

# **Rethinking the Design of Human-Centric AI Systems for Deployment in Transportation**

**Rex Chen**

CMU-S3D-25-107

August 2025

Software and Societal Systems Department  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

**Thesis Committee:**

Fei Fang (Co-chair)  
Norman Sadeh (Co-chair)  
Sean Qian  
Matteo Pozzi  
Ryan Shi (University of Pittsburgh)  
Peter Stone (The University of Texas at Austin)

*Submitted in partial fulfillment of the requirements  
for the degree of Doctor of Philosophy in  
Societal Computing.*

Copyright © 2025 Rex Chen

This research was supported by a NSERC Postgraduate Scholarship (Doctoral); NSF Grants IIS-1850477, CNS-1801316, and CNS-1914444; a Mobility21 grant; the Tang Family Endowed Innovation Fund; and a SCS Presidential Graduate Fellowship. The views and conclusions contained in this document are those of the author and should not be interpreted as representing the official policies, either expressed or implied, of any sponsoring institution, any government, or any other entity.

**Keywords:** Artificial intelligence, transportation, multi-agent learning, traffic signal control, gig driving, traffic simulation

## Abstract

Artificial intelligence (AI) has had a transformative impact in improving the efficiency, safety, and accessibility of transportation systems. Successes in existing deployments have led AI researchers to seek to develop and apply new AI algorithms — particularly those based on deep learning and multi-agent systems — to improve the performance and scalability of transportation. Many of these algorithms have shown strong performance in simulation-based evaluations. However, most of the state-of-the-art algorithms in the AI literature have never been physically deployed.

This thesis argues that a barrier to the deployment of many advanced AI technologies in transportation lies in that their designs are divorced from the key practical considerations of stakeholders. When these AI technologies are to be deployed in existing transportation systems, they face four key categories of challenges:

- (1) *uncertainty* in present and projected traffic conditions;
- (2) *heterogeneity* among users and deployment contexts;
- (3) *assurance* in terms of the understandability and safety of algorithms; and
- (4) *coordination* at the individual and system levels.

Using gig driving and traffic signal control (TSC) as representative problems, this thesis focuses on understanding how these challenges can be addressed by the design of AI systems, and proposes new designs and algorithms to improve the status quo. Specifically, my work involved:

- Understanding how designs that expose uncertainty in *gig driver schedule recommendation* can improve users’ trust over repeated interactions
- Evaluating the impact of heterogeneity in driver behaviour models and simulation scale on *traffic simulation* outcomes with statistically rigorous experiments
- Building an algorithmic pipeline for demand modelling in a *traffic simulation* to incorporate uncertain, heterogeneous detector data and stakeholder feedback
- Imposing safety constraints upon coordinated reinforcement learning (RL)-based *TSC* policies through a suite of action postprocessing techniques
- Developing a performant but scalable algorithm to distil RL-based *TSC* policies into coordinated, understandable decision trees

Based on my work, I show that designing AI technologies to better address these challenges also leads to algorithms with technical novelty. Various other problems within transportation (e.g. freight logistics and autonomous driving) and beyond transportation (e.g. robotic navigation and computer networking) are potential applications where my technical contributions could better align AI technologies with stakeholders’ needs and preferences. Regardless of the domain, I suggest that the path to the successful deployment of AI technologies lies in designing them with people in mind at every stage of the data-to-deployment pipeline.





## Acknowledgements

To my advisors — Fei Fang and Norman Sadeh — for leading me through five years of research, for opening doors and creating opportunities even in difficult times, for supporting my good ideas and pushing back against my bad ideas, for continually shaping me into a better communicator, and for coming from drastically different but complementary perspectives that, together, completed my PhD experience.

To my thesis committee — Sean Qian, Matteo Pozzi, Ryan Shi, and Peter Stone — for each approaching my interdisciplinary work from another unique angle of expertise, for bearing with us during many scheduling difficulties and logistical mishaps, and for asking hard questions and providing feedback that fundamentally interrogated, focused, and refined my research in the home stretch of my PhD.

To my mentees — Karen Wu and Nick Fettig — for showing initiative in exploring widely and thinking deeply as an indispensable part of my projects, for making the most of your talents under my imperfect leadership, and for instilling in me the confidence that, wherever you go with your career from this point onwards, you will help to create a bright future as the next generation of computer scientists.

To my other academic collaborators — Stephanie Milani, Zhicheng Zhang, Ruiyi Wang, Akshara Khare, Anagha Ravi Shankara, Chun Kai Ling, James Cunningham, Justin Kiefel, Daniel Smullen, Peter Story, Yuanyuan Feng, Yaxing Yao, Abhilasha Ravichander, and Aerin Zhang — for all of the brilliant ideas, thoughtful conversations, and productive working sessions that moved my research forward, whether or not the final product will see the light of day as part of this thesis.

To my stakeholder partners — Ryan Green, Brian Finnermore, Kunal Prajapati, and Karan Moudgil at Gridwise; John McCartney, Scott Morse, and Dave Palmer at Path Master; Eric Raamot, Tim Romashko, Kyle Love, and Jon Ringler at Econolite; Ken Mikula and Lori Daley at the City of Strongsville; Mike Schweikart at TMS Engineers; and Stan Caldwell, Lisa Kay Schweyer, Karen Lightman, and Chandani Sharma at Traffic21/Safety21 — for giving me the opportunity to understand and engage with the fascinating problems that constitute your day-to-day work, for sharing your time and resources with us to enable our research even as priorities and timelines shifted, and for moving my aspirations of benefitting real people closer to reality.

To my advisors' assistants — Linda Moreci and KC Turner — and S3D IT — Emanuel Bowes, Tom Pope, and Cole Jester — for all of their tireless dedication in responding to the logistical nightmares that underlaid many of my last-minute requests, the firm foundation without which none of my research would have been possible.

To my mentors at UBC — Kevin Leyton-Brown, Chris Cameron, and Jason Hartford — for introducing me to research, for inspiring my pursuit of social impact through my work, and for imparting ideas and influences that persist even in this thesis.

To my officemates — Jenny Liang, Manisha Mukherjee, Kush Jain, and Mei — for all of the companionship, the laughter, and the honesty through long days and nights, for opening me up to entirely new dimensions of computer science that reshaped my work, and for bearing witness to my growth in every facet of my life.

To my other friends made at CMU — Jenny Tang, Ally Nisenoff, Eli Claggett, Nadia Nahar, Chenyang Yang, Vasu Vikram, Kyle Liang, Kathleen Salazar, Kaia Newman, Rachel Kim, and Katie Vu — for making my experience here colourful and multifaceted, however long and however deeply I knew you in my time here.

To my friends and acquaintances made at FAccT 2025 — (in order of appearance) David Widder, Jessie J. Smith, Nari Johnson, Kimi Wenzel, Nedah Nemati, Jayesh Tripathi, Aditya Karan, Kaivalya Rawal, Joseph Summers, Jake Stein, Varun Rao, Quirine Smit, Nina van Liebergen, Cor Veenman, Min Kyung Lee, Lauren Olson, Emitzá Guzmán, and Alicia — for making my presentation of Chapter 3 a shining joy.

To my parents — Catherine Yu and Chi-Ming Chen — for the many endless sacrifices you have made to get me here, even when I have disappointed you.

To my fam, my coven, my band — Laurel Sipe, Connor Petrella, Nate Brooks and Emma Dawson, Josh Mateer, Anna Emmons, Michelle Bjornstad, Elanor Tang, Carlos and Kaitlyn Ortega, Kristin and Brandon Miller, Christy and Logan MacKenzie, Christina Cheng and Jonah Jeng, Sam, Amber, and Sophie Carlin, Asha Jacob, Martin Semwanga, Wesley Chen and Natalie Bloss, Vincent Ma, Annaliese Calzadilla, and Simon and Gabby Lingle — for being present in my highs and my lows, for seeing and embracing my most authentic self, for growing alongside me in ways I never thought possible, and for being the best friends I have ever had.

To Jed Whedon, Maurissa Tancharoen, Jeffrey Bell, Clark Gregg, and the cast and crew of *Marvel's Agents of S.H.I.E.L.D.*, for 136 episodes of life-changing television.

TO THE FIRST AND THE LAST, THE FAITHFUL AND TRUE, THE LIGHT THAT BROUGHT ME FORTH FROM DARKNESS, THE SHEPHERD THAT DEFENDED ME IN WEAKNESS, THE KING THAT PREPARED ME TO DO THIS GOOD WORK IN ADVANCE,  
THIS IS YOURS ALONE.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem Domains . . . . .	2
1.2	Thesis Statement . . . . .	3
1.3	Contributions . . . . .	4
1.3.1	Understanding the Longitudinal Impact of Designs that Expose Uncertainty on Gig Drivers' Trust in AI Decision Aids . . . . .	5
1.3.2	Evaluating the Impact of Heterogeneity in Driver Behaviour and Simulation Scale on Traffic Simulation Outcomes . . . . .	5
1.3.3	Building an AI-Enabled Pipeline for Traffic Simulation from Noisy, Multimodal Detector Data and Stakeholder Feedback . . . . .	6
1.3.4	Imbuing Reinforcement Learning Algorithms for Coordinated Traffic Signal Control with Operational Constraints . . . . .	6
1.3.5	Coordinating Decision Trees for Efficient Interpretable Multi-Agent Reinforcement Learning . . . . .	7
<b>2</b>	<b>The Real Deal</b>	
	<b>A Review of Gig Driving, Traffic Signal Control, Reinforcement Learning, and Deploying AI for Transportation</b>	<b>8</b>
2.1	Gig Driving . . . . .	8
2.1.1	Problem Formulation . . . . .	8
2.1.2	AI Algorithms for Gig Platforms . . . . .	10
2.1.3	AI Algorithms for Gig Drivers . . . . .	10
2.2	Traffic Signal Control . . . . .	11
2.2.1	Problem Formulation . . . . .	11
2.2.2	Detection Hardware . . . . .	12
2.2.3	Control Hardware and Algorithms . . . . .	13
2.2.4	Traffic Simulation . . . . .	15
2.3	Reinforcement Learning . . . . .	16
2.3.1	Markov Decision Processes . . . . .	16
2.3.2	Multi-Agent MDPs . . . . .	19
2.3.3	Partially Observable MDPs . . . . .	19
2.3.4	RL for TSC . . . . .	20
2.4	AI Design and Deployment Challenges . . . . .	21
2.4.1	Uncertainty . . . . .	21
2.4.2	Heterogeneity . . . . .	23

2.4.3	Assurance . . . . .	24
2.4.4	Coordination . . . . .	26
<b>3</b>	<b>Missing Pieces</b>	
	<b>Understanding the Longitudinal Impact of Designs that Expose Uncertainty on Gig Drivers' Trust in AI Decision Aids</b>	<b>28</b>
3.1	Introduction . . . . .	28
3.2	Related Work . . . . .	29
3.2.1	Trust in AI . . . . .	29
3.2.2	Recommendation Systems . . . . .	32
3.3	Pilot Interview Study . . . . .	32
3.3.1	Prototype Design . . . . .	32
3.3.2	Methodology . . . . .	33
3.3.3	Participants . . . . .	33
3.3.4	Results . . . . .	34
3.4	AI-Based Schedule Recommendation for Gig Driving . . . . .	37
3.4.1	Decision Aid Design . . . . .	37
3.4.2	Interface Design . . . . .	38
3.4.3	Interface Conditions . . . . .	38
3.5	Longitudinal User Study Design . . . . .	41
3.5.1	Participants and Data Sources . . . . .	41
3.5.2	User Study Activities . . . . .	42
3.5.3	Interview Procedure . . . . .	43
3.6	Quantitative Analysis . . . . .	43
3.6.1	Metrics of Trust and Reliance . . . . .	44
3.6.2	RQ1: Longitudinal Effects . . . . .	45
3.6.3	RQ2: Effects of Conditions . . . . .	47
3.7	Qualitative Analysis . . . . .	50
3.7.1	Motivations and Routines . . . . .	50
3.7.2	RQ1: Perceptions of Accuracy . . . . .	51
3.7.3	RQ2: Perceptions of Uncertainty . . . . .	51
3.8	Discussion . . . . .	52
3.8.1	Key Findings and Implications . . . . .	52
3.8.2	Limitations . . . . .	53
3.8.3	Recommendations for Future Work . . . . .	54
3.9	Conclusion . . . . .	54
<b>4</b>	<b>Purpose in the Machine</b>	
	<b>Evaluating the Impact of Heterogeneity in Driver Behaviour and Simulation Scale on Traffic Simulation Outcomes</b>	<b>55</b>
4.1	Introduction . . . . .	55
4.2	Related Work . . . . .	56
4.2.1	Validating Traffic Simulators . . . . .	56
4.2.2	Modelling Driver Behaviour . . . . .	57
4.3	Comparing SUMO and CityFlow . . . . .	57
4.4	Varying Driver Behaviour . . . . .	59

4.4.1	Car-Following Models . . . . .	59
4.4.2	Lane-Changing Models . . . . .	60
4.5	Experiments . . . . .	61
4.5.1	Setup . . . . .	62
4.5.2	Results . . . . .	63
4.5.3	Parameter Validity . . . . .	68
4.5.4	Acceleration-Speed Diagrams . . . . .	68
4.6	Conclusion . . . . .	70
<b>5</b>	<b>Out of the Past</b>	
	<b>Building an AI-Enabled Pipeline for Traffic Simulation from Noisy, Multimodal Detector Data and Stakeholder Feedback</b>	<b>71</b>
5.1	Introduction . . . . .	71
5.2	Related Work . . . . .	73
5.2.1	Demand Modelling for Traffic Simulation . . . . .	73
5.2.2	Computer Vision for Traffic Footage . . . . .	74
5.2.3	Large Language Models for Transportation Research . . . . .	74
5.3	Demand Modelling Pipeline . . . . .	74
5.3.1	Computer Vision-Based Vehicle Counting from Camera Footage . . . . .	74
5.3.2	Optimisation-Based Vehicle Route Generation from Multimodal Data . . . . .	75
5.3.3	LLM Agent Simulation Refinement from Natural Language Feedback . . . . .	77
5.4	Simulation Results: Strongsville, Ohio . . . . .	79
5.4.1	RQ1: Accuracy of Vehicle Counting . . . . .	79
5.4.2	RQ2: Accuracy of Generated Simulation . . . . .	81
5.4.3	RQ3: Accuracy of LLM Agent-Generated Constraints . . . . .	82
5.5	Conclusion . . . . .	87
<b>6</b>	<b>Self Control</b>	
	<b>Imbuing Reinforcement Learning Algorithms for Coordinated Traffic Signal Control with Operational Constraints</b>	<b>88</b>
6.1	Introduction . . . . .	88
6.2	Related Work . . . . .	90
6.2.1	Safety Constraints in Reinforcement Learning . . . . .	90
6.2.2	Safety Assessment in Traffic Signal Control . . . . .	92
6.3	Problem Formulation . . . . .	93
6.4	CycleLight . . . . .	95
6.4.1	Action Masking . . . . .	95
6.4.2	Imitation-Based Loss . . . . .	97
6.4.3	Action Projection . . . . .	98
6.5	Experiments . . . . .	99
6.5.1	Simulation Environment . . . . .	100
6.5.2	Simulations . . . . .	101
6.5.3	Algorithms and Baselines . . . . .	102
6.5.4	Evaluation Metrics . . . . .	102
6.5.5	Cologne Corridor Results . . . . .	103
6.5.6	Strongsville Network Results . . . . .	109

6.6	Conclusion	109
<b>7</b>	<b>Making Teams and Influencing Agents</b>	
	<b>Coordinating Decision Trees for Efficient Interpretable Multi-Agent Reinforcement Learning</b>	<b>111</b>
7.1	Introduction	111
7.2	Background	113
7.2.1	Decision Trees	113
7.2.2	VIPER	113
7.3	Related Work	114
7.4	HYDRAVIPER	114
7.4.1	Dataset Resampling: Centralised-Q Weighting	115
7.4.2	Training Rollouts: Adaptive Budget Allocation	117
7.4.3	Validation Rollouts: UCB Policy Selection	119
7.4.4	Agent Clustering: Scaling Up HYDRAVIPER	120
7.5	Experiments	122
7.5.1	Environments	122
7.5.2	Baselines and Setup	123
7.5.3	Results	124
7.5.4	Hyperparameter Sensitivity	130
7.5.5	Ablation	131
7.6	Conclusion	132
<b>8</b>	<b>End of the Beginning</b>	
	<b>Lessons, Perspectives, and Visions on AI Design and Deployment</b>	<b>134</b>
8.1	Introduction	134
8.2	Lessons on Designing and Deploying AI Systems	135
8.2.1	Incorporating Stakeholder Heterogeneity	135
8.2.2	Optimising for Assurances and Understandability	136
8.2.3	Managing Uncertainty and Missing Data	136
8.2.4	Coordinating with Human Interaction and Control	137
8.3	Perspectives on Goals and Possibilities of AI Deployments	138
8.3.1	Aligning Goals in Research and Deployment	138
8.3.2	Fostering Trust to Drive Progress	139
8.3.3	Generalising to Applications in Other Domains	139
8.4	Visions on Furthering the Capabilities of AI Technologies	143
<b>A</b>	<b>Appendix</b>	<b>145</b>
A.1	Chapter 3	145
A.1.1	Survey Questions	145
A.1.2	Interview Scripts	153
A.1.3	Full Quantitative Results	157
A.2	Chapter 5 — Full Claude Prompt	158

# Chapter 1

## Introduction

Welcome to Level Seven.

---

*Philip J. Coulson*

Emerging technologies based on *artificial intelligence* (AI) have revolutionised the state of modern transportation systems. During the past decade, AI technologies have been deployed on exponentially increasing scales to optimise traffic signals to reduce congestion and emissions [219, 376]; to dispatch gig drivers to improve the mobility of people and goods [159, 229]; to streamline public transport and freight logistics to improve utilisation and efficiency [165, 299]; and to control connected and autonomous vehicles to make individual transport safer and more accessible [132, 237]. In these applications, AI has mitigated systematic inefficiencies caused by suboptimal practices and insufficient capacity in existing infrastructure [318, 354].

The key to these successes has been the proliferation of *deep learning*, which has allowed AI algorithms to achieve performance beyond that of humans or even non-deep learning algorithms [346]. Concomitant advances in computing infrastructure have also allowed AI technologies to adapt to rapid change by ingesting data at massive scales [310]. These results have fuelled the growth of transportation as an application of AI, along with optimistic attitudes among forward-thinking stakeholders who believe in AI's capabilities to improve the status quo [314].

Consequently, transportation has also attracted attention from the AI literature as an application domain. Researchers have sought to develop even more performant AI algorithms for various problems derived from transportation systems [366]. These problems have ranged from prediction, forecasting, and pattern recognition to control, planning, and optimisation [174], based on applications in ground (including basic infrastructure, vehicle control, fleet management, and public transit) [120, 245, 307], air [96, 330], and maritime [293] transportation. Many of these novel AI algorithms have demonstrated strong capabilities in simulation-based evaluations (e.g. [15, 247]).

However, there is a disconnect between theory and practice. Most of the state-of-the-art algorithms developed in the AI literature have never been deployed in practice, despite their seemingly strong performance. Even when they have been deployed, their potentials have largely remained unrealised. In particular, there are many stakeholders who remain concerned about the possible impacts of these AI algorithms on the efficiency, equity, and ethics of transportation systems [314]. In this thesis, I ask: Why does the wheel of progress turn so slowly? What challenges prevent

stakeholders from being willing to deploy AI technologies for transportation? What technical advances in AI are needed to pave the way towards successful deployments, so that the promising performance in experimental evaluations can be translated into tangible benefits?

## 1.1 Problem Domains

Within the broad space of AI for transportation, I address these motivating questions with a focus on two representative problems to which AI has been applied:

- *Gig driving* (reviewed in Section 2.1), where platforms dispatch a pool of drivers to provide on-demand transportation of passengers or orders
- *Traffic signal control* (TSC; reviewed in Section 2.2), where controllers set the light sequences of a series of traffic signals to enable vehicles to pass through intersections

In Table 1.1, I summarise the characteristics of gig driving and traffic signal control. Although every transportation problem involves an idiosyncratic set of objectives and challenges, some of these characteristics are shared with other problems. My goal is accordingly to generalise the insights that can be derived from these two representative problems.

Gig driving	Traffic signal control
<b>Similarities</b>	
<b>Multi-agent:</b> Usually involves coordinating many drivers / intersections	
<b>Longitudinal:</b> Need to keep taking gigs / signalling phases	
<b>Differences</b>	
<b>Individual-level control:</b> Each driver and platform makes self-interested decisions with partial supply and demand information	<b>System-level control:</b> Control could be centralised but is distributed for technical reasons such as efficiency
<b>Mixed competitive-cooperative:</b> Various different drivers and platforms may collaborate or conflict with each other	<b>Cooperative:</b> All intersection agents seek to optimise global performance metrics across the road network
<b>Opacity:</b> Platforms are revenue maximisers and hide information from their users	<b>Transparency:</b> Stakeholders need traceability in order to pursue social benefits

Table 1.1: Comparison of problem characteristics for gig driving and traffic signal control

Both of these problems have been the subject of vigorous study in the AI literature. Thus, many AI-based solutions have been proposed, but few of them have been deployed. For gig driving, dispatching mechanisms with various theoretical guarantees have been designed, but their implementations are more often based on heuristics [304]. For TSC, reinforcement learning (RL) algorithms easily control and coordinate between dozens, hundreds, and even thousands of traffic signals in simulations [65]. However, after a decade of advances in RL, I am only aware of one deployment of these algorithms, which is limited to three intersections considered in isolation [85].



## 1.2 Thesis Statement

*This thesis proposes methods to integrate the key practical considerations of stakeholders into the design of AI systems intended for deployment in transportation.* In doing so, I address what I consider the most critical barrier to the practical applicability of AI technologies in transportation: not their performance, but rather their ability to be embedded in complex, human-centric deployment contexts. By optimising for gains in performance rather than deployability, many state-of-the-art AI algorithms have become divorced from the practical challenges faced by stakeholders. In this thesis, I address research questions relating to four design and deployment challenges:

- **Uncertainty.** How can AI technologies in transportation achieve robust levels of performance when longitudinal variability in the movement of people makes it difficult to observe or predict traffic conditions? This challenge arises in the transition from simulation to deployment: through varying traffic patterns and sensor imperfections in traffic signal control, and through the on-demand nature of gigs and information asymmetry in gig driving. Stakeholders may be unable to trust the performance of AI technologies trained in uncertainty-devoid simulations when they are deployed in the real world.
- **Heterogeneity.** How can AI technologies in transportation be broadly applied in the presence of variation between end-users and deployment contexts, and how might this be hindered by assumptions made during their design? This challenge arises during the generalisation of AI technologies: road users or gig drivers have unique personalities, motivations, and vehicle types that lead them to react to AI-driven decisions differently, and road networks and cities differ in scale and traffic characteristics. Stakeholders who represent minorities in datasets may be short-changed when AI technologies abstract out such variation.
- **Assurance.** How can AI technologies in transportation present their decisions to stakeholders in a way that is understandable and safe? This challenge arises during design and after deployment: deployed AI technologies must provide assurances to stakeholders about the decisions that they make (e.g. traffic signal plans or gig driver work schedules). In particular, they should align with stakeholders' mental models and provide guarantees against risks. Stakeholders will hesitate to delegate control to opaque AI technologies that lack these guarantees. User-centred evaluations are needed to thoroughly address this challenge.
- **Coordination.** How can AI technologies in transportation coordinate between interactions between individual agents (both human and AI) to achieve optimal system-level outcomes? This challenge arises inherently in the design process: traffic signal control requires coordinated timings between intersections, and gig driving requires coordination between drivers' movements. Contrary to assumptions in previous work, such coordinated behaviour is not guaranteed to be emergently learnt by AI technologies that distribute control between agents. Stakeholders will take issue with misalignment between individual and system-level goals.

However, theoretical and practical contributions in gig driving and TSC are not mutually exclusive. This thesis explores many dimensions along which stakeholder-centric considerations can be incorporated into the design and construction of AI technologies. Such dimensions include learning from their considerations in existing deployments (Chapter 2), fostering their perceptions of trust and reliance on AI (Chapter 3), accounting for their distinct patterns of behaviour (Chapter 4),

using their knowledge to fill in unknowns (Chapter 5), enforcing constraints that they can depend on for safety (Chapter 6), and creating interpretable model structures that they can understand (Chapter 7). My solutions to these challenges lead to technical novelties that generalise to other AI applications within transportation (e.g. freight logistics and autonomous driving) and beyond transportation (e.g. robotic navigation and computer networking), which I review in Chapter 8.

I argue that these practical considerations about people cannot be an afterthought but instead must be integrated throughout design, implementation, and evaluation. By being involved, stakeholders in transportation can be assured that the purpose of AI technologies is not to substitute or frustrate them, but to complement and enhance their work [234, 331]. Only then will AI technologies in transportation move smoothly from the cutting edge of the literature to deployments, thus achieving their potential to improve our fast-changing transportation systems.

## 1.3 Contributions

My main contributions within this thesis are summarised in Table 1.2.

Ch.	Domain	Challenges	Contribution	Citation
3	Gig driving	Uncertainty	User study showing that design factors impact longitudinal trust in driver schedule recommendation	[71]
4	Traffic simulation	Heterogeneity	Experiment showing that driver behaviour models significantly affects traffic simulation outcomes	[68]
5	Traffic simulation	Uncertainty Heterogeneity	Demand modelling pipeline that incorporates detection error and stakeholder feedback	[72]
6	Traffic signal control	Assurance Coordination	Suite of techniques for guaranteeing that RL-based TSC policies follow safety constraints	
7	Traffic signal control	Assurance Coordination	Performant, efficient algorithm that distils RL-based TSC policies into simpler decision trees	[70]

Table 1.2: Summary of contributions in this thesis. The last two chapters are not yet published.

Chapter 2 begins with an overview of gig driving and traffic signal control (including traffic simulation), the two representative problems that I consider in this thesis. In addition, I also review common formulations and algorithms in reinforcement learning, which are central to Chapters 6 and 7. Lastly, I consider how the challenges of uncertainty, heterogeneity, assurance, and coordination arise in these two problems; why these challenges are relevant considerations for stakeholders;

what progress has been made by the AI literature in addressing them; and where open, stakeholder-centric problems remain to be addressed as a prerequisite to deployment.

### **1.3.1 Understanding the Longitudinal Impact of Designs that Expose Uncertainty on Gig Drivers' Trust in AI Decision Aids**

The majority of this thesis is based on collaborations that I conducted with stakeholders in gig driving and traffic signal control. First, I focus on my collaboration with Gridwise, the developer of an assistant app for gig drivers. While the current Gridwise app mainly aims to empower gig drivers through presenting insights, this collaboration explored the potential of prescribing behavioural nudges to drivers so as to improve their earnings.

Chapter 3 focuses on designing an AI decision aid for gig drivers: a tool that recommends a schedule of when to drive and also predicts the earnings that drivers could make by following the schedule. The task of estimating drivers' earnings in the complex gig driving ecosystem is fraught with uncertainty. Due to natural variance in supply and demand, as well as in driver behaviour, these estimates can easily diverge from reality. My work examines how this uncertainty can be presented to users through different designs and how uncertainty impacts their trust and reliance on the tool. By doing so, I draw generalisable conclusions about the innate psychological factors that underpin stakeholders' willingness to engage with AI technologies in transportation.

I created a schedule recommendation tool that estimates weekly and hourly earnings from historical data and elicits constraints from users to recommend a schedule. Next, I conducted a longitudinal user study where 51 gig drivers recruited through Gridwise interacted with the tool for 7 days. Some drivers were exposed to uncertainty through designs that showed range-based estimates and hedging text. Based on a combination of quantitative statistical analyses and qualitative interviews, I show that ranges improved trust and reliance over repeated interactions, but that hedging had the opposite effect. Overall, my results suggest that the difference between these conditions is task alignment and personalisation, which are important for building trust in AI systems.

### **1.3.2 Evaluating the Impact of Heterogeneity in Driver Behaviour and Simulation Scale on Traffic Simulation Outcomes**

Second, in the remainder of my thesis, I focus on a collaboration aimed at deploying RL-based TSC in the city of Strongsville, Ohio. This work has involved the City's Traffic Management Centre; Path Master, a distributor of traffic control technologies; Econolite, a technology company that develops intelligent transportation systems; and PTV, a technology company that develops traffic simulation and planning software. In Chapters 4 to 7, I move through the deployment pipeline for RL-based TSC policies. These chapters expound on how deployment challenges arise at each stage of the pipeline: choosing a realistic simulator, building a data-driven simulation, training and evaluating the RL policy, and lastly deploying the RL policy via an interpretable surrogate. At each stage, a precondition for deployment has been that my modelling choices must adhere as closely as possible to real-world conditions in Strongsville.

Chapter 4 begins with the problem of choosing a realistic simulator and simulation parameters. Stakeholders can only trust the results of RL training if the traffic simulator is able to realistically represent the impacts of different factors on traffic. For this chapter, I conducted a comparison of how heterogeneity is modelled by two traffic simulators commonly used for RL-based TSC, SUMO and CityFlow. Relative to SUMO, CityFlow abstracts away many details in how it implements driver behaviour — such as car-following and lane-changing models — and traffic flows, in exchange for a significant speedup. Previous work had used these simulators interchangeably, based on a limited evaluation that found no significant differences in their simulation outcomes.

I contend that the evaluation and selection of traffic simulators must be based on a variety of simulations and metrics. Thus, I performed a statistically rigorous evaluation of how the similarity in simulation outcomes between the two simulators changes as heterogeneity is introduced. SUMO and CityFlow show statistically significant differences ( $p \ll 0.001$ ) in system-level (travel time), lane-level (vehicle counts), and vehicle-level (speed and acceleration) metrics. Generally, these differences increased as drivers were modelled to be more aggressive, and as the traffic volume or the road network scale were increased. Although the tradeoff between simulators ultimately depends on research needs, I rely on the more realistic SUMO in the remainder of this thesis.

### **1.3.3 Building an AI-Enabled Pipeline for Traffic Simulation from Noisy, Multimodal Detector Data and Stakeholder Feedback**

Chapter 5 uses the SUMO simulator to create a data-driven traffic simulation pipeline, specifically to create a simulation of Strongsville. Modelling traffic demand is a crucial task in traffic simulation. Given traffic detector data from Strongsville, a typical approach would be to apply route generation procedures that reconstruct this data into origin-destination vehicular flows; to run a simulation using these flows; and then to train an RL policy using the simulation. However, stakeholders noted significant discrepancies between simulated and real-world traffic conditions. A closer examination revealed uncertainty and heterogeneity in the data: it originates from multiple detector types, each with varying error rates that depend on traffic volumes.

I show that existing traffic detectors and simulation pipelines fail to properly account for these errors. To replace them, I contribute a generalisable demand modelling pipeline that includes several technical novelties: (1) an accurate computer vision-based method for vehicle counting directly from raw camera footage; (2) an uncertainty-aware quadratic integer program for integrating multiple detector data sources into a single simulation; (3) and a large language model agent that directly uses stakeholders’ natural language feedback to further refine the accuracy of the simulation. Using this pipeline, I generated a 36-intersection traffic simulation with more than 200 000 vehicles, which stakeholders in Strongsville found to be more trustworthy than previous iterations.

### **1.3.4 Imbuing Reinforcement Learning Algorithms for Coordinated Traffic Signal Control with Operational Constraints**

Chapter 6 uses the SUMO simulation of Strongsville to train RL-based TSC policies. For such policies to be deployable, stakeholders need them to be coordinated to optimise traffic flow glob-

ally. Stakeholders also need assurance that these policies will not violate safety constraints, such as a minimum green light duration that allows vehicles to safely traverse an intersection. Existing RL algorithms for TSC do well in coordination, but poorly in enforcing constraints. By contrast, Strongsville’s current TSC algorithm outputs structured, constraint-satisfying signal plans, but they are optimised heuristically. My work combines the strengths of these approaches by imbuing RL policies with Strongsville’s signal plan format and constraints.

I introduce CycleLight, a suite of postprocessing techniques that can extend state-of-the-art RL algorithms for TSC. The primary techniques are action masking, which limits the set of valid actions for policies that make periodic signalling decisions; imitation learning, which outputs a surrogate policy that mimics periodic policies while committing to signal durations in advance; and action projection, which explicitly trains the surrogate policy to follow the constraints through differentiable optimisation. Evaluations in a benchmark simulation show that CycleLight allows RL-based TSC algorithms to outperform fixed signal plans by 10% while adhering to operational constraints; evaluations in the Strongsville simulation are ongoing but promising.

### **1.3.5 Coordinating Decision Trees for Efficient Interpretable Multi-Agent Reinforcement Learning**

Chapter 7 applies imitation learning to distil RL-based TSC policies into a more understandable format. Specifically, I represent RL policies using decision trees, which are much more understandable than the opaque deep neural network architectures at the core of state-of-the-art RL algorithms for TSC. More interpretable surrogate policies can be deployed in place of neural policies to improve trust among stakeholders and can allow stakeholders to modify and control them. However, these surrogates must be both performant and computationally efficient. Existing interpretable multi-agent RL algorithms fail to balance performance and computational efficiency, either being uncoordinated or relying on costly training procedures to achieve good performance.

I improve upon this prior work with HYDRAVIPER, an efficient, interpretable multi-agent RL algorithm that extracts coordinated decision trees. HYDRAVIPER jointly samples environment interactions to obtain shared training datasets for teams of cooperative agents, thus improving performance; it also adaptively allocates environment interaction budgets for training and evaluating surrogates, thus improving computational efficiency. For TSC environments, HYDRAVIPER reduces runtime relative to baseline algorithms by nearly an order of magnitude while retaining comparable performance. Although HYDRAVIPER shows promise as a means of fostering trust among stakeholders in Strongsville, user evaluations are still needed to ensure its practical utility.

Lastly, in Chapter 8, I take a bird’s eye view of my work to distil insights into how to advance AI technologies for transportation. I ask: What have I learnt about designing AI systems to meet the needs of these stakeholders? What have I learnt about moving these AI systems from designs to concrete deployments? What have I learnt about the fundamental goals of working with stakeholders to deploy AI systems? In doing so, I aim to provide insight for future research that aims to narrow the gap between theory and practice in AI for transportation and beyond.

# Chapter 2

## The Real Deal

### A Review of Gig Driving, Traffic Signal Control, Reinforcement Learning, and Deploying AI for Transportation

*Domain:* Gig driving, traffic signal control  
*Challenges:* Uncertainty, heterogeneity,  
interpretability, coordination

To wait even one second to do this is  
playing with fire.

---

*Philip J. Coulson*

In Chapter 1, I introduced *gig driving* and *traffic signal control* (TSC) as two transportation problems that I focus on as deployment contexts of AI technologies. This chapter provides background on both problems. Section 2.1 and Section 2.2 begin with an overview of gig driving and TSC from the perspective of stakeholders — gig platforms, gig drivers, and traffic engineers — rather than from the perspective of the AI literature. I show that pain points do exist in the state of practice: information asymmetry leads to unpredictability and inequity in gig driving, while imperfect detection and control reduce throughput in TSC. Next, in Section 2.3, I introduce *reinforcement learning* (RL), the class of AI algorithms for TSC that I consider in Chapters 6 and 7. Lastly, in Section 2.4, I describe how the challenges of uncertainty, heterogeneity, interpretability, and coordination from Chapter 1 occur in gig driving and TSC. I also discuss how existing methods from the AI literature have addressed these challenges for both problems, how they have been inadequate, and what work still remains to be done.

This chapter was adapted from text in various publications, but a significant portion was published as a review paper that surveyed deployment challenges of RL for TSC, at the International Workshop on Agents in Traffic and Transportation (ATT) at IJCAI in 2022 [69].

## 2.1 Gig Driving

### 2.1.1 Problem Formulation

*Gig drivers* are independent contractors who use their personal vehicles to provide on-demand transportation for passengers (*ridesourcing*) and restaurant or grocery orders (*food delivery*) from

app platforms such as Uber, Lyft, DoorDash, and Instacart. I consider ridesourcing and food delivery to be a unified problem, not only due to similarities in their problem formulations [167], but also because drivers view them interchangeably [218] (e.g. Uber and its food delivery business, Uber Eats, are mutually cannibalistic platforms [78]).

In gig driving, there are interactions between three different groups of stakeholders: the customers who request rides or deliveries, the drivers who service these requests, and the platforms that mediate their interactions (Figure 2.1). Customers initiate by sending their requests and payments directly to the platforms. The platforms then attempt to match the set of requests with the set of available drivers. Each driver is then shown a list of requests that the platform recommends, and is free to accept or decline each one. If they accept, they fulfil the request; if they decline, the platform falls back to an alternative assignment. Lastly, platforms distribute payments to drivers. From the platforms’ perspective, this is an iterative weighted matching problem [167].

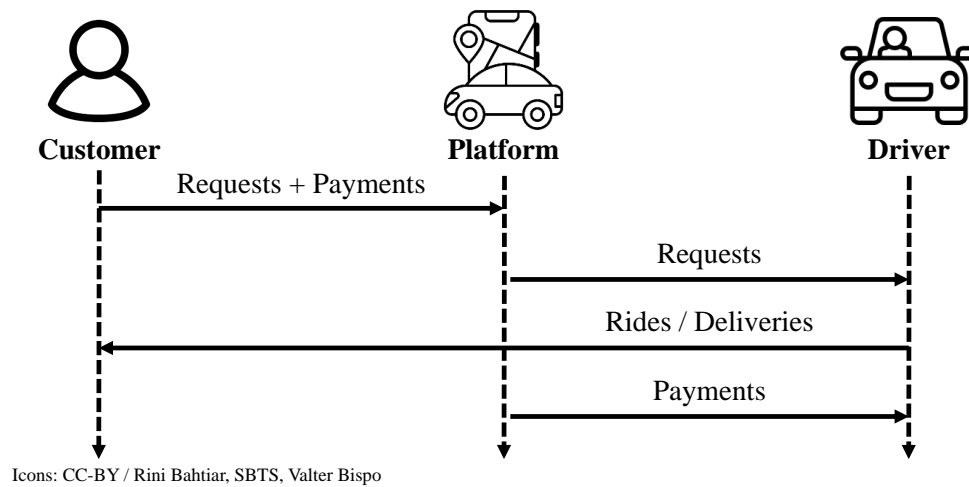


Figure 2.1: Flow of transactions between customers, platforms, and drivers in gig driving.

Although the realisation of this model is a relatively recent phenomenon, gig driving has experienced explosive growth [156] due to the benefits it provides. For customers, gig driving provides a convenient alternative to both private vehicles and public transport [76, 172, 349, 360]. For drivers, gig driving offers them the flexibility and autonomy to choose when and where to work [67, 218]. However, gig driving has also introduced new challenges. Negative interactions between customers and drivers degrade the experiences of both parties [160, 335, 384]. For drivers, the promise of autonomy is also hindered by the opacity of platforms’ assignment, pricing, and evaluation mechanisms [144, 218, 372, 401, 448], leading to volatility and systemic inequities in driver outcomes [95, 271]. Lastly, gig driving has created societal externalities. At least in the short term, gig driving will continue to contribute to increased congestion: despite the efforts of platforms to manage their supply and demand, they have not only shifted demand away from other transport modes but also created new demand altogether [159, 349].

### 2.1.2 AI Algorithms for Gig Platforms

In practice, gig platforms extend the basic problem formulation to account for and to shape fluctuations in supply and demand. One of the most prevalent strategies adopted by these platforms is *dynamic pricing*, referred to as “surge pricing” by Uber, “Prime Time” by Lyft, “Peak Pay” by DoorDash, and “Blitz Pricing” by Postmates [218, 220]. By setting higher prices during peak hours, platforms balance demand by disincentivising customers from making requests, while also incentivising drivers to reposition to high-demand regions and accept requests in exchange for bonuses [244, 304]. Precise implementation details for these mechanisms are generally scant. Ong et al. [304] describe Lyft’s “escrow” mechanism, in which convex optimisation problems are solved to prospectively allocate customers’ expected Prime Time payments equitably among drivers. However, deployed algorithms likely incorporate various heuristics based on business requirements.

Just as external researchers have little visibility into the inner workings of dynamic pricing algorithms, drivers must contend with both information asymmetries and gamification mechanisms that some perceive to be unfair [477]. Within their limited control, drivers have responded to the algorithmic management exerted by platforms by switching between platforms [218], participating in discussion forums [253, 477], and even colluding to induce artificial surges [410]. Nevertheless, incentives created by these opaque pricing mechanisms still create tangible effects on drivers, including inequitable increases in the level of competition and effort required by drivers [271].

### 2.1.3 AI Algorithms for Gig Drivers

One line of work in AI for gig driving has focused on gig platforms by designing mechanisms for driver dispatching and payment. These mechanisms aim to achieve various desirable properties, such as welfare maximisation and incentive compatibility (i.e. guaranteeing that drivers will accept dispatched trips) [41, 131, 247, 332]. In particular, incentive compatibility is enforced by ensuring that prices offset drivers’ opportunity costs [131, 247]. However, these mechanisms generally make unrealistic simplifying assumptions. For instance, drivers in reality rarely have perfect information about even dispatched trips (e.g. drop-off locations are often obscured [477]), much less the opportunity costs of declining them. Even works grounded in particular platforms (such as Ong et al. [304]) fail to consider the complexities that arise from interactions between platforms: how do drivers’ opportunity costs change when they can switch to a different platform?

Simultaneously, a complementary line of work has focused on understanding and responding to the needs of drivers; this line of work is closer to the focus that I adopt. Zhang et al. [477] conducted focus groups with ridesourcing drivers to envision possible improvements to gig platforms; one of their core findings was a need for native in-platform data-driven insights that would eliminate the need for third-party tools (e.g. mileage tracking apps). Based on this, Zhang et al. [476] designed “data probes” to help drivers understand how their work patterns and positionalities interact with the management practices of platforms. Khan et al. [190] created a measurement suite that quantitatively analyses the dynamic pricing strategies of platforms. Yet, given this abundance of information, there is no work on how to help drivers use this information to make better decisions.



## 2.2 Traffic Signal Control

### 2.2.1 Problem Formulation

*Traffic signal control* (TSC) aims to allocate green time at one or more intersections in a road network to traffic moving in different directions. In this thesis, I consider each intersection to be controlled by an individual decision-making agent. At each intersection, every *approach* (roadway, e.g. northbound or southbound) is divided into *lanes*. Vehicles traversing the intersection can follow different *movements*, each of which is defined by one incoming and one outgoing lane (e.g. northbound left turn or straight through) [141, 436, 495]. A *phase*  $\phi \in \Phi = \{\phi_1, \dots, \phi_P\}$  is a combination of movements that is signalled as a single unit [203]. For efficiency, pairs of non-conflicting movements are often signalled simultaneously (e.g. westbound/eastbound left turn, northbound left turn/straight through) [203, 302, 435].

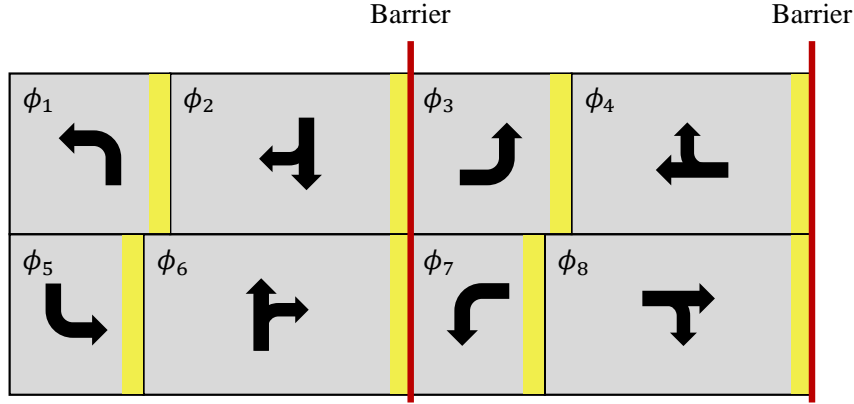


Figure 2.2: An eight-phase NEMA ring barrier diagram for a four-legged intersection.

A *signal plan* for an intersection is defined by a sequence  $((\phi_1, t_1) \dots (\phi_P, t_P))$  of phases and the time allocated to each phase. In a signal plan, the time allocated for each phase is divided into three intervals: *green*  $G_i$  (i.e. a green light, where vehicles can pass freely), *yellow change*  $Y_i$  (i.e. a yellow light, where vehicles are warned of a coming change in right of way), and *red clearance*  $R_i$  (i.e. a red light, which occurs before the start of the conflicting green interval to ensure that the intersection clears properly) [203]. A typical four-legged intersection follows the eight-phase *NEMA phasing* sequence, which is shown in the *ring barrier diagram* in Figure 2.2. The phase sequence consists of two rings of four phases, each of which serves one road (i.e.  $\phi_{1,2,5,6}$  serve the north-south road, while  $\phi_{3,4,7,8}$  serve the east-west road). Pairs of non-conflicting phases are usually signalled simultaneously [302, 435]. For instance,  $\phi_1$  does not create traffic conflicts with either  $\phi_5$  or  $\phi_6$ , and thus can overlap with either one; I denote these combinations  $\phi_{1,5}$  and  $\phi_{1,6}$ . However, the two rings conflict with each other, so the ends of  $\phi_{2,6}$  and  $\phi_{4,8}$  must be synced.

The task of TSC is to find a signal plan for each intersection that optimises certain *signal performance measures* (SPMs). Since vehicles move through the entire road network and improving throughput requires coordinating all of the signals they encounter, TSC is a centralised problem. Some typical SPMs used by stakeholders in TSC include (see also Koohy et al. [202]):

- **Queue length**  $q$  (vehs), the number of vehicles that are stopped at an intersection on red.
- **Queue density**  $q$  (vehs / m), the number of stopped vehicles normalised by road length.
- **Pressure**  $p$  (vehs), the difference in the queue length (or density) between an intersection's incoming and outgoing lanes. Maximising pressure provably maximises throughput [433].
- **Delay**  $d$  (veh-hrs), the amount of time that vehicles cumulatively spend stopped at an intersection. This is equal to the difference in cumulative travel time between the theoretical optimum (all vehicles moving at free-flow speed) and the actual time.
- **Flow rate** or **throughput**  $s$  (vehs / hr), the rate at which vehicles traverse an intersection.
- **Percentage on green**  $g$ , the proportion of traffic that arrives at an intersection during a green light and thus does not incur signalling-induced delay.

## 2.2.2 Detection Hardware

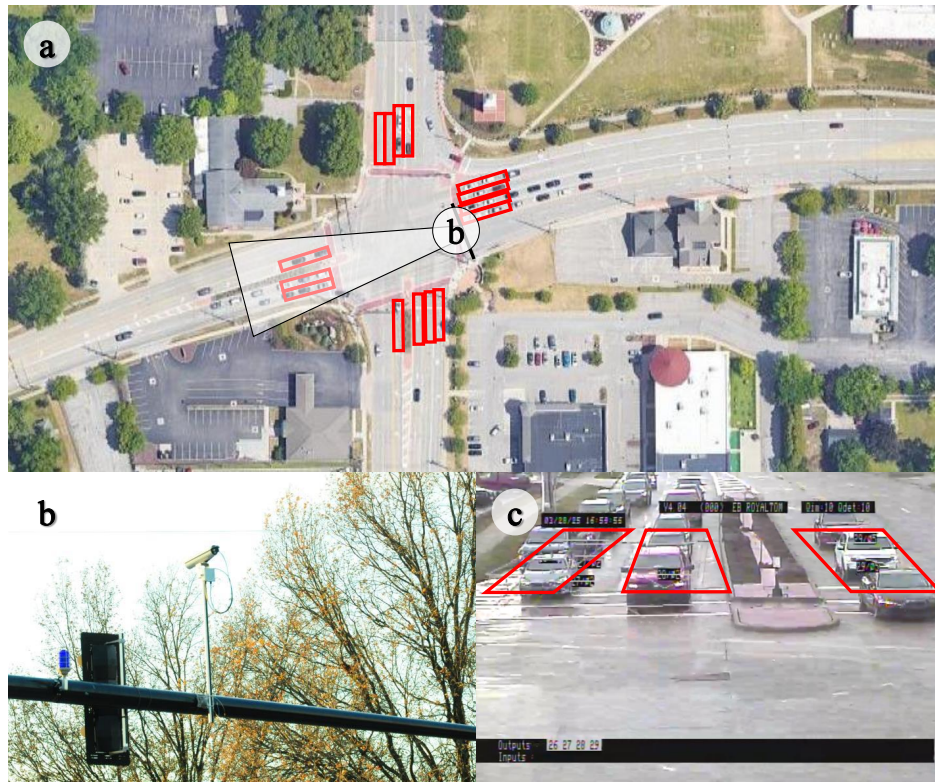


Figure 2.3: Traffic conditions in Strongsville, Ohio, at the intersection of US 42 (north-south) and SR 82 (east-west). (a) Satellite image with shapes of camera detection zones marked in red; (b) photograph of eastbound camera detector; (c) footage from detector with detection zones in red.

How do TSC algorithms interface with the environment? They must first perceive the road state through *detection* hardware, and then alter the signal state accordingly through *control* hardware. Detection hardware usually consists of three types of detectors in the United States.

**Loop Detectors** An induction loop detector consists of a loop of wire embedded in the pavement (i.e. it is an *intrusive* detector), which is actuated when a vehicle passes over it. Loop detectors are relatively robust to the environment. However, actuation depends on the detector’s sensitivity, which is difficult to configure accurately and may result in undercounting or overcounting [215]. Overcounting can occur due to excessive sensitivity to adjacent lanes (splashover) or detector interference (chattering) [216]. Furthermore, loop detectors are highly vulnerable to wear and tear [136]. The maintenance costs of induction loops are the highest among different detectors, and their maintenance induces additional externalities by preventing the use of the road [196]. For this reason, many municipalities are replacing failing loop detectors that have exceeded their lifespans with *non-intrusive* detectors, such as camera and radar detectors [319, 388].

**Camera Detectors** A camera detector is typically mounted in a fixed position above the roadway, and detection zones are placed on the camera’s field of view (Figure 2.3). In the United States, the AutoScope vehicle detection algorithm [273] is used for many cameras. It extracts features to label each detection zone as being in one of three discrete states: “background”, “uncertain”, and “vehicle”. The vehicle counts generated by this algorithm do not reflect actual vehicle volumes — they simply are the number of times each detection zone was actuated. Detection zones must be carefully configured based on their intended purpose. In most cases, they only need to detect vehicles’ presence or absence. However, if their purpose is accurate counting, the detection zones cannot be large. Otherwise, consecutive vehicles may continuously actuate a detection zone, leading to undercounting. Vehicles that are far from the camera may also be indistinguishable [275].

Inclement environmental conditions also contribute to inaccuracy in camera detector counts. Darkness (due to nighttime or fog) and precipitation (such as rain or snow, which cause glare) obfuscate the visual signal of vehicles, making it more difficult for vehicle detection algorithms to isolate them from the background [255, 265, 274]. Even under nominal conditions, shadows can also result in false detections [333]. Finally, detections may be missed on high-speed roads [393].

**Radar Detectors** Radar detectors, which are also mounted above the road surface, can be used to supplement camera detectors. They mitigate some of the limitations of camera detectors and require less maintenance [196]. However, they are still affected by adverse environmental conditions, including oscillation caused by wind [267].

### 2.2.3 Control Hardware and Algorithms

Control hardware is usually installed in a controller cabinet, in which a microprocessor-based *controller* ingests detector inputs and sends control commands to signals. Historically, many controller standards have existed in the United States, but in the early 2000s they were unified by the US Department of Transportation into the Advanced Transportation Controller (ATC) Type 2070 controller standard [327]. Type 2070 controllers are widely used throughout the United States today. In addition to controllers, cabinets typically also include detector processing hardware, load switches that interface between controllers and signals, and a malfunction management unit (or conflict monitor) that performs postprocessing checks as a fail-safe [300].

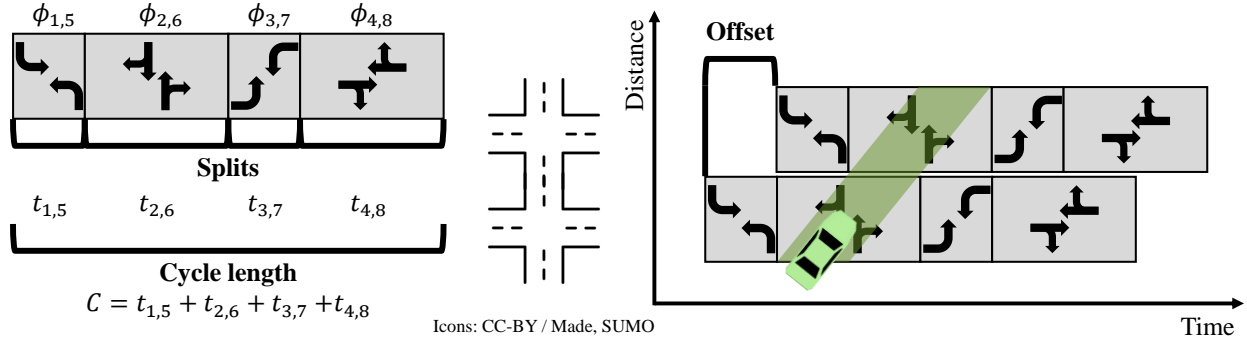


Figure 2.4: Illustration of cycle-offset-split plans for intersections sharing a four-phase signal plan. (L) Diagram of cycle length and splits. (R) Time-space diagram showing effect of offsets. If the offset between two intersections corresponds exactly to a vehicle’s free-flow travel time between them, green progression can be achieved.

The most common signal plans implemented by controllers are cyclic, looping through fixed phase sequences in cycles. According to the NTCIP standard, such *cycle-offset-split* (COS) plans are defined by three sets of parameters (Figure 2.4): the overall *cycle length*  $C$ , *splits*  $t_\phi$  (proportions of time allocated to each phase), and *offsets*  $o_\phi$  (timing offsets between adjacent intersections to achieve coordination). For offsets, one intersection’s controller is usually selected as the master clock. In coordinated control, all intersections in a traffic corridor must share the same cycle length  $C$  and splits  $t_\phi$ , with  $\sum_\phi t_\phi = C$ . The offset must also correspond to the travel time between adjacent intersections at free-flow speed. This has the effect of creating *green progression*, where a vehicle that arrives at the first signal in a corridor on green can continue uninterrupted [203].

Controllers support three main types of control algorithms:

- In *fixed-time control*, a single signal plan or a small number of time-of-day plans are optimised by traffic engineers based on historical traffic volumes. Cycle length optimisation typically depends on setting a desired critical intensity ratio (i.e. the portion of capacity that should be used) for the phase with the highest traffic intensity. The splits, in turn, are selected to balance this ratio between the remaining phases [259]. Software such as Synchro and TRANSYT are used to perform these computations offline [326].
- In *actuated control*, a signal plan is adjusted based on detector inputs through a fixed set of logical rules. An actuated signal plan typically assigns a *minimum green time* to each phase. Each *call* from a detected vehicle lengthens the amount of time available for a phase up to a *maximum green time*, at which point the phase *maxes out*. Alternatively, if no vehicles are detected after an interval known as the *passage time*, the phase *gaps out* [203]. These parameters are, again, usually optimised through offline computation.
- In *adaptive control*, a signal plan is adjusted based on more complex optimisation algorithms, such that the mapping from detector inputs to plans can vary over time [115, 141]. These algorithms are usually too complex to be implemented in controllers directly; instead, they are executed in the cloud, and the resulting signal plans are sent to controllers. Software such as SCOOT [382, 383] and Edaptive [112] are used to perform these computations.

Deployed adaptive algorithms are largely based on local search and optimisation techniques to make incremental changes to cycles, splits, and offsets of predefined time-of-day plans [92, 93], and focus more on leveraging domain-specific heuristics than searching for a globally optimal solution. For example, many deployed offset optimisation methods are based on the link-pivoting combination method of Day and Bullock [92], which greedily chooses offsets for additional intersections in an ordering based on the road network topology. Although these algorithms are designed to minimally disrupt existing workflows for fixed-time control [363], it is not clear whether they leverage the full potential of adaptive signal control.

## 2.2.4 Traffic Simulation

Before an optimised signal plan can be deployed, its quality must be evaluated. However, iterative evaluation and refinement of a signal plan is usually not possible in the real world, because deploying a suboptimal signal plan would result in efficiency and safety costs [128]. *Traffic simulations* provide a safe sandbox within which the performance of signal plans can be assessed and refined. If appropriately constructed based on real-world data, they can serve as digital twins capable of closely capturing the impacts of signal plans on physical traffic.

Traffic simulations exist on a spectrum of how granularly they emulate reality [25, 108]. *Microscopic* simulators provide the most granular simulations, as they simulate the behaviour of individual vehicles (including acceleration, deceleration, and lane changing). This behaviour is generally based on car-following models, which model the acceleration of a vehicle as a time series that depends on its own speed and other vehicles' speeds. *Macroscopic* simulators, by contrast, aggregate vehicles into flows, with time series describing the volume, speed, and density of flows between different points in a road network. *Mesoscopic* simulators provide an intermediate solution between microscopic and macroscopic simulators that balances detail and computational efficiency; they model individual vehicles as flows, or organise groups of vehicles into platoons.

What traffic simulators are used in the AI literature and in industry practice? Noaeen et al. [302] reviewed 160 papers on reinforcement learning-based TSC; they found that the most popular simulators are SUMO [10], an open-source simulator; VISSIM, a proprietary simulator; PARAMICS, a proprietary simulator; GLD [446], an open-source simulator; and AIMSUN, a proprietary simulator. Some of these simulators allow for mixed microscopic/mesoscopic simulations. All of these simulators are designed to be realistic; proprietary simulators generally model more features [341, 412]. Researchers have also created traffic simulators dedicated to training AI methods. On the one hand, Zhang et al. [482]'s CityFlow simplifies traffic simulations greatly in exchange for a 20-fold speedup over SUMO; Chapter 4 compares these two simulators. On the other hand, Garg et al. [128]'s Traffic3D simulates environmental perturbations in a detailed 3D environment to increase the robustness of trained agents. Based on my interaction with industry professionals, traffic engineers typically use VISSIM; VISUM, its macroscopic counterpart; and SimTraffic, a microscopic simulator used by Synchro to perform macroscopic calculations.

What about the traffic simulations executed in these simulators? Noaeen et al. [302]'s review found that 62% of papers relied exclusively on synthetically generated simulations, while another 34% of papers included simulations generated based on real-world data. Chapter 5 discusses different methods for generating data-driven traffic simulations.

## 2.3 Reinforcement Learning

One emerging approach to adaptive TSC is *reinforcement learning* (RL). RL is a paradigm for sequential decision making in which agents learn how to act through trial-and-error interactions with an environment. The goal of RL is to learn *policies*, which describe how agents should act conditioned on the current state of the environment.

Early work in RL during the 1980s and 1990s, including the seminal *Q*-learning algorithm [430], relied on tabular enumeration of environment states and agent actions. RL remained relatively difficult to scale until the emergence of neural networks and other function approximation methods in the 2010s, which led to *deep RL* [281]. Since then, the popularity and complexity of RL have experienced explosive growth. Game-playing deep RL agents have achieved superhuman performance in card games and video games with high-dimensional state and action spaces and real-time decision making, such as AlphaGo (Go) [371], Libratus (heads-up no-limit poker) [50], and AlphaStar (StarCraft II) [418]. Deep RL has also found novel applications in practical domains such as robotics, natural language processing, finance, and healthcare [226]. Transportation has been one of the most significant applications of deep RL, with tasks including autonomous driving [195], vehicle dispatching [315] and routing [297], and TSC. I refer the reader to Sutton and Barto [396] for an in-depth review of the history of RL.

### 2.3.1 Markov Decision Processes

The most common sequential decision-making problem formulation for RL is the *Markov decision process* (MDP), which can be described as a tuple  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$  [394]. In an MDP, a single agent interacts with an environment, which usually consists of a single intersection in the TSC setting, over a number of discrete timesteps  $t \in \{0..T\}$ . At each timestep  $t$ :

**States** The agent receives a representation of the current environment *state*  $s_t \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of all possible states. Before the agent has taken any actions, the first state in which it finds itself follows an *initial state distribution*  $\rho_0 : \mathcal{S} \rightarrow [0, 1]$  [394]. In the TSC setting, the state usually entails an abstract, numerical representation of a single intersection. As reviewed by Noaeen et al. [302], five of the most common state features in RL for TSC are (1) the queue length in each lane (38%), (2) the current phase (11%), (3) the total vehicle count in each lane (10%), (4) the positions of vehicles (6%), and (5) the speeds of vehicles (6%). A minority of algorithms (3%) directly apply deep image processing techniques to frames from camera detectors or simulators.

**Actions** Based on  $s_t$ , the agent picks an *action*  $a_t \in \mathcal{A}$ , where  $\mathcal{A}$  is the set of all possible actions.

In the TSC setting, the action is a signalling decision. Noaeen et al. [302] found that, in a majority of algorithms (62%), the action space consists of the index and/or duration of the next phase — which is more myopic than the COS plans discussed in Section 2.2.3. Other algorithms (32%) are based more explicitly on COS plans, with the action space involving the splits (i.e. the length and sequence of phases) or the cycle length.

**Transitions** The action  $a_t$  affects the environment immediately through a probabilistic transition to the next state. This is modelled by the *state transition function*  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ . Given  $s_t, a_t$ , the environment transitions to each possible next state  $s_{t+1} \in \mathcal{S}$  with probability

$$P(s_t, a_t, s_{t+1}) = \Pr(s_{t+1} \mid s_t, a_t),$$

where  $\sum_{s_{t+1} \in \mathcal{S}} P(s_t, a_t, s_{t+1}) = 1, \forall s_t, a_t$ . RL methods can be divided into two categories based on how they handle state transitions. (1) *Model-based RL* first learns a model of  $P$  before training the agent, usually by fitting it to observed transition probabilities. (2) *Model-free RL* does not learn a model of  $P$ ; instead,  $P$  is considered a notional component of the MDP, and the agent is optimised in expectation over trajectories sampled from  $P$  [396].

In the TSC setting, most RL algorithms follow a model-free approach [435]. Given that queue lengths and other common state features are effectively unbounded, the typical sizes of the state and action spaces make the explicit representation and learning of  $P$  prohibitively costly.

**Rewards** After the agent takes action  $a_t$ , the environment also gives the agent a numerical *reward*  $r_t = R(s_t, a_t) \in \mathbb{R}$ , following a *reward function*  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . It uses rewards to learn how good actions are in various states, so that it learns to take the best action in all states [394].

In the TSC setting, the reward typically denotes the effect of the signalling action on the state (per Noaen et al. [302], 30% use queue lengths; 6% use vehicle counts), or on vehicle-specific quality metrics (13% use the delays of vehicles, in terms of increase in travel time; 9% use the waiting times of vehicles; 4% use the throughput of intersections).

**Policies** The goal of the agent is to learn an optimal *policy*  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maps the current state to the action that it should take,  $a_t = \pi(s_t)$  (for a deterministic policy), or a distribution over actions,  $a_t \sim \pi(a_t \mid s_t)$  (for a stochastic policy). As this is a sequential problem, the agent cannot greedily choose actions to maximise estimated rewards at every timestep, because its actions may have persistent effects (for instance, deciding not to clear a queue in timestep  $t$  may lead to congestion in timestep  $t + 1$ ). Therefore, the agent’s objective is to optimise its *expected return*

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, \pi(s_t))} \left[ \sum_{t=0}^T \gamma^t R(s_t, \pi(s_t)) \right] := \mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, \pi(s_t))} \mathcal{R},$$

where  $\gamma \in [0, 1]$  is a *discount factor*. The higher the discount factor, the more the agent optimises for rewards in future timesteps; a completely myopic agent has  $\gamma = 0$ , and an agent that weights rewards from all timesteps equally has  $\gamma = 1$  [394].

Across many trajectories, a policy  $\pi$  converges to a discounted *state visitation distribution*, which is given by  $\rho_{\pi}(s) = \sum_{t=0}^T \gamma^t \Pr(s_t = s)$ ; this is not normalised. For a deterministic policy, the state visitation distribution is determined solely by  $\rho_0$  (the initial state distribution) and  $\pi$ .



**Q-Values** RL agents need to evaluate the quality of a state, or of an action in some state. These notions are encoded by the *value function*  $V(s)$  and the *state-action value function*  $Q(s, a)$ :

$$V^\pi(s) = \mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, a_t)} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$Q^\pi(s, a) = \mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, a_t)} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a, a_t = \pi(s_t), \forall t \in \{1..T\} \right]$$

Intuitively, the value function represents the expected return of beginning in state  $s_t$ , and then following  $\pi$  from then on; the  $Q$ -function represents the expected return of beginning in state  $s_t$ , *taking action*  $a_t$ , and then following  $\pi$  from then on [394]. At optimality, the value function and  $Q$ -function follow the *Bellman equations*:

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} P(s, a, s') V^\pi(s')$$

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} P(s, a, s') \max_{a'} Q(s', a').$$

**RL Algorithms** RL algorithms can be broadly divided into two categories.

*Q-learning* algorithms iteratively learn  $Q$  by minimising the loss of the Bellman equation:

$$Q_{(t)}^\pi(s, a) = (1 - \alpha) Q_{(t-1)}^\pi(s, a) + \alpha \left( R(s, a) + \gamma \sum_{s'} P(s, a, s') \max_{a'} Q_{(t-1)}^\pi(s', a') \right),$$

where  $\alpha$  is a learning rate. Over many iterations, the algorithm converges to a fixed point  $Q^* = T(Q^*)$  (where  $T$  is an operator for the right hand side of the Bellman equation). Then, the optimal policy is simply given by maximising the learnt  $Q$ -function,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  [397].

*Policy gradient* algorithms explicitly learn a parameterised form for the policy  $\pi_\theta$  by optimising its gradient with respect to the  $Q$ -value, which is given by the *policy gradient theorem*:

$$J_\theta = \mathbb{E}_{s, a \sim \rho^{\pi_\theta}} Q^{\pi_\theta}(s, a)$$

$$\nabla_\theta J_\theta = \mathbb{E}_{s, a \sim \rho^{\pi_\theta}} \nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a).$$

The policy gradient theorem facilitates optimisation by moving the expectation outside the gradient. However, the task of estimating  $Q^{\pi_\theta}(s, a)$  remains. In the basic *REINFORCE* algorithm,  $Q^{\pi_\theta}(s, a)$  is directly estimated from rewards observed in trajectories. Meanwhile, *actor-critic* algorithms combine policy gradient and  $Q$ -learning algorithms by training a *critic* — an estimator of  $Q$  — alongside the policy, or the *actor* [395].



### 2.3.2 Multi-Agent MDPs

The *Markov game* generalises MDPs to the setting of *multi-agent reinforcement learning* (MARL). In MARL, multiple agents exist and interact with each other in an environment; in the TSC setting, this can represent a multi-intersection road network, where each intersection is controlled by one agent [434, 435]. A Markov game has a set of agents  $\{1..N\}$ , and each agent has an action space  $\mathcal{A}_i$  with a *joint action* space  $\mathcal{A} = \prod_i \mathcal{A}_i$ . Likewise, each agent has its own policy  $\pi_i$ , which can be composed into a *policy profile*  $\pi = (\pi_1, \dots, \pi_N)$ . I refer to a policy profile excluding agent  $i$  as  $\pi_{-i}$ . At each timestep, each agent  $i$  simultaneously executes its own action  $a_t^{(i)}$  and receives its own reward  $r_t^{(i)} = R_i(s_t, a_t)$ . Thus, a Markov game can be described using a tuple  $(N, \mathcal{S}, \mathcal{A}, P, R, \gamma)$ .

In MARL, the optimal policy for an agent will depend on knowledge of the global state, and of the actions of other agents in the road network. It is possible for all policies to be learnt in a *centralised* fashion, such that the problem reduces to a single “meta-agent” observing the global state and taking joint actions. They can also be learnt in a fully *decentralised* fashion for each agent, where the other agents are viewed as part of the environment. Centralised training can be difficult to scale due to the high dimensionality of the state and action spaces, but decentralised training can result in suboptimal policies where agents fail to account for each other’s behaviour. *Centralised training, decentralised execution* (CTDE) actor-critic algorithms help mitigate these issues: the critic is centralised, but is only used during training to help coordinate all agents, while the actors for each agent are decentralised during execution [125, 243].

Coordination can also be introduced as part of the problem formulation. In Markov games where the agents are *cooperative* — as is the case in TSC — a common modification is to share the reward function between all agents:  $R = R_1 = \dots = R_N$ ;  $V$  and  $Q$  are thus also shared [485]. For such Markov games, *value decomposition* algorithms still train individual  $Q$ -functions  $\tilde{Q}_i$  to learn each agent’s contributions to the overall value, but constrain  $Q = \sum_i \tilde{Q}_i$  [325].

In *mixed cooperative-competitive* Markov games, there may be teams of cooperating agents, but not all agents share the same reward function [485]. Chapter 7 considers this setting. Here, the agent set is partitioned into a set of disjoint teams  $\mathcal{T}_1, \dots, \mathcal{T}_L \subseteq [N]$ , and the environment yields a joint reward to each team  $\mathcal{T}_\ell$  according to  $R_\ell : \mathcal{S} \times \prod_{i \in \mathcal{T}_\ell} \mathcal{A}_i \rightarrow \mathbb{R}$ . Each agent maximises its team’s total expected return  $\mathcal{R}_\ell = \sum_{t=0}^T \gamma^t r_\ell^t$ . I refer to the joint policy profile of team  $\mathcal{T}_\ell$  as  $\pi_\ell = (\pi_i, \forall i \in \mathcal{T}_\ell)$ , and a joint policy profile excluding all agents in  $\mathcal{T}_\ell$  as  $\pi_{-\ell} = (\pi_i, \forall i \notin \mathcal{T}_\ell)$ .

### 2.3.3 Partially Observable MDPs

In *partially observable MDPs* (POMDPs), the agent does not directly observe the state  $s_t$ . Instead, there is a space of *observations*  $\mathcal{O}$ . The agent’s observations  $o_t \in \mathcal{O}$  are assumed to be samples from some state-dependent probability distribution  $O(s_t, o_t) = \Pr(o_t \mid s_t)$ , with  $\sum_{o_t \in \mathcal{O}} O(s_t, o_t) = 1, \forall s_t$ . The agent never knows exactly what state it is in; instead, it maintains a probability distribution over states known as a *belief state*,  $b_t \in \Delta^{|\mathcal{S}|}$ , which is based on an initial belief  $b_0 \in \Delta^{|\mathcal{S}|}$ . Thus, a POMDP can be described using a tuple  $(\mathcal{S}, \mathcal{A}, P, R, \mathcal{O}, O, b_0, \gamma)$ .

Markov games and POMDPs can also be combined into *partially observable Markov games* or *dec-POMDPs*, where there are multiple agents each having partial observability of the state.

Assuming that the agents are mixed competitive-cooperative, this problem can be described using a tuple  $(N, \mathcal{T}, \mathcal{S}, \mathcal{A}, P, R, \mathcal{O}, O, b_0, \gamma)$ . Generally, partial observability is useful for when the same state may be observed differently by agents depending on roles or on randomness. In the TSC setting, there is a global road network state but not every intersection has access to it; instead, each intersection agent can only observe the state that is local to that intersection [6, 250].

### 2.3.4 RL for TSC

TSC has historically been a popular problem for RL practitioners, as it involves scaling up coordination to a relatively large number of agents. RL algorithms for TSC have either focused on improving performance through explicit, centralised coordination, or on improving scalability through implicit, decentralised coordination. This dichotomy already existed before the advent of deep learning in 2015, as reviewed by Yau et al. [463]. For instance, Medina and Benekohal [266] learnt coordinated  $Q$ -functions using the *max-plus* algorithm (a precursor to value decomposition methods), which exchanges messages between agents over a pre-defined *coordination graph*; El-Tantawy et al. [113]’s algorithm involved agents training envisioned policies for their neighbours and best-responding accordingly. Meanwhile, Prabuchandran et al. [313] induced coordination simply by sharing queue length observations between agents.

In 2016, van der Pol and Oliehoek [416] were the first to apply deep RL algorithms to TSC. They used a deep  $Q$ -network (DQN) implementation of the max-plus algorithm with an image-based state representation, and evaluated it on a four-intersection traffic simulation in SUMO. In 2018, Nishi et al. [301] applied *graph neural networks* (GNNs) to RL for TSC; they constructed a graph of traffic volumes for different movements at an intersection and used a GNN to process this state representation. Out of these works arose an abundance of deep RL algorithms for TSC. Among these, two lines of research are particularly notable.

- Wei, Li, and others developed a series of RL algorithms for TSC from 2019 onwards. Some were explicitly coordinated algorithms, such as CoLight [434], which uses a GNN with attention weights; and CoSLight [339], which uses a Transformer to learn a “collaborator matrix” of probabilistic team compositions.

Another line of their work focused on learning scalable representations of traffic states, with coordination as a secondary concern. These included IntelliLight, a DQN algorithm that conditions on the phase [437]; LIT, a similar algorithm that uses queue length as a reward [496]; PressLight, a DQN algorithm that uses pressure as a reward [433]; FRAP, an intersection topology-agnostic state representation based on “competition” between phases [495]; and MPLight, a DQN algorithm that shared the parameters of a FRAP network between all agents and used a pressure-based reward [65]. As the culmination of this line of work, MPLight’s parameter sharing allowed it to scale up to a 2 510-agent simulation and thus attracted significant attention. I make use of MPLight policies in Chapter 7.

In parallel, they developed the CityFlow simulator [482] and a version with vehicle behaviour models based on machine learning, CityFlowER [89]; LibSignal, a library that integrates APIs for SUMO, CityFlow, and various benchmark simulation environments [268]; and LibSignal++, a physical test environment with toy vehicles and traffic signals [492].

- Garg, Chli, and Vogiatzis focused less on algorithm development than on the practicalities of deploying RL for TSC. They first introduced the Traffic3D simulator [128]. Then, they designed a policy architecture that directly receives visual input from the simulator and processes it using *convolutional neural networks* (CNNs), which they trained using a REINFORCE-based RL algorithm [126, 127]. After adopting a parameter-shared actor-critic framework [129], they trained the policy using domain randomisation (see Section 2.4.1) and footage from real camera detectors [130]. Finally, they achieved a small-scale physical deployment in Coventry, UK [85].

In addition to LibSignal, Ault and Sharon [15] created RESCO, another benchmark dataset of various traffic simulations and RL-based TSC algorithms. I use this benchmark in Chapters 4 and 7. Alegre et al. [6] implemented an API that wraps around SUMO for RL algorithms; I modified and extended this wrapper for Chapters 6 and 7.

## 2.4 AI Design and Deployment Challenges

In Chapter 1, I introduced four practically significant challenges that hinder the design and deployment of AI technologies in transportation. In the following sections, I review how uncertainty (Section 2.4.1), heterogeneity (Section 2.4.2), assurance (Section 2.4.3), and coordination (Section 2.4.4) arise in the domains of gig driving and TSC.

### 2.4.1 Uncertainty

One of the central challenges for AI technologies in transportation is that the level of demand and supply is not predictable over time, even by experienced stakeholders. Transportation systems are not closed [425], as the number of users can vary according to recurrent (e.g. rush hours) and non-recurrent (e.g. special events) factors. When the outcomes of AI technologies' decisions do not correspond to reality, stakeholders may lose trust in them.

- For TSC, only some, not all, impacts of signalling actions remain consistent under different traffic conditions. Alegre et al. [6] showed that the performance of fixed RL policies degrades in previously unseen contexts, but also has a high variance. Although time-of-day plans can handle routine traffic variations, incidents such as lane closures and accidents can lead to inefficiency if not addressed rapidly. Yao and Qian [462] showed that traffic incidents can be predicted based on their traffic impacts up to 30 minutes before they are reported.
- For gig driving, the flexibility of driver schedules [67] and the on-demand nature of customer requests [9] introduce volatility into platform dynamics, even if drivers can experientially acquire intuitions and patterns for their work [218]. Platforms impose another challenge: the destinations of trips are withheld from drivers, which leads to significant difficulties in planning due to not knowing the level of demand at the destination [476]. Some drivers even circumvent the platform to directly communicate with passengers [357].

Leaving aside the challenge of predicting future traffic, even observing the current traffic level can be difficult for AI technologies without a complete view of the transportation system.

- For TSC, although most RL algorithms assume perfect state observability [302], I showed in Section 2.2.2 that detectors are often inaccurate. What is the impact of this inaccuracy on signal performance measures? In simulations, Sunkari et al. [393] significant differences in delay and queue length resulting from different detector types and passage times. Alegre et al. [6] also found that RL performance degrades when the road state is coarsely discretised. This results in *state aliasing*, where states are indistinguishable by the policy.
- For gig driving, the information asymmetry imposed by gig platforms limits visibility into real-time system conditions for both drivers and third-party tools alike [190]. Even Khan et al. [190]’s tool, which analyses web requests from gig platform apps to infer supply and demand in real time, is constrained by the data that the platforms choose to provide to their users. Tenured drivers have become accustomed to this lack of information and have found ways to adapt, but it still diminishes their sense of control [477].

Some strategies that AI technologies have used to address uncertainty include:

- **Robust training.** Injecting uncertainty directly into the training process allows AI algorithms to adapt to their presence. *Domain randomisation* methods optimise the performance of RL algorithms in expectation over different environmental parameters sampled from a distribution. These parameters can include ambient weather conditions [130] or the probability of vehicle detection [292]. *Transfer learning* and *meta-RL* methods go further: they model uncertainty as tasks sampled from an underlying distribution, and generalise policy parameters between different tasks [292, 303]. *Distributional RL* methods learn all quantiles of the return distribution, and optimise for them at the same time [364]. In solving for driver repositioning strategies that maximise earnings, Chaudhari et al. [63] optimise for worst-case expected earnings given the uncertainty in the probabilities of trip origins and destinations.
- **State prediction.** When this uncertainty causes data to be not just noisy but also missing, AI algorithms can be conditioned to behave differently depending on whether data is missing. For instance, Jiang et al. [180]’s TSC algorithm BlindLight predicts the probability that an intersection is “blinded”, and uses it to choose between  $Q$ -networks for normal and blinded intersections. Alternatively, missing data can be imputed. In the TSC setting, traffic prediction models can be integrated with RL-based control. Mei et al. [269] train models to impute missing states and rewards for RL policies from neighbouring intersections. External knowledge can also be used: Guo et al. [149] leverage ontological data, while Da et al. [90] use a *large language model* (LLM) to reason about unknown transition dynamics. In the gig driving setting, models exist to accurately predict short-term surge pricing [30], demand [66], and other information unavailable to drivers.

