1	A Scalable Model-Free Deep Reinforcement Learning-Based Perimeter Metering Control
2	Method for Multi-Region Urban Networks
3	
4	Dongqin Zhou
5	Department of Civil and Environmental Engineering
6	The Pennsylvania State University, University Park, PA, 16802
7	Email: dongqin.zhou@psu.edu
8	
9	Vikash V. Gayah <sup>*</sup>
10	Department of Civil and Environmental Engineering
11	The Pennsylvania State University, University Park, PA, 16802
12	Email: gayah@engr.psu.edu
13	* Corresponding author.
14	
15	Word Count: 6868 words + 1 table (250 words per table) = 7118 words
16	
17	
18	Submitted [07/07/2022]

# 1 ABSTRACT

- 2 Perimeter metering control is a convenient strategy to mitigate urban congestion by manipulating vehicular
- 3 movements around regional perimeters without modeling the detailed behaviors and interactions. Multi-
- region perimeter metering control holds promise for efficient management of urban traffic in large-scale
   networks. However, most existing methods for multi-region control require knowledge of either the traffic
- 6 dynamics or network properties (i.e., the critical accumulations), and completely model-free techniques are
- still lacking in the literature. To fill this gap, this paper proposes a novel control scheme based upon model-
- 8 free multi-agent reinforcement learning. The scheme features value function decomposition in the paradigm
- 9 of centralized training with decentralized execution, coupled with critical advances of single-agent deep
- 10 reinforcement learning and specialized problem reformulation. The effectiveness of the scheme is
- 11 demonstrated with numerical experiment conducted on a seven-region urban network.
- 12 Keywords: Macroscopic Fundamental Diagram (MFD); multi-region perimeter metering control; model-
- 13 free multi-agent reinforcement learning (MARL)
- 14

# 1 INTRODUCTION

2 The macroscopic fundamental diagram (MFD) has been recently shown promising in terms of the aggregate 3 traffic dynamics models and network-level control schemes it facilitates (1-5), both of which are critical 4 for large-scale urban traffic management. The initial theoretical investigation of the MFD dates back to the 5 1960s (6), but its existence was not verified until recently (1, 2). These seminal works have since inspired 6 a large number of research efforts on the existence analysis (7-9) and estimation of MFDs (10-17). Other 7 than the derivations, properties of well-defined MFDs have also been examined extensively (3, 18-20). 8 These references have shown that urban networks are subject to instability, hysteresis, and bifurcation 9 phenomena with heterogeneous distribution of vehicle presence. Fortunately, network partitioning can be 10 applied to divide a large heterogeneous area into several smaller regions such that congestion homogeneity 11 is maintained for each region (21-24).

12 Well-defined MFDs enable low-complexity modeling of traffic dynamics by focusing on aggregate 13 vehicular movements within and across homogeneous regions. This elegant modeling paradigm has led to 14 the development of numerous regional control schemes, e.g., congestion pricing (25-27), route guidance 15 (28-30), and others. The most common application utilizing MFDs is perimeter metering control (PMC), 16 which entails regulating the inter-regional transfer flows using traffic signals residing on the boundary of 17 the neighboring regions. The first examination of perimeter control was presented in (1) for a single region, 18 which formulated the traffic dynamics with MFDs and proposed the optimal Bang-Bang control policy. 19 Perimeter control for two-region networks, as first formulated in (31), has also attracted substantial research 20 interests over the years. For example, analytical and data-driven approaches have been adopted to design 21 solution schemes (4, 32-35), while stability and modeling uncertainty are examined in (36-40).

22 Another line of PMC research pertains to the efficient operations of traffic flows in a multi-region 23 setting. Early endeavors in this vein include (41, 42) where the traffic dynamics are formulated for a multi-24 reservoir and a mixed network. However, in these works the receiving capacity constraint was neglected; 25 and it was later integrated in (43) where a region-based and subregion-based MFD models were proposed. 26 To further enhance traffic dynamics modeling, numerous works have been conducted to consider: boundary 27 queue dynamics (37, 44, 45), time-delay effects (46), demand stochasticity (47), and parameter uncertainty 28 in MFDs (48). Multi-region PMC embodies great potential for city-level traffic management, for which 29 various solution methods have been proposed in the literature. Examples include linear quadratic regulator 30 (41, 44), model predictive control (29, 43), model-free adaptive control (49, 50), and others. Importantly, (49, 50) proposed solution schemes that are data-driven and model-free, yet the critical accumulation is still 31 32 explicitly blended into the controller designs. Other mentioned methods, on the other hand, are heavily 33 dependent on knowledge of the environment dynamics. Therefore, it is a high research priority to develop 34 a completely model-free method for multi-region perimeter control.

This paper bridges this gap by proposing a model-free scheme based on multi-agent reinforcement learning that features centralized training with decentralized execution and value function decomposition. Moreover, the scheme adopts the Bang-Bang type action design, which was corroborated as the optimal action form for PMC problems (1, 32, 44). The effectiveness of the proposed scheme is demonstrated via comparison with the model predictive control by conducting perimeter control on a seven-region network.

The remainder of the paper is structured as follows. The next section provides the traffic dynamics modeling of multi-region urban networks. Then, the proposed scheme is explained in detail, followed by the numerical experiment results. Finally, concluding remarks are presented in the last section.

43

# 44 TRAFFIC DYNAMICS OF MULTI-REGION URBAN NETWORKS

The general traffic dynamics for an *R*-region urban network are introduced here; see Fig. 1 for an illustration of a network with R = 7. The regions are assumed to be homogenous in terms of congestion distribution; however, if this assumption does not hold, network partitioning can be applied to maintain homogeneity (21-23). As such, a well-defined MFD function  $f_i(n_i(t))$  that relates trip completion rate to the regional accumulation  $n_i(t)$  is assumed to be able to model each region. The evolution of accumulations in region *i* can then be expressed as (28, 43):

5 
$$n_i(t) = \sum_{j \in \mathcal{R}} n_{ij}(t)$$
 (1)

6 
$$\dot{n}_{ii}(t) = q_{ii}(t) - M_{ii}(t) + \sum_{h \in N_i} u_{hi}(t) \cdot M_{hii}(t)$$
(2)

$$\dot{n}_{ij}(t) = q_{ij}(t) + \sum_{h \in N_i; h \neq j} u_{hi}(t) \cdot M_{hij}(t) - \sum_{h \in N_i} u_{ih}(t) \cdot M_{ihj}(t)$$
(3)

