# Dynamic Incentives for Efficient Traffic Management: A Reinforcement Learning Approach

Germán Pardo-González[1], Shaghayegh Vosough[2], Katerina Papadaki[1], and Claudio Roncoli[2,3]

[1]Department of Mathematics, London School of Economics and Political Science, UK
[2]Department of Built Environment, Aalto University, Finland
[3]Centre for Industrial Management / Traffic and Infrastructure, KU Leuven, Belgium

## SHORT SUMMARY

Traffic management literature often overlooks encouragement-based strategies in favour of road pricing. However, road pricing often raises concerns about accessibility and public dissatisfaction, leading to its prohibition in some places such as Finland. Similar to pricing, incentives can push the User-Equilibrium (UE) flow pattern towards the System Optimum (SO). We formulate and solve a problem to determine the allocation of incentives, encouraging users to reroute onto an alternative (potentially longer) path, to achieve overall social benefit. The proposed approach uses multi-agent reinforcement learning to assign incentives to drivers, in which travel times are estimated dynamically using a traffic simulation software (SUMO). We also introduce a dynamic pseudo-SO as a benchmark for evaluating the incentives' effectiveness. Under an unlimited budget, the incentives scheme achieves virtually identical performance to the dynamic pseudo-SO. As the budget decreases, the solution gradually degrades, reaching the dynamic UE. The numerical results demonstrate that unlimited incentives can reduce the total travel time by an average of 23%. However, employing a budget time equivalent to around 6.3% of the UE total travel time can achieve a 10% reduction in total travel time.

**Keywords**: Traffic management, Incentives, Multi-agent reinforcement learning, Q-learning, Traffic simulation.

## 1 INTRODUCTION

*Motivation*

Nowadays, large cities need reliable, robust, and efficient transport. To achieve this, the implementation of management strategies to reduce traffic congestion and emissions is paramount. A well-known strategy applied to congested cities such as London (TfL, 2024) is tolling. While area-based road pricing, like in London, primarily reduces demand by discouraging car use in specific zones, route- or link-based pricing mechanisms can directly influence travel behaviours to align traffic flows. Route- or link-based pricing can steer the user equilibrium (UE) flow pattern towards system optimal (SO), namely where the total social benefit reaches the highest level. However, road pricing usually causes public dissatisfaction, inequitable welfare distribution, and legal/policy restrictions. That is why incentivising schemes with voluntary participation have recently gained more popularity (Niroumand et al., 2025; Luan et al., 2023).

By implementing an incentive scheme, we aim to reduce the network's overall congestion by rerouting drivers and steering them to follow SO routes. As an example, imagine you are deciding whether to take route A or B, where the latter has a longer travel time. Your motivation to choose route B lies in the incentive the road authority offers for it. Opting for route B would help to alleviate congestion on A and potentially reduce the total travel time (TTT) in the whole network. Due to budget limitations, the efficient assignment of incentives is important. Under an unlimited budget, an incentives scheme minimising the TTT in the network can achieve SO in the same way as pricing schemes. However, given a limited budget, there is no guarantee to achieve SO and the aim is to assign incentives in an optimal way so that the TTT is the minimum (Niroumand et al., 2025).

*Background*

To achieve SO traffic flow, some drivers need to take routes (slightly) longer than their shortest paths (Van Essen et al., 2019). Therefore, a strong stimulus, such as an incentive, can be employed to encourage drivers to choose routes that may be less desirable than their preferred (e.g., faster) options to steer traffic flow towards SO (Vosough & Roncoli, 2024; Klein & Ben-Elia, 2018; Djavadian et al., 2014). While many studies have examined the impacts of incentive schemes on network flow and total travel time under static traffic conditions (Niroumand et al., 2024; Luan et al., 2023), it is essential to recognise the dynamic nature of traffic. This dynamic aspect underscores the importance of developing incentive schemes within a dynamic framework.

Developing effective traffic management strategies, e.g., incentive schemes, within a dynamic framework is challenging due to the inherent complexity of dynamic traffic assignment. To address these challenges, simulation tools such as SUMO (Simulation of Urban MObility) can be valuable. Gawron (1998) first introduces a basic yet common algorithm for computing dynamic UE. However, the study reveals that even with a two-link network, the algorithm is unstable and does not always converge. Hence, the author presents a modification that stores a subset of paths for every O-D pair, and for each, it computes a probability distribution that assigns the probability of assignment to those paths. This concept is implemented in SUMO's duaIterate tool, which provides an approximate solution to dynamic UE by iteratively refining path probabilities.

