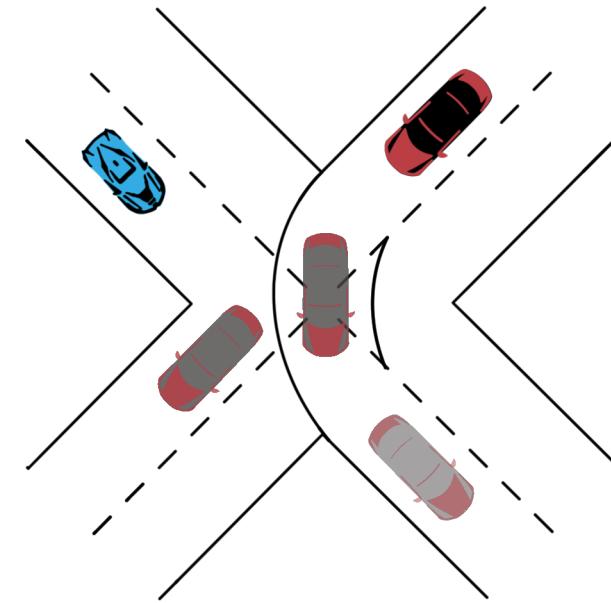




To Be Written

TOMMY TRAM • Learning When To Drive in Uncertain Environments • 2022



# Learning When To Drive in Uncertain Environments

A Reinforcement Learning approach



TOMMY TRAM



DEPARTMENT OF ELECTRICAL ENGINEERING

CHALMERS UNIVERSITY OF TECHNOLOGY  
Gothenburg, Sweden 2022  
[www.chalmers.se](http://www.chalmers.se)

THESIS FOR THE DEGREE OF DOCTOR OF PHILOSOPHY

---

# Learning When To Drive in Uncertain Environments

*A reinforcement learning approach*

TOMMY TRAM

Department of Electrical Engineering  
Chalmers University of Technology  
Gothenburg, Sweden, 2022

# **Learning When To Drive in Uncertain Environments**

*A reinforcement learning approach*

TOMMY TRAM

ISBN 978-91-7905-623-0

© 2022 TOMMY TRAM

All rights reserved.

Doktorsavhandlingar vid Chalmers tekniska högskola

Ny serie nr 5089

ISSN 0346-718X

Department of Electrical Engineering

Chalmers University of Technology

SE-412 96 Gothenburg, Sweden

Phone: +46 (0)31 772 1000

Email: [tommy.tram@chalmers.se](mailto:tommy.tram@chalmers.se); [tommy.tram@gmail.com](mailto:tommy.tram@gmail.com)

Cover:

An illustration of two vehicles approaching an intersection. The blue vehicle has to make a decision to yield or drive while the intention of the red vehicle is uncertain.

Printed by Chalmers Reproservice

Gothenburg, Sweden, December 2022

*To my loving family*



# **Learning When To Drive in Uncertain Environments**

*A reinforcement learning approach*

TOMMY TRAM

Department of Electrical Engineering  
Chalmers University of Technology

## **Abstract**

To Be Written

**Keywords:** Autonomous driving, reinforcement learning, uncertain environments, deep Q-learning 



## List of Publications

This thesis is based on the following publications:

- [A] Tommy Tram, Anton Jansson, Robin Grönberg, Mohammad Ali, and Jonas Sjöberg, “Learning Negotiating Behavior Between Cars in Intersections using Deep Q-Learning”. *Published in 2018 21st International Conference on Intelligent Transportation Systems (ITSC)*.
- [B] Tommy Tram, Ivo Batković, Mohammad Ali, and Jonas Sjöberg, “Learning When to Drive in Intersections by Combining Reinforcement Learning and Model Predictive Control”. *Published in 2019 IEEE Intelligent Transportation Systems Conference (ITSC)*.
- [C] Carl-Johan Hoel, Tommy Tram, and Jonas Sjöberg, “Reinforcement Learning with Uncertainty Estimation for Tactical Decision-Making in Intersections”. *Published in 2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*.
- [D] Tommy Tram, Maxime Bouton, Jonas Sjöberg, and Mykel Kochenderfer, “Belief State Reinforcement Learning for Autonomous Vehicles in Intersections”. *Published in TBD.* 
- [E] Hannes Eriksson, Tommy Tram, Debabrota Basu, Jonas Sjöberg, and Christos Dimitrakakis, “Maximum Likelihood Estimation for Transfer Reinforcement Learning Problems in Autonomous Driving”. *Artificial Intelligence and Statistics 2023 (AISTATS)*.



## Acknowledgments

Jonas Sjoberg

Volvo Cars, Zenuity, Zenseact. Mohammad Ali, Beliving in me and chosing me for this phd project Carl Lindberg coaching, allowing me to work close to my family during though times. och Mats NordLund. Mats Lestander. Managers enabling me to do this research and giving me all the resources.

Carl-Johan Hoel Ivo Batkovic Arian Ranjbar Maxime Bouton Anton och robin, Started as a master thesis student ended up as friends Hannes and debba

WASP friends and colleagues Krook, John, Daniel, Dapeng, Erik, Emil, Anton, Rui, Lars, Linnea, Oskar, Christian, Viktor, and the rest of my colleagues from Chalmers, Zenseact, and WASP.

SISL Mykel Kochenderfer Mattias, shushman, Jayesh, Jeremy, Ayan Wallenberg stanford

*Tommy Tram,  
Göteborg, December, 2022.*



## **Acronyms**

AD:	Autonomous driving
ADAS:	Advanced driving assistance systems
AV:	Autonomous vehicles
CNN:	Convolutional neural network
DQN:	Deep Q-network
EQN:	Ensabmle quantile networks
IDM:	Intelligent driver
LSTM:	Long short-term memory
MDP:	Markov descision process
MPC:	Model predictive control
POMDP:	Partially observable Markov descision process
RL:	Reinforcement learning





---

## Contents

---

<b>Abstract</b>	i
<b>List of Papers</b>	iii
<b>Acknowledgements</b>	v
<b>Acronyms</b>	vii
<b>I Overview</b>	1
<b>1 Introduction</b>	3
1.1 Problem formulation . . . . .	4
1.2 Scope and limitations . . . . .	4
1.3 System architecture . . . . .	5
1.4 Contributions . . . . .	5
1.5 Thesis outline . . . . .	6
<b>2 Related work</b>	7
2.1 Learning-based methods . . . . .	9
	ix

<b>3 Technical background</b>	 13
3.1 Partially observable Markov decision process . . . . .	13
3.1.1 Partially observable Markov decision process . . . . .	14
3.2 Markov decision processes . . . . .	14
3.3 Reinforcement learning . . . . .	16
<b>4 Model-free RL approaches</b> 	<b>19</b>
4.1 State representation . . . . .	19
4.1.1 Observable states . . . . .	19
 4.1.2 Unobservable states . . . . .	19
4.2 Approach . . . . .	20
4.3 Simulated experiments 	20
4.4 Results and discussion . . . . .	20
<b>5 Combining MPC and RL</b> 	<b>21</b>
 5.1 Action space, options . . . . .	21
 5.2 MPC . . . . .	21
 5.3 Approach . . . . .	21
5.4 Simulated experiments . . . . .	22
5.5 Results and discussion . . . . .	22
<b>6 Estimating the uncertainty</b>	<b>23</b>
6.1 Uncertainty of the decision . . . . .	23
6.1.1 Approach . . . . .	23
6.1.2 Simulated experiments . . . . .	23
6.1.3 Results and discussion . . . . .	23
6.2 Uncertrainty of the intention . . . . .	23
6.2.1 Approach . . . . .	23
6.2.2 Simulated experiments . . . . .	23
6.2.3 Results and discussion . . . . .	23
<b>7 Generalize over different scenarios</b> 	<b>25</b>
7.1 Approach . . . . .	25
7.2 Simulated experiments . . . . .	25
7.3 Results and discussion . . . . .	25
<b>8 Discussion</b> 	<b>27</b>

<b>9 Concluding remarks and future work</b>	<b>29</b>
9.1 Future work 	29
<b>10 Summary of included papers</b>	<b>31</b>
 10.1 Paper A . . . . .	31
10.2 Paper B . . . . .	32
10.3 Paper C . . . . .	32
10.4 Paper D . . . . .	33
10.5 Paper E . . . . .	34
<b>References</b>	<b>35</b>



<b>II Papers</b>	<b>37</b>
<b>A Learning Negotiating Behavior Between Cars in Intersections using Deep Q-Learning</b>	<b>A1</b>
1 Introduction . . . . .	A3
2 Overview . . . . .	A4
3 Problem formulation . . . . .	A5
3.1 System architecture . . . . .	A5
3.2 Actions as Short Term Goals . . . . .	A5
3.3 Observations that make up the state . . . . .	A7
3.4 Partially Observable Markov Decision Processes . . . . .	A8
4 Finding the optimal policy . . . . .	A8
5 Method . . . . .	A9
5.1 Deep Q-learning . . . . .	A9
5.2 Experience Replay . . . . .	A10
5.3 Dropout . . . . .	A10
5.4 Long short term memory . . . . .	A10
6 Implementation . . . . .	A11
6.1 Simulation environment . . . . .	A11
6.2 Reward function tuning . . . . .	A12
6.3 Neural Network Setup . . . . .	A12
7 Results . . . . .	A14
7.1 Effect of using Experience replay . . . . .	A15
7.2 Effect of using Dropout . . . . .	A16

7.3	Comparing DQN and DRQN . . . . .	A16
7.4	Effect of sharing weights in the network . . . . .	A16
8	Conclusion . . . . .	A17
<b>B Learning When to Drive in Intersections by Combining Reinforcement Learning and Model Predictive Control</b>		<b>B1</b>
1	Introduction . . . . .	B3
2	System . . . . .	B5
3	Problem formulation . . . . .	B6
3.1	Partially Observable Markov Decision Process . . . . .	B6
3.2	Deep Q-Learning . . . . .	B7
4	Agents . . . . .	B7
4.1	MPC agent . . . . .	B7
4.2	Sliding mode agent . . . . .	B11
4.3	Intention agents . . . . .	B12
5	Implementation . . . . .	B12
5.1	Deep Q-Network . . . . .	B12
5.2	Q-masking . . . . .	B14
5.3	Simulation environment . . . . .	B14
5.4	Reward function tuning . . . . .	B15
6	Results . . . . .	B17
7	Discussion . . . . .	B17
8	Conclusion . . . . .	B19
<b>C Reinforcement Learning with Uncertainty Estimation for Tactical Decision-Making in Intersections</b>		<b>C1</b>
1	Introduction . . . . .	C3
2	Approach . . . . .	C5
2.1	Reinforcement learning . . . . .	C5
2.2	Bayesian reinforcement learning . . . . .	C6
2.3	Confidence criterion . . . . .	C7
3	Implementation . . . . .	C8
3.1	Simulation setup . . . . .	C8
3.2	MDP formulation . . . . .	C9
3.3	Fallback action . . . . .	C11
3.4	Network architecture . . . . .	C12
3.5	Training process . . . . .	C12

3.6	Baseline method . . . . .	C13
4	Results . . . . .	C13
4.1	Within training distribution . . . . .	C14
4.2	Outside training distribution . . . . .	C17
5	Discussion . . . . .	C17
6	Conclusion . . . . .	C20
	References . . . . .	C20

## **D Belief State Reinforcement Learning for Autonomous Vehicles in Intersections**

		<b>D1</b>
1	Introduction . . . . .	D3
2	Background . . . . .	D6
2.1	Partially observable Markov decision process . . . . .	D6
2.2	Driver model . . . . .	D7
3	Proposed approach . . . . .	D8
3.1	Problem formulation . . . . .	D8
3.2	Belief state estimation using a particle filter . . . . .	D11
3.3	Belief state reinforcement learning . . . . .	D13
4	Experiments . . . . .	D16
4.1	Simulator setup . . . . .	D16
4.2	Neural network architecture . . . . .	D18
4.3	Training procedure . . . . .	D20
5	Results / Results and discussion . . . . .	D21
5.1	Baseline DQN results . . . . .	D22
5.2	QMDP and QMDP-IE results . . . . .	D23
5.3	QPF and QID results . . . . .	D24
5.4	Discussion . . . . .	D25
6	Conclusion . . . . .	D25
	References . . . . .	D26

## **E Maximum Likelihood Estimation for Transfer Reinforcement Learning Problems in Autonomous Driving**

		<b>E1</b>
1	Introduction . . . . .	E3
1.1	Related Work . . . . .	E4
1.2	Contribution . . . . .	E4
2	Background . . . . .	E5
2.1	Markov Decision Process . . . . .	E5

2.2	Likelihoods . . . . .	E6
2.3	Feature Representations . . . . .	E6
3	Transfer Reinforcement Problem Formulation . . . . .	E6
3.1	Finite and realisable plausible models . . . . .	E7
3.2	Infinite and realisable plausible models . . . . .	E7
3.3	Infinite and non-realisable bounded plausible models . . . . .	E8
3.4	Infinite and arbitrary plausible models . . . . .	E9
4	Maximum Likelihood Estimation of MDPs in Transfer Reinforcement Learning . . . . .	E9
4.1	Discrete MDPs . . . . .	E9
4.2	Linear MDPs . . . . .	E12
5	Planning in Markov Decision Processes . . . . .	E13
5.1	Discrete MDPs . . . . .	E13
5.2	Linear MDPs . . . . .	E13
6	Algorithms . . . . .	E14
7	Bounds of TRLBI . . . . .	E16
7.1	Value Estimation Error . . . . .	E18
7.2	Model Estimation Error . . . . .	E18
7.3	Realizability Error . . . . .	E18
8	Experiments . . . . .	E18
8.1	Discrete MDPs . . . . .	E18
8.2	Continous Space and Action . . . . .	E18
9	Autonomous Driving Scenarios with Knowledge Transfer . . . . .	E19
9.1	Setting 1: Transfer knowledge from similar tasks, e.g. cut-in to overtake . . . . .	E20
9.2	Setting 2: Transfer knowledge from similar cars, e.g. sedan to combi . . . . .	E21
9.3	Setting 3a: Transfer knowledge from similar driving locations, e.g. USA to China, using traffic rules, left-side / right-side traffic . . . . .	E22
9.4	Setting 3b: Transfer knowledge from similar driving locations, e.g. USA to China, using driver interaction . . . . .	E23

# **Part I**

# **Overview**



# CHAPTER 1

---

## Introduction

---

The way of transportation is currently evolving and autonomous driving technology is expected to have a big impact on this transformation. Over one million people are killed in traffic-related accidents each year, where the vast majority of the accidents are caused by human mistakes [1], [2]. By helping humans with perception, prediction and decision making autonomous driving could significantly improve traffic safety and make way for new innovative road infrastructure. Without the requirement of a present human for each transportation vehicle, the efficiency of traffic can be improved by scheduling commercial transports outside of rush hours [3].

**ToDo:** cite future of mixed-autonomy

**ToDo:** level 2: volvo pilot assist and tesla auto pilot. Level 4: Waymo.  
Woven automated city

Thanks to the rapid success of deep learning during the last decades, major progress towards deploying autonomous vehicles in the real world.

**Tommy:** rewrite

Major progress towards deploying autonomous vehicles has been made during the last decade. The perception systems have been remarkably improved, largely due to the success of deep learning techniques [4]. The low-level control of the vehicle is a mature research area and can be solved with classical control theory methods [5]. However, how to approach the high-level decision-making in complex traffic situations is less explored and forms the main topic of this thesis.

## 1.1 Problem formulation

The work presented in this thesis investigate the following research questions:

 **Q1.** How can RL be used to create a decision-making agent for autonomous driving, that can handle different intersection and roundabouts (complex urban scenarios)? Learn a scalable policy that is able to handle different scenarios. relative coordinate system. Action space. (specificera for att komma undan varför har du inte kollat pa andra metoder. How can we use RL for AD )

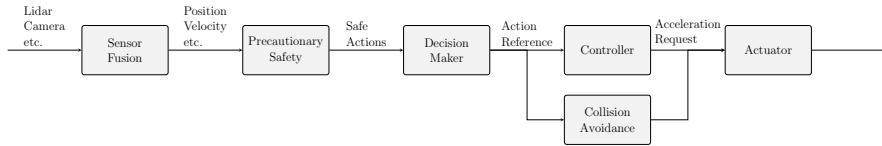
**Q2.** How can AD domain knowlage (and models) be used to improve the action and state space for a RL agent? MPC for actions, Particle filter for intention distrubution. How can AD domain knowlage be used to create a state and action space that improves the RL agent?

**Q3.** How can the quality of a RL agent be imporved by accounting for uncertainty? (How can the uncertainty of the RL agent be utilized?) (RPF-in the output and PF-in the input space)

1. We want to drive through intersections.
2. The intersection can be of different shapes. We assume we have a map of the intersection.
3. There will be other cars crossing the same intersection.

## 1.2 Scope and limitations

FILL



**Figure 1.1:** Representation of the system architecture.

1. We have access to sensors on-board the ego vehicle.
2. We don't assume any knowledge of traffic signs or traffic lights.
3. We don't have v2v, or v2x communication.

To compensate for not having v2v or v2x communication, we have to predict what other driver will do.



## 1.3 System architecture

outline system architecture and limit this paper to comfortable decisions reinforcement learning (RL)

## 1.4 Contributions

**ToDo:** rewrite once thesis is in a better state

The main contributions of this thesis are:

1. General approach to creating a decision making agent for driving in interactions.
2. A neural network architecture that is invariant to permutations of the order of which surrounding traffic participants are observed, which speeds up training and improves the quality of the trained agent.
3. A belief state representation using a particle filter and a comparison and analysis of different algorithms that utilize the belief state.
4. Two approaches to solving a POMDP with hidden intention state. LSTM layer and belief state.

5. General state space representation that is invariant to permutations of the intersection design.
6. Extension of RL methods that provide an estimate of the epistemic uncertainty and use it to create a confidence criteria that can identify situations with high uncertainty.

## **1.5 Thesis outline**

FILL

# CHAPTER 2

---

## Related work

---

**ToDo:** Urban challenge, winner Carnegie Mellon University John Dolan

**ToDo:** Rule based methods

**ToDo:** Motion planning, Predicting motion of surrounding vehicles, reactive (not interactive)

MPC cite ivo: A robust scenario MPC approach for uncertain multi-modal obstacles, Real-Time Constrained Trajectory Planning and Vehicle Control for Proactive Autonomous Driving With Road Users.

**ToDo:** POMDP model, online and offline solvers

Cite Maxime: Cooperation-aware reinforcement learning for merging in dense traffic, Belief state planning for autonomously navigating urban intersections.

**ToDo:** online solvers, MCTS

MCTS cite CJ: Combining planning and deep reinforcement learning in tactical decision making for autonomous driving, Tactical decision-making in autonomous driving by reinforcement learning with uncertainty estimation.

**ToDo:** cite Bo wahlberg, Morteza Haghir Chehreghani, Fernandez Llorca David, Christian Berge

**Tommy:** rewrite text

Early approaches to tactical decision-making for autonomous vehicles often used rule-based methods, commonly implemented as handcrafted state machines. For example, during the DARPA Urban Challenge, a rule-based approach was adopted by the winning team from the Carnegie Mellon University, where different modules handled the behavior for the different driving scenarios that were encountered **darpaCMU**. While successful for a limited and controlled environment such as the Urban Challenge event, it is unlikely that rule-based approaches could scale to handle the complexity and diversity of real-world driving.

Another group of algorithms treats the decision-making task as a motion planning problem. Commonly, a prediction model is used to predict the motion of the other agents, and then the behavior of the vehicle that is being controlled, henceforth referred to as the ego vehicle, is planned accordingly. This results in a reactive behavior, since the predictions are independent of the ego vehicle plan. Therefore, interaction between the ego vehicle and other agents is not explicitly considered, but may happen implicitly by frequent replanning. Search-based planners typically discretize the state space and then apply Dijkstra's algorithm **Dijkstra1959** or one of the algorithms from the A\* family **Hart1968**. These techniques were also commonly used during the DARPA Urban Challenge **Bacha2008**, **darpaStanford**. However, since real-time performance can be hard to achieve with graph-search algorithms, sampling-based algorithms such as rapidly-exploring random trees **Lavalle1998** have been used for motion planning, e.g., by Karaman et al. **Karaman2011**. A third approach to solve the motion planning problem is to use optimization-based techniques, for example optimal control, which was applied to highway driving scenarios by Werling et al. **Werling2010**. Since human behavior is complex and varies between individuals, some algorithms use a probabilistic prediction as input to the motion planning. This is for example shown in a study by Damerow et al. **Damerow2015**, which aims to minimize the risk during an intersection scenario. Additional approaches to motion planning for autonomous driving are provided in the surveys by González et al. **Gonzales2016** and Paden et al. [5].

It is common to model decision-making problems as Markov decision processes or partially observable Markov decision processes (POMDPs) [6]. This mathematical framework allows modeling of uncertainty in the current state, uncertainty in the future development of the traffic scene, and modeling of an interactive behavior. The task of finding the optimal policy for a POMDP is most often intractable, but many approximate methods exist. One way to group these methods is in offline and online methods. There are powerful offline algorithms for planning in POMDPs, which can solve complex situations. One example is shown by Brechtel et al., which proposes a solution to how measurement uncertainty and occlusions in an intersection can be handled **Brechtel2014**. In their work, an offline planner precomputes the policy by using a state representation that is learned for the specific scenario. A similar approach is adopted by Bai et al. for an intersection scenario **Bai2014**. The main drawback of these offline methods is that they are designed for specific scenarios. Due to the large number of possible real-world scenarios, it is challenging to precalculate a policy that is generally valid.

Online methods compute a policy during execution, which makes them more versatile than offline methods. However, the limited available computational resources require a careful problem formulation and limit the solution quality. Another online approach for solving a POMDP is the family of Monte Carlo tree search algorithms **Browne2012**, which is used by Sunberg et al. to make lane changing decisions on a highway **Sunberg2017**. Furthermore, other drivers' intentions are estimated with a particle filter. A hybrid approach between offline and online planning is pursued in a study by Sonu et al., where a hierarchical decision-making structure is used **Sonu2018**. The decision-making problem is modeled on two levels as MDPs, since full observability is assumed. The high-level MDP is solved offline by value iteration and the low-level MDP is solved online with MCTS.

