

EnduRL: Enhancing Safety, Stability, and Efficiency of Mixed Traffic Under Real-World Perturbations Via Reinforcement Learning

Bibek Poudel¹, Weizi Li¹, Kevin Heaslip²

Abstract—Human-driven vehicles (HVs) amplify naturally occurring perturbations in traffic, leading to congestion – a major contributor to increased fuel consumption, higher collision risks, and reduced road capacity utilization. While previous research demonstrates that Robot Vehicles (RVs) can be leveraged to mitigate these issues, most such studies rely on simulations with simplistic models of human car-following behaviors. In this work, we analyze real-world driving trajectories and extract a wide range of acceleration profiles. We then incorporate these profiles into simulations for training RVs to mitigate congestion. We evaluate the safety, efficiency, and stability of mixed traffic via comprehensive experiments conducted in two mixed traffic environments (Ring and Bottleneck) at various traffic densities, configurations, and RV penetration rates. The results show that under real-world perturbations, prior RV controllers experience performance degradation on all three objectives (sometimes even lower than 100% HVs). To address this, we introduce a reinforcement learning based RV that employs a congestion stage classifier to optimize the safety, efficiency, and stability of mixed traffic. Our RVs demonstrate significant improvements: safety by up to 66%, efficiency by up to 54%, and stability by up to 97%.

I. INTRODUCTION

Every major city in the world faces traffic congestion, from infrastructure-induced bottlenecks on bridges to more subtle phenomena like phantom jams on highways. Unsteady traffic flow during congestion increases travel time, energy consumption, and collision risk [1]. While the causes of congestion are not fully comprehended, a prominent explanation attributing to congestion is the asymmetric driving theory [2]. This theory posits that human driving is characterized by frequent under- and over-reactions, stemming from heterogeneous driving styles and individual differences in estimation, reaction, and actuation times. These factors will intensify in dense traffic and eventually lead to congestion [3].

As more vehicles with varying levels of autonomy are introduced into our transportation system, the concept of *mixed traffic control*, i.e., using Robot Vehicles (RVs) to address perturbations produced and amplified by Human-driven Vehicles (HVs), has emerged [4]–[7]. Various control strategies are introduced such as those based on heuristics [8], models of longitudinal dynamics [9], [10], and machine learning [11]. These techniques have been proven

effective in reducing congestion in scenarios such as intersections [12], highways [13], bottlenecks [14], and large networks [15] at RV penetration rates as low as 5%. Most of these studies are carried out in simulation using software such as Simulation of Urban MObility (SUMO) [16], where model-based methods dominate in representing car-following behaviors of HVs [17], [18].

Among various car-following models for enabling longitudinal control of a vehicle, Intelligent Driver Model (IDM) [19] is a popular choice with minimal parameters. To simulate heterogeneous driving behaviors that account for drivers’ individual differences, stochasticity is often integrated with the model, and imperfections in perception, processing, and actuation are accounted for in the simulation [20]. However, even with the adjustments, car-following behaviors of HVs are largely constrained: 92% of the accelerations derived using the stochastic IDM model with the default parameters [17] lie within the range $[-0.5, 0.5] m/s^2$. In contrast, only 68% of real-world accelerations (extracted from the I-24 MOTION dataset [21]) fall within the same range, resulting in a **24% difference**. Further, real-world accelerations have a *long tail* that extends up to $[-3, 3] m/s^2$. This observation reveals that simulated car-following mainly depicts safe or timid behaviors, lacking the representation of more aggressive real-world behaviors such as sudden braking and rapid accelerations [22], [23]. Subsequently, *RVs developed and validated in improving traffic conditions with only timid driving behaviors may struggle to handle real-world perturbations*, thereby impeding their effectiveness in achieving the goals of safety, stability, and efficiency.

Limited research to date bridges the gap between IDM and real-world driving [24]. Enhancements to IDM, including the integration of random noise and calibration with real-world driving data, are proposed to improve simulation accuracy. However, they fall short in capturing the full traffic variability and often lack broad applicability due to restrictive assumptions and artificial constraints on vehicle behaviors [25]–[30]. We introduce EnduRL, a framework overcomes these limitations by periodically sampling accelerations from the real-world dataset during car-following. This produces increased variability in driving behavior with more aggressive accelerations and decelerations. In addition, we propose a reinforcement learning (RL)-based RV that leverages downstream traffic information to forecast congestion stages and takes preemptive actions to improve traffic conditions.

We compare EnduRL with RVs developed in prior studies through comprehensive experiments in two mixed traffic environments: Ring [31] and Bottleneck [14] (see Fig. 1)

¹Bibek Poudel and Weizi Li are with Min H. Kao Department of Electrical Engineering and Computer Science at University of Tennessee, Knoxville, TN, USA bpoudel13@vols.utk.edu, weizili@utk.edu

²Kevin Heaslip is with Department of Civil and Environmental Engineering at University of Tennessee, Knoxville, TN, USA kheaslip@utk.edu

under 1500+ simulation runs. To evaluate safety, we use two surrogate measures, time to collision (TTC) and deceleration rate to avoid a crash (DRAC); for efficiency we measure fuel economy and throughput; for stability, we measure acceleration variation and wave attenuation. These metrics are widely adopted in traffic engineering [10], [32], [33]. Our results show that in Ring, our RV improves both safety and efficiency by up to **54%**, and stability by up to **97%**. Whereas in Bottleneck, our RV can improve safety by up to **66%**, efficiency up to **41%**, and stability up to **34%**. To the best of our knowledge, EnduRL is among the few studies that address the crucial gap between car-following behaviors in simulation and the real world and significantly improves the safety, stability, and efficiency of mixed traffic compared to previous studies. The project code can be found in the repository: <https://github.com/poudel-bibek/EnduRL>

II. RELATED WORK

To improve simulation fidelity, enhancements to IDM such as incorporating random noise [19] and calibrating IDM using real-world trajectories [34], [35] are introduced. However, they are limited in reproducing the variability found in real-world traffic [36] and have limited generalizability [37]. Prior studies also impose artificial bounds that limit the acceleration range of HVs [25]–[30]. More recent studies have adopted machine learning for approximating human driving behaviors. These techniques include applying supervised learning to features extracted from real-world driving data, replacing IDM with deep neural networks [38], and adopting deep RL with reward function tuned based on real-world data [39], [40]. Other studies propose equipping HVs with bilateral information, i.e., from both leader and follower vehicles to improve the traffic condition [41]. The key difference of our work to most existing studies is that we impose no artificial approximation or bounds. Instead, we directly extract and sample car-following dynamics from a real-world traffic dataset [21].

Regarding traffic control, numerous RVs (as controllers) have been proposed [25]–[27] and benchmarked [28]–[30]. We broadly classify them into model-based, heuristic-based, and learning-based RVs. We test two model-based RVs, Bilateral Control Module (BCM) [9] and Linear Adaptive Cruise Control (LACC). BCM is a linear acceleration model that uses relative speed and headway information from both the leader and follower vehicles, whereas LACC considers only the vehicle in front. For LACC, we use the constant time headway model described by Rajamani [10] with the calibration parameters chosen to mimic the speed of IDM vehicles at equilibrium. Moreover, we evaluate two heuristic-based RVs, FollowerStopper (FS) [8] and Proportional-integral with saturation (PIwS) [8]. FS is a velocity controller that tracks a set average velocity and travels slightly below that velocity to open up a gap, allowing the RV to dampen oscillations and brake smoothly when needed. PIwS estimates the average equilibrium velocity of the vehicles in the network and then drives at the estimated velocity. Both FS and PIwS require a

calibration on desired (or equilibrium) velocity and may fail to stabilize traffic if the equilibrium velocity is set high.

III. METHODOLOGY

We introduce Intelligent Driver Model and detail the components of EnduRL.

