# Cloud-to-edge continuum



| | Users | Network Operator | | | Service Providers |
|---|---|---|---|---|---|

| | ON-DEVICE | ON-PREMISE | FAR EDGE | NEAR EDGE | | CLOUD |
|---|---|---|---|---|---|---|
| | | | | | EDGE \| CLOUD | |
| Typical distance | | <1 km | 1-100 km | 100-1000 km | | >1000 km |
| Average latency* | | 1 ms | 2-5 ms | 10-20ms | | > 20 ms |
| | Millions | 100 000s | 1 000s    100s | 10s | | <10 |

Cell Site Edge

Aggregation node

Near-premise

Central Office

Mini DCs

In-country Data Center

Central Data Center

Backgroung figure from European industrial technology roadmap for the next generation cloud-edge offering, p12, The European Commission, 2021, link
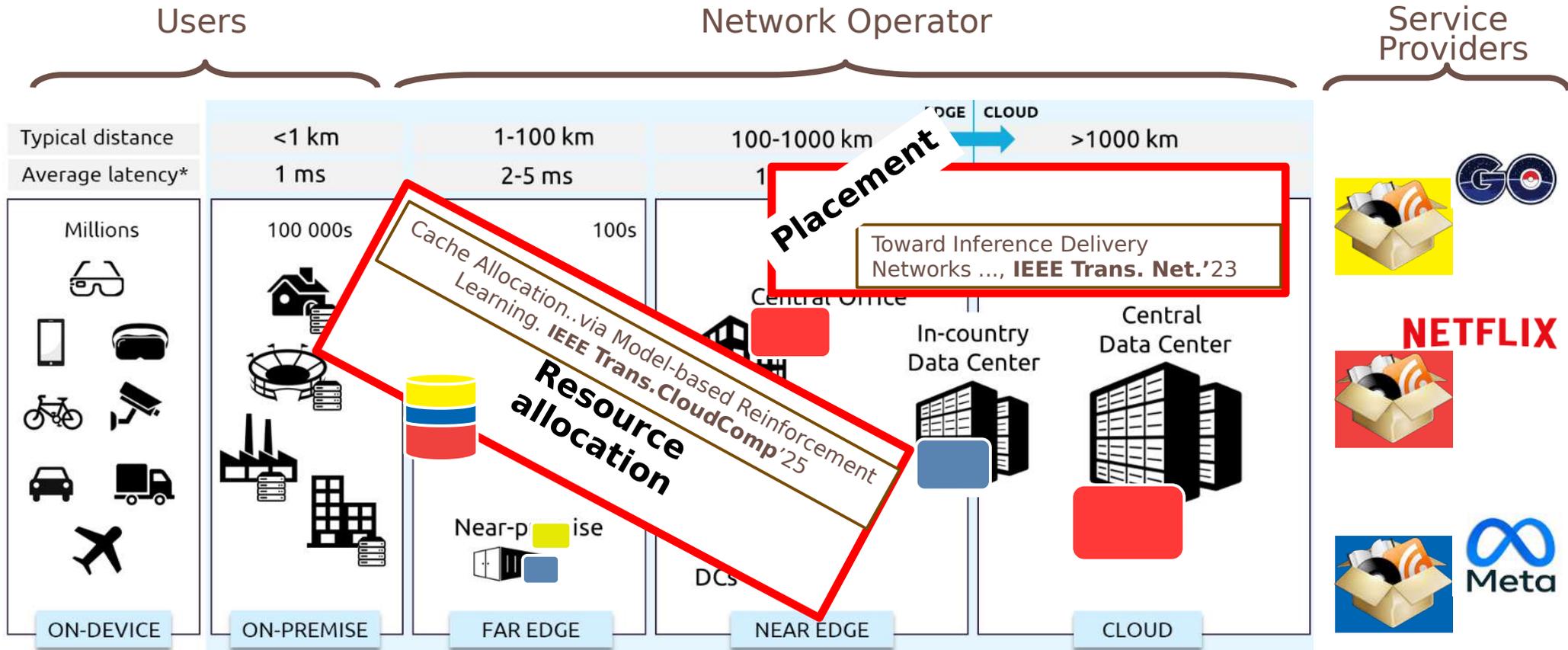
# Cloud-to-edge continuum
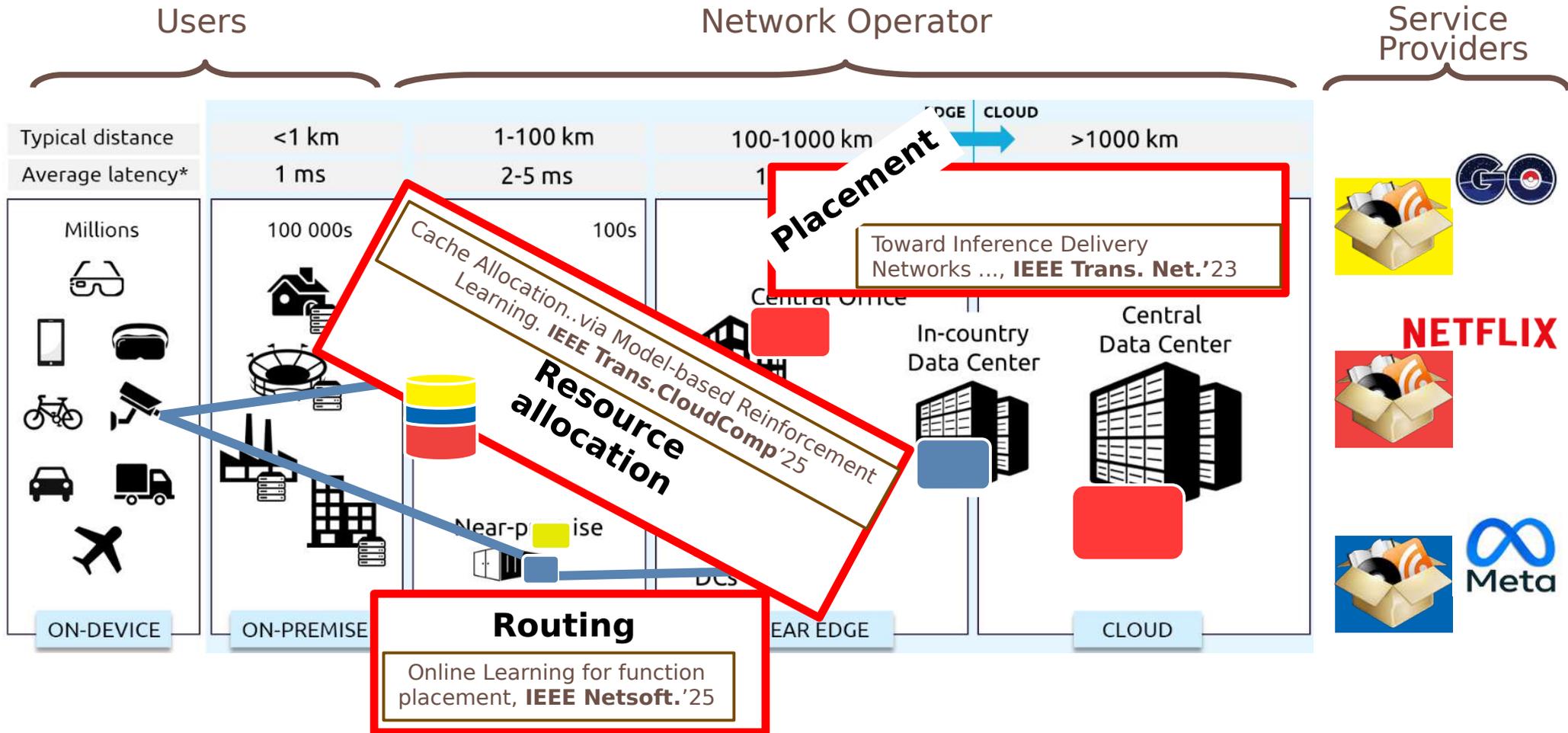


Backgroung figure from  European industrial technology roadmap for the next generation cloud-edge offering, p12, The European Commission, 2021, link

# Cloud-to-edge continuum



Users

Network Operator

Service Providers

| | | | | EDGE | CLOUD | |
| --- | --- | --- | --- | --- | --- | --- |
| Typical distance | <1 km | 1-100 km | 100-1000 km | | >1000 km | |
| Average latency* | 1 ms | 2-5 ms | | | | |

Placement

Toward Inference Delivery Networks ..., **IEEE Trans. Net.**'23

Cache Allocation..via Model-based Reinforcement Learning. **IEEE Trans.CloudComp**'25

Resource allocation

ON-DEVICE   ON-PREMISE   FAR EDGE   NEAR EDGE   CLOUD

Central Office   In-country Data Center   Central Data Center
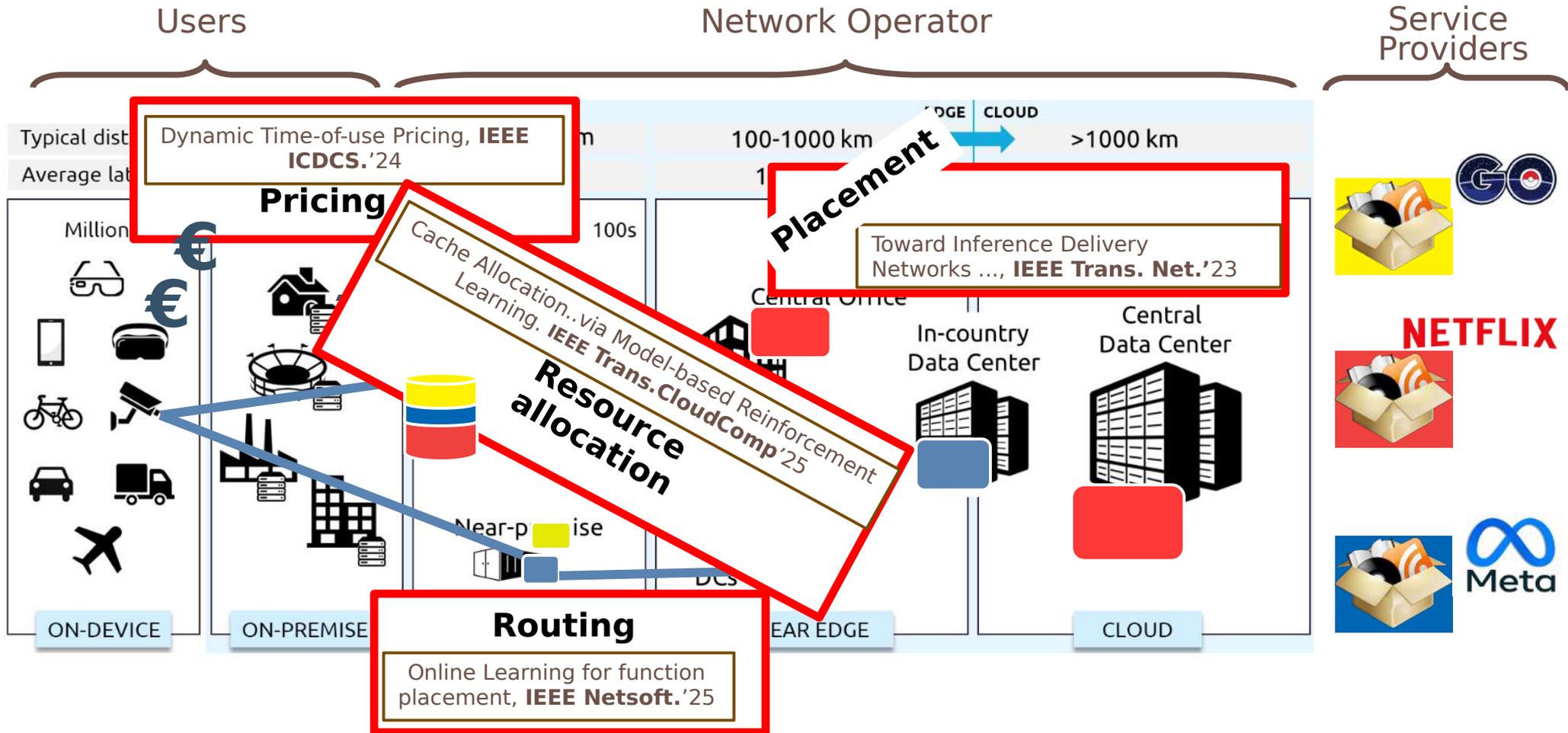
Near-premise   DCs

