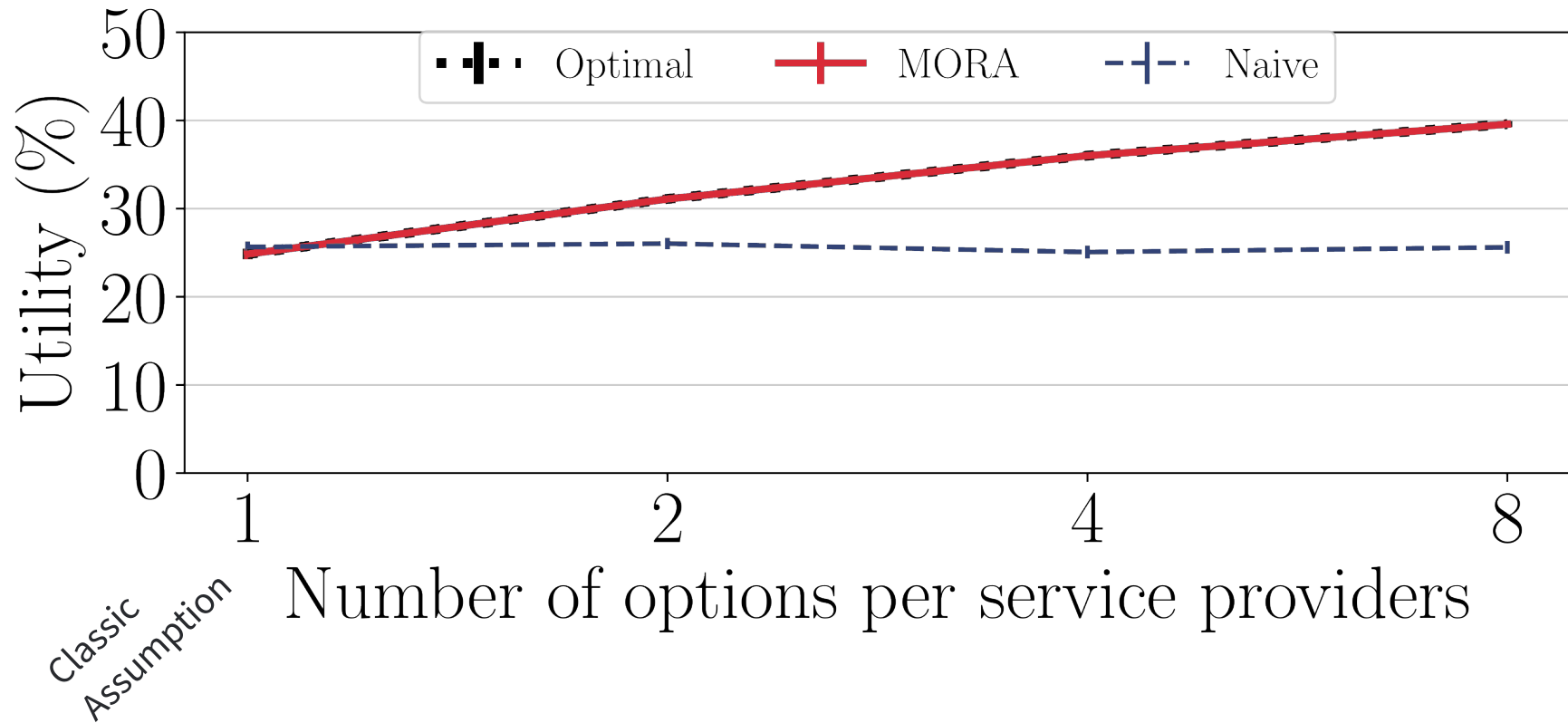
- Goal of the Netw. Operator: max utility
- Decisions:
 - **Option Selection**
 - **Container Placement**



Heuristic: performance



Araldo et Al, *Resource Allocation for Edge Computing with Multiple Tenant Configurations*, **ACM SAC 2020**




Utility is normalized with the maximum one (when selecting, for each SP, the option with the largest utility)

Current effort

- Proof of Concept (Docker and Kubernetes)
 - <https://github.com/mora-resource-allocation-edge-cloud/mora>

Conclusion and perspectives

- **Low latency services**
 - ← Need to open the edge to 3rd party service providers
- **Multi-tenant Edge Computing**
 - Virtualization
- **Data-driven resource allocation**
 - Stochastic Perturbation
 - Reinforcement Learning
- **Multiple Options Resource Allocation (MORA)** 
 - Service elasticity

Backup

- Backup

- **Subgradient** [1]

