Orchestrating a brighter world

NEC

5GTRANSFORMER

# Overbooking Network Slices through Yield-driven End-to-End Orchestration
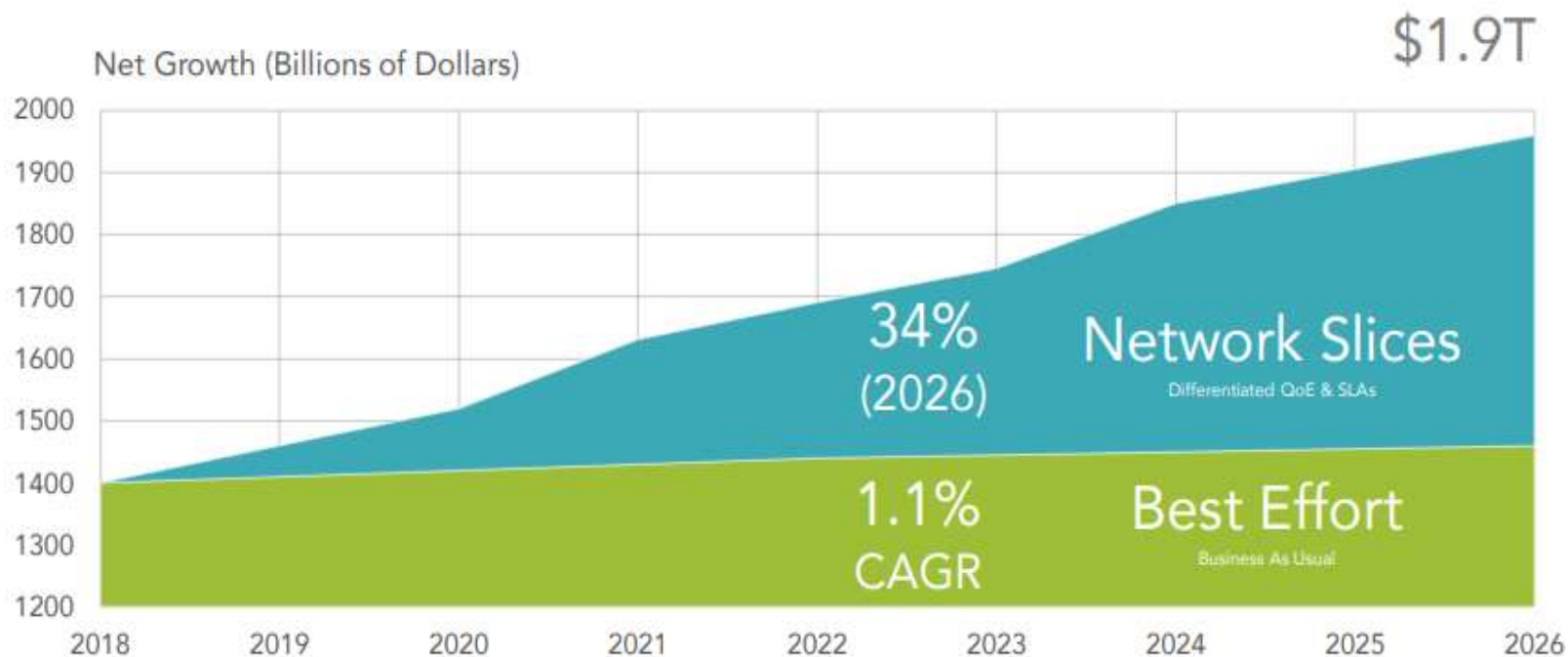
Josep Xavier Salvat

Lanfranco Zanzi

Andres Garcia-Saavedra

Vincenzo Sciancalepore

Xavier Costa-Perez

**NEC Laboratories Europe**

# Introduction: Network Slicing

▌ (possibly hyper-hyped but...) Telcos and vendors expect that **NS to unlock around $300bn in business opportunities** with verticals/private enterprises
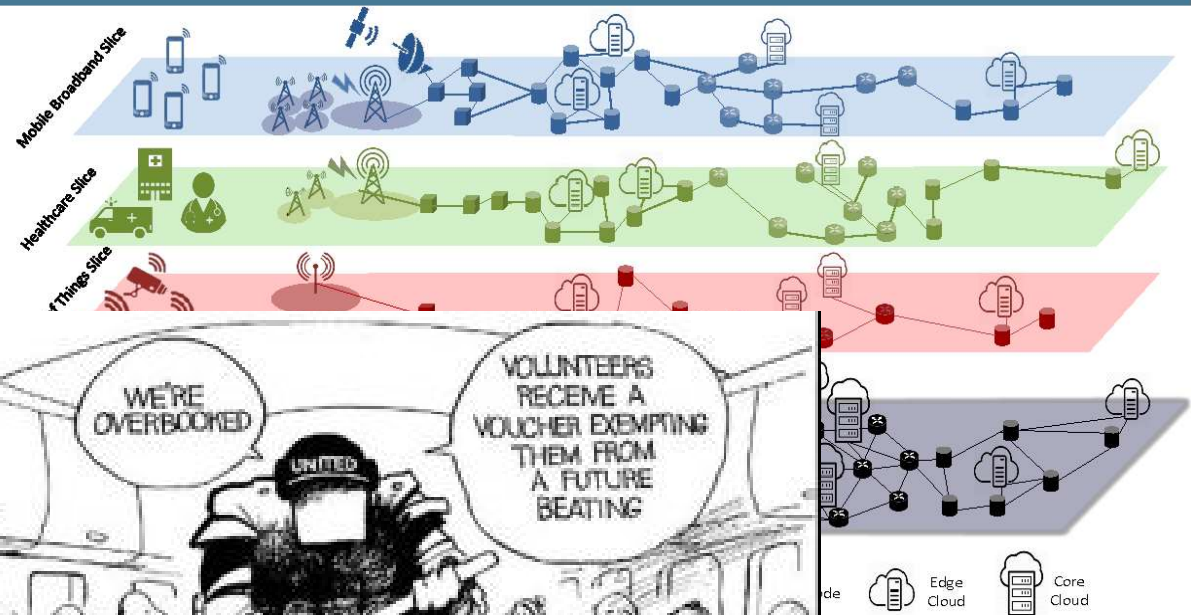
## 34% GROWTH TO 2026 FROM MOVING OFF "BEST EFFORT"

Net Growth (Billions of Dollars)

$1.9T

34% (2026)