8 where  $\mathcal{R} = \{1, 2, \dots, R\}$  denotes the network,  $n_{ij}$  and  $q_{ij}$  are the number of vehicles and traffic demands in 9  $R_i$  destined for  $R_j$ , with  $n_{ii}$  and  $q_{ii}$  defined similarly.  $u_{ih}$  is the perimeter controller ( $u_{ih} \in [u_{min}, u_{max}]$ 10 with  $0 \le u_{min} < u_{max} \le 1$ ) that specifies the allowable ratio of transfer flow from  $R_i$  to  $R_h$ , with h11 belonging to the neighboring regions of  $R_i$ ,  $N_i$ .  $M_{ihj}(t)$  represents the transfer flow from  $R_i$  to  $R_j$  via the 12 next region h, while  $M_{ii}(t)$  is the exit flow of region i, as calculated by:

13 
$$M_{ihj}(t) = \theta_{ihj}(t) \cdot \frac{n_{ij}(t)}{n_i(t)} \cdot f_i(n_i(t)), i \neq j, h \in N_i$$
(4)

14 
$$M_{ii}(t) = \frac{n_{ii}(t)}{n_i(t)} \cdot f_i(n_i(t))$$
(5)

15 where  $\theta_{ihj}(t) \in [0, 1]$  is the route choice term for vehicles traveling from  $R_i$  to  $R_j$  via the next region *h* 16 (hence  $\sum_{h \in N_i} \theta_{ihj}(t) = 1$ ) and is inversely related to the travel time of paths utilizing  $R_h$ .

17 The receiving capacity of regions with high accumulations might be insufficient to contain all 18 inflow vehicles, thus restraining the full penetration of transfer flows. As such, the capacity-restrained 19 transfer flows  $\hat{M}_{ihj}(t)$  are defined as (28, 43):

20 
$$\widehat{M}_{ihj}(t) = \min\left(M_{ihj}(t), C_{ih}(n_h(t)) \cdot \frac{M_{ihj}(t)}{\sum_{k \in R, k \neq i} M_{ihk}(t)}\right)$$
(6)

21 where  $C_{ih}(n_h(t))$  is the boundary capacity between  $R_i$  and  $R_h$  and is a function of  $n_h(t)$  as in:

22 
$$C_{ih}(n_h(t)) = \begin{cases} C_{ih}^{max}, & 0 \le n_h(t) \le \alpha \cdot n_h^{jam} \\ \frac{C_{ih}^{max}}{1-\alpha} \cdot (1 - \frac{n_h(t)}{n_h^{jam}}), & \alpha \cdot n_h^{jam} \le n_h(t) \le n_h^{jam} \end{cases}$$
(7)

where  $C_{ih}^{max}$  is the maximum boundary capacity between region *i* and *h*,  $n_h^{jam}$  is the accumulation value of region *h* where gridlock arises, and  $\alpha \in (0,1)$  is a parameter that signals the decrease of receiving capacity with the increase of accumulation.

26



Fig. 1. A seven-region urban network. The dash lines represent the perimeter controllers.

# 4 METHODOLOGY

5 This section reformulates the seven-region control problem in the context of multi-agent reinforcement 6 learning (MARL), followed by detailed explanations of the proposed scheme.

7

1 2

3

# 8 **Problem Reformulation**

9 The multi-region PMC problem can be viewed as a cooperative multi-agent task where a group of *n* agents 10  $(\mathcal{N} = \{1, \dots, n\})$  learn to achieve a common goal via individualized interactions with the same environment. 11 In this work, each agent is supposed to regulate two inter-regional vehicle movements by selecting values 12 for a pair of perimeter controller on a regional boundary. Formally, the multi-region PMC problem is 13 presented as a decentralized partially observable Markov decision process defined by a tuple < 14  $\mathcal{S}, \mathcal{O}, \mathcal{U}, \mathcal{P}, r, \pi, \gamma, \mathcal{N} >:$ 

- State space, S, and observation space, O. The state contains the global information about the entire network. However, due to partial observability, the agents can only observe local instances of the state and act based on these observations. In this work, the state  $s_t$  consists of all regional accumulations, traffic demands, and a binary congestion indicator that denotes whether the regions are congested or not. The observation  $o_t$  includes similar information for a pair of regions, i.e., two regional accumulations, traffic demands between the two regions, and the congestion indicator.
- Action space, U. The Bang-Bang control policy, i.e., to alternate the controller between the minimum and maximum values, has been shown as the optimal form of actions for perimeter control (1, 32, 44). Motivated by this, the agents assume the Bang-Bang type actions, i.e., either u<sub>min</sub> or u<sub>max</sub> for each perimeter controller.
- **Transition dynamics**,  $\mathcal{P}$ . The selected actions of the individual agents form a joint action, which is executed in the environment and leads a transition to a new state, according to the dynamics

 $\mathcal{P}(s_{t+1}|s_t, u_t): S \times U \to S$ . Note that, the proposed scheme is model-free and thus internalizes 1 such dynamics through the learning process without explicit modeling. 2 3 **Reward function**, *r*. After executing the joint action, the environment returns a real-time scalar 4 reward back to the agents as a quality assessment. The reward  $r(s_t, u_t)$  helps guide the agents to 5 achieve the control objective, i.e., to maximize the cumulative trip completion; and therefore, it is 6 defined as the trip completion in a time step. 7 **Policy**,  $\pi$ , and **discount factor**,  $\gamma$ . The agents select actions for the perimeter controllers based 8 upon the local observation  $o_t^a$  according to the policy  $\pi^a(u^a|o^a)$ . To differentiate immediate 9 rewards from delayed ones, a discount factor  $\gamma \in [0,1]$  is utilized. Collectively, the agents learn via 10 trial and error to maximize the expected total discounted reward, i.e., the return, as calculated by  $G_t = \sum_{\tau=t}^T \gamma^{\tau-t} r_{\tau+1}$  where *T* is the total number of time steps. 11 12

# 13 Algorithms

14 This section first introduces a canonical single-agent deep reinforcement learning method, which provide 15 theoretical background for the proposed scheme that is to be explained subsequently.

- 16
- 17 Double Deep Q Networks (Double DQN)

18 As a foundational reinforcement learning technique, Q-learning (51) has received sustained interests over 19 the years. Using a tabular form, it stores the long-term quality measurements of distinct state-action pairs, 20 i.e., the Q value  $Q(s_t, u_t)$  that denotes the expected return from the environment after taking action  $u_t$  at 21 state  $s_t$ . During the learning process, the Q values are updated iteratively, according to:

22 
$$Q(s_t, u_t) \leftarrow Q(s_t, u_t) + \kappa \cdot \left(r_{t+1} + \gamma \cdot \max_u Q(s_{t+1}, u) - Q(s_t, u_t)\right)$$
(8)

where  $\kappa$  is the learning rate. With sufficient learning updates, the Q values tend towards invariant, and the final learned policy can be derived in a greedy manner with respect to the Q values.