Definition 9. Given a function $\bar{L} : \mathbb{R}^P \rightarrow \mathbb{R}$, a function $\bar{g} : \mathcal{C} \subseteq \mathbb{R}^P \rightarrow \mathbb{R}^P$ is a **subgradient** of \bar{L} over \mathcal{C} iff

$$\bar{L}(\theta') - \bar{L}(\theta) \geq \bar{g}(\theta)^T \cdot (\theta' - \theta), \forall \theta, \theta' \in \mathcal{C}.$$

- **Supermartingale** [2]

Definition 10.2.1 Let $(M_n : n \geq 0)$ be an integrable sequence of random variables that is adapted to $(Z_n : n \geq 0)$. If for $n \geq 0$,

$$E[M_{n+1} | Z_0, \dots, Z_{n-1}] \leq M_n,$$

then $(M_n : n \geq 0)$ is said to be a **supermartingale** with respect to $(Z_n : n \geq 0)$.

- **Convergence theorem** [3]:

- All supermartingales with finite expectation converge almost surely

[1] Andrea Araldo ; György Dán ; Dario Rossi. (2018). Caching Encrypted Content Via Stochastic Cache Partitioning. IEEE/ACM Transactions on Networking, 26(1), 548–561.

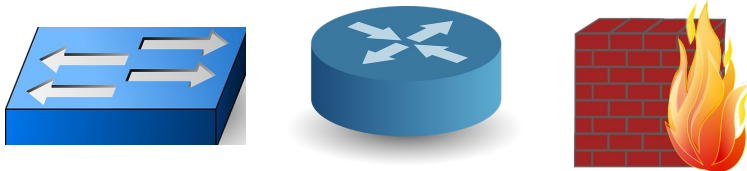
[2] Glynn, P. W. (2013). Martingales.

[3] https://en.wikipedia.org/wiki/Doob%27s_martingale_convergence_theorems

Slicing / Edge Computing

- **Slicing**

- Several logical networks on top of a real network
- Software network components



- Network operations

- **Edge Computing (EC)**

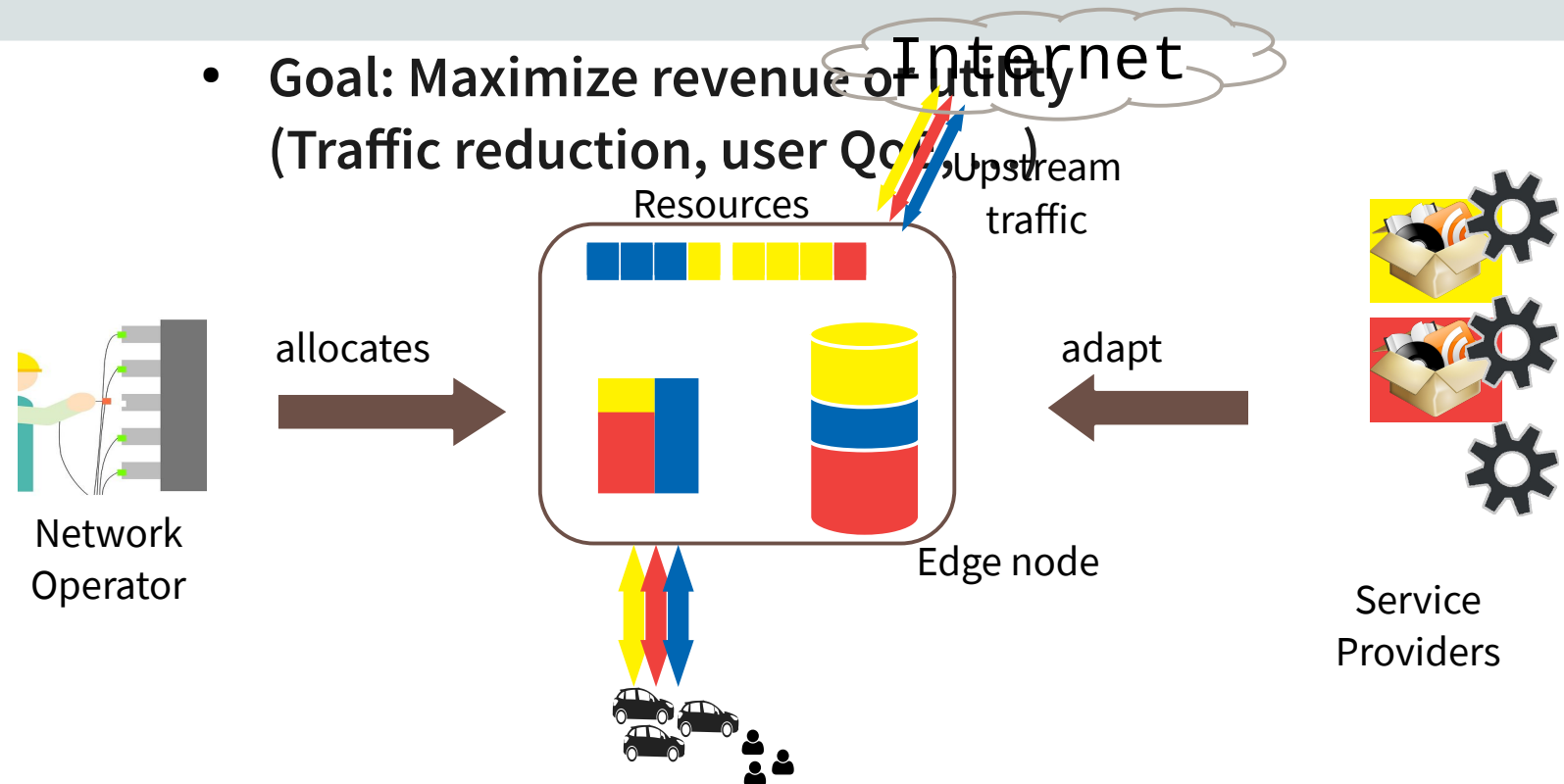
- Performing service computation at the Edge