However, more fundamental questions exist. What part of this uncertainty is useful to model from the perspective of stakeholders? How can AI systems properly convey uncertainty to stakeholders? For TSC, most work has arbitrarily defined distributions for noisy or missing data without empirical validation. In Chapter 5, I investigate just how much uncertainty in traffic volume arises from detector data. For gig driving, the data that is unavailable to drivers is better defined, but it is unclear how this missing data impacts drivers’ decision-making processes. In Chapter 3, I investigate how exposing the presence of uncertainty impacts gig drivers’ trust in AI.

## 2.4.2 Heterogeneity

Variation between individual stakeholders can impact their responses to AI technologies in transportation, and the effects of this variation on system-level outcomes must be accounted for. Otherwise, stakeholders who do not fit the typical profile of training data may be treated inequitably.

- For TSC, classes of vehicles such as trucks [142], public transit vehicles [439], and gig drivers' vehicles [159] all have disparate impacts on traffic patterns. Within vehicle classes, factors such as driving speed, reaction time, and braking sharpness affect how individual drivers respond to signals, such as the "dilemma zone" for yellow lights [186]. However, most traffic simulations used for RL training abstract away inter-vehicle variation by treating all vehicles as passenger cars, and many of them ignore pedestrian traffic [302, 482].
- For gig driving, drivers have a variety of motivations and habits. Some rely on it as a primary source of income, while others view it as supplementary income [252, 452]; some drive full-time, while others drive part-time [37, 218, 252]. They also respond differently to the information asymmetry and algorithmic management of gig platforms [452, 477]. Although some mechanisms incentivise drivers to report their preferences while maintaining theoretical guarantees [332], they do not account for the diversity of drivers in the real world.

Heterogeneity exists not just at the level of individuals, but also at the level of markets. The performance of AI technologies in transportation can vary between deployments, especially if they make assumptions about environmental properties that hold in some contexts but not in others.

- For TSC, the most common strategy for RL is to train and evaluate policies on a single road network, often compared to only a single benchmark algorithm [268, 302]. Furthermore, as I discussed in Section 2.2.4, most of the simulations used in the literature have been synthetic, not grounded in realistic traffic data [302]. Although the RESCO [15] and Libsignal [268] benchmarks represent useful steps toward homogenising evaluation practices, they have received limited attention outside of their originating research groups.
- For gig driving, it has been consistently recognised that the market entry of platforms has disparate impacts in different cities depending on their socioeconomic makeup [95, 133] and existing transportation infrastructure [133, 360]. Beyond the spatiotemporal locality considered by Ma et al. [247], there has been insufficient consideration of how these dimensions could influence drivers' response to gig platforms across multiple cities, even though the empirical data of Ong et al. [304] suggests that they do have an impact.

Some strategies that AI technologies have used to address heterogeneity include:

- **Inclusive policies.** Optimising for distinct subgroups of users is one way for AI algorithms to account for their presence. In the TSC setting, various RL algorithms assign different weights to different vehicle classes in reward functions. Some prioritise pedestrians [464, 489], while others prioritise public service vehicles such as buses, ambulances, and firetrucks [130, 210, 347]. In the latter category, some RL methods implement *preemption*, where normal signalling patterns are overridden in the presence of emergency vehicles [361, 386]. In the gig driving setting, Liu and Jiang [238] and Zhou et al. [500] personalise route recommendations based on drivers' features and expressed preferences. Their work is complemented by that of Di et al. [101], who cluster drivers by shift locations and times.

- **Cross-context generalisation.** Incorporating and differentiating heterogeneous environmental contexts into training can allow AI algorithms to better adapt to them. In the TSC setting, generalised state representations such as those of FRAP [495] and MetaLight [473] allow policies to transfer easily between contexts, while Zhang et al. [483]’s meta-RL algorithm groups traffic flows into tasks based on patterns in states, actions, rewards, and travel time. Additional state features for railways [270], trucks [400], and trams [486] also help to capture heterogeneity. In the gig driving setting, Guo et al. [150, 151] predict gig prices by integrating data from gig platforms themselves, taxi services, and public transport; thus, their models capture relationships between transportation modes.

Similar questions as uncertainty arise for heterogeneity. What is the right way to model heterogeneity, and just how much of an impact does it have? When integrating multiple sources of data or multiple objectives into the training process of AI algorithms, how can stakeholders’ perspectives on their relative importance be captured? In this thesis, I explore the impact of heterogeneity on traffic simulations, upstream of the training process of RL for TSC: Chapter 4 models the heterogeneous behaviour of drivers, while Chapter 5 integrates heterogeneous sources of detector data.

### 2.4.3 Assurance

*Assurance* combines two related notions: *interpretability* and *safety*. First, to ensure that AI-driven decisions are *interpretable* by stakeholders, they must be in alignment with stakeholders’ mental models. This is because AI technologies will be applied to established systems in transportation where stakeholders will have acquired mental models experientially. Without such assurances, stakeholders will be hesitant to collaborate with or delegate decisions to AI technologies.

- For TSC, the acyclic signal plans common in the RL literature [302] are a far cry from the COS signal plans used by traffic engineers, because they only myopically optimise for the next phase. In Section 2.2.3, I noted that a minority of existing RL methods instead optimise cycles, offsets, or splits [183, 254, 365, 458, 474]. However, these methods generally suffer from one of two limitations. Some are heuristic methods that postprocess signalling actions externally to the learning process [183, 458, 474], while others are not scalable because they enumerate the space of valid actions [254, 365].
- For gig driving, the needs and mental models of drivers are known to be a significant influence on the extent to which they trust platforms’ management mechanisms. When drivers’ mental models of these mechanisms misalign with reality, they reject platforms’ behavioural nudges (such as to surge chase) [61, 218, 476]. Although platforms have no incentive to reveal their exact mechanisms, I noted in Section 2.1.3 that drivers engage in sense-making activities [253] that mitigate the impact of this information asymmetry. In doing so, they have developed desiderata for how they want platforms to explain their processes [323, 477].

Second, *safety* constraints must be placed as safeguards against adverse outcomes. Both TSC and gig driving are domains with elevated stakes, where wrong decisions made by AI technologies can have significant impacts on stakeholders. When such impacts occur, mechanisms for tracing failures back to their root causes must also exist as a means of recourse. There is a trade-off between guaranteeing safety (by enforcing constraints) and performance (by loosening constraints).

- For TSC, failing to set appropriate constraints on signal plans can lead to safety issues. Tort cases show that concrete regulatory and mortal risks arise from improper signalling [403], and thus documentation and justification of signal planning are necessary if these records must be subpoenaed in proceedings. Ample scholarship exists regarding how signalling decisions affect crash rates, and these effects have been quantified as *crash modification factors* [249, 378, 450]. However, safety-related objectives are usually incorporated as ad hoc components of the reward function in RL-based TSC, if at all [117, 140, 232, 239, 465].
- For gig driving, failing to accurately estimate potential demand can lead to opportunity costs in terms of drivers’ monetary earnings and time. Guda and Subramanian [148] and Ma et al. [247] have examined opportunity costs in mechanisms that are incurred through mispredictions. However, communicating assurances about potential earnings to drivers remains an open problem. This is because guaranteeing earnings is a difficult task in the presence of uncertainty, even for the gig platforms themselves [263]. When platforms provide guarantees, their mechanisms must include driver payments to compensate for any shortfall [452].

Some strategies that AI technologies have used to address assurance include:

- **Pre-hoc interpretability.** One way to improve interpretability is to ensure that the decision-making processes of AI algorithms are inherently understandable by people. In the TSC setting, RL policies can be designed to explicitly balance multiple objectives in the reward function [57, 140, 479, 480]; to possess monotonicity in state variables [14] and other desirable properties; or even to consist of syntactically correct code, generated through combinatorial search procedures [147]. In the gig driving setting, Zhang et al. [477] and Rao et al. [323] collected suggestions from drivers about specific information disclosures and pricing structures that they would like platforms to implement.
- **Post-hoc explainability.** When it is not possible to make the AI algorithm itself interpretable, different explanation methods can be applied to summarise the algorithm’s decisions instead. In the TSC setting, many machine learning explainability methods can be applied. *Decision trees*, which branch based on logical rules (see Section 7.2), are more interpretable than neural network-based RL policies [282, 370]; Jayawardana et al. [177] train a decision tree to imitate an RL policy for TSC. *SHAP values*, which measure the contributions of features to model outputs, have also been applied to identify important state features in RL for TSC [182, 355, 481]. Decision trees and SHAP values have also been applied to crash modification factors [442]. In the gig driving setting, Li et al. [225] and Zhang et al. [476] used formative studies to design tools that inform and empower gig drivers. In particular, Li et al. [225] envisioned a data-sharing network that would quantitatively surface inequalities in task assignment and compensation.
- **Optimisation constraints.** Explicit guardrails on the behaviour of AI algorithms can help guarantee safety. In the TSC setting, unsafe actions can be masked out from selection by RL policies [291]; see Section 6.2.1 for a longer review of safety constraints in RL. Such methods are more rigorous alternatives to approaches based on reward shaping [206]. In the gig driving setting, Raman et al. [321] and Kumar et al. [209]’s income distribution schemes for gig drivers guarantee minimum earnings, although they considered simplified settings.

For assurance, the state of the art of the AI literature in TSC and gig driving have opposite problems. Many interpretability and safety methods have been implemented for TSC, but they do not satisfy the practical desiderata of scalability and alignment with human mental models. I focus on addressing these limitations in RL for TSC by constraining RL policies with practical signalling constraints (Chapter 6) and by improving the scalability of decision tree imitation learning methods (Chapter 7). On the other hand, driver-centred studies have unearthed a variety of desiderata and potential system designs to improve the interpretability of gig driving, but limited progress has been made towards realising these possibilities as deployed systems.

## 2.4.4 Coordination

To achieve equitable system-level outcomes in transportation, AI technologies must be able to coordinate interactions between many individual agents, all of whom may be self-interested. Although various AI algorithms in transportation assume decentralisation between agents, the resulting behaviour lacks guarantees. Specifically, it is not guaranteed that agents can global coordination by acting locally, except when a target state visitation distribution is given [421]. By contrast, stakeholders will be reluctant to adopt AI technologies for transportation systems if their benefits cannot be realised at the levels of both individual users and the overall system.

- For TSC, the signalling actions taken by individual intersections are not independent, but will affect the traffic state of upstream and downstream intersections. RL algorithms such as CoLight can emergently learn green progression, and examining the weights of the learnt policies shows that they capture intuitive patterns of importance along arterials [432, 433]. However, despite their complexity, these RL-based policies are acyclic and still less synchronised than the COS signal plans that are deployed in practice [1], and the presence of coordination cannot be reliably guaranteed for stakeholders.
- For gig driving, the choices taken by individual drivers to accept or reject rides will also shift the overall distribution of demand and supply. For example, a route planning tool that recommends many drivers to go to the same high-earnings location may instead decrease potential earnings by inducing oversupply. Hence, drivers avoid chasing surges [323]. While mechanisms such as those of Garg and Nazerzadeh [131] and Ma et al. [247] are designed to be incentive-compatible for individual utility-maximising drivers, they fail to provide concrete incentives for drivers to coordinate in a way that satisfies demand [304].

Distributed hierarchical structures of control can help stakeholders reason about how coordination can be achieved, and also help AI agents coordinate on objectives [5, 459]. However, control dynamics across such structures can be unpredictable, and are not always practical.

- For TSC, various RL algorithms have used the graph structure of road networks for coordination — either emergently learning hierarchical representations using graph neural networks [301, 422, 434, 475], or explicitly introducing higher-level controller agents that coordinate intersection agents by setting subgoals or reward signals [2, 94, 250, 474, 493]. However, this work has not focused on the practicalities of deploying hierarchical systems. If these computations must take place in the cloud, communication delays will likely exist [246], and therefore appropriate points of synchronisation must be identified.



- For gig driving, platforms can themselves be viewed as controller agents that direct the actions of driver agents. Is what is best for a driver also what is best for the platform? Zhang et al. [478] studied this problem using a cognitive hierarchy approach, and found that individual-level and system-level outcomes are not always in alignment — especially when the strategic sophistications of agents differ. In turn, competitive dynamics between platforms [78, 181] influence their dispatching mechanisms. This breaks the typical assumption for mechanisms that the platform has complete information on supply and demand [247].

Some strategies that AI technologies have used to address coordination include:

- **Mediation mechanisms.** With which other agents should coordinating agents exchange information? The introduction of a mediation mechanism, implicit or explicit, can help agents better account for others. In the TSC setting, regional *manager agents* can be used to influence the policies of lower-level agents [2, 145, 250]. More implicitly, *coordination graphs* [416, 501] and other collaborator selection mechanisms based on graph attention [339, 362, 434] allow agents to dynamically determine which others to coordinate with. In the gig driving setting, mechanisms have been designed in which a third-party, social welfare-optimising integrator dispatches gigs centrally to individual platforms [24, 429].
- **Communication protocols.** What information should coordinating agents exchange? Standardised communication protocols allow coordinating agents to learn from each other without centralised coordination. In the TSC setting, messages exchanged between agents range from observation embeddings [230, 454, 490] and policy representations [77, 250] to structures that are more explicitly optimised for compactness [45, 498]. In the gig driving setting, d’Orey et al. [106] and Yu et al. [466] designed distributed asynchronous messaging protocols for ridesourcing, in which passengers directly broadcast requests to nearby drivers, and drivers exchange messages to determine which passengers to service.
- **Spatiotemporal partitioning.** When and where should coordinating agents exchange information? In the TSC setting, several RL algorithms have accounted for the fact that not all agents act synchronously, due to factors such as detection and communication latency. These include algorithms based on distributed computation [451], state prediction [305], and delayed observation sharing [454]. In the gig driving setting, Jin et al. [184] modelled a hierarchy of dispatching agents located in geographically adjacent regions. Other mechanisms have modelled driver repositioning as a spatiotemporally localised problem, leveraging the intuition that coordination is only necessary under periodic peaks in demand [64, 162].

Although coordination in outcomes is a desirable property for stakeholders, exactly how it can be achieved is challenging to reason about. How can coordination be guaranteed instead of exhibited as an emergent behaviour? How can misalignment in individual and global objectives be managed? In this thesis, while Chapters 6 and 7 mainly focus on *retaining* the performance of coordinated RL algorithms, I also impose coordination by adding interpretable structure to policies. For example, in Chapter 6, I constrain cycle lengths to be shared by TSC policies; in Chapter 7, I constrain teams of cooperative agents to be trained on the same datasets for imitation learning.

# Chapter 3

## Missing Pieces

### Understanding the Longitudinal Impact of Designs that Expose Uncertainty on Gig Drivers’ Trust in AI Decision Aids

*Domain:* Gig driving  
*Challenges:* Uncertainty

I don’t ask twice.

---

*Pachakutiq (as Philip J. Coulson)*

#### 3.1 Introduction

In this chapter, I use gig driving as a case of how *uncertainty* impacts people’s interactions with AI technologies in transportation. As I discussed in Section 2.1, gig driving is a complex ecosystem. It does not operate on a strict timetable: drivers are free to drive when they wish, and customers are free to submit requests when they wish. It is also deeply interlinked: drivers and customers all reside in a common pool, and their actions can easily influence each other. This means that the level of demand and supply is never perfectly predictable, even if all of these parties can experientially gain knowledge of common patterns [218, 476].

All of these factors make gig driving an ideal context for studying the impact of uncertainty on interactions between AI technologies and their users. In this chapter, I consider a fundamental question that underlies all such interactions. If the goal of AI technologies should be to benefit human stakeholders, what factors would lead a stakeholder to want to use (or not want to use) AI technologies as part of their workflows? More specifically, how does uncertainty impact the way people *trust* and *rely* on AI technologies? To answer this question, I focus on *AI decision aids*, which function by (1) recommending decisions and (2) predicting how good the outcomes of following those decisions will be. When uncertainty impacts the predictability of AI decision aids — as is the case when a driver’s earnings does not meet expectations created by an AI’s predictions — past work has hypothesised that users’ trust in the decision aid will be eroded [176, 377].

Prior literature on trust in AI decision aids under uncertainty can be organised into two complementary lines of work. One line of work has studied specific factors that influence trust in AI decision aids, using laboratory experiments in simulated, single-shot, and low-stakes scenarios that

require limited domain expertise [20, 55, 175, 199, 491]. However, context is an important factor for trust in AI [194, 377]. Therefore, the contrived nature of these experiments limits their generalisability to the real-world use contexts of AI decision aids. Another line of work has studied trust in real-world AI decision aids [33, 194, 422], using qualitative, observational studies to assess how existing users interact with decision aids. These studies are not quantitative assessments of design factors and provide limited insight into how to design new and trustworthy AI decision aids.

Here, I provide a deeper exploration of trust in AI decision aids by combining the strengths of these two lines of work. I contribute the first in situ study of how exposing the uncertainty of an AI decision aid longitudinally impacts users' trust and reliance on the decision aid. Using gig driving as a testbed, I study trust in a real-world, medium-to-high-stakes decision-making scenario where users have existing expertise. Specifically, I comparatively evaluate different designs that expose the potential for misprediction in an AI decision aid. I address the following research questions:

**RESEARCH QUESTION 3.1.** *How do users' trust and reliance on an AI decision aid depend longitudinally on their perception of its predictive accuracy?*

**RESEARCH QUESTION 3.2.** *How do different designs that expose the inherent uncertainty in predictive performance impact users' trust and reliance on an AI decision aid?*

I addressed these questions by conducting a longitudinal user study where  $n = 51$  gig drivers used a schedule recommendation tool, which did not leverage sophisticated AI methods but was designed to emulate interactions with tools that do. By measuring the trust and reliance of participants over repeated interactions, I tested the effects of exposing uncertainty in the tool's predictions through range-based earnings estimates and hedging text. My quantitative and qualitative findings show that participants' initial perceptions of the tool's accuracy improved their trust in it over time. Range-based uncertainty not only improved trust and reliance in single-shot settings, but also strengthened it over repeated interactions; meanwhile, hedging had the opposite effect.

This chapter was published at the ACM Conference on Fairness, Accountability, and Transparency (FAccT) in 2025 [71].

## 3.2 Related Work

### 3.2.1 Trust in AI

#### Defining and Measuring Trust

*Trust* is a critical factor that enables individuals and organisations to collaborate productively, especially in uncertain and volatile situations [104]. As automated systems increasingly replace human roles in collaborative tasks, trust has become a central focus of research on human-AI interaction. Empirical studies in various domains have highlighted the role of trust in the acceptance of AI systems by their users [75, 192, 488]. Trust in AI is particularly challenging due to the opacity of many AI models [334]. The literature on AI explainability arose out of these concerns; conversely, AI explainability methods have been empirically found to enhance trust [122, 368].

A lack of uniformity exists in the human-AI interaction literature on how to define and evaluate

trust in AI systems. Kohn et al. [200] and Ueno et al. [411] both conducted reviews of trust measures, the latter specifically for human-AI interaction papers. Both reviews found a great diversity of measures and contexts, but also that experimental design choices were seldom informed by or validated against models of trust. I follow Mayer et al. [262] in operationalising and distinguishing the constructs of trust and *reliance*: in the context of AI decision aids, they hypothesise that trust is “the willingness of a person to be vulnerable to the actions of an [AI decision aid], based on the expectation that it will perform a [decision-making task] important to the trustor”, and that reliance is the external, behavioural expression of that internal attitude [377].

One point of broad consensus in the literature is that the design of AI systems impacts expectations of their performance, and thus the trust of their users [52, 91, 175, 199, 211, 350, 467]. However, the majority of this work has been based on controlled laboratory experiments. Compared to real-world use contexts of AI decision aids, such experiments have two main limitations.

### Trust Over Repeated Interactions

First, past experimental evaluations have largely been *single-shot*, involving only a single session of AI use with no temporal separation between decision points. Ueno et al. [411] reported that most measures of trust were evaluated using questionnaires, and most works evaluated trust only once. However, in the real world, users rarely interact with AI systems only once. As users gain experience with systems, the dynamics of trust between them will also change. Lee and Moray [217] identified three bases of trust that influence the perceived trustworthiness of automated systems by users: *performance* (i.e. perceived task competency), *process* (i.e. perceived transparency of behaviour), and *purpose* (perceived helpfulness and identity of benefactors). Solberg et al. [377] hypothesised that trust in AI decision aids is initially purpose-based; this gives way to performance-based and process-based trust as users experientially develop their perceptions of the AIs. Factors external to the interaction process may also impact trust dynamics, including changes in the environment and in the AI systems themselves (e.g. through model updates [428]).

Most studies that measured trust across multiple interactions have either ignored or attempted to control for learning effects (e.g. through randomisation of scenario ordering) [91, 294, 447]. Several studies have explicitly compared *two* temporally separated trust measurements. Mou and Cohen [287] and Mou et al. [288] measured user trust in health e-services during two sessions separated by five weeks; Kunkel et al. [211] compared trust in human-generated and item similarity-based movie recommendations after two recommendations separated by two weeks. However, these studies still provide a limited perspective on the evolution of trust dynamics. Unlike all of these works, I explicitly test the effects of different designs within a more complicated setting, which includes multiple interactions as well as a changing environmental context.

### Trust in High-Stakes Situations

Second, past experimental evaluations have largely been *low-stakes*, involving contrived or hypothetical decision-making scenarios in which an element of *risk* and thus vulnerability is largely absent [194]. Under Solberg et al. [377]’s hypothesis, the level of risk in a given context influences the extent to which a user trusts an AI decision aid — even if the user trusts the decision

aid in general, they may still feel that trusting it in risky scenarios may lead to adverse outcomes. Schoeffer et al. [353] and Yurrita et al. [471] provided empirical support for this. Scharowski et al. [350] showed that the influence of risk on trust can be modulated by the design of AI systems; they found that user trust in a loan approval AI system significantly increased in the presence of a “certification label”, with a greater increase under a high-stakes context.

In past work, some experiments have studied AI decision aids for decision-making domains entailing high stakes in the real world, such as the medical [52, 55, 33, 175, 422, 374, 453, 461] and financial [312, 350, 353] contexts. However, for practical and ethical reasons, participants in these studies cannot receive feedback from the real world for their decisions. A minority of work has evaluated trust in AI decision aids within their real-world contexts [33, 194, 422, 444], but these have been limited to observational studies that did not compare multiple designs. I contribute a mixed-methods study that assesses the effects of different designs on trust under financial risk, one of the nine different risk domains in a taxonomy introduced by Stuck et al. [385].

## Uncertainty and Trust

*Uncertainty* encapsulates sources of variability that make it difficult for users to reason about the outcomes of relying on an AI, thus increasing the risk of this reliance [385]. AI systems are rarely guaranteed to make perfect decisions, but inherent opacity in their designs makes it difficult for users to ascertain whether good performance can be expected. Various experimental studies have tested designs that improve users’ awareness of the presence and impact of uncertainty on AI [20, 60, 199, 295, 427, 457]. This work is grounded in *trust calibration* from the AI explainability literature, which aims to give users realistic expectations regarding when and why AI systems may or may not perform well [39, 407, 491].

Trust calibration allows users to modulate their reliance on AI decision aids, potentially by rejecting their outputs in some contexts [377]. By contrast, *trust enhancement* aims to uniformly improve reliance, which may lead to negative outcomes [491]. An important consideration in trust calibration is whether the imperfection of AI systems arises from *aleatoric* uncertainty (i.e. environmental variability) or *epistemic* uncertainty (i.e. model limitations) [39, 407]. I consider a context where users have existing mental models of how aleatoric uncertainty affects them, but are experiencing the impacts of epistemic uncertainty anew. For both trust calibration and enhancement, prior experiments have lacked real-world context, being limited to hypothetical or simple tasks where participants require little background knowledge.

One thus far underexplored dimension in trust calibration is the use of *lexical hedging*: verbiage that expresses uncertainty. Kim et al. [193] assessed the effects of lexical hedging in a large language model’s medical answers on trust but not on reliance. Zhang et al. [487] measured trust and reliance on an AI decision aid with lexical hedging for a contrived shape identification task. I perform a real-world evaluation of two designs for presenting uncertainty in an AI decision aid: presenting scalar ranges for estimates, and qualifying estimates with lexical hedging.

### 3.2.2 Recommendation Systems

Recommendation systems are mediators of the interfaces between people and the technologies they use. They provide value to their users by mitigating the burden of decision overload and helping them make better choices [47]; however, they have also benefitted technology providers by driving user engagement and financial opportunities (e.g. [309]). As is the case for AI decision aids, users' innate propensity for trust and domain knowledge also modulate their trust in recommendation systems [58, 197]. Although their primary goal is to learn and accurately model the preferences of users, recommendation systems are also persuasive in that their outputs can feed back into user preferences [143]. Much work has focused on different longitudinal effects that arise from such feedback loops. On the negative side, recommendation systems can create bias in user ratings [499], concentration in the set of recommended items [123], and degradation in system accuracy [484]. On the positive side, temporal diversity objectives [214] and nudging-based designs [231] can create trust [367] and lead to further exploration of content [231]. However, item preference is ultimately a subjective measure that can vary considerably between users [317]. In this chapter, I explore the feedback loop that arises from a more objective measure, the earnings of gig drivers.

## 3.3 Pilot Interview Study

To perform an ecologically valid study of trust and reliance in an AI decision aid among gig drivers, I needed to first design a decision aid with practical utility. What kind of AI decision aid would be most relevant to gig drivers? From Section 2.1, one of the biggest challenges faced by gig drivers is the volatility of platform dynamics, which makes it difficult for them to plan their driving activity so as to maximise their profits. I focus on one aspect of planning for gig drivers: choosing *when* to work, subject to their constraints and preferences. There is considerable variation in drivers' habits along this dimension [37, 218, 251]. Choosing *where* to work is another key aspect of planning [476], but I limited my study of uncertainty in this multi-objective problem to a single dimension.

Based on these insights, I concluded that gig drivers would find practical utility in a schedule recommendation tool. Like other AI decision aids [377], such a tool would (1) recommend a set of decisions (i.e. a schedule) to achieve a set objective (maximising their profits), and (2) predict the outcomes (i.e. estimated earnings) of following those decisions. I began my study by creating a prototype of such a tool. Next, I conducted a series of pilot interviews to understand gig drivers' needs and how well the prototype aligned with them. Through this process, I refined the design of the decision aid so that gig drivers would be more likely to find it useful in their daily decision making. The interview methodology was approved by our Institutional Review Board (IRB).

### 3.3.1 Prototype Design

Following common practice in UX design [121], I used Figma [124] to create the prototype design. The prototype interface consists of two pages.

First, a **constraint page** (Figure 3.3) elicits scheduling constraints from the user, i.e. when they are available or unavailable during the week on an hourly basis. These constraints would need to

be imposed on any schedule generated by the tool. Unlike Zhang et al. [476], I account for the fact that users' availability may change between days. Additionally, the page also asks users to set a target for themselves — either a maximum number of hours to work, or a minimum amount of money to earn, on a daily basis or on a weekly basis. I considered these as likely goals for drivers to have, based on prior driver studies [476, 477]. Due to technical limitations, these questions were implemented as dropdowns, with discretised options used in place of text boxes. However, I consider the impact of this on ecological validity to be minor since I controlled the page.

Second, a **schedule page** (Section 3.4.3) presents the tool's recommended schedule to the user. It was kept static to gather more uniform feedback from participants. The page includes two components. First, it displays a range of estimated weekly earnings, consisting of mean ("On an average week..."), pessimistic ("On a bad week..."), and optimistic ("On a good week...") earnings. Second, it shows a tabular schedule with estimated earnings for each hour of the week, and highlights time slots that the tool recommends for the user to drive during.

### 3.3.2 Methodology

Participants began by completing a web-based consent form and a demographic survey that collected their age, gender, and education level. After completing the web form, participants were invited via email to complete 20–30-minute audio-recorded Zoom interviews conducted by me. Participants were compensated with a \$10 Amazon gift card.

In the first 5–10 minutes, the interview focused on *formatively* understanding drivers' needs and motivations. In the last 15–20 minutes, the interview focused on *evaluatively* understanding how well the tool met drivers' needs. To ensure a consistent experience, the Figma prototype was opened in my browser and shown to participants in a screensharing session. First, on the constraint page, the participant was asked to work with me to interact with the page, entering the constraints as if they were using the tool for their actual planning. Then, on the schedule page, the participant was shown a schedule with mocked earnings estimates. Finally, the participant was asked about their overall opinions of the tool.

Transcripts for the interviews were generated by Zoom. I reviewed and corrected these transcripts, then performed structural coding [343]. Afterwards, I and a co-author separately used QualCoder 3.3 [88] to perform open coding [343] and axial coding [344]. We met to reconcile their codes and construct a unified codebook. Finally, I re-applied the updated codes.

### 3.3.3 Participants

Participants were recruited from the user base of Gridwise, a mobile assistant app for gig drivers, in July 2023. I chose to recruit from this user base to access a relatively large and diverse sample of both historical data and participants. Gridwise distributed recruitment messages to 500 users, but otherwise did not interact with participants. Recipients were sampled from Gridwise users in the United States who had completed at least one gig in a platform linked to the Gridwise app over the week preceding recruitment. I recruited 4 interview participants:

- **P1:** A 39-year-old female with a professional degree who drives exclusively for delivery platforms

- **P2:** A 55-year-old male with less than a high school degree who drives exclusively for delivery platforms
- **P3:** A 29-year-old female with an undergraduate degree who drives more frequently for ridesharing platforms
- **P4:** A 53-year-old male with a professional degree who drives more frequently for delivery platforms

### 3.3.4 Results

The four pilot interview participants reported a diversity of motivations and routines for driving. While all four participants had specific earning goals, P1, P2, and P4 considered their goals to be important and valued the time flexibility of gig work, whereas P3 was more motivated by the opportunity for human interaction. P1, P2 and P3 had typical times that they drive at; however, P1, P2, and P4 also adjusted their schedules based on demand. All four participants had encountered difficulties in planning due to the unpredictability of demand and/or supply (with P1, P2, and P3 feeling that gig platforms provide insufficient information), and indicated that they would find schedule recommendations to be useful.

All four participants found the initial design of the tool to be generally understandable, and felt that it would be useful for drivers in planning their activity. P1 and P2 liked the fact that the tool presents information to them in a way that reduces the need for guesswork while driving. In particular, P1 suggested that the tool would help mitigate a catch-22: it is not possible to view gig demand information in DoorDash without exiting their Dash (scheduled work period), but doing so seemingly deprioritises them.

On the constraint page, P1 and P4 indicated that the questions aligned well with their goals. P1 and P3 suggested that they would not set the constraints to perfectly align with their routines, so as to receive more information from the tool. On the schedule page, all four participants liked the estimated hourly earnings, with P1, P2, and P3 indicating that they would be helpful in deciding whether or not to work at particular times of day. Yet, P2, P3, and P4 acknowledged that the estimates would only be guesses. P1 and P4 also liked the range of weekly earnings, but P3 felt it assumed they would follow the recommended schedule perfectly. P2 and P3 noted that ranges for hourly earnings would be useful to display.

P1, P2, and P3 all felt that it was better for the tool to have a simple, easy-to-use design. All three indicated that the prototype fulfilled this requirement, although P3 suggested that wording and design improvements would be necessary (in particular, they felt that the monotone colour scheme of the tool was confusing). P2 and P3 felt that the design needed to be mobile-friendly. The participants also mentioned other desiderata:



Shift Recommendation Tool

Would you like to set a maximum number of hours that you would like to drive to earn as much money as possible, or a minimum amount of money that you would like to earn in as few hours as possible?

☐ A maximum number of hours to earn as much money as possible

☒ A minimum amount of money to earn in as few hours as possible

Would you like to set these constraints for each day individually, or for the entire week?

☒ For each day individually

☐ For the entire week

What is the minimum amount of money that you would like to earn for each day?

Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
\$140-\$150 ▾	\$140-\$150 ▾	\$140-\$150 ▾	\$140-\$150 ▾	\$140-\$150 ▾	\$140-\$150 ▾	\$140-\$150 ▾

When are you available during the week?

To select an entire day, click on the column header. To select an hour on every day, click on the row header.

Time	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am - 01 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
01 am - 02 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
02 am - 03 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
03 am - 04 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
04 am - 05 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
05 am - 06 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
06 am - 07 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
07 am - 08 am	Available	Available	Available	Available	Available	Available	Available
08 am - 09 am	Available	Available	Available	Available	Available	Available	Available
09 am - 10 am	Available	Available	Available	Available	Available	Available	Available
10 am - 11 am	Available	Available	Available	Available	Available	Available	Available
11 am - 12 pm	Available	Available	Available	Available	Available	Available	Available
12 pm - 01 pm	Available	Available	Available	Available	Available	Available	Available
01 pm - 02 pm	Available	Available	Available	Available	Available	Available	Available
02 pm - 03 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
03 pm - 04 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
04 pm - 05 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
05 pm - 06 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
06 pm - 07 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
07 pm - 08 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
08 pm - 09 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
09 pm - 10 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
10 pm - 11 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
11 pm - 12 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable

Save and Next

Figure 3.1: Screenshot of Figma prototype for the constraint page. On behalf of participants, I clicked on radio buttons, input fields, and table cells to set their constraints. For technical reasons, input fields for the numerical constraints were implemented as dropdown menus.

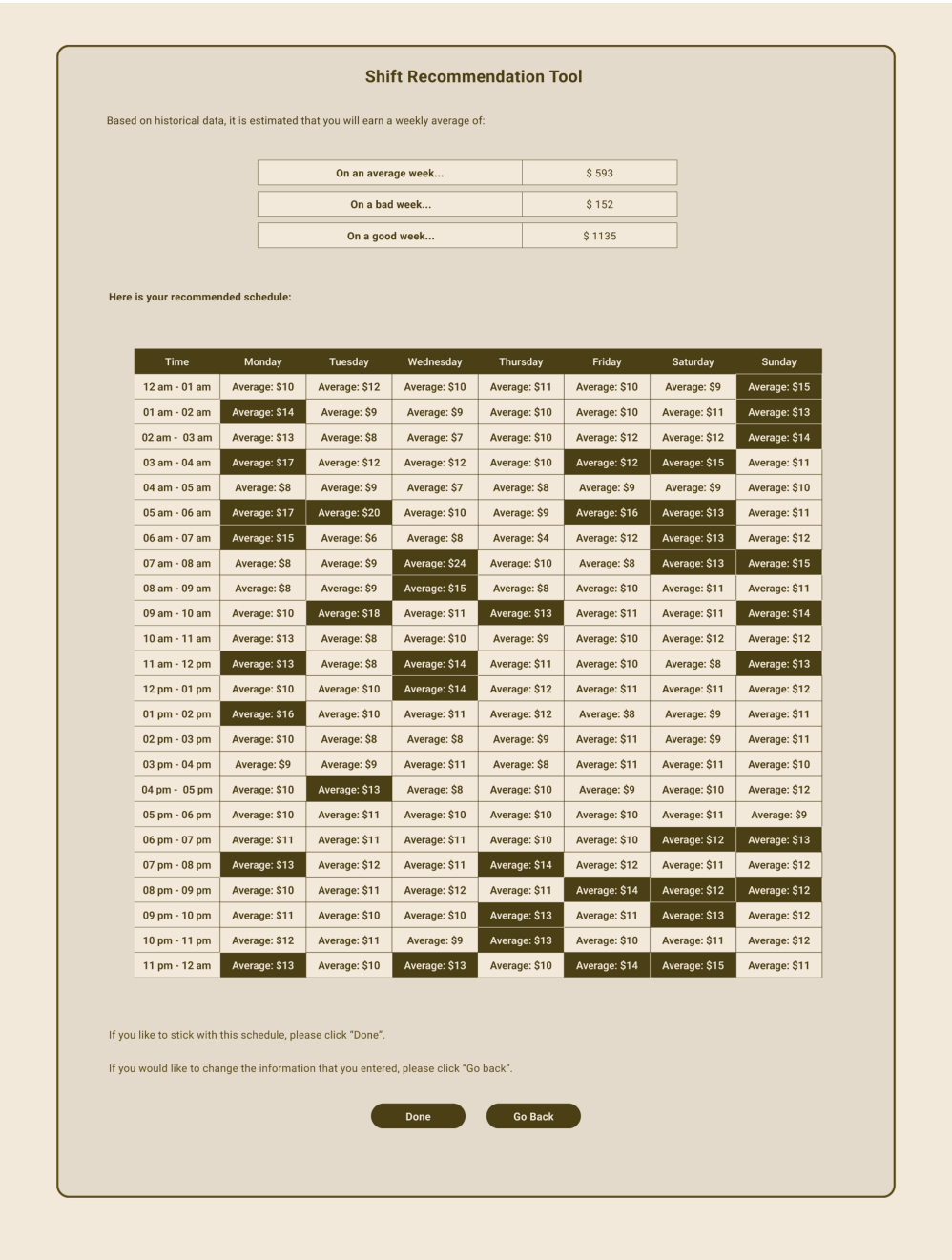


Figure 3.2: Screenshot of Figma prototype for the schedule page. A summary of mean, pessimistic, and optimistic weekly earnings is shown at the top of the page, followed by an hourly schedule where recommended cells are highlighted in darker colours. For technical reasons, this was shown as a static page not depending on previously-entered constraints.

- **More granular constraints.** P1, P3, and P4 all suggested ways to limit the scope of the historical gigs used to estimate their earnings. P3 wanted the tool to clarify that the historical data was limited to the city they drive in, and also how recent the data was. P1 and P3 wanted to limit the maximum distance of the historical gigs from their starting point. P4, who works for less popular delivery platforms, wanted to select which platforms the historical gigs came from, and indicated that this would improve their perceived control.
- **Feedback on performance.** P1 and P2 both wanted to compare the tool’s estimates with their actual earnings. Regardless of how the estimates compared to reality, P1 suggested that this feedback would be motivating; both P1 and P2 have gamified their gig-driving experiences to compare against either themselves or others. P2 also wanted to compare the estimates with their expenses, and P3 wanted a way to view the overall supply of drivers.

## 3.4 AI-Based Schedule Recommendation for Gig Driving

In this section, I describe my final design for a AI-based schedule recommendation tool for gig drivers. I use this tool as an exemplary AI decision aid to study the longitudinal relationship between the framing of uncertainty in outcomes and trust.

### 3.4.1 Decision Aid Design

My tool consists of two modules. First, an **estimation module** prospectively predicts  $e_{ij}$ , the earnings that drivers can expect during a specific hour  $j$  on a specific weekday  $i$ . These could be computed by a machine learning model or averaged from historical data. Second, for each driver, a **scheduling module** uses the estimated earnings and the driver’s constraints as inputs to produce an optimal set of working times. To do so, it solves a constrained optimisation problem to set variables  $x_{ij}$  to 1 or 0, denoting whether the driver is recommended to work in time slot  $(i, j)$ . Variants for these constraints were retained from the designs I tested in the pilot (Section 3.3.1).

- Some drivers wish to maximise their earnings while minimising their driving hours. For these drivers, the tool maximises the *objective function*: the sum of the estimated earnings  $e_{ij}$  for all recommended time slots  $(i, j)$  from the estimation module, i.e.  $\sum_{i,j} e_{ij}x_{ij}$ .

To set the *constraints*, the tool disallows time slots when the driver is not available, i.e.  $x_{ij} \leq a_{ij}$  where  $a_{ij}$  is an indicator of whether the driver is available during time slot  $(i, j)$ . It also places an upper bound on the total hours of recommended time slots per day by  $b_i$ , and per week by  $b_{tot}$ , i.e.  $\sum_j x_{ij} \leq b_i, \forall i; \sum_{i,j} x_{ij} \leq b_{tot}$ .

- Some drivers who value earnings to a greater extent set minimum targets for their hourly or daily earnings instead of restricting their driving hours. For these drivers, it minimises the *objective function*: the total hours of recommended time slots throughout the week, i.e.  $\sum_{i,j} x_{ij}$ .

To set the *constraints*, the tool disallows time slots when the driver is not available, i.e.  $x_{ij} \leq a_{ij}$ . It also places a lower bound on the estimated earnings per day by  $c_i$ , and per week by  $c_{tot}$ , i.e.  $\sum_j e_{ij}x_{ij} \geq c_i, \forall i; \sum_{i,j} e_{ij}x_{ij} \geq c_{tot}$ .

### 3.4.2 Interface Design

Next, the front-end interface of the decision aid, based on the prototype from Section 3.3.1, allows drivers to interact with the tool. The interface was implemented as an HTML/CSS/JavaScript website using Django 4.1 [103] and a PostgreSQL database. To mitigate potential biases, I designed my tool to be visually generic and distinct from apps or websites associated with any gig platforms. As with the prototype, the interface consists of two pages.

First, the **constraint page** (Figure 3.3) prompts users to select an optimisation objective: whether to maximise earnings or minimise hours on a daily or weekly basis ( $b_i$ ,  $b_{tot}$ ,  $c_i$ , and  $c_{tot}$ ). To maximise perceived control over the tool, I allowed users to choose these options freely rather than assigning them as conditions. The page also elicits hourly availability information ( $a_{ij}$ ). Compared to the prototype, the constraint page has a more varied colour scheme and clearer instructions, which were improved based on feedback from the pilot (Section 3.3.4).

Second, the **schedule page** (Section 3.4.3) shows the optimal schedule by highlighting the recommended time slots, i.e. the ones that lead to the highest earnings. Again, in the interests of maximising perceived control over the tool, I allowed users to revisit the constraint page until they were satisfied with the schedule. Like the constraint page, the prototype differed from the final design in its colour scheme and clarity of wording.

### 3.4.3 Interface Conditions

The schedule page uses the outputs of the estimation module to predict how much a driver following the recommended schedule would make per hour and per week. However, uncertainty inherently exists in these predictions, as they are based on historical data, and their realisation is contingent upon which gigs are offered to and accepted by drivers. To address Research Question 3.2, I varied the design of the schedule page between four conditions (Figure 3.4):

- (B) **Base condition.** Users were only shown their mean estimated earnings for the week and for each hour in the week.
- (D) **Daily estimates.** To assess the effect of introducing additional information irrelevant to uncertainty, users were shown their mean estimated earning for each day instead of their mean estimated weekly earning. As in (B), the schedule still showed mean estimated earnings for each hour.
- (R) **Ranged estimates.** To assess the effect of exposing uncertainty through range-based estimates (similar to Prabhudesai et al. [312]), users were shown mean, pessimistic, and optimistic estimates for hourly and weekly earnings. The prototype was most similar to Condition (R). Unlike the final design of Condition (R), however, ranges were not shown for hourly estimates; this was added based on feedback from the pilot (Section 3.3.4).
- (RH) **Ranged and hedged estimates.** To assess the effect of exposing uncertainty through lexical hedging (similar to Kim et al. [193] and Zhang et al. [487]), the textual description of the estimates was changed from (R). Instead of “Based on historical data, it is estimated that you will earn”, (RH) states “On average, based on historical data, a driver following this schedule will earn”.

### Tool Interaction

We will now give you access to a **tool that can recommend schedules of when you may want to drive**. To complete this study, you will need to interact with this tool for 7 days. Please take some time to enter information into the tool below about **your availability, preferences, and goals** when you drive for rideshare/delivery services.

The tool will use this information to recommend a schedule for you. **Please feel free to alter your answers and explore different options.** When you are done, you may scroll to the bottom and click to proceed. Overall, this should take about 5 minutes of your time.

### Shift Recommendation Tool

To enter your goals, answer the two following questions.

Would you like to specify a maximum number of hours that you would like to drive, or would you prefer to specify a minimum amount of money that you would like to earn?

☐ A maximum number of hours to drive ☒ A minimum amount of money to earn

Would you like to set these constraints for each day individually, or for the entire week?

☐ For each day individually ☒ For the entire week

What is the minimum amount of money that you would like to earn for the entire week?

\$

When are you available during the week?

To indicate that you are available during every hour of an entire day, click on the column header corresponding to that day.

To indicate that you are available during a particular hour on every day of the week, click on the row header corresponding to that hour.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
1 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
2 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
3 am	Available	Available	Available	Available	Available	Available	Available
4 am	Available	Available	Available	Available	Available	Available	Available
5 am	Available	Available	Available	Available	Available	Available	Available
6 am	Available	Available	Available	Available	Available	Available	Available
7 am	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
8 am	Available	Available	Available	Available	Available	Available	Available
9 am	Available	Available	Available	Available	Available	Available	Available
10 am	Available	Available	Available	Available	Available	Available	Available
11 am	Available	Available	Available	Available	Available	Available	Available
12 pm	Available	Available	Available	Available	Available	Available	Available
1 pm	Available	Available	Available	Available	Available	Available	Available
2 pm	Available	Available	Available	Available	Available	Available	Available
3 pm	Available	Available	Available	Available	Available	Available	Available
4 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
5 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
6 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
7 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
8 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
9 pm	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable	Unavailable
10 pm	Available	Available	Available	Available	Available	Available	Available
11 pm	Available	Available	Available	Available	Available	Available	Available

Next

[Discontinue Participation](#)

Figure 3.3: Screenshot of final constraint page, showing constraints entered by interview participant P1.

## Base condition (B)

Based on historical data, it is estimated that you will earn a weekly average of **\$1412**.

Here is your recommended schedule. Every cell shows an estimate of how much you will earn over that hour. The highlighted blue cells are hours that are recommended for driving. The other cells with darker text are hours for which you indicated you are available.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am	Average: \$13	Average: \$15	Average: \$14	Average: \$12	Average: \$12	Average: \$14	Average: \$15
1 am	Average: \$12	Average: \$12	Average: \$11	Average: \$11	Average: \$12	Average: \$15	Average: \$13
2 am	Average: \$12	Average: \$13	Average: \$12	Average: \$14	Average: \$13	Average: \$14	Average: \$14
3 am	Average: \$16	Average: \$16	Average: \$14	Average: \$15	Average: \$15	Average: \$12	Average: \$13

## Daily estimates (D)

Based on historical data, it is estimated that you will earn:

An average of... \$210 Monday \$194 Tuesday \$175 Wednesday \$173 Thursday \$184 Friday \$206 Saturday \$237 Sunday

Here is your recommended schedule. Every cell shows an estimate of how much you will earn over that hour. The highlighted blue cells are hours that are recommended for driving. The other cells with darker text are hours for which you indicated you are available.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am	Average: \$13	Average: \$15	Average: \$14	Average: \$12	Average: \$12	Average: \$14	Average: \$15
1 am	Average: \$12	Average: \$12	Average: \$11	Average: \$11	Average: \$12	Average: \$15	Average: \$13
2 am	Average: \$12	Average: \$13	Average: \$12	Average: \$14	Average: \$13	Average: \$14	Average: \$14
3 am	Average: \$16	Average: \$16	Average: \$14	Average: \$15	Average: \$15	Average: \$12	Average: \$13

## Ranged estimates (R)

Based on historical data, it is estimated that you will earn:

On an average week... \$1412  
On a bad week... \$391  
On a good week... \$2029

Mean  
10<sup>th</sup> percentile  
90<sup>th</sup> percentile

Here is your recommended schedule. Every cell shows an estimate of how much you will earn over that hour. The highlighted blue cells are hours that are recommended for driving. The other cells with darker text are hours for which you indicated you are available.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am	Average: \$13 Low: \$4 High: \$27	Average: \$15 Low: \$3 High: \$35	Average: \$14 Low: \$4 High: \$28	Average: \$12 Low: \$3 High: \$22	Average: \$12 Low: \$3 High: \$26	Average: \$14 Low: \$3 High: \$29	Average: \$15 Low: \$3 High: \$32
1 am	Average: \$12 Low: \$4 High: \$24	Average: \$12 Low: \$4 High: \$28	Average: \$11 Low: \$3 High: \$26	Average: \$11 Low: \$4 High: \$24	Average: \$12 Low: \$4 High: \$24	Average: \$15 Low: \$4 High: \$29	Average: \$13 Low: \$4 High: \$28
2 am	Average: \$12 Low: \$4 High: \$24	Average: \$13 Low: \$4 High: \$33	Average: \$12 Low: \$4 High: \$28	Average: \$14 Low: \$3 High: \$29	Average: \$13 Low: \$4 High: \$23	Average: \$14 Low: \$4 High: \$30	Average: \$14 Low: \$3 High: \$29
3 am	Average: \$16 Low: \$4 High: \$41	Average: \$16 Low: \$4 High: \$33	Average: \$14 Low: \$5 High: \$28	Average: \$15 Low: \$4 High: \$31	Average: \$15 Low: \$4 High: \$34	Average: \$12 Low: \$4 High: \$28	Average: \$13 Low: \$4 High: \$29

## Ranged + hedged estimates (RH)

In general, every driver may encounter different conditions while they are driving, which leads to differences in earnings from individual to individual and from one day to the next. On average, based on historical data, a driver following this schedule will earn:

On an average week... \$1412  
On a bad week... \$391  
On a good week... \$2029

Here is your recommended schedule. Every cell shows an estimate of how much the average driver will earn over that hour. The highlighted blue cells are hours that are recommended for driving. The other cells with darker text are hours for which you indicated you are available.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
12 am	Average: \$13 Low: \$4 High: \$27	Average: \$15 Low: \$3 High: \$35	Average: \$14 Low: \$4 High: \$28	Average: \$12 Low: \$3 High: \$22	Average: \$12 Low: \$3 High: \$26	Average: \$14 Low: \$3 High: \$29	Average: \$15 Low: \$3 High: \$32
1 am	Average: \$12 Low: \$4 High: \$24	Average: \$12 Low: \$4 High: \$28	Average: \$11 Low: \$3 High: \$26	Average: \$11 Low: \$4 High: \$24	Average: \$12 Low: \$4 High: \$24	Average: \$15 Low: \$4 High: \$29	Average: \$13 Low: \$4 High: \$28
2 am	Average: \$12 Low: \$4 High: \$24	Average: \$13 Low: \$4 High: \$33	Average: \$12 Low: \$4 High: \$28	Average: \$14 Low: \$3 High: \$29	Average: \$13 Low: \$4 High: \$23	Average: \$14 Low: \$4 High: \$30	Average: \$14 Low: \$3 High: \$29
3 am	Average: \$16 Low: \$4 High: \$41	Average: \$16 Low: \$4 High: \$33	Average: \$14 Low: \$5 High: \$28	Average: \$15 Low: \$4 High: \$31	Average: \$15 Low: \$4 High: \$34	Average: \$12 Low: \$4 High: \$28	Average: \$13 Low: \$4 High: \$29

Figure 3.4: Comparison of the four design conditions for the earnings estimates and recommended schedules on the schedule page, with abbreviations following Section 3.4.3. Boxes highlight differences between conditions in three areas: (1) weekly earnings estimates, (2) hourly earnings estimates, and (3) textual description of estimates.

## 3.5 Longitudinal User Study Design

To address Research Question 3.1, I conducted a longitudinal, in situ user study in which gig drivers repeatedly interacted with my schedule recommendation tool, and I measured their trust and reliance over these interactions. My participants used the tool for 7 days over a 14-day period, with the longer time window meant to accommodate variability in participants' availability. The methodology for this study was approved by our Institutional Review Board (IRB). Figure 3.5 illustrates the flow of the user study; I detail each day's study activities in Section 3.5.2.

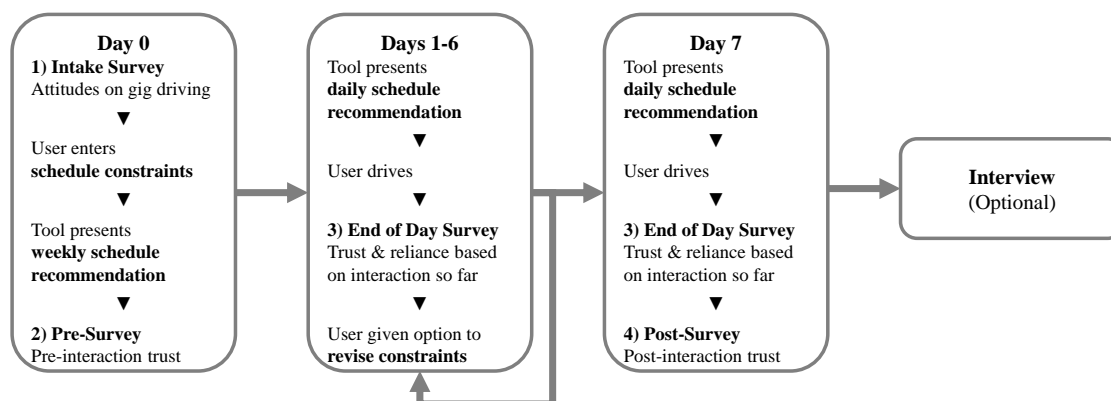


Figure 3.5: Flow of activities for the longitudinal user study. To be compensated, participants needed to complete Day 0 activities.

Based on a pilot conducted with 7 participants in August 2023, I determined that the Intake Survey, Pre-Survey, and tool interaction on Day 0 took an average of 14 minutes and 27 seconds, the End-of-Day Survey took an average of 2 minutes and 23 seconds per day, and the Post-Survey took an average of 2 minutes and 21 seconds. Based on van Berkel and Kostakos [415]'s recommendation of micro-compensation, this led me to set the compensation as an Amazon gift card with \$6 for the Day 0 surveys, \$2 for each daily survey, and a \$20 completion bonus (\$40 for full study completion). I made one payment upon study completion or the passage of 14 days.

### 3.5.1 Participants and Data Sources

As with the pilot interview study (Section 3.3.3), participants were recruited from the user base of Gridwise in September 2023. Gridwise distributed recruitment messages to users who (1) had completed at least one gig in DoorDash, Grubhub, Instacart, Lyft, Uber, or Uber Eats over the month preceding recruitment, and (2) resided in one of the four cities with the historical data used to generate the tool's earnings estimates: Los Angeles, New York, Chicago, and Houston. These were the platforms and cities for which historical data was available.

Accordingly, I generated estimates using gig data from August 2023 in each of these four cities. For each city, the data included approximately 100 000–300 000 gig records distributed evenly across times and weekdays. Hourly earnings were estimated by the mean of what drivers historically earned in this slot, filtered to the participant's city and platforms. I used this static

estimator to focus on the effects of exposing uncertainty in the estimates. Accordingly, I did not emphasise to participants that the schedule page was generated using AI or optimisation, and did not provide any additional information about the data used to generate the estimates.

### 3.5.2 User Study Activities

**Day 0: Pre-Interaction** Participants received a link to the study website from a recruitment message distributed by email. After the consent form, they completed the first of four surveys, the **Intake Survey** (Appendix A.1.1). This 12-question survey asked about their needs and motivations as gig drivers, along with demographics. These were the same questions as the formative portion of the pilot interviews (Section 3.3.2).

Next, participants were directed to interact with the tool, which I displayed in an iframe to mitigate response bias [98, 198]. They entered their constraints on the constraint page, and received the tool’s recommended schedule for the entire week on the schedule page. Participants were assigned to one of the four conditions for the schedule page (Section 3.4.3) uniformly at random, such that each condition had an approximately equal number of participants.

Lastly, participants completed the second of four surveys, the **Pre-Survey** (Appendix A.1.1), which was a 5-question survey measuring trust before interaction with the tool (Section 3.6.1).

**Days 1–7: Interaction** Next, participants began their 7 days of interaction with the tool, beginning on the next day of the week for which they indicated they were available to drive.

On each day, participants first received their recommended schedule for that day, sent via an email scheduled for 30 minutes before the start of their indicated availability. Thus, the tool’s outputs were displayed right as they were deciding their driving schedules. During the day, participants independently made decisions about their driving activity; I emphasised that compliance with the recommended schedule was not a condition of full participation.

At the end of each participant’s indicated availability for the day, a second scheduled email sent them a link to the **End-of-Day Survey** (Appendix A.1.1). This was an 8-question survey that measured their trust in the tool for that day, and their intention to rely on the tool for the next day (Section 3.6.1). If the participant intended to continue relying on the tool, they were then presented with a daily variant of the schedule page. Here, they could review the recommended schedule for the following day, and revise their constraints for the day as desired. Updated schedules were generated by fixing the recommended time slots for previous days using equality constraints and then re-solving the optimisation problem. However, if the participant intended to pause their interaction for one day, an email was sent on the next day, which prompted them to either review the next day’s schedule or to pause for an additional day.

**Day 7: Post-Interaction** On the final day, we removed the last question measuring reliance from the End-of-Day Survey, and added the **Post-Survey** (Appendix A.1.1). This was a 10-question survey that retrospectively measured participants’ trust and distrust in the tool over the entire user study (Section 3.6.1). After completing the Post-Survey, participants were sent a final email that invited them to participate in an optional **Exit Interview** (Section 3.5.3).



### 3.5.3 Interview Procedure

For participants who indicated their desire to be interviewed, an audio-recorded Zoom interview of 20–30 minutes was conducted by me. The interview focused on assessing dimensions of participants’ experiences that were not evident from the surveys. I began with questions about participants’ *motivations and routines*, which led into questions assessing the *constraint page*’s alignment with their decision-making process. Next, I asked participants about the *schedule page*, including how the recommended schedules factored into their decision-making and how it impacted the outcomes of their driving. Further questions focused on the *earnings estimates*, including perceptions of their accuracy and whether participants would’ve preferred another condition. Then, I asked participants to recall a specific day of interaction in terms of how the tool affected their behaviour for that day and for the following day. Finally, participants were asked for their *overall thoughts* on the tool. The full interview script is shown in Appendix A.1.2.

## 3.6 Quantitative Analysis

Among the 51 participants in the study, 25 (49%) were from Los Angeles, 10 (19.6%) were from New York, 8 (15.7%) were from Chicago, and 8 (15.7%) were from Houston; 4 (7.8%) were aged 18–24, 15 (29.4%) were aged 25–34, 22 (43.13%) were aged 35–44, 8 (15.7%) were aged 45–54, and 2 (3.9%) were aged over 55; 34 (66.7%) were male, 15 (29.4%) were female, and 1 (2%) was non-binary; 5 (9.8%) had a graduate degree, another 16 (31.4%) had an undergraduate degree, another 11 (21.6%) had a professional degree, and another 19 (37.3%) had a high school degree.

Out of these 51 participants, 44 completed at least one day of interaction with the tool, and 34 completed all 7 days of interaction. Starting from Day 0, Day 7 was reached by 6 (46%) of the base Condition (B) participants; 10 (71%) of the daily estimate Condition (D) participants; 7 (58%) of the ranged estimate Condition (R) participants; and 11 (92%) of the ranged and hedged Condition (RH) participants. I show full retention statistics in Figure 3.6.

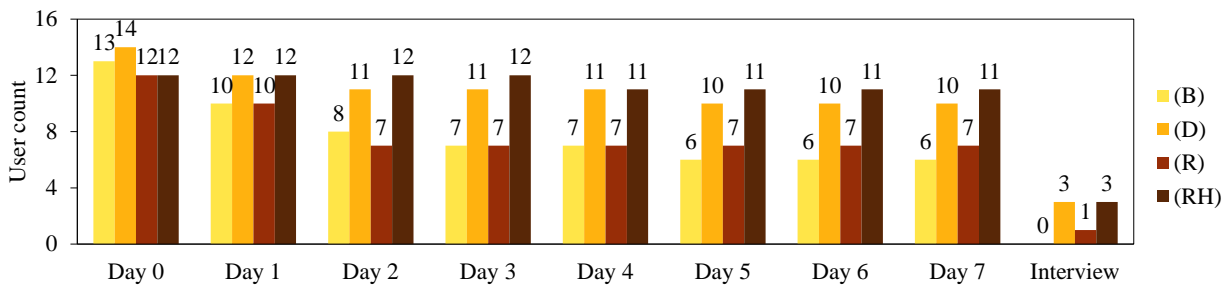


Figure 3.6: Retention statistics for the longitudinal user study, decomposed by design condition. Each day is labelled with the number of participants who completed all study activities for that day. Note the higher retention for Conditions (D) and (RH).

In the following sections, I first describe the metrics that I used to measure the participants’ trust and reliance (Section 3.6.1). Then, I analyse my findings from statistical models for these

metrics, specifically those relating to longitudinal effects (Section 3.6.2), and the effects of specific conditions (Section 3.6.3).

### 3.6.1 Metrics of Trust and Reliance

I measured the **trust** of participants using self-reported measures, following common practice [200]. I used two widely-used instruments for self-reported trust: the Human-Computer Trust Questionnaire (HCT) [257] and the Trust in Automation Scale (TiA) [179]. HCT measures 5 facets of trust using 5 questions each, while TiA measures both trust and distrust with 12 questions. On **Day 0** (pre-interaction), I included 5 items, one taken from each of the HCT's 5 facets of trust, in the Pre-Survey (Appendix A.1.1). On **Days 1–7** (during interaction), I included 3 items taken from 3 of the HCT's 5 facets of trust, in the End-of-Day Survey (Appendix A.1.1). On **Day 7** (post-interaction), I also included 5 items from the TiA in the Post-Survey (Appendix A.1.1), with 3 measuring trust and 2 measuring distrust. I chose not to include questions from this instrument earlier in the user study, as I felt that questions measuring distrust could have biased participants' reliance. All of these questions were presented to participants as 5-point Likert-type scales [79, 80]. From each survey, I computed an overall trust score by first inverting items measuring distrust, if any, and then averaging the Likert-scale responses.

I measured the **reliance** of participants, i.e. the external behavioural expression of trust, using both self-reported measures (End-of-Day Survey, Question 8; Appendix A.1.1) and their actual behaviour of discontinuing study participation. Specifically, I computed it as an ordinal variable with three levels: 1, if the participant indicated in the End-of-Day Survey that they intended to rely on the tool *more* tomorrow; 0, if the participant indicated that they intended to rely on the tool *about the same* tomorrow; and -1, if the participant indicated that they intended to rely on the tool *less* tomorrow, or did not complete the next day's study activities.

I use this notation to describe my statistical models:

- `pre_trust_score`: The **Day 0** (Pre-Survey) trust score.
- `trust_score`: The **current day**'s (End-of-Day) trust score.
- `reliance`: The **current day**'s (End-of-Day) reliance score.
- `post_trust_score`: The **Day 7** (Post-Survey) trust score.
- `day`: The day of interaction with the schedule recommendation tool (1–7).
- `user_id`: A randomly-assigned UUID for each participant, used as random effects.
- `condition`: The participant's schedule page condition.
- `estimate_accurate`: A binary indicator of whether the participant perceived their earnings to be about the same as *the tool's estimate* (End-of-Day Survey, Question 4).

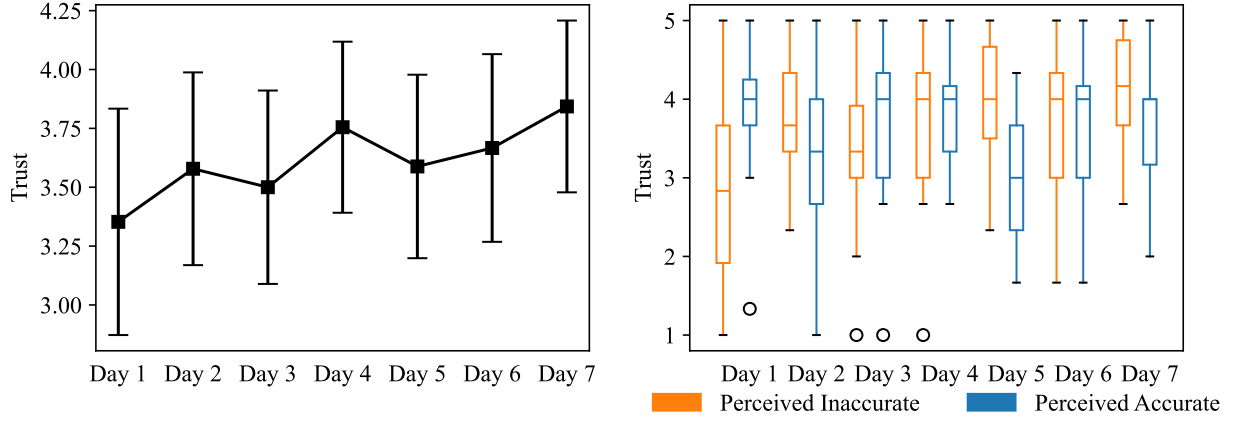


Figure 3.7: (Left) Means and 95% CIs of trust scores for the schedule recommendation tool on Days 1–7 among retained participants. Full statistics are shown in Table A.1 in Appendix A.1.3. (Right) Boxplots of trust scores on Days 1–7, decomposed by perceived accuracy.

### 3.6.2 RQ1: Longitudinal Effects

#### Effects on Trust

Participants reported a moderately high level of trust in the schedule recommendation tool ( $\mu = 3.631, \sigma^2 = 0.936$ ). To begin, I analysed how participants' trust in the schedule recommendation tool changed over time. For the 33 retained participants who completed all 7 days of interaction, Figure 3.7 (left) shows an upward trend in the mean trust score. To address Research Question 3.1, I then grouped each day's trust scores based on whether or not participants perceived the tool's estimates as being accurate. Figure 3.7 (right) shows that, on Day 1, perceived accuracy was positively correlated with trust; the interquartile ranges of the trust scores did not overlap between the two groups. This effect was less clear for Days 2–7, where trust scores for the two groups overlapped more extensively.

To further explore the longitudinal effects of perceived accuracy on trust, I fitted a linear mixed model (LMM) for `trust_score` using the R packages `lme4` 1.1-35.5 [29] and `lmerTest` 3.1-3 [212]. In this model, these longitudinal effects were modelled by the inclusion of the day, perceived accuracy (`estimate_accurate`), and their interaction as independent variables. I also included the `pre_trust_score` to adjust for participants' baseline level of trust in the tool (not on the same scale), and participant IDs as random effects to account for individual variance.

```
trust_score ~ pre_trust_score + day * estimate_accurate + (1 | user_id)
```

My model (Table 3.1) found that participants' pre-interaction trust was significantly and positively correlated with daily trust (`pre_trust_score`:  $\beta = 0.471, SE = 0.119, p = 0.00029$ ). Therefore, **participants' baseline trust persisted throughout their interactions with the tool**. Consistent with Figure 3.7, trust also increased significantly with each passing day (`day`:  $\beta = 0.130, SE = 0.027, p < 0.00001$ ). Also consistent with Figure 3.7, perceived accuracy was

Factor	Without Condition			With Condition		
	$\beta$	$SE$	$p$	$\beta$	$SE$	$p$
Intercept	1.447	0.468	0.00337**	1.994	0.549	0.00068***
pre_trust_score	0.471	0.119	0.00029***	0.411	0.128	0.00256**
day	0.130	0.027	0.00001***	0.052	0.046	0.25896
estimate_accurate	0.415	0.167	0.01357*	0.392	0.167	0.02010*
day:estimate_accurate	-0.121	0.037	0.00126**	-0.119	0.037	0.00168**
condition (D)				-0.494	0.341	0.15139
condition (R)				-0.169	0.360	0.64038
condition (RH)				-0.487	0.341	0.15755
day:condition (D)				0.111	0.053	0.03548*
day:condition (R)				0.103	0.056	0.06976†
day:condition (RH)				0.075	0.052	0.14971
Random intercept SD	0.605			0.613		

Table 3.1: Factors and coefficients ( $\beta$  with standard error  $SE$ ) for my linear mixed model of daily trust scores, without and with the `condition` as an independent variable. Statistically significant coefficients are denoted as † (0.1), \* (0.05), \*\* (0.01), \*\*\* (0.001).

significantly and positively correlated with trust (`estimate_accurate`:  $\beta = 0.415$ ,  $SE = 0.167$ ,  $p = 0.01357$ ), but it had less of an impact on trust with each passing day of the user study (`day:estimate_accurate`:  $\beta = -0.121$ ,  $SE = 0.037$ ,  $p = 0.00126$ ). This suggests that, **by the end of the user study, participants' trust was based less explicitly on perceived accuracy.**

### Effects on Reliance

Most participants indicated their desire to maintain their level of reliance on the schedule recommendation tool, corresponding to a reliance score of 0 ( $\mu = 0.038$ ,  $\sigma^2 = 0.643$ ). Trust and reliance were not strongly correlated ( $R^2 = 0.099$ ); some participants consistently expressed high reliance but also lower trust. The mean reliance score appeared to decrease over time, with the mean being lowest on Day 4, but I could discern no clear dependence on perceived accuracy (Figure 3.8). To clarify the nature of these longitudinal effects, I fitted another LMM using `lme4` and `lmerTest`. This model was similar to the model for trust, except the `reliance` score was the dependent variable, and I included the `trust_score` as an independent variable:

$$\text{reliance} \sim \text{pre\_trust\_score} + \text{day} * (\text{estimate\_accurate} + \text{trust\_score}) + (1 \mid \text{user\_id})$$

My model (Table 3.2) did not find significant effects for either the tool's perceived accuracy (`estimate_accurate`) or the `pre_trust_score`. However, two effects were significant: a negative effect from the `day`, supporting my initial observation ( $\beta = -0.220$ ,  $SE = 0.089$ ,  $p = 0.01444$ ), and a positive effect from the `day:trust_score` interaction ( $\beta = 0.058$ ,  $SE = 0.024$ ,  $p = 0.01640$ ). The latter suggests that reliance depended on perceived accuracy indirectly

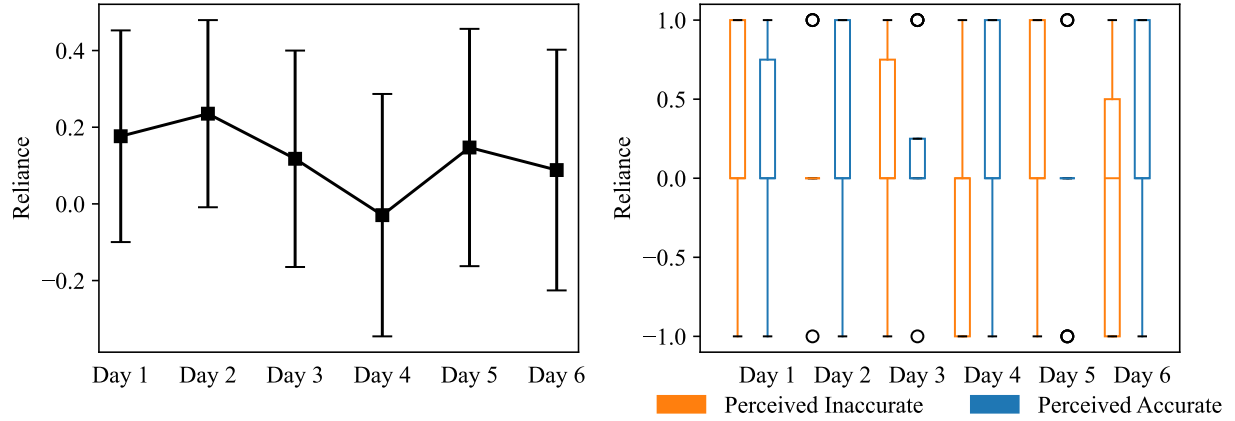


Figure 3.8: (Left) Means and 95% CIs of reliance scores for the schedule recommendation tool on Days 1–6 among retained participants. Full statistics are shown in Table A.1. (Right) Boxplots of reliance scores on Days 1–6, decomposed by perceived accuracy.

through trust. **Participants who trusted the tool more were more likely to continue relying on it;** this effect strengthened over interactions even as overall reliance weakened.

### 3.6.3 RQ2: Effects of Conditions

#### Pre-Interaction Trust

Next, I analysed the effects of the tool’s design conditions on trust and reliance, beginning with pre-interaction trust. Conditions (B)/(D)/(R)/(RH) had mean pre-interaction trust scores of 3.338, 3.629, 4.183, and 3.367; Condition (B) had the lowest, and Condition (R) had the highest.

To verify these initial observations, I used the Python package `statsmodels` 0.14.2 [356] to fit an ordinary least squares (OLS) model for pre-interaction trust, with the `condition` as an independent variable. Relative to Condition (B), the daily estimate Condition (D) did not significantly differ in pre-interaction trust (contrast (D) – (B):  $\beta = 0.290$ ,  $SE = 0.369$ ,  $p = 0.43188$ ); neither did the ranged and hedged estimate Condition (RH) (contrast (RH) – (B):  $\beta = 0.028$ ,  $SE = 0.384$ ,  $p = 0.94139$ ). Yet, Condition (R) had significantly higher pre-interaction trust relative to Condition (B) (contrast (R) – (B):  $\beta = 0.845$ ,  $SE = 0.384$ ,  $p = 0.02764$ ) and Condition (RH) (contrast (RH) – (R):  $\beta = -0.817$ ,  $SE = 0.391$ ,  $p = 0.03685$ ). Therefore, **exposing uncertainty through range-based estimates initially improved participants’ trust.**

#### Longitudinal Trust and Reliance

Next, I assessed the longitudinal effects of the schedule page design condition on trust and reliance. In Figure 3.9, I show the mean trust and reliance scores of participants in each of the four conditions. The means of the conditions were in most cases similar to the overall mean, with two exceptions: (1) the mean trust and reliance scores for Condition (R) were the highest of all con-

Factor	Without Condition			With Condition		
	$\beta$	$SE$	$p$	$\beta$	$SE$	$p$
Intercept	-0.024	0.422	0.95494	-0.129	0.496	0.79607
pre_trust_score	0.104	0.086	0.23316	0.103	0.092	0.27397
day	-0.220	0.089	0.01444*	-0.234	0.102	0.02280*
estimate_accurate	0.124	0.186	0.50517	0.125	0.187	0.50469
trust_score	-0.115	0.099	0.24285	-0.088	0.101	0.38470
day:estimate_accurate	0.012	0.048	0.81042	0.008	0.049	0.86786
day:trust_score	0.058	0.024	0.01640*	0.045	0.025	0.07127†
condition(D)				0.028	0.294	0.92461
condition(R)				-0.145	0.308	0.63782
condition(RH)				0.139	0.291	0.63414
day:condition(D)				0.062	0.067	0.35761
day:condition(R)				0.137	0.072	0.05958†
day:condition(RH)				0.043	0.065	0.50832
Random intercept SD	0.366			0.382		

Table 3.2: Factors and coefficients ( $\beta$  with standard error  $SE$ ) for my linear mixed model of daily reliance scores, without and with the `condition` as an independent variable. Statistically significant coefficients are denoted as † (0.1), \* (0.05), \*\* (0.01), \*\*\* (0.001).

ditions and showed a generally increasing trend; and (2) the trust scores for Condition (D) were lower on Days 5 and 6. To validate these trends, I added the `condition` and its interaction with the `day` as independent variables to the LMMs that I fitted in Section 3.6.2:

```
trust_score ~ pre_trust_score + day * (estimate_accurate + condition)
              + (1 | user_id)
reliance ~ pre_trust_score + day * (estimate_accurate + trust_score
              + condition) + (1 | user_id)
```

Condition (D) decomposed estimated earnings on a daily basis, thus providing information irrelevant to uncertainty. My models (Tables 3.1 and 3.2) indicate that this did not significantly improve either trust (`condition=(D):  $\beta = -0.494, SE = 0.341, p = 0.15139$` ) or reliance (`condition=(D):  $\beta = 0.028, SE = 0.294, p = 0.92461$` ) over the base Condition (B). While my trust model found a significant, positive longitudinal effect in Condition (D) ( $\beta = 0.111, SE = 0.053, p = 0.03548$ ), Figure 3.9 suggests that this does not represent a practically significant trend.

Condition (R) displayed uncertainty in predicted earnings using ranges of pessimistic and optimistic earnings. Again, my models did not find significant marginal effects for Condition (R) over Condition (B) in trust (`condition=(R):  $\beta = -0.169, SE = 0.360, p = 0.64038$` ) or reliance (`condition=(R):  $\beta = -0.145, SE = 0.308, p = 0.63782$` ). However, Condition (R) had nearly significant longitudinal effects for trust (`day:condition=(R):  $\beta = 0.103, SE = 0.057, p = 0.06976$` ) and also reliance (`day:condition=(R):  $\beta = 0.137, SE = 0.072, p =$`

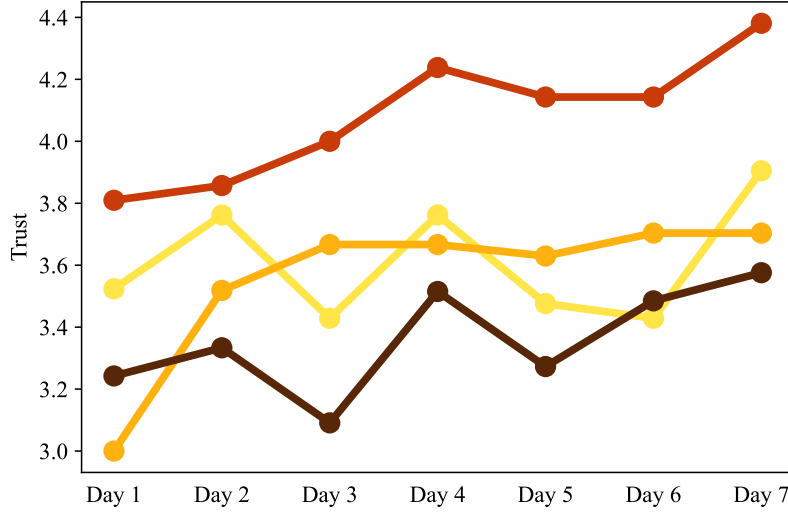


Figure 3.9: Mean trust (Left) and reliance (Right) scores for the schedule recommendation tool among retained participants, decomposed by condition. Note the higher mean scores for Condition (R). Full statistics are shown in Table A.1 in the appendix.

0.05958). This aligns with my observations based on the mean scores in Figure 3.9 as well as my findings for pre-interaction trust (Section 3.6.3). Therefore, despite exposing uncertainty in the tool’s earnings estimates, **the ranges of Condition (R) improved participants’ initial trust and then led them to maintain their trust and reliance** over daily interactions.

Condition (RH) added lexical hedging to the range-based earnings estimates in Condition (R). This condition was not significantly different from Condition (B) in marginal or longitudinal effects on trust and reliance. On Day 6, participants in Condition (RH) reported significantly lower reliance than participants in Condition (R) (Figure 3.8;  $\mu = 0.714, -0.182$ ; 95% CIs = (0.240, 1.188),  $(-0.649, 0.286)$ ), the only such significant pairwise difference on a daily basis (see Table A.1 in Appendix A.1.3). Combined with the pre-interaction trust of Condition (RH) being significantly lower than Condition (R) (Section 3.6.3), I conclude that **the addition of lexical hedging in Condition (RH) reversed the gains in trust and reliance from Condition (R)’s range-based uncertainty**.

