

Can ChatGPT Enable ITS? The Case of Mixed Traffic Control via Reinforcement Learning

Michael Villarreal, Bibek Poudel, Weizi Li

Abstract—The surge in Reinforcement Learning (RL) applications in Intelligent Transportation Systems (ITS) has contributed to its growth as well as highlighted key challenges. However, defining objectives of RL agents in traffic control and management tasks, as well as aligning policies with these goals through an effective formulation of Markov Decision Process (MDP), can be challenging and often require domain experts in both RL and ITS. Recent advancements in Large Language Models (LLMs) such as GPT-4 highlight their broad, general knowledge, reasoning capabilities, and commonsense priors across various domains. In this work, we conduct a large-scale user study involving 70 participants to investigate whether novices can leverage ChatGPT to solve complex mixed traffic control problems. The participants’ task is to develop the state space and reward function for three RL mixed traffic control environments, including ring road, bottleneck, and intersection. We find ChatGPT has mixed results. For intersection and bottleneck, ChatGPT increases number of successful policies by 150% and 136% compared to solely beginner capabilities, with some of them even outperforming experts. However, ChatGPT does not provide consistent improvements across all scenarios.

I. INTRODUCTION

Large Language Models (LLMs) represent a significant advancement in artificial intelligence. Their usefulness as a general-purpose tool is anticipated to have a profound societal impact. Experiments with LLMs such as GPT-4 [1], [2], LLaMA [3], and PALM2 [4] demonstrate their strong reasoning and common sense abilities across domains such as math, science, and coding [1], [2], [3], [4], [5], [6]. To achieve this success, LLMs leverage Reinforcement Learning from Human Feedback (RLHF) [7], [8] to solve the alignment problem, i.e., follow user intent by fine-tuning them on human feedback. While there is a strong focus on improving LLMs through Reinforcement Learning (RL), leveraging LLMs to assist RL problems is in nascent stages.

The RL framework is inherently challenging because it demands an understanding of Markov Decision Process (MDPs). Any RL problem must be formulated to a MDP through designing the state, action, reward, and discount-factor among other components [9]. This is usually a tedious task that requires numerous experimentation and careful analysis, often by a domain expert. Specifically, in research areas such as Intelligent Transportation Systems (ITS), it is challenging for novices to understand and design effective MDPs. Moreover, no standard technique exists for creating general-purpose MDPs that will work across environments and satisfy various goals.

Michael Villarreal, Bibek Poudel, and Weizi Li are with the Min H. Kao Department of Electrical Engineering and Computer Science, University of Tennessee, Knoxville, TN 37996, USA {tvillarr, bpoudel13}@vols.utk.edu, weizili@utk.edu

Over the last decade, RL has been adopted to address complex control problems in ITS such as traffic management and autonomous driving [10], [11]. As autonomous agents, including vehicles and traffic lights, become more prevalent [12] in ITS, they introduce new challenges and opportunities. One emerging topic is mixed traffic control that uses RL-empowered robot vehicles (RVs) to regulate human-driven vehicles (HVs), thus improving the overall traffic flow [13], [14], [15]. This surge in research interest has attracted a broader audience to participate in the topic, resulting in a growing demand for creative decision-making and control strategies enabled by RL. However, the initial technical barrier, specifically formulating MDP components to align with a control strategy, poses a challenge. LLMs, with their broad knowledge, commonsense priors, and creative capacity, show promise in reducing these barriers and simplifying the process.

In this project, we explore whether ChatGPT (with GPT-4 backend) can assist non-experts in ITS research to solve mixed traffic control problems. We conduct a large-scale user study with 70 participants who have no prior experience in ITS research. The participants are tasked to develop MDP components (state and reward) for three mixed traffic scenarios: a ring road, an intersection, and a bottleneck as shown in Fig. 1. We split participants into a control group, where participants attempt to solve problems solely based on common sense and prior knowledge, and a study group, where participants can prompt ChatGPT unrestricted. Participants are provided a manuscript with a general overview of RL, a few examples of MDP components, a reference bank of metrics related to traffic, and images and descriptions of the mixed traffic control environments. The formulated MDPs are then used to train policies and evaluate performance. From the study, we find that:

- 1) In the intersection environment, using ChatGPT can lead to a performance better than an expert.
- 2) In the bottleneck and intersection environments, using ChatGPT results in an increase in number of successful policies by 150% and 136%, respectively.
- 3) ChatGPT’s creativity enables a 363% increase in utilization of new metrics, although the use of these new metrics does not always result in a successful policy.
- 4) In the ring environment, ChatGPT’s help does not increase the policy success rate.