**Network Slices**
Differentiated QoE & SLAs

1.1% CAGR

**Best Effort**
Business As Usual

Source: Ericsson, Arthur D. Little

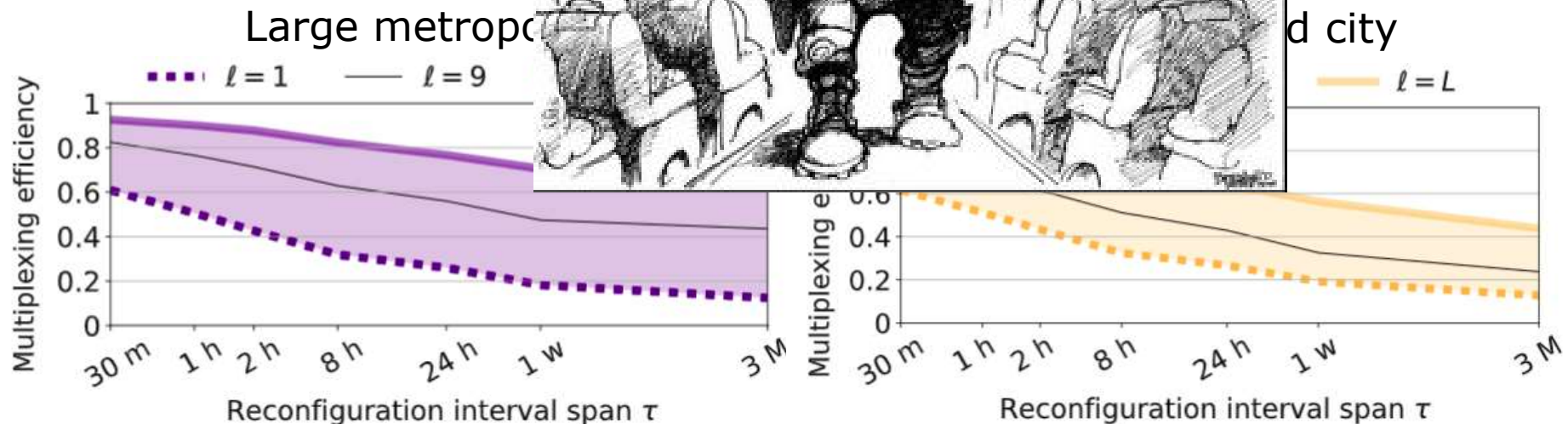\Orchestrating a brighter world  NEC

**1. End-to-End Orchestration**

**2. Network Slice Overbooking**
(Admission Control and Resource Reservation mechanism)



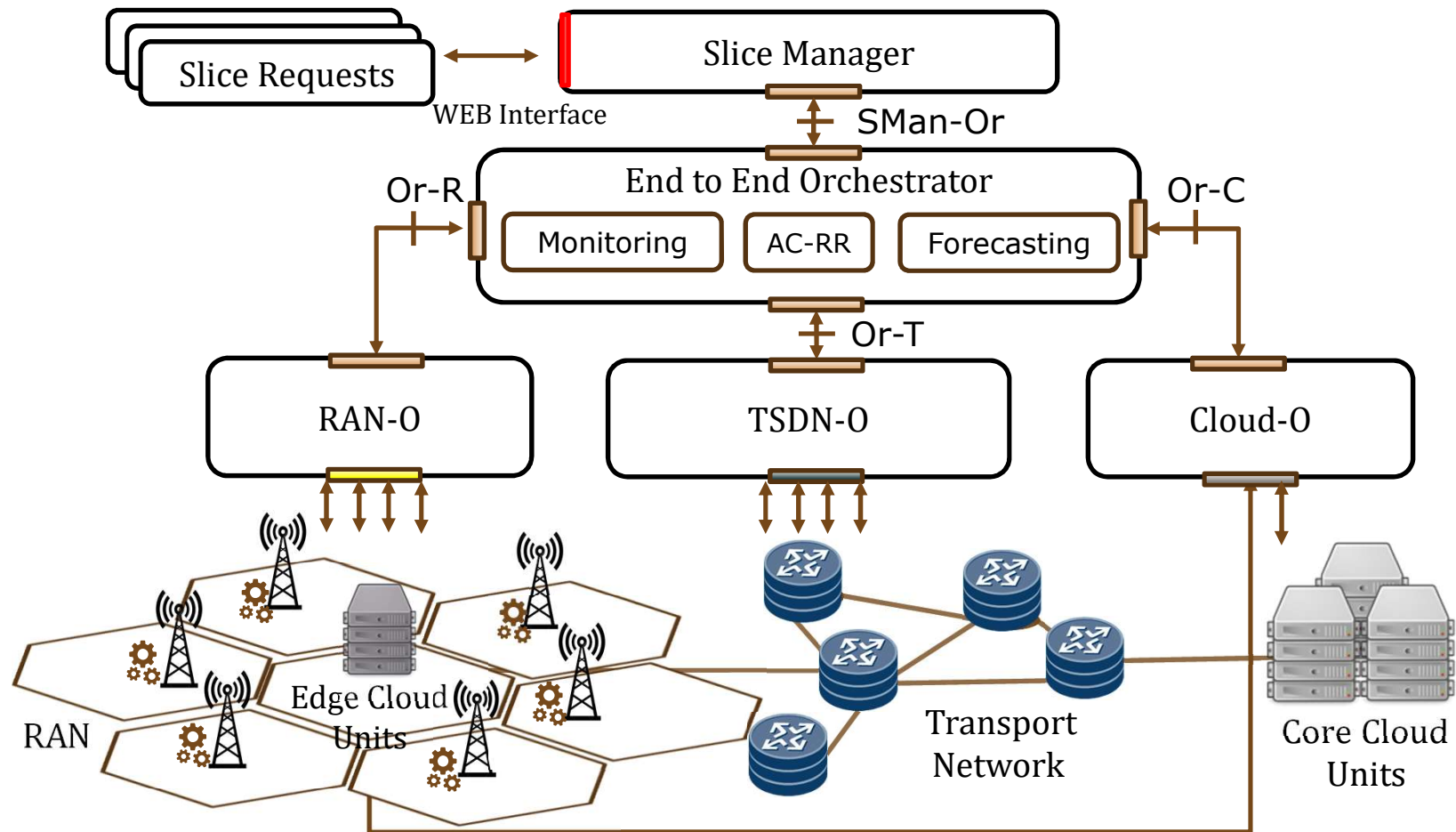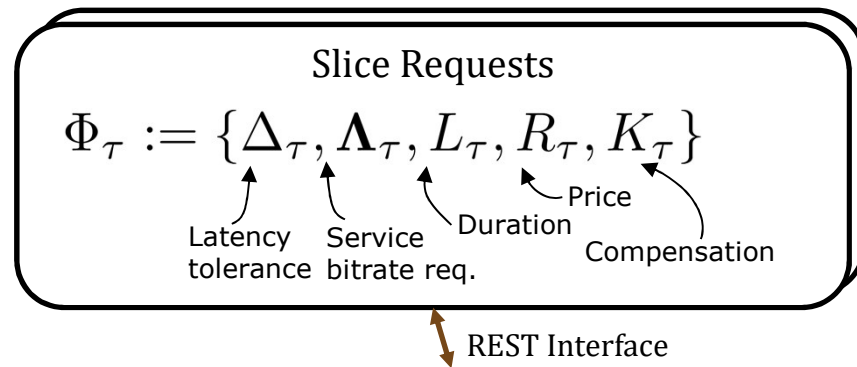NFV: Concepts, Architectures, pp. 80-87, May 2017.

