

Perimeter Control Using Deep Reinforcement Learning: A Model-free Approach towards Homogeneous Flow Rate Optimization

Xiaocan Li¹, Ray Coden Mercurius², Ayal Taitler¹, Xiaoyu Wang², Mohammad Noeen³, Scott Sanner¹, and Baher Abdulhai²

Abstract—Perimeter control maintains high traffic efficiency within protected regions by controlling transfer flows among regions to ensure that their traffic densities are below critical values. Existing approaches can be categorized as either model-based or model-free, depending on whether they rely on network transmission models (NTMs) and macroscopic fundamental diagrams (MFDs). Although model-based approaches are more data efficient and have performance guarantees, they are inherently prone to model bias and inaccuracy. For example, NTMs often become imprecise for a large number of protected regions, and MFDs can exhibit scatter and hysteresis that are not captured in existing model-based works. Moreover, no existing studies have employed reinforcement learning for homogeneous flow rate optimization in microscopic simulation, where spatial characteristics, vehicle-level information, and metering realizations — often overlooked in macroscopic simulations — are taken into account. To circumvent issues of model-based approaches and macroscopic simulation, we propose a model-free deep reinforcement learning approach that optimizes the flow rate homogeneously at the perimeter at the microscopic level. Results demonstrate that our model-free reinforcement learning approach without any knowledge of NTMs or MFDs can compete and match the performance of a model-based approach, and exhibits enhanced generalizability and scalability.

I. INTRODUCTION

Perimeter control manages inter-regional transfer flows by gating the feeder links of protected regions, ensuring that each region’s vehicle accumulation does not surpass a critical value that would negatively impact traffic efficiency. Urban traffic networks frequently become oversaturated during rush hour, necessitating perimeter control. The delineation of a congested area is called perimeter identification, which can be executed statically or dynamically [1]. Once the perimeter is identified, perimeter control operates at the regional level [2], [3] and is particularly effective in oversaturated situations where adaptive traffic signal control alone is insufficient [4].

Recent works can be divided into reinforcement learning (RL) based control and non-RL-based control. Non-RL control comprises classical feedback control methods, such as Proportional Integral (PI) controller [5], [6], which aims

to maintain densities near set-points determined by MFDs. Analytical control approaches such as optimal control [7], model predictive control [8], and switching algorithm on top of a linear quadratic regulator [3] require an NTM and/or an MFD. Adaptive and learning control methods such as [2], [9], [10] need a model for tuning and are limited by the policy class they can approximate. For RL-based control, [11] employs reinforcement learning to heterogeneously redistribute a predetermined total regional flow rate among different feeder links at the microscopic level. However, the prior knowledge of the total rate is still obtained from a model-based non-RL controller. In [12], a deep RL approach is utilized to regulate the transfer flow, using the NTM as a plant to facilitate the macroscopic environment. However, the macroscopic simulation has its limitations.

The macroscopic simulation based on the NTM lacks consideration for: 1) spatial characteristics, which include road geometry, the spatial distribution of congestion, and traffic network topology; 2) vehicle-level information, which comprises stochastic route choices, driving aggressiveness, physical attributes such as car length and width, and generation times; 3) most importantly, metering realizations, such as traffic lights and flow metering, cannot be simulated at a macroscopic level, even though they are crucial for field implementation. Macroscopic simulations do not account for these factors, resulting in an environment that deviates from reality. Consequently, we are motivated to propose a model-free controller that learns through interaction with a microscopic environment that are more detailed and realistic.

In this work, we formulate perimeter control as a model-free RL problem, employing the Proximal Policy Optimization (PPO) algorithm. To validate the proposed method, we create a traffic network and demand profile in a microscopic traffic simulator, allowing the RL agent to interact and collect training and evaluation data. Then, we provide analysis of state designs, policy visualization, and various traffic metrics. Finally, we conduct a comprehensive evaluation of different demand scenarios for the PPO and baseline controllers.