Cloud and Edge: differences

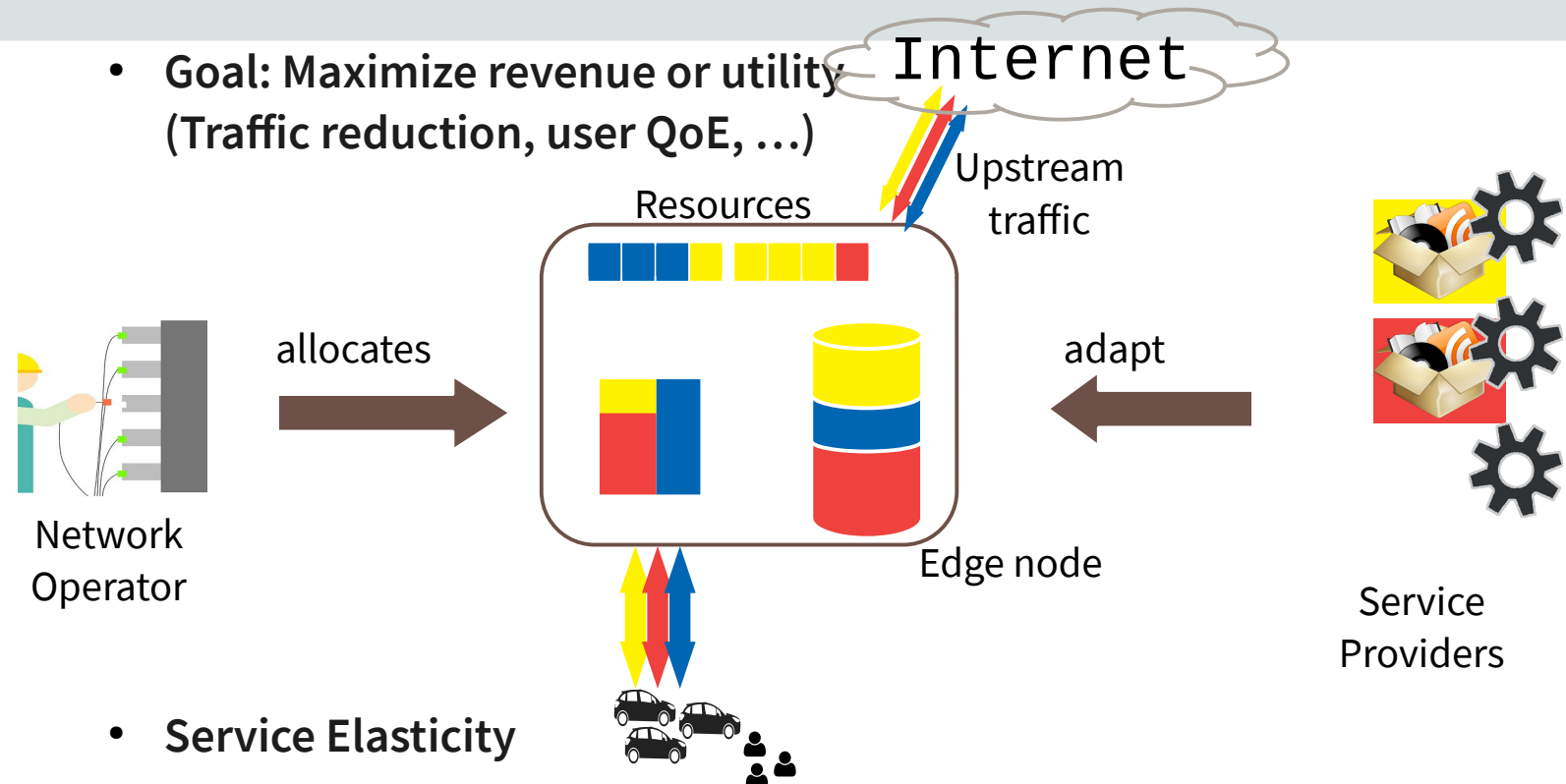
	Cloud	Edge
Amount of resources	Infinite	Limited → Contention
Objective	Maximize profit	Maximize profit (from rent) Maximize benefit (QoE, Traffic reduction)
Decision	Pricing strategy	Pricing Strategy + Allocation

