

# Reinforcement Learning for Road Pricing: A Review and Future Directions

Otto Vermeulen<sup>1\*</sup>, Arno Siebes<sup>1</sup> and Yannis Velegrakis<sup>1</sup>

<sup>1</sup>Department of Information and Computing Sciences, Utrecht  
University, Princetonplein 5, Utrecht, 3584 CC, The Netherlands.

\*Corresponding author(s). E-mail(s): [o.a.vermeulen@uu.nl](mailto:o.a.vermeulen@uu.nl);  
Contributing authors: [a.p.j.m.siebes@uu.nl](mailto:a.p.j.m.siebes@uu.nl); [i.velegrakis@uu.nl](mailto:i.velegrakis@uu.nl);

## Abstract

Demand for mobility is growing, and traffic on roads has increased substantially, leading to suboptimal traffic flow and congestion. Road pricing can encourage vehicles to change their behavior by charging for road use. Because traffic is not static, dynamic road pricing can help dynamically control traffic. Reinforcement Learning is an effective approach to optimizing the performance of a system. It has already been applied to control traffic signals and has recently found an application in dynamic road pricing for traffic optimization. We survey recent solutions and find that the methods proposed demonstrate the usefulness of reinforcement learning for road pricing. We compared how common challenges in reinforcement learning were approached in the works. Challenges which remain little explored are generalizability and scalability of solution approaches. Approaches to partial observability, credit assignment and non-stationarity are not in all cases taking full account of existing solutions for these common challenges. We further note the need for standardized benchmarks to allow comparisons between the performance of the provided solutions.

**Keywords:** Reinforcement Learning, Pricing and Resource Allocation, Optimization, Simulation, Transportation

## 1 Introduction

As traffic on road networks grows, so too does traffic congestion ([Mokbel et al., 2024](#); [Prieto Curiel et al., 2021](#)). Congestion occurs when vehicles moving on a road network are forced to reduce their speed as a result of the presence of other vehicles. Congestion

047 has significant productivity, environmental, and health consequences, therefore it is  
048 important to be able to manage it.

049 Intelligent Transportation Systems (Haydari and Yilmaz, 2020) (ITS) are systems  
050 that have been designed to optimize traffic by reducing the presence of congestions.  
051 One type of such systems are those imposing constraints to which vehicles must adhere.  
052 Examples include traffic signals, ramp metering, or speed limits. Another type is those  
053 that provide incentives for vehicles to follow a specific behavior. Examples include  
054 route guidance and road pricing.

055 Road pricing is the mechanism to determine the tolls that a vehicle must pay  
056 to use a specific road at a specific time. It is a mechanism that can contribute to  
057 alleviating congestion. This is based on the insight that people tend to make choices  
058 to optimize their utility, taking into account costs and benefits of their actions. In  
059 contrast to the other ITS mechanisms, road pricing uses a price mechanism, with the  
060 benefits of clarity, universality and efficiency (Lindsey and Verhoef, 2001). Its potential  
061 to influence vehicles' route choice and travel mode were advocated early on from an  
062 economic perspective (Beckmann et al., 1956; Knight, 1924; Pigou, 1924; Vickrey,  
063 1969). Farias et al. (2024) provide a recent comprehensive review of dynamic pricing  
064 for toll roads in the U.S. In the literature, there have been variations of this idea  
065 under the term of congestion pricing, congestion charging, tolling, toll charging, or  
066 road pricing. We use the term road pricing. With road pricing, roads with higher tolls  
067 tend to be less preferable to some of vehicle drivers. This means that, by controlling  
068 the tolls, one may indirectly control the flow of vehicles on the road network.

069 Effects of road pricing may be unanticipated if vehicles change their behavior.  
070 Other ITS, like ramp metering and variable speed limits, may also influence the effec-  
071 tiveness of road pricing. Assuming that vehicles change their behavior if roads are  
072 priced and that other mechanisms, like traffic signals and variable speed limits, are at  
073 play in road networks as well, we note that the effects of road pricing may be unan-  
074 ticipated. Because of this, we require a system that is able to learn from the effects of  
075 road pricing decisions and adjust this accordingly. Given the inherent dynamic aspect  
076 of road pricing, it seems ideal for a Reinforcement Learning approach.

077 Reinforcement Learning (RL) is a machine learning technique that learns agents to  
078 take actions to optimize their performance. In RL, agents interact with an environment  
079 that can be in a number of states. These agents want to achieve some objective.  
080 Agents take actions according to a policy, which leads to a change in the state of the  
081 environment. Agents learn by observing the effect of their actions in the environment  
082 through reward signals, to improve the policy to an optimal one (Sutton and Barto,  
083 2018). RL methods can use tables or functions. If these functions are neural networks  
084 we associate this with deep learning and correspondingly use the term Deep RL.

085 Despite its potential, RL has received limited attention within the research com-  
086 munity for road pricing applications. This is a notable gap, as RL presents a strong  
087 candidate for optimizing the tolling process, alongside offering numerous compelling  
088 research challenges. Traditional approaches like fixed tolling, rule-based heuristics or  
089 optimization-based techniques often rely on access to tractable traffic models and  
090 assume stationary demand patterns. However, in reality, traffic systems are dynamic,  
091 stochastic and are subject to incidents and human behavior, which evolves over time  
092

in response to (toll-) pricing but also other policies and control measures. These limitations hinder the effectiveness of conventional methods in practice.

The complex and analytically intractable nature of traffic propagation, influenced by a multitude of factors, makes the model-free learning capabilities of RL particularly well-suited for dynamic tolling optimization. By learning from real-time traffic data, RL agents can develop effective tolling policies that incentivize behavioral changes in road users through price adjustments. RL mechanisms can also adapt online to evolving traffic patterns and non-stationary conditions like incidents or changes in vehicle preferences.

Furthermore, the application of RL to road pricing raises significant research questions concerning scalability, generalization, and off-line learning. The interactions between multiple road pricing agents, diverse road user types (including human-driven and autonomous vehicles), and agents controlling other traffic elements (such as traffic signals, ramp meters, and variable speed limits) warrant in-depth investigation.

While RL has been extensively explored in the context of traffic signal control, its application to road pricing remains relatively underdeveloped. We attribute this to principal differences between the two domains. Traffic signal control typically benefits from more immediate and localized feedback, as well as greater availability of data and standardized interfaces for widely used simulation platforms. In contrast, RL for road pricing involves challenges such as delayed and indirect feedback, stochastic vehicle behavior, limited access to suitable datasets, and a lack of standardized interfaces. These factors have likely contributed to the slower adoption of RL in this domain. However, they also highlight opportunities for impactful research.

By addressing these, and the other challenges noted above, this survey aims to underscore and advance the use of RL in road pricing and stimulate further research in this promising domain.

## 1.1 Recent Related Surveys

Various works in the wider Intelligent Transportation Systems domain exist which use RL to solve traffic problems. Illustrative recent surveys covering such works are listed in Table 1. As motivated in the previous section, road pricing solutions cannot occur in isolation, but ultimately need to take account of other ITS mechanisms. We distinguish therefore two categories: surveys providing an overview of the use of RL in the wider ITS domain and surveys focusing on the use of RL in a specific ITS domain like traffic signal control.

## 1.2 Contribution of this Survey

Our contributions can be summarized as follows:

- We describe a traffic model and give an introduction to road pricing. We provide a high level overview of RL and different algorithms that have been used in the surveyed works.
- We formulate the problem how traffic can be managed by road pricing using RL and identify important specific research challenges. At the best of our knowledge this structured approach is unique.

**Table 1:** Surveys on ITS solutions through Reinforcement Learning on multiple aspects of a road network (top) or on specific aspects of the road network collectively (bottom). A survey for RL in road pricing does not yet exist.

Ref.	Features and relevance
Han et al. (2023)	Provides a comprehensive overview of RL based traffic control strategies, highlighting that few studies separate learning and testing. Identifies challenges for field implementation.
Haydari and Yilmaz (2020)	Categorizes application types, control methods and RL algorithms. Both for Single Agent and Multi Agent RL. Strong focus on traffic signal controls.
Farazi et al. (2021)	Broad overview of RL in the wider transportation domain. Includes mapping of problems to RL methods. Highlights simultaneous decision making of both vehicles and infrastructure.
Yan et al. (2022)	A comprehensive review of RL in logistics and supply chain management, including the Vehicle Routing Problem and urban logistics. Highlights the need to use different agents in complex cases.
Schmidt et al. (2022)	A review of MARL in autonomous mobility in four application areas: traffic control, autonomous vehicles, resource optimization and unmanned aerial vehicles.
Saharan et al. (2020)	Dynamic pricing techniques for ITS, importance of evaluation parameters, limitations of techniques.
Lombardi et al. (2021)	Overview of design, simulation, implementation and evaluation of tolling schemes.
Wei et al. (2021)	Recent advances in the use of RL for traffic signal control, including analysis of environments, experiment settings and evaluation of approaches.
Kumar and Raubal (2021)	Traffic congestion alleviation using Deep Learning. Systematic overview of recurring and non-recurring congestion types including RL solutions.
Qin et al. (2022)	RL in ride sharing business processes. Amongst other for pricing and dynamic routing.
Boggyrbayeva et al. (2024)	Machine Learning to Solve the Vehicle Routing Problems. Covers RL next to DL, heuristics and combinations.

- We analyze how RL is applied to solve the road pricing problem and what objectives are achieved and which challenges are addressed.
- Based on our model and the scope of the various works we identify the most important challenges. We provide additional areas for future work.

### 1.3 Methodology and Structure of the Survey

Research papers were identified through searches for reinforcement learning, tolling, road pricing, congestion pricing and congestion charging within prominent conference proceedings and journals (as indexed by Scopus, Google Scholar, DBLP and similar sources). We excluded works that were not sufficiently specific, such as those using pricing for electric vehicle charging as an incentive for changing user behavior or those considering the relationship with autonomous vehicles. This process yielded 19 works, specifically addressing RL and road pricing, for inclusion in this survey with 17 selected for detailed analysis. We used as source for the different theories the original and seminal papers as they are referenced in traffic and road pricing published works. For reinforcement learning theory, we consulted and included original works and existing surveys. The total number of works included in this survey is 144.

The remainder of this survey is structured as follows. Preliminaries and background are provided in three sections. Section 2 introduces modeling traffic and three traffic models which provide a basis for analyzing the traffic models in the works on RL and road pricing. In Section 3 we introduce road pricing. Section 4 provides a high level introduction to RL and a detailed introduction is provided in Appendix A. Section 5 formulates the traffic congestion problem in a reinforcement learning setup and, in the context of pricing, important research challenges. In the sequel, Section 6 discusses recent works applying reinforcement learning as solutions to the problem defined in Section 5. Section 7 describes how the challenges have been addressed and future research directions for reinforcement learning in the road pricing domain. Section 8 concludes the survey.

## 2 Modeling Traffic

### 2.1 Traffic Model - Network

We assume that we operate on a network (Dafermos and Sparrow, 1969). A network  $G(N, L)$  is a graph where  $N$  is a set of nodes and  $L$  is a set of directed edges between nodes in  $N$ , referred to as links. A link is typically denoted as a pair of nodes, the first being referred to as the start-point and the second as the end-point,  $(s, e)$ , or we may write  $l^{(s,e)}$ , instead. A network models a real life network, where the nodes correspond to real road junctions and the links to the real roads that connect the different junctions. In some cases synthetic networks are used, that are more abstract, to illustrate a concept or idea. A network is commonly governed by a network authority, which can impose a range of centralized and/or decentralized controls. A path is a sequence of links where the end-point of each link (apart from the one of the last) is the start-point of the subsequent link. The start-point of the first link and the end-point of the last are the origin and destination of the path, respectively. A path from a node  $o$  to a node  $d$  is denoted as  $p^{o \rightarrow d}$ , that is,  $p^{o \rightarrow d} = \{l^{(o,x_1)}, l^{(x_1,x_2)}, \dots, l^{(x_{n-1},x_n)}, l^{(x_n,d)}\}$ . Note that there may be more than one path from an origin node  $o$  to a destination node  $d$ . The set of such paths is denoted as  $P^{o \rightarrow d}$ . We denote by  $\mathbb{P}^G$  the set of all possible paths in a network  $G$ .

### 2.2 Traffic Model - Vehicles - State and Demand

Vehicles move on the network. Let  $\mathcal{V}$  be the universe of vehicles. At time  $\underline{t}$  a vehicle  $v$  is in position  $s$  on a link, moving towards its destination node  $o$ . We assume the final destination to be always a node, that is they never end in the middle of a road. The state of a vehicle  $v$  on the network at a specific time  $\underline{t}$  can be modeled as a tuple  $\langle \underline{t}, v, l^{(b,o)}, s, p^{o \rightarrow d} \rangle$ , where  $l^{(b,o)}$  is the link on which the vehicle is found,  $s$  is the position of the vehicle on the link  $l$ , and  $p^{o \rightarrow d}$  is the path that the vehicle intends to follow after it reaches the end-point of the link  $l^{(b,o)}$  in order to arrive to its final destination  $d$ . Although time is a continuous quantity, we consider discrete ordered time points  $\underline{t} = 0, 1, \dots, T$ . Let  $\mathcal{T}$  denote the set of these time points. Note that we use  $\underline{t}$  as  $t$  is reserved for episode steps, to be introduced in Section A.

The set of all the states of all vehicles at all different time points, denoted as  $\mathbb{D}$ , represents the amount of traffic on the network.  $\mathbb{D}$  is referred to as traffic states set. We will use the notation  $\mathbb{D}[cond]$  to select the states that satisfy certain conditions *cond*. For instance,  $\mathbb{D}_{[t \leq 5 \wedge l(b,o)=l]}$  denotes the states of the vehicles before time  $t \leq 5$  on the link  $l$ . If we assume an analysis period of traffic  $[t_a, t_b)$  and consider an origin  $o$  and destination  $d$  we define demand  $D_{od}(t_a, t_b)$  as the number of vehicles willing to travel from  $o$  to  $d$  during that period.

### 2.3 Link and Vehicle Properties

Every link in the network has a set of properties. There are three types of such properties. First, there are properties related to physical characteristics of the link that are naturally invariant. An example is the length  $le$  of the link. Second, there are properties that are invariant as the traffic progresses, but the network authority has the ability to explicitly set them to a different value. Examples of these adaptable properties include the maximum capacity of vehicles allowed on a link, the maximum speed limit  $u_{max}$ , or a possible toll price  $\tau$ . The network authority that controls these values may adapt them when it sees fit, or may have them also change based on some function of time, for example having higher tolls during business hours and much lower tolls after midnight. Finally there are characteristics that depend on the traffic at each specific moment, which means that these characteristics are dynamically changing over time. Illustrative examples are the density  $k$  (vehicles/m), the average speed of vehicles  $u$  (m/s) and the flow  $q$  (vehicles/s). The values of these dynamic properties of the links over time are referred to as a traffic pattern.

A vehicle also has some invariant properties, for example its maximum speed, maximum acceleration/deceleration, and length. Dynamic properties are driving speed, steering direction, and position.

### 2.4 Cost Function and Route and Departure Time Choice

When a vehicle is about to start moving from a point in the network towards a final destination node, it selects a route. A number of different parameters are taken into consideration, like the number of vehicles on the network, their speed, as well as other preferences of the vehicle.

To model this situation, we assume the existence of a cost function which a vehicle wants to minimize. It provides each vehicle a quantitative value, the predicted cost, for each possible route from the current position until the destination. The function takes into consideration the physical characteristics of the network, e.g. the number of links to traverse, the length or travel time at maximum allowed speed of links, in order to opt for shortest routes. It may also consider the dynamic characteristics like the observed traffic in order to opt for less congested links. It is also possible that the function takes the experienced costs into account which is their history of costs experienced on the network. In most cases the elements contributing to the cost are transformed to a single number. Thus, in its most general form the cost function for a vehicle  $v$  is  $c_v: \mathbb{P}^G \times \mathbb{D} \rightarrow \mathbb{R}$ . While this function, consisting of two arguments, is defined on all possible paths  $\mathbb{P}^G$  and all history is kept in  $\mathbb{D}$ , it is in practice always applied

to a small subset. For example, at decision time, a vehicle could assume the cost of traversing a link to be the free flow travel time (time required when no other vehicles are present). Often the formula is assumed the same for the vehicles, but varies in a set of parameters that identify how much the various factors affect the final costs. For instance a multiplication factor converting travel time to money, commonly called value of time (vot) is usually not homogeneous across vehicles. We illustrate the use of the abstract cost function  $c_v$  by relating it to the link cost function (11, Section 6.2.1) from Chen et al. (2018) which provides practical context. Recall the cost function:

$$c_v = \mathbb{P}^G \times \mathbb{D} \rightarrow \mathbb{R}. \quad (1)$$

Assume that vehicle  $v$  chooses a path  $p \in \mathbb{P}^G$  at decision time  $t$ . The path  $p$  is decomposed in its constituent links:  $p = \{l^{(o, x_1)}, l^{(x_1, x_2)}, \dots, l^{(x_{n-1}, x_n)}, l^{(x_n, d)}\}$ .

We now consider the per-link cost (as defined in Section 6.2.1):

$$c_v(l, \underline{t}) = tt(l, \underline{t}) * vot + \tau(l, \underline{t}). \quad (11)$$

This depends on the toll  $\tau(l, \underline{t})$  and  $tt(l, \underline{t})$  which is itself a function of  $\mathbb{D}_{[\underline{t}, l]}$  the state of traffic on that link at that time, which depends on the paths chosen by other vehicles, which were based on their evaluations, based on  $\mathbb{D}$ . We now evaluate the cost for the full path by aggregating the per-link costs:

$$c_v(p, \mathbb{D}) = \sum_{l \in p} c_v(l, t). \quad (2)$$

Substituting the expression for  $c_v(l, \underline{t})$  we get:

$$c_v(p, \mathbb{D}) = \sum_{l \in p} tt(l, \underline{t}) * vot + \tau(l, \underline{t}). \quad (3)$$

Each  $tt(l, \underline{t})$  is estimated based on vehicle speeds on link  $l$  at time  $\underline{t}$ , which recursively depends on the network-wide vehicle distribution on prior decisions of other vehicles. Hence the cost function incorporates, albeit implicitly, the full network traffic pattern encoded in  $\mathbb{D}$ .

The computation of the cost and the planning of the route to follow can take place either once when starting its route and keep the route that was decided, or may be recomputed every time the vehicle is at a node. The choice of the path based on the cost may be deterministic or stochastic. In case of stochastic path choice, a softmax function is commonly used to translate path costs to probabilities:

$$p_{v_i} = \frac{e^{-\theta c_{v_i}}}{\sum_j e^{-\theta c_{v_j}}}, \quad (4)$$

where  $\theta$  is a scaling parameter and  $c_{v_i}$  represents the cost for vehicle  $v$  when taking path  $i$ . In some cases the departure time is not fixed, meaning the vehicle can opt to

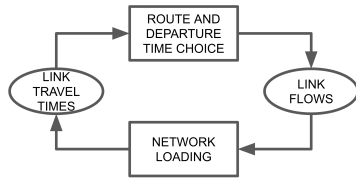
depart earlier or later, or even refrain from departing. We denote this process route and departure time choice as  $\mathcal{R}$ . For a network  $G$  and a set of vehicles  $V$  with a traffic demand set  $\mathbb{D}$  we refer to the tuple  $\langle G, V, \mathbb{D} \rangle$  as a traffic network.

## 2.5 Network Loading

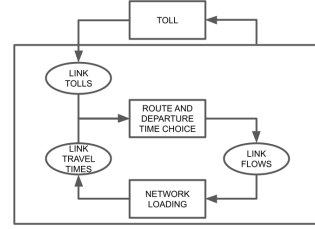
At any new time point, the vehicles that are on the network need to move towards their planned destinations as dictated by their planned paths. This process is known as network loading, and is denoted as  $\mathcal{M}$ . Given route choices for vehicles, network loading provides the resulting route travel times. After a network loading step, the network has new characteristics, and as such, a path refresh phase follows where the planned routes of the different vehicles get recalculated. After the network loading and the planned route refresh phase, the network gets into a new state:

$$\mathbb{D}_{[t=k+1]} = \mathcal{R}(\mathcal{M}(\mathbb{D}_{[t=k]})). \quad (5)$$

See Figure 1 (modified from [Boyles et al. \(2023\)](#)) for an illustration.



**Fig. 1:** Traffic Model without Toll



**Fig. 2:** Traffic Model with Toll

There are different types of model for network loading, to represent the movement of vehicles on a network of which we will describe three types. In a microscopic model, every vehicle is considered an individual entity and its behavior, characteristics and interactions are studied independently. Characteristics like position, speed, and distance from other vehicles are computed for each vehicle individually. In each network loading step, these characteristics are updated for every vehicle. For instance, in the Intelligent Driver Model (IDM) ([Treiber et al., 2000](#)), the acceleration of a vehicle is computed using vehicle and other vehicles' characteristics, leading to its position update. Other microscopic models ([Ahmed et al., 2021](#); [Chowdhury et al., 2000](#); [Diallo et al., 2021](#); [Knorr, 2013](#); [Nguyen et al., 2021](#)) exist but the main principle to calculate vehicle movement remains the same. Traffic simulators which use IDM are among other SUMO ([Lopez et al., 2018](#)), CityFlow ([Zhang et al., 2019](#)) and CBLab ([Liang et al., 2022](#)). Widely used commercial software for microscopic traffic simulation using different models are PTV VISSIM (Wiedemann model), Paramics (Fritzsche model) and AIMSUN (Gipps model) ([Ahmed et al. \(2021\)](#)).

To study traffic in a dynamic macroscopic representation, the Cell Transmission Model (CTM) ([Daganzo, 1994](#)) is often used. A link is divided into a finite number of



sections (called cells), and the traffic is studied by evaluating each cell at every time step, by transferring part of the traffic from each cell to the next. A key distinction from the microscopic model is that this model propagates vehicles not individually but in the aggregate. The CTM has its origins in a hydrodynamic model of traffic flow, the Lighthill-Whitham-Richards (LWR) model (Lighthill and Whitham, 1955; Richards, 1956). It uses three concepts: traffic flow  $q$ , speed  $u$ , density  $k$ , and their relationship on a link, the fundamental diagram,  $q = uk$ . A corresponding node model for CTM facilitates connections, merges and diverges (Daganzo, 1995; Tampère et al., 2011).