Millions   100 000s   100s

Background figure from  European industrial technology roadmap for the next generation cloud-edge offering, p12, The European Commission, 2021, link

# Cloud-to-edge continuum



Users

Network Operator

Service Providers

| Typical distance | <1 km | 1-100 km | 100-1000 km | >1000 km |
|---|---|---|---|---|
| Average latency* | 1 ms | 2-5 ms | | |

Millions    100 000s    100s

EDGE    CLOUD

Placement

Toward Inference Delivery Networks ..., **IEEE Trans. Net.**'23

Cache Allocation..via Model-based Reinforcement Learning. **IEEE Trans.CloudComp**'25

Resource allocation

Central Office

In-country Data Center

Central Data Center

Near-premise

DCs

**Routing**

Online Learning for function placement, **IEEE Netsoft.**'25

ON-DEVICE    ON-PREMISE    NEAR EDGE    CLOUD

NETFLIX

Meta

Backgroung figure from  European industrial technology roadmap for the next generation cloud-edge offering, p12, The European Commission, 2021, link

# Cloud-to-edge continuum

**Users**

**Network Operator**

**Service Providers**

Dynamic Time-of-use Pricing, **IEEE ICDCS.**'24

**Pricing**

**Placement**

Toward Inference Delivery Networks ..., **IEEE Trans. Net.**'23

Cache Allocation..via Model-based Reinforcement Learning. **IEEE Trans.CloudComp**'25

**Resource allocation**

**Routing**

Online Learning for function placement, **IEEE Netsoft.**'25

Backgroung figure from  European industrial technology roadmap for the next generation cloud-edge offering, p12, The European Commission, 2021, link

# Learning with no long training

## No long training is possible in a real system

- Model-based q-learning[1]
  - Calibrate a *model* to match observations
  - *Simulate* x times more transitions than the real one
  - Guarantee: convergence in probability boundedly close to the optimum
- Generalized Hidden Parameter Markov Decision Processes[2]
  - 2 Bayesian Neural Networks
    - 1 to learn state transitions
    - 1 to learn reward
  - Train on different synthetic scenarios
    - Keep some parameters fixed when changing scenarios (general dynamics)
    - Let other parameters change at each scenario (scenario-specific dynamics)
  - When applied in a real case, only few parameters need to be adapted
  - Guarantee: none

- Online learning[3]
  - Take a placement action
  - Assume worst case for the requests
  - Estimate a subgradient
  - Apply gradient ascent
  - Guarantee: sublinear regret
- Multi-armed bandits[4]
  - ~KL-UCB algorithm
  - Estimate bounds for the unknown parameters
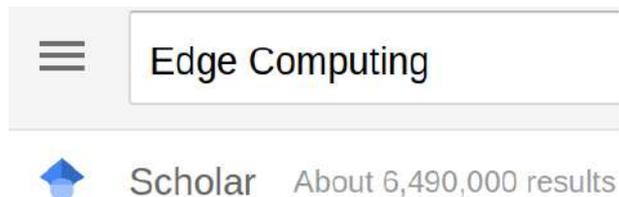  - Compute policy using bounds
  - Guarantee: sublinear regret

[1] Cache Allocation..via Model-based Reinforcement Learning. **IEEE Trans.CloudComp**'25
[2] Dynamic Time-of-use Pricing, **IEEE ICDCS.**'24
[3] Toward Inference Delivery Networks ..., **IEEE Trans. Net.**'23
[4] Online Learning for function placement, **IEEE Netsoft.**'25

# Why is not Edge Computing implemented?

Edge Computing

Scholar    About 6,490,000 results

- Unprecedented **business opportunity** for network operators

    - Network operators own the Edge

    - Service Providers must pass through network operators to run at the edge

# Why is not Edge Computing implemented?

Edge Computing

Scholar    About 6,490,000 results

- Unprecedented **business opportunity** for network operators

    - Network operators own the Edge

    - Service Providers must pass through network operators to run at the edge

But today, Edge Computing is not deployed

# Why is not Edge Computing implemented?

Edge Computing

Scholar    About 6,490,000 results

- Unprecedented **business opportunity** for network operators
  - Network operators own the Edge
  - Service Providers must pass through network operators to run at the edge

> But today, Edge Computing is not deployed

> Network operators are reluctant to bear the high cost and risk all alone

# Cooperative strategies for large technological infrastructures

- **Actors**
  - Network operators
  - Service Providers

- **Decisions**
  - Investment
    - Dimensioning resources
    - Timing of investments of each player
  - Sharing
    - of revenues, cost, risk
    - Dynamic resource allocation

- **Questions**
  - Is the coalition *stable*?
  - Is the coinvestment *profitable*?

- **Uncertainty**
  - User engagement, energy and resource availability are random processes
  - Bounds on probability of stability and profitability

- **Aim**
  - Propose multi-agent decision strategies ensuring stability or profitability with at least 99% probability

- **Coalitional (stochastic, robust) game theory**

Sakr, Araldo, Chahed, Patanè & Kofman . Coalitional game-theoretical approach to coinvestment with application to edge computing. **IEEE ICC** 2025

Sakr, Chahed, Patanè & Kofman. Co-Investment under Revenue Uncertainty Based on Stochastic Coalitional Game Theory, major revision in The Annals of Operations Research, 2026

# Sustainability

- How to reduce the environmental footprint of the computation continuum cloud→ edge→ devices?
- Moving computation close to users requires hardware

  -- Externality for its production, depletion of rare resources

  +Potential proximity to renewable energy
- Moving computation to the cloud

  -- Concentration of externalities

  -- Huge heating requirements

  +Consolidation
- Offload on available computation resources whenever possible
  - Vehicles,[1] network edge nodes, …

**RollingStone** **DARK SIDE OF AI**

Amazon has come to the state's eastern farmland, worsening a water pollution problem that's been linked to cancer and miscarriages

**The Guardian** Eur

Water levels across the Great Lakes are falling – just as US data centers move in

Region struggling with drought now threatened by energy-hungry facilities – but some residents are fighting back

**The Verge**

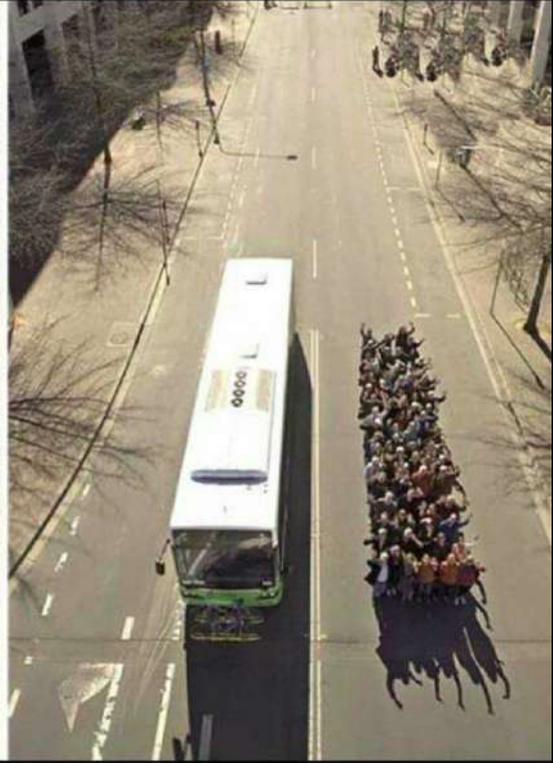/ Amazon could be accelerating the dangerous levels of nitrates in Morrow County's drinking water.

Patané, Araldo, Achir, Boukhatem, Vehicular Cloud Computing: A cost-effective alternative to Edge Computing in 5G networks, **Computer Networks**, 2025
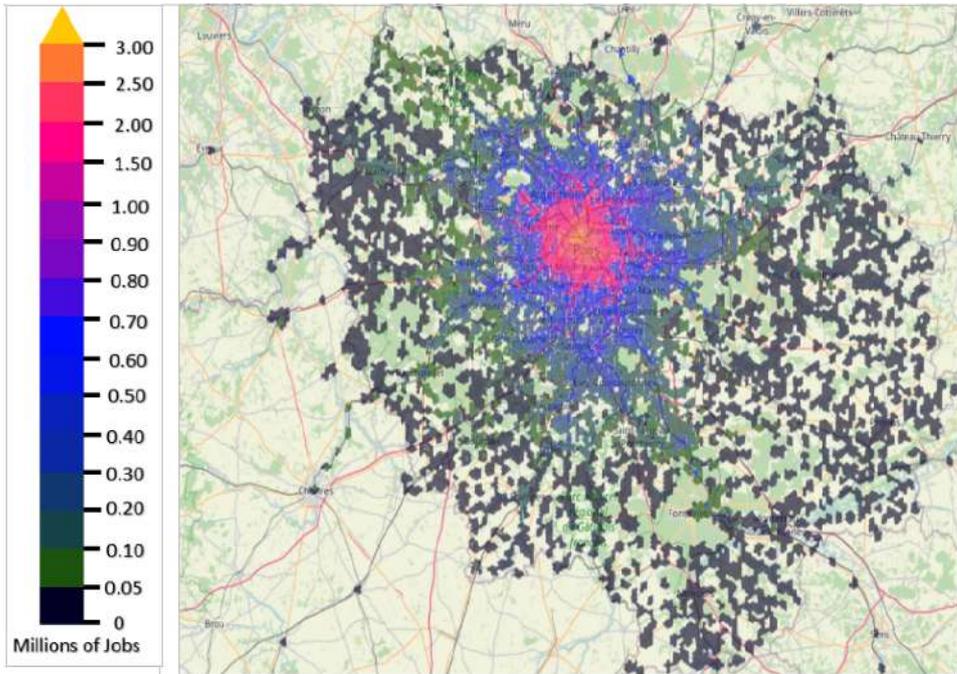
link