Despite its popularity, the tabular form of Q-learning limits its applicability to large problems that feature an abundance of state-action pairs. To mitigate this, research efforts have long been performed on value function approximation and its stability analysis (52-54), with the first success presented in the Deep Q-Networks (DQN) algorithm (55). This work has revealed the significant potential of deep reinforcement learning and has since inspired the development of more advanced learning techniques (56-60). Despite its success, however, the DQN method is prone to overestimation of the Q values (56). In the latter reference, an improved algorithm named Double DQN is proposed, which revises the learning target of DQN by using

32 the Q-network for action selection and target network for evaluation, as follows:

33 
$$Y_t = r_{t+1} + \gamma Q\left(s_{t+1}, \arg\max_u Q(s_{t+1}, u; \boldsymbol{\theta}_t); \boldsymbol{\theta}_t^-\right)$$
(9)

where  $Q(:,:;\boldsymbol{\theta}_t)$  and  $Q(:,:;\boldsymbol{\theta}_t^-)$  represent the Q- and target networks. The target network is a periodic copy of the Q-network, and its utilization provides relatively static targets and helps with learning stability.

36

## 37 Reinforcement learning controller design for multi-region perimeter control (MR-RL)

38 The success of single-agent RL has significantly boosted its extension to multi-agent systems. However,

39 directly applying single-agent techniques to multi-agent tasks is generally infeasible due to the curse of

- 40 dimensionality which renders it difficult to estimate the joint Q value. The most common paradigm to
- 41 alleviate this issue is centralized training with decentralized execution (CTDE, (61)), which adopts full
- 42 centralization conditioning on the global state during learning and decentralization conditioning on local

observations during action taking. Despite notable experimental results, the CTDE paradigm has a major scalability limitation due to fully centralized training. As such, value function decomposition has been proposed (62) as a mitigation strategy. Specifically, these methods factorize the centralized Q value as a function of the local Q values estimated by the agents conditioning on the local observations and actions. Importantly, this factorization ensures scalability as the joint action information is no longer required. As a representative method, QMIX (63) decomposes the joint Q value as a nonlinear but monotonic composition of the local Q values, and this decomposition has been widely adopted in later efforts, e.g., (64).

8 The multi-region control problem is a fully cooperative task where all agents collaborate to achieve 9 the highest trip completion. Naturally, higher local trip completions lead to higher total trip completion. 10 Hence, the monotonicity assumption holds and the QMIX is adopted to devise the proposed scheme, as 11 denoted by MR-RL that stands for Multi-Region Reinforcement Learning. The learning algorithm for the 12 proposed MR-RL is shown in Fig. 2, and its building components are explained in the following.

13



### 14 15 16

Fig. 2. A diagram of the learning algorithm for the MR-RL scheme.

17 The MR-RL scheme holds a group of agents for multi-region perimeter control, and each agent is 18 constructed as a multi-layer perceptron, a structure widely used in the literature (59, 63, 65). To improve 19 training efficiency, parameters of the agent network are shared. Hence, the agents with shared parameters 20 can be represented as  $Q(o^a, u^a; \theta^Q)$ , where  $\theta^Q$  is the network parameters. The agents receive as input the

local observations  $o^a$  and estimate the 4-dimensional local Q values for the two perimeter controllers (each 1 controller has two options,  $u_{min}$  and  $u_{max}$ ). The local actions can then be derived with the  $\epsilon$ -greedy 2 strategy regarding these values, i.e., the greedy action  $\arg \max_{u^a} Q(o^a, u^a; \theta^Q)$  is chosen with probability 3

4  $1 - \epsilon$  and a random action otherwise. To better balance exploration and exploitation, the  $\epsilon$  value is decayed 5 through time, with the decay schedule to be presented shortly.

6 The mixing network, as denoted by  $m(\cdot)$ , combines the estimated local O values to provide the 7 joint Q value and is central to value decomposition methods. The QMIX algorithm enforces monotonic 8 decomposition by utilizing non-negative weights for the mixing network, and separate hypernetworks are 9 exploited to produce such weights. Specifically, the hypernetworks take into the global state  $s_t$  and generate 10 the weights in a feed-forward structure with an absolute activation function. The hypernetworks also create biases for the mixing network, but these are not restricted to be non-negative. 11

12 The Double DQN update rule, along with the QMIX type value decomposition, is used to construct 13 learning targets (see Fig. 2) for the proposed MR-RL scheme, as follows:

14 
$$Y_{t} = r_{t+1} + \gamma \cdot m \left( s_{t+1}, \left\{ Q \left( o_{t+1}^{a}, \arg \max_{u^{a}} Q \left( o_{t+1}^{a}, u^{a}; \theta_{t}^{Q} \right); \theta_{t}^{Q-} \right) \right\}_{a=1}^{n}; \theta_{t}^{m-} \right)$$
(10)

15

where  $\arg \max Q(\cdot, \cdot; \theta_t^Q)$  is the local action selection using the shared agent network,  $Q(\cdot, \cdot; \theta_t^{Q^-})$  is the action evaluation with the target agent network, and  $m(\cdot, \cdot; \theta_t^{m^-})$  represents the target mixing network 16 whose inputs include the global state for non-negative weights construction. The major distinction between 17 18 this and the Double DQN target (Eq. (9)) is the mixing network which involves a group of local Q values. 19 Importantly though, this additional complexity precipitates significantly improved scalability to larger 20 multi-agent systems. The network parameters of the MR-RL scheme can be updated by minimizing the 21 following loss:

22 
$$\mathcal{L}(\theta_t^Q, \theta_t^m) = \sum_{i=1}^{b} \left[ Y_t^i - m \left( s_t^i, \{ Q(o_t^{a,i}, u_t^{a,i}; \theta_t^Q) \}_{a=1}^n; \theta_t^m \right) \right]^2$$
(11)

where b is the batch size of sampled transitions from the replay buffer used for network updates, and  $Y_t^i$  is 23 the learning target for the *i*-th transition. 24

25 The proposed MR-RL scheme is model-free in that it does not require a priori knowledge of the 26 environment dynamics. Instead, it learns the control policy from pure interactions with the environment, 27 and the interactions are stored in a replay buffer in the form of state-action-reward pairs, i.e., the transitions 28 in Fig. 2. The use of a replay buffer was initially presented in (66) and later consolidated in (55) as a critical 29 component of deep reinforcement learning. Specifically, the replay buffer is first utilized to store the 30 collected transitions; then during training, minibatches of transitions are randomly sampled from the buffer 31 to update the network parameters i.e., the shared agent network, hypernetworks, and the mixing network. 32 The replay buffer has been shown helpful to stabilize the learning process as the random sampling helps 33 remove correlations between the transitions. Further, to guarantee effective learning for the MR-RL scheme, 34 the Ape-X distributed architecture (65) is adopted. Concretely, the architecture maintains numerous 35 instantiations of the environment in parallel, with which the MR-RL interacts to collect an increased number 36 of transitions. These derived transitions are then pooled together in the replay buffer for future updates of 37 the network parameters. With enough training updates, the final learned control strategy can be obtained by applying the greedy policy on the fully trained agent network, i.e.,  $u_t^a = \pi(o_t^a) = \arg \max_u Q(o_t^a, u; \theta_t^Q)$ . 38

39 With these expositions, the proposed MR-RL scheme built with the learning algorithm and the 40 Ape-X architecture is formalized in Algorithm 1. Note that,  $\theta_t^m$  expresses the weights of the mixing network which include weights of the hypernetworks as a constituent element. In addition, the generator 41

refers to the instantiated environment, i.e., a transition generator. Finally, the list of hyperparameters along
 with their values is presented in Table 1.

