# UNITSA: A UNIVERSAL REINFORCEMENT LEARNING FRAMEWORK FOR V2X TRAFFIC SIGNAL CONTROL

**Maonan Wang** *, **Xi Xiong** †, **Yuheng Kan, Chengcheng Xu** ‡, **Man-On Pun** §¶

## ABSTRACT

Traffic congestion is a persistent problem in urban areas, which calls for the development of effective traffic signal control (TSC) systems. While existing Reinforcement Learning (RL)-based methods have shown promising performance in optimizing TSC, it is challenging to generalize these methods across intersections of different structures. In this work, a universal RL-based TSC framework is proposed for Vehicle-to-Everything (V2X) environments. The proposed framework introduces a novel agent design that incorporates a junction matrix to characterize intersection states, making the proposed model applicable to diverse intersections. To equip the proposed RL-based framework with enhanced capability of handling various intersection structures, novel traffic state augmentation methods are tailor-made for signal light control systems. Finally, extensive experimental results derived from multiple intersection configurations confirm the effectiveness of the proposed framework. The source code in this work is available at `https://github.com/wmn7/Universal_Light`

***Keywords*** Traffic signal control · Universal models · Reinforcement learning · Traffic state augmentation

## 1 Introduction

Traffic congestion presents a critical challenge in urban areas worldwide, leading to wasted time for individuals, excessive fuel consumption, and increased greenhouse gas emissions [1]. To alleviate congestion, conventional traffic signal control (TSC) methods such as fixed-cycle traffic control [2], the Webster method [3], and Self-Organizing Traffic Light Control (SOTL) [4] have been developed. However, as cities continue to grow, these traditional traffic management approaches often prove insufficient to handle increasing traffic volumes and dynamic road conditions [5]. The emergence of V2X technologies in transportation systems offers a promising solution by facilitating communication and data exchange among vehicles, infrastructure, and other road users [6]. By leveraging real-time data derived from vehicles, the traffic management infrastructure can efficiently regulate vehicle and pedestrian movements at intersections by dynamically adapting to the changing traffic conditions [7].

To control signal lights based on real-time traffic conditions, various RL-based methods have been proposed. These methods employ three primary approaches for adjusting traffic lights, namely "Choose the next phase", "Keep or change the phase" and "Set the phase duration". More specifically, in "Choose the next phase", the RL agent determines the next phase to activate, allowing for a flexible phase sequence rather than a predetermined one [8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22]. While this approach offers flexibility, it has the potential to confuse drivers since the phase selection may appear random, which could incur increased risks of traffic accidents. For the "Keep or change the phase" approach, the RL agent decides whether to maintain or change the current phases [23, 24, 25]. Finally, the "Set the phase duration" approach selects the optimal duration for the current phase from a set of predefined options [26, 27, 28]. Through direct interaction with the environment, the RL agent effectively learns to adapt to traffic condition changes based on real-world experiences.

---

*School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, China and Shanghai AI Laboratory, Shanghai, China.

†Key Laboratory of Road and Traffic Engineering, Ministry of Education, Tongji University, Shanghai, China.

‡SenseTime Group Limited and Shanghai AI Laboratory, Shanghai, China.

§School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, Shenzhen 518172, China.

¶Corresponding author: SimonPun@cuhk.edu.cn.

Despite the significant improvements achieved by the aforementioned RL-based methods, a major limitation is that they are designed for interactions of designated structures. In other words, these RL models have to be redesigned and re-trained from scratch when dealing with intersections of different structures including approaching roads, lanes, and phases, which incurs substantial resources in terms of data collection, model development, and testing [29]. Thus, it is crucial to develop universal models that can be easily adapted and deployed across a wide range of intersections. As a result, the implementation of V2X can be efficiently scaled up without requiring extensive customization or redevelopment for each individual intersection [30]. In the literature, several generalized models were proposed for different junctions [11, 18, 10, 12]. Despite their good performance, these generalized models are only applicable to those specific intersection configurations considered in their design. Furthermore, these models exhibit performance degradation when encountering intersections of unseen configurations.

Motivated the discussions above, we present a **Uni**versal **T**raffic **S**tate **A**ugmentation RL-based framework (UniTSA). The proposed framework enables the training of a universal agent using augmented data for TSC. In order to handle intersections of diverse configurations, a junction matrix is developed to describe the intersection state. By using this newly proposed matrix, intersections of different structures can be characterized by matrices of the same size. In addition, the "Keep or change the current phase" approach is employed as the action design in this work, ensuring consistent model structures across intersections of different configurations. To cope with unseen intersection structures, five traffic state augmentation methods are developed to enrich the agent's data collection during the training process. These augmentation methods provide more comprehensive training, which ultimately leads to improved performance in intersections of configurations not available in the training set. Furthermore, the Low-Rank Adaptation (LoRA) [31] is employed to further fine-tune the model for crucial intersections. Finally, extensive experiments using the Simulation of Urban MObility (SUMO) [32] platform are conducted by taking into account intersections of various approaching roads, lanes, and phases. The experimental results demonstrate that the proposed UniTSA model achieves excellent performance even in unseen intersections.

Our contributions can be summarized as follows:

- An adaptive TSC framework called UniTSA is proposed for V2X by leveraging a universal RL model built upon novel agent designs that can handle diverse intersection structures. In addition, a fine-tuning mechanism is devised to further enhance the performance in key intersections;
- Traffic state augmentation methods are developed for the proposed TSC framework, enhancing the agent's understanding of diverse intersections with improved performance in both training and testing sets;
- Extensive experiments on 12 intersections of diverse structures confirm that the proposed UniTSA model substantially outperforms conventional universal models. Moreover, given a new intersection, UniTSA can fine-tune its pre-trained model to achieve comparable or even better performance with significantly reduced training time as compared to training a new model from scratch.

The remainder of the paper is structured as follows: Section 2 provides a summary on the TSC-related work whereas Section 3 introduces the terminology related to road and traffic signals. After that, Section 4 presents the proposed UniTSA framework and the five traffic state augmentation methods before Section 5 elaborates on the experimental setup and results on UniTSA. Finally, Section 6 concludes the paper.

## 2 Related Work

Extensive research has been conducted in the field of transportation to study TSC. Conventionally, fixed-time signal control is one of the earliest and widely used approaches in TSC [33]. These methods rely on pre-determined signal timings based on historical traffic patterns or engineering guidelines. Various optimization techniques have been proposed to determine optimal fixed-time plans for specific intersection configurations, where the Webster [3] method is one of the most successful TSC methods for the single intersection case. It can calculate the cycle length and phase split for a single intersection according to the traffic volume during a certain period (i.e., past 15 minutes or 30 minutes). However, such fixed-time control methods often incur suboptimal performance due to their lack of adaptability to dynamically changing traffic conditions. To cope with this problem, actuated control methods such as Sydney Coordinated Adaptive Traffic System (SCATS) [34], Max-pressure control [35] and Self-Organizing Traffic Light Control (SOTL) [4] were designed to adaptively adjust signal timings based on real-time traffic demand. Despite their many advantages, these actuated control methods are handicapped by the necessity that expert settings are required for each intersection. Furthermore, the performance of these actuated control methods degrades in complex scenarios.

Recently, RL-based TSC methods have attracted substantial attention due to their outstanding adaptability to real-time traffic conditions and impressive capability of learning the optimal control policies in complex scenarios [36]. Generally speaking, these RL-based TSC methods can be categorized into three categories, namely the valued-based

methods [8, 23, 25, 9, 10, 11, 12, 13, 14], the policy-based methods [15, 16, 17, 18, 19] and the actor-critic methods [20, 28, 21, 22]. Despite their good performance, most existing RL-based TSC methods focus on training models for specific intersection configurations or scenarios. A few attempts on training generalized TSC models have been made in the literature. For instance, MetaLight [11] trains a more universal model by incorporating the meta-learning strategy proposed in [37]. However, MetaLight requires re-training of its model parameters for each new intersection encountered. To overcome this shortcoming, [18, 10, 12] established universal models through parameter sharing. However, these methods do not preserve the original phase structure of traffic lights. In contrast, our proposed method can maintain the original signal light structure while fine-tuning crucial intersections to achieve significantly improved performance.
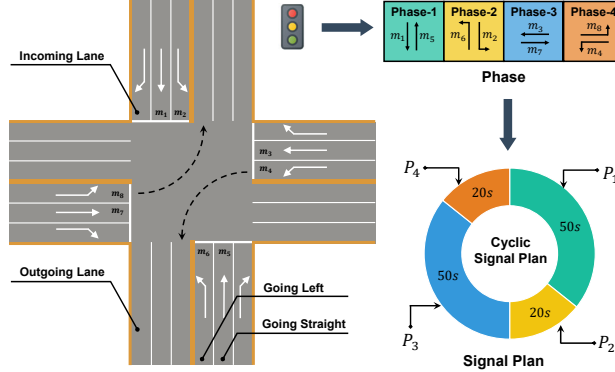


Figure 1: A standard 4-way intersection illustration.

## 3   Preliminary

In this section, some concepts used in this work are defined using a standard four-way intersection as depicted in Fig. 1 as an example. These concepts can be easily extended to the intersections of different structures.