DIFFERENCE BETWEEN

PRIVATE VEHICLE     PUBLIC BUS

# Inequity of accessibility distribution



[1] A. Badeanlou, **A. Araldo**, M. Diana, *Assessing transportation accessibility equity via open data*, subm. to hEART'22
[2] Biazzo et al. (2019). General scores for accessibility and inequality measures in urban areas. Royal Society Open Science
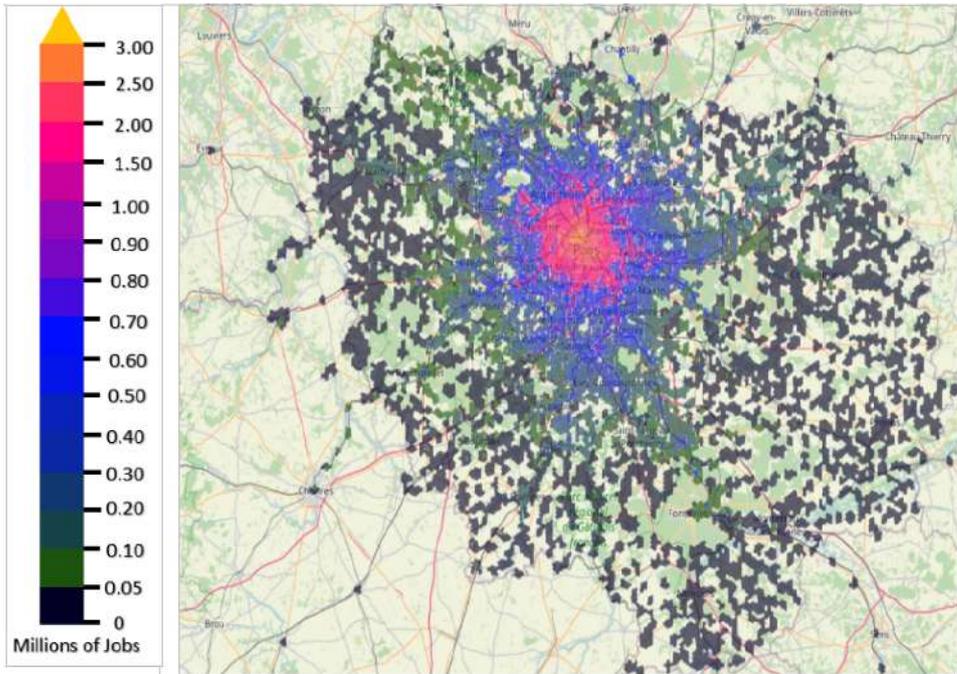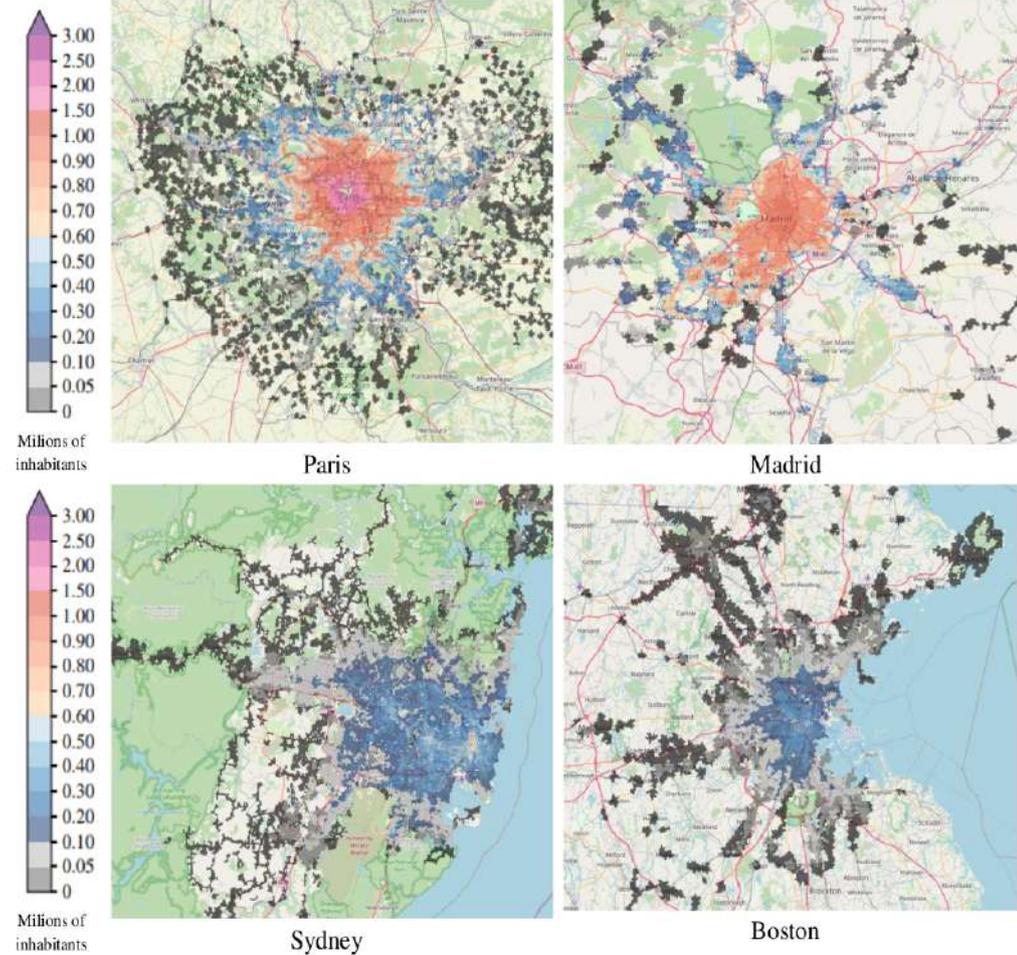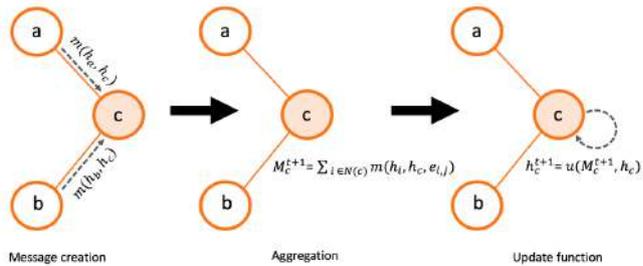
# Inequity of accessibility distribution



Figure 4-8 Sociality Score of four considered cities (Millions of Inhabitants)

Paris

Madrid

Sydney

Boston

[1] A. Badeanlou, **A. Araldo**, M. Diana, *Assessing transportation accessibility equity via open data*, subm. to hEART'22
[2] Biazzo et al. (2019). General scores for accessibility and inequality measures in urban areas. Royal Society Open Science

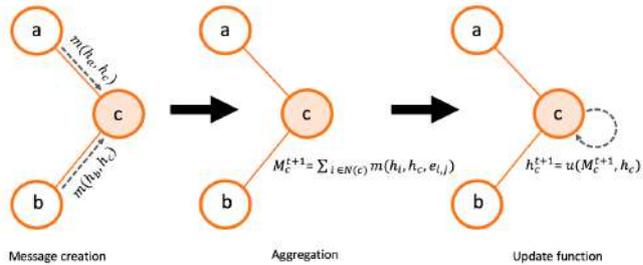# More equitable Public Transport

- Redesign bus lines for equality



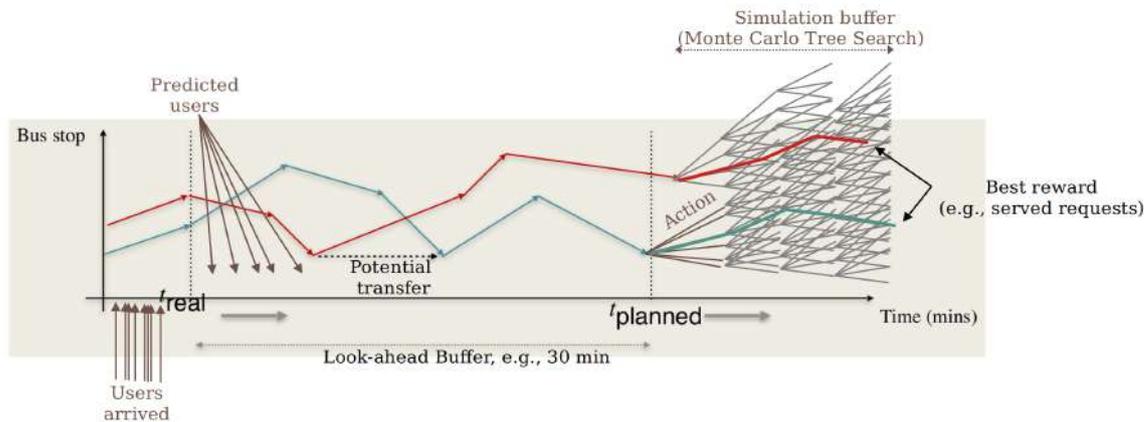Message Passing Neural Network (MPNN)

# More equitable Public Transport
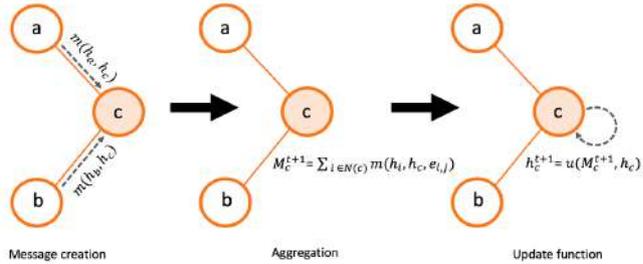
- Redesign bus lines for equality



Message Passing Neural Network (MPNN)

- Design a bus network in real time

# More equitable Public Transport

- **Redesign bus lines for equality**
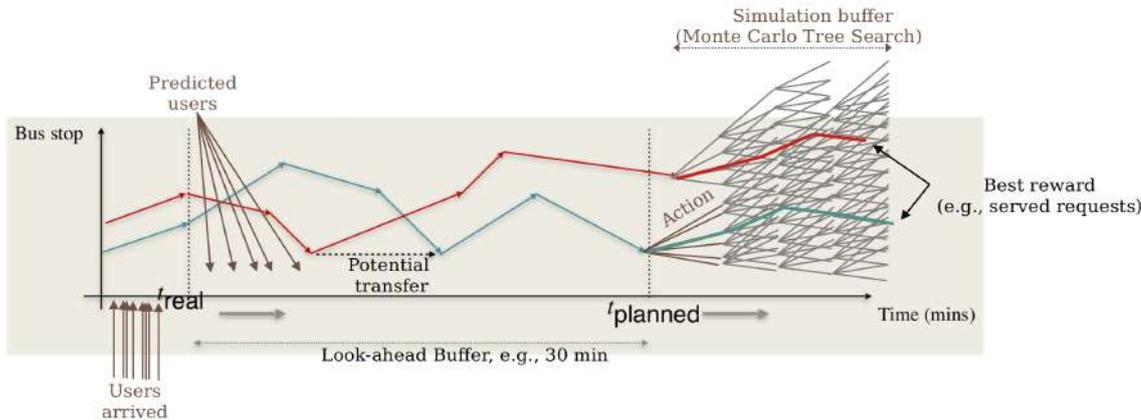


Message Passing Neural Network (MPNN)

- **Design a bus network in real time**



- **Mobility on Demand**
  - Estimate accessibility
  - Plan for equality
  - ?Stochastic geometry for strategic-level modeling?



**Number of jobs reachable within 30 min travel time**

PT only | PT+DRT



Accessibility score (number of jobs reachable)
0  500  1000  1500  2000  2500  3000  >3000

# Backup

# Toward **Inference Delivery Networks**:
## Distributing Machine Learning with Optimality Guarantees

### **IEEE Transactions on Networking** 2023

Tareq Salem, Gabriele Castellano, Giovanni Neglia, Fabio Pianese, Andrea Araldo

# Motivation

# Motivation

Research question:

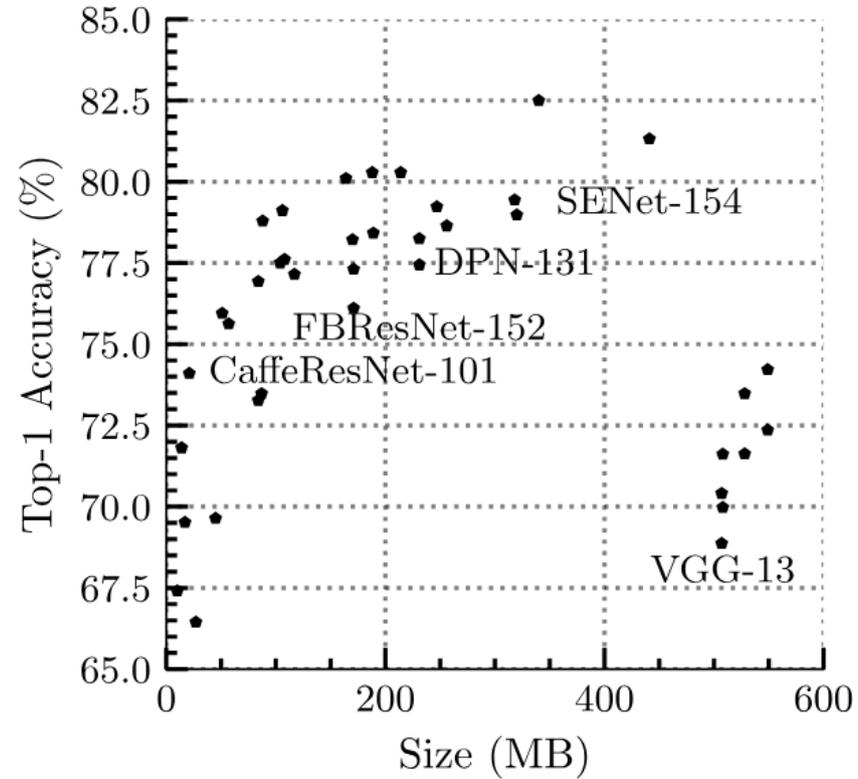In the cloud-to-edge continuum, **where** to place models and **which version**?