The wide range of results observed, from no success to outperforming the expert, across the mixed traffic control environments indicates that the effectiveness of utilizing

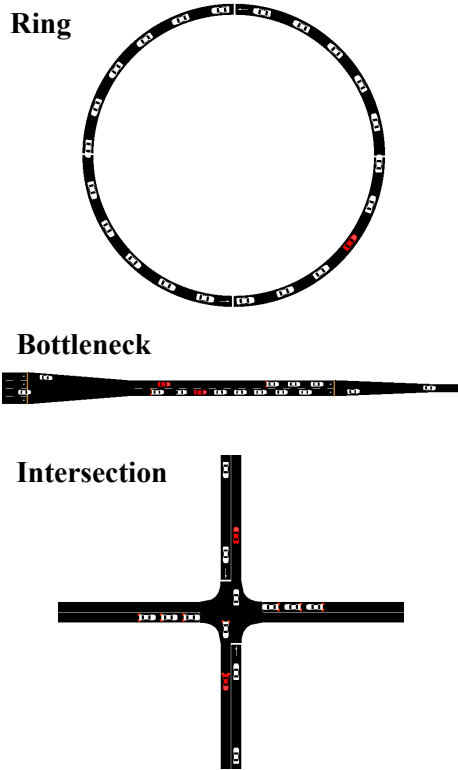


Fig. 1. Three mixed traffic control environments [13] (a deep reinforcement learning framework for traffic management), Ring, Bottleneck, and Intersection, are provided to the study participants. Robot vehicles are red and are controlled by learnt RL policies. Human-driven vehicles are white and are modeled by the Intelligent Driver Model [16].

ChatGPT in ITS depends on the complexity of the specific problem, as well as the extent and manner in which ChatGPT is utilized.

II. RELATED WORK

Recent works have studied how Large Language Models (LLMs) can be leveraged in a variety of RL tasks. In robotics, LLMs are utilized in planning and navigation by representing the robotic agent’s input (state) in natural language, additionally incorporating the input with visual or raw sensor data for grounding. This approach demonstrates high data efficiency and generalization to unseen environments [17], [18]. Similarly, using LLMs to guide exploration during training can also help achieve greater sample efficiency [19].

Limited research has been found in minimizing the human effort required to generate an effective policy. Some works use LLMs to substitute components of MDPs, such as making LLM a surrogate or proxy reward function. By using the in-context learning capability of LLMs, hard-to-specify reward functions (e.g., versatility, fairness) have been attempted [20]. Meanwhile, other studies have used LLMs to replace a policy entirely [21], [22]. However, to the best of our knowledge, no previous work has directly evaluated the effectiveness of LLMs (with and without them) as a tool to reduce the effort in designing MDP components by novices, comparing the obtained policies with those from experts.

III. PRELIMINARIES

In the following, we introduce the formulation of mixed traffic control as reinforcement learning (RL) tasks and discuss the corresponding test environments.

A. Reinforcement Learning

We model mixed traffic control as a Partially Observable Markov Decision Process (POMDP) represented by a tuple $(S, A, P, R, p_0, \gamma, T, \Omega, O)$ where S is the state space; A is the action space; $P(s'|s, a)$ is the transition probability function; R is the reward function; p_0 is the initial state distribution; $\gamma \in (0, 1]$ is the discount factor; T is the episode length (horizon); Ω is the observation space; and O is the probability distribution of retrieving an observation $\omega \in \Omega$ from a state $s \in S$. At each timestep $t \in [1, T]$, a robot vehicle (RV) uses its policy $\pi_\theta(a_t|s_t)$ to take an action $a_t \in A$, given the state $s_t \in S$. The RV’s environment provides feedback from taking action a_t by calculating a reward r_t and transitioning the agent into the next state s_{t+1} . The RV’s goal is to learn a policy π_θ that maximizes the discounted sum of rewards, i.e., return, $R_t = \sum_{i=t}^T \gamma^{i-t} r_i$. Proximal Policy Optimization [23] is used to learn π_θ .

B. Mixed Traffic Control Environments

1) *Ring*: The ring environment (shown in Fig. 1 top) consists of a single-lane circular road network and 22 vehicles (21 HVs and one RV). It simulates how perturbations due to imperfections in human driving behavior can amplify and propagate, leading to an eventual standstill for some vehicles. This situation, known as ‘stop-and-go traffic’, acts as a wave that propagates continually through the ring, opposite the direction of travel. The RV’s goal is to prevent the formation of these waves. Ring is a widely used benchmark in traffic control [24]. An expert’s [13] state space (major metrics given in Table I) is:

$$s = \left\{ \frac{v_{RV}}{v_{\max}}, \frac{v_{\text{lead}} - v_{RV}}{v_{\max}}, f(x_{\text{lead}} - x_{RV}) \right\}. \quad (1)$$

The difference in x_{lead} and x_{RV} is passed through a normalization function f . An expert’s reward function encourages high average velocity and low control actions (acceleration) though a weighted combination given by:

$$r = \frac{1}{n} \sum_i v_i - \alpha * |a_{RV}|, \quad (2)$$

where $n = 22$ and α is chosen empirically.