A third model is the static macroscopic model. In a static macroscopic model the congestion properties on a link are usually described using a Link Performance Function (LPF). This function expresses the average or steady state travel time on a link as a function of the traffic flow on the link. The most commonly used LPF is from the Bureau of Public Roads (United States Bureau of Public Roads, 1964). Detailed information on dynamic macroscopic models (Chiu et al., 2011) and on static models (Verhoef, 1999) can be found elsewhere. Other works make comparisons between the various models (Hoogendoorn and Bovy, 2001; Storani et al., 2021; van Wageningen-Kessels et al., 2015).

## 2.6 Objectives and Equilibria

One objective of the network authority is to optimize the performance, measured by some metric, of a traffic network. There are different metrics that can be considered for that purpose. One is the total travel time, the total time needed for all the vehicles to arrive at their destination, another the overall distance traveled by all vehicles, yet another the rate in which the vehicles arrive at their intended destination. The metric is calculated by some objective function  $J$ . We assume that the network authority can choose to apply means to influence the traffic. Now assume that every vehicle selects the route that minimizes its cost function  $c_v$ . Over time, traffic is then likely to move to a situation in which no vehicle can reduce its cost by choosing another path (Wardrop, 1952). This situation is called a User Equilibrium or Dynamic User Equilibrium (Chiu et al., 2011) in case of dynamic traffic assignment. It may imply that some links are congested, while others are underutilized.

Assume now that the network authority has the objective of minimizing the total cost and has full control over the vehicle route choices. In that case, diverting some vehicles to other paths, non-optimal for themselves, can help alleviate the congestion and improve the traffic in the network resulting in minimization of the total travel time. This situation in which the overall network traffic is optimized is known as the (Dynamic) Social Optimal equilibrium (Wardrop, 1952). However, in most cases, a network authority has not full control and will need to apply other means to influence vehicles to change their paths.

There are various means by which traffic can be controlled, implying vehicles to follow routes that are not always the most optimal for them. One is through hard restrictions, which divert vehicles towards specific links or restrict their speed (traffic signaling, traffic lights or speed limits). For this to materialize, the network authority needs to have a global view of the network and the ability to implement these controls. Once implemented, these controls can either operate decentralized, cooperatively or

not, or centralized. Another way to affect traffic is to let the final decision be made to individual vehicles, but at the same time providing incentives for them to choose routes that will improve overall network traffic, even if for the individual vehicle the route is not the most beneficial. The charging of a toll for road use is an example of such an incentive.

420

## 421 3 Road Pricing

422

### 423 3.1 Why Road pricing?

424

Road pricing (Tsekeris and Voß, 2009) can be used to manage travel demand, raise revenues for funding transport investment or a combination of both. Just as other demand management tools noted in Section 2, road pricing builds on the assumption that people weigh the cost and benefits of their actions and take their actions to maximize their utility.

429

### 430 3.2 Modeling Road Pricing

431

To illustrate the workings of road pricing we describe first best pricing, time dependent road pricing and second best pricing. First best road pricing is termed first best because it relies on a highly idealized scenario, assuming, among other things, that tolls are charged on all roads with perfectly differentiated prices. It is usually illustrated with a single road (Verhoef, 2002). All vehicles have the same cost function  $c_v(q)$  for this road and the demand function  $d(q)$  is also dependent on  $q$ . Traffic flow, speed and density are independent of time and uniform along the road. The User Equilibrium which will materialize in this case does not take account of the costs vehicles impose on each other. To mitigate this, a Marginal Cost Toll can be derived and applied, leading to a Social Optimal equilibrium. This toll is also known as a Pigouvian tax (Knight, 1924; Lindsey and Verhoef, 2001; Pigou, 1924). A detailed analysis using variations of the cost curve (for example, when it is extended with a backward bending part because of continuous congestion) is given by Verhoef (1999). Furthermore, marginal cost pricing for a single road, can also be applied to a general road network (Yang and Huang, 1998). We do not cover these cases here.

In a time dependent model, the time-independent model is extended with a time-dependent travel demand function and specifies how the flow changes over space and time. Demand is assumed fixed (price independent), but vehicles are heterogeneous with respect to their trip-timing as well as value of time. The private cost of a trip is (Arnott et al., 1993):

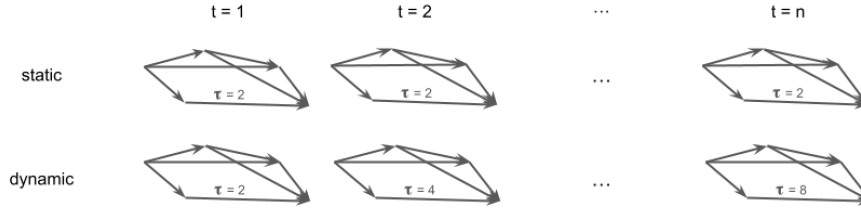
447

$$448 \quad c_v(\underline{t}) = \alpha T(\underline{t}) + \beta(\text{time early}) + \gamma(\text{time late}), \quad (6)$$

454

where  $\alpha$  is the *vot*,  $\beta$  is the unit cost for arriving early and  $\gamma$  is the unit cost for arriving late. The cost of arriving early or late is referred to as schedule delay cost and  $T$  as travel time cost. We note the desired arrival time as  $\underline{t}^*$ . Time early is  $\max[0, \underline{t}^* - \underline{t} - T(\underline{t})]$  and time late is  $\max[0, \underline{t} + T(\underline{t}) - \underline{t}^*]$ . Under certain conditions, vehicles cannot overtake each other and tolls can vary freely over time (Arnott et al.,

460



**Fig. 3:** Static and dynamic tolling of a single road

1993; Lindsey and Verhoef, 2001), a User Equilibrium and Social Optimal equilibrium can be calculated and tolls can be applied to achieve this Social Optimal equilibrium. More variants of this model exist (Lindsey and Verhoef, 2001). First best pricing and time dependent models are often not feasible given practical constraints (not all roads can be tolled), or information limitations (incomplete information for vehicles).

Second best pricing methods (Arnott, 2007; Ekström, 2014; Tsekeris and Voß, 2009; Verhoef, 2000, 2002), are more realistic, as various constraints are taken into account. Examples of second best pricing methods are pay lanes and cordon pricing in contrast to pricing every road, the use of step tolls as opposed to smoothly time varying tolls and a fixed daily tolling schedule rather than a schedule dependent on traffic conditions during the day. In relation to this, we also distinguish static and dynamic road pricing (Fig. 3). In static road pricing, toll prices are determined by analysis of historical and economic data, without taking into account current traffic. In dynamic road pricing, tolls are determined by also taking into account current or anticipated congestion levels, and hence can vary, based on the number of vehicles, their speed, and their location on the network (Cole et al., 2003; Como and Maggistro, 2021; Eliasson, 2017; Genser and Kouvelas, 2019, 2022; Maheshwari et al., 2024; Nohekhan et al., 2021; Paccagnan et al., 2021; Pandey and Boyles, 2019; Vickrey, 1963).

Dynamic road pricing is a challenging task, as driver behavior on the network needs to be continuously monitored and tolls need to be adjusted, with each having an immediate effect on the other. Determining and setting prices is performed by toll agents

### 3.3 Toll Agents

Toll agents ( $ta$ ) are charged with observing traffic on the network and determining tolls. A model representing a traffic network with a toll agent is given in Figure 2 (modified from (Boyles et al., 2023)). The figure illustrates that toll agents obtain information from the traffic system as a basis to determine the tolls to charge for link use to vehicles. Toll agents can be given objectives by the network authority, for example with respect to the flow on a link. If provided with a license the toll agents determine their own objective within the constraints of the license. For example to maximize profit on a tolled link in which they invested.

We detail baseline objectives for toll agents. We use the following equation:

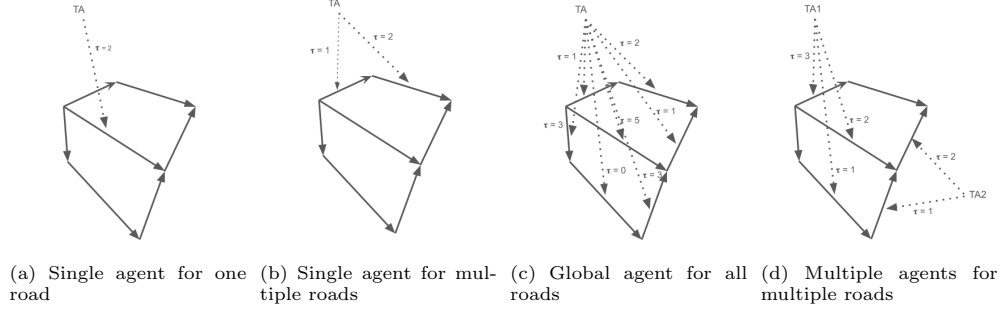
$$J = \sum_{v \in \mathcal{V}} (\phi * tt(v) + \psi * \sum_{l' \in p(v)} \tau_{l'}). \quad (7)$$

Here  $\phi$  is a scaling parameter, often representing the *vot*,  $tt(v)$  the experienced travel time of the vehicle,  $p(v)$  is the path which the vehicle has traversed,  $l'$  a tolled link and  $\tau_{l'}$  the toll incurred when traversing that link. The baseline objectives are given in Table 2. When considering Social Welfare, tolls are not taken into account as tolls are considered transfer payments which remain internal to society (Mirzaei et al. (2018a)). If the average is taken over the number of vehicles  $|\mathcal{V}|$  we use ASTT and ASC.

**Table 2:** Baseline objectives for toll agents

	$\phi$	$\psi$
Total System Travel Time (TSTT)	1	0
Total System Cost (TSC)	<i>vot</i>	1
Social Welfare (SW)	<i>vot</i>	0
Maximum Revenue (MR)	0	1

On a network, one (global) or more (local) toll agents can operate (Fig. 4). In the latter case each one covering a part of the network. Toll agents can work together to achieve the same objective or pursue their own objectives, leading to different levels of cooperation or competition. When starting to set tolls, vehicles may take other



**Fig. 4:** Example toll agent configurations

decisions with respect to their mode choice, departure time and route choice. This causes different traffic patterns which the toll agents may not have anticipated causing them to either adjust their prices or their tolling model.

In Section 3.2 we described how a toll is determined in the case of a first best pricing situation on a single link. To determine tolls for road pricing on networks, different methods and technologies have been proposed (De Palma and Lindsey, 2011)

using various coverage and toll differentiation mechanisms. [Joksimovic et al. \(2005\)](#) identifies a mechanism for optimal toll design in a dynamic network with route and departure time choice. Dynamic pricing for ITS surveyed in [\(Saharan et al., 2020\)](#) describe among other dynamic programming, evolutionary optimization, swarm optimization and game theory based techniques. A survey work [\(Lombardi et al., 2021\)](#) on model-based dynamic toll pricing presents an overview of methods for price definition, simulation techniques and technology applications. Reinforcement Learning, which we cover next, is another mechanism to determine toll prices for toll agents.

## 4 Reinforcement Learning

Reinforcement learning (RL) is a computational framework wherein an agent learns to take optimal decisions through direct interaction with an environment. The core interactive loop consists of the agent observing the environment’s state, selecting an action and receiving a reward signal. For the tolling problem, an action is, e.g., raising or lowering the toll on a given road, while the reward could, e.g., be (inversely) related to the time it takes vehicles to go from origin to destination. The agent’s goal is to develop a policy (a mapping from states to actions) that maximizes the long-term cumulative reward.

The standard formalism for this sequential decision-making problem is the Markov Decision Process (MDP), formally defined as a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$  that represents the set  $\mathcal{S}$  of all states  $s$ , the set  $\mathcal{A}$  of actions  $a$ , a transition probability function for the environment  $\mathcal{P}$ , and a reward function  $\mathcal{R}$ .

The central challenge in RL is to find an optimal policy  $\pi_*$  by learning value functions ( $V(s)$  or  $Q(s, a)$ ) which estimate expected cumulative rewards, or by directly optimizing the policy itself. Solutions to MDPs are based on the Bellman equations, which provide recursive relationships for value functions under a given policy [\(Sutton and Barto, 2018\)](#). In road pricing contexts specifically, RL enables the modeling of adaptive tolling strategies in dynamic, uncertain traffic environments. For a detailed treatment of these foundational concepts and algorithms, readers are directed to [Appendix A](#).

RL algorithms are typically categorized into three main categories, all three have been used for the tolling problem. Value-based methods (e.g., Q-learning), which focus on learning an accurate value function and are well-suited to problems with discrete action spaces. In contrast, policy-based methods (e.g., policy gradient), directly search a policy which make them suitable for handling continuous action spaces and learning stochastic policies. Actor-Critic methods, which represent a powerful hybrid architecture where an actor (policy) learns which actions to take, and the critic (value function) evaluates the quality of that action. Many state-of-the-art algorithms are based on this structure.

We further distinguish between model-free and model-based approaches. Model-free methods learn directly from experience without ever creating an explicit model of the environment’s dynamics ( $\mathcal{P}$  and  $\mathcal{R}$ ). In contrast, model-based methods either use a known model of the environment or first learn a model from experience and then use this model for planning and policy optimization.

599 While traditional tabular methods for the above approaches were feasible in  
600 domains with (relatively) small state and action spaces, they failed to scale to larger,  
601 complex problems. Deep Reinforcement Learning overcomes this limitation. Deep  
602 RL leverages deep neural networks as universal function approximators to represent  
603 policies and/or value functions. The integration of deep learning with reinforcement  
604 learning has enabled the field to solve larger problems with high-dimensional state- and  
605 action spaces. This makes it an essential tool to address modern real-world challenges  
606 such as the tolling problem analyzed in this survey.

607

## 608 5 Problem Formulation and Research Challenges

609

610 We will now formulate our problem statement and frame this in reinforcement learn-  
611 ing terminology.

612

613 **Problem.** *Given  $\langle G, V, \mathbb{D} \rangle$ , a horizon  $H$  and an objective  $J$  find a policy  $\pi$  for a*  
614 *Toll Agent  $ta$  to set toll prices on the links  $l$  so that  $J$  is optimized.*

615

### 616 5.1 The Road Pricing Problem in RL Terminology

617

618 The road pricing problem can be formalized as:

- 619 • *Environment:* the Traffic Network  $\langle G, V, \mathbb{D} \rangle$ .
- 620 • *Agent:* the Toll Agent,  $ta$ .

621

622 In this road pricing problem the  $ta$  has an objective  $J$  (see Table 2) that it needs to  
623 achieve.

624

The road pricing problem is further formalized as a MDP:

- 625 •  $\mathcal{S}$ : is the set of all states  $s$ . This is a representation of the current situation in the  
626 environment: speed, travel times, flow, density on the links and so on. It may also  
627 include information on nodes and information from other artefacts like speed limits,  
628 ramp metering and traffic lights. State information is collected by the agent. If it con-  
629 sists of a single number, without applying some function, we call it direct, otherwise  
630 composite. The state  $s$  or observation  $o$  could also be subject to preprocessing. Pre-  
631 processing can involve deep learning mechanisms like Graph Convolutional Networks  
632 (GCN) (Kipf and Welling, 2017) or transformers (Vaswani et al., 2017).
- 633 •  $\mathcal{A}$ : is the set of all actions  $a$ . On the basis of state  $s$  the agent takes an action  $a$ .  
634 This action  $a$  consists of setting new toll prices for one or more links in the Traffic  
635 Network respectively. The actions, can either be the toll price chosen from a discrete  
636 set of prices, from a certain price range, or adjustments to the toll price for one or  
637 more links. These adjustments are made in accordance with the  $ta$ 's policy  $\pi$ .
- 638 •  $\mathcal{P}$ : is the transition probability function mapping from any state  $s$  to a next state  
639  $s'$  after taking action  $a$ :  $\mathcal{P}(s'|s, a)$ . If a toll price is set, vehicles  $v$  will act on this  
640 price (and traffic conditions), propagating the traffic flow by network loading (5).  
641 This will lead to a new state. It is unlikely that  $\mathcal{P}$  is known in a traffic network so  
642 the agent needs to learn from experience.

643

644

- $\mathcal{R}$ : the reward function provides a reward from the environment to the agent. The reward signal could for example be the number of vehicles arriving periodically at a certain node, the average speed on a link, the density on a link (direct) or some combination of other values (composite). As the reward, and therefore return, are the signals to evaluate for the  $ta$  if its actions contribute in a positive or negative way to its objective  $J$ , it is therefore crucial that the reward signal is properly defined. Reward signals are combined to produce the return, using a discount factor  $\gamma \in [0, 1]$ , averaging or other mechanism. As the objective  $J$  can not always be observed by the  $ta$ , it needs to design its reward structure so that its return can act as a proxy for  $J$ .

The  $ta$  has a stochastic policy  $\pi(a|s)$  or a deterministic policy  $\pi(s)$ . Based on the state and this policy, the  $ta$  sets the toll price or toll price adjustment on one or more links. In case of more than one  $ta$  the problem needs to be generalized to a Markov Game (Section A.3) and, when other elements of the environment (vehicles, traffic signals) become learning agents, as well.

## 5.2 Research challenges

Common challenges in (Multi Agent) Reinforcement Learning as surveyed in (Albrecht et al., 2024; Du and Ding, 2021; Dulac-Arnold et al., 2021; Gronauer and Diepold, 2022; Patterson et al., 2024; Wong et al., 2022; Yuan et al., 2023) are outlined below with a specific focus on the implications for the road pricing problem.

Partial Observability Partial Observability (Section A.1), can occur if a toll agent has only information on the links for which it needs to provide the toll prices but not on all links of the network. It may also be that state information is distorted because of sensor failures, or simply because information is not available to the toll agent.

Credit Assignment Credit Assignment (Section A.2) in road pricing may be challenging due to spatio-temporal dependencies in the environment and actions of toll agents. Depending on the route and departure time model of vehicles and the resulting traffic dynamics, it can be difficult to specifically identify which action contributed to which (part of the) reward. Note that credit assignment in single agent settings is already challenging due to spatio-temporal dependencies; this is compounded in situations where multiple agents act simultaneously.

Non-stationarity Non-stationarity (Section A.3) may occur if multiple toll agents are learning and updating their policies simultaneously. When the state of the environment is dependent on the joint actions (toll prices) of all toll agents, the toll agents need to adapt to these new policies, which violates the MDP assumption and convergence may no longer be guaranteed (Wong et al., 2022).

Scalability When the problem is extended from a single agent to a multi agent setting this is accompanied by an expansion of the state and action dimension, which may lead to an exponential rise in the joint action dimension (Du and Ding, 2021; Yuan et al., 2023). Solutions with more agents may also increase the credit-assignment and non-stationarity challenges (Albrecht et al., 2024).



691 Generalizability This challenge addresses whether trained toll agents can perform  
692 well in new and unseen situations. In the road pricing problem, this would encom-  
693 pass among other different traffic demand patterns, changes in network topology and  
694 disturbances in data provided (for example inaccurate measurements, sensor noise,  
695 missing data).

696 In addition to the challenges outlined above, the definition of the state, action and  
697 the identification of a reward function is part of the research in road pricing problems  
698 as well. This contrasts to benchmark problems (Bettini et al., 2024; Bellemare et al.,  
699 2013; Sutton and Barto, 2018; Todorov et al., 2012) where these are explicitly defined.

700 The surveyed works in the next Section (Section 6) detail how these challenges are  
701 addressed in the road pricing problem and are summarized in Section 7.

## 703 6 Road Pricing Approaches using Reinforcement 704 Learning 705

706 In this section a detailed analysis is provided of the ways in which the problem and  
707 challenges, formulated in Section 5 has been addressed. The works are clustered in sub-  
708 sections by problem domain and RL solution approach. A conscious attempt has been  
709 made to align the notation between the various works and the notation introduced in  
710 Section 2 and Section A. The traffic model (Section 2) is the basis for describing the  
711 most relevant characteristics, specific for the works. The description of the toll mecha-  
712 nism and application of RL builds on the background information (Sections 3 and A).  
713 Additional detail is provided where necessary. The most important characteristics of  
714 the works are listed in Table 3.

715 We emphasize that in the models we both have episode steps  $t \in 0, 1, \dots, H$  and  
716 time steps  $\underline{t} \in 0, 1, \dots, T$  where there are  $\Gamma$  time steps in an episode step. This implies  
717 that if a vehicle uses the toll price of a link ( $\tau_l(\underline{t})$ ) it uses the toll price set by the  
718 agent at the latest episode step  $t$  prior to  $\underline{t}$ . Time can be indicated both by  $t$  and  $\underline{t}$ .

### 720 6.1 Policy Gradient for Enhanced $\Delta$ -Tolling in Road Networks

#### 722 6.1.1 Enhanced $\Delta$ -Tolling with Finite Difference Policy Gradient

723  $\Delta$  (Delta)-tolling (Sharon et al., 2016, 2017)(Sharon et al., 2017) dynamically tolls  
724 all links in a network optimizing Social Welfare using fixed parameters. In Enhanced  
725  $\Delta$ -tolling (Mirzaei et al., 2018a,b) ( $E\Delta$ -tolling) the objective is to optimize Social  
726 Welfare (See Table 2) by optimizing the parameters of  $\Delta$ -tolling. They therefore make  
727 the strong assumption that for each vehicle the  $vot$  is known. This is achieved by using  
728 a Finite Difference Policy Gradient Reinforcement Learning (FD PGRL) (Kohl and  
729 Stone, 2004) algorithm to learn the parameters that determine the toll.

