

Towards Inference Delivery Networks: Distributing Machine Learning with Optimality Guarantees

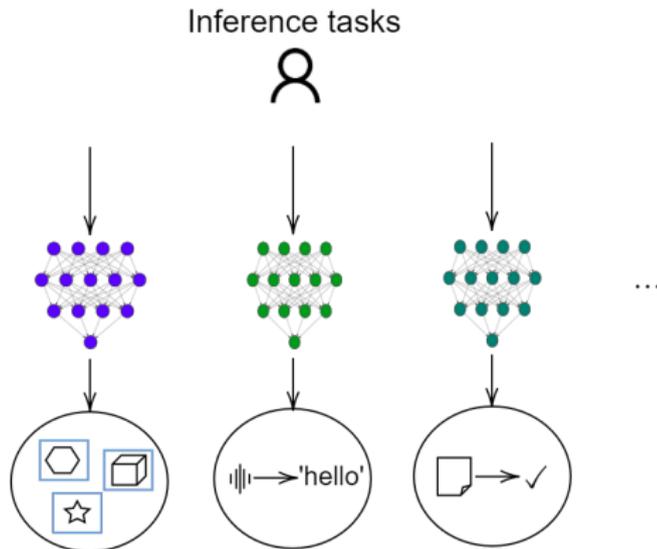
Tareq Si Salem¹, Gabriele Castellano^{1,2}, Giovanni Neglia¹, Fabio Pianese², and Andrea Araldo³

¹Inria, Université Côte d'Azur, France

²Nokia Bell Labs, France

³Télécom SudParis - Institut Polytechnique de Paris, France

Machine Learning Tasks



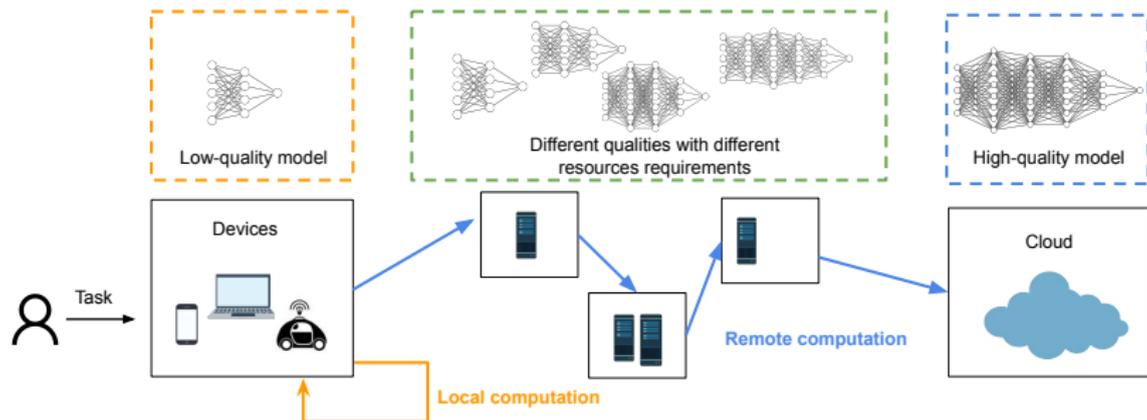
Inference tasks are generated by users and executed through pretrained inference models.

Inference Delivery Networks



- Simpler models available locally have low accuracy
- Complex models at the cloud may not meet latency constraints