Large metropo... ...d city

C. Marquez et al. "How Should I Slice My Network?: A Multi-Service Empirical Evaluation of Resource Sharing Efficiency". In ACM MobiCom 2018
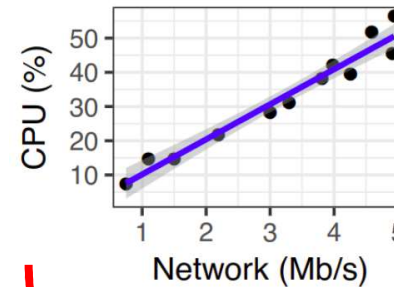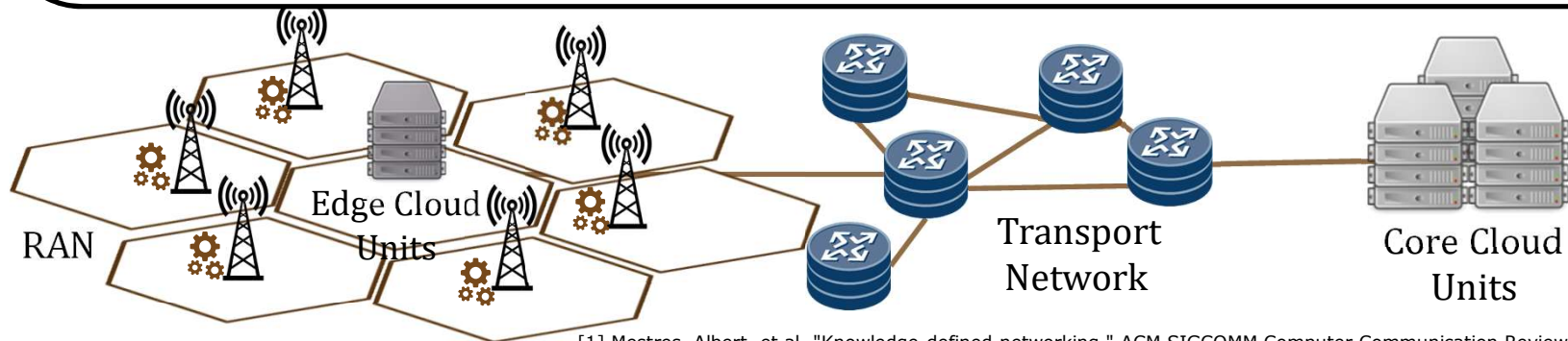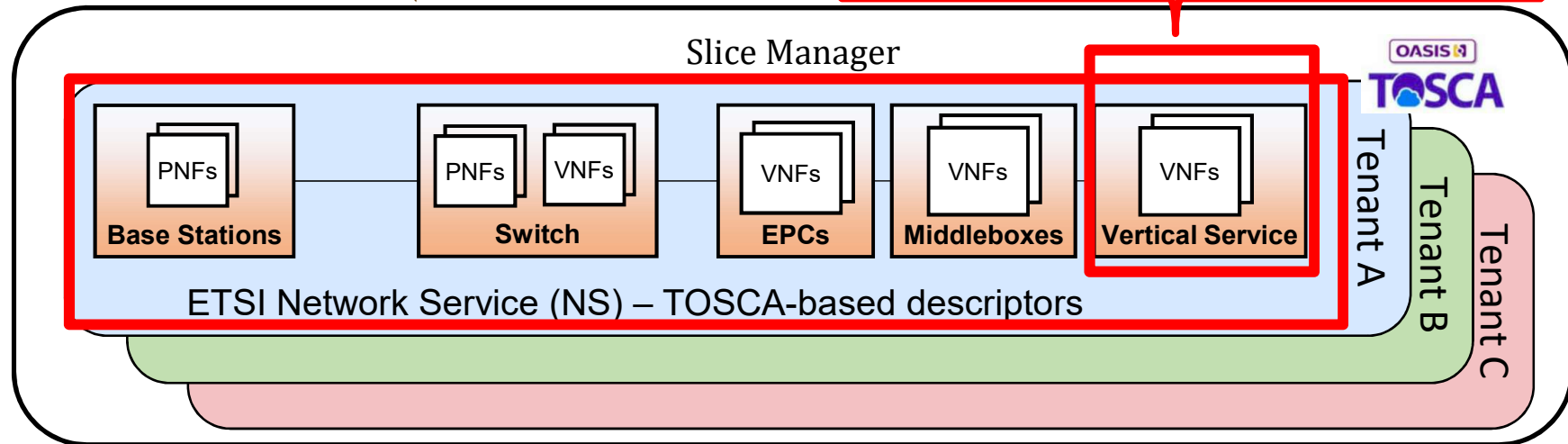
© NEC Laboratories Europe GmbH 2018

\Orchestrating a brighter world **NEC**

# System Design and Model

Slice Requests

$$\Phi_\tau := \{\Delta_\tau, \Lambda_\tau, L_\tau, R_\tau, K_\tau\}$$

Latency tolerance · Service bitrate req. · Duration · Price · Compensation

REST Interface

e.g., OpenFace:

e.g., snort [1]