- Lane: A lane refers to a roadway that provides a defined path for vehicles to travel in a specific direction. At a typical intersection, there are two types of lanes: incoming lanes $l_{in}$ (where vehicles enter) and outgoing lanes $l_{out}$ (where vehicle exit);

- Traffic Movement: A traffic movement refers to the connection between an incoming lane $l_{in}$ to an outgoing lane $l_{out}$. For the common 4-way intersection in the left side of Fig. 1, there are 12 movements in total, including right turns, left turns, and through movements in each of the four directions;

- Movement Signal: A movement signal is defined on the traffic movement. The green signal means the corresponding movement is allowed whereas the red signal is prohibited. As the right-turn traffic can move regardless of the signal, only eight movement signals out of the 12 possible movements in a 4-way intersection are used. More specifically, these eight movements denoted by $M_1, M_2, \cdots, M_8$ are Northbound (N), Northbound Left-turn (NL), Eastbound (E), Eastbound Left-turn (EL), Westbound(W), Westbound Left-turn (WL), Southbound (S), Southbound Left-turn (SL) movements, respectively. For instance, $m_8$ indicates that the vehicles can travel from west to north;

- Phase: A phase is a combination of movement signals. Each phase allows a specific group of traffic movements to occur while restricting others. The top-right portion of Fig. 1 illustrates the four phases of a 4-way intersection. For instance, Phase-1 involves $m_1$ and $m_5$, enabling vehicles traveling north-south to proceed while simultaneously prohibiting other movements;

- Signal Plan: A signal plan represents a prearranged sequence and duration of phases used to control the traffic signals at an intersection. Mathematically, a signal plan is denoted by $\{(p_1, t_1), (p_2, t_2), \cdots, (p_i, t_i), \cdots\}$, where $p_i$ and $t_i$ represent a phase and its duration, respectively. Usually, the phase sequence is in a cyclic order. For instance, the bottom-right portion of Fig. 1 shows a cycle-based signal plan and the duration of each phase is $t_1 = 50$, $t_2 = 20$, $t_3 = 50$ and $t_4 = 20$ as an example.
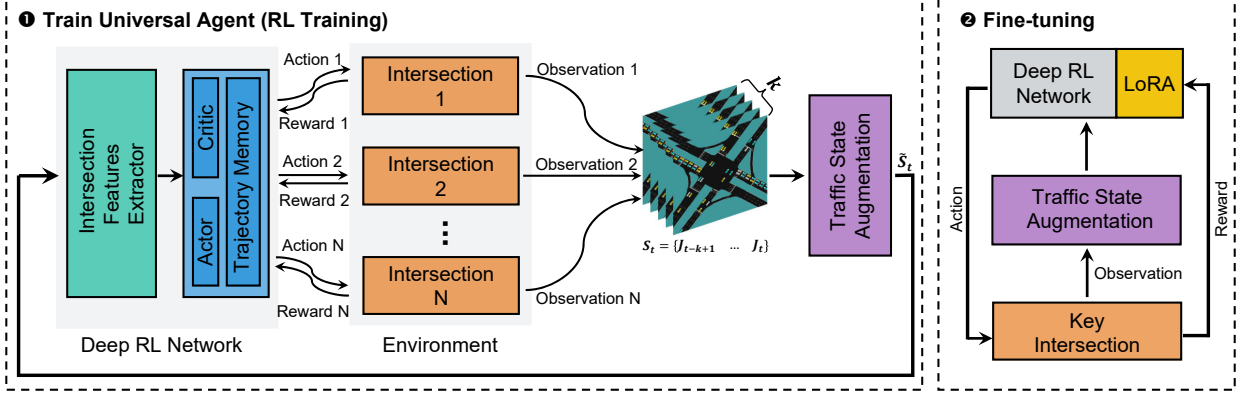
3

Figure 2: The overall framework of UniTSA.

# 4 Methodology

## 4.1 Framework

As shown in Fig. 2, the proposed UniTSA framework consists of two modules, one designed to train a universal agent through different intersections and one for fine-tuning the model for key intersections. More specifically, the first module capitalizes on a novel RL agent design to allow the model to characterize intersections of various topologies and signal schemes using the same structure by incorporating features on movements and actions with *next or not* while exploiting five novel traffic state augmentation methods. After that, the second module fine-tunes the model derived by the first module for some specific key intersections. More details about each step above will be elaborated in the following sections.

## 4.2 Agent Design and Junction Matrix

**State**: Intersections may possess different numbers of lanes, which results in state space of different dimensions when recording features in lane units. As discussed in Section 3, it is observed that there are only eight valid movement signals in a 4-way intersection, regardless of the number of lanes. Inspired by this observation, we propose to represent the information of an intersection at time $t$ as a junction matrix denoted as

$$\boldsymbol{J}_t = \left[\boldsymbol{m}_1^t, \boldsymbol{m}_2^t, \cdots, \boldsymbol{m}_8^t\right]^T, \tag{1}$$

where $[\cdot]^T$ stands for the transpose of the enclosed matrix while vector $\boldsymbol{m}_i^t$ of length eight represents information extracted from the $i$-th movement at time $t$ for $i = 1, 2, \cdots, 8$. More specifically, $\boldsymbol{m}_i^t$ encompasses three components: *traffic characteristics*, *movement characteristics*, and *traffic signal characteristics*. The *traffic characteristics* quantify the congestion level of the movement using parameters such as the average traffic flow $F^{i,t}$, the maximum occupancy $O_{max}^{i,t}$, and the average occupancy $O_{mean}^{i,t}$ within two consecutive actions. In addition, the *movement characteristics* provide specific details about the movement itself, including the direction ($I_s^i$) indicating whether it is a straight movement or not and the number of lanes ($L_i$) it occupies. Finally, the *traffic signal characteristics* comprise three binary parameters, namely $I_{cg}^{i,t}$, $I_{ng}^{i,t}$ and $I_{mg}^{i,t}$. These three parameters indicate whether the movement signal is currently green or not, whether it will be green in the next signal phase or not, and whether the current green duration has reached the minimum duration or not, respectively. These eight features can be readily obtained, making this design suitable for practical deployment. Therefore, vector $\boldsymbol{m}_i^t$ is defined as follows:

$$\boldsymbol{m}_i^t = \left[F^{i,t}, O_{max}^{i,t}, O_{mean}^{i,t}, I_s^i, L_i, I_{cg}^{i,t}, I_{ng}^{i,t}, I_{mg}^{i,t}\right]^T. \tag{2}$$

It has been observed that $\boldsymbol{J}_t$ at one time instance is not sufficient to provide a comprehensive understanding of the traffic dynamics for RL-based signal light controllers. To circumvent this obstacle, we propose to incorporate temporal traffic information by exploiting the $K$ latest observations, which enables the agent to better capture traffic patterns and trends in traffic behaviors. As a result, the agent can more effectively adapt its decision-making process in response to evolving traffic conditions. Mathematically, we define the state $\boldsymbol{S}_t \in \mathbb{R}^{K \times 8 \times 8}$ of the proposed agent at time $t$ as shown in Eq. (3):
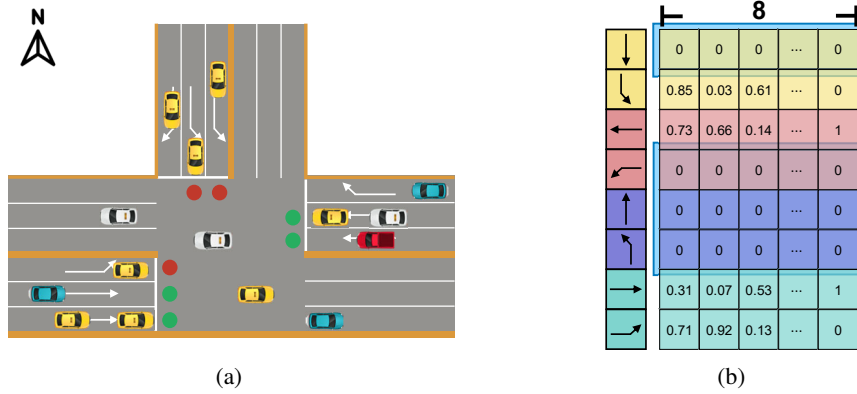
Figure 3: (a) A 3-way intersection with (b) its junction matrix with zero padding.

$$\boldsymbol{S}_t = [\boldsymbol{J}_{t-K+1}, \boldsymbol{J}_{t-K+2}, \cdots, \boldsymbol{J}_t].  \tag{3}$$

Finally, it is worth pointing out that zero padding is employed when the number of movements at an intersection is less than 8, e.g. a 3-way intersection. For instance, Fig. 3a shows a common 3-way intersection in which only E, EL, W, and SL movement signals are in use. As shown in Fig. 3b, zero padding is applied on the rows corresponding to those 4 unused movement signals, maintaining the matrix size identical to those of the 4-way intersection.

**Action**: A realistic and implementable action design has to take into account the safety of all traffic participants. While the action design *choose next phase* [10, 18, 12] can significantly improve the intersection efficiency, it disrupts the original cycle sequence of signal lights, thereby compromising driver safety. In sharp contrast, this work adopts the action design of "Keep or change" [23, 25, 9]. This design adheres to the concept of cycles, executing each phase sequentially (e.g., Phase 1, Phase 2, Phase 3, Phase 4, Phase 1, Phase 2, and so on). The agent determines whether to keep the current phase or change to the next phase based on the state $\boldsymbol{S}_t$. Due to the availability of current phase information $I_{cg}$ and next phase information $I_{ng}$ in the junction matrix, it is later shown that this action design exhibits excellent scalability for intersections with different signal plans.

**Reward**: The negative of the average queue length in each movement $q_i$ is adopted as the reward. Metrics such as waiting time, travel time, and delay are not used since it is impractical to obtain these metrics from real-world traffic detection devices. Consequently, the proposed reward function is defined as follows:

$$r_t = \frac{\left(-\sum_{i=1}^{8} q_i\right) - \mu}{\sigma + \epsilon},  \tag{4}$$

where $\epsilon$ is a small number to prevent division by zero. Furthermore, $\mu$ and $\sigma$ represent the mean and standard deviation of the first $R$ rewards, respectively. Mathematically, $\mu$ and $\sigma$ take the following form:

$$\mu = \frac{1}{R-1} \sum_{j=1}^{R-1} r_j,  \tag{5}$$

$$\sigma = \sqrt{\frac{1}{R-1} \sum_{j=1}^{R-1} (r_j - \mu)^2}.  \tag{6}$$

The reward is normalized to facilitate a faster training process.