Assumptions



- Service Elasticity

Assumptions



- **Goal: Maximize revenue or utility**
(Traffic reduction, user QoE, ...)

- **Service Elasticity**
- **Confidentiality and Isolation of Service Providers**

(traffic encryption, memory encryption)

- **They must be treated as black boxes**

==> Optimal allocation?

- **Two approaches**

(i) Data driven (ii) Multiple Options

Challenges and opportunities

- **Confidentiality and Isolation of Service Providers (SPs)**
 - ==> SPs are black boxes
 - ==> Allocation Benefit must be learned from measures
- **Dynamicity**
 - Allocation must evolve over time
- **Service Elasticity**

Challenges and opportunities

- **Confidentiality and Isolation of Service Providers (SPs)**

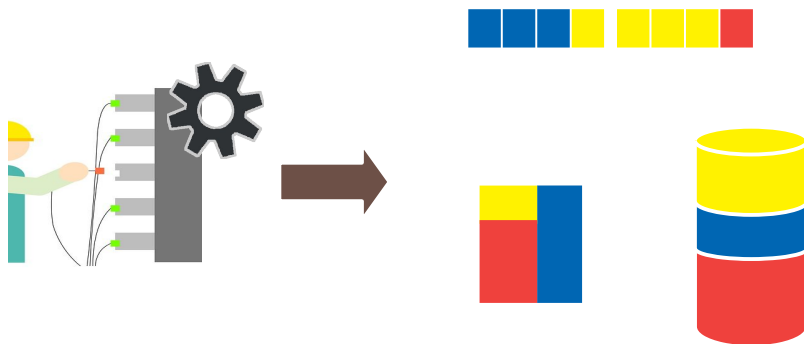
==> SPs are black boxes

==> Allocation Benefit must be learned from measures