Algo	rithm 1. Reinforcement Learning controller for Multi-Region perimeter control (MR-RL)			
1:	Randomly initialize shared agent network $\theta_0^Q$ and mixing network $\theta_0^m$ (hypernetworks included)			
	Initialize target agent and mixing networks $\theta_0^{Q^-} = \theta_0^Q, \theta_0^{m^-} = \theta_0^m$			
	Initailize replay buffer, buffer size B, sample size b, iteration number I, and genetaor number G			
2:	for $iter = 1$ to $I$ do			
3:	Compute the decayed $\epsilon$ value for $\epsilon$ –greedy exploration			
4:	for generator = 1 to $G$ do			
5:	Load the shared agent network $\boldsymbol{\theta}_{iter}^Q = \boldsymbol{\theta}_{iter-1}^Q$			
6:	$s_0, o_0 \leftarrow \text{Environment.Reset}()$			
7:	for $t = 1$ to $T$ do			
8:	$u_{t-1}^a = \arg \max_{i} Q(o_{t-1}^a, u; \theta_{iter}^Q)$ with probability $1 - \epsilon$			
	a random action with proability $\epsilon$			
	$\boldsymbol{u}_{t-1} = \{u_{t-1}^a\}_{a=1}^n$			
9:	$(r_t, s_t, \boldsymbol{o}_t) \leftarrow \text{Environment.Step}(s_{t-1}, \boldsymbol{o}_{t-1}, \boldsymbol{u}_{t-1})$			
10:	Store $(s_{t-1}, \boldsymbol{o}_{t-1}, \boldsymbol{u}_{t-1}, r_t, s_t, \boldsymbol{o}_t)$ into the replay buffer			
11:	end for			
12:	end for			
13:	if the number of stored transitions exceeds the buffer size B then			
14:	Remove outdated transitions			
15:	end if			
16:	Training samples $\leftarrow$ a batch of b transitions randomly drawn from the buffer			
17:	Periodically target networks $\theta_{iter}^{Q^-} = \theta_{iter-1}^Q, \theta_{iter}^{m^-} = \theta_{iter-1}^m$			
18:	$\boldsymbol{\theta}_{iter}^{Q}, \boldsymbol{\theta}_{iter}^{m} \leftarrow$ Update the network parameters by minimizing the loss as in Eq. (11)			
19:	end for			

29 30

Table 1. List of hyperparameters and their values

Hyperparameter	Value	Description
Iteration number (I)	250	The number of training iterations
Generator number $(G)$	6	The number of experiment instantiations to collect transitions
Replay buffer size $(B)$	10000	The storage capacity of the replay buffer
Sample size ( <i>b</i> )	1000	The number of transitions sampled for network updates
Initial $\epsilon$	0.90	The initial value of $\epsilon$ in $\epsilon$ – greedy exploration
$\epsilon$ decay	0.98	The exponential decay factor for the $\epsilon$ value
Final $\epsilon$	0.01	The final value of $\epsilon$ in $\epsilon$ – greedy exploration
Update epoch	5	The times to update the network parameters at each iteration
Initial learning rate	0.003	The initial learning rate used by RMSprop for the network updates
Learning rate decay	0.95	The exponential learning rate decay factor at each iteration
Minimum learning rate	0.0001	The minimum learning rate used by RMSprop
Discount factor	0.8	The discount factor used to compute the learning targets (Eq. (10))
Target networks lifetime	10	The number of iterations to update the target networks

## 1 **EXPERIMENTS**

In this section, the proposed MR-RL is applied for perimeter control on a seven-region network (with configurations shown in Fig. 1), and its performance is compared with two benchmarking methods to evaluate its effectiveness. Note that, there are 24 perimeter controllers for 12 pairs of neighboring regions in the network (see again Fig. 1), hence 12 agents are utilized.

6

# 7 Experiment Setup

8 In this work, a unit MFD consistent with the one observed in Yokohama (2) is utilized, with critical and 9 jam values of 8,240 veh and 34,000 veh, respectively (35, 67). Note that, the unit MFD assumes a piecewise 10 functional form rather than a 3-rd polynomial for the traffic dynamics to be more realistic. For all 11 experiments, each region is modeled with a slightly scaled (within  $\pm 10\%$ ) version of unit MFD, as similarly 12 done in (29). In addition, the parameters for the boundary capacity constraints are set to  $C_{ih}^{max} = 4.6$  veh/s 13 and  $\alpha = 0.48$ .

The traffic demand profiles adopted for the numerical experiments are shown in Fig. 3. A two-hour control period is simulated with high inflows to region 4 and relatively small demands among the periphery regions, which mimics traffic conditions during a morning peak. The duration of a time step is set as  $\Delta t =$ 60s, which is a realistic cycle length for the signalized intersections on the regional boundaries that implement perimeter control. In addition, the boundary values for the perimeter controllers are  $u_{min} =$ 0.1,  $u_{max} = 0.9$ . Finally, region 4 assumes a congested initial state with an accumulation value of 8,750 veh while all other regions are assumed to be uncongested initially with accumulations of 3,850 veh.