Grunitzki & Bazzan (2017) provided a first insight into the two possible ways of modelling dynamic traffic assignment using multi-agent reinforcement learning (MARL). The first one is the link-based approach where the agent needs to decide the link to traverse at each time step. The second approach is route-based meaning that the agent, given precomputed routes, decides the route to take. In another study, Shou et al. (2022) provided a survey of previous approaches to model the dynamic traffic assignment, showing the benefits and drawbacks of each one, although they do not compare their approach computationally. Their algorithm, based on deep Q-learning and bi-level optimisation, is the first one to include advanced MARL techniques. They tested their algorithm in a real-world network using SUMO.

While there are no studies directly investigating dynamic incentive schemes in urban transportation networks, some research has utilised simulation-based approaches to optimise tolls for traffic management. For instance, Stewart & Ge (2014) developed a solution framework based on simulation to perform toll optimisation in a dynamic UE setting. Also, Ramos et al. (2019) presented a formulation for the route-based traffic assignment by incorporating tolls determined through MARL. In their tolling mechanism, the reward function is designed as a weighted combination of route travel time and imposed tolls, allowing the user to prioritise between these objectives based on their preferences.

The objective of this study is motivated by the incentive scheme designed by Niroumand et al. (2024), where the authors investigated the effectiveness of static link and path incentive schemes on the TTT of the transportation networks under a budget limit. Building on this foundation, our work extends the analysis to a dynamic setting, introducing significant methodological differences.

*Research Contributions*

Our contributions to the field are as follows:

1. Developing a MARL-based algorithm for traffic assignment in SUMO to compute the dynamic pseudo-SO in the network, providing a benchmark for evaluating the proposed incentive scheme.

2. Designing a MARL-based algorithm to implement a dynamic incentive scheme for traffic management, aiming to push the UE towards the SO within the SUMO simulation environment.

3. Applying the proposed methods to a medium-sized urban network, i.e., the Kamppi area in Helsinki city centre, Finland, while incorporating a dynamic incentives scheme.

The remainder of this paper presents the methodology we employ to benchmark and investigate the incentives scheme in Section 2; the case study, outputs of the algorithms and the results' interpretations in Section 3; finally, conclusions are drawn in Section 4.

## 2   Methodology

***Markov Decision Process Formulation: Dynamic Pseudo-SO***

We formulate a Markov Decision Process (MDP) following the universal modelling framework from Powell (2019). Also, the static traffic assignment MARL algorithm by Grunitzki & Bazzan (2017) is adapted to a dynamic setting to compute a dynamic pseudo-SO solution. This will be used as a benchmark to assess the quality of the solution given by the incentives methodology. The MDP problem is solved by employing a solution algorithm based on Q-learning adapted to the multi-agent setup, namely the Independent Q-learning algorithm (Albrecht et al., 2024). Each traveller, with its origin and destination, represents an agent and has a corresponding Q-function.

*State Variable.* The state represents all the information needed to decide what happens after an action is taken. The goal is to find an optimal policy, i.e., a set of rules that tells the agent the best action to take. In this case, the "state" of an agent is simple: their origin and destination. Once an agent selects a route (an action), it reaches its destination without perturbations and thus it is not necessary to store the state. This is called a one-step MDP because the agent's decision is made at the start, and there are no further transitions or decisions along the way.

*Decision Variable.* The agent decides the path to take at the start of the trip. Let $\mathcal{W}$ be the set of agents and $\mathcal{P}_w$ the set of available paths for agent $w \in \mathcal{W}$. Also, for a given agent $w \in \mathcal{W}$, $X_w$ is a vector in $\mathbb{R}^{|\mathcal{P}_w|}$ with binary entries that represent the decision to take path $p \in \mathcal{P}_w$. It follows that $X_w = \{x_{wp}\}_{p \in \mathcal{P}_w}$ where $x_{wp}$ is a scalar that equals 1 if traveller $w \in \mathcal{W}$ takes path $p \in \mathcal{P}_w$ (and 0 otherwise). The decision variable for agent $w \in \mathcal{W}$ is defined as the vector $X_w$ whose dimension is the number of paths available. The action space for each agent looks as follows.

