# Learning to Control and Coordinate Mixed Traffic Through Robot Vehicles at Complex and Unsignalized Intersections

**Dawei Wang**
The University of Hong Kong
dawei@connect.hku.hk

**Weizi Li**
University of Tennessee, Knoxville
weizili@utk.edu

**Lei Zhu**
University of North Carolina at Charlotte
lzhu14@charlotte.edu

**Jia Pan**
The University of Hong Kong
jpan@cs.hku.hk

## Abstract

Intersections are essential road infrastructures for traffic in modern metropolises. However, they can also be the bottleneck of traffic flows as a result of traffic incidents or the absence of traffic coordination mechanisms such as traffic lights. Recently, various control and coordination mechanisms that are beyond traditional control methods have been proposed to improve the efficiency of intersection traffic. Amongst these methods, the control of foreseeable mixed traffic that consists of human-driven vehicles (HVs) and robot vehicles (RVs) has emerged. In this project, we propose a decentralized multi-agent reinforcement learning approach for the control and coordination of mixed traffic at real-world, complex intersections–a topic that has not been previously explored. Comprehensive experiments are conducted to show the effectiveness of our approach. In particular, we show that using 5% RVs, we can prevent congestion formation inside a complex intersection under the actual traffic demand of 700 vehicles per hour. In contrast, without RVs, congestion starts to develop when the traffic demand reaches as low as 200 vehicles per hour. When there exist more than 60% RVs in traffic, our method starts to achieve comparable or even better performance to traffic signals on the average waiting time of all vehicles at the intersection. Our method is also robust against both blackout events and sudden RV percentage drops, and enjoys excellent generalizablility, which is illustrated by its successful deployment in two unseen intersections.

## 1 Introduction

Traffic flow is the beating heart of a city, driving economic growth and ensuring daily lives. Despite the implementation of various traffic control methods, including traffic signals, ramp meters, and tolls, traffic congestion continues to be a global issue, with external expenses amounting to over $100 billion annually [Schrank et al., 2021]. With urbanization and motorization trends set to persist in the coming decades [Nations, 2019, Trouve et al., 2020], the need for better traffic system design and management is urgent.

Traffic is an interplay between vehicles and infrastructure. Modern urban road networks largely consist of linearly-coupled roads interconnected by intersections. The key to this design's functionality is *intersection*, which enables traffic flows from different directions to interchange and disperse. Any intersection blockage can disrupt traffic from all directions, leading to traffic spillover and even city-wide gridlock in severe cases. Unfortunately, intersections are vulnerable to traffic incidents due

to their varied (and potentially complex) topologies and conflicting traffic streams. In the U.S., more than 45% of all crashes take place at intersections [Choi, 2010]. Furthermore, extreme weather and energy shortages can disable traffic signals, leaving intersections without control for days or even weeks, and causing traffic to become stranded and congested [Press, 2022, Winck, 2022, Ramirez, 2022]. This raises the question: *how can we ensure uninterrupted traffic flows at intersections?*

While current transport policies and control methods have limited effectiveness in mitigating traffic delays and congestion, connected and autonomous vehicles (CAVs) offer us new opportunities. Recent studies [Sharon and Stone, 2017, Yang and Oguchi, 2020] have demonstrated the possibilities of using self-driving robot vehicles to enhance traffic throughput at intersections. However, these studies presume universal connectivity and centralized control of all vehicles, a scenario that may not materialize soon. The transition to varying levels of autonomous vehicles will be *gradual*, with a prolonged period of mixed traffic that includes both human-driven vehicles (HVs) and robot vehicles (RVs). Despite the challenges in modeling and controlling mixed traffic due to the diversity and suboptimality of human drivers, it is possible to regulate it by algorithmically determining the behaviors of RVs and using them to influence nearby HVs [Wu et al., 2022]. While progress has been made (see Sec. 1.1 for details), to the best of our knowledge, no evidence exists to demonstrate the feasibility of controlling mixed traffic with RVs at *real-world, complex intersections* where a large number of vehicles may potentially conflict Nevertheless, real-world intersection is an essential step towards achieving city-wide traffic control and realizing the full potential of self-driving vehicles for society.

In this project, we study the control of mixed traffic at real-world intersections. The intersection layouts and reconstructed traffic are shown in Fig. 1 LEFT. To test the limit of mixed traffic control and explore the envisioned benefits of RVs to our traffic system [Urmson et al., 2008], we further assume these intersections are unsignalized, with the flow of traffic entirely controlled by RVs and no other means. This project presents many challenges:

- How to design a model-free approach to control and coordinate large-scale traffic at complex intersections, considering that no models capture the behaviors of mixed traffic;

- How to design a generic representation of traffic conditions to ensure generalizablility. This representation has to account for intersection topologies, conflicting traffic streams, and real-world fluctuating traffic demands;

- How to generate mixed traffic and construct high-fidelity simulations based on real-world traffic data so that the RVs have an environment to interact with and learn from;

- How to ensure safety, given that developing effective and provably safe autonomous systems remains an open challenge; and

- How to evaluate such a control approach in terms of its effectiveness, generalizablility, robustness, and the design components.
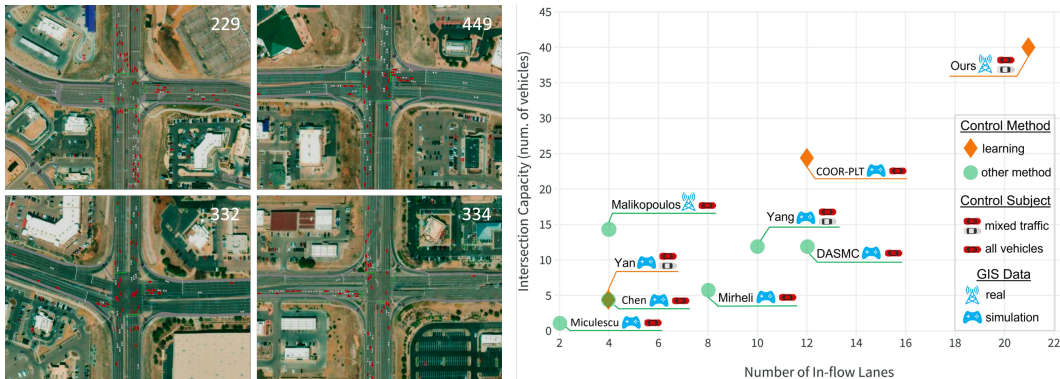


Figure 1: LEFT: We study real-world, complex intersections. These intersections are located at Colorado Springs, CO, USA. The traffic is reconstructed using the actual traffic data collected at these intersections. RIGHT: Comparison of example state-of-the-art studies on intersection traffic control.

We propose a model-free reinforcement learning (RL) method to mixed traffic control at complex intersections. Within the intersection, each RV will collect the status of surrounding vehicles and

encode them into a fixed-length representation for the mixed traffic control policy. Then, the policy outputs a high-level decision of whether the RV should enter or not enter the intersection.

Our approach falls into the paradigm of centralized training, decentralized execution. RVs are centrally trained with a reward function that accounts for traffic efficiency as well as the potential conflicts. During execution, all RVs make independent decisions while collectively ensuring smooth traffic flow at the intersection with virtually no centralized coordination. In addition, to guarantee safety, a fail-safe mechanism is implemented to eliminate potential conflicts within intersection areas which boosts the efficiency of the training process as well as ensures the safety of model inference.

We conduct comprehensive experiments, under high-fidelity traffic reconstruction and simulation. The real world traffic data provided by the city of Colorado Spring, CO, USA is used to reconstruct the simulated traffic flow, validating that our training environments and evaluation experiments closely resemble real-world conditions.

The results show that **with 60% or more RVs, our method outperforms traffic lights in traffic efficiency in most scenarios**. For example, the average waiting time of all vehicles is reduced by 30.9% and 42.7% compared to employing traffic lights at intersection 229, when the RV penetration rate is 70% and 90%, respectively. With 100% RVs, our method can reduce the average waiting time of the entire intersection traffic up to 81% compared to traffic light control and 97% compared to the traffic light absence baseline. We further explore the relationship between traffic demands, congestion, and RV penetration rates. Our findings show that with just 5% RVs, we can prevent congestion at the intersection under the actual traffic demand of 700 v/h. In contrast, without RVs, congestion emerges when the traffic demand is higher than 200 v/h.

We test the robustness and generalizablility of our approach. For robustness, we conduct a 'blackout' experiment when traffic lights suddenly stop working. During such an event, the RVs act as self-organized movable 'traffic lights' to coordinate the traffic and prevent congestion. Next, we examine the impact of sudden RV rate drops. The results demonstrate that even with 40% drop (from 90% to 50%), our method can still maintain stable and efficient traffic flows at the intersection. For generalizablility, we deploy our method in two unseen intersections without refinement: not only does our method prevent congestion, but also with at least 60% RVs, our method can surpass traffic light control on saving the average waiting time of all vehicles in one intersection. For another intersection, it is a tie for our method and traffic signal control.