- **Dynamicity**

– Allocation must evolve over time

- **Service Elasticity**



Challenges and opportunities

- **Confidentiality and Isolation of Service Providers (SPs)**

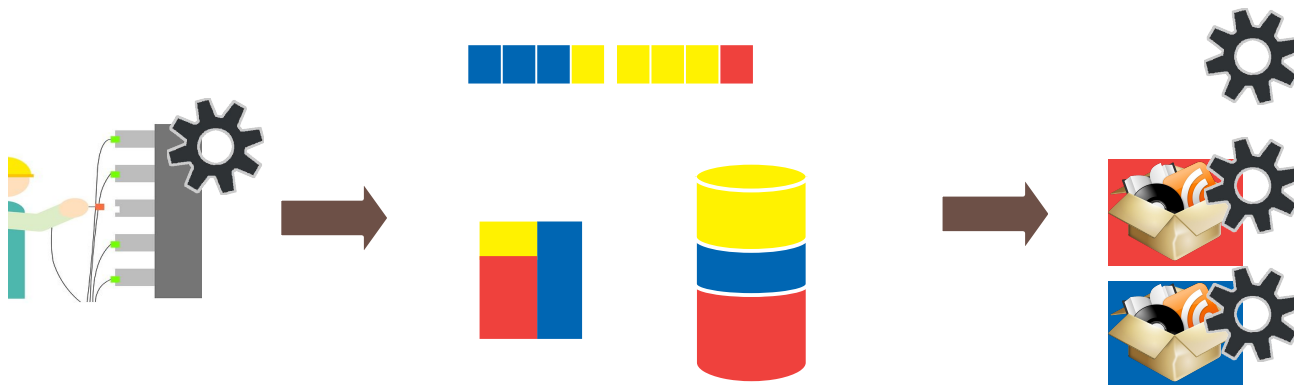
==> SPs are black boxes

==> Allocation Benefit must be learned from measures

- **Dynamicity**

– Allocation must evolve over time

- **Service Elasticity**



Challenges and opportunities

- **Confidentiality and Isolation of Service Providers (SPs)**

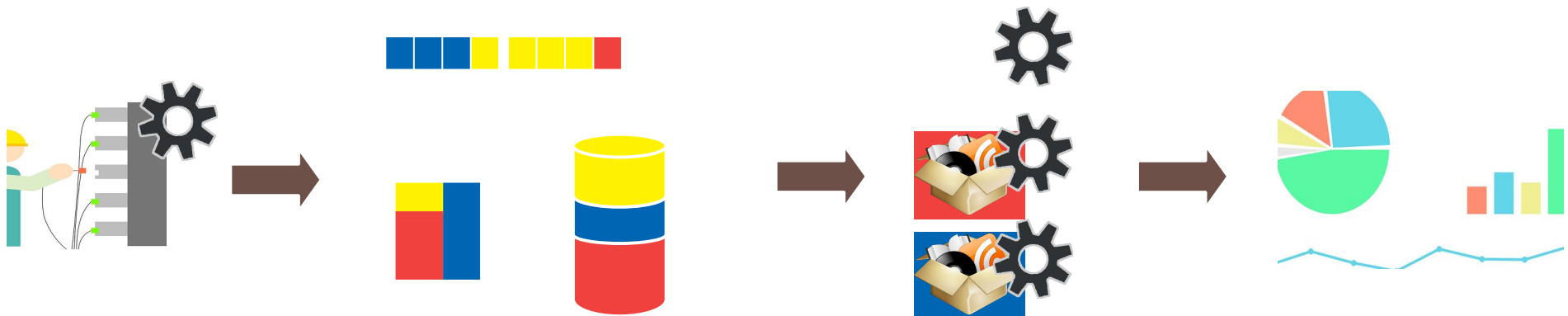
==> SPs are black boxes

==> Allocation Benefit must be learned from measures

- **Dynamicity**

– Allocation must evolve over time

- **Service Elasticity**



Challenges and opportunities

- **Confidentiality and Isolation of Service Providers (SPs)**

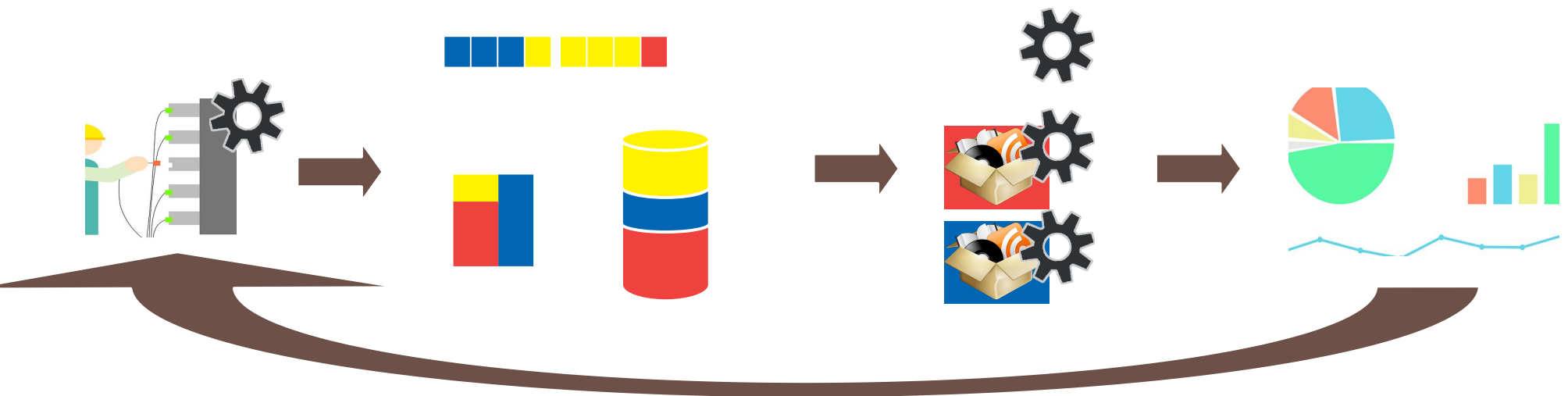