A. Intelligent Driver Model (IDM)

IDM [18] assumes that drivers aim to maintain a safe distance from their leader while trying to reach their desired speed. IDM vehicles accelerate when the headway to the leader is large, and decelerate below a set maximum deceleration. The acceleration is given by:

$$a_{IDM} = a \left[1 - \left(\frac{v}{v_0} \right)^\delta - \left(\frac{s^*(v, \nabla v)}{s} \right)^2 \right],$$

where s^* is the desired headway:

$$s^*(v, \nabla v) = s_0 + \max(0, v \cdot T + \frac{v \nabla v}{2\sqrt{ab}}),$$

where a is the maximum acceleration, v is the velocity, v_0 is the desired velocity, δ is the acceleration exponent, s is the headway, ∇v is relative velocity to the leader vehicle, s_0 is the minimum gap, T is the desired time headway, and b is maximum deceleration. The parameter values are set according to Treiber and Kesting [17] as $a = 1$, $b = 1.5$, $T = 1$, $\delta = 4$, and $s_0 = 2$.

B. Sampling Real-world Perturbations

We apply a car-following-filter [37] to the I-24 MOTION [21] dataset and extract 172,000 instantaneous accelerations. Analyzing the frequency and duration of these accelerations reveals a negative correlation. To ensure that these real-world perturbations are accurately mirrored in the simulation, we adopt a probabilistic approach. First, we uniformly sample the frequency of the perturbations within the observed range of [10, 30] per HV, for every 6 minutes of car-following. Then, we sample acceleration intensity uniformly within $[-3, 3] m/s^2$. For each selected intensity A_i , we find the most common duration \hat{T}_{A_i} by linearly mapping A_i within the minimum (T_{\min}) and maximum (T_{\max}) observed durations ($T_{\text{range}} = T_{\max} - T_{\min}$). We then sample T_{A_i} , the duration of A_i , from a piecewise triangular distribution using the following conditional probability density function:

$$P(T_{A_i} | \hat{T}_{A_i}) = \begin{cases} \frac{2(T_{A_i} - T_{\min})}{T_{\text{range}}(\hat{T}_{A_i} - T_{\min})}, & T_{\min} \leq T_{A_i} < \hat{T}_{A_i}, \\ \frac{2(T_{\max} - T_{A_i})}{T_{\text{range}}(T_{\max} - \hat{T}_{A_i})}, & \hat{T}_{A_i} \leq T_{A_i} \leq T_{\max}. \end{cases}$$

Lastly, we randomly assign each sampled acceleration to a HV during the testing of RVs.

C. Reinforcement Learning with Congestion Stage Classifier

Our RL-based approach leverages Congestion Stage Classifier (CSC), a neural network trained with supervised learning on position and velocity from preceding cars inside a sensing zone (Fig. 1). The data is collected at various densities and categorized into six classes, namely ‘Forming’, ‘Leaving’, ‘Congested’, ‘Free flow’, ‘Undefined’, and

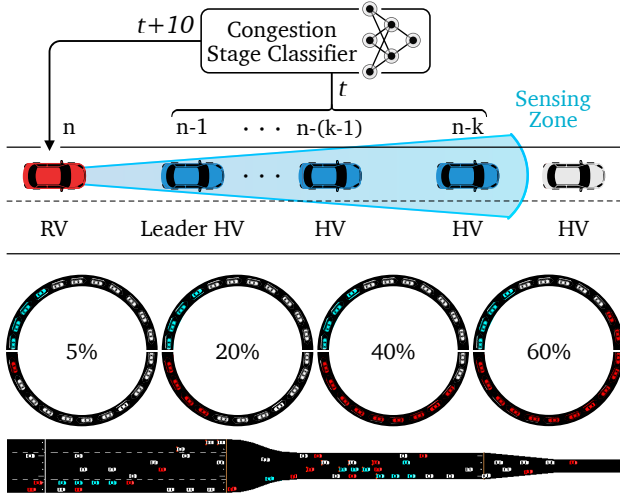


Fig. 1: TOP: The congestion stage classifier takes the position and velocity of the leader HVs of RV in the sensing zone to predict the congestion stage 10 timesteps into the future, enabling pro-active responses of the RV. MIDDLE: Our RVs deployed at various penetration rates in Ring. When penetration rates $> 5\%$, RVs are arranged as a platoon. BOTTOM: Our RVs deployed at 40% penetration rate in Bottleneck (truncated version shown).

‘No Vehicle’, inspired by asymmetric driving theory [2]. According to the theory, human drivers underestimate the required space headway during deceleration: when congestion is forming, availability of space headway decreases from one vehicle to the next as we move downstream within the sensing zone. Hence, a monotonic decrease in space headway is labeled congestion ‘Forming.’ The theory also suggests that human drivers overestimate the space headway during acceleration: when a congestion is relieved, availability of space headway increases from one vehicle to the next as we move downstream within the sensing zone. Such a monotonic increase in space headway is labeled congestion ‘Leaving.’ When all distances are above a threshold without a clear pattern, ‘Free-Flow’ is labeled; if all distances fall below a threshold, ‘Congested’ is labeled. With this labeling scheme, CSC predicts and classifies the congestion stage 10 timesteps in advance. We then incorporate the predictions by CSC into the observations and reward function of the RV. One CSC is trained for each environment: in Ring at density range 70–200 *veh/km* for 50 epochs with test accuracy of 95.5% and in Bottleneck at inflow range 3000–4500 *veh/hr* for 100 epochs, achieving test accuracy of 85.2%.

Since the objectives, safety, stability, and efficiency, have intrinsic conflicts [39] (e.g., optimizing for throughput may lead to aggressive driving behaviors, thus compromising safety and stability), we propose two types of RVs: the *safety + stability* RV, which prioritizes safety and stability, and the *efficiency* RV, which emphasizes efficiency. Both types of RVs leverage CSC, and share the action and observation space in both Ring and Bottleneck. We train our RVs using the PPO algorithm [42] with the following formulation:

- **Observation.** RV observes its leader vehicle’s position and

velocity, as well as the output of CSC:

$$\langle v_n, r_p(n, n-1), r_v(n, n-1) \rangle \oplus f_{\text{CSC}}(\{r_{p(i)}, r_{v(i)}\}_{i \in Z}),$$

where v_n is the velocity of the RV, $r_p(n, n-1)$ and $r_v(n, n-1)$ are the relative position and velocity respectively of RV with its immediate leader. For all $|Z|$ vehicles in the sensing zone (set to 50 *m*), the relative position $r_{p(i)}$ and velocity $r_{v(i)}$ are input to CSC to predict the congestion stage $f_{\text{CSC}}(\cdot)$.

- **Action.** RV controls its acceleration within $[-3, 3]$ *m/s*².
- **Reward.** Our reward is a linear combination of either RV velocity or average velocity of all vehicles, RV acceleration penalty, and a shaping component based on CSC output. For the two types of RVs, respective rewards are:

Reward Functions

Ring and Bottleneck: *efficiency* RV

```

reward ← 0.75 × v* − 2 × |an|
if fCSC(·) = Congested and sign(an) > 0 then
  reward ← reward + min(−1, λ1 · |an|)
end if
if fCSC(·) = Leaving and sign(an) < 0 then
  reward ← reward + λ2 · |an|
end if

```

Ring: *safety+stability* RV

```

reward ← 0.15 × v* − 4 × |an|
if fCSC(·) = Forming then
  reward ← reward + min(−1, λ3 · |an|)
end if

```

Bottleneck: *safety+stability* RV

```

reward ← 0.5 × v* − 4 × |an|
if fCSC(·) = Forming | Congested | Undefined then
  reward ← reward + λ4 · vn
  if sign(an) ≥ 0 then
    reward ← reward + min(−2, λ5 · |an|)
  end if
end if

```

where $v^* = \frac{4}{3n} \sum_{i=1}^n v_i$ in Ring with v_i as velocity of i^{th} vehicle and $v^* = v_n$ in Bottleneck, $f_{\text{CSC}}(\cdot)$ is the CSC output, a_n and v_n are the acceleration and velocity of the RV, respectively. $\lambda_1 = -10$, $\lambda_2 = -10$, $\lambda_3 = -5$, $\lambda_4 = -4$, and $\lambda_5 = -20$ are set empirically.