## 2.1 Learning-based methods

**Tommy:** rewrite

In planning-based methods, different algorithms are used to find a policy that maximizes the utility of the agent's behavior by using a model of the system. However, data-driven approaches are fundamentally different, since with this

class of methods, an agent instead learns how to behave from observing data. This section describes different types of data-driven approaches that have been explored for autonomous driving.

An intuitive approach is to collect data from when an expert is performing a task and then use supervised learning to imitate the behavior of the expert. This method is often referred to as behavioral cloning and was first applied to autonomous driving in the ALVINN project **Pomerleau1989**. Unfortunately, for many cases, behavioral cloning suffers from compounding errors, which refers to the problem when small mistakes gradually push the agent further away from the training distribution, into states from which the agent does not know how to recover. This problem can be mitigated by an active learning approach, where an expert can be queried during the training process **Ross2011**, which for example has been applied to autonomous driving by Kelly et al. **Kelly2019**. An alternative is to synthesize data by perturbing the expert’s driving, which was done in a study by Bansal et al. **Bansal2019**. Generative adversarial imitation learning **Ho2016** provides another method to handle the compounding error problem and has showed promising results in different highway driving scenarios **Kuefeler2017**.

Reinforcement learning is conceptually different from supervised learning, since labeled input-output samples are not available. Instead, the agent learns how to make decisions from interacting with the environment. RL methods are versatile, and have proven successful in various domains, such as playing Atari games **Mnih2015**, in continuous control **Lillicrap2015**, reaching a super human performance in the game of Go **Silver2017**, and beating the best chess computers **Silver2017chess**. One advantage of RL methods, compared to planning based methods, is that a model of the environment is not required, i.e., the transition probabilities between different states are not assumed to be known. Furthermore, many RL methods provide a general framework and an agent could, in theory, learn how to behave correctly in all possible driving situations. During the last few years, many papers have addressed the task of applying RL approaches to autonomous driving. For example, DQN-based agents were used by Isele et al. **Isele2018** for navigating through different intersection scenarios, with varying driver intentions and occlusions. Commonly, a high-level action space is used together with the DQN algorithm. Other studies use policy gradient RL methods to directly control the speed and the steering angle of the vehicle, for example in lane

changing and urban scenarios **Wang2019\_ddpg**, **Chen2019**. Safety of the RL-based agents has been addressed by restricting dangerous actions, either by heuristics **Mukadam2017**, linear temporal logic **Bouton2019**, or using an underlying motion planner with hard constraints **Shalev2016**.

A majority of these studies perform both the training and evaluation in simulated environments, whereas some train the agent in a simulator and then apply the trained agent in the real world **Pan2017**, **Bansal2019**, or for some limited scenarios, the training itself is also performed in the real world **Kendall2019**. Overviews of RL-based studies for autonomous driving are given by Kiran et al. **Kiran2021** and by Ye et al. **Ye2021**.

Both planning-based and RL-based methods use a reward function to find a policy that maximizes the cumulative future reward. However, for complex tasks such as autonomous driving, the design of the reward function is in itself a complicated task. For limited scenarios, the reward function can be manually specified and tuned until the agent finds a desired behavior, which is referred to as reward shaping **Ng1999**. However, for realistic scenarios, the number of possible reward features is massive and how to balance rewards related to, e.g., safety and efficiency is a complex issue. A practical approach is to instead learn the reward function from the behavior of human drivers by inverse reinforcement learning (IRL) **Ng2000**. An IRL approach was for example used by Kuderer et al. to learn the individual preferences of human drivers with different driving styles **Kuderer2015**. Sharifzadeh et al. combined IRL with DQN and Wang et al. used an adversarial IRL approach to simultaneously obtain both the reward function and the policy for different lane-changing scenarios **Sharifzadeh2016**, **Wang2019**. Zhu et al. **Zhu2021** provide an overview of IRL applied to autonomous driving. Except for performing planning and RL, the learned reward function can also be used to predict the behavior of other traffic participants. IRL for predictions was for example used by Ziebart et al. for pedestrians **Ziebart2009**, by Sun et al. for human drivers in intersections **Sun2019**, and by Sadigh for highway driving situations **Sadigh2016**.



# CHAPTER 3

---

## Technical background

---

### 3.1 Partially observable Markov decision process

- POMDP
- State - realative features, scalable to fidderent types of intersection design
- Action - IDM - MPC
- Observation - unknown intentions, red lights, stop signs yield.
- Transistion function - RL
- Reward Function - Goal, Comfort, crash
- discount factor

This chapter briefly introduce the Partially Observable Markov Decision Process (POMDP) framework and reinforcement learning.

### 3.1.1 Partially observable Markov decision process

A POMDP is a mathematical framework used for modeling sequential decision making problems under uncertainty. It is a generalization of the Markov Decision Process (MDP). It consists of a set of states  $\mathcal{S}$ , actions  $\mathcal{A}$ , observations  $\Omega$ , a transition model  $T$ , observation model  $O$ , reward function  $R$ , and a discount factor  $\gamma$ . Together they create the tuple  $(\mathcal{S}, \mathcal{A}, \Omega, T, O, R, \gamma)$  that formally defines a POMDP [6]. While in a state  $s \in \mathcal{S}$  and executing an action  $a \in \mathcal{A}$  the probability of transitioning to a future state  $s'$  is described by the transition model  $T(s' | s, a)$  and the reward  $r$  is given by the reward function  $R(s, a)$ . The agent only has access to partial information contained in its observation  $o$  distributed according to  $\Pr(o | s', a) = O(o, s', a)$ .

**Tommy:** CJ

This chapter provides a brief introduction to Markov decision processes and reinforcement learning. The purpose of the chapter is to summarize the most important concepts and introduce the notation that are used in the subsequent chapters. A comprehensive overview of MDPs and RL is given in the books by Kochenderfer [6] and Sutton and Barto [7], upon which this chapter is based.

## 3.2 Markov decision processes

Sequential decision-making problems in stochastic environments are commonly modeled as MDPs. Importantly, an MDP satisfies the Markov property, which requires that the probability distribution of the next state only depends on the current state and the action taken by the agent, i.e., not on the history of previous states or actions. An MDP is formally defined as the tuple  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$ , which is described in the following list [6, Ch. 4]:

- The state space  $\mathcal{S}$  represents the set of all possible states of the environment. This set could consist of both discrete and continuous states.
- The action space  $\mathcal{A}$  represents the set of all valid actions the agent can take. This set could also consist of both discrete and continuous actions. However, since this thesis focuses on high-level decision-making, only discrete actions are considered here.

- The state transition model  $T(s'|s, a)$  describes the probability that the system transitions to state  $s' \in \mathcal{S}$  from state  $s \in \mathcal{S}$  when action  $a \in \mathcal{A}$  is taken.
- The reward function  $R(s, a)$  returns a scalar reward  $r$  for each state-action pair.
- The discount factor  $\gamma \in [0, 1]$  is a scalar that discounts the value of future rewards. For a finite horizon MDP,  $\gamma$  could also take the value 1.

A policy  $\pi$  is a mapping from a state to an action, which could either be deterministic  $a = \pi(s)$  or probabilistic  $a \sim \pi(a|s)$ . The value of being in a state while following a policy is described by the value function

$$V^\pi(s) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, a_k) | s_0 = s, \pi \right]. \quad (3.1)$$

The goal of the agent is to find a policy which maximizes the value of each state.

In many decision-making problems, the agent does not have direct access to the state of the environment. Such a problem is commonly modeled as a partially observable Markov decision process, which is an extension to the MDP framework that also models state uncertainty. A POMDP is formally defined as the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{O}, T, O, R, \gamma)$ , where the state space, action space, transition model, reward function, and discount factor are defined as for an MDP. A POMDP has two additional components [6, Ch. 6]:

- The observation space  $\mathcal{O}$ , which is the set of possible observations.
- The observation model  $O(o|s', a)$ , which describes the probability of observing  $o \in \mathcal{O}$  in state  $s'$  after action  $a$  has been taken.

Since the agent does not have direct access to the current state in a POMDP, the agent needs to reason about the history of observations and actions. This history is often merged in a belief state  $b$ , which represents a probability distribution over the state space. In this case, the policy is a mapping from beliefs to actions  $\pi(b)$ .

For many real-world problems, it is not possible to represent the probability distributions  $T$  or  $O$  explicitly. For some solution approaches, only samples are needed, and then it is sufficient to define a generative model  $G$  that samples a

new state or observation from a given state and action, i.e.,  $s' \sim G(s, a)$  for the MDP case [6, Ch. 4] and  $(s', o) \sim G(s, a)$  for the POMDP case [6, Ch. 6].

### 3.3 Reinforcement learning

If all the elements  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$  of an MDP are known, an agent can use this model to directly compute an optimal policy. Such a problem is often considered a planning problem. For small MDPs, dynamic programming<sup>1</sup> techniques can provide an exact solution, which is calculated offline, i.e., before the agent is deployed in the environment. For example, in value iteration [6, Ch. 4], the Bellman operator is iteratively applied to the value function for all states,

$$V_{n+1}(s) = \max_a \left[ R(s, a) + \gamma \sum_{s'} T(s'|, a, s) V_n(s') \right]. \quad (3.2)$$

As  $n$  goes to infinity,  $V_n$  converges to the unique optimal value function  $V^*$ , and an optimal policy (not necessarily unique) is extracted by

$$\pi(s) = \operatorname{argmax}_a \left[ R(s, a) + \gamma \sum_{s'} T(s'|, a, s) V^*(s') \right]. \quad (3.3)$$

However, for many real-world problems with high dimensional state spaces, it is intractable to compute and store a policy offline. Contrarily to offline methods, online search methods perform planning from the current state up to some horizon, when the agent has been deployed. Thereby, the agent can limit the computation to states that are reachable from the current state, which is often significantly smaller than the full state space [6, Ch. 4].

In many problems, the state transition probabilities or the reward function are not known. These problems can be solved by reinforcement learning techniques, in which the agent learns how to behave from interacting with the environment [6, Ch. 5], see Figure ???. Compared to supervised learning, reinforcement learning presents some additional challenges. Since the data that are available to an RL-agent depends on its current policy, the agent must balance exploring the environment and exploiting the knowledge it has

---

<sup>1</sup>Dynamic programming refers to simplifying a complex problem by breaking it down into smaller sub-problems, often in a recursive manner.

already gained. Furthermore, a reward that the agent receives may depend on a crucial decision that was taken earlier in time, which makes it important to assign rewards to the correct decisions.

RL algorithms can be divided into model-based and model-free approaches [6, Ch. 5]. In the model-based versions, the agent first tries to estimate a representation of the state transition function  $T$  and then use a planning algorithm to find a policy. On the contrary, as the name suggests, model-free RL algorithms do not explicitly construct a model of the environment to decide which actions to take. The model-free approaches can be further divided into value-based and policy-based techniques. Value-based algorithms, such as  $Q$ -learning, aim to learn the value of each state and thereby implicitly define a policy. Policy-based techniques instead search for the optimal policy directly in the policy space, either by policy gradient methods or gradient-free methods, such as evolutionary optimization. There are also hybrid techniques that are both policy and value-based, such as actor critic methods.

RL algorithms generally assume that the environment is modeled as an MDP, i.e., that the state of the environment is known by the agent. However, in many cases of interest, only partial information about the state of the environment is available, which is modeled in the POMDP framework. For such cases, it is common to approximate the state by either the observation or a finite history of observations [7, Ch. 17]. The latter is referred to as a  $k$ -Markov approximation, where  $k$  defines the length of the included history. For a sufficiently long history, the Markov property is assumed to approximately hold, even though the environment is partially observable.



# CHAPTER 4

---

## Model-free RL approaches



### 4.1 State representation

This section describe the general state representation used in this research that enables these methods to be generalizable for different type of intersection and crossings. By describing the state space as a set of distances to intersection points, we can abstract the map layout of different intersections and the same algorithms would work for intersections variations that we havent specifically trained on.

**ToDo:** add image of intersection scenario with 90 degree entry point and 45 degree entry point.

#### 4.1.1 Observable states

position, velocity and acceleration.

#### 4.1.2 Unobservable states

intentions. traffic lights and traffic signs.

## **4.2 Approach**

## **4.3 Simulated experiments**

## **4.4 Results and discussion**

# CHAPTER 5

---

## Combining MCP and RL

---

### 5.1 Action space, options

### 5.2 MPC

### 5.3 Approach

MPC has a mixed integer problem, calculating the optimal path for all possible action is very computationally heavy. RL DQN. Only has descret actions. can not garantee safety. but is good at choosing actions with the best utility (value). The reward function takes in the predicted outcome of the model in the MPC and can penalize the choice of action. but if experience show that the outcome is better than the model, it can choose to take a bad action that would lead to a better total reward compared to only following a conservative model.

## **5.4 Simulated experiments**

We show the difference in a multi crossing scenario where the MPC can plan a path for both intersections while our previous DQN only handles one at a time.

## **5.5 Results and discussion**

# CHAPTER 6

---

## Estimating the uncertainty

---

**ToDo:** rewrite chapter name

### **6.1 Uncertainty of the decision**

**6.1.1 Approach**

**6.1.2 Simulated experiments**

**6.1.3 Results and discussion**

### **6.2 Uncertainty of the intention**

**6.2.1 Approach**

**6.2.2 Simulated experiments**

**6.2.3 Results and discussion**



# CHAPTER 7

---

Generalize over different scenarios

---

## **7.1 Approach**

## **7.2 Simulated experiments**

## **7.3 Results and discussion**

FILL



# CHAPTER 8

---

## Discussion

---

FILL



# CHAPTER 9

---

Concluding remarks and future work

---

FILL

## 9.1 Future work

FILL



# CHAPTER 10

---

## Summary of included papers

---

This chapter provides a summary of the included papers.

### 10.1 Paper A

Tommy Tram, Anton Jansson, Robin Grönberg, Mohammad Ali, and Jonas Sjöberg

Learning Negotiating Behavior Between Cars in Intersections using Deep Q-Learning

*Published in 2018 21st International Conference on Intelligent Transportation Systems (ITSC),*

pp. 3169–3174, Nov. 2018.

©2018 IEEE DOI: 10.1109/ITSC.2018.8569316.

This paper concerns automated vehicles negotiating with other vehicles, typically human driven, in crossings with the goal to find a decision algorithm by learning typical behaviors of other vehicles. The vehicle observes distance and speed of vehicles on the intersecting road and use a policy that adapts its speed along its pre-defined trajectory to pass the crossing efficiently. Deep Q-

learning is used on simulated traffic with different predefined driver behaviors and intentions. The results show a policy that is able to cross the intersection avoiding collision with other vehicles 98% of the time, while at the same time not being too passive. Moreover, inferring information over time is important to distinguish between different intentions and is shown by comparing the collision rate between a Deep Recurrent Q-Network at 0.85% and a Deep Q-learning at 1.75%.

## **10.2 Paper B**

Tommy Tram, Ivo Batković, Mohammad Ali, and Jonas Sjöberg  
Learning When to Drive in Intersections by Combining Reinforcement Learning and Model Predictive Control  
*Published in 2019 IEEE Intelligent Transportation Systems Conference (ITSC),*  
pp. 3263–3268, Oct. 2019.  
©2019 IEEE DOI: 10.1109/ITSC.2019.8916922.

In this paper, we propose a decision making algorithm intended for automated vehicles that negotiate with other possibly non-automated vehicles in intersections. The decision algorithm is separated into two parts: a high-level decision module based on reinforcement learning, and a low-level planning module based on model predictive control. Traffic is simulated with numerous predefined driver behaviors and intentions, and the performance of the proposed decision algorithm was evaluated against another controller. The results show that the proposed decision algorithm yields shorter training episodes and an increased performance in success rate compared to the other controller.

## **10.3 Paper C**

Carl-Johan Hoel, Tommy Tram, and Jonas Sjöberg  
Reinforcement Learning with Uncertainty Estimation for Tactical Decision-Making in Intersections  
*Published in 2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC),*  
pp. 1-7, Sep. 2020.

©2020 IEEE DOI: 10.1109/ITSC45102.2020.9294407.

This paper investigates how a Bayesian reinforcement learning method can be used to create a tactical decision-making agent for autonomous driving in an intersection scenario, where the agent can estimate the confidence of its decisions. An ensemble of neural networks, with additional randomized prior functions (RPF), are trained by using a bootstrapped experience replay memory. The coefficient of variation in the estimated Q-values of the ensemble members is used to approximate the uncertainty, and a criterion that determines if the agent is sufficiently confident to make a particular decision is introduced. The performance of the ensemble RPF method is evaluated in an intersection scenario and compared to a standard Deep Q-Network method, which does not estimate the uncertainty. It is shown that the trained ensemble RPF agent can detect cases with high uncertainty, both in situations that are far from the training distribution, and in situations that seldom occur within the training distribution. This work demonstrates one possible application of such a confidence estimate, by using this information to choose safe actions in unknown situations, which removes all collisions from within the training distribution, and most collisions outside of the distribution.

## 10.4 Paper D

Tommy Tram, Maxime Bouton, Jonas Sjöberg, and Mykel Kochenderfer  
Belief State Reinforcement Learning for Autonomous Vehicles in Intersections

*Published in TBD,*  
pp. 3169–3174, Nov. 2018.  
©2018 IEEE DOI: TBD.

This paper concerns automated vehicles negotiating with other vehicles, typically human driven, in crossings with the goal to find a decision algorithm by learning typical behaviors of other vehicles. The vehicle observes distance and speed of vehicles on the intersecting road and use a policy that adapts its speed along its pre-defined trajectory to pass the crossing efficiently. Deep Q-learning is used on simulated traffic with different predefined driver behaviors and intentions. The results show a policy that is able to cross the intersection avoiding collision with other vehicles 98% of the time, while at the same time

not being too passive. Moreover, inferring information over time is important to distinguish between different intentions and is shown by comparing the collision rate between a Deep Recurrent Q-Network at 0.85% and a Deep Q-learning at 1.75%.

## 10.5 Paper E

Hannes Eriksson, Tommy Tram, Debabrota Basu, Jonas Sjöberg, and Christos Dimitrakakis

Maximum Likelihood Estimation for Transfer Reinforcement Learning Problems in Autonomous Driving

*Artificial Intelligence and Statistics 2023 (AISTATS)*,  
pp. 3169–3174, Nov. 2023.

©2018 IEEE DOI: 10.1109/ITSC.2018.8569316.

For decision-problems with insufficient data, it is imperative to take into account not only what you know but also what you do not know. In this work, ways of transferring knowledge from known, existing tasks to a new setting is studied. In particular, for tasks such as autonomous driving, the optimal controller is conditional on things such as, the physical properties of the vehicle, the local and regional traffic rules and regulations and also on the specific scenario trying to be solved. Having separate controllers for every combination of these conditions is intractable. By assuming problems with similar structure, we are able to leverage knowledge attained from similar tasks to guide learning for new tasks. We introduce a maximum likelihood estimation procedure for solving Transfer Reinforcement Learning (TRL) of different types. This procedure is then evaluated over a set of autonomous driving settings, each of which constitutes an interesting scenario for autonomous driving agents to make use of external information. We prove asymptotic regret bounds for proposed method for general structured probability matrices in a specific setting of interest.

---

## References

---

- [1] N. Tran *et al.*, “Global status report on road safety 2018,” World Health Organization, Tech. Rep., 2018.
- [2] S. Singh, “Critical reasons for crashes investigated in the national motor vehicle crash causation survey,” National Highway Traffic Safety Administration, Tech. Rep. DOT HS 812 506, 2018.
- [3] D. J. Fagnant and K. Kockelman, “Preparing a nation for autonomous vehicles: Opportunities, barriers and policy recommendations,” *Transportation Research Part A: Policy and Practice*, vol. 77, pp. 167–181, 2015, ISSN: 0965-8564.
- [4] J. Janai, F. Güney, A. Behl, and A. Geiger, *Computer Vision for Autonomous Vehicles: Problems, Datasets and State of the Art*. Now Publishers Inc., 2020.
- [5] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli, “A survey of motion planning and control techniques for self-driving urban vehicles,” *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 1, pp. 33–55, 2016.
- [6] M. J. Kochenderfer, *Decision Making Under Uncertainty: Theory and Application*. MIT Press, 2015, ISBN: 0262029251, 9780262029254.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press, 2018.



# **Part II**

# **Papers**



# PAPER A

## Learning Negotiating Behavior Between Cars in Intersections using Deep Q-Learning

Tommy Tram, Anton Jansson, Robin Grönberg, Mohammad Ali, and Jonas Sjöberg

*Published in 2018 21st International Conference on Intelligent Transportation Systems (ITSC),  
pp. 3169–3174, Nov. 2018.  
©2018 IEEE DOI: 10.1109/ITSC.2018.8569316.*

*The layout has been revised.*

## Abstract

This paper concerns automated vehicles negotiating with other vehicles, typically human driven, in crossings with the goal to find a decision algorithm based on learning typical behavior of other vehicles. The vehicle observes distance and speed of vehicles on the intersecting road and use a policy that adapts its speed along its pre-defined trajectory to pass the crossing efficiently. Deep Q-learning is used on simulated traffic and the results show that policies can be trained to successfully drive comfortably through an intersection, avoiding collision with other cars and not being too passive. The policies generalize over different types driver behaviors and intentions. Moreover, to enable inferring information over time, a Deep Recurrent Q-Network is tested and compared to the Deep Q-learning. The results show that a Deep Recurrent Q-Network succeeds in three out of four attempts where a Deep Q-Network fails.