Our contributions are the following: 1) to the best of our knowledge, we are the first to propose a model-free deep RL approach for homogeneous flow rate optimization of perimeter control at the microscopic level; 2) the incorporation of microscopic environment facilitates a trustworthy simulation for future field implementation; 3) extensive analysis of the proposed approach’s robustness and generalizability towards unseen demand scenarios is provided.

¹Xiaocan Li, Ayal Taitler and Scott Sanner are with the Department of Mechanical & Industrial Engineering, University of Toronto, Canada hsiaotsan.li@mail.utoronto.ca, ataitler@gmail.com, ssanner@mie.utoronto.ca

²Ray Coden Mercurius, Xiaoyu Wang, Baher Abdulhai are with the Department of Civil & Mineral Engineering, University of Toronto, Canada ray.mercurius@mail.utoronto.ca, cnxiaoyu.wang@mail.utoronto.ca, baher.abdulhai@utoronto.ca

³Mohammad Noeen is with the Dalla Lana School of Public Health, University of Toronto, Canada m.noeen@utoronto.ca

II. PROBLEM STATEMENT

A single protected region perimeter control problem is formulated for its simplicity. In Figure 1, the abstract traffic network consists of one protected region and several feeder links. A fixed control plan of traffic signals is used for intersections in the protected region. During high traffic demand periods, such as morning rush hour, incoming vehicles from feeder links are regulated by meterings to protect the region. Only incoming vehicles are restricted, while outgoing vehicles are unaffected. At each control timestep, the controller sets the allowed flow rate for each flow metering. To simplify the problem, we assume a homogeneous perimeter control, implying that all feeder links employ an identical flow metering rate. Flow meterings are used instead of intersection timings for gating due to their simpler control nature, which eases the complexity of validation.

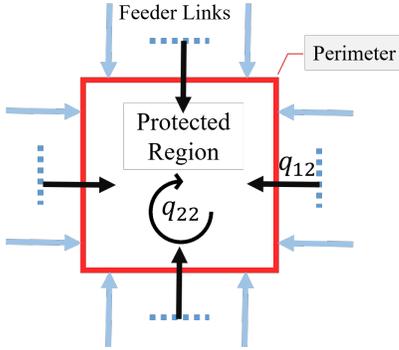


Fig. 1: A generic traffic network with single protected region.

A generic demand profile with a high peak of trips exceeding the protected region's capacity is required. Destinations are within the protected region, while origins are in both the protected region and feeder links. Demand flows are denoted by $q_{ij}(\text{veh/s}), i, j \in \{1, 2\}$, with 1 representing links outside the protected region and 2 representing the protected region. Therefore, q_{12} denotes flow from feeder links to the protected region, also known as exogenous flow; q_{22} denotes endogenous flow within the protected region; q_{21} denotes flow from protected region to exit links; and q_{11} denotes flow from feeder links to exit links. To simplify the demand profile, q_{11} and q_{21} are set to zero, as is in [11].

The objective of perimeter control is to minimize total delay. However, measuring delay is challenging due to difficulties in capturing stopping time and speed decrease for each vehicle. Instead, total time spent (TTS) is used as a surrogate, calculated as the difference between trip exit and generation times. Since TTS comprises both delay and travel time without delay, minimizing TTS is equivalent to minimizing total delay.

III. METHODOLOGY

A. Formulation as Markov Decision Processes

The perimeter control problem is formulated as a model-free reinforcement learning problem: At each timestep t , an agent takes an action a_t based on the current state s_t of

an environment and receives a scalar reward r_{t+1} from the environment. The reinforcement learning is formalized by the Markov Decision Process (MDP) described by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$:

1) State space \mathcal{S} : The design of the state is crucial for developing an agent that is robust and generalizable across various scenarios. Three traffic metrics serve as state variables: *I*) inner density D_p represents protected region congestion level defined as Eq. (1); *II*) feeder link density D_f defined as Eq. (2) is correlated with queue length outside the perimeter; *III*) future demand $n_{ij}, i, j \in \{1, 2\}$, derived from historical data, denotes the number of vehicles generated in the next control cycle T , defined as Eq. (3). Variables L_i, V_i represent the length of link i and the number of vehicles on link i , and \mathcal{P}, \mathcal{F} are the set of links for the protected region and feeder links, respectively. Three candidate designs of state space to be validated are defined in Eq. (4), (5), (6), where each design contains more information than previous ones. The definition in Eq. (6) is adopted from [12], while Eq. (4), (5) are our novel designs. Utilizing future demand suggests that the controller has access to a demand prediction model, albeit at a significantly lower level than a comprehensive model of environment dynamics;