Novel idea:

~~Content~~ Inference Delivery Networks

# Problem: Inference Delivery Networks



| | ON-DEVICE | ON-PREMISE | FAR EDGE | NEAR EDGE | | CLOUD |
|---|---|---|---|---|---|---|
| | | | | | EDGE \| CLOUD | |
| Typical distance | <1 km | 1-100 km | | 100-1000 km | | >1000 km |
| Average latency* | 1 ms | 2-5 ms | | 10-20ms | | > 20 ms |
| | Millions | 100 000s | 1 000s    100s | 10s | | <10 |

Cell Site Edge

Aggregation node

Near-premise

Central Office

Mini DCs

In-country Data Center

Central Data Center

# Problem: Inference Delivery Networks

Typical distance

Average latency*

<1 km
1 ms

1-100 km
2-5 ms

100-1000 km
10-20ms

Millions

100 000s

1 000s

100s

10s

EDGE

Cell Site Edge

Aggregation node

Near-premise

Central Office

In-country Data Center

Central Data Center

Mini DCs

**Repository**

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

App1: Conversational agent

App2: Augmented Reality

App3: Augmented Reality (other)

# Problem: Inference Delivery Networks

**Already trained models**

| | App1: Conversational agent |
| --- | --- |
| | App2: Augmented Reality |
| | App3: Augmented Reality (other) |

Typical distance

Average latency*

<1 km

1 ms

1-100 km

2-5 ms

100-1000 km

Millions

100 000s

1 000s

100s

**Repository**

Cell Site Edge

Aggregation node

Central Office

In-country Data Center

Central Data Center

Near-premise

Mini DCs

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

# Problem: Inference Delivery Networks

**Already trained models**

App1: Conversational agent

App2: Augmented Reality

App3: Augmented Reality (other)

Typical distance

Average latency*

<1 km

1 ms

1-100 km

ms

100-1000 km

Millions

100 000s

1 000s

100s

Cell Site Edge

Aggregation node

Central Office

In-country Data Center

Central Data Center

**Repository**

Near-premise

Mini DCs

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

# Problem: Inference Delivery Networks

**Already trained models**

App1: Conversational agent
App2: Augmented Reality
App3: Augmented Reality (other)

Typical distance    <1 km    1-100 km    100-1000 km

Average latency*    1 ms    ms    100s

Millions

Cell Site Edge

Central Office

In-country Data Center

Central Data Center

**Repository**

Aggregation node

Near-premise

Mini DCs

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

# Problem: Inference Delivery Networks

**Already trained models**

App1: Conversational agent

App2: Augmented Reality

App3: Augmented Reality (other)

Typical distance
Average latency*

<1 km
1 ms

1-100 km

100-1000 km

Millions

Central Office

In-country Data Center

Central Data Center

Cell Site Edge

**Repository**

Subgradient-based online learning algorithm

minimize
cost of inaccuracy + cost of latency

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

# Results

**Theorem**
The regret (worst-case deviation from optimum) grows as √T



mAP: mean average precision (area under the precision-recall curve averaged across all classes)

Ideal baselines:

- INFIDA$_{offline}$: ideal basecase (in hindsight)

- SG: static greedy

- OLAG: online greedy

Although request arrival is not known our performance is close to that of an ideal omniscient decision-maker

# Cache Allocation in Multi-tenant Edge Computing: an Online Model-based Reinforcement Learning Approach

Ayoub Ben-Ameur, Andrea Araldo, Tijani Chahed

# Resource allocation at the Edge

Edge node

Network Operator

Service Providers

Users

# Resource allocation at the Edge

Edge node

Network Operator

Service Providers

Users

# Resource allocation at the Edge

Edge node

Network Operator

Service Providers

Users

# Resource allocation at the Edge

Edge
node

Network Operator

Service
Providers

Users

Upstream
traffic

# Resource allocation at the Edge



Edge node

Network Operator

Service Providers

Users

Upstream traffic

# Resource allocation at the Edge



Edge node

Network Operator

Service Providers

Users

Upstream traffic

Reinforcement Learning

# Resource allocation at the Edge

**Edge node**

**Network Operator**

**Service Providers**

**Users**

**Upstream traffic**
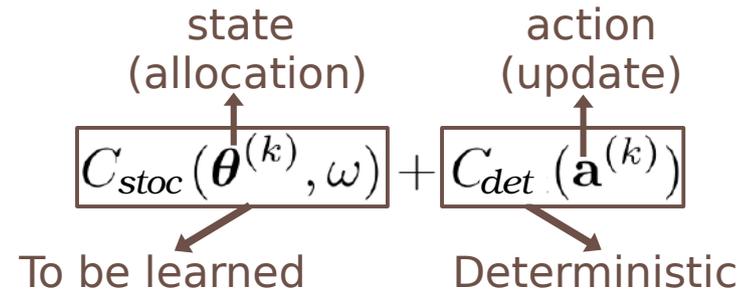
**Reinforcement Learning**

Reviewer's refrain: "Why don't you apply deep reinforcement learning (DRL)?"

- **Online learning**

  → No long training is possible

  → Experiences come from the running system

  → **Sample efficiency** (only
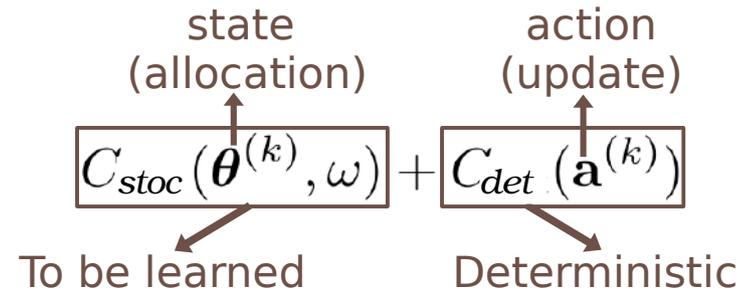
# Resolution: Model-Based Q-Learning

Instantaneous cost:

state
(allocation)

action
(update)

$$C_{stoc}\left(\boldsymbol{\theta}^{(k)},\omega\right) + C_{det}\left(\mathbf{a}^{(k)}\right)$$

To be learned          Deterministic

# Resolution: Model-Based Q-Learning

Instantaneous cost:

state
(allocation)

action
(update)

$$\boxed{C_{stoc}(\boldsymbol{\theta}^{(k)}, \omega)} + \boxed{C_{det}(\mathbf{a}^{(k)})}$$

To be learned       Deterministic

At each timeslot:

- Train a supervised model $\hat{C}_{stoc}(\cdot)$ on previous observations
- Apply Q-learning to model $\hat{C}_{stoc}(\cdot) + C_{det}$
- "Simulate" many transitions (200/sec)
    - Without perturbing the real system
    - We can extrapolate $\hat{C}_{stoc}(\boldsymbol{\theta})$ also at states $\boldsymbol{\theta}$ never visited before
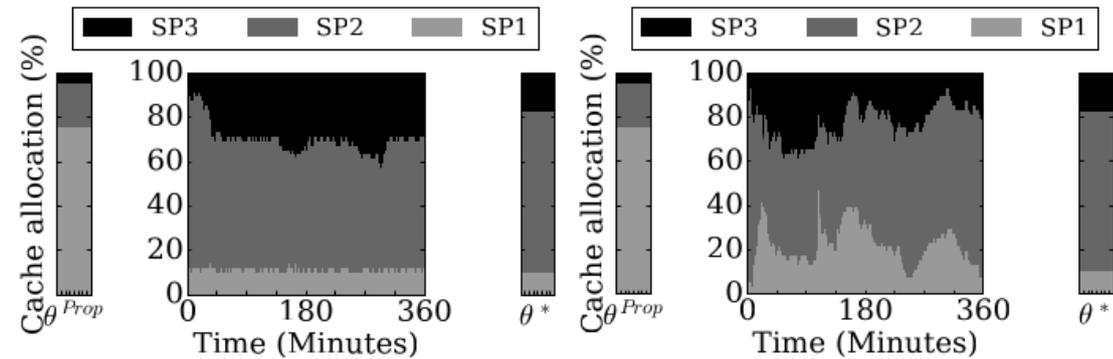
A "good" Q-table is learned very fast
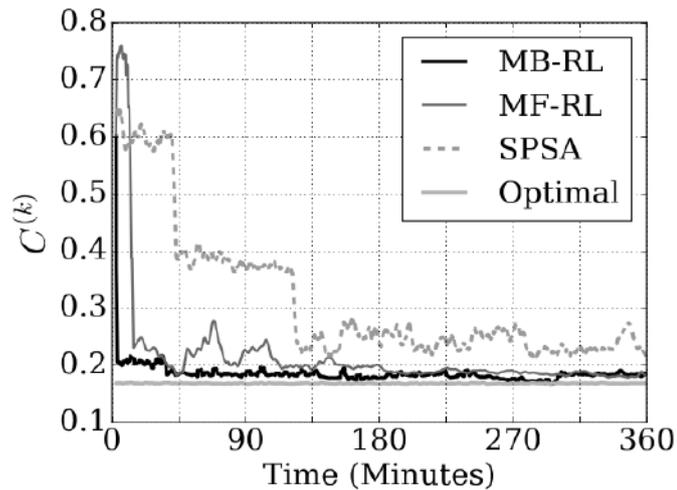
# Results

**Theorem V-B.1.** *If the discount factor $\gamma$ is sufficiently close to 1 and if the supervised model is an unbiased estimator of the nominal cost*

$$\lim_{k \to \infty} \boldsymbol{\theta}^{(k)} = \hat{\boldsymbol{\theta}}^* \text{ with probability 1.}$$



(a) MB-RL                                    (b) MF-RL

Fig. 8: Evolution of the allocation over time

# Conclusion

- AI is effective in managing systems with stochastic and unknown behavior

  – Good practical performance

  – Analytical guarantees

- **TODO**: Test these solutions beyond simulated environments

- Concern

  – "You should compare with deep learning!"

    - Hard to reproduce

    - Deep Learning can become a barrier to nice ideas

# Toward **Inference Delivery Networks**:
Distributing Machine Learning with Optimality Guarantees

Tareq Salem, Gabriele Castellano, Giovanni Neglia, Fabio Pianese, Andrea Araldo

# Decision making under uncertainty

- Uncertainty environment
  - User requests,
  - Cost resulting from certain decisions

# Decision making under uncertainty

- Uncertainty environment
  - User requests,
  - Cost resulting from certain decisions

  - How to make decisions?
    - If we know probability distributions
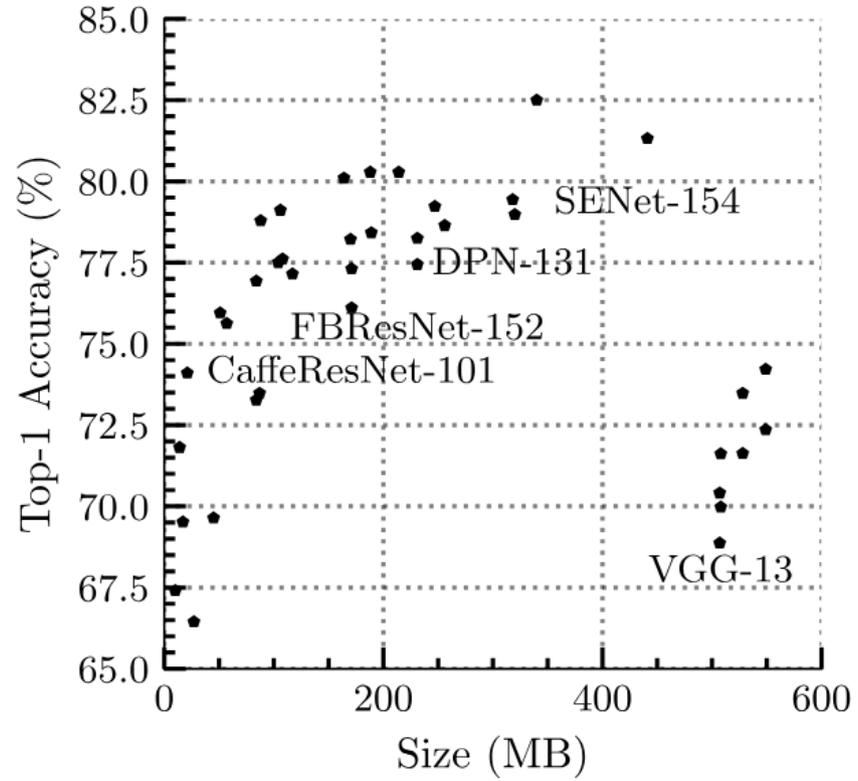      - Stochastic / Robust optimization

- Uncertainty environment
  - User requests,
  - Cost resulting from certain decisions

- How to make decisions?
  - If we know probability distributions
    - Stochastic / Robust optimization
  - Otherwise
    - Data-driven approaches,
      e.g. AI

# Definition of online learning

Personal definition inspired in particular from
N. Cesa Bianchi et al.,
Online Learning with Switching Cost and other Adaptive Adversaries
Neurips 2013

- Sequential decision making setting

- At each round:
  - The environment associates costs to each action
  - The learner selects an action
  - The learner observes the cost of such an action

- The learner aims to minimize the cost
  - The learner implements a strategy
  - Usually, the environment associates the cost so as to minimize the cost under the learner's strategy

# Motivation

# Motivation



Big model

Distillation

Smaller versions

# Motivation



Research question:

In the cloud-to-edge continuum, **where** to place models and **which version**?