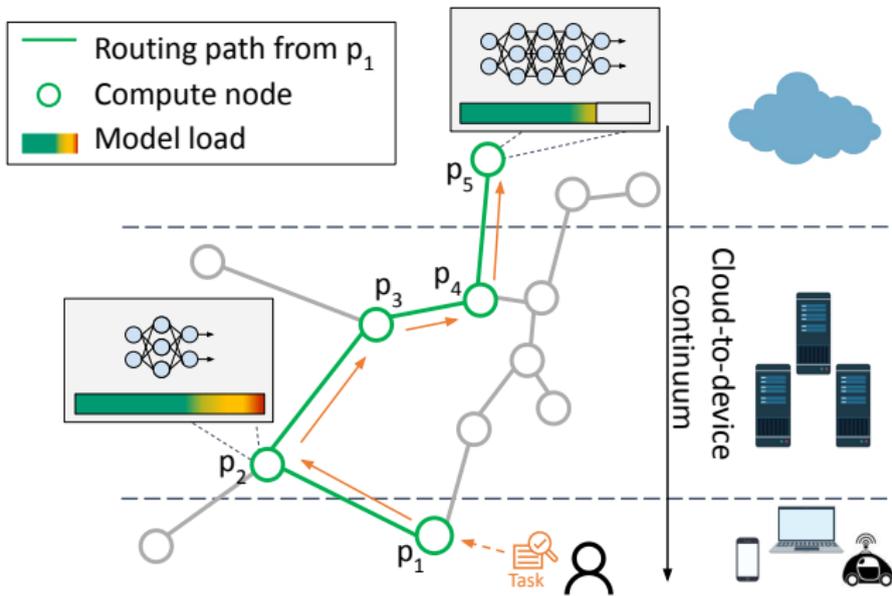
Inference Delivery Networks



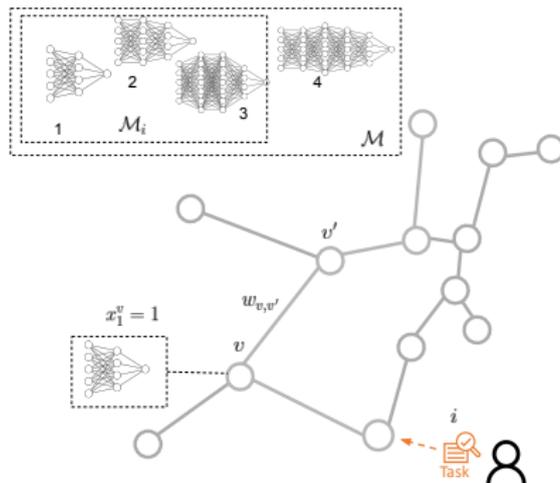
Integrate ML inference in the continuum between end-devices and the cloud.

- Present inference delivery networks
- Propose INFIDA, a distributed online allocation algorithm for IDNs with strong guarantees even w.o. a prior on the request process
- Evaluate INFIDA experimentally with greedy heuristics

System Overview

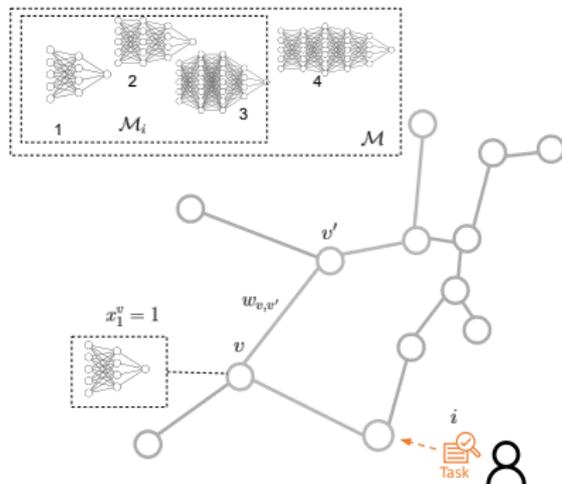


System Description - Compute Nodes and Models (1/2)



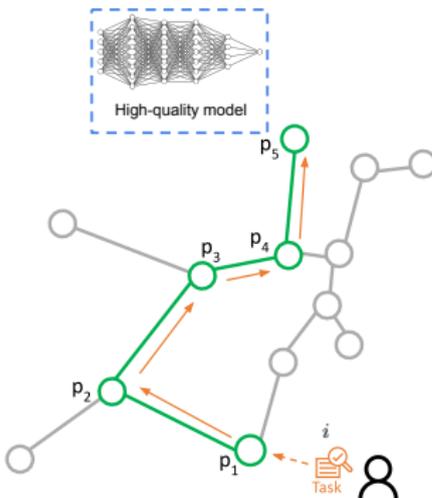
- We represent the inference delivery network (IDN) as a weighted graph $G(\mathcal{V}, \mathcal{E})$
- $\mathcal{N} = \{1, 2, \dots, N\}$ is the set of tasks the system can serve

System Description - Compute Nodes and Models (2/2)