- **Scaling laws.** In Ring, for RV penetration rates $> 5\%$, both *safety + stability* and *efficiency* RVs are configured as platoons with a single leader and multiple followers (see Fig. 1 MIDDLE). For example, our *efficiency* RV platoon at 40% penetration consists of 9 RVs ($\lfloor 22 \times 0.4 \rfloor = 9$) including a leader *efficiency* RV trained at 5% penetration rate and 8 follower RVs. The follower RVs observe the position and velocity of the entire platoon and optimize the following reward:

$$\lambda_6 \cdot r_p(n, n+j) + \lambda_7 \cdot r_v(n, n+j) + \lambda_8 \cdot |a_{n+j}| + \lambda_9,$$

where $r_p(n, n+j)$ and $r_v(n, n+j)$ are relative position and relative velocity of the j^{th} follower with the platoon leader, respectively; and a_{n+j} is the acceleration of the j^{th} follower. $\lambda_6 = -2$, $\lambda_7 = 4$, $\lambda_8 = -4$, and $\lambda_9 = 10$ are chosen empirically. Whereas in Bottleneck, our RVs are dispersed and share the same RL policy.

| Controller | Min. no. of Stabilizing vehicles | Time to Stabilize (s) | Average Velocity (m/s) |
|------------|----------------------------------|-----------------------|------------------------|
| IDM | Unstable | Unstable | 3.58 |
| FS | 1 | 108 | 5.20 |
| PIwS | 1 | 119 | 5.33 |
| RL+L | 1 | 74 | 5.25 |
| Ours (5%) | 1 | 130 | 5.28 |
| BCM | 4 | 146 | 5.26 |
| Ours (20%) | 4 | 59 | 4.86 |
| LACC | 9 | 690 | 5.25 |
| Ours (40%) | 9 | 185 | 5.53 |

TABLE I: RVs’ stabilization metrics in Ring at 81 *veh/hr* with HVs enabled by IDM. Stability is established when the standard deviation of the average velocity is below the IDM noise threshold 0.2 [28]. Metrics are averaged over 10 randomized simulation runs.

We benchmark EnduRL with two RL-based RVs. In Ring, we use RV(s) with only microscopic/local observations such as the relative position and velocity to its leader vehicle [11], referred to as RL+L hereafter. Whereas in Bottleneck, we use RV(s) with macroscopic/global observations such as traffic density [29], referred to as RL+G hereafter.

IV. EXPERIMENTS

We introduce the mixed traffic environments, the evaluation metrics, the experimental setup, and finally the results. To begin with, we test on two mixed traffic environments.

- **Ring**: a single-lane circular road network with 22 vehicles (see Fig. 1 MIDDLE). This classical scene illustrates congestion development without external disturbances.
- **Bottleneck**: a straight road where the number of lanes reduces from 8 to 4 and then to 2. This simulates vehicles experiencing a capacity drop [43] (see Fig. 1 BOTTOM).

Our evaluation metrics include the following.

- **Time to Collision (TTC)**: the time interval between two vehicles that will collide if they maintain a relative speed difference [44]:

$$TTC = \begin{cases} \frac{s-l}{v_f-v_l}, v_f > v_l, \\ \infty, v_f \leq v_l. \end{cases}$$

s is the space headway, l is the vehicle length, and v_f, v_l are the velocities of the RV and its HV leader, respectively. A lower TTC indicates a higher risk.

- **Deceleration Rate to Avoid a Crash (DRAC)**: the force experienced by a vehicle in an emergent braking in order to avoid a front-end collision:

$$DRAC = \begin{cases} \frac{(v_f-v_l)^2}{s-l}, v_f > v_l, \\ 0, v_f \leq v_l. \end{cases}$$

s is the space headway, l is the vehicle length, and v_f, v_l are the velocities of the RV and its HV leader, respectively. A lower DRAC represents a safer situation [45].

- **Fuel Economy (FE)**: the average fuel consumption of all vehicles in the network measured in miles per gallon

(*mpg*) using the Handbook Emission Factors for Road Transport 3 Euro 4 passenger car emission model [46].

- **Throughput**: the network flow rate of vehicles (*veh/hr*).
- **Controller Acceleration Variation (CAV)**: the standard deviation of RV acceleration. A higher CAV meaning an RV is overly sensitive to its inputs, causing perturbations in upstream traffic [41].
- **Wave Attenuation Ratio (WAR)**: measures the effectiveness of an RV in dampening perturbations [33]. For all RVs, a standard velocity perturbation is applied to their immediate leader HV:

$$WAR = 1 - \frac{\Delta v_{\text{follow_HV}}}{\Delta v_{\text{lead_HV}}}.$$

$\Delta v_{\text{follow_HV}}$ is the velocity drop in the HV that immediately follows RV(s) and $\Delta v_{\text{lead_HV}}$ is the velocity drop in the HV that immediately leads RV(s). A higher WAR indicates better dampening effect.

To set up the experiments, we use FLOW [11] and SUMO [16]. Various vehicle configurations are adopted: RVs are platooned in Ring when penetration rate is $> 5\%$ and dispersed in Bottleneck at all penetration rates. Both configurations have shown to stabilize traffic [28]. We select penetration rates 5%, 20%, 40%, and 60% to align with the minimum rates required for stabilizing traffic by different controllers as shown in Table I. To pursue rigorous evaluations of safety (TTC and DRAC) and CAV, when multiple RVs are present, we report the worst-case values of EnduRL. For the IDM baseline, worst-case values of TTC and DRAC are again reported, whereas CAV is measured for a single randomly selected vehicle. For each experiment, an RV is first given sufficient time to stabilize traffic (as reported in Table I), followed by applications of acceleration perturbations for six minutes. Under these acceleration perturbations, we measure the safety of RVs. Since only RVs are responsible for stabilizing the system, we also measure the stability of RVs. All vehicles in the network are accounted when measuring traffic efficiency.

We next report results. Fig. 2 illustrates both the progressive amplification of longitudinal perturbation in the absence of RVs (100% IDM vehicles), and the strategies used by various RVs to attenuate such perturbations and stabilize traffic. We use Ring with density 81 *veh/km* to demonstrate. Worth noting, the phenomenon shown in Fig. 2 is environment-agnostic and can be replicated in any environment where longitudinal perturbation applies. Specifically, at the 1140th s for LACC, 150th s for IDM, and 800th s for the others, an identical velocity perturbation is applied to the HV that immediately leads the RV(s) (a randomly chosen HV in case of IDM), in which the velocity of the leader HV is abruptly decreased to 3 m/s for 2 s . For IDM, the perturbation is introduced before stop-and-go waves, while for the others, it is applied after stop-and-go waves, and the RV(s) are given sufficient time to stabilize the traffic. We observe that in the absence of RVs, IDM vehicles oscillate within the range [2, 8] m/s , periodically increasing and decreasing their velocity, and eventually forming a stop-and-go wave.