21



22 23

24

Two benchmarking methods, model predictive control (MPC) and no control (NC), are adopted to compare the performances with the MR-RL, in terms of the cumulative trip completion (CTC) they achieve. The NC method does not impose limitations on the transfer flows and instead uses the maximum value for all perimeter controllers; it is usually adopted as a baseline method that provides the lower-bound control performances. In contrast, the MPC is a model-based rolling horizon optimization scheme that has achieved

- 1 state-of-the-art control performances. However, one major disadvantage of the MPC is that it builds upon
- 2 full knowledge of the environment dynamics (i.e., the MFDs and dynamic equations governing vehicle
- 3 movement between regions), which are generally difficult to obtain in the first place. In this paper, the MPC
- 4 is implemented as per the perimeter control-only scheme in (29) with a control horizon of 2 and a prediction
- 5 horizon of 3, as similarly adopted in numerous other works (5, 49, 50). Note that for the seven-region
- 6 perimeter control problem, a larger prediction horizon does not necessarily lead to better control 7 performance since the solution space of the formulated nonconvex optimization problem becomes so large
- 8 that finding the global optimum is increasingly difficult.
- 9

### 10 **Experiment Results**

11 For the experiment scenario considered, the traffic dynamics assumed by the MPC in the prediction model

12 are the same as those in the plant. The MR-RL is trained with five fixed random seeds and its performance

- 13 curves are shown in Fig. 4, where the darker line and shad ed area respectively represent the mean and 95%
- 14 confidence interval of the control gains (in terms of CTC). The MPC and NC are also run five times to
- 15 report their performance curves, but the curves are relatively invariant as they are not learning-based
- methods. Note that there is built-in randomness associated with the travel time calculation of different paths 16 used to determine the route choice term (see Eq. (4)); hence these two methods also exhibit a (negligibly
- 17
- 18 small) range of CTC values.
- 19





Fig. 4. Performance curves of different methods for the no uncertainty scenario.

22 23 As shown in Fig. 4, the NC method realizes the lower CTC value, which is expected since unlimited 24 vehicle inflow into region 4 aggravates the congestion therein and adversely impacts other inter-region 25 vehicular movements. More importantly, the proposed MR-RL can consistently learn and achieve control gains that are commensurate with (sometimes even better than) the MPC. This showcases the significant 26

1 potential of model-free reinforcement learning methods over model-based approaches. The MPC is an 2 optimization-based approach and derives control actions by solving a large nonlinear nonconvex program 3 that features a sizable solution space. As such, it may fail to find the global optimum, which leads to slight 4 underperformance to the proposed scheme, though the MPC could theoretically be the optimal control 5 technique with improved performances via guaranteed global optimum finding. However, implementing 6 this is not conceivably straightforward. Comparatively, the MR-RL learns the control policy via trial and 7 error, and through this process it can encounter better acting strategy than the MPC. Finally, note that 8 training performances of the MR-RL in the early period are noticeably worse than the NC method. This is 9 reasonable since during this period the MR-RL is principally exploring the environment. In this paper, the 10 training process is presumed to be completed with numerical simulations. Therefore, the initial control 11 performances are not important as only the fully trained MR-RL scheme will be applied.

12 To further demonstrate the effectiveness of the MR-RL, its control outcomes are examined more 13 carefully in the following. Fig. 5 presents the evolution of accumulations for each region, as achieved by 14 different control methods. The critical accumulations are also provided in dotted lines which helps 15 determine the congestion situation for the regions. As can be observed, under the NC method, region 4 16 becomes extremely congested in the end while the accumulations in other regions are generally smaller 17 than realized by the MPC or the MR-RL. This is understandable as the region 4-bounded flows are much 18 larger than the others. However, notable congestion in region 4 leads to small trip completion therein, which 19 also makes inter-regional travel time-consuming. For example, region 6-bounded vehicles in region 2 that 20 normally would travel via region 4 might need to take a longer route to reach their destinations. 21 Consequently, the trip completions in other regions will be negatively influenced and the NC method ends 22 up with the lowest CTC. Comparatively, both the MPC and MR-RL can significantly reduce the congestion 23 in region 4, while in the meantime keeping the accumulations in other regions under the critical values. 24 This implies that these methods can indeed perform effective perimeter control as the most destination-25 loaded region (i.e., region 4) are protected from over-congestion, which is consistent with the AB strategy 26 proposed in (1). Finally, the similarity of accumulations in all regions between the MPC and MR-RL 27 indicates great comparability between them.





Fig. 5. Accumulation plots for all regions. The dotted lines represent the critical accumulations.

1

16

12

2 Fig. 6 presents the control actions  $u_{i4}$  of the MPC and MR-RL, while all other controllers are 3 omitted from the presentation here. The selective presentation is done intentionally since other controllers 4 are nearly inactive, i.e., they all adopt the maximum value  $u_{max}$ . This is expected since the implementation 5 of perimeter control here is mostly designed at protecting region 4 from severe congestion, for which  $u_{i4}$ 6 being active is sufficient. Likewise, the NC actions are not included for comparison either since they are all 7 equal to the maximum value. Fig. 6 reveals that the MR-RL imposes stricter limitation on the transfer flows 8 to region 4 from regions 5, 6, and 7; hence the accumulations in these three regions are generally larger 9 than those resulted from the MPC actions. Similarly, the accumulations in regions 1, 2, and 3 are smaller 10 than those realized by the MPC since more transfer flows can be completed under the MR-RL policy that 11 exhibits looser control (i.e., more  $u_{max}$  values in the actions). In addition, both methods select the 12 maximum value for all controllers in the initial period, which is sensible as there does not exist pronounced 13 congestion within the network (even region 4 is only moderately congested). Overall, these control actions 14 help explain the resulting evolutions of accumulations, which also illustrates how the proposed MR-RL is 15 comparable to the MPC.



- 17 18
- 19 20

21 This paper presents a novel MR-RL scheme for multi-region perimeter control based on model-free multi-22 agent reinforcement learning. Specifically, the proposed MR-RL features value function decomposition, 23 which significantly helps with scalability, recent breakthrough of single-agent deep reinforcement learning 24 (such as the Ape-X architecture, double Q-learning update rule, experience replay, target networks, etc.), 25 and suitable problem reformulation governed by domain expertise (e.g., the Bang-Bang type action space 26 design). The control efficacy of the MR-RL is demonstrated with numerical experiments on a simulated 27 seven-region urban network, and the results suggest that the scheme can consistently learn and converge to 28 final control performances that are comparable to the MPC method. It is worth reiterating that, the proposed 29 MR-RL is completely model-free which does not require a priori information about the environment. Hence, 30 it is not affected by a wide range of modeling mismatch, e.g., scaling errors and the time-changing feature

**CONCLUDING REMARKS** 

1 of MFDs, to which recent data-driven approaches are liable (49, 50). Further, note that this paper presents 2 the first examination of completely model-free methods on seven-region perimeter control.

3 Future works to this paper could consider integrating low-level intra-regional signal control, which 4 holds promise for realistic city-wide traffic management. In addition, it would be interesting to investigate 5 whether the pretrained scheme from numerical simulations can be transferred to a microsimulation platform 6 and continue its learning course with real-time interactions. This could help evaluate an extra level of 7 transferability for the MR-RL and its ability to keep learning in a different platform.