- Each node $v \in \mathcal{V}$ has an *allocation budget* $b^v \in \mathbb{R}_+$, and $s_m^v \in \mathbb{R}_+$ is the size of model $m \in \mathcal{M}$
- The budget constraints: $\sum_{m \in \mathcal{M}} x_m^v s_m^v \leq b^v, \forall v \in \mathcal{V}$

System Description - Inference Requests



- We assume that every node has a predefined routing path towards a suitable repository node for each task $i \in \mathcal{N}$
- A request is determined by the pair (i, \mathbf{p})

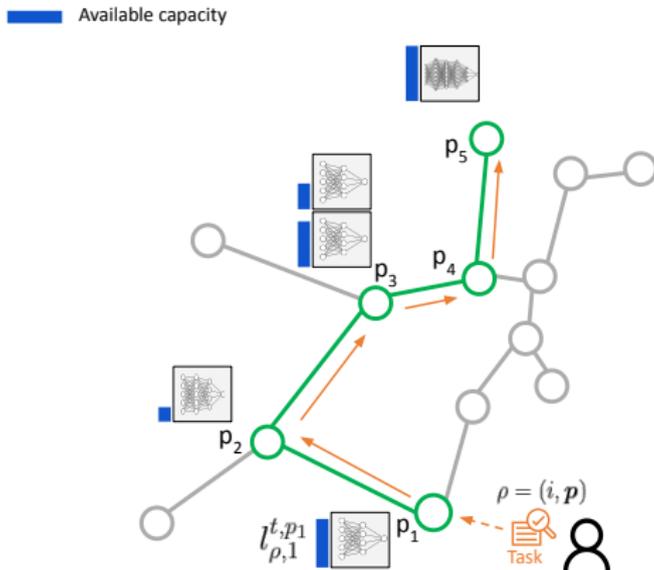
System Description - Cost Model

When serving request $\rho=(i, \mathbf{p}) \in \mathcal{R}$ on node p_j using model m , the system experiences a cost $C_{\mathbf{p},m}^{p_j} \in \mathbb{R}_+$.

Our theoretical results hold under this very general cost model, but to be concrete we refer to the following simpler model:

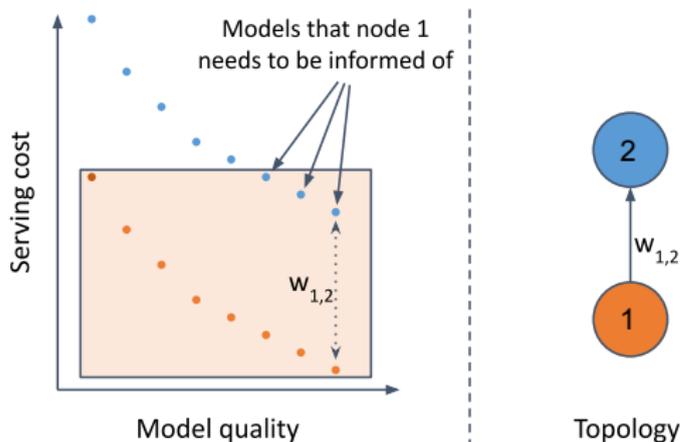
$$C_{\mathbf{p},m}^{p_j} = \underbrace{\sum_{j'=1}^{j-1} w_{p_{j'}, p_{j'+1}}}_{\text{round-trip latency}} + \underbrace{d_m^{p_j}}_{\text{inference delay}} + \underbrace{\alpha}_{\text{trade-off parameter}} \underbrace{(1-a_m)}_{\text{inaccuracy}}.$$

System Description - Request Load and Serving Capacity



- During a slot t the system receives a batch of requests $\mathbf{r}_t = [r_{\rho}^t]_{\rho \in \mathcal{R}} \in (\mathbb{N} \cup \{0\})^{\mathcal{R}}$
- Each model has an *available capacity* $l_{\rho,m}^{t,v} \in \mathbb{N} \cup \{0\}$

System Description - Serving Model (1/3)



For a model m allocated at v having the k -th smallest service cost, we denote by:

$$\underbrace{\gamma_{\rho}^k = C_{\mathbf{p},m}^v}_{\text{model service cost}}, \quad \underbrace{\lambda_{\rho}^k(I_t) = I_{\rho,m}^{t,v}}_{\text{potential available capacity}}, \quad \underbrace{z_{\rho}^k(I_t, \mathbf{x}) = x_m^v I_{\rho,m}^{t,v}}_{\text{effective available capacity}}.$$