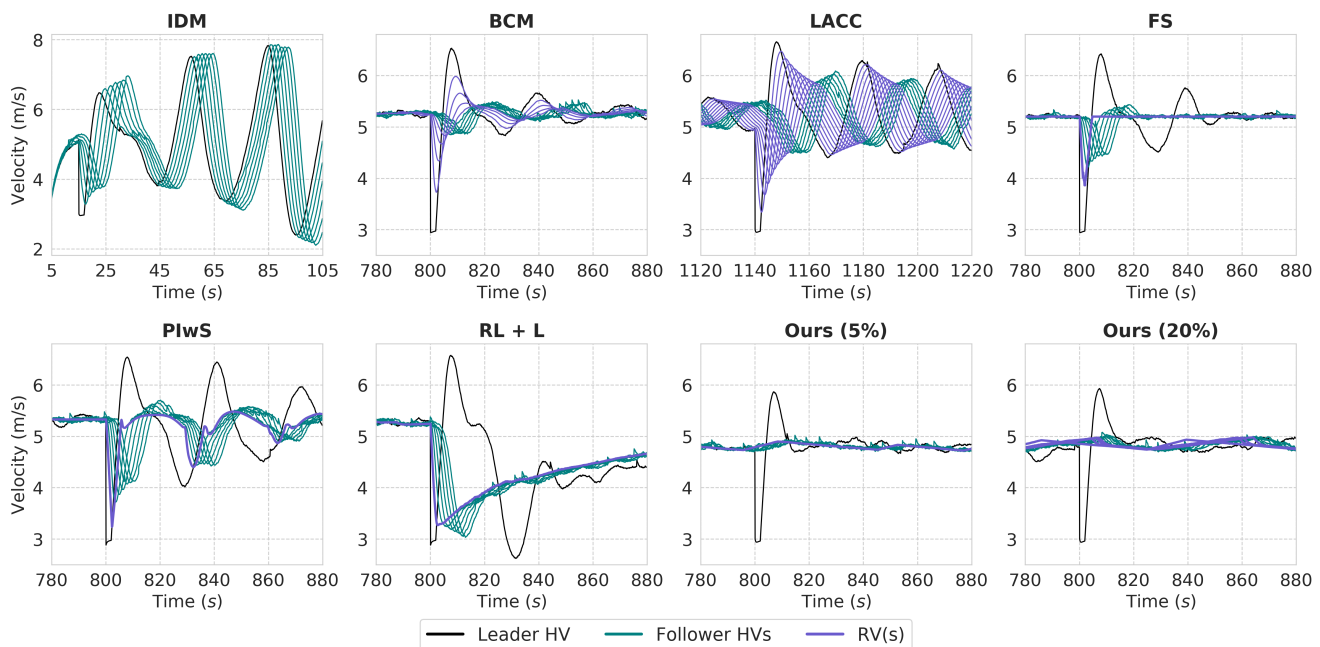


Fig. 2: Progressive amplification of longitudinal perturbation in the absence of RVs, and the strategies adopted by various RVs to attenuate such perturbations and stabilize traffic (6 HVs following the RVs are shown). Wave attenuation characteristics of various RVs at 81 veh/km are provided. The HV immediately in front of RV(s) produces a shock by applying a standard velocity perturbation of 3 m/s . With IDM (100% HVs), perturbation amplifies over time, whereas RVs dampen the perturbation over time.

In contrast, RVs attenuate the standard perturbation of the leader HV and stabilize traffic by adopting various strategies. For example, PIwS drives at an estimated average velocity and BCM maintains equal distance to its leader and follower HVs. The attenuation characteristics of our *safety + stability* RVs are shown in Fig. 2 at 5% and 20% penetration rates denoted by Ours (5%) and Ours (20%) accordingly. Compared to RL+L, our RVs open up a wider gap in front and track their HV leader at a slightly lower velocity allowing them to dampen incoming perturbations. For Ours (20%), multiple follower RVs closely track the leader RV. In Table II, the results of the *safety + stability* RV are shown in the Safety and Stability columns, while the results of the *efficiency* RV are presented in the Efficiency column. All values in Table II are averaged over 5 randomized simulation runs.

Table II LEFT presents the results of RVs in Ring at 85 veh/km density. Across all penetration rates, our *safety + stability* RV delivers the best safety measures. Notably, only our RV manages to exceed the critical 4 s threshold for TTC, as recommended by earlier studies [44], [47]. Additionally, DRAC is reduced by 54% at 5% penetration rate in comparison to IDM, and the CAV stays below 0.19 m/s^2 across all penetration rates. This is because the gap opened by our *safety + stability* RV dampens the perturbations caused by real-world disturbances. Furthermore, at 5%, 20%, 40%, and 60% penetration rates, our *efficiency* RV improves the throughput by 29%, 40%, 46%, and 54%, respectively, compared to IDM.

Table II RIGHT presents the results of RVs in Bottleneck at 150 veh/km peak density with an inflow rate of 3600 veh/hr . Unlike Ring, where multiple RVs form a platoon, compounding the individual dampening effects, in Bot-

tleneck, RVs are dispersed, with each RV aiming to dampen independently. Hence, the WAR for RVs remain constant at all penetration rates, among which our *safety + stability* RV maintains the highest WAR at 0.91, dampening majority of perturbations. Despite the high density in Bottleneck, our *safety + stability* RV exceeds the critical 4 s TTC threshold at 5% penetration rate, however, at higher penetration rates, TTC value is lower (but still the highest) among all other RVs. This is achieved by adopting the strategy of maintaining a front gap, consistent with its behavior in Ring. At all penetration rates, our *safety + stability* RV maintains the lowest DRAC (reduction of up to 66% compared to IDM).

Table II RIGHT also shows that our *efficiency* RV delivers the highest throughput at all penetration rates, up to 41% higher than IDM at 40% penetration rate. Since CSC does not encounter other RVs within the sensing zone during training, its accuracy decreases as more RVs appear in Bottleneck at higher penetration rates. This decline in performance is due to RVs behaving differently (usually lower accelerations and higher time gaps for *safety + stability* RV and higher accelerations and lower time gaps for *efficiency* RV) than IDM vehicles. Hence, CSC fails to anticipate the interactions between RVs and HVs within the sensing zone. Furthermore, zipper lanes (tapered sections that merge 8-4 and 4-2 lanes) break any achieved stability, e.g., a string of vehicles stabilized by an RV, when reaches the zipper lane, often brakes to allow vehicles from adjacent lanes to safely merge and cut-in-front. Lastly, the strategy adopted by our *efficiency* RV favors throughput over fuel economy as our RVs are not optimizing fuel economy via the reward function.

| RV Pen. Rate | RV Type | Ring | | | | | | Bottleneck | | | | | |
|--------------|----------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|
| | | Safety | | Efficiency | | Stability | | Safety | | Efficiency | | Stability | |
| | | TTC | DRAC | FE | Throughput | CAV | WAR | TTC | DRAC | FE | Throughput | CAV | WAR |
| | IDM* | 1.25 | 1.19 | 7.16 | 986 | 0.83 | Unstable | 1.23 | 5.40 | 15.84 | 1652 | 1.29 | Unstable |
| 5% | FS | 3.52 | 0.80 | 11.35 | 1214 | 0.47 | 0.54 | 1.40 | 5.07 | 16.09 | 1770 | 1.53 | 0.54 |
| | PIwS | 1.24 | 2.53 | 10.87 | 1260 | 0.34 | 0.61 | 1.35 | 5.00 | 16.37 | 1686 | 1.44 | 0.61 |
| | BCM | 1.03 | 2.87 | 7.58 | 1046 | 0.74 | Unstable | 1.51 | 4.93 | 16.06 | 1642 | 1.42 | 0.26 |
| | LACC | 1.08 | 1.40 | 7.26 | 1036 | 0.80 | Unstable | 1.44 | 5.65 | 15.31 | 1594 | 1.57 | 0.12 |
| | RL + L/G | 2.25 | 1.38 | 9.28 | 1060 | 0.27 | 0.47 | 1.24 | 4.74 | 15.24 | 1682 | 1.43 | 0.17 |
| | Ours | 5.07 | 0.54 | 11.44 | 1272 | 0.02 | 0.91 | 4.70 | 2.42 | 10.90 | 2086 | 0.85 | 0.91 |
| 20% | FS | 3.61 | 0.71 | 11.43 | 1318 | 0.39 | 0.92 | 1.32 | 5.81 | 16.02 | 1732 | 1.51 | 0.54 |
| | PIwS | 1.65 | 2.05 | 10.51 | 1296 | 0.22 | 0.92 | 1.25 | 5.86 | 16.69 | 1750 | 1.61 | 0.61 |
| | BCM | 1.12 | 1.49 | 12.28 | 1342 | 0.41 | 0.89 | 1.36 | 5.46 | 16.32 | 1666 | 1.61 | 0.26 |
| | LACC | 1.09 | 1.69 | 8.19 | 1148 | 0.82 | Unstable | 1.26 | 5.61 | 15.35 | 1704 | 1.75 | 0.12 |
| | RL + L/G | 1.97 | 1.49 | 4.91 | 654 | 0.12 | 0.47 | 1.42 | 5.01 | 16.45 | 1730 | 1.57 | 0.17 |
| | Ours | 5.38 | 0.53 | 11.99 | 1378 | 0.18 | 0.91 | 2.66 | 1.83 | 9.22 | 2270 | 0.94 | 0.91 |
| 40% | FS | 2.86 | 1.15 | 10.79 | 1346 | 0.38 | 0.97 | 1.26 | 5.70 | 17.16 | 1754 | 1.62 | 0.54 |
| | PIwS | 1.76 | 1.75 | 9.62 | 1284 | 0.21 | 0.90 | 1.25 | 5.89 | 16.86 | 1750 | 1.86 | 0.61 |
| | BCM | 1.96 | 1.02 | 10.93 | 1428 | 0.29 | 0.97 | 1.27 | 6.07 | 16.69 | 1726 | 1.78 | 0.26 |
| | LACC | 4.42 | 0.76 | 12.35 | 1444 | 0.28 | 0.72 | 1.19 | 5.48 | 16.76 | 1738 | 1.85 | 0.12 |
| | RL + L/G | 1.92 | 1.42 | 2.58 | 338 | 0.37 | 0.47 | 1.22 | 4.54 | 17.54 | 1694 | 1.44 | 0.17 |
| | Ours | 4.59 | 0.72 | 12.25 | 1448 | 0.18 | 0.91 | 1.86 | 2.28 | 8.68 | 2340 | 0.96 | 0.91 |
| 60% | FS | 2.50 | 0.95 | 8.86 | 1128 | 0.49 | 0.97 | 1.20 | 5.66 | 18.18 | 1826 | 1.96 | 0.54 |
| | PIwS | 1.60 | 1.94 | 9.67 | 1324 | 0.22 | 0.95 | 1.16 | 6.11 | 18.12 | 1884 | 1.96 | 0.61 |
| | BCM | 2.25 | 0.58 | 10.36 | 1382 | 0.23 | 0.97 | 1.17 | 6.15 | 18.26 | 1886 | 1.95 | 0.26 |
| | LACC | 2.90 | 0.85 | 11.77 | 1454 | 0.28 | 0.86 | 1.33 | 5.61 | 18.63 | 1868 | 1.84 | 0.12 |
| | RL + L/G | 1.88 | 1.57 | 2.03 | 280 | 0.25 | 0.47 | 1.42 | 4.41 | 19.90 | 1924 | 1.61 | 0.17 |
| | Ours | 4.42 | 0.81 | 12.58 | 1524 | 0.18 | 0.91 | 1.55 | 2.57 | 7.82 | 2034 | 1.04 | 0.91 |