## 1 Introduction

The development of autonomous driving vehicles is fast and there are regularly news and demonstrations of impressive technological progress, see eg **Bojarski2016EndCars**. However, one of the largest challenges does not have to do with the autonomous vehicle itself but with the human driven vehicles in mixed traffic situations. Human drivers are expected to follow traffic rules strictly, but in addition they also interact with each other in a way which is not captured by the traffic rules, **Liebner2012DriverModel**, **Lefevre2012EvaluatingIntentions**. This *informal* traffic behavior is important, since the traffic rules alone may not always be enough to give the safest behavior. This motivates the development of control algorithms for autonomous vehicles which behave in a "human-like" way, and in this paper we investigate the possibilities to develop such behavior by training on simulated vehicles.

In **Shalev-ShwartzSafeDriving** they raises two concerns when using Machine learning, specially Reinforcement learning, for autonomous driving applications: ensuring functional safety of the Driving Policy and that the Markov

Decision Process model is problematic, because of unpredictable behavior of other drivers. In the real world, intentions of other drivers are not always deterministic or predefined. Depending on their intention, different actions can be chosen to give the most comfortable and safe passage through an intersection. They also noted that in the context of autonomous driving, the dynamics of vehicles is Markovian but the behavior of other road users may not necessarily be Markovian. In this paper we solve these two concerns using a Partially Observable Markov Decision Process (POMDP) as a model and Short Term Goals (STG) as actions. With a POMDP the unknown intentions can be estimated using observations and that has shown promising results for other driving scenarios **BrechtelProbabilisticPOMDPs**. The POMDP is solved using a model-free approach called Deep (Recurrent) Q-Learning. With this approach a driving policy can be found using only observations without defining the states. Since we do not train on human driven vehicles, the results presented here cannot be considered human-like, but the general approach, to train the algorithms using traffic data, is shown working, and a possible next step could be to start with the pre-tuned policies from this work, and to continue the training in real traffic crossings.

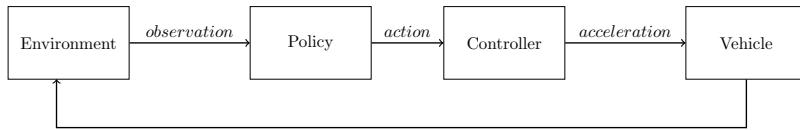
## 2 Overview

This paper starts by introducing the system architecture and defining the actions, observations and POMDP in Section 3. The final strategy of what action to take at a given situation is called a policy and is described in Section 4. The method used for finding this policy is called Q-learning, which uses a neural network to approximate a Q-value and is described in Section 5 together with techniques used to improve the learning, such as Experience replay, Dropout and a recurrent layer called Long short term memory (LSTM). We then present the simulation, reward function and neural network configurations in Section 6. The results are then presented in Section 7 comparing the effect of the methods mentioned in Section 5. Finally the conclusion and brief discussion is presented in Section 8.

### 3 Problem formulation

The objective is to drive along a main road that has one or two intersections with crossing traffic and control the acceleration in a way that avoids collisions in a comfortable way. All vehicles are assumed to drive along predefined paths on the road where they can either speed up or slow down to avoid collisions in the crossings. In this section the system architecture is defined along with the environment, observation and actions.

#### 3.1 System architecture



**Figure 1:** Representation of the architecture

Environment is defined as the world around the ego vehicle, including all vehicles of interest and the shape/type of the intersection. The environment can vary in different ways, e.g. number of vehicles and intersection or the distance to intersections. The environment is defined by the simulation explained in section 6.1. We assume that the ego vehicle receives observations from this environment at each sampling instant, as shown in Fig. 1. A policy then takes these observations and chooses a high level action that is defined in more detail in section 3.3. These actions are sent to a controller that calculates the appropriate acceleration request given to the ego vehicle, which will influence the environment and impact how other cars behave.

#### 3.2 Actions as Short Term Goals

Motivated by the insight that the ego vehicle has to drive before or after other vehicles when passing the intersection, decisions on the velocity profile is modeled by simply keeping a distance to other vehicles until they pass. This is done by defining the actions as Short Term Goals (STG), eg. keep set speed or yield for crossing car. This allows the properties of comfort on actuation

and safety to be tuned separately, making the decision a classification problem. The actions are then as followed:

- *Keep set speed:* Aims to keep a specified maximum speed  $v_{\max}$ , using a simple P-controller

$$a_p^e = K(v_{\max} - v^e) \quad (\text{A.1})$$

where  $a_p^e$  is the acceleration request and  $v^e$  is the velocity of ego vehicle towards the center of the intersection, while  $K$  is a proportional constant.

- *Keep distance to vehicle N:* Will control the acceleration in a way that keeps a minimum distance to a chosen vehicle  $N$ , a *Target Vehicle*, and can be done using a sliding mode controller, where the acceleration request is:

$$a_{sm}^e = \frac{1}{c_2}(-c_1 x_2 + \mu sign(\sigma(x_1, x_2))) \quad (\text{A.2})$$

$$\text{where } \begin{cases} x_1 = p^t - p^e \\ x_2 = v^t - v^e \end{cases}$$

where  $p^e$  and  $p^t$  is the position of ego and target vehicle respectively, shown in Fig. 2, and  $v^t$  is the velocity of target vehicle.  $c_1$  together with  $c_2$  are calibration parameters that can be set to achieve wanted performance with a surface plane  $\sigma$

$$\sigma = c_1 x_1 + c_2 x_2 \quad (\text{A.3})$$

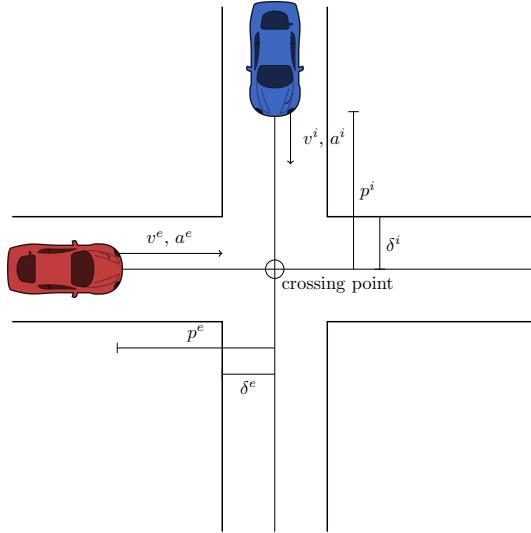
The final acceleration request  $a^e$  is achieved by a min arbitration between eq. A.1 and A.2

$$a^e = \min(a_{sm}^e, a_p^e) \quad (\text{A.4})$$

For more detailed information about sliding mode see **MemonAnalysisManoeuvres**. To distinguish between different cars to follow, each other vehicle will have its own action.

- *Stop in front of intersection:* Stops the car at the next intersection. Using the same controller as eq. A.4 while setting  $v^t = 0$  and  $p^t$  to start of intersection, the controller can bring ego vehicle to a comfortable stop before the intersection.

### 3.3 Observations that make up the state



**Figure 2:** Observations that makes the state

A human driver is, generally, good at assessing a scenario and it is hard to pin-point what information is used in their assessment. Therefore some assumptions are made on which features that are interesting to observe. The observation  $o_t$  at time  $t$  is defined as:

$$o_t = [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^i \quad v_t^i \quad a_t^i \quad \delta^i \quad a_{t+1}^{e,A} ]^T \quad (\text{A.5})$$

With notations as follows: consider Fig. 2, position of ego  $p_t^e$  and other vehicle  $p_t^i$  are defined as distance to common reference point, called *crossing point*, where  $i$  is an index of the other vehicle. The start of intersection for ego  $\delta^e$  and other vehicle  $\delta^i$  also uses the crossing point as reference. These are relevant in case a driver would choose to yield for other vehicles, then

they would most likely stop before the start of intersection. The velocity  $v^e$  and acceleration  $a^e$  of ego vehicle and velocity  $v^i$  and acceleration  $a^i$  of the other vehicles are observed to include the dynamics of different actors. The last feature in the observation,  $a_{t+1}^{e,A}$ , is the ego vehicle's predicted acceleration for each possible action  $A$ , which can be used to account for comfort in the decision.

### 3.4 Partially Observable Markov Decision Processes

The decision making process in the intersection is modeled as a POMDP. A POMDP works like a Markov Decision Process (MDP) **BellmanMDP** in most aspects, but the full state is not observable.

At each time instant, an action,  $a_t \in \mathcal{A}$ , is taken, which will influence to which new state vector,  $s_{t+1}$ , the system evolves and changes the environment. Each action  $a_t$  from a state  $s_t$  has a value called the reward  $r_t$ , which is given by a reward function  $\mathcal{R}_t$ .

One of the unobservable states could be the intentions of other drivers approaching the intersection. The state can only be perceived partially through observations  $o_t \in \Omega$  with the probability distribution of receiving observation  $o_t$  given an underlying hidden state  $s_t : o_t \leftarrow \mathcal{O}(s_t)$ , where  $\mathcal{O}(s_t)$  is the probability distribution.

## 4 Finding the optimal policy

Assuming the states are not known, we want a model-free method of finding a policy, and for this we use reinforcement learning. The goal is to have an agent learn how to maximize the future reward by taking different actions in a simulated environment. Details on the simulation environment used is described in Section 6.1. The strategy of which action to take given a state is called a policy  $\pi$  and can be modeled in two ways:

- As stochastic policy  $\pi(a|s) = \mathcal{P}[\mathcal{A} = a | \mathcal{S} = s]$
- As deterministic policy  $a = \pi(s)$

The standard assumption is made that the future reward is discounted by a factor  $\gamma$  per time step, making the discounted future reward  $\mathcal{R}_t = \sum_t^\tau \gamma^{t-1} r_t$ ,

where  $\tau$  is the time step where the simulation ends, e.g. when the agent crosses an intersection safely.

Similar to **MnihPlayingLearning**, the optimal action-value function  $Q^*(s_t, a_t)$  is defined as the maximum expected reward achievable by following a policy  $\pi$  given the state  $s_t$  and taking an action  $a_t$ :

$$Q^*(s_t, a_t) = \max_{\pi} \mathbb{E}[R_t | s_t, a_t, \pi] \quad (\text{A.6})$$

Using the Bellman equation,  $Q^*(s_t, a_t)$  can be defined recursively. If we know  $Q^*(s_t, a)$  for all actions  $a$  that can be taken in state  $s_t$ , then the optimal policy will be one that takes the action  $a_t$  that gives the highest immediate and discounted expected future reward  $r_t + \gamma Q^*(s_{t+1}, a_{t+1})$ . This gives us:

$$Q^*(s_t, a_t) = \mathbb{E}[r_t + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) | s_t, a_t] \quad (\text{A.7})$$

The optimal policy  $\pi^*$  is then given by taking actions according to an optimal  $Q^*(s_t, a_t)$  function:

$$\pi^*(s_t) = \arg \max_{a_t} Q^*(s_t, a_t) \quad (\text{A.8})$$

## 5 Method

From eq. A.8, the optimal policy is defined by taking an action that has the highest expected Q-value. Because the Q-value is not known, a non linear function approximation, such as a neural network, is used to estimate the Q-function. The method is known as Deep Q-Learning **MnihPlayingLearning** and in this section we will briefly describe Q-learning and methods used to improve the learning such as, Experience replay, dropout and Long-, Short-Term Memory (LSTM).

### 5.1 Deep Q-learning

Deep Q-Learning uses a neural network to approximate the  $Q$ -function. This neural network is called Deep Q-Network (DQN). The  $Q$ -function approximated by the DQN is denoted as  $Q(s_t, a_t | \theta^\pi)$ , where  $\theta^\pi$  are the neural network parameters for policy  $\pi$ . The state  $s_t$  is the input to the DQN and the output is the  $Q$ -value for each action  $\mathcal{A}$ .

## 5.2 Experience Replay

A problem with Deep Q-Learning, looking at eq. A.8, is that with a small change in the  $Q$ -function, can effect the policy drastically. The distribution of future training samples will be greatly influenced by the updated policy. If the network only trains on recent experiences, a biased distribution of samples is used. Such behavior can cause unwanted feedback loops which can lead to divergence of the trained network **Tsitsiklis1997AnApproximation**.

Proposed by **MnihPlayingLearning** in order to make DQN more robust, all observations  $o$ , together with taken actions  $a$  and their rewards  $r$ , are stored in an experience memory  $E$  and the agent can sample experiences  $E'$  from the experience memory. The sampled experiences are fed to the gradient descent algorithm as a mini-batch. Thus, the agent learns on an average of the experiences in  $E$  and is likely to train on the same experiences multiple times, which can speed up the network's convergence and reduce oscillations **Lin1992Self-ImprovingTeaching**.

## 5.3 Dropout

Overfitted neural networks have bad generalization performance **Hinton2012ImprovingDe** and to help reduce overfitting a technique called dropout was used. The idea with dropout is to temporarily remove random hidden neurons with their connections from the network, before each training iteration. This is done by, independently, for each neuron, setting its value to 0 with a probability  $p$ . For more details, see **Srivastava2014Dropout:Overfitting**.

## 5.4 Long short term memory

The effect of changes over time is explored in this paper and is done by adding a recurrent layer to the DQN making it a Deep Recurrent Q-Network (DRQN). A regular Recurrent Neural Network struggle to remember longer sequences due to vanishing gradients, and in **Hochreiter1997LONGMEMORY** LSTM is used to solve this problem. An LSTM is a recurrent network constructed for tasks with long-term dependencies. Instead of storing all information from previous time sample, LSTM stores information in a memory cell and modify this memory by using insert and forget gates. These gates decides if a memory cell should be kept or cleared, and during training, the network learns how to control these gates. As a result of this, both newly seen observations and

observations seen a long time ago can be stored and used by the network. A sequence length of 4 is used when training the LSTM, where the first 3 observations are only used to build the internal memory state of the LSTM cells, as described in **LamplePlayingLearning**.

## 6 Implementation

In this section we go through the experiment implementation. A simulation environment was set up to model the interactions. From section 3, both the number of observations and actions are dependent on the maximum number of cars. In this paper we consider up to 4 cars. The Deep Q Network can then also be fully defined with the help of observations from section 3 and finally we go through the reward function that defines our behavior.

### 6.1 Simulation environment

The simulation environment is set up as an intersection described in section 3.3. The number of other cars that are observable at the same time can vary from 1-4, while their intentions can vary between an aggressive *take way*, passive *give way* or a cautious driver. The take way driver does not slow down or yield for crossing traffic in an intersection, while the give way driver will always yield for other vehicles before continuing through the intersection. The cautious driver on the other hand, will slow down for crossing traffic but not come down to a full stop. With a maximum number of other cars set to 4 all possible actions the ego vehicle can take are:

- $\alpha_1$ : Keep set speed.
- $\alpha_2$ : Stop in front of intersection.
- $\alpha_3$ : Keep distance to vehicle 1.
- $\alpha_4$ : Keep distance to vehicle 2.
- $\alpha_5$ : Keep distance to vehicle 3.
- $\alpha_6$ : Keep distance to vehicle 4.

At the start of an episode, the ego vehicle's position and velocity, the number of other vehicles and their intentions are randomly generated. The episode only end when the ego vehicle fulfills one out of three conditions: 1. Crossing the intersection and reaching the other side, 2. Colliding with another vehicle. or 3. Running out of time  $\tau_m$ . Each car follows the control law from eq. A.4, trying to keep a set speed while keeping a set distance to the vehicle in front of its own lane.

All cars including the ego vehicle in these scenarios have a maximum acceleration set to  $5m/s^2$  and was set after comfort and normal driving conditions.

## 6.2 Reward function tuning

When using a DQN, the reward function is optimally distributed around  $[-1, 1]$ . If the defined reward values are too large, the  $Q_\pi$ -values can become large and cause the gradients to grow **VanHasseltLearningMagnitude**. The reward function is defined as follows:

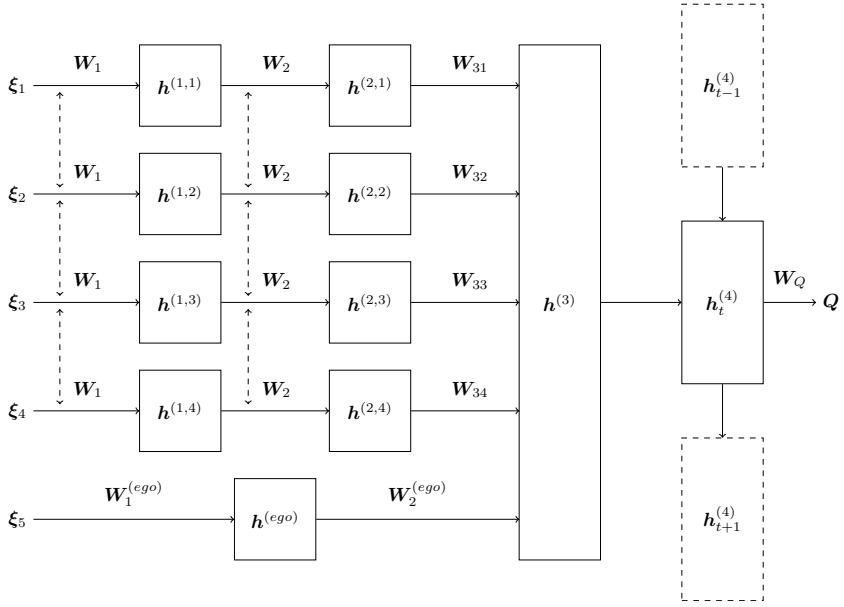
$$r_t = \hat{r}_t + \begin{cases} 1 - \frac{\tau}{\tau_m} & \text{on success,} \\ -2 & \text{on collision} \\ -0.1 & \text{on timeout, i.e. } \tau \geq \tau_m \\ -\left(\frac{j_t^e}{j_{\max}}\right)^2 \frac{\Delta\tau}{\tau_m} & \text{on non-terminating updates} \end{cases}$$

where  $\hat{r}_t = \begin{cases} -1 & \text{if chosen } a_t \text{ is not valid} \\ 0 & \text{otherwise} \end{cases}$

The actions  $\alpha_3, \dots, \alpha_6$  described should only be selected when a vehicle is visible and a valid target to follow. To learn when these action are valid, the agent gets punished on invalid choices using  $\hat{r}_t$ . Accelerations returned by the controller for different STG can vary, which increases jerk and can make the ride uncomfortable. Therefore the reward function also punishes the agent when acceleration jerk  $j_t^e$  is large, where  $\tau$  is the elapsed time sense the episode started and  $\tau_m$  is the maximum time has before a timeout.

## 6.3 Neural Network Setup

The DRQN structure is defined in Fig. 3. Where  $\mathbf{h}$  are the hidden layers of the network with weights  $\mathbf{W}$ . Because the observations  $o_t$  from section 2, are



**Figure 3:** Deep Recurrent Q Network layout with shared weights and a LSTM

used as input to the DRQN, the number of features must be fixed. With up to four other cars the input vectors  $\xi$  are as follows:

- $\xi_1 = [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^1 \quad v_t^1 \quad a_t^1 \quad \delta^1 ]^T$
- $\xi_2 = [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^2 \quad v_t^2 \quad a_t^2 \quad \delta^2 ]^T$
- $\xi_3 = [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^3 \quad v_t^3 \quad a_t^3 \quad \delta^3 ]^T$
- $\xi_4 = [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^4 \quad v_t^4 \quad a_t^4 \quad \delta^4 ]^T$
- $\xi_5 = [ a_{t+1}^{e,1} \quad a_{t+1}^{e,2} \quad a_{t+1}^{e,3} \quad a_{t+1}^{e,4} \quad a_{t+1}^{e,5} \quad a_{t+1}^{e,6} ]^T$

In case a vehicle is not visible, the input vector  $\xi$  are set to -1. All vehicle features are scaled to be a value between or close to  $[-1, 1]$  by using a car's sight range  $p_{\max}$ , maximum speed  $v_{\max}$  and maximum acceleration  $a_{\max}$ .

The output  $Q$  should be independent of which order other vehicles was observed in the input  $\xi_i$ . In other words, whether a vehicle was fed into  $\xi_1$  or

into  $\xi_4$ , the network should optimally result in the same decision, only based on the features' values. The network is therefore structured such that input features of one car, for instance  $\xi_1$ , are used as input to a sub-network with two layers  $\mathbf{h}^{(1,i)}$  and  $\mathbf{h}^{(2,i)}$ . Each other vehicle has a copy of this sub-network, resulting in them sharing weights ( $\mathbf{W}_1$  and  $\mathbf{W}_2$ ), as shown in Fig. 3. The first hidden layers are then given by:

$$\mathbf{h}^{(1,i)} = \tanh(\mathbf{W}_1 \boldsymbol{\xi}_i + \mathbf{b}_1) \quad (\text{A.9})$$

$$\mathbf{h}^{(2,i)} = \tanh(\mathbf{W}_2 \mathbf{h}^{(1,i)} + \mathbf{b}_2) \quad (\text{A.10})$$

$$\mathbf{h}^{(ego)} = \tanh(\mathbf{W}_1^{(ego)} \boldsymbol{\xi}_5 + \mathbf{b}^{(ego)}) \quad (\text{A.11})$$

The output of each sub-network,  $\mathbf{h}^{(2,i)}$  and  $\mathbf{h}^{(ego)}$ , is fed as input into a third hidden layer  $\mathbf{h}^{(3)}$ . The different sub-networks'  $\mathbf{h}^{(2,i)}$  outputs are multiplied with different weights  $\mathbf{W}_{31}, \dots, \mathbf{W}_{34}$  in order to distinguish different cars for different follow car actions. The ego features are also fed into layer 3 with its own weights  $\mathbf{W}_2^{(ego)}$ . The neurons in layer  $\mathbf{h}^{(3)}$  combine the inputs by adding them together:

$$\mathbf{h}^{(3)} = \tanh\left(\mathbf{W}_2^{(ego)} \mathbf{h}^{(ego)} + \sum_{i=1}^4 \mathbf{W}_{3i} \mathbf{h}^{(2,i)} + \mathbf{b}_3\right) \quad (\text{A.12})$$

The final layer  $\mathbf{h}^{(4)}$  uses the LSTM, described in section 5. This layer handles the storage and usage of previous observations, making it the recurrent layer of the network.

$$\mathbf{h}_t^{(4)} = \text{LSTM}\left(\mathbf{h}^{(3)} | \mathbf{h}_{t-1}^{(4)}\right) \quad (\text{A.13})$$

The output of the neural network is then the approximated  $\mathbf{Q}$ -value:

$$\mathbf{Q} = \mathbf{W}_Q \mathbf{h}^{(4)} + \mathbf{b}_4 \quad (\text{A.14})$$

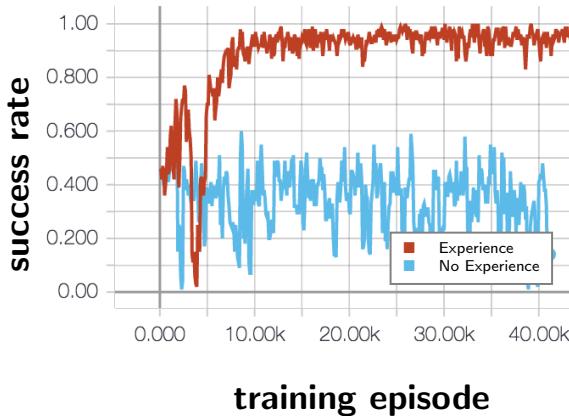
## 7 Results

Three metrics were selected for measurement of a training session: success rate, average episodic reward and collision to timeout ratio. With these metrics, the distribution between the three terminating states can be analyzed. Success

rate is defined as the success to fail ratio averaged over the last 100 episodes, where both collisions and timeouts are considered as failures. Average episodic reward is the summed reward over a whole episode, then averaged over 100 episodes. The collision to fail ratio displays the ratio between the number of collisions and unsuccessful episodes for the last 100 episodes. From the success rate and collision to fail ratio, a final collision rate is computed, which is the amount of episodes resulting in a collision, averaged over 100 episodes. The graphs presented are only using evaluation episodes, with a deterministic policy. Every 300 episode, the trained network is evaluated over 300 evaluation episodes.