$$D_p = \sum_{i \in \mathcal{P}} V_i / \sum_{i \in \mathcal{P}} L_i \quad (1)$$

$$D_f = \sum_{i \in \mathcal{F}} V_i / \sum_{i \in \mathcal{F}} L_i \quad (2)$$

$$n_{ij}(t) = \int_t^{t+T} q_{ij}(t') dt', \quad i, j \in \{1, 2\} \quad (3)$$

$$s_t^1 = [D_p] \quad (4)$$

$$s_t^2 = [D_p \quad D_f]^\top \quad (5)$$

$$s_t^3 = [D_p \quad D_f \quad n_{12} \quad n_{22} \quad n_{21} \quad n_{11}]^\top \quad (6)$$

2) Action space \mathcal{A} : The action is defined as the flow metering rate, which represents the number of vehicles allowed to pass the meter within a unit of time. The action space is one-dimensional due to the homogeneous control assumption. For practicality, the applied action is bounded by a minimum of 50 veh/h and a maximum of 300 veh/h per meter;

3) Transition dynamics: The dynamics from state s_t to new state s_{t+1} by taking action a_t . In a model-free setting, the dynamics is provided by the traffic simulator Aimsun and remains unknown to the agent;

4) Reward \mathcal{R} : The trip completion rate (veh/s) for each control cycle as a measure of traffic efficiency;

5) Discount factor γ : A scalar value ranged in $[0, 1]$ used to prioritize the current reward over future rewards. For our specific reward, $\gamma = 1$ cannot further stimulate the agent to learn when all trips are completed, as the total 'discounted' reward equals a constant of the total number of trips completed. In contrast, $\gamma < 1$ encourages the agent to complete all trips earlier;

The state value $V_\pi(s)$ defined in Eq. (7) as the expected

cumulative discounted reward over horizon H under policy π . The horizon is finite as the simulated demand shown in Figure 3 is time-limited. The objective of RL is to maximize the expected state value over a particular state distribution μ , and the optimal policy is found in Eq. (8).

$$V_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{H-1} \gamma^t r_t \mid s_t = s \right] \quad (7)$$

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{s \sim \mu} [V_{\pi}(s)] \quad (8)$$

B. Proximal Policy Optimization

Proximal Policy Optimization [13] is chosen as our RL algorithm due to its well-established reputation and ability to handle continuous action spaces such as flow metering rate. PPO is a policy gradient method that consists of an actor and a critic network. The actor network receives the state and generates action distributions, while the critic network receives the state and evaluates the state value. With these features, PPO is suitable for learning the optimal control policy in a perimeter control setting.

IV. EMPIRICAL RESULTS & DISCUSSION

A. Experimental Setup

Traffic Simulation Setup: Figure 2 depicts the testbed traffic network consisting of a protected region with 24 feeder links. Each traffic link is 170 meters long with two lanes per link, except for links connected to an origin or destination within the protected region have only one lane. In terms of the spatial distribution of origins and destinations, 100 points inside the protected region serve as both origins and destinations, while 24 points connected to each feeder link serve as origins only as $q_{21} = 0$ and $q_{11} = 0$.

In Figure 3, the traffic demand spanning 135 minutes is divided into nine 15-minute intervals. All origin-destination pairs share the same number of trips at each timestep, while the numbers of trips for nine intervals are determined by a vector of weights: $[r^0, r^1, r^2, r^3, r^4, r^3, r^2, r^1, r^0]$, where $r = 2$ serves as a hyperparameter to regulate the peakedness of the demand profile. The same weights are applied to both a total endogenous demand of 11000 vehicles and a total exogenous demand of 6000 vehicles.

During training, the random seed is fixed to 1 for a stable environment. However, during evaluation, distinct random seeds are used to assess the controller's performance more comprehensively and robustly under varying conditions.