TABLE II: Evaluations of various RVs and our method. IDM* denotes 100% HVs as baseline. LEFT: In Ring, under density 85 *veh/km*, our *safety + stability* RV uniquely achieves the highest safety measure (TCC > 4 s) and low oscillations (CAV < 0.19 *m/s²*) at all penetration rates, with a constant WAR (0.91). In contrast, WAR compounds for FS, PIwS, BCM, and LACC as penetration rate increases with BCM and LACC stabilizing traffic only at higher penetration rates (20%+ and 40%+, respectively). Our efficiency RV significantly improves the throughput (up to 54%) and fuel economy (up to 75%), while RL+L's efficiency declines as the number of RVs in traffic increases. RIGHT: In Bottleneck, under peak density 150 *veh/km*, our *safety + stability* RV maintains the highest WAR (0.91) at all penetration rates with the highest reduction in DRAC (up to 66%). Our *efficiency* RV experiences a significant increase in throughput at all penetration rates (up to 41%). Whereas at lower penetrations (5% and 20%) PIwS improves fuel economy by 3% and 5%, respectively; at higher penetration rates (40% and 60%), RL+G offers the highest fuel economy improvements at 10% and 25%, respectively.

V. CONCLUSION AND DISCUSSION

In this project, we incorporate real-world driving profiles into two mixed traffic environments, Ring and Bottleneck, with the goal to enhance the safety, efficiency, and stability of mixed traffic through RVs. For this purpose, we develop EnduRL with two classes of RVs: *safety + stability* RV and *efficiency* RV. Both types of RVs leverage the congestion stage classifier to optimize their objectives. In Ring, our RVs are able to increase the time to collision (TTC) above the critical threshold of 4 s, reduce the deceleration rate avoid a crash (DRAC) up to 54%, and increase the throughput up to 54%. This is achieved by dampening nearly all perturbations and maintaining the acceleration variation as low as 0.02 *m/s²*. In Bottleneck, our RVs are able to maintain the highest safety measure, i.e., TTC > 4 s at 5% penetration rate and DRAC < 2.42 at all penetration rates; and the highest throughput, i.e., improvements by up to 41%. Importantly, our RVs demonstrate a capacity to generalize and enhance performance when encountering more intricate dynamics,

such as merges and cut-ins, within the zipper lanes.

While we acknowledge that passenger comfort alongside theoretical analyses such as stabilizability, controllability, and reachability [48], [49] are significant and provide valuable insights, they fall beyond the scope of this work. We focus on directly measurable metrics, including enhancements in stability and reductions in oscillations. These not only offer quantitative assessments of performance but also directly correlate with other desired properties such as comfort. Given communicating global traffic states using vehicle-to-vehicle or vehicle-to-infrastructure techniques is in nascent stages and may require costly upgrades to current traffic infrastructure, our approach offers a practical solution because it solely relies on observations of individual RVs. The limitations of our work include not using additional traffic features such as lane-changing and heterogeneous vehicle types, and assuming perfect sensing without temporal delays. These issues serve as interesting future directions that we plan to explore.