Slice Manager

OASIS TOSCA

PNFs — Base Stations

PNFs VNFs — Switch

VNFs — EPCs

VNFs — Middleboxes

VNFs — Vertical Service

Tenant A / Tenant B / Tenant C

ETSI Network Service (NS) – TOSCA-based descriptors
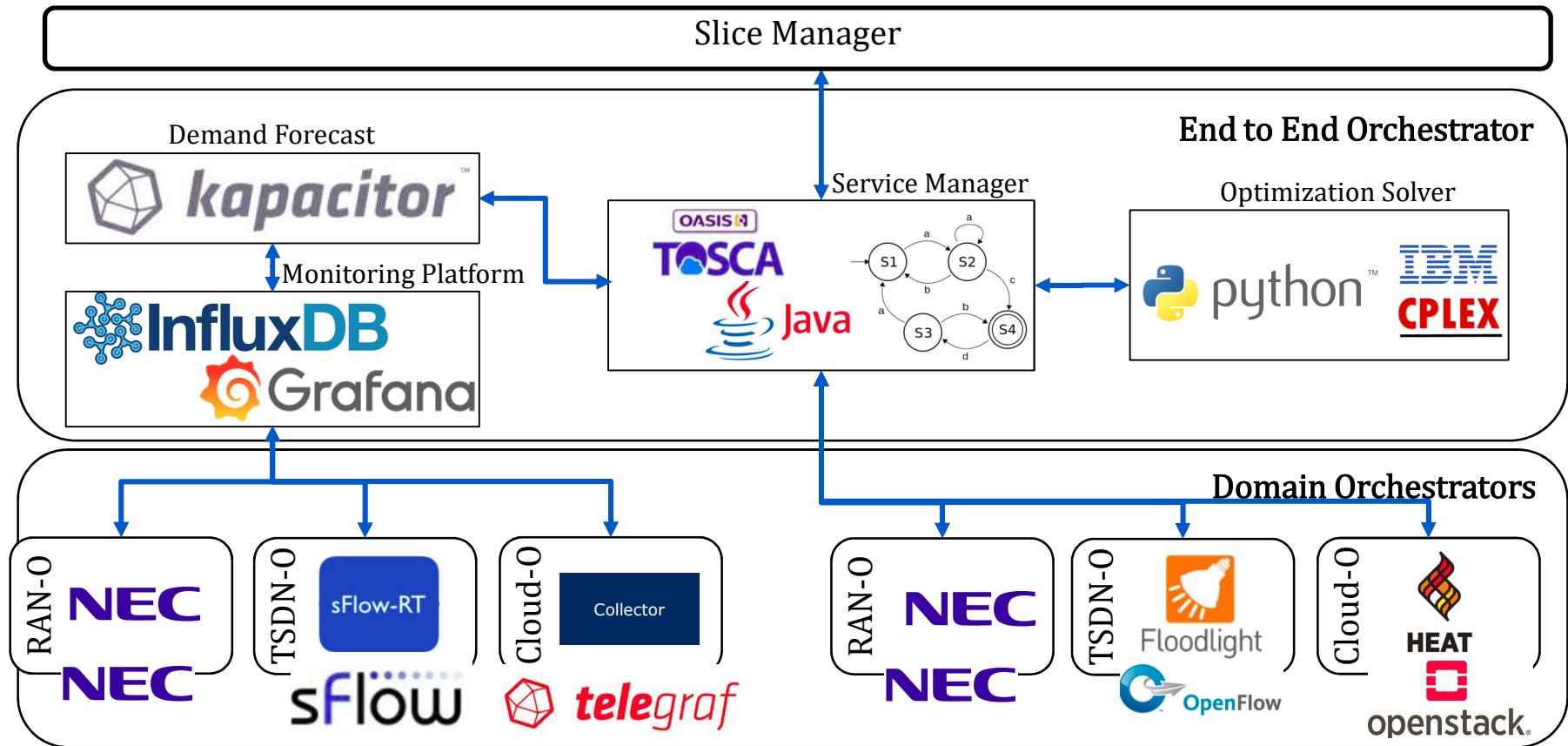
RAN · Edge Cloud Units · Transport Network · Core Cloud Units

[1] Mestres, Albert, et al. "Knowledge-defined networking." ACM SIGCOMM Computer Communication Review 47.3 (2017): 2-10.

Orchestrating a brighter world    NEC
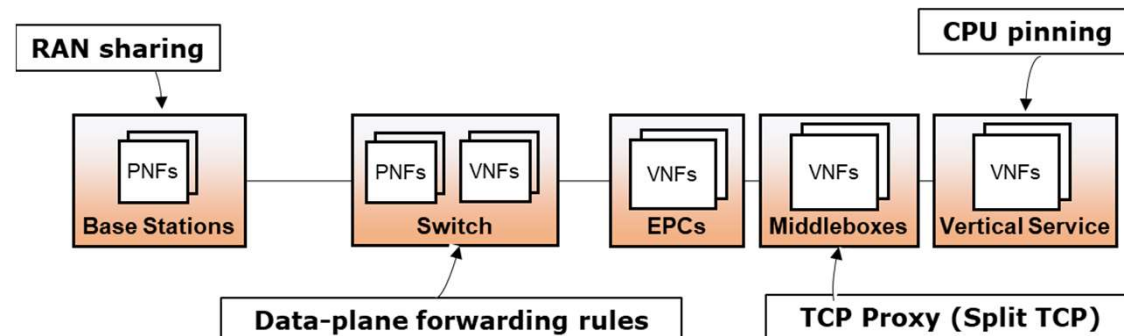
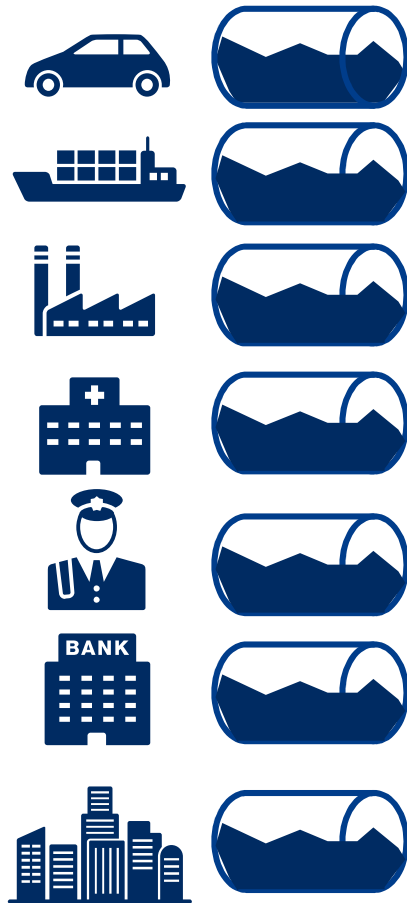**"Local" orchestrators**