RL Setup: The control cycle T is 96 seconds, same as the cycle of traffic signals in the protected region. Table I provides a comprehensive summary of PPO algorithm parameters and architecture for both actor and critic networks and key hyperparameters. State variables and rewards are scaled to the range $[0, 1]$ to ensure numerical stability.

B. Simulation Results

Three perimeter controllers are evaluated: No Perimeter Control (NPC) as a naïve baseline, PI controller as a model-based baseline, and our proposed PPO controller. Comparative analysis between NPC and PI controllers highlights the

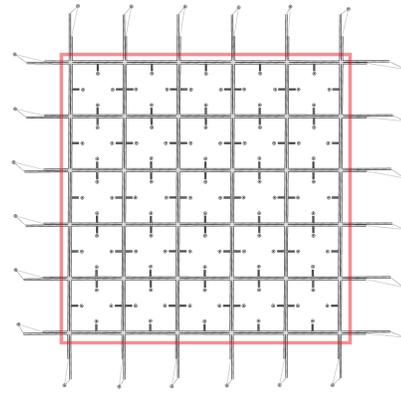


Fig. 2: Testbed: 5×5 grid network with single protected region and 24 feeder links.

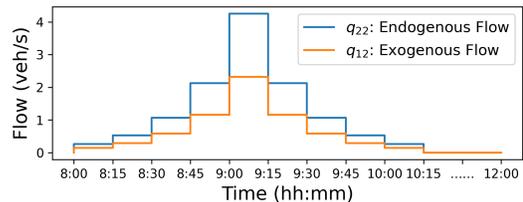


Fig. 3: Traffic demand profile.

TABLE I: Summary of PPO Parameters.

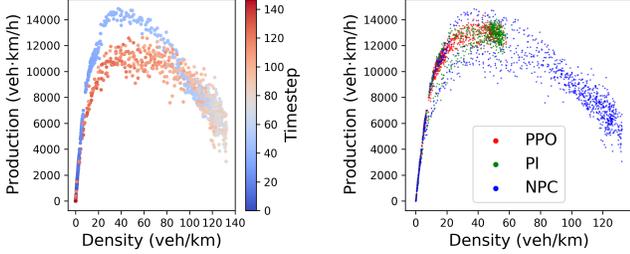
Parameter	Value
Actor Input Layer	$ S = 1, 2, 3$ units
Actor Hidden Layers	[256, 256] (Tanh activation)
Actor Output Layer	$ A = 1$ (ReLU activation)
Critic Input Layer	$ S = 1, 2, 3$ units
Critic Hidden Layers	[256, 256] (Tanh activation)
Critic Output Layer	1 (No activation)
Training Episodes	1200
Batch Size	25 episodes
Learning Rate	$5e-4$
Entropy Coefficient	0.01
Discount Factor γ	0.95, 0.9, 0.85

importance of perimeter control, while comparative analysis between PI and PPO emphasizes the benefits of the model-free RL approach. Evaluations are conducted using 10 replications with random seeds from 1 to 10, accounting for the inherent stochasticity of microscopic simulations. The total number of vehicles remains consistent across replications, with a maximum difference of 4%.

No Perimeter Control: Figure 4a illustrates the MFD of the protected region under the processes of loading and unloading of demand, labeled with cool and warm colors, respectively. In the loading process, the traffic production increases monotonically when the protected region is unsaturated. However, when the density of the protected region exceeds the critical density of 35 veh/km , traffic production declines, reaching a minimum production of only one-third of the maximum production when the maximal attained density is four times the critical density. In the unloading process, hysteresis is present in the MFD indicating the capacity drop, which occurs when the traffic network recovers

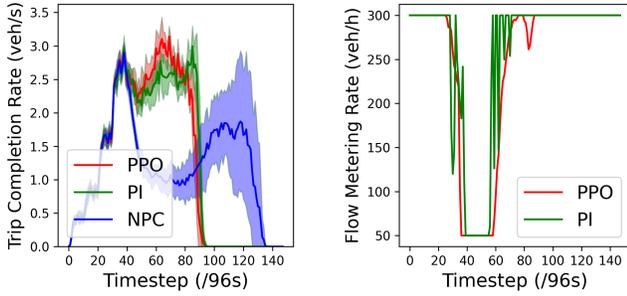
from extreme congestion.