8

### 9 **ACKNOWLEDGEMENTS**

- 10 This research was supported by NSF Grant CMMI-1749200.
- 11

### 12 **AUTHOR CONTRIBUTIONS**

13 The authors confirm contribution to the paper as follows: study conception and design: VG, DZ; analysis

- 14 and interpretation of results: VG, DZ; draft manuscript preparation: VG, DZ. All authors reviewed the 15 results and approved the final version of the manuscript.
- 16

### 17 REFERENCES

- 18 1. Daganzo, C. F. Urban Gridlock: Macroscopic Modeling and Mitigation Approaches. Transportation 19 Research Part *B*: Methodological, Vol. 41, No. 1, 2007, 49-62. pp. 20 https://doi.org/10.1016/j.trb.2006.03.001.
- 21 2. Geroliminis, N., and C. F. Daganzo. Existence of Urban-Scale Macroscopic Fundamental Diagrams: 22 Some Experimental Findings. Transportation Research Part B: Methodological, Vol. 42, No. 9, 23 2008, pp. 759–770.
- 24 3. Daganzo, C. F., V. V. Gayah, and E. J. Gonzales. Macroscopic Relations of Urban Traffic Variables: 25 Bifurcations, Multivaluedness and Instability. Transportation Research Part B: Methodological, 26 Vol. 45, No. 1, 2011, pp. 278–288. https://doi.org/10.1016/j.trb.2010.06.006.
- 27 4. Geroliminis, N., J. Haddad, and M. Ramezani. Optimal Perimeter Control for Two Urban Regions 28 with Macroscopic Fundamental Diagrams: A Model Predictive Approach. IEEE Transactions on 29 **Transportation** 14, Intelligent Systems, Vol. No. 1. 2013, pp. 348-359. 30 https://doi.org/10.1109/TITS.2012.2216877.
- 31 5. Yildirimoglu, M., I. I. Sirmatel, and N. Geroliminis. Hierarchical Control of Heterogeneous Large-32 Scale Urban Road Networks via Path Assignment and Regional Route Guidance. Transportation 33 Research Part *B*: Methodological, Vol. 118, 2018, 106–123. pp. 34 https://doi.org/10.1016/j.trb.2018.10.007.
- 35 Godfrey, J. W. The Mechanism of a Road Network. Traffic Engineering & Control, Vol. 11, No. 7, 6. 36 1969, pp. 323-327.
- 37 Fu, H., Y. Wang, X. Tang, N. Zheng, and N. Geroliminis. Empirical Analysis of Large-Scale 7. 38 Multimodal Traffic with Multi-Sensor Data. Transportation Research Part C: Emerging 39 Technologies, Vol. 118, 2020, p. 102725. https://doi.org/10.1016/j.trc.2020.102725.
- 40 8. Geroliminis, N., and J. Sun. Properties of a Well-Defined Macroscopic Fundamental Diagram for 41 Urban Traffic. Transportation Research Part B: Methodological, Vol. 45, No. 3, 2011, pp. 605–617. 42 https://doi.org/10.1016/j.trb.2010.11.004.
- 43 9. Paipuri, M., Y. Xu, M. C. González, and L. Leclercq. Estimating MFDs, Trip Lengths and Path 44 Flow Distributions in a Multi-Region Setting Using Mobile Phone Data. Transportation Research 45 Part C: Emerging Technologies, Vol. 118. 2020, p. 102709.

- https://doi.org/10.1016/j.trc.2020.102709.
   Ambühl, L., and M. Menendez. Data Fusion Algorithm for Macroscopic Fundamental Diagram Estimation. *Transportation Research Part C: Emerging Technologies*, Vol. 71, 2016, pp. 184–197.
   https://doi.org/10.1016/J.TRC.2016.07.013.
- 5 11. Buisson, C., and C. Ladier. Exploring the Impact of Homogeneity of Traffic Measurements on the 6 Existence of Macroscopic Fundamental Diagrams. Transportation Research Record: Journal of the 7 **Transportation** Research Board, Vol. 2124, No. 1, 2009, pp. 127–136. 8 https://doi.org/10.3141/2124-12.
- 9 12. Nagle, A. S., and V. V. Gayah. Accuracy of Networkwide Traffic States Estimated from Mobile
  10 Probe Data. *Transportation Research Record: Journal of the Transportation Research Board*, No.
  11 2421, 2014, pp. 1–11. https://doi.org/10.3141/2421-01.
- Du, J., H. Rakha, and V. V Gayah. Deriving Macroscopic Fundamental Diagrams from Probe Data: Issues and Proposed Solutions. *Transportation Research Part C: Emerging Technologies*, Vol. 66, 2016, pp. 136–149.
- 15 14. Daganzo, C. F., and L. J. Lehe. Traffic Flow on Signalized Streets. *Transportation Research Part B: Methodological*, Vol. 90, 2016, pp. 56–69. https://doi.org/10.1016/J.TRB.2016.03.010.
- Laval, J. A., and F. Castrillón. Stochastic Approximations for the Macroscopic Fundamental Diagram of Urban Networks. *Transportation Research Part B: Methodological*, Vol. 81, 2015, pp. 904–916. https://doi.org/10.1016/J.TRB.2015.09.002.
- 20 16. Leclercq, L., and N. Geroliminis. Estimating MFDs in Simple Networks with Route Choice. 21 Social and Behavioral Sciences. Vol. 80. 2013. 99–118. Procedia pp. 22 https://doi.org/10.1016/j.sbspro.2013.05.008.
- Tilg, G., S. Amini, and F. Busch. Evaluation of Analytical Approximation Methods for the
   Macroscopic Fundamental Diagram. *Transportation Research Part C: Emerging Technologies*, Vol.
   114, 2020, pp. 1–19. https://doi.org/10.1016/J.TRC.2020.02.003.
- Gayah, V. V., and C. F. Daganzo. Clockwise Hysteresis Loops in the Macroscopic Fundamental
   Diagram: An Effect of Network Instability. *Transportation Research Part B: Methodological*, Vol.
   45, No. 4, 2011, pp. 643–655. https://doi.org/10.1016/j.trb.2010.11.006.
- Mahmassani, H. S., M. Saberi, and A. Zockaie. Urban Network Gridlock: Theory, Characteristics,
  and Dynamics. *Transportation Research Part C: Emerging Technologies*, Vol. 36, 2013, pp. 480–
  497. https://doi.org/10.1016/j.trc.2013.07.002.
- Mazloumian, A., N. Geroliminis, and D. Helbing. The Spatial Variability of Vehicle Densities as
  Determinant of Urban Network Capacity. Vol. 368, No. 1928, 2010, pp. 4627–4647.
  https://doi.org/10.1098/rsta.2010.0099.
- Ji, Y., and N. Geroliminis. On the Spatial Partitioning of Urban Transportation Networks.
   *Transportation Research Part B: Methodological*, Vol. 46, No. 10, 2012, pp. 1639–1656.
   https://doi.org/10.1016/j.trb.2012.08.005.
- Saeedmanesh, M., and N. Geroliminis. Clustering of Heterogeneous Networks with Directional
   Flows Based on "Snake" Similarities. *Transportation Research Part B: Methodological*, Vol. 91,
   2016, pp. 250–269. https://doi.org/10.1016/j.trb.2016.05.008.
- 41 23. Saeedmanesh, M., and N. Geroliminis. Dynamic Clustering and Propagation of Congestion in
  42 Heterogeneously Congested Urban Traffic Networks. *Transportation Research Part B:*43 *Methodological*, Vol. 105, 2017, pp. 193–211.
- Lopez, C., P. Krishnakumari, L. Leclercq, N. Chiabaut, and H. van Lint. Spatiotemporal Partitioning 44 24. 45 of Transportation Network Using Travel Time Data. Transportation Research Record: Journal of 46 the *Transportation* Board, Vol. 2623, Research No. 1. 2017. pp. 98–107. 47 https://doi.org/10.3141/2623-11.
- 48 25. Geroliminis, N., and D. M. Levinson. Cordon Pricing Consistent with the Physics of Overcrowding.
   49 Transportation and Traffic Theory 2009: Golden Jubilee, 2009, pp. 219–240.
- 50 26. Daganzo, C. F., and L. J. Lehe. Distance-Dependent Congestion Pricing for Downtown Zones.