1. Create data-plane abstractions
2. Monitoring data-plane usage
3. Enforce orchestration decisions

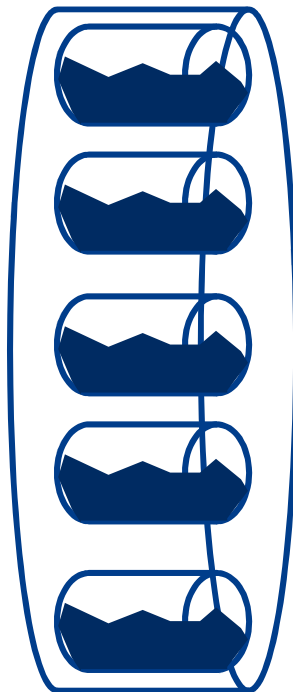Orchestrating a brighter world NEC
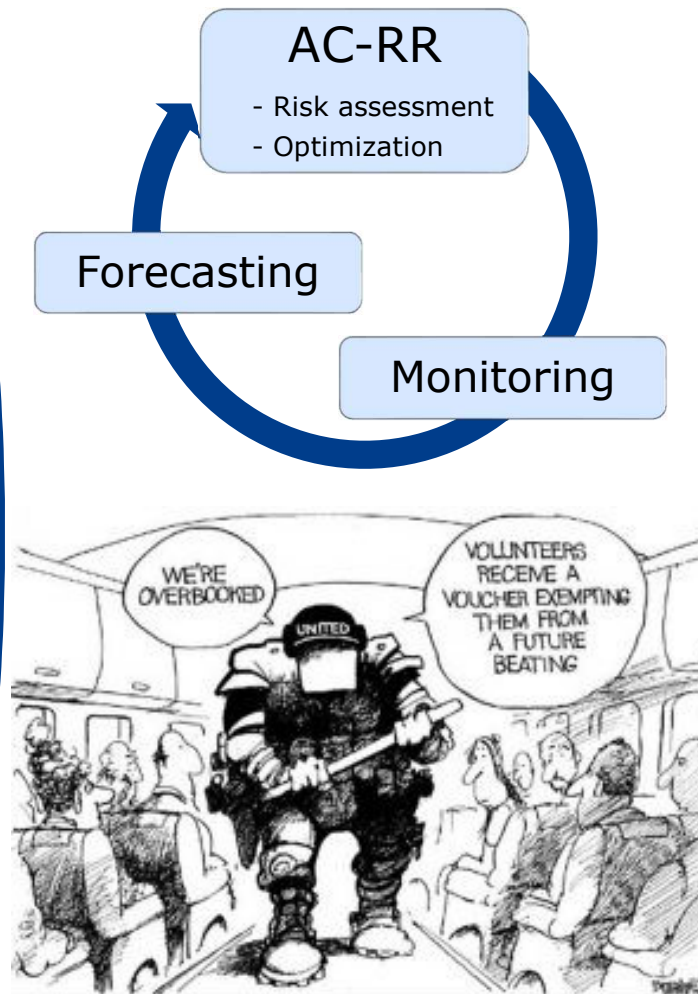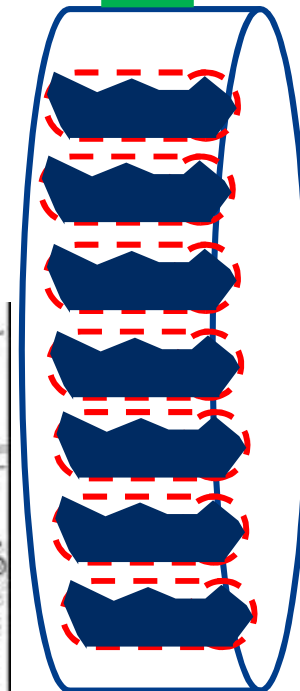
**Verticals Reqs**

**5G System**

**5G System**

AC-RR

- Risk assessment
- Optimization

Forecasting

Monitoring

Rate  Latency  CPU  Storage

\Orchestrating a brighter world  NEC

**Random processes hard to characterize**

Individual penalty of the slice

Individual reward of the slice

$$\min_{\boldsymbol{x}^{(t)} \in \{0,1\}^{\mathcal{S}}, \boldsymbol{z}^{(t)} \in \mathbb{R}_+^{\mathcal{S}}} \quad \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \sum_{\tau \in \mathcal{T}^{(t)}} \sum_{\substack{p \in \mathcal{P}_{b,c} \\ \forall b \in \mathcal{B}, c \in \mathcal{C}}} \overbrace{K_\tau x_{\tau,p}^{(t)} \Pr\left[z_{\tau,p}^{(t)} < \lambda_{\tau,p}^{(t)}\right]}^{\text{Expected penalty}} - \overbrace{R_\tau x_{\tau,p}^{(t)}}^{\text{Reward}}$$

Admission Control decisions

Resource Reservation decisions

Set of the slice requests at epoch t

Probability that the reservation is less that the usage (SLA violation)

subject to (1) capacity/delay/system constraints (linear and decoupled)

$$(2) \quad z_{\tau,p}^{(t)} \le x_{\tau,p}^{(t)} \Lambda_\tau, \qquad \forall \tau \in \mathcal{T}, \forall \mathcal{P}_{b,c}, \forall b \in \mathcal{B}, \forall c \in \mathcal{C}$$

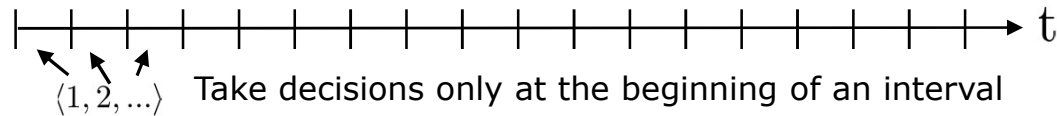**No more resources than SLA if access is granted..**
**... but nothing if access is rejected**

\Orchestrating a brighter world  NEC

## Practical simplifications

- $L_\tau$ is small compared to the system's time horizon.