### Post-Interaction Trust

Lastly, I fitted an OLS model for the `post_trust_score` with `statsmodels`. This model included the `condition` together with all previous trust (`pre_trust_score` and the daily `trust_score`) and reliance measurements.

$$\text{post\_trust\_score} \sim \text{condition} + \text{pre\_trust\_score} + \sum_{i=1}^7 (\text{trust\_score}_i + \text{reliance}_i)$$

For Conditions (B)/(D)/(R)/(RH), the mean post-interaction trust scores were 4.000, 3.720, 3.857, and 3.618. None of these conditions were significantly different from each other, and the coefficients for previous trust and reliance scores were not statistically significant either. The questions I adapted from the TiA asked participants to consider the entire duration of their interaction with the schedule recommendation tool. This broad, retrospective reflection may have failed to capture more nuanced longitudinal changes in trust and reliance like those I described in Section 3.6.

## 3.7 Qualitative Analysis

Overall, 7 participants completed the exit interview after they completed all seven days of the longitudinal user study. Three of these were from Condition (D), one was from Condition (R), and three were from Condition (RH):

- **P1:** A 47-year-old male with a high school degree driving for Instacart and Lyft in Los Angeles, Condition (D)
- **P2:** A 42-year-old female with a high school degree driving for Uber Eats in Los Angeles, Condition (D)
- **P3:** A 29-year-old female with a graduate degree driving for Lyft in Chicago, Condition (RH)
- **P4:** A 49-year-old male with an undergraduate degree driving for Lyft and Uber in Los Angeles, Condition (RH)
- **P5:** A 46-year-old male with an undergraduate degree driving for DoorDash, GrubHub, and Uber Eats in Chicago, Condition (R)
- **P6:** A 39-year-old male with a graduate degree driving for Lyft, Uber, and Uber Eats in Los Angeles, Condition (D)
- **P7:** A 39-year-old male with a professional degree delivering for DoorDash in Los Angeles by bike, Condition (RH)

To analyse these interviews, I used the same methodology as the pilot (Section 3.3.2). Now, I discuss my findings in relation to participants' motivations for using the tool (Section 3.7.1), perceptions of its accuracy (Section 3.7.2), and perceptions of its uncertainty based on the design conditions (Section 3.7.3).

### 3.7.1 Motivations and Routines

**Participants reported a diversity of motivations and routines for gig driving, which impacted their perceptions of the schedule recommendation tool's usefulness.** P1 and P2 viewed gig driving as a primary source of income, and thus found more value in the tool's earnings estimates:

[The tool was] definitely worthwhile, just because it gave me a number, a projection.

[...] They definitely motivate me to keep going the next day. (P1 (D))

Meanwhile, P3–P7 used gig driving to supplement other sources of income, but P4 still viewed his earning goals as important. Unlike other participants, P7 delivered with a bicycle in his spare time. He felt that his current commitment was insufficient to want to use the tool more:



If I take this job to a full time, take it seriously? I would [want to use it more].  
(P7 (RH))

Nevertheless, drivers found value in the tool regardless of their level of motivation. P1–P6 all reported challenges in estimating their potential earnings as a consequence of unpredictability in gig demand, pay, or location, or of gig platforms providing insufficient information. For instance:

Uber’s details that they offer to drivers through their interfaces are sorely lacking. So I’m grateful for the opportunity to interact with this tool.  
(P6 (D))

### 3.7.2 RQ1: Perceptions of Accuracy

When evaluating the tool’s accuracy, participants weighed its recommendations against their own routines and intuitions. For P1, P3, P6, and P7, the tool was a reference for how well they could perform in their existing routines, rather than something to reshape their routines:

I still would’ve followed my routine. [...] I was fortunate enough to at least have the tool make me a schedule based on the routine that I currently do.  
(P3 (RH))

However, P5 suggested that the tool could use a question-answering approach to nudge users into altering their routines, by first understanding their activity patterns and then suggesting modifications. When the tool was inaccurate, participants reacted in different ways. P1, P2, and P4 observed that instances of the tool being inaccurate decreased their desire to comply with the tool’s recommendations:

If I was making more than what it said,  
I would have done it more consistently on the schedule.  
(P2 (D))

P4’s reactions to inaccuracies were influenced by his expectations. He was motivated on one instance by the tool’s estimates exceeding his goals:

My target’s [...] \$30 an hour. Because those [estimated earnings] were consistently below \$30, [...] I wasn’t motivated to study it. But when I saw the 4 to 6 am,  
that kind of piqued my interest.  
(P4 (RH))

**Maintaining consistent perceptions of accuracy over time was important for building trust in this context.** P1, P2, and P5 indicated that the outcomes of their first one or two days of interaction impacted their willingness to follow the recommendations for the rest of the study. P4 and P5 indicated that their use of the tool would be strengthened longitudinally if they consistently perceived its predictions as being accurate:

Once I learned that it was accurate, and I had trust in it, and it was really helping, then I’d probably use it more and more.  
(P4 (RH))

### 3.7.3 RQ2: Perceptions of Uncertainty

Some participants recognised that the accuracy of the tool’s earnings estimates would be impacted by both their own decisions (P1, P3) and other environmental factors (P6):

It also depends, too, on the rides that I accept.  
(P3 (RH))

It might be true that I might earn the forecasted average earnings.  
But surges can definitely make a difference.  
(P6 (D))

Note that the tool's uncertainty was not exposed to P6, suggesting that this observation originated from their innate mental model. Recognising the effect of their own agency led P1 and P3, as well as P5, to adopt the tool's estimates as goals for their own earnings:

[...] setting daily goals of how much money I would like to make [...] was definitely something that I wasn't really doing prior to doing this study or using this tool.

(P3 (RH))

Also, P1 suggested that being able to compare their earnings to the tool's estimates in an hourly breakdown would be helpful for goal-setting.

All participants in Conditions (R) and (RH) (P3–P5, P7) appreciated the presence of ranges. P5 compared the tool's range to his own experiences:

I was always over the average.

So, to me, I was kind of in my head using that as a low.

(P5 (R))

For participants in the other conditions, P1 and P2 indicated that they would've preferred to have had ranges. However, P6 suggested that ranges may lead to disappointment when they are used for goal-setting:

If they earn less than [the higher number], then they probably might feel disappointed in the tool through no fault of its own, right? If you say \$12 to \$18, and it comes in at \$14, [...] I could understand how folks might look at that as a let down.

(P6 (D))

Thus, **ranged-based uncertainty was useful for decision-making, but needed to be calibrated against expectations.**

Both P4 and P5 struggled with the idea that the uncertainty in the tool's estimates could have originated from drivers with habits different to themselves:

Obviously no one's ever gonna work if it's just \$4 or \$5 an hour.

(P4 (RH))

This was in spite of the lexical hedging presented to P4. At least for this participant, the verbiage in the hedges may thus have failed to achieve its goal of leading him to consider potential sources of uncertainty more carefully.

## 3.8 Discussion

### 3.8.1 Key Findings and Implications

**Trust in AI decision aids is built both initially and over time.** I found that participants' pre-interaction trust in my schedule recommendation tool significantly impacted their trust during interaction (Sections 3.6.2 and 3.6.3). This is consistent with findings in the medical domain that practitioners [52, 56] and patients [287, 288] prefer to gauge their trust prior to interaction. My interviews similarly showed that perceptions of the tool's accuracy in the first two days influenced subsequent trust (Section 3.7.2). Yet, I also found that trust and reliance increased across interactions with the tool. While perceived accuracy had diminishing impacts on trust in later stages of the user study (Section 3.6.2), P4's experience (Section 3.7.2) shows that critical incidents where estimates differ significantly from expectations can cause catastrophic losses of trust [380]. P4 found this difficult to recover from. However, losses of trust could be mitigated prospectively by

calibrating perceptions of accuracy, e.g. by emphasising that drivers' outcomes are also a function of their own decisions. This could foster *appropriate reliance* by helping users decide when or when not to rely on the tool [351].

**Interactivity could help to maintain trust over time.** Losses of trust can also be mitigated retrospectively based on trust repair strategies [306]. Based on my qualitative findings, I hypothesise two mechanisms by which interactivity could help to maintain and repair trust. First, interactivity may enhance perceptions of control. The modes of interaction suggested by participants, such as hourly breakdowns and question-answering, would assure users that the AI has the intent and agency to capture and learn from their preferences [306]. Second, interactivity could help users to better recall their experiences and decisions. Until I probed further, most interview participants could not recall whether their earnings significantly differed from the tool's estimates.

**The impact of exposing uncertainty on trust in AI decision aids depends on task alignment.** Prior work has reached mixed conclusions on how exposing uncertainty in AI impacts trust and reliance. On similar tasks, Zhang et al. [491] found that confidence scores improve reliance, whereas Prabhudesai et al. [312] found that distribution plots dampen trust and reliance. Yang et al. [457] found that the effects of these designs depended on individual characteristics, but my results suggest another dimension: task alignment in the designs themselves. Task-aligned uncertainty representations, i.e. scalar ranges as opposed to distributions, allowed my participants to incorporate uncertainty directly into their decision-making (Section 3.7.3), thus improving trust (Section 3.6.3). This is consistent with findings in the AI explainability literature that domain-aligned explanations are more persuasive [60, 295]. I hypothesise that task alignment also underlies the negative effect of hedging I observed (Section 3.6.3): thinking about other drivers is not helpful when drivers are trying to reason about their own outcomes (Section 3.7.3).

**How uncertainty is exposed should be adapted to user subpopulations.** My results did not find that a one-size-fits-all approach exists to fostering trust. Even within the same condition, participants exhibited variability in how they reacted to the outcomes of their reliance. As is the case in recommendation systems, the same design that sustained trust in one user degraded it in another. Nevertheless, I hypothesise that it may be helpful to adapt uncertainty displays to subgroups within the gig driver population, such as those found by Di et al. [101]. Specifically, my qualitative results point to differing perceptions of accuracy and uncertainty between highly motivated drivers (e.g. P1 and P2) and less motivated drivers (e.g. P7). How can these subgroups be helpful? On one hand, drivers wanted to receive estimates from those with habits similar to themselves (e.g. P4 and P5). This aligns with the bandwagon effect from recommendation systems research [233, 392]. On the other hand, drivers would find more value in a personalised system. This aligns with the Barnum effect for recommendation systems [201, 390, 398]. Subgroups that are neither fully personalised nor fully generalised could help to strike a balance between these two cognitive biases. A future large-scale study could help to confirm my hypothesis.

### 3.8.2 Limitations

My work has two primary limitations. First, I cannot claim that the design of my tool was optimal for engendering trust. My focus was on testing how designs for exposing uncertainty would impact trust and reliance. Thus, I attempted to isolate the effect of this design dimension by refining the

tool through a formative pilot study (Section 3.3). Nevertheless, further improvements may have been possible. For instance, I could not provide retrospective breakdowns of participants' earnings due to data availability limitations. Thus, design choices orthogonal to the exposure of uncertainty may have impacted participants' trust and reliance.

Second, despite my best efforts, my sample of drivers was limited. These individuals were at least aware, if not active users, of the Gridwise app, and thus they may have been more focused on their outcomes than the general gig driver population. The trust of users in an AI decision aid is contingent upon their domain knowledge [194, 453], and — as I demonstrate (Section 3.7.1) — the extent to which they integrate the decision aid into their existing routines. A future study aimed at a broader population of gig drivers could uncover additional insights by explicitly controlling for factors such as full-time status and driver tenure. I was also unable to reach participants who discontinued the user study. Future studies that follow up with such participants would be a valuable source of data on mechanisms of trust loss and repair in longitudinal settings.

### **3.8.3 Recommendations for Future Work**

The paucity of similar longitudinal, in-situ studies in prior work is understandable given the logistical challenges I encountered. I stress the importance of observational studies to improve domain understanding as a basis for longitudinal interventional studies. My pilot interviews helped me to design a tool that was task-aligned, which led participants to find value in it over repeated interactions. Furthermore, my study design aimed to increase perceived control while reducing user burden through flexibility in the scheduling of participation; customisability of the constraints on the AI decision aid; shorter survey instruments; and incremental compensation.

## **3.9 Conclusion**

In this chapter, I considered how uncertainty impacts the deployment of AI technologies in gig driving. I assessed how gig drivers' trust and reliance on an AI decision aid is influenced by designs that expose uncertainty, in the context of an in-situ deployment unlike the laboratory experiments used by previous work. Specifically, I conducted a longitudinal, in situ user study of an AI-based schedule recommendation tool with  $n = 51$  gig drivers. Their interactions with my tool impacted their actual earnings. My findings demonstrate that trust can be built by (1) maintaining perceptions of accuracy over repeated interactions and (2) displaying uncertainty in a task-aligned fashion. More generally, they point to the importance of designing with stakeholders in mind as a way to improve the trustworthiness of AI technologies in transportation. Although my work in this chapter constitutes an initial foray into improving the deployability of AI technologies for gig drivers, broader challenges surrounding uncertainty, heterogeneity, interpretability, and coordination remain to be addressed (Section 2.4).

# Chapter 4

## Purpose in the Machine

### Evaluating the Impact of Heterogeneity in Driver Behaviour and Simulation Scale on Traffic Simulation Outcomes

*Domain:* Traffic signal control

*Challenges:* Heterogeneity

One must accommodate the times, or things get messy.

---

*Philip J. Coulson*

#### 4.1 Introduction

In this chapter, I begin a series of studies of design and deployment challenges for reinforcement learning (RL) in the domain of traffic signal control (TSC). As I outlined in Section 2.2.4, traffic simulators are an important tool to evaluate the quality of a signal plan.

When these signal plans are based on deep RL, traffic simulators improve the efficiency of training and evaluation in two ways. First, it must be possible to repeatedly evaluate signal plans, as they may need to be refined iteratively based on complex and potentially conflicting requirements from many stakeholders [139]. These include both end-users (i.e. road users) and decision-makers (i.e. traffic engineers and city planners). Second, deep RL algorithms require particularly large quantities of data for training. Real-world traffic studies cannot be continuously-occurring and large-scale due to efficiency and safety costs [128], but traffic simulators can generate large quantities of realistic data for both training and evaluation [482].

At the same time, traffic simulators also allow *heterogeneity* in the deployment contexts of signal plans to be addressed more thoroughly during training and evaluation. First, heterogeneity between human stakeholders can be incorporated: traffic simulators can model different types of road users, who have unique needs and make unique choices that must be addressed individually and equitably. Second, heterogeneity between environmental characteristics can be incorporated: traffic simulators can produce a diverse set of scenarios that RL algorithms can learn from.

This chapter explores how these two types of heterogeneity are modelled by two different traffic simulators, SUMO and CityFlow, which I introduced in Section 2.2.4. Released in 2001, SUMO [10] is the most popular traffic simulator in RL for TSC [302]. It supports an expansive

framework for simulation definitions as well as efficient programmatic API access. However, it is single-threaded and thus scales less well to very large road networks. Motivated by this limitation, in 2019 Zhang et al. [482] introduced CityFlow, a traffic simulator designed for applications of RL to traffic signal control. It is multithreaded and thus has been reported to achieve a speedup of  $>20\times$  over SUMO. However, its framework for simulation definitions is more restricted and focuses on aspects which are essential for RL training. Its adoption is limited but increasing [302].

If traffic simulators are to serve as training environments for RL algorithms, the simulators' modelling assumptions must be sufficiently realistic that the resulting signal plans can learn to adapt to a variety of scenarios during deployment. Thus, real-world validation is crucial. However, granular validation with real-world data is usually not possible due to the aforementioned challenges of data collection. Comparisons between simulators take a partial step towards this goal by verifying that different simulators lead to equivalent outcomes. In this chapter, I compare CityFlow against the more granular SUMO. Similar comparisons were conducted by Zhang et al. [482] and Mei et al. [268], but their evaluation used the expected travel time of vehicles — a long-term, system-level outcome — given trained RL policies. Such indirect evaluations do not capture the effects of simulator design and simulation heterogeneity on the *inputs* of RL policies: queue lengths and other instantaneous, low-level features (Section 2.3.1). I conducted a more comprehensive comparison to answer the following research questions:

**RESEARCH QUESTION 4.1.** *Do the low-level simulation outcomes of CityFlow and SUMO have a statistically significant level of distributional equivalence?*

**RESEARCH QUESTION 4.2.** *How is this distributional equivalence affected by heterogeneity between vehicles in terms of different driver behavioural models?*

**RESEARCH QUESTION 4.3.** *How is this distributional equivalence affected by heterogeneity in the simulation's traffic characteristics, in terms of traffic demand and road network size?*

This chapter was published at the Winter Simulation Conference (WSC) in 2023 [68].

## 4.2 Related Work

### 4.2.1 Validating Traffic Simulators

Validation is an important yet challenging aspect of the development of traffic simulators that ensures their fidelity to the real world. A multitude of road networks have been used to validate SUMO itself [10, 32] and to calibrate its car-following models [205] in comparison to detector data. For instance, Lobo et al. [242] compared traces for vehicle counts and crossing times between SUMO and detector data for a road network in Ingolstadt, Germany; their simulation is used in this work. Meanwhile, CityFlow was validated by comparing average travel times under various traffic volumes to that of SUMO [268, 482].

A number of studies have compared outcomes from multiple simulators; my work in this chapter falls in this setting. Maciejewski [256] compared vehicle counts under different traffic demand and driver behaviour settings for SUMO, the commercial simulator VISSIM, and TRANSIMS.

Several studies involved SUMO and the commercial simulator AIMSUN. Leksono and Andriyana [221] compared travel times and queue lengths in these simulators under different traffic management interventions for a roundabout in Norrköping, Sweden. Ronaldo and Ismail [336] used *t*-tests to compare flows and speeds in these simulators to observations from a highway interchange in Stockholm, Sweden. Baza-Solares et al. [31] applied *t*-tests to vehicle counts from the two simulators for a road network from Bucaramanga, Colombia. This work differs from prior approaches in that: (1) the measures used relate to the distribution of outcomes across the network, not just at individual points; and (2) these measures were evaluated across multiple road network scales.

### 4.2.2 Modelling Driver Behaviour

There has been a significant body of literature on building realistic models of driver behaviour. Car-following behaviour was the first to be modelled, with the early Gazis-Herman-Rothery model being over 50 years old. Subsequent work has yielded optimal velocity, fuzzy logic, collision avoidance, action point, and cellular automaton models [227, 342]. I implemented collision avoidance and action point models. Lane-changing behaviour has been modelled by a newer, separate line of work [284, 497], which has produced models based on rules, discrete choice, game theory, and cellular automata. Some work has built unified models of car-following, lane-changing, and other driver behaviour [260, 406]. Here, car-following and lane-changing models were considered separately but co-varied in experiments to elucidate the effects of their interactions. The most relevant prior work is Capela Dias et al. [59], who analysed the impact of driver behaviour on system-level travel time in SUMO; by contrast, I focus on distributional equivalence between two simulators in terms of lower-level outcomes.

## 4.3 Comparing SUMO and CityFlow

Both CityFlow and SUMO are microscopic traffic simulators. They represent road networks as graphs, with intersections as nodes (“junctions” in SUMO) and roads as edges between nodes [22]. Vehicles are generated by flows; all vehicles within a flow share similar attributes, including behavioural parameters and routes [21]. Routes are defined as sequences of edges. In order to follow its route, a vehicle will proceed down a road edge until it arrives at an intersection. As shown in Figure 4.1, permissible movements through an intersection — through, right, and left — are defined by specific pairs of lanes (“roadlinks” in CityFlow, “connections” in SUMO).

When executing the simulation loop, both simulators perform the same procedures in each timestep, although the two simulators perform them in different orders (Figure 4.2):

- **Routing.** Vehicles decide whether to keep or alter their current routes.
- **Traffic signal logic.** The states of signals are updated according to signal plans.
- **Car-following.** Vehicles decide how to follow vehicles in front of them.
- **Intersection navigation.** Vehicles decide how to traverse roadlinks, if necessary.
- **Lane-changing.** Vehicles decide whether and how to change lanes.

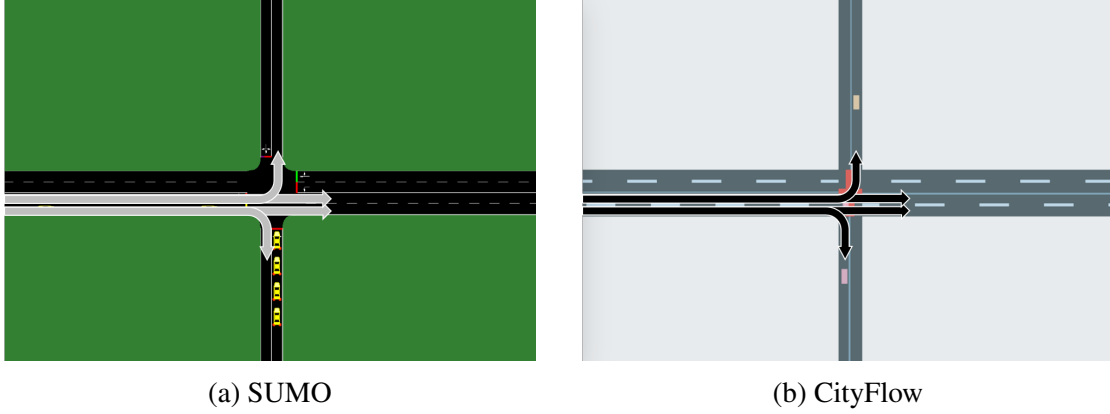


Figure 4.1: Screenshot of an intersection in the arterial4x4 road network in SUMO and CityFlow, showing connections in SUMO/roadlinks in CityFlow.

- **Environment update.** After all of these decisions have been made, all simulation objects are updated with their current state.

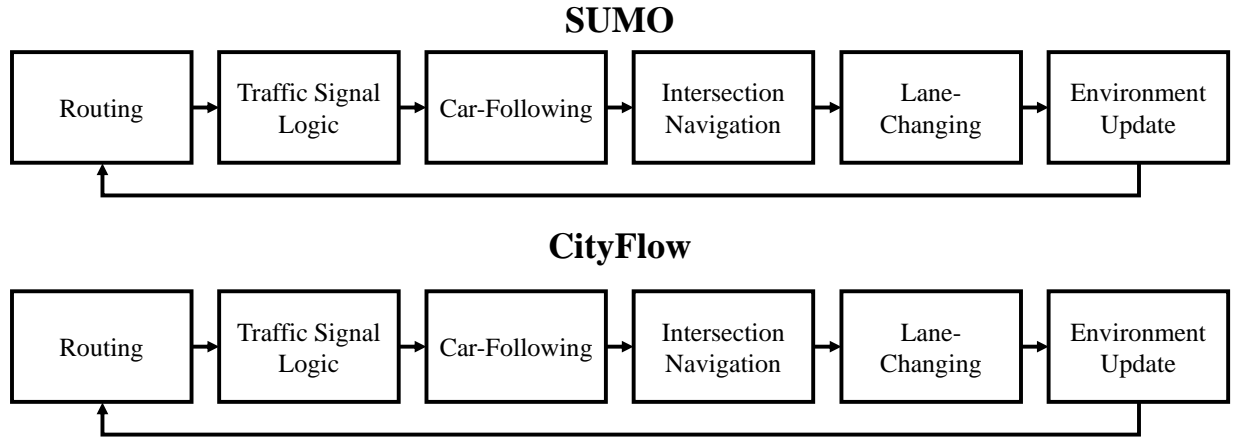


Figure 4.2: Diagram of simulation loop for SUMO and CityFlow.

However, as shown in Table 4.1, key differences between the two simulators exist in finer details of implementation [268]. Additional differences exist in how these simulators model heterogeneity in vehicle behaviour; I discuss these aspects in Section 4.4.

Where possible, I designed the experiments to control or co-vary differences between the simulators. In particular, I converted flows in both CityFlow and SUMO to be single-vehicle and deterministic, and to share the same routes. However, differences remain, particularly aspects of driver behaviour (Section 4.4) driven by randomness in SUMO; by contrast, CityFlow only uses random generation for vehicle priority. Such uncontrolled factors may account for their differences.



Feature	SUMO	CityFlow
<i>Performance</i>	Single-threaded if using <code>libsumo</code> , relatively slow due to socket communication if using <code>traci</code>	Multithreaded, relatively fast
<i>Vehicles</i>	Heterogeneous classes supported, including pedestrians and cyclists	Homogeneous vehicles other than size and behavioural parameters
<i>Roads</i>	Complex roads including sidewalks, bike lanes, ramps, roundabouts	Simpler roads with lanes being homogeneous other than width and speed
<i>Traffic signals</i>	Detailed signal plans including yellow change and pedestrian intervals	Binary green/red for roadlinks to determine status of traffic lights
<i>Randomness</i>	Incorporates random perception error into vehicle behaviour	Largely deterministic other than vehicle priority determination

Table 4.1: Comparison of implementations of SUMO and CityFlow.

## 4.4 Varying Driver Behaviour

Addressing Research Question 4.2 requires distributional equivalence to be assessed under different variations of driver behaviour; as established in Section 4.2.2, these are implemented through car-following and lane-changing models in traffic simulators. In both CityFlow and SUMO, car-following logic is used to maintain a safe gap to the leading vehicle, while lane changes are used to switch vehicles between lanes so that they can take appropriate roadlinks at intersections to continue their routes. However, in CityFlow, car-following and lane-changing are handled in separate threads; in SUMO, they are handled in sequence.

### 4.4.1 Car-Following Models

*Car-following models* determine the speeds at which vehicles travel unobstructed (free speed), follow a lead vehicle (following speed), and stop at an obstacle (stopping speed). Two types of numerical integration can be used to calculate vehicle speeds: a Euler update, which solves for the speed at discrete timesteps, and a ballistic update, which solves for the acceleration at discrete timesteps and applies it to the speed. I controlled both simulators to use ballistic updates.

Several shared parameters have varying effects on different car-following models. CityFlow assumes that vehicles have usual and maximum accelerations and decelerations; SUMO assumes that vehicles have a maximum possible *acceleration* and *deceleration* (which are used as the usual values), and a maximum *emergency deceleration*. I used the latter formulation here, although this may have resulted in more aggressive behaviour than if lower usual accelerations and decelerations were used. Additionally, both simulators model vehicles as having minimum desired following distances in terms of space (i.e. the *minimum gap*) and time (i.e. the *minimum headway*). Capela Dias et al. [59] co-varied these parameters to model the effect of driver aggressiveness on travel time. I adopted their taxonomy of aggressiveness types, but excluded gender effects due to disagreement in the literature on their significance [399, 402]. This gave six parameter settings (Table 4.2).

Table 4.2: Parameter settings for six aggressiveness types based on Capela Dias et al. [59]. I set maximum emergency deceleration to  $-9.0 \text{ m/s}^2$ , as they did not specify this parameter.

Type	Max. accel.	Max. decel.	Max. emerg. decel.	Min. gap	Min. headway
Aggressive young	$3.1 \text{ m/s}^2$	$-5.5 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	1.2 m	1.0 s
Courteous young	$2.5 \text{ m/s}^2$	$-4.5 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	2.5 m	1.0 s
Aggressive middle-aged	$2.9 \text{ m/s}^2$	$-5.0 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	2.0 m	1.3 s
Courteous middle-aged	$2.4 \text{ m/s}^2$	$-4.1 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	2.5 m	1.5 s
Aggressive old	$2.6 \text{ m/s}^2$	$-4.5 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	2.0 m	1.7 s
Courteous old	$2.3 \text{ m/s}^2$	$-3.8 \text{ m/s}^2$	$-9.0 \text{ m/s}^2$	2.5 m	1.9 s

The default car-following models in CityFlow and SUMO are both modified from the collision avoidance model of Krauß [207]. For free speed, CityFlow’s implementation uses the maximum speed, while SUMO’s implementation modulates this by the visible lookahead distance. For stopping speed, both simulators solve somewhat different quadratic equations to determine the deceleration needed to stop within a fixed distance. For following speed, a target distance is maintained, which is computed by the desired minimum headway as well as the speed and maximum deceleration of the lead vehicle. I added SUMO’s variant of the Krauß model to CityFlow. SUMO implements several other car-following models, of which I also re-implemented several in CityFlow to introduce variation in driver behaviour:

- The collision avoidance/action point model of Wagner [420], which probabilistically combines Krauß’s model with the action point model of Todosiev [405].
- The collision avoidance model of Wiedemann [445], a behavioural model that varies between free, approaching, following, and emergency modes based on the gap to the lead vehicle.
- The adaptive cruise control (ACC) model of Milanés and Shladover [276], which determines acceleration using a “speed control” method if the gap to the lead vehicle is large and using a “gap control” method if the gap is small.

#### 4.4.2 Lane-Changing Models

*Lane-changing models* are rule-based in CityFlow and SUMO. In both CityFlow and SUMO, a vehicle initiates a lane change by signalling the lead and lag vehicles on the destination lane. If the gaps to the lead and lag vehicles are sufficiently large (Figure 4.3), the vehicle makes the lane change. Upon completion of the lane change, the vehicle is instantly teleported. (SUMO also supports a sublane model [21, 359] that models lateral movement, which I excluded from consideration as a control in my experiments.) Despite these similarities, the two simulators’

lane-changing models differ irreconcilably in implementation details. Vehicles in CityFlow only perform lane changes to follow their routes, and lane-changing is based on the explicit insertion of a copy of the vehicle (the “shadow vehicle”) on the target lane that the lag vehicle will follow. Vehicles in SUMO follow a parametrised hierarchy of motivations, which includes strategic lane changes to follow routes as well as tactical lane changes for overtaking [116].

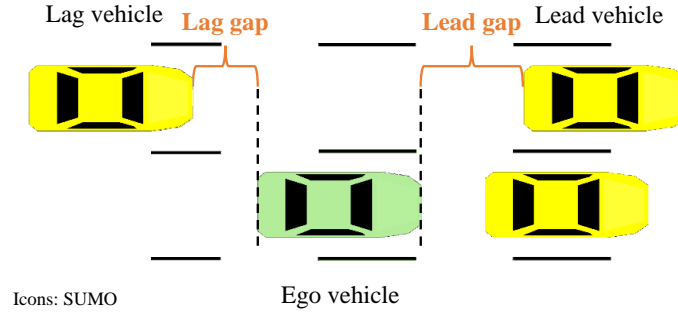


Figure 4.3: Diagram showing lead and lag gaps in lane-changing.

I introduced one of SUMO’s parameters, the gap tolerance factor, to CityFlow. When vehicles are changing lanes, the minimum gap required to initiate a lane change is given by the necessary gap for collision avoidance divided by this constant factor. A factor of 1.0 represents the default behaviour. My experiments varied this factor between 0.5, 0.82, 1.0, 1.18, and 1.5, representing decreases/increases of 18% and 50% in tolerance. This is based on Sun and Kondyli [389], who observed the mean gaps for forced, cooperative, and free lane change manoeuvres to be 45 ft, 53 ft, and 109 ft in a video dataset.

## 4.5 Experiments

To address the research questions, I designed two experiments, the basic structure of which is summarised in Figure 4.4. Both address Research Question 4.1 and Research Question 4.2 by assessing CityFlow and SUMO’s distributional equivalence under different settings of driver behaviour. The two experiments also address separate aspects of Research Question 4.3: Experiment 1 assesses distributional equivalence while varying the simulation scale by traffic demand, and Experiment 2 while varying it by road network size. I accomplished this using road networks from the benchmark dataset of Ault and Sharon [15]. All of these networks were initially defined using SUMO syntax; I used the SUMO-to-CityFlow network converter of Zhang et al. [482] to generate their CityFlow counterparts, and created a flow converter to map between flows in the two simulators.

I considered four independent variables in both experiments: the car-following model (Section 4.4.1, 5 levels), the car-following aggressiveness parameters (Section 4.4.1, 6 levels), the lane-changing gap tolerance (Section 4.4.2, 5 levels), and the road network (2 levels). For the car-following model, I compared the default and SUMO-based implementations of the Krauß model in CityFlow to the Krauß model in SUMO, while I compared the other models’ CityFlow implementations to their respective SUMO counterparts. I also held constant the fundamental lane-changing

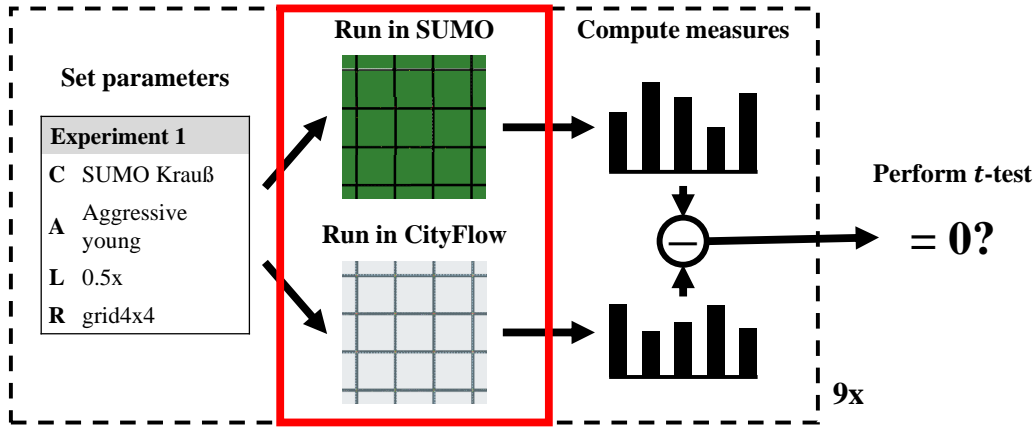


Figure 4.4: Diagram of experimental flow.

models for the two simulators, as was the traffic signal program — a simple fixed-time program retained from the original dataset.

I used eight instantiations of two types of measures to assess distributional equivalence. First, the root mean squared error (RMSE) quantifies the point-to-point difference in individual outcome measures. I computed this as the mean RMSE of the total travel time and waiting time (defined as the time that a vehicle spends queued with a speed  $< 0.1^{\text{m/s}}$ ) over all vehicles and timesteps; the mean RMSE of per-lane total vehicle counts and queued vehicle counts over all lanes and timesteps; and the mean RMSE of the speed and acceleration over all individual vehicles and timesteps. Second, the Kullback-Liebler (KL) divergence measures the difference in the distribution of outcomes over the entire road network. I computed this for the distributions of vehicle counts and queued vehicle counts as the mean over all timesteps.

### 4.5.1 Setup

#### Experiment 1: Traffic Demand

For this experiment, I used the road networks arterial4x4 [250] and grid4x4 [65] (Figure 4.5) from Ault and Sharon [15]’s RESCO benchmark. Both are synthetic grid networks with similar topologies, but they vary in the level of congestion. grid4x4 has six-lane roads with a uniformly distributed demand of 1,473 vehicles. arterial4x4 has major (four-lane) and minor (two-lane) roads, and a demand of 2,484 vehicles that alternates between major and minor roads. arterial4x4’s traffic pattern leads to congestion and degraded RL performance in simulations [15].

The following power analysis uses these variable names:  $C$  for car-following models,  $A$  for car-following aggressiveness,  $L$  for lane-changing gap tolerance, and  $R$  for road network. The second-order linear multiple regression included 5 ( $C$ ) + 6 ( $A$ ) + 1 ( $L$ ) + 1 ( $L^2$ ) + 2 ( $R$ ) + 10 ( $C \cdot R$ ) + 12 ( $A \cdot R$ ) + 2 ( $L \cdot R$ ) + 30 ( $C \cdot A$ ) + 5 ( $C \cdot L$ ) + 6 ( $A \cdot L$ ) = 80 variables. Using G\*Power 3.1.9.7’s power calculation for ordinary linear multiple regression with a fixed model

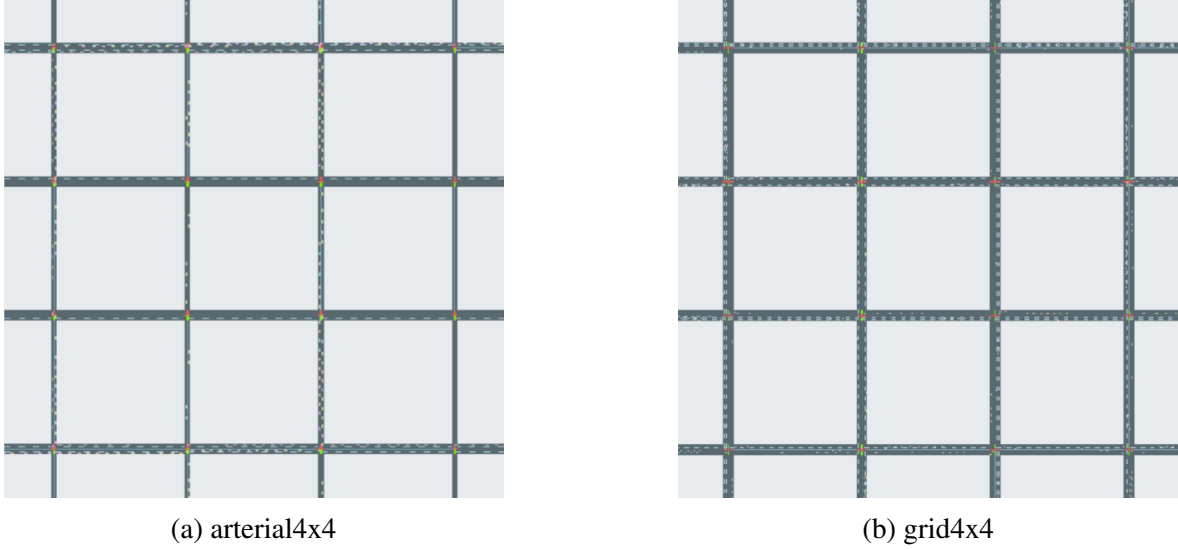


Figure 4.5: Screenshots in CityFlow of the arterial4x4 and grid4x4 road networks.

and  $R^2$  increase, a small effect size of 0.02, and  $\alpha = \beta = 0.95$ , I computed the total necessary sample size as 2,646. Divided by the number of cells,  $5 \cdot 6 \cdot 5 \cdot 2 = 300$ , I computed the number of replications per cell as  $\lceil \frac{2,646}{300} \rceil = 9$ . I executed all replications with Python 3.9.16, SUMO 1.12.0, and a modified version of CityFlow 0.1 on a shared server with four cores, two 4.2GHz Intel i7-7700K processors per core, and 62 GiB of RAM.

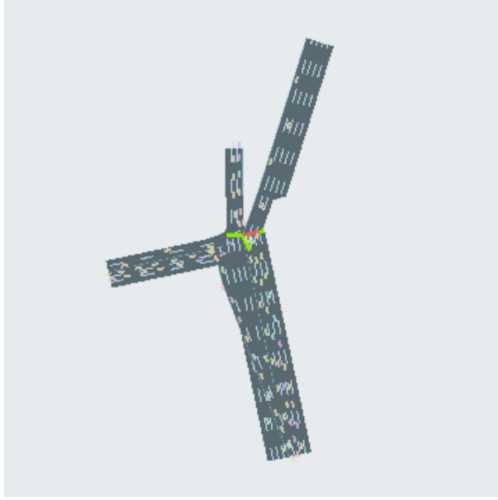
## Experiment 2: Network Scale

For this experiment, I used the road networks ingolstadt1 and ingolstadt7 (Figure 4.6) from Ault and Sharon [15]’s RESCO benchmark. Both are subsets of the Ingolstadt road network that was simulated by Lobo et al. [242]. They respectively contain 1 and 7 signalised intersections, representing a single busy intersection and a larger arterial road; ingolstadt7 is a superset of ingolstadt1. The total demands of the two road networks are respectively 1,716 and 3,031 vehicles.

Based on the analysis in Section 4.5.1, I computed the number of replications per cell as  $\lceil \frac{2,646}{300} \rceil = 9$ . I executed all replications using Python 3.9.16, SUMO 1.12.0, and a modified version of CityFlow 0.1 on a shared server with four cores, two 4.2GHz Intel i7-7700K processors per core, and 62 GiB of RAM.

### 4.5.2 Results

One-sample  $t$ -tests indicated that all RMSE and KL divergence measures were significantly different from 0, with a  $p$ -value  $\ll 0.001$  for all cells in Experiments 1 and 2. This suggests a lack of distributional equivalence between CityFlow and SUMO. The following subsections explore the results of second-order linear multiple regression for each of the measures. Notably, coefficients involving the difference between the SUMO and CityFlow Krauß model implementations were



(a) ingolstadt1



(b) ingolstadt7

Figure 4.6: Screenshots in CityFlow of the ingolstadt1 and ingolstadt7 road networks.

generally not significant.

In the following tables, abbreviations of variables follow Section 4.5.1. Subscripts denote the values of the levels. Car-following models are abbreviated as “dK” = default Krauß, “SK” = SUMO Krauß, “Wa” = Wagner, “Wi” = Wiedemann. Aggressiveness types are abbreviated as in Table 4.2. Road networks are abbreviated as “G” = grid4x4, “I7” = ingolstadt7.

### Experiment 1: Traffic Demand

For total time and waiting time RMSEs, the road network, car-following model, and aggressiveness generally had significant effects, as did various pairwise interactions between them. The uncongested grid4x4 network had significantly lower RMSEs (coefficients:  $-890.61/-1035.46$ ) than the congested arterial4x4 network (intercepts:  $1718.61/1815.46$ ), suggesting that congestion worsened the distributional equivalence of these measures. In arterial4x4, the Wagner (coefficients relative to SUMO Krauß:  $546.8/466.65$ ) and Wiedemann models (coefficients:  $100.96/208.99$ ) had significantly higher RMSEs; these differences were smaller for grid4x4 (coefficients relative to SUMO Krauß:  $49.5/33.75$  for Wagner;  $48.91/34.85$  for Wiedemann).

For total and queued vehicle count RMSEs and KL divergences, I generally found significant effects for the road network, car-following model, aggressiveness, and gap tolerance, along with various pairwise interactions between them. The RMSEs and KL divergences showed distinct patterns: the RMSEs were much lower for grid4x4 than arterial4x4 (coefficients:  $-5.15/-4.93$ ), but the KL divergences had less variation (coefficients:  $-0.621/-0.013$ ). Yet, the KL divergences had low enough standard deviations that the road network’s effects remained significant. High aggressiveness in arterial4x4 generally yielded higher RMSEs (coefficients for aggressive young relative to aggressive middle-aged:  $2.16/2.46$ ) and KL divergences (coefficients:  $0.228/0.139$ ). While the same was true for the RMSEs in grid4x4, its KL divergences were lower for more aggressive set-

tings (coefficients:  $-0.204/-0.122$ ). Despite higher time measures, the Wagner model led to lower measures for total vehicle count (coefficients relative to SUMO Krauß:  $-1.16/-0.193$ ).

For vehicle speed and acceleration RMSEs, the road network, car-following model, and aggressiveness generally had significant effects, as did various pairwise interactions between them. Greater equivalence in vehicle distributions did not always correspond to more similar vehicle-level measures. Both speed and acceleration RMSEs increased for grid4x4 (coefficients:  $2.62/0.418$ ) even though the other measures were lower on average. Also, unlike its vehicle count measures but like its time measures, the Wagner model had significantly higher speed and acceleration RMSEs (coefficients relative to SUMO Krauß:  $2.71/1.26$ ).

## Experiment 2: Network Scale

For total time and waiting time RMSEs, the aggressiveness and its interactions generally had significant effects, along with the Wiedemann model and its interactions. The best-fitting model for total time did not include a road network-gap tolerance interaction, whereas the model for waiting time did. The smaller ingolstadt1 network had lower but more variable RMSEs (intercepts:  $899.29/1170.3$ ), while the larger ingolstadt7 network had significantly higher but more uniform RMSEs (coefficients:  $1264.59/1046.07$ ). For ingolstadt1, the most aggressive parameter settings led to significantly higher RMSEs (coefficients of aggressive young relative to aggressive middle-aged:  $408.36/568.65$ ). Likewise, the Wiedemann model had significantly higher RMSEs than other car-following models in ingolstadt1 (coefficients relative to SUMO Krauß:  $192.29/1917.31$ ), but this effect was reversed for ingolstadt7 (coefficients:  $-349.13/-464.13$ ).

For total and queued vehicle count RMSEs and KL divergences, the road network and aggressiveness generally had significant effects, as did various pairwise interactions between them and with the gap tolerance. Unlike Experiment 1, the RMSE measures were lower in ingolstadt7 than in ingolstadt1 (coefficients:  $-2/-1.67$ ), but the KL divergence measures were higher (coefficients:  $0.461/2.48$ ). More aggressive parameter settings again led to significant increases in the RMSEs and KL divergences, with a larger increase in RMSEs (coefficients of aggressive young relative to aggressive middle-aged:  $1.35/2.29$ ) than in KL divergences (coefficients:  $0.595/0.717$ ). However, the increase in both measures was smaller in ingolstadt7 (coefficients:  $0.764/1.507$  for RMSE,  $0.283/0.413$  for KL divergences). Both the Wagner and ACC models had RMSEs that significantly increased with gap tolerance (coefficients relative to SUMO Krauß per unit of gap tolerance:  $0.151/0.056$  for Wagner,  $0.165/0.092$  for ACC), but KL divergences that significantly decreased with it (coefficients:  $-0.285/-0.511$  for Wagner,  $-0.354/-0.642$  for ACC).

For vehicle speed and acceleration RMSEs, the Wagner and Wiedemann models along with the aggressiveness had significant effects, as did various pairwise interactions of the road network with the car-following models and aggressiveness. Again, the road network had significant effects on both speed and acceleration RMSEs, with these measures being higher for ingolstadt7 (coefficients:  $5.51/0.173$ ). The Wagner car-following model had significantly higher RMSEs (coefficients relative to SUMO Krauß:  $0.363/0.461$ ), whereas the Wiedemann model had significantly lower RMSEs (coefficients:  $-1.94/-0.957$ ).

Table 4.3: Fitted regression coefficients for Experiment 1. Default levels are  $C = \text{ACC}$ ,  $A = \text{AM}$ ,  $R = \text{arterial4x4}$ . N/A denotes variables not included in the response surface. Significance levels: \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

Dep. var.	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Intercept	970.12***	923.27***	4.05***	2.58***	3.27***	2.43***	2.82***	4.92***
$R_G$	-940.8***	-928.6***	-3.4***	-0.289***	-2.96***	0.098***	2.1***	0.021
$C_{dK}$	-189.49***	-260.91***	2.61***	0.098***	1.94***	0.028*	0.821***	0.1***
$C_{SK}$	-192.08***	-263.95***	2.6***	0.102***	1.93***	0.029*	0.822***	0.119***
$C_{Wa}$	-788.82***	-734.96***	-2.03***	-0.291***	-2.05***	-0.224***	2.11***	1.48***
$C_{Wi}$	523.87***	528.37***	0.964***	0.446***	1.21***	0.394***	-0.309***	-0.073***
$A_{AO}$	15.1	-9.09	-0.026	0.082***	-0.182***	0.084***	-0.265***	-0.136***
$A_{AY}$	114.09***	121.5***	0.512***	0.002	0.631***	0.006	0.113***	0.067**
$A_{CM}$	-26.02	-62.04*	-0.143***	0.024**	-0.211***	0.038***	-0.291***	-0.175***
$A_{CO}$	-34.31	-121.33***	-0.154***	0.056***	-0.463***	0.074***	-0.408***	-0.194***
$A_{CY}$	-3.25	-7.94	-0.076*	-0.018*	0.056	-0.013	-0.135***	-0.046*
$R_G C_{dK}$	154.88***	229.28***	-2.61***	-0.081***	-1.94***	-0.012	-0.766***	-0.125***
$R_G C_{SK}$	155.5***	230.53***	-2.61***	-0.084***	-1.94***	-0.01	-0.764***	-0.125***
$R_G C_{Wa}$	775.41***	702.88***	1.81***	0.241***	1.79***	0.167***	-2.15***	-1.45***
$R_G C_{Wi}$	-267.66***	-144.34***	-0.775***	-0.315***	-0.564***	-0.265***	0.344***	0.008
$R_G A_{AO}$	-48.38***	-17.35	-0.087**	-0.087***	0.088*	-0.038***	0.179***	0.072***
$R_G A_{AY}$	-60.32***	-75.11***	-0.43***	-0.004	-0.521***	0.069***	-0.234***	-0.15***
$R_G A_{CM}$	-27.02*	-2.7	0.088**	-0.031***	0.162***	0.036***	0.278***	0.11***
$R_G A_{CO}$	-61.15***	17.82	-0.012	-0.058***	0.26***	0.001	0.418***	0.157***
$R_G A_{CY}$	-26.89*	-45.34**	0.134***	0.045***	-0.05	0.075***	0.136***	0.021
$C_{dK} A_{AO}$	45.19**	10.91	-0.038	-0.007	-0.159*	-0.028*	-0.096***	0.067***
$C_{SK} A_{AO}$	47.08**	12.14	-0.034	-0.012	-0.158*	-0.032**	-0.13***	0.058**
$C_{Wa} A_{AO}$	-4.14	10.12	0.542***	0.069***	0.601***	0.095***	0.147***	-0.007
$C_{Wi} A_{AO}$	40.85*	47.8*	0.103*	-0.006	0.194**	-0.023*	0.056*	0.086***
$C_{dK} A_{AY}$	-64.3***	-52.04*	0.058	-0.033**	0.072	-0.001	0.192***	-0.003
$C_{SK} A_{AY}$	-61.46***	-48.95*	0.064	-0.035***	0.076	-0.004	0.178***	-0.011
$C_{Wa} A_{AY}$	-71.42***	-74.05**	-0.486***	-0.007	-0.556***	-0.005	-0.165***	-0.034*
$C_{Wi} A_{AY}$	-56.73**	-50.16*	-0.087	-0.004	-0.163*	0.015	-0.001	0.003
$C_{dK} A_{CM}$	30.41	28.05	-0.087	-0.019*	-0.166**	-0.041***	-0.104***	0.073***
$C_{SK} A_{CM}$	34.88*	32.75	-0.08	-0.023*	-0.16*	-0.045***	-0.113***	0.066***
$C_{Wa} A_{CM}$	13.33	42.11	0.479***	0.127***	0.533***	0.158***	0.168***	0.031
$C_{Wi} A_{CM}$	80.26***	96.42***	0.005	0.032**	0.072	-0.007	0.001	0.118***
$C_{dK} A_{CO}$	116.82***	90.82***	-0.074	-0.005	-0.153*	-0.021	-0.193***	0.074***
$C_{SK} A_{CO}$	115.25***	89.26***	-0.074	-0.008	-0.155*	-0.022	-0.2***	0.072***
$C_{Wa} A_{CO}$	68.58***	124.53***	0.991***	0.133***	1.17***	0.181***	0.237***	-0.015
$C_{Wi} A_{CO}$	81.46***	98.69***	0.039	0.005	0.195**	-0.025*	0.041	0.134***
$C_{dK} A_{CY}$	-0.113	19.28	-0.041	-0.044***	0.098	-0.025*	-0.018	0.004
$C_{SK} A_{CY}$	3.37	22.38	-0.031	-0.044***	0.108	-0.023*	-0.022	0.006
$C_{Wa} A_{CY}$	-3.99	0.919	-0.089*	0.003	-0.128*	0.017	-0.013	0.033
$C_{Wi} A_{CY}$	39.96*	54.56*	-0.081	0.026**	-0.068	0.024*	-0.029	0.068***
$L$	-87.49***	-120.24***	-0.038	-0.019**	-0.167***	-0.022**	-0.062**	-0.034*
$R_G L$	123.75***	180.39***	0.068**	0.041***	0.286***	0.05***	0.041**	0.026**
$C_{dK} L$	13.24	15.42	0.028	0.003	0.049	0.001	-0.021	-0.014
$C_{SK} L$	13.38	15.46	0.034	0.001	0.054	0.001	-0.009	-0.025
$C_{Wa} L$	12.96	14.77	-0.021	-0.006	-0.015	-0.019	-0.028	-0.034*
$C_{Wi} L$	-287.21***	-425.27***	-0.187***	-0.142***	-0.68***	-0.129***	-0.047*	-0.004
$A_{AO} L$	8.21	9.99	N/A	N/A	N/A	N/A	0.08**	0.013
$A_{AY} L$	-5.87	-3.05	N/A	N/A	N/A	N/A	0.056*	0.068***
$A_{CM} L$	25.87	26.58	N/A	N/A	N/A	N/A	0.048*	0.033*
$A_{CO} L$	24.55	24.98	N/A	N/A	N/A	N/A	0.05*	0.021
$A_{CY} L$	26.82	35.96	N/A	N/A	N/A	N/A	0.047	0.034*



Table 4.4: Fitted regression coefficients for Experiment 2. Default levels are  $C = \text{ACC}$ ,  $A = \text{AM}$ ,  $R = \text{ingolstadt1}$ . N/A denotes variables not included in the response surface. Significance levels: \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

Dep. var.	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Intercept	7.6	22.81	2.28***	3.07***	1.21***	2.68***	5.8***	1.75***
$R_{17}$	426.11***	367.35***	-0.545***	1.71***	0.091*	2.44***	-0.476***	-0.025
$C_{dK}$	61.93	68.9	0.924***	-0.018	1.05***	0.099*	-0.202*	-0.033
$C_{SK}$	47.15	51.86	0.87***	-0.021	0.993***	0.104*	-0.126	-0.023
$C_{Wa}$	-15.79	1.18	-0.637***	-0.616***	-0.19**	-0.378***	0.545***	0.473***
$C_{Wi}$	1389.63***	1917.31***	2.24***	0.035	3.38***	-0.169***	0.246**	-0.302***
$A_{AO}$	28.69	31.65	0.317***	0.058	0.222**	0.031	0.035	-0.129***
$A_{AY}$	-84.5*	-90.65*	-0.384***	-0.104**	-0.419***	-0.194***	0.315***	0.157***
$A_{CM}$	12.43	16.06	0.366***	-0.177***	0.383***	-0.17***	-0.719***	-0.321***
$A_{CO}$	34.07	33.49	0.44***	-0.085*	0.36***	-0.073	-0.579***	-0.339***
$A_{CY}$	-47.2	-48.54	-0.104	-0.014	-0.054	0.135**	-0.553***	-0.171***
$R_{17}C_{dK}$	-92.99***	-79.78***	-0.913***	-0.209***	-0.889***	-0.291***	0.28***	0.075***
$R_{17}C_{SK}$	-113.04***	-102.86***	-0.936***	-0.219***	-0.917***	-0.295***	0.284***	0.087***
$R_{17}C_{Wa}$	-104.46***	-99.8***	0.051	0.423***	-0.257***	0.248***	1.01***	0.092***
$R_{17}C_{Wi}$	171.07***	-49.3*	-0.812***	0.156***	-1.44***	-0.11***	0.564***	-0.071***
$R_{17}A_{AO}$	-39.49	-54.69*	-0.3***	-0.07***	-0.221***	-0.114***	0.055	0.04***
$R_{17}A_{AY}$	-143.27***	-182.13***	-0.066*	0.092***	-0.164***	0.141***	-0.099*	0.093***
$R_{17}A_{CM}$	214.11***	267.24***	-0.04	0.059**	0.093*	0.056*	-0.343***	-0.162***
$R_{17}A_{CO}$	153.16***	201.83***	-0.224***	0.042*	-0.032	0.033	-0.489***	-0.144***
$R_{17}A_{CY}$	210.73***	260.83***	0.151***	0.047**	0.203***	-0.033	-0.228***	-0.157***
$C_{dK}A_{AO}$	-24.91	-29.01	0.125*	0.03	0.037	-0.021	-0.023	-0.044*
$C_{SK}A_{AO}$	1.49	0.955	0.157**	0.045	0.074	-0.015	0.017	-0.041*
$C_{Wa}A_{AO}$	-40.91	-29.49	-0.141**	-0.048	-0.103	0.068	-0.219**	-0.056**
$C_{Wi}A_{AO}$	11.65	39.36	-0.083	-0.056*	-0.036	-0.051	0.253***	0.024
$C_{dK}A_{AY}$	62.86	68.72	-0.086	-0.057*	0.013	-0.064	-0.09	0.037*
$C_{SK}A_{AY}$	92.72**	107.09**	-0.023	-0.04	0.082	-0.044	-0.06	0.033
$C_{Wa}A_{AY}$	46.29	52.08	0.072	-0.064*	0.105	-0.155***	-0.208**	-0.023
$C_{Wi}A_{AY}$	148.65***	213.29***	0.608***	0.052	0.712***	-0.081*	0.085	-0.052**
$C_{dK}A_{CM}$	68.22*	69.16	0.359***	0.113***	0.309***	0.054	0.469***	0.039*
$C_{SK}A_{CM}$	87.22**	92.12*	0.385***	0.135***	0.334***	0.074*	0.545***	0.048**
$C_{Wa}A_{CM}$	79.84*	78.72*	0.049	0.231***	0.04	0.297***	0.618***	0.074***
$C_{Wi}A_{CM}$	-164***	-211.64***	-0.402***	-0.072**	-0.483***	0.047	-0.238**	0.035*
$C_{dK}A_{CO}$	21.83	19.98	0.407***	0.07*	0.297***	-0.003	0.448***	0.015
$C_{SK}A_{CO}$	34.23	35.11	0.424***	0.068*	0.322***	-0.017	0.45***	0.013
$C_{Wa}A_{CO}$	43.95	47.56	-0.042	0.156***	-0.02	0.234***	0.457***	0.021
$C_{Wi}A_{CO}$	-164.22***	-204.38***	-0.464***	-0.105***	-0.508***	0.04	-0.257***	0.023
$C_{dK}A_{CY}$	0.075	-4.09	-0.055	0.055*	-0.029	0.037	0.4***	0.104***
$C_{SK}A_{CY}$	21.65	21.13	-0.043	0.037	-0.011	0.012	0.454***	0.115***
$C_{Wa}A_{CY}$	53	44.29	0.034	0.114***	0.009	0.012	0.57***	0.056**
$C_{Wi}A_{CY}$	-188.69***	-256.56***	-0.374***	-0.128***	-0.476***	-0.048	-0.495***	-0.02
$L$	52.1	32.01	-0.081	0.011	-0.105*	-0.031	-0.031	0.011
$R_{17}L$	-105.17***	-48.53*	0.092**	-0.02	0.155***	0.022	N/A	0.039***
$C_{dK}L$	-49.19	-58.58	-0.429***	0.06*	-0.459***	0.14***	0.047	0.027
$C_{SK}L$	-53.48	-64.42*	-0.414***	0.051*	-0.445***	0.121***	-0.048	0.014
$C_{Wa}L$	-25.45	-33.98	0.007	-0.091***	0.002	-0.07*	-0.103	-0.023
$C_{Wi}L$	-694.07***	-923.7***	-0.568***	-0.042	-0.817***	0.17***	-0.059	0.2***
$A_{AO}L$	-13.86	-27.5	-0.084	-0.017	-0.105	-0.032	-0.001	0.002
$A_{AY}L$	34.56	39.67	0.166**	0.012	0.17**	0.01	-0.022	-0.036*
$A_{CM}L$	-24.34	-32.66	-0.178***	0.119***	-0.227***	0.147***	0.302***	0.05**
$A_{CO}L$	-12.05	-21.61	-0.068	0.084**	-0.107	0.099**	0.235***	0.016
$A_{CY}L$	51.39	54.83	0.164**	0.032	0.152**	0.036	0.24***	-0.002

### 4.5.3 Parameter Validity

I conducted parameter validation by comparing the parameter settings for car-following and lane-changing models used in the experiments with prior literature. In general, the dependence of these parameters on external factors such as traffic density and speed suggests that the settings used in traffic simulators should be calibrated to specific road networks and conditions. However, the settings I used remain reasonable considering the variation reported in the literature.

For maximum acceleration and deceleration, my experiments respectively varied them from  $2.3 \text{ m/s}^2$  to  $3.1 \text{ m/s}^2$  and from  $-3.8 \text{ m/s}^2$  to  $-5.5 \text{ m/s}^2$ . These settings were based on Capela Dias et al. [59], with smaller values for older and less aggressive drivers (Section 4.4.1); however, they considered gender effects to be negligible. Similar values have been reported in prior work [119, 204, 296]. However, among driving simulator studies, Körber et al. [204] demonstrated an age effect opposite to that assumed by Capela Dias et al. [59]. Among real-world studies, Moon and Yi [283] found a dependence of the 95th percentile of braking decelerations on speed, and Nakagawa et al. [296] reported a significant interaction between age and gender.

For minimum gap and headway time, my experiments respectively varied them from 1.2 m to 2.5 m and from 1.0 s to 1.9 s. This again followed Capela Dias et al. [59], with larger values for older and less aggressive drivers (Section 4.4.1). Similar values have been reported in prior work [99, 272, 348]. However, among driving simulator studies, Körber et al. [204] demonstrated an age effect opposite to that assumed by Capela Dias et al. [59]. Among real-world studies, Moon and Yi [283] found a dependence of the 95th percentile of braking decelerations on speed, and Nakagawa et al. [296] reported a significant interaction between age and gender.

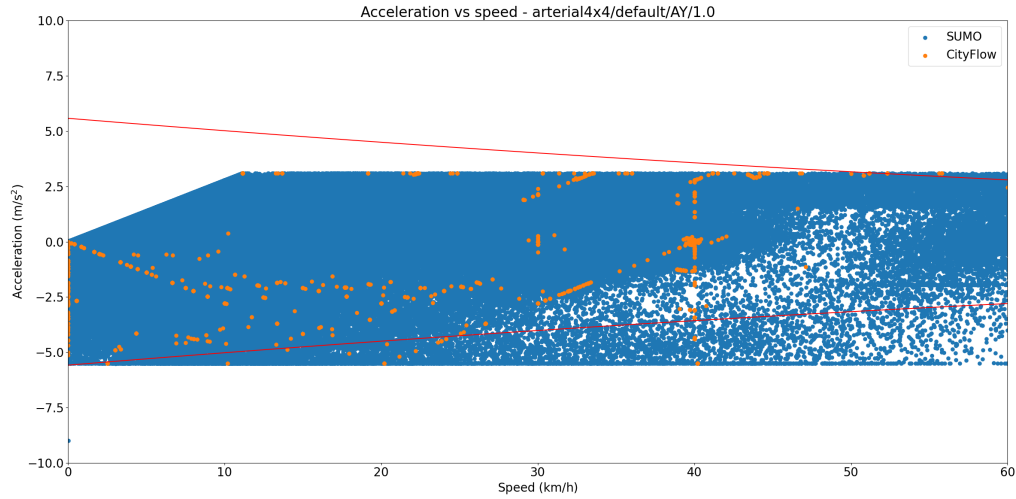
Lane-changing gap tolerance can be approximately quantified by variance in accepted gap widths from empirical lane-changing behaviour. In my experiments, I varied it upward and downward by 18% and 50%. As described in Section 4.4.2, Sun and Kondyli [389]’s real-world study reported the mean gaps for forced, cooperative, and free lane changes to be 45 ft, 53 ft, and 109 ft, representing increments of 18% and 105%. These settings are consistent with standard deviations in lead and lag gaps as reported in prior studies [19, 285]. Ali et al. [7]’s driving simulator study identified significant factors that impact gap tolerance: relative to the average male middle-aged driver, gaps are smaller for younger drivers, larger for female drivers, and smaller as speed increases. Likewise, Hill et al. [164]’s real-world study reported that the mean and standard deviation of lag gaps depended on congestion. Future work could use these factors to create a taxonomy of lane-changing behaviour similar to Capela Dias et al. [59].

### 4.5.4 Acceleration-Speed Diagrams

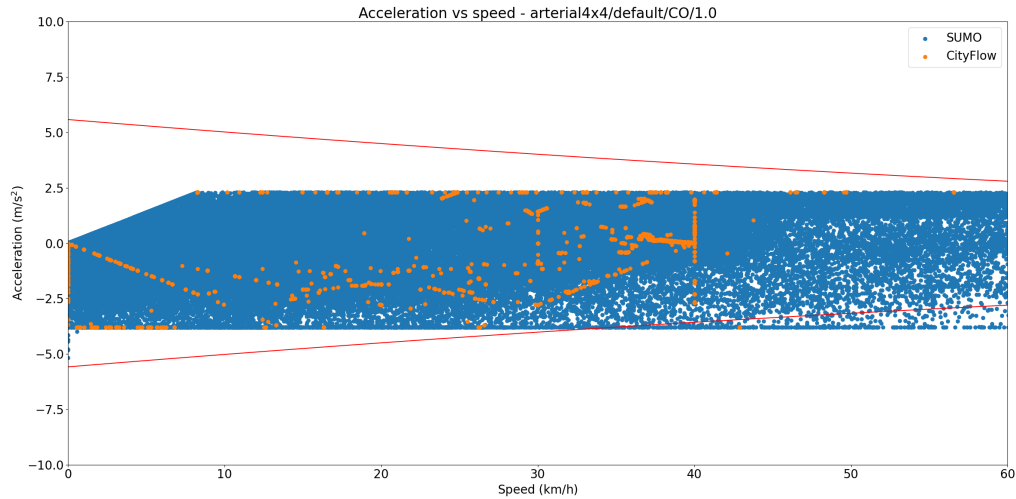
Based on movement and friction experienced by vehicles, Eboli et al. [111] derived the following upper bound on the safe magnitude of acceleration and deceleration as a function of speed:

$$|a| = g \cdot \left[ 0.198 \left( \frac{v}{100} \right)^2 - 0.592 \left( \frac{v}{100} \right) + 0.569 \right]$$

where  $a$  is the acceleration in  $\text{m/s}^2$ ,  $v$  is the speed in  $\text{km/h}$ , and  $g = 9.81 \text{ m/s}^2$  is the gravitational acceleration. The bound yields the point at which movement and friction forces are in balance for a vehicle. I plot this bound (in red) against actual acceleration-speed points from the two simulators in Figure 4.7. The most aggressive setting of car-following parameters clearly differs from the least aggressive in that it causes vehicles to drive closer to their maximum accelerations and exceed safe deceleration limits more often at low speed. The diagrams also show that perceptual randomness in SUMO likely yielded a much wider spread of points than in CityFlow.



(a) Aggressive young (AY) drivers



(b) Courteous old (CO) drivers

Figure 4.7: Acceleration-speed diagrams comparing outcomes of CityFlow (orange) and SUMO (blue) for aggressive young (AY) and courteous old (CO) drivers under the Krauß car-following model in the arterial4x4 network.

## 4.6 Conclusion

In this chapter, I designed experiments to compare the low-level simulation outcomes of two traffic simulators, CityFlow and SUMO. To capture the effects of modelling real-world heterogeneity, various parameters of driver behaviour and road network scale were varied. The results indicate a lack of distributional equivalence between the simulators, with certain parameter settings worsening distributional equivalence. However, as I noted in Section 4.3, these experiments were insufficient to provide a complete characterisation of what the critical differences between these simulators are. Many aspects of CityFlow and SUMO that were not controlled — simulation control-flow, other aspects of driver behaviour, the effects of traffic signals, and randomness — could all have contributed to the observed discrepancies. Future work should perform more comprehensive, controlled evaluations of these two simulators.

Regardless, researchers in RL for TSC must not take traffic simulators for granted as a *deus ex machina* for training, and must recognise that they may not be interchangeable. Which simulator, then, should be chosen? I do not aim to answer this question with my work in this chapter, but some observations can be made:

- **SUMO** provides a detailed simulation that models real-world heterogeneity, and captures additional aspects of traffic management and driver behaviour
- **CityFlow** provides an efficient simulation that abstracts out and homogenises various details, reducing the number of parameters that need to be tuned

Mei et al. [268] showed that, if the `libsumo` library is used as an API, SUMO’s runtime can be reduced tenfold. However, they also found that it remains slower than CityFlow for various environments. Therefore, the core trade-off between these two simulators (and others) involves realism and efficiency. There is no one best simulator; researchers must decide whether using a coarser abstraction of the environment is acceptable in exchange for faster training.

But how exactly should researchers make this decision? Crucially, RL-based signal plans may not necessarily perform better when they are trained with more granular simulators. Both introducing unnecessary complexity [496] and removing needed complexity [134] in the observation space may harm performance. Mei et al. [268] found no consistent ordering in performance across environments between RL policies trained with CityFlow or SUMO. Following their example, researchers should compare training results from different simulators and use them to design RL formulations in a principled way. One strategy may be to train a baseline using an efficient simulator (e.g. CityFlow), and then to finetune it by further training with a realistic one (e.g. SUMO).

Ultimately, the true goal is to ensure that RL-based signal plans can perform well in real-world traffic conditions, which requires that their training and validation environments are close to reality. As I previously noted, distributional equivalence between simulators is a proxy measure of this goal. Unfortunately, a chicken-and-egg problem exists in that traffic simulators are intended to replace real-world data collection, yet cannot be validated without it. For now, simulations should still be developed in collaboration with stakeholders to ensure that they meet acceptable standards of fidelity. However, the future holds promise for both realistic and efficient traffic simulations: the increasing prevalence of connected vehicles [261] means that the collection of granular real-world data for validation may be within reach.

# Chapter 5

## Out of the Past

### Building an AI-Enabled Pipeline for Traffic Simulation from Noisy, Multimodal Detector Data and Stakeholder Feedback

*Domain:* Traffic signal control  
*Challenges:* Uncertainty, heterogeneity

When it's playing out right in front of you, it feels like you can't escape it.

---

*Philip J. Coulson*

#### 5.1 Introduction

In Chapter 4, I studied how heterogeneity is modelled by different traffic simulators. I also argued that realistic, granular simulators are necessary to ensure that traffic signal plans created by reinforcement learning (RL) policies can perform well under real-world traffic conditions. However, a realistic *simulator* alone is insufficient to achieve this goal, because there must also be a realistic *simulation* that is executed in the simulator to generate traffic for the policy. The construction of simulations grounded in data from physical traffic systems is an understudied problem, but it is crucial to ensuring that the environment with which the RL algorithm interacts is close to reality – and, accordingly, that the results of RL training and evaluation are trustworthy for stakeholders.

Existing approaches to creating road network-scale traffic simulations have a number of limitations that hamper their realism, and thus their practical applicability. I focus on limitations surrounding the central task of *demand modelling*, or the modelling of traffic volumes within the simulation. Demand modelling methods that construct origin-destination matrices from activity data are unrealistic and fail to make use of traffic detector data. Meanwhile, detector data-driven approaches to demand modelling have relied on suboptimal heuristics. All of these approaches also consider the source data to be the ground truth; they do not perform any calibration to account for sources of *uncertainty* (noise) or *heterogeneity* (multimodality) in the data. In this chapter, I address these two challenges by seeking to answer the following research questions:

**RESEARCH QUESTION 5.1.** *Is it possible to correct for error in camera detection so that it can provide accurate vehicle counts for traffic simulation?*

**RESEARCH QUESTION 5.2.** *Can different sources of traffic detector data, each having different error rates and assumptions, be integrated to produce a single unified traffic simulation?*

**RESEARCH QUESTION 5.3.** *How can unstructured feedback from human stakeholders be used for the calibration of traffic simulations?*

To address these research questions, I contribute a detailed, systematic pipeline for modelling demand in a traffic simulation from noisy, multimodal detector data (Figure 5.1). Starting from raw detector data, my pipeline consists of three steps: (1) I apply a vehicle tracking-based computer vision method directly to camera footage, to obtain more accurate vehicle counts than the camera detectors themselves. (2) I solve a quadratic optimisation program to populate my simulation with vehicle routes. In doing so, I account for multimodality by imposing multiple sets of optimisation constraints based on different sources of vehicle counts. (3) I incorporate feedback from stakeholders to refine the simulation, using a large language model (LLM) agent that encodes natural language feedback into optimisation constraints.

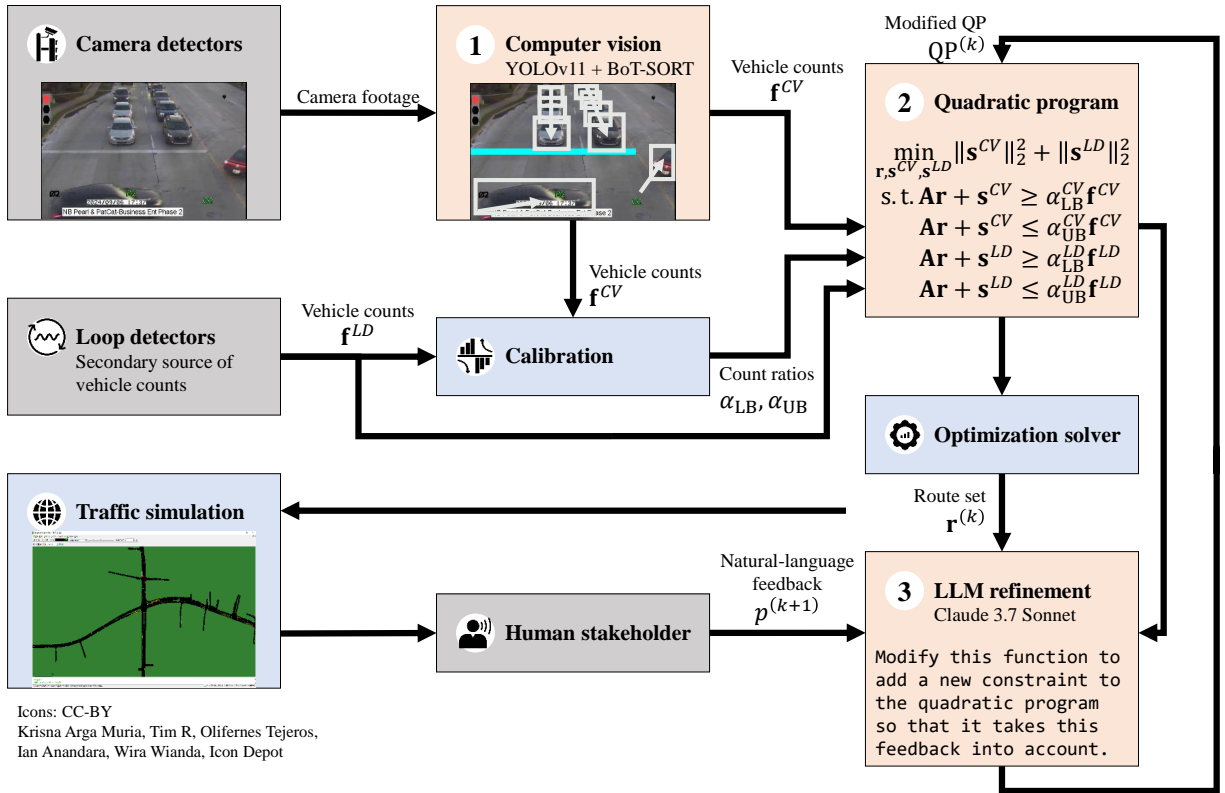


Figure 5.1: My pipeline for generating a traffic simulation from multimodal detector data.

As a proof of concept, I apply my pipeline to simulate a high-traffic road network from the city of Strongsville, Ohio. Beginning with 24 hours of recorded camera footage and detector data from 36 intersections, I created a fully realised traffic simulation. This chapter shows that: (1) My vehicle tracking-based computer vision method rectifies undercounting in camera detector data.

- (2) My optimisation method is able to generate a set of vehicle routes that is consistent with counts from both computer vision and loop detector counts, while still accounting for error in these counts.
- (3) My LLM agent is able to synthesise code representing sensible, quantified constraints based on qualitative stakeholder feedback.

This chapter was published at the Winter Simulation Conference (WSC) in 2025 [72].

## 5.2 Related Work

### 5.2.1 Demand Modelling for Traffic Simulation

*Demand modelling*, or the modelling of trips taken by individual vehicles from one point to another in a road network, is a central but difficult aspect of constructing data-driven traffic simulations [25]. One popular approach is *activity modelling*, where trips are extrapolated from censuses of the daily activities taken by a sample of households in the study region [44, 82, 83, 161, 222, 414]. While many municipalities collect this data for transportation planning, the locations of these activities are usually coarsely discretised. Furthermore, they represent a small, not necessarily representative sample of the population. This means that activity modelling-based simulations are prone to significant error [213]; I do not consider activity modelling in this chapter.

An alternative approach to demand modelling directly uses data from traffic detectors. The types of detectors used for simulations include induction loops [40] and video cameras [434, 455, 495] (see further discussion in Section 2.2.2). Detectors provide granular vehicle counts local to individual intersections, but converting them to fully-realised routes through a road network is nontrivial. Wei et al. [434], Xu et al. [455], and Zheng et al. [495] all provided no details on how they generated vehicle routes. The route generation procedures that have been specified in prior work have relied on suboptimal heuristics. Lobo et al. [242] and Rapelli et al. [324] used detector data to adjust activity models. Bieker et al. [40] generated routes probabilistically by using turn ratios to define distributions over movements at intersections. This approach only leads to correct simulation outcomes in expectation. Finally, Qiu et al. [316]’s approach, which uses scripts included with the traffic simulator SUMO [10], is most similar to mine: they applied a two-step process of first sampling routes randomly, and then solving a linear program (LP) to approximate how many times each route should be used to match the detector counts as closely as possible. Unlike them, I solve the problem exactly as a quadratic integer program (QIP) without intermediate approximations.

One further limitation of previous detector-based approaches is that they rely on a single source of detector data, which is assumed to be generally error-free. Among the works cited previously, only Bieker et al. [40] reported detectors that failed to report vehicle counts and were removed from their dataset. When these works validate their simulations, they treat the detector data as the ground truth to compare their simulations against [40, 82, 222, 242, 414]. By contrast, my approach (Section 5.3.2) integrates multiple methods of processing detector data, and I adopt a semi-automated approach to validation that combines detector data with manual verification. This approach allowed me to obtain higher realism in the vein of traffic simulations that rely exclusively on manual counting [290], but at a much larger scale.

## 5.2.2 Computer Vision for Traffic Footage

*Vehicle counting* can be decomposed into two distinct but related problems: *vehicle detection*, the identification of vehicles in footage; and *vehicle tracking*, the identification of these vehicles’ trajectories across frames [404]. These methods can be divided into two distinct waves of research. First, in the 1990s, advances in image processing led to vehicle detection methods based on extracting heuristically designed features [84]. Second, in the 2010s, the advent of *convolutional neural networks* (CNNs) led to various deep object detection algorithms capable of automatically extracting relevant features for vehicle detection [423]. While these lines of work have more recently overlapped methodologically, they have not been comparatively evaluated to my knowledge. I provide an in-situ evaluation of AutoScope, a widely-deployed image processing method, against CNN-based counting.

## 5.2.3 Large Language Models for Transportation Research

*Large language models* (LLMs) are useful for aligning AI systems with human intuition. As such, they have been increasingly applied to the domain of transportation. One line of work has generated simulation scenarios [62, 224, 391] and reward functions for vehicular agents [157, 504] based on natural language prompts. Another line of work has used LLMs to align simulations with reality based on general knowledge; Da et al. [90] used an LLM to infer how actions taken in a simulated environment would affect a real-world environment differently.