2) *Bottleneck*: The bottleneck environment (shown in Fig. 1 middle) simulates vehicles experiencing a capacity drop [25] where a road’s outflow significantly decreases after the road’s inflow surpasses a threshold. The RVs’ goal is to improve outflow. Bottleneck represents a bridge with lanes decreasing from $4 \times l$ to $2 \times l$ to l (where l is a scaling factor and is one for our work). The RV penetration rate is 10%. An expert’s [26] state space is:

$$s = \{\bar{X}_{HV}, \bar{V}_{HV}, \bar{X}_{RV}, \bar{V}_{RV}, o_{20}\}, \quad (3)$$

Type	Metrics
Vehicle	Position (x), Velocity (v), Acceleration (a), Gap to leader/follower (g_{lead}/g_{fol}), Fuel consumption (fc)
Road Environment	Outflow over last i seconds (o_i), Density (d), Speed limit (v_{max}), Average speed (\bar{v}), Average acceleration (\bar{a}), Minimum speed (v_{min}), Average time standstill (\bar{ss}_τ)
Driving Behavior	Distance/Time to complete stop (st_{dist}/st_τ), Time to collide (ttc), Jerk (j)

TABLE I

THE BANK OF METRICS PROVIDED TO PARTICIPANTS. CAPITALIZATION OF A METRIC INDICATES ITS VECTOR.

where the mean positions and velocities of both vehicle types are considered across user-defined segments of the road network. An expert’s reward function rewards increasing bottleneck outflow:

$$r = o_{10}. \quad (4)$$

3) *Intersection*: The intersection environment (shown in Fig. 1 bottom) represents an unsignalized intersection where east/westbound traffic flow is less than north/southbound traffic flow. This flow discrepancy leads to east/westbound traffic queues as crossing the intersection would be unsafe otherwise. RVs drive in the north/south directions with a 20% penetration rate. The RVs’ objective is minimizing east/west queues and increased average vehicle velocity. An expert’s [26] state space is:

$$s = \{V_{all}, I_{all}, E_{all}, D_{edge}, \bar{V}_{edge}\}, \quad (5)$$

where I_{all} is all vehicle’s distance to intersection, E_{all} is all vehicle’s edge number, D_{edge} is density of each edge, and \bar{V}_{edge} is average vehicle velocity of each edge. There are eight edges for each direction on both sides of the intersection. An expert’s reward function penalizes vehicle delay and vehicle standstills in traffic:

$$r = -\frac{t * \sum((V_{max} - V_{all})/V_{max})}{n + eps} - (gain * ss_n), \quad (6)$$

where t is current timestep, n is number of vehicles, eps prevents zero division, $gain$ is 0.2, and ss_n is the number of standstill vehicles.

IV. USER STUDY

A. Participants and Cohorts

We recruit 70 graduate students as participants (54 male, 16 female). The male/female participants have a median age of 23/22 and average 4.93/4.26 years of driving experience. Participants come from a machine learning or artificial intelligence course with minimal coverage of reinforcement learning (RL) taught by the same professor. At the study’s date, artificial intelligence participants have slightly more RL knowledge; however, we find no discernible difference between the two courses’ results. We attribute this to the manuscript and ChatGPT’s assistance.

The 70 participants are split into two cohorts: a control group (38 participants) that only use prior non-expert knowledge and the manuscript, and a study group (32 participants) that additionally has access to ChatGPT [27], [1]. The

	Control	Study	Control	Study	Control	Study
Valid	25	28	24	27	21	24
Invalid	13	4	14	5	17	8
Total	38	32	38	32	38	32

TABLE II

NUMBER OF VALID AND INVALID ANSWERS FOR EACH TRAFFIC ENVIRONMENT FOR THE CONTROL AND STUDY GROUPS.

control group helps assess whether ChatGPT can assist non-experts in solving traffic problems by providing baseline capabilities. We find the control group heavily uses collision-prevention metrics, including g_{lead}/g_{fol} , ttc , or st_{dist}/st_τ , with consistent use of x_{RV} , v_{RV} , and a_{RV} across all environments. Additionally, the control group tends to include objective-oriented metrics such as o or \bar{ss}_τ for bottleneck and metrics about queue length/time standstill in the east/west directions for intersection. These metrics are straightforward and intuitive, aligning with our expectations for non-experts in ITS research.

We poll the study group’s ChatGPT [27] use frequency with the following descending frequencies (number of participant selections): daily (2), several times a week (13), once a week (6), several times a month (5), once a month or less (2), and never (3).

B. ChatGPT Setup

The study group participants use GPT-4 with 8k context length and temperature = 0.7, a static model with no additional fine-tuning. The March 14th, 2023 version of GPT-4 is used. To simulate the practical use of ChatGPT, a fresh chat interface is provided to participants using TypingMind [28]. The participants use the same conversation with ChatGPT for the study to enable in-context learning.