As reward design is essential to an RL task, we justify our reward function by comparing it to the reward of the state-of-the-art method by Yan and Wu [Yan and Wu, 2021]. We show that our reward reflects intersection traffic conditions swiftly. In summary, our work is **the first to demonstrate the feasibility of controlling and coordinating mixed traffic at unsignalized intersections with complex topologies and real-world traffic demands**. As many challenges are addressed for the first time in mixed traffic control, we hope that our design can provide insights into these questions and stimulate future endeavors in the field. Our code is available at https://github.com/daweidavidwang/MixedTrafficControl

## 1.1 Related Work

Existing studies have demonstrated the potential of mixed traffic control in scenarios such as ring roads, figure-eight roads [Wu et al., 2022], highway bottleneck and merge [Vinitsky et al., 2018, Feng et al., 2021], two-way intersections [Yan and Wu, 2021], and roundabouts [Jang et al., 2019]. However, these scenarios typically lack real-world complexity and involve only a limited number of vehicles that may be in conflict.

Besides using RL, two approaches are developed to control and coordinate traffic at intersections [Rios-Torres and Malikopoulos, 2016]. The first approach is traffic signal control [Hunt et al., 1981, Hadi and Wallace, 1993, 1994], which has been extensively studied. However, our work differs from this line of research since we assume that intersections are unsignalized. The second is intersection management system [Miculescu and Karaman, 2019, Malikopoulos et al., 2018], which requires all vehicles to be centrally controlled and thus is not applicable to mixed traffic.

The comparison of our work and example studies of intersections is presented in Fig. 1 RIGHT. As these studies do not provide all measurements, we offer our best estimates. For complete information, we refer readers to COOR-PLT [Li et al., 2022], DASMC [Zhou et al., 2022], Yang [Yang and Oguchi,
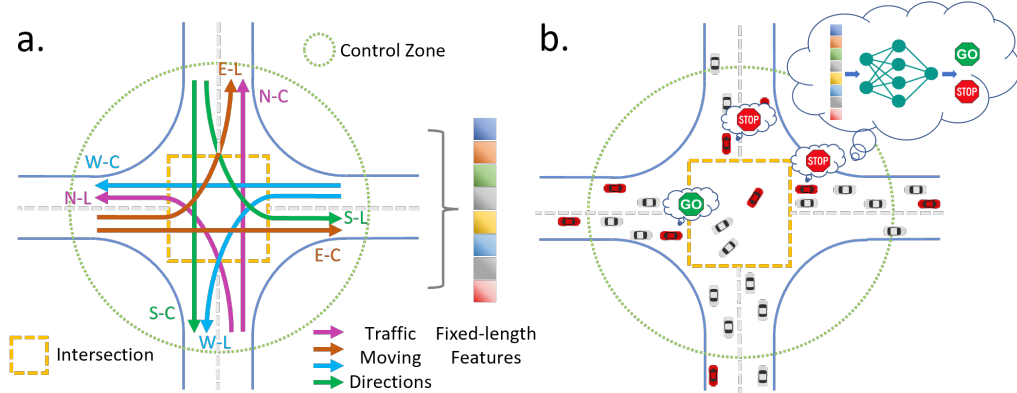
Figure 2: The pipeline of our approach. **a.** Each RV within the control zone encodes the intersection's traffic condition as a fixed-length representation. This includes both macroscopic traffic features such as queue length and waiting time, and microscopic traffic features such as vehicle locations in each moving direction (E, W, N, and S represent east, west, north, and south, respectively; C means cross and L means left-turn). **b.** The traffic-condition representation is then used by each RV to decide Stop or Go at the intersection entrance.

2020], Malikopoulos [Malikopoulos et al., 2018], Mirheli [Mirheli et al., 2019], Chen [Chen et al., 2022], Yan [Yan and Wu, 2021], and Miculescu [Miculescu and Karaman, 2019].

## 2 Methodology

The pipeline of our approach is shown in Fig. 2. Each RV entering the control zone employs our method and observes the traffic condition within. The RV subsequently encodes the traffic condition into a fixed-length representation (Fig. 2a), and then uses it to make a high-level decision (Stop or Go) at the intersection entrance (Fig. 2b).

### 2.1 Intersection Traffic

A standard four-way intersection comprises four moving directions: eastbound (E), westbound (W), northbound (N), and southbound (S); and three turning options: left (L), right (R), and cross (C). As an example, we use E-L and E-C to denote left-turning traffic and crossing traffic that travel eastbound, respectively. The complete notation is shown in Fig. 2a. We further define 'conflict' as two moving directions intersecting each other, e.g., E-C and N-C. In total, we consider eight traffic streams that may lead to conflicts: E-L, E-C, W-L, W-C, N-L, N-C, S-L, and S-C; and we define the conflict-free movement set $\mathcal{C}$ as (S-C, N-C), (W-C, E-C), (S-L, N-L), (E-L, W-L), (S-C, S-L), (E-C, E-L), (N-C, N-L), (W-C, W-L). For the pairs of traffic streams that are not in $\mathcal{C}$, conflicts may arise.

It is worth noting that in our scenarios, right-turning traffic is not considered since most intersections have an independent right-turn lane. As a result, right-turning traffic does not enter the intersection or only briefly occupies it. Additionally, in many countries such as the U.S., right-turn vehicles are typically not required to wait for the green light due to traffic regulations. Therefore, we do not coordinate right-turning traffic with traffic from other directions. Our experiments confirm that this empirical choice has minimal impact on controlling and coordinating intersection traffic.

To enable RVs to interact with HVs under real-world traffic conditions, we reconstruct traffic using actual traffic data and then conduct high-fidelity simulations. To create mixed traffic, we randomly assign newly spawned vehicles that enter the simulation to be either RV or HV based on a predetermined RV penetration rate. These procedures are detailed in Appendix A.1.

### 2.2 Decentralized RL For Mixed Traffic

We formulate mixed traffic control as POMDP, which consists of a 7-tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma)$, where $\mathcal{S}$ is a set of states ($\mathbf{s} \in \mathcal{S}$), $\mathcal{A}$ is a set of actions ($\mathbf{a} \in \mathcal{A}$), $\mathcal{T}$ is the transition probabilities between states $\mathcal{T}(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$, $\mathcal{R}$ is the reward function ($\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$), $\Omega$ is a set of observations $\mathbf{o} \in \Omega$, $\mathcal{O}$ is the set of conditional observation probabilities, and $\gamma \in [0, 1)$ is a discount factor. At each time $t$, when an RV $i$ enters the control zone, its action $a_i^t$ is determined based on the current traffic

4

condition $o_i^t$, which is a partial observation of the traffic state $s_i^t$ of the intersection. We present the policy $\pi_\theta$ as a neural network trained using the following loss:

$$\left(R_{t+1} + \gamma_{t+1}q_{\bar{\theta}}\left(S_{t+1}, \arg\max_{a'} q_\theta(S_{t+1}, a)\right) - q_\theta(S_t, A_t)\right)^2, \tag{1}$$

where $q$ denotes the estimated value from the value network, while $\theta$ and $\bar{\theta}$ respectively represent the value network and the target network [Hessel et al., 2018]. The target network is a periodic copy of the value network, which is not directly optimized during training.

### 2.2.1 Action Space

As our focus is on exploring mixed traffic control via RVs over traffic lights, we restrict the action space of RV to high-level decisions $A = \{\text{Stop}, \text{Go}\}$. An RV's action $a_i^t \in A$ determines whether the RV $i$ should enter the intersection or stop at the intersection entrance to hold its following vehicles.

The longitudinal acceleration of an RV is computed using Intelligent Driver Model (IDM) [Treiber et al., 2000] when the vehicle is outside the control zone. Within the control zone, if the RV decides Go, it accelerates using the maximum acceleration $a^t = a_{\max}$; conversely, if the RV decides Stop, it decelerates and comes to a halt at the intersection via $a^t = -v^2/2d_{\text{front}}$, where $d_{\text{front}}$ is the distance to the intersection. Note that alternative formulas can be used to calculate the acceleration.

### 2.2.2 Observation Space

To empower an RL policy capable of generalizing across diverse intersection topologies, we encode the traffic conditions observed by each RV into a fixed-length representation. The observation for each RV within the control zone (commencing from $30\ m$ before the intersection) encompasses three elements.

- The status of RV: The distance from RV $i$'s current position to the intersection, denoted as $d_i^t$.

- Traffic conditions within the control zone but outside the intersection: The queue length $l^{t,j}$ and the average waiting time $w^{t,j}$ of each of the eight traffic moving directions (introduced in Sec. 2.1). This is to quantify the anisotropic congestion levels of an intersection.

- Traffic condition inside the intersection: We design an occupancy map $m^{t,j}$ for each moving direction. As depicted in Fig. 3 LEFT, we divide the inner lane along a moving direction into 10 equally sized segments. A segment is labeled 'occupied' with a value of 1 if a vehicle's position is located within it, or labeled 'free' with a value of 0 if otherwise.