731 Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (c_v)$  For each link in a network the toll is  
732 calculated in two steps: first,  $\Delta$  is calculated every time step  $\underline{t}$  based on the difference  
733 between the travel time on a link  $tt(l, \underline{t})$  and the free flow travel time  $T(l)$  on the same  
734 link:

$$\Delta(\underline{t}) = tt(l, \underline{t}) - T(l). \quad (8)$$



**Table 3:** Overview of important characteristics in surveyed works

	6.1.1 (Mirzaei et al. (2018b))	6.2.1 (Chen et al. (2018))	6.2.2 (Qiu et al. (2019))	6.2.3 (Jin et al. (2021))	6.2.4 (Wang et al. (2022))	6.2.5 (Lu et al. (2024))	6.2.6 (He et al. (2024))	6.3.1 (Zhu and Ukusuri (2015))	6.3.2 (Pandey and Boyles (2018))	6.4.1 (Pandey et al. (2020))	6.4.2 (Zhang et al. (2023))	6.5.1 (Sato et al. (2021))	6.5.2 (Sato et al. (2022))	6.6.1 (Tavares and Bazzan (2014))	6.6.2 (Chakravarty et al. (2024))	6.7.1 (Ramos et al. (2018))	6.7.2 (Ramos et al. (2020))
<b>Toll agent objective</b>																	
Minimal Total System Travel Time																	
Minimal Average System Travel Time																	
Maximal Social Welfare	✓																
Minimal Total System Cost																	
Maximum Flow		✓															
Minimal Delay			✓	✓													
Maximum Revenue																	
<b>Vehicle departure time and route choice</b>																	
Departure time																	
Deterministic route choice	✓																
Stochastic route choice		✓	✓	✓													
Dynamic route choice	✓	✓	✓	✓													
Departure time and route choice																	
<b>Reinforcement learning approach</b>																	
Value based	✓																
Policy based		✓	✓	✓	✓	✓	✓										
Actor-critic																	
<b>Architecture if multi agent</b>																	
DTDE																	
CTCE																	
CTDE			✓	✓	✓	✓											
<b>State or observation</b>																	
Direct	✓																
Composite		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓				
Preprocessed			✓		✓	✓											
<b>Action</b>																	
Toll price road		✓	✓	✓													
Toll price road adjustment	✓																
Toll price route					✓												
Toll price region							✓										
<b>Reward</b>																	
Direct		✓	✓	✓			✓										
Composite					✓	✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

After that,  $\Delta$  is multiplied by a proportionality factor  $\beta$  so the product is proportional to the difference between travel time and free flow travel time. Then this and the previous toll  $\tau$  are smoothed with a weighting factor  $R$ :

$$\tau(l, t+1) = R(\beta\Delta(t)) + (1-R)\tau(l, t). \quad (9)$$

The calculated toll is then set to the link in real time. The corresponding cost function is:

$$c_v(l, \underline{t}) = tt(l, \underline{t}) * vot + \tau(l, \underline{t}). \quad (10)$$

( $\mathcal{R}$ ) Route choice: A vehicle optimizes its route dynamically at every node for minimal cost, based on actual travel times and actual tolls. ( $\mathcal{M}$ ) Network loading is done using the CTM.

RL set-up: In a FD PGRL algorithm the parameters  $\theta$  of the policy  $\pi$  are optimized with gradients as usual. The negative total cost (10) in an episode, exclusive of toll  $c_v(l, \underline{t}) = tt(l, \underline{t}) * vot$ , is used as the reward signal. However, these gradients are calculated after performing an episode with a number of policies with slightly perturbed parameters (Peters and Schaal, 2006) which lead to different performance. Using the differences in the parameters and the differences in performance, the gradients are estimated. Here,  $\beta$  and  $R$  are optimized by a FD PGRL-algorithm. The specific version used here (Kohl and Stone, 2004) updates the policy parameters with an adjustment vector. The variants in this work consist of different Local/Global combinations of  $\beta$  and  $R$  in the network. Local means every link has its own parameter and Global means all links have the same parameter. Variants are: Global  $R$ , Local  $\beta$  ( $E\Delta$ -tolling $_{\beta}$ ), Local  $R$ , Global  $\beta$  ( $E\Delta$ -tolling $_R$ ), and Local  $R$ , Local  $\beta$  ( $E\Delta$ -tolling $_{\beta, R}$ ).

Results: Experiments are conducted on various networks. All variants show improvements up to 45% over no tolls.  $E\Delta$ -tolling $_{\beta}$  underperforms the other variants and performs only slightly better than regular  $\Delta$ -tolling. Another drawback is that  $E\Delta$ -tolling is limited in its convergence rate; it is optimized for a specific traffic pattern and might be too slow to learn a new traffic pattern.

## 6.2 Actor-Critic Dynamic Toll Collection for Road Networks

### 6.2.1 Dynamic Electronic Toll Collection - PG- $\beta$

A Dynamic Electronic Toll Collection (DyETC) system in a network, consisting of tolled and non-tolled links, uses Policy Gradient with actor-critic (Chen et al., 2018). The objective is to optimize traffic flow by dynamically adjusting tolls. The policy is approximated by a parameterized function and optimized by a modified Policy Gradient (PG)-algorithm.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ : ( $c_v$ ) Link costs are based on travel time on a link and toll. The travel time is calculated with a LPF. The toll  $\tau(l, \underline{t})$  in this case is equal to the action  $A_t^l$ :

$$c_v(l, \underline{t}) = tt(l, \underline{t}) * vot + \tau(l, \underline{t}). \quad (11)$$

Route costs consist of the sum of its link costs. ( $\mathcal{R}$ ) Route choice of vehicles is based on a Stochastic User Equilibrium and route choice for vehicles is updated every time step. ( $\mathcal{M}$ ) Network loading is done using a modification of a macroscopic static model.

RL set-up: A new algorithm, Policy Gradient- $\beta$  (PG- $\beta$ ) with a separate critic is proposed, with a continuous action space. This improves traditional PG algorithms by using: a) time-dependent value and policy functions, b) use of the  $\beta$ -distribution as an alternative to the normal distribution to obtain a bounded action space, and c) state

abstraction: an assumption that vehicles on the same link have almost equal effects on tolls. The problem is formulated as a time-dependent Markov Decision Process on the traffic network, with time horizon  $H$ , and discrete episode steps  $t = 0, 1, \dots, H$ . maintaining time-dependent value and policy functions. State is defined as  $S_t^{lf}$ , the number of vehicles on link  $l$  with destination node  $f$  at time  $t$ . The state vector of a link as  $\mathbf{S}_t^l$  and state matrix of the network  $G$  as  $\mathbf{S}_t$ . An action  $\mathbf{A}_t$  is defined as the toll vector which consists of the tolls  $A_t^l$  applied to link  $l$ . Once an action is applied at time  $t$ , the number of vehicles on each link is updated:

$$S_{t+1}^{lf} = S_t^{lf} - S_t^{lf,out} + S_t^{lf,in}. \quad (12)$$

The immediate reward function  $R_t(\mathbf{S}_t)$  is defined as the number of vehicles that arrived at their destination during episode step  $t$ . Time-dependent policy and value functions imply that there is a set of policies and a set of value functions; for each episode step  $t$  a policy and value function. The value function at episode step  $t$  is defined as:

$$v_t(\mathbf{S}_t) = \sum_{t'=t}^H \gamma^{t'-t} R_{t'}(\mathbf{S}_{t'}). \quad (13)$$

The optimal policy for the toll agent at episode step  $t$  is:

$$\pi_{*,t}(\mathbf{A}_t|\mathbf{S}_t) = \arg \max_{\pi_t} v_t(\mathbf{S}_t). \quad (14)$$

A policy function approximator parameterized by  $\boldsymbol{\theta}_t, \pi_t(\mathbf{A}_t|\mathbf{S}_t, \boldsymbol{\theta}_t)$  and a value function approximator parameterized by  $\boldsymbol{\vartheta}_t, \hat{v}_t(\mathbf{S}_t|\boldsymbol{\vartheta}_t)$  are used in an actor-critic algorithm to train the policies. Different from the regular actor-critic is that each episode step has its own policy and value function for which the parameters are learned. Secondly, where usually in PG algorithms a Normal probability density function (pdf) is used for MDPs with continuous action spaces, here they use a  $\beta$ -pdf to adapt to a bounded action space.

Results: Experiments are on a synthetic traffic network and a real network. PG- $\beta$  outperforms algorithms where the episode step is incorporated in the state or not considered at all. A similar algorithm using the Normal distribution does not learn an effective policy. Using PG- $\beta$  with state abstraction is 75% more efficient than without. PG- $\beta$ -abs is then compared in five scenarios with four other tolling mechanisms: Fixed toll: proportional to average demand, DyState: dynamic toll proportional to state,  $\Delta$ -tolling (see Section 6.1.1) and no toll. In all scenarios PG- $\beta$ -abs outperforms the other tolling schemes. In a second experiment on a real network, PG- $\beta$ -abs outperforms all other mechanisms, and compared to the Fixed tolling scheme, which was second best, the realized flow was 8% higher and TSTT 14.6% lower. Scalability to larger networks is identified in the work as a challenge.

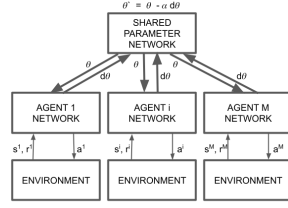
### 6.2.2 Dynamic Electronic Toll Collection - MARL

An extension of DyETC (Chen et al., 2018) (Section 6.2.1) with MARL, to enable larger networks and overcome the scalability problems, led to DyETC-MARL using

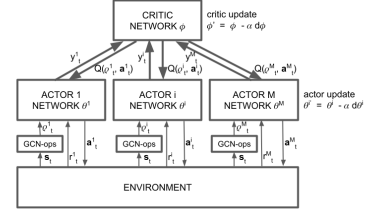
edge-based Graph Convolutional Networks (GCN) (Qiu et al., 2019). The work proposes the use of an edge-based GCN to extract spatio-temporal correlations of the network status and to apply cooperative MARL where each toll agent serves a part of the network. The objective of the toll agents is to maximize the total number of vehicles arriving at their destination. This means that each agent gets the same reward. The agents share a common value network.

Traffic model set-up: Similar to Section 6.2.1. In addition, the network  $G$  is split in partitions, that each cover one or more links.

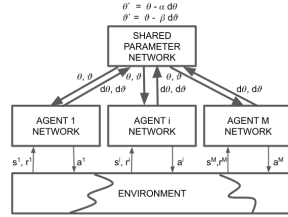
RL set-up: The RL problem is formulated as in Section 6.2.1. The optimal policy for each agent  $ta^i, i \in 1, \dots, N$  covers its  $(ta^i)$  partition of the network. The agents adapt the CTDE framework using an actor-critic method for learning. Figure 6 presents the high-level architecture of the proposed system. Each  $ta$  learns its own policy (actor),



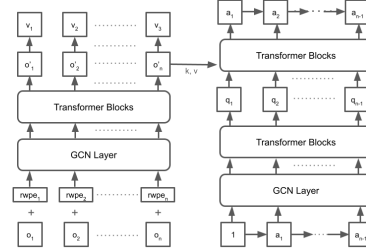
**Fig. 5:** Asynchronous updating of parameters by multiple agents (Mnih et al., 2016)



**Fig. 6:** DyETC MARL with GCN (Qiu et al., 2019)



**Fig. 7:** DADO (Jin et al., 2021)



**Fig. 8:** MAGT-toll (Lu et al., 2024)

and together they use the same value network (critic). By sharing the critic and getting the same rewards, the agents have a common goal; this implies a natural way of cooperating between them. Because of the finite horizon of the MDP, the parameters of the actor and critic are time-dependent. The captured state information is preprocessed using a GCN (Kipf and Welling, 2017). With a GCN the graph nature of  $G$  is combined with the state information of each link to generate a richer input to each agent. The GCN's output  $\rho_i$ , is specific for each  $ta$ , using information from the

partition of that toll agent. So agent  $ta^i$  obtains  $\rho_i$  and generates the action for every link in its partition.

Results: Experiments are conducted on two real networks. MARL-eGCN, DPG- $\beta$ -eGCN (not using the multi-agent approach but PG- $\beta$  with a neural network) and PG- $\beta$  are compared with each other on both networks. MARL-eGCN and DPG- $\beta$ -eGCN, significantly outperform PG- $\beta$ . In addition MARL-eGCN and DPG- $\beta$ -eGCN converge both faster than PG- $\beta$ .

In the next set of experiments, various ablation studies and robustness tests were conducted, among other different traffic settings and different pricing schemes. An important finding in this work is that performance decreases when more than 8 agents are used, and best performance is achieved with 6 agents. The use of more agents suffers from feature redundancy and more coordination costs. Fewer agents (than 6) cannot leverage the power of multi-agent learning (Qiu et al., 2019).

The solution demonstrates that decomposing the state and action spaces per agent and cooperation between agents helps solve the DyETC problem. Preprocessing state with a GCN exploits correlations between the links in the network.

### 6.2.3 Dynamic and Deadline Oriented Road Pricing Mechanism

DADO, considering vehicles' time requirements takes deadlines into account (Jin et al., 2021). For vehicles with a deadline, route cost is appreciated differently if the time interval  $x$  between the calculated arrival time and the deadline is below or above a threshold value  $D$ . The objective for the agents is to optimize the number of vehicles arriving prior to their deadline, using different reward functions.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (c_v)$  For vehicles with  $x > D$  remaining time is abundant and their cost is based on road tolls and the remaining time, increasing slowly if  $x$  decreases. If  $x < D$  the time cost will increase dramatically. For vehicles without a deadline the regular way of calculating the cost per link is used (Equation (11)). Each vehicle is provided with a deadline and updates its route at every junction based on Stochastic User Equilibrium.  $(\mathcal{R})$  and  $(\mathcal{M})$  network loading are similar to the DyETC model (Section (6.2.1)) The objective of toll agents is to maximize the number of vehicles that arrive at the destination before the deadline.

RL set-up: State is defined as  $S_t^{lfd}$  where the deadline  $d$  is an additional index to the state defined in DyETC; for vehicles with no deadline,  $d = 0$ . The index  $d$  is also used in calculating the fraction of traffic flow in the Stochastic User Equilibrium calculation, as well as in the state transition calculation. Actions (tolls) are as defined in DyETC. Three different reward functions  $R$  are defined: R1 to maximize the number of vehicles arriving at their destination before the deadline, R2 to minimize the number of vehicles arriving after the deadline and R3 to minimize the time gap between time of arrival if beyond the deadline and the deadline. The value function is defined as:

$$v_t(\mathbf{S}_t) = \sum_{t'=t}^H \gamma^{t'-t} r_{t'}(\mathbf{S}_{t'}). \quad (15)$$

The policy is as in DyETC time-dependent:  $\pi_t^i(\mathbf{A}_t|\mathbf{S}_t, \boldsymbol{\theta})$ . The optimal policy by training the agent maximizes the value function, which is determined by the critic,

parameterized by  $\boldsymbol{\vartheta}$ ,  $v(\mathbf{S}_t, \boldsymbol{\vartheta})$ . The multi-agent actor-critic algorithm, using the CTDE concept (Figure 7), is based on three key components: a) The local agents for each tolled link determine the toll prices by the local actor and collect the rewards; b) After completing an episode, the local agent determines the accumulated gradients for the local actor ( $\boldsymbol{\theta}$ ) and critic ( $\boldsymbol{\vartheta}$ ); and c) A global actor-critic is updated asynchronously by the accumulated gradients of the local agents, whose parameters are pulled by the local agents after all updates by local agents are complete.

Results: Experiments are performed on a real network. The results are compared with: a) fixed tolling (fee proportional to average demand); b)  $\Delta$ -tolling (Sharon et al., 2017) (Section 6.1.1) and c) DyETC (Chen et al., 2018) (Section 6.2.1). When evaluated against the objective of the reward functions R1, R2 and R3, DADO outperforms the other tolling mechanisms. When evaluated without taking account of deadlines, DADO slightly underperforms DyETC. When reward function R1 is used and, respectively the initial state on the links, initial demand, cost sensitivity and maximum toll price are varied, DADO and DyETC outperform fixed tolling and  $\Delta$ -tolling and DADO is still optimal. Fixed tolling and  $\Delta$ -tolling have problems with performance under changing traffic conditions.

#### 6.2.4 Toll Pricing with Attention Network and Soft Actor-Critic

Another variation (Cooperative Tolling with Reinforcement Learning, CTRL) uses an attention mechanism and SAC (Section A.5.2) to determine toll prices for routes (Wang et al., 2022). Attention (Vaswani et al., 2017) allows modeling of dependencies without regard to their distance in input or output sequences, in this case the information contained in the route between upstream and downstream links. The quality score of a route (Q-value) is also derived using an attention network. This aims to produce comparable tolls for routes taking both state and action (of other routes) into account.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (c_v)$  The price (toll) of a route is determined by the toll agent. ( $\mathcal{R}$ ) Route choice: For every OD pair three routes are predefined and route choice by vehicles is performed at the start of their route based on lowest toll. Vehicles can not change their route after departure. ( $\mathcal{M}$ ) Network loading is done with a microscopic simulation model. The objective of the toll agent is to minimize the Average System Travel Time for all vehicles that arrive at their destinations.

RL set-up: First the delays are calculated. Delays are defined as a scaled difference between actual travel time and free flow travel time on all links, denoted as  $\mathbf{d}_l^t$ . State  $S_t$  for route  $p^{o \rightarrow d}$  is defined as the output of the attention mechanism denoted by  $Agg$  and based on the concatenated delays of the links comprising that route:

$$S_t^{p^{o \rightarrow d}} = Agg(\mathbf{d}_l^t), l \in p^{o \rightarrow d}. \quad (16)$$

The action  $A_t^{p^{o \rightarrow d}}$  is the route price. Reward per link is defined as the cumulative distance covered by vehicles on a link during a time step. Reward per route is defined as the average reward of all links per route. SAC (Haarnoja et al., 2018) is the algorithm to train the agent determining the route prices. In this set-up the actor network calculates the action (toll price) for each path of an OD-pair. The Q values for each

state action pair are calculated by the critic network. Part of this critic network is an attention mechanism to ensure that Q-values are not only based on the state/action per route but also that other routes are taken into account. The parameters of the critic network are updated, after which the policy is updated.

Results: In experiments the set-up is tested on three real networks. Experiments are conducted to compare with other tolling mechanisms: a) No Change: no toll; b) Random: random route choice by vehicles; c) Formula: formula based toll price depending on vehicles on link; d)  $\Delta$ -tolling:  $\Delta$ -tolling (Sharon et al., 2017) (see: 6.1.1); e) Indi-SAC: Soft-Actor-Critic per link; f) Share-SAC: shared Soft-Actor-Critic; and g) MARL-eGCN: MARL-eGCN with one partition (Qiu et al., 2019) (see: 6.2.2 hereafter). The comparison is on throughput, Average System Travel Time and return. Its best result is a reduction of 44% in Average System Travel Time of vehicles that arrived at their destination. Overall CTRL outperforms the other toll mechanisms. In only one network  $\Delta$ -tolling shows comparable results. An ablation study shows that the Q-attention mechanism provides better tolls (defined as differentiated for all three routes and fair) and State-attention increases throughput and decreases Average System Travel Time. Overall CTRL increases throughput and decreases Average System Travel Time.

### 6.2.5 Dynamic Toll Collection using Transformers and Graph Neural Networks

Further building on the ideas described in Sections 6.2.2 and 6.2.4 leads to a model (Lu et al., 2024) where a Multi Agent Transformer (Wen et al., 2022) architecture is combined with Random Walk Positional Encoding (RWPE) (Dwivedi et al., 2022) with a GCN for structural Position Encoding (PE) (Ye et al., 2023). The objective of the *ta's* is to reduce the TSTT. For a high-level architecture, see Figure 8.

Traffic model set up: ( $\mathcal{R}$ ) Before departure, vehicles select the lowest cost route, which they do not update after departure. A number of links in the network is tolled. In the graph for the GCN, the links are modeled as vertices, and the junctions as edges. ( $\mathcal{M}$ ) For network loading a microscopic model is used (Zhang et al., 2019).

RL set up: State is defined as the number of vehicles that entered the link since the previous episode step combined with the current number of vehicles on the link. The action is a selection of one of five toll rates per link. Reward is given by the negative average queue length of vehicles on the link during the decision interval, started with  $A_t(l)$  which consists of  $\Gamma$  time steps:

$$R_{t+1}(l) = -\frac{1}{\Gamma} \sum_{i=1}^{\Gamma} q_i(l). \quad (17)$$

The approach in this model is to first pre-calculate the PE for the vertices. When taking an episode step, the observations per vertex are concatenated with the PE and fed into the GCN. The output of the GCN is input for the Multi Agent Transformer, which expresses the sequential decision process for each agent (gives the tolls for each link) using the Multi-Agent Advantage Decomposition theorem (Wen et al., 2022).



Results: The model is tested on two city networks using a microscopic traffic model (Zhang et al., 2019). The experiment encompassed three traffic scenarios simulating low, medium and high vehicle demand. The MAGT solution was compared with a MAT solution (without PE/GCN), no toll, adaptive toll (based on  $\Delta$ -tolling (Sharon et al., 2016)) and EGCN toll (Qiu et al., 2019). The MAGT toll solution outperforms all other mechanisms except in a low demand scenario, where EGCN and MAGT perform almost equally. Its best result was a reduction of 18% in Average System Travel Time. The performance is attributed to the combination of positional encoding and the Multi Agent Transformer.

1068

## 1069 6.2.6 Dynamic Toll Pricing in Regions

1070 An agent is trained to set tolls for a region-based approach (He et al., 2024) using  
1071 a similar type of day-to-day and intraday dynamics as in Section 6.5.1. The  $ta$  has  
1072 the objective of minimizing the Total System Travel Time. The performance is also  
1073 evaluated on speed profiles and total system cost.