When responding to prompts, LLMs can use external tools. For instance, Li et al. [224]’s agent generates and executes command-line calls to the SUMO simulator; Wang et al. [424]’s agent chooses between different perception and decision tools to perform reinforcement learning for traffic signal control. I leverage a strength of LLMs that has not been explored for transportation research to my knowledge: the synthesis of syntactically and semantically correct programs [16]. I use an LLM within an iterative framework in which the *only* external input required is qualitative natural language feedback.

## 5.3 Demand Modelling Pipeline

### 5.3.1 Computer Vision-Based Vehicle Counting from Camera Footage

In Section 2.2.2, I noted that the camera detectors deployed in most municipal transportation systems — such as those based on AutoScope — are primarily meant to detect the presence or absence of vehicles, and do not necessarily provide accurate counts. To detect individual vehicles more accurately than the camera detectors themselves, I use the YOLOv11 object detection model [185]. In each frame of footage, the model predicts a set of bounding boxes, each of which encloses an object. For each box, the model outputs a class  $b_c$ , the  $x$  and  $y$  coordinates of the centre of the bounding box ( $b_x$  and  $b_y$ ), and the width and height of the box ( $b_w$  and  $b_h$ ).

How can YOLO’s detections of vehicles be converted to counts for each lane? I manually annotate each frame of the traffic footage with the stop bar’s  $y$ -coordinate ( $S_y$ ), and with  $x$ -coordinates for each lane ( $S_L^\ell, S_R^\ell$ ). The most straightforward method to perform counting is to verify whether



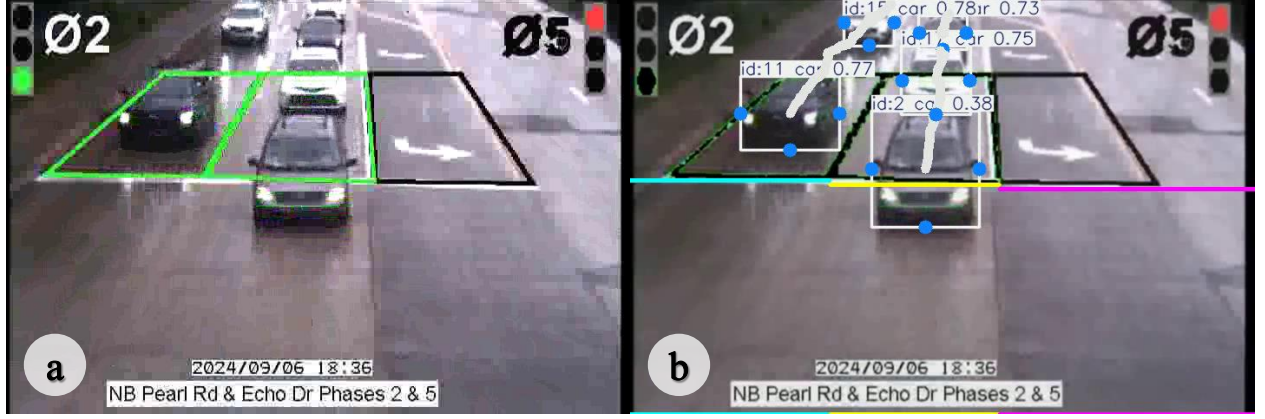


Figure 5.2: Demonstration of my vehicle tracking method on camera detector footage from intersection 8 (US 42 & Echo Rd) in Strongsville, Ohio. (a) Raw footage, showing two vehicles actuating a detection zone in the centre lane. (b) Footage with preprocessing filters applied, and bounding boxes and tracks for counted vehicles annotated. The coloured lines represent manually labelled stop bar positions.

the bounding box for an object identified as a vehicle has crossed the position of the stop bar, i.e. I increment the count for lane  $\ell$  when a detected vehicle has  $S_L^\ell \leq b_x \leq S_R^\ell$ , and  $|b_y - S_y| \leq \epsilon$  for some predefined threshold  $\epsilon$ . When I applied this method in practice, I encountered two issues. (1) Due to instability in the real-time streaming (RTSP) connection over which the detector footage is retrieved, frames are frequently dropped. For many vehicles, this leads to the absence of the frames in which their bounding boxes' borders  $b_y$  are close to the stop bar  $S_y$ . These vehicles first appear far above the stop bar ( $b_y \ll S_y$ ), and then far below the stop bar ( $b_y \gg S_y$ ) once the footage resumes. (2) Pixelation artifacts, particularly around the detection zones marked on the footage, also obfuscate the bounding boxes.

To address dropped frames, I use the BoT-SORT algorithm [4] to perform vehicle tracking. BoT-SORT reidentifies each bounding box across consecutive frames to provide a consistent ID  $b_t$ . In each frame, for each bounding box, I verify whether  $b_t$  has already been counted. To ensure that I do not capture traffic in other directions, I only consider  $b_t$  if  $b_y < S_y$  initially. If  $b_t$  has not yet been counted and  $b_y > S_y$ , I increment the vehicle count and mark  $b_t$  as counted. Even if the footage is missing the frame where  $|b_y - S_y| \leq \epsilon$  for a vehicle, the vehicle will be counted in a later frame. I also apply filters to smooth the footage as a preprocessing step, including a non-local means filter [51], a spatiotemporal denoising filter (hqdn3d), a frame-blending motion interpolation filter (minterpolate), and a filter to remove detector actuation overlays. The results are shown in Figure 5.2b.

### 5.3.2 Optimisation-Based Vehicle Route Generation from Multimodal Data

My vehicle tracking-based computer vision method from Section 5.3.1 outputs vehicle counts  $f_j^{CV}$  for a set of counting locations  $j$ . These represent the traffic volumes at the eastbound, northbound, southbound, and westbound approaches for each intersection (if they are available). I also have a

set of counts from loop detectors  $f_j^{LD}$ , which overlap with the computer vision counts at a subset of counting locations. But how can these multimodal vehicle counts be integrated to generate vehicle routes for a fully-realised simulation?

I make two key assumptions to identify the set of feasible routes. (1) Given an origin and destination in the road network, I assume that vehicles perform shortest-path (Dijkstra) routing. (2) I assume that most of the traffic in the road network originates at the fringes of the network, and few routes begin and end in the middle of a road edge (representing traffic from unmodelled driveways). This assumption holds as long as all major sources and sinks of traffic are modelled. Based on these assumptions, I enumerate the full set of routes between all counting locations, instead of randomly sampling them as in prior work.

Given this set of feasible routes, I aim to solve for the number of times each route should be used, so that the number of times they pass through the counting locations match the given vehicle counts as closely as possible. I do so by dividing 24 hours of count data into 15-minute time segments, which are indexed as  $t \in \{0, \dots, 95\}$ . For each time segment, I match the total counts at each counting location, as well as counts for dedicated left-turn and right-turn lanes if they exist.

How closely should the counts be matched? I assume that, for some counting locations  $j$  and time segments  $t$ , there are ground truth counts  $f_{jt}^M$ , and that error exists in both my computer vision counts  $f_{jt}^{CV}$  and loop detector counts  $f_{jt}^{LD}$ . Let  $M$ ,  $CV$ , and  $LD$  denote the sets of locations for these sources. Based on how much computer vision overcounts or undercounts for locations shared with the ground truth, I extrapolate conservative lower and upper bounds for the true counts:

$$\alpha_{LB}^{CV} = \min_t \min_{j \in M \cap CV} \frac{f_{jt}^M}{f_{jt}^{CV}}, \alpha_{UB}^{CV} = \max_t \max_{j \in M \cap CV} \frac{f_{jt}^M}{f_{jt}^{CV}}.$$

I also derive bounds for loop detectors,  $\alpha_{LB}^{LD}$  and  $\alpha_{UB}^{LD}$ , in a similar fashion. Then, I assume that, for time segments  $t \in \{0, \dots, 95\}$ ,  $\alpha_{LB}^{CV} f_{jt}^{CV} \leq f_{jt}^M \leq \alpha_{UB}^{CV} f_{jt}^{CV}$ ,  $\forall j \in CV$ , and  $\alpha_{LB}^{LD} f_{jt}^{LD} \leq f_{jt}^M \leq \alpha_{UB}^{LD} f_{jt}^{LD}$ ,  $\forall j \in LD$ .

Now, for each 15-minute time segment  $t$ , I used the solver Gurobi to solve the following quadratic integer program (QIP), where the decision variable is the number of usages  $r_{it}$  for each route  $i \in \{1, \dots, n\}$ :

$$\min_{\mathbf{r}_t, \mathbf{s}_t^{CV}, \mathbf{s}_t^{LD}} \|\mathbf{s}_t^{CV}\|_2^2 + \|\mathbf{s}_t^{LD}\|_2^2 + \lambda_{\text{nonfringe}} \sum_{i \in \text{nonfringe}} r_{it} + \lambda_{\text{temporal}} \|\mathbf{r}_t - \mathbf{r}_{t-1}\|_2^2 \quad (5.1a)$$

$$\text{s.t. } \mathbf{A}\mathbf{r}_t + \mathbf{s}_t^{CV} \geq \alpha_{LB}^{CV} \mathbf{f}_t^{CV} \quad (5.1b)$$

$$\mathbf{A}\mathbf{r}_t + \mathbf{s}_t^{CV} \leq \alpha_{UB}^{CV} \mathbf{f}_t^{CV} \quad (5.1c)$$

$$\mathbf{A}\mathbf{r}_t + \mathbf{s}_t^{LD} \geq \alpha_{LB}^{LD} \mathbf{f}_t^{LD} \quad (5.1d)$$

$$\mathbf{A}\mathbf{r}_t + \mathbf{s}_t^{LD} \leq \alpha_{UB}^{LD} \mathbf{f}_t^{LD} \quad (5.1e)$$

$$\mathbf{r}_t \in (\mathbb{Z}^{\geq 0})^n,$$

where  $\mathbf{A} \in \{0, 1\}^{n \times m}$  is a binary matrix denoting which counting locations are used by routes:  $\mathbf{A}_{ij}$  is 1 if route  $i$  passes counting location  $j$ , and is 0 otherwise, such that  $\mathbf{A}\mathbf{r}_t$  gives the number of times the generated routes collectively pass each counting location  $j \in \{1, \dots, m\}$ ;  $\mathbf{s}_t^{CV}, \mathbf{s}_t^{LD} \in \mathbb{R}^m$  are

slack variables that represent the error between the generated routes' counts  $\mathbf{Ar}_t$  and the actual counts  $\mathbf{f}_t^{CV}$  and  $\mathbf{f}_t^{LD}$ ; nonfringe is the set of indices for routes where the start or the end of the route are interior edges in the road network;  $\lambda_{\text{nonfringe}}$  is a hyperparameter for weighting the objective function penalty for these routes; and  $\lambda_{\text{temporal}}$  is a hyperparameter for penalising discrepancies between the generated routes of adjacent time segments.

In the QIP, constraints (5.1b) and (5.1c) specify that  $\mathbf{Ar}_t$  should lie within a probable range of counts extrapolated from computer vision counts. The lower bound  $\alpha_{LB}^{CV} \mathbf{f}_t^{CV}$  assumes that computer vision is overcounting, and the upper bound  $\alpha_{UB}^{CV} \mathbf{f}_t^{CV}$  assumes that it is undercounting. The sum-of-squares of the error  $\|\mathbf{s}_t^{CV}\|_2^2$  is minimised in the objective function (5.1a). The two following constraints, (5.1d) and (5.1e), are analogous constraints for loop detector counts. Again, the error  $\|\mathbf{s}_t^{LD}\|_2^2$  is minimised in the objective. Because all of these bounds may be relatively loose, the problem is underconstrained. Within the possible space of solutions, the final two terms in the objective function optimise for two heuristics of simulation realism: first, the  $\lambda_{\text{nonfringe}}$  term minimises the number of nonfringe routes, which are rare under my assumptions; and second, the  $\lambda_{\text{temporal}}$  term encodes the intuition that traffic flow should be relatively continuous over time.

The solution  $\mathbf{r}$  represents a set of vehicles that should enter the road network within the 15-minute time segment indexed as  $t$ . How can these vehicles be distributed within the time segment? The simplest strategy is to uniformly distribute them. However, under this strategy, two adjacent traffic segments with similar volumes may have very different traffic patterns. Instead, I solve the following QP that enforces similarity between the traffic patterns of adjacent time segments:

$$\begin{aligned} \min_{\mathbf{c}_t} & \|\mathbf{c}_t - \mathbf{c}_{t-1}\|_2^2 \\ \text{s.t.} & \sum_{m=1}^{15} c_{i,(t,m)} = r_{it}, \forall i \in \{1, \dots, n\}, \\ & \mathbf{c}_t \in (\mathbb{Z}^{\geq 0})^{n \times 15} \end{aligned}$$

where  $\mathbf{c}_t$  is the minute-by-minute distribution of routes  $\mathbf{r}_t$  within time segment  $t$ :  $c_{i,(t,m)}$  is the number of times a vehicle with route  $i$  appears in the simulation during minute  $m$  of the 15-minute time segment  $t$ . As a base condition, I uniformly distribute traffic for the first time segment  $t = 0$ , but I solve this QP for time segments  $t \in \{1, \dots, 95\}$ .

More sophisticated methods could be used to enforce temporal continuity in traffic flow. Instead of a uniform distribution, the knowledge of stakeholders could be used to initialise the traffic flow of  $t = 0$ . Dynamics models could also be used to compute the expected travel time of vehicles through the road network. This would mitigate cases where a vehicle appears in time segment  $t$ , but does not reach an intersection where it is counted in  $\mathbf{f}_t$  until time segment  $t + 1$  or later.

### 5.3.3 LLM Agent Simulation Refinement from Natural Language Feedback

The QIP (5.1) remains fundamentally underconstrained. For each counting location, the generated counts could lie anywhere between the lower bounds (constraints (5.1b) and (5.1d)) and the upper bounds (constraints (5.1c) and (5.1e)). Additionally, most municipalities do not install camera detectors at every intersection, meaning that my vehicle tracking method (Section 5.3.1) does not

generate counts or impose bounds for the entire road network. Not all possible ways of assigning routes to match these counts are equally realistic. The domain knowledge of stakeholders, such as traffic engineers, can be leveraged to ensure that the traffic simulation is aligned with downstream use cases such as traffic analysis. Yet, without experience in optimisation, it is difficult for these stakeholders to directly modify QIP (5.1) to align with their intuition.

My problem formulation is as follows. There are  $K$  pieces of structured natural language feedback  $p^{(k)} = (t^{(k)}, j^{(k)}, l^{(k)})$ , where each piece consists of a time, intersection, and a natural language description of what corrections (if any) should be made to the simulated traffic state at this intersection. There is also code which solves the original problem  $\text{QP}^{(0)}$ , and the route counts  $\mathbf{Ar}^{(0)}$  obtained from solving  $\text{QP}^{(0)}$ . The objective is to produce an updated problem  $\text{QP}^{(k)}$ , which has been modified so that it will produce a new route set  $\mathbf{r}^{(k)}$  that addresses  $\{p^{(1)}, \dots, p^{(K)}\}$ . The core difficulty in this problem is converting the natural language feedback into concrete optimisation constraints, which cannot be accomplished by traditional optimisation methods. Instead, I solve this problem by using an LLM agent to answer prompts containing  $(p^{(k)}, \text{QP}^{(k-1)}, \mathbf{Ar}^{(k-1)})$ , and leveraging its code generation capabilities to generate  $\text{QP}^{(k)}$ .

Notably, I do not provide the LLM agent with any handcrafted information beyond the time segment  $t^{(k)}$  and intersection  $j^{(k)}$  that the feedback is targeted at; what is already available from  $\text{QP}^{(k)}$ ; and a list of intersections and main roads. Based on the set of route counts  $\mathbf{Ar}^{(k-1)}$  from the previous simulation, the LLM agent must automatically extract concrete, quantitative constraints that are aligned with the qualitative feedback. To solve this task, I prompt the LLM agent using a chain of thought [438] to:

- (1) Extract the relevant counts by formulating a call to a `get_counts` tool, which retrieves the previous counts  $\mathbf{Ar}_{jt}^{(k-1)}$  for a particular location and time segment  $(j, t)$ ;
- (2) Write a constraint corresponding to the feedback  $p^{(k)}$ , using the counts from the previous step to make subjective judgments on how to set undetermined coefficients;
- (3) Translate this constraint to Python code for the package `cvxpy` 1.5.3 [102];
- (4) For the time segment  $t^{(k)}$  specified in the feedback, add this constraint to the optimisation function while minimally modifying the rest of the code;
- (5) For adjacent time segments  $(t^{(k)} - 1, t^{(k)} + 1)$ , add relaxed constraints to ensure temporal continuity.

I use this LLM agent within an iterative simulation refinement framework like that of Behari et al. [34]. For each of  $K$  pieces of feedback, I first use the LLM agent to generate a program. Then, I apply a rapid verification procedure to the generated program based on three criteria:

- (1) *Syntactic correctness*. I attempt to execute the program in a Python interpreter to ensure it represents syntactically correct Python. If not, then it cannot generate an updated simulation.
- (2) *Feasibility*. I attempt to solve the new QIP for the time segment  $t^{(k)}$ , as well as for adjacent time segments  $(t^{(k)} - 1, t^{(k)} + 1)$ . Assuming that the feedback is internally consistent, and given the underconstrained nature of the problem, I expect that the solver should be able to quickly find (see runtime results in Section 5.4.2) at least one feasible solution  $\hat{\mathbf{r}}^{(k)}$  for the LLM-generated QIP.

- (3) *Semantic correctness.* I attempt to verify that the LLM agent’s modification to the simulation actually corresponds to the feedback given, based on the solution  $\hat{\mathbf{r}}^{(k)}$  to the feasibility check. To do so, I use the LLM agent to perform *self-reflection* [369]. It uses the `get_counts` tool to first retrieve relevant counts  $\mathbf{A}\mathbf{r}^{(k-1)}$  from the previous solution, and then the same counts  $\mathbf{A}\hat{\mathbf{r}}^{(k)}$  from the candidate solution. Then, the LLM agent compares these counts while taking into account the feedback  $p^{(k)}$  to return a binary signal of whether the modification is semantically correct.

If the program fails any of these three verification criteria, I discard the program and prompt the LLM agent to generate a new one. This process repeats until the LLM agent generates a correct program  $\mathbf{QP}^{(k)}$  for feedback  $p^{(k)}$ . In the next iteration, I prompt the LLM agent to directly modify  $\mathbf{QP}^{(k)}$  to produce  $\mathbf{QP}^{(k+1)}$ . After at least  $K$  generations, I obtain a single program  $\mathbf{QP}^{(K)}$ , which I execute for all timesteps  $t$  to obtain the final solution  $\mathbf{r}^{(K)}$  and a corresponding simulation.

## 5.4 Simulation Results: Strongsville, Ohio

I applied my demand modelling pipeline from Section 5.3 to simulate a large road network from the city of Strongsville, Ohio. The Strongsville road network experiences heavy through traffic due to its connection to two interstates, I-71 and I-80; the ramps of these interstates respectively connect to two intersecting arterials, SR 82 (Royalton Road) and US 42 (Pearl Road). The daily traffic volumes of both of these roads have exceeded their designed capacities, leading to the implementation of various countermeasures to improve throughput [118]. As part of these countermeasures, Strongsville installed an adaptive traffic signal control system on SR 82 and US 42.

This system uses three types of traffic detectors. (1) Camera detectors are used for the main roads at each intersection (i.e. along SR 82 and US 42) and on some side roads. (2) Loop detectors are used for most side roads and some turning movements. (3) Radar detectors are used for detection upstream and downstream of intersections. As my goal is to match the traffic state at the intersections themselves, I do not consider data from Strongsville’s radar detectors.

My simulation, as shown in Figure 5.3, covers the intersections along SR 82 and US 42 for which the city has installed adaptive signal control. I first converted OpenStreetMap data to a SUMO [10] road network. Next, on Friday, September 6, 2024, I captured 24 hours of footage from 74 out of 86 counting locations where AutoScope detectors are installed. I used counts from AutoScope and loop detectors to fill missing counts from vehicle tracking. After I applied my pipeline, I randomly assigned vehicles to different vehicle classes [441], following a survey conducted by the Ohio DOT in September 2022. Finally, I implemented the traffic signal patterns that were in use (Section 6.5.1).

In the rest of this section, I evaluate the accuracy of my pipeline steps for this simulation to address my research questions from Section 5.1.

### 5.4.1 RQ1: Accuracy of Vehicle Counting

To evaluate the accuracy of my vehicle tracking-based counting method (Section 5.3.1), I manually counted traffic from camera detector footage. As doing so would be infeasible for the entire



Figure 5.3: Screenshot of traffic simulation of Strongsville, Ohio. Counts are available for 36 intersections, either from vehicle tracking on camera detector footage (approaches in blue) or from loop/AutoScope detectors (approaches in orange). Each intersection is labelled with its number.

simulation, I selected footage from four different intersections that are important to stakeholders (from the south, centre, east, and north of the road network), and two one-hour time segments (12 pm, an off-peak hour, and 5 pm, a peak hour) for each intersection. Figure 5.2 shows a screenshot from one of these pieces of footage.

In Figure 5.4, I compare the counts generated by my vehicle tracking method and by AutoScope to the ground truth from manual counting. Based on my evaluation, my method was able to faithfully capture the traffic state of Strongsville. For all 16 of the footage excerpts that I and another co-author manually counted, my method had an error of less than 2 vehicles per minute (120 vehicles per hour). I used the same set of hyperparameters for preprocessing and detection/tracking across all counting locations; tuning these hyperparameters for individual counting locations could yield further gains.

Meanwhile, AutoScope exhibited a persistent pattern of undercounting across all of the approaches that I manually counted. In fact, it had an error of *more* than 3 vehicles per minute (180 vehicles per hour) for 14 out of 16 footage excerpts, with the exception being intersection 17's northbound approach (where I observed that lane switching resulted in duplicated actuations). The primary cause of this undercounting was the continuous actuation of detection zones by consecutive vehicles, as I discussed in Section 2.2.2. Consistent with this, AutoScope was generally more accurate under intermittent traffic during the 12 pm time segment, and its accuracy degraded under increased traffic levels during the 5 pm time segment.

Table 5.1: Comparison of vehicle counts obtained from my method (CV) and AutoScope (AS) with manual counting.

	8 — Pearl Rd & Echo Rd				17 — Pearl Rd & Business Entrance			
	Northbound		Southbound		Northbound		Southbound	
	12:00	17:00	12:00	17:00	12:00	17:00	12:00	17:00
Manual	823	959	781	1041	1042	952	1106	1329
CV	872	905	839	1071	1065	950	1114	1322
AS	582	637	594	688	1036	1026	647	895

	33 — Royalton Rd & I-71 Ramps				52 — Pearl Rd & Sprague Rd			
	Eastbound		Westbound		Northbound		Southbound	
	12:00	17:00	12:00	17:00	12:00	17:00	12:00	17:00
Manual	1889	2088	1059	726	979	1048	1166	1441
CV	1779	1962	1089	801	1096	1178	1191	1553
AS	1405	1406	587	220	442	578	853	960

An ablation experiment indicated that the preprocessing filters I applied to the footage had a significant impact on the accuracy of my vehicle counting method. The original footage included green detector actuation overlays that disrupted YOLO’s ability to detect vehicles (Figure 5.2a). Because these overlays overlapped with the position of the stop bar, this resulted in a failure to count any vehicles for some counting locations. Replacing the green pixels with black successfully remedied this issue. Frame-blending motion interpolation also mitigated the issue of BoT-SORT failing to associate vehicle tracks before and after dropped frames.

I used these results to estimate bounds for the ratios between my vehicle tracking counts and the ground truth:  $\alpha_{LB}^{CV} = 0.94$ ,  $\alpha_{UB}^{CV} = 1.12$ . As there is insufficient overlap between manual and loop detector counts, I extrapolated bounds for loop detector counts from the ratio between vehicle tracking and loop detector counts. These bounds are loose due to error in the loop detectors:

$$\alpha_{LB}^{LD} = \alpha_{LB}^{CV} \min_t \min_{j \in CV \cap LD} \frac{f_{jt}^{CV}}{f_{jt}^{LD}} = 0.02, \alpha_{UB}^{LD} = \alpha_{UB}^{CV} \max_t \max_{j \in CV \cap LD} \frac{f_{jt}^{CV}}{f_{jt}^{LD}} = 19.06.$$

## 5.4.2 RQ2: Accuracy of Generated Simulation

Next, I solved QIP (5.1) to generate a set of routes consistent with these counts (Section 5.3.2). Owing to the underconstrained nature of the problem, obtaining a feasible solution for the 18 496-variable QIP in each time segment required less than 0.5 seconds; thus, my approach scales well to moderately-sized road networks. However, to optimise solution quality in the final simulation, I ran the QIP for each time segment for 60 seconds. The QP for ensuring temporal continuity in route distributions was slower, requiring approximately 3 minutes per time segment. The objective



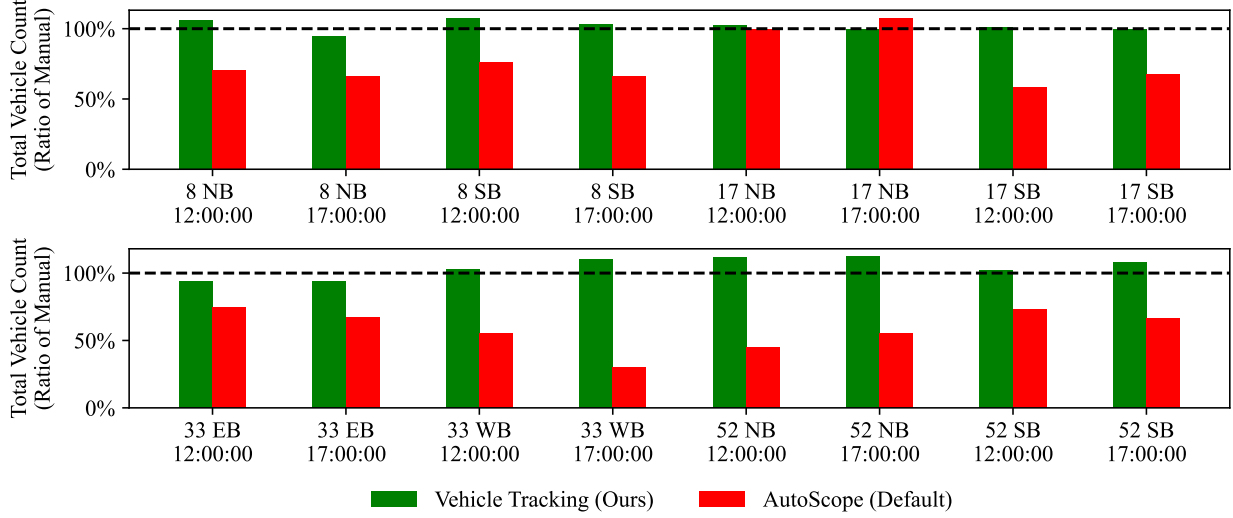


Figure 5.4: Plot of counts from my vehicle tracking method (green) and AutoScope (red), as ratios relative to manual counts. If a method perfectly matches the manual counts, it has a ratio of 100%.

function’s value was mainly determined by bound violations for vehicle tracking counts  $\|\mathbf{s}_t^{CV}\|_2^2$ , which were two orders of magnitude larger than violations for loop detector counts  $\|\mathbf{s}_t^{LD}\|_2^2$  and the fringe route penalty  $\sum_{i \in \text{nonfringe}} r_{it}$ . To balance the objective function, I set  $\lambda_{\text{nonfringe}} = \lambda_{\text{temporal}} = 10$ . I obtained a simulation with a total volume of 182 230 vehicles over 24 hours. Among these vehicles, 72.64% had routes that started and ended on the fringes of the road network.

In Figure 5.5, I focus on the accuracy of my simulation for the time interval between sunrise (6:59 am) and sunset (7:51 pm) on September 6, 2024. Across counting locations on average, the simulation was accurate to my computer vision counts  $f_j^{CV}$ , with no overflow or underflow. This can be attributed to the relatively narrow range of  $[\alpha_{LB}^{CV}, \alpha_{UB}^{CV}]$ . However, there are individual counting locations where the simulation has substantial overflow or underflow, especially so for loop detector counts. I attribute these to counting locations with few vehicles where the detected traffic flow is inconsistent. Violating the expected bounds of these counts results in a small penalty compared to the rest of the objective function. The underconstrained nature of the problem points to a need for stakeholder-driven refinement.

### 5.4.3 RQ3: Accuracy of LLM Agent-Generated Constraints

As my LLM for simulation refinement (Section 5.3.3), I used Claude 3.7 Sonnet. An earlier iteration of this model achieved state-of-the-art performance on code generation benchmarks [12, 503]. I performed two rounds of evaluation: one on synthetic feedback (where a ground truth exists for constraint correctness), and one on real stakeholder feedback (where there is no ground truth). My evaluation focused on the three criteria used by the LLM agent to verify generated code (Section 5.3.3): syntactic correctness, feasibility, and semantic correctness.

First, I randomly generated  $K = 20$  pieces of structured feedback in the form of (intersection,



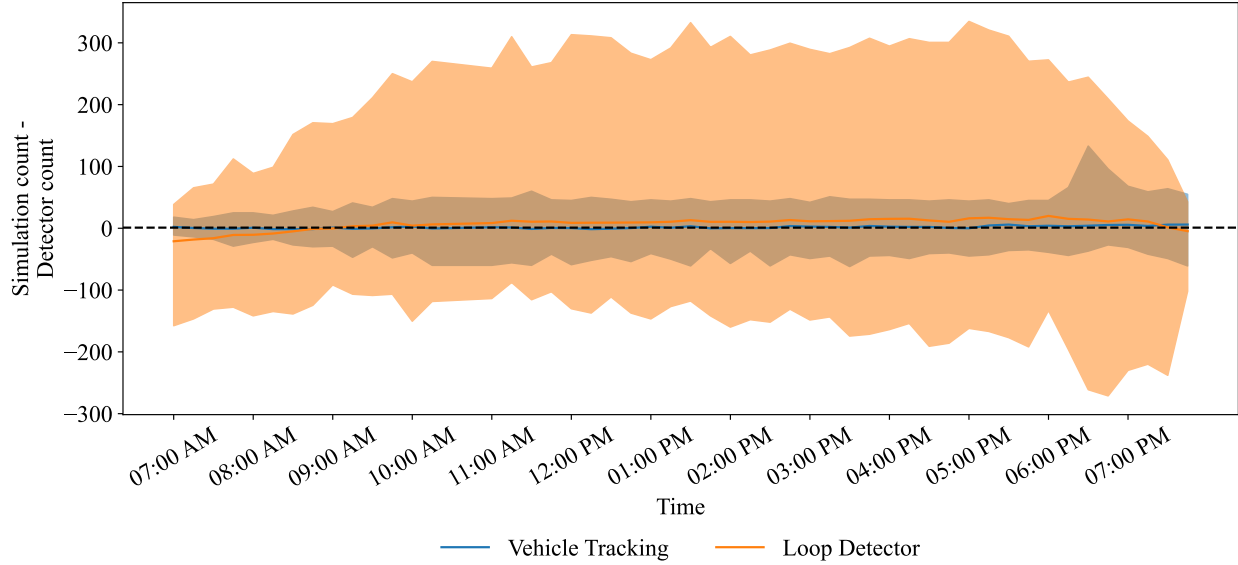


Figure 5.5: Plot of difference between counts in the QIP-generated simulation, counts from my vehicle tracking-based computer vision method, and counts from loop detectors. For each time segment, the solid line is the mean across counting locations, while the shaded region is the range.

direction, approach, increase/decrease) tuples. I used Claude to rephrase this feedback to match the style of stakeholder feedback. Here, I did not use reflection to assess semantic correctness, but instead directly verified the traffic volume in the updated simulation against the structured feedback. I generated ten programs for each piece of feedback with a temperature of 0.8.

As shown in Figure 5.6, the LLM agent always generated valid Python code, giving a syntactic correctness rate of **100%**. The feasibility rate was **87%**. I found that tool use was important to prevent hallucination of counts. Not all generated programs were feasible due to two issues in the added constraints: (1) they included a slack variable term  $s_t$ , which conflicted with the slack variable constraints from the original optimisation problem, or (2) they were formulated in terms of vehicle tracking-based computer vision counts  $f_t^{CV}$ , which were not always available. Lastly, the semantic correctness rate was **87%** — whenever the generated program was feasible, the result was also correct. This gave me confidence in continuing to use my approach, after adding the reflection procedure and modifying the prompt to prevent the two aforementioned issues.

Second, I collected  $K = 20$  pieces of feedback by presenting my simulation for the 5:00 pm time segment to a stakeholder familiar with Strongsville’s traffic conditions. Among the 20 pieces of feedback, 12 pointed to intersections that were true to real life (particularly those with counts based on vehicle tracking-based computer vision), while 8 pointed to intersections where the simulated traffic needed improvement. With reflection in place, I generated a single program for each piece of feedback, again with a temperature of 0.8. The LLM agent achieved a syntactic correctness rate of **100%**, a feasibility rate of **100%**, and a semantic correctness rate of **100%** from reflection, including for the 8 pieces of feedback that indicated changes.

Now, I demonstrate the capabilities of my framework with examples of the constraints gener-

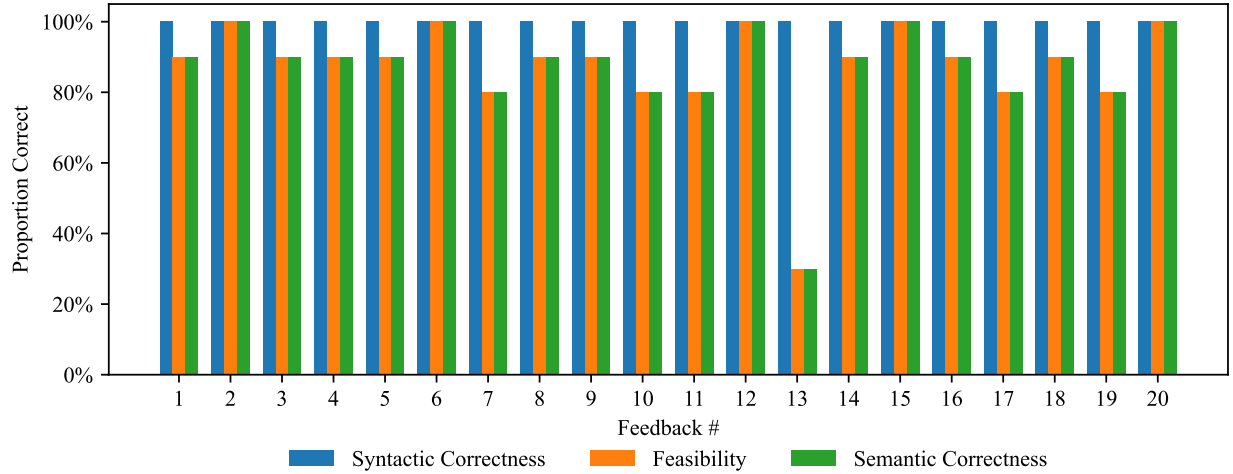


Figure 5.6: Evaluation of  $B = 10$  optimisation programs generated by Claude for  $K = 20$  synthetically generated pieces of feedback. I used three criteria to evaluate the programs: syntactic correctness (perfect), feasibility, and semantic correctness (both high but not perfect).

ated by the LLM agent. Intersection 25 (US 42 & SR 82) was an intersection with camera detectors for which I was unable to capture footage. The stakeholder commented:

I would expect to see more cars both eastbound and westbound at this time of day. Especially at this intersection at this time of day, I would expect every approach to max out. Each left turn, each through, main line, side street, I would expect it to be pretty packed at this time of day. There might be just a couple of lingering [vehicles], but I'd expect pretty much each queue to get its full allocated time.

The LLM agent queried the total eastbound, northbound, southbound, and westbound counts for the old simulation during the 5:00 pm time segment (119, 157, 281, and 334), as well as turn counts for the eastbound and westbound approaches. It then generated the following code:

```
# Add constraints for intersection 25 during peak hour (5:00 PM)
# Indices for intersection 25: 40 (EB), 41 (NB), 42 (SB), 43 (WB)
if hour == "05:00 PM":
    # Enforce higher traffic at intersection 25 based on
    # traffic engineer's feedback
    # Increase eastbound traffic (was 119)
    flow_constraints.append((cat_locs[:, 40, 0].T @ route_count)
                           >= 200)
    # Increase westbound traffic (was 334)
    flow_constraints.append((cat_locs[:, 43, 0].T @ route_count)
                           >= 450)
    # Increase northbound traffic (was 157)
    flow_constraints.append((cat_locs[:, 41, 0].T @ route_count)
                           >= 250)
```

```

# Increase southbound traffic (was 281)
flow_constraints.append((cat_locs[:, 42, 0].T @ route_count)
    >= 400)

# Also increase left and right turn volumes
# Eastbound left turns (was 7)
flow_constraints.append((cat_locs[:, 40, 1].T @ route_count)
    >= 30)
# Eastbound right turns (was 43)
flow_constraints.append((cat_locs[:, 40, 2].T @ route_count)
    >= 80)
# Westbound left turns (was 178)
flow_constraints.append((cat_locs[:, 43, 1].T @ route_count)
    >= 250)
# Westbound right turns (was 138)
flow_constraints.append((cat_locs[:, 43, 2].T @ route_count)
    >= 200)
# Northbound right turns (was 108)
flow_constraints.append((cat_locs[:, 41, 2].T @ route_count)
    >= 150)
# Southbound left turns (was 141)
flow_constraints.append((cat_locs[:, 42, 1].T @ route_count)
    >= 200)
# Southbound right turns (was 17)
flow_constraints.append((cat_locs[:, 42, 2].T @ route_count)
    >= 40)

# Add interpolated constraints for adjacent time segments
# (4:45 PM and 5:15 PM)
elif hour == "04:45 PM" or hour == "05:15 PM":
    # Apply milder increases (about 75% of the 5:00 PM increase)
    # for adjacent time periods
    # Eastbound traffic
    flow_constraints.append((cat_locs[:, 40, 0].T @ route_count)
        >= 180)
    # Westbound traffic
    flow_constraints.append((cat_locs[:, 43, 0].T @ route_count)
        >= 400)
    # Northbound traffic
    flow_constraints.append((cat_locs[:, 41, 0].T @ route_count)
        >= 220)

```

```
# Southbound traffic
flow_constraints.append((cat_locs[:, 42, 0].T @ route_count)
    >= 350)
```

where `cat_locs` is the counting location matrix  $\mathbf{A}$ , and `route_count` is the solution  $\mathbf{r}_t^{(k)}$ .

This snippet demonstrates that the LLM agent is able to correctly identify the indices of the counting locations (40 to 43); quantify the stakeholder's expectation of additional traffic by adding lower bound constraints on the total flow and on the turning flows; and add looser constraints to the adjacent time segments for continuity, all in valid Python code.

For other pieces of feedback where the stakeholder indicated that the simulation was accurate, the LLM agent added constraints to reinforce the existing traffic patterns in the simulation. For intersection 46 (Pearl Rd & Valley View Pkwy), the stakeholder commented:

There's not a whole lot of traffic there on the side streets, what you have going there looks accurate.

This led to the following code, which imposes both lower and upper bounds on traffic volumes:

```
# Add constraints for intersection 46 (Valley View Pkwy)
# based on traffic engineer's feedback
# Intersection 46 corresponds to index 28,
# so counting locations are 112-115
if hour == "05:00 PM":
    # Maintain low side street volumes
    # as engineer confirmed they're accurate
    # Eastbound constraint (location 112)
    flow_constraints.append(cat_locs[:, 112, 0].T @ route_count
        >= 45) # Lower bound
    flow_constraints.append(cat_locs[:, 112, 0].T @ route_count
        <= 57) # Upper bound

    # Westbound constraint (location 115)
    flow_constraints.append(cat_locs[:, 115, 0].T @ route_count
        >= 42) # Lower bound
    flow_constraints.append(cat_locs[:, 115, 0].T @ route_count
        <= 52) # Upper bound

# Add similar but relaxed constraints for adjacent time periods
elif hour == "04:45 PM" or hour == "05:15 PM":
    # Slightly more relaxed bounds for adjacent time periods
    # Eastbound constraint (location 112)
    flow_constraints.append(cat_locs[:, 112, 0].T @ route_count
        >= 40) # Lower bound
    flow_constraints.append(cat_locs[:, 112, 0].T @ route_count
        <= 62) # Upper bound
```

```
# Westbound constraint (location 115)
flow_constraints.append(cat_locs[:, 115, 0].T @ route_count
    >= 37) # Lower bound
flow_constraints.append(cat_locs[:, 115, 0].T @ route_count
    <= 57) # Upper bound
```

My final simulation of Strongsville was created by executing the simulation refinement procedure in sequence for all  $K = 20$  pieces of feedback from the stakeholder. For each modification to the original simulation, the LLM agent’s reflection procedure indicated that it accurately captured the feedback. When the stakeholder was presented with the final simulation, they concurred with the LLM agent regarding the improvements that had been made. The simulation had a total volume of 200 332 vehicles over 24 hours, with 64.72% having fringe routes. This increase can be in part attributed to increased volume at intersections without vehicle tracking counts.

## 5.5 Conclusion

In this chapter, I presented an end-to-end pipeline for modelling demand in traffic simulations with three steps: computer vision-based vehicle counting, combinatorial optimisation-based vehicle route generation, and LLM-based iterative simulation refinement from natural language feedback. I applied my pipeline to a high-traffic road network in Strongsville, Ohio. Based on my evaluation results, my demand modelling methodology adheres more faithfully to real-world traffic conditions than the approaches used in past work, and it holds promise in generalising to road networks from other municipalities with similarly multimodal detector data. For Strongsville, I could generate simulations quickly even when I exhaustively enumerated the route set. However, the number of routes increases exponentially with the number of intersections. To improve the scalability of my pipeline for even larger road networks, I suggest that the route set should not be sampled, but instead clustered into geographic subregions connected at boundaries by major roads.

My pipeline represents an initial, offline proof-of-concept of how LLM agents can enable interactive simulation generation. One line of future work is to convert it into a streaming pipeline capable of near-real-time use. Streaming capabilities would allow simulations to be updated based on live traffic. They would also enable the creation of an interactive interface where stakeholders can iterate on detection and optimisation parameters using natural language feedback, while reviewing the results of their feedback instantaneously. I also envision that other sources of data could be incorporated into my demand modelling framework. Road state reports (e.g. Waze), weather data, and business information can all be indicative of factors that impact traffic. As LLMs’ capabilities improve, they hold promise for integrating data from these heterogeneous sources [62]. This frontier of AI-enabled possibilities can help traffic simulations to better reflect heterogeneity in real-world traffic conditions and to better serve their users.

# Chapter 6

## Self Control

### Imbuing Reinforcement Learning Algorithms for Coordinated Traffic Signal Control with Operational Constraints

*Domain:* Traffic signal control  
*Challenges:* Assurance, coordination

The simple life. Now we can have that.

---

*Philip J. Coulson (Life Model Decoy)*

#### 6.1 Introduction

In Chapter 5, I created a realistic traffic simulation for a road network from the city of Strongsville, Ohio. I also studied how AI techniques can be used to integrate data from different traffic detector types with inherent uncertainty and heterogeneity. This chapter focuses on the intended downstream application of my simulation: as a training and evaluation environment for traffic signal control (TSC) policies based on reinforcement learning (RL). As I outlined in Chapter 1, the goal of training such a policy should be to deploy it in a physical transportation system, so that it can improve the efficiency of traffic signals over existing plans for the benefit of various stakeholders. What characteristics, then, would make an RL-based signal plan trained using my simulation suitable for deployment in Strongsville’s real-world infrastructure?

On the one hand, signal plans should be *coordinated*. Transportation agencies aim to optimise traffic broadly across an entire road network, not just locally at intersections. The key system-level objectives of TSC, such as throughput and green progression [100], involve flows between individual intersections, and therefore can only be achieved by coordinating the behaviour of their controllers. State-of-the-art RL algorithms for TSC such as CoLight [434] perform well in coordination, being able to replicate green progression and other desirable traffic patterns.

On the other hand, signal plans should also provide *assurances* to stakeholders. In practice, stakeholders impose various constraints on how they would like signal plans to operate. Drivers seek predictability and consistency in traffic signals; they experience stress [191] and express complaints [49] about suboptimal signal patterns that differ from what they are used to. Traffic engineers also seek to understand and explain the impact of signal changes using their experience-based

mental models; black-box signal plans (such as those based on RL) hinder these goals [381]. For all of these stakeholders, *safety* is an essential consideration that underlies their operational constraints, and it cannot be compromised in exchange for better performance [470]. Misalignment in mental models often leads to crashes and other safety issues.

The city of Strongsville currently uses Econolite’s Edaptive algorithm for adaptive signal control. As I discussed in Section 2.2.3, Edaptive applies heuristic optimisation methods to adjust base time-of-day patterns according to pre-defined signal plan constraints. Based on discussions with the city’s Traffic Management Centre and other relevant stakeholders, I learnt that it is critical for deployed RL-based signal plans to follow the constraints and format of Edaptive’s signal plans. An example of such a plan is shown in Figure 6.1). Every two cycles, Edaptive uses the last 15 minutes of traffic data to increase or decrease splits. Minimum and maximum green constraints are also enforced but are not shown in this interface. Edaptive may also alter the cycle length and offsets every 3 runs; for simplicity, I consider offsets to be fixed in this chapter.

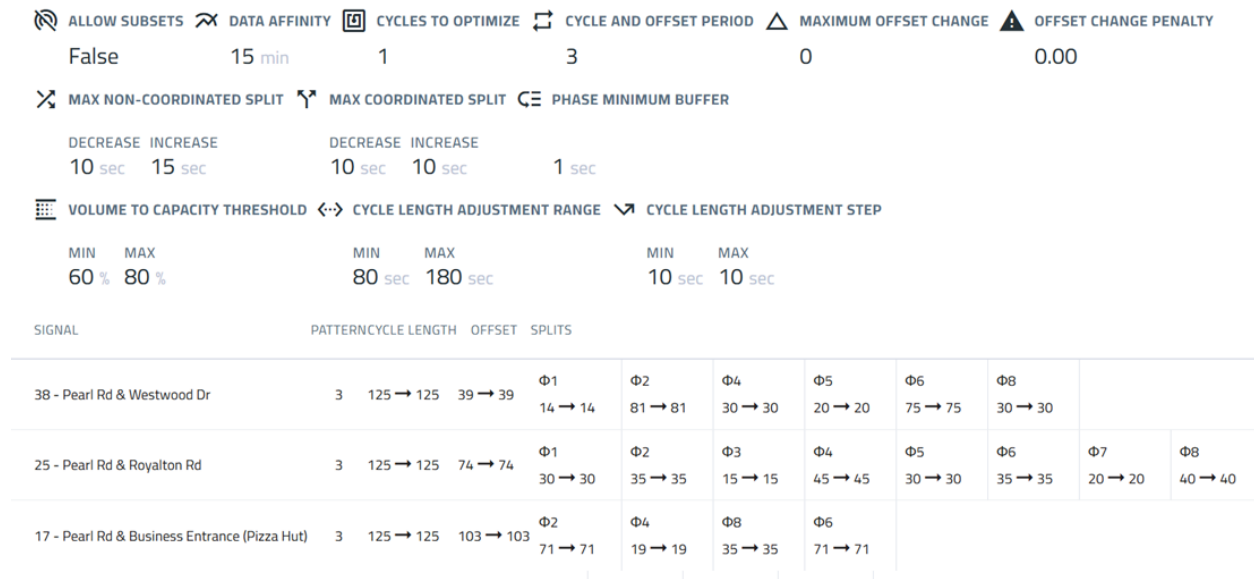


Figure 6.1: Edaptive configuration and output for the central US 42 corridor in Strongsville, Ohio on September 6, 2024 (the date corresponding to the simulation created in Chapter 5).

(Top) Constraints imposed on the Edaptive optimisation procedure.

(Bottom) Sample adjustment to signal plan offsets and splits made by Edaptive at 5:17 pm.

Edaptive’s signal plans obey the following constraints, which I formally define in Section 6.3:

- **Minimum and maximum green:** These constraints on the duration of each phase ensure safety by giving road users from all approaches enough time to traverse an intersection [203]. Maximum green constraints also mitigate the frustration of drivers stuck in traffic [191].
- **Bounded cycle length:** This constraint on the total duration of all phases in a cycle improves the predictability of signal plans for drivers, who use information about the time remaining in phases to determine their behaviour in response to phase changes [328].

- **Fixed phase sequence:** This constraint, which requires a consistent ordering of phases within a cycle, likewise improves the predictability of signal plans for drivers, and also mitigates losses in efficiency from discontinuous changes between different phases [158].
- **Smooth signal plan transition:** This constraint on the allowable change in splits between adjacent cycles improves the verifiability of signal plans, and also mitigates the losses in efficiency caused by transitions in splits between inconsistent plans [49].
- **Split predetermination:** This constraint, which requires the entire sequence of splits for a cycle to be specified in advance, improves the verifiability of signal plans by traffic engineers.

In this chapter, I introduce CycleLight, a suite of techniques that can be used to control the behaviour of signalling policies generated by state-of-the-art RL algorithms for TSC. This adds a dimension of assurance to algorithms that already exhibit strong coordination capabilities. First, I address the first four of these constraints by applying *action masking*, where phase-based acyclic policies are prevented from selecting invalid actions. Next, I address these constraints in conjunction with split predetermination using an *imitation-based loss* with performant acyclic policies. Lastly, to ensure that these cyclic policies follow the constraints, I apply *action projection* to project split-based cyclic policies back to safe sets through differentiable optimisation.

For my experimental evaluation, I seek to address the following research questions:

**RESEARCH QUESTION 6.1.** *Can CycleLight improve safety and adherence to signal plan constraints when it is applied to RL algorithms for TSC?*

**RESEARCH QUESTION 6.2.** *Can CycleLight improve signal performance metrics relative to a static signal plan when it is applied to RL algorithms for TSC?*

**RESEARCH QUESTION 6.3.** *Which signalling constraints among Equation (6.1) to (6.7), as implemented in CycleLight, have the greatest impact on the performance of RL algorithms for TSC?*

I apply these techniques to the MA2C [77] algorithm, which I use to train signal control policies for a three-intersection benchmark and the Strongsville simulation from Chapter 5. I assess how the operational constraints listed above impact both of the most important desiderata in TSC: *performance* and *safety*. To assess safety, I also implement proxy metrics of crash rates as novel reward functions for RL-based TSC. In the benchmark environment, action masking and projection produce policies that improve queue lengths over the fixed-time plan by 10%, while also adhering to all signalling constraints and improving safety-related metrics. While training is ongoing for the Strongsville environment, preliminary results suggest that similar gains can be achieved.

## 6.2 Related Work

### 6.2.1 Safety Constraints in Reinforcement Learning

What does it mean for an RL agent’s behaviour to be “safe”? Various notions of safety have been proposed based on different components of the RL problem formulation. In the most general form, a safety specification  $\varphi$  defines the sets of states, actions, or rewards that are safe or unsafe for



an RL agent [8], such that  $\varphi(s, a, r) = 1$  if the agent is safe and 0 if the agent is unsafe (which should be avoided) [206]. Different methods have been developed for each of these aspects of the RL problem formulation, although they are complementary to each other. In this chapter, I enforce safety in actions, but I also briefly discuss safety in states and safety in rewards.

For safety in actions, the safety specification defines a set of safe actions in each state,  $\mathcal{A}_\varphi(s) = \{a \in \mathcal{A} \mid \varphi(s, a) = 1\}$ . There are three broad categories of methods that have been used in past work to guarantee the safety of actions taken by RL policies [206]:

- **Action masking** methods apply most directly to  $Q$ -learning and other RL algorithms for environments with discrete actions. For a deterministic  $Q$ -learning policy, the masked policy selects the safe action with the highest  $Q$ -value,  $\tilde{a} = \operatorname{argmax}_{a \in \mathcal{A}_\varphi(s)} Q(s, a)$ . For a stochastic policy with action probabilities  $\pi(a \mid s)$ , the masked policy assigns probability 0 to unsafe actions and then renormalises the probabilities of safe actions:  $\tilde{\pi}(a \mid s) = \frac{\varphi(s, a)\pi(a \mid s)}{\sum_{a' \in \mathcal{A}_\varphi(s)} \pi(a' \mid s)}$  [38, 170, 419]. These two approaches were respectively used to implement signal plan constraints for TSC by Liu et al. [239] and Müller and Sabatelli [291]. Unlike them, I extend action masking to account for cycle-level (not just phase-level) constraints for many agents.
- **Action replacement** methods replace an unsafe action  $a$  having  $\varphi(s, a) = 0$  with some fallback action  $\tilde{a}$ , which can be obtained by sampling the safety set, as in *shielded RL* [8], or by using an alternative failsafe controller [27, 171, 358]. Although Krasowski et al. [206] found that action replacement with sampling worked well for control tasks with continuous action spaces, the safety set is typically much smaller for TSC. I use the next phase in a sequence as a failsafe action so that constraints can be satisfied.
- **Action projection** methods are similar to action replacement methods, but also minimise a notion of distance between the replacement action  $\tilde{a}$  and the original action  $a$ :

$$\begin{aligned} \tilde{a} = \operatorname{argmin}_{a'} \operatorname{dist}(a, a') \\ \text{s.t. } \varphi(s, a') = 1 \end{aligned}$$

Two common families of approaches are based on model predictive control: *control barrier functions* (CBFs) [74, 440] and *safe set algorithms* (SSAs) [248, 494]. These approaches are rigorous, but also difficult to apply due to challenges in formulating RL as control problems [350]. In this chapter, I take a more direct optimisation approach to project vector-valued TSC actions, and I also combine action projection with imitation learning.

For safety in states, the formulation is similar to that of safety in actions. There is a set of safe states in each of which *some* safe action can be taken,  $\mathcal{S}_\varphi = \{s \in \mathcal{S} \mid \exists a \in \mathcal{A}, \varphi(s, a) = 1\}$ . Due to the similar problem formulation, many action projection methods also apply to state safety [494]. Both CBFs and SSAs can be used to constrain actions in a way that provides theoretical guarantees about the overall state distributions encountered by agents, either asymptotically or during training. However, to provide these guarantees, these methods require perfect knowledge or specific structural properties of the environment’s transition dynamics [494], which are not known in TSC. Therefore, I do not consider such global notions of state safety.

For safety in rewards, a *cost function* is typically used to penalise unsafe behaviour. *Reward shaping* is the most basic framework for using cost functions: the cost function is directly added

as a penalty to the reward function. This is the most common approach to enforce safety in RL-based TSC, where penalties for safety violations partially [140, 232, 465] or completely [117] constitute the reward function. However, a critical issue of reward shaping is that agents are no longer optimising for the original reward function [322]. This may lead to *reward hacking*, where this fundamental misalignment degrades performance [375]. Even with approaches that address reward hacking [35, 169], these hybrid reward functions still fail to provide concrete assurances while also failing to optimise directly for coordinated performance.

A more rigorous formulation that separates reward functions and cost functions is the *constrained MDP*  $(\mathcal{S}, \mathcal{A}, P, R, \gamma, C, d)$ . In this formulation,  $C : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the cost function,  $\mathcal{C} = \sum_{t=0}^T \gamma^t C(s_t, \pi(s_t))$  is the *cost return*, and  $d$  is an upper bound on the cost return. The objective is to maximise the expected reward return while bounding the cost return [146], i.e.

$$\begin{aligned} \pi^* &= \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, \pi(s_t))} \mathcal{R} \\ \text{s.t. } &\mathbb{E}_{s_0 \sim \rho, s_{t+1} \sim P(s_t, \pi(s_t))} \mathcal{C} \leq d. \end{aligned}$$

Constrained MDPs are typically solved using *Lagrangian methods* [154, 329, 460], where the constrained MDP is converted to and solved as an unconstrained MDP. Specifically, the cost function is added to the reward function as a penalty and adaptively weighted by a Lagrange multiplier  $\lambda$ , giving  $\operatorname{argmax}_{\pi} \min_{\lambda} (\mathcal{R} - \lambda \mathcal{C})$  as the equivalent unconstrained objective. Although Lagrangian methods are a more rigorous alternative to reward shaping, the same challenges of appropriate definition and proper optimisation exist for cost functions as they do for reward functions. Huang and n3n [170] also found that the performance of action-based methods scales better than reward-based methods. Therefore, I do not use the constrained MDP formulation.

## 6.2.2 Safety Assessment in Traffic Signal Control

What does it mean for a traffic signal plan to be “safe”? The metric of most direct concern is the rate of crashes or *conflicts* (i.e. road user interactions that would lead to crashes if not averted) [135]. For safety reasons, it is not possible to deploy a plan and evaluate the crash rate in the real world, and crashes are rare events either way. Therefore, there is virtually no physical data on how different signalling policies affect crash rates [340]. While crash prediction models exist, they have mostly been designed for freeways and other simpler, limited-access roadway designs [168]. Crash prediction models dedicated to arterial roadways have focused on indirect features such as speed, volume, and weather [469], and are of limited utility in signal plan design.

As I established in Chapter 4, traffic simulation provides an alternative means of systematically analysing the safety of signal plans. Although simulations can be used to directly compute crash rates, a more efficient workflow is to use them to evaluate *safety performance functions*. These functions are regression models that rely on proxy metrics of crash rates [235], and are fit using large datasets from simulations or traffic monitoring systems. One of these is the US Federal Highway Administration’s Surrogate Safety Assessment Model (SSAM) [135]. The most relevant application of SSAM was by Sabra et al. [340], who used it to model how different modifications to a Synchro-optimised signal plan’s cycle length, offset, and splits impact crash rates.

Among other proxy metrics, SSAM uses the post-encroachment time and the time to collision of conflicts [73] to estimate crash rates. Although such proxy metrics for crash rate analysis are practically significant, they are virtually unknown to researchers in RL for TSC [302]. I know of only two exceptions: Mirbakhsh and Azizi [279] evaluated the performance of RL for a single intersection with and without safety-based reward components, while Karbasi et al. [187] performed a similar evaluation to assess the impact of connected and automated vehicles on safety. In this chapter, I use the rate of constraint violations to evaluate safety for RL-based traffic signal plans. Unlike previous smaller-scale environments in the RL literature, my evaluation is partially based on the high-fidelity, 36-agent simulation of Strongsville from Chapter 5.

### 6.3 Problem Formulation

In this chapter, I follow the problem setup of TSC from Section 2.2.1: a set of intersections in a road network are each controlled by one agent, and the goal of each agent is to take signalling actions for its intersection to optimise various signal performance measures. Thus, I also model TSC as a multi-agent RL problem using the *Markov game* formulation I introduced in Section 2.3.2.

Within this formulation, I assume that each timestep represents some atomic unit of real time, e.g. 1 second. Each intersection has a set of valid phases  $\Phi_i = \{\phi_i^{(1)}, \dots, \phi_i^{(P_i)}\}$ . In timestep  $t$ , each intersection  $i$  is in phase  $\phi_{i,t} \in \Phi_i$ , and its agent  $i$  takes a signalling action  $a_{i,t} = \phi_{i,t+1} \in \Phi_i$ . (Henceforth, for simplicity, I drop the suffix  $i$  and consider each intersection separately.) In a given policy rollout, a phase  $\phi$  may be signalled for  $K_\phi$  blocks of contiguous timesteps, which can be considered analogously to *options* or other temporally extended action abstractions [18]. For each block  $k \in \{1, \dots, K_\phi\}$ , the starting timestep is  $t_{\phi,k}^{\text{start}}$ , the ending timestep is  $t_{\phi,k}^{\text{end}}$ , and the duration of the block is the split for that phase,  $t_{\phi,k} := t_{\phi,k}^{\text{end}} - t_{\phi,k}^{\text{start}}$ .

Then, the signal operation constraints outlined in Section 6.1 can be formulated as follows:

- **Minimum and maximum green.** Given *minimum green times*  $t_\phi^{\min}$  that lower bound the duration for each phase:

$$t_{\phi,k} \geq t_\phi^{\min}, \forall \phi \in \Phi, k \in \{1, \dots, K_\phi\}. \quad (6.1)$$

Similarly, given *maximum green times*  $t_\phi^{\max}$  that upper bound phase durations:

$$t_{\phi,k} \leq t_\phi^{\max}, \forall \phi \in \Phi, k \in \{1, \dots, K_\phi\}. \quad (6.2)$$

In Figure 6.2, each phase  $\phi^{(p)}$  is shown with a green range on the time axis, denoting its allowable duration under these two constraints. Any shorter, and it would violate the minimum green constraint; any longer, and it would violate the maximum green constraint.

- **Bounded cycle length.** Intuitively, the splits of all phases should sum to a constant duration for predictability. However, this is often relaxed in practice; see the cycle length adjustment parameters in Figure 6.1. Under the relaxed constraint, given a total *cycle length*  $C_k$ , and a distance  $\delta_k$ , the splits of all phases must sum to within  $\delta_k$  of  $C_k$ :

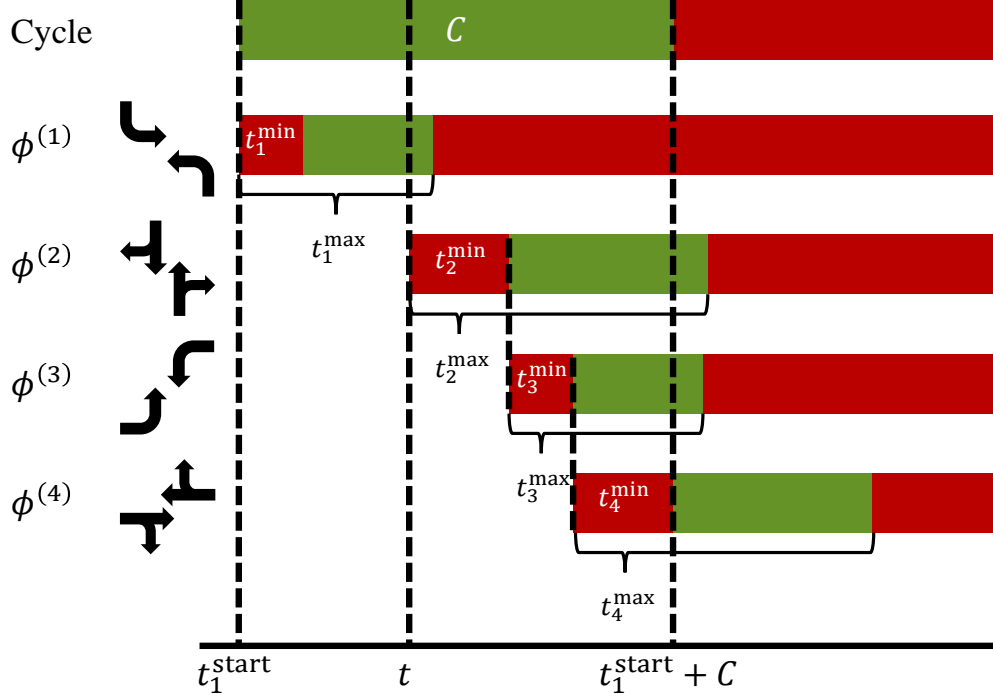


Figure 6.2: Diagram of signalling constraints for a four-phase intersection. For each phase, green denotes permissible durations, while red denotes prohibited durations. The intersection has an overall cycle length  $C$ , which the splits for all phases should sum to. Phases  $\phi^{(1)}$  to  $\phi^{(4)}$  each have a minimum duration,  $t_{\phi}^{\min}$ , and a maximum duration,  $t_{\phi}^{\max}$ .  $\phi^{(1)}$  cannot be signalled up to its maximum duration; otherwise, the cycle is too short for the other phases' minimum durations.

$$\left| \sum_{\phi} t_{\phi,k} - C_k \right| \leq \delta_k, \forall k \in \{1, \dots, K\}, \quad (6.3)$$

where  $K = K_{\phi}, \forall \phi \in \Phi$  now represents the index of the cycle. For example,  $C_k = 120$  s,  $\delta_k = 5$  s gives  $[115, 125]$  as the range of allowable cycle lengths;  $\delta_k = 0$  enforces a constant cycle length.  $C_k$  is not required to be constant across intersections and cycles, although this is often true in practice. Note that this requires the RL policy to signal all phases  $\phi \in \Phi$  exactly once during a cycle, although it does not constrain their order. Thus, there is an implicit constraint on the ordering of start and end timesteps for phases in adjacent cycles:

$$\max_{\phi \in \Phi} t_{\phi,k}^{\text{end}} < \min_{\phi \in \Phi} t_{\phi,k+1}^{\text{start}}, \forall k \in \{1, \dots, K-1\}. \quad (6.4)$$

In Figure 6.2, the overall cycle length  $C$  is shown in green; all four phases,  $\phi^{(1)}, \dots, \phi^{(4)}$ , must be signalled in this time. Beyond this duration, a new cycle (in red) must begin.

- **Fixed phase sequence.** Across all cycles  $k \in \{1, \dots, K\}$ ,  $K = K_\phi, \forall \phi \in \Phi$ , the start and end times of the phases always must occur in the same fixed order,  $\{\phi^{(1)}, \dots, \phi^{(P)}\}$ :

$$\begin{aligned} t_{1,k}^{\text{start}} < t_{2,k}^{\text{start}}, t_{1,k}^{\text{end}} < t_{2,k}^{\text{end}}, \\ \forall \phi_1 < \phi_2 \in \Phi, k \in \{1, \dots, K\}. \end{aligned} \quad (6.5)$$

- **Smooth split transition.** Between two temporally adjacent cycles, the splits of a given phase can change by no more than some threshold  $\epsilon_k$ :

$$\max_{\phi \in \Phi} |t_{\phi,k} - t_{\phi,k+1}| \leq \epsilon_k, \forall k \in \{1, \dots, K-1\}. \quad (6.6)$$

My formulation allows the threshold to vary over time to encapsulate two different types of constraints: (1) *maintaining* the same splits for multiple cycles (in which case  $\epsilon_k = 0$ ), corresponding to the time between Edaptive’s optimisation runs, and (2) *limiting* the increase or decrease of splits during an optimisation run. To simplify the formulation from that of Edaptive (Figure 6.1), I assume the same bound for increases and decreases in splits.

- **Split predetermination.** Let the current cycle index of intersection  $i$  (i.e. the index of the last cycle that has begun in this intersection) be  $k_t$ . Instead of each action  $\mathbf{a}_t$  only deciding the next phase  $\phi_{t+1}$ , it should decide the splits for the entire next cycle  $k_t + 1$ :

$$\mathbf{a}_t = \{t_{1,k_t+1}, \dots, t_{P,k_t+1}\}, \forall t \in \{1, \dots, T\}. \quad (6.7)$$

## 6.4 CycleLight

How can these signal operation constraints be imposed on RL-based policies? In this section, I describe CycleLight, a generalisable suite of postprocessing techniques for signalling actions taken by RL policies. CycleLight can be combined with different RL algorithms for TSC to train policies that satisfy these operational constraints.

### 6.4.1 Action Masking

As I noted in Section 2.3.1, state-of-the-art RL algorithms for TSC are acyclic in that their actions consist only of the phase for the next timestep, i.e.  $a_t = \phi_{t+1}$  [302]. For these algorithms with discrete action spaces, CycleLight uses action masking to implement all of the constraints defined in Section 6.3 except for split predetermination.

Based on these constraints, CycleLight uses action masking to limit the set of phases  $\Phi_{t+1}$  that can be signalled next. Let the current phase be  $\phi^{(p)}$ . The constraints that are enforced at timestep  $t$  depends on the duration of the current phase so far,  $t_{p,k}^{\text{curr}}$ , and the duration of the current cycle  $k$  so far,  $\sum_{\phi \in \Phi} t_{\phi,k}^{\text{curr}}$ . The following cases are possible:

- **Phase cannot be changed.** If the minimum green constraint Equation (6.1) is enforced and  $t_{p,k}^{\text{curr}} < t_p^{\text{min}}$ , then only the current phase can be signalled:  $\Phi_{t+1} = \{\phi^{(p)}\}$ . This can also occur if the split transition constraint Equation (6.6) is enforced, if the phase has not been signalled for long enough relative to the previous cycle:  $t_{p,k}^{\text{curr}} < t_{p,k-1} - \epsilon_k$ .

- **Phase can be changed.** The current phase has been signalled for at least its minimum green time. If both the phase ordering and the phase sequence constraints (Equations (6.4) and (6.5)) are enforced, only the current phase or the next phase in the sequence is valid:  $\Phi_{t+1} = \{\phi^{(p)}, \phi^{((p+1) \bmod P)}\}$ . If only the former is enforced, the current phase or any other phase not yet signalled in this cycle is valid:  $\Phi_{t+1} = \{\phi \in \Phi \mid \phi = \phi^{(p)} \vee t_{\phi, k_t}^{\text{curr}} = 0\}$ . Otherwise, any phase is valid:  $\Phi_{t+1} = \Phi$ .
- **Phase must be changed.** The current phase has been signalled for too long. Several constraint violations can lead to this case. For example, if the maximum green constraint Equation (6.2) is enforced, then this case occurs if  $t_{p,k}^{\text{curr}} > t_p^{\text{max}}$ . This can also occur if the split transition constraint Equation (6.6) is enforced, and the phase has not been signalled for too long relative to the previous cycle:  $t_{p,k}^{\text{curr}} > t_{p,k-1} + \epsilon_k$ .

Lastly, Figure 6.2 illustrates a case where the cycle length constraint Equation (6.3) can lead to this case. Let  $\Phi_k^\emptyset \subseteq \Phi$  be the subset of phases that have not been signalled in cycle  $k$ . If Equation (6.5) is enforced, then  $\Phi_k^\emptyset = \{\phi^{((p+1) \bmod P)}, \dots, \phi^{(P)}\}$ ; otherwise,  $\Phi_k^\emptyset = \{\phi \in \Phi \mid t_{\phi, k_t}^{\text{curr}} = 0\}$ . Under Equation (6.3), every phase must be signalled in a cycle, and each phase must be signalled for its minimum green time. This constraint is violated if not enough time in the cycle is left for the minimum greens of the remaining phases:

$$C + \delta - \sum_{\phi \in \Phi - \Phi_k^\emptyset} t_{\phi, k}^{\text{curr}} < \sum_{\phi \in \Phi_k^\emptyset} t_{\phi}^{\text{min}}.$$

In Figure 6.2,  $\phi^{(p)}$  is  $\phi^{(1)}$ . Although  $\phi^{(1)}$  has not been signalled up to its maximum time, there is not enough time in the cycle  $C$  to signal  $\phi^{(2)}$  for  $t_{\phi^{(2)}}^{\text{min}}$ ,  $\phi^{(3)}$  for  $t_{\phi^{(3)}}^{\text{min}}$ , and  $\phi^{(4)}$  for  $t_{\phi^{(4)}}^{\text{min}}$ .

If both the phase ordering and the phase sequence constraints (Equations (6.4) and (6.5)) are enforced, only the next phase of the sequence is valid:  $\Phi_{t+1} = \{\phi^{((p+1) \bmod P)}\}$ . If only the former is enforced, any other phase not signalled yet in this cycle is valid:  $\Phi_{t+1} = \{\phi \in \Phi \mid t_{\phi, k_t}^{\text{curr}} = 0\}$ . Otherwise, any phase is valid:  $\Phi_{t+1} = \Phi$ .

A phase that is valid to signal in the current timestep may still lead to constraint violations in future timesteps. For example, in Figure 6.2, signalling  $\phi^{(1)}$  up to time  $t$  is valid; however, if the policy signals  $\phi^{(2)}$  for more than  $t_{\phi^{(2)}}^{\text{min}}$ , the cycle length constraint Equation (6.3) will still be violated. Therefore, CycleLight solves integer linear programs (ILPs) to find a subset of phases  $\Phi_{t+1}^{\text{ILP}}$  which will definitely lead the acyclic policy to generate splits that satisfy the signalling constraints.

Specifically, assume that the policy has signalled up to phase  $\phi^{(p)}$  at timestep  $t$ , it is currently in the  $k$ th cycle, and its chosen actions so far have induced splits  $\mathbf{t}_k$  for this cycle. Suppose that the policy signals for  $t_\Delta$  (e.g. 5) seconds at a time. To verify whether signalling  $\phi^{(p)}$  for  $t_\Delta$  additional seconds will lead to feasible splits, CycleLight solves the ILP

$$\min_{\hat{\mathbf{t}}} 1 \quad (6.8)$$

$$\text{s.t. } \sum_{p'=p}^P t_{\Delta} \hat{\mathbf{t}}^{(p')} = C + \delta - \sum_{p'=1}^{p-1} t_{\Delta} \mathbf{t}_k^{(p')} \quad (6.9)$$

$$t_{p'}^{\min} \leq t_{\Delta} \hat{\mathbf{t}}^{(p')} \leq t_{p'}^{\max}, \forall p' \in \{p, \dots, P\} \quad (6.10)$$

$$\left| t_{\Delta} \hat{\mathbf{t}}^{(p')} - \mathbf{t}_{k-1}^{(p')} \right| \leq \epsilon_k, \forall p' \in \{p, \dots, P\} \quad (6.11)$$

$$\hat{\mathbf{t}}^{(p)} \geq \mathbf{t}_k^{(p)} + t_{\Delta} \quad (6.12)$$

$$\hat{\mathbf{t}} \in (\mathbb{Z}^+)^{P-p+1}, \quad (6.13)$$

where  $\hat{\mathbf{t}}$  are the splits for phases  $\phi^{(p)}, \dots, \phi^{(P)}$  divided by  $t_{\Delta}$ ;  $\mathbf{t}_k$  are the undivided splits induced by the acyclic policy *so far* in the current cycle; and  $\mathbf{t}_{k-1}$  are the splits induced in the previous cycle. Equation (6.9) satisfies the cycle length constraint (6.3); Equation (6.10) satisfies the phase duration constraints (6.1) and (6.2); Equation (6.11) satisfies the split transition constraint (6.6); Equation (6.12) constrains the split for the current phase  $\Phi^{(p)}$  to be equal to its total duration if it were to be signalled; and Equation (6.13) constrains all phase times to be multiples of  $t_{\Delta}$ . Verifying whether  $\phi^{((p+1) \bmod P)}$  is feasible involves a similar ILP, except the phases are divided into  $\{1, \dots, p\}$  and  $\{p+1, \dots, P\}$ , and Equation (6.12) is replaced by  $\hat{\mathbf{t}}^{((p+1) \bmod P)} = t_{\Delta}$ .

This ILP is similar to the LP solved by my schedule recommendation tool in Chapter 3 to reoptimise schedules while fixing existing shifts. As a simplifying assumption, I assume that the cycle length is a perfect multiple of the timestep increment,  $C \bmod t_{\Delta} = \delta \bmod t_{\Delta} = 0$ . If not, I allow the last phase to have an arbitrary length to satisfy the cycle length requirement:  $\hat{\mathbf{t}}^{(P)} \in \mathbb{R}$ .

For a deterministic policy, CycleLight selects the highest-weighted (e.g. by  $Q$ -value) valid phase from the ILP-filtered set  $\Phi_{t+1}^{\text{ILP}}$  as the action for timestep  $t+1$ . For a stochastic policy, given the action probabilities for each phase  $\pi(a = \phi \mid s)$ , CycleLight masks out invalid phases and renormalises the probabilities of the remaining valid phases, which is equivalent to applying Dropout deterministically to the policy's logits [206]:

$$\tilde{\pi}(a = \phi \mid s) = \frac{\mathbb{1}[\phi \in \Phi_{t+1}] \pi(a = \phi \mid s)}{\sum_{\phi' \in \Phi_{t+1}} \pi(a = \phi' \mid s)}.$$

Then, it samples from  $\tilde{\pi}$  to produce the action for timestep  $t+1$ .

## 6.4.2 Imitation-Based Loss

Cyclic RL algorithms for TSC, which output the splits for an entire cycle at once, satisfy the cycle length and predetermination constraints — Equations (6.3), (6.4) and (6.7) — by construction. Policies trained by these algorithms output logits  $\sigma^{(p)}$  for each phase  $p$ , which are normalised into a distribution over the cycle time. Given a fixed cycle length  $C$ , for cycle  $k$  of intersection  $i$ , such a policy would output a length- $P$  vector of the form

$$\mathbf{a}_k = \left\{ \frac{C \sigma^{(1)}}{\sum_{p=1}^P \sigma^{(p)}}, \dots, \frac{C \sigma^{(P)}}{\sum_{p=1}^P \sigma^{(p)}} \right\}.$$

(For cyclic algorithms, I set the cycle length distance  $\delta = 0$ .) If these splits are always executed in the order  $\{\phi^{(1)}, \dots, \phi^{(P)}\}$ , Equation (6.5) is satisfied. At the start of cycle  $k$ ,  $t_{\text{start}} = t_{1,k}^{\text{start}}$ , the policy executes  $\mathbf{a}_t = \mathbf{a}_k$ . Until the end of the cycle at  $t_{\text{end}} = t_{\text{start}} + C$ , it takes no action.

Due to the reduced flexibility in actions from having to commit to a full sequence of splits before observing traffic conditions, cyclic RL algorithms are likely to be less performant than their acyclic counterparts given the same amount of training time. Therefore, cyclic RL trades off a loss in performance for gains in interpretability and safety. However, I propose to use *imitation learning* to make cyclic RL more performant.

In imitation learning, a *student* or *learner* policy  $\hat{\pi}$  is trained based on demonstrations from a *teacher* or *expert* policy  $\pi^*$  [228]. The goal is to minimise some loss between the student and the teacher,  $J(\hat{\pi}) - J(\pi^*)$ , so that the student achieves a level of performance (e.g. in terms of returns) similar to the teacher [320]. For the TSC setting, the student  $\hat{\pi}$  is a cyclic policy, and the teacher  $\pi^*$  is an acyclic policy trained with all relevant signalling constraints imposed.