The improvement of using Dropout and Experience replay, from Section 5, are clearly shown in Fig. 4 and 5. Studying the red curve in Fig. 4, with all methods included, the best policy had a success rate of 99%, average episodic reward 0.8 and collision to timeout ratio at 40%.

## 7.1 Effect of using Experience replay

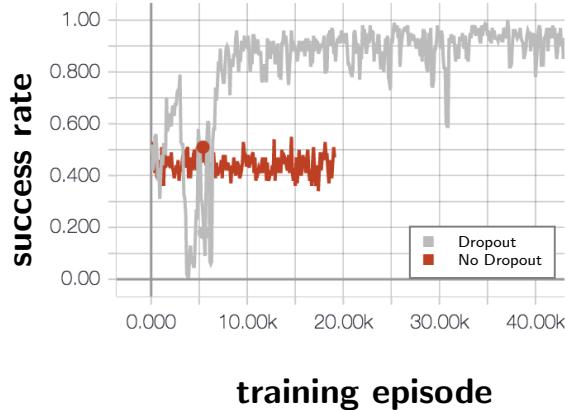


**Figure 4:** Success rate trend comparing using experience replay (red) and not using experience replay (blue)

Without either method the success rate does not converge to a value higher than 60%. When experience replay was not used, the highest success rate was 53%, average episodic reward  $-0.1$  and collision to timeout ratio at 90%.

Compared to not using dropout, not using experience replay has a higher lower variation on the success rate.

## 7.2 Effect of using Dropout



**Figure 5:** Success rate trend comparing using dropout (grey) and not using dropout (red)

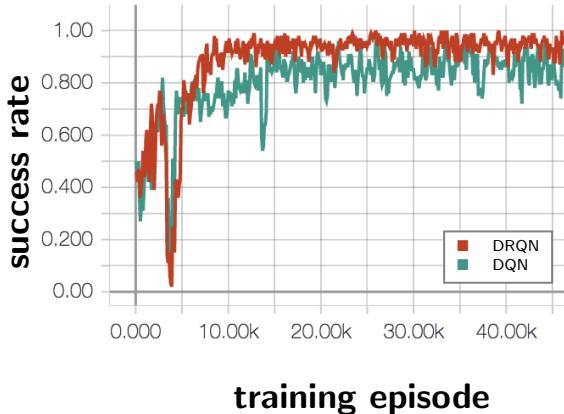
In the case of not using Dropout, the best policy had a success rate of 58%, average episodic reward  $-0.7$  and a collision to timeout ratio at 90%.

## 7.3 Comparing DQN and DRQN

In Fig. 6, we can see the effect off having a recurrent layer by comparing a DQN, without a LSTM layer, and with a DRQN, with LSTM. The faded colored line show actual sampled values and the thick line acts as a trend line which for DRQN converges towards a success rate of around 97.2% and a 0.85% collision rate, compared to a success rate of 87.5% and a collision rate of 1.75% for DQN.

## 7.4 Effect of sharing weights in the network

When introducing multiple cars in the scenario, the success rate converged considerably slower, as shown in Fig. 7. Using the network structure that share



**Figure 6:** Graphs presenting the performance of a DRQN (red) compared to a DQN with a single observation (green), running on scenarios with cars that have different behaviors. When compared a DRQN succeeds in 3 out of 4 attempts, where a DQN fails.

weights between cars, significantly improved how fast the network converged compared to a fully connected DRQN.

## 8 Conclusion

In this paper, Deep Q-Learning was presented in the domain of autonomous vehicle control. The goal of the ego agent is to drive through an intersection, by adjusting longitudinal acceleration using short-term goals. Short-term goals also allowed a smoother and more human-like behavior by controlling the acceleration and comfort with a separate controller. Instead of finding a policy with a continuous control output the problem became a classification problem.

The results also show that trained policies can generalize over different types of driver behaviors. The same policy is able to respond to other vehicles' actions and handle traffic scenarios with a varied number of cars, without knowing traffic rules or the type of intersection it drives in. Multiple observations are needed in order to recognize cars' behaviors, and can be utilized by for instance a DRQN.



**Figure 7:** The figure shows that the success rate for a network with shared weights (brown line) converge faster than the fully connected network structure which do not share weights (turquoise line).

There was a significant performance improvement in using a DRQN instead of a DQN with a single observation. In other words, the environment for these scenarios is better modeled as a POMDP instead of an MDP and the agent needs multiple observations in order to draw enough conclusions about other cars' behaviors. Convergence of the neural network was shown to be improved by sharing weights between the first layers to which the target car features are fed, compared to a fully connected neural network structure. The results are still limited to the tested traffic scenarios and driver behaviors, and expanding the domain beyond a simulator is a natural next step. The selected features are also limited to intersections.

The results show a success rate of around 99% for recognizing behaviors. However, the ego vehicle still collides. The ego vehicle is limited to drive comfortably, meaning that in some cases, the ego agent is not allowed to break hard enough. In a complete system, a collision avoidance procedure, which does not have comfort constraints, would need to take over the control to ensure a safe ride. A collision in this paper is defined by two areas overlapping,

and in a real world implementation this does not have to mean an actual collision. Instead this could be interpreted as an intervention from a more safety critical system. This way, in the low chances a good action could not be found, the safety of the vehicles can still be guaranteed.

In section 3.2, a sliding mode controller was chosen, but this can be replaced by any controller. One other option could be a Model Predictive Controller, where safer actuation can be achieved by using constraints. Also, the actions in this paper used the same controller tuning for all actions, this does not have to be the case. An action can have the same STG but only differ by the controller's tuning parameters. This way, the agent gains more flexibility while the comfort can remain intact, possibly increasing the success rate.



# PAPER B

**Learning When to Drive in Intersections by Combining Reinforcement Learning and Model Predictive Control**

Tommy Tram, Ivo Batković, Mohammad Ali, and Jonas Sjöberg

*Published in 2019 IEEE Intelligent Transportation Systems Conference (ITSC),  
pp. 3263–3268, Oct. 2019.  
©2019 IEEE DOI: 10.1109/ITSC.2019.8916922.*

*The layout has been revised.*

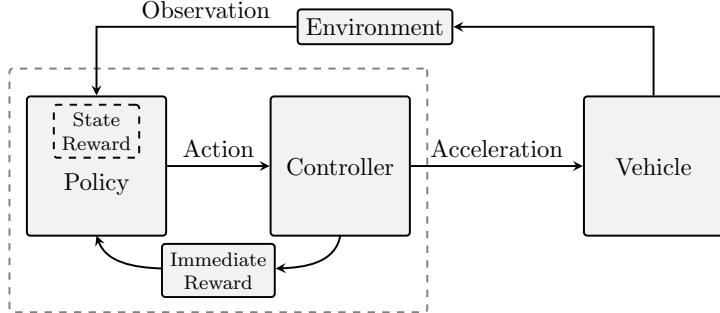
## Abstract

In this paper, we propose a decision making algorithm intended for automated vehicles that negotiate with other non-automated vehicles in crossings. The decision algorithm is separated into two parts: a high-level decision module based on reinforcement learning, and a low-level planning module based on model predictive control. Traffic is simulated with different predefined driver behaviors and intentions, and evaluate the performance of the proposed decision algorithm and benchmark it against using a sliding mode controller. The results show that the proposed decision algorithm yields faster training times and an increased performance compared to the sliding mode controller.

## 1 Introduction

Self driving cars is a fast advancing field with Advanced Driver-Assistance Systems becoming a requirement in modern day cars. Decision making for self driving cars can be difficult to solve with simple rule based system in complex scenarios like intersections, while human drivers have a good intuition about when to drive and how to drive comfortably. Sharing the road with other road users requires interaction, which can make rule based decision making complex **Liebner2012DriverModel**. Many advancements aim to bring self driving Level 4 to the market by trying to imitate human drivers **Bansal2018ChauffeurNet:Worst** or predicting what other drivers in traffic are planning to do **Zyner2017LongPrediction**.

Previous research **Tram2018LearningQ-Learning** showed that reinforcement learning (RL) can be used to learn a negotiation behaviour between cars without vehicle to vehicle communication when driving in an intersection. The method found a policy that learned to drive though an intersection, with crossing traffic, where other vehicles have different intentions and avoided collision. The previous method could use the same algorithm, but trained on different type of intersections and still find a general policy that would get to the other side of the intersection while avoiding collision. By modeling the decision process as a partially observable Markov decision process, uncertainty in the environment or sensing can be accounted for **BrechtelProbabilisticPOMDPs**



**Figure 1:** Representation of the decision making architecture

and still be safe **BoutonReinforcementDriving**.

Because the decision policy is separated from the control, high level decision making can focus on the task, when to drive, while the low level control that handles the comfort of passengers in the car by generating a smooth acceleration profile. Previous work showed how this worked for intersections with a single crossing point, where Short Term Goal (STG) actions could choose one car to follow. This gives the RL policy a flexibility to choose actions that can safely transverse through the intersection by switching between different STG.

A simple controller holds well when intersection crossing points are far away from each other, but when there are several crossing points in close succession, a simple controller would have a hard time handling it. In this paper, we instead propose using MPC that can consider multiple vehicles at the same time and generate an optimal trajectory. In contrast to **hult**, where they prove stability and recursive feasibility using an MPC approach, and assuming that agents can cooperate, we restrict ourselves to non-cooperative scenarios. The MPC is used to plan trajectories around other vehicles using available predictions from other vehicles, and in the worst case, use the prediction as early detection whenever a dangerous situation, e.g., a collision, may appear.

Applying MPC directly to the problem, could lead to a growing complexity with the number of vehicles in the intersection, e.g., which vehicle do we yield for and which vehicle do we drive in front. Therefore, we propose to separate the problem into two parts: the first being a high-level decision maker, which structures the problem, and the second being a low level planner, which

optimizes a trajectory given the traffic configuration.

For the high-level decision maker RL is used to generate decisions for how the vehicle should drive through the intersection, and MPC is used as a low-level planner to optimize a safe trajectory. Compared to **decentralizedMPC** where all vehicles are controlled using MPC to stay in a break-safe set based on a model of other vehicles future trajectory, this can be perceived as too conservative for a passenger. By combining RL and MPC, the decision policy will learn which action is optimal by using feedback from the MPC controller into the reward function. Since MPC uses predefined models, e.g. vehicle models and other obstacle prediction models, the performance relies on their accuracy and assumptions. To mitigate this, we use Q-learning, which is a model-free RL approach, to optimize the expected future reward based on its experience during an entire episode which is able to compensate to some extent for model errors, which is explained more in section 5.1.

This paper is structured as follows. Section 2 introduces the system architecture of our framework. The problem formulation, along with the two-layers of the decision algorithm is presented in Section 3. Section 4 present three different used agents for simulation and validation. Implementation details is presented in Section 5 and the results are shown in Section 6 followed by discussion in Section 7. Finally, conclusions and future research are presented in Section 8.

## 2 System

A full decision architecture system shown in Fig. 1, would include a precautionary safety layer that limits which acceleration values the system can actuate in order to stay safe. Followed by a decisions making system that makes a high level decision for when to drive in order to avoid collision and be comfortable. The policy maker later send that action to a controller that actuates the action turning it into a control signal, e.g., an acceleration request. The environment state, together with the new acceleration request, is sent though a collision avoidance system that checks if the current path has a collision risk, and mitigate if needed. With such a structure, the comfort that is experienced by the passenger is fully controlled by the low-level controller and partially affected by the decision. This paper focuses on the integration between policy

and actuation, by having an MPC controller directly giving feedback to the decision maker through immediate actions. This allows the policy to know how comfortably the controller can handle the action and give feedback sooner if the predicted outcome may be good or bad.

### 3 Problem formulation

The goal of the ego-vehicle is to drive along a predefined route that has one or two intersections with crossing traffic, where the intent of the other road users is unknown. Therefore, the ego-vehicle needs to assess the driving situation and drive comfortably, while avoiding collisions with any vehicle<sup>1</sup> that may cross. In this section, we define the underlying Partially Observable Markov Decision Process (POMDP) and present how the problem is decomposed using RL for decision making and MPC for planning and control.

#### 3.1 Partially Observable Markov Decision Process

The decision making process is modeled as a Partially Observable Markov Decision Process (POMDP). A POMDP is defined by the 7-tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  an action space that is defined in section 4.1,  $\mathcal{T}$  the transition function,  $\mathcal{R}$  the reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is defined in 5.4,  $\Omega$  an observation space,  $\mathcal{O}$  the probability of being in state  $s_t$  given the observation  $o_t$ , and  $\gamma$  the discount factor.

A POMDP is a generalization of the Markov Decision Process (MDP) **BellmanMDP** and therefore works in the same way in most aspects. At each time instant  $t$ , an action,  $a_t \in \mathcal{A}$ , is taken, which will change the environment state  $s_t$  to a new state  $s_{t+1}$ . Each transition to a state  $s_t$  with an action  $a_t$  has a reward  $r_t$  given by a reward function  $\mathcal{R}$ . The Key difference from a regular MDP is that the environment state  $s_t$  is not entirely observable because the intention of other vehicles are not known. In order to find the optimal solution for our problem, we need to know the future intention of other drivers. Instead we can only partially perceive the state though observations  $o_t \in \Omega$ .

---

<sup>1</sup>Although our approach can be extended to other road users, for convenience of exposition we'll refer to vehicles.

### 3.2 Deep Q-Learning

In the reinforcement learning problem, an agent observes the state  $s_t$  of the environment, takes an action  $a_t$ , and receives a reward  $r_t$  at every time step  $t$ . Through experience, the agent learns a policy  $\pi$  in a way that maximizes the accumulated reward  $\mathcal{R}$  in order to find the optimal policy  $\pi^*$ . In Q-learning, the policy is represented by a state action value function  $Q(s_t, a_t)$ . The optimal policy is given by the action that gives the highest Q-value.

$$\pi^*(s_t) = \arg \max_{a_t} Q^*(s_t, a_t) \quad (\text{B.1})$$

Following the Bellman equation the optimal Q-function  $Q^*(s_t, a_t)$  is given by:

$$Q^*(s_t, a_t) = \mathbb{E}[r_t + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) | s_t, a_t] \quad (\text{B.2})$$

Because Q-learning is a model-free algorithm and it does not make any assumptions on the environment, even though models are used to simulate the environment, this can be useful when the outcome does not match the prediction models.

## 4 Agents

This section explains the different agents types. MPC and sliding mode (SM) for the ego vehicle. Observed target vehicles has different intention agents based on SM agent.

### 4.1 MPC agent

We model the vehicle motion with states  $\mathbf{x} \in \mathbb{R}^3$  and control  $\mathbf{u} \in \mathbb{R}$ , defined as

$$\mathbf{x} := [p^e \quad v^e \quad a^e]^\top, \quad \mathbf{u} := j^e, \quad (\text{B.3})$$

where we denote the position along the driving path in an Frenet frame as  $p^e$ , the velocity as  $v^e$ , the acceleration as  $a^e$ , and the jerk as  $j^e$ , see Fig. 2. In addition, we assume that measurements of other vehicles are provided through an observation  $\mathbf{o}$ . We limit the scope of the problem to consider at most four vehicles, and define the observations as

$$\mathbf{o} := [p^1 \quad v^1 \quad p_{\text{ego}}^{\text{cross},1} \quad \dots \quad p^4 \quad v^4 \quad p_{\text{ego}}^{\text{cross},4}]^\top, \quad (\text{B.4})$$

where we denote the position along its path as  $p$ , the velocity as  $v$ , and  $p_{\text{ego}}^{\text{cross},j}$  for  $j \in [1, 4]$ , as the distance to the ego-vehicle from the intersection point, see Fig. 2.

In this paper, we assume that there exists a lateral controller that stabilizes the vehicle along the driving path. To that end, we only focus on the longitudinal control. Given the state representation, the dynamics of the vehicle is then modeled using a triple integrator with jerk as control input.

The objective of the agent is to safely track a reference, e.g. follow a path with a target speed, acceleration, and jerk profile, while driving comfortably and satisfying constraints that arise from physical limitations and other road users, e.g. not colliding in intersections with crossing vehicles. Hence, we formulate the problem as a finite horizon, constrained optimal control problem

$$J = \min_{\bar{\mathbf{x}}, \bar{\mathbf{u}}} \sum_{k=0}^{N-1} \left[ \begin{array}{c} \bar{\mathbf{x}}_k - \mathbf{r}_k^{\mathbf{x}} \\ \bar{\mathbf{u}}_k - \mathbf{r}_k^{\mathbf{u}} \end{array} \right]^T \left[ \begin{array}{cc} Q & S^T \\ S & R \end{array} \right] \left[ \begin{array}{c} \bar{\mathbf{x}}_k - \mathbf{r}_k^{\mathbf{x}} \\ \bar{\mathbf{u}}_k - \mathbf{r}_k^{\mathbf{u}} \end{array} \right] + \left[ \begin{array}{c} \bar{\mathbf{x}}_N - \mathbf{r}_N^{\mathbf{x}} \end{array} \right]^T P \left[ \begin{array}{c} \bar{\mathbf{x}}_N - \mathbf{r}_N^{\mathbf{x}} \end{array} \right] \quad (\text{B.5a})$$

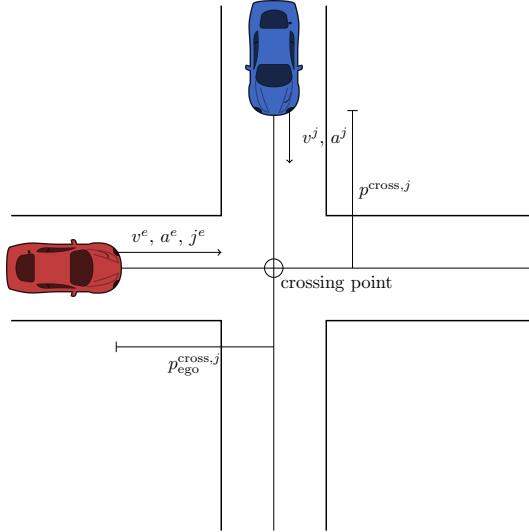
$$\text{s.t. } \bar{\mathbf{x}}_0 = \hat{\mathbf{x}}_0, \quad (\text{B.5b})$$

$$\bar{\mathbf{x}}_{k+1} = A\bar{\mathbf{x}}_k + B\bar{\mathbf{u}}_k, \quad (\text{B.5c})$$

$$h(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k, \bar{\mathbf{o}}_k, a_k) \leq 0, \quad (\text{B.5d})$$

where  $k$  is the prediction time index,  $N$  is the prediction horizon,  $Q$ ,  $R$ , and  $S$  are the stage costs,  $P$  is the terminal cost,  $\bar{\mathbf{x}}_k$  and  $\bar{\mathbf{u}}_k$  are the predicted state and control inputs,  $\mathbf{r}_k^{\mathbf{x}}$  and  $\mathbf{r}_k^{\mathbf{u}}$  are the state and control input references,  $\bar{\mathbf{o}}_k$  denotes the predicted state of vehicles in the environment which need to be avoided, and  $a$  is the action from the high-level decision maker. Constraint (B.5b) enforces that the prediction starts at the current state estimate  $\hat{\mathbf{x}}_0$ , (B.5c) enforces the system dynamics, and (B.5d) enforces constraints on the states, control inputs, and obstacle avoidance.

The reference points,  $\mathbf{r}_k^{\mathbf{x}}$ ,  $\mathbf{r}_k^{\mathbf{u}}$  are assumed to be set-points of a constant velocity trajectory, e.g. following the legal speed-limit on the road. Therefore, we set the velocity reference according to the driving limit, and the acceleration and jerk to zero.



**Figure 2:** Observations of a scenario

### Obstacle prediction

In order for the vehicle planner in (B.5) to be able to properly avoid collisions, it is necessary to provide information about the surrounding vehicles in the environment. Therefore, similarly to **batkovic2019**, we assume that a sensor system provides information about the environment, and that there exists a prediction layer which generates future motions of other vehicles in the environment. The accuracy of the prediction layer will heavily affect the performance of the planner, hence, it is necessary to have computationally inexpensive and accurate prediction methods.