1074 Traffic model set up: Given  $\langle G, V, \mathbb{D} \rangle$ :  $G$  is divided in a number of regions.  
1075 The traffic characteristics in the regions are based on a (Macroscopic) Fundamental  
1076 Diagram (2.5). These characteristics are determined by a trained neural network.  
1077 Vehicles travel from one region to another region. ( $c_v$ ) Vehicles take an experienced  
1078 cost on (day, time) =  $(j - 1, \underline{t})$  and the predicted cost  $(j, \underline{t})$  based on actual travel  
1079 times and tolls. Their perceived cost  $(j, \underline{t})$ , is a convex combination of experienced  
1080 and predicted cost. ( $\mathcal{R}$ ) The perceived cost is input to a softmax function to provide  
1081 the route choice probability. ( $\mathcal{M}$ ) Network loading is performed using a macroscopic  
1082 model.

1083 RL model set up: A2C (Section A.5) is used to determine the tolls. State is defined  
1084 as:

$$1085 S_t = [(n_i(t), \tau_i(t), t, i) | i \in \mathcal{I}], \quad (18)$$

1086 where  $n_i(t)$  indicates the number of vehicles in  $i$  at time  $t$ , and  $\mathcal{I}$  is the number of  
1087 regions. The action  $A_t \in [\tau_{min}, \tau_{max}]$  is the adjusted toll price, modeled as a continuous  
1088 variable. The reward is defined as the negative sum of all vehicles in the region in the  
1089 toll adjustment interval.

1090 Results: The  $ta$  is able to improve the speed profile and decrease both Total System  
1091 Cost and Total System Travel Time with a reduction of 6% as its best result. The model  
1092 also reaches an equilibrium state. When the system is tested to robustness against  
1093 input variability, the agent reaches an improved equilibrium state comparable to the  
1094 non-varied input. The model is further tested with different lengths of  $\Gamma$ . For  $\Gamma$  there  
1095 is a level (15 minutes) that provides a good balance between accuracy and feasibility  
1096 with values above (30, 60 minutes) and below (5 minutes) with lower performance.

1098

## 1099 6.3 Value Based Dynamic Toll Collection for Tolloed Lanes

### 1100 6.3.1 Distance Based Tolling with R-MART

1101 In (Zhu and Ukkusuri, 2015) a traffic network is considered in which some links have  
1102 both a non-tolled lane and a tolled lane. The objective is to dynamically determine the  
1103 toll prices for the tolled lanes so that the Total System Travel Time on these links is  
1104



minimized. Tabular Q-learning is used to find the optimal policy. One agent operates on each link with a tolled lane, which optimizes the performance only for that link, without cooperation or communication with other agents (DTDE).

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ :  $(G)$ . The links  $l$  are divided in  $N(l)$  cells of equal length. The lanes of a link are further indexed by  $i$  and, if tolled, by  $i'$  like  $l_i$  and  $l_{i'}$ .  $(c_v)$  Tolls are calculated based on the remaining distance from cell  $m$  to the exit of the tolled lane and actual toll rate  $\sigma_e, e \in 1, 2, 3, 4$ . The remaining distance is determined as the number of cells until the exit of the lane  $(N(l_{i'}) - m)$  multiplied by cell length  $C$ .

$$\tau(l_{i'}, m, \sigma_e, \underline{t}) = (N(l_{i'}) - m) * C * \sigma_e \quad (19)$$

The cost for a vehicle on the tolled lane, adding remaining travel time to the end of the lane, therefore is:

$$c_v(l_{i'}, m, \underline{t}) = tt(l_{i'}, m, \underline{t}) + (N(l_{i'}) - m) * C * \sigma_e. \quad (20)$$

And for a non-tolled lane:

$$c_v(l_i, m, \underline{t}) = tt(l_i, m, \underline{t}). \quad (21)$$

A vehicle determines the cost of using the non-tolled lane (based on remaining travel time to the end of the lane) and the tolled lane (based on the remaining travel time to the end of the lane and tolls).  $(\mathcal{R})$  Vehicles choose to either stay on a non-tolled lane, or change to a tolled lane. Route choice is implemented using a binomial logit model with lane cost as input. Tolled lanes can be entered at any position on the non-tolled lane; vice versa is not allowed.  $(\mathcal{M})$  Network loading is done using the CTM.

RL set-up: The environment is modeled as a MDP. The state  $S$  at lane  $l_i$  at time  $t$  is denoted as  $S_t^{l_i}$  and the next state as  $\bar{S}_t^{l_i}$ , and similar for  $l_{i'}$ . The states in lane  $l_i$  are discretized and based on the comparison between the density  $k$  at time  $t$  denoted as  $k(l_i, t)$  and the jam density  $k_{j, l_i}$  as below:

$$S_t^{l_i} = \begin{cases} 1, & \text{if } k(l_i, t) \leq 0.25k_{j, l_i} \\ 2, & \text{else if } k(l_i, t) \leq 0.50k_{j, l_i} \\ 3, & \text{else if } k(l_i, t) \leq 0.75k_{j, l_i} \\ 4, & \text{else if } k(l_i, t) \leq k_{j, l_i} \end{cases} \quad (22)$$

and similar for lane  $l_{i'}$ . The state on the link, comprising both the tolled and non-tolled lanes, is denoted  $S_t^l$ . The actions  $a$  in lane  $l_{i'}$  depend on the cell  $m$  and time  $t$  denoted as  $A_t^{l_{i'}}(m)$ . Action selection is  $\epsilon$ -greedy on the maximum Q-values. The toll rates corresponding to these actions are  $\sigma_b$  if  $A_t^{l_{i'}}(m) = b$  ( $b \in 1 \dots 4$ ) with  $\sigma_b$  the threshold values of the toll rate. Reward  $R_t^{l_{i'}}(S_t^l, A_t^{l_{i'}}(m), \bar{S}_t^l)$  is calculated based on total travel time on a link summing the travel times  $tt$  at lanes  $l_i$  and  $l_{i'}$  from cell  $m$  for all vehicles  $x(l, m, t)$  to the exit of the lane at time  $t$  for all lanes and all cells:

$$R_t^{l_{i'}}(S_t^l, A_t^{l_{i'}}(m), \bar{S}_t^l) = - \sum_{j \in l_i, l_{i'}} \sum_m tt(j, m, t) * x(j, m, t). \quad (23)$$

1151 A R-Markov Average Reward Technique (R-MART (Sutton and Barto, 2018)) is used  
 1152 which applies to continuing problems. In continuing problems, the interaction between  
 1153 agent and environment goes on and on forever. The quality of a policy is defined as  
 1154 the average rate of reward  $\rho$ . The Q-values for toll lane  $i$  for all cells  $m$  on a link  $l$  are  
 1155 updated as:

$$1156 \quad Q(S_t^l, A_t^{l'}(m)) \leftarrow Q(S_t^l, A_t^{l'}(m)) + \alpha \left[ R_t^{l'}(\cdot) - \rho + \max_a Q(\bar{S}_t^l, a) - Q(S_t^l, A_t^{l'}(m)) \right]. \quad (24)$$

1159 Results: In an experiment, the Total System Travel Time on the links with both a  
 1160 non-tolled and tolled lane decreases with 25% on a real network. Varying the number  
 1161 of states and actions demonstrates that increasing the number of states and/or actions  
 1162 does not lead per sé to an improvement in travel time; there is even a trend that the  
 1163 algorithm worsens when the number of states and actions are growing. In the work  
 1164 they highlight finding the best combination of states and actions as a problem for  
 1165 further study.

### 1167 6.3.2 MARL for Distributed Dynamic Pricing of Managed Lanes

1169 A distributed Multi Agent model is used to manage dynamic pricing for managed  
 1170 lanes with multiple entrances and exits (Pandey and Boyles, 2018). To this effect, the  
 1171 sum of each agents' local value is used as an approximation for the total state value.  
 1172 Furthermore, in determining the action (toll price) the agents collaborate by taking  
 1173 the actions of downstream agents into account. Therefore we categorize the training  
 1174 model as CTDE. The objective of the toll agents is to maximize revenue.

1175 Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ : ( $G$ ) Managed lane networks are considered,  
 1176 consisting of one general-purpose lane and one tolled lane, with one or more entries  
 1177 and one or more exits. ( $\mathbb{D}$ ) vehicles are categorized in different classes according to  
 1178 a set of different *vot*. ( $\mathcal{R}$ ) vehicles choose according to minimal cost of the available  
 1179 routes. ( $\mathcal{M}$ ) Network loading is done using CTM.

1180 RL set-up: State  $\mathbf{S}_t$  is defined as the vector that contains the number of vehicles  
 1181 of each *vot* class per cell. Action  $\mathbf{A}_t$  is the action vector with the actions of all agents  
 1182 (toll prices). Reward  $R_t(\mathbf{S}_t, \mathbf{A}_t)$  is the product of the number of vehicles choosing the  
 1183 managed lane times the (toll rate/mile) times the length of of the managed lane. A  
 1184 penalty is applied if the managed lane becomes congested. Because the assumption  
 1185 is that toll agents collaborate and need to coordinate their actions with only a few  
 1186 neighboring agents, the MDP is relaxed. This is done by approximating the value of  
 1187 the state as the sum of the individual value functions of the agents' state:

$$1188 \quad V_*(\mathbf{S}_t) = \sum_{n \in N} V_*^n(\mathbf{S}_t). \quad (25)$$

1192 For each toll agent the Bellman equation (A8) at optimality is:

$$1193 \quad V_*^n(\mathbf{S}_t^n) = \max_{\mathbf{A}_t} R^n(\mathbf{S}_t, \mathbf{A}_t) + V_*^n(\mathbf{S}^n(t+1)) \quad \forall n \in N, \quad (26)$$

where also the reward is decomposed into the rewards for each agent. To solve this relaxed MDP some modifications are applied to Q-learning. The first modification is that the actions of downstream agents are taken into account in action selection. The second modification is that the value function is used based on cooperative Q-learning (Kok and Vlassis, 2006) as opposed to the action-value (Q) function. The resulting algorithm is called SparseV.

Results: SparseV is compared with three other algorithms using heuristics: a) Density: if the density downstream in the managed lane differs from the required density, toll price is adjusted up or down; b) Ratio: the ratio, between the density in the cells downstream in the managed lane and the density in the cells in the general-purpose lane, determines adjustment of the toll price, and c) Random: a number of policies with random action is simulated, the best performing random policy is chosen. Experiments are performed on a synthetic and real network.

SparseVs outperforms the Density and Ratio heuristic which produce a revenue that is 70-75% lower. The Random policy outperforms SparseV by approximately 9 % in the first network and is outperformed by approximately 24 % in the other network. The better performance of the SparseV and Random policies is explained by the 'jam-and-harvest' nature of these policies; by setting initial tolls such that the General Purpose Lane becomes congested and subsequently charging higher tolls when demand increases. The work highlights convergence to sub-optimal values, attributed to exploration in relation to the aggregation level of the state space and the jam-and-harvest nature, unwanted in practice, as weaknesses.

## 6.4 Actor-Critic Dynamic Toll Collection for Tolled Lanes

### 6.4.1 Deep RL for Dynamic Pricing of Express lanes

Deep RL methods are tested in a wide range of use cases (access locations, OD-pairs, VOT heterogeneity, partial observability) where links have both a toll lane and a General Purpose Lane (GPL) (Pandey et al., 2020). A2C, PPO and SAC (Section A) are compared against a feedback control heuristic. The agents are trained for two different objectives (not in the same experiment): Maximum Revenue and minimum Total System Travel Time.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ :  $(G)$ . All links are denoted as  $l \in L$ , links which are tolled are denoted as  $l' \in L$ .  $(c_v)$  Toll lanes have multiple entrances and exits, where tolls are charged.  $(\mathbb{D})$ . Vehicles are divided in different classes based on their *vot* and on their destination.  $(\mathcal{R})$  Lane choice: Two models for lane choice are used, a multi-class binary logit model (two routes, stochastic choice, multiple *vot*) if there are only two routes, and a multi-class decision route model (multiple routes, deterministic choice, multiple *vot*) if more routes are available.  $(\mathcal{M})$  Network loading is done using CTM.

RL set-up: As the state is not fully observable the problem is defined as a POMDP for which the policy  $\pi(\mathbf{a}|\mathbf{o}(s))$  is learned. The work uses a finite horizon  $H$  and incorporates time (the toll update step number) in state observations. This is because of the temporal dependence of the congestion pattern. The observation vector  $\mathbf{o}(s)$ , in which they also model noise, for state  $s$  comprises the number of vehicles on each link,

1243 using  $\mathbf{1}_l(v) = 1$  if  $v$  is on link  $l$  and 0 if not:

1244

1245

1246

1247

$$\mathbf{o}(s) = \left\{ \sum_{i=1}^I \mathbf{1}_l(d_i) | l \in L \right\}. \quad (27)$$

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

The action  $\mathbf{a}$  in state  $s$  is the vector of tolls  $\tau(l')$  which are being applied on the tolled links  $l'$ . Tolls are modeled as a continuous variable. Two types of reward are analysed: rewards based on revenue maximization  $R_{MR}(s, \mathbf{a})$  which is the sum of the product of toll rate and vehicles entering a toll lane during a toll update time period. In the other type rewards are based on travel time minimization  $R_{TSTT}(s, \mathbf{a})$  which counts the number of vehicles every toll update time step. The algorithms to find the policy are A2C, PPO and SAC.

Results: A series of experiments is conducted on synthetic and real networks, for various objectives, variations of the observation space and variations on demand, route choice and *vot*. All three algorithms with a single objective converge for all networks. SAC outperforms PPO and A2C, which is attributed to the entropy regularization used by SAC. In an experiment on the real network it appears that reducing the observation space of the network (reducing 15 links to 4 links to 1 link under observation) does not reduce the convergence speed. As this was an unexpected result it was speculated that this happens due to spatial correlation of congestion patterns on a corridor. When conducting experiments on the same network, with changed demand and *vot* to test generalization and transferability, the key observations are that a) the algorithms learn from the new input situation at equal speed as the original situation, and b) that the transferred policies from earlier experiments (PPO, A2C), perform within a bandwidth of 5 - 12% of the newly trained policies in their experiments. This is not the case if the lane choice model is changed; in that case the transferred policy performs badly, attributed to the fact that the lane choice model has a large impact on the evolution of congestion. Comparison with an industry-based heuristic shows that for all networks for both objectives, DRL outperforms this heuristic.

1273

1274

#### 6.4.2 A Priori Link Selection and Dynamic Tolling of Expressways

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

In case of a high-speed road network, where each link is tolled, the distance difference between routes can be very large, and route adjustment is relatively difficult. Building on the work in Section 6.4.1 another model is designed (Zhang et al., 2023) in which vehicles weigh distance, travel time and cost and apply multi-objective optimization. Each *ta* is responsible for one toll lane.

In experiments, the model demonstrates its effectiveness. However they note that it does not take account of temporal correlation of traffic flow between different road sections in the network.

## 6.5 Dynamic Tolling to Manage Departure Time Choice

### 6.5.1 Dynamic Congestion Pricing for Departure Time Choice

Another model is proposed for dynamic toll pricing for departure time choice (Sato et al., 2021). The objective is to minimize the total delay by influencing the departure time of the vehicles. Q-learning is applied to learn the optimal policy, using a Q-value function approximator.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (G)$ . All links in the network have one bottleneck.  $(c_v)$  The cost of a link with a bottleneck is based on the day  $j$ , the time  $t$ , the delay  $w$ , the required arrival time  $\underline{t}^*$ , a vehicle's *vot*, toll and a penalty for arriving early or late at the destination.

$$c_v(l, j, \underline{t}) = \tau(l, j, \underline{t}) + \text{vot} * w(l, j, \underline{t}) + \begin{cases} \beta(\underline{t}^* - \underline{t})(\underline{t} < \underline{t}^*), \\ \gamma(\underline{t} - \underline{t}^*)(\text{otherwise}). \end{cases} \quad (28)$$

In this equation,  $\beta$  is the early arrival penalty factor and  $\gamma$  the late arrival penalty factor (see Section 3). Tolls  $\tau(l, j, t)$  consist of a stabilization toll  $\tau^s(l, j, t)$ , a parameterized function based on departure rate (Seo, 2020), and a congestion toll  $\tau^c(l, j, t)$  based on the action of the policy.  $(\mathcal{R})$  The model studies the morning commute. The day-to-day dynamics for departure flow in the morning commute are modeled using replicator dynamics (Schuster and Sigmund, 1983). This mechanism gives for every OD-pair the time  $\underline{t}$  and day  $j$  dependent departure rate  $\varsigma(p, j, \underline{t})$  for all paths  $p \in p^{o \rightarrow d}$ . All paths  $p^{o \rightarrow d}$  going through link  $l$  are denoted by  $P_l^{o \rightarrow d}$ . The departure rates per route per OD pair at time  $\underline{t}$  are given by a logit model. Both have travel cost per route per departure time as input. The departure time of the vehicle is therefore not fixed.  $(\mathcal{M})$  Network loading is performed with an extension of the bottleneck model. The bottleneck model can be seen as a special case of the LWR model (Section 2). In a bottleneck model (Arnott et al., 1990; Seo, 2020; Vickrey, 1969) the bottleneck has a fixed capacity (service rate) and if the arrival rate exceeds this capacity, a queue begins to form. In the model used here, the queue occupies no space.

A vehicle arrives at the bottleneck of a link after the free flow travel time from the begin of the link to the bottleneck of the link, and departs the link immediately after leaving the bottleneck. The time spent at the bottleneck is the waiting time  $w$  (also: delay).

RL set-up: Q-learning is used to determine the congestion tolls. The effect of tolls on departure time should lead to a reduction in waiting time at the bottleneck. A spatially shared reward structure is to take account of tolls on other links. State  $\mathbf{S}$  is a vector of the departure rate from an origin in an OD-pair, waiting time at the bottleneck and the stabilisation toll (T for transposed):

$$\mathbf{S}_t^l = \left( \sum_{p' \in P_l^{o \rightarrow d}} \varsigma(p', t), w(l, t), \tau^s(l, t) \right)^T. \quad (29)$$

1335 Action  $A_t^l$  is the update to the congestion toll, therefore:

$$1336 \quad \tau^c(l, t + 1) = \tau^c(l, t) + A_t^l. \quad (30)$$

1339 Reward  $R$  is the weighted sum (by a factor  $K$ ) of the inverse difference between the  
1340 bottleneck capacity and traffic flow at the bottleneck in link  $l$  lower bounded by a  
1341 constant  $E$  (denoted  $\xi_t^l$ ) and the bottleneck capacity and traffic flow for all other links  
1342  $l'$  with bottlenecks collectively also lower bounded by a constant  $E$  (denoted by  $\xi_t^{l'}$ ):

$$1344 \quad R_t^l = (\xi_t^l)^{-1} + K(\xi_t^{l'})^{-1}. \quad (31)$$

1346 More reward is given if departure rates are close to bottleneck capacity and it takes  
1347 account of this for the other bottlenecks. This mechanism implements the spatially  
1348 shared reward structure. The policy finds the adjustment tolls (the actions), by choos-  
1349 ing actions  $\epsilon$ -greedy, leading to minimal total waiting time. The parameter updates  
1350 are learned by Q-value function approximation for which Radial Basis Functions are  
1351 used (Sutton and Barto, 2018).

1352 Results: Experiments are performed on two synthetic networks: one where they  
1353 apply the method to a single link single bottleneck model, the other with a single  
1354 OD-pair and three links with three bottlenecks. After training, the agent succeeds  
1355 in reducing the waiting time to almost zero and outperforms a trial and error con-  
1356 gestion pricing scheme (Seo, 2020). Performance in the three bottleneck model was  
1357 inefficient compared to the single bottleneck model. This is attributed to the reward  
1358 function in the sense that the spatially shared reward function did not provide sufficient  
1359 coordination between the bottlenecks.

## 1361 6.5.2 Pricing for Departure Time Choice with DDPG

1362 Extending Sato et al. (2021) outlined in the departure choice model above  
1363 (Section 6.5.1) with DRL, and temporally switching learning leads to Distributed  
1364 Pricing-Deep Deterministic Policy Gradient (DP-DDPG) (Sato et al., 2022).  
1365 DDPG (Lillicrap et al., 2015), (Section A), is used as RL mechanism. A modifica-  
1366 tion is that an action is fixed at zero if the moving average of the waiting time is less  
1367 than a threshold (temporally switching learning). This leads to more efficiency and  
1368 less excessive increases in tolls. Distributed control, using spatially shared rewards  
1369 (Section 6.5.1), is implemented to enable cooperation.

1371 Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ :  $(G)$  and  $(c_v)$  are as per (Sato et al., 2021)  
1372 detailed above (Section 6.5.1).  $(\mathcal{R})$  Different from (Section 6.5.1), the departure time  
1373 and route choice are based on a multinomial logit model. Furthermore, vehicles have  
1374 bounded rationality, meaning that departure time and route will not be changed if  
1375 the difference between expected cost of current choices compared with the expected  
1376 costs of alternatives remain below a certain threshold.  $(\mathcal{M})$  For network loading the  
1377 bottleneck model is used. The objective of the toll agents is to reduce delay on the  
1378 link under its responsibility.

1379 RL set-up: State is defined in a similar way as above (Section 6.5.1), adapted for  
1380 input to a neural network. Action is an adjustment to the toll price for a bottleneck.

Action  $A_t^l$  is the update to the congestion toll:

$$\tau(l, j, t) = \tau(l, j, t) + A_t^l. \quad (32)$$

The spatially shared reward  $R$  is a weighted average of the waiting time at bottleneck  $i$  of agent  $ta^i$  and the bottlenecks of the other agents  $ta^{-i}$ , slightly modified from (Sato et al., 2021) above. With this reward, each agent  $ta^i$  takes into account the rewards of the other agents  $ta^{-i}$  and therefore their objectives are aligned. Each agent  $ta^i$  uses its own instance of the DDPG algorithm to observe state, leading to an action and obtaining a reward and observation of the next state. This data  $(\mathbf{S}_t^l, A_t^l, R_t^l, \mathbf{S}_{t+1}^l)$  for each agent is stored in an experience replay memory, from which data is sampled to update the actor and critic networks of the agents. If the temporally switching learning criterion has been met, no action is taken (so no toll updates) and no data is added to the experience replay buffer, so this data is not used for learning.