System Description - Serving Model (2/3)

The aggregate cost incurred by the system at time slot t is

$$C(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x}) = \sum_{\rho \in \mathcal{R}} \sum_{k=1}^{K_\rho} \gamma_\rho^k \cdot \underbrace{\min \left\{ r_\rho^t - \sum_{k'=1}^{k-1} z_\rho^{k'}(\mathbf{l}_t, \mathbf{x}), z_\rho^k(\mathbf{l}_t, \mathbf{x}) \right\}}_{\text{(a)}} \cdot \underbrace{\mathbb{1} \left\{ \sum_{k'=1}^{k-1} z_\rho^{k'}(\mathbf{l}_t, \mathbf{x}) < r_\rho^t \right\}}_{\text{(b)}}$$

- **(a)** is the number of requests served by the k -th ranked model
- **(b)** is zero when all the requests can be served before the k -th ranked model

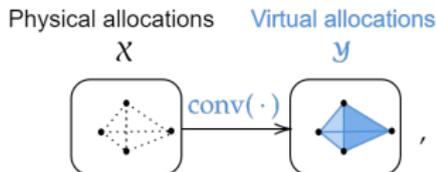
Equivalently, our objective is to maximize the *allocation gain* defined as $G(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x}) = C(\mathbf{r}_t, \mathbf{l}_t, \boldsymbol{\omega}) - C(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x})$.

System Description - Serving Model (3/3)

The allocation gain has the following equivalent expression:

$$G(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x}) = \sum_{\rho \in \mathcal{R}} \sum_{k=1}^{K_\rho-1} \underbrace{(\gamma_\rho^{k+1} - \gamma_\rho^k)}_{\text{cost saving}} \underbrace{(Z_\rho^k(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x}) - Z_\rho^k(\mathbf{r}_t, \mathbf{l}_t, \boldsymbol{\omega}))}_{\text{additional requests}},$$

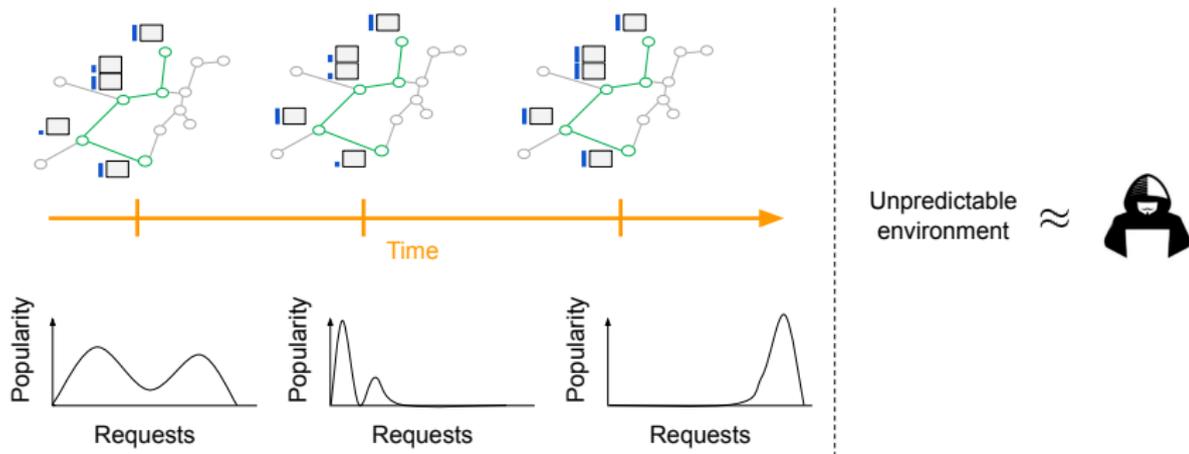
where $Z_\rho^k(\mathbf{r}_t, \mathbf{l}_t, \mathbf{x}) \triangleq \min \left\{ r_\rho^t, \sum_{k'=1}^k z_\rho^{k'}(\mathbf{l}_t, \mathbf{x}) \right\}$.