$$\mathcal{X}_w \in \left\{ x_w \in \{0,1\}^{|\mathcal{P}_w|} : x_w \mathbb{1} = 1 \right\}, \quad \forall w \in \mathcal{W}$$

for which $X_w$ takes values in $\mathcal{X}_w$. Moreover, the set that encompasses the decisions of all agents is $X = \{X_w\}_{w \in \mathcal{W}}$. Finally, we define the complete action space for all agents.

$$\mathcal{X} = \{\mathcal{X}_w\}_{w \in \mathcal{W}}.$$

*Reward Function.* The following reward function is defined

$$R_w(X) = \omega C_w(X) + (1 - \omega)\text{TTT}(X), \quad \forall w \in \mathcal{W} \tag{1}$$

where $\omega$ is a weight that ranges from 0 to 1. Each vehicle is treated as an autonomous agent, with specific characteristics, such as speed and acceleration. These rewards correspond to a weighted combination of the agent's travel time and the TTT, which are known only after all agents have made their routing decision and the behaviour of each individual driver on the network has been implemented. We define $C_w(X)$ to be the travel time experienced by agent $w \in \mathcal{W}$, given the choice of the other agents stored in $X$. Also, $\text{TTT}(X)$ is the estimated total travel time of the network that also depends on the actions of every agent.

*Objective Function.* For this part, we consider the classical action-value function (Q-function) which corresponds to the value of taking a given action while being in a given state. But, as we do not store the state, this reduces to an action function that represents the value of taking an action regardless of the state. We use a different Q-function for each agent:

$$Q_w^*(X_w) = \min_{X' \in \mathcal{X}} R_w(X'), \qquad \forall w \in \mathcal{W}. \tag{2}$$

The following equation is solved iteratively at every iteration $n$ using the following update rule, where there is no future state as the destination is a terminal:

$$Q_w(X_w) \leftarrow (1 - \alpha_n)Q_w(X_w) + \alpha_n R_w(X), \qquad \forall p \in \mathcal{P}_w, w \in \mathcal{W} \tag{3}$$

where $\alpha_n$ is the learning rate, that influences the learning performance of the algorithm. One typically starts with a big $\alpha_n$ so that the algorithm gives more importance to new information. At the end, we want a small $\alpha_n$ so that it converges.

*Solution Algorithm.* We employ a simulation-based independent multi-agent Q-learning method, as follows, to find the dynamic pseudo-SO. Note that the following algorithm does not consider incentives.

---

**Initialise**

    1. Compute $k$ shortest paths for every agent $w \in \mathcal{W}$ and store them in $\mathcal{P}_w$.

    2. Set $Q_w(X_w) \leftarrow 0$ for all agents $w \in \mathcal{W}$.

**For each episode** $n = 1, ..., N$**:**

    3. Set $\epsilon_n = 1 - \frac{n}{N}$ and $\alpha_n = \frac{a}{(b+n)}$.

  **For each agent** $w \in \mathcal{W}$**:**

    4. Take action $X_w$ as follows:

$$X_w = \begin{cases} \text{Random action } X_w \in \mathcal{X}_w & \text{with probability } \epsilon_n \\ \arg\min_{X_w \in \mathcal{X}_w} Q_w(X_w) & \text{with probability } 1 - \epsilon_n \end{cases}$$

    5. Add $X_w$ to $X$, the set of all actions taken.

  6. After all agents have completed their trips, retrieve $\text{TTT}(X)$ and travel time $C_w(X)$ for each agent.

  **For each agent** $w \in \mathcal{W}$**:**

    7. Compute the reward using Eq. (1) and update the Q-function with Eq. (3).

---

### *Markov Decision Process Formulation: Dynamic UE with Incentives*

We formulate a MDP similar to the previous section, using the universal modelling framework from Powell (2019) as well. It is based on MARL, uses simulation, and is novel in the literature. It aims at finding the optimal incentive for each driver under a limited budget. Then, the solution to this problem is compared to the solution provided by the dynamic pseudo-SO to show the effectiveness of the incentive scheme.

*State Variable.* This is a one-step MDP for which the state is not stored, as before.