This efficiency loss is also evident in the trip completion rate over time shown in Figure 5a. Before timestep 40, the rate gradually increases to around 3 *veh/s*, as the number of vehicles in the network rises. Between timestep 40 and 80, the rate decreases to approximately 1 *veh/s*. After timestep 80, the rate recovers from the loss; however, it does not recover to the previous rate, consistent with the production loss observed during the recovery process in the MFD.



(a) Loading and unloading of demand in NPC case. (b) MFD with and without perimeter control.

Fig. 4: MFDs under different controllers.



(a) Trip completion rate. (b) Flow metering rate.

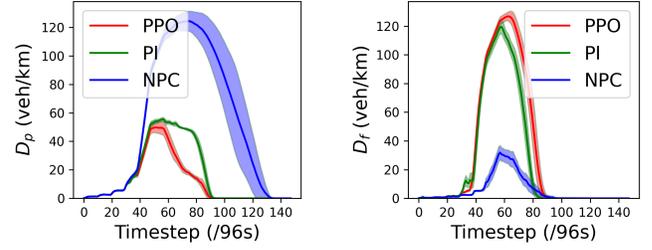
Fig. 5: Traffic metrics over time.

PI Controller: The PI controller is reproduced from [5] based on our testbed with coefficients of proportional and integral terms being $K_P = 36.88$ and $K_I = 1.24$, and the set-point being $\hat{N}_2 = 2700veh$. Variables $q_{in}(k)$, $N_2(k)$ in Eq. (9) are the total flow metering rate at timestep k , number of vehicles in the protected region at timestep k , respectively. The flow metering rate is clipped with the same lower and upper bounds as the PPO controller. Note that the reproduced PI controller may not be optimal due to parameter design being vulnerable to uncertainties in the choice of set-point from the MFD and slope estimation.

$$q_{in}(k) = q_{in}(k-1) - K_P[N_2(k) - N_2(k-1)] + K_I[\hat{N}_2 - N_2(k)] \quad (9)$$

Compared to NPC, the PI controller can keep the inner density near the critical density resulting in near-optimal traffic production, as is illustrated in Figure 4b of MFD and Figure 6a of inner density over time. For the trip completion rate, the PI controller maintains the same rate as NPC during

the early simulation, while the PI controller has a higher rate than NPC between timestep 40 and 90. Note that the area under the curve is the total number of finished trips, the rate reaches zero earlier under the PI control than NPC, which implies the PI controller clears all trips earlier.



(a) Inner density. (b) Feeder link density.

Fig. 6: Congestion level over time.

PPO Controller: In this part, the analysis of state space designs, and visualization of 1D, 2D and 6D state policies are provided. A comparison of controllers across various traffic metrics is also conducted to evaluate their effectiveness.

State Space Designs: Table II presents the results of three state designs across 10 replications. Our states evolve in the order of increasing variable complexity and attainment difficulty from a traffic engineering perspective. By calculating the performance improvement relative to using D_p only, including D_f in the state enhances the performance in 8 out of 10 replications, with an average improvement of 2.8%. Moreover, the incorporation of future demand into the state yields the best average performance among all state designs, with an average improvement of 7.8%.

1D State Policy: The simplest state design incorporates only the inner density D_p as a one-dimensional state. The learnt 1D state policy is visualized in Figure 7a exhibiting a three-piecewise linear function with 3 characteristics:

- 1) When D_p is lower than 8 *veh/km* being an undersaturated state, the policy employs the maximal metering rate to fully utilize the traffic network;
- 2) When D_p exceeds 25 *veh/km*, signifying a near or oversaturated state, the policy applies the minimal metering rate to alleviate the internal congestion;
- 3) For D_p between 8 *veh/km* and 25 *veh/km*, the flow metering rate decreases smoothly and linearly as D_p increases.