Overall, the observation space of RV $i$ at $t$ is

$$o_i^t = \oplus_j^J \langle l^{t,j}, w^{t,j}\rangle \oplus_j^J \langle m^{t,j}\rangle \oplus \langle d_i^t\rangle, \tag{2}$$

where $\oplus$ is the concatenation operator and $J = 8$ is the total number of traffic moving directions.
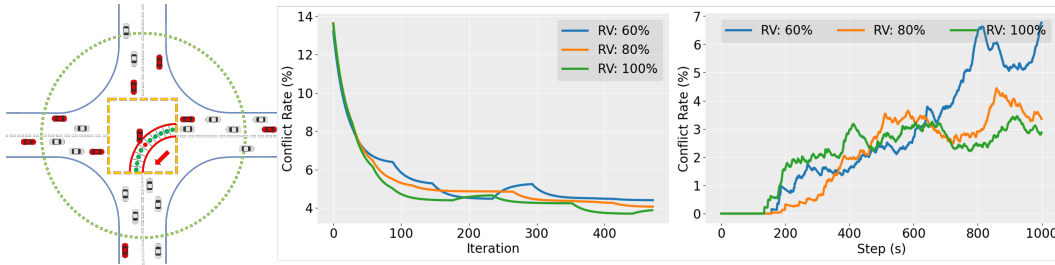


Figure 3: LEFT: The occupancy map along the moving direction W-L. Each of the 10 segments is labeled with either free (green dot) or occupied (red dot). MIDDLE: As learning progresses, the frequency of conflicting decisions decreases and stabilizes at a low level over all three RV penetration rates. RIGHT: Regardless of the RV penetration rate, the conflict rate (calculated as the number of conflict decisions divided by the total number of RVs' decisions) stays low, e.g., $\sim$6% for 60%, $< 4\%$ for 80% RVs, and $< 3\%$ for 100% RVs.

5

### 2.2.3 Conflict-Aware Reward

To encourage the RV to consider not only its own efficiency but also the potential conflicts within the intersection, we design a conflict-aware reward function for the RV:

$$r(s^t, a^t, s^{t+1}) = \lambda_L r_L + p_c, \tag{3}$$

where $r_L$ is the local reward, $p_c$ is the conflict punishment, and $\lambda_L$ is the coefficient.

The local reward $r_L$ is

$$\begin{cases} -w^{t+1,j}, & \text{if } a^t = \text{Stop}; \\ w^{t+1,j}, & \text{otherwise.} \end{cases} \tag{4}$$

$w^{t+1,j}$ is the average waiting time of all vehicles in the $j$th direction, which is normalized to $[0, 1]$, $p_c$ denotes the punishment for conflicts. If the RV decides Stop, the local reward is the negative waiting time $-w^{t+1,j}$; otherwise, it is positive $w^{t+1,j}$.

The conflict punishment $p_c$ is

$$\begin{cases} -1, & \text{if conflict}; \\ 0, & \text{otherwise.} \end{cases} \tag{5}$$

If the RV's movement conflicts with other vehicles in the intersection, it incurs a penalty of $-1$.

The reward design is inspired by the observation that waiting time [Zhang et al., 2020, Gregurić et al., 2020] has been used to measure traffic congestion. We demonstrate the effectiveness of our reward design in improving intersection traffic through extensive experimentation. The analysis of the reward function can be found in Sec. 3.3.

### 2.2.4 RL Algorithm

We employ the Rainbow DQN algorithm [Hessel et al., 2018], which is a state-of-the-art technique that combines six extensions of the original DQN algorithm [Mnih et al., 2015]. We use it and the reward to centrally train all RVs. During execution, each RV makes independent decisions by utilizing the same policy—a neural network with three fully connected (FC) layers and each FC layer contains 512 hidden units and uses ReLU as the activation function. We train the policy with a learning rate of 0.0005, a discount factor of 0.99, and the minibatch size of 32. Training takes approximately 48 hours using Intel i9-13900K and NVIDIA GeForce RTX 4090.

## 2.3 Conflicting Traffic Streams

The primary cause of intersection congestion and accidents is conflicting directions of movement. Despite penalizing RVs for conflicting decisions, our reward function cannot entirely eliminate conflicts, i.e., simultaneous Go decisions from RVs at conflicting directions. Coordinating a multi-agent system using RL is difficult [Lowe et al., 2017], and developing autonomous systems that are both effective and provably safe remains an open challenge [Gu et al., 2022].

Although conflicting decisions among RVs at the intersection entrance may not result in crashes, to prevent congestion and ensure safety, we implement a fail-safe coordination mechanism to post-process RL decisions. To begin with, if there are no vehicles on the conflicting stream or inside the intersection, and no conflicting decisions among the RVs, an ego RV who decides Go will enter the intersection. However, if there are vehicles inside the intersection, particularly on the conflicting stream of the ego RV, it is not permitted to enter. When multiple RVs on conflicting streams arrive at the intersection entrance and all decide Go, the RV with the highest priority score, calculated by averaging waiting time and queue length, is granted entry, while the others must wait.

We evaluate our reward function in avoiding conflicts and the use of the fail-safe mechanism. As depicted in Fig. 3 MIDDLE, the number of conflicts decreases as training progresses, stabilizing at a low level. We further investigate the conflict rate by calculating it as the ratio of conflicts to the number of RVs within the control zone. As shown in Fig. 3 RIGHT, the conflict rate for either 60% RVs or 80% RVs tends to converge around 4%, while the conflict rate for 100% RVs is less than 3% after 500 steps. These results demonstrate the effectiveness of our reward in coordinating intersection traffic and the infrequent usage of the fail-safe coordination mechanism.

## 2.4 Assumptions of Robot Vehicles

We define robot vehicles (RVs) in this project to differ from conventional autonomous vehicles (AVs). While AVs usually include a full suite of perception-to-planning modules, our RVs focus on high-level decisions (Stop/Go) and only require basic communication to obtain the positions and decisions of other vehicles within the control zone. Other sensors for perception and motion planning of a vehicle, such as cameras and lidars, are unnecessary. Therefore, the learning process is different than a typical AV training process, end-to-end or otherwise. Furthermore, our design promotes human-AI cooperation by allowing human drivers to execute low-level controls while receiving high-level suggestions from AI. Incentives can be designed to encourage humans to follow these suggestions, thus contributing to a more efficient traffic system. Even if the suggestions are not followed regularly, the proposed control mechanism can still benefit mixed traffic as the RV penetration rate increases in the expected future (see Sec. 3.1). These features render our method applicable to all levels of vehicle autonomy, making it a more practical solution for mixed traffic control.

## 3 Experiments and Results

We evaluate our method by comparing it to four baseline models: 1) **TL**: the traffic signal program deployed in the city of Colorado Spring, CO; 2) **NoTL**: no traffic lights; 3) **Yan** [Yan and Wu, 2021]: the state-of-the-art RL traffic controller with 100% RV penetrate rate[1]; and 4) **Yang** [Yang and Oguchi, 2020]: the state-of-the-art CAV control method for unsignalized intersections.

### 3.1 Overall Performance

Table 1 shows the main results measured with reduced average waiting time in percentage at the four intersections shown in Fig. 1 LEFT. The waiting time is defined in Appendix A.2.1. We test RV penetration rates from 40% to 100%, conducting ten experiments at each rate and reporting the averaged results. Traffic is reconstructed using real-world data (see Appendix A.1 for details). Each experiment lasts 1000 steps (1000 $s$ in simulation). The location and behaviors of HVs in each experiment are stochastic. The performance of our method between different intersections are varied. With only 40% RVs, our method can considerably surpass the traffic light control in the intersection where we perform best. For other intersections, our method with 60% can achieve comparable or even better performance to the traffic signal control baseline. An example comparing our approach with using traffic lights on all moving directions at intersection 229 is shown in Fig. 5a. In the absence of traffic lights and with 100% RVs, we can achieve up to 97% reduction in average waiting time. These findings suggest that our approach can scale to various RV penetration rates and result in efficient and coordinated mixed traffic.

| | Reduced Average Waiting Time (%) | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Compared to TL | | | | | | Compared to NoTL |
| Intersection 229 | -30.67% | 5.06% | 30.94% | 49.78% | 42.68% | 81.73% | 97.88% |
| Intersection 449 | 67.55% | 77.47% | 76.87% | 75.97% | 77.10% | 79.35% | 83.89% |
| Intersection 332 | -13.05% | -24.60% | 20.41% | 30.71% | 25.09% | 63.37% | 88.54% |
| Intersection 334 | -33.59% | 19.39% | 3.69% | 61.09% | 45.96% | 69.71% | 89.01% |
| RV Rate | 50% | 60% | 70% | 80% | 90% | 100% | 100% |

Table 1: Reduced average waiting time in percentage at each intersection under different RV penetration rates. Our method outperforms traffic signals in all cases when the RV penetration rate is 60% or higher. Moreover, the time saved generally increases with higher RV penetration rates. In comparison to scenarios without traffic lights, with 100% RVs, we can achieve up to 97% reduction in average waiting time.