### 4.3 Traffic State Augmentation

In recent studies, data augmentation techniques have demonstrated their effectiveness in enhancing the generalization capabilities of RL models [38, 39, 40]. By training on a more diverse set of augmented samples, RL agents can improve
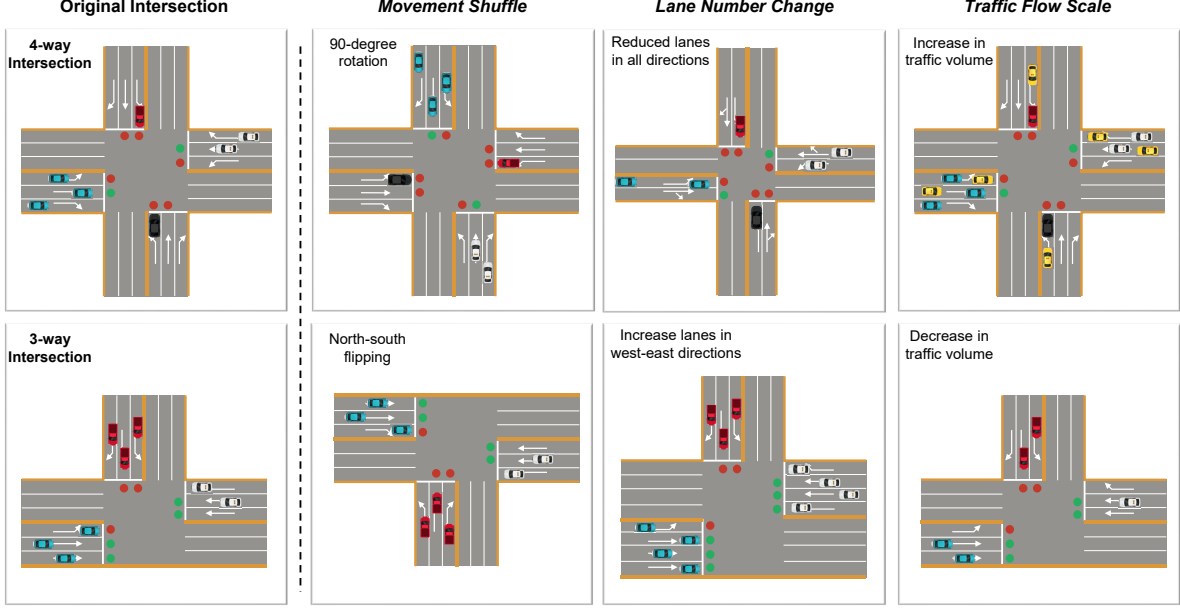
Figure 4: Illustration of three traffic state augmentation methods applied to both 4-way and 3-way intersections.
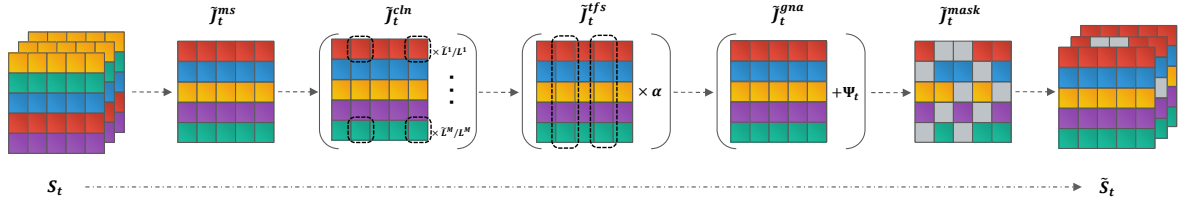


Figure 5: Detailed steps in traffic state augmentation block.

their capability of handling unseen tasks. The most common data augmentation techniques reported in the literature are *Gaussian noise addition*, and *masking*. In this work, we propose three additional novel traffic state augmentation methods specifically designed for RL-based TSC tasks, namely *movement shuffling*, *change of lane numbers* and *traffic flow scaling* as illustrated in Fig. 4.

Fig. 5 provides a detailed illustration of the step-by-step process involving the traffic state $S_t$ as it undergoes the five data augmentation methods, leading to the final output $\tilde{S}_t$. During training, a minibatch of data is randomly sampled from the replay buffer or recently augmented trajectories. While augmentation across the minibatch is stochastic, it is consistent across the stacked frames. It is also worth noting that these five traffic state augmentation methods can be directly applied to the junction matrix $J$, enabling the agent to learn and adapt to different traffic scenarios and intersection structures.

**Movement Shuffling**: This method involves shuffling the rows of the junction matrix, simulating different rotations, and flipping at the same intersection. Intuitively, shuffling can be interpreted as an effective rotation of the original intersection as depicted in Fig. 4. The assumption behind shuffling is that the action taken by the agent should not change after the rotation of the intersection. Mathematically, the movement shuffling operation can be modeled as follows:

$$\tilde{\boldsymbol{J}}_t^{\mathrm{ms}} = \boldsymbol{P} \cdot \boldsymbol{J}_t, \tag{7}$$

where $\boldsymbol{J}_t$ and $\tilde{\boldsymbol{J}}_t^{\mathrm{ms}}$ represent the original and augmented junction matrices, respectively. Furthermore, $\boldsymbol{P}$ is a permutation matrix designed to exchange rows in $\boldsymbol{J}_t$. After "movement shuffling" is applied to all junction matrices in $\tilde{\boldsymbol{S}}_t$ in Eq. (3), the new state can be represented as $\tilde{\boldsymbol{S}}_t^{\mathrm{ms}}$:

$$\tilde{\boldsymbol{S}}_t^{\mathrm{ms}} = \left[ \tilde{\boldsymbol{J}}_{t-K+1}^{\mathrm{ms}}, \tilde{\boldsymbol{J}}_{t-K+2}^{\mathrm{ms}}, \cdots, \tilde{\boldsymbol{J}}_t^{\mathrm{ms}} \right]. \tag{8}$$

**Change of Lane Numbers**: To expose the agent to a wider range of road structure combinations, we propose to randomly modify the number of lanes $L_i$ in each movement vector $\tilde{\boldsymbol{m}}_i^{\mathrm{ms}} \in \tilde{\boldsymbol{J}}_t^{\mathrm{ms}}$. This augmentation method allows the agent to encounter various lane configurations during training, enhancing its capability of handling diverse intersection layouts. In addition, the traffic characteristics (e.g., traffic flow and occupancy) are multiplied by the corresponding coefficients to maintain relative values.

An example of this traffic state augmentation method is shown in Fig. 4. In this example, a 4-way intersection was modified to two in all directions by reducing the number of lanes. Similarly, for the 3-way intersection, we increase the number of west and east-approaching lanes to 4. While modifying the number of lanes, we maintain the relative number of vehicles in each direction. Thus, it is reasonable to assume that the actions taken by the agent before and after such modifications should be identical. Mathematically, the operation to change the lane number can be expressed as:

$$\tilde{\boldsymbol{m}}_i^{\mathrm{cln}} = f(\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}, \tilde{L}_i), \ i = 1, 2, \cdots, 8, \tag{9}$$

where $L_i$ and $\tilde{L}_i$ are uniformly distributed random variables representing the original and modified numbers of lanes, respectively. Furthermore, $\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}$ stands for an individual movement vector within the junction matrix after the "movement shuffling" method whereas $f(\cdot, \cdot)$ denotes the function that adjusts the traffic characteristics. Specifically, $f(\cdot, \cdot)$ can take the following form:

$$f(\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}, L') = \begin{cases} [\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}]_k \times \left( \dfrac{\tilde{L}_i}{L_i} \right), & \text{if } k = 1, 2, 3, 5 \\ [\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}]_k, & \text{otherwise} \end{cases} \tag{10}$$

where $[\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}]_k$ represents the $k$-th entry of $\tilde{\boldsymbol{m}}_i^{\mathrm{ms}}$. Note that $f$ is designed to ensure that the traffic characteristics $F$, $O_{max}$, and $O_{mean}$ maintain relative values regardless of the variations in the lane configuration. After applying the "change of lane numbers" method in all junction matrices in $\tilde{\boldsymbol{S}}_t^{\mathrm{ms}}$, the new state can be represented as $\tilde{\boldsymbol{S}}_t^{\mathrm{cln}}$:

$$\tilde{\boldsymbol{S}}_t^{\mathrm{cln}} = \left[ \tilde{\boldsymbol{J}}_{t-K+1}^{\mathrm{cln}}, \tilde{\boldsymbol{J}}_{t-K+2}^{\mathrm{cln}}, \cdots, \tilde{\boldsymbol{J}}_t^{\mathrm{cln}} \right]. \tag{11}$$

where

$$\tilde{\boldsymbol{J}}_t^{\mathrm{cln}} = \left[ \tilde{\boldsymbol{m}}_1^{\mathrm{cln}}, \tilde{\boldsymbol{m}}_2^{\mathrm{cln}}, \cdots, \tilde{\boldsymbol{m}}_8^{\mathrm{cln}} \right]^T. \tag{12}$$

**Traffic Flow Scaling**: To shift the agent's focus from absolute car numbers to relative car distributions, we introduce a flow scaling factor. By multiplying the flow and occupancy values in the junction matrix with a uniformly distributed random number $\alpha$, we can create variations in the relative traffic volume for each movement. Notably, $\alpha$ remains consistent across all movements within the same traffic state, ensuring that the relative vehicle proportions between movements remain unchanged. This augmentation method facilitates the agent to prioritize the relative significance of different movements, thereby reducing its reliance on absolute values. It is worth noting that traffic flow scaling does not alter the number of lanes in each movement.

Fig. 4 illustrates two plausible scaling methods by proportionally increasing or decreasing the number of vehicles on each approach or adding small changes to the traffic volume in two scenarios. Mathematically, the flow scaling operation can be defined as:

$$\tilde{\boldsymbol{m}}_i^{\mathrm{tfs}} = g(\tilde{\boldsymbol{m}}_i^{\mathrm{cln}}, \alpha), \ i = 1, 2, \cdots, 8, \tag{13}$$

where $\alpha$ is the flow scaling factor. Furthermore, $g(\cdot, \cdot)$ denotes a function that scales the flow and occupancy values and takes the following form:

$$g(\tilde{\boldsymbol{m}}_i^{\mathrm{cln}}, \alpha) = \begin{cases} [\tilde{\boldsymbol{m}}_i^{\mathrm{cln}}]_k \times \alpha, & \text{if } k = 1, 2, 3 \\ [\tilde{\boldsymbol{m}}_i^{\mathrm{cln}}]_k. & \text{otherwise} \end{cases} \tag{14}$$

After replacing $\tilde{\boldsymbol{m}}_1^{\mathrm{cln}}$ in Eq. (12) with $\tilde{\boldsymbol{m}}_i^{\mathrm{tfs}}$, the resulting state from the "traffic flow scaling" method is denoted as $\tilde{\boldsymbol{S}}_t^{\mathrm{tfs}}$.