Big model

Distillation

Smaller versions

Novel idea:

~~Content~~ Inference Delivery Networks

# Problem: Inference Delivery Networks

| | EDGE | CLOUD | |
|---|---|---|---|
| Typical distance | <1 km | 1-100 km | 100-1000 km | >1000 km |
| Average latency* | 1 ms | 2-5 ms | 10-20ms | > 20 ms |

Millions — 100 000s — 1 000s — 100s — 10s — <10

Cell Site Edge
Aggregation node
Near-premise
Central Office
Mini DCs
In-country Data Center
Central Data Center

ON-DEVICE — ON-PREMISE — FAR EDGE — NEAR EDGE — CLOUD

| Typical distance | <1 km | 1-100 km | 100-1000 km | >1000 km |
| Average latency* | 1 ms | 2-5 ms | 10-20ms | > 20 ms |
| | Millions | 100 000s | 1 000s | 100s | 10s | <10 |

EDGE | CLOUD

Cell Site Edge

Aggregation node

Near-premise

Central Office

In-country Data Center

Central Data Center

Mini DCs

Repository

ON-DEVICE | ON-PREMISE | FAR EDGE | NEAR EDGE | CLOUD

# Problem: Inference Delivery Networks

EDGE

| Typical distance | <1 km | 1-100 km | 100-1000 km | 100-1000 km |
| --- | --- | --- | --- | --- |
| Average latency* | 1 ms | 2-5 ms | 10-20ms | |
| Millions | 100 000s | 1 000s | 10s | |

Cell Site Edge

Aggregation node

Near-premise

Central Office

In-country Data Center

Central Data Center

Mini DCs

**Repository**

ON-DEVICE  ON-PREMISE  FAR EDGE  NEAR EDGE  CLOUD

Task1: Conversational agent

Task2: Augmented Reality

Task3: Augmented Reality (other)

**Already trained models**

Task1: Conversational agent

Task2: Augmented Reality

Task3: Augmented Reality (other)

Typical distance

Average latency*

<1 km

1 ms

1-100 km

2-5 ms

100-1000 km

Millions

100 000s

1 000s

100s

Cell Site Edge

Aggregation node

Central Office

In-country Data Center

Central Data Center

Near-premise

Mini DCs

**Repository**

ON-DEVICE

ON-PREMISE

FAR EDGE

NEAR EDGE

CLOUD

**Already trained models**

**Legend:**
- Task1: Conversational agent
- Task2: Augmented Reality
- Task3: Augmented Reality (other)

| | Typical distance | Average latency* | | | |
|---|---|---|---|---|---|
| | <1 km | 1 ms | | | |
| | 1-100 km | 2-5 ms | | | |
| | 100 000s | | | 100-1000 km | |

Millions

Cell Site Edge

Aggregation node

Central Office

In-country Data Center

Central Data Center

**Repository**

Near-premise

Mini DCs

ON-DEVICE    ON-PREMISE    FAR EDGE    NEAR EDGE    CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

# Problem: Inference Delivery Networks



**Already trained models**

| | |
|---|---|
| 🟨 | Task1: Conversational agent |
| 🟥 | Task2: Augmented Reality |
| 🟦 | Task3: Augmented Reality (other) |

Typical distance | <1 km | 1-100 km | 100-1000 km
Average latency* | 1 ms | ms | 100s

Millions | 100 000s | 1 000s

Cell Site Edge
Aggregation node
Near-premise
Central Office
In-country Data Center
Central Data Center
Mini DCs

**Repository**

ON-DEVICE | ON-PREMISE | FAR EDGE | NEAR EDGE | CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

**Already trained models**

Task1: Conversational agent

Task2: Augmented Reality

Task3: Augmented Reality (other)

Typical distance | <1 km | 1-100 km | 100-1000 km
Average latency* | 1 ms | ms | 

Millions

Cell Site Edge

Central Office

In-country Data Center

Central Data Center

Aggregation node

Near-premise

Mini DCs

**Repository**

ON-DEVICE | ON-PREMISE | FAR EDGE | NEAR EDGE | CLOUD

lower delay
lower capacity
smaller model
worse precision

higher delay
higher capacity
larger model
better precision

a request for task $i$ going through path $p$, when the request is served via model $m \in \mathcal{M}_i$ on node $p_j \in p$ :
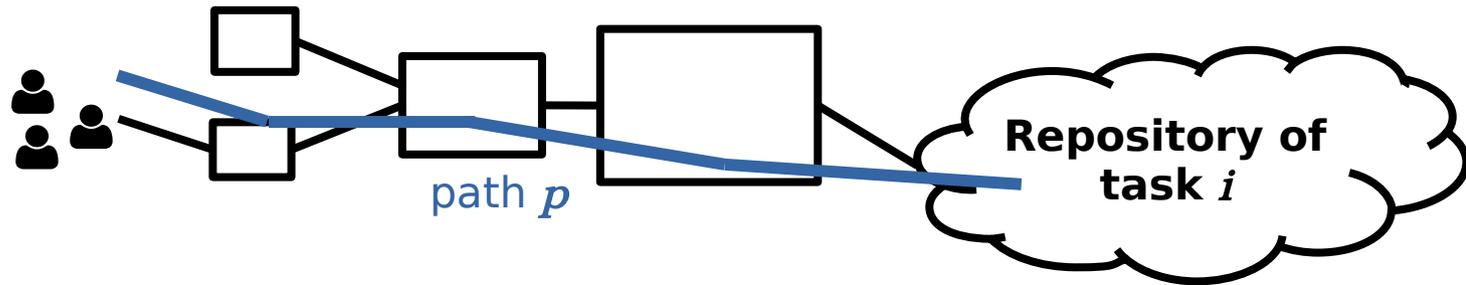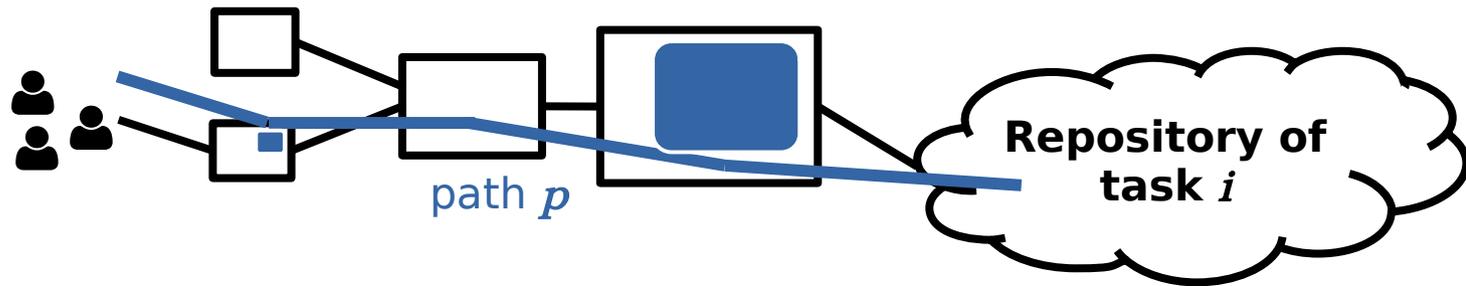


path $p$

**Repository of task $i$**

# Model

a request for task $i$ going through path $p$, when the request is served via model $m \in \mathcal{M}_i$ on node $p_j \in p$ :



path $p$
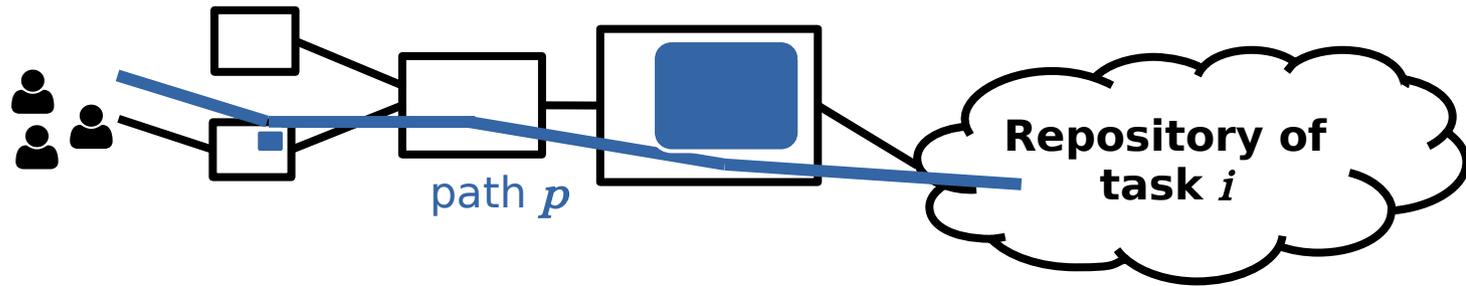
Repository of task $i$

# Model

- Cost suffered by a request for task $i$ going through path $p$, when the request is served via model $m \in \mathcal{M}_i$ on node $p_j \in p$ :

$$C_{\boldsymbol{p},m}^{p_j} = \sum_{j'=1}^{j-1} \underbrace{w_{p_{j'},p_{j'+1}}}_{\text{Round trip time}} + \underbrace{d_m^{p_j}}_{\text{Elaboration time}} + \alpha(1-\underbrace{a_m}_{\text{Accuracy}})$$



path $p$

**Repository of task $i$**

- Cost suffered by a request for task $i$ going through path $p$, when the request is served via model $m \in \mathcal{M}_i$ on node $p_j \in p$ :
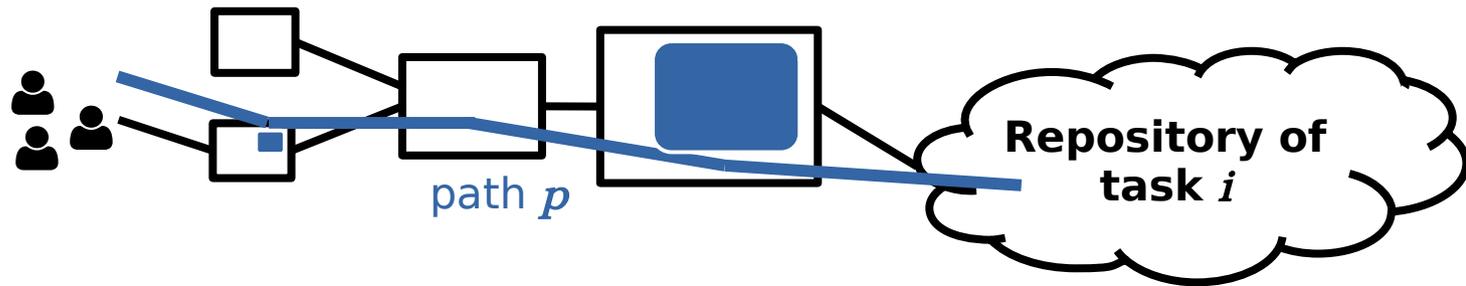
$$C_{p,m}^{p_j} = \underbrace{\sum_{j'=1}^{j-1} w_{p_{j'},p_{j'+1}}}_{\text{Round trip time}} + \underbrace{d_m^{p_j}}_{\text{Elaboration time}} + \alpha \underbrace{(1-a_m)}_{\text{Accuracy}}$$



path $p$

Repository of task $i$

- Input: requests arriving at timeslot $t$

$$r_t = [r_\rho^t]_{\rho \in \mathcal{R}} \qquad \forall \rho = (i,p),\ r_\rho^t \text{ is the number of requests for model } i \text{ following path } p$$