In this paper, for simplicity the future motion of other agents is estimated by a constant velocity prediction model. The motion is predicted at every time instant for prediction times  $k \in [0, N]$ , and is used to form the collision avoidance constraints, which we describe in the next section. Even though more accurate prediction methods do exist, e.g. **lefeuvre2014survey**, **batkovic2018**, we use this simple model to show the potential of the overall framework.

### Collision avoidance

We denote a vehicle  $j$  with the following notation  $\mathbf{x}^j := [p^j \ v^j \ a^j]^\top$ , and an associated crossing point at position  $p^{\text{cross},j}$  in its own vehicle frame, which translated into the ego-vehicle frame is denoted as  $p_{\text{ego}}^{\text{cross},j}$ . With a predefined road topology, we assume that the vehicles will travel along the assigned paths, and that collisions may only occur at the crossing points  $p^{\text{cross},j}$  between an obstacle and the ego vehicle. Hence, for collision avoidance, we use the predictions of the future obstacle states  $\bar{\mathbf{x}}_k^j$  for times  $k \in [0, N]$ , provided by a prediction layer outside of the MPC framework. Given the obstacle measurements, the prediction layer will generate future states throughout the prediction horizon. With this information, it is possible to identify the time slots when an obstacle will enter the intersection.

Whenever an obstacle  $j$  is predicted to be within a threshold of  $p^{\text{cross},j}$ , e.g. the width of the intersecting area, the ego vehicle faces a constraint of the following form

$$\bar{p}_k^e \geq p_{\text{ego}}^{\text{cross},j} + \Delta, \quad \underline{p}_k^e \leq p_{\text{ego}}^{\text{cross},j} - \Delta,$$

where  $\Delta$  ensures sufficient padding from the crossing point that does not cause a collision. The choice of  $\Delta$  must be at least such that  $p_k$  together with the dimensions of the ego-vehicle does not overlap with the intersecting area.

### Take way and give way constraint

Since the constraints from the surrounding obstacles become non-convex, we rely on the high-level policy maker to decide through action  $a$  how to construct the constraint (B.5d) for Problem (B.5). The take-way action implies that the ego-vehicle drives first through the intersection, i.e., it needs to pass the intersection before all other obstacles. This implies that for any vehicle  $j$  that reaches the intersection during prediction times  $k \in [0, N]$ , the generated constraint needs to lower bound the state  $p_k$  according to

$$\max_j p^{\text{cross},j} + \Delta \leq p_k^e. \quad (\text{B.6})$$

Similarly, if the action is to give way, then the position needs to be upper bounded by the closest intersection point so that

$$p_k^e \leq \min_j p_{\text{ego}}^{\text{cross},j} - \Delta, \quad (\text{B.7})$$

for all times  $k$  that the obstacle is predicted to be in the intersection.

### Following an obstacle

For any action  $a$  that results in the following of an obstacle  $j$ , the ego-vehicle position is upper bounded by  $p_k^e \leq p_{\text{ego}}^{\text{cross},j}$ . We construct constraints for obstacles  $i \neq j$  according to

- if  $p^{\text{cross},i} < p^{\text{cross},j}$  then  $p^{\text{cross},i} + \Delta \leq p_k^e$ , which implies that the ego-vehicle should drive ahead of all obstacles  $i$  that are approaching the intersection;
- if  $p^{\text{cross},i} > p^{\text{cross},j}$  then  $p_k^e \leq p^{\text{cross},i} - \Delta$ , which implies that the ego-vehicle should wait to pass obstacle  $j$  and other obstacles  $i$ ;
- if  $p^{\text{cross},i} = p^{\text{cross},j}$  then the constraints generated for obstacle  $i$  becomes an upper or lower bound depending on if obstacle  $i$  is ahead or behind the obstacle  $j$  into the intersection.

## 4.2 Sliding mode agent

To benchmark the performance of using MPC, a SM controller that was used in **Tram2018LearningQ-Learning** is introduced.

$$a_{\text{sm}}^e = \frac{1}{c_2}(-c_1 x_2 + \mu \text{sign}(\sigma(x_1, x_2))), \quad (\text{B.8a})$$

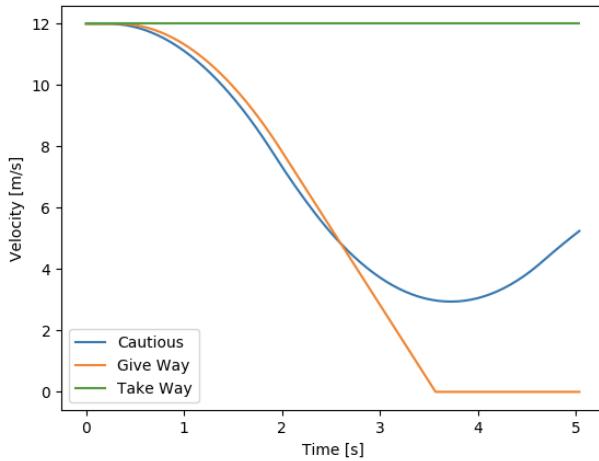
$$\text{where } \begin{cases} x_1 = p^t - p^e, \\ x_2 = v^t - v^e, \end{cases} \quad (\text{B.8b})$$

$$\sigma = c_1 x_1 + c_2 x_2, \quad (\text{B.8c})$$

$$a_p^e = K(v_{\max} - v^e), \quad (\text{B.8d})$$

$$a^e = \min(a_{\text{sm}}^e, a_p^e). \quad (\text{B.8e})$$

The SM controller aims to keep a minimum distance to a target car with a velocity of  $v^e$ , by controlling the acceleration  $a_{\text{sm}}^e$ .  $c_1$ ,  $c_2$ , and  $\mu$  are tuning parameters to control the comfort of the controller. In case there is no target car, the controller maintains a target velocity  $v_{\max}$  with a regular p-controller from (B.8d) with a proportional constant  $K$ . The final acceleration is given by (B.8e). For more details about the SM agent see **Tram2018LearningQ-Learning**.



**Figure 3:** Example of how velocity profile of the different intention agents can look like. All agents has the same starting velocity of 12m/s and are approaching the same intersection

### 4.3 Intention agents

There are three different intention agents for crossing traffic with predetermined intentions, three velocity profiles are shown as an example in Fig. 3. All agents were implemented with a SM controller with different target values. The take way intention does not yield for the crossing traffic and simply aim to keep its target reference speed. Give way intention however, slow down to a complete stop at the start of the intersection until crossing traffic has passed before continuing through. The third intention is cautious, slowing down but not to a full stop. This makes it difficult for a constant velocity or acceleration model to predict what other agents will do.

## 5 Implementation

### 5.1 Deep Q-Network

The deep Q-network is structured as a three layer neural network with shared weights and a Long Short-Term Memory based on previous work

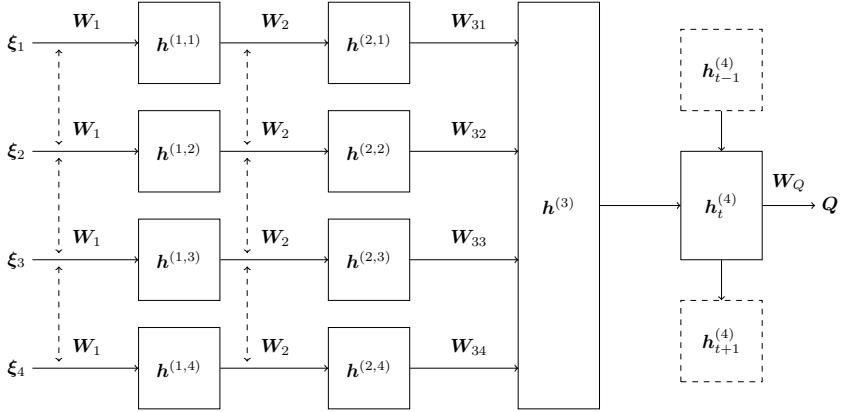


Figure 4: Representation of the network structure

**Tram2018LearningQ-Learning** and shown in Fig. 4. The input features  $\xi_n$  are composed of observations  $o_t$ , introduced in section 3.1 and shown in Fig. 2, with up to four observed vehicles:

$$\begin{aligned}\xi_1 &= [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^1 \quad v_t^1 \quad a_t^1 \quad \delta^1 ]^T \\ \xi_2 &= [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^2 \quad v_t^2 \quad a_t^2 \quad \delta^2 ]^T \\ \xi_3 &= [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^3 \quad v_t^3 \quad a_t^3 \quad \delta^3 ]^T \\ \xi_4 &= [ p_t^e \quad v_t^e \quad a_t^e \quad \delta^e \quad p_t^4 \quad v_t^4 \quad a_t^4 \quad \delta^4 ]^T\end{aligned}\tag{B.9}$$

Normalization of the input features is done by scaling the features down to values between  $[-1, 1]$  using the maximum speed  $v_{\max}$ , maximum acceleration  $a_{\max}$  and a car's sight range  $p_{\max}$ . Empty observation of other vehicles  $[p_t^n \quad v_t^n \quad a_t^n \quad \delta^n]$  has a default value of  $-1$ . The input vectors are sent through two hidden layers  $\mathbf{h}^{(1,i)}$  and  $\mathbf{h}^{(2,i)}$  with shared weights  $\mathbf{W}_1$  and  $\mathbf{W}_2$  respectively

$$\mathbf{h}^{(1,i)} = \tanh(\mathbf{W}_1 \xi_i + \mathbf{b}_1) \tag{B.10}$$

$$\mathbf{h}^{(2,i)} = \tanh(\mathbf{W}_2 \mathbf{h}^{(1,i)} + \mathbf{b}_2). \tag{B.11}$$

A similar study for lane changes on a highway confirmed the importance of having equals weights for inputs that describe the state of interchangeable objects **Hoel**. The output of each sub-network is then sent through a fully

connecting layer

$$\mathbf{h}^{(3)} = \tanh \left( \sum_{i=1}^4 \mathbf{W}_{3i} \mathbf{h}^{(2,i)} + \mathbf{b}_3 \right). \quad (\text{B.12})$$

That is then connected to an Long Short-Term Memory (LSTM) **Hochreiter1997LONGM** that can store and use previous features

$$\mathbf{h}_t^{(4)} = \text{LSTM} \left( \mathbf{h}^{(3)} | \mathbf{h}_{t-1}^{(4)} \right). \quad (\text{B.13})$$

The approximated Q-value is then

$$\mathbf{Q}_{approx} = \mathbf{W}_Q \mathbf{h}^{(4)} + \mathbf{b}_4 \quad (\text{B.14})$$

The Q-value is then masked using Q-masking, explain in section 5.2

$$\mathbf{Q} = \mathbf{Q}_{approx} \mathbf{Q}_{mask} \quad (\text{B.15})$$

the optimal policy  $\pi^*$  is then given by taking the action that gives the highest Q-value

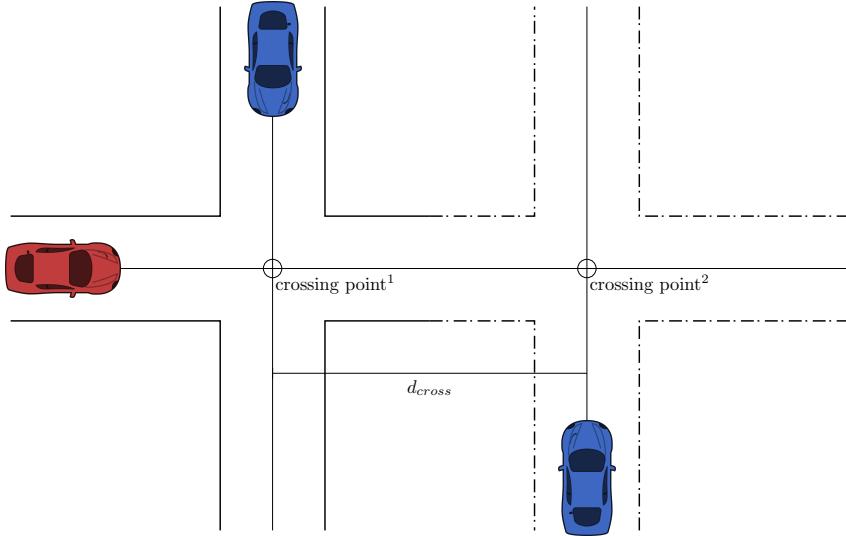
$$\pi^*(s_t) = \arg \max_{a_t} Q^*(s_t, a_t) \quad (\text{B.16})$$

## 5.2 Q-masking

Q-masking **Mukadam2017** helps the learning process by reducing the actions space by disabling actions the agent does not need to explore. If there are less than  $N$  cars, it would then be meaningless to choose to follow a car that does not exist. Which motivates masking off cars that does not exist. In previous work **Tram2018LearningQ-Learning**, a high negative reward was given when an action to follow a car that did not exist was chosen, while the algorithm continued with a default action take way. The agent quickly learned to not choose cars that did not exist, but with Q-masking, the agent does not even have to explore these options. For other details about the training see **Tram2018LearningQ-Learning**.

## 5.3 Simulation environment

All agents are spawned with random initial speed  $v_0 \in [10, 30]\text{m/s}$ , position  $p_0^i \in [10, 55]\text{m}$  and intention. The cars dimensions are 2 m wide and 4 m long. The ego car operates within comfort bounds and therefore has a limited



**Figure 5:** Illustration of an intersection scenario, where the solid line is a single crossing and together with the dashed line creates a double crossing.

maximum acceleration and deceleration of  $5 \text{ m/s}^2$ . Two main types of crossing were investigated. One and two crossing points as shown in Fig. 5, where the distance between crossing points  $d_{cross}$  vary between  $[4, 8, 12, 25, 30, 40] \text{ m}$  with different scenarios.

The MPC agent was discretized at 30Hz, with a prediction horizon of  $N = 100$  and cost tuning of

$$Q = \text{blockdiag}(0.0, 1.0, 1.0), \quad R = 1, \quad S = \mathbf{0}. \quad (\text{B.17})$$

## 5.4 Reward function tuning

There are three states that terminates an episode; success, failure, and timeout. Success is when the ego agent reaches the end of the road defined by the scenario. Failure is when the frame of the ego agent overlaps with another road users frame, e.g., in a collision, this frame can be the size of the vehicle or a safety boundary around a vehicle. The final terminating state is timeout and that is simply when the agent can't reach the two previous terminating states before the timeout time  $\tau_m$ . According to **VanHasseltLearningMagnitude**, the

$Q_\pi$  values and gradient can grow to be very large if the total reward values are too large. All rewards are therefore scaled with the episode timeout time  $\tau_m$ , which is set to 25s, to keep the total reward  $r_t \in [-2, 1]$ . The reward function is defined as follows:

$$r_t = \begin{cases} 1 & \text{on success,} \\ -1 & \text{on failure,} \\ 0.5 & \text{on timeout, i.e. } \tau \geq \tau_m, \\ f(p_{\text{crash}}, p_{\text{comf}}) & \text{on non-terminating updates,} \end{cases}$$

where  $f(p_{\text{crash}}, p_{\text{comf}})$  consists of

$$f(p_{\text{crash}}, p_{\text{comf}}) = \alpha p_{\text{crash}} \frac{\tau_m}{\tau - t_{\text{pred}}} + \beta p_{\text{comf}} \frac{\tau_m}{\tau}, \quad (\text{B.18})$$

with  $\alpha \in [0, 1]$ ,  $\beta \in [0, 1]$  being weight parameters, and  $\alpha + \beta = 1$ . The first term in the function corresponds to a feasibility check of Problem (B.5), which to a large extent depends on the validity of the accuracy of the prediction layer. The high-level decision from the policy-maker affects how the constraints are constructed, and may turn the control problem infeasible, e.g. if the decided action is to take way, while not being able to pass the intersection before all other obstacles. Therefore, whenever the MPC problem becomes infeasible we set  $p_{\text{crash}} = 1$  to indicate that the selected action most likely will result in a collision with the surrounding environment.

The second term  $p_{\text{comf}}$  relates to the comfort of the planned trajectory, which is estimated by computing and weighting the acceleration and jerk profiles as

$$p_{\text{comf}} = \frac{1}{\sigma N} \left( \sum_{k=0}^{N-1} \bar{a}_k^2 Q^a + \bar{j}_k^2 R^j + a_N^2 Q^a \right),$$

where  $\bar{a}$ , and  $\bar{j}$  are the acceleration and jerk components of the state and control input respectively,  $Q^a$  and  $R^j$  are the corresponding weights, and  $\sigma$  is a normalizing factor which ensures that  $p_{\text{comf}} \in [0, 1]$ . For the simulation we used  $Q^a = 1$  and  $R^j = 1$ .

The timeout reward 0.5 was set to be higher than the average accumulated reward from  $p_{\text{comf}}$ , so that the total accumulated reward would be positive in case of timeouts. Because  $p_{\text{crash}}$  usually only triggers close to a potential collision, that is why  $t_{\text{pred}}$  is set to the first time a crash prediction is triggered. This will scale the negative reward higher in collision episodes.

**Table 1:** Average success rates and collision to timeout rates.

Controller	Success Rate		Timeout Ratio	
	Single	Double	Single	Double
SM	96.1%	90.9%	72%	93%
MPC	97.3%	95.2%	45%	76%

## 6 Results

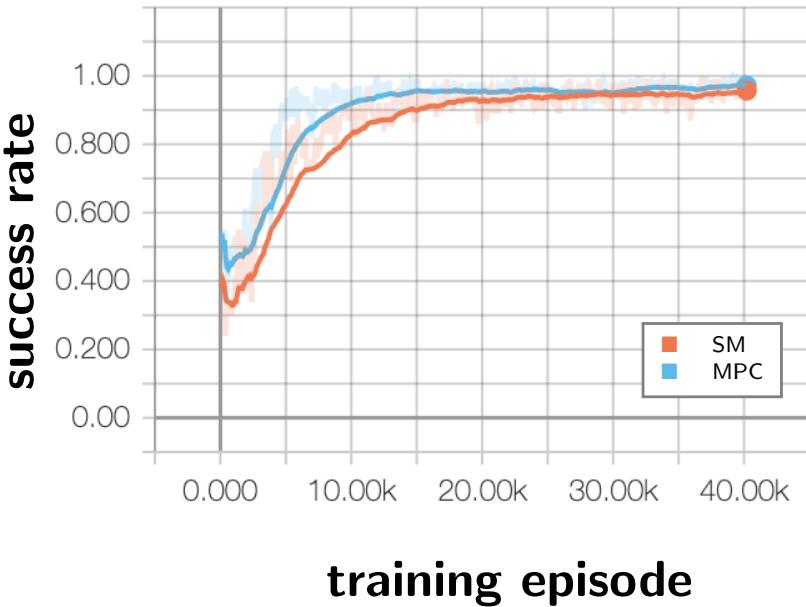
For evaluation we compared the success rate of the decision-policy together with a collision to timeout ratio (CTR). The success rate is defined as the number of times the agent is able cross the intersections without colliding with other obstacles, or exceeding the time limit to cross. Since we define a time-out to be a failure, we use the CTR to separate potential collisions with the agent being too conservative.

Fig. 6 shows a comparison in success rate between the proposed MPC architecture and the previous SM agent for scenarios with only one crossing. In this scenario, the MPC agent converges after  $10^4$  training episodes, while the previous SM agent converges after  $4 \cdot 10^4$  training episodes. In addition, comparing the CTR metric, Fig. 7 shows that the MPC agent has 0.45 CTR while the SM agent has 0.72 CTR. Evidently, it is visible that the MPC is able to leverage future information into its planning horizon in order to achieve faster training, and also avoiding collisions as a result.

We evaluate the performance of the MPC and SM agents for the more difficult double intersection problem, where we vary the distance between the intersection points. Table 1 shows the performance of the MPC and SM agent for both the single and double scenarios. The performance drops for both agents for the double crossing scenario. However, it is visible that the MPC agent suffers less performance degradation compared to the SM agent. The CTR however more than doubles for the MPC agent for the double crossing, while the already high CTR rate for the SM agent increases above rates of 0.9.

## 7 Discussion

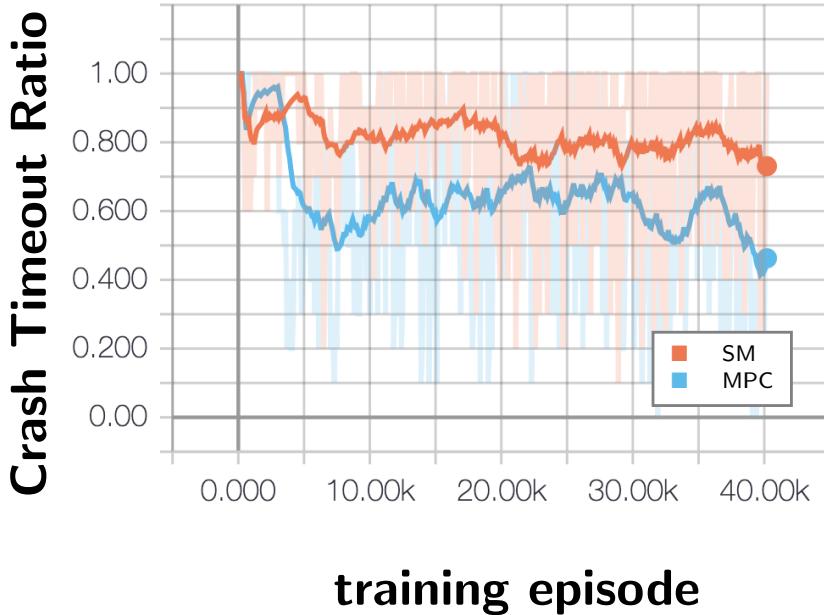
The benefit of being able to use a prediction horizon for the MPC is shown to mostly impact the training time for the traffic scenarios compared to the



**Figure 6:** Average MPC and SM success rate for a single crossing after evaluating the policy 300 episodes.

SM agent. This allows the RL decision-policy to get feedback early in the training process to see whether an action most likely will lead to a collision. In addition, the lower CTR also implies that the use of a prediction horizon also makes the decision-policy more conservative, since it rather times out than risk collisions.