1

2

*Transportation Research Part B: Methodological*, Vol. 75, 2015, pp. 89–99. https://doi.org/10.1016/j.trb.2015.02.010.

- Zheng, N., R. A. Waraich, K. W. Axhausen, and N. Geroliminis. A Dynamic Cordon Pricing
  Scheme Combining the Macroscopic Fundamental Diagram and an Agent-Based Traffic Model. *Transportation Research Part A: Policy and Practice*, Vol. 46, No. 8, 2012, pp. 1291–1303.
  https://doi.org/10.1016/j.tra.2012.05.006.
- https://doi.org/10.1016/j.tra.2012.05.006.
  Yildirimoglu, M., M. Ramezani, and N. Geroliminis. Equilibrium Analysis and Route Guidance in Large-Scale Networks with MFD Dynamics. *Transportation Research Part C: Emerging Technologies*, Vol. 59, 2015, pp. 404–420. https://doi.org/10.1016/j.trc.2015.05.009.
- 10 29. Sirmatel, I. I., and N. Geroliminis. Economic Model Predictive Control of Large-Scale Urban Road 11 Networks via Perimeter Control and Regional Route Guidance. IEEE Transactions on Intelligent 12 *Transportation* Systems, Vol. 19. No. 4. 2018. pp. 1112-1121. 13 https://doi.org/10.1109/TITS.2017.2716541.
- Menelaou, C., S. Timotheou, P. Kolios, and C. G. Panayiotou. Joint Route Guidance and Demand
   Management for Real-Time Control of Multi-Regional Traffic Networks. *IEEE Transactions on Intelligent Transportation Systems*, 2021. https://doi.org/10.1109/TITS.2021.3077870.
- Haddad, J., and N. Geroliminis. On the Stability of Traffic Perimeter Control in Two-Region Urban
  Cities. *Transportation Research Part B: Methodological*, Vol. 46, No. 9, 2012, pp. 1159–1176.
  https://doi.org/10.1016/j.trb.2012.04.004.
- Aalipour, A., H. Kebriaei, and M. Ramezani. Analytical Optimal Solution of Perimeter Traffic Flow
   Control Based on MFD Dynamics: A Pontryagin's Maximum Principle Approach. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 20, No. 9, 2019, pp. 3224–3234.
   https://doi.org/10.1109/TITS.2018.2873104.
- 33. Haddad, J. Optimal Perimeter Control Synthesis for Two Urban Regions with Aggregate Boundary
  Queue Dynamics. *Transportation Research Part B: Methodological*, Vol. 96, 2017, pp. 1–25.
  https://doi.org/10.1016/j.trb.2016.10.016.
- 27 Su, Z. C., A. H. F. Chow, N. Zheng, Y. P. Huang, E. M. Liang, and R. X. Zhong. Neuro-Dynamic 34. 28 Programming for Optimal Control of Macroscopic Fundamental Diagram Systems. Transportation 29 Technologies, Research Part C: Emerging Vol. 116. 2020. p. 102628. 30 https://doi.org/10.1016/j.trc.2020.102628.
- 31 35. Zhou, D., and V. V. Gayah. Model-Free Perimeter Metering Control for Two-Region Urban
   32 Networks Using Deep Reinforcement Learning. *Transportation Research Part C: Emerging* 33 *Technologies*, Vol. 124, 2021, p. 102949.
- 36. Haddad, J. Robust Constrained Control of Uncertain Macroscopic Fundamental Diagram Networks.
   *Transportation Research Part C: Emerging Technologies*, Vol. 59, 2015, pp. 323–339.
   https://doi.org/10.1016/J.TRC.2015.05.014.
- 37 37. Li, Y., M. Yildirimoglu, and M. Ramezani. Robust Perimeter Control with Cordon Queues and
   38 Heterogeneous Transfer Flows. *Transportation Research Part C: Emerging Technologies*, Vol. 126,
   39 2021, p. 103043. https://doi.org/10.1016/j.trc.2021.103043.
- 40 38. Mohajerpoor, R., M. Saberi, H. L. Vu, T. M. Garoni, and M. Ramezani. H∞ Robust Perimeter Flow
  41 Control in Urban Networks with Partial Information Feedback. *Transportation Research Part B:*42 *Methodological*, Vol. 137, 2020, pp. 47–73. https://doi.org/10.1016/j.trb.2019.03.010.
- 39. Zhong, R. X., C. Chen, Y. P. Huang, A. Sumalee, W. H. K. Lam, and D. B. Xu. Robust Perimeter
  Control for Two Urban Regions with Macroscopic Fundamental Diagrams: A Control-Lyapunov
  Function Approach. *Transportation Research Part B: Methodological*, Vol. 117, 2018, pp. 687–707.
  https://doi.org/10.1016/j.trb.2017.09.008.
- 47 40. Sirmatel, I. I., and N. Geroliminis. Stabilization of City-Scale Road Traffic Networks via
  48 Macroscopic Fundamental Diagram-Based Model Predictive Perimeter Control. Control
  49 Engineering Practice, Vol. 109, 2021, p. 104750.
  50 https://doi.org/10.1016/j.conengprac.2021.104750.