C. Manuscript

We develop two versions of manuscripts¹ (for different cohorts) with the following four sections. Study duration is 95 minutes. Within the 95 minutes, participants are required to spend the first 25 minutes reading the first three sections outlined below, and then given the remaining time to look over the entire manuscript. This is to ensure a higher probability of participants reading over the provided material, potentially leading to better quality state spaces and reward functions for the mixed traffic environments.

1) General Instructions and Background Questions:

The general instructions outline participant conduct for the study’s duration. and warn participants that rule violations comprise their answers’ integrity and result in their dismissal. Participants complete the background questions discussed in Sec. IV-A.

2) Reinforcement Learning Overview:

The RL overview section provides a brief explanation of the state space and reward function RL components. We provide transportation-related examples and explanations to ensure common understanding of RL across all participants.

¹<https://github.com/tmv11r1rl/its-study>

Novice A Prompting with ChatGPT Responses

Novice B Prompting with ChatGPT Responses

<p>Novice A: if a robot car is moving in a circular motion to avoid stop and go traffic is state space= (speed limit, gap to vehicle's front) is good state space</p> <p>ChatGPT: In the case of a robot car moving in a circular motion to avoid stop-and-go traffic, the suggested state space of (speed limit, gap to vehicle's front) may be insufficient to model the complete behavior required for safe navigation. . . .</p> <p style="text-align: center;">⋮</p> <p>Novice A: if a robot car is moving in a circular motion to avoid stop and go traffic what is a good reward function</p> <p>ChatGPT: When designing a reward function for a robot car that aims to avoid stop-and-go traffic while moving in a circular motion, the function should encourage safe and efficient navigation. . . .</p>	<p>Novice B: In reinforcement learning, for car driving in loop, what metric could best describe the state space. You are given following metrics: Position, speed, acceleration/deceleration, gap to vehicle in front, fuel consumption. You are allowed to use statistical computation like mean, median, mode, standard deviation, variance.</p> <p>ChatGPT: In a reinforcement learning problem for car driving in a loop, the goal is to find the best policy that maps the given state space to the actions to be performed. . . .</p> <p style="text-align: center;">⋮</p> <p>Novice B: how do you define the reward function for the above scenario with TTC?</p> <p>ChatGPT: Defining a reward function for the car driving scenario considering Time-To-Collision (TTC) requires balancing multiple objectives . . .</p>
--	--

Fig. 2. Example prompting by two novices, Novice A and Novice B, with ChatGPT response excerpts for the ring environment. The novices' prompts are given exactly. Only conversation portions are provided for brevity purposes. Novice A, unlike Novice B and most other participants, starts the conversation by asking if a certain state space is good to which ChatGPT responds the state space might be "insufficient." Novice B follows more closely participants' trend to completely rely on ChatGPT for the state space after providing context of the ring environment. Both novices heavily rely on ChatGPT for the reward function; however, Novice B prompts ChatGPT to include time-to-collision based on prior conversation with ChatGPT.

State Space, Novice

Reward Function, Novice + ChatGPT

<p>1. Give the state space as a tuple of metrics. If using metrics that are not in the bank of metrics, explain what they mean with the space provided.</p> <p>$s = (\text{Gap to vehicle's front}, \text{Gap to vehicles back}, \text{speed}, \text{acceleration}, \text{Position (robot)})$.</p> <p>Let, $G_f = \text{Gap to vehicle's front}$ } Can be measured using sensor. $G_b = \text{Gap to vehicles back}$ }</p> <p>→ Acceleration can be both positive & negative.</p> <p>Note! The state space is about Robot vehicle.</p>	<p>2. Give the reward function as equation that evaluates to a single numerical value. If using metrics that are not in the bank of metrics, explain what they mean in the space provided.</p> <p>$R = \text{maintain speed} + \text{penalize abrupt changes in velocity} + \text{maintain safe distance} + \text{maintain safe Time to collision (TTC)}$</p> <p>$R = \alpha * (1 - \text{Target speed} - v_{\text{robot}} / \text{target speed}) - \beta * \Delta v_{\text{robot}} - \gamma * \max(0, \text{min-TTW} + K * v_{\text{robot}} - \text{THW}) - \delta * (\max(0, \text{safe-TTC} - \text{TTC}_{\text{rear}})) + \max(0, \text{safe-TTC} - \text{TTC}_{\text{front}})$</p> <p>$\alpha, \beta, \gamma, \delta = \text{weight of each reward comping}$ $v_{\text{robot}} = \text{robot velocity}$ $\text{THW} = \text{time headway (time to reach front vehicle at current speed)}$</p>
--	---

Fig. 3. Example state space from novice (left) and reward function from novice with ChatGPT (right) for ring. The novice without ChatGPT's reward function heavily uses provided metrics, while using ChatGPT allows the reward function to contain complex, generated terms. ChatGPT provides sound reasoning for using the reward function metrics (see Sec. V-A for details).