**Gaussian Noise Addition**: Gaussian noise is added directly to the junction matrix to introduce randomness into the training data. The additive noise can affect all components of the junction matrix, including traffic characteristics,

7

movement characteristics, and traffic signal characteristics. This augmentation method allows the agent to adapt to noisy and uncertain traffic conditions, improving its robustness during inference. Mathematically, the Gaussian noise addition operation can be modeled as:

$$\tilde{\boldsymbol{J}}_t^{\text{gna}} = \tilde{\boldsymbol{J}}_t^{\text{tfs}} + \boldsymbol{\Psi}_t, \tag{15}$$

where $\boldsymbol{\Psi}_t \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ denotes the Gaussian noise matrix, and $\tilde{\boldsymbol{J}}_t^{\text{tfs}} \in \tilde{\boldsymbol{S}}_t^{\text{tfs}}$ corresponds to the junction matrix after applying the "traffic flow scaling" method. After applying "Gaussian noise addition", the new state can be represented as $\tilde{\boldsymbol{S}}_t^{\text{gna}}$.

**Masking**: To encourage the agent to learn the traffic flow changes, we randomly set the values in the junction matrix to zero at a specific time instance. By masking certain components of the junction matrix, we create situations where the agent must rely on the information before and after the masked period to infer traffic dynamics. This augmentation method promotes the agent's capability of understanding and responding to traffic fluctuations.

Finally, a novel state denoted as $\tilde{\boldsymbol{S}}_t$ is obtained after the application of the five data augmentation methods discussed above.
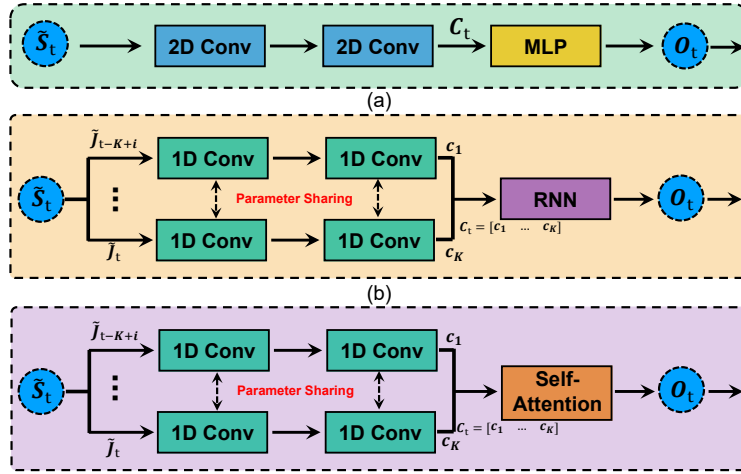


Figure 6: Three kinds of intersection feature extraction block (a) CNN-based Structure; (b) RNN-based Structure; (c) Transformer-based Structure.

## 4.4 Intersection Feature Extraction

In this section, three neural network structures are utilized to extract intersection information from the augmented traffic states, namely Convolutional Neural Network (CNN), Recurrent Neural Network (RNN) and transformers.

**CNN-based Structure**: As shown in Fig. 6(a), a CNN-based structure equipped with two 2D convolutional layers is utilized to extract time series information within a road junction, resulting in a hidden representation $\boldsymbol{C}_t$ given as:

$$\boldsymbol{C}_t = \text{ReLU}\left(\mathbf{W}_2^{2d}\mathbf{W}_1^{2d}\tilde{\boldsymbol{S}}_t\right), \tag{16}$$

where $\mathbf{W}_1^{2d}$ and $\mathbf{W}_2^{2d}$ represent the learnable parameters of the two 2D convolutional blocks, respectively. More specifically, the first convolutional block, $\mathbf{W}_1^{2d}$, extracts information regarding the traffic movement, and subsequently, the second convolutional block, $\mathbf{W}_2^{2d}$, captures information specific to the junction based on the movement information. Furthermore, ReLU$(\cdot)$ denotes the ReLU function.

Next, the resulting $\boldsymbol{C}_t$ is passed through a multilayer perceptron (MLP) layers, producing a feature vector $\boldsymbol{O}_t$ as follows:

$$\boldsymbol{O}_t = \mathbf{W}_u\boldsymbol{C}_t + \mathbf{b}_u, \tag{17}$$

where $\mathbf{W}_u$ and $\mathbf{b}_u$ are learnable parameters of the MLP layer.

**RNN-based Structure**: In contrast to the CNN-based approach, the RNN-based structure illustrated in Fig. 6(b) employs a parameter-sharing 1D convolutional layer to extract information from each junction matrix. The 1D

convolutional layer operates on each junction matrix $\tilde{\boldsymbol{J}}_{t-K+i}$ for $i = 1, 2, \cdots, K$, which is the augmented state at a particular time step within the history window. The resulting outputs denoted by $\boldsymbol{c}_1, \boldsymbol{c}_2, \ldots, \boldsymbol{c}_k$, capture the information from the intersection at $K$ different time instances and can be expressed as:

$$\boldsymbol{c}_i = \text{ReLU}\left(\mathbf{W}_2^{1d}\mathbf{W}_1^{1d}\tilde{\boldsymbol{J}}_{t-K+i}\right), \; i = 1, 2, \cdots, K, \tag{18}$$

where $\mathbf{W}_1^{1d}$ and $\mathbf{W}_2^{1d}$ represent the learnable parameters of the two 1D convolutional layers. These outputs are then fed into the RNN module:

$$\boldsymbol{h}_i = \tanh(\mathbf{W}_x\boldsymbol{c}_i + \boldsymbol{h}_{i-1}\mathbf{W}_h + \mathbf{b}_h), \; i = 1, 2, \cdots, K, \tag{19}$$

where $\mathbf{W}_x$, $\mathbf{W}_h$, and $\mathbf{b}_h$ represent the weights for the hidden layer in the RNN structure. The final hidden state $\boldsymbol{h}_K$ derived from the last RNN output is used to calculate the features of the entire intersection over a period of time, denoted as $\boldsymbol{O}_t$:

$$\boldsymbol{O}_t = \mathbf{W}_v\boldsymbol{h}_K + \mathbf{b}_v, \tag{20}$$

where $\mathbf{W}_v$ and $\mathbf{b}_v$ are the weights for the output layer in the RNN structure.

**Transformer-based Structure**: Fig. 6(c) shows the transformer-based approach. we employ a weight-shared CNN network to extract features from the junction matrix at each time step, as mentioned in Eq. (18). However, instead of using an RNN block, we utilize a transformer encoder to capture temporal dependencies in the sequence of features $\tilde{\boldsymbol{C}}_t = [\boldsymbol{c}_1, \boldsymbol{c}_2, \ldots, \boldsymbol{c}_K]$. To incorporate timing information, we insert a learnable embedding denoted as $\boldsymbol{c}_{\text{class}}$ to $\tilde{\boldsymbol{C}}_t$ as follows:

$$\boldsymbol{C}_t = \left[\boldsymbol{c}_{\text{class}}, \tilde{\boldsymbol{C}}_t\right] = [\boldsymbol{c}_{\text{class}}, \boldsymbol{c}_1, \boldsymbol{c}_2, \ldots, \boldsymbol{c}_K]. \tag{21}$$

The output state of the transformer encoder, based on this modified input sequence, serves as the traffic state representation $\boldsymbol{O}_t$. In the transformer encoder block, self-attention is calculated as follows:

$$\boldsymbol{Q}_C = \boldsymbol{C}_t\boldsymbol{W}_Q, \boldsymbol{K}_C = \boldsymbol{C}_t\boldsymbol{W}_K, \boldsymbol{V}_C = \boldsymbol{C}_t\boldsymbol{W}_V, \tag{22}$$

and

$$\boldsymbol{Z}_t = \phi\left(\frac{\boldsymbol{Q}_C\boldsymbol{K}_C^T}{\sqrt{d}}\right)\boldsymbol{V}_C, \tag{23}$$

where $\boldsymbol{Q}_C$, $\boldsymbol{K}_C$, and $\boldsymbol{V}_V$ represent the projected query, key and value features, respectively, while $\boldsymbol{W}_Q$, $\boldsymbol{W}_K$, and $\boldsymbol{W}_V$ the corresponding parameter metrics. In addition, $\boldsymbol{Z}_t$ represents the output of self-attention, and $\phi(\cdot)$ denotes the softmax function. The final output $\boldsymbol{O}_t$ is the first value of $\boldsymbol{Z}_t$.

### 4.5 RL Training and Fine-tuning

The Proximal Policy Optimization (PPO) algorithm [41] is adopted to train our UniTSA model. As depicted in Fig. 2, the agent gathers trajectories consisting of observations, actions, and rewards during the interactions with diverse traffic scenarios of different structures. These trajectories serve as the basis for computing the policy loss and value loss utilized to update the weights of the Actor and Critic networks.

The policy loss quantifies the difference between the current policy and the updated policy derived from the collected trajectories. It encourages the agent to increase the probabilities of actions that lead to higher rewards while reducing the probabilities of actions that lead to lower rewards. Mathematically, the policy loss can be formulated as:

$$\mathcal{L}_{\text{pf}}(\theta) = \hat{\mathbb{E}}_t\left[\min\left(r_tA_t, \text{clip}\left(r_t, 1 - \epsilon, 1 + \epsilon\right)A_t\right)\right], \tag{24}$$

where $\hat{\mathbb{E}}_t$ stands for the expectation operator and $r_t$ represents the ratio between the new policy $\pi_\theta(a_t|s_t)$ and the old policy $\pi_{\theta_{\text{old}}}(a_t|s_t)$:

$$r_t = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}. \tag{25}$$

Furthermore, $A_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ denotes the advantage function whereas the clip function can ensure stable policy updates.