Results: Two experiments are performed, one on a network with three links for one OD-pair with a bottleneck on every link. Another experiment involved a real network to evaluate the performance where multiple OD pairs and multiple routes exist. In the first experiment, the deadline for all vehicles is the same and DP-DDPG is compared with centralized DDPG, fully distributed DDPG and Q-Learning (Sato et al., 2021). DP-DDPG outperforms centralized DDPG and fully distributed DDPG as these could not decrease waiting time. Q-Learning however succeeded in decreasing waiting time. In a simplified real network, toll is charged on four links with a bottleneck and there are three links with a bottleneck without charging tolls. The same arrival time is set for all vehicles. The results demonstrate that the waiting time reduces, meaning that with a trained agent, the vehicles also adapt their behavior leading to lower waiting times.

## 6.6 Value Based Dynamic Tolling Including Vehicles as Learning Agents

### 6.6.1 Multi Agent Road Pricing With Learning vehicles And Toll Agents

A traffic network is modeled as a Multi Agent system with two types of agents (Tavares and Bazzan, 2014). For the demand side, vehicles that are traversing the network, with the objective to optimize route choice. For the supply side, toll agents that are setting tolls on links in the network, with the objective to maximize flow. For both types of agents DTDE applies.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (c_v)$ . Vehicles have knowledge of the network based on experience, which means they only know the toll  $\tau$  and travel time  $tt$  for a link if they travelled it in the past and recorded it as  $\langle \underline{t}, v, l, tt, \tau \rangle$ . If vehicles have never used a link, the travel time is based on the free flow travel time  $T(l)$  and for the toll, the maximum toll  $\tau'(l)$  divided by 2:  $\tau = \tau'(l)/2$ . The cost for a vehicle is the sum of link costs. Link costs  $c_v(l, tt, \tau)$  are vehicle specific, and based on a weighting factor  $\eta(v) \in [0, 1]$ :

$$c_v(l, tt, \tau) = \eta(v) * tt(l, v) + (1 - \eta(v)) * \tau(l, v), \quad (33)$$



where travel times and tolls are based on prior recorded experiences or alternatively the default values. ( $\mathcal{R}$ ) Prior to each episode, vehicles calculate their optimal (lowest cost) route based on calculated link costs and do not update their route while traversing the network. ( $\mathcal{M}$ ) Network loading is done using a microscopic traffic model.

RL set-up: The toll agents are trained as independent learners to limit the problem size (DTDE). The toll agents' actions consist of setting a toll price  $\tau(l)$  as a fraction of the maximum toll price  $\tau'(l)$  for a link. The model is stateless and only a reward  $R_l$  is observed after an action; the reward  $R_l$  is the number of vehicles  $\nu_l$  which have entered a link during the episode. The action values  $Q(\tau(l))$  are stored in a Q-table and are learned according to:

$$Q_l(\tau(l)) \leftarrow (1 - \alpha)Q_l(\tau(l)) + \alpha R_l. \quad (34)$$

where  $\alpha$  is the learning rate. The objective is to find the  $\epsilon$ -greedy policy which maximizes the flow, where  $\epsilon$  is decreasing over time. Vehicles initially only know the values of free flow travel time for all links and maximum price for all links. Vehicles can update these values by acquiring local knowledge, that is when they travel a route learning the actual travel times of the links and tolls they pay.

Results: Under different vehicle preference scenarios, the performance of the *ta* is compared with a fixed toll scenario and a Dynamic User Equilibrium scenario which assumes global knowledge for the vehicles. The DUE solution outperforms the MARL solution, which again outperforms the fixed toll scenario. While the average performance of the *ta*'s increases, large differences exist between best and worst performers. In summary, the case where both vehicles and toll agents learn outperforms the fixed-toll solution in these experiments. Global performance increased (more vehicles completed trips in the network) and vehicle costs decreased. These outcomes were achieved without an explicit coordination mechanism between the Toll Agents. With global information DUE led to better results.

## 6.6.2 Deep Multi Agent Road Pricing with Learning Vehicles and Toll Agent

A variation on 6.6.1 uses DQN (Mnih et al., 2015) both for a *ta* and for vehicles making their lane choice (Chakravarty et al., 2024). The objective of the *ta* is to maximize revenue, generated by the tolled links ( $l'$ ) defined as  $rev(t) = k_{l'}(t) \cdot v_{l'}(t) \cdot \tau_{l'}(t)$ , the product of density, speed and toll price. The objective of the vehicles is to be in the fast lane, but taking into account the toll price.

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ : In the network are two parallel links, one tolled ( $l'$ ), one non-tolled ( $l$ ) where the difference in speed is noted as  $\Delta v = v_{l'} - v_l$ . Shallow neural networks are pre-trained using real-world data to relate observed density to speed ( $f_1(k) = u$ ) and toll price and speed to provide the fraction  $\xi$  choosing the tolled lane ( $f_2(\tau, k) = \xi$ ) and the non-tolled lane ( $1 - \xi$ ). ( $c_v$ ) The vehicle cost function is not made explicit. ( $\mathcal{R}$ ) Route choice is made by an agent using DQN. ( $\mathcal{M}$ ) We categorize this as a macroscopic model.



RL Set up: The state for vehicles is defined as:

$$S_t = [\tau, \Delta v, \int \Delta v dt], \quad (35)$$

and the reward function, with  $F = 0, 1$  for being in the slower or faster lane as:

$$R_t = -a_1 \tau^2 + a_2 F. \quad (36)$$

Vehicle characteristics are modeled by  $a_1, a_2$  to indicate how much they value speed over toll price. The action  $a_t$  consists of choosing the tolled or non-tolled lane.

For the toll agent state is defined as:

$$S_t = [v_T, \Delta v, \int \Delta v dt, k_{in}, rev(t)], \quad (37)$$

with  $k_{in}$  total incoming density. The action space is a set of toll prices. Finally, the  $ta$  reward combines revenue and  $\Delta v$ :

$$R_t = b_1 * rev(t)^2 - 2 * B, \quad (38)$$

where  $B = 0$  if  $\Delta v \geq 0$  and  $B = 1$  else.

Results: The driver agent succeeded in improving its reward. Training the agent ceased when the average reward was above a certain threshold for five episodes. Another outcome was that vehicles were never willing to pay more than a certain toll price, although the speed in the tolled lane was higher. The  $ta$  succeeded in increasing its reward as well. Training for this agent was stopped if the average reward was above a certain threshold.

The authors note that both the  $ta$  and the vehicle agent were able to improve their performance. However they indicate the vehicle agent could be improved by a more complex reward function and that their set-up of the environment showed some weaknesses in calculating the density fraction.

## 6.7 Value Based Marginal Cost Tolling

### 6.7.1 Multi Agent Reinforcement Learning - Only Vehicles

In [Ramos et al. \(2018, 2020\)](#) vehicles are represented by a Q-learning agent whose objectives are to find the routes with the optimal Q-value. The tolling scheme, based on Marginal Cost Tolling (MCT), is distributed. The objective of the tolling scheme is to bring the User Equilibrium closer to the System Optimum (Section 2).

Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle: (c_v)$ . The toll on a link for a vehicle is the difference between the free flow travel time  $T(l)$  and the actual driving time on the link  $tt(l)$ , multiplied by a constant  $\beta$  which is specific for the Link Performance Function (Section 2).

$$\tau(v, l, tt(l, \underline{t}), \underline{t}) = \beta(tt(l, \underline{t}) - T(l)). \quad (39)$$

1519 The cost for a link is the sum of the toll and actual driving time:

1520

$$1521 \quad c_v(l, tt(l, \underline{t}), \underline{t}) = tt(l, \underline{t}) + \tau(v, l, tt(l, \underline{t}), \underline{t}). \quad (40)$$

1522

1523 The self-calculated cost by the vehicles for a path  $p^{o \rightarrow d}$  is the sum of the cost of the  
 1524 links  $l \in p^{o \rightarrow d}$  and is denoted as  $c_v(p^{o \rightarrow d})$ . ( $\mathcal{R}$ ) Vehicles know a subset of their routes  
 1525 a priori, and the only action is to select the route every day, based on the Q-value of  
 1526 the routes in the subset. ( $\mathcal{M}$ ) Network loading is done by a microscopic model.

1527 RL set-up: This problem is modeled as a stateless MDP, as the vehicle takes an  
 1528 action  $A$ , based on the highest Q-value of that action and then obtains a reward  $R(A)$   
 1529 based on that action. The reward  $R(A)$  is defined as:

1530

$$1531 \quad R(A) = -c_v(p^{o \rightarrow d}). \quad (41)$$

1532

1533 The vehicles want to take a route  $p^{o \rightarrow d}$  which minimizes their costs. Therefore, this  
 1534 entails finding the policy  $\pi$ , which provides the route to take and maximises the agent's  
 1535 average reward. The algorithm used is Q-learning (Watkins and Dayan, 1992) with  
 1536 independent Q-Learners (Claus and Boutilier, 1998) thereby implementing a DTDE  
 1537 framework. The Q-learning algorithm is modified to a stateless version:

1538

$$1539 \quad Q(A) = (1 - \alpha)Q(A) + \alpha R(A). \quad (42)$$

1540

1541 The learning rate  $\alpha$  is exponentially decaying per episode. Every episode an agent takes  
 1542 an action (chooses a route), using an  $\epsilon$ -greedy exploration strategy with exponentially  
 1543 decreasing exploration rate per episode. Then it updates its Q-value.

1544 Results: The algorithm is tested on 13 different (synthetic and real) networks  
 1545 and compared against three other algorithms namely  $\Delta$ -tolling (Sharon et al., 2017)  
 1546 (see 6.1.1), toll-free Q-Learning where they apply the same model without tolls and  
 1547 one other mechanism. Their method provides similar results as  $\Delta$ -tolling and better  
 1548 performance than toll-free Q-learning to approach SO. In comparison to  $\Delta$ -tolling this  
 1549 method has a learning scheme, a fairer tolling scheme and convergence guarantees. The  
 1550 work however highlights that speed of convergence is a major area for improvement.  
 1551 A similar work is (Ramos et al., 2020) in which it is also highlighted that speed of  
 1552 convergence is a major area for improvement.

1553

## 1554 6.7.2 Multi Agent Reinforcement Learning - Heterogeneous Agents

1555

1556 An algorithm capable of realigning agents' heterogeneous preferences over travel time  
 1557 and monetary expenses to obtain a system efficient equilibrium (SO) is Generalised  
 1558 Toll-based Q-Learning (Ramos et al., 2020). It also includes a mechanism to enforce  
 1559 agents to truthfully report their preferences. The key difference with (Ramos et al.,  
 1560 2018, 2020) (see 6.7.1) is the calculation of the cost function. The objective of the  
 1561 tolling scheme is to bring the User Equilibrium closer to the System Optimum.

1562 Traffic model set-up: Given  $\langle G, V, \mathbb{D} \rangle$ : ( $c_v$ ) In this model with heterogeneous  
 1563 preferences for a vehicle  $v$ , travel time and toll are weighted by  $(1 - \eta(v))$  and  $\eta(v)$

1564

respectively. Furthermore the cost includes side payments indicated by  $\rho_{\psi(v)}$ :

$$c_v(p^{o \rightarrow d}) = -\rho_{\psi(v)} + \sum_{l \in p^{o \rightarrow d}} ((1 - \eta(v)) * tt(l, \underline{t}) + \eta(v) * \tau(v, l, tt(l, \underline{t}), \underline{t})). \quad (43)$$

The side payments are a redistribution of a part of the collected tolls to the vehicles as not all collected tolls are required to operate the toll system.  $(\mathcal{R}, \mathcal{M})$  as in Section 6.7.1

RL set-up: The reward  $R$  is as Equation (41) and Q-learning algorithm as Equation (42). The difference is that the reward also includes the side payments, described above. To enforce truthful reporting (to prevent misuse of the side payments) the number of times an agent has not chosen the least-cost action based on their preferences is logged. If the number of inconsistent actions, which means reporting action  $a$  and executing action  $a'$ , exceeds a threshold the agent is expected to be cheating which leads to a penalty.

Results: The algorithm is tested on 15 different (synthetic and real) networks. The results are compared with the algorithm in (Ramos et al., 2018, 2020) (see 6.7.1) and with  $\Delta$ -tolling from (Sharon et al., 2017) (see Section 6.1.1). In general, GTQ-learning was able to converge to a SO efficient equilibrium, leading to a reduction of 30% in Average System Travel Time. Furthermore, when vehicle preferences were applied, GTQ-learning outperformed in almost all cases the other two algorithms. It was also experimentally demonstrated that side payments did not deteriorate the equilibrium. However the other algorithms also achieved reasonable results. The mechanism to prevent misreporting, neutralized agent misreporting, restoring system optimality. Although the mechanisms (Ramos et al., 2018, 2020; Sharon et al., 2017) produced similar results, results in this study were obtained under more realistic assumptions. The work compares to (Ramos et al., 2018, 2020) which does not take heterogeneous preferences into account and to (Sharon et al., 2017) which does not work in a decentralized way. All three do ignore misreporting and do not provide for tax returns.

## 6.8 Other Works

One work (Nisha et al., 2024) describes an approach to deploy RL for road pricing using a comprehensive state, action and reward set. In their experiments, they compare their solution with three traditional models (Decision Tree, SVM and logistic regression) which are outperformed. The main part of the work describes the design process and the components of a complete solution. Although it highlights the potential of RL for road pricing, we exclude it from further analysis as it does not provide details on the (RL-) design itself.

Another work (Chiou, 2024) combines a stochastic mathematical program with equilibrium constraints with a Q-learning performance index (KAQPI). This value function takes a combination of the rate of delay, the number of stops in the network (caused by traffic signals) and the sum of link tolls which it strives to minimize. The main focus in this work is stochastic optimization and modeling, which we exclude in this survey. We refer to the work for details.

## 1611 6.9 An assessment of the potential of RL for road pricing

1612 From the works analyzed in this Section, it is a challenge to specifically identify  
 1613 what contributes the most to the demonstrated performance improvements. As can be  
 1614 inferred from Table 3, there is a wide variety in toll agent objectives, RL design choices,  
 1615 and vehicle characteristics. In addition, most of the works use different simulators  
 1616 (see Table B1) and perform their experiments on different road networks (Table B2),  
 1617 which also may affect the results. This further complicates drawing general conclu-  
 1618 sions whether the performance is attributable to the agent, the specific traffic network  
 1619 settings, or the choice of suitable (hyper) parameters.

1620 Nevertheless, across all works in this Section, seven provide a percentage for the  
 1621 reduction of Average/Total System Travel time (see Table 5). We observe that the  
 1622 bandwidth of reduction ranges from 6% (Section 6.2.6) to 44%, 45% (Sections 6.1.1  
 1623 and 6.2.4). The other works had results between these maxima and minimum. While  
 1624 these results suggest that using RL for road pricing can lead to meaningful reductions  
 1625 of Average/Total System Travel time we caution against interpreting them as directly  
 1626 comparable. The works differ significantly in terms of traffic models, simulation envi-  
 1627 ronments, road networks, demand patterns, agent designs and RL algorithms. These  
 1628 variations make it difficult, maybe even impossible, to attribute these performance  
 1629 gains to specific factors or draw generalizable conclusions. However, the diversity across  
 1630 the works demonstrates that various combinations of environment, agent design and  
 1631 objectives can achieve beneficial outcomes. We consider even a 6% improvement in  
 1632 Average/Total System Travel time to be a significant achievement, underscoring RL  
 1633 as a relevant contribution to solving the road pricing problem.

1634

1635

1636 **Table 4:** Scaling parameter softmax

Section	$\theta$
6.2.1	0.5
6.2.3	0.5
6.2.6	0.03
6.3.1	1.0
6.5.1	0.015

1644

1645

1646

## 1647 7 Challenges and Future Research Directions

1648 The works covered in Section 6 of this survey demonstrate the usefulness of Reinforce-  
 1649 ment Learning for various road pricing scenarios, ranging from tolled lanes to larger  
 1650 road real world networks and corresponding demand data.

1651

### 1652 7.1 How Identified Challenges are Addressed

1653

1654 Either explicit or implicit, the challenges related to partial observability, credit assign-  
 1655 ment and non-stationarity have been addressed in most of the surveyed works.

1656

**Table 5:** Reduction in Travel Time

Section	Reduction
6.1.1	45%
6.2.1	15%
6.2.4	44%
6.2.5	18%
6.2.6	6%
6.3.1	25%
6.7.2	30%

However, given the limited number of works, we cannot conclude which factors contribute the most to driving performance and effectively addressing these challenges. On the other hand, a CTDE architecture in multi agent cases and the use of pre-processing networks provide positive indicators of their usefulness. Scalability and generalizability have received little attention in the surveyed works in Section 6.

Partial observability In some of the approaches (Section 6.4.1), partial observability was explicitly noted, in which case the observation was provided to the learning algorithm. It also modeled noise effects of sensors. In the other single agent approaches, the choice (and implicit simplification) to frame the problem as an MDP enabled the use of simpler algorithms and demonstrate their feasibility. In the solution approaches involving multiple agents, the CTDE architecture (Sections 6.2.2, 6.2.3, 6.2.4, 6.2.5, 6.3.2 and 6.5.2) and the use of communication between agents (Section 6.2.2), address this challenge. In the single case where DTDE is compared with CTDE (Section 6.5.2) the CTDE architecture outperforms DTDE. We observe that in the other DTDE cases, the toll agents were not using state (Sections 6.6.2, 6.7.1, 6.7.2) or no comparison with a CTDE architecture was performed (Section 6.3.1).

Credit Assignment The approaches in Sections 6.2.2, 6.2.4, 6.2.5 address the spatio-temporal dependencies by using a GCN (Kipf and Welling, 2017) and transformer (Vaswani et al., 2017). The multi agent credit assignment challenge has various kinds of solution mechanisms surveyed in the works (Du and Ding, 2021; Gronauer and Diepold, 2022; Wong et al., 2022). In Section 6.3.2 one of these mechanisms (value decomposition) is used, in other works this challenge is not explicitly addressed. We can not conclude that Credit Assignment, for RL applied to road pricing, is solved. Credit Assignment in this domain is complicated because routing decisions, as observed from the works, also depend on effects from other traffic. A clear approach to link the toll price, taken into account by individual vehicles, and its effect on future traffic on that road has not emerged from the works.

Non-stationarity A CTDE architecture and communication (Sections 6.2.2, 6.2.3, 6.2.4, 6.2.5, 6.3.2 and 6.5.2) partially address this challenge. However, the challenges when more complex vehicle behavior is modeled (Section 6.6.1), and when assumptions about vehicle behavior may not hold in the real world (Section 6.1.1) remain. If in addition aspects like mode choice behavior (other means of traffic) and modeling (temporary) disruptions of roads in the network are considered, the non-stationarity challenge is further exacerbated.

Scalability In most of the works, learning is first demonstrated on a synthetic network followed by experiments on real world networks. As noted (Sections 6.1.1, 6.2.2 and 6.2.3) a larger state and action space is used than for example in classic RL problems like cartpole or mountain car (Sutton and Barto, 2018) but overall the state and action space remain small. The impact of state and action space on performance (Section 6.3.1) and of the traffic model used (Section 6.4.1) was mentioned as related challenges in the analyzed works. The works (Sections 6.2.1, 6.5.2) specifically highlighted scalability as an issue.

Generalizability This aspect was partially highlighted (Section 6.4.1) to consider how algorithms could transfer across different traffic models. In the works, different topologies or origin-destination pairs were not considered. As generalizability would

increase acceptance in real-world implementations, this challenge has in our view only minimally been addressed and therefore also remains. We do note that existing, known approaches for the above challenges, surveyed by, among others, (Du and Ding, 2021; Gronauer and Diepold, 2022; Wong et al., 2022) are only to a limited extent utilized. With respect to the design of state, action, and reward, in each of the works, tailor-made state/observations were created for each environment. The actions varied from discrete to continuous toll prices, for roads and routes. Reward functions also showed a great variety.

#### 7.1.1 Additional Findings

In the simulated environments, covering both synthetic and real-world networks, agent learning has been demonstrated. Nevertheless, specific learning issues noted in the works remain. These were finding the balance between exploration and exploitation (Sections 6.3.1, 6.3.2 and 6.6.1), sub-optimal performance because of the aggregation level of the state space (Section 6.3.2) and convergence challenges (Section 6.1.1) related to changing traffic patterns.

Furthermore, experiments demonstrated that in a multi-agent setting an optimal number of agents exist (Section 6.2.2) and that more fine-grained observation and action spaces not necessarily lead to better performance (Section 6.3.1). However, a detailed analysis of these findings was left as future work by the authors. Another observation concerns the choice of the scaling parameter  $\theta$  used in stochastic route choice models (see Equation 4 and Table 4). The value of  $\theta$  varies considerably across the surveyed works, affecting the resulting traffic distributions. A detailed analysis of the selection or calibration of  $\theta$ , and its impact on agent performance, was generally not provided. In contrast to learning in simulated environments, none of the works covered off-line learning based on real data, or learning in the real world. Unfortunately, despite being present and of major importance, existing literature on RL for road pricing has not elaborated enough on deployment barriers. There are technical challenges, like legacy infrastructures, data acquisition, integration, quality and exploration. Financial obstacles are definitely present, like deployment, operational and maintenance cost. Furthermore, there are societal issues, like acceptance and use by the general population, or diversion of traffic through narrow roads or quiet neighborhoods. Last, but not least, there may be physical bottlenecks, like the ability to install high tech equipment or communication coverage issues. To the best of our knowledge, no cities or transport agencies have yet implemented full-fledged reinforcement learning-based road pricing schemes.

Although all solution approaches demonstrated convergence in learning, the variability in selection of among other the environment (including the use of different traffic simulators and different networks), multi-agent architectures, agent algorithms and evaluation mechanisms, makes comparison of solutions hard. The inherent variability in RL experiments, where evaluation presents significant challenges (Patterson et al., 2024), is aggravated by differences in setups in the works described above. This points to the need for a common benchmark.

## 7.2 Future Research Directions

In addition to the challenges which remain, as noted above, we give two future research directions based on what we think has not, or only to a limited extent, been covered.