To evaluate the generalizability of our approach, we test it on two previously unseen intersections shown in Fig. 5b, one of which is a three-way intersection. We apply our policy directly to the unseen four-way intersection without refinement, and it functions well. As shown in Fig. 13, our policy can achieve comparable performance to conventional traffic signal control baseline when the RV

---

[1]To apply this approach, we extend the network input to the maximum number of incoming lanes in our study to accommodate real-world intersections.
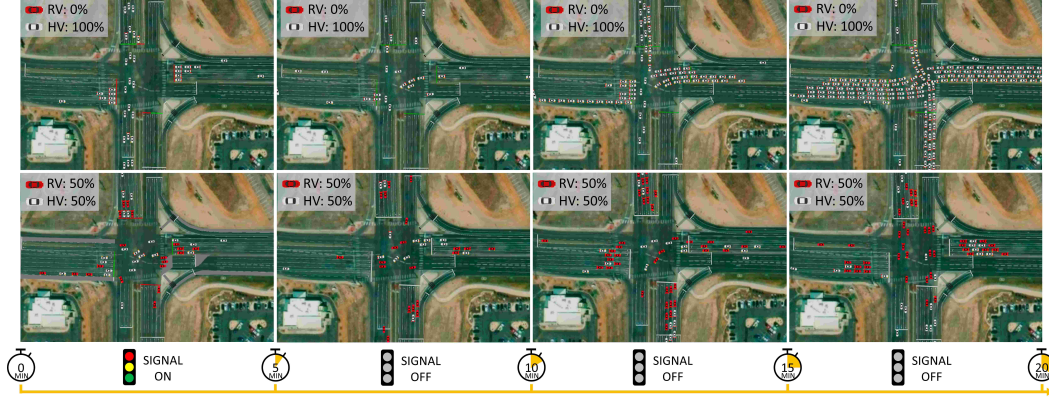
Figure 4: Comparison between traffic conditions with and without RVs during a blackout event at intersection 229. The blackout event occurs at the 5-minute mark. Congestion forms rapidly within 15 minutes in traffic without RVs. Conversely, traffic regulated with 50% RVs does not result in congestion.

penetration rate is 60% or higher. Our policy is also deployed without refinement on the unseen three-way intersection, which requires coordination among only four directions: S-C, S-L, W-L, and N-C. To adapt our policy, we set the input values of other directions that appear in four-way intersections to zero. Despite never having encountered this intersection topology and its traffic demand, our policy coordinates traffic well and prevents congestion. As shown in Fig. 14, our approach outperforms the traffic light control baseline when RV penetration is at 40% or higher. At 90% RV penetration, our method reduces average waiting time by approximately 62.5% compared to using traffic lights. These results demonstrate the excellent generalizablility of our approach. The statistics of the six intersections are given in Appendix A.2.2 and the detailed performance of our method at these intersections are provided in Appendix A.3.1.

To demonstrate our approach's robustness, we simulate blackout events in which all traffic signals are off. Fig. 4 shows the results for no RVs and 50% RVs. If no RVs are present, a gridlock forms at the intersection within 15 minutes after the traffic lights go off (starting from the 5th minute). In contrast, with 50% RVs, no congestion is observed. The quantitative results associated with Fig. 4 and more discussion about this evaluation are provided in Appendix A.3.2.

We also examine the effect of sudden RV rate drops due to unforeseeable reasons such as unstable vehicle-to-vehicle communication, software failures, and humans taking over the control. The 'offline' RVs are simulated using the IDM model [Treiber et al., 2000], which is used for all HVs. The results are shown in Fig. 5c. All drops occur at the 100th step. While the average waiting time of all vehicles at the intersection is increased, our approach rapidly stabilizes the system and prevent congestion by maintaining the average waiting time under certain thresholds.

## 3.2 Traffic Demands and Congestion

We further analyze the relationship of traffic demands and congestion. The results using intersection 229 as the testbed are shown in Fig. 5d. By increasing the traffic demand from 150 v/h to 300 v/h with no traffic lights and no RVs, we observe congestion starting to form at 200 v/h. (indicated by a low average speed of all vehicles at the intersection. In this paper, we define that congestion occurred when average speed inside the intersection is lower than 1 m/s.) In contrast, with the actual traffic demand at 700 v/h, congestion does not form with just 5% RVs deployed in traffic. Fig. 5d also demonstrates that the minimum RV penetration rate required to avoid congestion under the real-world traffic demand is 5%.

## 3.3 Analysis of the Reward Function

Designing effective rewards for controlling mixed traffic at complex intersections is challenging. Varied intersection topology, conflicting traffic streams, and the use of real-world traffic data which can lead to unpredictable and unstable inflow/outflow all complicate the task. To address conflicts within intersections and prevent traffic jams from negatively affecting the learning process, our insight
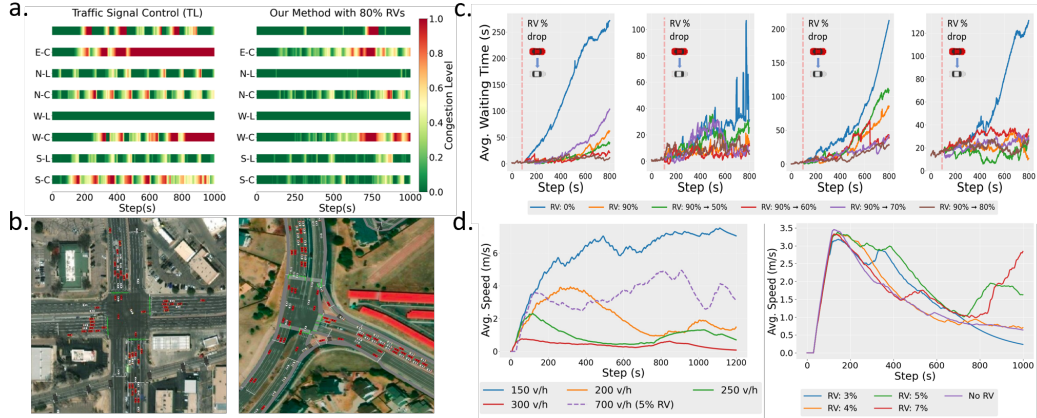
Figure 5: **a.** Our method results in efficient traffic flows in all moving directions compared to using traffic lights at intersection 229. **b.** Two unseen intersections used in testing our approach. **c.** Our method ensures stable and uncongested traffic, even the RV penetration rate abruptly drops. The sub-figures correspond to intersections 229, 332, 449, and 334 (from left to right). **d.** LEFT: The solid lines represent scenarios with no traffic lights and no RVs. Congestion starts to occur when the demand exceeds 200 v/h. However, the real-world demand (denoted by the dashed line), 700 v/h, does not lead to congestion due to the deployment of 5% RVs in the traffic. **d.** RIGHT: A minimum RV penetration rate of 5% is required to prevent congestion in traffic.

is to consider the impact of each RV's actions on traffic flow in its own direction while penalizing conflicting decisions among RVs.

Ideally, the designed reward should promptly reflect the congestion severity and traffic conditions of an intersection. In Fig. 6, we illustrate two traffic conditions (congestion formed in a. and b., traffic flow remains constant in c. and d.) with corresponding average vehicle speed in the intersection, our reward and Yan and Wu's reward [Yan and Wu, 2021]. The average speed of vehicles are used as an intersection congestion indicator. In this paper, if the average speed is constantly lower than 1 m/s, the intersection is seemed as congested.

As shown in Fig. 6a., with no traffic lights and 100% HVs, congestion arises, causing the average speed of all vehicles to decrease rapidly (orange curve). Our reward (blue curve) swiftly captures this trend, making it a timely indicator for the learning process. In comparison, the reward of the state-of-the-art method by Yan and Wu [Yan and Wu, 2021] in Fig. 6 b. fails to do so. In addition, we also show the counterexample of congestion traffic condition. As shown in Fig. 6c. and d., with 80% RVs, the traffic flow remain constant, our reward curve in c. is much higher than in a., which demonstrate the effectiveness of our reward function. However, the Yan and Wu's reward function keeps the same trend in both traffic scenarios. The further discussion can be found in Appendix A.3.3.
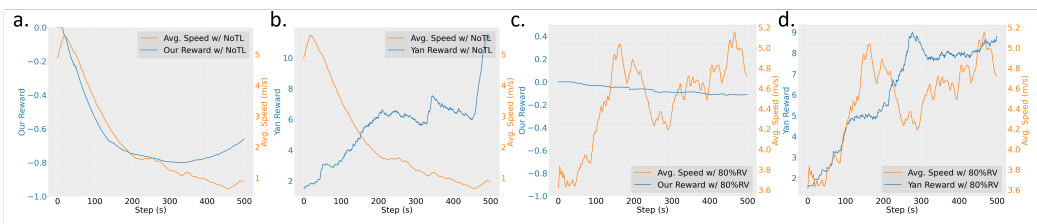


Figure 6: a. b. Congestion and deadlock occurred as the average speed inside the intersection is decreasing to near zero (the orange line) Our reward timely reflects congestion forming at the intersection. (as shown in a.) By contrast, The reward by Yan and Wu [Yan and Wu, 2021] fails to do the same. (as shown in b.) c. d. When 80% RVs are activated, the traffic flow remain uninterrupted, as the average speed maintains at a decent level. Comparing to a., our reward has a significant increase since the traffic condition is improved. However, the difference of the reward by Yan and Wu between two traffic conditions are very minor.