- Cost suffered by a request for task $i$ going through path $p$, when the request is served via model $m \in \mathcal{M}_i$ on node $p_j \in p$ :

$$C_{\boldsymbol{p},m}^{p_j} = \underbrace{\sum_{j'=1}^{j-1} w_{p_{j'},p_{j'+1}}}_{\text{Round trip time}} + \underbrace{d_m^{p_j}}_{\text{Elaboration time}} + \alpha\underbrace{(1-a_m)}_{\text{Accuracy}}$$
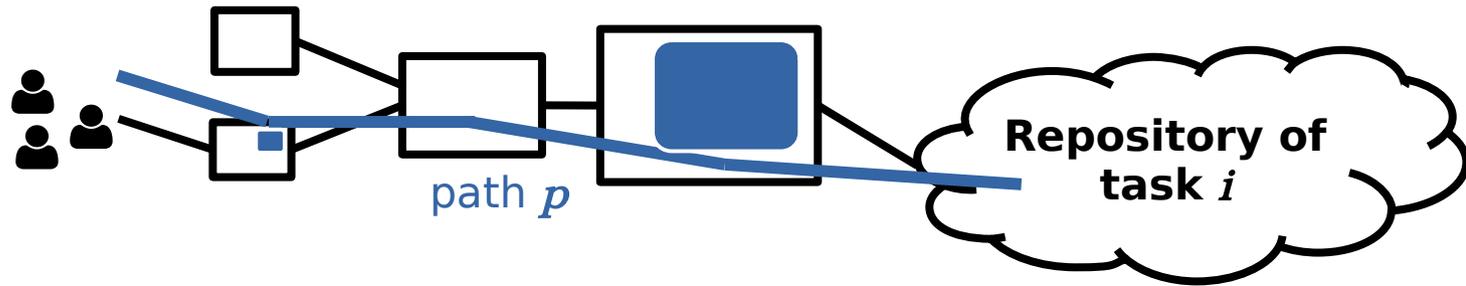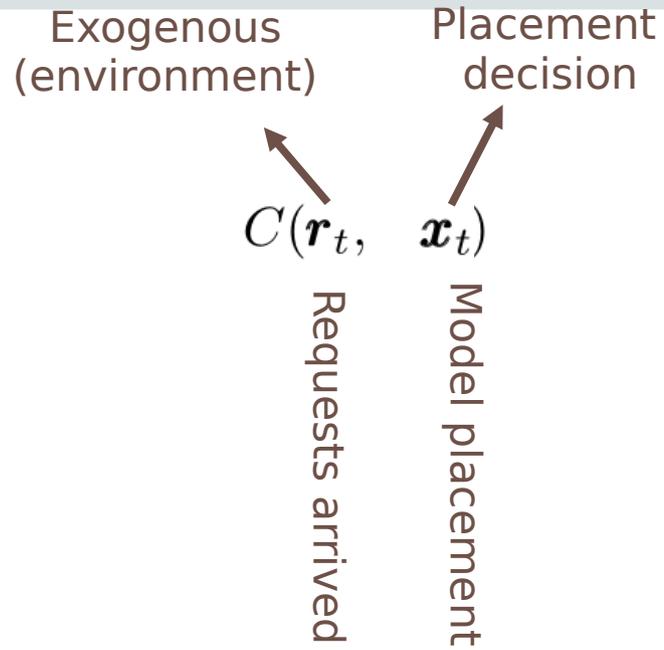


path $p$

Repository of task $i$

- Input: requests arriving at timeslot $t$

$$\boldsymbol{r}_t = \left[r_\rho^t\right]_{\rho \in \mathcal{R}} \qquad \forall \rho=(i,\boldsymbol{p}),\ r_\rho^t \text{ is the number of requests for model } i \text{ following path } \boldsymbol{p}$$
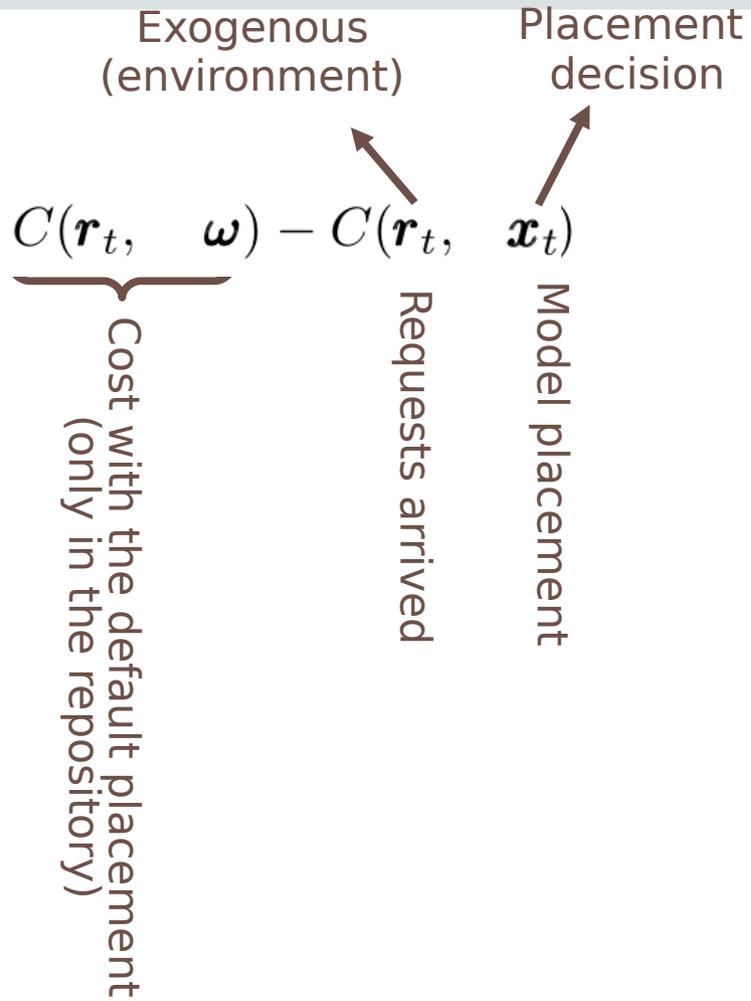
- Decision at timeslot $t$

$$\mathbf{x}_t = \left[x_{t,m}^v\right]_{t,v,m} \qquad\qquad x_{t,m}^v=1 \text{ iff model } m \in \mathcal{M}_i \text{ is placed on node } v.$$

# Online learning formulation

Exogenous
(environment)

Placement
decision

$$C(\boldsymbol{r}_t, \quad \boldsymbol{x}_t)$$

Requests arrived

Model placement

# Online learning formulation

Exogenous
(environment)

Placement
decision

$$C(\boldsymbol{r}_t, \quad \boldsymbol{\omega}) - C(\boldsymbol{r}_t, \quad \boldsymbol{x}_t)$$

Cost with the default placement
(only in the repository)

Requests arrived

Model placement

# Online learning formulation

Exogenous (environment)

Placement decision

$$C(\boldsymbol{r}_t, \quad \boldsymbol{\omega}) - C(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) = G(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) : \text{Instantaneous gain}$$

Cost with the default placement (only in the repository)

Requests arrived

Model placement

# Online learning formulation

Exogenous
(environment)

Placement
decision

$$C(\boldsymbol{r}_t, \quad \boldsymbol{\omega}) - C(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) = G(\boldsymbol{r}_t, \quad \boldsymbol{x}_t)$$ : Instantaneous gain

Cost with the default placement
(only in the repository)

Requests arrived

Model placement

$$\psi \in (0, 1]$$

What we lose compared to the best
placement in hindsight

$$\psi\text{-Regret}_{T,\mathcal{X}}$$

$$\triangleq \sup_{\{\boldsymbol{r}_t \quad \}_{t=1}^{T} \in \mathcal{A}^T} \left\{ \psi \sum_{t=1}^{T} G(\boldsymbol{r}_t, \quad \boldsymbol{x}_*) - \mathbb{E}\left[ \sum_{t=1}^{T} G(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) \right] \right\}$$

Best placement in hindsight

$$\boldsymbol{x}_* \in \arg\max_{\boldsymbol{x} \in \mathcal{X}} \sum_{t=1}^{T} G(\boldsymbol{r}_t, \quad \boldsymbol{x})$$

# Online learning formulation

Exogenous (environment)

Placement decision

$$C(\boldsymbol{r}_t, \quad \boldsymbol{\omega}) - C(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) = G(\boldsymbol{r}_t, \quad \boldsymbol{x}_t)$$ : Instantaneous gain

Cost with the default placement (only in the repository)