*Decision Variable.* The agent should decide whether a path should be incentivised or not with the restriction that at most one can be incentivised. Note that the agent could have the decision that no paths are incentivised. Then, our decision variable is a vector $Y_w \in \mathbb{R}^{|\mathcal{P}_w|}$ with binary entries where 1 means that the path will be incentivised and 0 otherwise. Consider the action space, that represents set of all possible decisions for agent $w \in \mathcal{W}$:

$$\mathcal{Y}_w \in \left\{ y_w \in \{0,1\}^{|\mathcal{P}_w|} : y_w \mathbb{1} \leq 1 \right\}, \quad \forall w \in \mathcal{W}$$

for which $Y_w$ takes values in $\mathcal{Y}_w$.

Let the vector of path travel times of all paths of agent $w \in \mathcal{W}$ be $\tau_w \in \mathbb{R}^{|\mathcal{P}_w|}$. Let $\delta$ be the amount we reduce the time of the incentivised path. Thus, the travel times are adjusted as follows.

$$\tau_w \leftarrow \tau_w - Y_w \delta.$$

The above update only modifies a single entry of $\tau_w$, the one that corresponds to the incentivised path. Since we want the incentivised path to have the lowest travel time, we set the time reduction $\delta$ to be the difference between the incentivised path's travel time and the minimum path's travel time, plus an extra term $\phi > 0$:

$$\delta = \tau_w^T Y_w - \min\{\tau_w\} + \phi.$$

This will make the incentivised path's travel time to be $\phi$ units below the minimum travel time. Finally, the path with minimum cost is selected, which corresponds to the incentivised one, namely, $\arg\min_{p \in \mathcal{P}_w}\{\tau_w\}$. If there are no incentivised paths, all entries of $Y_w$ will be 0, and no modification will be made.

*Budget Limitation.* Let $B$ be the total budget available for incentivising drivers and $b$ the budget used so far, for convenience the budget is measured in time units. The algorithm checks whether there is enough budget to complete the action, i.e., $\delta + b \leq B$. If there is not enough budget, it

does not modify the travel times and selects the original shortest path. The budget is continuously tracked throughout the process, representing the state of the central authority responsible for assigning incentives.

*Reward Function.* To compute the reward we need $X$, as it gives the paths used by all agents. Even though the decision variable is $Y_w$ we still use $X_w$ and $X$ to calculate the reward. In the case that a path is incentivised, it will be used, and thus $X_w = Y_w$. However, when $Y_w = 0$ (no path is incentivised), then $X_w$ will have entry 1 for the path given by $\arg\min_{p \in \mathcal{P}_w}\{\tau_w\}$, which is the shortest path.

The definition of reward $R_w(X)$ is described as in Eq. (1), like in the previous MDP.

*Objective Function.* We define the following Q-function for each agent:

$$Q_w^*(Y_w) = \min_{X' \in \mathcal{X}} R_w(X'). \tag{4}$$

The above equation is solved iteratively for every episode $n$ using the following update rule.

$$Q_w(Y_w) \leftarrow (1 - \alpha_n)Q_w(Y_w) + \alpha_n R_w(X), \qquad \forall w \in \mathcal{W} \tag{5}$$

where the learning rate is $\alpha_n$.

*Solution Algorithm.* Consider the solution algorithm next.

**Initialise**

1. Compute $k$ shortest paths for every agent $w \in \mathcal{W}$ and store them in $\mathcal{P}_w$.

2. For each $w \in \mathcal{W}$, set the entries of $\tau_w$ to be the free-flow travel times of $w$'s chosen paths from SUMO.

3. Set $Q_w(Y_w) \leftarrow 0$ for all agents $w \in \mathcal{W}$ and set $\phi$ as a small number.

**For each episode** $n = 1, ..., N$**:**

4. Set $\epsilon_n = 1 - \frac{n}{N}$, $\alpha_n = \frac{a}{(b+n)}$ and $b = 0$.

**For each agent** $w \in \mathcal{W}$**:**

5. Take action $Y_w$ as follows:
$$Y_w = \begin{cases} \text{Random action } Y_w \in \mathcal{Y}_w & \text{with probability } \epsilon_n \\ \arg\min_{Y_w \in \mathcal{Y}_w} Q_w(Y_w) & \text{with probability } 1 - \epsilon_n \end{cases}$$

6. Calculate the amount of the incentive $\delta$:
$$\delta = \tau_w^T Y_w - \min\{\tau_w\} + \phi$$

7. If there is enough budget left, i.e., $\delta + b \leq B$:

   Modify travel times $\tau_w$ and update budget $b$:
$$\tau_w \leftarrow \tau_w - Y_w \delta$$
$$b \leftarrow b + \delta$$

   Else:

   Do not apply the incentive or modify the travel times $\tau_w$.