3) *Answer Instructions and Bank of Metrics:* The answer instructions cover how their answers should be formatted with examples. The bank of metrics is provided in Table I. The participants can use any new metric with explanation.

4) *Mixed Traffic Environment Descriptions and Questions:* Each of the three traffic environments (shown in Fig. 1; details in Sec. III-B) receives a general description, problem explanation, and the mixed traffic objective of the RVs. The general description includes details such as the number of RVs present, general flow behavior of traffic, ratio of RVs to human-driven vehicles. We provide a supplementary video² demonstrating the mixed traffic environments. We ask the participants three questions per environment: the state space, the reward function, and briefly explain the rationale behind the reward function. We ask participants to explain their reward function rationale to encourage deeper understanding from the participants. The state spaces participants provide are observation spaces; however, we solely use "state space" to prevent additional complexity/avoid confusion.

²<https://youtu.be/qTqgfl76FAo>

V. RESULTS

Next, we analyze the participants' answers. Then, we briefly discuss experiment setup and the results in the three mixed traffic control environments.

A. Participant Answers Analysis

Responses	Ring		Bottleneck		Intersection	
	Control	Study	Control	Study	Control	Study
Valid	25	28	24	27	21	24
Invalid	13	4	14	5	17	8
Total	38	32	38	32	38	32

TABLE III
NUMBER OF VALID AND INVALID ANSWERS FOR EACH TRAFFIC ENVIRONMENT FOR THE CONTROL AND STUDY GROUPS.

We find that ChatGPT [27], [1] impacts provided state spaces/reward functions by significantly decreasing invalid answers over using prior knowledge. We consider an answer valid if it contains metrics provided in the metrics bank or have well-defined explanations. Table III presents the

number of valid/invalid answers for all three environments. On average, only 61% of answers from the control group are valid, while 82% are valid for the study group, a 21% increase. This illustrates ChatGPT’s capabilities in guiding participants to create valid state spaces/reward functions.

Another impact of ChatGPT is that the study group uses (in ring, bottleneck, intersection order) 35, 63, and 59 new metrics compared to the control group’s 8, 17, and 21 metrics that do not exist in the metrics bank. On average, this is a 363% increase in new metrics. This significant increase implies that ChatGPT can provide new perspectives to solving the mixed traffic problems.

We provide example prompting by two novices (labeled Novice A and Novice B) from the study group for the ring environment’s state space and reward function in Fig. 2. Novice A deviates from the general trend of participants by asking ChatGPT if a provided state space is good for the ring environment (to which ChatGPT correctly assumes the state space may be insufficient). The general trend of participants is showcased with Novice B where Novice B provides context for the ring environment, then asks what a state space might be for the environment. Additionally, both novices more closely follow participants’ general trend of heavily relying on ChatGPT by asking ChatGPT for a good reward function. Novice B demonstrates a common behavior with their reward function prompting by wanting to have a specific metric (for Novice B, time-to-collision) in the function.

We provide example state space and reward function in Fig. 3. The left image is a control group (Novice) state space, while the right image is a reward function from the study group (Novice + ChatGPT). The novice state space heavily uses existing metrics (a trend with the control group), while the ChatGPT reward shows the intricate metrics ChatGPT generates. ChatGPT also offers explanations that appear reasonable with the terms. For example, for the “maintain safe time to collision” in Fig. 3, ChatGPT states, “Since your state space includes the time-to-collision (TTC) for the front and rear vehicles, you can encourage the agent to maintain a safe TTC with both vehicles.... The reward is applied only when the actual TTC is less than the determined *safeTTC*.”

B. Experiment Setup

We train an RL policy for each valid answer using Proximal Policy Optimization (PPO) [23] with default RLlib hyperparameters [29]. The HVs use the Intelligent Driver Model (IDM) [16] with the stochastic noise range $[-0.2, 0.2]$ added to account for heterogeneous driving behaviors. Policies are trained for 200 episodes. A fully-connected neural network with 2 hidden layers of size 8 are used for the ring and intersection environments and size 16 are used for the bottleneck environment. Experiments are conducted with Intel i7-12700k CPU and 32G RAM.

C. Experiment Results

1) *Training*: We supply training curves for control (Novices) and study (Novices + ChatGPT) groups in Fig. 4.

The curves are normalized to $[0,1]$ with the mean reward values averaged across the two respective groups at each episode. We also supply the expert’s training curve. For all three environments and both groups, the curves show reward improving during training, validating participants were able to develop trainable policies as a result of the RVs’ actions. Increasing rewards does not guarantee the RV achieves the environment’s objective, as the RV may pursue actions that enhance rewards not in line with the goal.