The value loss measures the discrepancy between the estimated value function and the actual rewards obtained during the interaction. It makes the value function to better approximate the expected cumulative rewards. The value loss can be defined as:

$$\mathcal{L}_{\mathrm{vf}}(\theta) = \hat{\mathbb{E}}_t \left[ \left( V_\theta(s_t) - \hat{R}_t \right)^2 \right], \tag{26}$$

where $V_\theta(s_t)$ is the estimated value function under policy $\theta$, and $\hat{R}_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ denotes the reward-to-go. Once the policy loss and the value loss have been calculated, the final objective function is expressed as Eq. (27):

$$\mathcal{L}(\theta) = -\mathcal{L}_{\mathrm{pf}}(\theta) + \lambda \mathcal{L}_{\mathrm{vf}}(\theta), \tag{27}$$

where $\lambda$ is the coefficient of value loss.

This work proposes an effective universal model to control the traffic signals by optimizing this objective function through PPO. Furthermore, since certain intersections are more critical than others in practice, LoRA [31] is adopted in the proposed model to fine-tune the performance on these important intersections. More specifically, the LoRA modules are added to the weights of dense layers in the Actor and Critic networks as shown in Fig. 2. During the fine-tuning process, the original pretrained weights are kept constant while the LoRA modules is being updated. This design allows the model to adapt to the intersection-specific features without significantly increasing the number of parameters, while maintaining training efficiency.
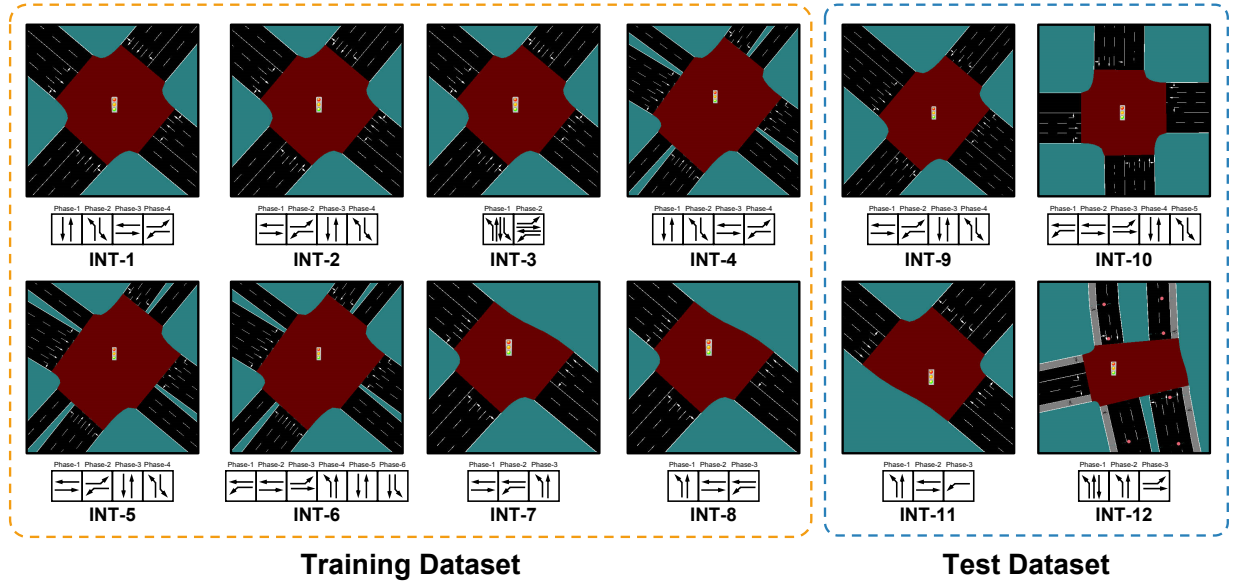


Figure 7: All intersection topologies with their available phases which used for the training and testing.

Let us consider a pretrained weight matrix $\boldsymbol{W} \in \mathbb{R}^{n \times m}$ in the network accompanied by a LoRA module $\Delta \boldsymbol{W} = \boldsymbol{W}_A \boldsymbol{W}_B^T$, where $\boldsymbol{W}_A \in \mathbb{R}^{n \times d}$, $\boldsymbol{W}_B \in \mathbb{R}^{m \times d}$ with $d \ll n$. The output of this layer can be obtained as:

$$z = \boldsymbol{W}x + \Delta \boldsymbol{W}x = \boldsymbol{W}x + \frac{\alpha}{r} \boldsymbol{W}_A \boldsymbol{W}_B^T x, \tag{28}$$

where $\boldsymbol{W}_A$ and $\boldsymbol{W}_B$ are initialized as a zero matrix and a zero-mean Gaussian distribution matrix, respectively. Furthermore, $\alpha$ is a constant scale hyperparameter whereas $r$ is the rank of the LoRA module. Through the combination of RL training using PPO and fine-tuning with LoRA, the proposed universal model can effectively address the challenges posed by important intersections of varying structures.

Table 1: Hyperparameter Setting

| Hyper-parameter | Value |
|---|---|
| Learning rate | 0.0001 |
| Trajectory memory size | 3000 |
| Clipping range $\epsilon$ | 0.2 |
| Discount factor $\gamma$ | 0.99 |
| Value function coefficient $\lambda$ | 0.9 |
| Scale hyperparameter $\alpha$ | 1 |
| Rank of LoRA module | 8 |

## 5 Experiments

### 5.1 Experiment Settings

Extensive experiments were conducted to validate the proposed model using the SUMO software package [32] in this section. SUMO is an open-source microscopic traffic simulation tool designed for handling large networks. It provides the Traffic Control Interface (TraCI) for controlling traffic lights and retrieving traffic condition information for intersections. We calculate the flow and occupancy of each movement by analyzing the positions and trajectories of vehicles on the road. It is important to note that, in order to simulate real-world conditions, we consider only vehicles within proximity of $150$ m to the intersection, in lieu of all vehicles along the entire road. Moreover, a green light was followed by a yellow light of $3$ s before transitioning to a red light to ensure driver safety. The waiting time per vehicle was used as a performance metric to evaluate the effectiveness of the different methods. A low waiting time indicates that vehicles spent less time passing through the intersection.

We utilized the Proximal Policy Optimization (PPO) implementation provided by the Stable Baselines3 library [42]. To accelerate training, we employed 30 parallel processes, and the total number of training environment steps was set to 10M. The state representation included the previous $K = 8$ snapshots of the junction matrix. The interval between two consecutive actions was $5$ s.

The hyper-parameters are configured as shown in Table 1. Furthermore, the actor and critic networks were designed as two-layer fully connected networks. The input sizes were $\{64, 32\}$, and the output sizes were $\{32, 2\}$ and $\{32, 1\}$ respectively for the actor and critic networks.

Table 2: All intersection configurations

| Intersection ID | Training Dataset | | | | | | | | Test Dataset | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | INT-1 | INT-2 | INT-3 | INT-4 | INT-5 | INT-6 | INT-7 | INT-8 | INT-9 | INT-10 | INT-11 | INT-12 |
| roads | 4 | 4 | 4 | 4 | 4 | 4 | 3 | 3 | 4 | 4 | 3 | 3 |
| lanes per road | (3,3,3,3) | (3,3,3,3) | (3,3,3,3) | (3,4,4,5) | (3,4,4,5) | (3,4,4,5) | (3,3,3) | (3,3,3) | (3,4,3,4) | (3,3,3,3) | (4,3,3) | (2,3,2) |
| phases | 4 | 4 | 2 | 4 | 4 | 6 | 3 | 3 | 4 | 5 | 3 | 3 |

### 5.2 Datasets

This study considers intersections with diverse structures, aiming to employ one single universal model to predict actions for all intersections. Specifically, 12 intersections of varying numbers of phases, lanes on each road, and approaching roads (i.e., 3-way or 4-way intersections) are constructed and used for experiments. Among these 12 intersections, eight are used for training, while the remaining four are reserved for testing. The topologies and phases of the 12 intersections are depicted in Fig. 7, and Table 2 provides a summary on their configurations. For instance, INT-4 shown in Fig. 7 consists of four bi-directional approaching roads with three lanes in the north-south direction, five lanes in the west-east direction, and four lanes in each of the other two directions. INT-4 includes four phases, each combining two different movement signals. Consequently, the configurations for INT-4 in Table 2 indicate the presence of four roads, lanes per road specified as $(3, 4, 4, 5)$ in clockwise order, and four phases.

There are three different intersection topologies in the training dataset. INT-1, INT-2 and INT-3 represent a regular 4-way intersection scenario, with each road consisting of three lanes; INT-4, INT-5 and INT-6 represent a large 4-way intersection scenario, featuring more than four lanes per road; INT-7, INT-8 and INT-9 depict the 3-way intersection scenario. For the intersections with the same topology, we generate new intersections by altering the sequence of phases or the number of phases. For instance, INT-1 and INT-2 have identical configurations, but the sequence of phases differs. Similarly, INT-1 and INT-3 vary in terms of the number of phases.

To assess the performance of the proposed model on unseen intersections, four testing scenarios are formed, namely INT-9, INT-10, INT-11 and INT-12. For instance, INT-9 and INT-10 also represent 4-way intersections, but INT-9 features different lane configurations compared to the training set, while INT-10 differs in the number and combination of phases from the intersections in the training set. INT-11 modifies the lane and phase of each road based on INT-8. Finally, INT-12 simulates real-world traffic within the city of Ingolstadt, Germany [43].

In addition to considering intersections with different structures, we generate 100 unique pieces of the route for each intersection. Three-quarters of these routes were utilized for training, while the remaining quarter was reserved for evaluation. Each route has a duration of $30,000$ seconds, equivalent to approximately 8 hours.