REFERENCES

- [1] Colin Buchanan. *Traffic in Towns: A study of the long term problems of traffic in urban areas*. Routledge, 2015.
- [2] Hwasoo Yeo. *Asymmetric microscopic driving behavior theory*. University of California, Berkeley, 2008.
- [3] Hwasoo Yeo and Alexander Skabardonis. Understanding stop-and-go traffic in view of asymmetric traffic theory. In *Transportation and Traffic Theory 2009: Golden Jubilee: Papers selected for presentation at ISTTT18, a peer reviewed series since 1959*, pages 99–115. Springer, 2009.
- [4] Xuan Di and Rongye Shi. A survey on autonomous vehicle control in the era of mixed-autonomy: From physics-based to ai-guided driving policy learning. *Transportation research part C: emerging technologies*, 125:103008, 2021.
- [5] Michael Villarreal, Bibek Poudel, Jia Pan, and Weizi Li. Mixed traffic control and coordination from pixels. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [6] Michael Villarreal, Dawei Wang, Jia Pan, and Weizi Li. Analyzing emissions and energy efficiency in mixed traffic control at unsignalized intersections. In *IEEE Forum for Innovative Sustainable Transportation Systems (FISTS)*, 2024.
- [7] Michael Villarreal, Bibek Poudel, and Weizi Li. Can chatgpt enable its? the case of mixed traffic control via reinforcement learning. In *IEEE International Conference on Intelligent Transportation Systems (ITSC)*, pages 3749–3755, 2023.
- [8] Raphael E Stern, Shumo Cui, Maria Laura Delle Monache, Rahul Bhadani, Matt Bunting, Miles Churchill, Nathaniel Hamilton, Hannah Pohlmann, Fangyu Wu, Benedetto Piccoli, et al. Dissipation of stop-and-go waves via control of autonomous vehicles: Field experiments. *Transportation Research Part C: Emerging Technologies*, 89:205–221, 2018.
- [9] Berthold KP Horn. Suppressing traffic flow instabilities. In *16th International IEEE Conference on Intelligent Transportation Systems (ITSC 2013)*, pages 13–20. IEEE, 2013.
- [10] Rajesh Rajamani. *Vehicle dynamics and control*. Springer Science & Business Media, 2011.
- [11] Cathy Wu, Abdul Rahman Kreidieh, Kanaad Parvate, Eugene Vinitzky, and Alexandre M Bayen. Flow: A modular learning framework for mixed autonomy traffic. *IEEE Transactions on Robotics*, 38(2):1270–1286, 2021.
- [12] Dawei Wang, Weizi Li, Lei Zhu, and Jia Pan. Learning to control and coordinate mixed traffic through robot vehicles at complex and unsignalized intersections. *arXiv preprint arXiv:2301.05294*, 2023.
- [13] Mustafa Yildirim, Sajjad Mozaffari, Luc McCutcheon, Mehrdad Dianati, Alireza Tamaddon-Nezhad, and Saber Fallah. Prediction based decision making for autonomous highway driving. In *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pages 138–145. IEEE, 2022.
- [14] Eugene Vinitzky, Kanaad Parvate, Aboudy Kreidieh, Cathy Wu, and Alexandre Bayen. Lagrangian control through deep-rl: Applications to bottleneck decongestion. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 759–765. IEEE, 2018.
- [15] Dawei Wang, Weizi Li, and Jia Pan. Large-scale mixed traffic control using dynamic vehicle routing and privacy-preserving crowdsourcing. *IEEE Internet of Things Journal*, 11(2):1981–1989, 2024.
- [16] 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.
- [17] Martin Treiber and Arne Kesting. Traffic flow dynamics. *Traffic Flow Dynamics: Data, Models and Simulation*. Springer-Verlag Berlin Heidelberg, pages 983–1000, 2013.
- [18] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000.
- [19] Martin Treiber and Arne Kesting. The intelligent driver model with stochasticity-new insights into traffic flow oscillations. *Transportation research procedia*, 23:174–187, 2017.
- [20] Sumo driver state. Accessed: 2023-09-02.
- [21] Derek Gloudemans, Yanbing Wang, Junyi Ji, Gergely Zachar, William Barbour, Eric Hall, Meredith Cebelak, Lee Smith, and Daniel B Work. I-24 motion: An instrument for freeway traffic science. *Transportation Research Part C: Emerging Technologies*, 155:104311, 2023.
- [22] Bryan Higgs and Montasir Abbas. A two-step segmentation algorithm for behavioral clustering of naturalistic driving styles. In *16th International IEEE Conference on Intelligent Transportation Systems (ITSC 2013)*, pages 857–862. IEEE, 2013.
- [23] Guanghan Peng, Hongdi He, and Wei-Zhen Lu. A new car-following model with the consideration of incorporating timid and aggressive driving behaviors. *Physica A: Statistical Mechanics and its Applications*, 442:197–202, 2016.
- [24] Saleh Albeaik, Alexandre Bayen, Maria Teresa Chiri, Xiaoqian Gong, Amaury Hayat, Nicolas Kardous, Alexander Keimer, Sean T McQuade, Benedetto Piccoli, and Yiling You. Limitations and improvements of the intelligent driver model (idm). *SIAM Journal on Applied Dynamical Systems*, 21(3):1862–1892, 2022.
- [25] Cathy Wu, Aboudy Kreidieh, Kanaad Parvate, Eugene Vinitzky, and Alexandre M Bayen. Flow: Architecture and benchmarking for reinforcement learning in traffic control. *arXiv preprint arXiv:1710.05465*, 10, 2017.
- [26] Abdul Rahman Kreidieh, Cathy Wu, and Alexandre M Bayen. Dissipating stop-and-go waves in closed and open networks via deep reinforcement learning. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 1475–1480. IEEE, 2018.
- [27] Lokesh Chandra Das and Myounggyu Won. Saint-acc: Safety-aware intelligent adaptive cruise control for autonomous vehicles using deep reinforcement learning. In *International Conference on Machine Learning*, pages 2445–2455. PMLR, 2021.
- [28] Fang-Chieh Chou, Alben Rome Bagabaldo, and Alexandre M Bayen. The lord of the ring road: a review and evaluation of autonomous control policies for traffic in a ring road. *ACM Transactions on Cyber-Physical Systems (TCPS)*, 6(1):1–25, 2022.
- [29] Eugene Vinitzky, Aboudy Kreidieh, Luc Le Flem, Nishant Kheterpal, Kathy Jang, Cathy Wu, Richard Liaw, Eric Liang, and Alexandre M Bayen. Benchmarks for reinforcement learning in mixed-autonomy traffic. In *Conference on robot learning*, pages 399–409. PMLR, 2018.
- [30] Mayuri Sridhar and Cathy Wu. Piecewise constant policies for human-compatible congestion mitigation. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 2499–2505. IEEE, 2021.
- [31] Yuki Sugiyama, Minoru Fukui, Macoto Kikuchi, Katsuya Hasebe, Akihiro Nakayama, Katsuhiro Nishinari, Shin-ichi Tadaki, and Satoshi Yukawa. Traffic jams without bottlenecks—experimental evidence for the physical mechanism of the formation of a jam. *New journal of physics*, 10(3):033001, 2008.
- [32] Douglas Gettman and Larry Head. Surrogate safety measures from traffic simulation models. *Transportation Research Record*, 1840(1):104–115, 2003.
- [33] Sehyun Tak, Sunghoon Kim, and Hwasoo Yeo. A study on the traffic predictive cruise control strategy with downstream traffic information. *IEEE Transactions on Intelligent Transportation Systems*, 17(7):1932–1943, 2016.
- [34] Arne Kesting and Martin Treiber. Calibrating car-following models by using trajectory data: Methodological study. *Transportation Research Record*, 2088(1):148–156, 2008.
- [35] Li Li, Xiquan Micheal Chen, and Lei Zhang. A global optimization algorithm for trajectory data based car-following model calibration. *Transportation Research Part C: Emerging Technologies*, 68:311–332, 2016.
- [36] MN Sharath and Nagendra R Velaga. Enhanced intelligent driver model for two-dimensional motion planning in mixed traffic. *Transportation Research Part C: Emerging Technologies*, 120:102780, 2020.
- [37] Meixin Zhu, Xuesong Wang, Andrew Tarko, et al. Modeling car-following behavior on urban expressways in shanghai: A naturalistic driving study. *Transportation research part C: emerging technologies*, 93:425–445, 2018.
- [38] Xiao Wang, Rui Jiang, Li Li, Yilun Lin, Xinhui Zheng, and Fei-Yue Wang. Capturing car-following behaviors by deep learning. *IEEE Transactions on Intelligent Transportation Systems*, 19(3):910–920, 2017.
- [39] Meixin Zhu, Yin Hai Wang, Ziyuan Pu, Jingyun Hu, Xuesong Wang, and Ruimin Ke. Safe, efficient, and comfortable velocity control based on reinforcement learning for autonomous driving. *Transportation Research Part C: Emerging Technologies*, 117:102662, 2020.

- [40] Meixin Zhu, Xuesong Wang, and Yin Hai Wang. Human-like autonomous car-following model with deep reinforcement learning. *Transportation research part C: emerging technologies*, 97:348–368, 2018.
- [41] Tianyu Shi, Yifei Ai, Omar ElSamadisy, and Bahar Abdulhai. Bilateral deep reinforcement learning approach for better-than-human car following model. *arXiv preprint arXiv:2203.04749*, 2022.
- [42] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [43] Meead Saberi and Hani S Mahmassani. Hysteresis and capacity drop phenomena in freeway networks: Empirical characterization and interpretation. *Transportation research record*, 2391(1):44–55, 2013.
- [44] Katja Vogel. A comparison of headway and time to collision as safety indicators. *Accident analysis & prevention*, 35(3):427–433, 2003.
- [45] Dale F Cooper and N Ferguson. Traffic studies at t-junctions. 2. a conflict simulation record. *Traffic Engineering & Control*, 17(Analytic), 1976.
- [46] Peter De Haan and Mario Keller. Modelling fuel consumption and pollutant emissions based on real-world driving patterns: the hbeqa approach. *International journal of environment and pollution*, 22(3):240–258, 2004.
- [47] TJ Ayres, L Li, David Schleunig, and D Young. Preferred time-headway of highway drivers. In *ITSC 2001. 2001 IEEE Intelligent Transportation Systems. Proceedings (Cat. No. 01TH8585)*, pages 826–829. IEEE, 2001.
- [48] Yang Zheng, Jiawei Wang, and Keqiang Li. Smoothing traffic flow via control of autonomous vehicles. *IEEE Internet of Things Journal*, 7(5):3882–3896, 2020.
- [49] Jiawei Wang, Yang Zheng, Qing Xu, Jianqiang Wang, and Keqiang Li. Controllability analysis and optimal control of mixed traffic flow with human-driven and autonomous vehicles. *IEEE Transactions on Intelligent Transportation Systems*, 22(12):7445–7459, 2020.