To train each intersection agent’s cyclic policy  $\hat{\pi}$  with split actions  $\mathbf{a}_k$ , I backpropagate a reward  $r_t = R(s_t, \mathbf{a}_k)$  for each cycle  $k$  after the action  $\mathbf{a}_k$  has been fully executed. The agent’s loss function  $\mathcal{L}$  consists of three components. In addition to the base RL loss  $\mathcal{L}_{\text{base}}$ , I add a penalty of  $\lambda_{\text{viol}}$  for each constraint from Equation (6.1) to Equation (6.6) that is violated by the cyclic policy’s splits  $\mathbf{a}_k$ , as well as a Euclidean distance-based imitation learning loss  $\mathcal{L}_{\text{imit}}$ :

$$\begin{aligned} \mathcal{L}(s_t, \mathbf{a}_k) &= \mathcal{L}_{\text{base}}(s_t, \mathbf{a}_k) + \lambda_{\text{viol}} \mathcal{L}_{\text{viol}}(s_t, \mathbf{a}_k) + \lambda_{\text{imit}} \mathcal{L}_{\text{imit}}(\hat{\mathbf{t}}_k, \mathbf{a}_k) \\ &= \mathcal{L}_{\text{base}}(s_t, \mathbf{a}_k) + \lambda_{\text{viol}} \left[ \left( \sum_{p=1}^P \mathbf{a}_k^{(p)} - C \right)^2 + \sum_{p=1}^P \left( \min(0, t_p^{\min} - \mathbf{a}_k^{(p)}) \right) \right. \\ &\quad \left. + \sum_{p=1}^P \left( \min(0, \mathbf{a}_k^{(p)} - t_p^{\max}) \right) + \sum_{p=1}^P \mathbb{1} \left[ \left| \mathbf{a}_k^{(p)} - \mathbf{a}_{k-1}^{(p)} \right| > \epsilon_k \right] \left( \mathbf{a}_k^{(p)} - \mathbf{a}_{k-1}^{(p)} \right)^2 \right] \\ &\quad + \lambda_{\text{imit}} \left\| \mathbf{a}_k - \hat{\mathbf{t}}_k \right\|_2^2, \end{aligned} \tag{6.14}$$

where  $\hat{\mathbf{t}}_k$  are the splits that would have been induced by the acyclic policy in the same cycle. This loss function penalises the cyclic policy for constraint violations, as well as based on how much the its splits deviate from the splits that the acyclic policy executed given the exact same observations. As shown in Figure 6.3, the imitation learning loss is at the level of cycles: rather than at individual timesteps, the cyclic policy imitates the acyclic policy’s behaviour over entire cycles.

### 6.4.3 Action Projection

By itself, imitation learning does not guarantee that a cyclic policy will follow all the signalling constraints that I consider. Although the penalty  $R_{\text{viol}}$  incentivises the cyclic policy to follow the signalling constraints, no hard constraints are placed on the actual cyclic policy’s behaviour. To follow the remaining constraints, Equations (6.1), (6.2) and (6.6), CycleLight uses action projection



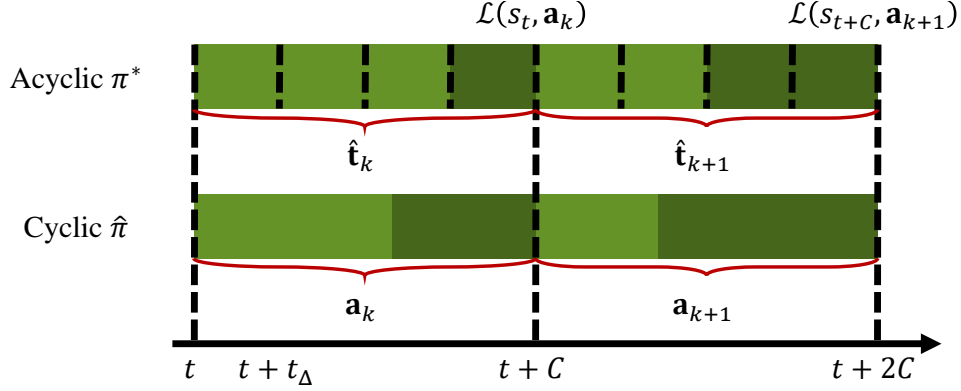


Figure 6.3: Diagram illustrating computation of imitation learning loss. Starting at time  $t$ , a constrained acyclic policy  $\pi^*$  and a cyclic policy  $\hat{\pi}$  both signal phases over cycle length  $C = 4t_\Delta$ .  $\hat{\pi}$  signals splits  $\mathbf{a}_k$ , while the phases signalled by  $\pi^*$  every  $t_\Delta$  timesteps induce splits  $\hat{\mathbf{t}}_k$ . The loss for  $\hat{\pi}$  is taken at the end of the cycle between  $\mathbf{a}_k$  and  $\hat{\mathbf{t}}_k$ , using the state  $s_t$  at the start of the cycle.

by solving the following quadratic program (QP) to postprocess  $\mathbf{a}_k$ :

$$\hat{\mathbf{a}}_k = \underset{\hat{\mathbf{a}}}{\operatorname{argmin}} \|\hat{\mathbf{a}} - \mathbf{a}_k\|_2^2 \quad (6.15)$$

$$\text{s.t. } \sum_{p=1}^P \hat{\mathbf{a}}^{(p)} = C \quad (6.16)$$

$$t_p^{\min} \leq \hat{\mathbf{a}}^{(p)} \leq t_p^{\max}, \forall p \in \{1, \dots, P\} \quad (6.17)$$

$$\left| \hat{\mathbf{a}}^{(p)} - \mathbf{a}_{k-1}^{(p)} \right| \leq \epsilon_k, \forall p \in \{1, \dots, P\} \quad (6.18)$$

$$\hat{\mathbf{a}} \in \mathbb{R}^P,$$

where  $\mathbf{a}_{k-1}$  is the action, i.e. splits, taken by the cyclic policy in the previous cycle. Equation (6.15) minimises the Euclidean distance between the original and projected splits, while Equations (6.16) to (6.18) enforce the cycle length, phase duration, and split transition constraints as before. These optimisation constraints are disabled if the corresponding constraints are not enforced.

One issue of a naïve split postprocessing approach is that policy gradients are lost during the projection step. Nevertheless, there is a straightforward way to obtain gradients from this optimisation problem. Using the differentiable optimisation package `cvxpylayers` [3], CycleLight is able to directly backpropagate through the action projection step. The use of `cvxpylayers` in CycleLight eliminates the need for reward shaping and hyperparameter tuning for the constraint violation penalty from Equation (6.14), and always produces a valid action.

## 6.5 Experiments

In my experimental evaluation of CycleLight, I apply it to train TSC policies for a three-intersection benchmark simulation and the Strongsville simulation from Chapter 5.

### 6.5.1 Simulation Environment

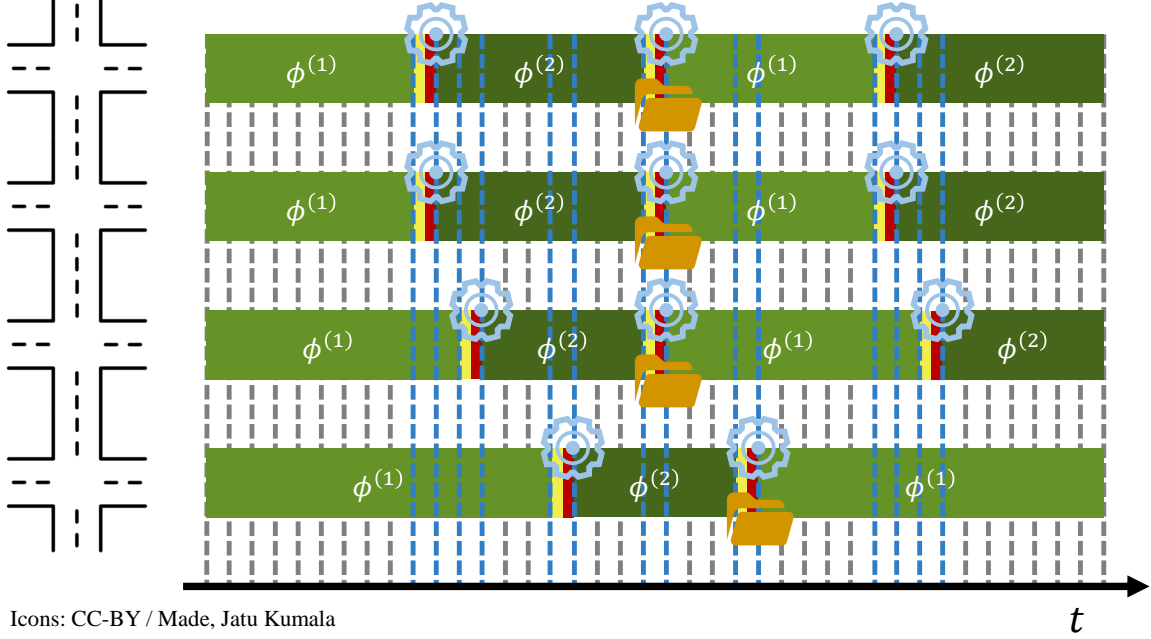


Figure 6.4: Diagram showing execution of a four-intersection traffic simulation in my event queue-based modification of the `sumo-rl` TSC environment, where each intersection has two phases and is controlled by a cyclic policy. The simulation updates the reward in each timestep (shown as dark dashed lines) *without* interfacing with intersection controllers, until it reaches a queued change in signal state for any intersection (shown as blue dashed lines). Blue gears denote these timesteps, at which a cyclic policy changes the intersection to a new phase. Yellow folders denote timesteps at which a cyclic policy adds an experience (based on the entire cycle) to its training dataset.

To create an interface between RL policies and the traffic simulator SUMO [10], I initially used the `sumo-rl` environment of Alegre et al. [6]. However, this environment was designed for acyclic policies and thus has a number of limitations. While `sumo-rl` implements minimum green (6.1) and maximum green (6.2) constraints, it enforces the same limits for all phases at a given intersection, whereas this is usually not true of signal plans in practice. Furthermore, it does not implement the execution of splits-based actions, and therefore provides no way to implement the cycle length (6.3), phase ordering (6.5), or split transition (6.6) constraints.

While retaining the basic framework of `sumo-rl`, I created a new wrapper environment for SUMO to address these limitations. From SUMO simulation files, the environment extracts individual yellow, red, and minimum and green durations for each phase, as well as the full set of signal plans and cycle lengths defined in the simulation. The environment also allows policies to execute actions as individual phases (in which case they are signalled for  $t_{\Delta} \geq 1$  timesteps) or as the full set of splits. It does so using an event queue: instead of each intersection’s controller state being updated instantaneously and independently, actions executed by controllers are enqueued into a centralised queue. During every timestep representing 1 second of simulated time, the envi-

ronment inspects the queue, and triggers the controllers to execute every signal state change event that occurs within that second. While the environment is capable of executing sub-second splits, I use an update granularity of 1 second to trade off simulation runtime and fidelity.

As shown in Figure 6.4, the environment does not enforce a one-to-one correspondence between simulation timesteps and agent actions. During every timestep, the environment updates the current reward of each agent to avoid biasing the average computed over the entire episode. The environment interacts only with lower-level controllers to update the phase of a signal when indicated by the event queue. If the RL policy is acyclic, it queries the policy for new actions and adds an experience to its buffer. Otherwise, if the policy is cyclic, it only queries for a new action and adds a new experience once the entire sequence of enqueued splits has finished executing. In either case, the RL algorithm observes each action’s cumulative effects on the reward during the training process: the observations stored into the buffer originate from the first timesteps of each action, while the rewards stored into the buffer originate from the last timesteps.

## 6.5.2 Simulations

Using this RL environment, I run experiments using two traffic simulations.

- The first simulation is the *Cologne corridor* [413], a three-intersection environment that represents a traffic corridor from the city of Cologne (Köln), Germany. It has a total volume of 4 494 vehicles in 7–8 am rush hour traffic. I used the implementation of this environment from the RESCO benchmark [15]. This implementation includes a default 90-second signal plan for each intersection, which has 3-second yellow change intervals but no red clearance intervals for each phase. It also constrains the minimum and maximum durations of each phase to 5 and 50 seconds for all intersections. I use this environment to iterate on the design of CycleLight, as it runs relatively quickly.
- The second simulation is the *Strongsville road network* from Chapter 5, using the final simulation of the 5–6 pm time segment after stakeholder feedback. For this environment, I replicated the city’s signalling constraints as faithfully as possible. I extracted the phasing sequence for each intersection, together with their yellow, red, minimum green, and maximum green times, from the signal plans for US 42 (Pearl Road) and SR 82 (Royalton Road). However, to reduce the size and complexity of the action space, I removed overlap phases (e.g. in Figure 2.2, I combined the northbound and southbound left turn phases into one action, and straight through/right turn phases into another action), and I set the minimum and maximum green durations to the maxima over all combined phases.

To set cycle lengths and split transition constraints, I replicated the base time-of-day patterns and optimisation constraints from the city’s Edaptive configuration (Figure 6.1). In particular, I set the cycle length of each intersection based on the signal plans they executed during the simulated time segment: 130-second plans for US 42, and 112-second plans for SR 82. For the maximum split distance  $\epsilon$ , I used the coordinated split limit of 10 seconds.

### 6.5.3 Algorithms and Baselines

I use CycleLight to impose signalling constraints on MA2C [77], an observation-sharing policy gradient algorithm that exchanges policy “fingerprints” between adjacent intersections in the road network. The implementation of this algorithm was modified from the RESCO benchmark [15]. MA2C with action masking is denoted Mask-MA2C.

For imitation learning, I needed a cyclic policy. I experimented with modifying the vector of logits output by MA2C, which are normalised into a probability distribution  $\pi(a \mid s)$ . To produce splits, I multiplied this distribution by the target cycle length to obtain  $\mathbf{a} = C\pi(a \mid s)$ . However, this straightforward modification of MA2C did not perform well compared to the acyclic algorithm, and I also encountered difficulties in backpropagating gradients through `cvxpylayers`.

This led me to implement a new cyclic architecture based on MADDPG [243], which I call CyclePG. CyclePG shares an embedding network between the actor and the critic, which consists of three feedforward layers with 64 output units and leaky ReLU activations. The actor network consists of a 64-unit LSTM layer, followed by a leaky ReLU activation, a feedforward layer with 1 output unit, and a softmax activation. The critic network consists of 2 feedforward layers with 64 and 1 output units. A CyclePG policy that both imitates MA2C and implements action projection is denoted Proj-MA2C. As baselines, I compare Mask-MA2C and Proj-MA2C to the default fixed-time signal plan from both environments, a random acyclic policy, and unconstrained MA2C.

All experiments were run in parallel on a server with 56 2.75GHz AMD EPYC 7453 processors, 252 GiB of RAM, and 2 each of 48 GiB NVIDIA RTX A6000 and NVIDIA RTX 6000 Ada Generation GPUs. RL algorithms were trained for 1 000 episodes with 600 timesteps (corresponding to 10 minutes of simulation time) in each episode, using a learning rate of  $1 \times 10^{-3}$ , a batch size of 32, 10 timesteps between each policy update, and 100 timesteps between each target policy update. To solve the ILPs and QPs, I respectively use the packages `cvxpy` 1.5.3 [102] and `cvxpylayers` 0.1.6 [3], using the Gurobi 12.0.2 [153] and ECOS 2.0.12 [105] solvers.

### 6.5.4 Evaluation Metrics

For the reward function and the primary performance evaluation metric of all algorithms, I use the *queue length*, taken as a sum across all intersections and as a mean across all timesteps in an episode. Accordingly, for all RL algorithms, the observations consist of the current phase and the queue lengths for each phase. To compute the queue length, I count the number of vehicles with a speed less than 0.1 m/s in each incoming lane, and then group lanes into phases. I also compute the mean *waiting time* incurred by vehicles at a given intersection, defined as the time elapsed since it last moved; and the mean *speed* of vehicles at the intersection. I do not consider speed as a safety metric; there is not necessarily a one-to-one relationship between lower speeds and safety [87].

In addition, I compute *constraint violation rates* for each constraint enforced. For acyclic algorithms, these are mostly computed every time the policy takes an action. The only exceptions are the cycle length and split transition constraints, which are only well-defined if the policy’s actions induce splits  $\hat{\mathbf{t}}_k$  in each cycle. For cyclic algorithms, these are computed with the splits  $\mathbf{a}_k$  generated by the policy for each cycle. The violation rate for each intersection is computed as the number of violations over 600 timesteps. The overall violation rate is summed across intersections.

Lastly, I use SUMO’s SSM devices [23] to compute two safety metrics that are used in SSAM and other common traffic analysis methods. In SUMO, vehicles with these virtual devices equipped monitor their trajectories for potential car-following, merging, or crossing conflicts between two vehicles, each of which is defined by a “conflict area”. The two safety metrics are:

- *Post-encroachment time*. This is the time between when the first vehicle leaves the conflict area and the second vehicle enters the conflict area.
- *Time to collision*. This is the time for a vehicle to enter a conflict area if it does not change its behaviour, as given by the distance to the conflict area divided by the vehicle’s speed.

I track these metrics for values below 100 s. Both of these metrics should ideally be maximised.

### 6.5.5 Cologne Corridor Results

**Research Question 6.1 RL policies trained with CycleLight eliminate signal plan violations and maintain safety-related metrics.** I begin by assessing the ability of CycleLight to enforce the signal plan constraints I outlined in Section 6.1. As shown in Table 6.1, when trained with action masking (Mask-MA2C) or action projection with an imitation-based loss (Proj-MA2C), MA2C-based policies always follow all of the signal plan constraints by construction. By contrast, the unconstrained MA2C policy frequently violates the minimum green and phase ordering constraints; the cycle length and split transition constraints are almost always violated due to the low probability of an unconstrained policy signalling phases in a fixed sequence. The constraint violation rate of unconstrained MA2C approaches that of the random policy in both environments. In fact, in the Cologne corridor environment, unconstrained MA2C violates the maximum green constraint, while the random policy does not.

For the remainder of this section, I focus solely on the Cologne corridor environment; I return to discussing results for the Strongsville network in Section 6.5.6.

Figures 6.5 and 6.6 compare the TTC and PET of MA2C-based policies trained with CycleLight to the baselines. For both of these metrics, Proj-MA2C has a higher mean value than unconstrained MA2C. In particular, Proj-MA2C has a significantly higher mean PET than all other policies except the high-variance Mask-MA2C. Meanwhile, both unconstrained MA2C and Proj-MA2C have TTCs intermediate between fixed-time and random policies. Although Mask-MA2C has a lower mean TTC than the random policy, I note that TTC is a relatively noisy metric, and none of the policies have significantly different means. Thus, while both Mask-MA2C and Proj-MA2C are able to prevent explicit violations of signalling constraints, the cyclic policy output by Proj-MA2C also leads the simulated vehicles to behave in a more consistently safe manner.

**Research Question 6.2 RL policies trained with CycleLight improve signal performance metrics over strong fixed-time plans.** Next, I investigate the impact of CycleLight on signal performance metrics in the Cologne corridor. In terms of the reward function, i.e. the queue length, Figure 6.7 shows that both Mask-MA2C and Proj-MA2C significantly outperform non-RL baselines. The two CycleLight-based policies both achieve a total queue length of approximately 33, compared to 37.02 for the fixed-time policy ( $\approx 10\%$  improvement,  $p \approx 0.005$ ) and 50.93 for the random policy ( $\approx 35\%$  improvement,  $p \approx 4 \times 10^{-5}$ ). Unconstrained MA2C outperforms

Algorithm	Min Green	Max Green	Phase Ordering
<b>Cologne Corridor</b>			
Fixed-Time	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
Random	187.14 $\pm$ 3.07	0.0 $\pm$ 0.0	158.56 $\pm$ 2.92
MA2C	98.00 $\pm$ 6.92	0.42 $\pm$ 0.25	144.48 $\pm$ 2.63
Mask-MA2C	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
Proj-MA2C	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
<b>Strongsville Network</b>			
Fixed-Time	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
Random	2 982.64 $\pm$ 14.63	67.38 $\pm$ 3.83	1878.40 $\pm$ 5.07
MA2C	1562.40 $\pm$ 29.61*	4.40 $\pm$ 1.33*	1049.00 $\pm$ 28.80*
Mask-MA2C	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0
Proj-MA2C	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0

Table 6.1: Number of violations of signalling constraints by different TSC policies over a 600-second episode (means and 95% confidence intervals across 5 different randomly-seeded runs). The bounded cycle length and smooth split transition constraints are not well-defined for acyclic policies such as random and MA2C, and they are always 0 for the other algorithms by construction. Preliminary results are denoted with an asterisk.

CycleLight-constrained MA2C by a very large margin, with a total queue length of 10.19 ( $\approx 72\%$  improvement over the fixed-time policy). However, it achieves this by breaking all of the signal plan constraints, as shown in Table 6.1. By contrast, CycleLight is able to simultaneously enforce signalling constraints (thus improving assurance) while improving performance relative to the fixed-time policy (through improving coordination). I hypothesise that CycleLight-constrained MA2C has reached the lower bound of queue length that is achievable while respecting all constraints; in later results, I explore the effect of individual constraints on MA2C’s performance.

Figure 6.8 shows that Proj-MA2C leads to significantly shorter waiting times for vehicles than the fixed-time policy ( $\approx 8\%$  improvement from 10.52 s to 9.65 s,  $p \approx 0.03$ ). By contrast, Mask-MA2C leads to longer waiting times ( $\approx 20\%$  degradation from 10.52 s to 12.58 s), although this difference is only nearly significant (two-tailed  $p \approx 0.06$ ). I hypothesise that the more myopic nature of Mask-MA2C leads to signalling decisions that reduce queue length without also reducing waiting time. However, both Mask-MA2C and Proj-MA2C significantly outperform the random baseline. Again, unconstrained MA2C significantly outperforms all other signal control policies in terms of waiting time, but it accomplishes this through a policy that violates signal plan constraints.

Lastly, Figure 6.9 compares the mean speeds of vehicles under CycleLight-constrained MA2C policies and baseline policies. Neither Mask-MA2C or Proj-MA2C significantly differ from the fixed-time policy in speed, with a mean speed of approximately 3 m/s (two-tailed  $p \approx 0.48, 0.11$ ). The random policy has a significantly lower speed of 1.85 m/s. This is sensible, given that the

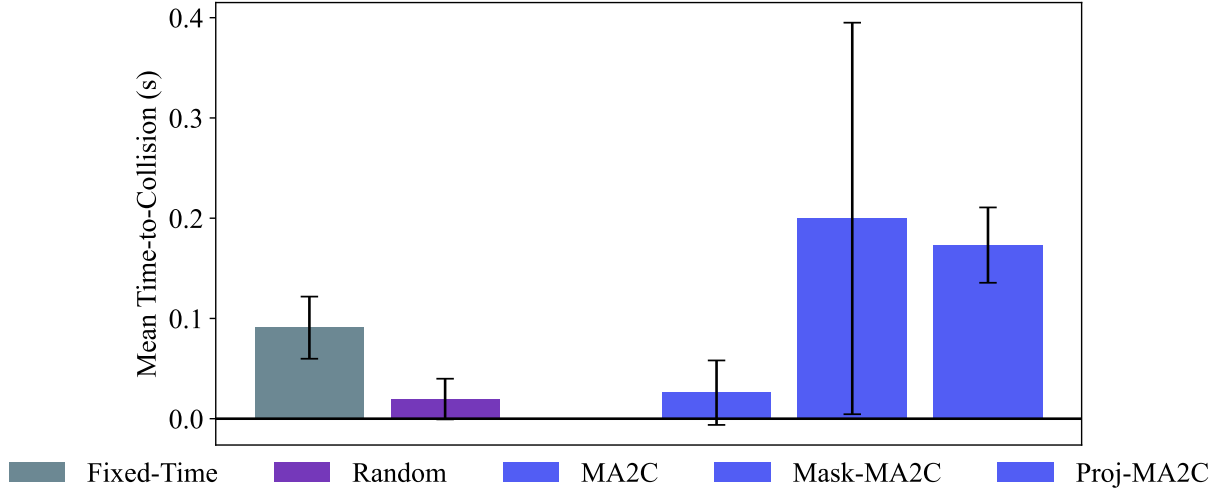


Figure 6.5: Mean post-encroachment time (PET) of vehicles in the Cologne corridor under different signal control policies. Bars show 95% confidence intervals based on 5 randomly-seeded runs.

policy does not attempt to minimise queued vehicles and also leads to irregular transitions between phases. Meanwhile, unconstrained MA2C has the highest mean speed at 4.61 m/s. On the one hand, higher speeds generally improve the throughput of vehicles; on the other hand, higher speeds may be associated with less conservative driving behaviour, as evidenced by the PET metric in Figure 6.5. Thus, the retention of vehicle speeds similar to the status quo by CycleLight-constrained policies may be preferable to stakeholders.

**Research Question 6.3 Relaxing the phase ordering constraint could allow RL policies trained with CycleLight to achieve the most substantial performance improvements.** Which of the signal plan constraints from has the greatest impact on the performance of MA2C? In Figure 6.10, I show the effect of adding each constraint individually to Mask-MA2C, where these constraints must be enforced directly rather than by construction of the policy (as in Proj-MA2C).

Each of these constraints on their own leads to rewards (i.e. queue lengths) that are intermediate between fully constrained MA2C (mean queue length 33.26) and fully unconstrained MA2C (mean queue length 10.19). The phase ordering constraint (Equation (6.5)) has the greatest impact; among MA2C policies that follow a single constraint, this leads to the longest mean queue length of 28.60. The split transition constraint (Equation (6.6)), which is only defined if the phase ordering constraint is, only leads to a marginal increase in the mean queue length to 30.10. This suggests that a different phase sequence, or possibly a phase sequence that changes over time, would help to improve signalling performance in the Cologne corridor.

Meanwhile, the maximum green constraint (Equation (6.2)) leads to the next largest increase in the mean queue length to 16.69. This suggests that the maximum phase durations in this environment may be too short. Finally, the minimum green (Equation (6.1)) and cycle length (Equation (6.3)) constraints minimally impact the mean queue length (respectively, 10.65 and 10.26).

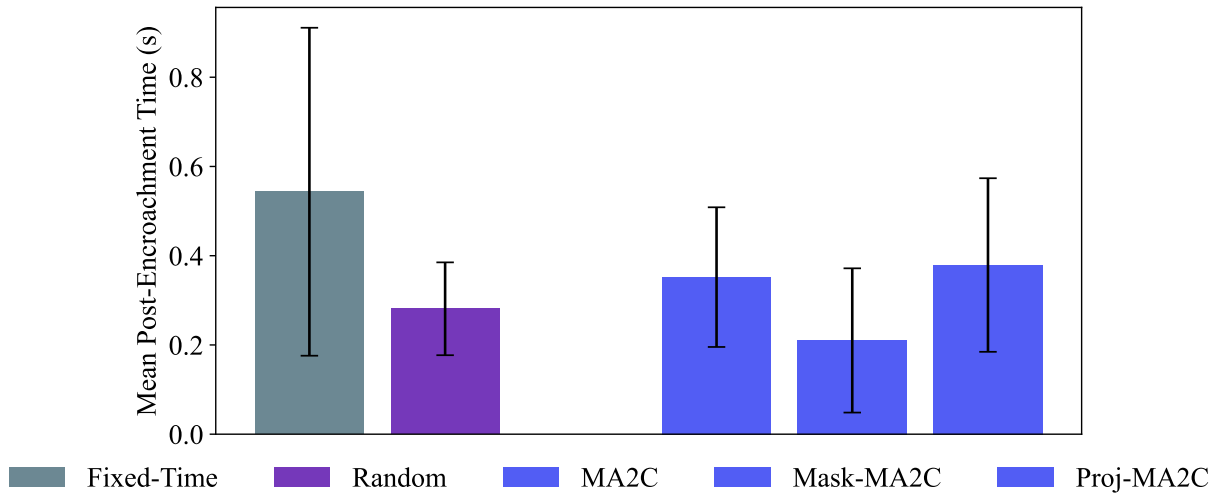


Figure 6.6: Mean time to collision (TTC) of vehicles in the Cologne corridor under different signal control policies. Bars show 95% confidence intervals based on 5 randomly-seeded runs.

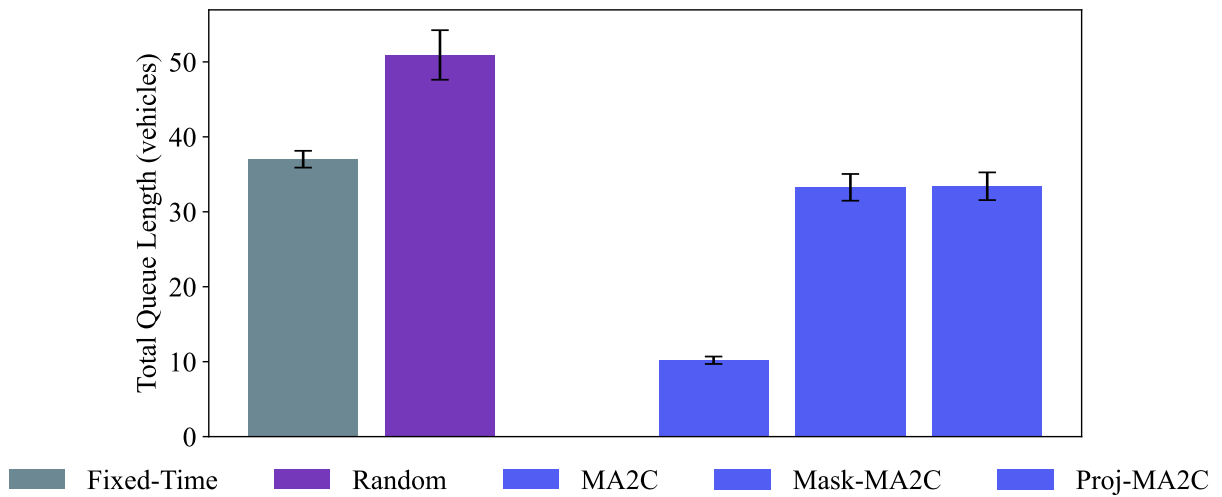


Figure 6.7: Total queue length in the Cologne corridor under different signal control policies. Bars show 95% confidence intervals based on 5 randomly-seeded runs.



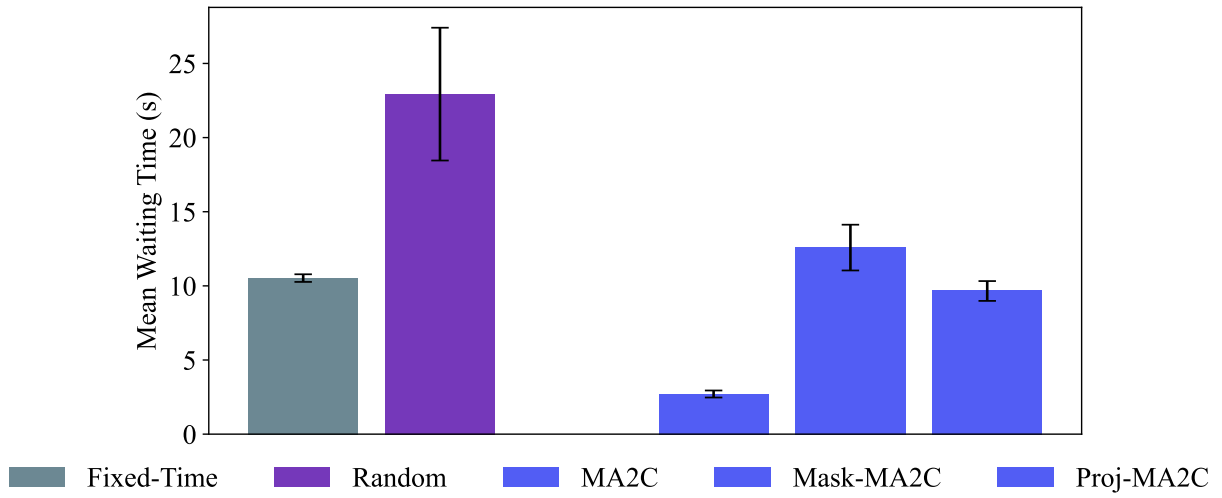


Figure 6.8: Mean waiting time of vehicles in the Cologne corridor under different signal control policies. Bars show 95% confidence intervals based on 5 randomly-seeded runs.

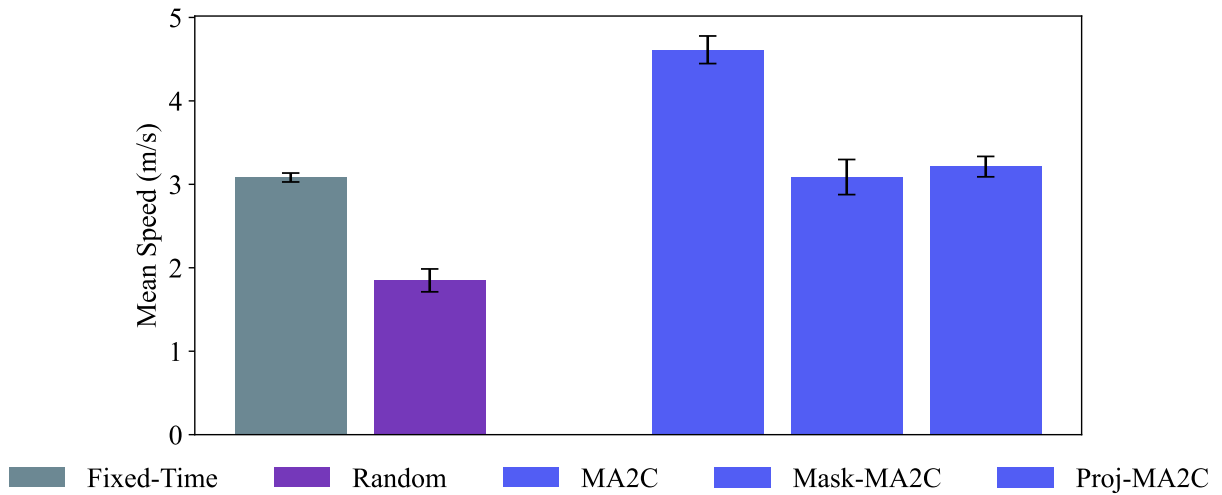


Figure 6.9: Mean speed of vehicles in the Cologne corridor under different signal control policies. Bars show 95% confidence intervals based on 5 randomly-seeded runs.

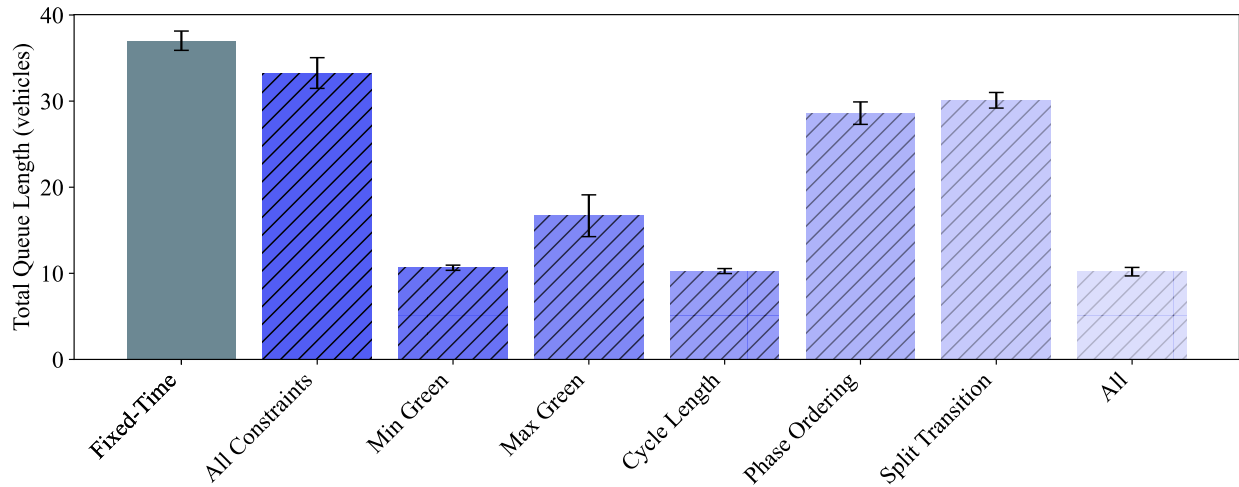


Figure 6.10: Effect of adding individual constraints on total queue length in the Cologne corridor under an MA2C policy trained with action masking (Mask-MA2C). Bars show 95% confidence intervals based on 5 randomly-seeded runs.

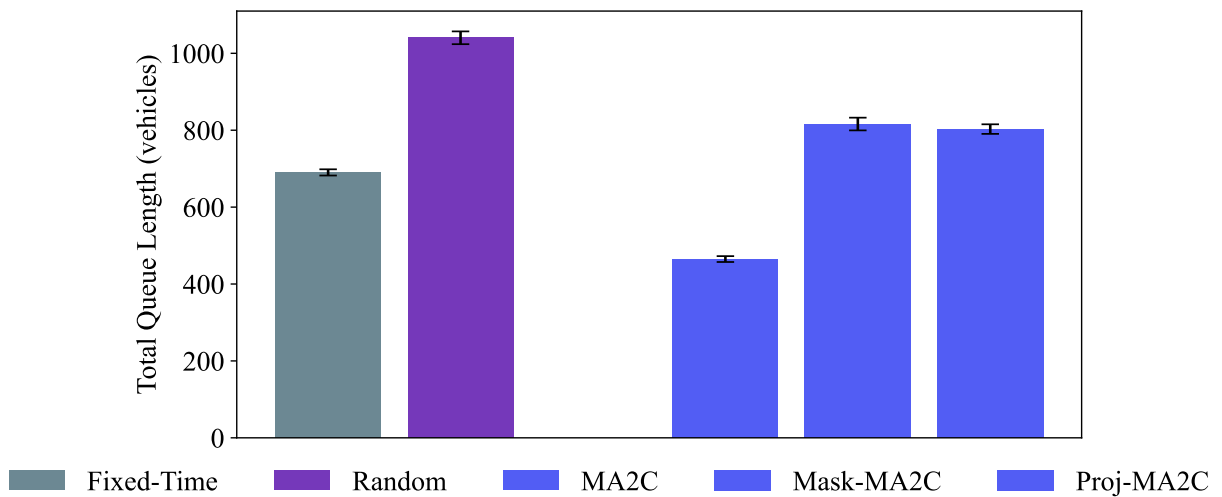


Figure 6.11: Total queue length in the Strongsville network under different signal control policies, based on preliminary training results. Bars show 95% confidence intervals based on 5 randomly-seeded runs.

### 6.5.6 Strongsville Network Results

In Section 6.5.5, I did not discuss results on the Strongsville network for most of the signal performance metrics that I assessed. At the time of writing, training was not yet complete in this environment. The primary cause of this was the long execution time of the Strongsville simulation in my event queue-based environment; a single episode of Proj-MA2C requires approximately 32 minutes to execute for the Strongsville network, compared to 22 seconds for the Cologne corridor. Since the training process did not appear to have fully converged after 600 episodes, and further hyperparameter tuning was infeasible, I am only able to report preliminary results which represent a lower bound on the true performance of CycleLight for the Strongsville network.

Figure 6.11 shows that unconstrained MA2C reduces the mean queue length relative to the fixed-time policy from 690.26 to 464.90 ( $\approx 33\%$  improvement). Although the training of the policy was incomplete, unconstrained MA2C had reached a 70% improvement over the fixed-time policy at this point in the training process for the Cologne corridor. This suggests that performant policies are more difficult to learn in this environment. Thus, it is not unexpected that neither Mask-MA2C nor Proj-MA2C (mean queue lengths 816.10 and 802.94) have been able to achieve the performance of the fixed-time policy at this point. Nevertheless, both have significantly improved over the random policy (mean queue length 1040.40), having closed approximately  $2/3$  of the performance gap between the random and fixed-time policies. When Proj-MA2C is trained without the RL policy loss, only the violation and imitation losses, it exceeds the performance of the fixed-time policy (mean queue length 670.06) after only 200 episodes. I envision that, with proper hyperparameter tuning and acceleration of the simulation environment, CycleLight could provide similar benefits for the Strongsville network as for the Cologne corridor.

## 6.6 Conclusion

In this chapter, I introduced CycleLight, a suite of action postprocessing techniques that enforce operational constraints on RL-based TSC policies. CycleLight combines (1) action masking to restrict the action set of acyclic TSC policies; (2) an imitation-based loss to train cyclic TSC policies from performant acyclic policies; and (3) action projection with differentiable optimisation to directly incorporate signalling constraints into the end-to-end training of cyclic TSC policies. In a three-intersection benchmark simulation, applying CycleLight to the state-of-the-art RL algorithm MA2C leads to policies that outperform fixed-time control policies while adhering perfectly to all constraints. With action projection, CycleLight not only optimises directly for the reward of queue length, but also leads to improvements in related signal performance and safety metrics including waiting time, speed, and post-encroachment time. My incomplete results suggest that the same could be true for the Strongsville network environment.

Could CycleLight deliver comparable performance benefits for all TSC environments, and how should its performance be assessed? Some environments, such as the Strongsville network, are more difficult to optimise than others. They may also have existing signal plans that have already been optimised. I suggest that signal performance measures for TSC should not be assessed in absolute terms, but rather in terms of relative improvement between a completely uninformative

baseline (e.g. a random policy) and a near-optimal baseline (e.g. a policy optimised by Edaptive, which I was unable to access for my experimental evaluation). Consultation with stakeholders is necessary to determine what threshold of improvement can be considered practically significant.

I also suggest that the performance of RL-based signal plans should be evaluated under a variety of traffic scenarios, including generalisation to previously unseen traffic conditions [6]. The lack of such an evaluation is a limitation of my work, but also more generally of the literature on RL for TSC. Nevertheless, given that I simulated 24 hours of traffic in Strongsville in Chapter 5, it should be straightforward to extend my evaluation to cover traffic at different times of the day.

Simulation and RL training were both significant bottlenecks in my experimental evaluation. A single training episode in the Cologne corridor took 7 seconds with unconstrained MA2C, 17 seconds with action masking, and 22 seconds with action projection. How could these runtimes be improved? Surprising speedups can be achieved in many RL environments through straightforward changes to data structures [387]. Likewise, there is likely significant room to accelerate my wrapper environment for SUMO, such as using the `libsumo` library rather than socket communication for API calls [268]. Meanwhile, more efficient but less well-supported alternatives to `cvxpylayers` also exist for differentiable optimisation [54].

Ultimately, the constraints imposed by CycleLight on RL policies in this chapter still represent an overly rigorous abstraction of those followed by Edaptive and other deployed adaptive signal control methods. For example, Edaptive optimises splits more frequently than cycle lengths and offsets, whereas CycleLight does not capture the notion of offsets and optimises the other two signal plan components simultaneously. As demonstrated by my evaluation of the impacts of individual signal plan constraints, there is an inherent tradeoff between performance and adherence to constraints. The choice of which constraints to model and which ones to relax should, again, depend on stakeholder requirements. Regardless, CycleLight introduces a rigorous framework that can flexibly incorporate operational constraints into RL policies for TSC by design.

# Chapter 7

## Making Teams and Influencing Agents

### Coordinating Decision Trees for Efficient Interpretable Multi-Agent Reinforcement Learning

*Domain:* Traffic signal control  
*Challenges:* Assurance, coordination

You're an important part of this team, and  
I don't want you to forget that.

---

*Philip J. Coulson*

#### 7.1 Introduction

In Chapter 6, I addressed one aspect of *assurance* in the deployment of AI technologies for transportation: *safety*. I designed an approach to imbue traffic signal plans based on reinforcement learning (RL) with operational constraints. These constraints align the format of the policy with traditional cycle-offset-split plans, making the resulting signal plans safer and more predictable by stakeholders. In this chapter, I focus on the other aspect of assurance: *interpretability*.

Despite these safety improvements, two challenges still make it difficult for stakeholders to interpret how deep RL policies map states to signalling actions. First, the deep neural network (NN) architectures used for RL policies typically have thousands to millions of parameters. Second, the behaviour of RL policies is difficult to predict and verify over long time horizons. Past literature has found that human stakeholders understand and trust RL policies less than their simpler rule-based counterparts, even if RL yields superior performance [373]. At the same time, as I noted in Chapter 3, interpretability techniques can help stakeholders trust AI systems more [109, 491].

In traffic signal control (TSC) and other applications where the safety and verifiability of RL policies is critical [137, 177], users may deploy interpretable surrogate policies instead of NN policies. Such surrogate policies should be *performant* — capable of achieving high returns. In multi-agent RL, coordinating the training of surrogates is critical for performance: if multiple surrogates are deployed simultaneously, they cannot assume that they are interacting with performant experts, as their performance may be influenced by other agents' suboptimal behaviour.

At the same time, surrogate policies should be *computationally efficient* — it should be possible to generate them with minimal environment interactions and runtime. As I noted in Chapter 4,

running a high-fidelity traffic simulator can be computationally intensive. In human-in-the-loop frameworks where users, such as traffic engineers, provide oversight to correct undesirable policy behaviour [258], the ability to iterate on surrogate policies in a rapid fashion is also critical [449]. The desiderata of performance and computational efficiency exist in tension: more complex models capable of stronger performance and coordination capabilities are also less efficient [277].

Decision trees (DTs) are an attractive model class for interpretable RL due to their comprehensibility [370]. They also enable the design of responsible AI systems, as their branching rules can be easily verified and constrained by human experts or automated processes [42]. DTs lie at the core of the imitation learning framework VIPER [28], which has been applied to distil deep RL policies into DTs in domains such as TSC [177, 502], autonomous vehicles [352], and robotics [338]. However, generalising VIPER to the MARL setting is challenging. Past work [278] introduced two multi-agent VIPER algorithms, IVIPER and MAVIPER, which are both impractical for deployment. IVIPER fails to coordinate the training of DTs, thus sacrificing performance; MAVIPER trains DTs in a coordinated but computationally inefficient manner.

To this end, I introduce HYDRAVIPER, an efficient method to extract coordinated DT policies for cooperative MARL. My method makes three key algorithmic contributions: (1) HYDRAVIPER coordinates agent training by jointly resampling the training dataset for each team of cooperative agents. (2) When interacting with the environment to collect a training dataset, HYDRAVIPER adaptively collects critical trajectories closer to convergence. (3) When interacting with the environment for evaluation, HYDRAVIPER uses a multi-armed bandit-based evaluation strategy to identify promising sets of trained surrogates. Through experimental evaluations, I seek to answer the following questions about HYDRAVIPER:

**RESEARCH QUESTION 7.1.** *Is HYDRAVIPER both performant and efficient (in terms of environment interactions and runtime)?*

**RESEARCH QUESTION 7.2.** *Does HYDRAVIPER maintain performance optimality as the environment interaction budget decreases?*

**RESEARCH QUESTION 7.3.** *Can HYDRAVIPER maintain performance optimality while scalability is improved through agent clustering?*

Experiments demonstrate that HYDRAVIPER is able to balance performance and computational efficiency. HYDRAVIPER also improves the applicability of DT-based interpretable MARL policies: users can exchange training time for performance by altering its environment interaction budgets, but its performance remains optimal at different budget levels. Lastly, HYDRAVIPER’s efficiency on large environments can be improved while maintaining coordination by dividing the agent set into mutually influential teams.

This chapter was published at the AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES) in 2025 [70].

## 7.2 Background

In this chapter, I model the multi-agent cooperation problem as a team-based, mixed competitive-cooperative, *partially observable Markov game*, following the formulation from Section 2.3.3. My algorithm assumes access to value and  $Q$ -functions that take the global observations of all agents, not states, as input:  $V^{\pi_i}(\mathbf{o})$  and  $Q^{\pi_i}(\mathbf{o}, \mathbf{a})$ . I define *mean* value functions and state-action value functions for each team:  $\bar{V}^{\pi_\ell}(\mathbf{o}) := \frac{1}{|\mathcal{T}_\ell|} \sum_{i \in \mathcal{T}_\ell} V^{\pi_i}(\mathbf{o})$ ,  $\bar{Q}^{\pi_\ell}(\mathbf{o}, \mathbf{a}) := \frac{1}{|\mathcal{T}_\ell|} \sum_{i \in \mathcal{T}_\ell} Q^{\pi_i}(\mathbf{o}, \mathbf{a})$ .

### 7.2.1 Decision Trees

A *decision tree* (DT) recursively partitions an input space  $\mathcal{X}$  through functions  $f_j : \mathcal{X} \rightarrow \mathbb{R}$  and thresholds  $\tau_j$  at each internal node  $j$ . Each internal node induces a partition of samples,  $P_j = \{x \in \mathcal{X} : f_j(x) \leq \tau_j\}$ . In a DT policy, internal nodes  $(f_j, \tau_j)$  encode observation-dependent decision criteria, while leaf nodes  $l \in \mathcal{L}$  map partitioned observations to actions:  $\hat{\pi}_i : \mathcal{O}_i \rightarrow \mathcal{A}_i, \forall i \in [N]$ .

### 7.2.2 VIPER

Icons: CC-BY-NC-ND / Milani et al. (2024), Savannah Logan

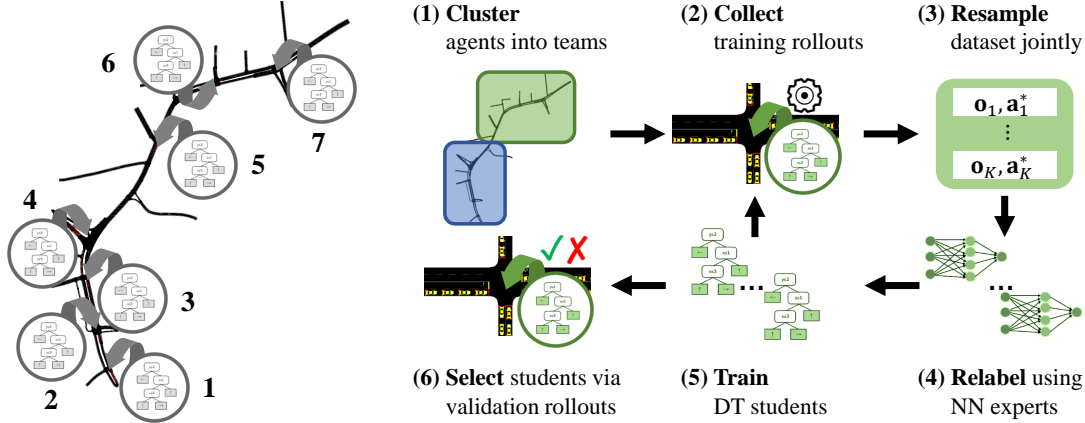


Figure 7.1: (Left) Imitation learning in traffic signal control, where a decision tree must be learnt to imitate the RL-based policy of each intersection’s signal controller agent (the seven-intersection Ingolstadt corridor TSC environment from Section 7.5.1 is shown, with intersections numbered). (Right) The HYDRAVIPER framework, in which DT students are trained *independently* using a *jointly resampled* dataset of environment trajectories and relabelled by an NN expert.

As I introduced in Section 6.4.2, the goal of *imitation learning* is to train a *student* policy  $\hat{\pi}$  based on demonstrations from a *teacher* policy  $\pi^*$  [228], so that the student and the teacher achieve similar performance [320]. *VIPER* [28] is an imitation learning framework, adapted from the more general DAGGER [337], that trains DTs as student policies. Given a trained expert (NN) policy  $\pi^*$ , VIPER iteratively generates student (DT) policies  $\hat{\pi}^m$ . Specifically, in each iteration  $m$ , VIPER:

- (1) *Collects*  $K$  new rollouts  $\{\mathbf{o}, \hat{\pi}^{m-1}(\mathbf{o})\}$  using the previous students from iteration  $m - 1$  (where  $\hat{\pi}^0 := \pi^*$ )

- (2) *Resamples* a dataset  $\mathcal{D}$  from all trajectories collected so far, based on upweighting critical states where taking a suboptimal action may be costly in terms of  $Q$ -values:

$$p_k \propto V^{\pi^*}(\mathbf{o}_k) - \min_{\mathbf{a}} Q^{\pi^*}(\mathbf{o}_k, \mathbf{a})$$

- (3) *Relabels* the dataset with the expert actions  $\pi^*(\mathbf{o}_k)$

- (4) *Trains* new DT students  $\hat{\pi}^m$  on  $\{\mathbf{o}_k, \pi^*(\mathbf{o}_k) \mid \mathbf{o}_k \in \mathcal{D}\}$

After  $M$  iterations, VIPER (5) *selects* a student through validation on an additional set of rollouts. Ross et al. [337] showed that such a procedure is guaranteed to find a student which is performant on the distribution of states that it induces.

## 7.3 Related Work

**Interpretable Multi-Agent Learning** Past methods for interpretable MARL have focused on using feature importances to construct saliency maps [163, 173, 240, 286], building logical structures [189, 178, 426], and defining domain concepts [472]. Each of these categories of methods has limitations. Feature importances and saliency maps are visually clear, but only highlight aspects of the state space without showing how policies use them. Policies based on logical rules and concepts allow users to align the execution of these policies with domain knowledge, but require extensive feature engineering. By contrast, I learn simple policy representations grounded directly in the environment feature space.

**Decision Trees for Reinforcement Learning** Relative to deep NNs, shallow DT policy representations are intrinsically [282] and empirically [370] more comprehensible. One line of work in DT-based RL directly trains DT policies [86, 241, 370, 408] using relaxations amenable to direct optimisation. However, these methods suffer from training instability and performance degradation. Another line of work follows the *VIPER* framework [28], in which a surrogate DT is trained by imitation learning of a performant expert. Although this approach has achieved success in various single-agent settings [177, 338, 352, 502], there exist only two *VIPER*-based algorithms for the multi-agent setting: *IVIPER* and *MAVIPER* [278]. *IVIPER* takes a decentralised approach in which each agent views other agents as a stationary part of the environment, and thus independently trains DTs for each agent. This approach is efficient, but lack of coordination hinders its performance. I show the pseudocode of *IVIPER* in Algorithm 1. Meanwhile, *MAVIPER* takes a centralised approach in which each agent accounts for the potential impact that its joint actions with others may have, and thus jointly trains the DTs by projecting the actions that other agents' DTs would predict. This approach is performant, but the joint DT growth procedure is inefficient. Thus, neither algorithm achieves a balance between performance and computational efficiency.

## 7.4 HYDRAVIPER

In this section, I present *HYDRAVIPER* (Algorithm 2), my algorithm for performant and efficient interpretable MARL. As shown in Figure 7.1, *HYDRAVIPER* builds on the *DAGGER* and *VIPER*



---

**Algorithm 1** IVIPER

---

**Input:**  $(\mathcal{S}, \mathcal{A}, P, r_i, o_i), \pi^*, Q^{\pi^*}, K_{\text{train}}, K_{\text{valid}}, M$   
**Output:**  $\hat{\pi}$

- 1: **Initialise** dataset  $\mathcal{D} \leftarrow \emptyset$  and policies  $\hat{\pi}_i^0 \leftarrow \pi_i^*, \forall i \in N$
- 2: **for** iteration  $m \in \{1, \dots, M\}$  **do**
- 3:   **for** each agent  $i \in \{1, \dots, N\}$  **do**
- 4:     **Collect and relabel**  $K_{\text{train}}$  training rollouts:  $\mathcal{D}_i^m \leftarrow \{(\mathbf{o}, \pi_i^*(o_i)) \sim d^{(\hat{\pi}_i^{m-1}, \pi_{-i}^*)}\}$
- 5:     **Aggregate** dataset  $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \mathcal{D}_i^m$
- 6:     **Set** weights for each  $(\mathbf{o}_k, \mathbf{a}_k) \in \mathcal{D}_i$ :  $p_{ik} \leftarrow V^{\pi_i^*}(\mathbf{o}_k) - \min_{a_i} Q^{\pi_i^*}(\mathbf{o}_k, a_i, \pi_{-i}^*(\mathbf{o}_{-ik}))$
- 7:     **Resample** dataset  $\mathcal{D}'_i \leftarrow \{(\mathbf{o}_k, \mathbf{a}_k) \sim p_{ik}\}$
- 8:     **Train** DT  $\hat{\pi}_i^m \leftarrow \text{TrainDT}(\mathcal{D}'_i)$
- 9:     **Average** return of  $K_{\text{valid}}$  validation rollouts:  $\hat{\mu}_i^m \leftarrow \frac{1}{K_{\text{valid}}} \sum_{k=1}^{K_{\text{valid}}} \bar{r}_{ik}, \bar{r}_{ik} \sim d^{(\hat{\pi}_i^m, \pi_{-i}^*)}$
- 10: **for** each agent  $i \in \{1, \dots, N\}$  **do**
- 11:   **Select** best student  $\hat{\pi}_i \leftarrow \arg\max_m \hat{\mu}_i^m$
- 12: **return** best set of agents  $\hat{\pi} = (\hat{\pi}_1, \dots, \hat{\pi}_N)$

---

frameworks by iteratively collecting data from environment rollouts to train DT policies. HYDRAVIPER first (1) *partitions agents* into clusters for scalability (line 4). Next, in each of  $M$  iterations, HYDRAVIPER: (2) collects a dataset of rollouts from the environment, using an *adaptive procedure* (lines 6–7); (3) resamples the dataset to prioritise learning the correct actions in critical states, using *team-based Q-values* (lines 9–10); (4) and trains DTs based on these datasets (lines 11–12). After it completes all  $M$  training iterations, HYDRAVIPER (5) identifies the best-performing student for each agent, using a *multi-armed bandit algorithm*, and returns them as a policy profile (lines 13–14). Now, I describe each of these algorithm components in detail.

### 7.4.1 Dataset Resampling: Centralised-Q Weighting

VIPER-based algorithms include a dataset resampling step (Algorithm 2, lines 9–10) so that students can focus their learning on more critical states. At a high level, they construct a training dataset by computing sample weights on the aggregated dataset of environment rollouts, typically using some notion of value based on the expert  $Q$ -functions. Measuring value is straightforward in the single-agent setting, but — as I have mentioned — a key obstacle in multi-agent learning is efficient coordination among agents. To address this challenge, HYDRAVIPER induces coordination in the resampling step using a team-based notion of value (Algorithm 3), but trains DTs independently for each agent.

Specifically, HYDRAVIPER resamples the dataset for DT construction based on weights  $p_{\ell k}$ , which represent the relative importance of each sample for each team of agents  $\mathcal{T}_\ell$  (Algorithm 3, line 1). Past work computed this importance based on *individual*  $Q$ -functions, meaning that each agent must maintain its own dataset and induce coordination through (typically computationally expensive) joint training procedures. By contrast, I propose an intuitive change: HYDRAVIPER uses the mean of the expert  $Q$ -functions within each *team* of coordinated agents,

---

**Algorithm 2** HYDRAVIPER
 

---

**Input:** Markov game  $(\mathcal{S}, \mathcal{A}, P, R_i, O_i)$ , experts  $\pi^*$ , expert  $Q$ -functions  $Q^{\pi^*}$ , per-iteration rollout count  $K_{\text{train}}$ , rollout budgets  $(B_{\text{train}}, B_{\text{valid}})$ , threshold  $\epsilon$ , iteration count  $M$ , scaling factor  $c$ , agent distance function  $d$

**Output:** Trained students  $\hat{\pi}$

- 1: **Initialise** dataset  $\mathcal{D} \leftarrow \emptyset$ , policies  $\hat{\pi}_i^0 \leftarrow \pi_i^*, \forall i \in N$
- 2: **Initialise** rollout count  $n_{\text{train}} \leftarrow 0$
- 3: **Initialise** dropped rollout count  $K_{\text{drop}} \leftarrow \infty$ 
  - ▷ **Section 7.4.4: Agent Clustering**
- 4: **Cluster** agents  $\mathcal{T}_1, \dots, \mathcal{T}_L \leftarrow \text{Partition}(\Gamma, \pi^*, d)$
- 5: **for**  $m \in \{1, \dots, M\}$  **do**
  - ▷ **Section 7.4.2: Training Rollouts**
- 6:  $\mathcal{D}, n_{\text{train}} \leftarrow \text{TR-A}(\mathcal{D}, \hat{\pi}^{m-1}, m, K_{\text{train}}, B_{\text{train}}, K_{\text{drop}}, n_{\text{train}})$
- 7: **Reinitialise** dropped rollout count  $K_{\text{drop}} \leftarrow \infty$
- 8: **for each team**  $\mathcal{T}_\ell \in \{\mathcal{T}_1, \dots, \mathcal{T}_L\}$  **do**
  - ▷ **Section 7.4.1: Dataset Resampling**
- 9:  $\mathcal{D}'_\ell, K'_{\text{drop}} \leftarrow \text{C-Q}(\mathcal{D}_\ell, \mathcal{T}_\ell, \pi^*, Q^{\pi^*}, \epsilon)$
- 10:  $K_{\text{drop}} \leftarrow \min(K_{\text{drop}}, K'_{\text{drop}})$
- 11: **for each agent**  $i \in \mathcal{T}_\ell$  **do**
- 12:  $\hat{\pi}_i^m \leftarrow \text{TrainDT}(\mathcal{D}'_\ell)$
- 13: **for each team**  $\mathcal{T}_\ell \in \{1, \dots, L\}$  **do**
  - ▷ **Section 7.4.3: Validation Rollouts**
- 14:  $\hat{\pi}_i, \forall i \in \mathcal{T}_\ell \leftarrow \text{VR-UCB}(\{\hat{\pi}_\ell^m\}_{m=1}^M, \mathcal{T}_\ell, B_{\text{valid}}, c)$
- 15: **return**  $\hat{\pi} = (\hat{\pi}_1, \dots, \hat{\pi}_N)$

---

$\bar{Q}^{\pi_\ell^*} := \frac{1}{|\mathcal{T}_\ell|} \sum_{j \in \mathcal{T}_\ell} Q^{\pi_j^*}$ , to prioritise samples according to their value to the team. Then, I compute the weights as the difference in value between the optimal joint team action and the worst-case joint team action. Intuitively, highly-weighted samples are those where coordinating on joint actions matters for performance. The weights are defined as:

$$\begin{aligned}
 p_{\ell k} &\propto \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \pi^*(\mathbf{o}_k)) - \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})) \\
 &= \bar{V}^{\pi_\ell^*}(\mathbf{o}_k) - \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})).
 \end{aligned} \tag{7.1}$$

For further gains in sample efficiency, HYDRAVIPER does not compute  $p_{\ell k}$  by enumerating joint actions over all agents in the environment. Instead, it only enumerates possible joint actions  $\mathbf{a}_\ell$  over the *team* and uses expert actions  $\pi_{-\ell}^*(\mathbf{o}_{-\ell})$  for the opponent agents. This novel resampling procedure eliminates the need for per-agent datasets in IVIPER and MAVIPER, allowing agents to prioritise the same critical states without computationally expensive joint training.

HYDRAVIPER uses each team’s jointly sampled dataset  $\mathcal{D}_\ell$  to independently train DTs for each agent  $i$  (Algorithm 2, lines 8–9). The DT  $\hat{\pi}_i$  uses individual observations  $o_i$  to fit  $\pi_i^*$ ’s actions in the dataset. Modifying the input dataset rather than the training procedure provides HYDRAVIPER with flexibility in the choice of DT learning algorithm. We use CART [48], but more

---

**Algorithm 3** Centralised-Q Resampling (C-Q)

---

**Input:** Team dataset  $\mathcal{D}_\ell$ , team  $\mathcal{T}_\ell$ , experts  $\pi^*$ , expert  $Q$ -functions  $Q^{\pi^*}$ , threshold  $\epsilon$   
**Output:** Resampled dataset  $\mathcal{D}'_\ell$ , dropped rollout count  $K_{\text{drop}}$

- 1: **Set** weights for each  $(\mathbf{o}_k, \mathbf{a}_k) \in \mathcal{D}_\ell$ :  $p_{\ell k} \leftarrow \bar{V}^{\pi_\ell^*}(\mathbf{o}_k) - \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k}))$
- 2: **Update**  $K_{\text{drop}} \leftarrow \min_\ell \left\lceil \frac{1}{T} |\{(\mathbf{o}_k, \mathbf{a}_k) \in \mathcal{D}_\ell \mid p_{\ell k} \leq \epsilon\}| \right\rceil$
- 3: **Resample** dataset:  $\mathcal{D}'_\ell \leftarrow \{(\mathbf{o}_k, \mathbf{a}_k) \sim p_{\ell k}\}$
- 4: **return**  $\mathcal{D}'_\ell, K_{\text{drop}}$

---

advanced models such as random forests or mixtures of DTs [417] can also be used to improve performance.

### 7.4.2 Training Rollouts: Adaptive Budget Allocation

---

**Algorithm 4** Adaptive Training Rollouts (TR-A)

---

**Input:** Team dataset  $\mathcal{D}_\ell$ , students  $\hat{\pi}^{m-1}$ , iteration  $m$ , per-iteration rollout count  $K_{\text{train}}$ , training rollout budget  $B_{\text{train}}$ , dropped rollout count  $K_{\text{drop}}$ , total rollout count  $n_{\text{train}}$   
**Output:** Updated dataset  $\mathcal{D}$ , total rollout count  $n_{\text{train}}$

- 1: **Set**  $K_{\text{train}}^m \leftarrow \min(K_{\text{drop}}, K_{\text{train}}) \mathbb{1}[n_{\text{train}} \leq B_{\text{train}}]$
- 2: **Update**  $n_{\text{train}} \leftarrow n_{\text{train}} + K_{\text{train}}^m \mathbb{1}[m > 1]$
- 3: **for** each team  $\mathcal{T}_\ell \in \{1, \dots, L\}$  **do**
- 4:     **Collect and relabel**  $K_{\text{train}}^m$  rollouts:  $\mathcal{D}_\ell^m \leftarrow \{(\mathbf{o}_\ell, \pi_\ell^*(\mathbf{o}_\ell)) \sim d^{(\hat{\pi}^{m-1})}\}$
- 5:     **Aggregate** dataset:  $\mathcal{D}_\ell \leftarrow \mathcal{D}_\ell \cup \mathcal{D}_\ell^m$
- 6: **return**  $\mathcal{D}, n_{\text{train}}$

---

Thus far, I have assumed that HYDRAVIPER has access to a dataset of observation-action pairs for training. To collect this dataset, HYDRAVIPER follows the DAGGER-style iterative procedure of collecting a dataset at each iteration  $m$  by rolling out the current student policies  $\hat{\pi}^{m-1}$  (Algorithm 2, line 6–7). The next set of students are trained on the aggregate of all collected datasets, therefore building up the set of inputs likely to be encountered by the student policies during execution. However, collecting training rollouts is computationally expensive. Past work has employed an inefficient static allocation strategy that uniformly performs  $K_{\text{train}}$  rollouts in each iteration. This strategy is problematic because the students are far from convergence early in training, so the distribution of trajectories collected earlier in training potentially diverges from those that converged students would encounter. HYDRAVIPER addresses this challenge through an adaptive rollout strategy that dynamically allocates the training budget at each iteration and prioritises critical states encountered later in training.

Recall that, for each team of cooperative agents  $\mathcal{T}_\ell$ , HYDRAVIPER follows Equation (7.1) to compute weights  $p_{\ell k}$  for resampling the training dataset. I show the following:

**Theorem 7.1.** *Given a dataset of observation-action pairs for team  $\mathcal{T}_\ell$  in iteration  $m$  of HYDRAVIPER,  $\mathcal{D}_\ell = \{(\mathbf{o}_\ell, \mathbf{a}_\ell)\}$ , assume there exists a pair  $(\mathbf{o}_{\ell k}, \mathbf{a}_{\ell k})$  that receives the weight  $p_{\ell k}^{(m)} = 0$ . Then, in iteration  $m + 1$  of HYDRAVIPER, this pair also receives the weight  $p_{\ell k}^{(m+1)} = 0$ .*

*Proof.* If  $p_{\ell k}^{(m)} \propto \bar{V}^{\pi_\ell^*}(\mathbf{o}_k) - \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})) = 0$ , then by definition

$$\bar{V}^{\pi_\ell^*}(\mathbf{o}_k) := \max_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})) = \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})), \quad (7.2)$$

i.e. joint team actions  $\mathbf{a}_\ell$  have no effect on the value given observation  $\mathbf{o}_k$ . When HYDRAVIPER resamples the dataset in iteration  $m$  (Algorithm 3, line 3),  $(\mathbf{o}_{\ell k}, \mathbf{a}_{\ell k})$  will not be part of the resampled dataset  $\mathcal{D}'_\ell$ . However, the resampled dataset  $\mathcal{D}'_\ell$  does not replace the original dataset  $\mathcal{D}_\ell$ .

In iteration  $m + 1$ ,  $\mathcal{D}_\ell$  is aggregated with a newly-collected dataset of observation-action pairs  $\mathcal{D}_\ell^m$  (Algorithm 4, line 5), and  $(\mathbf{o}_{\ell k}, \mathbf{a}_{\ell k})$  continues to be part of this dataset. A new set of weights  $p_\ell^{(m+1)}$  are computed using this expanded dataset (Algorithm 2, line 8). Assume that  $p_{\ell k}^{(m+1)} \neq 0$ . Without loss of generality, let  $p_{\ell k}^{(m+1)} > 0$ . Then by definition

$$\bar{V}^{\pi_\ell^*}(\mathbf{o}_k) := \max_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})) > \min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k})).$$

None of  $\mathbf{o}_k$ ,  $\pi_{-\ell}^*$ ,  $\bar{V}^{\pi_\ell^*}$ , or  $\bar{Q}^{\pi_\ell^*}$  changed between iterations  $m$  and  $m + 1$ , since HYDRAVIPER uses experts and expert value functions to compute the weights. This contradicts Equation (7.2). Thus, if  $p_{\ell k}^{(m)} = 0$ , then  $p_{\ell k}^{(m+1)} = 0$ .  $\square$

A similar proof holds if the strict equality is replaced by a defined threshold  $\epsilon$ . As a result, samples  $(\mathbf{o}_{\ell k}, \mathbf{a}_{\ell k})$  with  $p_{\ell k} \leq \epsilon$  are effectively *removed* from the dataset  $\mathcal{D}$ . This intuition serves as the motivation behind HYDRAVIPER’s adaptive training rollout budget allocation (Algorithm 4): after samples are dropped during the resampling procedure, HYDRAVIPER performs rollouts to replenish the dataset.

Specifically, I treat the first iteration as a warm-up period, in which the experts collect a predefined number of  $K_{\text{train}}$  rollouts (Algorithm 4, lines 3–5). This leads to an initial dataset of  $T \cdot K_{\text{train}}$  observation-action pairs. Each team  $\mathcal{T}_\ell$  discards non-critical samples from its dataset (Algorithm 3, line 2), i.e. those where the range in the  $Q$ -value is at most a predefined threshold  $\epsilon$ . With the goal of efficiency in mind, HYDRAVIPER computes the minimum number of such discarded samples across all teams of cooperative agents. This then determines the minimum number of rollouts required to collect at least this many samples in the next iteration. The expected number of dropped rollouts, and therefore the budget for the next iteration, is:

$$K_{\text{drop}} = \min_{\ell} \left\lceil \frac{1}{T} \left| \{(\mathbf{o}_k, \mathbf{a}_k) \in \mathcal{D}_\ell \mid p_{\ell k} \leq \epsilon\} \right| \right\rceil.$$

During the remaining  $M - 1$  iterations, HYDRAVIPER continues to collect rollouts using students until it exhausts its total budget of  $B_{\text{train}}$  training rollouts (Algorithm 4, line 1). Choosing different rollout budgets allows performance and efficiency to be traded off. A higher budget is likely to lead to superior performance, as more rollouts will be collected from students closer to convergence before the budget is exhausted, but it also requires more computation time.

### 7.4.3 Validation Rollouts: UCB Policy Selection

---

**Algorithm 5** UCB Validation Rollouts (VR-UCB)

---

**Input:** Policies  $\{\hat{\pi}_\ell^m\}_{m=1}^M$ , team  $\mathcal{T}_\ell$ , validation rollout budget  $B_{\text{valid}}$ , scaling factor  $c$   
**Output:** Selected policies  $\hat{\pi}_i, \forall i \in \mathcal{T}_\ell$

- 1: **Initialise**  $n_m \leftarrow 0$  for all  $m \in \{1, \dots, M\}$
- 2: **Initialise**  $n_{\min} \leftarrow \lceil 2 \ln B_{\text{valid}} \rceil$
- 3: **Initialise** return estimates:  $\mu_\ell^m \leftarrow \frac{1}{C_{\min}} \sum_{k=1}^{C_{\min}} \bar{R}_{\ell k}, \bar{R}_{\ell k} \sim d(\hat{\pi}_\ell^m, \pi_{-\ell}^*)$
- 4: **for** rollout  $k \in \{1, \dots, (B_{\text{valid}} - mn_{\min})\}$  **do**
- 5:     **Set**  $m^* \leftarrow \operatorname{argmax}_m \hat{\mu}_\ell^m + \sqrt{\frac{c \ln B_{\text{valid}}}{n_m}}$
- 6:     **Collect** mean return:  $\bar{R}_{\ell k} \sim d(\hat{\pi}_\ell^{m^*}, \pi_{-\ell}^*)$
- 7:     **Update** rollout count:  $n_{m^*} \leftarrow n_{m^*} + 1$
- 8:     **Update** running average of mean return:  $\hat{\mu}_\ell^{m^*} \leftarrow \frac{n_{m^*}-1}{n_{m^*}} \hat{\mu}_\ell^{m^*} + \frac{1}{n_{m^*}} \bar{R}_{\ell k}$
- 9: **return**  $\hat{\pi}_\ell^{m^*}, \forall i \in \mathcal{T}_\ell$

---

Following  $M$  iterations, HYDRAVIPER produces  $M$  joint policy profiles for each team. It then must select the best-performing policy profile (Algorithm 2, lines 13–14). HYDRAVIPER iterates through the policy profiles to estimate the team performance of each using a set of validation rollouts. The performance metric it uses is the undiscounted mean return of the team,  $\bar{R}_\ell^m = \frac{1}{T} \sum_{t=0}^T R_\ell(s_t, \hat{\pi}_\ell^m(\mathbf{o}_t))$ .

As is the case for training, collecting validation rollouts is computationally intensive, so these rollouts also need to be efficiently allocated. However, the problem setting differs here. The goal is not to collect a *diverse* set of training rollouts, but rather to identify the *most performant* policy profiles using as few rollouts as possible. The mean return of each policy profile is unknown *a priori*; it must be estimated by selecting policy profiles and performing rollouts with noisy returns. Again, a fixed allocation strategy of  $K_{\text{valid}}$  environment rollouts for each policy profile is wasteful. The rollouts assigned to clearly poorly performing policy profiles could be reallocated to reduce the variance in the estimated returns of promising policy profiles. This motivation aligns with that of multi-armed bandit (MAB) problems.

Given a limited budget of  $B_{\text{valid}}$  rollouts, I represent the task of selecting the best-performing policy profile as a MAB problem. For each team  $\mathcal{T}_\ell$ , the policy profile  $\hat{\pi}_\ell^m$  from each iteration  $m$  is an arm, and its return is a random variable  $\bar{R}_\ell^m$  with unknown mean  $\mu_\ell^m$ . Each rollout samples from one such random variable, which captures the distribution of returns from environment and policy randomness. The objective is to identify the best arm  $m_\ell^* = \operatorname{argmax}_m \mu_\ell^m$  in as few rollouts as possible, i.e. to minimise the regret with respect to the policy that selects  $m_\ell^*$  for every rollout.

In this work, I use a modification of the UCB1 algorithm [13]. This allows me to achieve logarithmic regret given a readily satisfiable assumption: that the returns  $\bar{R}_\ell^m$  of the arms are bounded. Although UCB assumes that the arms are bounded in  $[0, 1]$ , it can be modified in a manner equivalent to rescaling the rewards to remain in  $[0, 1]$ . HYDRAVIPER relies on the general form of the Chernoff-Hoeffding bound:

**Theorem 7.2.** (Theorem 2 of Hoeffding [166]) For independent random variables  $X_1, \dots, X_n$  with mean  $\mu$ ,  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ , and  $a_i \leq X_i \leq b_i, \forall i \in \{1, \dots, n\}$ , then for  $\alpha > 0$

$$\Pr(\bar{X} \geq \mu + \alpha) \leq e^{-\frac{2n^2\alpha^2}{\sum_{i=1}^n (b_i - a_i)^2}}.$$

**Corollary 7.1.** Assume that  $\bar{R}_\ell^m$  is bounded by  $[a, b]$  with  $\Delta = b - a$  for all  $i, m$ . For  $c = 2\Delta^2$ , Theorem 7.2 shows that, for the empirical mean  $\hat{\mu}_\ell^m$ , (notation simplified for clarity)

$$\Pr\left(\hat{\mu}_\ell^m \geq \mu_\ell^m + \sqrt{\frac{c \ln B}{n_m}}\right) \leq e^{-\frac{2c \ln B}{\Delta^2}} = B^{-4}.$$

This is the same bound as demonstrated for UCB1, and the same holds for the lower confidence bound  $\mu_\ell^m - \sqrt{\frac{c \ln B}{n_m}}$ . Overall, this choice of  $c$  leads to the same  $O(\log B)$  regret bound as UCB1. However, HYDRAVIPER can also be extended to use other MAB algorithms. If the mean returns of each policy profile are assumed to be normally distributed, the UCB1-NORMAL algorithm [13] could be used; it also achieves logarithmic regret. This algorithm effectively chooses  $c$  to be proportional to the arms' sample variance; the greater the variance, the wider the confidence bound. An offline (but biased) estimate of the sample variance can also be obtained by performing expert rollouts before running UCB.

Given a total budget of  $B_{\text{valid}}$  validation rollouts, HYDRAVIPER performs them as follows. For each policy profile, it first performs  $n_{\min} = \lceil 2 \ln B_{\text{valid}} \rceil$  rollouts to generate initial estimates of the mean returns (Algorithm 5, lines 2–3). To allocate the remainder of the budget (lines 4–8), HYDRAVIPER follows UCB1 to select the policy profile index for the  $k$ th validation rollout as

$$m_{\ell k}^* = \operatorname{argmax}_m \left( \hat{\mu}_\ell^m(k) + \sqrt{\frac{c \ln B_{\text{valid}}}{n_m(k)}} \right),$$

where  $n_m(k) = \sum_{k'=1}^k \mathbb{1}[m_{\ell k'}^* = m]$  is the number of rollouts that have used policy profile  $m$  thus far,  $\hat{\mu}_\ell^m(k) = \frac{\sum_{k'=1}^k \bar{R}_{\ell k'} \mathbb{1}[m_{\ell k'}^* = m]}{n_m(k)}$  is the empirical mean of the returns  $\bar{R}_{\ell k}$  from policy profile  $m$ ,  $B_{\text{valid}}$  is the total budget of rollouts, and  $c$  is a scaling constant for the confidence bound (see Section 7.5.4). HYDRAVIPER maintains a running average for the mean return of each policy profile, which it updates using the mean return  $\bar{R}_{\ell k}$  of each rollout (line 8).