## 5.3 Compared Methods

To evaluate the performance of the proposed UniTSA, we compare the resulting universal model with several classic and state-of-the-art RL-based methods for TSC.

**FixTime** [2]: Fixed-time control utilizes a predetermined cycle and phase duration plan, which is widely used in situations with steady traffic flow. We consider two versions of FixTime, FixTime-30 (Fix-30) and FixTime-40 (Fix-40). These variants correspond to fixed-time signal control plans where each phase has a duration of 30 seconds and 40 seconds, respectively.

**Webster** [3]: The Webster method determines the cycle length and phase split based on traffic volume during a specific period. It has been proven that when the traffic is uniform, the Webster method minimizes the travel time of all vehicles passing the intersection or maximizes the intersection capacity [36]. Additionally, the Webster method can be adapted for real-time applications. For fairness, we employ the Webster method to adjust the traffic lights based on real-time traffic in this study.

**SOTL** [4]: Self-Organizing Traffic Light Control (SOTL) is an actuated signal control method that dynamically adjusts signal durations based on a manually determined threshold for the number of waiting vehicles. In this experiment, we set the threshold based on [44].

**MPLight** [10]: MPLight incorporates the FRAP structure proposed in [9] and a pressure-based reward mechanism inspired by [45]. Furthermore, MPLight employs a sharing-parameter multilayer perceptron (MLP) to enhance adaptability across different intersection configurations.

**AttendLight** [18]: This method adopts the attention mechanism to train a universal model capable of handling intersections with diverse structures and traffic flow distributions. It employs two attention models: the first attention model addresses variations in the number of roads and lanes, while the second attention model enables decision-making across intersections with different numbers of phases.

In addition to the baseline methods, we also consider several variations of our method:

**UniTSA (Single)**: This method focuses on training the model within a single environment. It utilizes the agent design described in Section 4.2 and employs an RNN-based structure for extracting information from the traffic state.

**UniTSA (Multi)**: In contrast to UniTSA (Single), this method trains the model across multiple environments simultaneously. We explore the performance of UniTSA (Multi) using various neural network designs, as discussed in Section 4.4, including UniTSA (Multi+CNN), UniTSA (Multi+RNN), and UniTSA (Multi+Trans).

**UniTSA (Multi+TSA)**: This variant of UniTSA enhances the training process by incorporating five traffic state augmentation methods into UniTSA (Multi). It results in UniTSA (Multi+CNN+TSA), UniTSA (Multi+RNN+TSA), and UniTSA (Multi+Trans+TSA).

## 5.4 Results of the training intersections

In this section, the performance of the proposed UniTSA is compared against that derived from several existing approaches, including the FixTime approach, Webster model, SOTL, MPLight, and AttendLight, on the training intersections. In particular, UniTSA models of different network structures were examined.

Table 3 shows the average waiting time per vehicle achieved by UniTSA and other baseline methods on the training intersections. It is clear that the RL-based approaches demonstrated superior performance as compared to the conventional approaches in most intersection scenarios. Despite that SOTL achieved the shortest waiting time in the INT-2, it required manually defined thresholds for different environments, limiting its generalization in large-scale scenarios.

Among the RL-based universal methods, namely MPLight, AttendLigh and UniTSA, UniTSA demonstrated significant performance improvements as compared to other RL-based methods. On average, UniTSA achieved 15% and 12%

Table 3: Quantitative results (average waiting time per vehicle) of training intersections for universal models. A lower value indicates better performance and the lowest values are highlighted in bold.

|  | INT-1 | INT-2 | INT-3 | INT-4 | INT-5 | INT-6 | INT-7 | INT-8 |
|---|---|---|---|---|---|---|---|---|
| Fix-30 [2] | 39.458 | 38.862 | 8.323 | 35.123 | 35.031 | 51.923 | 18.920 | 18.595 |
| Fix-40 [2] | 50.696 | 52.108 | 10.667 | 44.888 | 45.540 | 61.743 | 23.411 | 24.759 |
| Webster [3] | 26.466 | 26.889 | 5.751 | 25.413 | 24.541 | 40.128 | 12.755 | 13.574 |
| SOTL [4] | 16.048 | 16.561 | **2.764** | 28.169 | 27.902 | 27.404 | 8.639 | 7.925 |
| MPLight [10] | 19.111 | 14.659 | 4.469 | 16.067 | 19.925 | 19.115 | 6.654 | 7.523 |
| AttendLight [18] | 16.483 | 13.893 | 3.860 | 16.903 | 18.915 | 20.795 | 6.532 | 8.104 |
| UniTSA (Multi+CNN) | 13.776 | 13.580 | 3.265 | 14.790 | 15.437 | 18.751 | 6.592 | 6.894 |
| UniTSA (Multi+RNN) | 13.692 | 13.437 | 3.495 | 15.198 | 14.896 | 18.612 | 6.393 | 6.625 |
| UniTSA (Multi+Trans) | 14.976 | 20.459 | 2.811 | 16.489 | 16.481 | 23.793 | 7.552 | 7.175 |
| UniTSA (Multi+CNN+TSA) | 14.071 | 14.150 | 3.036 | 15.990 | 16.472 | 24.383 | 6.329 | 6.328 |
| UniTSA (Multi+RNN+TSA) | 13.450 | 13.578 | 3.007 | **14.470** | **14.462** | 18.689 | **6.242** | **6.164** |
| UniTSA (Multi+Trans+TSA) | **13.335** | **13.314** | 3.311 | 14.648 | 14.566 | **18.579** | 6.643 | 6.571 |

performance improvement over MPLight and AttendLight, respectively, across the eight intersections evaluated. When compared to MPLight, the proposed method incorporates not only parameter sharing techniques before replacing the MLP with RNN or Transformer block to capture the temporal information of the traffic state. In contrast to AttendLight, UniTSA simplifies the state design by representing the traffic state at a specific time instance using the junction matrix. Furthermore, five methods of traffic state enhancement are introduced, allowing the agent to observe a wider variety of traffic intersection states during training, which contributes to the improved performance of our model.

Finally, we explored the impact of different network structures and traffic state augmentation. Inspection of Table 3 reveals that an RNN-based structure yielded superior results compared to a CNN-based structure. Furthermore, in the *absence* of traffic state augmentation, UniTSA (Multi+CNN) and UniTSA (Multi+RNN) outperformed UniTSA (Multi+Trans) across most intersections, which is rather surprising. This can be attributed to the fact that the Transformer-based approach necessitates a larger volume of training data. However, in the *presence* of traffic state augmentation, the model can interact with a broader range of intersections with varying structures. As a result, UniTSA (Multi+Trans+TSA) outperformed UniTSA (Multi+RNN) and UniTSA (Multi+RNN+TSA) in many scenarios.

Table 4: Quantitative results of test intersections for universal models.

|  | INT-9 | INT-10 | INT-11 | INT-12 |
|---|---|---|---|---|
| Fix-30 | 40.341 | 38.360 | 17.136 | 17.440 |
| Fix-40 | 60.043 | 50.580 | 23.504 | 20.570 |
| Webster | 26.191 | 27.766 | 11.981 | 13.280 |
| SOTL | 23.031 | 28.066 | 7.452 | 8.070 |
| MPLight | 23.674 | 21.475 | 8.447 | 15.047 |
| AttendLight | 18.080 | 18.501 | 7.393 | 12.982 |
| UniTSA (Multi+CNN) | 17.877 | 18.244 | 7.063 | 12.820 |
| UniTSA (Multi+RNN) | 16.771 | 15.450 | 6.807 | 10.140 |
| UniTSA (Multi+Trans) | 17.640 | 16.269 | 6.598 | 9.630 |
| UniTSA (Multi+CNN+TSA) | 13.075 | 14.245 | 6.794 | 7.550 |
| UniTSA (Multi+RNN+TSA) | **12.677** | **14.208** | **5.914** | 6.730 |
| UniTSA (Multi+Trans+TSA) | 13.054 | 16.186 | 5.983 | **6.190** |

## 5.5   Results of the test intersections

Next, we evaluate the key feature of UniTSA to see whether it can be utilized for *unseen* intersections. Four intersections, namely INT-9, INT-10, INT-11 and INT-12, are specifically used for the testing purposes. As discussed in Section 5.2, these intersections differ from the scenarios in the training set in terms of the number of lanes or the number of phases. Table 4 summarizes the performance achieved by different methods, including baseline approaches and various variants of our proposed UniTSA method. Consistent with the results obtained on the training set, the RL-based algorithms significantly outperformed the traditional traffic control algorithms.