$\langle 1, 2, \ldots \rangle$  Take decisions only at the beginning of an interval

- We do forecasting. And weight by a deterministic risk-cost function.

| **Risk of resource deficit** |
| :---: |
| $P_{\tau,p} := \dfrac{\Lambda_{\tau,p} - z_{\tau,p}}{\Lambda_{\tau,p} - \hat{\lambda}_{\tau,p}}, \quad 0 \le P_{\tau,p} \le 1$ |

**Risk of wrong predictions**
(forecast uncertainty amplified by slice duration)

$$\xi_{\tau,p} := \hat{\sigma}_{\tau,p} L_\tau, \quad 0 < \xi_{\tau,p} \le L_\tau$$

$$\min_{\boldsymbol{x}^{(t)} \in \{0,1\}^{\mathcal{S}}, \boldsymbol{z}^{(t)} \in \mathbb{R}_+^{\mathcal{S}}} \lim_{T} \frac{1}{T} \sum_{t=1}^{T} \sum_{\substack{\tau \in \mathcal{T}^{(t)}}} \sum_{\substack{p \in \mathcal{P}_{b,c} \\ \forall b \in \mathcal{B} \; c \in \mathcal{C}}} \overbrace{K_\tau x_{\tau,p}^{(t)} \Pr\left[z_{\tau,p}^{(t)}\right]}^{\text{Expected penalty}} - \overbrace{R_\tau x_{\tau,p}^{(t)}}^{\text{Reward}}$$

subject to   (1) Capacity/delay/system constraints (linear and decoupled)

(2) $x_{\tau,p}^{(t)} \hat{\lambda}_{\tau,p}^{(t)} \le z_{\tau,p}^{(t)} \le x_{\tau,p}^{(t)} \Lambda_\tau$, $\qquad \forall \tau \in \mathcal{T}, \forall \mathcal{P}_{b,c}, \forall b \in \mathcal{B}, \forall c \in \mathcal{C}$
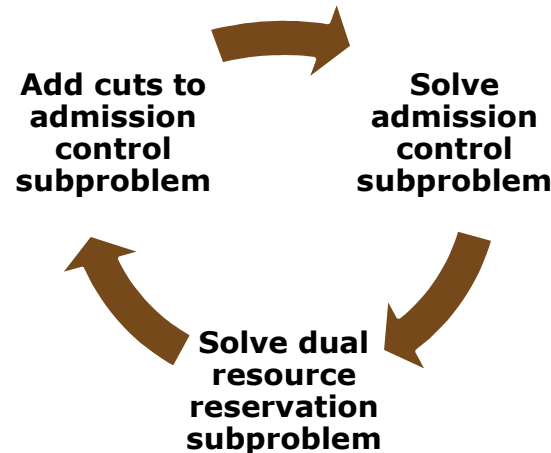
**No less resources than forecasted peak demand**

**No more resources than SLA if access is granted**

# Problem formulation

## We can linearize our problem easily

- … but still a MILP with coupled constraints (NP-hard)

**Add cuts to admission control subproblem** → **Solve admission control subproblem** → **Solve dual resource reservation subproblem** → (cycle)
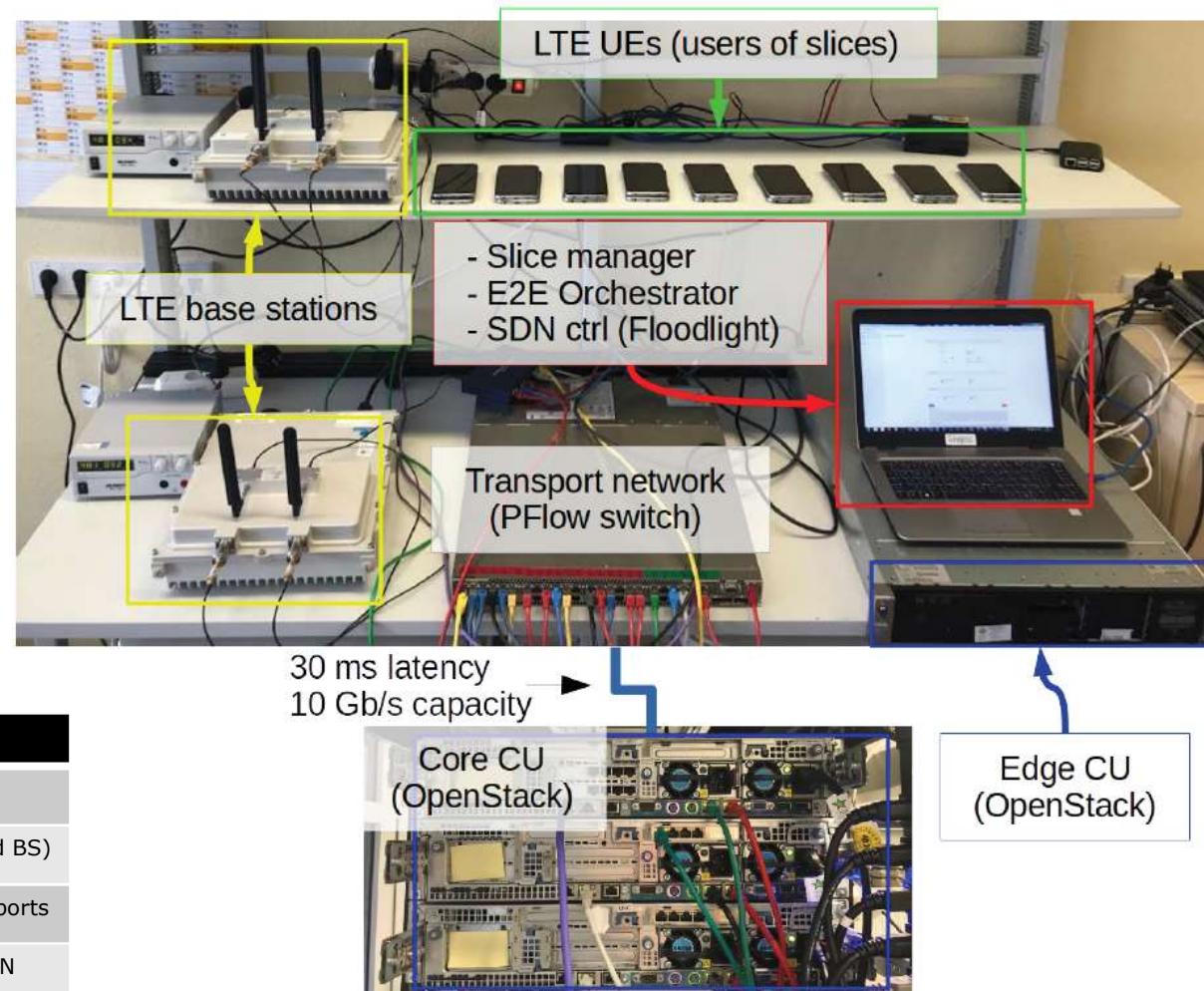
## We use two methods to solve this problem

- **Bender's decomposition**
  - Developed in 1960s, increasingly popular last years.
  - Finds an exact solution, but no guarantees for the time needed