==> SPs are black boxes

==> Allocation Benefit must be learned from measures

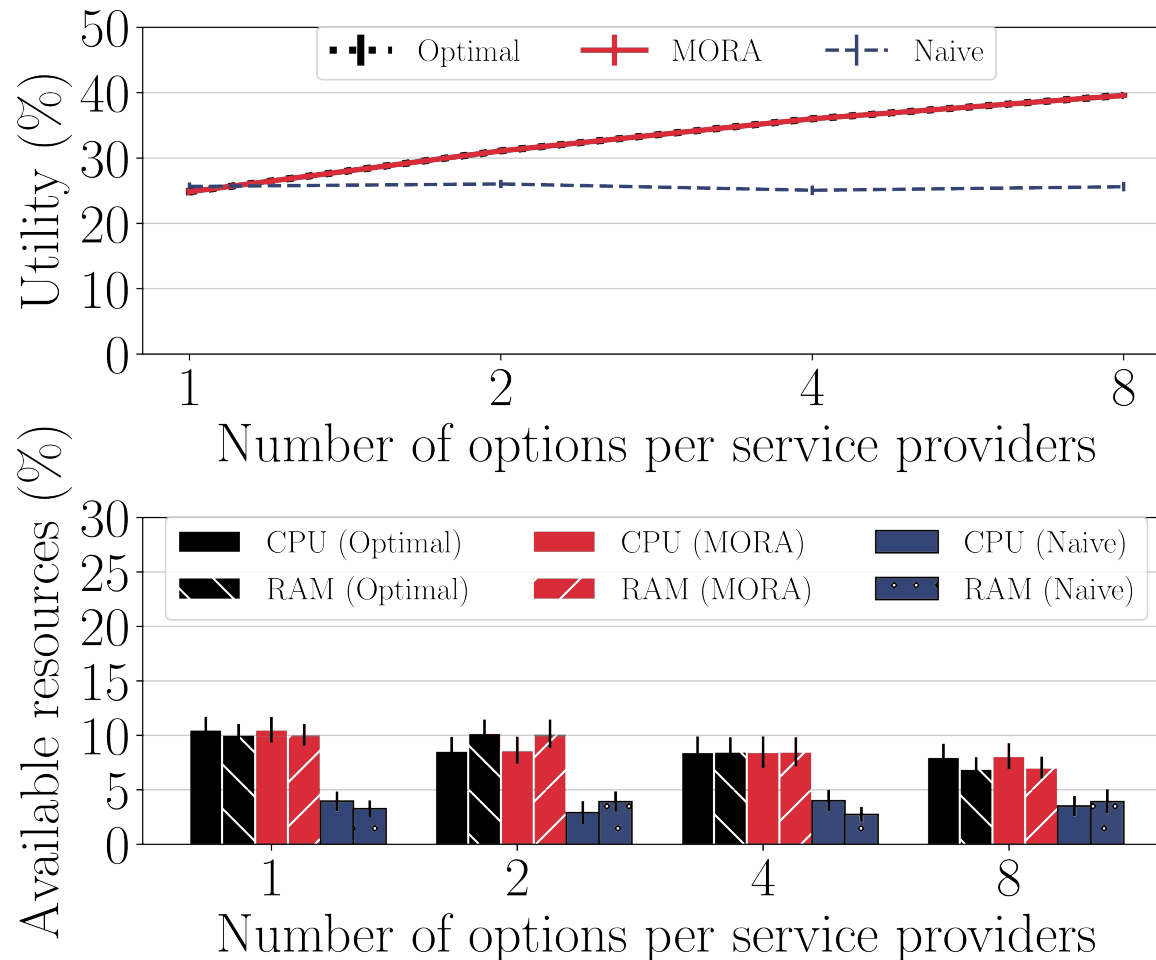
- **Dynamicity**

- Allocation must evolve over time

- **Service Elasticity**



Heuristic: performance



Utility is normalized with the maximum one (when selecting, for each SP, the option with the largest utility)

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $\mathbf{a}=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $a=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			



Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $\mathbf{a}=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

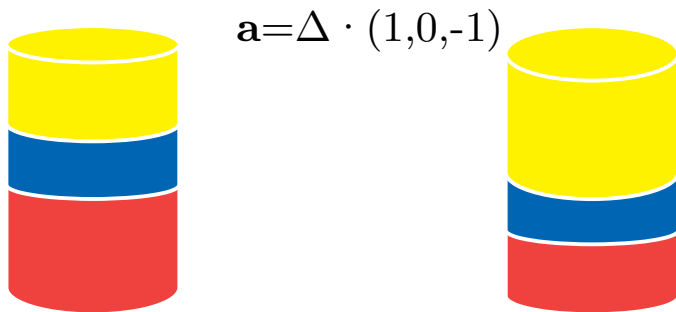


$$\mathbf{a}=\Delta \cdot (1,0,-1)$$

Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $\mathbf{a}=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