# 4 Conclusion and Future Work

We propose a decentralized RL approach for the control and coordination of mixed traffic at real-world and unsignalized intersections. Our approach consists of many novel designs for mixed traffic control, including a conflict-aware reward function that is suitable for coordinating large-scale traffic and a generic representation for encoding intersection traffic conditions. Our method is the first to control mixed traffic under real-world traffic conditions at complex intersections. Various experiments are conducted to show the effectiveness, robustness, and generalizablility of our approach. Detailed analysis are also pursued to justify the design of the components of our method.

In the future, we would like to further improve our method in several aspects. First, the learning algorithm could use a hierarchical design so that the low-level control (e.g., longitudinal and lateral acceleration) can be incorporated into the RL policy's output. Second, we want to ease the coordination mechanism so that vehicles have more freedom to move inside the intersection. Nevertheless, we anticipate that certain prior knowledge remains necessary for ensuring no-conflict movements. Lastly, we aim to integrate our approach with traffic flow prediction to enhance large-scale traffic coordination. Since real-world traffic demands vary over time, accurate flow predictions are invaluable for optimizing intersection traffic control and coordination.

# References

Michael Behrisch, Laura Bieker, Jakob Erdmann, and Daniel Krajzewicz. SUMO–simulation of urban mobility: an overview. In *International Conference on Advances in System Simulation*, 2011.

Xiaolong Chen, Manjiang Hu, Biao Xu, Yougang Bian, and Hongmao Qin. Improved reservation-based method with controllable gap strategy for vehicle coordination at non-signalized intersections. *Physica A: Statistical Mechanics and its Applications*, page 127953, 2022.

Eun-Ha Choi. Crash factors in intersection-related crashes: An on-scene perspective. *National Highway Traffic Safety Administration, U.S. Department of Transportation*, 2010.

Shuo Feng, Xintao Yan, Haowei Sun, Yiheng Feng, and Henry X Liu. Intelligent driving intelligence test for autonomous vehicles with naturalistic and adversarial environment. *Nature communications*, 12(1):1–14, 2021.

Martin Gregurić, Miroslav Vujić, Charalampos Alexopoulos, and Mladen Miletić. Application of deep reinforcement learning in traffic signal control: An overview and impact of open traffic data. *Applied Sciences*, 10(11):4011, 2020.

Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*, 2022.

Mohammed A Hadi and Charles E Wallace. Hybrid genetic algorithm to optimize signal phasing and timing. *Transportation Research Record*, (1421):104–112, 1993.

Mohammed A Hadi and Charles E Wallace. Optimization of signal phasing and timing using cauchy simulated annealing. *Transportation Research Record*, 1456:64–71, 1994.

Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2018.

PB Hunt, DI Robertson, RD Bretherton, and RI Winton. Scoot-a traffic responsive method of coordinating signals. Technical report, Transport and Road Research Laboratory (TRRL), United Kingdom, 1981.

Kathy Jang, Eugene Vinitsky, Behdad Chalaki, Ben Remer, Logan Beaver, Andreas A Malikopoulos, and Alexandre Bayen. Simulation to scaled city: zero-shot policy transfer for traffic control via autonomous vehicles. In *ACM/IEEE International Conference on Cyber-Physical Systems*, pages 291–300, 2019.

Duowei Li, Jianping Wu, Feng Zhu, Tianyi Chen, and Yiik Diew Wong. Coor-plt: A hierarchical control model for coordinating adaptive platoons of connected and autonomous vehicles at signal-free intersections based on deep reinforcement learning. *arXiv preprint arXiv:2207.07195*, 2022.

Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.

Andreas A Malikopoulos, Christos G Cassandras, and Yue J Zhang. A decentralized energy-optimal control framework for connected automated vehicles at signal-free intersections. *Automatica*, 93: 244–256, 2018.

David Miculescu and Sertac Karaman. Polling-systems-based autonomous vehicle coordination in traffic intersections with no traffic signals. *IEEE Transactions on Automatic Control*, 65(2): 680–694, 2019.

Amir Mirheli, Mehrdad Tajalli, Leila Hajibabai, and Ali Hajbabaie. A consensus-based distributed trajectory control in a signal-free intersection. *Transportation research part C: emerging technologies*, 100:161–176, 2019.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

United Nations. World urbanization prospects: The 2018 revision (st/esa/ser.a/420). *Department of Economic and Social Affairs, Population Division, New York: United Nations*, 2019.

Associated Press. Power still out to 50k customers, days after memphis storm. https://www.usnews.com/news/best-states/tennessee/articles/2022-02-07/power-still-out-to-60k-customers-days-after-memphis-storm, February 2022.

Rachel Ramirez. Power outages are on the rise, led by texas, michigan and california. here's what's to blame. https://www.cnn.com/2022/09/14/us/power-outages-rising-extreme-weather-climate/index.html, September 2022.

Jackeline Rios-Torres and Andreas A Malikopoulos. A survey on the coordination of connected and automated vehicles at intersections and merging at highway on-ramps. *IEEE Transactions on Intelligent Transportation Systems*, 18(5):1066–1077, 2016.

David Schrank, Bill Eisele, Tim Lomax, and Jim Bak. Urban mobility scorecard. *Texas A&M Transportation Institute and INRIX*, 2021.

Guni Sharon and Peter Stone. A protocol for mixed autonomous and human-operated vehicles at intersections. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 151–167, 2017.

Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000.

Mallory Trouve, Gaele Lesteven, and Fabien Leurent. Worldwide investigation of private motorization dynamics at the metropolitan scale. *Transportation Research Procedia*, 48:3413–3430, 2020.

Chris Urmson et al. Self-driving cars and the urban challenge. *IEEE Intelligent Systems*, 23(2):66–68, 2008.

Eugene Vinitsky, Kanaad Parvate, Aboudy Kreidieh, Cathy Wu, and Alexandre Bayen. Lagrangian control through deep-rl: Applications to bottleneck decongestion. In *IEEE International Conference on Intelligent Transportation Systems*, pages 759–765, 2018.

Ben Winck. Get ready for blackouts from london to la, as the global energy crisis overwhelms grids and sends energy prices skyrocketing. https://www.businessinsider.com/global-europe-energy-crisis-power-electricity-outages-blackouts-energy-grid-2022-9?op=1, September 2022.

Cathy Wu, Abdul Rahman Kreidieh, Kanaad Parvate, Eugene Vinitsky, and Alexandre M Bayen. Flow: A modular learning framework for mixed autonomy traffic. *IEEE Transactions on Robotics*, 38(2):1270–1286, 2022.

Zhongxia Yan and Cathy Wu. Reinforcement learning for mixed autonomy intersections. In *IEEE International Intelligent Transportation Systems Conference*, pages 2089–2094, 2021.

Hao Yang and Ken Oguchi. Intelligent vehicle control at signal-free intersection under mixed connected environment. *IET Intelligent Transport Systems*, 14(2):82–90, 2020.

Rusheng Zhang, Akihiro Ishikawa, Wenli Wang, Benjamin Striner, and Ozan K Tonguz. Using reinforcement learning with partial vehicle detection for intelligent traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 22(1):404–415, 2020.

Anye Zhou, Srinivas Peeta, Menglin Yang, and Jian Wang. Cooperative signal-free intersection control using virtual platooning and traffic flow regulation. *Transportation research part C: emerging technologies*, 138:103610, 2022.

## A  Appendix

### A.1  Mixed Traffic

#### A.1.1  Reconstruction and Simulation

In order for RVs to interact with HVs under real-world traffic conditions, we need to first reconstruct traffic using actual traffic data and then carry on high-fidelity simulations. We reconstruct the intersection traffic using turning count data at each intersection provided by the city of Colorado Springs, CO, USA[2]. The turning count data records the number of vehicles moving in a particular direction at the intersection and is collected via in-road sensors such as infrastructure-mounted radars.

Given the GIS data (traffic data and digital map), we pursue traffic simulations in SUMO [Behrisch et al., 2011], a widely-adopted traffic simulation platform. In SUMO, a directed graph is used to describe the simulation area: each edge of the graph represents a road segment with an ID and a vehicle's route is defined by a list of edge IDs.