2D State Policy: The learnt policy is shown in Figure 7b with a color scheme representing D_f . The difference between 1D and 2D state policies is the 2D state policy enforces a higher flow metering rate with larger D_f when the protected region is uncongested, allowing a less restricted perimeter control to fully utilize the traffic network.

6D State Policy: The learnt policy is visualized in Figure 7c with a color scheme representing the future demand. For the ease of visualization, the feeder link density D_f is fixed to the average value across the simulation. The policy projection consists of several parallel piecewise linear functions resembling the 1D state policy, but shifted by future demand:

TABLE II: Performance* of different state designs for the PPO controller.

State\Seed	1	2	3	4	5	6	7	8	9	10	Mean
Inner Density Only	5050	5640	5143	5181	4495	5356	5326	5968	5429	4998	5259
Inner & Feeder Link Densities	4894 (3.1%)	5833 (-3.4%)	5064 (1.6%)	4793 (7.5%)	4573 (-1.7%)	4934 (7.9%)	5273 (1.0%)	5801 (2.8%)	5134 (5.4%)	4802 (3.9%)	5110 (2.8%)
Two Densities & Future Demand	4360 (13.7%)	5151 (8.7%)	5043 (1.9%)	4901 (5.4%)	4218 (6.1%)	4607 (14.0%)	5069 (4.8%)	5541 (7.2%)	4801 (11.6%)	4771 (4.5%)	4846 (7.8%)

*The performance is measured by the TTS (hour) and the percentage is the relative improvement to the 1D State case.

TABLE III: Performance* comparison of perimeter controllers on different replication seeds.

Controller\Seed	1	2	3	4	5	6	7	8	9	10	Mean
NPC	8614	14083	12471	12380	9537	13437	11800	14440	12359	11620	12074
PI	4840 (43.8%)	5706 (59.5%)	5458 (56.2%)	5328 (57.0%)	5332 (44.1%)	5410 (59.7%)	5664 (52.0%)	5680 (60.7%)	5459 (55.8%)	5299 (54.4%)	5417 (55.1%)
PPO (Ours)	4360 (49.4%)	5151 (63.4%)	5043 (59.6%)	4901 (60.4%)	4218 (55.8%)	4607 (65.7%)	5069 (57.0%)	5541 (61.6%)	4801 (61.1%)	4771 (58.9%)	4846 (59.9%)

*The performance is measured by the TTS (hour) and the percentage is the relative improvement to the NPC case.

For a fixed D_p , a larger future demand corresponds to a lower flow metering rate, logically resulting in stricter perimeter control when more vehicles are anticipated to be generated in both the protected region and feeder links. The most restrictive perimeter control, with a minimal flow metering rate of 50 *veh/h*, is implemented when the protected region is highly congested, even if the future demand is low.

In summary, the incorporation of feeder link density and future demand as state variables enhances the adaptiveness of our control policy as well as the performance in general.

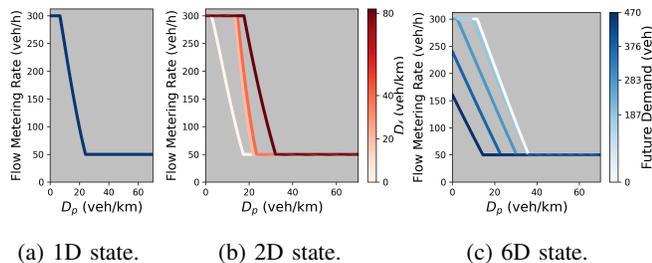


Fig. 7: Visualization of policies with various state designs.

Comparison of Controllers: The 6D state design for the PPO controller is chosen to compare with other controllers due to its outstanding performance among all state designs. Figure 8 displays the performance of all 10 replications: The scatter plot demonstrates that the PPO controller outperforms both the PI controller and the NPC in every replication. Furthermore, the kernel density estimation and box plot reveal a higher variability (low consistency) for the NPC case among all 10 replications. In contrast, the PI controller exhibits the lowest variability, although its average performance is slightly inferior to that of the PPO controller. The numerical results for Figure 8 are presented in Table III. On average, the PPO controller reduces the TTS by 59.9% and 10.5% compared to the NPC and PI controllers, respectively.