8. Select the path to be taken: $p = \arg\min_{p \in \mathcal{P}_w}\{\tau_w\}$. Let $X_w$ be the decision vector of choosing path $p$.

9. Add $X_w$ to $X$, the set of all paths taken by agents.

10. After all agents have completed their trips, retrieve $\text{TTT}(X)$, travel time $C_p(X)$ from SUMO, and update the travel times $\tau_w$ for each agent.

**For each agent** $w \in \mathcal{W}$**:**

11. Compute the reward using Eq. (1), and update the Q-function with Eq.(5).

## 3 CASE STUDY AND RESULTS

### Case study description

We apply our proposed methodology to the Kamppi area, Helsinki, Finland, on the network depicted in Fig. 1. This area is known for its congestion during peak hours and central location.

The data, for simulation in SUMO, including a set of trips with departure times for each traveller and the network definition is obtained from Bochenina et al. (2023). The network contains 235 edges and 152 nodes and there are 1100 trips (same as O-D pairs).

### Numerical results

We demonstrate the efficiency of the incentive scheme. When operating with an unlimited budget, the incentive method achieves performance comparable to the dynamic pseudo-SO, as illustrated in Fig 2, with the TTT being nearly 23% lower than the UE's average over the last 100 episodes (simulation iterations). This effectively highlights how, in a realistic scenario, incentives can significantly alleviate congestion in a heavily congested network. Achieving a 23% reduction in TTT is a substantial milestone, capable of making a significant difference in real-world networks.
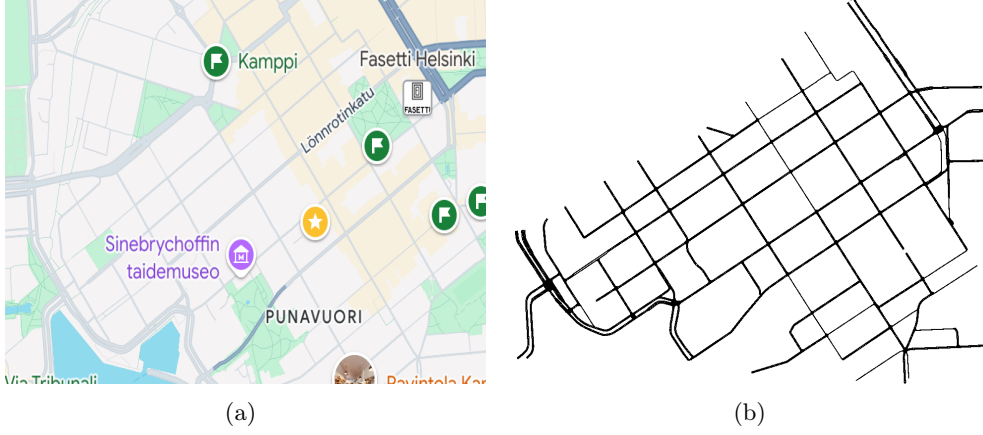
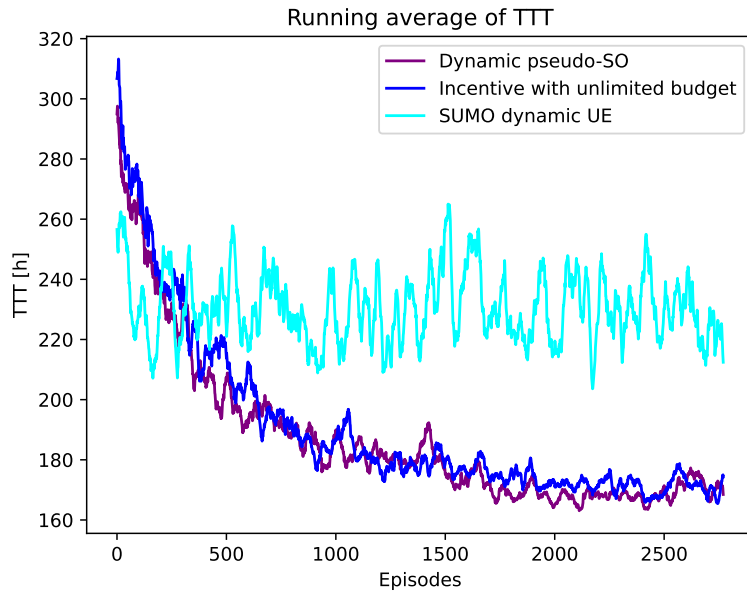Figure 1: Kamppi, Helsinki Area transportation network.