Among the RL-based universal methods, UniTSA models demonstrated superior performance across all test intersections as shown in Table 4. Among all the methods, UniTSA (Multi+RNN+TSA) and UniTSA (Multi+Trans+TSA) excelled
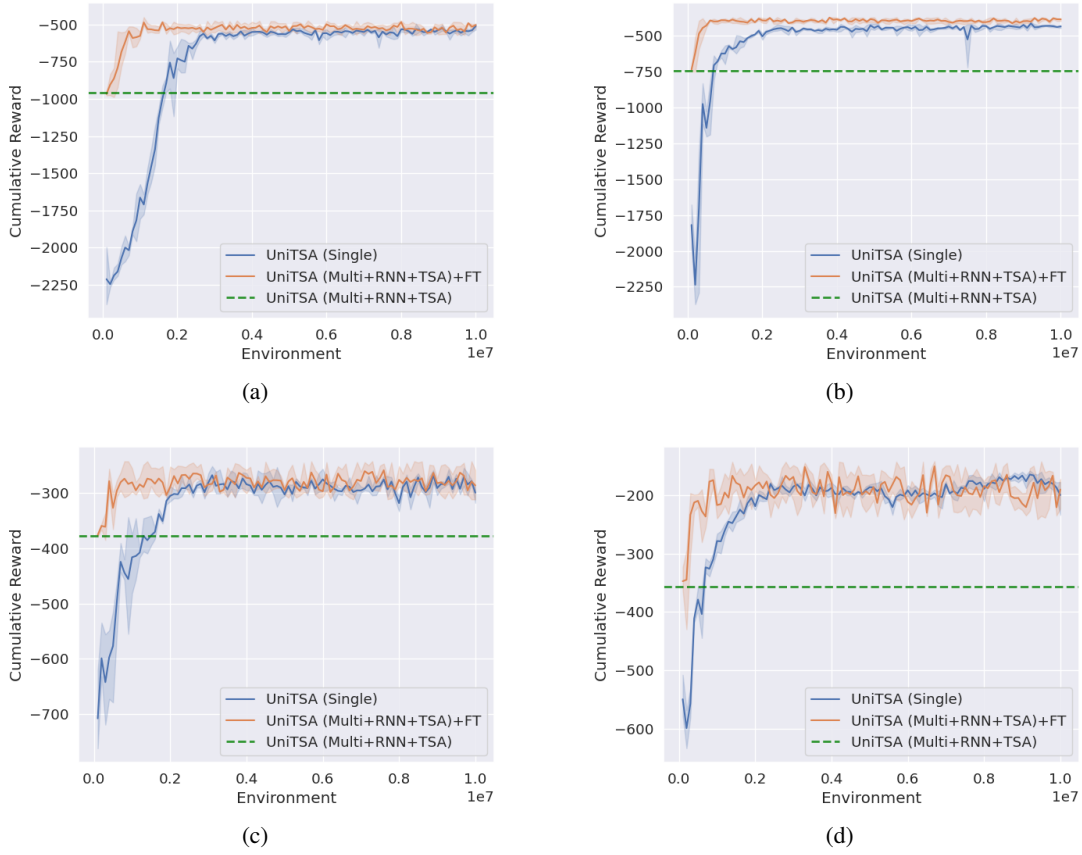
Figure 8: The environment steps of different methods at the four test intersections. (a) INT-9. (b) INT-10. (c) INT-11. (d) INT-12.

in terms of reducing average travel time for vehicles passing through the intersections. UniTSA (Multi+RNN+TSA) showcases an average travel time reduction of approximately $32.9\%$ and $41.3\%$ as compared to MPLight and Attend-Light, respectively. These results confirm the effectiveness of our approach in optimizing TSC and enhancing traffic flow efficiency.

Notably, incorporating traffic state augmentation techniques leads to improved performance among the different UniTSA variants. For instance, UniTSA (Multi+RNN+TSA), UniTSA (Multi+RNN+TSA), and UniTSA (Multi+RNN+TSA) exhibit enhancements of $23.4\%$, $19.8\%$, and $17.9\%$, respectively, when compared to UniTSA (Multi+RNN), UniTSA (Multi+RNN), and UniTSA (Multi+RNN). This improvement can be attributed to the inclusion of a greater variety of intersection scenarios within the training data through traffic state augmentation techniques, such as the "Lane number change" method, which enables the generation of diverse combinations of lane configurations.

Table 5: Quantitative results of fine-tuning in test intersections.

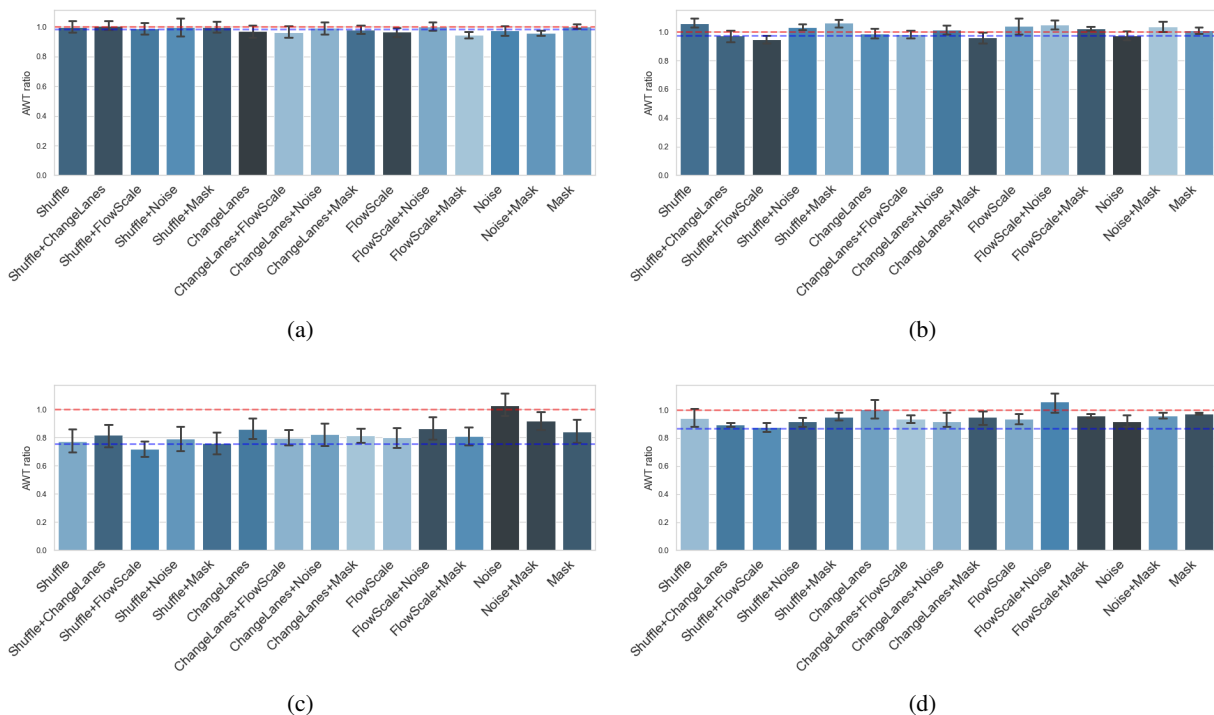|  | INT-9 | INT-10 | INT-11 | INT-12 |
|---|---|---|---|---|
| UniTSA (Multi+RNN) | 16.771 | 15.450 | 6.807 | 10.140 |
| UniTSA (Multi+RNN+TSA) | 12.677 | 14.208 | 5.914 | 6.730 |
| *1M Environment Steps* | | | | |
| UniTSA (Single) | 23.683 | 14.863 | 8.417 | 5.024 |
| UniTSA (Multi+RNN+TSA) + FT | **10.995** | **12.553** | **4.475** | **3.534** |
| *10M Environment Steps* | | | | |
| UniTSA (Single) | 10.353 | 12.254 | 4.359 | **3.089** |
| UniTSA (Multi+RNN+TSA) + FT | **10.292** | **11.081** | **4.153** | 3.110 |

14

Figure 9: Comparative analysis of traffic state augmentation methods on the selected training and test intersections (a) INT-1. (b) INT-7. (c) INT-9. (d) INT-11.

## 5.6 Results of fine-tuning in test intersections

In practical scenarios, certain intersections require special attention due to their significance. To address this, we begin with the universal model trained by UniTSA (Multi+RNN). The RNN-based UniTSA is chosen because it demonstrates superior performance across most intersections compared with both CNN-based and Transformer-based structures. In comparison to the UniTSA (Single), which is trained on a single scenario, the resulting model can quickly reach or even surpass the performance of the single-environment model after only a few training steps.

Fig. 8 shows the change of cumulative rewards over training steps for different models in the test intersections. The green dashed line represents the result of applying the universal model directly to new intersections without any fine-tuning. Notably, the model already exhibits promising results without any additional fine-tuning or transfer learning. The blue line represents the model trained from scratch whereas the orange line the fine-tuned model based on the universal model. It is observed that the single-environment model converges at approximately 3M training steps. In sharp contrast, the fine-tuning model achieves comparable performance with only around 1M training steps, resulting in approximately $66\%$ reduction in computation time while maintaining similar performance.

Table 5 provides a detailed analysis of the model's performance after fine-tuning. At 1M training steps, the fine-tuning model demonstrated an average performance improvement of $36\%$ as compared to UniTSA (Single) across the four test intersections. Even after 10M training steps, the fine-tuning models continued to outperform the models trained from scratch by $3\%$. This aspect is particularly appealing for real-time applications. For instance, in a road network with over 1000 junctions, it is possible to significantly reduce the number of interactions with the environment while maintaining comparable performance, thereby greatly enhancing training efficiency.

## 5.7 Comparative analysis on traffic state augmentations

In this section, we analyze the effectiveness of traffic state augmentation methods in improving the performance of the UniTSA (Multi+RNN) model. We conducted experiments by applying different combinations of two traffic state augmentation techniques on both the training and test intersections. To evaluate the impact of these methods, we propose the Average Waiting Time (AWT) ratio between the models with and without traffic state augmentations with an AWT ratio of less than 1 indicating an improvement achieved by the corresponding traffic state augmentation methods.

15

Fig. 9 showcases the outcomes of the comparative analysis on INT-1, INT-7, INT-9, and INT-11 intersections, which represent common intersection structures encountered in real-world scenarios. Each bar in the figure corresponds to the average AWT ratio for a specific combination of traffic state augmentations, and the error bars represent the $95\%$ confidence interval. The blue dashed line represents the AWT ratio of UniTSA (Multi+RNN+TSA), which employs all available traffic state augmentation methods.

Fig. 9a and Fig. 9b depict the results obtained from the training intersections INT-1 and INT-7, respectively. Inspection of these figures reveals that the average performance improvement achieved through traffic state augmentations was around $2\%$. Furthermore, the inclusion of noise and mask methods incurred performance degradation in INT-7, which can be attributed to the fact that the model has already captured the underlying patterns and characteristics of the training intersections to a large extent. As a result, introducing additional variations through traffic state augmentation may not provide substantial benefits in the training set. However, traffic state augmentation methods demonstrated significant improvements when confronted with unseen intersection structures.