#### 7.4.4 Agent Clustering: Scaling Up HYDRAVIPER

When resampling the dataset, HYDRAVIPER calculates sample weights using Equation (7.1). This computation requires enumerating joint actions  $\mathbf{a}_\ell$  for each team  $\mathcal{T}_\ell$ , in order to find the worst-case joint action that minimises the team's mean  $Q$ -function,  $\min_{\mathbf{a}_\ell} \bar{Q}^{\pi_\ell^*}(\mathbf{o}_k, \mathbf{a}_\ell, \pi_{-\ell}^*(\mathbf{o}_{-\ell k}))$ . The complexity of this step scales with the size of the joint action space and thus exponentially with the size of the team. Some mixed competitive-cooperative environments (see Section 7.5.1) have an

---

**Algorithm 6** Agent Graph Clustering (Partition)

---

**Input:** Markov game  $(\mathcal{S}, \mathcal{A}, P, R_i, O_i)$ , experts  $\pi^*$ , agent distance function  $d$   
**Output:** Agent teams  $\mathcal{T}_1, \dots, \mathcal{T}_L$

- 1: **Construct** graph  $G = (V = \{1, \dots, N\}, E, w = 0)$
- 2: **for** each agent  $i \in \{1, \dots, N\}$  **do**
- 3:     **for** each agent  $j \in \{1, \dots, N\}$  **do**
- 4:         **Assign** edge weight  $w_{ij} \leftarrow \frac{1}{d(i,j)}$
- 5: **Partition** graph  $\mathcal{T}_1, \dots, \mathcal{T}_L \leftarrow \text{METIS}(G, L)$
- 6: **return**  $\mathcal{T}_1, \dots, \mathcal{T}_L$

---

inherent team structure that can reduce this complexity. In cooperative environments such as TSC, HYDRAVIPER clusters the agent set into teams to improve training efficiency (Algorithm 4).

The goal is to find a clustering of the agent set into teams  $\mathcal{T}_1 \dots \mathcal{T}_L$  so that HYDRAVIPER-trained DT students have *performance* similar to those trained on the full agent set, but improved *scalability* in that the number of actions to enumerate per team is much smaller than the full agent set:  $\prod_{i \in \mathcal{T}_\ell} \mathcal{A}_i \ll \prod_{i \in \{1, \dots, N\}} \mathcal{A}_i$ . I leverage the intuition that agents distant from each other (in terms of environmental distance, trajectory similarity, or other metrics) are unlikely to be influential on each other in most environments.

Suppose that there is a function  $d(i, j)$  that computes this distance between a pair of agents. In my clustering procedure (Algorithm 6), I first construct a complete graph  $G = (V, E) = K_N$  where the nodes represent agents, and the weight between node  $i$  and node  $j$ ,  $w_{ij}$ , is inversely proportional to  $d(i, j)$  (lines 1–4). Then, I perform *graph partitioning* to divide  $G$  into  $L$  contiguous, connected node clusters of approximately equal size (line 5), such that the sum of the weights of inter-cluster edges is minimised. I use the hierarchical METIS algorithm [188] to accomplish this. Note that I solve a graph partitioning problem instead of a min-cut problem to prevent the clusters from being imbalanced. Otherwise, in the worst case, the largest cluster could have size  $O(N)$ , thus yielding minimal gains in scalability.

How can the distance metric  $d$  be defined? For an environment that has an inherent team structure (such as the physical deception environment in Section 7.5.1), I define the graph  $G$  as a complete subgraph for each team, which corresponds to the distance function

$$d(i, j) = \begin{cases} 1 & \text{if } i \text{ and } j \text{ share a team} \\ \infty & \text{otherwise} \end{cases}$$

For traffic signal control environments, I note that the road network inherently forms a graph  $G_{env}$ , which can be partitioned to obtain sets of spatially proximal agents that correspond to neighbouring intersections. In this case,  $G = G_{env}$ , or equivalently

$$d(i, j) = \begin{cases} 1 & e_{ij} \in G_{env} \\ \infty & e_{ij} \notin G_{env} \end{cases}$$

Environment-agnostic distance metrics can also be designed. I consider a measure of proximity that aligns with the VIPER framework: the influence of agents on each other’s  $Q$ -values. Recall from Section 7.4.1 that HYDRAVIPER measures the importance of samples by the mean  $Q$ -function  $\bar{Q}^{\pi_i^*}$ . In a cooperative setting, if one agent’s actions have a significant impact on another agent’s  $Q$ -values, then the joint actions of these agents are likely to have a significant impact on the overall agent set’s mean  $Q$ -function. Including both agents on the same team would allow HYDRAVIPER to capture the effects of these joint actions. I define  $G$  using the distance function

$$d(i, j) = \mathbb{E}_{(\mathbf{o}, \mathbf{a}) \in \mathcal{D}_{\text{dist}}} \frac{2}{\delta_{ij} + \delta_{ji}}, \delta_{ij} = Q^{\pi_j^*}(\mathbf{o}, \pi^*(\mathbf{o})) - \min_{a_i} Q^{\pi_j^*}(\mathbf{o}, a_i, \pi_{-i}^*(\mathbf{o}_{-i})),$$

where  $\delta_{ij}$  is the range in agent  $j$ ’s  $Q$ -values induced by agent  $i$ , and the distance  $d_{ij}$  is the inverse of the average of  $\delta_{ij}$  and  $\delta_{ji}$ . This symmetrisation of influence is a simplifying assumption to obtain a single weight for each edge. Alternative distance metrics could be designed to better capture agent pairs where one is significantly more influential than the other. Since METIS requires integral edge weights, I rescale  $\delta_{ij}$  to percentiles between  $\min_{i,j} \delta_{ij}$  and  $\max_{i,j} \delta_{ij}$ .

## 7.5 Experiments

Now, I demonstrate the utility of HYDRAVIPER for interpretable MARL using experiments in various benchmark environments. In doing so, I perform a functionally grounded evaluation of interpretability [107], where I assess the quality of the generated DTs in terms of *performance* and *computational efficiency*. As the DTs would be used directly in place of NN-based policies in deployment, I consider these to be good proxy metrics for their practical applicability.

### 7.5.1 Environments

I evaluate HYDRAVIPER in four environments: two environments in the *multi-agent particle world* (MPE) benchmark [243], and two *traffic signal control* (TSC) environments in the RESCO benchmark [15]. In MPE environments, agents must navigate in a 2D space to achieve a coordinated objective, making these environments ideal for assessing coordination capabilities. These benchmark environments can also be considered as abstractions of robotic navigation tasks.

**Cooperative navigation (CN)** In this environment, a team of three agents must coordinate to split up and cover three different targets while avoiding collisions with each other.

**Physical deception (PD)** In this environment, a team of two defender agents must cooperate to protect two targets from an adversary agent. One of the two targets is the “goal” of the adversary; this is not known to the adversary, which can only observe the positions of the targets and defenders. I train the two defender agents against an NN adversary.



In TSC environments, each agent controls a single intersection by selecting different signal phases; each phase allows vehicles from a subset of lanes to pass through the intersection. Both environments are based on real-world road corridors reproduced in the traffic simulator SUMO [10]. To interface with the simulator, I use the OpenAI Gym-style wrapper `sumo-rl` [6]. I focus on imitating experts for all agents as a team.

**Cologne corridor (CC)** [413] This environment is the same three-agent environment from the experimental evaluation in Chapter 6.

**Ingolstadt corridor (IC)** [242] This environment is the same seven-agent environment from the experimental evaluation in Chapter 4.

## 7.5.2 Baselines and Setup

I compare HYDRAVIPER with IVIPER and MAVIPER, which represent the state of the art in interpretable multi-agent RL with DT surrogate policies. In addition, I compare with *expert* policies — MADDPG [243] for MPE and MPLight [65] for TSC — and an additional baseline, *imitation DT*. Imitation DT does not use students to collect rollouts, nor does it perform dataset resampling; it collects the same number of training rollouts as the other algorithms and trains DTs on the collected dataset. As imitation DT performs worse than the other algorithms by a wide margin, I only show its performance in Table 7.3.

For MPE environments, I use a horizon of 25 timesteps per episode, and I trained MADDPG for 60 000 episodes as the expert for the DT students to imitate. For TSC environments, I use a horizon of 125 timesteps per episode (each timestep represents 20 seconds of simulation time), and I trained MPLight for 500 episodes as the expert. All imitation learning algorithms were run for 100 iterations to produce DTs with a maximum depth of 4. IVIPER and MAVIPER ran  $K_{\text{train}} = K_{\text{valid}} = 50$  training and validation rollouts per iteration for MPE (including for the initial iteration where rollouts are collected by the experts), and 10 rollouts per iteration for TSC. Imitation DT ran the same number of training rollouts. I set these to equalise the number of environment interactions per iteration.

I repeated all experiments 10 times with different random seeds, and I report the mean and 95% confidence interval of the reward over 10 rollouts performed with the final student policy profiles generated from these runs. Most experiments were run in parallel on a server with 56 2.75GHz AMD EPYC 7453 processors and 252 GiB of RAM. For these experiments, I report the *number of rollouts collected*, not *runtimes*, as the rollout time is roughly constant. However, I also report runtimes for the execution of IVIPER, MAVIPER, and HYDRAVIPER on all four environments. For these experiments, I use the `kernprof` profiler (v4.1.3) to run them in sequence, with no other concurrent processes running except system routines. These experiments were run on another server with 8 4.2GHz Intel i7-7700K processors and 62 GiB of RAM.

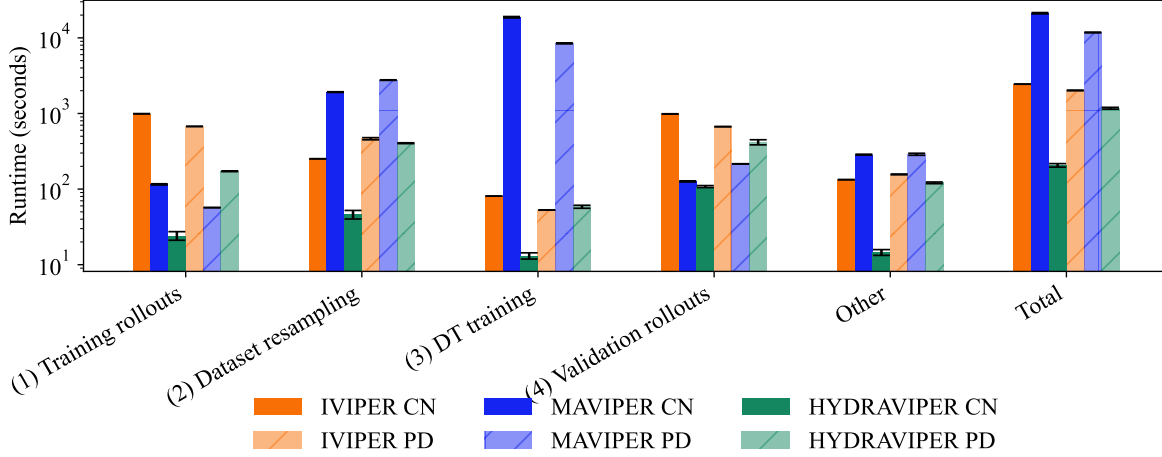


Figure 7.2: Runtime decomposition for IVIPER, MAVIPER, and HYDRAVIPER (with full environment interaction budget) on the multi-agent particle world environments, cooperative navigation and physical deception. Error bars show 95% confidence intervals based on 10 random seeds.

### 7.5.3 Results

**Research Question 7.1 HYDRAVIPER achieves strong, coordinated performance in a computationally efficient manner.** First, I assess HYDRAVIPER’s performance as I vary it between two environment interaction budget levels, high (5 000 training/5 000 validation rollouts for MPE, 1 000 training/1 000 validation rollouts for TSC) and low (500 training/1 500 validation rollouts for MPE, 100 training/100 validation rollouts for TSC). As shown in Table 7.1, HYDRAVIPER students perform better than or comparable to students trained by the most performant DT baseline (MAVIPER for MPE, IVIPER for TSC) in all environments at both budget levels. HYDRAVIPER’s performance is also better than or comparable to the NN experts for all environments except cooperative navigation, in which all DT-based algorithms cannot achieve expert-level performance. In physical deception, although neither MAVIPER nor HYDRAVIPER substantially outperforms IVIPER given the considerable stochasticity in the environment, HYDRAVIPER achieves a level of performance much closer to MAVIPER, while its training time is an order of magnitude shorter than MAVIPER.

In TSC environments, HYDRAVIPER is the best performing algorithm at both the high and low interaction budget levels. Notably, HYDRAVIPER at the high budget level substantially outperforms the expert on the Ingolstadt corridor. By contrast, MAVIPER fails to coordinate the intersection agents and is in general the worst-performing algorithm. HYDRAVIPER more than halves the runtime of both IVIPER and MAVIPER on both TSC environments.

In Figure 7.5, I show a decision tree generated by HYDRAVIPER for one agent in the Ingolstadt corridor at the low budget level (100 training/100 validation rollouts). Even with this limited environment interaction budget, HYDRAVIPER generates DTs that are not only performant, but also intuitively sensible. Each internal node of the DT compares the number of queueing vehicles for some turning movement to a threshold; the left branch includes all samples where there are

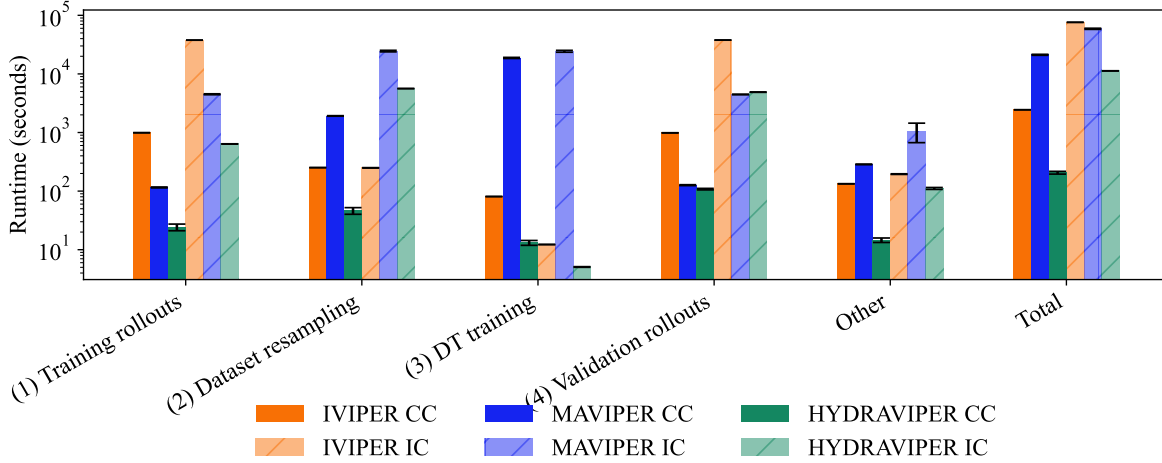


Figure 7.3: Runtime decomposition for IVIPER, MAVIPER, and HYDRAVIPER (with full environment interaction budget) on the traffic signal control environments, Cologne corridor and Ingolstadt corridor. Error bars show 95% confidence intervals based on 10 random seeds.

fewer vehicles than the threshold, while the right branch includes all samples where there are more than the threshold. The root assesses traffic from the westbound side street. If it is low, the DT coordinates north-south traffic on the main road; if it is high, the DT coordinates turn traffic from the side street. Note the similar structure of the depth-3 subtrees, where the left branch splits on a feature, and the right branch predicts a single action. As traffic intensity increases, the action set that the policy can use to optimise traffic shrinks. Stakeholders, such as traffic engineers, can follow such a workflow to visualise, reason about, and supervise DTs generated by HYDRAVIPER.

**Research Question 7.2 As the environment interaction budget decreases, HYDRAVIPER still outperforms baselines.** Now, I investigate the ability of HYDRAVIPER to adapt to increasing budget constraints for environment interaction, as would be imposed by users who wish to quickly iterate on DT policy training. As shown in Figure 7.4, HYDRAVIPER’s performance in all four environments does not change substantially as the training and validation rollout budgets are individually reduced. Furthermore, in all four environments, HYDRAVIPER achieves performance on par with or better than MAVIPER at all budget levels. Therefore, HYDRAVIPER is able to maintain a Pareto frontier in the tradeoff between performance and computational efficiency.

In cooperative navigation, the performance of both HYDRAVIPER and MAVIPER remains similar as the training and validation budgets are reduced individually. However, when both budgets are reduced simultaneously (shown in Table 7.1), the performance of HYDRAVIPER but not MAVIPER remains essentially unchanged. In physical deception, HYDRAVIPER still performs well even as its training budget is reduced by a factor of 10, whereas MAVIPER performs substantially worse. Furthermore, the 95% confidence intervals of HYDRAVIPER’s rewards are smaller than those of MAVIPER at all validation budget levels. Thus, HYDRAVIPER is able to identify performant policy profiles more consistently than MAVIPER.

In the Cologne corridor, HYDRAVIPER’s performance consistently remains within the ex-

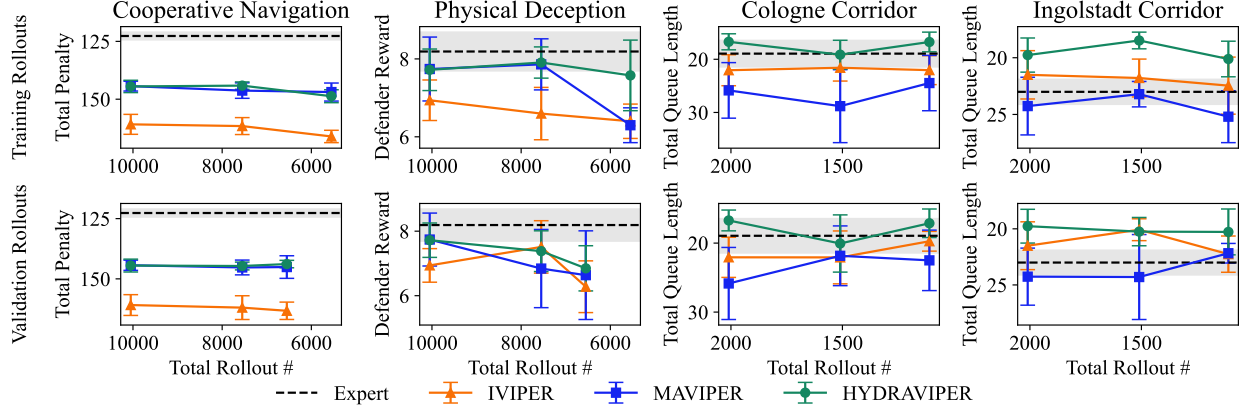


Figure 7.4: Performance of HYDRAVIPER and baselines as the number of rollouts decreases. Top shows decreasing training rollouts; bottom shows decreasing validation rollouts. HYDRAVIPER’s performance stays consistent as the number of rollouts decreases. For physical deception, higher rewards are better; for all other environments, lower rewards are better. Bars show 95% confidence intervals based on 10 randomly-seeded runs. Full results are shown in Table 7.3.

pert’s 95% confidence interval at all environment interaction budget levels, whereas the same is not true of MAVIPER. Meanwhile, the performance of HYDRAVIPER on the Ingolstadt corridor substantially exceeds the expert at all budget levels, whereas MAVIPER and IVIPER (except for the 500 validation rollout setting) remain in the expert’s 95% confidence interval.

**Research Question 7.3 Even when the agent set is decomposed through clustering, HYDRAVIPER maintains its performance.** Finally, I evaluate the effect of agent clustering on the performance of HYDRAVIPER. For the Ingolstadt corridor environment in the high budget setting (1 000 training/1 000 validation rollouts), I evaluate two strategies from Section 7.4.4: (1) clustering the agent set into two teams based on the road network graph  $G_{env}$  (graph-metis), and (2) using pairwise  $Q$  values to identify mutually impactful agents, and either  $k$ -means clustering (marginal-kmeans) or METIS (marginal-metis) for partitioning.

As shown in Figure 7.6 and Table 7.2, these clustering strategies allow HYDRAVIPER to retain its performance even when the size of the agent set is approximately halved for each team. Clustering reduces the largest team’s joint action set in size from 1 944 to 72, and the total runtime of HYDRAVIPER by up to 47%. The best-performing strategy combines pairwise  $Q$ -value weights with METIS for partitioning, but they all perform similarly to the unpartitioned algorithm.

These clustering methods also outperform two baselines. First, the `random` baseline randomly assigns each agent to one of two teams; this baseline has very high variance in performance. Second, the `contiguous` baseline uses a handcrafted division of the agent set into two subsets; `graph-metis` recovers this division automatically, but `marginal-metis` improves further by grouping agents that are not adjacent in the environment.

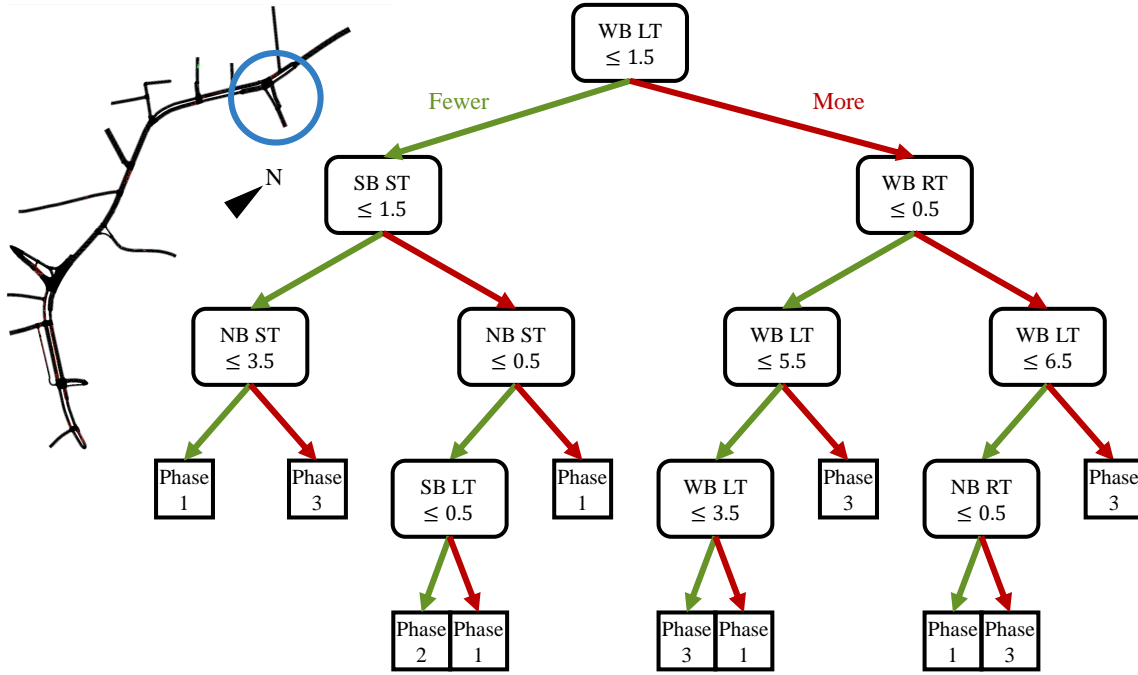


Figure 7.5: Decision tree generated by HYDRAVIPER (100 training/100 validation rollouts) for the Ingolstadt corridor (IC). I have replaced feature indices with human-readable interpretations as turn movements, and I have pruned subtrees that predict a single action.

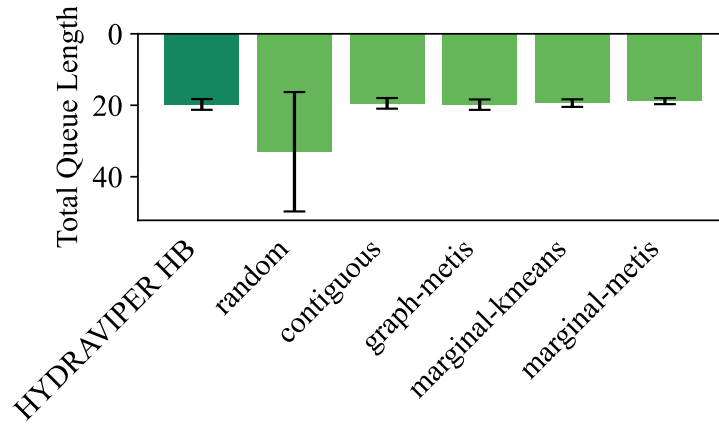


Figure 7.6: Performance of HYDRAVIPER on the Ingolstadt corridor (IC) under different agent set clustering methods at the high budget (HB) level. *graph-metis* applies METIS to the environment graph; *marginal-kmeans* applies  $k$ -means clustering to pairwise  $Q$ -values; and *marginal-metis* applies METIS to the same.

Environment		Expert	IVIPER	MAVIPER
Cooperative Navigation	Total Penalty	$122.67 \pm 1.67$	$160.87 \pm 4.31$	$144.35 \pm 2.12$
	Runtime (s)	N/A	$2\,444.6 \pm 9.1$	$21\,188.7 \pm 408.6$
Physical Deception	Defender	$8.19 \pm 0.50$	$6.94 \pm 0.52$	$7.74 \pm 0.82$
	Reward			
	Runtime (s)	N/A	$2\,017.2 \pm 21.3$	$11\,782.4 \pm 137.8$
Cologne Corridor	Queue Length	$18.94 \pm 2.49$	$22.06 \pm 2.91$	$25.85 \pm 5.22$
	Runtime (s)	N/A	$33\,841.4 \pm 441.3$	$37\,503.8 \pm 834.8$
Ingolstadt Corridor	Queue Length	$23.01 \pm 1.10$	$21.51 \pm 2.13$	$24.26 \pm 2.54$
	Runtime (s)	N/A	$75\,709.6 \pm 441.3$	$52\,316.9 \pm 6\,280.2$
Environment		HYDRAVIPER	HYDRAVIPER LB	
Cooperative Navigation	Total Penalty	$144.48 \pm 2.67$	$144.84 \pm 2.12$	
	Runtime (s)	$206.2 \pm 11.1$	$180.5 \pm 9.3$	
Physical Deception	Defender	$7.72 \pm 0.53$	$7.12 \pm 0.84$	
	Reward			
	Runtime (s)	$1\,173.5 \pm 21.6$	$388.4 \pm 5.6$	
Cologne Corridor	Queue Length	$16.72 \pm 1.51$	$18.77 \pm 3.69$	
	Runtime (s)	$13\,651.6 \pm 254.2$	$1\,865.4 \pm 26.1$	
Ingolstadt Corridor	Queue Length	$19.77 \pm 1.51$	$21.87 \pm 1.59$	
	Runtime (s)	$11\,205.8 \pm 40.2$	$6\,490.0 \pm 98.0$	

Table 7.1: Performance and runtimes (means and 95% confidence intervals) for HYDRAVIPER and baselines. All algorithms are given the same environment interaction budget, except for low budget (LB) HYDRAVIPER (which uses 20% of the rollouts for MPE, 10% of the rollouts for TSC). HYDRAVIPER achieves or exceeds the performance of MAVIPER using a fraction of the runtime, and still performs well in the low budget setting. For physical deception, higher rewards are better; for all other environments, lower rewards are better. Figure 7.2 and Figure 7.3 show runtimes for individual algorithm steps.

Clustering Method	Runtime (s)	Worst Clustering
HYDRAVIPER HB	11 263.8 $\pm$ 55.8	[[1,2,3,4,5,6,7]]
+ random	6 925.5 $\pm$ 897.0	[[1,2,3,6,7],[4,5]]
+ contiguous	6 731.3 $\pm$ 214.3	[[1,2,3],[4,5,6,7]]
+ graph-metis	6 026.4 $\pm$ 39.1	[[1,2,3],[4,5,6,7]]
+ marginal-kmeans	8 538.0 $\pm$ 643.0	[[1,2,3,4,6,7],[5]]
+ marginal-metis	6 067.5 $\pm$ 57.0	[[1,2,4,7],[3,5,6]]

Table 7.2: Runtimes and worst-performing clusterings across 10 different random seeds of HYDRAVIPER under the clustering methods in Figure 7.6. The intersection agent numbers follow those shown in Figure 7.1.

Environment	Training	Validation	Expert	Imitation DT	IVIPER	MAVIPER	HYDRAVIPER
<b>Cooperative Navigation</b>	5000	5000	122.67 $\pm$ 1.67	221.19 $\pm$ 8.58	160.87 $\pm$ 4.31	144.35 $\pm$ 2.12	144.48 $\pm$ 2.67
	2500	5000		211.84 $\pm$ 6.32	161.62 $\pm$ 3.66	146.28 $\pm$ 3.34	144.13 $\pm$ 1.59
	500	5000		218.22 $\pm$ 5.90	166.11 $\pm$ 2.67	146.91 $\pm$ 3.90	148.71 $\pm$ 2.80
	5000	2500			161.94 $\pm$ 4.93	145.27 $\pm$ 3.07	144.66 $\pm$ 1.62
	5000	1500			163.25 $\pm$ 3.60	145.10 $\pm$ 4.69	143.86 $\pm$ 1.54
	500	1500			168.31 $\pm$ 5.52	153.66 $\pm$ 4.07	144.84 $\pm$ 2.12
<b>Physical Deception</b>	5000	5000	8.19 $\pm$ 0.50	6.27 $\pm$ 0.43	6.94 $\pm$ 0.52	7.74 $\pm$ 0.82	7.72 $\pm$ 0.53
	2500	5000		5.73 $\pm$ 0.31	6.60 $\pm$ 0.67	7.84 $\pm$ 0.66	7.91 $\pm$ 0.40
	500	5000		5.32 $\pm$ 0.61	6.40 $\pm$ 0.44	6.30 $\pm$ 0.45	7.58 $\pm$ 0.91
	5000	2500			7.51 $\pm$ 0.81	6.84 $\pm$ 1.21	7.38 $\pm$ 0.63
	5000	1500			6.28 $\pm$ 0.80	6.64 $\pm$ 1.37	6.85 $\pm$ 0.70
	500	1500			6.03 $\pm$ 0.70	7.36 $\pm$ 0.99	7.12 $\pm$ 0.84
<b>Cologne Corridor</b>	1000	1000	18.94 $\pm$ 2.49	137.67 $\pm$ 0.64	22.06 $\pm$ 2.91	25.85 $\pm$ 5.22	16.72 $\pm$ 1.51
	500	1000		211.84 $\pm$ 8.00	21.60 $\pm$ 2.50	28.82 $\pm$ 6.86	19.13 $\pm$ 2.70
	100	1000		218.22 $\pm$ 10.29	22.05 $\pm$ 2.72	24.47 $\pm$ 5.21	16.75 $\pm$ 1.85
	1000	500			22.07 $\pm$ 3.83	21.84 $\pm$ 4.33	20.06 $\pm$ 4.15
	1000	100			19.73 $\pm$ 1.43	22.50 $\pm$ 4.38	17.12 $\pm$ 2.07
	100	100			23.40 $\pm$ 3.53	28.91 $\pm$ 7.27	18.77 $\pm$ 3.69
<b>Ingolstadt Corridor</b>	1000	1000	23.01 $\pm$ 1.10	169.43 $\pm$ 5.26	21.51 $\pm$ 2.91	24.26 $\pm$ 2.54	19.77 $\pm$ 1.51
	500	1000		170.55 $\pm$ 3.41	21.79 $\pm$ 1.67	23.21 $\pm$ 1.12	18.48 $\pm$ 0.74
	100	1000		166.99 $\pm$ 3.67	22.46 $\pm$ 2.52	25.20 $\pm$ 2.29	20.11 $\pm$ 1.57
	1000	500			20.08 $\pm$ 0.97	24.30 $\pm$ 3.79	20.25 $\pm$ 1.26
	1000	100			22.26 $\pm$ 1.62	22.18 $\pm$ 0.87	20.28 $\pm$ 2.04
	100	100			27.23 $\pm$ 5.31	23.31 $\pm$ 1.44	21.87 $\pm$ 1.59

Table 7.3: Performance (means and 95% confidence intervals) for IVIPER, MAVIPER, and HYDRAVIPER at different training and validation budget levels. For physical deception, higher rewards are better; for all other environments, lower rewards are better.

### 7.5.4 Hyperparameter Sensitivity

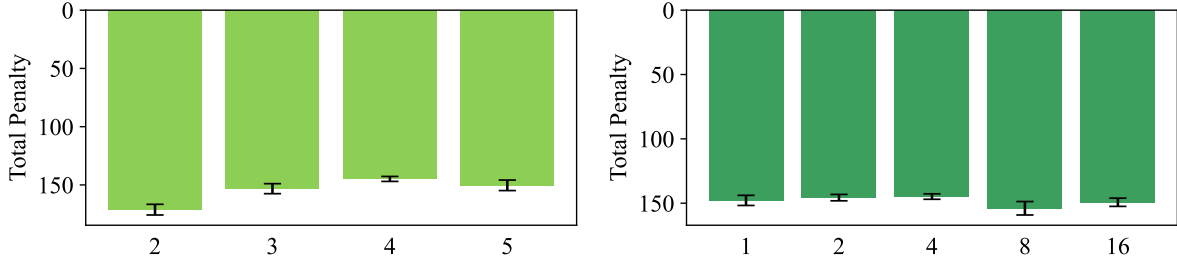


Figure 7.7: Sensitivity of HYDRAVIPER to two hyperparameters, DT depth and the UCB scaling constant  $c$ , on the cooperative navigation (CN) environment. Lower rewards are better.

To understand the effects of HYDRAVIPER’s hyperparameters on its performance, I conduct experiments to vary the depth of the DT students, and the scaling constant  $c$  for UCB policy selection (Section 7.4.3), on the cooperative navigation environment. I choose this environment due to its relatively low level of randomness. For these experiments, I use HYDRAVIPER at the low budget level (500 training/1 500 validation rollouts) as the baseline algorithm, and fix all hyperparameters other than those of interest.

By default, I use a DT depth of 4; my results show that DTs of this depth provide a good tradeoff between expressiveness and computational efficiency. As shown on the left of Figure 7.7, depth-4 DTs outperform depth-2 and depth-3 DTs on cooperative navigation. This is an intuitive result; the optimal agent policy for this environment cannot be represented with such shallow DTs, as they must condition on the positions of the other agents and the landmarks. However, I find that depth-4 DTs also marginally outperform depth-5 DTs. This same pattern exists in all of the environments that I use for evaluation. I hypothesise that the amount of data collected by HYDRAVIPER at the low budget level is insufficient to coordinate between depth-5 DTs.

My default value for  $c$  is also 4. As I outlined in Corollary 7.1 in Section 7.4.3, the value that theoretically allows logarithmic regret to be achieved is based on the range of mean returns in the environment. Since the agents navigate in a  $2 \times 2$  square environment, the maximum distance of an agent to a target is  $2\sqrt{2}$ . The reward in this environment is the negation of the minimum agent distance to each landmark, plus a penalty of -1 for each agent that the ego agent collides with. Therefore, the maximum possible penalty is  $\Delta = 3 \cdot 2\sqrt{2} + 2$ , which requires  $c \approx 219.88$  to achieve the guarantee of Corollary 7.1. However, as shown on the right of Figure 7.7,  $c = 4$  empirically performs best;  $c = 8$  and  $c = 16$  are already excessively conservative given the low randomness in the environment. Since the agents are already trained, I hypothesise that the potential range of returns is less useful in practice for finding a good policy profile than the typical range of returns. This can be approximated by expert rollouts, as I described in Section 7.4.3.

Sensitivity results for the three other environments are also shown below. The best DT depths for these environments are all 4, as with cooperative navigation, while the best tested values of  $c$  for physical deception and the Cologne corridor are, respectively, 2 and 16. In physical deception, the reward is the  $L_2$  distance of the adversary to the goal, minus the minimum  $L_2$  distance of



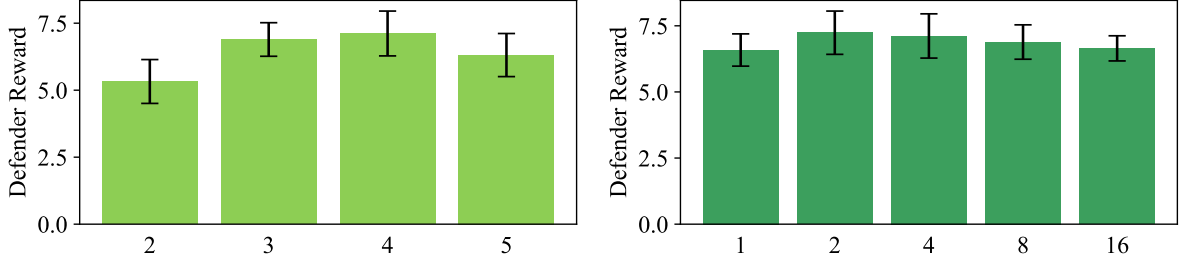


Figure 7.8: Sensitivity of HYDRAVIPER to two hyperparameters, DT depth and the UCB scaling constant  $c$ , on the physical deception (PD) environment. Higher rewards are better.

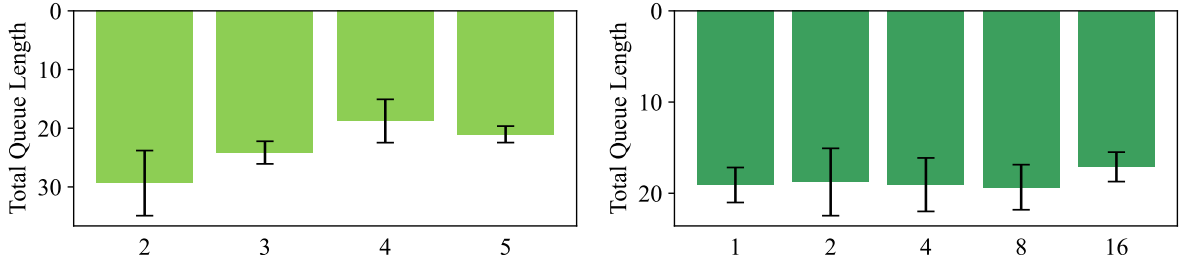


Figure 7.9: Sensitivity of HYDRAVIPER to two hyperparameters, DT depth and the UCB scaling constant  $c$ , on the Cologne corridor (CC) environment. Lower rewards are better.

any defender to the goal. Again, the maximum distance from an agent to the goal is  $\Delta = 2\sqrt{2}$ . Therefore,  $c = 16$  is required to achieve the guarantee of Corollary 7.1. In the TSC environments, the reward in these environments is the total queue length, which is effectively unbounded.

### 7.5.5 Ablation

Lastly, to understand which components of HYDRAVIPER are responsible for its success, I conduct an ablation study for HYDRAVIPER at two budget levels in all four environments. For the high budget level, I use 5 000 training/5 000 validation rollouts for MPE environments, and 1 000 training/1 000 validation rollouts for TSC environments; for the low budget level, I use 500 training/1 500 validation rollouts for MPE environments, and 100 training/100 validation rollouts for TSC environments. I compare HYDRAVIPER’s centralised- $Q$  resampling with IVIPER’s independent resampling (HYDRAVIPER - CQ). In addition, I study the impact of removing adaptive training budget allocation (HYDRAVIPER - TR-A) and UCB-based validation budget allocation (HYDRAVIPER - VR-UCB).

Figure 7.11 shows my ablation results in multi-agent particle world environments, while Figure 7.12 shows my ablation results in traffic signal control environments. In all environments and at both budget levels, centralised- $Q$  resampling outperforms the IVIPER resampling scheme, although the performance gap is more pronounced for cooperative navigation due to randomness in the other environments. This result suggests that sampling the training dataset independently

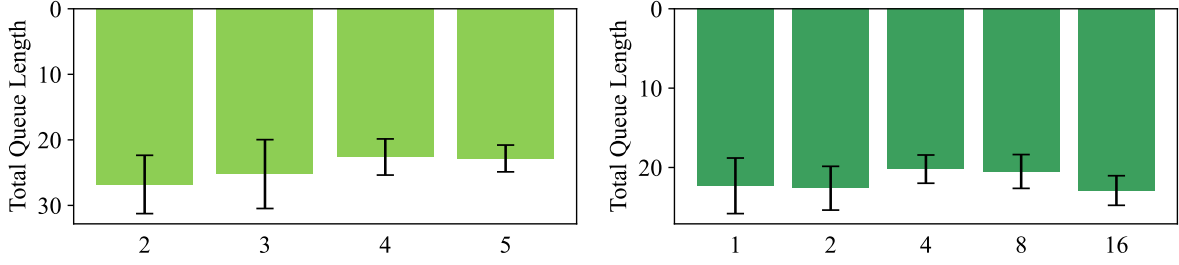


Figure 7.10: Sensitivity of HYDRAVIPER to two hyperparameters, DT depth and the UCB scaling constant  $c$ , on the Ingolstadt corridor (IC) environment. Lower rewards are better.

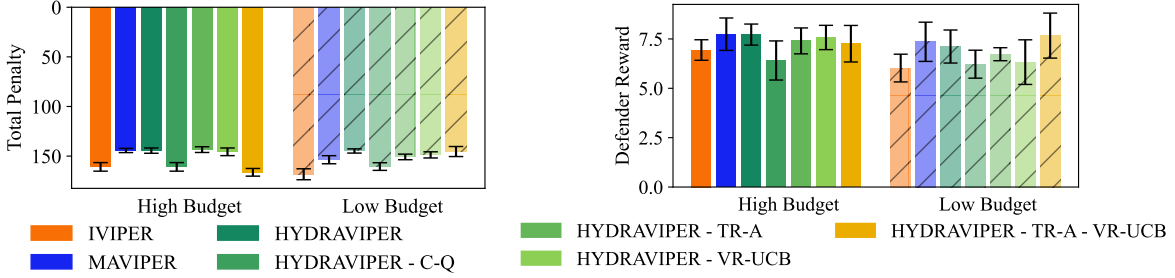


Figure 7.11: Ablation of HYDRAVIPER on the cooperative navigation (CN, left, lower rewards are better) and physical deception (PD, right, higher rewards are better) environments.

for each agent, instead of according to team performance, is insufficient to achieve coordinated behaviour in the resulting students.

Meanwhile, removing the budget allocation methods degrades HYDRAVIPER’s performance. In most cases, having both budget allocation methods is necessary to improve HYDRAVIPER’s reward; having a single method further degrades the reward, especially in low budget settings. Another notable benefit of the budget allocation methods is in reducing the variance of HYDRAVIPER’s reward. This is shown most clearly by the Cologne corridor, where the variance is greater both when *only* centralised- $Q$  resampling is present and also when it is *removed*. For physical deception and Ingolstadt corridor, training budget allocation is more effective at variance reduction than validation budget allocation. These results suggest that the primary benefit of the two rollout budget allocation mechanisms is to stabilise HYDRAVIPER’s learning process, especially in the low budget setting when extracting the most information from each rollout is critical.

## 7.6 Conclusion

In this chapter, I introduced a new DT-based interpretable MARL method, HYDRAVIPER. HYDRAVIPER addresses several limitations of prior multi-agent methods in the VIPER framework: (1) it improves performance by using a joint dataset resampling scheme based on team  $Q$ -values, and (2) it improves computational efficiency by adaptively allocating fixed budgets of environment

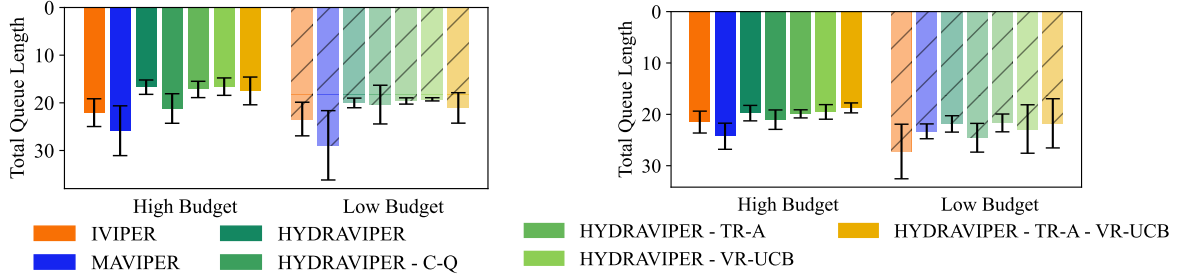


Figure 7.12: Ablation of HYDRAVIPER on the Cologne corridor (CC, left) and Ingolstadt corridor (IC, right) environments. Lower rewards are better.

interactions for training and validation, as well as by dividing agents into jointly-trained teams. Based on experiments in benchmark environments for multi-agent coordination and traffic signal control, I showed that HYDRAVIPER achieves performance comparable with MAVIPER (a centralised method) and even NN experts, all with a runtime less than IVIPER (a decentralised method). I also demonstrated HYDRAVIPER’s sample efficiency in its ability to retain a similar level of performance using a fraction of the environment interactions.

Through my experiments in the Ingolstadt corridor environment, I scaled up the VIPER framework to seven agents. To my knowledge, this is the largest team of coordinated agents to which interpretable MARL has been applied so far. However, environments based on real-world domains can have many more agents than the environments that I studied. For example, the review of Noaen et al. [302] showed that TSC environments of dozens or even hundreds of agents are used in the RL literature. In the most extreme case, Chen et al. [65] used parameter-shared MP-Light policies as controller agents for an extremely large simulation of 2 510 traffic lights. My agent clustering approach shows promise in scaling up to larger environments while retaining performance comparable to that of expert policies. In scaling up, other methods for assigning edge weights  $e_{ij}$  to agent pairs could be developed to leverage structure present in MARL environments, including weights based on observation-action trajectories. I envision that the flexibility of the HYDRAVIPER framework will allow it to adapt to characteristics of different MARL environments while maintaining Pareto optimality in the performance-computational efficiency tradeoff.

Ultimately, HYDRAVIPER and other algorithms in the VIPER framework aim to improve the interpretability of RL policies by converting deep NNs to DTs. In this chapter, I have shown that this conversion can be done in a performant, computationally efficient manner. However, my functionally grounded evaluation has not shown whether DTs are sufficient to help stakeholders better understand these policies. An application-grounded evaluation [107] of HYDRAVIPER that includes user studies would be necessary to assess its practical utility. Indeed, in Chapter 6, I learnt that stakeholders in transportation prefer familiar policy representations. Although DTs are popular in the interpretable AI literature, they are rarely used in traffic engineering. Combining the general framework of HYDRAVIPER with alternative policy structures, including those based on natural language, may help make RL-based traffic signal control more understandable and controllable.

# Chapter 8

## End of the Beginning

### Lessons, Perspectives, and Visions on AI Design and Deployment

The way you see the world... it's different than the rest of us.

---

*Philip J. Coulson*

#### 8.1 Introduction

In this chapter, I conclude by reflecting on my experiences with designing and deploying AI systems in transportation. During the past five years, my two collaborations with stakeholders elucidated their day-to-day operations and their perceptions, preferences, pain points, and priorities. One focused on gig driving; involved *Gridwise*, the developer of a gig driver assistant app; and led to Chapter 3. The other focused on traffic signal control (TSC); involved the *Traffic Management Centre* (TMC) of Strongsville, Ohio; *Econolite*, a leading traffic technology company; and *Path Master*, a traffic technology distributor; and led to Chapters 4 to 7. Had I studied these problems using sanitised formulations, these insights would have been unobtainable.

The lessons I have learnt have driven my technical contributions in Chapters 3 to 7. However, the work presented in these chapters represents the culmination of much broader research efforts. The same is true of many other socially impactful applications of AI technologies [366]: many contributions to AI for social good are non-technical, and methods do not necessarily generalise beyond the immediate application context [46]. For this reason, important, generalisable knowledge about the process of conducting stakeholder-driven research is under-reported.

Section 8.2 begins with takeaways from the challenges that I faced during the design and deployment of AI systems for this thesis. Next, Section 8.3 reflects more generally on the fundamental goals of these collaborations and how my research priorities intersected with the needs of stakeholders. I also envision future deployments that could generalise my work to other problem domains. Finally, Section 8.4 outlines a number of directions for future scholarship to build on the technical contributions that I have made in this thesis.

## 8.2 Lessons on Designing and Deploying AI Systems

Both of my stakeholder collaborations involved extensive problem scoping and formulation. During this process, I gained an appreciation for the importance of stakeholder-centred design. It is insufficient to translate a sanitised problem — along with algorithms designed for it — directly into a deployment. Doing so surfaced fundamental misalignments in mental models and assumptions between the AI systems that I designed and the variety of stakeholders they served. Thus, *heterogeneity* and *assurance* were the most critical challenges during the design process.

Even after thorough design work, unforeseen challenges arise when AI systems are deployed. AI systems are designed in decontextualised environments that represent abstractions of real-world deployment contexts. Unexpected outcomes can occur when unavailable resources and logistical impracticalities cause deployment contexts to diverge from assumptions made during the design and training process. In such cases, AI systems fall short of the potential of their designs. Thus, *uncertainty* and *coordination* were the most critical challenges during the deployment process.

This section reflects on the lessons I have learnt through addressing stakeholders' key considerations related to these design and deployment challenges. Although neither collaboration resulted in the fully realised deployments that I initially envisioned, my hope is that these takeaways will be instructive for future research that aims to design and deploy AI systems for social impact.

### 8.2.1 Incorporating Stakeholder Heterogeneity

Stakeholders are not monolithic. AI systems interact with a wide variety of stakeholders, and unintended consequences can arise when AI systems do not account for *heterogeneity* in their diverse preferences and needs. **To mitigate such issues during deployments, as many stakeholders as possible should be involved in the AI system design process as early and deeply as possible.** For example, while collaborating directly with Gridwise gave me a good understanding of the typical planning problem faced by gig drivers, I could not fully appreciate the diversity of motivations and routines among gig drivers until I conducted my user studies. In interviews, I found that some drivers are more motivated by social interaction than profit (Gridwise pilot P3), while others use unconventional modes of transportation (Gridwise interview P7). These individualised insights would have been difficult to identify from the aggregated, anonymised datasets I used to estimate earnings. **Formative user studies for AI systems that engage with individual stakeholders are thus crucial**, as I noted in Chapter 3. The lack of such studies is a crucial shortcoming of my evaluation in Chapters 6 and 7, and indeed more generally of the AI literature.

More broadly, AI systems for transportation are deployed in complex socio-technical contexts. **Multifaceted interactions between different groups of stakeholders lead to complexity in AI systems' design requirements.** The AI deployment process is not a linear transfer of control from developers to end-users, but also involves organisational structures and communities that may be impacted by their decisions. However, these broader groups of stakeholders are often not considered in the design of responsible AI systems [208]. For the Strongsville project, the city's TMC has internal regulatory requirements for signal plans, but their day-to-day operations are more often informed by complaints from the general public. This influenced their focus when they reviewed the traffic simulation from Chapter 5. Furthermore, their role as an intermediary

between developers and the public led to their relative reluctance to conduct field tests or make configuration changes (as in the case of their camera detectors) compared to Econolite.

### 8.2.2 Optimising for Assurances and Understandability

If stakeholders are so diverse, can a single AI system be designed to provide *assurances* regarding the social welfare of many stakeholders? One-size-fits-all solutions are unlikely to exist. **Problem objectives, solution constraints, and information displays in AI systems should instead be flexibly adapted to the expectations and mental models of stakeholders.** What are these expectations? I learnt from speaking with stakeholders in gig driving and TSC that some expect AI systems to deliver maximal improvements, while others consider a minimum threshold of performance to be satisfactory. However, incorrect assumptions made during the design process may not become evident until users have experienced the AI system in action. **Alignment between AI systems' assumptions and stakeholders' understanding must be thoroughly and repeatedly interrogated to ensure practical utility.** For instance, it was not until the exit interviews that the gig drivers participating in the Gridwise user study reported confusion about whether my tool's earnings estimates should be interpreted as including tips and bonuses (Gridwise interview P2/P6). Likewise, HYDRAVIPER generates performant decision trees in an efficient manner, but they remain insufficiently aligned with the needs of traffic engineers without user-centred evaluations.

Are out-of-the-box AI algorithms sufficient to optimise for stakeholders' specialised needs? **Generic objective functions commonly used in the literature encode restrictive assumptions, which provide insufficient assurances for stakeholders about the behaviour of AI systems.** For gig driving, I initially created an earnings estimation model that optimised for accuracy across the entire dataset of historical trips, but user studies showed that this was not the correct goal. Instead, gig drivers would have preferred higher-confidence estimates, based on a narrower and potentially more predictive set of historical data from other drivers similar to themselves (Gridwise interview P4). For TSC, I found that stakeholders are equally concerned with optimising performance (e.g. queue length) and safety (e.g. crash rate proxies) metrics, but only the former category of metrics are used as RL rewards. Furthermore, the lack of constraints imposed on RL-generated signal plans meant that the Strongsville stakeholders could not trust their good performance, which required me to re-implement both existing algorithms and environments.

### 8.2.3 Managing Uncertainty and Missing Data

While the gap between simulations and reality can be partially bridged by principled environment and algorithm design [345], the dynamic and open nature of transportation systems means that AI technologies will often encounter an element of *uncertainty* that cannot be anticipated during the design process. **Starting with benchmarks and sanity checks can help assess and mitigate the impacts of edge cases during deployments of AI systems.** For both CycleLight (Chapter 6) and HYDRAVIPER (Chapter 7), I iterated on algorithms using small, lightweight benchmark environments. Although these environments were less similar to the deployment context of Strongsville, they allowed me to run larger numbers of replications to assess the algorithm's performance in the presence of randomness. Meanwhile, my work on the Gridwise user study is a negative example

of this lesson. The short one-month timeline of the user study, combined with the difficulty of recruiting participants, motivated me to directly proceed to full deployments with drivers in four cities. However, the scale of this study made it infeasible to individually sanity-check earnings estimates for different time slots and driver types in each city, thus resulting in some estimates that drivers did not find useful. **Keeping pilot tests for AI systems small-scale and easy to run is important to ensure that the uncertainty under assessment remains manageable.**

AI algorithms from the literature often assume that the input data is complete and noise-free. **Engineering the input data, in practice, can be much more time-consuming than designing an AI system, and may even reshape research questions.** The data acquisition process in both my collaborations was fraught with logistical difficulties. For the Gridwise project, the lack of real-time access to gig driver records meant that I could not assess the objective accuracy of my tool’s earnings estimates. Instead, I had to rely on participants’ self-reported perceptions of accuracy, which reshaped my research questions to focus on how subjective perceptions impact trust and reliance. Similarly, for the Strongsville project, my inability to reconfigure Strongsville’s detectors to obtain more accurate vehicle counts led me to engineer a novel computer vision approach. This greatly expanded my demand generation pipeline from what was originally an intermediate step in the RL training process into a standalone contribution (Chapter 5).

## 8.2.4 Coordinating with Human Interaction and Control

While AI algorithms are highly capable of achieving *coordination* between AI-controlled agents, in many deployments, they serve as decision aids rather than primary decision agents. Just as humans cannot be sure that AI systems will perform optimally, AI systems also cannot assume that humans will comply with (i.e. perfectly follow) recommendations in a fully rational manner [264]. **Imperfect compliance with AI systems may lead to unexpected outcomes during deployments.** For instance, during the Gridwise user study, participants did not perfectly follow the recommended schedules, and this led many to report perceived discrepancies in earnings from my tool’s estimates. **When full compliance with AI systems is not guaranteed, information about outcomes under counterfactual decisions can help to assure or even nudge users.** As I remarked in Chapter 3, estimating potential earnings under alternate schedules would have been helpful in maintaining trust, but this was not possible given data limitations. Estimating counterfactual signal performance measures was also a potential contribution that I explored with Econolite.

People will not always run on the same timetables as AI systems. **Judicious decisions need to be made about when AI systems and their deployers, should interact with stakeholders.** User burden must be balanced with appropriate delegation to ensure that these interactions maximise the value of information [152] for decision makers. For the Gridwise user study, I timed email reminders of recommended schedules to coincide with the start of gig drivers’ daily activity, which participants found to be useful (e.g. Gridwise interview P5). But I was also coordinating with Gridwise’s growth team, who needed to manually upload trip records on a daily basis. Miscoordination on this front led to the aforementioned data limitations. For the Strongsville project, iterating directly on the traffic simulation with the Strongsville TMC would not have been productive. Even still, I needed to repeatedly seek feedback from Path Master. This motivated me to develop the LLM agent from Chapter 5 as a way to reduce the user burden for all parties involved.

## 8.3 Perspectives on Goals and Possibilities of AI Deployments

Now, I take a step back from the design and deployment challenges of uncertainty, heterogeneity, assurance, and coordination to consider my collaborations in a more general way. What is the goal of deploying an AI system in transportation? How can a deployment of an AI system be sustained? What other deployments are possible based on the work that I have done in this thesis? My findings from Chapter 3 are a surprisingly critical component of the answer to all three questions: there must be *alignment* between the purposes of AI systems and the needs of stakeholders, and there must be *trust* throughout the interactions between AI systems and stakeholders. The gamut of stakeholder needs within and beyond transportation leads to a variety of potential applications for my work.

### 8.3.1 Aligning Goals in Research and Deployment

When designing AI systems, there is often a gap in complexity and novelty between solutions that address stakeholders' needs and solutions that are technically interesting. **Deployment challenges and research contributions can mutually feed into each other, but for a given project one often takes precedence over the other.** For the Gridwise project, the initial project goal was to develop an accurate earnings prediction model that Gridwise could use for schedule recommendations. However, once our priorities fell out of alignment, I decided to pursue a user study to further knowledge about human-AI interaction. This required the tool's estimates to be significantly simplified. For the Strongsville project, the priority of the Strongsville TMC was to ensure the accuracy of the simulation. My contribution in Chapter 5 mainly integrated existing state-of-the-art methods into a unified framework. However, this framework had the greatest impact in moving the project closer to deployment. By contrast, our priority as researchers was developing novel algorithms, in which case Strongsville was more of a motivating example. Chapter 7 advanced the state of the art in interpretable MARL, but had limited relevance to the actual deployment. **Blockers encountered by deployment-oriented projects are opportunities to make progress in research-oriented projects, and vice versa.** This was certainly the case for my work.

Data privacy and intellectual property are crucial interests in collaborative projects, especially between academic and corporate entities. **Logistics related to these interests delay the deployment of AI systems, but establishing a clear, standard precedent can help reduce overhead.** My collaborations encountered three types of administrative roadblocks: data licencing agreements (DLAs), technology transfer agreements, and Institutional Review Boards (IRBs).

- First, for both Gridwise and Econolite, we needed to finalise DLAs before we could receive any data to work on the respective projects. In the case of Gridwise, repeated iteration and discussion on the DLA meant that, although we finalised the goals of the project in October 2021, I was unable to begin working with any data until March 2022.
- Second, we considered two viable options for companies to make use of my final research products: (1) we would make standard technology transfer agreements with companies, where our work products would be licenced to the companies; or (2) I would re-implement my contributions within the company as an intern. We favoured the former option, but this again required a separate agreement to be made with these companies.



- Lastly, for the Gridwise collaboration, we had to communicate with our IRB to assess the potential of different user study designs. This two-month delay to the project caused Gridwise to lose their availability and their initial sample of drivers.

Had we iterated more quickly on all of these procedures, it is possible that more could have been accomplished from the Gridwise collaboration before their priorities shifted.

### 8.3.2 Fostering Trust to Drive Progress

Success and failure in both of these multi-year collaborative projects were driven by a single, central theme. **Over repeated interactions, trust is the lynchpin that determines the willingness of stakeholders to continuously engage with AI systems.** Again, no one-size-fits-all solution exists. During my collaborations, some stakeholders generally trusted the AI systems that I developed (e.g. Gridwise interview P1 and Econolite), while others were more sceptical and sought to better understand the AI’s inner workings first (e.g. Gridwise interview P4 and the Strongsville TMC). I learnt from Econolite that this is also how their customers respond to existing systems such as Edaptive. Regardless, Econolite emphasised that traffic engineers seek the flexibility to modify algorithmic decisions based on their own domain knowledge, and that this need would be strengthened with RL. **Control and agency over outcomes can help stakeholders trust AI systems.** This lesson was also reflected in feedback that I received from gig drivers about my tool (e.g. Gridwise interview P5, who wanted it to understand their preferences through dialogue).

What is the basis of trust? In Chapter 3, I adopted a definition of trust that entailed the “expectation that [AI systems] will perform [decision-making tasks] important to the trustor” [262]. **Ultimately, the best way to build trust is to assure stakeholders that the AI system is achieving its stated purpose.** To ensure alignment between AI systems and the stakeholders they serve, trust must be evaluated and enforced in a clear, contextualised manner (e.g. my adaptation of the HCT and TiA instruments to the gig driving setting), grounded in the goals and needs communicated by stakeholders during the design process. For both projects, jumping directly into a full deployment (e.g. the Gridwise user study, or initial RL results on Strongsville) did not accomplish this. Instead, I needed to progressively build up to a complete deployment over time, using sanity checks and pilot tests to assure stakeholders of each component’s functionality. For the Strongsville project, I did so by separating the simulation environment from the RL training procedure. Meanwhile, the Gridwise user study remained a negative example of this lesson. Participants wanted to receive feedback on their earnings (Gridwise interview P1), but were instead unable to rationalise critical misalignments between their outcomes and their expectations.

### 8.3.3 Generalising to Applications in Other Domains

Although neither of my collaborations reached full deployment, they yielded both practical lessons and a number of technical contributions. Through pursuing opportunities and responding to challenges, I advanced the frontiers of AI technologies designed to incorporate stakeholders’ needs. My resulting insights generalise to other transportation problems and to the broader AI literature.

Meanwhile, the frontiers of the deployment contexts of AI systems are also evolving rapidly. Especially in the last decade, technological innovations have continuously created both new possi-

bilities and new pitfalls, and society is still learning to adapt to these disruptive changes. However, the core of these socio-technical systems has never changed: the human stakeholders who use and are impacted by AI technologies. Above all else, this thesis sought to establish a framework for designing AI systems to address the deployment challenges that arise from these considerations.

## Within Transportation

Our transportation systems face two critical challenges: they suffer from *systemic inefficiencies*, and they are *rapidly changing*. Although the COVID-19 pandemic temporarily alleviated both challenges, these effects are no longer evident. For example, Schrank et al. [354] measured the average delay of commuters in the United States and found that it more than doubled from 1982 to 2019; it diminished to near-1982 levels in 2020; and it resurged to 2019 levels in 2022. However, society has not returned to the status quo as much as it has established a new normal. My contributions in this thesis are poised to address several emergent trends:

- **Increasing personalisation of travel preferences.** As hybrid work schedules are increasingly replacing the 9-to-5 workday [26], commutes are spreading out from the typical morning and evening rush hours. This flexibility in working schedules has led to a broad reconsideration by commuters of their chosen transportation modes [354].

*Chapter 3* studies gig drivers' motivations, routines, and attitudes, using survey instruments adapted to the specific context of their day-to-day tasks. Similarly grounded formative studies can be used to understand stakeholder preferences in other transportation systems.

*Chapter 4* shows that the way traffic simulations model individual road users' decisions can significantly affect outcomes. This holds for the microscopic driver behaviour I considered in that chapter, but also for more macroscopic behaviour such as mode choice.

- **Increasing demand for ridesourcing.** Due to the personalisation of travel preferences, the growth of the gig driving industry rebounded rapidly after 2020. Although this growth has tapered, the industry remains optimistic about its future [379]. The continued growth of gig driving will increase economic opportunities but also opportunity costs for gig drivers [443].

*Chapter 3* contributes not just a predictive tool to estimate gig drivers' earnings, but also a prescriptive tool to recommend work schedules for them. The tool can be flexibly modified to account for other objectives, e.g. travel time, that drivers may be concerned about.

- **Increasing freight traffic volumes.** Compared to physical stores, online retailers and other e-commerce platforms have grown in popularity, even amid global economic volatility [280]. The resulting movement of goods has also increased the proportion of traffic — and thus congestion — that can be attributed to freight vehicles [354].

*Chapter 3*'s schedule recommendation tool is intended for gig drivers, but my findings can be generalised to present uncertainty in delivery schedules to freight platforms. These schedules would need to account for not just temporal availability but also spatial movement.

*Chapter 5*'s demand modelling pipeline includes a simple heuristic for computing heavy vehicle traffic volumes and incorporating them into simulations to assess their impact. My vehicle tracking-based counting method could also be modified to distinguish vehicle classes.

- **Increasing awareness of the need to replace basic infrastructure.** Particularly in the United States, a significant amount of infrastructure is no longer able to meet current demand [318]. This is exemplified by the collapse of the 50-year-old Fern Hollow Bridge in 2022 [43]. Substantial investments have been made, but room for improvement remains [431].

*Chapter 4* focuses on traffic simulators, which are an important part of the policy analysis toolkit in transportation. My experimental methodology from that chapter can be used to statistically compare the effects of different infrastructural improvements (e.g. road redesigns).

*Chapter 5* further builds upon the foundation of Chapter 4 with an LLM agent-based framework that allows stakeholders to iteratively refine traffic simulations. This proof of concept could be extended to enable the real-time assessment of policy interventions.
- **Increasing shift of vehicular traffic from urban roads to freeways.** Generally, congestion-related delays have grown faster for freeways in the United States than they have grown for city arterials [354]. Whatever the cause, this shift implies an increase in long-distance travel. More freeway capacity can be beneficial, but the effects are unlikely to be immediate [11].

*Chapter 5* improves the accuracy of demand modelling, and also scales it up to a large road network. Balancing accuracy and computational efficiency is the key to network-level traffic analysis in a city like Strongsville, which has heavy traffic at junctions with freeways.

*Chapter 6* integrates operational constraints into RL-based TSC, providing stakeholders with greater understandability and controllability over signalling decisions. Such assurances are important as the scale of traffic management expands beyond municipal boundaries.
- **Increasing maturity of connected and automated vehicles (CAVs).** Despite persistent safety and regulatory concerns [289], makers of CAVs are aggressively expanding into new markets worldwide [114]. However, full penetration is unlikely to occur for decades; until then, CAVs must operate in a hybrid traffic environment that includes human drivers [237].

*Chapter 5* leverages error calibration to integrate multiple sources of detector data into a single, more accurate simulation. Vehicle-to-infrastructure communications are a promising source of low-error data that could be used as part of my demand modelling pipeline.

*Chapter 7* lessens the cognitive burden of stakeholders in understanding complex deep learning policies by distilling them into decision trees. If applied to CAV algorithms, decision trees could help stakeholders interpret behaviour and assess counterfactual outcomes.
- **Increasing ownership and use of private vehicles.** Instead of making use of transit or using ridesourcing services, individuals are increasingly owning their own vehicles despite increasing car prices [354]. This is also partially in anticipation of economic instability [81]. As is the case for trucks, however, this increased ownership will increase congestion.

*Chapter 6* jointly optimises for both the performance and the safety of RL-based TSC policies. Since individual drivers are both less efficient and less safe than buses or rail transit [236], it will be important to optimise these criteria centrally through traffic signals.

*Chapter 7*, just as with CAVs, can be applied to make RL-based TSC policies more understandable. In particular, stakeholders may increasingly need to communicate policy decisions to individual members of the general public who are impacted by heavy traffic.

## Beyond Transportation

What other domains can my contributions and insights apply to? In this thesis, I viewed gig driving and traffic signal control as representative problems for studying the design and deployment of AI technologies. As I outlined in Table 1.1, both problems have characteristics — including their multi-agent, longitudinal natures — that generalise to a variety of domains beyond transportation.

- **Chapter 3** focuses on trust and reliance on AI technologies, in an uncertain, longitudinal environment with elevated stakes. These characteristics are hardly exclusive to gig driving. AI systems — especially LLMs — have penetrated every dimension of society, yet trust in AI is decreasing [138]. My findings on how to foster trust through task alignment and personalisation apply broadly across deployments of AI decision aids in healthcare, education, and criminal justice [53]. However, the impact of uncertainty on trust could differ between contexts. For example, gig drivers receive near-instant feedback through earnings, but medical practitioners may not observe the impacts of their decisions until long afterwards.
- **Chapter 4** calibrates two traffic simulators against each other as an intermediate step to ensure the ecological validity of simulation results. The sim-to-real problem, which aims to align simulations with reality, is a well-known challenge in robotics [345]. However, many robotic systems have well-established physical dynamics, leaving limited room for variation in simulator design. Other embodied systems where simulations are used to train AI technologies have less predictable dynamics, including homes [97], hospitals [308], and battlefields [223]. In these domains, experimental studies such as mine can help verify simulators’ modelling assumptions, and assess tradeoffs between simulation fidelity and efficiency.
- **Chapter 5** solves a specific instance of a flow prediction problem on a network, and it does so by integrating data from multiple error-prone sources. This framework can be generalised to other flow prediction tasks, such as those involving computer networks [311], electrical circuits [36], or river systems [456]. More generally, my LLM agent-based simulation refinement and reflection procedure can be adapted to various deployment contexts for aligning the outputs of AI systems with stakeholders’ expectations. For example, Behari et al. [34] applied this framework to healthcare: they used an LLM agent to learn a multi-armed bandit reward function that encapsulates desired policy outcomes for particular sub-populations.
- **Chapter 6** augments action postprocessing methods for RL with an imitation-based loss and differentiable optimisation. Both of these methods are well-suited to the enforcement of operational constraints in RL application domains that involve cyclical control actions, such as robotics (especially locomotion control tasks) [17, 155] and sensor networking [110]. More generally, action-constrained RL can also be applied to other safety-critical contexts such as power systems [468], cyber-physical systems [409], and financial investments [298].
- **Chapter 7** distils RL policies into decision trees using the HYDRAVIPER algorithm. This interpretable RL workflow readily generalises to all of the aforementioned RL applications, but especially to large-scale problems where many agents act independently and where human oversight is needed. In addition to TSC, I evaluated HYDRAVIPER in two coordination environments from Lowe et al. [243]. These environments can be considered as abstractions of robotic navigation tasks; others in the benchmark are based on computer security.

## 8.4 Visions on Furthering the Capabilities of AI Technologies

My work has filled a number of critical gaps in the AI literature. However, it was also limited by design and deployment challenges. There is room for future research to build on the capabilities of my contributions and thus improve the status quo of transportation systems. In this section, I conclude by summarising potential improvements to my work.

- In **Chapter 3**, I developed a tool that recommends working schedules for gig drivers, and estimates the potential earnings of drivers from following these schedules. The tool’s estimation and scheduling modules both have substantial room for improvement, since the focus of my study was testing designs of the tool’s user interface, not improving the tool’s accuracy. The earnings estimates were directly based on averages across historical trips. Earlier prototypes had used a richer set of features, including platforms and geographic locations, to predict earnings. Based on feedback from the user study, I also envision that predictions could be tailored to specific clusters of drivers. To this end, data from a larger and more diverse sample of drivers would provide further insight into the diversity of driver experiences. Beyond working times, drivers must also choose different platforms to drive for and regions to drive in. Because trips reposition drivers to regions with different levels of supply and demand, adding a spatial component makes the problem more sequential. Another extension that I previously considered is the coordination of recommendations between drivers, which can mitigate the failure mode of too many drivers crowding a potentially lucrative region.
- In **Chapter 4**, I conducted an experimental comparison of simulation outcomes from two traffic simulators, SUMO and CityFlow. Although certain aspects of simulation heterogeneity, such as driver behaviour and traffic scale, increased the divergence in outcomes between the simulators, my experiments were insufficient to identify the source of the discrepancies. SUMO and CityFlow make a variety of different assumptions and modelling choices that are not mutually compatible. I made several changes to CityFlow’s driver behaviour models to align them to SUMO, but many other differences remain. An ablation-type experiment could help isolate the factors that contribute the most to these discrepancies. It is unclear whether either simulator replicates real-world traffic with sufficient realism. Data from CAVs could be used to validate their driver behaviour models, as well as the input parameters of the models. My experiments used somewhat arbitrary values for car-following and lane-changing parameters, which could be better tailored to specific road networks.
- In **Chapter 5**, I built a algorithmic pipeline for demand modelling in traffic simulations that is more rigorous than existing approaches. However, several aspects of the pipeline still relied on heuristics. These included the methods I used to process unstable camera detector footage, and to distribute vehicular flows within time segments and between vehicle classes. One of my main contributions was the use of an LLM agent to refine the simulation from natural language feedback. Although the agent uses reflection to identify errors in generated simulations, it is still not designed to handle inconsistent or vague feedback from stakeholders. Implementing a dialogue-based system for interacting with the LLM agent and integrating external data sources would make my pipeline more usable by stakeholders.

Another limitation of my pipeline is its scalability. The Strongsville road network was small enough that I could exhaustively enumerate possible routes, but this approach would not scale to larger geographic regions. Clustering routes within subregions could be helpful for scaling up. Improvements to optimisation routines — particularly for the QP that enforces temporal continuity in traffic flow — would also accelerate simulation generation.

- In **Chapter 6**, I combined a suite of action postprocessing techniques — CycleLight — to impose operational constraints on RL-based signal plans. The main limitation of my work in this chapter is the incompleteness of my evaluation on the Strongsville network; the CycleLight-constrained policies had not converged at the time of writing.

My inability to complete the training process for Strongsville can be attributed to two bottlenecks: the simulation environment and the training process. Although my wrapper environment for SUMO is flexible, the efficiency of its implementation is likely suboptimal. The differentiable optimisation subroutine for action projection could also be sped up.

Variability in traffic flow significantly impacts the performance of TSC policies. I evaluated CycleLight using only two simulations with pre-defined traffic patterns, and I also did not assess the generalisation performance of the policies. In addition, I was unable to replicate Strongsville’s Edaptive controllers, and had to evaluate against fixed time-of-day plans.

- In **Chapter 7**, I designed a performant, scalable algorithm — HYDRAVIPER — to distil RL policies into decision trees. However, my evaluation did not address the fundamental question of whether these decision trees improve the understandability of RL policies for stakeholders. User studies would be necessary to ensure their practical utility.

Nevertheless, HYDRAVIPER is ultimately more of an improvement upon existing algorithms than one designed especially for stakeholders. Even if decision trees can be visualised and used to assess counterfactual scenarios, stakeholders in TSC prefer to directly evaluate cycle-offset-split plans. The HYDRAVIPER framework could be extended to imitate adjustments to cyclic policies as atomic actions, or to output policies based on natural language.

Within the HYDRAVIPER framework, various algorithmic subroutines could be substituted. Potential improvements could involve the decision tree training algorithm, the training and validation budget allocation methods (e.g., respectively, importance sampling-based approaches that discard unlikely observation-action pairs, and alternative variants of the UCB algorithm), and the agent clustering strategy (e.g. using trajectory embeddings).

Although the end of this thesis represents a stopping point for my work, it is by no means an apotheosis for the design and deployment of AI systems in the real world. However, through this thesis, I have argued that technical improvements alone are not and cannot be sufficient. The crux of this forward movement must be the stakeholders who will be impacted. Only by working with stakeholders to deeply understand, to collaboratively build, and to responsibly use AI technologies can we manifest a better tomorrow through deploying AI systems in transportation and beyond.

# Appendix A

## Appendix

### A.1 Chapter 3

#### A.1.1 Survey Questions

##### Survey 1 — Intake Survey

1. Please tell us which of the following regions you primarily drive in.
    - Los Angeles
    - New York
    - Chicago
    - Houston
  2. Please tell us which of the following services you currently drive for.
    - DoorDash
    - Grubhub
    - Instacart
    - Lyft
    - Uber
    - Uber Eats
    - Other
  3. In 2 or 3 sentences, please tell us what you like most about driving for ridesharing and/or delivery services.
  4. In 2 or 3 sentences, please tell us what you like least about driving for ridesharing and/or delivery services.
- How much do you agree or disagree with the following statements?*
5. When I drive, it's important to me that I make some minimum amount of money.
    - Strongly disagree

- Disagree
  - Not sure
  - Agree
  - Strongly agree
6. When I drive, I have an accurate sense of how much money I will make.
- Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree
7. When I don't earn the amount that I expect to from driving, it causes difficulties for me.
- Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree
8. I try to stick to a regular routine for times and places to drive.
- Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree
9. I am happy with how I currently decide when and where to drive.
- Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree
10. Please tell us your age.
11. Please tell us what gender you identify as.
- Male
  - Female
  - Non-binary



- Prefer not to disclose
- Other

12. Please tell us your highest education level.

- Less than high school
- High school
- Some two-year professional degree
- Some undergraduate degree
- Some graduate degree (MS, PhD, JD, or MD)

## Survey 2 — Pre-Survey

*How much do you agree or disagree with the following statements?*

1. I understand how the tool used my answers to generate this recommended schedule.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item U2 from the Perceived Understandability questions of the HCT [257], “I understand how the system will assist me with decisions I have to make.” As the constraint page is intended to encapsulate the user’s decision-making process, I consider the generation of the recommended schedule to be how the tool assists the user with their decisions.*

2. I feel that I can rely on the tool to produce recommendations which accommodate the things that matter most to me.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item R4 from the Perceived Reliability questions of the HCT [257], “I can rely on the system to function properly.” I consider the tool to be properly functioning if its recommendations account for the user’s goals and preferences.*

3. I feel that the driving times recommended by the tool are as good as what an experienced driver would recommend to me.

- Strongly disagree

- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item T3 from the Perceived Technical Competence questions of the HCT [257], “The advice the system produces is as good as that which a highly competent person could produce.” My tool’s advice is its recommended schedule, and to my participants a competent individual would be an experienced driver.*

- I feel that the times suggested by the tool are good even if I don’t know for certain that they will maximise my earnings / minimise my hours *[depending on the constraints selected]*.
  - Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree

*This item is based on item F1 from the Faith questions of the HCT [257], “I believe advice from the system even when I don’t know for certain that it is correct.” Again, my tool’s advice is its recommended schedule of driving times. Since no clear notion of correctness applies to continuous estimates of earnings, I reworded this question to focus on alignment with the user’s objectives.*

- I would like to use the tool to decide my driving hours in the future.
  - Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree

*This item is based on item P4 from the Personal Attachment questions of the HCT [257], “I like using the system for decision making.” I reworded it to better assess participants’ level of intended future reliance on the tool.*

### **Survey 3 — End-of-Day Survey**

- How often did you follow the times in the recommended schedule today?
  - I did not follow the recommendations at all
  - I followed the recommendations for one hour during the day
  - I followed the recommendations for two or three hours during the day
  - I followed the recommendations for four or more hours during the day

2. How satisfied do you feel you are with your earnings from today?
  - Very dissatisfied
  - Somewhat dissatisfied
  - Neither satisfied nor dissatisfied
  - Somewhat satisfied
  - Very satisfied
3. As far as you remember, how did your earnings today compare to your expectations?
  - Lower
  - About the same
  - Higher
  - Not sure
4. As far as you remember, how did your earnings today compare to the tool's estimate?
  - Lower
  - About the same
  - Higher
  - Not sure

*How much do you agree or disagree with the following statements?*

5. I felt that the recommended schedule provided by the tool was easy to follow.
  - Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree

*This item is based on item U4 of the Perceived Understandability questions of the HCT [257], "It is easy to follow what the system does." Instead of asking the user about the tool's operation generally, I focused the question on the interpretability of its recommended schedule for that day.*

6. I felt that the recommended schedule provided all of the information that I needed to decide when to drive.
  - Strongly disagree
  - Disagree
  - Not sure
  - Agree
  - Strongly agree

*This item is based on item R1 of the Perceived Reliability questions of the HCT [257], “The system always provides the advice I require to make my decision.” Again, my tool’s advice is its recommended schedule. I focused the question on the user’s decisions for that particular day.*

7. When I was unsure of when to drive today, I followed the recommended schedule.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item F2 from the Faith questions of the HCT [257], “When I am uncertain about a decision I believe the system rather than myself.” I focused the question on the user’s decisions for that particular day, and reworded “believe” to “follow” to assess compliance more clearly.*

*I left out questions based on Perceived Technical Competence and Personal Attachment for length.*

8. Which of the following statements do you agree with most?

- I intend to rely on the tool less tomorrow than I did today
- I intend to rely on the tool about the same tomorrow as I did today
- I intend to rely on the tool more tomorrow than I did today
- I intend to pause my interaction with the tool for one day tomorrow

#### **Survey 4 — Post-Survey**

*How much do you agree or disagree with the following statements?*

1. I feel that I have become familiar with how to use the tool.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item 12 of the TiA [179], “I am familiar with the system.” I reworded the question in light of the fact that users did not have any existing experience with using the tool before the study.*

2. When I am using a navigation app that suggests routes to me, I feel like I would want to follow the suggestions more if the app asked me questions about my preferences (like this tool did) before giving its suggestions.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is an original question that prompts the participant to consider their interactions with other types of recommendation systems. It assesses the extent to which participants would appreciate granular controls based on their preferences in such systems.*

3. Were there questions that you wanted the tool to ask you that it didn't? If so, please tell us about them in 2 or 3 sentences.