2) *Ring*: Results are given in Fig. 5 (left). Due to the task’s complexity, we plot a policy’s best performance over 10 tests for each trained ring policy. All vehicle average speed (x-axis; meters/second) and minimum speed of any vehicle in the ring (y-axis) for the last 100 seconds (total testing period is 600 seconds) are considered. We consider policies successful (outlined) if their minimum speed is greater than zero while maintaining relatively good average speed. An expert’s performance [13] is also provided for comparison.

When using ChatGPT’s assistance, five policies are successful, while seven policies are successful using only non-expert knowledge. This result defies expectations given ChatGPT’s ability to increase valid answers and inject new metrics into the state spaces/reward functions for the ring environment. One explanation is despite the addition of new metrics, the metrics do not improve the robot vehicle’s ability to prevent stop-and-go traffic. Another conjecture is while the study group is encouraged to use ChatGPT, the level of usage in participants varies from asking ChatGPT a few questions to completely relying on it. This impacts ChatGPT’s ability to help the participants. Additionally, the high number of non-expert successful policies is unexpected. While none of the given policies reach the expert’s level, the anticipated number is close to zero given the tasks’ complexity. Non-experts have more capability than originally hypothesized.

3) *Bottleneck*: Fig. 5 (middle) shows the bottleneck results. Each trained RL policy is tested 10 times with the average reported. Outflow (x-axis; vehicles/hour) is considered, and an expert’s performance is given as a pink, vertical line. We consider a policy successful if the policy’s outflow is greater than 1400.

While the outflow range between novices and novices with ChatGPT is similar, we observe an increase in successful policies when using ChatGPT’s aid. The control group has 14 successful policies, while the study group has 19 successful policies, a 136% increase. Similar to ring, no participant-given policy outperforms the expert being near 100 shy.

4) *Intersection*: Intersection results are given in Fig. 5 (right). The trained RL policy is evaluated 10 times with average results reported. All vehicle average speed (x-axis; meters/second) and east/westbound queue length (y-axis; number of waiting vehicles) are considered. We plot an expert’s performance [13] and consider the nearest neighbors as successes (outlined). Four policies are successful with only non-expert knowledge, but six policies are successful (a 150% increase) when using ChatGPT aid. Additionally, of the six with-ChatGPT-help policies, four of them outperform

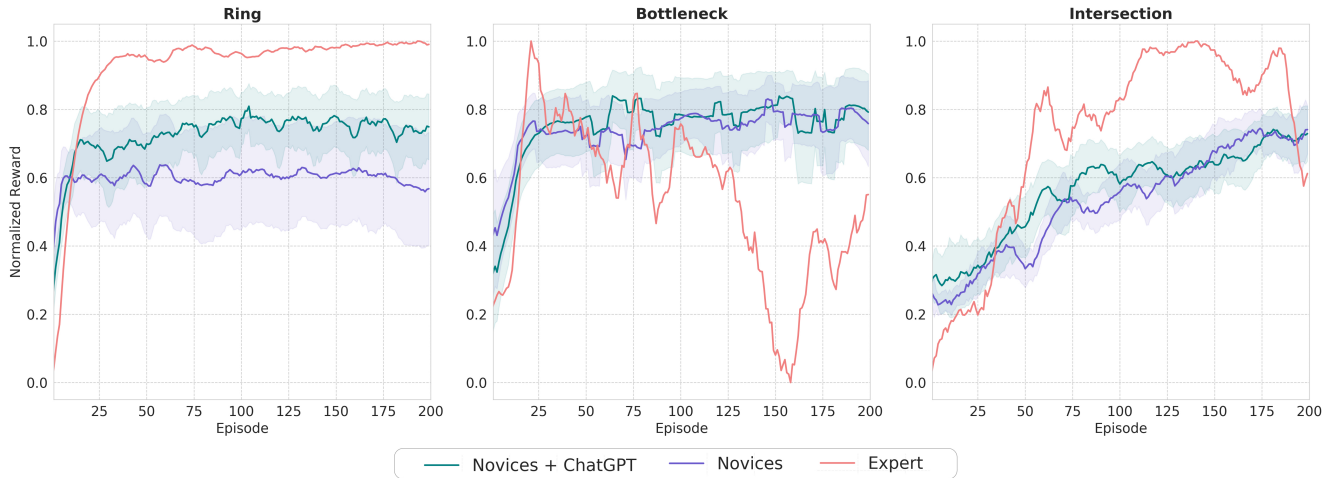


Fig. 4. Average normalized reward curves during training are shown for the three traffic control environments. Each environment consists of three curves for the control group (Novices), the study group (Novices + ChatGPT), and an expert. For Novices and Novices + ChatGPT, the solid line indicates the average for all participants, with the shaded region representing variance. In both groups across all networks, average rewards increase over the course of training. This validates that both groups are able to develop state spaces and reward functions that are trainable using RL.