$G(\mathbf{r}_t, \mathbf{l}_t, \mathbf{y})$ is concave over \mathbf{y}

$$= \text{[3D plot 1]} + \text{[3D plot 2]} + \text{[3D plot 3]} + \dots$$

System Description - Adversarial Setting



- We consider a “pessimistic” scenario where both requests and available capacities are selected by an adversary

Subgradients computation

For every $v \in \mathcal{V}$ the gain function has a subgradient \mathbf{g}_t^v at point $\mathbf{y}_t^v \in \mathcal{Y}^v$ given by

$$\mathbf{g}_t^v = \left[\sum_{\rho \in \mathcal{R}} l_{\rho,m}^{t,v} \left(\gamma_{\rho}^{K_{\rho}^*(\mathbf{y}_t)} - C_{\rho,m}^v \right) \mathbb{1}_{\{\kappa_{\rho}(v,m) < K_{\rho}^*(\mathbf{y}_t)\}} \right]_{m \in \mathcal{M}},$$

where $K_{\rho}^*(\mathbf{y}_t) = \min \{k \in [K_{\rho} - 1] : \sum_{k'=1}^k z_{\rho}^{k'}(\mathbf{l}_t, \mathbf{y}_t) \geq r_{\rho}^t\}$.

INFIDA distributed allocation

- 1: **procedure** INFIDA($\mathbf{y}_1^v = \arg \min_{\mathbf{y}^v \in \mathcal{Y}^v \cap \mathcal{D}^v} \Phi^v(\mathbf{y}^v)$, $\eta^v \in \mathbb{R}_+$)
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: $\mathbf{x}_t^v \leftarrow \text{DepRound}(\mathbf{y}_t^v)$ \triangleright Sample a physical allocation
- 4: Compute $\mathbf{g}_t^v \in \partial_{\mathbf{y}^v} G(\mathbf{r}_t, \mathbf{I}_t, \mathbf{y}_t)$
- 5: $\hat{\mathbf{y}}_t^v \leftarrow \nabla \Phi^v(\mathbf{y}_t^v)$ \triangleright Map state to the dual space
- 6: $\hat{\mathbf{h}}_{t+1}^v \leftarrow \hat{\mathbf{y}}_t^v + \eta^v \mathbf{g}_t^v$ \triangleright Take gradient step in the dual space
- 7: $\mathbf{h}_{t+1}^v \leftarrow (\nabla \Phi^v)^{-1}(\hat{\mathbf{h}}_{t+1}^v)$ \triangleright Map dual state back to the primal space
- 8: $\mathbf{y}_{t+1}^v \leftarrow \mathcal{P}_{\mathcal{Y}^v \cap \mathcal{D}^v}^{\Phi^v}(\mathbf{h}_{t+1}^v)$ \triangleright Project new state onto the feasible region
- 9: **end for**
- 10: **end procedure**

Theoretical Guarantees

We provide the optimality guarantees of INFIDA in terms of the ψ -regret ($\psi = 1 - \frac{1}{e}$).

ψ -Regret

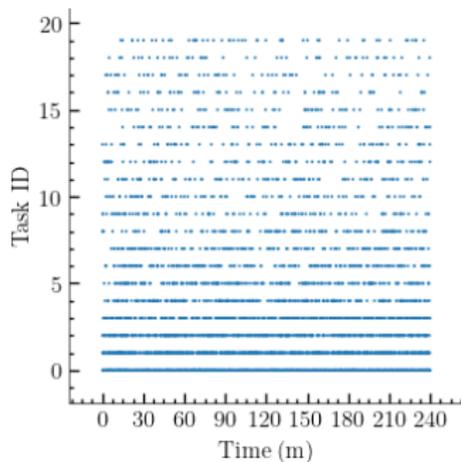
$$= \underbrace{\psi}_{\text{discount}} \text{ Total gain of the static optimum} - \text{Expected total gain of INFIDA}$$

$$\leq \psi \frac{RL_{\max} \Delta_C}{s_{\min}} \sqrt{2s_{\max} |\mathcal{V}| |\mathcal{M}| \sum_{v \in \mathcal{V}} \min\{b^v, \|\mathbf{s}^v\|_1\} \log \left(\frac{\|\mathbf{s}^v\|_1}{\min\{b^v, \|\mathbf{s}^v\|_1\}} \right) T}$$