Vehicles are routed using jtcrouter[3] based on the turning count data. By default, jtcrouter will select edges that are close to the intersection as the starting and ending edges of a route, which can result in extremely short routes, thus affecting the simulation fidelity. To mitigate this issue, we adjust vehicle routes by proposing more suitable edges for vehicles to enter and leave the network. Specifically, for the traffic stream on the main road that connects the four intersections examined in this project, we assign the starting and ending edges of vehicle routes to the boundary of the main road. For the traffic stream on other roads, the starting edges are moved to the successor upstream intersection and the ending edges are moved to the successor downstream intersection. Note that the route planning of a vehicle is determined during traffic reconstruction. There is no re-planning of a vehicle's route during simulation when RVs and HVs are interacting.

After re-assigning the starting edge and ending edge of each route, 'extra traffic counts' can occur. For example, a vehicle traveling through intersection 334 from northbound can also travel through intersections 229, 449, and 332, contributing to the northbound count for all four intersections. To alleviate this problem, we consider the coordination of traffic flows among adjacent intersections to avoid traffic double-counting, and then refine the number of routes to ensure the turning counts in the simulation match the actual turning count data. To evaluate whether the simulated flow resembles the real-world flow in terms of turning counts, we adopt the absolute percentage error (APE):

$$APE = \frac{|TC_{\text{real}} - TC_{\text{sim}}|}{TC_{\text{real}}}, \tag{6}$$

where $TC_{\text{real}}$ and $TC_{\text{sim}}$ are the turning counts from the real-world traffic and simulated traffic, respectively. As a result, the APE scores for intersections 229, 499, 332, and 334 are 0.22, 0.21,

---

[2]https://coloradosprings.gov/
[3]https://sumo.dlr.de/docs/jtrrouter.html

0.16, and 0.17, respectively. Because APE = 0 represents a perfect match between simulated and real-world traffic, these low APE scores serve as validation for our simulation.

### A.1.2  Generating Mixed Traffic

To create a mixture of RVs and HVs, at each time step, newly spawned vehicles are randomly assigned to be either RV or HV according to a pre-specified RV penetration rate. For an HV, the longitudinal acceleration is computed using Intelligent Driver Model (IDM) [Treiber et al., 2000]. For an RV, when it is outside the control zone, IDM is again used to determine the longitudinal acceleration; if it is inside the control zone, the high-level decisions Stop/Go are determined by the RL policy, while its low-level longitudinal acceleration is determined by the formulas introduced in Sec. 2.2.1.

## A.2  Experiment Details

### A.2.1  Evaluation Metrics

One of our evaluation metrics is the average waiting time. We define the waiting time for each vehicle as the total consecutive time it remains stationary in the control zone after entering it. To obtain the waiting times, we utilize the SUMO API [Behrisch et al., 2011]. The average waiting time for a specific direction is the mean of the waiting times for all vehicles moving in that direction, while the average waiting time for an intersection is the mean of the waiting times for all vehicles present at that intersection.

### A.2.2  Intersections

Our RVs are trained at four intersections: 229, 449, 332, and 334. For testing, we evaluate their performance not only on the four intersections used in training, but also on two previously unseen intersections, 140 and 205. It is worth noting that intersection 205 is different from the others, as it is a three-way intersection. The details of these intersections are provided in Table 2.

| Intersection | Num. incoming lanes | Num. non-empty lanes | Traffic demand (v/h * lane) |
|---|---|---|---|
| 229 (four-way) | 21 | 19 | 1157 |
| 449 (four-way) | 19 | 18 | 1089 |
| 332 (four-way) | 18 | 17 | 928 |
| 334 (four-way) | 16 | 14 | 789 |
| 140 (four-way) | 24 | 24 | 987 |
| 205 (three-way) | 10 | 10 | 1115 |

Table 2: Intersection details. RVs are trained on intersections 229, 449, 332, and 334. They are tested on all six intersections. The traffic demand is extracted from real-world traffic data. However, during the data capture period, some lanes of the intersections do not have any cars passing through them. Thus, we list the number of non-empty lanes of each intersection.

## A.3  Additional Results and Analyses

### A.3.1  Intersection Performance

In Fig. 9, we show the detailed performance of our method at intersection 229. The results include two parts. The first part (the top row of the four figures) shows the average waiting time along the eight traffic moving directions.

The second part (the bottom row of the four figures) reports the influence of different RV penetration rates on the average waiting time. Similarly, Fig. 10, 11, and 12 illustrate the detailed performance at intersections 449, 332, and 334, respectively.

**Intersection 229**. As shown in Fig. 9, for all moving directions, NoTL and Yan perform the worst and are excluded from the zoomed-in sub-figure on the top, RIGHT row. From the zoomed-in sub-figure, we can see that the average waiting time is continuously reduced when the RV penetration rate is increased from 40% to 80% (except for W-C and W-L). Overall, TL and Yang have similar performance. However, their performance along different directions are different, for S-C, S-L, W-L,
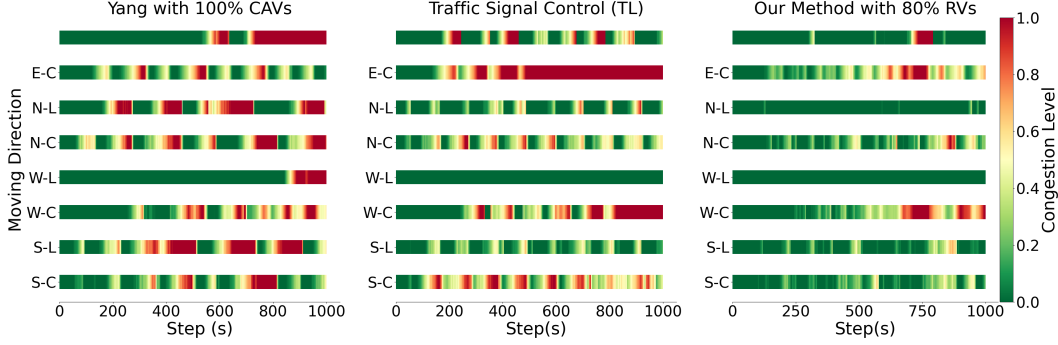
Figure 7: Traffic congestion levels at intersection 229 under different control mechanisms. Our approach with 80% RVs consistently achieves lower levels of congestion than Yang and TL. Unlike Yang and TL, which control intersection traffic using fixed phases, our method learns to use adaptive phases for control based on traffic conditions.

N-C, N-L, E-L, TL performs much better than Yang. In general, our method starts to outperform TL and Yang when the RV penetration rate is 60% or higher.

We further show traffic congestion levels of intersection 229 during our evaluation in Fig. 7. The congestion level is defined as AVT/Threshold, where AVT denotes the average waiting time of all vehicles of a moving direction, and Threshold is for normalization. For results shown in Fig. 7, Threshold is set to 40, which is the maximum average waiting time during our evaluation at intersection 229. The results illustrate that traffic controlled using our method achieves much lower congestion levels than Yang and TL. In addition, our method can flexibly coordinate conflicting moving directions based on varied traffic conditions, which is different than Yang and TL that employ fixed-phase coordination. These results hint that varied phases of control can positively influence the efficiency of intersection traffic.

**Intersection 449**. In general, similar results are observed as those of intersection 229 and are shown in Fig. 10. For most moving directions, the performances of TL are worse than ours, except for the direction S-L. Our method with 40% RVs or higher can overtake significantly Yang and TL in nearly all cases.

**Intersection 332**. The results are shown in Fig. 11. Similar to the intersections 229 and 449, Yan and NoTL are worse than Yang and TL, as well as our method with RV penetration rate 40% or higher. For Yang and TL, our method with at least 70% RVs can outperform them in general.

**Intersection 334**. The results are shown in Fig. 12. In general, the average waiting time decreases as the RV penetration rate increases. More specifically, the performances of our method is better than TL for most moving directions, except for the direction E-C and E-L. Note that in the actual traffic data, some directions do not have traffic (e.g., W-L) and thus leave blank in the figure. There is an interesting phenomenon where the waiting time of NoTL for some moving directions is lower than that of TL. This is because the traffic demands of different directions are not equal at intersection 334. For some directions, the traffic demands are low or even zero (e.g., W-L). Hence, the chance of congestion inside the intersection (and subsequently longer waiting time) is lower, which further attenuates the importance of traffic lights on efficient traffic flows.

Worth mentioning, across all results of all intersections, the average waiting time of all vehicles may not monotonically decrease when the RV penetration rate increases. The median waiting time of a higher RV percentage can be lower than the median waiting time of a lower RV percentage, e.g., 80% RVs vs 90% RVs at the intersection 334. This is because during repeated experiments, while traffic demands are matched between simulations, the actual data, behaviors, and positions of individual vehicles are stochastic. These unpredictable factors can lead to a large variance in performance. In addition, policy training convergence is also a possible reason for this phenomenon.