Figure 2: Total travel time comparison with and without incentives.

In Fig. 3, we demonstrate that as the budget becomes increasingly constrained, the TTT of the network intuitively rises. With fewer paths incentivised, travel time increases due to selfish routing behaviour, shifting the solution closer to the UE (or worse). However, this eventually converges after a substantial number of iterations. Our findings align with the static incentive scheme investigated by Niroumand et al. (2024).

The mean and standard deviation of the last 100 episodes are shown in Table 1. It is important to emphasise that our algorithms are more reliable than Gawron's algorithm of SUMO, as they exhibit a lower standard deviation, meaning that they have a better travel time reliability. Moreover, the incentive scheme with an unlimited budget achieves performance nearly identical to the dynamic pseudo-SO. However, it is worth noting that with no budget, our algorithm performs worse than SUMO's dynamic UE, as it ultimately resorts to greedily assigning the shortest paths.

## 4  CONCLUSIONS

We implement a dynamic incentives scheme aimed at reducing overall congestion in the network. With this, we demonstrate how encouragement can lead to social benefits by minimising total travel
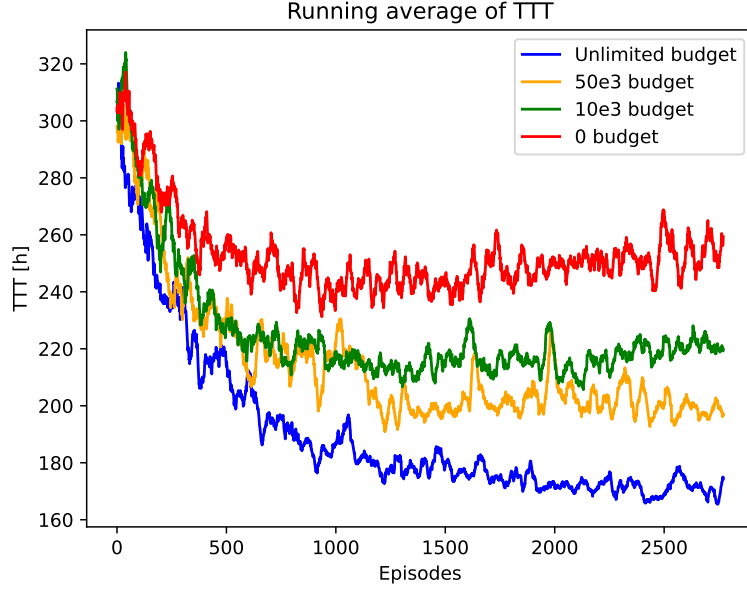
Figure 3: Total travel time comparison for different values of the budget available for incentives.

Table 1: Mean and standard deviation of TTT for the last 100 episodes.

| Scenario | Mean [h] | Standard deviation [h] |
|---|---|---|
| SUMO Dynamic UE | 219.77 | 30.69 |
| Dynamic Pseudo-SO | 169.23 | 10.49 |
| Unlimited budget | 170.02 | 11.60 |
| 50e3 budget | 198.11 | 11.38 |
| 10e3 budget | 220.55 | 14.82 |
| 0 budget (UE) | 257.96 | 28.38 |

time while offering a more publicly acceptable alternative to road pricing, which often causes public dissatisfaction and uneven welfare distribution.

To this end, we develop two MDP formulations and algorithms to perform traffic assignment with and without incentives. The latter that performs the SO assignment is called dynamic pseudo-SO as there is no way to prove the optimality of an actual SO in a dynamic situation. The dynamic pseudo-SO is used as a benchmark to show the efficiency of the incentive schemes in the network. The solution algorithms are based on MARL where SUMO is the environment to estimate travel times dynamically for more realistic results and we use simulated data from the Kamppi area in Helsinki, Finland. The results indicate that our algorithms are reliable as for an unlimited budget, the incentive method demonstrates comparable performance to the dynamic pseudo-SO. When the budget is limited the solutions move to the dynamic UE as the budget decreases. Additionally, our algorithms are more reliable than SUMO's as they have lower variances and converge faster.