It is important to note little effort was put into tuning the MPC agent, and that we used very primitive prediction methods that do not hold very well in crossing scenarios, e.g. the simulated agents did not keep constant speed profiles while approaching the intersections. However, under these circumstances, the decision algorithm still managed to obtain a success rate above 95% for the double crossings.



**Figure 7:** Average MPC and SM crash to timeout ratio for a single crossing after evaluating the policy in 300 episodes. A CTR of 0 means that all failures are timeouts, while a CTR of 1 means that all failures are collisions.

## 8 Conclusion

In this paper, we proposed a decision making algorithm for intersections which consists of two components: a high-level decision maker that uses Deep Q-learning to generate decisions for how the vehicle should drive through the intersection, and a low-level planner that uses MPC to optimize safe trajectories. We tested the framework in a traffic simulation with randomized intent of other road users for both single and double crossings. Results showed that the proposed MPC agent outperforms the previous SM agent by almost 5% in scenarios with double crossings and in cases of failure, more often timeout than colliding.



# PAPER C

## Reinforcement Learning with Uncertainty Estimation for Tactical Decision-Making in Intersections

Carl-Johan Hoel, Tommy Tram, and Jonas Sjöberg

*Published in 2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC),  
pp. 1-7, Sep. 2020.*

©2020 IEEE DOI: 10.1109/ITSC45102.2020.9294407.

*The layout has been revised.*

## Abstract

This paper investigates how a Bayesian reinforcement learning method can be used to create a tactical decision-making agent for autonomous driving in an intersection scenario, where the agent can estimate the confidence of its decisions. An ensemble of neural networks, with additional randomized prior functions (RPF), are trained by using a bootstrapped experience replay memory. The coefficient of variation in the estimated  $Q$ -values of the ensemble members is used to approximate the uncertainty, and a criterion that determines if the agent is sufficiently confident to make a particular decision is introduced. The performance of the ensemble RPF method is evaluated in an intersection scenario and compared to a standard Deep Q-Network method, which does not estimate the uncertainty. It is shown that the trained ensemble RPF agent can detect cases with high uncertainty, both in situations that are far from the training distribution, and in situations that seldom occur within the training distribution. This work demonstrates one possible application of such a confidence estimate, by using this information to choose safe actions in unknown situations, which removes all collisions from within the training distribution, and most collisions outside of the distribution.

## 1 Introduction

To make safe, efficient, and comfortable decisions in intersections is one of the challenges of autonomous driving. A decision-making agent needs to handle a diverse set of intersection types and layouts, interact with other traffic participants, and consider uncertainty in sensor information. The fact that around 40% of all traffic accidents during manual driving occur in intersections indicates that decision-making in intersections is a complex task **NHTSA**. To manually predict all situations that can occur and tailor a suitable behavior is not feasible. Therefore, a data-driven approach that can learn to make decisions from experience is a compelling approach. A desired property of such a machine learning approach is that it should also be able to indicate

how confident the resulting agent is about a particular decision.

Reinforcement learning (RL) provides a general approach to solve decision-making problems [1], and could potentially scale to all types of driving situations. Promising results have been achieved in simulation by applying a Deep Q-Network (DQN) agent to intersection scenarios **Isele2018**, **Tram2018**, and highway driving **Wang2018**, **Hoel2018**, or a policy gradient method to a lane merging situation **Shalev2016**. Some studies have trained an RL agent in a simulated environment and then deployed the agent in a real vehicle **Pan2017**, **Bansal2018**, and for a limited case, trained the agent directly in a real vehicle **Kendall2017**.

Generally, a fundamental problem with the RL methods in previous work is that the trained agents do not provide any confidence measure of their decisions. For example, if an agent that was trained for a highway driving scenario would be exposed to an intersection situation, it would still output a decision, although it would likely not be a good one. A less extreme example involves an agent that has been trained in an intersection scenario with nominal traffic, and then faces a speeding driver. McAllister et al. further discuss the importance of estimating the uncertainty of decisions in autonomous driving **McAllister2017**.

A common way of estimating uncertainty is through Bayesian probability theory [2]. Bayesian deep learning has previously been used to estimate uncertainty in autonomous driving for image segmentation **Kendall2017** and end-to-end learning **Michelmore2018**. Dearden et al. introduced Bayesian approaches to RL that balances the trade off between exploration and exploitation **Dearden1998**. In recent work, this approach has been extended to deep RL, by using an ensemble of neural networks **Osband2018**. However, these studies focus on creating an efficient exploration method for RL, and do not provide a confidence measure for the agents' decisions.

This paper investigates an RL method that can estimate the uncertainty of the resulting agent's decisions, applied to decision-making in an intersection scenario. The RL method uses an ensemble of neural networks with randomized prior functions that are trained on a bootstrapped experience replay memory, which gives a distribution of estimated  $Q$ -values (Sect. 2). The distribution of  $Q$ -values is then used to estimate the uncertainty of the recommended action, and a criterion that determines the confidence level of the agent's decision is introduced (Sect. 2.3). The method is used to train a decision-making agent

in different intersection scenarios (Sect. 3), in which the results show that the introduced method outperforms a DQN agent within the training distribution. The results also show that the ensemble method can detect situations that were not present in the training process, and thereby choose safe fallback actions in such situations (Sect. 4). Further characteristics of the introduced method is discussed in Sect. 5. This work is an extension to a recent paper, where we introduced the mentioned method, but applied to a highway driving scenario **Hoel2020**.

## 2 Approach

This section gives a brief introduction to RL, describes how the uncertainty of an action can be estimated by an ensemble method, and introduces a measure of confidence for different actions. Further details on how this approach was applied to driving in an intersection scenario follows in Sect. 3.

### 2.1 Reinforcement learning

Reinforcement learning is a branch of machine learning, where an agents explores an environment and tries to learn a policy  $\pi(s)$  that maximizes the future expected return, based on the agent's experiences [1]. The policy determines which action  $a$  to take in a given state  $s$ . The state of the environment will then transitions to a new state  $s'$  and the agent receives a reward  $r$ . A Markov Decision Process (MDP) is often used to model the reinforcement learning problem. An MDP is defined by the tuple  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $T$  is a state transition model,  $R$  is a reward model, and  $\gamma$  is a discount factor. At each time step  $t$ , the agent tries to maximize the future discounted return

$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}. \quad (\text{C.1})$$

In a value-based branch of RL called  $Q$ -learning **Watkins1992**, the objective of the agent is to learn the optimal state-action value function  $Q^*(s, a)$ . This function is defined as the expected return when the agent takes action  $a$  from state  $s$  and then follow the optimal policy  $\pi^*$ , i.e.,

$$Q^*(s, a) = \max_{\pi} \mathbb{E} [R_t | s_t = s, a_t = a, \pi]. \quad (\text{C.2})$$

The  $Q$ -function can be estimated by a neural network with weights  $\theta$ , i.e.,  $Q(s, a) \approx Q(s, a; \theta)$ . The weights are optimized by minimizing the loss function

$$L(\theta) = \mathbb{E}_M \left[ (r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta))^2 \right], \quad (\text{C.3})$$

which is derived from the Bellman equation. The loss is obtained from a mini-batch  $M$  of training samples, and  $\theta^-$  represents the weights of a target network that is updated regularly. More details on the DQN algorithm are presented by Mnih et al. **Mnih2015**.

## 2.2 Bayesian reinforcement learning

One limitation of the DQN algorithm is that only the maximum likelihood estimate of the  $Q$ -values is returned. The risk of taking a particular action can be approximated as the variance in the estimated  $Q$ -value **Garcia2015**. One approach to obtain a variance estimation is through statistical bootstrapping **Efron1982**, which has been applied to the DQN algorithm **Osband2016**. The basic idea is to train an ensemble of neural network on different subsets of the available replay memory. The ensemble will then provide a distribution of  $Q$ -values, which can be used to estimate the variance. Osband et al. extended the ensemble method by adding a randomized prior function (RPF) to each ensemble member, which gives a better Bayesian posterior **Osband2018**. The  $Q$ -values of each ensemble member  $k$  is then calculated as the sum of two neural networks,  $f$  and  $p$ , with equal architecture, i.e.,

$$Q_k(s, a) = f(s, a; \theta_k) + \beta p(s, a; \hat{\theta}_k). \quad (\text{C.4})$$

Here, the weights  $\theta_k$  of network  $f$  are trainable, and the weights  $\hat{\theta}_k$  of the prior network  $p$  are fixed to the randomly initialized values. A parameter  $\beta$  scales the importance of the networks. With the two networks, the loss function in Eq. C.3 becomes

$$\begin{aligned} L(\theta_k) = \mathbb{E}_M & \left[ (r + \gamma \max_{a'} (f_{\theta_k^-} + \beta p_{\hat{\theta}_k})(s', a') \right. \\ & \left. - (f_{\theta_k} + \beta p_{\hat{\theta}_k})(s, a))^2 \right]. \end{aligned} \quad (\text{C.5})$$

Algorithm 1 outlines the complete ensemble RPF method, which was used in this work. An ensemble of  $K$  trainable and prior neural networks are

---

**Algorithm 1** Ensemble RPF training process

---

```

1: for  $k \leftarrow 1$  to  $K$  do
2:   Initialize  $\theta_k$  and  $\hat{\theta}_k$  randomly
3:    $m_k \leftarrow \{\}$ 
4:    $i \leftarrow 0$ 
5:   while networks not converged do
6:      $s_i \leftarrow$  initial random state
7:      $\nu \sim \mathcal{U}\{1, K\}$ 
8:     while episode not finished do
9:        $a_i \leftarrow \text{argmax}_a Q_\nu(s_i, a)$ 
10:       $s_{i+1}, r_i \leftarrow \text{STEPENVIRONMENT}(s_i, a_i)$ 
11:      for  $k \leftarrow 1$  to  $K$  do
12:        if  $p \sim \mathcal{U}(0, 1) < p_{\text{add}}$  then
13:           $m_k \leftarrow m_k \cup \{(s_i, a_i, r_i, s_{i+1})\}$ 
14:         $M \leftarrow$  sample mini-batch from  $m_k$ 
15:        update  $\theta_k$  with SGD and loss  $L(\theta_k)$ 
16:       $i \leftarrow i + 1$ 

```

---

first initialized randomly. Each ensemble member is also assigned a separate experience replay memory buffer  $m_k$  (although in a practical implementation, the replay memory can be designed in such a way that it uses negligible more memory than a shared buffer). For each new training episode, a uniformly sampled ensemble member,  $\nu \sim \mathcal{U}\{1, K\}$ , is used to greedily select the action with the highest  $Q$ -value. This procedure handles the exploration vs. exploitation trade-off and corresponds to a form of approximate Thompson sampling. Each new experience  $e = (s_i, a_i, r_i, s_{i+1})$  is then added to the separate replay buffers  $m_k$  with probability  $p_{\text{add}}$ . Finally, the trainable weights of each ensemble member are updated by uniformly sample a mini-batch  $M$  of experiences and using stochastic gradient descent (SGD) to backpropagate the loss of Eq. C.5.

### 2.3 Confidence criterion

The agent's uncertainty in choosing different actions can be defined as the coefficient of variation<sup>1</sup>  $c_v(s, a)$  of the  $Q$ -values of the ensemble members. In

---

<sup>1</sup>Ratio of the standard deviation to the mean.

previous work, we introduced a confidence criterion that disqualifies actions with  $c_v(s, a) > c_v^{\text{safe}}$ , where  $c_{\text{safe}}$  is a hard threshold **Hoel2020**. The value of the threshold should be set so that  $(s, a)$  combinations that are contained in the training distribution are accepted, and those which are not will be rejected. This value can be determined by observing values of  $c_v$  in testing episodes within the training distribution, see Sect. 4.1 for further details.

When the agent is fully trained (i.e., not during the training phase), the policy chooses actions by maximizing the mean of the  $Q$ -values of the ensemble members, with the restriction  $c_v(s, a) < c_v^{\text{safe}}$ , i.e.,

$$\begin{aligned} \operatorname{argmax}_a \frac{1}{K} \sum_{k=1}^K Q_k(s, a), \\ \text{s.t. } c_v(s, a) < c_v^{\text{safe}}. \end{aligned} \quad (\text{C.6})$$

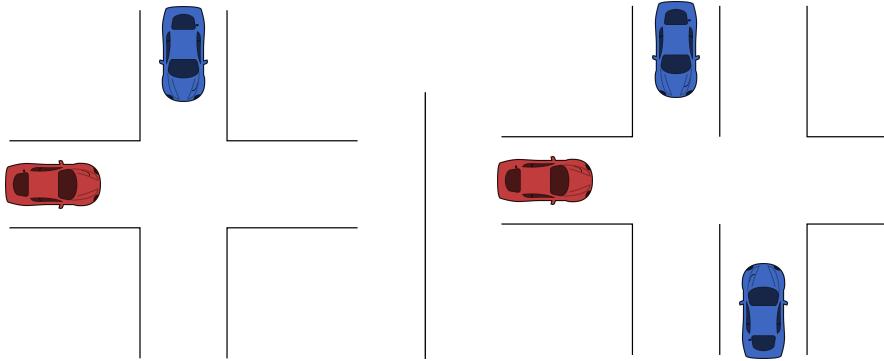
In a situation where no possible action fulfills the confidence criterion, a fallback action  $a_{\text{safe}}$  is chosen.

## 3 Implementation

The ensemble RPF method, which can obtain an uncertainty estimation of different actions, is tested on different intersection scenarios. In this work, the uncertainty information is used to reject unsafe actions and reduce the number of collisions. This section describes how the simulation of the scenarios is set up, how the decision-making problem is formulated as an MDP, the architecture of the neural networks, and the details on how the training is performed.

### 3.1 Simulation setup

The simulated environment consists of different intersection scenarios, and is based on previous work **tram2019**. For completeness, an overview is presented here. Each episode starts by randomly selecting a single or bi-directional intersection, shown in Fig. 1, and placing the ego vehicle to the left with a random distance  $p_e^{c,j}$  to the intersection and a speed of 10 m/s. A random number  $N$  of other vehicles are positioned along the top and bottom roads with a random distance  $p_o^{c,j}$  to the intersection, and a random desired speed  $v_d^j$ . The other vehicles follow the Intelligent Driver Model (IDM) **idm2000**, with



**Figure 1:** The two intersection scenarios considered in this work; single directional to the left and bidirectional to the right. The agent controls the red car.

**Table 1:** Parameters for simulator

Number of other vehicles, $N$	$\{1, 2, 3, 4\}$
Starting position ego, $p_e^{c,j}$	[50, 60] m
Starting position target, $p_o^{c,j}$	[10, 55] m
Desired velocity, $v_d^j$	[8, 12] m/s

a set time gap of  $t_d^j = 1$  s. One quarter of the vehicles stop at the intersection and three quarters continue through the intersection, regardless of the behavior of the ego vehicle. When a vehicle has passed the intersection and reached the end of the road, it is moved back to the other side of the intersection, which creates a constant traffic flow. The simulator is updated at 25 Hz, and decisions are taken at 4 Hz. The goal of the ego vehicle is to reach a position that is located 10 m to the right of the last crossing point.

### 3.2 MDP formulation

The following Markov decision process is used to model the decision-making problem. The full state is not directly observable, since the intentions of the surrounding vehicles are not known to the agent. Therefore, the problem is a Partially Observable Markov Decision Process (POMDP) [Kaelbling1998](#).

However, by using a  $k$ -Markov approximation, where the state consists of the  $k$  last observations, the POMDP can be approximated as an MDP **Mnih2015**. For the scenarios that were considered in this work, it proved sufficient to simply use the last observation.

### State space, $\mathcal{S}$

The design of the state of the system,

$$s = (p_e^g, v_e, a_e, \{p_e^{s,j}, p_e^{c,j}, p_o^{s,j}, p_o^{c,j}, v_o^j, a_o^j\}_{j=0,\dots,N}), \quad (\text{C.7})$$

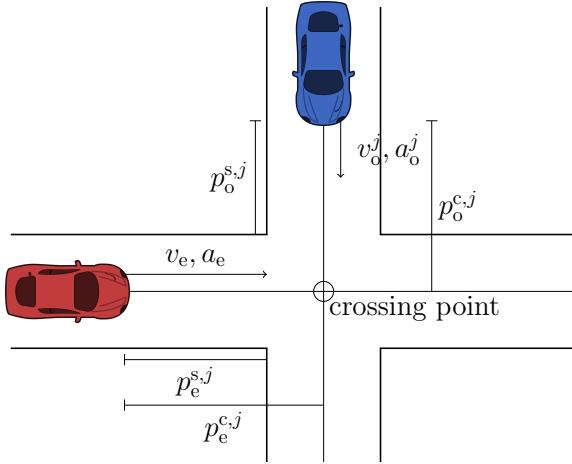
allows the description of intersections with different layouts **Tram2018**. The state, illustrated in Fig. 2, consists of the distance from the ego vehicle to the goal  $p_e^g$ , the velocity and acceleration of the ego vehicle,  $v_e$ ,  $a_e$ , and the other vehicles,  $v_o^j$ ,  $a_o^j$ , where  $j$  denotes the index of the other vehicles. Furthermore,  $p_e^{s,j}$  and  $p_e^{c,j}$  are the distances from the ego vehicle to the start of the intersection and crossing point, relative to target vehicle  $j$  respectively. The distances  $p_o^{s,j}$  and  $p_o^{c,j}$  are the distance from the other vehicles to the start of the intersection and the crossing point.

### Action space, $\mathcal{A}$

The action space consists of six tactical decisions:  $\{\text{'take way'}, \text{'give way'}, \text{'follow car } \{1, \dots, 4\}\}$ , which set the target of the IDM controller. The ‘take way’ action treats the situation as an empty road, whereas the ‘give way’ action sets a target distance of  $p_e^{s,j}$  and a target speed of 0 m/s. The ‘follow car  $j$ ’ actions sets the target distance to  $p_e^{c,j} - p_o^{c,j}$  and target speed to  $v_o^j$ . In cases where  $p_o^{c,j} > p_e^{c,j}$ , the target distance is set to a value that corresponds to timegap 0.5 s. The output of the IDM model is further limited by a maximum jerk  $j_{\max} = 5 \text{ m/s}^3$  and maximum acceleration  $a_{\max} = 5 \text{ m/s}^2$ . If less than four vehicles are present, the actions that correspond to choosing an absent vehicle are pruned by using Q-masking **Mukadam2017**.

### Reward model, $R$

The objective of the agent is to reach the goal on the other side of the intersection, without colliding with other vehicles and for comfort reasons, with as little jerk  $j_t$  as possible. Therefore, the reward at each time step  $r_t$  is



**Figure 2:** The state space definitions for a single crossing scenario, where subscript e and o denotes ego and other vehicle, respectively.

defined as

$$r_t = \begin{cases} 1 & \text{at reaching the goal,} \\ -1 & \text{at a collision,} \\ -\left(\frac{j_t}{j_{\max}}\right)^2 \frac{\Delta\tau}{\tau_{\max}} & \text{at non-terminating steps.} \end{cases}$$

The non-terminating reward is scaled with the maximum time of an episode,  $\tau_{\max}$ , and the step time  $\Delta\tau = 0.04$  s, to ensure  $\sum_{t=0}^{t=\tau_{\max}} \in [-1, 0]$ . Further details about the reward function can be found in previous research **tram2019**.

### Transition model, $T$

The state transition probabilities are not known to the agent. However, the true transition model is defined by the simulation model, described in Sect. 3.1.

### 3.3 Fallback action

As mentioned in Sect. 2.3, a fallback action  $a_{\text{safe}}$  is used when  $c_v > c_v^{\text{safe}}$  for all available actions. This fallback action is set to ‘give way’, with the

difference that no jerk limitation is applied and with a higher acceleration limit  $a_{\max} = 10 \text{ m/s}^2$ .

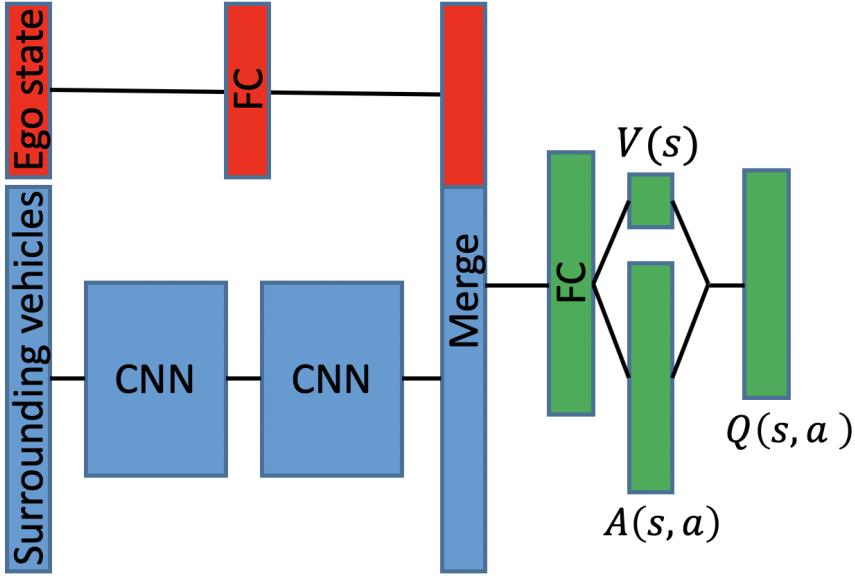
### 3.4 Network architecture

In previous studies, we have showed that a network architecture that applies the same weights to the input that describes the surrounding vehicles results in a better performance and speeds up the training process **Hoel2018, Tram2018**. Such an architecture can be constructed by applying a one-dimensional convolutional neural network (CNN) structure to the surrounding vehicles' input. The network architecture that is used in this work is shown in Fig. 3. The first convolutional layer has 32 filters, with size and stride set to six, which equals the number of state inputs of each surrounding vehicle, and the second convolutional layers has 16 filter, with size and stride set to one. The fully connected (FC) layer that is connected to the ego vehicle input has 16 units, and the joint fully connected layer has 64 units. All layers use rectified linear units (ReLUs) as activation functions, except for the last layer, which has a linear activation function. The final dueling structure of the network separates the estimation of the state value  $V(s)$  and the action advantage  $A(s, a)$  **Wang2016**. The input vector is normalized to the range  $[-1, 1]$ . The input vector contains slots for four surrounding vehicles, and if less vehicles are present in the traffic scene, the empty input is set to  $-1$ .