14

### A.3.2 Blackout

In Fig. 8, the blackout events occur at the 100th step. Without RVs, the absence of traffic lights leads to significant increases in average waiting time due to traffic jams. However, with the presence of 50% RVs, the average waiting time remains stable during the blackout events. Essentially, our method enables RVs to act as 'self-organized traffic lights,' effectively coordinating traffic at the intersection and preventing gridlocks.
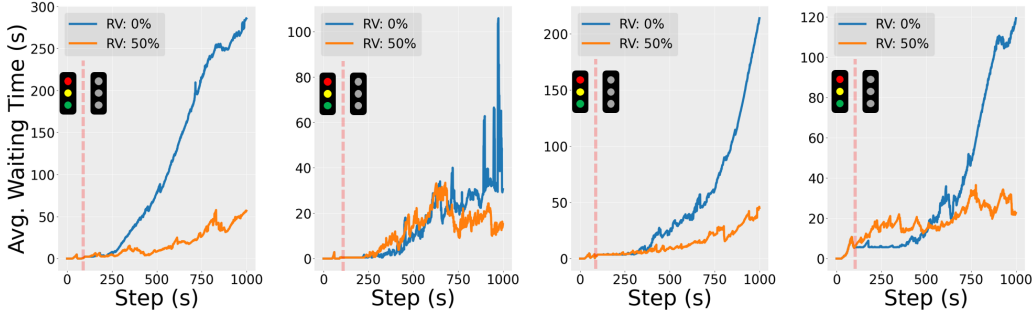


Figure 8: Blackout experiments. We simulate blackout events (traffic signals off) at intersections 229, 332, 449, and 334 (from left to right) since the 100th step. Without any RV, a gridlock will form at the intersection causing the average waiting time of all vehicles to increase rapidly. In contrast, with 50% RVs, no gridlock appears and the waiting time of all vehicles remain low and stable.

### A.3.3 Reward Analysis

The reward function by Yan and Wu [Yan and Wu, 2021] (Yan) takes the format $R_{Yan} = \texttt{outflow}(s_t, s_{t+1}) - \texttt{collision}(s_t, s_{t+1})$, where $\texttt{outflow}(s_t, s_{t+1})$ denotes the number of vehicles exiting the network from $t$ to $t + 1$, and $\texttt{collision}(s_t, s_{t+1})$ is the number of collisions in the network from $t$ to $t + 1$. We record this reward during the evaluation of the NoTL baseline (no traffic lights and 100% HVs) to analyze its characteristics. The result is shown in Fig. 6.

As anticipated, the absence of traffic lights leads to intersection congestion, as evidenced by the average speed of all vehicles rapidly approaches zero. However, the Yan Reward fails to capture changes in traffic conditions promptly. This is due to the outflow of a network being a delayed indicator; congestion within the intersection does not prevent previously cleared vehicles from contributing to the outflow. As a result, the delayed reward hinders the learning process, with episodes often terminating due to congestion before the reward can manifest it.
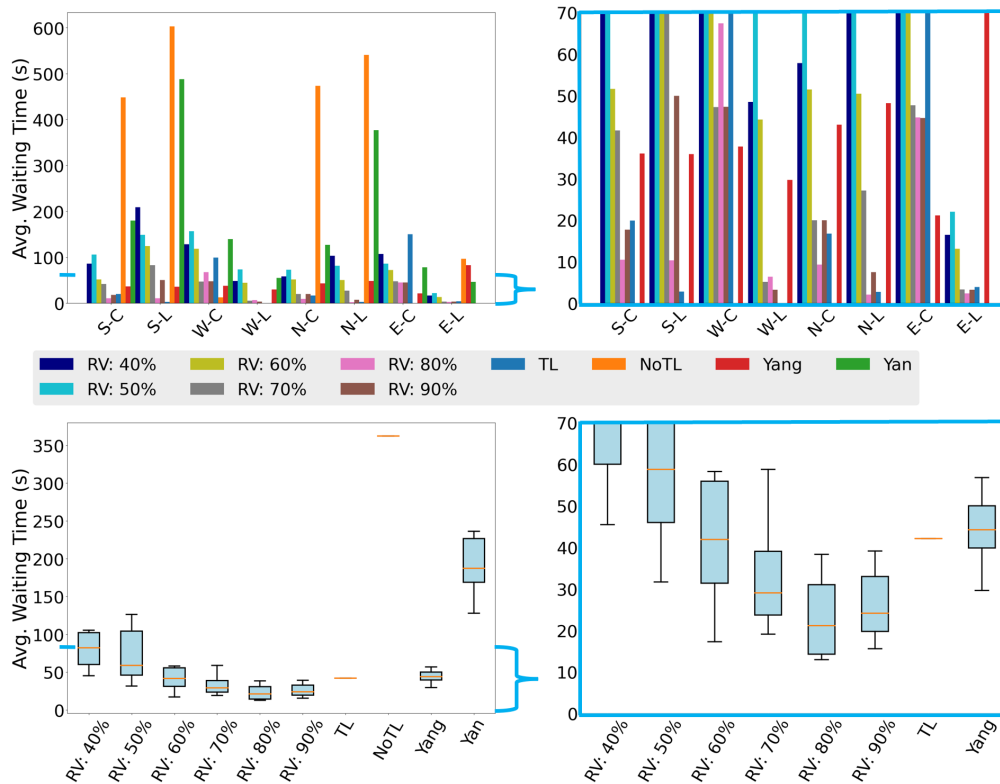
Figure 9: The overall results measured in average waiting time at the intersection 229. The RIGHT sub-figures are zoomed-in versions of the LEFT sub-figures by excluding NoTL and Yan. In general, as the RV penetration rate equals or passes 60%, our method achieves consistent better performance over the other four baselines.
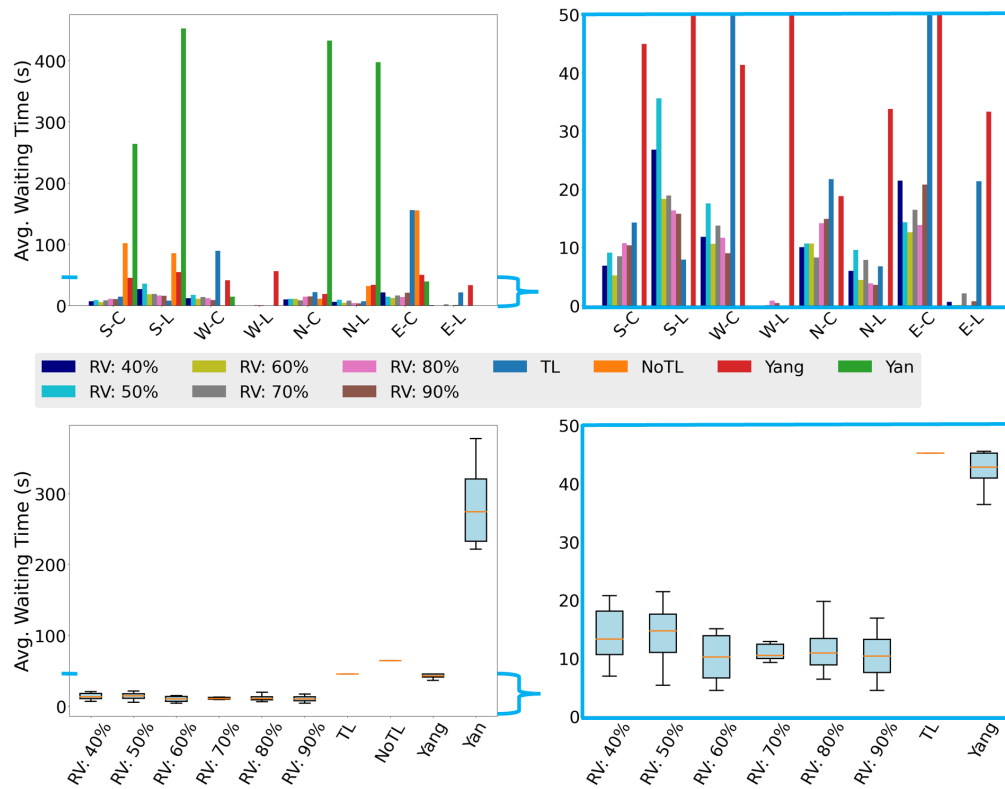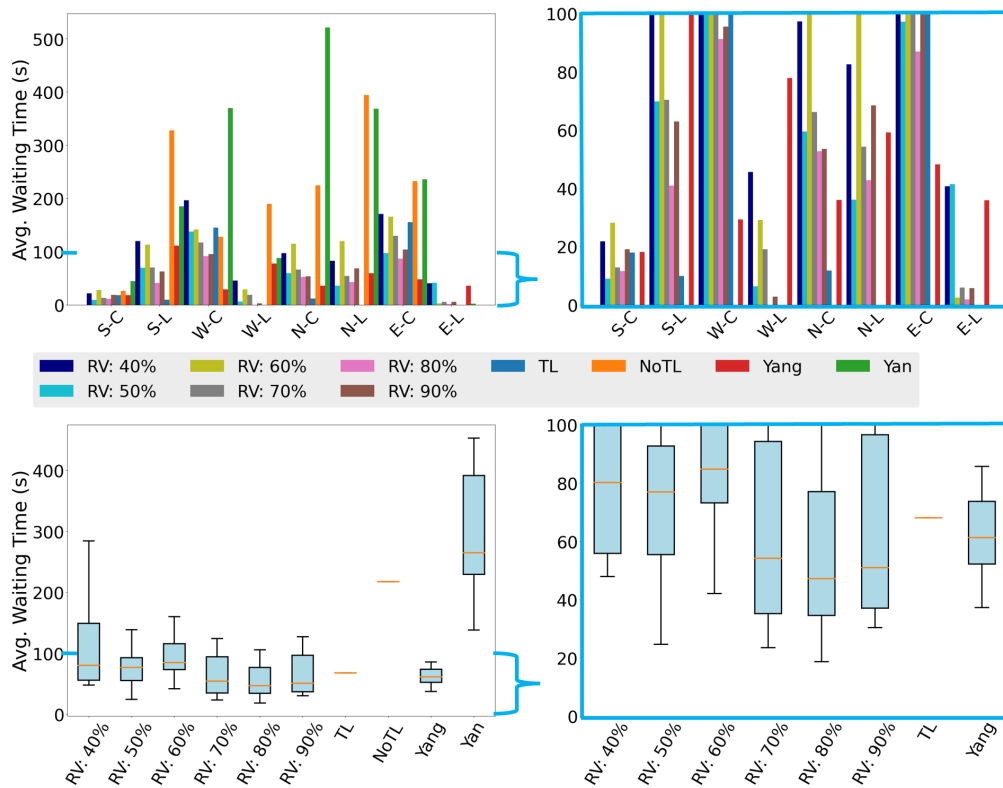
Figure 10: The overall results measured in average waiting time at the intersection 449. The RIGHT sub-figures are zoomed-in versions of the LEFT sub-figures by excluding NoTL and Yan. With 40% or more RVs, our method consistently outperforms all other baseline methods.