Figure 6a illustrates that the PPO controller maintains the lowest D_p throughout the simulation, while NPC has the highest D_p . Figure 6b indicates that the PPO controller exhibits the highest D_f as the controller retains incoming ve-

hicles outside the protected region to preserve a higher traffic production inside. As shown in Table IV, the reduced TTS inside the protected region is 9318 hours, which significantly surpasses the increased TTS on feeder links, totaling 2088 hours. Consequently, the PPO controller reduces the average TTS by 7220 hours compared to the NPC case.

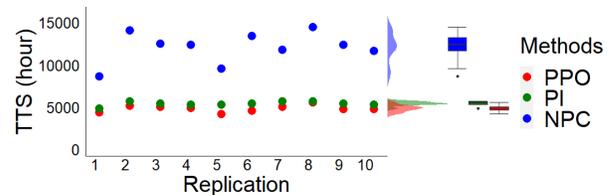


Fig. 8: Performance comparison on every replication.

TABLE IV: TTS (hour) inside & outside protected region.

	TTS inside	TTS outside	TTS
PPO	2425	2422	4847
NPC	11743	334	12077
Difference	-9318	2088	-7220

In Figure 5b, both the PPO and PI controllers exhibit highly dynamic actions over time for replication 1. The PI controller’s action oscillates considerably, while the PPO controller displays a much smoother action progression. Additionally, between timestep 20 and timestep 90, the PPO controller employs lower actions than the PI controller to have a more restricted control.

Figure 5a presents the trip completion rate under different perimeter controllers. Before timestep 38, the rates of the PPO and PI controllers are nearly identical to the NPC case since the entire network is loading demand from an empty state. Between timestep 38 and timestep 85, the PPO controller achieves a higher rate than the PI controller and reaches a zero rate earliest, indicating that all trips are completed. However, The NPC has a 50% lower rate compared to either PPO or PI controllers and all trips are only completed near timestep 140. Additionally, the variability for the NPC case is the highest among the compared controllers.

Generalizability Tests: The objective of the generalizability tests is to evaluate the controller’s performance under out-of-distribution scenarios not experienced during training. Two generalizability tests are performed: 1) a scalability test scales the total number of vehicles up or down while maintaining a fixed peakedness; 2) a demand peakedness test varies the demand shape from a flat one to a highly peaked one while keeping the total number of vehicles constant. Results are presented in figures 9a and 9b correspondingly.

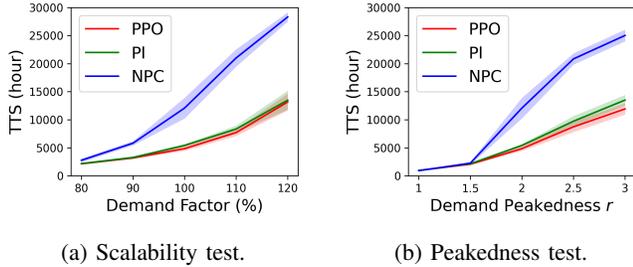


Fig. 9: Generalizability tests for robustness.

In the scalability test, the demand is adjusted with factors ranging from 80% to 120%, incremented by 10%. Each controller is trained exclusively on the 100% factor. In Figure 9a, the PPO controller achieves the lowest average TTS, followed by the PI controller and NPC. Moreover, the performance gap between the PPO controller and NPC widens as the demand factor increases, indicating that the benefits of perimeter control are more significant when faced with a more congested scenario. This suggests that perimeter control is beneficial only when a traffic network experiences production degradation due to congestion; otherwise, implementing perimeter control is unnecessary.

In the peakedness test, the demand profile’s peakedness is changed with the parameter r from 1.0 to 3.0 in increments of 0.5. All controllers are trained exclusively on demand peakedness $r = 2$. In Figure 9b, when the demand profile is entirely or nearly uniform over time ($r = 1.0$ and $r = 1.5$), the performances of the three controllers are quite similar, indicating that perimeter control may not be needed. As r increases beyond 1.5, the performance gap between the NPC and both the PPO and PI controllers widens with the peakedness. Moreover, the PPO controller consistently outperforms the PI controller in all tested scenarios, demonstrating that the PPO controller remains robust under changes in demand peakedness relative to the baselines.