- **Heuristic algorithm**
  1. Cast the admission control subproblem into a classical 0-1 Knapsack problem and use an heuristic to solve it.
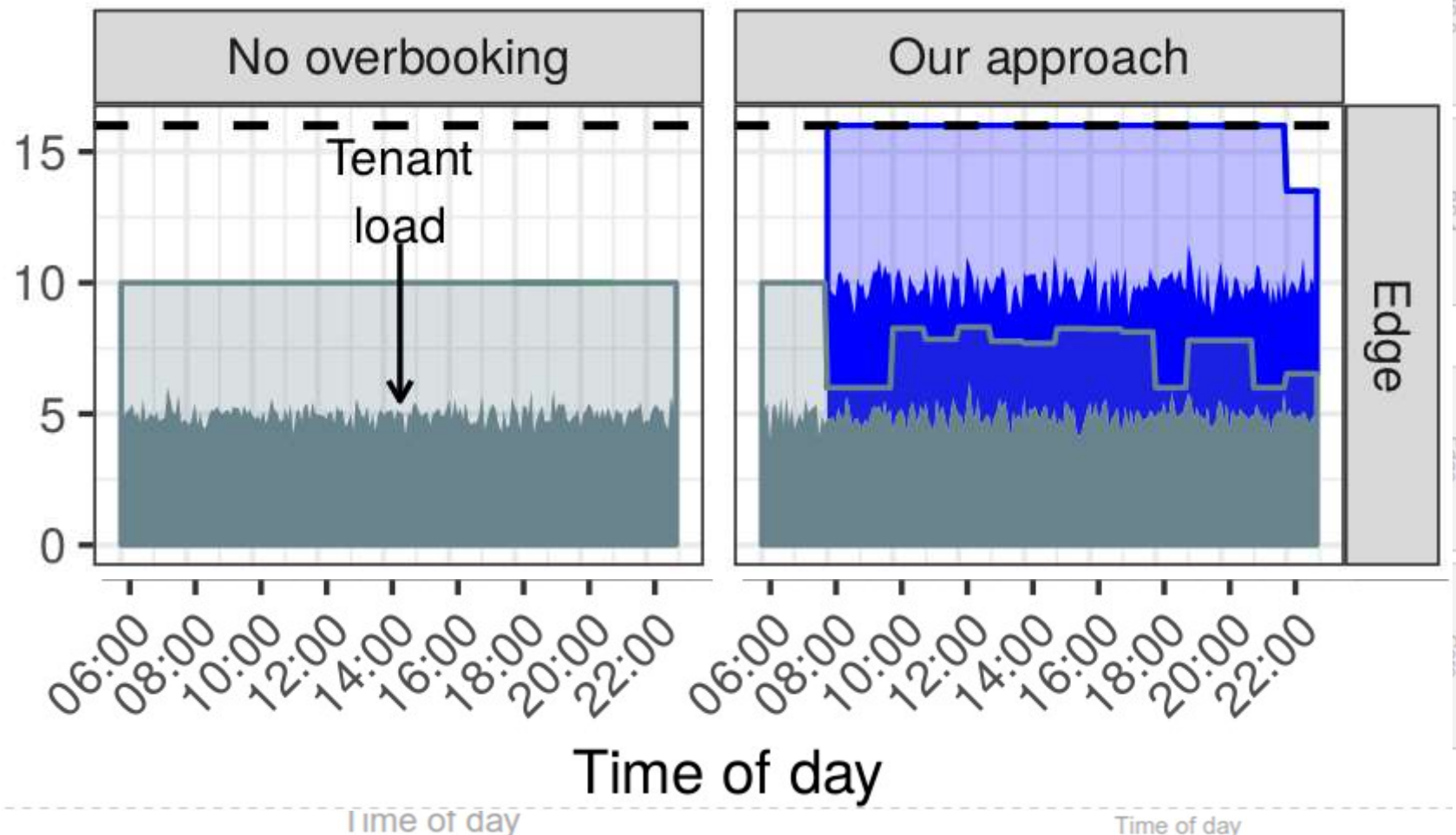  2. Add only feasibility cuts

\Orchestrating a brighter world **NEC**

| Device | Description |
|--------|-------------|
| vEPCs | OpenEPC Rel. 7 (1x per slice) |
| UEs | Samsung Galaxy 7 (1x per slice and BS) |
| BH | Openflow 1.5 switch w/ 48 1-Gb/s ports |
| RAN | 2x 20-MHz NEC small cell w/ RAN sharing |
| CUs | OpenStack Queens w/ 16 (Edge) and 64 (Core) cores |

Orchestrating a brighter world  NEC

# Evaluation: Experimental proof-of-concept

| Type | Delay req. (ms) | Service bitrate req. (Mb/s) | CPU req. (number of CPUS) | Price (monetary units) | Penalty (monetary units) |
|------|------|------|------|------|------|
| eMBB | 30 | 50 | 0 | 1 | 1/50 |
| mMTC | 30 | 10 | 2λ | 3 | 3/10 |
| uRLLC | 5 | 25 | 0.2λ | 2 | 2/25 |