### 3.5 Training process

Algorithm 1 is used to train the agent. The loss function of Double DQN is applied, which subtly modifies the maximization operation of Eq. C.3 to  $\gamma Q(s', \text{argmax}_{a'} Q(s', a'; \theta_i); \theta_i^-)$  **Hasselt2016**. The Adam optimizer is used to update the weights **Kingma2014**, and  $K$  parallel workers are used for the backpropagation step. The hyperparameters of the training process are shown in Table 2, and the values were selected by an informal search, due to the computational complexity.

If the current policy of the agent decides to stop the ego vehicle, an episode could continue forever. Therefore, a timeout time is set to  $\tau_{\max} = 20 \text{ s}$ , at which the episode terminates. The last experience of such an episode is not added to the replay memory. This trick prevents the agent to learn that an episode can end due to a timeout, and makes it seem like an episode can



**Figure 3:** The neural network architecture that was used in this work.

continue forever, which is important, since the terminating state due to the time limit is not part of the MDP **Hoel2018**.

### 3.6 Baseline method

The Double DQN method, hereafter simply referred to as the DQN method, is used as a baseline. For a fair comparison, the same hyperparameters as for the ensemble RPF method is used, with the addition of an annealing  $\epsilon$ -greedy exploration schedule, which is shown in Table 2. During test episodes, a greedy policy is used.

## 4 Results

The results show that the ensemble RPF method outperforms the DQN method, both in terms of training speed and final performance, when the resulting

**Table 2:** Hyperparameters of Algorithm 1 and baseline DQN.

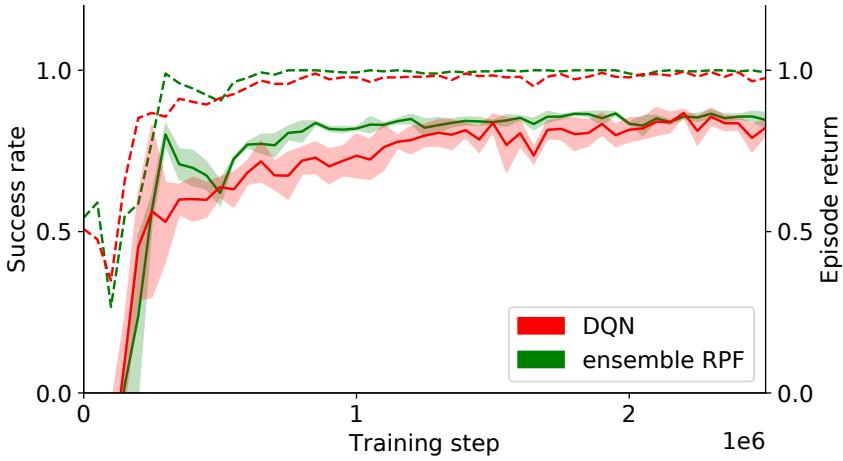
Number of ensemble members, $K$	10
Prior scale factor, $\beta$	1
Experience adding probability, $p_{\text{add}}$	0.5
Discount factor, $\gamma$	0.99
Learning start iteration, $N_{\text{start}}$	50,000
Replay memory size, $M_{\text{replay}}$	500,000
Learning rate, $\eta$	0.0005
Mini-batch size, $M_{\text{mini}}$	32
Target network update frequency, $N_{\text{update}}$	20,000
Huber loss threshold, $\delta$	10
Initial exploration constant, $\epsilon_{\text{start}}$	1
Final exploration constant, $\epsilon_{\text{end}}$	0.05
Final exploration iteration, $N_{\epsilon-\text{end}}$	1,000,000

agents are tested on scenarios that are similar to the training scenarios. When the fully trained ensemble RPF agent is exposed to situations that are outside of the training distribution, the agent indicates a high uncertainty and chooses safe actions, whereas the DQN agent collides with other vehicles. More details on the characteristics of the results are presented and briefly discussed in this section, whereas a more general discussion follows in Sect. 5.

The ensemble RPF and DQN agents were trained in the simulated environment that was described in Sect. 3. After every 50,000 training steps, the performance of the agents were evaluated on 100 random test episodes. These test episodes were randomly generated in the same way as the training episodes, but kept fixed for all the evaluation phases.

#### 4.1 Within training distribution

The average return and the average proportion of episodes where the ego vehicle reached the goal, as a function of number of training steps, is shown in Fig. 4, for the test episodes. The figure also shows the standard deviation

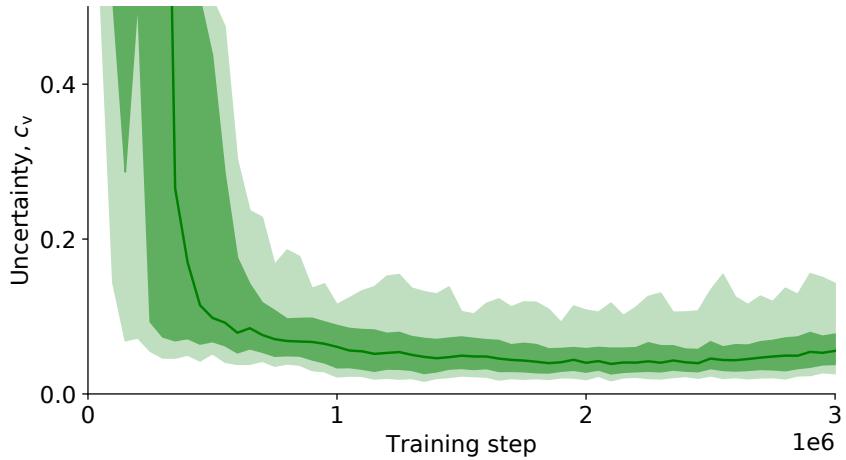


**Figure 4:** Proportion of test episodes where the ego vehicle reached its goal (dashed), and episode return (solid), over training steps for the ensemble RPF and DQN methods. The shaded areas show the standard deviation for 5 random seeds.

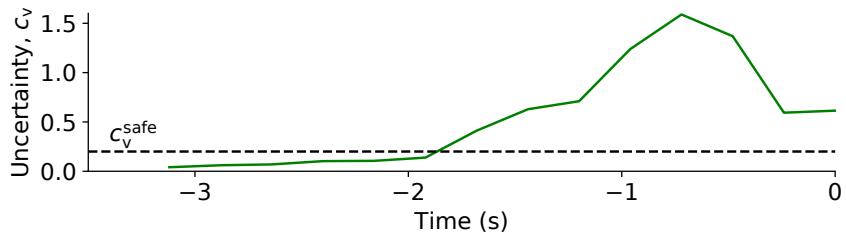
for 5 random seeds, which generates different sets of initial parameters of the networks and different training episodes, whereas the test episodes are kept fixed. The results show that the ensemble RPF method both learns faster, yields a higher return, and causes less collisions than the DQN method.

Fig. 5 shows how the coefficient of variation  $c_v$  of the chosen action varies during the testing episodes. Note that the uncertainty of actions that are not chosen can be higher, which is often the case. After around one million training steps, the average value of  $c_v$  settles at around 0.04, with a 99 percentile value of 0.15, which motivates the choice of setting  $c_v^{\text{safe}} = 0.2$ .

As shown in Fig. 4, occasional collisions still occur during the test episodes when deploying the fully trained ensemble RPF agent. The reasons for these collisions are further discussed in Sect. 5. In one particular example of a collision, the agent fails to brake early enough and ends up in an impossible situation, where it collides with another vehicle in the intersection. However, the estimated uncertainty increases significantly during the time before the collision, when the incorrect actions are taken, see Fig. 6. When applying the



**Figure 5:** Mean coefficient of variation  $c_v$  for the chosen action during the test episodes. The dark shaded area shows percentiles 10 to 90, and the bright shaded area shows percentiles 1 to 99.



**Figure 6:** Uncertainty  $c_v$  during the time steps before one of the collisions in the test episodes, within the training distribution. The collision occurs at  $t = 0$  s.

confidence criterion (Sect. 2.3), the agent instead brakes early enough, and can thereby avoid the collision. The confidence criterion was also applied to all the test episodes, which removed all collisions.

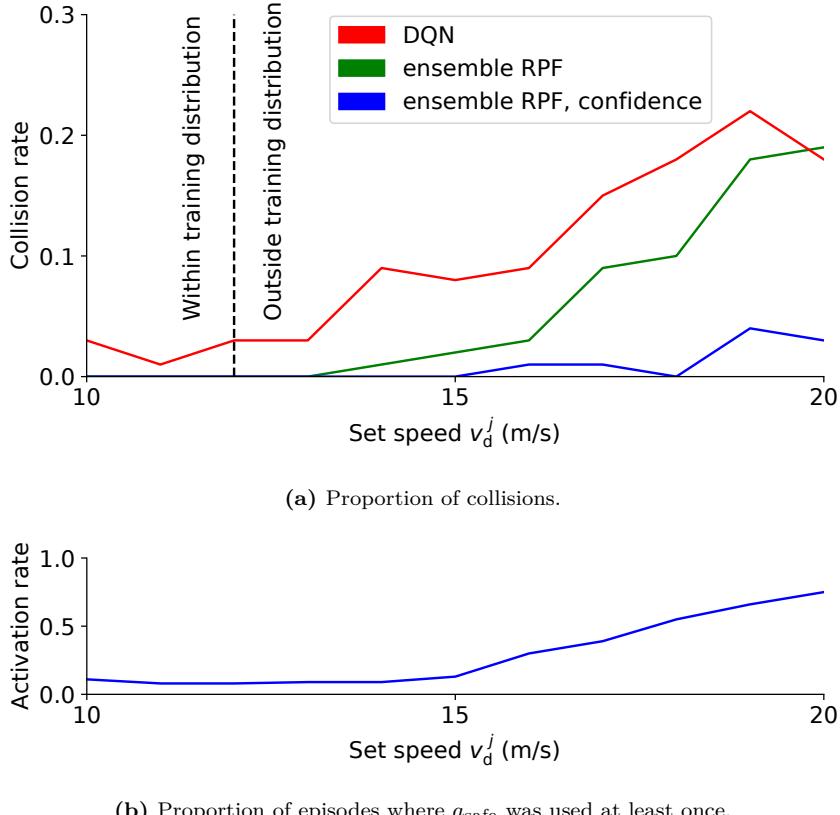
## 4.2 Outside training distribution

The ensemble RPF agent that was obtained after three million training steps was tested in scenarios outside of the training distribution, in order to evaluate the agent’s ability to detect unseen situations. The same testing scenarios as for within the distribution was used, with the exception that the speed of the surrounding vehicles was set to a single deterministic value, which was varied during different runs in the range  $v_d^j = [10, 20]$  m/s. The proportion of collisions as a function of set speed of the surrounding vehicles is shown in Fig. 7, together with the proportion of episodes where the confidence criterion was violated at least once. The figure shows that when the confidence criterion is used, most of the collisions can be avoided. Furthermore, the violations of the criterion increase when the speed of the surrounding vehicles increase, i.e., the scenarios move further from the training distribution.

An example of a situation that causes a collision is shown in Fig. 8, where an approaching vehicle drives with a speed of 20 m/s. The  $Q$ -values of both the trained ensemble RPF and DQN agents indicate that the agents expect to make it over the crossing before the other vehicle. However, since the approaching vehicle drives faster than what the agents have seen during the training, a collision occurs. When the confidence criterion is applied, the uncertainty rises to  $c_v > c_v^{\text{safe}}$  for all actions when the ego vehicle approaches the critical region, where it has to brake in order to be able to stop, and a collision is avoided by choosing action  $a_{\text{safe}}$ .

## 5 Discussion

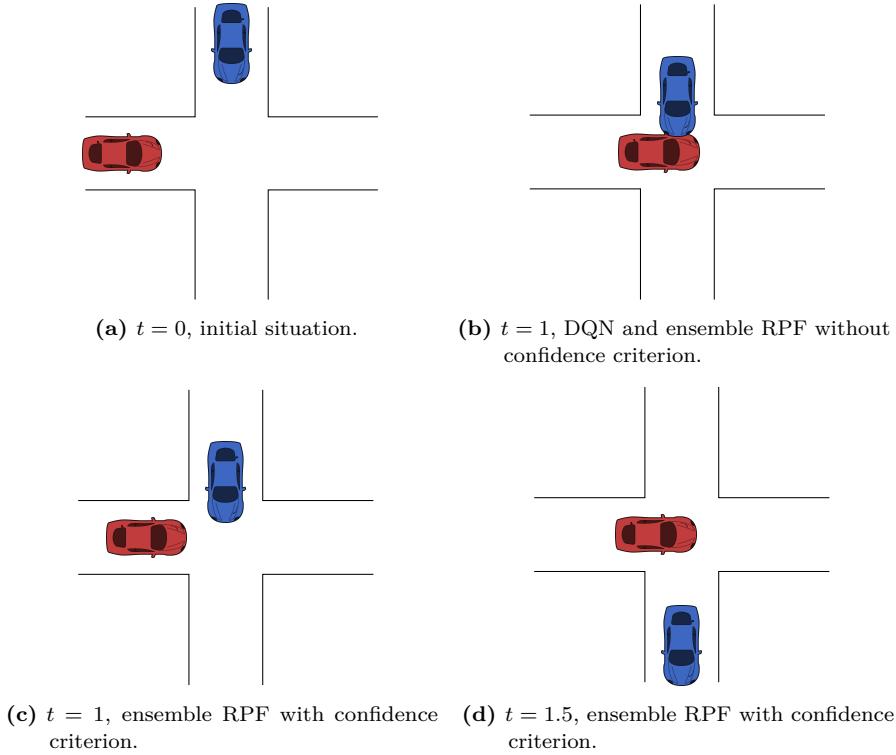
The results show that the ensemble RPF method can indicate an elevated uncertainty for situations that the agent has been insufficiently trained for, both within and outside of the training distribution. In previous work by the authors of this paper, we observed similar results when using the ensemble RPF method to estimate uncertainty outside of the training distribution in a highway driving scenario **Hoel2020**. In contrast, this paper shows that, in some cases, the ensemble RPF method can even detect situations with high uncertainty within the training distribution. Such situations include rare events that seldom or never occur during the training process, which makes it hard for the agent to provide an accurate estimate of the  $Q$ -values for the corresponding states. Since these states are seldom used to update the neural



**Figure 7:** Performance of the ensemble RPF agent, with and without the confidence criterion, and the DQN agent, in test episodes with different set speeds  $v_d^j$  for the surrounding vehicles.

networks of the ensemble, the weights of the trainable networks will not adapt to the respective prior networks, and the uncertainty measure  $c_v$  will remain high for these rare events. This information is useful to detect edge cases within the training set and indicate when the decision of the trained agent is not fully reliable.

In this work, the estimated uncertainty is used to choose a safe fallback action if the uncertainty exceeds a threshold value. For the cases that are



**Figure 8:** Example of a situation outside of the training distribution, where there would be a collision if the confidence criterion is not used. The vehicle at the top is here approaching the crossing at 20 m/s.

considered here, this confidence criterion removes all collisions within the training distribution, and almost all collisions when the speed of the surrounding vehicles is increased to levels outside of the training distribution. However, to guarantee safety by using a learning-based method is challenging, and an underlying safety layer is often used **Underwood2016**. The presented method could decrease the number of activations of such a safety layer, but possibly more importantly, the uncertainty measure could also be used to guide the training process to focus on situations that the current agent needs to explore further. Moreover, if an agent is trained in simulation and then deployed in real traffic, the uncertainty estimation of the agent could detect situations that

should be added to the simulated world, in order to better match real-world driving.

The results show that the ensemble RPF method performs better and more stable than a standard DQN method within the training distribution. The main disadvantage is the increased computational complexity, since  $K$  neural networks need to be trained. This disadvantage is somewhat mitigated in practice, since the design of the algorithm allows an efficient parallelization. Furthermore, the tuning complexity of the ensemble RPF and DQN methods are similar. Hyperparameters for the number of ensemble members  $K$  and prior scale factor  $\beta$  are introduced, but the parameters that control the exploration of DQN are removed.

## 6 Conclusion

The results of this paper demonstrates the usefulness of using a Bayesian RL technique for tactical-decision making in an intersection scenario. The ensemble RPF method can be used to estimate the confidence of the recommended actions, and the results show that the trained agent indicates high uncertainty for situations that are outside of the training distribution. Importantly, the method also indicates high uncertainty for rare events within the training distribution. In this work, the confidence information was used to choose a safe action in situations with high uncertainty, which removed all collisions from within the training distribution, and most of the collisions in situations outside of the training distribution.

The uncertainty information could also be used to identify situations that are not known to the agent, and guide the training process accordingly. To investigate this further is a topic for future work. Another subject for future work involves how to set the parameter value  $c_v^{\text{safe}}$  in a more systematic way, and how to automatically update the value during training.

## References

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press, 2018.
- [2] M. J. Kochenderfer, *Decision Making Under Uncertainty: Theory and Application*. MIT Press, 2015, ISBN: 0262029251, 9780262029254.

# PAPER D

## Belief State Reinforcement Learning for Autonomous Vehicles in Intersections

Tommy Tram, Maxime Bouton, Jonas Sjöberg, and Mykel Kochenderfer

*Published in TBD,*  
pp. 3169–3174, Nov. 2018.  
©2018 IEEE DOI: TBD.

*The layout has been revised.*

### Abstract

This paper investigates different approaches to find a safe and efficient driving strategy through an intersection with other drivers. Because the intentions of the other drivers to yield, stop, or go are not observable, we use a particle filter to maintain a belief state. We study how a reinforcement learning agent can use these representations efficiently during training and evaluation. This paper shows that an agent trained without any consideration of the intentions of others is both slower at reaching the goal and results in more collisions. Four algorithms that use a belief state generated by a particle filter are compared. Two of the algorithms have access to the intention only during training while the others do not. The results show that explicitly trying to predict the intention gave the best performance in terms of safety and efficiency.

## 1 Introduction

During the last decade, major progress has been made towards deploying autonomous vehicles (AV) and improving traffic safety. Structured traffic scenarios can often be solved by rule-based or planning-based methods, where the surrounding vehicles are assumed to follow a simple driving model. However, according to the Insurance Institute for Highway Safety **IIHS2019**, in 2019, an estimated 115 741 people were injured by drivers running a red light and 928 of them were killed. Hence, to give AV the capability to avoid accidents when other traffic participants do not follow the rules, they need the capability to consider the intentions and predict future actions of other participants. For example, in an intersection or a roundabout scenario, a human driver would observe other vehicles and form a belief of the driver's intention to yield or not and would then act according to this belief. Although AVs have the potential to simultaneously observe more objects than a human, there is a vast number of possible traffic situations, making it difficult for an engineer to anticipate every situation and design a suitable rule-based strategy. Furthermore, both traffic rules and driving styles vary between different countries. Therefore, we focus on approaches that instead learn from experience how to behave in

different situations and adapt to different driving styles. A desirable property of such a learning-based agent is that it accounts for the intentions of the surrounding traffic participants.

Using reinforcement learning (RL) to create a general decision-making agent for autonomous driving has been suggested in many studies during the last few years. For example, **Isele2018** demonstrated how the deep Q-network (DQN) algorithm could be used to navigate through occluded intersections. **chen2019** compared the performance of DQN with policy gradient and an actor critic method for different roundabout scenarios. An overview of RL-based studies for autonomous driving is given by **Ye2021** and **kiran2021**. Most of these studies both train and evaluate the agents in simulated environments, whereas **Bansal2018** and **Pan2017** focus on transferring an agent that has been trained in simulation to drive in the real world. **Kendall2019** trained in the real world. Game-theoretic approaches to unsignalized intersection scenarios have been proposed by **chen2020** and **li2021**. A drawback is that these methods are limited by their models of other driver intentions and the prediction of their future actions. Therefore, we propose using RL to learn a policy based on experience.

The decision-making task of an autonomous vehicle can be formulated as a Partially Observable Markov Decision Process (POMDP), e.g., where the intentions of other traffic participants are included in the state space but are not directly observable. **Sunberg2017** models freeway lane changes as a POMDP with latent internal states and found a policy using Monte Carlo Tree Search. Our previous work **Tram2018** showed that it was possible for an agent to solve such a POMDP by finding gaps between cars in an intersection by using a long short-term memory (lstm), but the intentions were not explicitly estimated. In this paper, we hypothesize that explicitly estimating the intention probability distribution of the other drivers can improve the performance of the policy by reducing the number of collisions and reach the other side of the intersection faster. **Hubmann2017** proposed a similar POMDP framework for solving the same intersection problem using a particle filter, but they did not use RL. By using RL, the method is more scalable to different environments and has the potential to adapt to real world behaviors through experience.

This hypothesis, that explicitly estimating the intention can reduce the number of collisions, is verified by training and comparing two DQN agents: one with full observability (**DQN FO**) including the true intention as an input

feature to the network; and one without intention, that is referred to as **DQN without intention**. While we show in our simulations that observing the true intention leads to performance improvements, it is not a realizable approach in practice since the intention is not observable. We will therefore investigate and answer the following questions: *How can the intention of the other vehicles approaching an intersection be estimated? With a belief state, is it better to train on the full distribution of the belief or use an estimate of the intention? Can the approximation of the neural network compensate for the inaccuracy of the particle filter?*

To answer the first question, a particle filter is proposed in 3.2, which uses a sequence of past observations to generate a belief state. This belief state can then be used, instead of the true intent, to create a policy. In addition to the two baseline DQN agents, four approaches that combine a particle filter with the DQN algorithm are presented in 3.3. The **Q Particle Filter (QPF)** approach uses the belief state, i.e., the full set of particles to estimate the  $Q$ -values both during the training and evaluation processes. If the intention is accessible during training the second approach, **QMDP**, is trained under full observability and then tested with the belief state as input. The third approach, **Q Intention Distribution (QID)**, only uses a probability distribution of the non-observable states as an input both during the training and testing processes. Finally, **QMDP-Intention Estimate (QMDP-IE)** is trained under full observability and tested with an estimate of the intention from the particle filter. The detail of each approach are explained in more detail in 3.3 and their performance is evaluated in an unsignalized intersection scenario, explained in 4.