*How much do you agree or disagree with the following statements?*

4. I felt that I was able to trust the schedules recommended by the tool.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item 11 of the TiA [179], "I can trust the system." I focused the scope of this question on the output of the tool, the recommended schedule, rather than the tool as a whole.*

5. I felt that I was able to depend on the schedules recommended by the tool for deciding when to drive.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item 9 of the TiA [179], "The system is dependable." Again, I focused the scope of this question on the output of the tool, the recommended schedule.*

6. When I am using a navigation app that suggests routes to me, I feel like I would want to follow the suggestions more if the app gave me information about the minimum and maximum possible time of the trip (similar to what this tool did).

- Strongly disagree
- Disagree
- Not sure
- Agree

- Strongly agree

*This item is an original question that prompts the participant to consider their interactions with other types of recommendation systems. It assesses the extent to which participants would appreciate increased exposure of uncertainty through range-based estimates in such systems.*

7. I felt that the recommended schedule was misleading.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item 1 of the TiA [179], “The system is deceptive.” In addition to focusing the scope of this question on the output of the tool, I also reworded “deceptive” to “misleading” to capture the broader possibility of the tool being perceived as unintentionally providing incorrect information.*

8. I felt that the recommended schedule harmed my earnings.

- Strongly disagree
- Disagree
- Not sure
- Agree
- Strongly agree

*This item is based on item 5 of the TiA [179], “The system’s actions will have a harmful or injurious outcome.” In this context, the outcome for the user is their earnings from driving while following the recommended schedule. I reworded the question to assess the outcome retrospectively.*

9. If there were any, please identify some of the driving times recommended by the tool that did not align with your expectations.

- When was the time?
- In 1 or 2 sentences, why did it not align with your expectations?

10. Do you have any other questions or comments regarding this tool that you would like to share with us?

## A.1.2 Interview Scripts

### Pilot Interviews

**Formative Questions** *As the interviews were semi-structured, the script below focuses on the guiding questions that I asked participants. I also probed participants further depending on their responses.*

- Please tell us what you like most about driving for ridesharing and/or delivery services.
- Please tell us what you like least about driving for ridesharing and/or delivery services.
- When you drive, how important is it to you that you make some minimum amount of money daily/weekly?
- When you drive, do you have an accurate sense of how much money you will make?
- Do you try to stick to a regular routine for times and places to drive?
- Are you happy with your current routines in terms of when and where you drive?
- Would getting recommendations for times to drive would be helpful to you?

**Evaluative Questions** We will show you a tool that can suggest personalised driving schedules. The tool will ask you some questions about your availability and preferences, as well as revenue targets that you might have. Different people might want to use the tool differently. However, we expect a typical user to use it as follows.

First, they would fill in some information about when they are available during the week, along with either how long they want to work or how much they want to make. The tool will then suggest a recommended schedule for the entire week. As they return to the tool every day to plan out their schedules, users will have the opportunity to interact with the tool, tweaking their availability and possible revenue targets to see how the recommendations change.

*I begin screensharing the constraint page prototype.*

- Here's the initial page of the tool that lets you specify your availability and goals.
- Do you believe you understand what is being shown on this page?
- Do you feel that this tool is asking you the right questions about your availability and goals?
- Are there other important questions that you wish the tool would ask?
- Think about your upcoming week. Using this screen, please tell us what information you think you would want to enter to get a useful recommendation. We will click on the page for you.

*I switch to the schedule page prototype.*

- Here is an example of a recommended schedule that the tool would generate based on the information you just provided.
- Do you feel that you understand what the recommended schedule is suggesting?
- What part of the recommended schedule do you feel is the most useful?

- What part of the recommended schedule do you feel is the least useful? Is there anything that's missing from the schedule?
- Do you feel the recommended schedule gives you enough information to decide whether you would want to follow it?

Finally, we'd like to ask about your overall opinion of the tool.

- What did you like about this tool?
- What did you dislike about this tool?
- Did you feel that interacting with the tool took too much time, or that it was too complicated or confusing for you? Why or why not?
- What sort of information would increase the chance that you want to use this tool and follow its recommendations? How big of a difference do you think that having this information would make?
- Do you believe that drivers would generally find a tool like this to be useful for when they're planning their driving? Why or why not?

## **User Study Interviews**

**Formative Questions** Let's start with talking about your driving for rideshare/delivery services in general.

- Could you start by telling us why you are driving?
  - Is it primary or supplemental income?
  - What other commitments do you balance it with (jobs, family, hobbies)?
- To what extent do you rely on making a target amount when you are driving?
- Can you talk through your typical process for deciding when to drive?

**Feedback on Constraint Design** Now, let's think back to the times when you were interacting with the tool, particularly when it asked you to enter your availability and goals.

- How similar or different were the tool's questions to the way you typically make these decisions?
- Did you feel like you were able to use the tool to adequately specify your main considerations for when you'd like to drive?
  - Were you ever unsure of what information the tool was asking for?
  - Would you have preferred the tool to ask for information differently, or to ask for different information?
- Did you feel like you were able to influence the recommended schedule that the tool generated for you?
  - Did you try to experiment with entering in different information?



- Did you feel that interacting with the tool took too much time, or that it was too complicated or confusing for you?
  - Could you see yourself spending more time interacting with the tool than you did (e.g. to enter more details)? Why or why not?

**Feedback on Schedules** Now, let's talk about your how the tool's recommended schedules may or may not have influenced your driving activity over the last few days.

- Did you find that the recommended schedules made sense?
  - Did you feel that you understood how the tool used your answers to generate schedules? Why or why not?
- To what extent did you rely on the email reminders of the schedules?
  - Did you ever miss the email reminders?
  - When did you typically check the schedule, if at all?
- If you saw the recommended schedules, how did they impact your process for deciding when to work?
  - To what extent did you follow the schedules?
  - Were there times at which you prioritised your own intuition over the schedules?
  - If so, were there times at which you wished you followed the schedule more closely? Why or why not?
- How did your response to recommended schedules change throughout the week, if at all?
  - Did you look at the schedules more or less as time went on?
  - Were there any particular days on which you wanted to check the schedule more? Why or why not?
- Did you feel that the tool gave you more or different information than you would otherwise get from the services that you drive for/from Gridwise? Why or why not?
- Are there some additional details which could have increased the chance that you followed the recommended schedules?
  - For example, would you have preferred to see the entire week's schedule on every day?
- Did the recommended schedules lead you to drive at different times and/or locations than before?
  - Did this happen early on or later?
  - At what times of day?
- When you followed the recommended schedules, did you feel that you ended up making more money, less money, or about the same relative to before?
  - How closely do you track your earnings in general?
  - Did you track your earnings more closely when using the tool?

### Feedback on Estimates

- How much did you focus on the tool's estimates for how much you could earn?
- Did you feel like you could rely on the estimates to achieve your earning goals?
  - Did you feel that these estimates were meant to be accurate projections of how much you could earn, or that they were rough ballpark figures?
- In general, did you feel that the estimated earnings had the right level of detail, or would you have liked to see additional information?
  - *[If participants were in Conditions (B) or (D)]* Would you have preferred to see a range for how much you could earn?
  - *[If participants were in Conditions (R) or (RH)]* Would you have preferred to see a single number for how much you could earn?

*I select a particular day on which the participant interacted with the tool. If earnings data was available, this was a day on which the participant earned more than the tool's estimate; otherwise, this was the sixth day of their interaction with the tool.*

- Let's talk about [weekday], Day [day] of your interaction with the tool.
- Do you recall the extent to which you looked at the recommended schedule?
  - If you did, do you remember how you decided whether you wanted to follow it? Was this influenced by how much the tool estimated your earnings to be? Why or why not? What did you think of the estimate that the tool gave you?
  - If you didn't, do you happen to remember why? Was this influenced by how much you earned on the previous day? Why or why not?
- Do you recall whether you made more or less than the tool estimated on that day?
  - If there were any differences, do you have any idea why?
- Did that influence your decision to look at the recommended schedule for the next day? Why or why not?
- We checked your records briefly and found that you earned \$xxx.xx, compared to the tool's estimate of \$xxx.xx. Does that change how you feel at all?

**Overall Thoughts** Now, we'd like to wrap up with a few general questions about the tool.

- Did you feel that the time you spent interacting with the tool was worthwhile or not worthwhile? Why or why not?
- If you had the option of using a tool like this one, what are the chances that you might actually use it to decide your driving schedule in the future? Why or why not?
- Beyond what you've mentioned already, is there anything else you believe might increase the chance that you would use this tool in the future?
- Do you have any other questions, comments, or concerns?

### A.1.3 Full Quantitative Results

In Table A.1, I report the means, standard errors, and 95% confidence intervals of the daily trust and reliance scores plotted in Figures 3.7 and 3.8.

Day	Condition	Trust			Reliance		
		$\mu$	$SE$	95% CI	$\mu$	$SE$	95% CI
1	Overall	3.353	0.196	(2.872, 3.834)	0.176	0.107	(-0.100, 0.453)
	(B)	3.556	0.444	(2.468, 4.643)	0.167	0.307	(-0.623, 0.957)
	(D)	3.033	0.390	(2.080, 3.986)	0.200	0.200	(-0.314, 0.714)
	(R)	3.810	0.348	(2.959, 4.660)	0.143	0.261	(-0.528, 0.813)
	(RH)	3.242	0.379	(2.315, 4.170)	0.182	0.182	(-0.286, 0.649)
2	Overall	3.578	0.167	(3.169, 3.988)	0.235	0.095	(-0.009, 0.479)
	(B)	3.722	0.416	(2.703, 4.741)	0.333	0.333	(-0.524, 1.190)
	(D)	3.567	0.205	(3.065, 4.069)	0.200	0.200	(-0.314, 0.714)
	(R)	3.857	0.397	(2.885, 4.830)	0.429	0.202	(-0.091, 0.948)
	(RH)	3.333	0.362	(2.447, 4.220)	0.091	0.091	(-0.143, 0.325)
3	Overall	3.500	0.168	(3.089, 3.911)	0.118	0.110	(-0.165, 0.400)
	(B)	3.833	0.331	(3.025, 4.642)	-0.167	0.307	(-0.957, 0.623)
	(D)	3.400	0.364	(2.508, 4.292)	0.100	0.180	(-0.361, 0.561)
	(R)	4.000	0.309	(3.245, 4.755)	0.571	0.202	( 0.052, 1.091)
	(RH)	3.091	0.270	(2.430, 3.752)	0.000	0.191	(-0.490, 0.490)
4	Overall	3.755	0.148	(3.392, 4.118)	-0.029	0.123	(-0.346, 0.287)
	(B)	3.778	0.306	(3.028, 4.527)	-0.167	0.307	(-0.957, 0.623)
	(D)	3.667	0.157	(3.282, 4.051)	-0.200	0.200	(-0.714, 0.314)
	(R)	4.238	0.347	(3.390, 5.086)	0.286	0.286	(-0.449, 1.020)
	(RH)	3.515	0.334	(2.697, 4.333)	0.000	0.234	(-0.600, 0.600)
5	Overall	3.588	0.159	(3.199, 3.978)	0.147	0.120	(-0.162, 0.457)
	(B)	3.278	0.475	(2.116, 4.439)	-0.167	0.307	(-0.957, 0.623)
	(D)	3.733	0.257	(3.104, 4.363)	0.100	0.180	(-0.361, 0.561)
	(R)	4.143	0.290	(3.434, 4.852)	0.429	0.297	(-0.336, 1.193)
	(RH)	3.273	0.273	(2.605, 3.940)	0.182	0.226	(-0.400, 0.764)
6	Overall	3.667	0.163	(3.268, 4.065)	0.088	0.122	(-0.226, 0.402)
	(B)	3.333	0.487	(2.142, 4.525)	-0.333	0.333	(-1.190, 0.524)
	(D)	3.733	0.276	(3.059, 4.408)	0.200	0.200	(-0.314, 0.714)
	(R)	4.143	0.340	(3.311, 4.975)	0.714	0.184	( 0.240, 1.188)
	(RH)	3.485	0.275	(2.813, 4.157)	-0.182	0.182	(-0.649, 0.286)
7	Overall	3.843	0.149	(3.478, 4.208)	N/A	N/A	N/A
	(B)	3.889	0.351	(3.029, 4.749)	N/A	N/A	N/A
	(D)	3.733	0.247	(3.128, 4.339)	N/A	N/A	N/A
	(R)	4.381	0.286	(3.682, 5.080)	N/A	N/A	N/A
	(RH)	3.576	0.292	(2.862, 4.289)	N/A	N/A	N/A

Table A.1: Statistics for daily trust and reliance, as measured by End-of-Day Surveys.

## A.2 Chapter 5 — Full Claude Prompt

*In the following prompt, the placeholders  $\{\text{intersection}\}$ ,  $\{\text{int\_idx}\}$ , and  $\{\text{feedback}\}$  are automatically filled in based on the feedback provided by the human stakeholder.  $\{\text{num\_intersections}\}$ ,  $\{\text{intersections\_list}\}$ , and  $\{\text{main\_roads\_list}\}$  are filled in based on JSON-structured data, which can be customised to different deployment contexts.*

Consider a road network with  $\{\text{num\_intersections}\}$ . The following list shows the intersections, in the format “zero-based intersection index. (intersection number) - intersection name”:

$\{\text{intersections\_list}\}$

The main roads are:

$\{\text{main\_roads\_list}\}$

You are trying to generate a traffic simulation where vehicles follow routes that are consistent with observed counts.

For each intersection in the road network, you are given a count of how many vehicles enter the intersection in the eastbound, northbound, southbound, and westbound directions. Denote these counts as  $f_j$  for each of  $m$  counting locations  $j$ .

You are also given a set of  $n$  possible routes that vehicles can follow. Denote the number of times that route  $i$  is used as  $r_i$ . You are given a  $n \times m$  matrix  $\mathbf{A}$ , where entry  $\mathbf{A}_{ij}$  is 1 if route  $i$  passes counting location  $j$ , and is 0 otherwise.

For some counting locations  $j$ , you are also given specific counts for left turn vehicles and right turning vehicles. Denote these counts as  $f_{LT,j}$  and  $f_{RT,j}$  respectively. You are also given  $n \times m$  matrices  $\mathbf{A}_{LT}$  and  $\mathbf{A}_{RT}$ , where entries  $\mathbf{A}_{LT,ij}$  and  $\mathbf{A}_{RT,ij}$  are 1 if route  $i$  passes counting location  $j$  and respectively turns left or right, and are 0 otherwise.

The counting locations  $j$  are indexed as follows. Each of the intersections listed above occupies four consecutive entries for eastbound, northbound, southbound, and westbound.

You solve a quadratic program to find a solution for the  $n$ -dimensional variable  $\mathbf{r}$ .

This is the first requirement: For each counting location  $j$ , you want the sum of all routes  $i$  that pass through  $j$  to be within some range of the given counts  $f_j$ . Let  $\mathbf{s}$  be an  $m$ -dimensional slack variable, and let  $\alpha_{LB}$  (e.g. 0.94) and  $\alpha_{UB}$  (e.g. 1.12) be lower and upper bounds on the given counts. This requirement can be represented as the constraints  $\mathbf{Ar} + \mathbf{s} \geq \mathbf{f}\alpha_{LB}$ , and  $\mathbf{Ar} + \mathbf{s} \leq \mathbf{f}\alpha_{UB}$ , and the objective function  $\min_{\mathbf{s}} \sum_j s_j^2$ .

This is the second requirement: For the counting locations  $j$  where turning vehicle counts are given, they should also be matched within some range of the given counts. Let  $\mathbf{s}_{LT}$  and  $\mathbf{s}_{RT}$  be  $m$ -dimensional slack variables, and let  $\alpha_{LB}$  and  $\alpha_{UB}$  be lower and upper bounds as before. This requirement can be represented as the constraints  $\mathbf{A}_{LT}\mathbf{r} + \mathbf{s}_{LT} \geq \mathbf{f}_{LT}\alpha_{LB}$ ,  $\mathbf{A}_{LT}\mathbf{r} + \mathbf{s}_{LT} \leq \mathbf{f}_{LT}\alpha_{UB}$ ,  $\mathbf{A}_{RT}\mathbf{r} + \mathbf{s}_{RT} \geq \mathbf{f}_{RT}\alpha_{LB}$ , and  $\mathbf{A}_{RT}\mathbf{r} + \mathbf{s}_{RT} \leq \mathbf{f}_{RT}\alpha_{UB}$ , and the objective function  $\min_{\mathbf{s}_{LT}, \mathbf{s}_{RT}} \sum_j s_{LT,j}^2 + s_{RT,j}^2$ .

This is the third requirement: Some of the routes begin and end in the middle of a road and are therefore less realistic. They can still be used to satisfy the counts, but as few as possible of them should be used. If “nonfringe” is the set of such unrealistic routes, this requirement can be represented as the objective function  $\lambda \min_{\mathbf{r}} \sum_{i \in \text{nonfringe}} r_i$ , where  $\lambda$  is a tunable parameter.

Altogether, this is the quadratic program that you solve:

$$\begin{aligned}
& \min_{s, s_{LT}, s_{RT}} \sum_j s_j^2 + s_{LT,j}^2 + s_{RT,j}^2 + \lambda \min_{\mathbf{r}} \sum_{i \in \text{nonfringe}} r_i \\
& \text{subject to } \mathbf{A}\mathbf{r} + \mathbf{s} \geq \mathbf{f}\alpha_{LB} \\
& \quad \mathbf{A}\mathbf{r} + \mathbf{s} \leq \mathbf{f}\alpha_{UB} \\
& \quad \mathbf{A}_{LT}\mathbf{r} + \mathbf{s}_{LT} \geq \mathbf{f}_{LT}\alpha_{LB} \\
& \quad \mathbf{A}_{LT}\mathbf{r} + \mathbf{s}_{LT} \leq \mathbf{f}_{LT}\alpha_{UB} \\
& \quad \mathbf{A}_{RT}\mathbf{r} + \mathbf{s}_{RT} \geq \mathbf{f}_{RT}\alpha_{LB} \\
& \quad \mathbf{A}_{RT}\mathbf{r} + \mathbf{s}_{RT} \leq \mathbf{f}_{RT}\alpha_{UB} \\
& \quad r \text{ is a vector of nonnegative integers}
\end{aligned}$$

You solve this problem once for every 15-minute interval over 24 hours.

Now, a traffic engineer has reviewed the simulation you generated from the solution. They have given some feedback on the traffic at the intersection with number  $\{\text{intersection}\}$ , which is at index  $\{\text{int\_idx}\}$  in the list given above. This is their feedback for the “ $\{\text{time}\}$ ” time segment: “ $\{\text{feedback}\}$ ”

The attached Python function, `solve_routes`, solves this quadratic program using the Python library `cvxpy`. Modify this function to add one or more new constraints to the quadratic program so that it takes this feedback into account.

The inputs and constants in the function are as follows:

- `input_locs`: A 3D array for  $\mathbf{A}_{\text{fringe}}$ , the counting locations passed by the subset of realistic (“fringe”) routes. The first dimension is the route, the second dimension is the counting location, and the third dimension is indexed as 0 = total flow, 1 = left turn flow, 2 = right turn flow. The indices of the counting locations follow the intersection numbering given previously.
- `eps_locs`: A 3D array for  $\mathbf{A}_{\text{nonfringe}}$ , the counting locations passed by the subset of unrealistic (“nonfringe”) routes. The first dimension is the route, the second dimension is the counting location, and the third dimension is indexed as 0 = total flow, 1 = left turn flow, 2 = right turn flow.
- `cat_locs`: A 3D array formed by concatenating `input_locs` and `eps_locs`. The first dimension is the route, the second dimension is the counting location, and the third dimension is indexed as 0 = total flow, 1 = left turn flow, 2 = right turn flow.
- `primary_counts`: A 2D array for  $\mathbf{f}$ , the counts at each of the counting locations. These counts are more reliable and should be given priority in matching. Each row is a time segment and each column is a counting location.
- `secondary_counts`: A 2D array for  $\mathbf{f}$ , the counts at each of the counting locations. These counts are less reliable and do not need to be matched as closely. Each row is a time segment and each column is a counting location.

- `hour_labels`: An array of string labels following the format “(number with no leading zeroes):(two-digit number) (AM/PM)”, e.g. “{time}”.
- `scale_prim_lb`: A float for  $\alpha_{LB}$ , denoting how much the traffic from more reliable primary counts should be scaled by as a lower bound, default 1.
- `scale_prim_ub`: A float for  $\alpha_{UB}$ , denoting how much the traffic from more reliable primary counts should be scaled by as an upper bound, default 1.
- `scale_sec_lb`: A float for  $\alpha_{LB}$ , denoting how much the traffic from less reliable secondary counts should be scaled by as a lower bound, default 1.
- `scale_sec_ub`: A float for  $\alpha_{UB}$ , denoting how much the traffic from less reliable secondary counts should be scaled by as an upper bound, default 1.
- `eps_lambda`: A float, a hyperparameter used to weight the penalty for unrealistic (“non-fringe”) routes.
- `time_lambda`: A float, a hyperparameter used to weight the squared loss between the solution from timestep  $t$  and timestep  $t + 1$ .
- `route_lambda`: A float, a hyperparameter used to weight the squared loss between the previous solution and the current solution.
- `start_time`: A datetime object representing a 24-hour timestamp of the day, before which secondary counts are used in place of primary counts (i.e. before sunrise).
- `end_time`: A datetime object representing a 24-hour timestamp of the day, after which secondary counts are used in place of primary counts (i.e. after sunset).
- `time_limit`: A float, a time limit for the quadratic program solver, default 60 seconds.
- `prev_route_counts`: A 2D array for  $\mathbf{r}^{\text{old}}$ , the solution from the previous iteration of the quadratic program. Each row is a time segment and each column is a route.
- `CONNECTIONS`: A dictionary, where each index is one of the intersections named above.
- `TURN_MAPS`: A dictionary, where the indices are “EB”, “NB”, “SB”, and “WB” for the incoming directions at each intersection.

Let's think step by step.

1. Compute the counts of vehicles that entered intersection  $\{\text{intersection}\}$  from each direction under the previous simulation during the " $\{\text{time}\}$ " time segment. The total counts ( $\mathbf{Ar}$ ), left-turn counts ( $\mathbf{A}_{LT}\mathbf{r}$ ), and right-turn counts ( $\mathbf{A}_{RT}\mathbf{r}$ ) for the previous solution  $\mathbf{r}^{\text{old}}$  are accessible through the `get_counts` tool. DO NOT use the tool any more than the minimum number of times that is necessary to answer this question.
  - a. Determine what counting location, time segment, and count type flag to pass into the `get_counts` tool.
  - b. Using the `get_counts` tool, determine how many vehicles entered intersection  $\{\text{intersection}\}$  during the " $\{\text{time}\}$ " segment from the eastbound, northbound, southbound, and westbound directions.
  - c. If the count for a particular direction is 0 for every time segment during the day, that means that the given direction does not exist at this intersection. You should not add a constraint for this particular direction at this intersection, or else the problem will be infeasible.
2. Based on the traffic engineer's feedback, determine if the current simulation is reflective of real-world traffic conditions at this intersection. If it is reflective, you should add a constraint to maintain the current level of traffic at this intersection. If it is not reflective, you should add a constraint to modify the current level of traffic at this intersection.
3. In mathematical notation, write out a constraint corresponding to the traffic engineer's feedback. This constraint may involve parameters with unspecified values, along with  $f_j$  and  $\mathbf{A}_{ij}$ . It should NOT involve the slack variables  $s_j$ , as the constraint may make the problem infeasible.
4. Use your best judgement to determine what values of the unspecified parameters in the constraint would best correspond to the traffic engineer's feedback.
5. Explain how your parameter values from step 4 are sensible, given the eastbound, northbound, southbound, and westbound traffic at intersection  $\{\text{intersection}\}$  that you calculated in step 2 for the previous simulation.
6. Determine the indices of counting locations ( $j$ ) corresponding to intersection  $\{\text{intersection}\}$ .
  - (a) What is the index of intersection  $\{\text{intersection}\}$  in the list of intersection names given above? This should be a 0-based index.
  - (b) If each intersection occupies 4 counting locations, what are the counting locations for intersection  $\{\text{intersection}\}$ ? These should be 0-based indices. For example:  
Intersection 4 is the first entry (index 0) in the list given above. It would occupy indices  $j = 0 \times 4 + 0 = 0, 0 \times 4 + 1 = 1, 0 \times 4 + 2 = 2$ , and  $0 \times 4 + 3 = 3$ .  
Intersection 6 is the second entry (index 1) in the list given above. It would occupy indices  $j = 1 \times 4 + 0 = 4, 1 \times 4 + 1 = 5, 1 \times 4 + 2 = 6$ , and  $1 \times 4 + 3 = 7$ .

7. Translate the constraint from step 4 to a single line of `cvxpy` code.

It should NOT involve the slack variables named `slack_prim`, `slack_lt_prim`, `slack_rt_prim`, `slack_sec`, `slack_lt_sec`, or `slack_rt_sec`, because otherwise the problem may be infeasible.

It should also NOT involve the count variables named `primary_counts` and `secondary_counts`, because they may not exist. Instead, the constraint should be in terms of hardcoded numerical values wherever possible; failing that, it should be in terms of the previous solution `prev_route_counts`; and failing that, it should be in terms of `total_count_flow`, `lt_count_flow`, or `rt_count_flow`.

8. For the “{time}” time segment, add the line of code from step 7 to the list of constraints in the `solve_routes` function. You should modify the rest of the code as little as possible.
9. For the 15-minute time segments adjacent to the “{time}” time segment, add a similar constraint that interpolates between normal traffic and traffic that follows the traffic engineer’s feedback. If possible, use the same constraint for both before and after. You should modify the rest of the code as little as possible.
10. Output the entirety of the modified `solve_routes` function. Do not skip any lines of code. Beyond the constraints you just added, you should modify the rest of the code as little as possible.



# Bibliography

- [1] Montasir Abbas, Darcy Bullock, and Larry Head. 2001. Real-Time Offset Transitioning Algorithm for Coordinating Traffic Signals. *Transportation Research Record* 1748 (2001), 26–39.
- [2] Monireh Abdoos and Ana L.C. Bazzan. 2021. Hierarchical traffic signal optimization using reinforcement learning and traffic prediction with long-short term memory. *Expert Systems with Applications* 171 (2021), 114580.
- [3] Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J. Zico Kolter. 2019. Differentiable Convex Optimization Layers. In *Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS '19)*. NeurIPS, Vancouver, Canada, 9562–9574.
- [4] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. 2022. BoT-SORT: Robust Associations Multi-Pedestrian Tracking. *arXiv:2206.14651*
- [5] Sanjeevan Ahilan and Peter Dayan. 2019. Feudal Multi-Agent Hierarchies for Cooperative Reinforcement Learning. *arXiv:1901.08492*
- [6] Lucas N. Alegre, Ana L.C. Bazzan, and Bruno C. da Silva. 2021. Quantifying the impact of non-stationarity in reinforcement learning-based traffic signal control. *PeerJ Computer Science* 7 (2021), e575.
- [7] Yasir Ali, Zuduo Zheng, Md. Mazharul Haque, Mehmet Yildirimoglu, and Simon Washington. 2020. Understanding the discretionary lane-changing behaviour in the connected environment. *Accident Analysis & Prevention* 137 (2020), 105463.
- [8] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. 2018. Safe Reinforcement Learning via Shielding. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI '18)*. AAAI, New Orleans, USA, 2669–2678.
- [9] Tal Altshuler, Yaniv Altshuler, Rachel Katoshevski, and Yoram Shiftan. 2019. Modeling and Prediction of Ride-Sharing Utilization Dynamics. *Journal of Advanced Transportation* 2019 (2019), 6125798.
- [10] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lücken, Johannes Rummel, Peter Wagner, and Evamarie Wießner. 2018. Microscopic Traffic Simulation using SUMO. In *Proceedings of the 21st International Conference on Intelligent Transportation Systems (ITSC '18)*. IEEE, Pis-

cataway, USA, 2575–2582.

- [11] Alex Anas. 2024. “Downs’s Law” under the lens of theory: Roads lower congestion and increase distance traveled. *Journal of Urban Economics* 139 (2024), 103607.
- [12] Anthropic. 2025. *The Claude 3 Model Family: Opus, Sonnet, Haiku*. Technical Report. Anthropic. 1–42 pages.
- [13] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. 2002. Finite-time Analysis of the Multiarmed Bandit Problem. *Machine Learning* 47 (2002), 235–256.
- [14] James Ault, Josiah P. Hanna, and Guni Sharon. 2020. Learning an Interpretable Traffic Signal Control Policy. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS ’20)*. IFAAMAS, Auckland, NZ, 88–96.
- [15] James Ault and Guni Sharon. 2021. Reinforcement learning benchmarks for traffic signal control. In *Proceedings of the 35th Conference on Neural Information Processing Systems, Datasets and Benchmarks Track (NeurIPS ’21)*. NeurIPS, Virtual, 1–11.
- [16] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Fine-Tuning Language Models from Human Preferences. arXiv:2108.07732
- [17] Shingo Ayabe, Takuto Otomo, Hiroshi Kera, and Kazuhiko Kawamoto. 2025. Robustness Evaluation of Offline Reinforcement Learning for Robot Control Against Action Perturbations. arXiv:2412.18781
- [18] Pierre-Luc Bacon, Jean Harb, and Doina Precup. 2017. The Option-Critic Architecture. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI ’17)*. AAAI, San Francisco, USA, 1726–1734.
- [19] Esmaeil Balal, Ruey Long Cheu, Thompson Gyan-Sarkodie, and Jessica Miramontes. 2014. Analysis of Discretionary Lane Changing Parameters on Freeways. *International Journal of Transportation Science and Technology* 3, 3 (2014), 277–296.
- [20] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the Whole Exceed its Parts? The Effect of AI Explanations on Complementary Team Performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI ’21)*. ACM, Yokohama, Japan, 1–16.
- [21] Angelo Banse and Jakob Erdmann. 2025. *Demand — Definition of Vehicles, Vehicle Types, and Routes*. German Aerospace Centre. [https://sumo.dlr.de/docs/Demand/Introduction\\_to\\_demand\\_modelling\\_in\\_SUMO.html](https://sumo.dlr.de/docs/Demand/Introduction_to_demand_modelling_in_SUMO.html)
- [22] Angelo Banse and Jakob Erdmann. 2025. *Introduction to SUMO Networks*. German Aerospace Centre. [https://sumo.dlr.de/docs/Networks/SUMO\\_Road\\_Networks.html](https://sumo.dlr.de/docs/Networks/SUMO_Road_Networks.html)
- [23] Angelo Banse and Jakob Erdmann. 2025. *SSM Device*. German Aerospace Centre. [https://sumo.dlr.de/docs/Simulation/Output/SSM\\_Device.html](https://sumo.dlr.de/docs/Simulation/Output/SSM_Device.html)
- [24] Yue Bao, Guangzhi Zang, Hai Yang, Ziyong Gao, and Jiancheng Long. 2023. Mathematical

- modeling of the platform assignment problem in a ride-sourcing market with a third-party integrator. *Transportation Research Part B: Methodological* 178 (2023), 102833.
- [25] Jaume Barceló. 2010. Models, Traffic Models, Simulation, and Traffic Simulation. In *Fundamentals of Traffic Simulation*, Jaume Barceló (Ed.). Springer, New York, USA, 1–62.
  - [26] Jose Maria Barrero, Nicholas Bloom, and Steven J. Davis. 2021. *Why Working from Home Will Stick*. Technical Report 28731. National Bureau of Economic Research. 68 pages.
  - [27] Osbert Bastani. 2021. Safe Reinforcement Learning with Nonlinear Dynamics via Model Predictive Shielding. In *Proceedings of the 2021 American Control Conference (ACC '21)*. ACM, New Orleans, USA, 3488–3494.
  - [28] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. 2018. Verifiable Reinforcement Learning via Policy Extraction. In *Proceedings of the 32nd Conference on Neural Information Processing Systems (NeurIPS '18)*. NeurIPS, Montréal, Canada, 2494–2504.
  - [29] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software* 67, 1 (2015), 1–48.
  - [30] Matthew Battifarano and Sean Qian. 2019. Predicting real-time surge pricing of ride-sourcing companies. *Transportation Research Part C: Emerging Technologies* 107 (2019), 444–462.
  - [31] Nelson Baza-Solares, Ruben Velasquez-Martínez, Cristian Torres-Bohórquez, Yerly Martínez-Estupiñán, and Cristian Poliziani. 2022. Traffic Simulation with Open-Source and Commercial Traffic Microsimulators: A Case Study. *Communications — Scientific Letters of the University of Zilina* 24, 2 (2022), E49–E62.
  - [32] Luca Bedogni, Marco Gramaglia, Andrea Vesco, Marco Fiore, Jérôme Härri, and Francesco Ferrero. 2015. The Bologna Ringway Dataset: Improving Road Network Conversion in SUMO and Validating Urban Mobility via Navigation Services. *IEEE Transactions on Vehicular Technology* 64, 12 (2015), 5464–5476.
  - [33] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamviboonsuk, and Laura M. Vardoulakis. 2020. A Human-Centered Evaluation of a Deep Learning System Deployed in Clinics for the Detection of Diabetic Retinopathy. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)*. ACM, Honolulu, USA, 1–12.
  - [34] Nikhil Behari, Edwin Zhang, Yunfan Zhao, Aparna Taneja, Dheeraj Nagaraj, and Milind Tambe. 2024. A Decision-Language Model (DLM) for Dynamic Restless Multi-Armed Bandit Tasks in Public Health. In *Proceedings of the 38th Conference on Neural Information Processing Systems (NeurIPS '24)*. NeurIPS, Vancouver, Canada, 1–38.
  - [35] Paniz Behboudian, Yash Satsangi, Matthew E. Taylor, Anna Harutyunyan, and Michael Bowling. 2022. Policy invariant explicit shaping: an efficient alternative to reward shaping. *Neural Computing and Applications* 34 (2022), 1673–1686.
  - [36] Dominik Beinert, Clara Holzhüter, Josephine M. Thomas, and Stephan Vogt. 2023. Power flow forecasts at transmission grid nodes using Graph Neural Networks. *Energy and AI* 14

(2023), 100262.

- [37] Thor Berger, Carl Benedikt Frey, Guy Levin, and Santosh Rao Danda. 2020. Uber happy? Work and well-being in the “Gig Economy”. *Economic Policy* 34, 99 (2020), 429–477.
- [38] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d.O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. 2019. Dota 2 with Large Scale Deep Reinforcement Learning. arXiv:1912.06680
- [39] Umang Bhatt, Javier Antorán, Yunfeng Zhang, Q. Vera Liao, Prasanna Sattigeri, Riccardo Fogliato, Gabrielle Melançon, Ranganath Krishnan, Jason Stanley, Omesh Tickoo, Lama Nachman, Rumi Chunara, Madhulika Srikumar, Adrian Weller, and Alice Xiang. 2021. Uncertainty as a Form of Transparency: Measuring, Communicating, and Using Uncertainty. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (AIES '21)*. ACM, Virtual, 401–413.
- [40] Laura Bieker, Daniel Krajzewicz, Antonio Pio Morra, Carlo Michelacci, and Fabio Carotolano. 2014. Traffic Simulation for All: A Real World Traffic Scenario from the City of Bologna. In *Proceedings of the 2014 SUMO Conference (SUMO '14)*. Springer, Berlin, Germany, 47–60.
- [41] Kostas Bimpikis, Ozan Candogan, and Daniela Saban. 2019. Spatial Pricing in Ride-Sharing Networks. *Operations Research* 67, 3 (2019), 744–769.
- [42] Hendrik Blockeel, Laurens Devos, Benoît Frénay, Géraldin Nanfack, and Siegfried Nijssen. 2023. Decision trees: from efficient prediction to responsible AI. *Frontiers in Artificial Intelligence* 6 (2023), 1124553.
- [43] National Transportation Safety Board. 2024. *Collapse of the Fern Hollow Bridge*. Technical Report HIR-24-02. National Transportation Safety Board. 116 pages.
- [44] Klavdiya Bochenina, Anton Taleiko, and Laura Ruotsalainen. 2023. Simulation-Based Origin-Destination Matrix Reduction: A Case Study of Helsinki City Area. In *Proceedings of the 2023 SUMO User Conference (SUMO '23)*. SUMO, Berlin, Germany, 1–7.
- [45] Rohit Bokade, Xiaoning Jin, and Christopher Amato. 2023. Multi-Agent Reinforcement Learning Based on Representational Communication for Large-Scale Traffic Signal Control. *IEEE Access* 11 (2023), 47646–47658.
- [46] Elizabeth Bondi, Lily Xu, Diana Acosta-Navas, and Jackson A. Killian. 2021. Envisioning Communities: A Participatory Approach Towards AI for Social Good. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (AIES '21)*. ACM, Virtual, 425–436.
- [47] P Bonhard and M A Sasse. 2006. “Knowing me, knowing you” — using profiles and social networking to improve recommender systems. *BT Technology Journal* 24, 3 (2006), 84–98.
- [48] Leo Breiman, Jerome Friedman, R.A. Olshen, and Charles J. Stone. 1984. Splitting Rules.

In *Classification and Regression Trees*. Taylor & Francis, New York.

- [49] W. Brilon and T. Wietholt. 2013. Experiences with Adaptive Signal Control in Germany. *Transportation Research Record* 2356 (2013), 9–16.
- [50] Noam Brown and Tuomas Sandholm. 2017. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science* 359, 6374 (2017), 418–424.
- [51] Antoni Buades, Bartomeu Coll, and Jean-Michel Morel. 2011. Non-Local Means Denoising. *Image Processing On Line* 1 (2011), 208–212.
- [52] Eleanor R. Burgess, Ivana Jankovic, Melissa Austin, Nancy Cai, Adela Kapuścińska, Suzanne T. Currie, J. Marc Overhage, Erika S. Poole, and Jofsh Kaye. 2023. Healthcare AI Treatment Decision Support: Design Principles to Enhance Clinician Adoption and Trust. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–19.
- [53] Anita Burgun, Virginie Do, Aïda Elamrani, Thierry Kirat, Jérôme Lang, and Thierry Poibeau. 2025. Challenges in AI-Assisted Decision-Making: Insights from the Medical and the Legal domains. *arXiv:2506.2847*
- [54] Andrew Butler. 2023. SCQPTH: an efficient differentiable splitting method for convex quadratic programming. *arXiv:2308.08232*
- [55] Carrie J. Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. “Hello AI”: Uncovering the Onboarding Needs of Medical Practitioners for Human-AI Collaborative Decision-Making. *Proceedings of the ACM on Human-Computer Interaction* 3 (2019), 1–24.
- [56] Carrie J. Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2021. Onboarding Materials as Cross-functional Boundary Objects for Developing AI Assistants. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems (CHI EA '21)*. ACM, Yokohama, Japan, 1–7.
- [57] Sibin Cai, Jie Fang, and Mengyun Xu. 2025. XLight: An interpretable multi-agent reinforcement learning approach for traffic signal control. *Expert Systems with Applications* 273 (2025), 126938.
- [58] Wanling Cai, Yucheng Jin, and Li Chen. 2022. Impacts of Personal Characteristics on User Trust in Conversational Recommender Systems. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. ACM, New Orleans, USA, 1–14.
- [59] José António Capela Dias, Penousal Machado, and Francisco Câmara Pereira. 2013. Simulating the Impact of Drivers’ Personality on City Transit. In *Proceedings of the 13th World Conference on Transport Research (WCTR '13)*. World Conference on Transport Research Society, Leeds, UK, 1–13.
- [60] Federico Maria Cau, Hanna Hauptmann, Lucio Davide Spano, and Nava Tintarev. 2023. Effects of AI and Logic-Style Explanations on Users’ Decisions under Different Levels of Uncertainty. *ACM Transactions on Interactive Intelligent Systems* 13, 4 (2023), 1–42.
- [61] Ngai Keung Chan and Lee Humphreys. 2018. Mediatization of Social Space and the Case

- of Uber Drivers. *Media and Communication* 6, 2 (2018), 29–38.
- [62] Cheng Chang, Siqi Wang, Jiawei Zhang, Jingwei Ge, and Li Li. 2024. LLMScenario: Large Language Model Driven Scenario Generation. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 54, 11 (2024), 6581–6594.
  - [63] Harshal A. Chaudhari, John W. Byers, and Evimaria Terzi. 2018. Putting Data in the Driver’s Seat: Optimizing Earnings for On-Demand Ride-Hailing. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining (WSDM ’22)*. ACM, Marina Del Rey, USA, 90–98.
  - [64] Harshal A. Chaudhari, John W. Byers, and Evimaria Terzi. 2020. Learn to Earn: Enabling Coordination Within a Ride-Hailing Fleet. In *Proceedings of the 2020 IEEE International Conference on Big Data (Big Data ’20)*. IEEE, Atlanta, USA, 1127–1136.
  - [65] Chacha Chen, Hua Wei, Nan Xu, Guanjie Zheng, Ming Yang, Yuanhao Xiong, Kai Xu, and Zhenhui Li. 2020. Toward A Thousand Lights: Decentralized Deep Reinforcement Learning for Large-Scale Traffic Signal Control. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI ’20)*. AAAI, New York, USA, 3414–3421.
  - [66] Long Chen, Piyushimita (Vonu) Thakuriah, and Konstantinos Ampountolas. 2021. Short-Term Prediction of Demand for Ride-Hailing Services: A Deep Learning Approach. *Journal of Big Data Analytics in Transportation* 3 (2021), 175–195.
  - [67] M. Keith Chen, Peter E. Rossi, Judith A. Chevalier, and Emily Oehlsen. 2019. The Value of Flexible Work: Evidence from Uber Drivers. *Journal of Political Economy* 127, 6 (2019), 2735–2794.
  - [68] Rex Chen, Kathleen M. Carley, Fei Fang, and Norman Sadeh. 2023. Purpose in the Machine: Do Traffic Simulators Produce Distributionally Equivalent Outcomes for Reinforcement Learning Applications?. In *Proceedings of the 2023 Winter Simulation Conference (WSC ’23)*. ACM, San Antonio, USA, 1842–1853.
  - [69] Rex Chen, Fei Fang, and Norman Sadeh. 2022. The Real Deal: A Review of Challenges and Opportunities in Moving Reinforcement Learning-Based Traffic Signal Control Systems Towards Reality. In *Proceedings of the 12th International Workshop on Agents in Traffic and Transportation (ATT ’22)*. CEUR, Vienna, Austria, 1–21.
  - [70] Rex Chen, Stephanie Milani, Zhicheng Zhang, Norman Sadeh, and Fei Fang. 2025. Making Teams and Influencing Agents: Efficiently Coordinating Decision Trees for Interpretable Multi-Agent Reinforcement Learning. In *Proceedings of the 2025 AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES ’25)*. ACM, Madrid, Spain, n. pag.
  - [71] Rex Chen, Ruiyi Wang, Fei Fang, and Norman Sadeh. 2025. Missing Pieces: How Do Designs that Expose Uncertainty Longitudinally Impact Trust in AI Decision Aids? An In Situ Study of Gig Drivers. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’25)*. ACM, Athens, Greece, 790–816.
  - [72] Rex Chen, Karen Wu, John McCartney, Norman Sadeh, and Fei Fang. 2025. Out of the Past: An AI-Enabled Pipeline for Traffic Simulation from Noisy, Multimodal Detector Data

- and Stakeholder Feedback. In *Proceedings of the 2025 Winter Simulation Conference (WSC '25)*. ACM, Seattle, USA, n. pag.
- [73] Xiwen Chen, Hao Wang, Abolfazl Razi, Brendan Russo, Jason Pacheco, and John Roberts. 2023. Network-Level Safety Metrics for Overall Traffic Safety Assessment: A Case Study. *IEEE Access* 11 (2023), 17755–17778.
  - [74] Richard Cheng, Gábor Orosz, Richard M. Murray, and Joel W. Burdick. 2019. End-to-End Safe Reinforcement Learning through Barrier Functions for Safety-Critical Continuous Control Tasks. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI '19)*. AAAI, Honolulu, USA, 3387–3395.
  - [75] Hyesun Choung, Prabu David, and Arun Ross. 2023. Trust in AI and Its Role in the Acceptance of AI Technologies. *International Journal of Human-Computer Interaction* 39, 9 (2023), 1727–1739.
  - [76] Peter Christensen and Adam Osman. 2023. *The Demand for Mobility: Evidence from an Experiment with Uber Riders*. Working Paper 31330. NBER. 48 pages.
  - [77] Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. 2020. Multi-Agent Deep Reinforcement Learning for Large-Scale Traffic Signal Control. *IEEE Transactions on Intelligent Transportation Systems* 21, 3 (2020), 1086–1095.
  - [78] Hyuck David Chung, Yue Maggie Zhou, and Christine Choi. 2022. When Uber Eats its Own Business, and That of its Competitors Too. *Academy of Management Proceedings* 2022, 1 (2022), 15263.
  - [79] Seung Youn Chyung, Megan Kennedy, and Ingrid Campbell. 2018. Evidence-Based Survey Design: The Use of Ascending or Descending Order of Likert-Type Response Options. *Performance Improvement* 57, 9 (2018), 9–16.
  - [80] Seung Youn Chyung, Katherine Roberts, Ieva Swanson, and Andrea Hankinson. 2017. Evidence-Based Survey Design: The Use of a Midpoint on the Likert Scale. *Performance Improvement* 56, 10 (2017), 15–23.
  - [81] Kris Van Cleave. 2025. *U.S. vehicle sales surge as buyers hope to beat tariff hikes*. CBS News. <https://www.cbsnews.com/news/u-s-new-vehicle-sales-surge-buyers-hope-to-beat-tariff-hikes/>
  - [82] Lara Codecà, Raphael Frank, and Thomas Engel. 2015. Luxembourg SUMO Traffic (LuST) Scenario: 24 hours of mobility for vehicular networking research. In *Proceedings of the 2015 Vehicular Networking Conference (VNC '15)*. IEEE, Kyoto, Japan, 1–8.
  - [83] Lara Codecà and Jérôme Härri. 2017. Towards multimodal mobility simulation of C-ITS: The Monaco SUMO traffic scenario. In *Proceedings of the 2017 Vehicular Networking Conference (VNC '17)*. IEEE, Turin, Italy, 97–100.
  - [84] Benjamin Coifman, David Beymer, Philip McLauchlan, and Jitendra Malik. 1998. A real-time computer vision system for vehicle tracking and traffic surveillance. *Transportation Research Part C: Emerging Technologies* 6, 4 (1998), 271–288.
  - [85] Christopher Court-Dobson and Jon Lawson. 2024. *Can AI control traffic*

- lights, signals and intersections?* Traffic Technology International. <https://www.traffichtechnologytoday.com/features/feature-should-ai-control-traffic-lights.html>
- [86] Marco Crespi, Andrea Ferigo, Leonardo Lucio Custode, and Giovanni Iacca. 2023. A population-based approach for multi-agent interpretable reinforcement learning. *Applied Soft Computing* 147 (2023), 110758.
  - [87] Jiaxun Cui, William Macke, Harel Yedidsion, Aastha Goyal, Daniel Urieli, and Peter Stone. 2021. Scalable Multiagent Driving Policies For Reducing Traffic Congestion. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. ACM, Virtual, 386–394.
  - [88] Colin Curtain. 2023. *QualCoder*. University of Tasmania. <https://github.com/ccbogel/QualCoder/releases/tag/3.3>
  - [89] Longchao Da, Chen Chu, Weinan Zhang, and Hua Wei. 2024. CityFlowER: An Efficient and Realistic Traffic Simulator with Embedded Machine Learning Models. In *Proceedings of the 2024 European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD '24)*. ACM, Vilnius, Lithuania, 368–373.
  - [90] Longchao Da, Minquan Gao, Hao Mei, and Hua Wei. 2024. Prompt to transfer: sim-to-real transfer for traffic signal control with prompt learning. In *Proceedings of the 38th AAAI Conference on Artificial Intelligence (AAAI '24)*. AAAI, Vancouver, Canada, 82–90.
  - [91] Simon Danner, Matthias Pfromm, and Klaus Bengler. 2020. Does Information on Automated Driving Functions and the Way of Presenting It before Activation Influence Users' Behavior and Perception of the System? *Information* 11, 1 (2020), 54.
  - [92] Christopher M. Day and Darcy M. Bullock. 2011. Computational Efficiency of Alternative Algorithms for Arterial Offset Optimization. *Transportation Research Record* 2259 (2011), 37–47.
  - [93] Christopher M. Day and Darcy M. Bullock. 2011. *Optimization of Offsets and Cycle Length Using High Resolution Signal Event Data*. Working Paper SPR-3409. Joint Transportation Research Program. 36 pages.
  - [94] Taylor de O. Antes, Ana L.C. Bazzan, and Anderson Rocha Tavares. 2022. Information upwards, recommendation downwards: reinforcement learning with hierarchy for traffic signal control. *Procedia Computer Science* 201 (2022), 24–31.
  - [95] Arjan de Ruijter, Oded Cats, and Hans van Lint. 2024. Ridesourcing platforms thrive on socio-economic inequality. *Scientific Reports* 14 (2024), 7371.
  - [96] Augustin Degas, Mir Riyanul Islam, Christophe Hurter, Shaibal Barua, Hamidur Rahman, Minesh Poudel, Daniele Ruscio, Mobyen Uddin Ahmed, Shahina Begum, Md Aquif Rahman, Stefano Bonelli, Giulia Cartocci, Gianluca Di Flumeri, Gianluca Borghini, Fabio Babiloni, and Pietro Aricó. 2022. A Survey on Artificial Intelligence (AI) and eXplainable AI in Air Traffic Management: Current Trends and Development with Future Research Trajec-



- tory. *Applied Sciences* 12, 3 (2022), 1295.
- [97] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve<sup>1</sup>, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark Yatskar, and Ali Farhadi. 2020. RoboTHOR: An Open Simulation-to-Real Embodied AI Platform. In *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR '20)*. CVPR, Seattle, USA, 3164–3174.
  - [98] Nicola Dell, Vidya Vaidyanathan, Indrani Medhi, Edward Cutrell, and William Thies. 2012. “Yours is better!”: participant response bias in HCI. In *Proceedings of the 2012 CHI Conference on Human Factors in Computing Systems (CHI '12)*. ACM, Austin, USA, 1321–1330.
  - [99] Delphine Delorme and Bongsob Song. 2001. *Human Driver Model for SmartAHS*. Technical Report UCB-ITS-PRR-2001-12. Institute of Transportation Studies, University of California, Berkeley, Berkeley, USA. 1–49 pages.
  - [100] Richard W. Denney, Larry Head, and Kevin Spencer. 2008. *Signal timing under saturated conditions*. Technical Report FHWA-HOP-09-008. U.S. Department of Transportation Federal Highway Administration. 76 pages.
  - [101] Yining Di, Meng Xu, Zheng Zhu, Hai Yang, and Xiqun Chen. 2022. Analysis of ride-sourcing drivers’ working Pattern(s) via spatiotemporal work slices: A case study in Hangzhou. *Transport Policy* 125 (2022), 336–351.
  - [102] Steven Diamond and Stephen Boyd. 2016. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research* 17, 83 (2016), 1–5.
  - [103] Django. 2023. *Django 4.1*. Django Software Foundation. <https://docs.djangoproject.com/en/4.1>
  - [104] Mark Dodgson. 1993. Learning, Trust, and Technological Collaboration. *Human Relations* 46, 1 (1993), 77–95.
  - [105] Alexander Domahidi, Eric Chu, and Stephen Boyd. 2013. ECOS: An SOCP solver for embedded systems. In *Proceedings of the 2013 European Control Conference (ECC '13)*. IEEE, Zurich, Switzerland, 3071–3076.
  - [106] Pedro M. d’Orey, Ricardo Fernandes, and Michel Ferreira. 2012. Empirical Evaluation of a Dynamic and Distributed Taxi-Sharing System. In *Proceedings of the 15th International IEEE Conference on Intelligent Transportation Systems (ITSC '12)*. IEEE, Anchorage, USA, 140–146.
  - [107] Finale Doshi-Velez and Been Kim. 2017. Towards A Rigorous Science of Interpretable Machine Learning. arXiv:1702.08608
  - [108] Richard Dowling, Alexander Skabardonis, and Vassili Alexiadis. 2004. *Traffic Analysis Toolbox, Volume III: Guidelines for Applying Traffic Microsimulation Modeling Software*. U.S. Department of Transportation Federal Highway Administration.
  - [109] Jeff Druce, Michael Harradon, and James Tittle. 2021. Explainable artificial intelligence (XAI) for increasing user trust in deep reinforcement learning driven autonomous systems. arXiv:2106.03775

- [110] Efi Dvir, Mark Shifrin, and Omer Gurewitz. 2024. Cooperative Multi-Agent Reinforcement Learning for Data Gathering in Energy-Harvesting Wireless Sensor Networks. *Mathematics* 12, 13 (2024), 2102.
- [111] Laura Eboli, Gabriella Mazzulla, and Giuseppe Pungillo. 2016. Combining speed and acceleration to define car users’ safe or unsafe driving behaviour. *Transportation Research Part C: Emerging Technologies* 68 (2016), 113–125.
- [112] Econolite. 2023. *Centracs Edaptive datasheet*. Document CNTRC-EDPTV 11.2023. Econolite Product Resource Library. 2 pages.
- [113] Samah El-Tantawy, Baher Abdulhai, and Hossam Abdelgawad. 2013. Review of Microscopic Lane-Changing Models and Future Research Opportunities. *IEEE Transactions on Intelligent Transportation Systems* 14, 4 (2013), 1942–1956.
- [114] Jennifer Elias. 2025. *Waymo to begin testing in Philadelphia with safety drivers behind the wheel*. CNBC. <https://www.cnbc.com/2025/07/07/waymo-to-begin-testing-in-philadelphia-with-drivers-behind-the-wheel.html>
- [115] Myungeun Eom and Byung-In Kim. 2020. The traffic signal control problem for intersections: a review. *European Transport Research Review* 12 (2020), 50.
- [116] Jakob Erdmann. 2014. Lane-Changing Model in SUMO. In *Proceedings of the 2014 SUMO User Conference (SUMO ’14)*. German Aerospace Center, Cologne, Germany, 77–88.
- [117] Mohamed Essa and Tarek Sayed. 2020. Self-learning adaptive traffic signal control for real-time safety optimization. *Accident Analysis & Prevention* 146 (2020), 105713.
- [118] Euthenics and TranSystems. 2023. *Preliminary Feasibility Study — Cuyahoga/Medina Traffic Study*. Technical Report PID 116069. City of Strongsville. 1–759 pages.
- [119] Karim Fadhloun, Hesham Rakha, Amara Loulizi, and Abdessattar Abdelkef. 2015. Vehicle Dynamics Model for Estimating Typical Vehicle Accelerations. *Transportation Research Record* 2491 (2015), 61–71.
- [120] Nahid Parvez Farazi, Bo Zou, Tanvir Ahamed, and Limon Barua. 2021. Deep reinforcement learning in transportation research: A review. *Transportation Research Interdisciplinary Perspectives* 11 (2021), 100425.
- [121] K. J. Kevin Feng, Tony W. Li, and Amy X. Zhang. 2023. Understanding Collaborative Practices and Tools of Professional UX Practitioners in Software publishers. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI ’23)*. ACM, Hamburg, Germany, 1–20.
- [122] Andrea Ferrario and Michele Loi. 2022. How Explainability Contributes to Trust in AI. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’22)*. ACM, Seoul, South Korea, 1457–1466.
- [123] Andres Ferraro, Dietmar Jannach, and Xavier Serra. 2020. Exploring Longitudinal Effects of Session-based Recommendations. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys ’20)*. ACM, Virtual, 474–479.
- [124] Figma, Inc. n.d.. *Figma*. Figma, Inc. <https://figma.com>

- [125] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual multi-agent policy gradients. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI '18)*. AAAI, New Orleans, USA, 2974–2982.
- [126] Deepeka Garg, Maria Chli, and George Vogiatzis. 2018. Deep Reinforcement Learning for Autonomous Traffic Light Control. In *Proceedings of the 2018 3rd International Conference on Intelligent Transportation Engineering (ICITE '18)*. IEEE, Singapore, 214–218.
- [127] Deepeka Garg, Maria Chli, and George Vogiatzis. 2019. A Deep Reinforcement Learning Agent for Traffic Intersection Control Optimization. In *Proceedings of the 2019 IEEE Intelligent Transportation Systems Conference (ITSC '19)*. IEEE, Auckland, NZ, 4222–4229.
- [128] Deepeka Garg, Maria Chli, and George Vogiatzis. 2019. Traffic3D: A Rich 3D-Traffic Environment to Train Intelligent Agents. In *Proceedings of the 19th International Conference on Computational Science (ICCS '19)*. Springer, New York, USA, 749–755.
- [129] Deepeka Garg, Maria Chli, and George Vogiatzis. 2020. Multi-Agent Deep Reinforcement Learning for Traffic optimization through Multiple Road Intersections using Live Camera Feed. In *Proceedings of the 2020 IEEE Intelligent Transportation Systems Conference (ITSC '20)*. IEEE, Rhodes, Greece, 1–8.
- [130] Deepeka Garg, Maria Chli, and George Vogiatzis. 2022. Fully-Autonomous, Vision-based Traffic Signal Control: from Simulation to Reality. In *Proceedings of the 21th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '22)*. IFAA-MAS, Auckland, NZ, 454–462.
- [131] Nikhil Garg and Hamid Nazerzadeh. 2022. Driver Surge Pricing. *Management Science* 68, 5 (2022), 3219–3235.
- [132] Timothy Geary and David Danks. 2019. Balancing the Benefits of Autonomous Vehicles. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (AIES '19)*. ACM, Honolulu, USA, 1–6.
- [133] Steven R. Gehrke. 2020. Uber service area expansion in three major American cities. *Journal of Transport Geography* 86 (2020), 102752.
- [134] Wade Genders and Saiedeh Razavi. 2018. Evaluating reinforcement learning state representations for adaptive traffic signal control. In *Proceedings of the 9th International Conference on Ambient Systems, Networks and Technologies (ANT '18)*. Procedia Computer Science, Porto, Portugal, 26–33.
- [135] Douglas Gettman, Lili Pu, Tarek Sayed, and Steve Shelby. 2008. *Surrogate Safety Assessment Model and Validation: Final Report*. Technical Report FHWA-HRT-08-051. U.S. Department of Transportation Federal Highway Administration. 322 pages.
- [136] David Gibson, Milton K. (Pete) Mills, and Doug Rekenthaler Jr. 1998. Staying in The Loop: The Search for Improved Reliability of Traffic Sensing Systems Through Smart Test Instruments. *Public Roads* 62, 2 (1998), 47–51.
- [137] Thomas Krendl Gilbert, Nathan Lambert, Sarah Dean, Tom Zick, Aaron Snoswell, and So-

- ham Mehta. 2023. Reward Reports for Reinforcement Learning. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society (AIES '23)*. ACM, Montréal, Canada, 84–130.
- [138] Nicole Gillespie, Steve Lockey, Alexandria Macdade, Tabi Ward, and Gerard Hassed. 2025. *Trust, attitudes and use of artificial intelligence: A global study 2025*. Technical Report. The University of Melbourne and KPMG. 115 pages.
  - [139] Yang Miang Goh and Peter E.D. Love. 2012. Methodological application of system dynamics for evaluating traffic safety policy. *Safety Science* 50, 7 (2012), 1594–1605.
  - [140] Yaobang Gong, Mohamed Abdel-Aty, Jinghui Yuan, and Qing Cai. 2020. Multi-Objective reinforcement learning approach for improving safety at intersections with adaptive traffic signal control. *Accident Analysis & Prevention* 144 (2020), 105655.
  - [141] Robert L. Gordon and Warren Tighe. 2005. *Traffic Control Systems Handbook*. U.S. Department of Transportation Federal Highway Administration.
  - [142] Lance R. Grenzeback, William R. Reilly, Paul O. Roberts, and Joseph R. Stowers. 1990. Urban Freeway Gridlock Study: Decreasing the Effects of Large Trucks on Peak-Period Urban Freeway Congestion. *Transportation Research Record* 1256 (1990), 16–26.
  - [143] Ulrike Gretzel and Daniel R. Fesenmaier. 2006. Persuasion in Recommender Systems. *International Journal of Electronic Commerce* 11, 2 (2006), 81–100.
  - [144] Kathleen Griesbach, Adam Reich, Luke Elliott-Negri, and Ruth Milkman. 2019. Algorithmic Control in Platform Food Delivery Work. *Socius* 5 (2019), 1–15.
  - [145] Hankang Gu, Shangbo Wang, Xiaoguang Ma, Dongyao Jia, Guoqiang Mao, and Eng Gee Lim. 2024. Large-Scale Traffic Signal Control Using Constrained Network Partition and Adaptive Deep Reinforcement Learning. *IEEE Transactions on Intelligent Transportation Systems* 25, 7 (2024), 7619–7632.
  - [146] Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, and Jun Wang. 2024. A Review of Safe Reinforcement Learning: Methods, Theories, and Applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46, 12 (2024), 11216–11235.
  - [147] Yin Gu, Kai Zhang, Qi Liu, Weibo Gao, Longfei Li, and Jun Zhou. 2024.  $\pi$ -Light: Programmatic Interpretable Reinforcement Learning for Resource-Limited Traffic Signal Control. In *Proceedings of the 38th AAAI Conference on Artificial Intelligence (AAAI '24)*. AAAI, Vancouver, Canada, 21107–21115.
  - [148] Harish Guda and Upender Subramanian. 2019. Your Uber Is Arriving: Managing On-Demand Workers Through Surge Pricing, Forecast Communication, and Worker Incentives. *Management Science* 65, 5 (2019), 1995–2014.
  - [149] Jiaying Guo, Saeedeh Ghanadbashi, Shen Wang, and Fatemeh Golpayegani. 2023. Urban Traffic Signal Control at the Edge: An Ontology-Enhanced Deep Reinforcement Learning Approach. In *Proceedings of the 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC '23)*. IEEE, Bilbao, Spain, 6027–6033.
  - [150] Suiming Guo, Chao Chen, Jingyuan Wang, Yaxiao Liu, Ke Xu, and Zhiwen Yu. 2019. ROD-

- Revenue: Seeking Strategies Analysis and Revenue Prediction in Ride-on-Demand Service Using Multi-Source Urban Data. *IEEE Transactions on Mobile Computing* 19, 9 (2019), 2202–2220.
- [151] Suiming Guo, Chao Chen, Jingyuan Wang, Yaxiao Liu, Ke Xu, Daqing Zhang, and Dah Ming Chiu. 2018. A Simple but Quantifiable Approach to Dynamic Price Prediction in Ride-on-demand Services Leveraging Multi-source Urban Data. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2, 3 (2018), 1–24.
  - [152] Ziyang Guo, Yifan Wu, Jason D. Hartline, and Jessica Hullman. 2024. A Decision Theoretic Framework for Measuring AI Reliance. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT '24)*. ACM, Rio de Janeiro, Brazil, 221–236.
  - [153] Gurobi Optimization, LLC. 2025. *Gurobi Optimizer Reference Manual*. Gurobi Optimization, LLC. <https://www.gurobi.com>
  - [154] Sehoon Ha, Peng Xu, Zhenyu Tan, Sergey Levine, and Jie Tan. 2021. Learning to Walk in the Real World with Minimal Human Effort. In *Proceedings of the 4th Annual Conference on Robot Learning (CoRL '21)*. CoRL, London, UK, 1110–1120.
  - [155] Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. 2019. Learning to Walk via Deep Reinforcement Learning. In *Proceedings of the 15th Conference on Robotics: Science and Systems (RSS '19)*. RSS, Freiburg im Breisgau, Germany, 11.
  - [156] Jonathan V. Hall and Alan B. Krueger. 2018. An Analysis of the Labor Market for Uber's Driver-Partners in the United States. *ILR Review* 71, 3 (2018), 705–732.
  - [157] Xu Han, Qiannan Yang, Xianda Chen, Zhenghan Cai, Xiaowen Chu, and Meixin Zhu. 2024. AutoReward: Closed-Loop Reward Design with Large Language Models for Autonomous Driving. *IEEE Transactions on Intelligent Vehicles* Early Access (2024), 1–13.
  - [158] Xiao Han, Xiangyu Zhao, Liang Zhang, and Wanyu Wang. 2023. Mitigating Action Hysteresis in Traffic Signal Control with Traffic Predictive Reinforcement Learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23)*. ACM, Long Beach, USA, 673–684.
  - [159] Chengzheng Hang, Zhenfei Liu, Yujing Wang, Caiyi Hu, Yuelong Su, and Zhenning Dong. 2019. Sharing diseconomy: impact of the subsidy war of ride-sharing companies on urban congestion. *International Journal of Logistics Research and Applications* 22, 5 (2019), 491–500.
  - [160] Benjamin V. Hanrahan, Ning F. Ma, and Chien Wen Yuan. 2017. The Roots of Bias on Uber. In *Proceedings of the 15th European Conference on Computer-Supported Cooperative Work (ECSCW '17)*. EUSSET, Sheffield, UK, 1–17.
  - [161] Brian Yueshuai He, Qinhua Jiang, Haoxuan Ma, and Jiaqi Ma. 2024. Multi-Agent Multimodal Transportation Simulation for Mega-cities: Application of Los Angeles. *Procedia Computer Science* 238 (2024), 736–741.

- [162] Suining He and Kang G. Shin. 2019. Spatio-Temporal Capsule-based Reinforcement Learning for Mobility-on-Demand Network Coordination. In *Proceedings of the 2019 World Wide Web Conference (WWW '19)*. ACM, New York, USA, 2806–2813.
- [163] Alexandre Heuillet, Fabien Couthouis, and Natalia Díaz-Rodríguez. 2022. Collective eXplainable AI: Explaining Cooperative Strategies and Agent Contribution in Multiagent Reinforcement Learning With Shapley Values. *IEEE Computational Intelligence Magazine* 17, 1 (2022), 59–71.
- [164] Corey Hill, Lily Elefteriadou, and Alexandra Kondyli. 2015. Exploratory Analysis of Lane Changing on Freeways Based on Driver Behavior. *Journal of Transportation Engineering* 141, 4 (2015), 1–11.
- [165] Gayang Ho, Clémence Morlet, Teik Soon Looi, Evan Gwee, Joel Teo, Peisi Keg, and Alok Jain. 2018. *Artificial Intelligence in Mass Public Transit*. Technical Report. UITP Asia-Pacific Centre for Transport Excellence. 13 pages.
- [166] Wassily Hoeffding. 1963. Probability Inequalities for Sums of Bounded Random Variables. *J. Amer. Statist. Assoc.* 58 (1963), 13–30.
- [167] Hannah Horner, Jennifer Pazour, and John E. Mitchell. 2021. Optimizing driver menus under stochastic selection behavior for ridesharing and crowdsourced delivery. *Transportation Research Part E: Logistics and Transportation Review* 153 (2021), 102419.
- [168] Moinul Hossain, Mohamed Abdel-Aty, Mohammed A. Quddus, Yasunori Muromachi, and Soumik Nafis Sadeek. 2019. Real-time crash prediction models: State-of-the-art, design pathways and ubiquitous requirements. *Accident Analysis & Prevention* 124 (2019), 66–84.
- [169] Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and Changjie Fan. 2020. Learning to Utilize Shaping Rewards: A New Approach of Reward Shaping. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS '20)*. NeurIPS, Vancouver, Canada, 15931–15941.
- [170] Shengyi Huang and Santiago Onta nón. 2022. A Closer Look at Invalid Action Masking in Policy Gradient Algorithms. arXiv:2006.14171
- [171] Nathan Hunt, Nathan Fulton, Sara Magliacane, Trong Nghia Hoang, Subhro Das, and Armando Solar-Lezama. 2021. Verifiably safe exploration for end-to-end reinforcement learning. In *Proceedings of the 24th International Conference on Hybrid Systems: Computation and Control (HSCC '21)*. ACM, Virtual, 1–11.
- [172] Hyeonjun Hwang, Clifford Winston, and Jia Yan. 2020. *Measuring the Benefits of Ridesharing Services to Urban Travelers: The Case of The San Francisco Bay Area*. Working Paper 70. Hutchins Center. 18 pages.
- [173] Shariq Iqbal and Fei Sha. 2019. Actor-Attention-Critic for Multi-Agent Reinforcement Learning. In *Proceedings of the 36th International Conference on Machine Learning (ICML '19)*. JMLR, Long Beach, USA, 2961–2970.
- [174] Lakshmi Shankar Iyer. 2021. AI enabled applications towards intelligent transportation.

*Transportation Engineering* 5 (2021), 100083.

- [175] Maia Jacobs, Jeffrey He, Melanie F. Pradier, Barbara Lam, Andrew C. Ahn, Thomas H. McCoy, Roy H. Perlis, Finale Doshi-Velez, and Krzysztof Z. Gajos. 2021. Designing AI for Trust and Collaboration in Time-Constrained Medical Decisions: A Sociotechnical Lens. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*. ACM, Yokohama, Japan, 1–14.
- [176] Alon Jacovi, Ana Marasović, Tim Miller, and Yoav Goldberg. 2021. Formalizing Trust in Artificial Intelligence: Prerequisites, Causes and Goals of Human Trust in AI. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. ACM, Chicago, USA, 624–635.
- [177] Vindula Jayawardana, Anna Landler, and Cathy Wu. 2021. Mixed Autonomous Supervision in Traffic Signal Control. In *Proceedings of the 2021 IEEE 24th International Conference on Intelligent Transportation Systems (ITSC '21)*. IEEE, Indianapolis, USA, 1767–1773.
- [178] Bokai Ji, Guangxia Li, and Gang Xiao. 2023. Enhancing the Interpretability of Deep Multi-agent Reinforcement Learning via Neural Logic Reasoning. In *ICANN '23*. Springer, Crete, Greece, 199–210.
- [179] Jiun-Yin Jian, Ann M. Bisantz, and Colin G. Drury. 2000. Foundations for an Empirically Determined Scale of Trust in Automated Systems. *International Journal of Cognitive Ergonomics* 4, 1 (2000), 53–71.
- [180] Qize Jiang, Minhao Qin, Hanyuan Zhang, Xinyu Zhang, and Weiwei Sun. 2024. BlindLight: High Robustness Reinforcement Learning Method to Solve Partially Blinded Traffic Signal Control Problem. *IEEE Transactions on Intelligent Transportation Systems* 25, 11 (2024), 16625–16641.
- [181] Shan Jiang, Le Chen, Alan Mislove, and Christo Wilson. 2018. On Ridesharing Competition and Accessibility: Evidence from Uber, Lyft, and Taxi. In *Proceedings of the 2018 World Wide Web Conference (WWW '18)*. ACM, Lyon, France, 863–872.
- [182] Xia Jiang, Jian Zhang, and Bo Wang. 2022. Energy-efficient driving for adaptive traffic signal control environment via explainable reinforcement learning. *Applied Sciences* 12, 11 (2022), 5380.
- [183] Junchen Jin and Xiaoliang Ma. 2019. A Multi-Objective Agent-Based Control Approach With Application in Intelligent Traffic Signal System. *IEEE Transactions on Intelligent Transportation Systems* 20, 10 (2019), 3900–3912.
- [184] Jiarui Jin, Ming Zhou, Weinan Zhang, Minne Li, Zilong Guo, Zhiwei Qin, Yan Jiao, Xiaocheng Tang, Chenxi Wang, Jun Wang, Guobin Wu, and Jieping Ye. 2019. CoRide: Joint Order Dispatching and Fleet Management for Multi-Scale Ride-Hailing Platforms. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM '19)*. ACM, Beijing, China, 1983–1992.
- [185] Glenn Jocher and Jing Qiu. 2024. *YOLO11*. Technical Report. Ultralytics. v11.0.0, <https://github.com/ultralytics/ultralytics>.

- [186] Min-Wook Kang, Moynur Rahman, and Joyoung Lee. 2020. Determination and Utilization of Dilemma Zone Length and Location for Safety Assessment of Rural High-Speed Signalized Intersections. *Transportation Research Record* 2674, 4 (2020), 272–280.
- [187] Amir Hossein Karbasi, Hao Yang, and Saiedeh Razavi. 2024. Exploring the impact of traffic signal control and connected and automated vehicles on intersections safety: A deep reinforcement learning approach. In *Proceedings of the 103rd Annual Meeting of the Transportation Research Board (TRB '24)*. TRB, Washington, DC, USA, 1–19.
- [188] George Karypis and Vipin Kumar. 1998. A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. *SIAM Journal on Scientific Computing* 20, 1 (1998), 359–392.
- [189] Dmitry Kazhdan, Zohreh Shams, and Pietro Lio. 2020. MARLeME: A Multi-Agent Reinforcement Learning Model Extraction Library. In *Proceedings of the 2020 International Joint Conference on Neural Networks (IJCNN '20)*. IEEE, Glasgow, UK, 1–8.
- [190] Hassan Ali Khan, Muhammad Shahzad, Hassan Iqbal, and Guoliang Jin. 2022. RMS: Removing Barriers to Analyze the Availability and Surge Pricing of Ridesharing Services. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. ACM, New Orleans, USA, 1–18.
- [191] Zulqarnain H. Khattak, Michael D. Fontaine, and Richard A. Boateng. 2018. Evaluating the impact of adaptive signal control technology on driver stress and behavior using real-world experimental data. *Transportation Research Part F: Traffic Psychology and Behaviour* 58 (2018), 133–144.
- [192] Jungkeun Kim, Marilyn Giroux, and Jacob C. Lee. 2021. When do you trust AI? The effect of number presentation detail on consumer trust and acceptance of AI recommendations. *Psychology & Marketing* 38 (2021), 1140–1155.
- [193] Sunnie S. Y. Kim, Q. Vera Liao, Mihaela Vorvoreanu, Stephanie Ballard, and Jennifer Wortman Vaughan. 2024. “I’m Not Sure, But...”: Examining the Impact of Large Language Models’ Uncertainty Expression on User Reliance and Trust. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT '24)*. ACM, Rio de Janeiro, Brazil, 822–835.
- [194] Sunnie S. Y. Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, and Andrés Monroy-Hernández. 2023. Humans, AI, and Context: Understanding End-Users’ Trust in a Real-World Computer Vision Application. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23)*. ACM, Chicago, USA, 77–88.
- [195] B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Yogamani, and Patrick Pérez. 2021. Deep Reinforcement Learning for Autonomous Driving: A Survey. *IEEE Transactions on Intelligent Transportation Systems* 23, 6 (2021), 4909–4926.
- [196] Lawrence A. Klein, Milton K. Mills, and David R.P. Gibson. 2006. *Traffic Detector Handbook: Third Edition — Volume II*. U.S. Department of Transportation Federal Highway Administration.