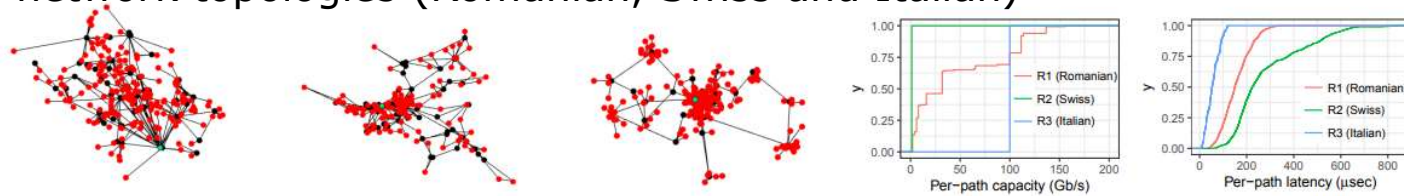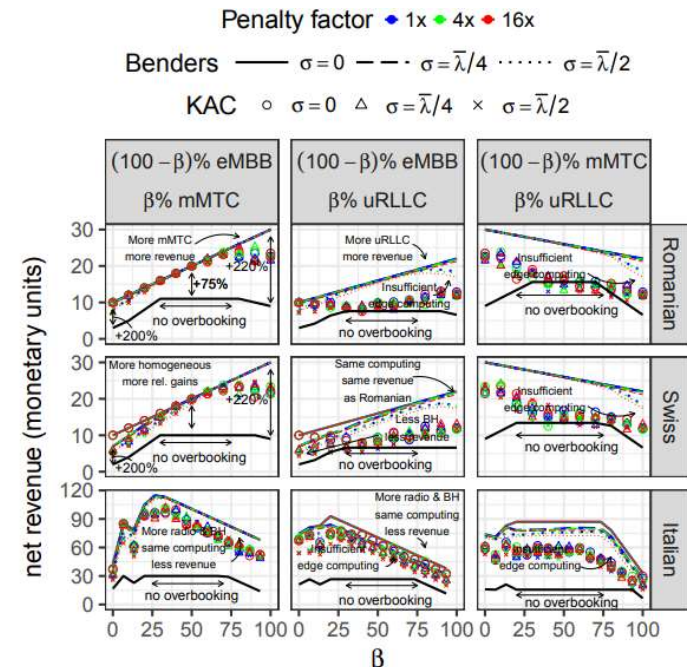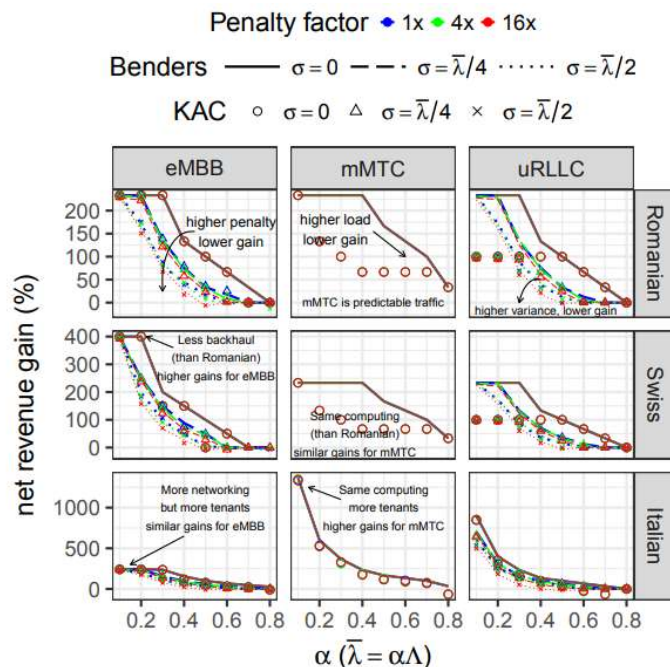© NEC Laboratories Europe GmbH 2018   \Orchestrating a brighter world **NEC**

# Overbooking brings significant gains

- 3 real network topologies (Romanian, Swiss and Italian)



(a) Romanian topology (N1).  (b) Swiss topology (N2).  (c) Italian topology (N3).  (d) Path Capacity Distribution  (e) Path Delay Distribution

- Evaluation: Wide set of penalty, traffic load and traffic variability
- Up to 200% revenue gains when load is low and is predictable



© NEC Laboratories Europe GmbH 2018     \Orchestrating a brighter world   NEC

# Conclusions

## Network Slicing will be a key technology for 5G

- New sources of revenue for mobile operators **and** vendors
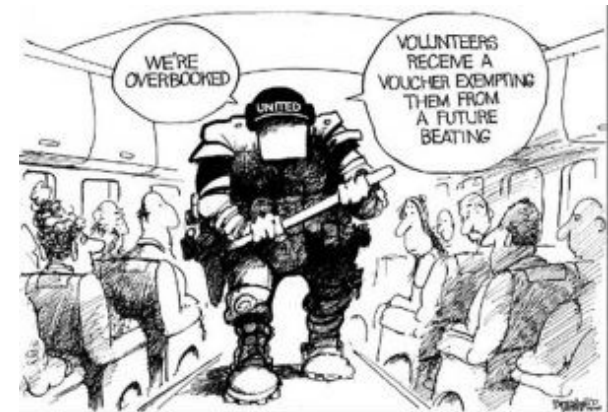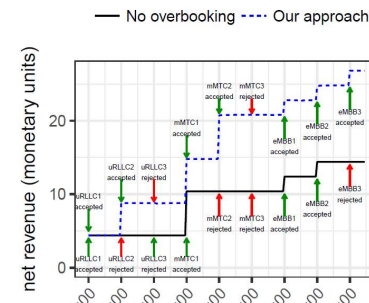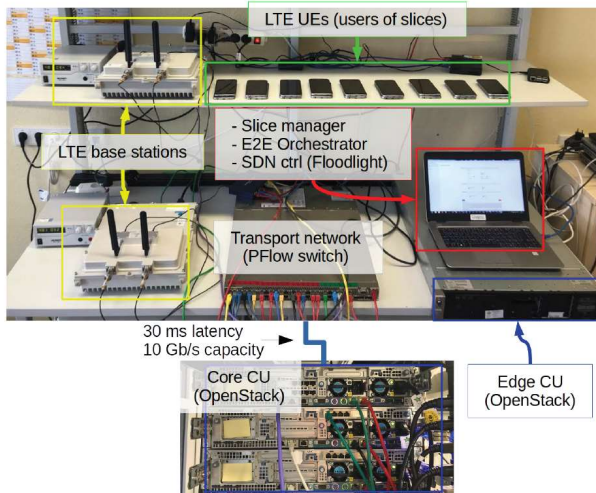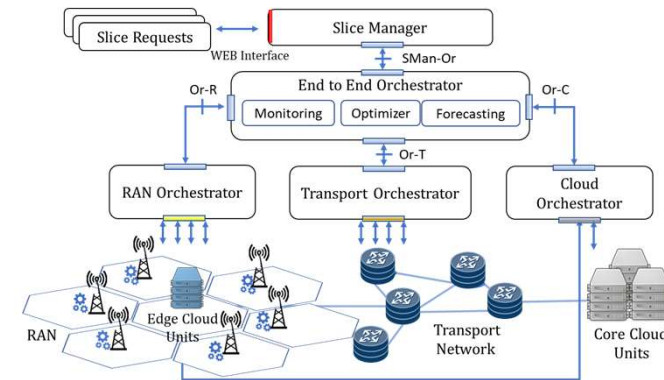
## Two main challenges:

- Service/Slice **Orchestration Platform**
  - Hierarchical orchestration for large-scale
  - Feedback loop -> monitoring/control
- Admission Control and Resource Reservation
  - We explore the concept of **slice overbooking**
  - 2 algorithms, 1 experimental PoC, 3 large-scale topologies







(a) Romanian topology (N1).  (b) Swiss topology (N2).  (c) Italian topology (N3).

\Orchestrating a brighter world  NEC

\Orchestrating a brighter world    **NEC**

# Overbooking Network Slices through Yield-driven End-to-End Orchestration

J. X. Salvat,

L. Zanzi,

Andres Garcia-Saavedra,

V. Sciancalepore,

X. Costa-Perez

**NEC Laboratories Europe**