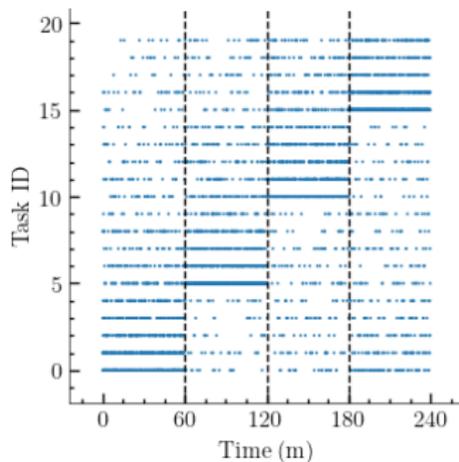
$$= \mathcal{O}(\sqrt{T}).$$

- The average ψ -regret goes to 0 when T is large enough; we perform as well as a ψ -approximation of the static optimum that knows the sequence of the models' capacities and requests in hindsight!
- ψ is the best approximation bound achievable for the problem, assuming $P \neq NP$

Experiments - Requests Traces



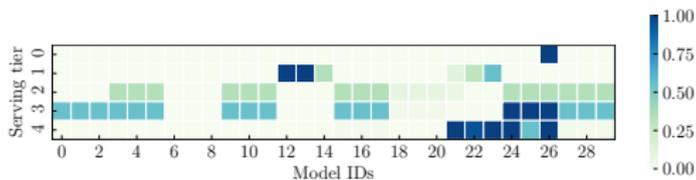
(a) Fixed Popularity Profile



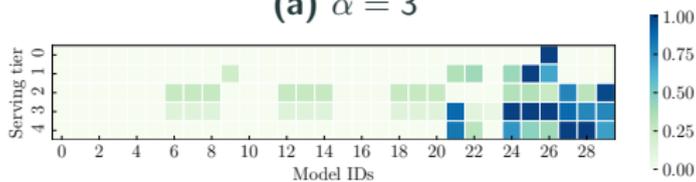
(b) Sliding Popularity Profile

Performance Metric. The performance of a policy is evaluated in terms of the time averaged gain normalized to the number of requests per second (NTAG).

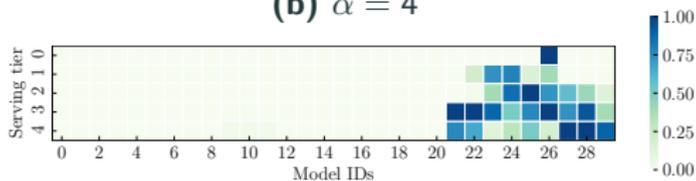
Experiments - Trade-off between Latency and Accuracy



(a) $\alpha = 3$



(b) $\alpha = 4$



(c) $\alpha = 5$

Figure 2: Fractional allocation decisions y_m^v of INFIDA on the various tiers of *Network Topology I* under *Fixed Popularity Profile*.

Experiments - Trade-off between Latency and Accuracy (1/3)

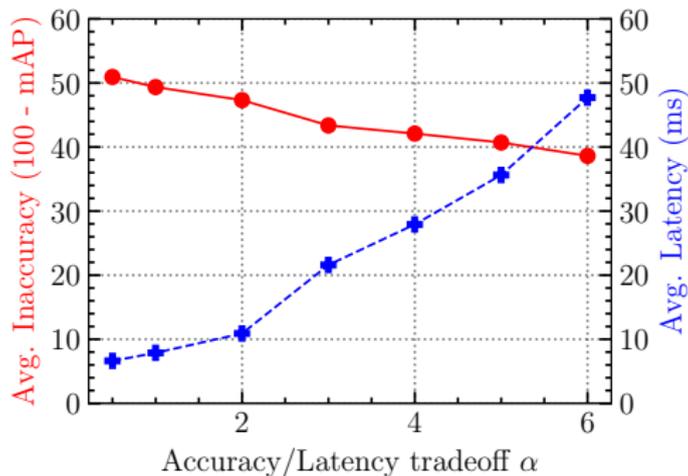


Figure 3: Average latency (dashed line) and inaccuracy (solid line) costs experienced with INFIDA for different values of α under *Network Topology I* and *Fixed Popularity Profile*.