Figure 11: The overall results measured in average waiting time at the intersection 332. The RIGHT sub-figures are zoomed-in versions of the LEFT sub-figures by excluding NoTL and Yan. Generally speaking, NoTL and Yan do not perform very well. Our method starts to outperform TL and Yang when the RV penetration rate is 70% or higher.
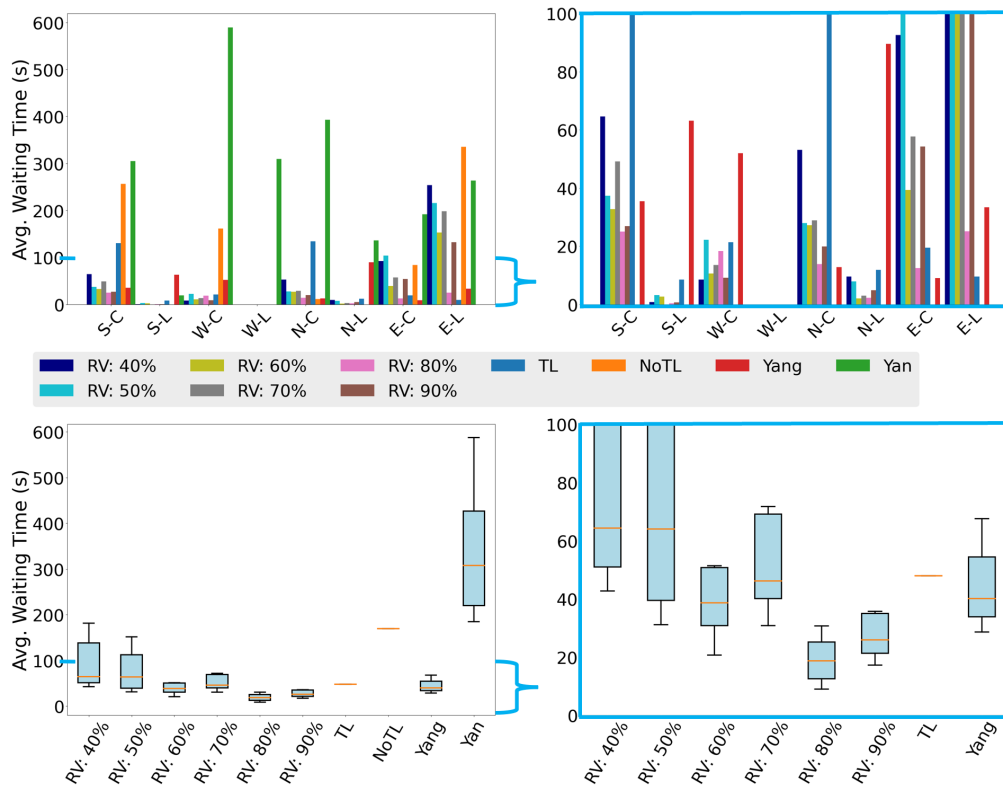
Figure 12: The overall results measured in average waiting time at the intersection 334. The RIGHT sub-figures are zoomed-in versions of the LEFT sub-figures by excluding NoTL and Yan. In general, our method with 60% RVs or more outperforms all four baselines.
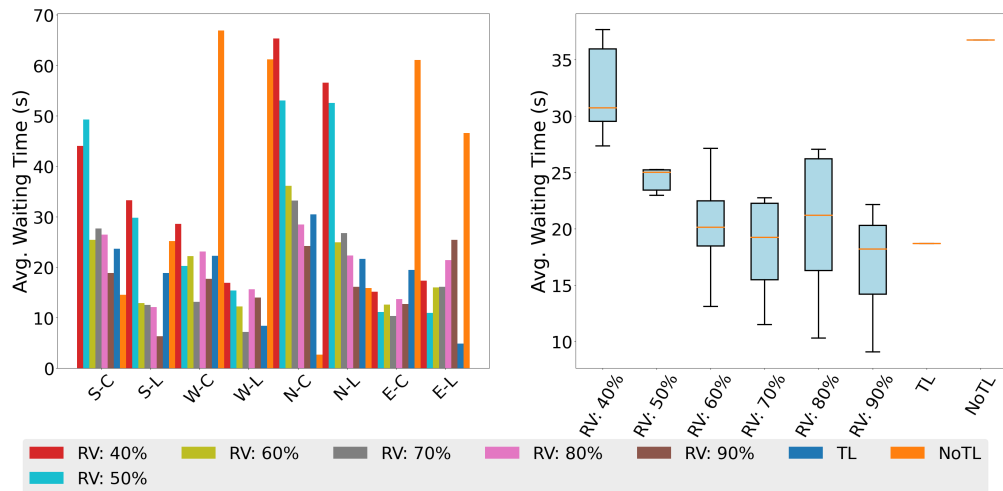


Figure 13: The overall results measured in average waiting time at the intersection 140 (an unseen four-way intersection during training). Starting from 60-70% RVs, our method beats the traffic lights (TL) baseline.
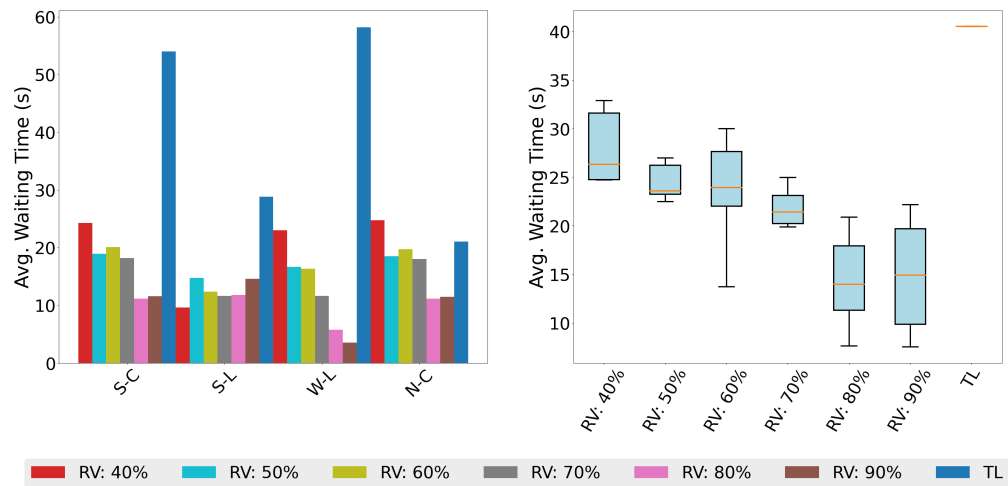
19

Figure 14: The overall results measured in average waiting time at the intersection 205. This is a three-way, unseen intersection during training and thus only four directions are shown. The RIGHT sub-figures are zoomed-in versions of the LEFT sub-figures by excluding NoTL. Our approach starts to outperform the TL baseline when RVs are 40% or more.