Requests arrived

Model placement

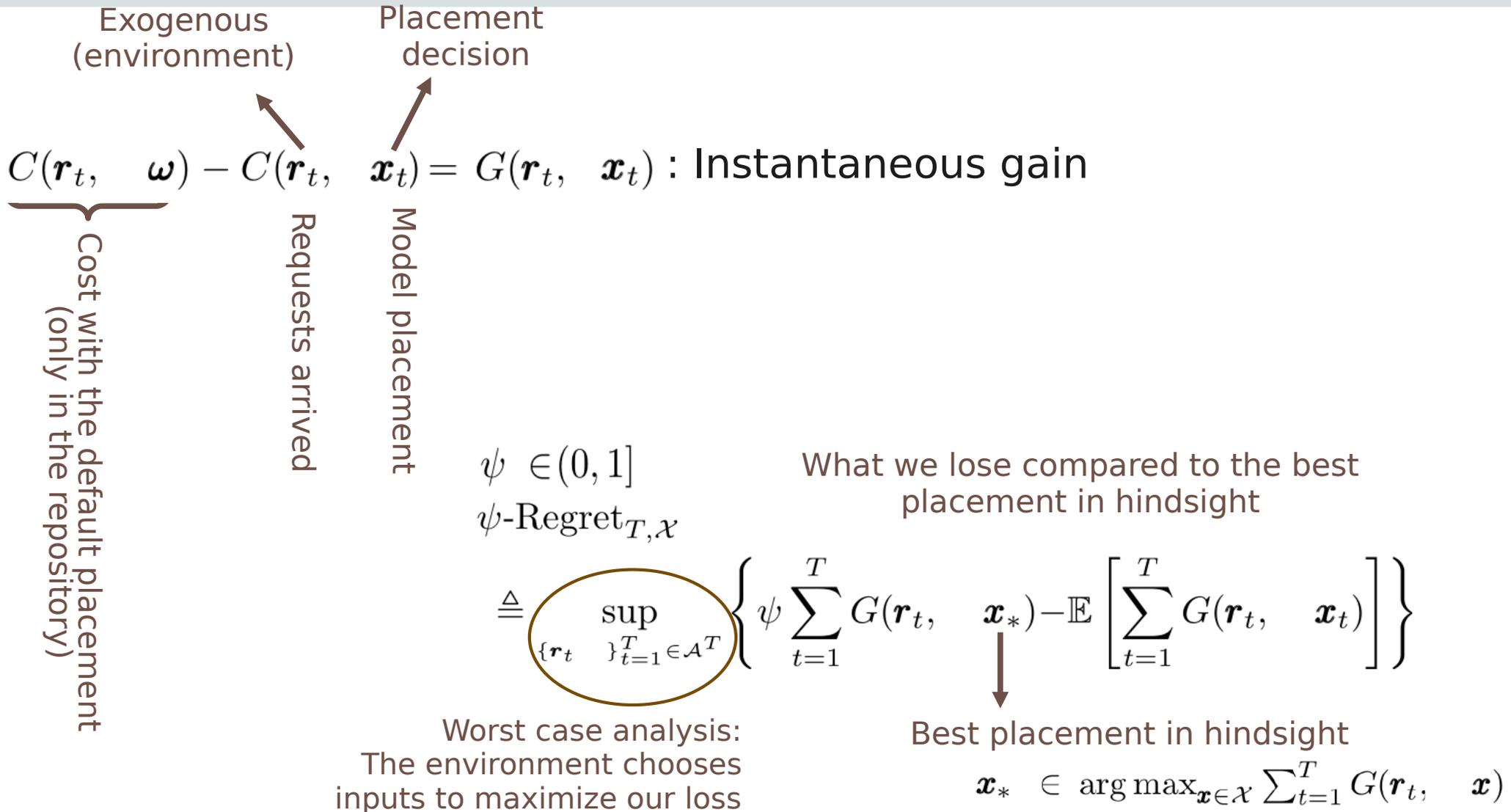$$\psi \in (0, 1]$$

$$\psi\text{-Regret}_{T,\mathcal{X}}$$

What we lose compared to the best placement in hindsight

$$\triangleq \sup_{\{\boldsymbol{r}_t \quad \}_{t=1}^T \in \mathcal{A}^T} \left\{ \psi \sum_{t=1}^T G(\boldsymbol{r}_t, \quad \boldsymbol{x}_*) - \mathbb{E}\left[ \sum_{t=1}^T G(\boldsymbol{r}_t, \quad \boldsymbol{x}_t) \right] \right\}$$

Worst case analysis:
The environment chooses
inputs to maximize our loss

Best placement in hindsight

$$\boldsymbol{x}_* \in \arg\max_{\boldsymbol{x} \in \mathcal{X}} \sum_{t=1}^T G(\boldsymbol{r}_t, \quad \boldsymbol{x})$$

**Algorithm 1** INFIDA Distributed Allocation on Node $v$

1: **procedure** INFIDA($\boldsymbol{y}_1^v = \underset{\boldsymbol{y}^v \in \mathcal{Y}^v \cap \mathcal{D}^v}{\arg\min} \Phi^v(\boldsymbol{y}^v)$, $\boldsymbol{x}_1^v = \text{DEPROUND}(\boldsymbol{y}_1^v)$,

$\eta \in \mathbb{R}_+$)

2:      **for** $t = 1, 2, \ldots, T$ **do**

3:          Compute $\boldsymbol{g}_t^v \in \partial_{\boldsymbol{y}^v} G(\boldsymbol{r}_t, \boldsymbol{l}_t, \boldsymbol{y}_t)$ through (18).

4:

5:          $\hat{\boldsymbol{h}}_{t+1}^v \leftarrow \hat{\boldsymbol{y}}_t^v + \eta \boldsymbol{g}_t^v$      ▷ Take gradient step

6:

7:          $\boldsymbol{y}_{t+1}^v \leftarrow \mathcal{P}_{\mathcal{Y}^v \cap \mathcal{D}^v}^{\Phi^v}(\boldsymbol{h}_{t+1}^v)$   ▷ Project new state onto the feasible region using Algorithm 2

8:          $\boldsymbol{x}_{t+1}^v \leftarrow \text{DEPROUND}(\boldsymbol{y}_{t+1}^v)$      ▷ Sample a discrete allocation

*Theorem 5.1:* INFIDA has a sublinear $(1 - 1/e)$-regret w.r.t. the time horizon $T$, i.e., there exists a constant $A$ such that:

$$(1 - 1/e)\text{-Regret}_{T,\mathcal{X}} \leq A\sqrt{T}, \qquad (21)$$

# Inference delivery network - proof

Salem, Castellano, Neglia, Pianese, Araldo, Toward Inference Delivery Networks: Distributing Machine Learning With Optimality Guarantees, **IEEE Trans. Net.** 2023

## Sketch of the proof

*Proof.* To prove the $\psi$-regret guarantee: *(i)* we first establish an upper bound on the regret of the INFIDA policy over its fractional allocations domain $\mathcal{Y}$ against a fractional optimum, then *(ii)* we use it to derive a corresponding $\psi$-regret guarantee over the integral allocations domain $\mathcal{X}$.

**Fractional domain regret guarantee.** To establish the regret guarantee of running Algorithm 1 at the level of each computing node $v \in \mathcal{V}$, we showed that the following properties hold:

1) The function $G$ is concave over its domain $\mathcal{Y}$ (Lemma F.1).
2) The mirror map $\Phi : \mathcal{D} \to \mathbb{R}$ is $\theta$-strongly convex w.r.t. the norm $\|\cdot\|_{l_1(s)}$ over $\mathcal{Y} \cap \mathcal{D}$, where $\theta$ is equal to Eq. (95) (Lemma F.2).
3) The gain function $G : \mathcal{Y} \to \mathbb{R}$ is $\sigma$-Lipchitz w.r.t $\|\cdot\|_{l_1(s)}$: the subgradients are bounded under the norm $\|\cdot\|_{l_\infty(\frac{1}{s})}$ by $\sigma$, i.e., the subgradient of $G(\boldsymbol{r}_t, \boldsymbol{l}_t, \boldsymbol{y})$ at point $\boldsymbol{y}_t \in \mathcal{Y}$ is upper bounded ($\|\boldsymbol{g}_t\|_{l_\infty(\frac{1}{s})} \leq \sigma$) for any $(\boldsymbol{r}_t, \boldsymbol{l}_t) \in \mathcal{A}$ (Lemma F.3).
4) $\|\cdot\|_{l_\infty(\frac{1}{s})}$ is the dual norm of $\|\cdot\|_{l_1(s)}$ (Lemma F.4).
5) The Bregman divergence $D_\Phi(\boldsymbol{y}_*, \boldsymbol{y}_1)$ in Eq. (63) is upper bounded by a constant $D_{\max}$ where $\boldsymbol{y}_* = \arg\max_{\boldsymbol{y} \in \mathcal{Y}} \sum_{t=1}^{T} G(\boldsymbol{r}_t, \boldsymbol{l}_t, \boldsymbol{y})$ and $\boldsymbol{y}_1 = \arg\min_{\boldsymbol{y} \in \mathcal{Y} \cap \mathcal{D}} \Phi(\boldsymbol{y})$ is the initial allocation (Lemma F.5).

We then apply results from [1]

[1] S. Bubeck, "Convex Optimization: Algorithms and Complexity," Foundations and Trends® in Machine Learning, Nov. 2015.

# Markov Decision Process

- State at timeslot $k$

$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)}, ..., \theta_P^{(k)})$$
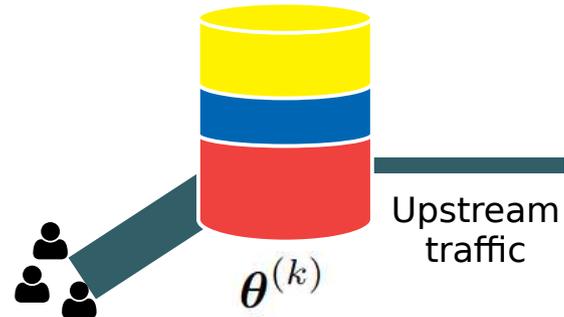
$\boldsymbol{\theta}^{(k)}$

# Markov Decision Process

- State at timeslot $k$
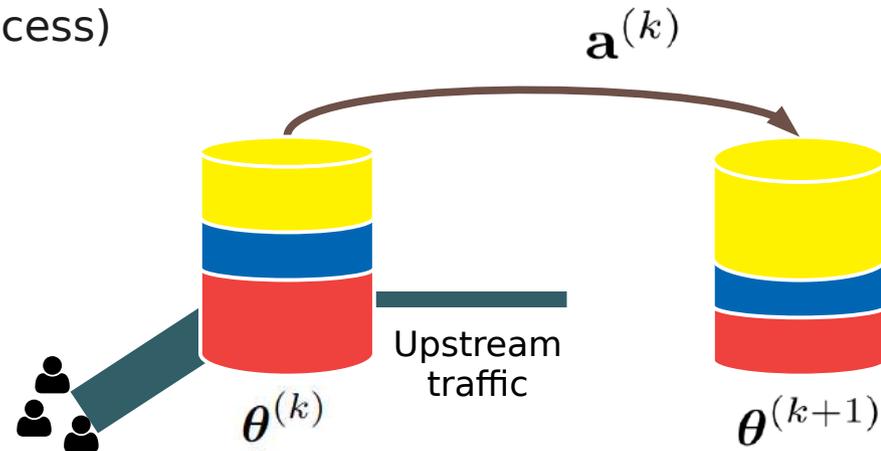
$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)}, ..., \theta_P^{(k)})$$

- Nominal cost (upstream traffic)

$$C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega)$$  Random (it depends on users' requests)



$\boldsymbol{\theta}^{(k)}$

Upstream traffic

# Markov Decision Process

- State at timeslot $k$

$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)}, ..., \theta_P^{(k)})$$
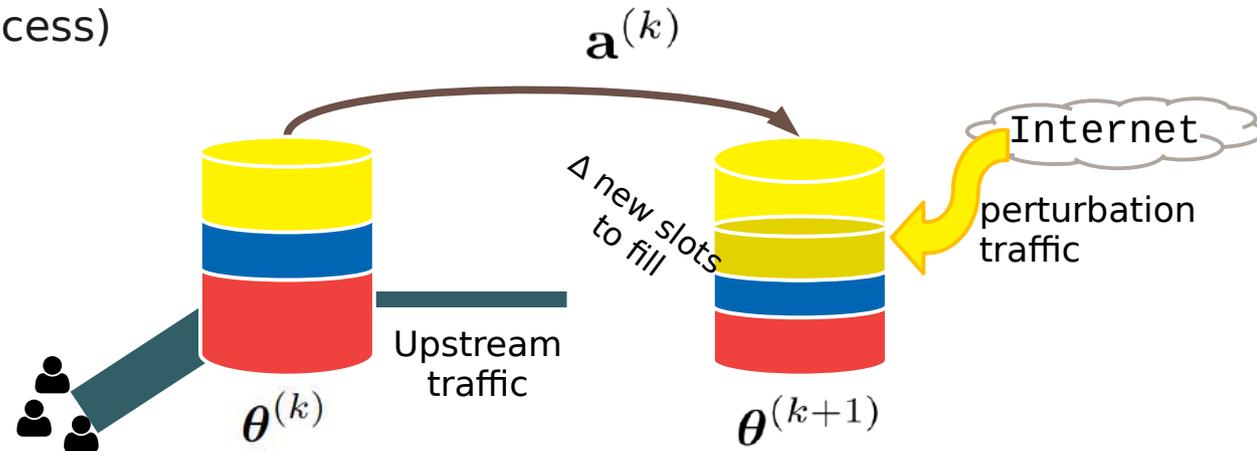
- Nominal cost (upstream traffic)

$$C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega)$$ Random (it depends on users' requests)

- Action

$$\mathbf{a}^{(k)} = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$$

(Deterministic Markov Decision Process)

$\mathbf{a}^{(k)}$

Upstream traffic

$\boldsymbol{\theta}^{(k)}$

$\boldsymbol{\theta}^{(k+1)}$

# Markov Decision Process

- State at timeslot $k$

$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)},...,\theta_P^{(k)})$$

- Nominal cost (upstream traffic)

$$C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega)$$   Random (it depends on users' requests)