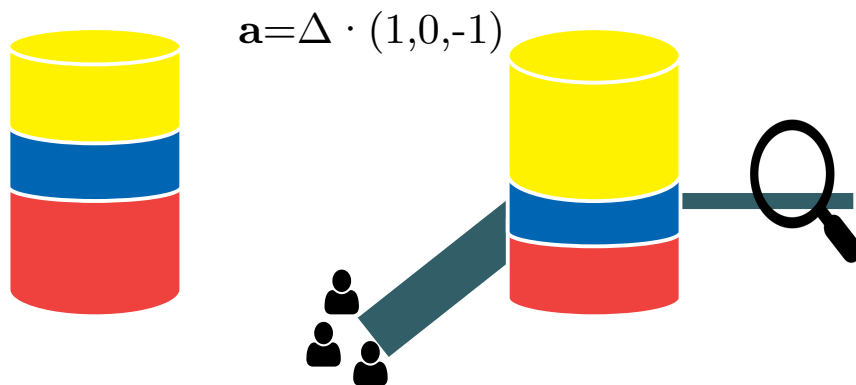
	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			



Reinforcement Learning for Edge Cache Allocation [1]

- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $\mathbf{a}=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Reinforcement Learning for Edge Cache Allocation [1]

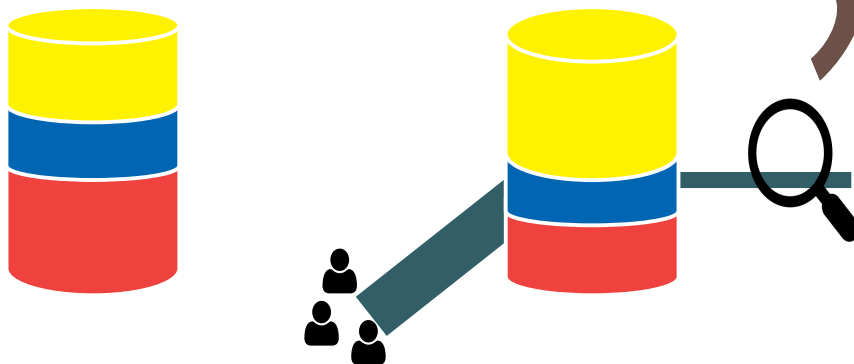
- State: allocation $\theta=(\theta_1,\dots,\theta_p)$
- Action: perturbation, e.g. $\mathbf{a}=\Delta \cdot (1,0,-1)$
- Instantaneous cost: upstream traffic in 1 s
- Q-table
(estimations of cumulative costs)

	action 1	action 2
allocation 1	C_{11}	C_{12}	
allocation 2	C_{21}	C_{22}	
....			

SARSA algorithm

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

- We learn a good Q-table by perturbing-and-observing the system



[1] T. Bouganim, A. **Araldo** et Al., “The Cost of Learning Fast with Reinforcement Learning for Edge Cache Allocation”, ITC PhD Workshop, 2020

Augmented Reality with EdgeAI Devices

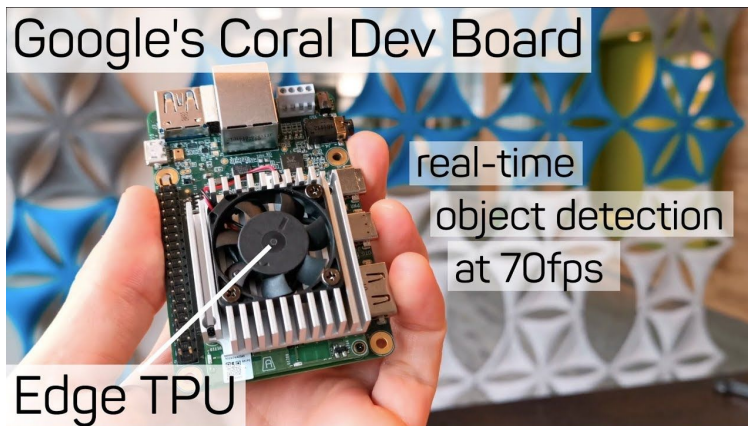
A. Ben-Ameur, A. Araldo, F. Bronzino,

*On the Deployability of Augmented Reality Using Embedded Edge
Devices,*

IEEE CCNC 2021

NOKIA Bell Labs

EdgeAI devices



~150\$



Pictures from:

<https://www.youtube.com/watch?v=bOYWx1jJCZo>

<https://www.phoronix.com/scan.php?page=article&item=nvidia-jetson-nano&num=1>

Centralized vs. Distributed Architecture

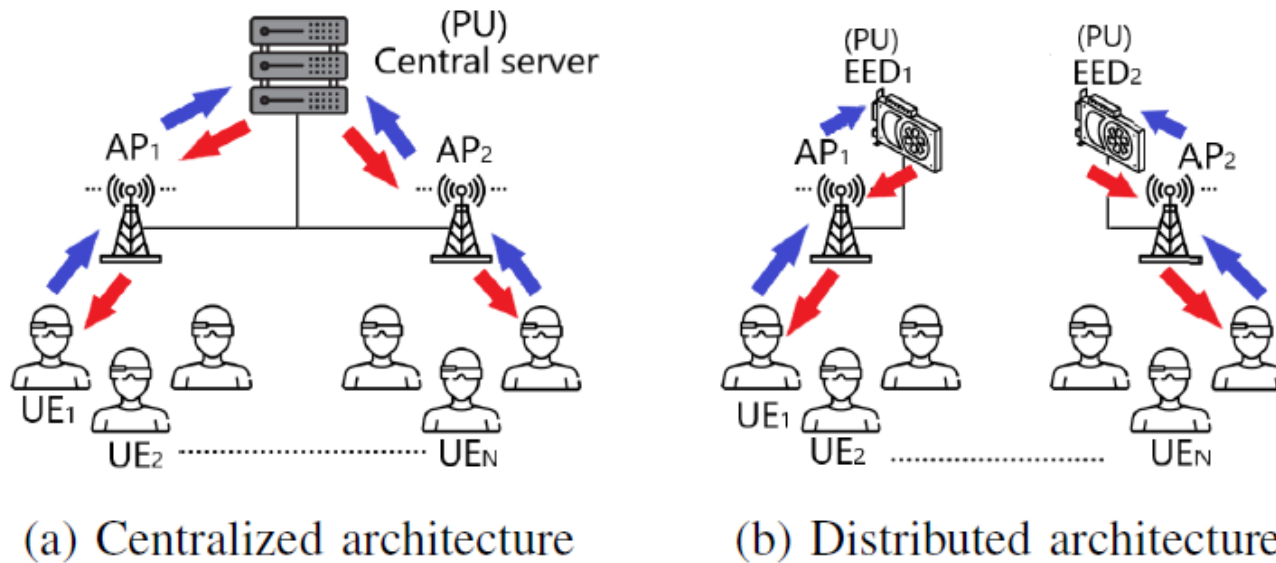


Fig. 1: The centralized vs. distributed architecture

- **Measurements**
 - Latency and recognition accuracy
 - Bare metal server, Google TPU Board, Jetson Nano
- **Analytical Model**
- **NS3 simulation**

TABLE I: Augmented Reality requirements

AR requirements	Latency
Low Responsiveness (LR)	500 ms [6]
Mid Responsiveness (MR)	100 ms
High Responsiveness (HR)	16 ms [11]

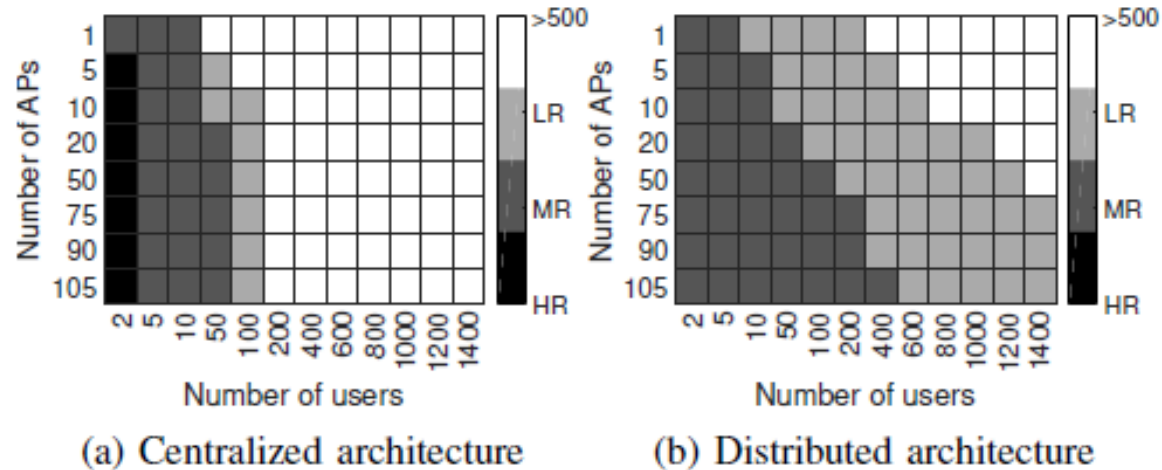


Fig. 5: Achievable requirements for $R = 450$ Mbps.