Experiments - Trade-off between Latency and Accuracy (2/3)

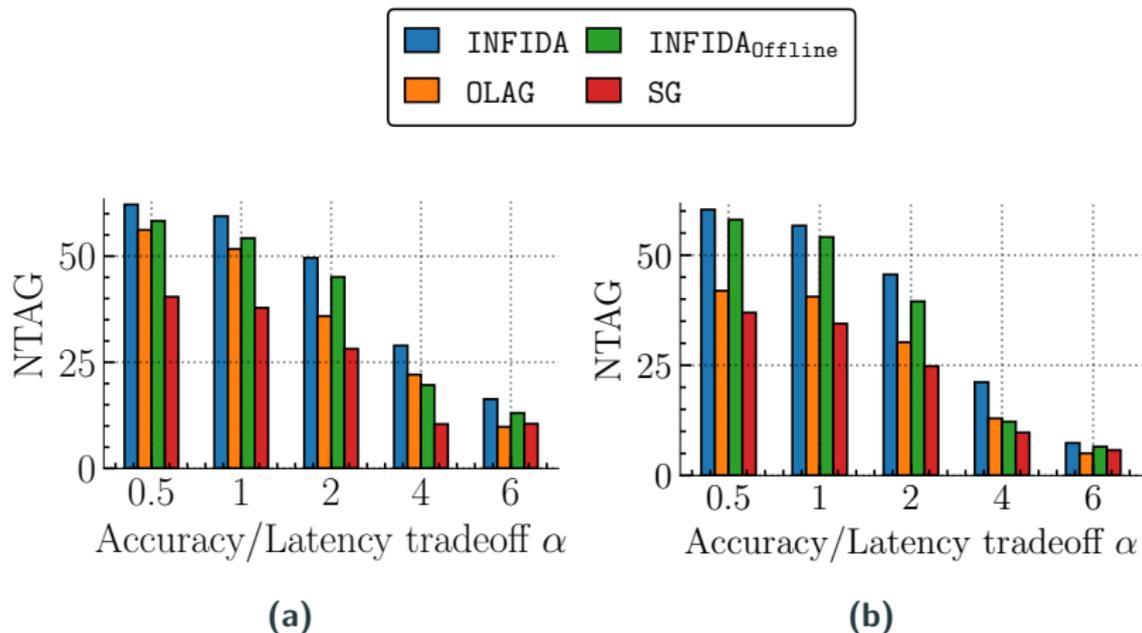


Figure 4: NTAG of the different policies under *Sliding Popularity Profile* and network topologies: (a) *Network Topology I*, and (b) *Network Topology II*.

Experiments - Trade-off between Latency and Accuracy (3/3)

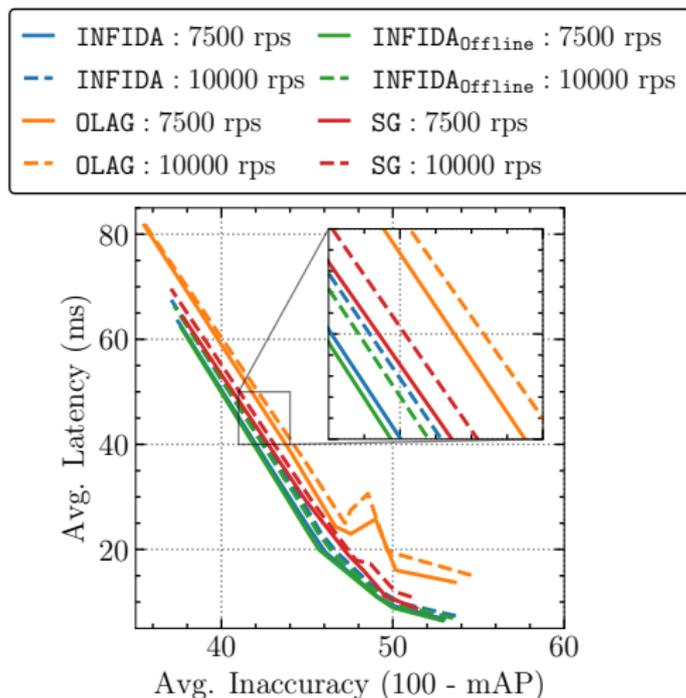


Figure 5: Average Latency vs. Average Inaccuracy obtained for different values of $\alpha \in \{0.5, 1, 2, 3, 4, 5, 6\}$ under *Fixed Popularity Profile* and *Network Topology II*.