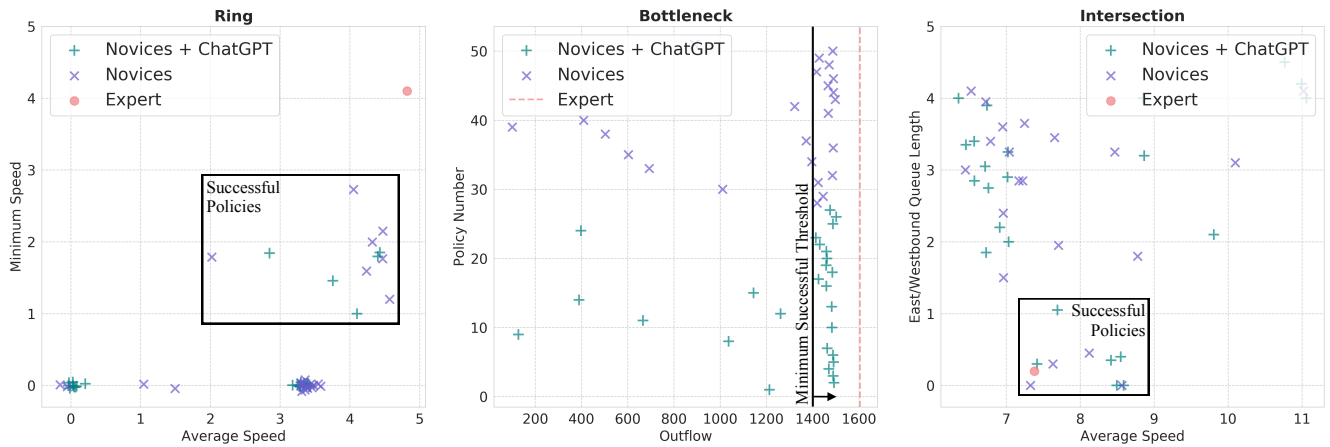


Fig. 5. Results for the three mixed traffic environments with successful RL policies denoted. For ring, five policies using ChatGPT’s help are successes, while seven policies are successful using only non-expert knowledge. Using ChatGPT sees a decrease in successful policies by two compared to only using non-expert knowledge, illustrating ChatGPT needs to be better prompted or further improvements to be useful in this task. For bottleneck, 14 policies are successful without ChatGPT’s assistance, while 19 are successful without ChatGPT, a 136% increase. For intersection, only using non-expert knowledge results in four successful policies, while ChatGPT increases successes to six (a 150% increase) with four (the rightmost green markers) of those outperforming the expert. The bottleneck and intersection increases illustrate how ChatGPT can enable more non-experts to solves complex mixed traffic control tasks. However, the number of increases is lower than expected, potentially showing ChatGPT needs better prompts or further improvement.

the expert policy by a significant margin. This result is significant and showcases how ChatGPT can give non-experts the ability to compete with ITS domain experts. Two non-expert policies outperform the expert policy, an unexpected outcome though half as much as using ChatGPT.

For both bottleneck and intersection, we observe a 136% and 150% increase, respectively, in successful policies. Examples of successful policies are provided in Fig. 6. ChatGPT is new with its training set not provided meaning ChatGPT could have not been trained on a sufficient amount of RL, ITS, or both training data to provide even more assistance. While the level of ChatGPT assistance is participant-determined, we observe a significant number of policies with extensive ChatGPT are not successful.

VI. CONCLUSIONS

In this work, we conduct a large-scale user study involving non-experts in intelligent transportation systems (ITS) research trying to provide quality reinforcement learning (RL) state spaces and reward functions for three mixed control traffic environments. Our study finds using ChatGPT can increase the number of successful polices by 150% and 135% in the intersection and bottleneck environments, respectively. However, using ChatGPT does not increase successes in the ring environment. Additionally, the improvement rate from using ChatGPT is less than originally theorized. This potentially means that an insufficient amount of RL ITS problems were provided to ChatGPT during training.

In the future, we intend to expand the study to include

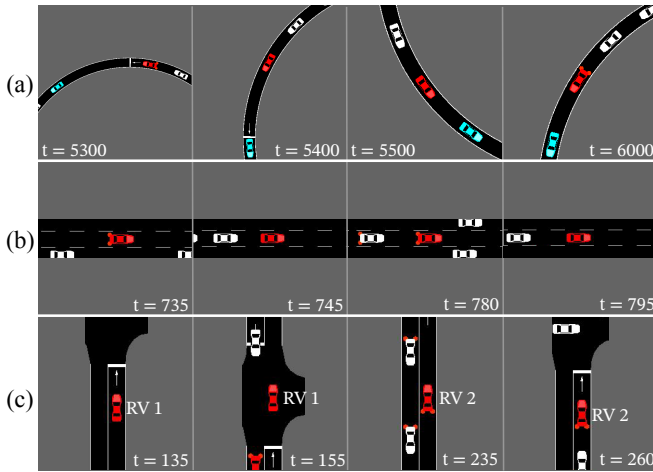


Fig. 6. Example successful policies with ChatGPT’s aid. (a) **Ring**. The RV first creates a gap to the leading vehicle (blue) to slow every vehicle down before gradually accelerating to stabilize the ring, preventing stop-and-go traffic. (b) **Bottleneck**. The RV gradually slows down (before speeding up again in last frame) to force following vehicles to slow down as well, easing congestion closer to the bottleneck’s end. (c) **Intersection**. The first RV travels safely through the intersection, while a second RV slows down so all east/westbound vehicles can cross without causing an incident.