- Action

$$\mathbf{a}^{(k)} = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$$

(Deterministic Markov Decision Process)

- Perturbation Cost

# Markov Decision Process

- State at timeslot $k$

$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)}, ..., \theta_P^{(k)})$$

- Nominal cost (upstream traffic)

$$C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega)$$  Random (it depends on users' requests)
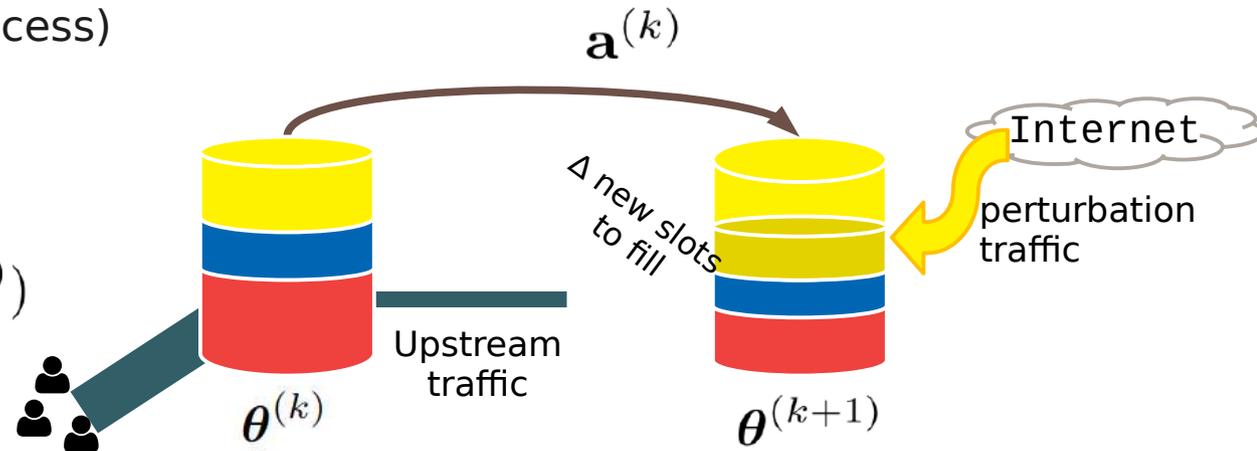
- Action

$$\mathbf{a}^{(k)} = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$$

(Deterministic Markov Decision Process)

- Perturbation Cost

- Instantaneous cost

$$C^{(k)} \triangleq C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega) + C_{\text{pert}}(\mathbf{a}^{(k)})$$

# Markov Decision Process

- State at timeslot $k$

$$\boldsymbol{\theta}^{(k)} = (\theta_1^{(k)},...,\theta_P^{(k)})$$

- Nominal cost (upstream traffic)

$$C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega)$$   Random (it depends on users' requests)

- Action

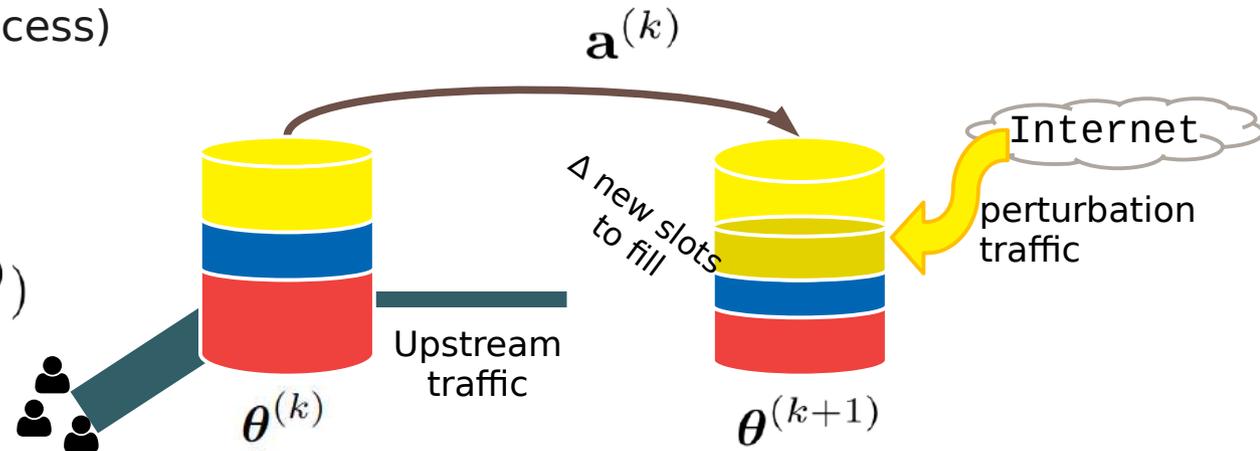$$\mathbf{a}^{(k)} = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$$

(Deterministic Markov Decision Process)

- Perturbation Cost

- Instantaneous cost

$$C^{(k)} \triangleq C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega) + C_{\text{pert}}(\mathbf{a}^{(k)})$$

$$C_{\text{cum}}^{\gamma} = \lim_{Z \to \infty} \mathbb{E} \left[ \sum_{k=0}^{Z} \gamma^{(k)} \cdot \underbrace{\left( C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega) + C_{\text{pert}}(\mathbf{a}^{(k)}) \right)}_{\text{Instantaneous cost } C^{(k)}} \right]$$



$\mathbf{a}^{(k)}$

Internet

Δ new slots to fill

perturbation traffic

Upstream traffic

$\boldsymbol{\theta}^{(k)}$

$\boldsymbol{\theta}^{(k+1)}$

**Algorithm 1:** $k$-th step of Model-based RL

1   $\alpha^{(k)} \leftarrow$ calculate the value of $\alpha$ ;    // formula (17)
2   $\epsilon^{(k)} \leftarrow$ calculate the value of $\epsilon$ ;    // formula (18)
3   with probability $\epsilon^{(k)}$: $\mathbf{a}^{(k)} \leftarrow$ random action ;
   // $\epsilon$-greedy policy
4   with probability $1 - \epsilon^{(k)}$: $\mathbf{a}^{(k)} \leftarrow$ best action from $Q^{(k)}(\boldsymbol{\theta}, \mathbf{a})$ ;
5   $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + \mathbf{a}^{(k)}$;
6   $C^{(k)} \leftarrow C_{\text{nom}}(\boldsymbol{\theta}^{(k)}, \omega) + C_{\text{pert}}(\mathbf{a}^{(k)})$;
7   $Q^{(k)}(\boldsymbol{\theta}^{(k)}, \mathbf{a}^{(k)}) \leftarrow (1 - \alpha^{(k)}) \cdot Q^{(k)}(\boldsymbol{\theta}^{(k)}, \mathbf{a}^{(k)}) +$
   $\alpha^{(k)} \cdot \left( C^{(k)} + \gamma \min_{\mathbf{a} \in \mathcal{A}_{\boldsymbol{\theta}^{(k+1)}}} Q^{(k)}(\boldsymbol{\theta}^{(k+1)}, \mathbf{a}) \right)$ ;
   // update $Q^{(k)}$
8   //////////////
9   /// Memory replay
10   $\mathcal{M}^{(k)} \leftarrow \mathcal{M}^{(k-1)} \cup \{(\boldsymbol{\theta}^{(k)}, \mathbf{a}^{(k)}, C_{\text{nom}}(\boldsymbol{\theta}^{(k)}))\}$;
11   **for** $N_{\text{memory}}$ times ;    // $N_{\text{memory}}$ is the size of the memory mini batch
12   **do**
13    $(\boldsymbol{\theta}^{\text{rd}}, \mathbf{a}^{\text{rd}}, C_{\text{nom}}^{\text{rd}}) \leftarrow$ random element from $\mathcal{M}^{(k)}$;
14    $\boldsymbol{\theta}'^{\text{rd}} \leftarrow \boldsymbol{\theta}^{\text{rd}} + \mathbf{a}^{\text{rd}}$;
15    $Q^{(k)}(\boldsymbol{\theta}^{\text{rd}}, \mathbf{a}^{\text{rd}}) \leftarrow (1 - \alpha^{(k)}) \cdot Q^{(k)}(\boldsymbol{\theta}^{\text{rd}}, \mathbf{a}^{\text{rd}}) + \alpha^{(k)} \cdot$
    $\left( C_{\text{nom}}^{\text{rd}} + C_{\text{pert}}(\mathbf{a}^{\text{rd}}) + \gamma \min_{\mathbf{a} \in \mathcal{A}_{\boldsymbol{\theta}'^{\text{rd}}}} Q^{(k)}(\boldsymbol{\theta}'^{\text{rd}}, \mathbf{a}) \right)$ ;
    // update $Q^{(k)}$
16   **end**

17   //////////////
18   /// Model training and inference
19   $\mathcal{D}_p^{(k)} \leftarrow \mathcal{D}_p^{(k-1)} \cup \{(\theta_p^{(k)}, C_{\text{nom,p}}(\theta_p^{(k)}))\}$;    // collect
   realization of $C_{\text{nom,p}}(\theta_p^{(k)})$ for each SP $p$
20   $\hat{C}_{\text{nom},p}^{(k)}(\theta_p) \leftarrow$ estimate model from $\mathcal{D}_p^{(k)}$ **for** $N_{\text{model}}$
   *times* ;    // $N_{\text{model}}$ is the size of the model mini batch
21   **do**
22    $\boldsymbol{\theta}^{\text{rd}} \leftarrow$ random state from $\mathcal{S}$;
23    $\mathbf{a}^{\text{rd}} \leftarrow$ random action from $\mathcal{A}_{\boldsymbol{\theta}^{\text{rd}}}$;
24    $\boldsymbol{\theta}'^{\text{rd}} \leftarrow \boldsymbol{\theta}^{\text{rd}} + \mathbf{a}^{\text{rd}}$;
25    Compute $\hat{C}_{\text{nom},p}^{(k)}(\theta_p^{\text{rd}})$ ;    // predict the nominal cost using the model
26    $\hat{C} \leftarrow \sum_{p=1}^P \hat{C}_{\text{nom,p}}(\theta_p^{\text{rd}})) + C_{\text{pert}}(\mathbf{a}^{\text{rd}})$;
27    $Q^{(k)}(\boldsymbol{\theta}^{\text{rd}}, \mathbf{a}^{\text{rd}}) \leftarrow (1 - \alpha^{(k)}) \cdot Q^{(k)}(\boldsymbol{\theta}^{\text{rd}}, \mathbf{a}^{\text{rd}}) + \alpha^{(k)} \cdot$
    $\left( \hat{C} + \gamma \min_{\mathbf{a} \in \mathcal{A}_{\boldsymbol{\theta}'^{\text{rd}}}} Q^{(k)}(\boldsymbol{\theta}'^{\text{rd}}, \mathbf{a}) \right)$ ;    // update $Q^{(k)}$
28   **end**

# RL for caching - proof

Ben-Ameur, Araldo, Chahed, Cache Allocation in Multi-tenant Edge Computing: An Online Model-based Reinforcement Learning Approach, **IEEE ICC** 2022 & **Major Revision in IEEE Trans. Cloud Comp.**

**Theorem V-B.1.** *If the discount factor $\gamma$ is sufficiently close to 1*

$$\lim_{k \to \infty} \boldsymbol{\theta}^{(k)} = \hat{\boldsymbol{\theta}}^* \text{ with probability } 1.$$

Sketch of the proof

- We prove that our Q-table $Q^{(k)}$ converges to the optimal Q-table $Q^*$ with probability 1.
- We prove that the sequence of actions and states induced by $Q^*$ has an absorbing state that is the discretely optimal state $\hat{\boldsymbol{\theta}}^*$.
- We prove that the sequence of actions and states induced by our Q-table $Q^{(k)}$ also follows $Q^*$.
- We prove that the sequence of actions and states that we take online converges with probability 1 to the sequence induced by our Q-table $Q^{(k)}$ (assuming no more exploration).
- Finally, we show that this sequence converges with probability 1 to a sequence induced by $Q^*$.