- 41. Aboudolas, K., and N. Geroliminis. Perimeter and Boundary Flow Control in Multi-Reservoir
   Heterogeneous Networks. *Transportation Research Part B: Methodological*, Vol. 55, 2013, pp. 265–281. https://doi.org/10.1016/j.trb.2013.07.003.
- 4 42. Haddad, J., M. Ramezani, and N. Geroliminis. Cooperative Traffic Control of a Mixed Network
  with Two Urban Regions and a Freeway. *Transportation Research Part B: Methodological*, Vol.
  54, 2013, pp. 17–36. https://doi.org/10.1016/j.trb.2013.03.007.
- 54, 2013, pp. 17–36. https://doi.org/10.1016/j.trb.2013.03.007.
  43. Ramezani, M., J. Haddad, and N. Geroliminis. Dynamics of Heterogeneity in Urban Networks: Aggregated Traffic Modeling and Hierarchical Control. *Transportation Research Part B: Methodological*, Vol. 74, 2015, pp. 1–19. https://doi.org/10.1016/j.trb.2014.12.010.
- 10 44. Ni, W., and M. Cassidy. City-Wide Traffic Control: Modeling Impacts of Cordon Queues.
   11 *Transportation Research Part C: Emerging Technologies*, Vol. 113, 2020, pp. 164–175.
   12 https://doi.org/10.1016/j.trc.2019.04.024.
- 13 45. Sirmatel, I. I., D. Tsitsokas, A. Kouvelas, and N. Geroliminis. Modeling, Estimation, and Control in 14 Large-Scale Urban Road Networks with Remaining Travel Distance Dynamics. Transportation 15 Emerging Research Part C: Technologies. Vol. 128. 2021. 103157. p. 16 https://doi.org/10.1016/J.TRC.2021.103157.
- Haddad, J., and Z. Zheng. Adaptive Perimeter Control for Multi-Region Accumulation-Based
  Models with State Delays. *Transportation Research Part B: Methodological*, Vol. 137, 2020, pp. 133–153. https://doi.org/10.1016/J.TRB.2018.05.019.
- Zhong, R. X., Y. P. Huang, C. Chen, W. H. K. Lam, D. B. Xu, and A. Sumalee. Boundary Conditions and Behavior of the Macroscopic Fundamental Diagram Based Network Traffic Dynamics: A Control Systems Perspective. *Transportation Research Part B: Methodological*, Vol. 111, 2018, pp. 327–355. https://doi.org/10.1016/J.TRB.2018.02.016.
- 48. Haddad, J., and B. Mirkin. Coordinated Distributed Adaptive Perimeter Control for Large-Scale
  Urban Road Networks. *Transportation Research Part C: Emerging Technologies*, Vol. 77, 2017,
  pp. 495–515. https://doi.org/10.1016/j.trc.2016.12.002.
- 49. Lei, T., Z. Hou, and Y. Ren. Data-Driven Model Free Adaptive Perimeter Control for Multi-Region
  Urban Traffic Networks With Route Choice. *IEEE Transactions on Intelligent Transportation*Systems, 2019, pp. 1–12. https://doi.org/10.1109/tits.2019.2921381.
- So. Ren, Y., Z. Hou, I. I. Sirmatel, and N. Geroliminis. Data Driven Model Free Adaptive Iterative
   Learning Perimeter Control for Large-Scale Urban Road Networks. *Transportation Research Part C: Emerging Technologies*, Vol. 115, 2020, p. 102618. https://doi.org/10.1016/j.trc.2020.102618.
- Watkins, C. J. C. H., and P. Dayan. Q-Learning. *Machine Learning*, Vol. 8, No. 3–4, 1992, pp. 279–
   https://doi.org/10.1007/bf00992698.
- 35 52. Sutton, R. S., and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- Tsitsiklis, J. N., and B. Van Roy. An Analysis of Temporal-Difference Learning with Function
   Approximation. 1997.
- van Hasselt, H., Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep Reinforcement
   Learning and the Deadly Triad. 2018.
- 40 Mnih, V., K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. 55. 41 Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, 42 D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-Level Control through Deep 43 Reinforcement Learning. Nature, Vol. 518, No. 7540, 2015, pp. 529-533. 44 https://doi.org/10.1038/nature14236.
- 45 56. van Hasselt, H., A. Guez, and D. Silver. Deep Reinforcement Learning with Double Q-Learning.
  46 30th AAAI Conference on Artificial Intelligence, AAAI 2016, 2015, pp. 2094–2100.
- 47 57. Hessel, M., J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M.
  48 Azar, and D. Silver. Rainbow: Combining Improvements in Deep Reinforcement Learning. 32nd
  49 AAAI Conference on Artificial Intelligence, AAAI 2018, 2017, pp. 3215–3222.
- 50 58. Schaul, T., J. Quan, I. Antonoglou, and D. Silver. Prioritized Experience Replay. 2016.

- 59. Lillicrap, T. P., J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
   Continuous Control with Deep Reinforcement Learning. 2016.
- Wang, Z., T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas. Dueling Network
  Architectures for Deep Reinforcement Learning. *33rd International Conference on Machine Learning, ICML 2016*, Vol. 4, 2015, pp. 2939–2947.
- 6 61. Oliehoek, F. A., M. T. J. Spaan, and N. Vlassis. Optimal and Approximate Q-Value Functions for
   7 Decentralized POMDPs. *Journal Of Artificial Intelligence Research*, Vol. 32, 2008, pp. 289–353.
   8 https://doi.org/10.1613/jair.2447.
- 9 62. Koller, D., and R. Parr. Computing Factored Value Functions for Policies in Structured MDPs .
  10 1999.
- Rashid, T., M. Samvelyan, C. S. de Witt, G. Farquhar, J. Foerster, and S. Whiteson. QMIX:
   Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. 2018.
- Peng, B., T. Rashid, C. A. S. de Witt, P.-A. Kamienny, P. H. S. Torr, W. Böhmer, and S. Whiteson.
   FACMAC: Factored Multi-Agent Centralised Policy Gradients. 2021.
- 15 65. Horgan, D., J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. van Hasselt, and D. Silver.
   16 Distributed Prioritized Experience Replay. 2018.
- 17 Lin, L.-J. Self-Improving Reactive Agents Based on Reinforcement Learning, Planning and 66. 18 Teaching. Machine Learning, Vol. 8, No. 3-4. 293-321. 1992, pp. 19 https://doi.org/10.1007/bf00992699.
- Gao, X. (Shirley), and V. V. Gayah. An Analytical Framework to Model Uncertainty in Urban
  Network Dynamics Using Macroscopic Fundamental Diagrams. *Transportation Research Part B: Methodological*, Vol. 117, 2018, pp. 660–675. https://doi.org/10.1016/j.trb.2017.08.015.
- Wang, Y., B. Han, T. Wang, H. Dong, and C. Zhang. Off-Policy Multi-Agent Decomposed Policy
   Gradients. 2020.

25