more participants, traffic control environments, and Large Language Models (LLMs). A comparative analysis of various LLMs could also be an interesting direction. Further, fine-tuning existing LLMs to be able to serve the ITS community by reducing the human effort required is another avenue that we want to explore.

ACKNOWLEDGEMENT

This research is supported by NSF IIS-2153426. The authors would like to thank NVIDIA and the University of Tennessee, Knoxville for their support. The authors would also like to thank Heath Mitchell for GPT-4 API access.

REFERENCES

- [1] OpenAI. Gpt-4 technical report. Technical report, OpenAI, 2023.
- [2] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [3] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [4] Google. Palm2 technical report. Technical report, Google, 2023.
- [5] Yousef Wardat, Mohammad A Tashtoush, Rommel AlAli, and Adeeb M Jarrah. Chatgpt: A revolutionary tool for teaching and learning mathematics. *Eurasia Journal of Mathematics, Science and Technology Education*, 19(7):em2286, 2023.
- [6] Katharina Jeblick, Balthasar Schachtner, Jakob Dexl, Andreas Mittermeier, Anna Theresa Stüber, Johanna Topalis, Tobias Weber, Philipp Wesp, Bastian Sabel, Jens Rieke, et al. Chatgpt makes medicine easy to swallow: An exploratory case study on simplified radiology reports. *arXiv preprint arXiv:2212.14882*, 2022.
- [7] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

- [8] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [9] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [10] Ammar Haydari and Yasin Yılmaz. Deep reinforcement learning for intelligent transportation systems: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(1):11–32, 2020.
- [11] Matthew Veres and Medhat Moussa. Deep learning for intelligent transportation systems: A survey of emerging trends. *IEEE Transactions on Intelligent transportation systems*, 21(8):3152–3168, 2019.
- [12] SAE International. Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles. https://www.sae.org/standards/content/j3016_202104/, Apr 2021.
- [13] 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, 2021.
- [14] Zhongxia Yan and Cathy Wu. Reinforcement learning for mixed autonomy intersections. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 2089–2094. IEEE, 2021.
- [15] Dawei Wang, Weizi Li, Lei Zhu, and Jia Pan. Learning to control and coordinate hybrid traffic through robot vehicles at complex and unsignalized intersections. *arXiv preprint arXiv:2301.05294*, 2023.
- [16] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000.
- [17] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Ayaan Wahid, Jonathan Tompson, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, et al. Palm-e: An embodied multimodal language model. *International Conference on Machine Learning (ICML)*, 2023.
- [18] Vishnu Sashank Dorbala, James F Mullen Jr, and Dinesh Manocha. Can an embodied agent find your” cat-shaped mug”? llm-based zero-shot object navigation. *arXiv e-prints*, pages ArXiv–2303, 2023.
- [19] Kolby Nottingham, Prithviraj Ammanabrolu, Alane Suhr, Yejin Choi, Hannaneh Hajishirzi, Sameer Singh, and Roy Fox. Do embodied agents dream of pixelated sheep?: Embodied decision making using language guided world modelling. In *Workshop on Reinventing Reinforcement Learning at ICLR 2023*, 2023.
- [20] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. *International Conference on Machine Learning (ICML)*, 2023.
- [21] Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. Foundation models for decision making: Problems, methods, and opportunities. *arXiv preprint arXiv:2303.04129*, 2023.
- [22] Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordination. *International Conference on Machine Learning (ICML)*, 2023.
- [23] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [24] 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.
- [25] Meead Saberi and Hani S Mahmassani. Empirical characterization and interpretation of hysteresis and capacity drop phenomena in freeway networks. *Transportation Research Record: Journal of the Transportation Research Board, Transportation Research Board of the National Academies, Washington, DC*, 2013.
- [26] Eugene Vinitsky, Aboudy Kreidieh, Luc Le Flem, Nishant Kheterpal, Kathy Jang, Cathy Wu, Fangyu 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.
- [27] OpenAI. Chatgpt: Optimizing language models for dialogue., 2023.
- [28] Typingmind. <https://www.typingmind.com>. Accessed: April 2023.
- [29] Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph E. Gonzalez, Michael I. Jordan, and Ion Stoica. RLlib: Abstractions for distributed reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2018.