First, in general, we observe that an accepted benchmark, to compare the solutions against, is lacking. In the works we observe that the features of state and observation differ, as well as the reward signals and actions. Furthermore, pre-processing of data, the (Reinforcement Learning) algorithms, and also the means of network loading vary per surveyed work. This makes it hard to specifically identify what causes the changes in behavior and performance of the models.

A significant next step for road pricing using Reinforcement Learning lies in establishing standardized benchmarks and fostering shared datasets. These are essential to enable fair comparison of diverse solutions and accelerate progress in the field. We envision two primary objectives for such benchmarks: First, to rigorously test and compare the performance of various toll agents and their algorithms, a dedicated simulation environment would be highly beneficial. Such an environment, ideally constructed using frameworks like Gymnasium (Towers et al. (2024)), PettingZoo (Terry et al. (2021)), RL4CO (Berto et al. (2023)) or BenchMARL (Bettini et al. (2024)), would need to encompass diverse network topologies, realistic demand patterns, and standardized settings for route- and departure time choice, alongside consistent network loading models. Second, for demonstrating scalability and the efficacy of agents in more realistic, complex traffic settings, future efforts should focus on meticulously detailing the specific simulators used, their full range of traffic settings, and comprehensive interface specifications between the agents and simulators. Crucially, this must be coupled with initiatives to facilitate the sharing of generated datasets (like demand and traffic patterns), which is paramount for enhancing reproducibility and comparability across studies.

By pursuing these avenues, the community can collectively move towards more robust, reproducible, and impactful research in RL-based road pricing. To expand further on this, performing real world experiments is also needed to test the validity of outcomes from simulations.

Next, road pricing is only one of a range of ITS mechanisms. As exemplified by other surveys (Section 1.1) mechanisms like traffic signal controls, variable speed limits and ramp metering are modeled using an MDP formulation as well and involve agents which take actions. Understanding the interaction between agents from other ITS and road pricing agents is essential. First by observing and analysing emergent behaviors. Second by creating and studying agents which can cope with these situations. Obviously the same could be said about the vehicles; in most of the works these are modeled as probabilistic actors who are not learning.

Moving beyond the strict scope of the works covered in this survey, future research could explore connections with adjacent fields such as economics, sociology, and issues of equity and fairness. There is also significant potential in multidisciplinary approaches, where theoretical solutions are tested and evaluated under real-world conditions. This would include technical integration with existing traffic management systems.



## 1795 8 Conclusions

1796  
1797 Road pricing is an effective and convenient method to address traffic congestion. It  
1798 poses various challenges given its traffic dynamics, the complex behavioral patterns  
1799 of vehicles and the need to periodically adjust the tolls in the traffic network by one  
1800 or more toll agents. Reinforcement Learning is a promising solution to tackle these  
1801 challenges as it lends itself well to model the interaction between traffic network and  
1802 toll agents. It provides algorithms, both for single and multi agent cases, to optimize  
1803 the performance of the traffic network to the benefit of vehicles and society. In this  
1804 work we presented a traffic model, relevant elements of Reinforcement Learning and  
1805 the different ways Reinforcement Learning can be used for effective toll pricing. While  
1806 the current solutions in the works make the potential of RL abundantly clear, research  
1807 challenges remain before solutions are ready for implementation.

1808 Based on our analysis, partial observability, non-stationarity and credit assign-  
1809 ment were partially addressed in the surveyed works. While the existing literature  
1810 in Reinforcement Learning on addressing partial observability, non-stationarity, and  
1811 credit assignment challenges provides valuable insights, opportunities remain to more  
1812 effectively leverage these in the context of road pricing. Scalability and generalizabil-  
1813 ity, critical challenges in Reinforcement Learning, and real-world experiments received  
1814 limited attention in the surveyed works. We also emphasized the need for a benchmark  
1815 to enable comparison between solution approaches. Finally, specific challenges related  
1816 to interaction between heterogeneous agents would benefit from more research in this  
1817 domain.

1818

## 1819 Declarations

1820

## 1821 Funding

1822 No funding was received for conducting this study.

1823

## 1824 Competing interests

1826 The authors have no competing interests to declare that are relevant to the content  
1827 of this article.

1828

## 1829 Acknowledgements

1830

1831 We would like to thank the editors and four anonymous reviewers for providing their  
1832 constructive feedback. Without their suggestions this manuscript would not have  
1833 looked as in this final version.

1834

## 1835 References

1836

1837 Albrecht, S.V., Christianos, F., Schäfer, L.: Multi-Agent Reinforcement Learn-  
1838 ing: Foundations and Modern Approaches. MIT Press, Cambridge (2024).  
1839 <https://www.marl-book.com>

1840



Arnott, R., De Palma, A., Lindsey, R.: Economics of a bottleneck. <i>Journal of Urban Economics</i> <b>27</b> (1), 111–130 (1990)	1841 1842 1843
Arnott, R., De Palma, A., Lindsey, R.: A structural model of peak-period congestion: A traffic bottleneck with elastic demand. <i>The American Economic Review</i> , 161–179 (1993)	1844 1845 1846 1847
Ahmed, H.U., Huang, Y., Lu, P.: A review of car-following models and modeling tools for human and autonomous-ready driving behaviors in micro-simulation. <i>Smart Cities</i> <b>4</b> (1), 314–335 (2021)	1848 1849 1850
Arnott, R.: Congestion tolling with agglomeration externalities. <i>Journal of Urban Economics</i> <b>62</b> (2), 187–203 (2007)	1851 1852 1853
Busoniu, L., Babuska, R., De Schutter, B.: A comprehensive survey of multiagent reinforcement learning. <i>IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)</i> <b>38</b> (2), 156–172 (2008)	1854 1855 1856 1857
Bellemare, M.G., Dabney, W., Munos, R.: A distributional perspective on reinforcement learning. In: <i>International Conference on Machine Learning</i> , pp. 449–458 (2017). PMLR	1858 1859 1860 1861
Bono, G., Dibangoye, J.S., Matignon, L., Pereyron, F., Simonin, O.: Cooperative multi-agent policy gradient. In: <i>Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018, Proceedings, Part I</i> 18, pp. 459–476 (2019). Springer	1862 1863 1864 1865 1866
Berto, F., Hua, C., Park, J., Kim, M., Kim, H., Son, J., Kim, H., Kim, J., Park, J.: RL4CO: a unified reinforcement learning for combinatorial optimization library. In: <i>NeurIPS 2023 Workshop: New Frontiers in Graph Learning</i> (2023). <a href="https://openreview.net/forum?id=YXSJxi8dOV">https://openreview.net/forum?id=YXSJxi8dOV</a>	1867 1868 1869 1870 1871
Boyles, S.D., Lownes, N.E., Unnikrishnan, A.: <i>Transportation Network Analysis</i> vol. 1. University of Texas, Austin (2023). edition 0.91	1872 1873 1874
Barth-Maron, G., Hoffman, M.W., Budden, D., Dabney, W., Horgan, D., Dhruva, T., Muldal, A., Heess, N., Lillicrap, T.: Distributed distributional deterministic policy gradients. In: <i>International Conference on Learning Representations</i> (2018)	1875 1876 1877
Bogrybayeva, A., Meraliyev, M., Mustakhov, T., Dauletbayev, B.: Machine learning to solve vehicle routing problems: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i> (2024)	1878 1879 1880 1881
Beckmann, M., McGuire, C.B., Winsten, C.B.: <i>Studies in the economics of transportation</i> . Technical report, Cowles Commission for Research in Economics, New Haven (1956)	1882 1883 1884 1885 1886

1887 Bellemare, M.G., Naddaf, Y., Veness, J., Bowling, M.: The arcade learning environ-  
1888 ment: An evaluation platform for general agents. *Journal of Artificial Intelligence*  
1889 *Research* **47**, 253–279 (2013)

1890

1891 Bettini, M., Prorok, A., Moens, V.: Benchmarl: Benchmarking multi-agent reinforce-  
1892 ment learning. *Journal of Machine Learning Research* **25**(217), 1–10 (2024)

1893

1894 Chen, H., An, B., Sharon, G., Hanna, J., Stone, P., Miao, C., Soh, Y.: Dyetc: Dynamic  
1895 electronic toll collection for traffic congestion alleviation. In: *Proceedings of the*  
1896 *AAAI Conference on Artificial Intelligence*, vol. 32 (2018)

1897

1898 Claus, C., Boutilier, C.: The dynamics of reinforcement learning in cooperative  
1899 multiagent systems. *AAAI/IAAI* **1998**(746-752), 2 (1998)

1900

1901 Chiu, Y.-C., Bottom, J., Mahut, M., Paz, A., Balakrishna, R., Waller, T., Hicks, J.:  
1902 Dynamic traffic assignment: A primer. *Transportation Research Circular (E-C153)*  
1903 (2011)

1904

1905 Cole, R., Dodis, Y., Roughgarden, T.: Pricing network edges for heterogeneous selfish  
1906 users. In: *Proceedings of the Thirty-fifth Annual ACM Symposium on Theory of*  
1907 *Computing*, pp. 521–530 (2003)

1908

1909 Chiou, S.-W.: A knowledge-assisted reinforcement learning optimization for road net-  
1910 work design problems under uncertainty. *Knowledge-Based Systems* **292**, 111614  
1911 (2024)

1912

1913 Como, G., Maggistro, R.: Distributed dynamic pricing of multiscale transportation  
1914 networks. *IEEE Transactions on Automatic Control* **67**(4), 1625–1638 (2021)

1915

1916 Chowdhury, D., Santen, L., Schadschneider, A.: Statistical physics of vehicular traffic  
1917 and some related systems. *Physics Reports* **329**(4-6), 199–329 (2000)

1918

1919 Chakravarty, S., Tanveer, M.H., Voicu, R.C., Banerjee, M.: Optimal-tolling using  
1920 reinforcement learning. In: *SoutheastCon 2024*, pp. 1317–1321 (2024). *IEEE*

1921

1922 Daganzo, C.F.: The cell transmission model: A dynamic representation of highway  
1923 traffic consistent with the hydrodynamic theory. *Transportation Research Part B: Methodological* **28**(4), 269–287 (1994)

1924

1925 Daganzo, C.F.: The cell transmission model, part ii: network traffic. *Transportation*  
1926 *Research Part B: Methodological* **29**(2), 79–93 (1995)

1927

1928 Dulac-Arnold, G., Levine, N., Mankowitz, D.J., Li, J., Paduraru, C., Goyal, S., Hester,  
1929 T.: Challenges of real-world reinforcement learning: definitions, benchmarks and  
1930 analysis. *Machine Learning* **110**(9), 2419–2468 (2021)

1931

1932 Du, W., Ding, S.: A survey on multi-agent deep reinforcement learning: from the  
1933 perspective of challenges and applications. *Artificial Intelligence Review* **54**(5),

3215–3238 (2021)	1933
Diallo, A.O., Lozenguez, G., Doniec, A., Mandiau, R.: Comparative evaluation of road traffic simulators based on modeler’s specifications: An application to intermodal mobility behaviors. In: ICAART (1), pp. 265–272 (2021)	1934 1935 1936 1937 1938
Dwivedi, V.P., Luu, A.T., Laurent, T., Bengio, Y., Bresson, X.: Graph neural networks with learnable structural and positional representations. In: International Conference on Learning Representations (2022). <a href="https://openreview.net/forum?id=wTTjnvGphYj">https://openreview.net/forum?id=wTTjnvGphYj</a>	1939 1940 1941 1942
De Palma, A., Lindsey, R.: Traffic congestion pricing methodologies and technologies. <i>Transportation Research Part C: Emerging Technologies</i> <b>19</b> (6), 1377–1399 (2011)	1943 1944 1945
Degrís, T., Pilarski, P.M., Sutton, R.S.: Model-free reinforcement learning with continuous action in practice. In: 2012 American Control Conference (ACC), pp. 2177–2182 (2012). IEEE	1946 1947 1948 1949
Dafermos, S.C., Sparrow, F.T.: The traffic assignment problem for a general network. <i>Journal of Research of the National Bureau of Standards B</i> <b>73</b> (2), 91–118 (1969)	1950 1951 1952
Degrís, T., White, M., Sutton, R.S.: Off-policy actor-critic. <i>Proceedings of the 29<sup>th</sup> International Conference on Machine Learning</i> (2012)	1953 1954 1955
Ekström, J.: Finding second-best toll locations and levels by relaxing the set of first-best feasible toll vectors. <i>European Journal of Transport and Infrastructure Research</i> <b>14</b> (1) (2014)	1956 1957 1958
Eliasson, J.: Congestion pricing. In: <i>The Routledge Handbook of Transport Economics</i> , pp. 209–226. Routledge, New York and London (2017)	1959 1960 1961
Foerster, J.N.: Deep multi-agent reinforcement learning. PhD thesis, University of Oxford (2018)	1962 1963 1964
Farazi, N.P., Zou, B., Ahamed, T., Barua, L.: Deep reinforcement learning in transportation research: A review. <i>Transportation research interdisciplinary perspectives</i> <b>11</b> , 100425 (2021)	1965 1966 1967 1968
Farias, A.V., Zhu, S., Mardan, A.: An overview of dynamic pricing toll roads in the united states: Pricing algorithms, operation strategies, equity concerns, and funding mechanism. <i>Case Studies on Transport Policy</i> , 101226 (2024)	1969 1970 1971 1972
Gronauer, S., Diepold, K.: Multi-agent deep reinforcement learning: a survey. <i>Artificial Intelligence Review</i> <b>55</b> (2), 895–943 (2022)	1973 1974 1975
Gupta, J.K., Egorov, M., Kochenderfer, M.: Cooperative multi-agent control using deep reinforcement learning. In: <i>Autonomous Agents and Multiagent Systems: AAMAS 2017 Workshops, Best Papers, São Paulo, Brazil, May 8-12, 2017, Revised</i>	1976 1977 1978

1979 Selected Papers 16, pp. 66–83 (2017). Springer

1980

1981 Genser, A., Kouvelas, A.: Dynamic congestion pricing for multi-region networks: A

1982 traffic equilibria approach. In: 19th Swiss Transport Research Conference (STRC

1983 2019) (2019). STRC

1984

1985 Genser, A., Kouvelas, A.: Dynamic optimal congestion pricing in multi-region urban

1986 networks by application of a multi-layer-neural network. *Transportation Research*

1987 *Part C: Emerging Technologies* **134**, 103485 (2022)

1988

1989 Hoogendoorn, S.P., Bovy, P.H.: State-of-the-art of vehicular traffic flow modelling.

1990 *Proceedings of the Institution of Mechanical Engineers, Part I: Journal of Systems*

1991 *and Control Engineering* **215**(4), 283–303 (2001)

1992

1993 Hernandez-Leal, P., Kaisers, M., Baarslag, T., De Cote, E.M.: A survey of learn-

1994 ing in multiagent environments: Dealing with non-stationarity. *arXiv preprint*

1995 *arXiv:1707.09183* (2017)

1996

1997 Hernandez-Leal, P., Kartal, B., Taylor, M.E.: A survey and critique of multiagent

1998 deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems* **33**(6),

1999 750–797 (2019)

2000

2001 He, Q., Ma, M., Li, C., Liu, W.: Learning and managing stochastic network traf-

2002 fic dynamics: an iterative and interactive approach. *Transportmetrica B: transport*

2003 *dynamics* **12**(1), 2303050 (2024)

2004

2005 Haarnoja, T., Tang, H., Abbeel, P., Levine, S.: Reinforcement learning with deep

2006 energy-based policies. In: *International Conference on Machine Learning*, pp. 1352–

2007 1361 (2017). PMLR

2008

2009 Han, Y., Wang, M., Leclercq, L.: Leveraging reinforcement learning for dynamic traf-

2010 fic control: A survey and challenges for field implementation. *Communications in*

2011 *Transportation Research* **3**, 100104 (2023)

2012

2013 Haydari, A., Yilmaz, Y.: Deep reinforcement learning for intelligent transportation

2014 systems: A survey. *IEEE Transactions on Intelligent Transportation Systems* (2020)

2015

2016 Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maxi-

2017 mum entropy deep reinforcement learning with a stochastic actor. In: *International*

2018 *Conference on Machine Learning*, pp. 1861–1870 (2018). PMLR

2019

2020 Joksimovic, D., Bliemer, M.C., Bovy, P.H.: Optimal toll design problem in dynamic

2021 traffic networks with joint route and departure time choice. *Transportation Research*

2022 *Record* **1923**(1), 61–72 (2005)

2023

2024 Jin, J., Zhu, X., Wu, B., Zhang, J., Wang, Y.: A dynamic and deadline-oriented road

pricing mechanism for urban traffic management. *Tsinghua Science and Technology*

27(1), 91–102 (2021)	2025
Kakade, S.M.: A natural policy gradient. <i>Advances in neural information processing systems</i> <b>14</b> (2001)	2026
Kaelbling, L.P., Littman, M.L., Cassandra, A.R.: Planning and acting in partially observable stochastic domains. <i>Artificial intelligence</i> <b>101</b> (1-2), 99–134 (1998)	2027
Kaelbling, L.P., Littman, M.L., Moore, A.W.: Reinforcement learning: A survey. <i>Journal of artificial intelligence research</i> <b>4</b> , 237–285 (1996)	2028
Knight, F.H.: Some fallacies in the interpretation of social cost. <i>The Quarterly Journal of Economics</i> <b>38</b> (4), 582–606 (1924)	2029
Knorr, F.: Applicability and application of microscopic traffic simulations. PhD thesis, Universitätsbibliothek Duisburg-Essen (2013)	2030
Kumar, N., Raubal, M.: Applications of deep learning in congestion detection, prediction and alleviation: A survey. <i>Transportation Research Part C: Emerging Technologies</i> <b>133</b> , 103432 (2021)	2031
Kohl, N., Stone, P.: Policy gradient reinforcement learning for fast quadrupedal locomotion. In: <i>IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA'04. 2004</i> , vol. 3, pp. 2619–2624 (2004). IEEE	2032
Kok, J.R., Vlassis, N.: Collaborative multiagent reinforcement learning by payoff propagation. <i>Journal of Machine Learning Research</i> <b>7</b> , 1789–1828 (2006)	2033
Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: <i>International Conference on Learning Representations</i> (2017). <a href="https://openreview.net/forum?id=SJU4ayYgl">https://openreview.net/forum?id=SJU4ayYgl</a>	2034
Lopez, P.A., Behrisch, M., Bieker-Walz, L., Erdmann, J., Flötteröd, Y.-P., Hilbrich, R., Lücken, L., Rummel, J., Wagner, P., Wießner, E.: Microscopic traffic simulation using sumo. In: <i>2018 21st International Conference on Intelligent Transportation Systems (ITSC)</i> , pp. 2575–2582 (2018). IEEE	2035
Liang, C., Huang, Z., Liu, Y., Liu, Z., Zheng, G., Shi, H., Du, Y., Li, F., Li, Z.: Cblab: Scalable traffic simulation with enriched data supporting. <i>arXiv preprint arXiv:2210.00896</i> (2022)	2036
Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. <i>arXiv preprint arXiv:1509.02971</i> (2015)	2037
Lu, J., Hong, C., Wang, R.: Magt-toll: A multi-agent reinforcement learning approach to dynamic traffic congestion pricing. <i>PloS one</i> <b>19</b> (11), 0313828 (2024)	2038

2071 Li, Y.: Deep reinforcement learning: An overview. arXiv preprint arXiv:1701.07274  
 2072 (2017)  
 2073  
 2074 Littman, M.L.: Markov games as a framework for multi-agent reinforcement learning.  
 2075 In: Machine Learning Proceedings 1994, pp. 157–163. Elsevier, Amsterdam (1994)  
 2076  
 2077 Levine, S., Kumar, A., Tucker, G., Fu, J.: Offline reinforcement learning: Tutorial,  
 2078 review, and perspectives on open problems. arXiv preprint arXiv:2005.01643 (2020)  
 2079  
 2080 Lombardi, C., Picado-Santos, L., Annaswamy, A.M.: Model-based dynamic toll  
 2081 pricing: An overview. Applied Sciences **11**(11), 4778 (2021)  
 2082  
 2083 Lindsey, R., Verhoef, E.: Traffic congestion and congestion pricing. In: Handbook of  
 2084 Transport Systems and Traffic Control. Emerald Group Publishing Limited, United  
 2085 Kingdom (2001)  
 2086  
 2087 Lighthill, M.J., Whitham, G.B.: On kinematic waves ii. a theory of traffic flow on long  
 2088 crowded roads. Proceedings of the Royal Society of London. Series A. Mathematical  
 2089 and Physical Sciences **229**(1178), 317–345 (1955)  
 2090  
 2091 Mnih, V., Badia, A.P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver,  
 2092 D., Kavukcuoglu, K.: Asynchronous methods for deep reinforcement learning. In:  
 2093 International Conference on Machine Learning, pp. 1928–1937 (2016). PMLR  
 2094  
 2095 Maheshwari, C., Kulkarni, K., Pai, D., Yang, J., Wu, M., Sastry, S.: Congestion pricing  
 2096 for efficiency and equity: Theory and applications to the san francisco bay area.  
 2097 arXiv preprint arXiv:2401.16844 (2024)  
 2098  
 2099 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G.,  
 2100 Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., *et al.*: Human-level  
 2101 control through deep reinforcement learning. nature **518**(7540), 529–533 (2015)  
 2102  
 2103 Mirzaei, H., Sharon, G., Boyles, S., Givargis, T., Stone, P.: Enhanced delta-tolling:  
 2104 Traffic optimization via policy gradient reinforcement learning. In: 2018 21st Inter-  
 2105 national Conference on Intelligent Transportation Systems (ITSC), pp. 47–52  
 2106 (2018). IEEE  
 2107  
 2108 Mirzaei, H., Sharon, G., Boyles, S., Givargis, T., Stone, P.: Link-based parameter-  
 2109 ized micro-tolling scheme for optimal traffic management. In: Proceedings of the  
 2110 17th International Conference on Autonomous Agents and MultiAgent Systems, pp.  
 2111 2013–2015 (2018)  
 2112  
 2113 Mokbel, M., Sakr, M., Xiong, L., Züfle, A., Almeida, J., Anderson, T., Aref, W.,  
 2114 Andrienko, G., Andrienko, N., Cao, Y., et al.: Mobility data science: Perspectives  
 2115 and challenges. ACM Transactions on Spatial Algorithms and Systems (2024)  
 2116