- [197] Bart P. Knijnenburg, Niels J.M. Reijmer, and Martijn C. Willemsen. 2011. Each to His Own: How Different Users Call for Different Interaction Methods in Recommender Systems. In *Proceedings of the 5th ACM Conference on Recommender Systems (RecSys '11)*. ACM, Chicago, USA, 3–13.
- [198] Amy J. Ko, Thomas D. LaToza, and Margaret M. Burnett. 2015. A practical guide to controlled experiments of software engineering tools with human participants. *Empirical Software Engineering* 20 (2015), 110–141.
- [199] Rafal Kocielnik, Saleema Amershi, and Paul N. Bennett. 2019. Will You Accept an Imperfect AI? Exploring Designs for Adjusting End-user Expectations of AI Systems. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. ACM, Glasgow, UK, 1–14.
- [200] Spencer C. Kohn, Ewart J. de Visser, Eva Wiese, Yi-Ching Lee, and Tyler H. Shaw. 2021. Measurement of Trust in Automation: A Narrative Review and Reference Guide. *Frontiers in Psychology* 12 (2021), 604977.
- [201] Sherrie Y. X. Komiak and Izak Benbasat. 2006. The Effects of Personalization and Familiarity on Trust and Adoption of Recommendation Agents. *MIS Quarterly* 30, 4 (2006), 941–960.
- [202] Behrad Koohy, Sebastian Stein, Enrico Gerding, and Ghaithaa Manla. 2022. Reward Function Design in Multi-Agent Reinforcement Learning for Traffic Signal Control. In *Proceedings of the 12th International Workshop on Agents in Traffic and Transportation (ATT '22)*. International Joint Conference on Artificial Intelligence, Vienna, Austria, 1–13.
- [203] Peter Koonce, Lee Rodegerdts, Kevin Lee, Shaun Quayle, Scott Beaird, Cade Braud, Jim Bonneson, Phil Tarnoff, and Tom Urbanik. 2008. *Traffic Signal Timing Manual*. U.S. Department of Transportation Federal Highway Administration.
- [204] Moritz Körber, Christian Gold, David Lechner, and Klaus Bengler. 2016. The influence of age on the take-over of vehicle control in highly automated driving. *Transportation Research Part F: Traffic Psychology and Behaviour* 39 (2016), 19–32.
- [205] Daniel Krajzewicz, Georg Hertkorn, C. Rössel, and Peter Wagner. 2002. SUMO (Simulation of Urban MObility) - an open-source traffic simulation. In *Proceedings of the 4th Middle East Symposium on Simulation and Modelling (MESM '02)*. SCS Europe, Dubai, UAE, 183–187.
- [206] Hanna Krasowski, Jakob Thumm, Marlon Müller, Lukas Schäfer, Xiao Wang, and Matthias Althoff. 2023. Provably Safe Reinforcement Learning: Conceptual Analysis, Survey, and Benchmarking. *Transactions on Machine Learning Research* November 2023 (2023), 1–38.
- [207] Stefan Krauß. 1998. *Microscopic Modeling of Traffic Flow: Investigation of Collision Free Vehicle Dynamics*. Ph. D. Dissertation. German Aerospace Center, Cologne, Germany.
- [208] Blaine Kuehnert, Rachel Kim, Jodi Forlizzi, and Hoda Heidari. 2025. The “Who”, “What”, and “How” of Responsible AI Governance: A Systematic Review and Meta-Analysis of (Actor, Stage)-Specific Tools. In *Proceedings of the 2025 ACM Conference on Fairness*,

- Accountability, and Transparency (FAccT '25)*. ACM, Athens, Greece, 2991–3005.
- [209] Ashwin Kumar, Yevgeniy Vorobeychik, and William Yeoh. 2023. Using Simple Incentives to Improve Two-Sided Fairness in Ridesharing Systems. In *Proceedings of the 33rd International Conference on Automated Planning and Scheduling (ICAPS '23)*. AAAI, Prague, Czech Republic, 227–235.
  - [210] Neetesh Kumar, Sarthak Mittal, Vaibhav Garg, and Neeraj Kumar. 2022. Deep Reinforcement Learning-Based Traffic Light Scheduling Framework for SDN-Enabled Smart Transportation System. *IEEE Transactions on Intelligent Transportation Systems* 23, 3 (2022), 2411–2421.
  - [211] Johannes Kunkel, Tim Donkers, Lisa Michael, Catalin-Mihai Barbu, and Jürgen Ziegler. 2019. Let Me Explain: Impact of Personal and Impersonal Explanations on Trust in Recommender Systems. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. ACM, Scotland, UK, 1–12.
  - [212] Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest Package: Tests in Linear Mixed Effects Models. *Journal of Statistical Software* 82, 13 (2017), 1–26.
  - [213] Min-Ah Kwak, Theo Arentze, Erik de Romph, and Soora Rasouli. 2012. Activity-based dynamic traffic modeling: Influence of population sampling fraction size on simulation error. In *Proceedings of the 13th International Conference on Travel Behaviour Research (IATBR '12)*. IATBR, Toronto, Canada, 1–17.
  - [214] Neal Lathia, Stephen Hailes, Licia Capra, and Xavier Amatriain. 2010. Temporal Diversity in Recommender Systems. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '10)*. ACM, Geneva, Switzerland, 210–217.
  - [215] Ho Lee and Benjamin Coifman. 2011. Quantifying Loop Detector Sensitivity and Correcting Detection Problems on Freeways. *Journal of Transportation Engineering* 138, 7 (2011), 871–881.
  - [216] Ho Lee and Benjamin Coifman. 2012. Identifying chronic splashover errors at freeway loop detectors. *Transportation Research Part C: Emerging Technologies* 24 (2012), 141–156.
  - [217] John Lee and Neville Moray. 1992. Trust, control strategies and allocation of function in human-machine systems. *Ergonomics* 10 (1992), 1243–1270.
  - [218] Min Kyung Lee, Danyel Kusbit, Evan Metsky, and Laura Dabbish. 2015. Working with Machines: The Impact of Algorithmic and Data-Driven Management on Human Workers. In *Proceedings of the 2015 CHI Conference on Human Factors in Computing Systems (CHI '15)*. ACM, Seoul, South Korea, 1–13.
  - [219] Lauren Leffer. 2024. *Can Google Make Stoplights Smarter?* Scientific American. <https://www.scientificamerican.com/article/googles-project-green-light-uses-ai-to-take-on-city-traffic/>
  - [220] Yanzhe (Murray) Lei, Stefanus Jasin, Jingyi Wang, Houtao Deng, and Jagannath Putrevu.

2020. Dynamic Workforce Acquisition for Crowdsourced Last-Mile Delivery Platforms. arXiv:3532844
- [221] Catur Yudo Leksono and Tina Andriyana. 2012. *Roundabout Microsimulation using SUMO: A Case Study in Idrottsparken Roundabout, Norrköping, Sweden*. Master’s thesis. Linköping University, Linköping.
  - [222] Jonas F. Leon, Francesca Giancola, Andrea Boccolucci, and Mattia Neroni. 2023. A Demand Modelling Pipeline for an Agent-Based Traffic Simulation of the City of Barcelona. In *Proceedings of the 2023 Winter Simulation Conference (WSC ’23)*. ACM, San Antonio, USA, 1777–1782.
  - [223] Hung-Xin Li. 2025. Exploration of Wargaming and AI Applications in Military Decision-Making. In *Proceedings of the 2025 International Conference on Military Technologies (ICMT ’25)*. IEEE, Brno, Czech Republic, 1–5.
  - [224] Shuyang Li, Talha Azfar, and Ruimin Ke. 2024. ChatSUMO: Large Language Model for Automating Traffic Scenario Generation in Simulation of Urban MObility. *IEEE Transactions on Intelligent Vehicles* Early Access (2024), 1–12.
  - [225] Toby Jia-Jun Li, Yuwen Lu, Jaylexia Clark, Meng Chen, Victor Cox, Meng Jiang, Yang Yang, Tamara Kay, Danielle Wood, and Jay Brockman. 2022. A Bottom-Up End-User Intelligent Assistant Approach to Empower Gig Workers against AI Inequality. In *Proceedings of the 2022 Symposium on Human-Computer Interaction for Work (CHIWORK ’22)*. ACM, Durham, USA, 1–10.
  - [226] Yuxi Li. 2018. Deep Reinforcement Learning. arXiv:1810.06339
  - [227] Yongfu Li and Dihua Sun. 2012. Microscopic car-following model for the traffic flow: the state of the art. *Journal of Control Theory and Applications* 10 (2012), 133–143.
  - [228] Yichen Li and Chicheng Zhang. 2022. On Efficient Online Imitation Learning via Classification. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS ’22)*. NeurIPS, New Orleans, USA, 32383–32397.
  - [229] Ziru Li, Chen Liang, Yili Hong, and Zhongju Zhang. 2022. How Do On-demand Ridesharing Services Affect Traffic Congestion? The Moderating Role of Urban Compactness. *Production and Operations Management* 31, 1 (2022), 239–258.
  - [230] Zhenning Li, Hao Yu, Guohui Zhang, Shangjia Dong, and Cheng-Zhong Xu. 2021. Network-wide traffic signal control optimization using a multi-agent deep reinforcement learning. *Transportation Research Part C: Emerging Technologies* 125 (2021), 103059.
  - [231] Yu Liang and Martijn C. Willemsen. 2022. Exploring the longitudinal effects of nudging on users’ music genre exploration behavior and listening preferences. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys ’22)*. ACM, Seattle, USA, 3–13.
  - [232] Lyuchao Liao, Jierui Liu, Xinke Wu, Fumin Zou, Jengshyang Pan, Qi Sun, Shengbo Eben Li, and Maolin Zhang. 2020. Time Difference Penalized Traffic Signal Timing by LSTM Q-Network to Balance Safety and Capacity at Intersections. *IEEE Access* 8 (2020), 80086–80096.

- [233] Mengqi Liao and S. Shyam Sundar. 2022. When E-Commerce Personalization Systems Show and Tell: Investigating the Relative Persuasive Appeal of Content-Based versus Collaborative Filtering. *Journal of Advertising* 51 (2022), 256–267.
- [234] Ulrich Lichtenthaler. 2018. Substitute or Synthesis: The Interplay between Human and Artificial Intelligence. *Research-Technology Management* 61, 5 (2018), 12–14.
- [235] In-Kyu Lim and Young-Jun Kweon. 2013. Identifying High-Crash-Risk Intersections: Comparison of Traditional Methods with the Empirical Bayes–Safety Performance Function Method. *Transportation Research Record* 2364, 1 (2013), 44–50.
- [236] Todd Litman. 2023. *Safer Than You Think! Revising the Transit Safety Narrative*. Technical Report. Victoria Transport Policy Institute. 32 pages.
- [237] Todd Litman. 2025. *Autonomous Vehicle Implementation Predictions: Implications for Transport Planning*. Technical Report. Victoria Transport Policy Institute. 49 pages.
- [238] Shan Liu and Hai Jiang. 2022. Personalized route recommendation for ride-hailing with deep inverse reinforcement learning and real-time traffic conditions. *Transportation Research Part E: Logistics and Transportation Review* 164 (2022), 102780.
- [239] Ying Liu, Lei Liu, and Wei-Peng Chen. 2017. Intelligent Traffic Light Control Using Distributed Multi-agent Q Learning. In *Proceedings of the 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC '17)*. IEEE, Yokohama, Japan, 1–8.
- [240] Zichuan Liu, Yuanyang Zhu, and Chunlin Chen. 2023. NA2Q: Neural Attention Additive Model for Interpretable Multi-Agent Q-Learning. In *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*. JMLR, Honolulu, USA, 22539–22558.
- [241] Zichuan Liu, Yuanyang Zhu, Zhi Wang, Yang Gao, and Chunlin Chen. 2025. MIXRTs: Toward Interpretable Multi-Agent Reinforcement Learning via Mixing Recurrent Soft Decision Trees. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 47, 5 (2025), 4090–4107.
- [242] Silas Lobo, Stefan Neumeier, Evelio M. G. Fernandez, and Christian Facchi. 2020. In-TAS - The Ingolstadt Traffic Scenario for SUMO. In *Proceedings of the 2020 SUMO User Conference (SUMO '20)*. SUMO, Berlin, Germany, 73–92.
- [243] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2017. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. In *Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPS '17)*. NeurIPS, Long Beach, USA, 6379–6390.
- [244] Alice Lu, Peter Frazier, and Oren Kislev. 2018. Surge Pricing Moves Uber’s Driver Partners. arXiv:3180246
- [245] Zhuoran Lu and Ming Yin. 2021. Human Reliance on Machine Learning Models When Performance Feedback is Limited: Heuristics and Risks. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*. ACM, Yokohama, Japan, 1–16.
- [246] Felipe Luyanda, Douglas Gettman, Larry Head, Steven Shelby, Darcy Bullock, and Pitu

- Mirchandani. 2003. ACS-Lite Algorithmic Architecture: Applying Adaptive Control System Technology to Closed-Loop Traffic Signal Control Systems. *Transportation Research Record* 1856, 1 (2003), 175–184.
- [247] Hongyao Ma, Fei Fang, and David C. Parkes. 2022. Spatio-Temporal Pricing for Ridesharing Platforms. *Operations Research* 70, 2 (2022), 1025–1041.
  - [248] Haitong Ma, Changliu Liu, Shengbo Eben Li, Sifa Zheng, Wenchao Sun, and Jianyu Chen. 2024. Learn Zero-Constraint-Violation Safe Policy in Model-Free Constrained Reinforcement Learning. *IEEE Transactions on Neural Networks and Learning Systems* 36, 2 (2024), 2327–2341.
  - [249] Jiaqi Ma, Michael D. Fontaine, Fang Zhou, and Jia Hu. 2016. Estimation of Crash Modification Factors for an Adaptive Traffic-Signal Control System. *Journal of Transportation Engineering* 142, 12 (2016), 04016061.
  - [250] Jinming Ma and Feng Wu. 2020. Feudal Multi-Agent Deep Reinforcement Learning for Traffic Signal Control. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '20)*. IFAAMAS, Auckland, NZ, 816–824.
  - [251] Ning F. Ma and Benjamin V. Hanrahan. 2019. Part-Time Ride-Sharing: Recognizing the Context in which Drivers Ride-Share and its Impact on Platform Use. *Proceedings of the ACM on Human-Computer Interaction* 3 (2019), 1–17.
  - [252] Ning F. Ma, Zheng Yao, Veronica A. Rivera, and Dongwook Yoon. 2022. “Brush it Off”: How Women Workers Manage and Cope with Bias and Harassment in Gender-agnostic Gig Platforms. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. ACM, New Orleans, USA, 1–13.
  - [253] Ning F. Ma, Chien Wen Yuan, Moojan Ghafurian, and Benjamin V. Hanrahan. 2018. Using Stakeholder Theory to Examine Drivers’ Stake in Uber. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. ACM, Montréal, Canada, 1–12.
  - [254] Zian Ma, Chengcheng Xu, Yuheng Kan, Maonan Wang, and Wei Wu. 2021. Adaptive Coordinated Traffic Control for Arterial Intersections based on Reinforcement Learning. In *Proceedings of the 2021 IEEE 24th International Conference on Intelligent Transportation Systems (ITSC '21)*. IEEE, Indianapolis, USA, 2562–2567.
  - [255] Carl A. MacCarley, Stephen L.M. Hockaday, Daniel Need, and Samuel Taff. 1992. Evaluation of Video Image Processing Systems for Traffic Detection. *Transportation Research Record* 1360 (1992), 46–49.
  - [256] Michal Maciejewski. 2010. A Comparison of Microscopic Traffic Flow Simulation Systems for an Urban Area. *Transport Problems* 5, 4 (2010), 27–38.
  - [257] Maria Madsen and Shirley Gregor. 2000. Measuring Human-Computer Trust. In *Proceedings of the 11th Australasian Conference on Information Systems (ACIS '00)*. AAIS, Brisbane, Australia, 6–8.
  - [258] Travis Mandel, Yun-En Liu, Emma Brunskill, and Zoran Popović. 2017. Where to Add

- Actions in Human-in-the-Loop Reinforcement Learning. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI '17)*. AAAI, San Francisco, USA, 2322–2328.
- [259] Fred L. Mannering and Scott S. Washburn. 2018. Traffic Control and Analysis at Signalized Intersections. In *Principles of Highway Engineering and Traffic Analysis*. Wiley, Hoboken, USA, 243–310.
- [260] Gustav Markkula, Ola Benderius, Krister Wolff, and Mattias Wahde. 2012. A Review of Near-Collision Driver Behavior Models. *Human Factors* 54, 6 (2012), 1117–1143.
- [261] Jijo Mathew, Jairaj Desai, Rahul Suryakant Sakhare, Woosung Kim, Howell Li, and Darcy M. Bullock. 2021. Big Data Applications for Managing Roadways. *Institute of Transportation Engineers Journal* 91, 2 (2021), 28–35.
- [262] Roger C. Mayer, James H. Davis, and F. David Schoorman. 1995. An Integrative Model of publisheral Trust. *The Academy of Management Review* 20, 3 (1995), 709–734.
- [263] Jay Mayfield. 2024. *FTC Takes Action to Stop Lyft from Deceiving Drivers with Misleading Earnings Claims*. Federal Trade Commission. <https://www.ftc.gov/news-events/news/press-releases/2024/10/ftc-takes-action-stop-lyft-deceiving-drivers-misleading-earnings-claims>
- [264] Bryce McLaughlin and Jann Spiess. 2024. Designing Algorithmic Recommendations to Achieve Human–AI Complementarity. arXiv:2405.01484
- [265] Juan Medina, Madhav Chitturi, and Rahim Benekohal. 2010. Effects of fog, snow, and rain on video detection systems at intersections. *Transportation Letters* 2, 1 (2010), 1–12.
- [266] Juan C. Medina and Rahim F. Benekohal. 2012. Traffic signal control using reinforcement learning and the max-plus algorithm as a coordinating strategy. In *Proceedings of the 15th International IEEE Conference on Intelligent Transportation Systems (ITSC '12)*. IEEE, Anchorage, USA, 596–601.
- [267] Juan C. Medina, Hani Ramezani, and Rahim (Ray) F. Benekohal. 2013. Evaluation of Microwave Radar Vehicle Detectors at a Signalized Intersection Under Adverse Weather Conditions. *Transportation Research Record* 2356 (2013), 100–108.
- [268] Hao Mei, Xiaoliang Lei, Longchao Da, Bin Shi, and Hua Wei. 2024. LibSignal: an open library for traffic signal control. *Machine Learning* 113 (2024), 5235–5271.
- [269] Hao Mei, Junxian Li, Bin Shi, and Hua Wei. 2023. Reinforcement learning approaches for traffic signal control under missing data. In *Proceedings of the 32nd International Joint Conference on Artificial Intelligence (IJCAI '23)*. ACM, Macau, 2261–2269.
- [270] Xutao Mei, Nijiro Fukushima, Bo Yang, Zheng Wang, Tetsuya Takata, Hiroyuki Nagasawa, and Kimihiko Nakano. 2023. Reinforcement Learning-Based Intelligent Traffic Signal Control Considering Sensing Information of Railway. *IEEE Sensors Journal* 23, 24 (2023), 31125–31136.
- [271] Wei Miao, Yiting Deng, Wei Wang, Yongdong Liu, and Christopher S. Tang. 2023. The effects of surge pricing on driver behavior in the ride-sharing market: Evidence from a quasi-experiment. *Journal of Operations Management* 69, 5 (2023), 794–822.

- [272] Paul G Michael, Frank C Leeming, and William O Dwyer. 2000. Headway on urban streets: observational data and an intervention to decrease tailgating. *Transportation Research Part F: Traffic Psychology and Behaviour* 3 (2000), 55–64. Issue 2.
- [273] Panos G. Michalopoulos. 1991. A real-time computer vision system for vehicle tracking and traffic surveillance. *IEEE Transactions on Vehicular Technology* 40, 1 (1991), 21–29.
- [274] Dan Middleton and Ricky Parker. 2000. *Initial Evaluation of Selected Detectors to Replace Inductive Loops on Freeways*. Technical Report FHWA/TX-00/1439-7. Texas Department of Transportation. 1–88 pages.
- [275] Dan Middleton, Ricky Parker, and Ryan Longmire. 2007. *Investigation of Vehicle Detector Performance and ATMS Interface*. Technical Report FHWA/TX-07/0-4750-2. Texas Department of Transportation. 1–144 pages.
- [276] Vicente Milanés and Steven E. Shladover. 2014. Modeling cooperative and autonomous adaptive cruise control dynamic responses using experimental data. *Transportation Research Part C: Emerging Technologies* 48 (2014), 285–300.
- [277] Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. 2024. Explainable Reinforcement Learning: A Survey and Comparative Review. *Comput. Surveys* 56, 7 (2024), 1–36.
- [278] Stephanie Milani, Zhicheng Zhang, Nicholay Topin, Zheyuan Ryan Shi, Charles Kamhoua, Evangelos E. Papalexakis, and Fei Fang. 2022. MAVIPER: Learning Decision Tree Policies for Interpretable Multi-agent Reinforcement Learning. In *Proceedings of the 2022 European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD '22)*. ACM, Grenoble, France, 251–266.
- [279] Shahin Mirbakhsh and Mahdi Azizi. 2024. Adaptive Traffic Signal’s Safety and Efficiency Improvement by MultiObjective Deep Reinforcement Learning Approach. *International Journal of Innovative Research in Multidisciplinary Education* 3, 7 (2024), 1245–1257.
- [280] Anuja Bharat Mistry and Shivani Tanna. 2025. *US online spending surges \$24.1 billion as steep discounts boost sales, Adobe says*. Reuters. <https://www.reuters.com/business/retail-consumer/us-online-spending-surges-241-billion-steep-discounts-boost-sales-adobe-says-2025-07-12/>
- [281] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumar, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518 (2015), 529–533.
- [282] Christoph Molnar. 2019. Taxonomy of Interpretability Methods. In *Interpretable Machine Learning*. Independent, Munich, Germany.
- [283] Seungwuk Moon and Kyongsu Yi. 2008. Human driving data-based design of a vehicle adaptive cruise control algorithm. *Vehicle System Dynamics* 46, 8 (2008), 661–690.
- [284] Sara Moridpour, Majid Sarvi, and Geoff Rose. 2010. Lane changing models: a critical

- review. *Transportation Letters* 2, 3 (2010), 157–173.
- [285] Sara Moridpour, Majid Sarvi, Geoff Rose, and Euan Ramsay. 2008. Variables influencing lane changing behaviour of heavy vehicles. In *Proceedings of the 31st Australasian Transport Research Forum (ATRF '08)*. Australasian Transport Research Forum, Canberra, Australia, 1–15.
  - [286] Yoshinari Motokawa and Toshiharu Sugawara. 2023. Interpretability for Conditional Coordinated Behavior in Multi-Agent Reinforcement Learning. In *Proceedings of the 2023 International Joint Conference on Neural Networks (IJCNN '23)*. IEEE, Gold Coast, Australia, 1–8.
  - [287] Jian Mou and Jason F. Cohen. 2017. Trust and online consumer health service success: A longitudinal study. *Information Development* 33, 2 (2017), 169–189.
  - [288] Jian Mou, Dong-Hee Shin, and Jason Cohen. 2017. Understanding trust and perceived usefulness in the consumer acceptance of an e-service: a longitudinal investigation. *Behaviour & Information Technology* 36, 2 (2017), 125–139.
  - [289] Brittany Moye. 2025. AAA: Fear in Self-Driving Vehicles Persists. AAA. <https://newsroom.aaa.com/2025/02/aaa-fear-in-self-driving-vehicles-persists/>
  - [290] Arthur Müller, Vishal Rangras, Georg Schnittker, Michael Waldmann, Maxim Friesen, Tobias Ferfers, Lukas Schreckenberger, Florian Hufen, Jürgen Jasperneite, and Marco Wiering. 2021. Towards Real-World Deployment of Reinforcement Learning for Traffic Signal Control. In *Proceedings of the 20th IEEE International Conference on Machine Learning and Applications (ICMLA '21)*. IEEE, Pasadena, USA, 507–514.
  - [291] Arthur Müller and Matthia Sabatelli. 2022. Safe and Psychologically Pleasant Traffic Signal Control with Reinforcement Learning using Action Masking. In *Proceedings of the 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC '22)*. IEEE, Macau, China, 951–958.
  - [292] Arthur Müller and Matthia Sabatelli. 2023. Bridging the Reality Gap of Reinforcement Learning based Traffic Signal Control using Domain Randomization and Meta Learning. In *Proceedings of the 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC '23)*. IEEE, Bilbao, Spain, 5271–5278.
  - [293] Ziaul Haque Munim, Mariia Dushenko, Veronica Jaramillo Jimenez, Mohammad Hassan Shakil, and Marius Imset. 2020. Big data and artificial intelligence in the maritime industry: a bibliometric review and future research directions. *Maritime Policy & Management* 47, 5 (2020), 577–597.
  - [294] Mohammad Naiseh, Reem S. Al-Mansoori, Dena Al-Thani, Nan Jiang, and Raian Ali. 2021. Nudging through Friction: An Approach for Calibrating Trust in Explainable AI. In *Proceedings of the 8th International Conference on Behavioral and Social Computing (BESC '21)*. IEEE, Virtual, 1–5.
  - [295] Mohammad Naiseh, Dena Al-Thani, Nan Jiang, and Raian Ali. 2023. How the different



- explanation classes impact trust calibration: The case of clinical decision support systems. *International Journal of Human-Computer Studies* 169 (2023), 102941.
- [296] Satoru Nakagawa, Dean Kriellaars, Christine Blais, Jeannette Montufar, and Michelle M. Porter. 2006. *Speed and acceleration patterns of younger and older drivers*. Technical Report. University of Manitoba, Winnipeg, Canada. 1–14 pages.
  - [297] Mohammadreza Nazari, Afshin Oroojlooy, Martin Takáč, and Lawrence V. Snyder. 2018. Reinforcement Learning for Solving the Vehicle Routing Problem. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems (NeurIPS '18)*. NeurIPS, Montréal, Canada, 9861–9871.
  - [298] Philip Ndikum and Serge Ndikum. 2024. Advancing Investment Frontiers: Industry-grade Reinforcement Learning for Portfolio Optimization. arXiv:2403.07916
  - [299] Gim Huay Neo, Robin Riedel, and Evgeni Kochman. 2025. *Intelligent Transport, Greener Future: AI as a Catalyst to Decarbonize Global Logistics*. Technical Report. World Economic Forum. 33 pages.
  - [300] Daiheng Ni. 2020. Controllers and Detectors. In *Signalized Intersections: Fundamentals to Advanced Systems*. Springer, New York, USA, 179–209.
  - [301] Tomoki Nishi, Keisuke Otaki, Keiichiro Hayakawa, and Takayoshi Yoshimura. 2018. Traffic Signal Control Based on Reinforcement Learning with Graph Convolutional Neural Nets. In *Proceedings of the 2018 IEEE 21st International Conference on Intelligent Transportation Systems (ITSC '18)*. IEEE, Maui, USA, 877–883.
  - [302] Mohammad Noaeen, Atharva Naik, Liana Goodman, Jared Crebo, Taimoor Abrar, Zahra Shakeri Hossein Abad, Ana L.C. Bazzan, and Behrouz Far. 2022. Reinforcement learning in urban network traffic signal control: A systematic literature review. *Expert Systems with Applications* 199 (2022), 116830.
  - [303] Mojtaba Norouzi, Monireh Abdoos, and Ana L. C. Bazzan. 2021. Experience classification for transfer learning in traffic signal control. *The Journal of Supercomputing* 77 (2021), 780–795.
  - [304] Hao Yi Ong, Daniel Freund, and Davide Crapis. 2021. Driver Positioning and Incentive Budgeting with an Escrow Mechanism for Ride-Sharing Platforms. *INFORMS Journal on Applied Analytics* 51, 5 (2021), 373–390.
  - [305] Aoyu Pang, Maonan Wang, Yirong Chen, Man-On Pun, and Michael Lepech. 2024. Scalable Reinforcement Learning Framework for Traffic Signal Control Under Communication Delays. *IEEE Open Journal of Vehicular Technology* 5 (2024), 330–343.
  - [306] Saumya Pareek, Eduardo Velloso, and Jorge Goncalves. 2024. Trust Development and Repair in AI-Assisted Decision-Making during Complementary Expertise. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT '24)*. ACM, Rio de Janeiro, Brazil, 546–561.
  - [307] Darsh Parekh, Nishi Poddar, Aakash Rajpurkar, Manisha Chahal, Neeraj Kumar, Gyanendra Prasad Joshi, and Woong Cho. 2022. A Review on Autonomous Vehicles: Progress,

- Methods and Challenges. *Electronics* 11, 14 (2022), 2162.
- [308] Jay J. Park, Jakov Tiefenbach, and Andreas K. Demetriades. 2022. The role of artificial intelligence in surgical simulation. *Frontiers in Medical Technology* 4 (2022), 1076755.
  - [309] Bhavik Pathak, Robert Garfinkel, Ram D. Gopal, Rajkumar Venkatesan, and Fang Yin. 2010. Empirical Analysis of the Impact of Recommender Systems on Sales. *Journal of Management Information Systems* 27, 2 (2010), 159–188.
  - [310] Konstantin F. Pilz, Lennart Heim, and Nicholas Brown. 2025. Increased Compute Efficiency and the Diffusion of AI Capabilities. In *Proceedings of the 39th AAAI Conference on Artificial Intelligence (AAAI '25)*. AAAI, Philadelphia, USA, 27582–27590.
  - [311] Pascal Poupart, Zhitang Chen, Priyank Jaini, Fred Fung, Hengky Susanto, and Yanhui Geng. 2016. Online flow size prediction for improved network routing. In *Proceedings of the 2016 IEEE 24th International Conference on Network Protocols (ICNP '16)*. IEEE, Singapore, 1–6.
  - [312] Snehal Prabhudesai, Leyao Yang, Sumit Asthana, Xun Huan, Q. Vera Liao, and Nikola Banovic. 2023. Understanding Uncertainty: How Lay Decision-makers Perceive and Interpret Uncertainty in Human-AI Decision Making. In *Proceedings of the ACM Conference on Intelligent User Interfaces 2023 (IUI '23)*. ACM, Sydney, Australia, 379–396.
  - [313] K.J. Prabuchandran, A.N. Hemanth Kumar, and Shalabh Bhatnagar. 2015. Decentralized learning for traffic signal control. In *Proceedings of the 7th International Conference on Communication Systems and Networks (COMSNETS '15)*. IEEE, Bangalore, India, 2529–2534.
  - [314] Yiheng Qian, Tejaswi Polimetla, Thomas W. Sanchez, and Xiang Yan. 2024. How do transportation professionals perceive the impacts of AI applications in transportation? A latent class cluster analysis. arXiv:2401.08915
  - [315] Zhiwei (Tony) Qin, Xiaocheng Tang, Yan Jiao, Fan Zhang, Zhe Xu, Hongtu Zhu, and Jieping Ye. 2020. Ride-Hailing Order Dispatching at DiDi via Reinforcement Learning. *INFORMS Journal on Applied Analytics* 50, 5 (2020), 272–286.
  - [316] Anjie Qiu, Prapul Alemada Sathish, Donglin Wang, and Hans D. Schotten. 2024. Advanced Traffic Demand Generation in SUMO: ML-based Prediction of Flow Rate based on Real-world Measured Datasets. In *Proceedings of the 2024 IEEE 99th Vehicular Technology Conference (VTC '24)*. IEEE, Singapore, 1–7.
  - [317] Filip Radlinski, Craig Boutilier, Deepak Ramachandran, and Ivan Vendrov. 2022. Subjective Attributes in Conversational Recommendation Systems: Challenges and Opportunities. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI '22)*. AAAI, Vancouver, Canada, 12287–12293.
  - [318] Md. Mokhlesur Rahman, Pooya Najaf, Milton Gregory Fields, and Jean-Claude Thill. 2022. Traffic congestion and its urban scale factors: Empirical evidence from American urban areas. *International Journal of Sustainable Transportation* 16, 5 (2022), 406–421.
  - [319] Ram Rajagopal and Pravin Varaiya. 2007. *Health of California's Loop Detector System*.

- Technical Report UCB-ITS-PRR-2007-13. California PATH Program. 59 pages.
- [320] Nived Rajaraman, Lin F. Yang, Jiantao Jiao, and Kannan Ramchandran. 2020. Toward the Fundamental Limits of Imitation Learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS '20)*. NeurIPS, Vancouver, Canada, 2914–2924.
  - [321] Naveen Raman, Sanket Shah, and John P. Dickerson. 2021. Data-Driven Methods for Balancing Fairness and Efficiency in Ride-Pooling. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI '21)*. IJCAI, Montréal, Canada, 363–369.
  - [322] Jette Randsløv and Preben Alstrøm. 1998. Learning to Drive a Bicycle Using Reinforcement Learning and Shaping. In *Proceedings of the 15th International Conference on Machine Learning (ICML '98)*. JMLR, Madison, USA, 463–471.
  - [323] Varun Nagaraj Rao, Samantha Dalal, Eesha Agarwal, Dana Calacci, and Andrés Monroy-Hernández. 2025. Rideshare Transparency: Translating Gig Worker Insights on AI Platform Design to Policy. In *Proceedings of the 28th ACM Conference on Computer-Supported Cooperative Work and Social Computing (CSCW '25)*. ACM, Bergen, Norway, 1–49.
  - [324] Marco Rapelli, Claudio Casetti, and Giandomenico Gagliardi. 2022. Vehicular traffic simulation in the city of Turin from raw data. *IEEE Transactions on Mobile Computing* 21, 12 (2022), 4656–4666.
  - [325] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning (ICML '18)*. JMLR, Stockholm, Sweden, 7234–7284.
  - [326] Nedal T. Ratrouf and Imran Reza. 2014. Comparison of Optimal Signal Plans by Synchro & TRANSYT-7F Using PARAMICS — A Case Study. *Procedia Computer Science* 32 (2014), 372–379.
  - [327] Robert Rausch and David Miller. 2007. Advanced Transportation Controller Standards Update. *Institute of Transportation Engineers Journal* 77, 5 (2007), 29–33.
  - [328] Bijul Raveendran, Tom V. Mathew, and Nagendra R. Velaga. 2025. Impact of countdown timer on drivers' anticipation at the onset of yellow at signalized intersections. *Transportation Letters* 17, 2 (2025), 356–368.
  - [329] Alex Ray, Joshua Achiam, and Dario Amodei. 2019. *Benchmarking Safe Exploration in Deep Reinforcement Learning*. Technical Report. OpenAI. 1–25 pages.
  - [330] Pouria Razzaghi, Amin Tabrizian, Wei Guo, Shulu Chen, Abenezer Taye, Ellis Thompson, Alexis Bregeon, Ali Baheri, and Peng Wei. 2022. A Survey on Reinforcement Learning in Aviation Applications. arXiv:2211.02147
  - [331] Yuqing Ren, Xuefei (Nancy) Deng, and K.D. Joshi. 2023. Unpacking Human and AI Complementarity: Insights from Recent Works. *ACM SIGMIS Database* 54, 3 (2023), 6–10.
  - [332] Duncan Rheingans-Yoo, Scott Duke Kominers, Hongyao Ma, and David C. Parkes. 2019. Ridesharing with driver location preferences. In *Proceedings of the 28th International Joint*

- Conference on Artificial Intelligence (IJCAI '19)*. ACM, Macau, 557–564.
- [333] Avery Rhodes, Darcy M. Bullock, James R. Sturdevant, and Zachary Thomas Clark. 2005. *Evaluation of Stop Bar Video Detection Accuracy at Signalized Intersections*. Technical Report FHWA/IN/JTRP-2005/28. Joint Transportation Research Program, Indiana Department of Transportation and Purdue University. 1–418 pages.
  - [334] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2006. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*. ACM, San Francisco, USA, 1135–1144.
  - [335] Brishen Rogers. 2015. The Social Costs of Uber. *University of Chicago Law Review* 82 (2015), 85–102.
  - [336] Aji Ronaldo and M. Taufiq Ismail. 2012. *Comparison of the two Micro-Simulation Software AIMSUN & SUMO for Highway Traffic Modelling*. Master’s thesis. Linköping University, Linköping.
  - [337] Stéphane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. 2011. A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS '11)*. PMLR, Fort Lauderdale, USA, 627–635.
  - [338] Aaron M. Roth, Jing Liang, Ram Sriram, Elham Tabassi, and Dinesh Manocha. 2023. MSVIPER: Improved Policy Distillation for Reinforcement-Learning-Based Robot Navigation. *Journal of the Washington Academy of Sciences* 109, 2 (2023), 27–58.
  - [339] Jingqing Ruan, Ziyue Li, Hua Wei, Haoyuan Jiang, Jiaming Lu, Xuantang Xiong, Hangyu Mao, and Rui Zhao. 2024. CoSLight: Co-optimizing Collaborator Selection and Decision-making to Enhance Traffic Signal Control. In *Proceedings of the 30th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '24)*. ACM, Barcelona, Spain, 2500–2511.
  - [340] Ziad A. Sabra, Douglas Gettman, R. David Henry, and Venkata Nallamothu. 2010. *Balancing Safety and Capacity in an Adaptive Signal Control System — Phase 1*. Technical Report FHWA-HRT-10-038. U.S. Department of Transportation Federal Highway Administration. 106 pages.
  - [341] Mustapha Saidallah, Abdeslam El Fergougui, and Abdelbaki Elbelrhiti Elalaoui. 2016. A Comparative Study of Urban Road Traffic Simulators. In *Proceedings of the 2016 International Conference on Civil, Transportation and Environment (ICCTE '16)*. MATEC, Guangzhou, China, 1–6.
  - [342] Mohammad Saifuzzaman and Zuduo Zheng. 2014. Incorporating human-factors in car-following models: A review of recent developments and research needs. *Transportation Research Part C: Emerging Technologies* 48 (2014), 379–403.
  - [343] Johnny Saldaña. 2015. First-Cycle Coding Methods. In *The Coding Manual for Qualitative Researchers* (3rd ed.). SAGE, San Diego, USA, 67–206.

- [344] Johnny Saldaña. 2015. Second-Cycle Coding Methods. In *The Coding Manual for Qualitative Researchers* (3rd ed.). SAGE, San Diego, USA, 233–268.
- [345] Erica Salvato, Gianfranco Fenu, Eric Medvet, and Felice Andrea Pellegrino. 2021. Crossing the Reality Gap: A Survey on Sim-to-Real Transferability of Robot Controllers in Reinforcement Learning. *IEEE Access* 9 (2021), 153171–153187.
- [346] Iqbal H. Sarker. 2021. Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and Research Directions. *SN Computer Science* 2 (2021), 420.
- [347] R M Savithramma and R. Sumathi. 2023. Intelligent traffic signal controller for heterogeneous traffic using reinforcement learning. *Green Energy and Intelligent Transportation* 2, 6 (2023), 100124.
- [348] J.R. Sayer, M.L. Mefford, P.S. Fancher, R.E. Ervin, and S.E. Bogard. 1997. An Experimental Design for Studying How Driver Characteristics Influence Headway Control. In *Proceedings of the 1st International Conference on Intelligent Transportation Systems (ITSC '97)*. IEEE, Piscataway, USA, 870–875.
- [349] Bruce Schaller. 2021. Can sharing a ride make for less traffic? Evidence from Uber and Lyft and implications for cities. *Transport Policy* 102 (2021), 1–10.
- [350] Nicolas Scharowski, Michaela Benk, Swen J. Kühne, Léane Wettstein, and Florian Brühlmann. 2023. Certification Labels for Trustworthy AI: Insights From an Empirical Mixed-Method Study. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23)*. ACM, Chicago, USA, 248–262.
- [351] Max Schemmer, Niklas Kühl, Carina Benz, Andrea Bartos, and Gerhard Satzger. 2023. Appropriate Reliance on AI Advice: Conceptualization and the Effect of Explanations. In *Proceedings of the 28th International Conference on Intelligent User Interfaces (IUI '23)*. ACM, Sydney, Australia, 410–422.
- [352] Lukas M Schmidt, Georgios Kontes, Axel Plinge, and Christopher Mutschler. 2021. Can you trust your autonomous car? Interpretable and verifiably safe reinforcement learning. In *Proceedings of the 2021 IEEE Intelligent Vehicles Symposium (IV '21)*. IEEE, Nagoya, Japan, 171–178.
- [353] Jakob Schoeffer, Niklas Kuehl, and Yvette Machowski. 2022. “There Is Not Enough Information”: On the Effects of Explanations on Perceptions of Informational Fairness and Trustworthiness in Automated Decision-Making. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*. ACM, Seoul, Korea, 1616–1628.
- [354] David Schrank, Luke Albert, Kartikeya Jha, and Bill Eisele. 2023. *2023 Urban Mobility Report*. Technical Report. Texas A&M Transportation Institute. 80 pages.
- [355] Lincoln V Schreiber, Lucas N Alegre, Ana LC Bazzan, and Gabriel de O Ramos. 2022. On the explainability and expressiveness of function approximation methods in RL-based traffic signal control. In *Proceedings of the 2022 International Joint Conference on Neural Networks (IJCNN '22)*. IEEE, Padua, Italy, 1–8.

- [356] Skipper Seabold and Josef Perktold. 2010. statsmodels: Econometric and statistical modeling with Python. In *Proceedings of the 9th Python in Science Conference (SciPy '10)*. Python in Science Conference, Austin, USA, 92–96.
- [357] Ujjwal Sehrawat, Namit Sawhney, Tejaswini Yeleswarapu, and Nimmi Rangaswamy. 2021. The Everyday HCI of Uber Drivers in India: A Developing Country Perspective. In *Proceedings of the 24th ACM Conference on Computer-Supported Cooperative Work and Social Computing (CSCW '21)*. ACM, Virtual, 1–22.
- [358] Mahmoud Selim, Amr Alanwar, Shreyas Kousik, Grace Gao, Marco Pavone, and Karl H. Johansson. 2022. Safe Reinforcement Learning Using Black-Box Reachability Analysis. *IEEE Robotics and Automation Letters* 7, 4 (2022), 10665–10672.
- [359] Marc Semrau, Jakob Erdmann, Bernhard Friedrich, and René Waldmann. 2016. Simulation framework for testing ADAS in Chinese traffic situations. In *Proceedings of the 2016 SUMO User Conference (SUMO '16)*. German Aerospace Center, Cologne, Germany, 103–114.
- [360] Susan Shaheen and Nelson Chan. 2016. Mobility and the Sharing Economy: Potential to Facilitate the First- and Last-Mile Public Transit Connections. *Built Environment* 42, 4 (2016), 573–588.
- [361] Mahboubeh Shamsi, Abdolreza Rasouli Kenari, and Roghayeh Aghamohammadi. 2022. Reinforcement learning for traffic light control with emphasis on emergency vehicles. *The Journal of Supercomputing* 78, 4 (2022), 4911–4937.
- [362] Jing Shang, Shunmei Meng, Jun Hou, Xiaoran Zhao, Xiaokang Zhou, and Rong Jiang. 2025. Graph-Based Cooperation Multiagent Reinforcement Learning for Intelligent Traffic Signal Control. *IEEE Internet of Things Journal* 12, 10 (2025), 14362–14374.
- [363] Steven G. Shelby, Darcy M. Bullock, Doug Gettman, Raj S. Ghaman, Ziad A. Sabra, and Nils Soyke. 2008. An Overview and Performance Evaluation of ACS Lite — A Low Cost Adaptive Signal Control System. In *Proceedings of the 87th Annual Meeting of the Transportation Research Board (TRB '08)*. TRB, Washington, DC, USA, 1–17.
- [364] Tianyu Shi, Francois Xavier Devailly, Denis Larocque, and Laurent Charlin. 2024. Improving the Generalizability and Robustness of Large-Scale Traffic Signal Control. *IEEE Open Journal of Intelligent Transportation Systems* 5 (2024), 2–15.
- [365] Yang Shi, Zhenbo Wang, Tim J. LaClair, Chieh (Ross) Wang, Yunli Shao, and Jinghui Yuan. 2023. A Novel Deep Reinforcement Learning Approach to Traffic Signal Control with Connected Vehicles. *Applied Sciences* 13, 4 (2023), 2750.
- [366] Zheyuan Ryan Shi, Claire Wang, and Fei Fang. 2020. Artificial Intelligence for Social Good: A Survey. arXiv:2001.01818
- [367] Donghee Shin. 2020. How do users interact with algorithm recommender systems? The interaction of users, algorithms, and performance. *Computers in Human Behavior* 109 (2020), 106344.
- [368] Donghee Shin. 2021. The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable AI. *International Journal of Human-Computer*

*Studies* 146 (2021), 102551.

- [369] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS '23)*. NeurIPS, New Orleans, USA, 8634–8652.
- [370] Andrew Silva, Taylor Killian, Ivan Rodriguez Jimenez, Sung-Hyun Son, and Matthew Gombolay. 2020. Optimization Methods for Interpretable Differentiable Decision Trees in Reinforcement Learning. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS '20)*. PMLR, Palermo, Italy, 1855–1865.
- [371] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 529 (2016), 484–489.
- [372] Anubha Singh, Patricia Garcia, and Silvia Lindtner. 2023. Old Logics, New Technologies: Producing a Managed Workforce on On-Demand Service Platforms. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–15.
- [373] Ho Chit Siu, Jaime Peña, Edenna Chen, Yutai Zhou, Victor Lopez, Kyle Palko, Kimberlee Chang, and Ross Allen. 2021. Evaluation of human-AI teams for learned and rule-based agents in Hanabi. In *Proceedings of the 35th Conference on Neural Information Processing Systems (NeurIPS '21)*. NeurIPS, Virtual, 16183–16195.
- [374] Venkatesh Sivaraman, Leigh A. Bukowski, Joel Levin, Jeremy M. Kahn, and Adam Perer. 2023. Ignore, Trust, or Negotiate: Understanding Clinician Acceptance of AI-Based Treatment Recommendations in Health Care. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–18.
- [375] Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krashennnikov, and David Krueger. 2022. Defining and Characterizing Reward Hacking. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS '22)*. NeurIPS, New Orleans, USA, 24824–24837.
- [376] Stephen Smith, Gregory Barlow, Xiao-Feng Xie, and Zack Rubinstein. 2013. Smart urban signal networks: Initial application of the SURTRAC adaptive traffic signal control system. In *Proceedings of the 23rd International Conference on Automated Planning and Scheduling (ICAPS '13)*. PKP, Rome, Italy, 434–442.
- [377] Elizabeth Solberg, Magnhild Kaarstad, Maren H. Rø Eitrheim, Rossella Bisio, Kine Reegard, and Marten Bloch. 2022. A Conceptual Model of Trust, Perceived Risk, and Reliance on AI Decision Aids. *Group & publisher Management* 47, 2 (2022), 187–222.
- [378] Raghavan Srinivasan, Craig Lyon, Bhagwant Persaud, Jongdae Baek, Frank Gross, Sarah Smith, and Carl Sundstrom. 2012. Crash Modification Factors for Changes to Left-Turn

- Phasing. *Transportation Research Record* 2279, 1 (2012), 108–117.
- [379] Akash Sriram. 2025. *Uber’s quarterly revenue miss turns spotlight on slowing ride-hailing growth*. Reuters. <https://www.reuters.com/business/autos-transportation/uber-forecasts-upbeat-second-quarter-steady-ride-hailing-delivery-demand-2025-05-07/>
  - [380] Jeff C. Stanley and Stephen L. Dorton. 2023. Exploring Trust With the AI Incident Database. In *Proceedings of the 2023 Human Factors and Ergonomics Society Annual Meeting (HFES ’23)*. ACM, Washington DC, USA, 489–494.
  - [381] Aleksandar Stevanonic. 2010. Chapter Six - Implementation Costs and Benefits. In *Adaptive Traffic Control Systems: Domestic and Foreign State of Practice*. NHCPR Synthesis, Vol. 403. TRB, Washington, DC, USA, 36–42.
  - [382] Aleksandar Stevanovic, Cameron Kergaye, and Peter T. Martin. 2009. SCOOT and SCATS: A Closer Look into Their Operations. In *Proceedings of the 88th Annual Meeting of the Transportation Research Board (TRB ’09)*. TRB, Washington, DC, USA, 1–17.
  - [383] Aleksandar Stevanovic and Peter T. Martin. 2008. Split-Cycle Offset Optimization Technique and Coordinated Actuated Traffic Control Evaluated Through Microsimulation. *Transportation Research Record* 2080 (2008), 48–56.
  - [384] Eroze Sthapit and Peter Björk. 2019. Sources of value co-destruction: Uber customer perspectives. *Tourism Review* 74, 4 (2019), 780–794.
  - [385] Rachel E. Stuck, Brittany E. Holthausen, and Bruce N. Walker. 2021. Chapter 8 — The role of risk in human-robot trust. In *Trust in Human-Robot Interaction*. Academic Press, Cambridge, UK, 179–194.
  - [386] Haoran Su, Yaofeng Desmond Zhong, Biswadip Dey, and Amit Chakraborty. 2022. EMV-Light: A Decentralized Reinforcement Learning Framework for Efficient Passage of Emergency Vehicles. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI ’22)*. AAAI, Virtual, 4593–4601.
  - [387] Joseph Suárez. 2024. PufferLib: Making Reinforcement Learning Libraries and Environments Play Nice. arXiv:2406.12905
  - [388] Dazhi Sun, Leslie Dodoo, Andres Rubio, Harsha Kalyan Penumala, Michael Pratt, and Srinivasa Sunkari. 2012. *Synthesis study of Texas signal control systems: technical report*. Technical Report FHWA/TX-13/0-6670-1. Texas A&M Transportation Institute. 1–65 pages.
  - [389] Daniel (Jian) Sun and Alexandra Kondyli. 2010. Modeling Vehicle Interactions during Lane-Changing Behavior on Arterial Streets. *Computer-Aided Civil and Infrastructure Engineering* 25 (2010), 557–571.
  - [390] Yuan Sun, Magdalayna Drivas, Mengqi Liao, and S. Shyam Sundar. 2023. When Recommender Systems Snoop into Social Media, Users Trust them Less for Health Advice. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI ’23)*. ACM, Hamburg, Germany, 1–14.



- [391] Yuan Sun, Navid Salami Pargoo, Peter J. Jin, and Jorge Ortiz. 2024. Optimizing Autonomous Driving for Safety: A Human-Centric Approach with LLM-Enhanced RLHF. In *Companion of the 2024 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp Companion '24)*. ACM, New York, USA, 76–80.
- [392] S. Shyam Sundar, Anne Oeldorf-Hirsch, and Qian Xu. 2008. The Bandwagon Effect of Collaborative Filtering Technology. In *Proceedings of the 2008 CHI Conference on Human Factors in Computing Systems (CHI '08)*. ACM, Florence, Italy, 3453–3458.
- [393] Srinivasa Sunkari, Apoorba Bibeka, Nadeem Chaudhary, and Kevin Balke. 2019. *Impact of Traffic Signal Controller Settings on the Use of Advanced Detection Devices*. Technical Report FHWA/TX-18/0-6934-R1. Texas A&M Transportation Institute. 1–67 pages.
- [394] Richard S. Sutton and Andrew G. Barto. 2018. Finite Markov Decision Processes. In *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, USA, 53–88.
- [395] Richard S. Sutton and Andrew G. Barto. 2018. Policy Approximation. In *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, USA, 257–264.
- [396] Richard S. Sutton and Andrew G. Barto. 2018. The Reinforcement Learning Problem. In *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, USA, 1–26.
- [397] Richard S. Sutton and Andrew G. Barto. 2018. Temporal-Difference Learning. In *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, USA, 143–166.
- [398] Pang Suwanaposee, Carl Gutwin, Zhe Chen, and Andy Cockburn. 2023. “Specially For You” — Examining the Barnum Effect’s Influence on the Perceived Quality of System Recommendations. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–11.
- [399] Meirav Taieb-Maimon and David Shinar. 2001. Minimum and Comfortable Driving Headways: Reality versus Perception. *Human Factors* 43, 1 (2001), 159–172.
- [400] Kai Liang Tan, Anuj Sharma, and Soumik Sarkar. 2020. Robust Deep Reinforcement Learning for Traffic Signal Control. *Journal of Big Data Analytics in Transportation* 2, 3 (2020), 263–274.
- [401] Zhi Ming Tan, Nikita Aggarwal, Josh Cowls, Jessica Morley, Mariarosaria Taddeo, and Luciano Floridi. 2021. The ethical debate about the gig economy: A review and critical analysis. *Technology and Society* 65 (2021), 101594.
- [402] Leo Tasca. 2000. *A review of the literature on aggressive driving research*. Technical Report. Ontario Ministry of Transportation, Toronto, Canada. 1–25 pages.
- [403] Larry W. Thomas. 2014. *Effect of MUTCD on Tort Liability of Government Transportation Agencies*. Legal Research Digest 63. National Cooperative Highway Research Program. 99 pages.
- [404] David Edmigio Valdivieso Tituana, Sang Guun Yoo, and Roberto O. Andrade. 2022. Vehicle Counting using Computer Vision: A Survey. In *Proceedings of the 2022 IEEE 7th International Conference for Convergence in Technology (I2CT '22)*. IEEE, Pune, India, 1–7.

- [405] Ernest Peter Todosiev. 1963. *The action point model of the vehicle-driver system*. Ph.D. Dissertation. The Ohio State University, Columbus, USA.
- [406] Tomer Toledo, Haris N. Koutsopoulos, and Moshe Ben-Akiva. 2007. Integrated driving behavior modeling. *Transportation Research Part C: Emerging Technologies* 15, 2 (2007), 96–112.
- [407] Richard Tomsett, Alun Preece, Dave Braines, Federico Cerutti, Supriyo Chakraborty, Mani Srivastava, Gavin Pearson, and Lance Kaplan. 2020. Rapid Trust Calibration through Interpretable and Uncertainty-Aware AI. *Patterns* 1, 4 (2020), 100049.
- [408] Nicholay Topin, Stephanie Milani, Fei Fang, and Manuela Veloso. 2020. Iterative Bounding MDPs: Learning Interpretable Policies via Non-Interpretable Methods. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI '20)*. AAAI, New York, USA, 9923–9931.
- [409] Hoang-Dung Tran, Feiyang Cai, Manzanar Lopez Diego, Patrick Musau, Taylor T. Johnson, and Xenofon Koutsoukos. 2019. Safety Verification of Cyber-Physical Systems with Reinforcement Learning Control. *ACM Transactions on Embedded Computing Systems* 18, 5 (2019), 1–22.
- [410] Manish Tripathy, Jiaru Bai, and H. Sebastian (Seb) Heese. 2022. Driver collusion in ride-hailing platforms. *Decision Sciences* 54, 4 (2022), 434–446.
- [411] Takane Ueno, Yuto Sawa, Yeongdae Kim, Jacqueline Urakami, Hiroki Oura, and Katie Seaborn. 2022. Trust in Human-AI Interaction: Scoping Out Models, Measures, and Methods. In *Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems (CHI EA '22)*. ACM, New Orleans, USA, 1–7.
- [412] M. Rehmat Ullah, Khuram S. Khattak, Zawar H. Khan, Mushtaq A. Khan, Nasru Minallah, and Akhtar N. Khan. 2021. Vehicular Traffic Simulation Software: A Systematic Comparative Analysis. *Pakistan Journal of Engineering and Technology* 4, 1 (2021), 66–78.
- [413] Sandesh Uppoor and Marco Fiore. 2011. Large-scale urban vehicular mobility for network-ing research. In *Proceedings of the 2011 IEEE Vehicular Networking Conference (VNC '11)*. IEEE, Amsterdam, Netherlands, 62–69.
- [414] Sandesh Uppoor and Marco Fiore. 2012. A large-scale vehicular mobility dataset of the Cologne urban area. In *Proceedings of the 14th French Conference on Algorithms and Telecommunications (AlgoTel '12)*. HAL, Hérault, France, 1–4.
- [415] Niels van Berkel and Vassilis Kostakos. 2021. Recommendations for Conducting Longitu-dinal Experience Sampling Studies. In *Advances in Longitudinal HCI Research*. Springer, Berlin, Germany, 59–78.
- [416] Elise van der Pol and Frans A. Oliehoek. 2016. Coordinated Deep Reinforcement Learners for Traffic Light Control. In *Proceedings of the 30th Conference on Neural Information Processing Systems (NeurIPS '16)*. NeurIPS, Barcelona, Spain, 1–8.
- [417] Marko Vasić, Andrija Petrović, Kaiyuan Wang, Mladen Nikolić, Rishabh Singh, and Sarfraz Khurshid. 2022. MoËT: Mixture of Expert Trees and its application to verifiable reinforce-

- ment learning. *Neural Networks* 151 (2022), 34–47.
- [418] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575 (2019), 350–354.
  - [419] Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, John Quan, Stephen Gaffney, Stig Petersen, Karen Simonyan, Tom Schaul, Hado van Hasselt, David Silver, Timothy Lillicrap, Kevin Calderone, Paul Keet, Anthony Brunasso, David Lawrence, Anders Ekermo, Jacob Repp, and Rodney Tsing. 2017. StarCraft II: A New Challenge for Reinforcement Learning. arXiv:1708.04782
  - [420] Peter Wagner. 2008. Action point models of human driving behaviour. In *Proceedings of the 2008 Traffic Simulation Workshop*. Institute of Highway Engineering and Transport Planning, Graz University of Technology, Graz, Austria, 1.
  - [421] Caroline Wang, Ishan Durugkar, Elad Liebman, and Peter Stone. 2023.  $DM^2$ : Decentralized Multi-Agent Reinforcement Learning via Distribution Matching. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence (AAAI '23)*. AAAI, Washington, DC, USA, 11699–11707.
  - [422] Dakuo Wang, Liuping Wang, Zhan Zhang, Ding Wang, Haiyi Zhu, Yvonne Gao, Xiangmin Fan, and Feng Tian. 2021. “Brilliant AI Doctor” in Rural Clinics: Challenges in AI-Powered Clinical Decision Support System Deployment. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*. ACM, Yokohama, Japan, 1–18.
  - [423] Hai Wang, Yijie Yu, Yingfeng Cai, Xiaobo Chen, Long Chen, and Qingchao Liu. 2019. A Comparative Study of State-of-the-Art Deep Learning Algorithms for Vehicle Detection. *IEEE Intelligent Transportation Systems Magazine* 11, 2 (2019), 82–95.
  - [424] Maonan Wang, Aoyu Pang, Yuheng Kan, Man-On Pun, Chung Shue Chen, and Bo Huang. 2024. LLM-Assisted Light: Leveraging Large Language Model Capabilities for Human-Mimetic Traffic Signal Control in Complex Urban Environments. arXiv:2403.08337
  - [425] Xiaoyu Wang, Baher Abdulhai, and Scott Sanner. 2023. A Critical Review of Traffic Signal Control and A Novel Unified View of Reinforcement Learning and Model Predictive Control Approaches for Adaptive Traffic Signal Control. In *Handbook on Artificial Intelligence and Transport*. Edward Elgar, Northampton, USA, Chapter 17, 482–532.
  - [426] Xinzhi Wang, Huao Li, Hui Zhang, Michael Lewis, and Katia Sycara. 2021. Explanation of Reinforcement Learning Model in Dynamic Multi-Agent System. arXiv:2008.01508
  - [427] Xinru Wang and Ming Yin. 2021. Are Explanations Helpful? A Comparative Study of

- the Effects of Explanations in AI-Assisted Decision-Making. In *Proceedings of the 26th International Conference on Intelligent User Interfaces (IUI '21)*. ACM, College Station, USA, 318–328.
- [428] Xinru Wang and Ming Yin. 2023. Watch Out for Updates: Understanding the Effects of Model Explanation Updates in AI-Assisted Decision Making. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–19.
  - [429] Yinquan Wang, Jianjun Wu, Huijun Sun, Ying Lv, and Junyi Zhang. 2024. Promoting Collaborative Dispatching in the Ride-Sourcing Market With a Third-Party Integrator. *IEEE Transactions on Intelligent Transportation Systems* 25, 7 (2024), 6889–6901.
  - [430] Christopher J. C. H. Watkins and Peter Dayan. 1992. Q-learning. *Machine Learning* 8 (1992), 279–292.
  - [431] Tammy Webber and Michael Phillis. 2025. *US infrastructure improved with Biden-era spending but there's a long way to go*. AP News. <https://apnews.com/article/infrastructure-report-card-civil-engineers-365edbaa0118b4f6af360f7b908a1d9e>
  - [432] Hua Wei, Chacha Chen, Kan Wu, Guanjie Zheng, Zhengyao Yu, Vikash Gayah, and Zhenhui Li. 2019. Deep Reinforcement Learning for Traffic Signal Control along Arterials. In *Proceedings of the 1st Workshop on Deep Reinforcement Learning for Knowledge Discovery (DRL4KDD '19)*. ACM, Anchorage, USA, 1–7.
  - [433] Hua Wei, Chacha Chen, Guanjie Zheng, Kan Wu, Vikash Gayah, Kai Xu, and Zhenhui Li. 2019. PressLight: Learning Max Pressure Control to Coordinate Traffic Signals in Arterial Network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '19)*. ACM, Anchorage, USA, 1290–1298.
  - [434] Hua Wei, Nan Xu, Huichu Zhang, Guanjie Zheng, Xinshi Zang, Chacha Chen, Weinan Zhang, Yanmin Zhu, Kai Xu, and Zhenhui Li. 2019. CoLight: Learning Network-level Cooperation for Traffic Signal Control. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM '19)*. ACM, Beijing, China, 1913–1922.
  - [435] Hua Wei, Guanjie Zheng, Vikash Gayah, and Zhenhui Li. 2019. A Survey on Traffic Signal Control Methods. arXiv:1904.08117
  - [436] Hua Wei, Guanjie Zheng, Vikash Gayah, and Zhenhui Li. 2021. Recent Advances in Reinforcement Learning for Traffic Signal Control: A Survey of Models and Evaluation. *ACM SIGKDD Explorations Newsletter* 22, 2 (2021), 12–18.
  - [437] Hua Wei, Guanjie Zheng, Huaxiu Yao, and Zhenhui Li. 2018. IntelliLight: A Reinforcement Learning Approach for Intelligent Traffic Light Control. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '18)*. ACM, London, UK, 2496–2505.
  - [438] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H.

- Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS '22)*. NeurIPS, New Orleans, USA, 24824–24837.
- [439] Keji Wei, Vikrant Vaze, and Alexandre Jacquillat. 2021. Transit Planning Optimization Under Ride-Hailing Competition and Traffic Congestion. *Transportation Science* 56, 3 (2021), 725–749.
- [440] Tianhao Wei and Changliu Liu. 2022. Safe Control with Neural Network Dynamic Models. In *Proceedings of The 4th Annual Learning for Dynamics and Control Conference (L4DC '22)*. PMLR, Palo Alto, USA, 739–750.
- [441] Herbert Weinblatt, Erik Minge, and Scott Petersen. 2013. Length-Based Vehicle Classification Schemes and Length Bin Boundaries. *Transportation Research Record* 2339, 1 (2013), 19–29.
- [442] Xiao Wen, Yuanchang Xie, Liming Jiang, Yan Li, and Tingjian Ge. 2022. On the interpretability of machine learning methods in crash frequency modeling and crash modification factor development. *Accident Analysis & Prevention* 168 (2022), 106617.
- [443] Tom Wenzel, Clement Rames, Eleftheria Kontou, and Alejandro Henao. 2019. Travel and energy implications of ridesourcing service in Austin, Texas. *Transportation Research Part D: Transport and Environment* 70 (2019), 18–34.
- [444] David Gray Widder, Laura Dabbish, James D. Herbsleb, Alexandra Holloway, and Scott Davidoff. 2021. Trust in Collaborative Automation in High Stakes Software Engineering Work: A Case Study at NASA. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*. ACM, Yokohama, Japan, 1–13.
- [445] Rainer Wiedemann. 1974. Simulation des straßenverkehrsflusses. *Publications of the Institute for Transportation, University of Karlsruhe* 8 (1974), 1–42.
- [446] M. Wiering, J. Vreeken, J. van Veenen, and A. Koopman. 2004. Simulation and optimization of traffic in a city. In *Proceedings of the 2004 Intelligent Vehicles Symposium (IV '04)*. IEEE, Parma, Italy, 453–458.
- [447] Philipp Wintersberger, Frederica Janotta, Jakob Peintner, Andreas Löcken, and Andreas Riener. 2021. Evaluating feedback requirements for trust calibration in automated vehicles. *it — Information Technology* 63, 2 (2021), 1–12.
- [448] Alex J Wood, Mark Graham, Vili Lehdonvirta, and Isis Hjorth. 2019. Good Gig, Bad Gig: Autonomy and Algorithmic Control in the Global Gig Economy. *Employment and Society* 33, 1 (2019), 56–75.
- [449] Jingda Wu, Zhiyu Huang, Zhongxu Hu, and Chen Lv. 2023. Toward Human-in-the-Loop AI: Enhancing Deep Reinforcement Learning via Real-Time Human Guidance for Autonomous Driving. *Engineering* 21 (2023), 75–91.
- [450] Lingtao Wu, Dominique Lord, and Yajie Zou. 2015. Validation of Crash Modification Factors Derived from Cross-Sectional Studies with Regression Models. *Transportation Research Record* 2514 (2015), 88–96.

- [451] Qiang Wu, Jianqing Wu, Jun Shen, Bo Du, Akbar Telikani, Mahdi Fahmideh, and Chao Liang. 2022. Distributed agent-based deep reinforcement learning for large scale traffic signal control. *Knowledge-Based Systems* 241 (2022), 108304.
- [452] Qingjun Wu, Hao Zhang, Zhen Li, and Kai Liu. 2019. Labor control in the gig economy: Evidence from Uber in China. *Journal of Industrial Relations* 61, 4 (2019), 574–596.
- [453] Oskar Wysocki, Jessica Katharine Davies, Markel Vigo, Anne Caroline Armstrong, Dónal Landers, Rebecca Lee, and André Freitas. 2023. Assessing the communication gap between AI models and healthcare professionals: Explainability, utility and trust in AI-driven clinical decision-making. *Artificial Intelligence* 316 (2023), 103839.
- [454] Donghan Xie, Zhi Wang, Chunlin Chen, and Daoyi Dong. 2020. IEDQN: Information Exchange DQN with a Centralized Coordinator for Traffic Signal Control. In *Proceedings of the 2020 International Joint Conference on Neural Networks (IJCNN '20)*. IEEE, Glasgow, UK, 1–8.
- [455] Bingyu Xu, Yaowei Wang, Zhaozhi Wang, Huizhu Jia, and Zongqing Lu. 2021. Hierarchically and cooperatively learning traffic signal control. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence (AAAI '21)*. AAAI, Virtual, 669–677.
- [456] Wei Xu, Yanan Jiang, Xiaoli Zhang, Yi Li, Run Zhang, and Guangtao Fu. 2020. Using long short-term memory networks for river flow prediction. *Hydrology Research* 51, 6 (2020), 1358–1376.
- [457] Fumeng Yang, Chloe Rose Mortenson, Erik Nisbet, Nicholas Diakopoulos, and Matthew Kay. 2024. In Dice We Trust: Uncertainty Displays for Maintaining Trust in Election Forecasts Over Time. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI '24)*. ACM, Honolulu, USA, 1–24.
- [458] Huan Yang, Han Zhao, Yu Wang, Guoqiang Liu, and Danwei Wang. 2022. Deep Reinforcement Learning Based Strategy For Optimizing Phase Splits in Traffic Signal Control. In *Proceedings of the 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC '22)*. IEEE, Macau, 2329–2334.
- [459] Jiachen Yang, Igor Borovikov, and Hongyuan Zha. 2020. Hierarchical Cooperative Multi-Agent Reinforcement Learning with Skill Discovery. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '20)*. IFAA-MAS, Auckland, NZ, 1566–1574.
- [460] Qisong Yang, Thiago D. Sim ao, Simon H. Tindemans, and Matthijs T. J. Spaan. 2023. Advanced Transportation Controller Standards Update. *Machine Learning* 112 (2023), 859–887.
- [461] Qian Yang, Yuexing Hao, Kexin Quan, Stephen Yang, Yiran Zhao, Volodymyr Kuleshov, and Fei Wang. 2023. Harnessing Biomedical Literature to Calibrate Clinicians' Trust in AI Decision Support Systems. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–14.
- [462] Weiran Yao and Sean Qian. 2020. Learning to Recommend Signal Plans under Incidents

- with Real-Time Traffic Prediction. *Transportation Research Record* 2674, 6 (2020), 45–59.
- [463] Kok-Lim Alvin Yau, Junaid Qadir, Hooi Ling Khoo, Mee Hong Ling, and Peter Komisarczuk. 2017. A Survey on Reinforcement Learning Models and Algorithms for Traffic Signal Control. *Comput. Surveys* 50, 3 (2017), 34.
  - [464] Mobin Yazdani, Majid Sarvi, Saeed Asadi Bagloee, Neema Nassir, Jeff Price, and Hossein Parineh. 2023. Intelligent vehicle pedestrian light (IVPL): A deep reinforcement learning approach for traffic signal control. *Transportation Research Part C: Emerging Technologies* 149 (2023), 103991.
  - [465] Bingquan Yu, Jinqiu Guo, Qinpei Zhao, Jiangfeng Li, and Weixiong Rao. 2020. Smarter and Safer Traffic Signal Controlling via Deep Reinforcement Learning. In *Proceedings of the 29th ACM International Conference on Information and Knowledge Management (CIKM '20)*. ACM, Virtual, 3345–3348.
  - [466] Haoxiang Yu, Vaskar Raychoudhury, and Shrawani Silwal. 2020. Dynamic Taxi Ride Sharing using Localized Communication. In *Proceedings of the 21st International Conference on Distributed Computing and Networking (ICDCN '20)*. ACM, Kolkata, India, 1–10.
  - [467] Kun Yu, Shlomo Berkovsky, Ronnie Taib, Jianlong Zhou, and Fang Chen. 2019. Do I trust my machine teammate? An investigation from perception to decision. In *Proceedings of the 24th International Conference on Intelligent User Interfaces (IUI '19)*. ACM, Los Angeles, USA, 460–468.
  - [468] Peipei Yu, Hongcai Zhang, Yonghua Song, Zhenyi Wang, Huiyu Dong, and Liang Ji. 2025. Safe reinforcement learning for power system control: A review. *Renewable and Sustainable Energy Reviews* 223 (2025), 116022.
  - [469] Jinghui Yuan, Mohamed Abdel-Aty, Ling Wang, Jaeyoung Lee, Xuesong Wang, and Rongjie Yu. 2018. Real-Time Crash Risk Analysis of Urban Arterials Incorporating Bluetooth, Weather, and Adaptive Signal Control Data. In *Proceedings of the 97th Annual Meeting of the Transportation Research Board (TRB '18)*. TRB, Washington, DC, USA, 1–8.
  - [470] Rui Yue, Guangchuan Yang, Yichen Zheng, Yuxin Tian, and Zong Tian. 2022. Effects of traffic signal coordination on the safety performance of urban arterials. *Computational Urban Science* 2, 3 (2022), 1–13.
  - [471] Mireia Yurrita, Tim Draws, Agathe Balayn, Dave Murray-Rust, Nava Tintarev, and Alessandro Bozzon. 2023. Disentangling Fairness Perceptions in Algorithmic Decision-Making: the Effects of Explanations, Human Oversight, and Contestability. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–21.
  - [472] Renos Zabounidis, Joseph Campbell, Simon Stepputtis, Dana Hughes, and Katia P. Sycara. 2023. Concept Learning for Interpretable Multi-Agent Reinforcement Learning. In *Proceedings of the 7th Annual Conference on Robot Learning (CoRL '23)*. CoRL, Atlanta, USA, 1828–1837.
  - [473] Xinshi Zang, Huaxiu Yao, Guanjie Zheng, Nan Xu, Kai Xu, and Zhenhui Li. 2020. MetaL-

- ight: Value-Based Meta-Reinforcement Learning for Traffic Signal Control. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI '20)*. AAAI, New York, USA, 1153–1160.
- [474] Jing Zeng, Jie Xin, Ya Cong, Jiancong Zhu, Yihao Zhang, Weihao Jiang, and Shiliang Pu. 2022. HALight: Hierarchical Deep Reinforcement Learning for Cooperative Arterial Traffic Signal Control with Cycle Strategy. In *Proceedings of the 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC '22)*. IEEE, Macau, 479–485.
  - [475] Zheng Zeng. 2021. GraphLight: Graph-based Reinforcement Learning for Traffic Signal Control. In *Proceedings of the 6th International Conference on Computer and Communication Systems (ICCCS '21)*. IEEE, Las Vegas, USA, 645–650.
  - [476] Angie Zhang, Alexander Boltz, Jonathan Lynn, Chun-Wei Wang, and Min Kyung Lee. 2023. Stakeholder-Centered AI Design: Co-Designing Worker Tools with Gig Workers through Data Probes. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM, Hamburg, Germany, 1–18.
  - [477] Angie Zhang, Alexander Boltz, Chun-Wei Wang, and Min Kyung Lee. 2022. Algorithmic Management Reimagined For Workers and By Workers: Centering Worker Well-Being in Gig Work. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. ACM, New Orleans, USA, 1–20.
  - [478] Dongping Zhang, Jason Hartline, and Jessica Hullman. 2024. Designing Shared Information Displays for Agents of Varying Strategic Sophistication. In *Proceedings of the 27th ACM Conference on Computer-Supported Cooperative Work and Social Computing (CSCW '24)*. ACM, San José, Costa Rica, 1–34.
  - [479] Gongquan Zhang, Fangrong Chang, Helai Huang, and Zilong Zhou. 2024. Dual-Objective Reinforcement Learning-Based Adaptive Traffic Signal Control for Decarbonization and Efficiency Optimization. *Mathematics* 12, 13 (2024), 2056.
  - [480] Gongquan Zhang, Fangrong Chang, Jieling Jin, Fan Yang, and Helai Huang. 2024. Multi-objective deep reinforcement learning approach for adaptive traffic signal control system with concurrent optimization of safety, efficiency, and decarbonization at intersections. *Accident Analysis and Prevention* 199 (2024), 107451.
  - [481] Gongquan Zhang, Jieling Jin, Fangrong Chang, and Helai Huang. 2024. Real-time traffic conflict prediction at signalized intersections using vehicle trajectory data and deep learning. *International Journal of Transportation Science and Technology* In Press (2024), 1–15.
  - [482] Huichu Zhang, Siyuan Feng, Chang Liu, Yaoyao Ding, Yichen Zhu, Zihan Zhou, Weinan Zhang, Yong Yu, Haiming Jin, and Zhenhui Li. 2019. CityFlow: A Multi-Agent Reinforcement Learning Environment for Large Scale City Traffic Scenario. In *Proceedings of the 2019 World Wide Web Conference (WWW '19)*. ACM, New York, USA, 3620–3624.
  - [483] Huichu Zhang, Chang Liu, Weinan Zhang, Guan jie Zheng, and Yong Yu. 2020. GeneralLight: Improving Environment Generalization of Traffic Signal Control via Meta Reinforcement Learning. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (CIKM '20)*. ACM, Virtual, 1783–1792.



- [484] Jingjing Zhang, Gediminas Adomavicius, Alok Gupta, and Wolfgang Ketter. 2020. Consumption and Performance: Understanding Longitudinal Dynamics of Recommender Systems via an Agent-Based Simulation Framework. *Information Systems Research* 31, 1 (2020), 76–101.
- [485] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. 2021. Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms. In *Handbook of Reinforcement Learning and Control*. Springer, New York, USA, 321–384.
- [486] Lun Zhang, Shan Jiang, and Zheng Wang. 2017. Schedule-Driven Signal Priority Control for Modern Trams Using Reinforcement Learning. In *Proceedings of the 17th COTA International Conference of Transportation Professionals (CICTP '17)*. ASCE, Shanghai, China, 2122–2132.
- [487] Qiaoning Zhang, Matthew L. Lee, and Scott Carter. 2022. You Complete Me: Human-AI Teams and Complementary Expertise. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. ACM, New Orleans, USA, 1–28.
- [488] Tingru Zhang, Da Tao, Xingda Qu, Xiaoyan Zhang, Jihong Zeng, Haoyu Zhu, and Han Zhu. 2020. Automated vehicle acceptance in China: Social influence and initial trust are key determinants. *Transportation Research Part C: Emerging Technologies* 112 (2020), 220–233.
- [489] Yunchang Zhang and Jon Fricker. 2021. Investigating Smart Traffic Signal Controllers at Signalized Crosswalks: A Reinforcement Learning Approach. In *Proceedings of the 7th International Conference on Models and Technologies for Intelligent Transportation Systems (MT-ITS '21)*. IEEE, Heraklion, Greece, 1–6.
- [490] Yifeng Zhang, Harsh Goel, Peizhuo Li, Mehul Damani, Sandeep Chinchali, and Guillaume Sartoretti. 2025. CoordLight: Learning Decentralized Coordination for Network-Wide Traffic Signal Control. *IEEE Transactions on Intelligent Transportation Systems* Early Access (2025), 1–16.
- [491] Yunfeng Zhang, Q. Vera Liao, and Rachel K. E. Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAccT '20)*. ACM, Barcelona, Spain, 295–305.
- [492] Yiran Zhang, Khoa Vo, Longchao Da, Tiejun Chen, Xiaoou Liu, and Hua Wei. 2025. Poster Abstract: Reproducible and Low-Cost Sim-to-Real Environment for Traffic Signal Control. In *Proceedings of the 16th ACM/IEEE International Conference on Cyber-Physical Systems (ICCPS '25)*. ACM, Irvine, USA, 1–2.
- [493] Pengqian Zhao, Yuyu Yuan, and Ting Guo. 2022. Extensible Hierarchical Multi-Agent Reinforcement-Learning Algorithm in Traffic Signal Control. *Applied Sciences* 12, 24 (2022), 12783.
- [494] Weiye Zhao, Tairan He, and Changliu Liu. 2023. Probabilistic Safeguard for Reinforcement Learning Using Safety Index Guided Gaussian Process Models. In *Proceedings of The 5th Annual Learning for Dynamics and Control Conference (LADC '23)*. PMLR, Philadelphia,

USA, 783–796.

- [495] Guanjie Zheng, Yuanhao Xiong, Xinshi Zang, Jie Feng, Hua Wei, Huichu Zhang, Yong Li, Kai Xu, and Zhenhui Li. 2021. Learning Phase Competition for Traffic Signal Control. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM '19)*. ACM, Beijing, China, 1963–1972.
- [496] Guanjie Zheng, Xinshi Zang, Nan Xu, Hua Wei, Zhengyao Yu, Vikash Gayah, Kai Xu, and Zhenhui Li. 2019. Diagnosing Reinforcement Learning for Traffic Signal Control. arXiv:1905.04716
- [497] Zuduo Zheng. 2014. Recent developments and research needs in modeling lane changing. *Transportation Research Part B: Methodological* 60 (2014), 16–32.
- [498] Bin Zhou, Qishen Zhou, Simon Hu, Dongfang Ma, Sheng Jin, and Der-Horng Lee. 2024. Cooperative Traffic Signal Control Using a Distributed Agent-Based Deep Reinforcement Learning With Incentive Communication. *IEEE Transactions on Intelligent Transportation Systems* 25, 8 (2024), 10147–10160.
- [499] Meizi Zhou, Jingjing Zhang, and Gediminas Adomavicius. 2021. *Longitudinal Impact of Preference Biases on Recommender Systems' Performance*. Technical Report 2021-10. Kelley School of Business Research Paper.
- [500] Xinyu Zhou, ZhuHua Liao, Yijiang Zhao, Yizhi Liu, and Aiping Yi. 2025. Ride-hailing pick-up area recommendation in a vehicle-cloud collaborative environment: a feature-aware personalized clustering federated learning approach. *Cluster Computing* 28 (2025), 32.
- [501] Hong Zhu, Fengmei Sun, Keshuang Tang, Tianyang Han, and Junping Xiang. 2024. A Coordination Graph Based Framework for Network Traffic Signal Control. *IEEE Transactions on Intelligent Transportation Systems* 25, 10 (2024), 14298–14312.
- [502] Yuanyang Zhu, Xiao Yin, and Chunlin Chen. 2022. Extracting decision tree from trained deep reinforcement learning in traffic signal control. *IEEE Transactions on Computational Social Systems* 10, 4 (2022), 1997–2007.
- [503] Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong Hoang, Armel Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, Binyuan Hui, Niklas Muennighoff, David Lo, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro von Werra. 2025. BigCodeBench: Benchmarking Code Generation with Diverse Function Calls and Complex Instructions. In *Proceedings of the 2025 International Conference on Learning Representations (ICLR '25)*. ICLR, Singapore, 1–55.
- [504] Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-Tuning Language Models from Human Preferences. arXiv:1909.08593