APPENDIX

I. CAR FOLLOWING FILTER

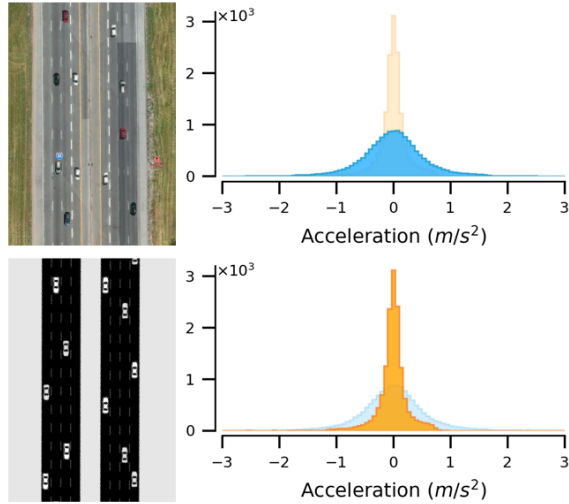


Fig. 3: Instantaneous accelerations observed during car-following behaviors at densities $[70, 150] \text{ veh/km}$. TOP: Real-world data from the I-24 MOTION dataset reveals a distribution having long tails extending to $[-3, 3] \text{ m/s}^2$. BOTTOM: IDM (in simulation) produces accelerations mostly within $[-0.5, 0.5] \text{ m/s}^2$, indicating much ‘timid’ driving behaviors than the real world.

We analyze the I-24 MOTION dataset [21] with study length = 6.75 km and study time = 4 h . The dataset contains various vehicle types such as semi-trailers, mid-sized trucks, motorbikes, and cars under different traffic conditions such as approaching standing traffic, lane changing, and free flow. To

extract car-following trajectories, we select data points that meet the following criteria:

- Ego car is following another car, i.e., has a leader.
- Leader and ego cars are in the same lane $\geq 5 \text{ s}$.
- Ego car’s speed is greater than 10% of the speed limit, i.e., not approaching stationary traffic.
- Ego car’s space headway is less than 124 m , applying 4 s rule at the speed limit to avoid free flow conditions.

II. MODEL-BASED ROBOT VEHICLES

Bilateral Control Module (BCM): BCM [9] uses information about both follower and leader vehicles to obtain a linear model whose acceleration is given by:

$$a = k_d \cdot \Delta_d + k_v \cdot (\Delta v_l - \Delta v_f) + k_c \cdot (v_{des} - v), \quad (1)$$

where Δ_d , Δv_l , Δv_f , v_{des} , and v , represent the difference in distance to the leader compared to the distance to the follower, the difference in velocity to the leader, the difference in velocity to the follower, the set desired velocity, and the current velocity of the vehicle, respectively. $k_d = 1$, $k_v = 1$, and $k_c = 1$ are gain parameters.

Linear Adaptive Cruise Control (LACC): LACC is an improvement on existing cruise control systems that allows vehicles to maintain a safe distance or speed without communication. The constant time-headway model by Rajamani [10] employs a first-order differential equation for approximation. The control acceleration at time t is given by:

$$a_t = \left(1 - \frac{\Delta t}{\tau}\right) \cdot a_{(t-1)} + \frac{\Delta t}{\tau} a_{cmd,(t-1)}, \quad (2)$$

$$a_{cmd} = k_1 \cdot e_x + k_2 \cdot \Delta v_l, \quad (3)$$

$$e_x = s - h \cdot v, \quad (4)$$

where $k_1 = 0.3$ and $k_2 = 0.4$ are design parameters, e_x is the gap error, s is the space headway, Δv_l is the velocity difference to the leader, $h = 1$ is the desired time gap, Δt is the control time-step, and $\tau = 0.1$ is the time lag of the system.

III. HEURISTIC-BASED ROBOT VEHICLES

FollowerStopper (FS): FS [8] is an RV that travels at a fixed command velocity (target) under safe conditions but when required, slightly lowers the target velocity, opening up a gap to the vehicle ahead. This allows it to dampen oscillations and brake smoothly when needed. The command velocity is given by:

$$v_{cmd} = \begin{cases} 0, & \text{if } \Delta x \leq \Delta x_1 \\ v \frac{\Delta x - \Delta x_1}{\Delta x_2 - \Delta x_1}, & \text{if } \Delta x_1 < \Delta x \leq \Delta x_2 \\ v + (U - v) \frac{\Delta x - \Delta x_2}{\Delta x_3 - \Delta x_2}, & \text{if } \Delta x_2 < \Delta x \leq \Delta x_3 \\ U, & \text{if } \Delta x_3 < \Delta x \end{cases} \quad (5)$$

where $v = \min(\max(v_{lead}, 0), U)$ is the speed of the leader vehicle, Δx is the headway of the RV, and U is the desired velocity. The thresholds $(\Delta x_1, \Delta x_2, \Delta x_3)$ are defined as

$$\Delta x_k = \Delta x_k^0 + \frac{1}{2d_k} (\Delta v_-)^2, \quad k = 1, 2, 3. \quad (6)$$

The model parameters Δx_k^0 , Δv_- , and d_k determine the spacing between vehicles and the RV’s responsiveness to changes in velocity.

Proportional-integral with saturation (PIwS): PIwS [8] estimates the desired average velocity (U) of the vehicles in the network using its historical average velocity. The PIwS RV calculates the target velocity as

$$v_{target} = U + v_{catch} \times \min \left(\max \left(\frac{\Delta x - g_l}{g_u - g_l}, 0 \right), 1 \right), \quad (7)$$

which is used to calculate the command velocity at $t + 1$ as

$$v_{cmd}^{t+1} = \beta_t (\alpha_t v_{target}^t + (1 - \alpha_t) v_{lead}^t) + (1 - \beta_t) v_{cmd}^t, \quad (8)$$

where v_{catch} is the catch-up velocity—a velocity higher than the average velocity allows the RV to catch up with its leader, Δx is the difference in position between the RV and its leader, g_l and g_u represent the lower and upper threshold distance, respectively; α_t and β_t represent the weight factors for target velocity v_{target} and command velocity v_{cmd} , respectively. Finally, v_{lead} represents the velocity of the leader vehicle.

IV. REINFORCEMENT LEARNING (RL) BENCHMARKS

To benchmark with other RL techniques, we reproduce their original policies by following the provided experiment parameters and closely matching the performance. Specifically, to obtain RL policy with only local observations (RL+L), we follow Wu et al. [11]; to obtain RL with global observations (RL+G), we follow Vinitzky et al. [14]. Our reproduced RL+L achieves the performance within 1% error (measured with stabilization time and average velocity during stabilization) of the original work. Whereas for RL+G, our reproduction achieves the performance within 3% error (measured with outflow) of the original work. The precise implementations of the other RL methods validate our benchmarking experiments.

V. CONGESTION STAGE CLASSIFIER (CSC)

One CSC each is trained in Ring and Bottleneck with independent datasets collected in the two environments. For each RV, the position and velocity of all vehicles in its local zone (set to 50 m) are collected. Fig. 4 shows the K-means clustering of collected data over all six classes (‘Forming’, ‘Leaving’, ‘Congested’, ‘Free flow’, ‘Undefined’, and ‘No Vehicle’) in both environments. As the data collected is sequential in nature and CSC predictions are made a number of time-steps into the future, a time offset of 10 time-steps is chosen to balance usefulness and accuracy (to illustrate, a prediction of congestion stage 100 time-steps in the future would be very useful, however, not very accurate; whereas a prediction of congestion state 1 time-step into the future can be very accurate but not very useful).