Experiments - Update Costs

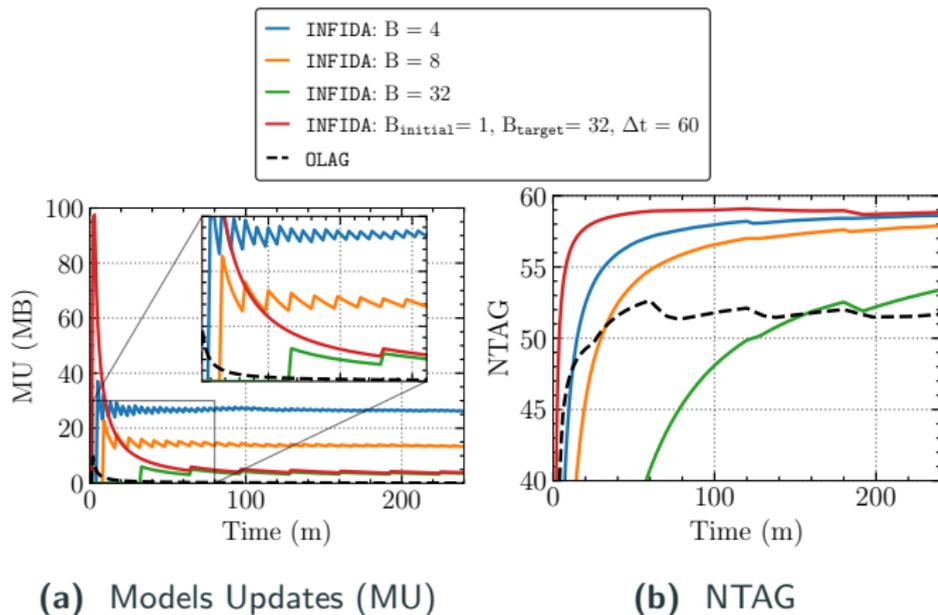
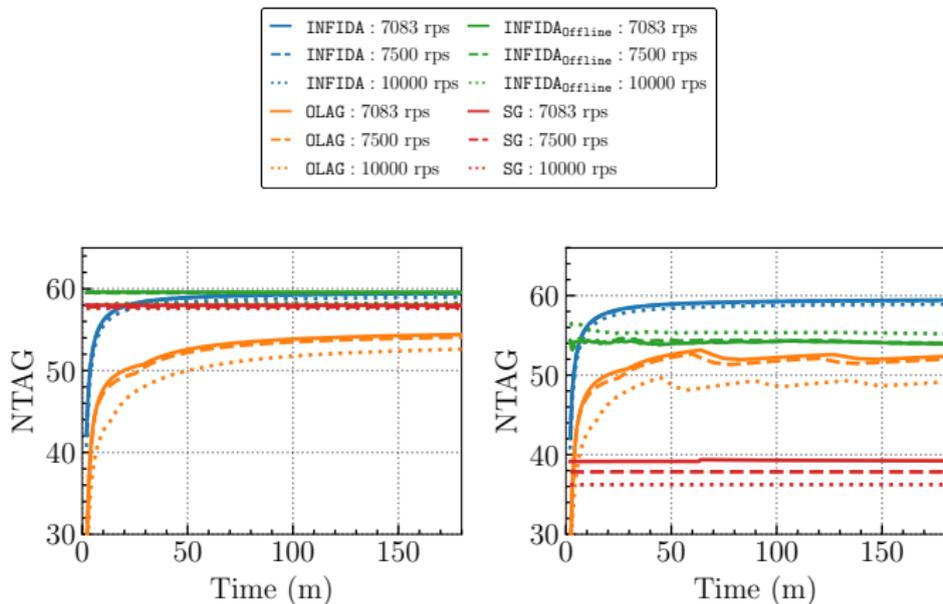


Figure 6: (a) Models Updates (MU) and (b) NTAG of Greedy and INFIDA for different values of refresh period $B \in \{4, 8, 16\}$, and for a dynamic refresh period. The experiment is run under *Network Topology I* and *Sliding Popularity Profile*.

Experiments - Scalability on Requests Load



(a) Fixed Popularity Profile

(b) Sliding Popularity Profile

Figure 7: NTAG of the different policies for different request rates under *Network Topology I*.

Thank you for your attention!