Fig. 9c and Fig. 9d illustrate the AWT ratios for the test intersections INT-9 and INT-11, respectively. These results clearly demonstrate that most traffic state augmentation methods enhanced the performance of the base policy in the test intersections, which consists of unseen intersections. The diverse training samples generated through traffic state augmentation contribute to the improved performance. Among the traffic state augmentation techniques, movement shuffle and traffic flow scale were particularly effective in enhancing the model's performance. These techniques enable the model to adapt and learn from a wider range of scenarios, resulting in improved performance on the test set.

## 6 Conclusion

In this paper, a universal RL-based TSC framework called UniTSA has been proposed for diverse intersection structures in V2X environments. More specifically, UniTSA offers the capability to train a universal RL agent by incorporating a junction matrix to characterize intersection states. To handle unseen intersections, new traffic state augmentation methods have been proposed to enrich the agent's data collection, resulting in improved performance and generalization for unseen intersection configurations. As a result, UniTSA eliminates the necessity of extensive customization and redevelopment for each individual intersection while offering a simple, efficient, and open-sourced implementation, which makes UniTSA a valuable framework for future research in data-efficient and generalizable RL-based TSC methods within the field of V2X. Extensive experimental results have demonstrated that UniTSA achieved the shortest average waiting time across various intersection configurations, surpassing the performance of the existing methods and outperforming the models trained from scratch with fine-tuning.

## Acknowledgments

## References

[1] Fehda Malik, Hasan Ali Khattak, and Munam Ali Shah. Evaluation of the impact of traffic congestion based on sumo. In *2019 25th International Conference on Automation and Computing (ICAC)*, pages 1–5. IEEE, 2019.

[2] Alan J Miller. Settings for fixed-cycle traffic signals. *Journal of the Operational Research Society*, 14(4):373–386, 1963.

[3] Thomas Urbanik, Alison Tanaka, Bailey Lozner, Eric Lindstrom, Kevin Lee, Shaun Quayle, Scott Beaird, Shing Tsoi, Paul Ryus, Doug Gettman, et al. *Signal timing manual*, volume 1. Transportation Research Board Washington, DC, 2015.

[4] Carlos Gershenson. Self-organizing traffic lights. *arXiv preprint nlin/0411066*, 2004.

[5] Ishu Tomar, S Indu, and Neeta Pandey. Traffic signal control methods: Current status, challenges, and emerging trends. *Proceedings of Data Analytics and Management: ICDAM 2021, Volume 1*, pages 151–163, 2022.

[6] Wang Tong, Azhar Hussain, Wang Xi Bo, and Sabita Maharjan. Artificial intelligence for vehicle-to-everything: A survey. *IEEE Access*, 7:10823–10843, 2019.

[7] Tamás Wágner, Tamás Ormándi, Tamás Tettamanti, and István Varga. Spat/map v2x communication between traffic light and vehicles and a realization with digital twin. *Computers and Electrical Engineering*, 106:108560, 2023.

[8] Yit Kwong Chin, Lai Kuan Lee, Nurmin Bolong, Soo Siang Yang, and Kenneth Tze Kin Teo. Exploring q-learning optimization in traffic signal timing plan management. In *2011 third international conference on computational intelligence, communication systems and networks*, pages 269–274. IEEE, 2011.

[9] Guanjie Zheng, Yuanhao Xiong, Xinshi Zang, Jie Feng, Hua Wei, Huichu Zhang, Yong Li, Kai Xu, and Zhenhui Li. Learning phase competition for traffic signal control. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pages 1963–1972, 2019.

[10] Chacha Chen, Hua Wei, Nan Xu, Guanjie Zheng, Ming Yang, Yuanhao Xiong, Kai Xu, and Zhenhui Li. Toward a thousand lights: Decentralized deep reinforcement learning for large-scale traffic signal control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3414–3421, 2020.

[11] Xinshi Zang, Huaxiu Yao, Guanjie Zheng, Nan Xu, Kai Xu, and Zhenhui Li. Metalight: Value-based meta-reinforcement learning for traffic signal control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 1153–1160, 2020.

[12] Enming Liang, Zicheng Su, Chilin Fang, and Renxin Zhong. Oam: An option-action reinforcement learning framework for universal multi-intersection control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 4550–4558, 2022.

[13] Azzedine Boukerche, Dunhao Zhong, and Peng Sun. A novel reinforcement learning-based cooperative traffic signal system through max-pressure control. *IEEE Transactions on Vehicular Technology*, 71(2):1187–1198, 2022.

[14] Liang Zhang, Qiang Wu, Jun Shen, Linyuan Lü, Bo Du, and Jianqing Wu. Expression might be enough: representing pressure and demand for reinforcement learning based traffic signal control. In *International Conference on Machine Learning*, pages 26645–26654. PMLR, 2022.

[15] Yuanhao Xiong, Guanjie Zheng, Kai Xu, and Zhenhui Li. Learning traffic signal control from demonstrations. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pages 2289–2292, 2019.

[16] Stefano Giovanni Rizzo, Giovanna Vantini, and Sanjay Chawla. Time critic policy gradient methods for traffic signal control in complex and congested scenarios. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '19, page 1654–1664, New York, NY, USA, 2019. Association for Computing Machinery.

[17] Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. Multi-agent deep reinforcement learning for large-scale traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 21(3):1086–1095, 2019.

[18] Afshin Oroojlooy, Mohammadreza Nazari, Davood Hajinezhad, and Jorge Silva. Attendlight: Universal attention-based reinforcement learning model for traffic signal control. *Advances in Neural Information Processing Systems*, 33:4079–4090, 2020.

[19] Zian Ma, Chengcheng Xu, Yuheng Kan, Maonan Wang, and Wei Wu. Adaptive coordinated traffic control for arterial intersections based on reinforcement learning. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 2562–2567. IEEE, 2021.

[20] Seyed Sajad Mousavi, Michael Schukat, and Enda Howley. Traffic light control using deep policy-gradient and value-function-based reinforcement learning. *IET Intelligent Transport Systems*, 11(7):417–423, 2017.

[21] Mohammad Aslani, Mohammad Saadi Mesgari, Stefan Seipel, and Marco Wiering. Developing adaptive traffic signal control by actor–critic and direct exploration methods. In *Proceedings of the Institution of Civil Engineers-Transport*, volume 172, pages 289–298. Thomas Telford Ltd, 2019.

[22] Haoran Su, Yaofeng D Zhong, Joseph YJ Chow, Biswadip Dey, and Li Jin. Emvlight: A multi-agent reinforcement learning framework for an emergency vehicle decentralized routing and traffic signal control system. *Transportation Research Part C: Emerging Technologies*, 146:103955, 2023.

[23] Elise Van der Pol and Frans A Oliehoek. Coordinated deep reinforcement learners for traffic light control. *Proceedings of learning, inference and control of multi-agent systems (at NIPS 2016)*, 8:21–38, 2016.

[24] Patrick Mannion, Jim Duggan, and Enda Howley. An experimental review of reinforcement learning algorithms for adaptive traffic signal control. *Autonomic road transport support systems*, pages 47–66, 2016.

[25] Hua Wei, Guanjie Zheng, Huaxiu Yao, and Zhenhui Li. Intellilight: A reinforcement learning approach for intelligent traffic light control. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2496–2505, 2018.

[26] Lun-Hui Xu, Xin-Hai Xia, and Qiang Luo. The study of reinforcement learning for traffic self-adaptive control under multiagent markov game environment. *Mathematical Problems in Engineering*, 2013, 2013.

[27] Mohammad Aslani, Mohammad Saadi Mesgari, and Marco Wiering. Adaptive traffic signal control with actor-critic methods in a real-world traffic network with different traffic disruption events. *Transportation Research Part C: Emerging Technologies*, 85:732–752, 2017.

[28] Mohammad Aslani, Stefan Seipel, Mohammad Saadi Mesgari, and Marco Wiering. Traffic signal optimization through discrete and continuous reinforcement learning with robustness analysis in downtown tehran. *Advanced Engineering Informatics*, 38:639–655, 2018.

[29] Halit Bugra Tulay and Can Emre Koksal. Road state inference via channel state information. *IEEE Transactions on Vehicular Technology*, pages 1–14, 2023.

[30] Mohamed MG Farag, Hesham A Rakha, Emadeldin A Mazied, and Jayanthi Rao. Integration large-scale modeling framework of direct cellular vehicle-to-all (c-v2x) applications. *Sensors*, 21(6):2127, 2021.

[31] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

[32] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lücken, Johannes Rummel, Peter Wagner, and Evamarie Wießner. Microscopic traffic simulation using sumo. In *2018 21st international conference on intelligent transportation systems (ITSC)*, pages 2575–2582. IEEE, 2018.

[33] Peter Koonce and Lee Rodegerdts. Traffic signal timing manual. Technical report, United States. Federal Highway Administration, 2008.

[34] Arthur G Sims and Kenneth W Dobinson. The sydney coordinated adaptive traffic (scat) system philosophy and benefits. *IEEE Transactions on vehicular technology*, 29(2):130–137, 1980.

[35] Pravin Varaiya. The max-pressure controller for arbitrary networks of signalized intersections. *Advances in dynamic network modeling in complex transportation systems*, pages 27–66, 2013.

[36] Hua Wei, Guanjie Zheng, Vikash Gayah, and Zhenhui Li. A survey on traffic signal control methods. *arXiv preprint arXiv:1904.08117*, 2019.

[37] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169, 2021.

[38] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *Advances in neural information processing systems*, 33:19884–19895, 2020.

[39] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.

[40] Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data augmentation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13611–13617. IEEE, 2021.

[41] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[42] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *The Journal of Machine Learning Research*, 22(1):12348–12355, 2021.

[43] James Ault and Guni Sharon. Reinforcement learning benchmarks for traffic signal control. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.

[44] Seung-Bae Cools, Carlos Gershenson, and Bart D'Hooghe. Self-organizing traffic lights: A realistic simulation. *Advances in applied self-organizing systems*, pages 45–55, 2013.

[45] Hua Wei, Chacha Chen, Guanjie Zheng, Kan Wu, Vikash Gayah, Kai Xu, and Zhenhui Li. Presslight: Learning max pressure control to coordinate traffic signals in arterial network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1290–1298, 2019.