Some of the limitations of this work include not considering the behaviour of the drivers as in realistic situations a driver may follow or not the incentivised route. This can be taken into account by assuming a participation/compliance rate. Furthermore, as shown in Niroumand et al. (2024), incentives on paths outperform incentives on links. None of them were tried in this study as the incentives are along trips, however, it would be insightful to compare results with path incentives.

## References

Albrecht, S. V., Christianos, F., & Schäfer, L. (2024). *Multi-agent reinforcement learning: Foundations and modern approaches.* MIT Press. Retrieved from `https://www.marl-book.com`

Bochenina, K., Taleiko, A., & Ruotsalainen, L. (2023, 06). Simulation-based origin-destination matrix reduction: A case study of helsinki city area. *SUMO Conference Proceedings*, *4*, 1–13. Retrieved from `https://www.tib-op.org/ojs/index.php/scp/article/view/197` doi: 10.52825/scp.v4i.197

Djavadian, S., Hoogendoorn, R. G., Van Arerm, B., & Chow, J. Y. (2014). Empirical evaluation of drivers' route choice behavioral responses to social navigation. *Transportation Research Record*, *2423*(1), 52–60.

Gawron, C. (1998). *Simulation-based traffic assignment. computing user equilibria in large street networks* (Unpublished doctoral dissertation). Universität zu Köln.

Grunitzki, R., & Bazzan, A. L. (2017). Comparing two multiagent reinforcement learning approaches for the traffic assignment problem. In *2017 brazilian conference on intelligent systems (bracis)* (pp. 139–144).

Klein, I., & Ben-Elia, E. (2018). Emergence of cooperative route-choice: A model and experiment of compliance with system-optimal atis. *Transportation research part F: traffic psychology and behaviour*, *59*, 348–364.

Luan, M., Waller, S. T., & Rey, D. (2023). A non-additive path-based reward credit scheme for traffic congestion management. *Transportation Research Part E: Logistics and Transportation Review*, *179*, 103291.

Niroumand, R., Vosough, S., Rinaldi, M., & Roncoli, C. (2024). Beyond links: The power of path incentives in alleviating congestion and emissions in urban networks. In *12th symposium of the european association for research in transportation.*

Niroumand, R., Vosough, S., Roncoli, C., Rinaldi, M., & Connors, R. (2025). Evaluating link and path incentives: Which is the most effective strategy for mitigating traffic congestion? In *Transportation research board 104th annual meeting, washington dc, usa.*

Powell, W. B. (2019). A unified framework for stochastic optimization. *European Journal of Operational Research*, *275*(3), 795-821. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0377221718306192` doi: https://doi.org/10.1016/j.ejor.2018.07.014

Ramos, G. D. O., Rădulescu, R., & Nowé, A. (2019). A budged-balanced tolling scheme for efficient equilibria under heterogeneous preferences. In *Proceedings of the adaptive and learning agents workshop (ala-19) at aamas.*

Shou, Z., Chen, X., Fu, Y., & Di, X. (2022). Multi-agent reinforcement learning for markov routing games: A new modeling paradigm for dynamic traffic assignment. *Transportation Research Part C: Emerging Technologies*, *137*, 103560.

Stewart, K., & Ge, Y.-E. (2014). Optimising time-varying network flows by low-revenue tolling under dynamic user equilibrium. *European Journal of Transport and Infrastructure Research*, *14*(1).

TfL. (2024). *Congestion Charge.* `https://tfl.gov.uk/modes/driving/congestion-charge#:~:text=The%20Congestion%20Charge%20is%20a,Sat%2DSun%20and%20bank%20holidays.` (Accessed: 2024-09-02)

Van Essen, M., Eikenbroek, O., Thomas, T., & Van Berkum, E. (2019). Travelers' compliance with social routing advice: Impacts on road network performance and equity. *IEEE Transactions on Intelligent Transportation Systems*, *21*(3), 1180–1190.

Vosough, S., & Roncoli, C. (2024). Achieving social routing via navigation apps: User acceptance of travel time sacrifice. *Transport Policy*, *148*, 246–256.