After windowing, the dataset includes instances where the congestion stage changes from t to $t + 10$, as well as instances where the congestion stage remains the same over the time window. To train CSC, we sample data to ensure a balanced

| Category | Parameter | Value |
|-----------------------|------------------------------|----------------|
| Ring Simulation | Time Step (Δt) | 0.1 |
| | Simulation Horizon (T) | 4500 |
| | Warmup Time-steps | 2500 |
| | Speed Limit (m/s) | 30 |
| | Initial Speed (m/s) | 0 |
| Bottleneck Simulation | Time Step (Δt) | 0.5 |
| | Simulation Horizon (T) | 1300 |
| | Warmup Time-steps | 100 |
| | Speed Limit (m/s) | 17 |
| | Initial Speed (m/s) | 6 |
| PPO Algorithm | Inflow Rate (veh/hr) | 3600 |
| | Learning Rate (α) | 0.00005 |
| | Discount Factor (γ) | 0.999 |
| | GAE Estimation (λ) | 0.97 |
| | KL Divergence Target | 0.02 |
| | Entropy Coefficient Initial | 0.1 |
| | Entropy Coefficient Final | 0.01 |
| | Value Function Clip Param | 20 |
| | SGD Iterations | 2 |
| | Ring CSC | Neural Network |
| Batch Size | | 32 |
| Learning Rate | | 0.01 |
| Bottleneck CSC | Epochs | 50 |
| | Neural Network | 32, 16, 16 |
| | Batch Size | 256 |
| Policy Networks | Learning Rate | 0.001 |
| | Epochs | 100 |
| | Our leader RV in Ring | 64, 32, 16 |
| | Our follower RV in Ring | 64, 32, 16 |
| | Our RV in Bottleneck | 32, 16, 8 |
| | RL+L | 32, 32, 32 |
| | RL+G | 256, 256 |

TABLE III: Detailed experiment parameters. We show the simulation parameters of Ring and Bottleneck, as well as the parameters of Proximal Policy Optimization (PPO) and Congestion Stage Classifier (CSC). The hidden layer dimensions of various policy networks are also shown.

representation of transition/non-transition instances as well as instances containing all six classes. Worth noting, the ‘No vehicle’ class presents a unique challenge. The collected data may contain instances changing from ‘No vehicle’ to another class after the 10 time-steps. However, based on the input corresponding to ‘No vehicle’ at t , we cannot predict the congestion stage at $t + 10$. Consequently, we discard data points where the ‘No Vehicle’ class transitions to another class after 10 time-steps. We replace this discarded data with synthetic examples that simulate various scenarios for the RV’s position and velocity without leader vehicles.

The accuracy of the trained CSC is 95.5% in Ring and 85.2% in Bottleneck. The confusion matrix is shown in Fig. 5. CSC only observes downstream HVs in the same lane. Thus, when facing zipper lanes of Bottleneck (where traffic merges from adjacent lanes), the CSC cannot anticipate the merging traffic, resulting in lower accuracy in Bottleneck.

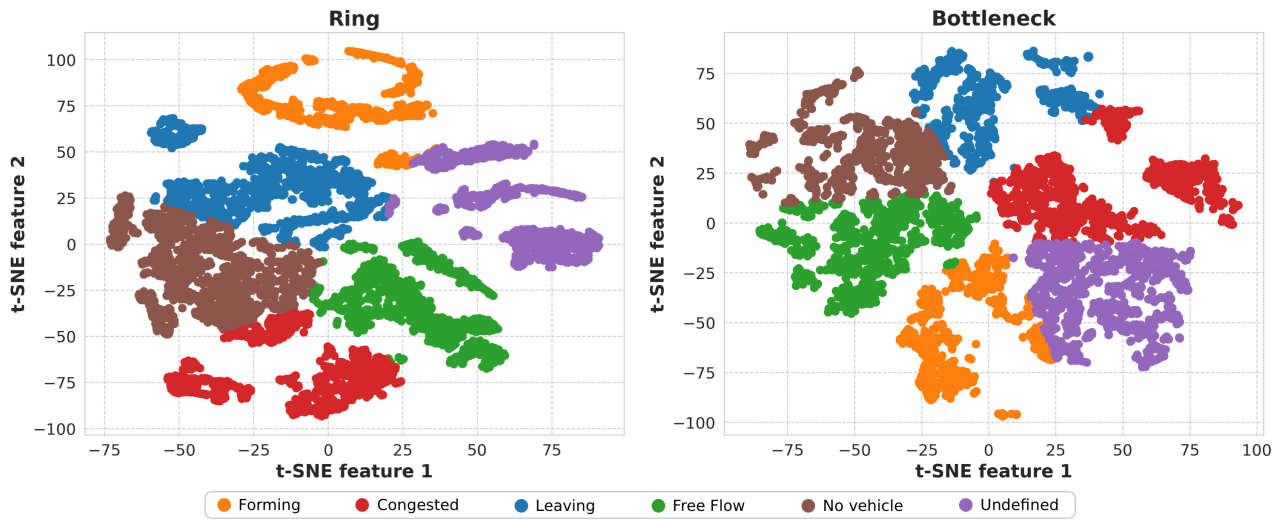


Fig. 4: The results of applying K-means clustering with t-SNE on a subset of CSC training data. LEFT: In Ring, the clusters are spread out, suggesting that the data is easily classifiable. RIGHT: In Bottleneck, overlapping clusters indicate that more complex interactions exist among the congestion stages, possibly due to the presence of zipper lanes causing vehicles abruptly merge.

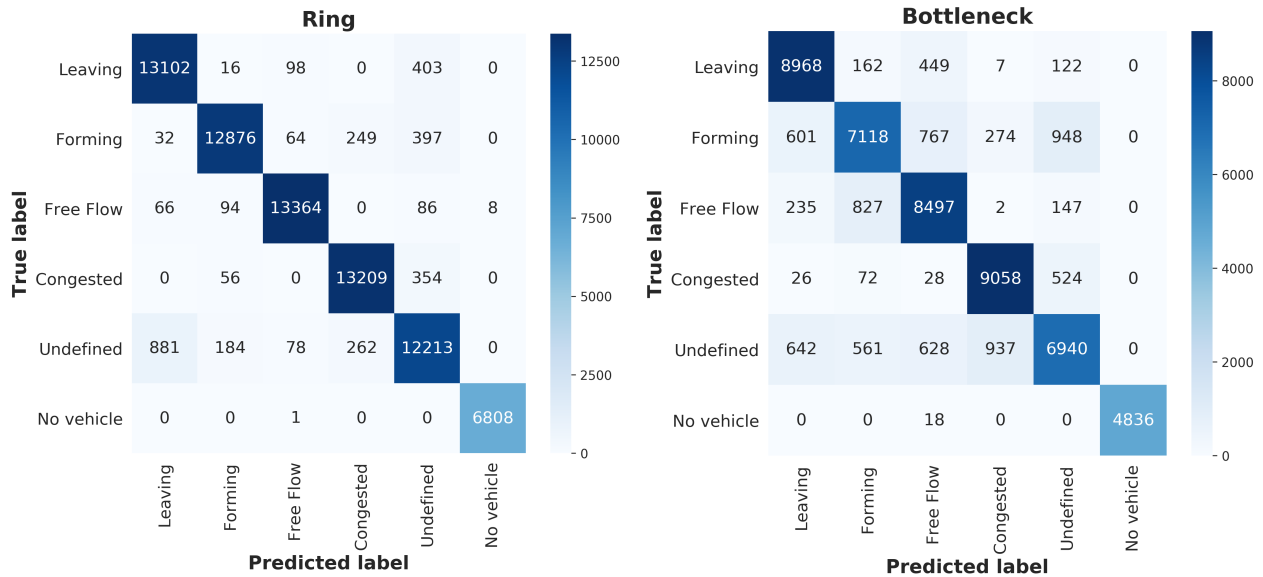


Fig. 5: Confusion matrix of a trained CSC in Ring (LEFT) and Bottleneck (RIGHT) on the validation set.

The CSC training parameters are provided in Table III.