Marbach, P., Tsitsiklis, J.N.: Simulation-based optimization of markov reward processes. <i>IEEE Transactions on Automatic Control</i> <b>46</b> (2), 191–209 (2001)	2117 2118 2119
Nash Jr, J.F.: Equilibrium points in n-person games. <i>Proceedings of the National Academy of Sciences</i> <b>36</b> (1), 48–49 (1950)	2120 2121 2122
Nguyen, J., Powers, S.T., Urquhart, N., Farrenkopf, T., Guckert, M.: An overview of agent-based traffic simulators. <i>Transportation research interdisciplinary perspectives</i> <b>12</b> , 100486 (2021)	2123 2124 2125 2126
Nisha, A.S.A., Rao, N.V., Venkatesh, G., Murugan, S., Meenakshi, B., <i>et al.</i> : Efficient congestion management through iot-driven road user charging systems with reinforcement learning. In: <i>2024 Second International Conference on Intelligent Cyber Physical Systems and Internet of Things (ICoICI)</i> , pp. 431–436 (2024). IEEE	2127 2128 2129 2130
Nohekhan, A., Zahedian, S., Sadabadi, K.F.: Investigating the impacts of i-66 inner beltway dynamic tolling system. <i>Transportation Engineering</i> <b>4</b> , 100059 (2021)	2131 2132 2133
Padakandla, S.: A survey of reinforcement learning algorithms for dynamically varying environments. <i>ACM Computing Surveys (CSUR)</i> <b>54</b> (6), 1–25 (2021)	2134 2135 2136
Pandey, V., Boyles, S.D.: Multiagent reinforcement learning algorithm for distributed dynamic pricing of managed lanes. In: <i>2018 21st International Conference on Intelligent Transportation Systems (ITSC)</i> , pp. 2346–2351 (2018). IEEE	2137 2138 2139 2140
Pandey, V., Boyles, S.D.: Comparing route choice models for managed lane networks with multiple entrances and exits. <i>Transportation Research Record</i> <b>2673</b> (10), 381–393 (2019)	2141 2142 2143 2144
Paccagnan, D., Chandan, R., Ferguson, B.L., Marden, J.R.: Optimal taxes in atomic congestion games. <i>ACM Transactions on Economics and Computation (TEAC)</i> <b>9</b> (3), 1–33 (2021)	2145 2146 2147 2148
Prieto Curiel, R., González Ramírez, H., Quiñones Domínguez, M., Orjuela Mendoza, J.P.: A paradox of traffic and extra cars in a city as a collective behaviour. <i>Royal Society open science</i> <b>8</b> (6), 201808 (2021)	2149 2150 2151 2152
Pigou, A.C.: <i>The Economics of Welfare</i> . Macmillan, London (1924)	2153 2154
Patterson, A., Neumann, S., White, M., White, A.: Empirical design in reinforcement learning. <i>Journal of Machine Learning Research</i> <b>25</b> (318), 1–63 (2024)	2155 2156
Peters, J., Schaal, S.: Policy gradient methods for robotics. In: <i>2006 IEEE/RSJ International Conference on Intelligent Robots and Systems</i> , pp. 2219–2225 (2006). IEEE	2157 2158 2159 2160
Pandey, V., Wang, E., Boyles, S.D.: Deep reinforcement learning algorithm for dynamic pricing of express lanes with multiple access locations. <i>Transportation</i>	2161 2162



2163 Research Part C: Emerging Technologies **119**, 102715 (2020)  
2164  
2165 Qiu, W., Chen, H., An, B.: Dynamic electronic toll collection via multi-agent deep  
2166 reinforcement learning with edge-based graph convolutional networks. In: IJCAI,  
2167 pp. 4568–4574 (2019)  
2168  
2169 Qin, Z.T., Zhu, H., Ye, J.: Reinforcement learning for ridesharing: An extended survey.  
2170 Transportation Research Part C: Emerging Technologies **144**, 103852 (2022)  
2171  
2172 Ramos, G.d.O., Da Silva, B.C., Rădulescu, R., Bazzan, A.L., Nowé, A.: Toll-based  
2173 reinforcement learning for efficient equilibria in route choice. The Knowledge  
2174 Engineering Review **35** (2020)  
2175  
2176 Ramos, G.d.O., Silva, B.C., Rădulescu, R., Bazzan, A.L.: Learning system-efficient  
2177 equilibria in route choice using tolls. In: Proceedings of the Adaptive Learning  
2178 Agents Workshop, vol. 2018 (2018)  
2179  
2180 Richards, P.I.: Shock waves on the highway. Operations research **4**(1), 42–51 (1956)  
2181  
2182 Ramos, G.d.O., Rădulescu, R., Nowé, A., Tavares, A.R.: Toll-based learning for min-  
2183 imising congestion under heterogeneous preferences. In: Proceedings of the 19th  
2184 International Conference on Autonomous Agents and Multiagent Systems (AAMAS  
2185 2020), pp. 1098–1106 (2020). IFAAMAS  
2186  
2187 Sutton, R.S., Barto, A.G.: Reinforcement Learning: An Introduction. MIT Press,  
2188 Cambridge (2018)  
2189  
2190 Saharan, S., Bawa, S., Kumar, N.: Dynamic pricing techniques for intelligent trans-  
2191 portation system in smart cities: A systematic review. Computer Communications  
2192 **150**, 603–625 (2020)  
2193  
2194 Schmidt, L.M., Brosig, J., Plinge, A., Eskofier, B.M., Mutschler, C.: An introduction  
2195 to multi-agent reinforcement learning and review of its application to autonomous  
2196 mobility. In: 2022 IEEE 25th International Conference on Intelligent Transportation  
2197 Systems (ITSC), pp. 1342–1349 (2022). IEEE  
2198  
2199 Storani, F., Di Pace, R., Bruno, F., Fiori, C.: Analysis and comparison of traffic flow  
2200 models: a new hybrid traffic flow model vs benchmark models. European Transport  
2201 Research Review **13**(1), 1–16 (2021)  
2202  
2203 Seo, T.: Trial-and-error congestion pricing scheme for morning commute problem with  
2204 day-to-day dynamics. Transportation Research Procedia **47**, 561–568 (2020)  
2205  
2206 Sharon, G., Hanna, J., Rambha, T., Albert, M., Stone, P., Boyles, S.D.: Delta-tolling:  
2207 Adaptive tolling for optimizing traffic throughput. In: ATT@ IJCAI (2016)  
2208  
2209 Sharon, G., Hanna, J.P., Rambha, T., Levin, M.W., Albert, M., Boyles, S.D., Stone,  
2210 P.: Real-time adaptive tolling scheme for optimized social welfare in traffic networks.

In: Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS-2017) (2017)	2209
	2210
	2211
Singh, S., Kearns, M.J., Mansour, Y.: Nash convergence of gradient dynamics in general-sum games. In: UAI, pp. 541–548 (2000)	2212
	2213
	2214
Shoham, Y., Leyton-Brown, K.: Multiagent Systems: Algorithmic, Game-theoretic, and Logical Foundations. Cambridge University Press, Cambridge (2008)	2215
	2216
	2217
Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., Riedmiller, M.: Deterministic policy gradient algorithms. In: International Conference on Machine Learning, pp. 387–395 (2014). PMLR	2218
	2219
	2220
	2221
Sharon, G., Levin, M.W., Hanna, J.P., Rambha, T., Boyles, S.D., Stone, P.: Network-wide adaptive tolling for connected and automated vehicles. Transportation Research Part C: Emerging Technologies <b>84</b> , 142–157 (2017)	2222
	2223
	2224
Sutton, R.S., McAllester, D., Singh, S., Mansour, Y.: Policy gradient methods for reinforcement learning with function approximation. Advances in neural information processing systems <b>12</b> (1999)	2225
	2226
	2227
	2228
Schuster, P., Sigmund, K.: Replicator dynamics. Journal of theoretical biology <b>100</b> (3), 533–538 (1983)	2229
	2230
	2231
Sato, K., Seo, T., Fuse, T.: A reinforcement learning-based dynamic congestion pricing method for the morning commute problems. Transportation Research Procedia <b>52</b> , 347–355 (2021)	2232
	2233
	2234
	2235
Sato, K., Seo, T., Fuse, T.: Dynamic network congestion pricing based on deep reinforcement learning. arXiv preprint arXiv:2206.12188 (2022)	2236
	2237
	2238
Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)	2239
	2240
	2241
Tan, M.: Multi-agent reinforcement learning: Independent vs. cooperative agents. In: Proceedings of the Tenth International Conference on Machine Learning, pp. 330–337 (1993)	2242
	2243
	2244
	2245
Tavares, A.R., Bazzan, A.L.: An agent-based approach for road pricing: system-level performance and implications for drivers. Journal of the Brazilian Computer Society <b>20</b> (1), 1–15 (2014)	2246
	2247
	2248
Terry, J., Black, B., Grammel, N., Jayakumar, M., Hari, A., Sullivan, R., Santos, L.S., Dieffendahl, C., Horsch, C., Perez-Vicente, R., <i>et al.</i> : Pettingzoo: Gym for multi-agent reinforcement learning. Advances in Neural Information Processing Systems <b>34</b> , 15032–15043 (2021)	2249
	2250
	2251
	2252
	2253
Tampère, C.M., Corthout, R., Cattrysse, D., Immers, L.H.: A generic class of first order	2254

node models for dynamic macroscopic simulation of traffic flows. *Transportation Research Part B: Methodological* **45**(1), 289–309 (2011)

Todorov, E., Erez, T., Tassa, Y.: Mujoco: A physics engine for model-based control. In: 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 5026–5033 (2012). IEEE

Treiber, M., Hennecke, A., Helbing, D.: Congested traffic states in empirical observations and microscopic simulations. *Phys. Rev. E* **62**, 1805–1824 (2000) <https://doi.org/10.1103/PhysRevE.62.1805>

Towers, M., Kwiatkowski, A., Terry, J., Balis, J.U., De Cola, G., Deleu, T., Goulão, M., Kallinteris, A., Krimmel, M., KG, A., et al.: Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032* (2024)

Tsekeris, T., Voß, S.: Design and evaluation of road pricing: state-of-the-art and methodological advances. *NETNOMICS: Economic Research and Electronic Networking* **10**, 5–52 (2009)

United States Bureau of Public Roads: Traffic Assignment Manual for Application with a Large, High Speed Computer vol. 2. US Department of Commerce, Bureau of Public Roads, Office of Planning, Urban Planning Division, Washington D.C. (1964)

Verhoef, E.T.: Time, speeds, flows and densities in static models of road traffic congestion and congestion pricing. *Regional Science and Urban Economics* **29**(3), 341–369 (1999)

Verhoef, E.T.: Second-best congestion pricing in general networks-algorithms for finding second-best optimal toll levels and toll points. Technical report, Tinbergen Institute Discussion Paper (2000)

Verhoef, E.T.: Second-best congestion pricing in general static transportation networks with elastic demands. *Regional Science and Urban Economics* **32**(3), 281–310 (2002)

Vickrey, W.S.: Pricing in urban and suburban transport. *The American Economic Review* **53**(2), 452–465 (1963)

Vickrey, W.S.: Congestion theory and transport investment. *The American Economic Review* **59**(2), 251–260 (1969)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)

Wageningen-Kessels, F., Van Lint, H., Vuik, K., Hoogendoorn, S.: Genealogy of traffic flow models. *EURO Journal on Transportation and Logistics* **4**(4), 445–473 (2015)

Wardrop, J.G.: Road paper. some theoretical aspects of road traffic research. Proceedings of the institution of civil engineers <b>1</b> (3), 325–362 (1952)	2301 2302 2303
Watkins, C.J.C.H.: Learning from delayed rewards. PhD thesis, King’s College, Cambridge United Kingdom (1989)	2304 2305 2306
Wong, A., Bäck, T., Kononova, A.V., Plaat, A.: Deep multiagent reinforcement learning: Challenges and directions. Artificial Intelligence Review, 1–34 (2022)	2307 2308 2309
Watkins, C.J., Dayan, P.: Q-learning. Machine learning <b>8</b> (3), 279–292 (1992)	2310 2311
Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine learning <b>8</b> (3), 229–256 (1992)	2312 2313 2314
Wang, Y., Jin, H., Zheng, G.: Ctrl: Cooperative traffic tolling via reinforcement learning. In: Proceedings of the 31st ACM International Conference on Information & Knowledge Management, pp. 3545–3554 (2022)	2315 2316 2317
Wen, M., Kuba, J., Lin, R., Zhang, W., Wen, Y., Wang, J., Yang, Y.: Multi-agent reinforcement learning is a sequence modeling problem. Advances in Neural Information Processing Systems <b>35</b> , 16509–16521 (2022)	2318 2319 2320 2321
Wei, H., Zheng, G., Gayah, V., Li, Z.: Recent advances in reinforcement learning for traffic signal control: A survey of models and evaluation. ACM SIGKDD Explorations Newsletter <b>22</b> (2), 12–18 (2021)	2322 2323 2324 2325
Yan, Y., Chow, A.H., Ho, C.P., Kuo, Y.-H., Wu, Q., Ying, C.: Reinforcement learning for logistics and supply chain management: Methodologies, state of the art, and future opportunities. Transportation Research Part E: Logistics and Transportation Review <b>162</b> , 102712 (2022)	2326 2327 2328 2329
Yang, H., Huang, H.-J.: Principle of marginal-cost pricing: how does it work in a general road network? Transportation Research Part A: Policy and Practice <b>32</b> (1), 45–54 (1998)	2330 2331 2332 2333
Ye, G., Song, J., Feng, M., Zhu, G., Shen, P., Zhang, L., Shah, S.A.A., Bennamoun, M.: Position and structure-aware graph learning. Neurocomputing <b>556</b> , 126581 (2023)	2334 2335 2336
Yuan, L., Zhang, Z., Li, L., Guan, C., Yu, Y.: A survey of progress on cooperative multi-agent reinforcement learning in open environment. arXiv preprint arXiv:2312.01058 (2023)	2337 2338 2339 2340
Zhang, H., Feng, S., Liu, C., Ding, Y., Zhu, Y., Zhou, Z., Zhang, W., Yu, Y., Jin, H., Li, Z.: Cityflow: A multi-agent reinforcement learning environment for large scale city traffic scenario. In: The World Wide Web Conference, pp. 3620–3624 (2019)	2341 2342 2343 2344
Zhu, F., Ukkusuri, S.V.: A reinforcement learning approach for distance-based dynamic tolling in the stochastic network environment. Journal of Advanced	2345 2346

2347 Transportation **49**(2), 247–266 (2015)

2348

2349 Zhang, X., Wang, W., Chen, J.: A priori lane selection strategy for reinforcement  
2350 learning of dynamic expressway tolling. In: 2023 International Conference on Pattern  
2351 Recognition, Machine Vision and Intelligent Algorithms (PRMVIA), pp. 143–154  
2352 (2023). IEEE

2353

## 2354 Appendix A Concise introduction to 2355 Reinforcement Learning

2356

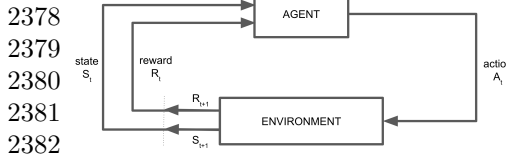
### 2357 A.1 Reinforcement Learning

2358

2359 In RL (Sutton and Barto, 2018), an agent interacts dynamically with an environment  
2360 at each of a sequence of discrete episode steps  $t = 0, 1, 2, 3, \dots$ . RL is about learning  
2361 from these interactions to achieve a goal. The agent is the learner and decision maker,  
2362 the environment is everything else. Framing the problem of learning from interaction to  
2363 achieve a goal can be done with a Markov Decision Process. A Markov Decision Process  
2364 (MDP) is formally defined as a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$  where  $\mathcal{S}$  is the set of all states  $s$ , a  
2365 complete description of the environment.  $\mathcal{A}$  is the set of all actions  $a$ ,  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow$   
2366  $[0, 1]$  is the transition probability function mapping from any state  $s \in \mathcal{S}$  to any next  
2367 state  $s' \in \mathcal{S}$  after execution of  $a \in \mathcal{A}$ , also denoted by  $\mathcal{P}(s'|s, a)$ .  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is  
2368 the reward function that provides the reward  $r$ , also denoted by  $\mathcal{R}(s, a, s')$ . The agent  
2369 seeks to maximize the total amount of reward over time through its choice of actions  
2370 using a policy  $\pi$ . Policy  $\pi(a|s)$  gives the probability of taking action  $a$  given state  
2371  $s$ . Policies can be stochastic ( $\pi(a|s)$ ) or deterministic ( $\pi(s)$ ). The MDP and agent,  
2372 following the policy  $\pi$ , give rise to a sequence of states  $S_t = s$ , actions  $A_t = a$  and  
2373 rewards  $R_t = r$  which is called a trajectory  $T$ :  $T = S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots$   
2374 An episode is a trajectory  $T$  that has reached a final episode step  $H$ . See Figure  
2375 A1 that illustrates agent environment interaction. The probability of a trajectory  $T$ ,

2376

2377



2383 **Fig. A1:** Agent environment interac-  
2384 tion in a MDP (Sutton and Barto,  
2385 2018)

2386

2387

2388 following policy  $\pi$  and  $S_0$  sampled from the start state distribution, is:

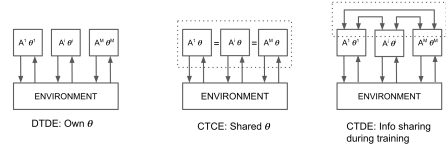
2389

2390

2391

2392

$$p(T|\pi) = p(S_0) \prod_{t=0}^{t=H-1} \pi(A_t|S_t) p(S_{t+1}|S_t, A_t). \quad (\text{A1})$$



2393 **Fig. A2:** Approaches to cooperative  
2394 MARL (Gronauer and Diepold, 2022;  
2395 Wong et al., 2022)

The return  $G_t$  after episode step  $t$  is defined as the cumulative discounted reward over time:

$$G_t = \sum_{i=t+1}^H \gamma^{i-t-1} R_i = R_{t+1} + \gamma G_{t+1}, \quad (\text{A2})$$

where  $\gamma \in [0, 1]$  is the discount factor, reflecting the importance of future rewards. The goal of the agent is defined to maximize the total expected discounted return per episode:

$$J = \mathbb{E}_{T \sim p(T)} G_0(T). \quad (\text{A3})$$

Any policy can be used by an agent for execution in an environment. RL algorithms specify how to find an optimal policy  $\pi_*$  using the experience gained by the agent. In such an algorithm, the agent's experience consists of one or more episodes. On the basis of this experience, the agent can learn (an estimate of) the optimal policy. Where the agent generates its own data we refer to this as online learning, while off-line learning refers to the situation where the data is collected by another policy (Levine et al., 2020). Experiences in Partially Observable MDPs, where the agent does not have perfect state information, lead to partial observations (Kaelbling et al., 1998, 1996). In that case we need and use a set of observations  $\mathcal{O}$  and an observation probability function  $\mathcal{Z} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{O}$ .

Value Functions. We measure the value of a policy, given a starting state, by a value function to compare the performance of different policies. The value function,  $v_\pi : \mathcal{S} \rightarrow \mathbb{R}$  is defined as the expected return when starting in state  $s$ , at any episode step  $t$  and following policy  $\pi$  thereafter. Similarly, the action value function,  $q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , also known as q-value function gives the value when taking action  $a$  when in state  $s$ , following policy  $\pi$  thereafter.

For MDPs this can be defined as:

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]. \quad (\text{A4})$$

Similarly, the action value function,  $q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , also known as q-value function or q-function, is defined as:

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a], \quad (\text{A5})$$

which gives the value when taking action  $a$  when in state  $s$ , following policy  $\pi$  thereafter.

The estimated value of a state  $S_t$  is denoted as  $V(S_t)$  and for q-value functions the estimated value is denoted as  $Q(S_t, A_t)$ . Approximate values, based on a function parameterized by weights  $\theta$ , are denoted as  $\hat{v}(s, \theta)$  or  $\hat{q}(s, a, \theta)$ . The relation between the value function  $v_\pi(s)$  and action-value function  $q_\pi(s, a)$  is that, given the same state, the value function is the expectation of the action-value function with respect to the actions:

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a). \quad (\text{A6})$$

The difference between the action-value function and value function for a state and action can be considered as an evaluation of the policy compared with the average

2439 policy and is called advantage:

2440

$$2441 \quad A(S_t, A_t) = Q(S_t, A_t) - V(S_t). \quad (A7)$$

2442

2443 To solve an RL problem usually involves estimating the value function which  
2444 describes, in terms of expected return, how good it is to be in a certain state  $s$ . We  
2445 then use the value function to find the optimal policy. If we substitute the return (A2)  
2446 in the value function (A4) we can derive:

2447

$$2448 \quad v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_\pi(s')], \quad \forall s \in \mathcal{S}. \quad (A8)$$

2449

2450  
2451 Equation (A8) is called the Bellman equation for  $v_\pi$  and provides the recursive relation  
2452 between the value of a state and its successor states. For the action-value function the  
2453 Bellman equation reads:

2454

$$2455 \quad q_\pi(s, a) = \sum_{s', r} p(s', r|s, a) \left[ r + \gamma \sum_{a'} \pi(a'|s') q_\pi(s', a') \right]. \quad (A9)$$

2456

2457  
2458 If we consider a policy  $\pi$ , for which the expected return is larger than from any  
2459 other policy  $\pi'$  for all states  $s$ , we call this the optimal policy  $\pi_*$ . The optimal value  
2460 function and optimal action-value function are defined as:  $v_*(s) = \max_\pi v_\pi(s)$  and  
2461  $q_*(s, a) = \max_\pi q_\pi(s, a)$ . If both the transition function  $\mathcal{P}$  and reward function  $\mathcal{R}$   
2462 are known, we can find the optimal policy by planning, using dynamic programming.  
2463 These are called model-based methods. However, when  $\mathcal{P}$  and  $\mathcal{R}$  are unknown, we  
2464 need to solve them using a model-free RL algorithm. These are based on learning;  
2465 interacting with the environment, obtaining rewards and states and improving the  
2466 policy based on this.