The performance disparity between the PI controller and the PPO controller is rooted in their characteristics. The PI controller is a regulator, while the PPO controller is an optimal controller. The PI controller does not explicitly strive for optimality and relies on an accurately identified set-point. If this set-point is not precise, the performance of the PI controller may not match that of the PPO controller.

V. CONCLUSION AND FUTURE WORK

We have explored an RL-based homogeneous flow rate optimization for perimeter control that does not require

NTMs or MFDs. The evaluation of control policies is performed with a grid-patterned traffic network. In addition, the choice of each state variable is supported by an in-depth policy visualization analysis. Moreover, our model-free approach matches a standard model-based PI controller and proves to be more robust under scalability and demand peakedness tests. The exploration of the characteristics and requirements of a model-free approach demonstrates that it is more adaptive to uncertainty and hence less susceptible to model inaccuracies.

Future research directions include model-free heterogeneous flow rate optimization for enhanced perimeter queue balancing and management, the replacement of flow meterings with traffic signals on the perimeter, and an in-depth comparison of model-based RL perimeter control methods against model-free ones.

REFERENCES

- [1] Y. Guo, L. Yang, S. Hao, and J. Gao, “Dynamic identification of urban traffic congestion warning communities in heterogeneous networks,” *Physica A: Statistical Mechanics and its Applications*, vol. 522, pp. 98–111, 2019.
- [2] Y. Ren, 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, p. 102618, 2020.
- [3] Y. Li, R. Mohajerpoor, and M. Ramezani, “Perimeter control with real-time location-varying cordon,” *Transportation Research Part B: Methodological*, vol. 150, pp. 101–120, 2021.
- [4] M. Keyvan-Ekbatani, X. Gao, V. V. Gayah, and V. L. Knoop, “Traffic-responsive signals combined with perimeter control: investigating the benefits,” *Transportmetrica B: Transport Dynamics*, vol. 7, no. 1, pp. 1402–1425, 2019.
- [5] M. Keyvan-Ekbatani, A. Kouvelas, I. Papamichail, and M. Papageorgiou, “Exploiting the fundamental diagram of urban networks for feedback-based gating,” *Transportation Research Part B: Methodological*, vol. 46, no. 10, pp. 1393–1403, 2012.
- [6] M. Keyvan-Ekbatani, M. Papageorgiou, and V. L. Knoop, “Controller design for gating traffic control in presence of time-delay in urban road networks,” *Transportation Research Procedia*, vol. 7, pp. 651–668, 2015.
- [7] A. Aalipour, 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, pp. 3224–3234, 2018.
- [8] M. Ramezani, J. Haddad, and N. Geroliminis, “Dynamics of heterogeneity in urban networks: aggregated traffic modeling and hierarchical control,” *Transportation Research Part B: Methodological*, vol. 74, pp. 1–19, 2015.
- [9] A. Kouvelas, M. Saeedmanesh, and N. Geroliminis, “Enhancing model-based feedback perimeter control with data-driven online adaptive optimization,” *Transportation Research Part B: Methodological*, vol. 96, pp. 26–45, 2017.
- [10] Z. Su, A. H. Chow, N. Zheng, Y. Huang, E. Liang, and R. Zhong, “Neuro-dynamic programming for optimal control of macroscopic fundamental diagram systems,” *Transportation Research Part C: Emerging Technologies*, vol. 116, p. 102628, 2020.
- [11] W. Ni and M. J. Cassidy, “Cordon control with spatially-varying metering rates: A reinforcement learning approach,” *Transportation Research Part C: Emerging Technologies*, vol. 98, pp. 358–369, 2019.
- [12] D. Zhou and V. V. Gayah, “Model-free perimeter metering control for two-region urban networks using deep reinforcement learning,” *Transportation Research Part C: Emerging Technologies*, vol. 124, p. 102949, 2021.
- [13] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.