2467 Apart from strictly adhering to the policy  $\pi(a|s)$  by taking action  $a$  when in state  
2468  $s$ , greedy and  $\epsilon$ -greedy policies can be distinguished. A policy  $\pi$  is called greedy if  
2469 it always selects the action with maximum estimated action value, in which case the  
2470 agent is exploiting the current knowledge of the value of its actions. A policy  $\pi$  is  
2471 called  $\epsilon$ -greedy if the agent acts greedily most of the time, and with probability  $\epsilon$ ,  
2472 it selects an action randomly. If one of the non-greedy actions is chosen, we call this  
2473 exploring because this exposes the agent potentially to new experiences. This enables  
2474 improvement of the estimate of the non-greedy action's value.

2475

## 2476 A.2 Solution Methods to Find the Policy

2477  
2478 We distinguish three main categories: methods based on value functions, methods  
2479 directly updating the policy and methods based on a combination of these, actor-critic.  
2480

### 2481 A.2.1 Value Based Methods

2482

2483 Value based methods search the optimal policy based on the (action-)value function.  
2484 This means that the value is estimated or approximated first. Based on these values the



policy is optimized. Model free Monte Carlo (MC) methods learn value functions and optimal policies from experience by sample episodes. MC methods update estimates of values of a state  $V(S_t)$  by waiting for the sample return  $G_t$  of an episode and use that as target for  $V(S_t)$ .

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]. \quad (\text{A10})$$

where  $\alpha \in (0, 1]$  is a constant step-size parameter. Model free Temporal Difference (TD) methods update estimates, based on other learned estimates. For example, where MC needs to obtain the sample return (A2), the one-step TD method, TD(0), uses the reward  $R_{t+1}$  and the discounted estimate of the value of the next state  $\gamma V(S_{t+1})$ , to update the estimate of  $V(S_t)$ :

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]. \quad (\text{A11})$$

in which we call  $R_{t+1} + \gamma V(S_{t+1})$  the TD(0)-target and  $\delta_t$  the TD-error:

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t). \quad (\text{A12})$$

In RL,  $\delta_t$  occurs in various forms depending on the number of steps (N-step returns) taken into account to determine the target. An illustration of a value based method is the Q-learning (Watkins and Dayan, 1992) algorithm. Q-learning, is an off-policy RL algorithm. Off-policy means that the selection of actions is based on another policy than the actions used to estimate the Q-value. The update function for Q-learning is:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]. \quad (\text{A13})$$

The learned Q-function then directly approximates the optimal q-function  $q_*(s, a)$  (Watkins, 1989; Watkins and Dayan, 1992).

### A.2.2 Policy Based Methods

A policy based method (Sutton et al., 1999), directly optimizes the parameterized policy of an agent with respect to the expected return by gradient descent. The stochastic policy  $\pi(a|s, \theta)$  (hereafter denoted as  $\pi_\theta$ ) is parameterized by  $\theta$  and differentiable with respect to its parameters. Policy gradient methods learn the policy parameters based on the gradient of some performance measure  $J(\pi_\theta)$  (hereafter noted as  $J(\theta)$ ) with respect to  $\theta$ . The parameter update is expressed as:

$$\theta_{t+1} = \theta_t + \alpha \hat{\nabla} J(\theta_t), \quad (\text{A14})$$

where  $\hat{\nabla} J(\theta_t)$  is an estimate of the gradient of the performance measure, with respect to  $\theta_t$  (Sutton and Barto, 2018). If we define, in an episodic case, the performance measure as the value at the start of the episode:

$$J(\theta) = v_{\pi_\theta}(s_0), \quad (\text{A15})$$

2531 then, with  $\mu(s)$  as the on-policy distribution of states under  $\pi$  it can be derived that:

2532

$$2533 \quad \nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a|s, \boldsymbol{\theta}). \quad (\text{A16})$$

2534

2535

2536 Equation (A16) is called the policy gradient theorem (Sutton et al., 1999)(Marbach  
2537 and Tsitsiklis, 2001) for stochastic policies. A deterministic case leading to  $\pi(s, \boldsymbol{\theta})$   
2538 can be derived as well (Silver et al., 2014). A classical example of a policy gradient  
2539 algorithm is REINFORCE (Williams, 1992). From the policy gradient theorem (A16),  
2540 it can be derived that:

2541

$$2542 \quad \nabla J(\boldsymbol{\theta}) \propto \mathbb{E}_\pi \left[ G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right], \quad (\text{A17})$$

2543

2544 where the expectation is taken with respect to the policy  $\pi$ . The expression between  
2545 the brackets can be sampled on each episode step and be used to update  $\boldsymbol{\theta}$ :

2546

$$2547 \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}. \quad (\text{A18})$$

2548

2549 REINFORCE uses samples of the full return, and therefore may have a high variance  
2550 which slows learning. To address this, the policy gradient theorem can be generalized  
2551 to include a baseline function. As long as the baseline function does not vary with  
2552 actions  $a$  any function  $b(s)$  can be used:

2553

$$2554 \quad \nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a (q_\pi(s, a) - b(s)) \nabla \pi(a|s, \boldsymbol{\theta}). \quad (\text{A19})$$

2555

2556 In addition to the Policy Gradient method discussed above there are Natural Policy  
2557 Gradient (Kakade, 2001) and Finite Difference Policy Gradient (Kohl and Stone, 2004)  
2558 methods as different approaches, also directly optimizing the policy.

2559

### 2560 **A.2.3 Actor-critic Methods.**

2561

2562 A frequently used architecture, actor-critic (Degris et al., 2012)(Degris et al., 2012),  
2563 combines the value-based methods and policy-based methods described above. Actor-  
2564 critic uses the policy gradient theorem (A16) and value functions (A4, A5). In actor-  
2565 critic algorithms the actor is the learned policy  $\pi(a|s, \boldsymbol{\theta})$  and the critic the learned  
2566 value function  $\hat{v}(s, \boldsymbol{\vartheta})$ , or the action value function  $\hat{q}(s, a, \boldsymbol{\vartheta})$ . In the example below a  
2567 one-step actor-critic method uses the TD-error  $\delta_t$ :

2568

$$2569 \quad \delta_t = r' + \gamma \hat{v}(s', \boldsymbol{\vartheta}) - \hat{v}(s, \boldsymbol{\vartheta}). \quad (\text{A20})$$

2570

2571 Replacing  $q_\pi(s, a)$  in (A16), with the TD-error (A20) leads to:

2572

$$2573 \quad \nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a (r' + \gamma \hat{v}(s', \boldsymbol{\vartheta}) - \hat{v}(s, \boldsymbol{\vartheta})) \nabla \pi(a|s; \boldsymbol{\theta}). \quad (\text{A21})$$

2574

The learning method to update the critic’s parameters  $\boldsymbol{\vartheta}$  can be performed using semi-gradient TD(0). Semi-gradient TD(0) takes the effect of changing the parameters  $\boldsymbol{\vartheta}_t$  on the learned state value ( $\hat{v}(S_t, \boldsymbol{\vartheta})$ ) into account, but not on the target in minimizing  $\delta_t$ .

#### A.2.4 Issues and Challenges.

Common issues in RL problems are credit assignment (which action leads to which outcome), sparse rewards (not delivering non-zero rewards frequently enough) and sample efficiency (Li, 2017). Important criteria, to evaluate the performance of algorithms, are whether they are stable (do they converge), their efficiency (how long does it take to converge), generalization (if they converge, can they be generalized) and scalability (can the solution scale efficiently with growing state and action space). Practical challenges exist as well. Due to the large number of design decisions in Reinforcement Learning experiments, much can go wrong illustrated by a listing of 20 common errors (Patterson et al., 2024).

### A.3 Multi Agent Reinforcement Learning

A Multi Agent system describes a system in which multiple agents take decisions based on their own policy and interact with the environment (Gronauer and Diepold, 2022; Wong et al., 2022). If this is done in the context of RL this is called Multi Agent Reinforcement Learning (MARL). Depending on the situation and objectives, the agents may compete, cooperate or demonstrate mixed behavior and may or may not communicate with each other.

A MARL system can be formalized as a Markov Game which is a framework that generalizes MDPs from one agent to Multiple Agents (Gronauer and Diepold, 2022; Littman, 1994; Shoham and Leyton-Brown, 2008). A Markov Game is defined as a tuple  $(\mathcal{N}, \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$  where:  $\mathcal{N}$  is the set of  $N > 1$  interacting agents  $n$ ,  $\mathcal{S}$  is the set of states observed by the agents  $s$ ,  $\mathcal{A}$  is the joint action space  $\mathcal{A}^1 \times \dots \times \mathcal{A}^N$  where  $\mathcal{A}^n$  is the action set of agent  $n$ ,  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability function mapping from any state  $s \in \mathcal{S}$  to any next state  $s' \in \mathcal{S}$  after execution of  $\mathbf{a} \in \mathcal{A}$ , also expressed as  $\mathcal{P}(s'|s, \mathbf{a})$ ,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function that provides the agent specific reward  $r^n$ , also expressed as  $\mathcal{R}(s, \mathbf{a}, s', n)$ .

The policy  $\pi$  taken by agent  $n$  is denoted as  $\pi^n$ . A joint policy  $\boldsymbol{\pi}$  where  $\boldsymbol{\pi} = \{\pi^1, \dots, \pi^N\}$ , is a mapping from states  $s$  to probabilities of selecting  $\mathbf{a}$ :  $\boldsymbol{\pi}(\mathbf{a}|s)$  is the probability of  $\mathbf{A}_t = \mathbf{a}$  given  $S_t = s$ . If all agents  $n$  follow a policy  $\pi^n \in \Pi^n$  at time  $t$  in state  $s$ , then joint action  $\mathbf{a}$  leads, with respect to  $\mathcal{P}$ , to a new state  $s'$  and each agent gets individual reward  $r^n$  from  $\mathcal{R}$ . As per convention  $-n$  means all agents except  $n$  so  $\boldsymbol{\pi}^{-n} = \{\pi^1, \dots, \pi^{n-1}, \pi^{n+1}, \dots, \pi^N\}$ .

As in the Single Agent case outlined before, the objective of each agent is to optimize its return. A (Nash) equilibrium (Nash Jr, 1950) exists in the situation that the set of policies  $\boldsymbol{\pi}_*$  consisting of one policy  $\pi_*^n$  for each agent, maximizes the return for each agent:

$$J^n(\pi_*^n, \boldsymbol{\pi}_*^{-n}) \geq J^n(\pi^n, \boldsymbol{\pi}_*^{-n}) \forall \pi^n \in \Pi^n. \quad (\text{A22})$$

This means, informally, that no agent, given policies of other agents fixed, can improve by unilaterally deviating from  $\pi_*$ . We note that multiple Nash equilibria may exist with different returns, making evaluation of Multi Agent learning harder than Single Agent learning. If in MARL all agents receive the same reward  $R$  this is a fully cooperative setting; if they are encouraged to cooperate but don't get an equal reward it is called cooperative. If the sum of rewards of the agents equals zero, it is a competitive setting; each agent in this case tries to maximize its own reward and minimize the rewards of others. Mixed settings are neither cooperative nor competitive. These are also called general sum games.

2632

## 2633 A.4 Solution methods

2634

MARL algorithms encompass similar solution methods as Single Agent RL; value based (e.g. (Tan, 1993)), policy based (e.g. (Singh et al., 2000)) or actor-critic (e.g. (Foerster, 2018)), which are combined with concepts from game theory (Busoniu et al., 2008; Hernandez-Leal et al., 2019). The three approaches below address different training and execution methods for agents. Architectures and more detail of the three approaches can be found in (Albrecht et al., 2024; Gronauer and Diepold, 2022; Wong et al., 2022).

2642

### 2643 A.4.1 Decentralized Training Decentralized Execution (DTDE)

An approach to find the policy for each of the agents in the Multi Agent setting would be to use the same approach for each individual agent as in the Single Agent setting. Each of the agents takes its input from the environment, take its own actions and get its own rewards. The agents do not share information with other agents (Bono et al., 2019; Gronauer and Diepold, 2022). An example is the Independent Learners approach (Claus and Boutilier, 1998; Kok and Vlassis, 2006) and is based on Q-learning. Each agent has an individual table for  $Q$ -values and the global  $Q$ -function is defined as a linear combination of individual contributions. The local  $Q$ -function is then updated as:

2653

$$2654 \quad Q^n(S_t, A_t^n) \leftarrow Q^n(S_t, A_t^n) + \alpha \left[ R_{t+1}^n(S_t, \mathbf{A}_t) + \gamma \max_{a^n} Q^n(S_{t+1}, a^n) - Q^n(S_t, A_t^n) \right].$$

(A23)

2656

Note that the global state  $S_t$  is used to determine the agents'  $Q$ -values as well as to determine the individual reward  $R_{t+1}^n(S_t, \mathbf{A}_t)$ . A disadvantage of this approach is that the environment appears non-stationary (Padakandla, 2021) to the agents. By non-stationarity we mean that the dynamics of the environment change over time. The agents do not have access to the other agents' actions, nor can they perceive the joint action. Finally, each agent is learning, that is updating its own policy independently, impacting the other agents' decision making.

2664

### 2665 A.4.2 Centralized Training Centralized Execution (CTCE)

2666

In another approach, agents learn their policies based on common or shared information, for example local observations or policies. The collection of observations from

2668

each of the agents is used to train a central policy. The central policy sends the actions to all agents for execution. In centralized Q-learning, the update equation (Albrecht et al., 2024) becomes:

$$Q(S_t, \mathbf{A}_t) \leftarrow Q(S_t, \mathbf{A}_t) + \alpha \left[ R_{t+1}(S_t, \mathbf{A}_t) + \gamma \max_{\mathbf{a}} Q(S_{t+1}, \mathbf{a}) - Q(S_t, \mathbf{A}_t) \right]. \quad (\text{A24})$$

A disadvantage of this approach is that the state-action spaces grow exponentially with the number of agents. This could be mitigated by policy factorization using the combined observations, but also by creating individual policies, leading to individual actions. Although this would significantly reduce the size of the action space, the exponential growth of the observation space would still be a problem (Gupta et al., 2017).

#### A.4.3 Centralized Training Decentralized Execution (CTDE)

Yet another approach would have homogeneous agents having an individual policy, using local observations to a distribution over individual actions but enable them to share information and/or resources during training. Agents execute their policies based on local observations. With this approach, non-stationarity and partial observability are mitigated even if other agents' policies are changing (Wong et al., 2022).

#### A.4.4 Issues and challenges.

MARL have additional challenges when compared to Single Agent RL. In addition to credit assignment, sparse rewards, stability and generalization these are partial observability, non-stationarity and the effectiveness of different training schemes (Wong et al., 2022). Except for the cooperative case, it is unlikely that agents share observations, policies or parameters. Additional challenges in that case are observing, analyzing and acting on the other agents without sharing the information as described in the cooperative approaches (Busoniu et al., 2008; Hernandez-Leal et al., 2017).

### A.5 Deep Reinforcement Learning

Deep Reinforcement Learning algorithms use neural networks as function approximators to solve problems which have continuous or high dimensional state and/or action spaces. We provide relevant examples below, which are selected as these or modifications are used as part of the solution architecture in road pricing. For details on the cited algorithms we refer to the original works.

#### A.5.1 Value Based Algorithms.

DQN (Deep Q Networks)(Mnih et al., 2015) is an adaptation of Q-Learning. A separate current Q-network  $Q(s, a, \theta)$  and a target Q-network  $Q(s, a, \theta')$  are used to determine the optimal policy. The experiences gathered and stored by taking actions based on  $Q(s, a, \theta)$  are used to learn the parameters  $\theta$  by minimizing a loss function using the current Q-network and target Q-network. Periodically the target parameters  $\theta'$  are updated using  $\theta$

### 2715 **A.5.2 Policy Based and Actor-Critic Algorithms.**

2716 In Proximal Policy Optimization (PPO) (Schulman et al., 2017) an update to a pol-  
2717 icy is maximized subject to a constraint on the size of the update. Maximization is  
2718 performed using clipped probability ratios, which compare the new policy  $\pi_{\theta}$  and old  
2719 policy  $\pi_{\theta_{old}}$ . In another solution approach, Soft Actor Critic (SAC) (Haarnoja et al.,  
2720 2018, 2017), the reward is augmented with an entropy term. Adding the entropy term  
2721 improves both exploration by acquiring diverse behaviors and in addition, the algo-  
2722 rithm is more robust for model and estimation errors. Yet another approach, Deep  
2723 Deterministic Policy Gradient (DDPG) (Lillicrap et al., 2015), combines ideas from  
2724 DPG (see Section A.2 and (Silver et al., 2014)) and DQN (see above). It is a model  
2725 free, off-policy actor-critic algorithm with neural networks as function approximators.  
2726

### 2727 **A.5.3 Multiple Actor/Learner for Value and Policy Based** 2728 **Algorithms.**

2729 These algorithms parallelize part of the learning experience. Each actor interacts with  
2730 its own version of the environment and experiences are shared. Mnih et al. (2016)  
2731 proposed a framework Asynchronous Advantage Actor Critic (A3C) that uses asyn-  
2732 chronous gradient descent for optimization. Multiple agents operate asynchronously,  
2733 in parallel, on multiple instances of the environment. Periodically the target network of  
2734 the actor is updated and an asynchronous update of the global parameters is performed  
2735 by sending the gradients. Multiple updates from the various learners are applied to  
2736 the central parameters and because of this, are likely to be less correlated, stabilizing  
2737 learning. See Figure 5 for a case with  $M$  agents (both actor and learner), indexed by  
2738  $i$ , asynchronously updating their local parameters to the global shared parameter net-  
2739 work. In case of a single agent it is referred to as A2C. Approaches to enhance DDPG  
2740 are D3PG and D4PG (Barth-Maron et al. (2018); Bellemare et al. (2017)).  
2741  
2742

### 2743 **A.5.4 Multi Agent Reinforcement Learning.**

2744 A range of Multi Agent RL algorithms using neural networks exist (Albrecht et al.,  
2745 2024; Gronauer and Diepold, 2022; Wong et al., 2022). The surveyed works below  
2746 mostly use an actor-critic variant using a CTDE architecture, specifically designed for  
2747 their use case and will be discussed there.  
2748

## 2749 **Appendix B Additional Tables**

2750  
2751  
2752  
2753  
2754  
2755  
2756  
2757  
2758  
2759  
2760

**Table B1:** Overview of simulators used in surveyed works

	6.1.1 (Mirzaei et al. (2018b))	6.2.1 (Chen et al. (2018))	6.2.2 (Qiu et al. (2019))	6.2.3 (Jin et al. (2021))	6.2.4 (Wang et al. (2022))	6.2.5 (Lu et al. (2024))	6.2.6 (He et al. (2024))	6.3.1 (Zhu and Ukkusuri (2015))	6.3.2 (Pandey and Boyles (2018))	6.4.1 Pandey et al. (2020))	6.4.2 (Zhang et al. (2023))	6.5.1 (Sato et al. (2021))	6.5.2 (Sato et al. (2022))	6.6.1 (Tavares and Bazzan (2014))	6.6.2 (Chakravarty et al. (2024))	6.7.1 (Ramos et al. (2018))	6.7.2 (Ramos et al. (2020))
<b>Simulation models</b>																	
Custom DTA Simulator (CTM)	✓																
Custom Simulator (BPR)		✓	✓														
Custom Simulator (BPR)				✓													
CityFlow					✓	✓											
Custom Simulator							✓										
Custom Simulator (CTM)								✓									
Custom Simulator (CTM)									✓	✓							
Custom Highway RL Environment											✓						
Custom Simulator												✓	✓				
SUMO														✓			
MATLAB/Simulink															✓		
Custom Simulator																✓	✓

This table is meant to illustrate the diversity in simulation models used in the surveyed works.



**Table B2:** Overview of road networks in surveyed works

	6.1.1 (Mirzaei et al. (2018b))	6.2.1 (Chen et al. (2018))	6.2.2 (Qiu et al. (2019))	6.2.3 (Jin et al. (2021))	6.2.4 (Wang et al. (2022))	6.2.5 (Lu et al. (2024))	6.2.6 (He et al. (2024))	6.3.1 (Zhu and Ukusuri (2015))	6.3.2 (Pandey and Boyles (2018))	6.4.1 Pandey et al. (2020))	6.4.2 (Zhang et al. (2023))	6.5.1 (Sato et al. (2021))	6.5.2 (Sato et al. (2022))	6.6.1 (Tavares and Bazzan (2014))	6.6.2 (Chakravarty et al. (2024))	6.7.1 (Ramos et al. (2018))	6.7.2 (Ramos et al. (2020))
<b>Road networks</b>																	
Sioux-Falls	✓																✓
Simplified Sioux-Falls								✓									
Downtown Austin	✓																
Uptown San Antonio	✓																
Singapore Central Region		✓	✓														
Singapore Whole Net			✓														
Qinhuai ext. network				✓													
Hangzhou					✓	✓											
Manhattan					✓	✓											
Porto Roadnet					✓												
Jinan City						✓											
Network of regions							✓										
DESE Network									✓	✓							
LBJ Texpress Toll Segment									✓	✓	✓						
Single Entry Single Exit										✓	✓						
MoPac Express Lane Network											✓						
Simple Expressway Network												✓					
GYExpress													✓				
Single Bottleneck														✓			
Three Bottleneck															✓		
Parallel Bottleneck Network																✓	
Test Network																	✓
Toll Lane, Regular Lane																	
B1-B7																✓	✓
BB1, BB3, BB5, BB7																✓	✓
OW																✓	✓
Anaheim																	✓
Eastern Massachusets																	✓

This table is meant to illustrate the diversity in road networks used in the surveyed works.