The results show that the intention state improves the policy when comparing the performance between a DQN with access to the true intention and a DQN without the intention state. The QMDP-IE algorithm found the safest policy with zero collisions in the evaluation set, while the QID found the fastest and most aggressive policy.

The main contributions of this paper are:

- Multiple DQN models that capture the belief state at training time either in the form of particle, point estimate of the intention, or the likelihood of the intention.
- A particle filter approach for estimating the intentions of surrounding traffic participants in an intersection.

- Empirical analysis of the performance of four approaches to explicitly consider the intentions of other traffic participants in an RL framework.

The structure of the paper is as follows. In 2, we review POMDPs and the intelligent driver model (IDM). In 3, we formulate the problem as a POMDP and introduce the particle filter used to generate the belief state along with the algorithms used to solve the POMDP. In 4, we describe the simulation setup, the neural network architecture, and training procedure. 5 compares the different algorithms and our conclusions are presented in 6.

## 2 Background

This section briefly describes the POMDP framework, the approximation method QMDP, deep Q-learning, and the IDM.

### 2.1 Partially observable Markov decision process

A POMDP is a mathematical framework used for modeling sequential decision making problems under uncertainty. It consists of a set of states  $\mathcal{S}$ , actions  $\mathcal{A}$ , observations  $\Omega$ , a transition model  $T$ , observation model  $O$ , reward function  $R$ , and a discount factor  $\gamma$ . Together they create the tuple  $(\mathcal{S}, \mathcal{A}, \Omega, T, O, R, \gamma)$  that formally defines a POMDP [1]. While in a state  $s \in \mathcal{S}$  and executing an action  $a \in \mathcal{A}$  the probability of transitioning to a future state  $s'$  is described by the transition model  $T(s' | s, a)$  and the reward  $r$  is given by the reward function  $R(s, a)$ . The agent only has access to partial information contained in its observation  $o$  distributed according to  $\Pr(o | s', a) = O(o, s', a)$ .

To accommodate for the missing information, the agent maintains a belief state. A belief state  $b$  is a probability distribution such that  $b(s) = \Pr(s | o_{1:t})$  is the probability of being in state  $s$  given observations  $o_{1:t}$ . At each time step  $t$ , the agent updates its belief using a Bayesian filtering approach given the previous belief and the current observation as follows:

$$b'(s') \propto O(o | s', a) \sum_{s \in S} T(s' | s, a) b(s). \quad (\text{D.1})$$

In a POMDP, a policy is a mapping from a belief state to an action. A value function  $Q^\pi(b, a)$  represents the expected discounted accumulated reward received by following policy  $\pi$  after taking action  $a$  from belief state  $b$ . The

goal is to find a policy, that maximizes the expected reward. Finding such a policy is generally intractable [1]. We can use an approximation referred to as QMDP:

$$Q(b, a) \approx \sum_s b(s)Q_{\text{MDP}}(s, a) \quad (\text{D.2})$$

where  $Q_{\text{MDP}}$  is the optimal value function of the Markov Decision Process (MDP) version of the problem that assumes that the agent has full observability.

While finding the optimal  $Q(b, a)$  is intractable, finding  $Q_{\text{MDP}}$  is usually easier. When the transition function can be written explicitly and the state space is finite, we can use dynamic programming to compute  $Q_{\text{MDP}}$ . Often, the transition function is only accessible in the form of a generative model (*e.g.* a simulator) and the state space is high dimensional and continuous. In such settings, we can use deep Q-learning to approximate  $Q_{\text{MDP}}$  **LITTMAN1995**.

In deep Q-learning, the value function is approximated by a neural network. The function approximating the optimal value function is given by minimizing the following loss function:

$$J(\theta) = E_{s'}[(r + \gamma \max_{a'} Q(s', a'; \theta) - Q(s, a; \theta))^2] \quad (\text{D.3})$$

where  $\theta$  represents the weights of the neural networks and  $(s, a, r, s')$  is obtained from one simulation step in the environment. The weights are updating through gradient updates. In addition, innovations such as the use of a replay buffer and a target network can greatly help the convergence **Mnih2015**.

## 2.2 Driver model

The general behaviour of cars are modeled using IDM **idm2000** and the acceleration is described by

$$\alpha = \alpha_{\max} \left( 1 - \left( \frac{v_\alpha}{v_0} \right)^\delta - \left( \frac{s^*(v_\alpha, \Delta v_\alpha)}{s_\alpha} \right)^2 \right)$$

with  $s^*(v_\alpha, \Delta v_\alpha) = s_0 + v_\alpha T_{\text{gap}} + \frac{v_\alpha \Delta v_\alpha}{2\sqrt{a_{\max} \alpha_b}}$

where  $v_0$  is the desired velocity,  $s_0$  is the minimum distance between cars,  $T_{\text{gap}}$  is the desired time gap with the vehicle in front,  $a_{\max}$  is the maximum vehicle acceleration, and  $\alpha_b$  and  $\delta$  are comfort braking deceleration and a

model parameter. The velocity difference  $\Delta v_\alpha$  between the two vehicles and the distance to the vehicle in front  $s_\alpha$  determine the acceleration for the next time step. The acceleration can be simplified into two terms, one for free road with no lead vehicle

$$a^{\text{free}} = \alpha_{\max} \left( 1 - \left( \frac{v_\alpha}{v_0} \right)^\delta \right) \quad (\text{D.4})$$

and an interaction term for when there is a vehicle in front

$$a^{\text{int}} = -a \left( \frac{(s_\alpha, \Delta v_\alpha)}{s_\alpha} \right)^2 \quad (\text{D.5})$$

$$= -a \left( \frac{s_0 + v_\alpha T_{\text{gap}}}{s_\alpha} + \frac{v_\alpha \Delta v_\alpha}{2\sqrt{a_{\max} \alpha_b s_\alpha}} \right)^2. \quad (\text{D.6})$$

These equations are used to both model behaviors of other drivers and generate the acceleration for the ego vehicle.

### 3 Proposed approach

In this section, we define the main problem that this paper addresses, as well as its components. Further, the particle filter used to estimate the belief state is introduced. Finally, we present the RL algorithms that are used to train the agents.

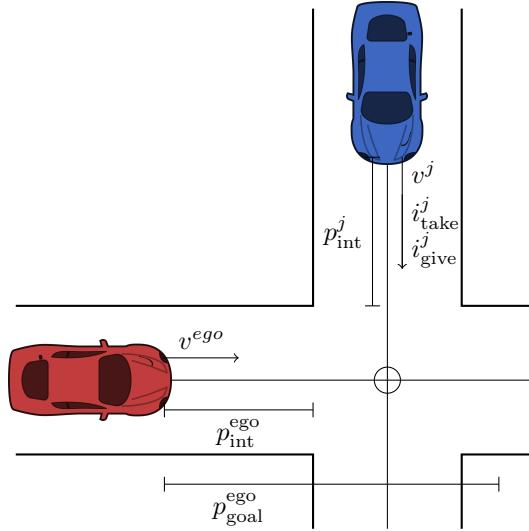
#### 3.1 Problem formulation

The problem setting is a four-way intersection shown in 1. The ego vehicle approaches an intersection with at least one other car on the perpendicular lane with about the same time to intersection assuming a constant velocity model. The goal for the ego vehicle is to reach the other side of the intersection as fast as possible without colliding with the other car. With only onboard sensors on the ego vehicle, it can observe the physical state of other vehicles but not their intention. Such an intersection can be described as a POMDP.

##### **State space, $\mathcal{S}$**

The state of the system,

$$s = (p_{\text{goal}}^{\text{ego}}, p_{\text{int}}^{\text{ego}}, v^{\text{ego}}, t_{\text{stop}}, \{p_{\text{int}}^j, v^j, i^j\}_{j=1}^N), \quad (\text{D.7})$$



**Figure 1:** State definitions for the intersection. The red vehicle is the ego vehicle.

consists of the ego vehicle state and the states of the surrounding vehicles. The ego vehicle state is described by the distance to the intersection  $p_{\text{int}}^{\text{ego}}$ , the distance to the goal  $p_{\text{goal}}^{\text{ego}}$ , and the velocity  $v^{\text{ego}}$ . Stop time  $t_{\text{stop}}$  is the time the ego vehicle is at stand still  $v^{\text{ego}} = 0$ . This state tracks the amount of time the ego vehicle has been standing still at the intersection and indicates the time to the terminal states *safe stop* or *deadlock*, which are explained later when we define the reward function. The states of the surrounding vehicles, indexed by  $j \in \{1, \dots, N\}$ , are described by the distance to the intersection  $p_{\text{int}}^j$ , the velocity  $v^j$ , and the intention  $i^j$ . The intentions are represented by a one-hot vector and can be either *give way*, which means that the vehicle will stop before the intersection, or *take way*, which means that the vehicle will drive through the intersection. These intentions control the behavior of the IDM model, described in 2.2.

### Observation space, $\Omega$

An observation  $o$  consists of the ego vehicle state and the physical state of the surrounding vehicles, while the intentions of the surrounding vehicles  $i^j$  are

not observed. An observation is described by

$$o = (p_{\text{goal}}^{\text{ego}}, p_{\text{int}}^{\text{ego}}, v^{\text{ego}}, t_{\text{stop}}, \{\hat{p}_{\text{int}}^j, \hat{v}^j\}_{j=1}^N). \quad (\text{D.8})$$

### Observation model

The agent observes the ego vehicle states precisely and takes noisy measurements of the positions and speeds of the surrounding vehicles, given by

$$\hat{p}_{\text{int}}^j = p_{\text{int}}^j + \epsilon_p, \quad (\text{D.9})$$

$$\hat{v}^j = v^j + \epsilon_v. \quad (\text{D.10})$$

Here,  $\epsilon_p \sim \mathcal{N}(0, \sigma_p)$  and  $\epsilon_v \sim \mathcal{N}(0, \sigma_v)$ .

### Action space, $\mathcal{A}$

The action space  $\mathcal{A}$  consists of two high level actions take way and give way, also known as options **SUTTON1999**. The motion of the ego vehicle is controlled by changing the acceleration using the IDM, introduced in 2.2. Depending on the action, the IDM parameters for target vehicle would be different. The take way action uses the free term from D.4, not following any car and just drives though the intersection. The give way action on the other hand uses the interaction term from D.6, with the variable distance  $s$  set to the distance to intersection and the velocity difference  $\Delta v_\alpha$  to 0. The action then controls the acceleration of the ego vehicle.

### Reward function, $R$

The reward function is designed to encourage safety and efficiency. The main goal of the agent is to reach the goal on the other side of the intersection without colliding with other traffic participants. There are four terminal states: goal, safe stop, collision, and deadlock. While goal and collision is self explanatory, safe stop and deadlock are reached when the agent chooses to stop at the intersection for  $t_{\text{stop}}$  consecutive seconds. To distinguish whether a stop was efficient, there are two different outcomes. If other vehicles are at stand still and waiting for the ego vehicle at the terminal state, the deadlock reward is received while the safe stop is received if all other vehicles are in

motion while the ego vehicle stopped. The reward function

$$r_t = \begin{cases} 8 & \text{reaching the goal,} \\ 0.4 & \text{safe stop,} \\ -10 & \text{collision,} \\ -0.6 & \text{deadlock,} \\ -0.01 & \text{otherwise} \end{cases}$$

is designed to have a large relative difference between reaching the goal and colliding, with a small incentive for reaching a terminal state faster. The reward for safe stop is positive to account for scenarios where there are many cars and none of them has the intention to yield. Combined with the negative reward for deadlock, the agent is incentivised to not stop unless another vehicle is yielding.

### Transition model, $T$

The state transition probabilities are implicitly defined by the generative simulation model and are not known to the agent. The simulation model is defined in 4.1.

## 3.2 Belief state estimation using a particle filter

To estimate the state of the environment D.7 from noisy observations, we use a particle filter algorithm that takes advantage of the IDM and assumptions of observation independence. The role of the particle filter in our method is to estimate the position and velocity of the observed vehicles along with their intention (give way or take way) from noisy measurements of position and velocity. This work does not focus on optimize the design of a particle filter but rather to demonstrate how it can be used jointly with a RL agent. In particular, we investigate whether it provides any benefits to perform state estimation at training time and whether the decision making agent can be made more efficient by reasoning about a distribution over intentions instead of just a state.

Previous work has shown that Bayesian filters can be used to infer driver intentions in merging scenarios **bouton2019**. One challenge in estimating driver intentions is that one might need to consider the joint distribution over

the intentions of all drivers interacting in the scenario. In this intersection scenario, a joint estimate of intentions for the four closest drivers is created. Each driver is modeled by a position  $p_{\text{int}}^j$ , a velocity  $v^j$ , and a binary intention  $i^j$  (give way or take way). Let  $s^j$  be the three dimensional state of a car. Considering four vehicles at the intersection, the state estimation algorithm must estimate a 12 dimensional distribution. A particle filter is used to estimate the state because the motion of a vehicle can easily be simulated given its intention.

We make independence assumptions about the other cars that allow us to factor the distribution joint transition model as:

$$\Pr(s_{t+1}^{1:4} | s_t^{1:4}) = \Pr(s_{t+1}^1 | s_t^1) \prod_{i=2}^4 \Pr(s_{t+1}^i | s_t^i, s_t^{i-1}). \quad (\text{D.11})$$

Without loss of generality, it is assumed that the order of vehicles is 1, 2, 3, and 4. Hence, the position and speed of vehicle 2 only depends on the front vehicle 1 according to the IDM. In addition, the observations of each vehicle are assumed to be independent from each other. Let  $o_t^i$  be the observation of vehicle  $i$  at time  $t$ . The joint observation distribution is:

$$\Pr(o_t^{1:4} | s_t^{1:4}) = \prod_{i=1}^4 \Pr(o_t^i | s_t^i). \quad (\text{D.12})$$

Given these two assumptions about the problem structure, we can define the full particle filter procedure. A set of  $M$  ordered particles for each observed vehicle are maintained, where a particle representing the joint state can be expressed as follows:  $s^{[m]} = (s^{1[m]}, \dots, s^{4[m]})$ . The set of joint particles can be represented by maintaining ordered sets of individual particles for each vehicle. The first particle associated with vehicle 1 will always represent the leading car to the first particle associated with vehicle 2. At each time step, the standard particle filter operations of prediction and measurement updates are performed with a few improvements.

- For all particle indices  $m$ ,  $s'^{[m]} \sim \text{simulate}(s^{[m]})$ . Predict the particle one step forward in time using the IDM. The prediction is performed sequentially by updating the leading vehicle first, then its following car, and so on, according to the process described by D.11.
- Compute observation weights for each particle. A Gaussian sensor model was used for the position and velocity. The weight  $w^{[m]}$  of a particle

is given by multiplying each of the four individual vehicle weights as described in D.12. The individual weights are assumed to follow a Gaussian observation model:

$$w^{j,[m]} \propto \mathcal{N}([p_{\text{int}}^j, v^j]^T; [\hat{p}_{\text{int}}^j, \hat{v}^j]^T, \text{diag}[\sigma_p^2, \sigma_v^2]). \quad (\text{D.13})$$

- A new set of particles is resampled from the set  $\{s'^{[m]}\}_{m=1}^M$  weighted by  $w^{[m]}$ , if the effective number of particles fall below a threshold.
- Noise is added to the resampled particles. This is done in two ways. By adding some noise to the acceleration in the prediction model and by letting the intention of each particle have a slight probability  $p_i$  to change intention. This noise models the fact that the transition model is not fully known.

At the end of these steps, the set of particles are updated to estimate the current state based on the current observation. The belief of our POMDP agent is represented by the set of resampled particles.

In order to make the particle filter efficient and avoid known issues such as particle depletion, two improvements were incorporated. When a vehicle is observed for the first time, a set of  $M$  particles are sampled, where each individual particle is associated with a joint particle. To respect the IDM, the position of a new vehicle is sampled uniformly around the first observed position and the position of its leading vehicle, while the velocity is sampled uniformly in the range described in 1. Intentions are initialized as 50% give way and 50% take way. Finally, the state of the individual particles for this new vehicle is appended to the joint particles of the other vehicles that have just been resampled.

The particle filter implementation is tailored to the problem of intersection navigation. We take advantage of observation independence and the structure of the car following model to efficiently update the joint particles. The particles are used to represent the belief state  $b$ , and the algorithms used train the RL agents are described in the next section.

### 3.3 Belief state reinforcement learning

This section describes the proposed approach to train an RL agent in the belief space. In classical deep RL methods for POMDPs, the learned value

	State $s$	Particles $b$	State Approximation $D(b)$
Training with intention	DQN FO	QMDP	QMDP-IE
Training without intention	DQN without intent	QPF	QID

**Figure 2:** Table showing the difference between proposed algorithms. The first row are algorithms that have access to the true intention during training, while the algorithms on the second row do not. The columns show the input to the neural network, where the first column contains the baseline algorithms that use the noisy observation from the sensors either with or without the intention  $i^j$ . The second column contains algorithms using all  $M$  particles as input and the final column uses an estimate of the intention.

function processes a sequence of observations using recurrent neural networks **HausknechtS15drqn** or by concatenating a finite history of observations and feeding it to a simple feed forward neural network **Mnih2015**. In this approach, a state estimation algorithm is used instead to compute a belief state based on the sequence of past observation and action. The resulting distribution is then fed to a deep RL agent. Instead of processing observations, our agent processes belief states. By learning in the belief space, our agent is more robust to state uncertainty and, contrary to the QMDP approach, it can compensate for some limitations of the state estimation module. This ability is key to learning intention aware policies for autonomously navigating intersections.

All six algorithms in this paper follow the same deep Q-learning structure, shown in 2, with the added subtlety that the architecture of the value function, in step 10 from 2, varies in order to process belief states. The loss function in D.3 takes the form of a mean-squared error with a target corresponding to  $r + \gamma \max_{a'} Q(s', a'; \theta)$  and a prediction corresponding to  $Q(s, a; \theta)$ . Depending on the algorithm presented, the prediction will be performed in a different way.

To confirm the hypothesis that the intention state is necessary, two baseline DQN algorithms are used.

**DQN FO:** This algorithm is trained with full observability and has access to the true intention at training. As mentioned earlier, intention can not be measured directly, and therefore this algorithm is just used to show the

best policy if the intention approximation is accurate. It also assumes perfect observability of the other state variables.

**DQN without intention:** Compared to the previous algorithm, this has access to the same states except the intention. By omitting the intention from the input state that is fed to the neural network, this baseline algorithm can be used to show the limit of a policy that does not have access to the intention.

**QMDP:** The first approach to use the belief state is the QMDP approach described in D.2. The belief state is used at test time but the agent is trained in an environment with full observability in order to estimate  $Q_{\text{MDP}}$ . The prediction used in the loss function is  $Q(s, a)$  and is learned from  $(s, a, r, s')$  tuples. The resulting policy can be sensitive to the performance of the state estimation algorithm used to generate the belief state  $b$ . To address this issue, we use the true state at training time so that the training procedure matches what the agent will experience at test time.

**QPF (Particle Filter):** With this algorithm, the agent is trained in an environment with partial observability and uses our proposed particle filter algorithm. The input to the agent is a set of  $M$  particles. The loss function is minimized from  $(b, a, r, b')$  experience samples and the Q-learning prediction is given by:

$$Q(b, a; \theta) = \frac{1}{M} \sum_{m=1}^M Q(s^{[m]}, a; \theta). \quad (\text{D.14})$$

This operation is differentiable and allows training  $Q(b, a)$  using deep Q-learning. One challenge with this approach is that the number of particles used to represent the belief can be very large in practice and cause the training to be unstable and very slow. It is challenging to update the individual Q-values by back propagating the loss function computed based on an aggregation of the values of many particles.

**QID (Intention Distribution):** The agent is also trained in an environment with partial observability and the loss is minimized using  $(b, a, r, b')$  samples. To simplify the training procedure compared to QPF, the probability distribution of the intention is used to compute the prediction:

$$Q(b, a; \theta) = Q(\mathbf{D}(b), a; \theta) \quad (\text{D.15})$$

where

$$\mathbf{D}(b) = (p_{\text{goal}}^{\text{ego}}, p_{\text{int}}^{\text{ego}}, v^{\text{ego}}, t_{\text{stop}}, \{\hat{p}_{\text{int}}^j, \hat{v}^j, \tilde{v}^j\}_{j=1}^N), \quad (\text{D.16})$$

and the probability distribution of the intention

$$\tilde{i}^j = \sum_{m=1}^M w^{j,[m]} [\tilde{i}_{\text{take}}, \tilde{i}_{\text{give}}]^{j,[m]} \quad (\text{D.17})$$

is given by the particle filter, where  $w^{j,[m]}$  is the weight of the particle  $m$  and  $[\tilde{i}_{\text{take}}, \tilde{i}_{\text{give}}]$  is the one-hot vector specifying intention. This operation does not affect the differentiability of the Q function that is trained using the same procedure as QPF. The advantage of this approach is that, at training time, the agent can experience changes in the estimated intention that would match what the state estimator would produce at test time. As a consequence, the agent becomes more robust to the error in intention estimation.

**QMDP-IE (Intention Estimate):** The last algorithm consists of training an agent on a fully observable MDP just like the QMDP approach. At test time, instead of averaging the Q-values of all the particle, a threshold  $i_{\text{threshold}}$  is set on the intention distribution and then use it as the true intention state. Similar to D.16, the particle filter is only used to generate the intention and is represented as a one-hot vector

$$\hat{i}^j = \begin{cases} [0 \ 1] & \text{if } \tilde{i}_{\text{give}}^j > i_{\text{threshold}} \\ [1 \ 0] & \text{otherwise,} \end{cases} \quad (\text{D.18})$$

where a high value on the threshold  $i_{\text{threshold}}$  makes the agent more conservative.

## 4 Experiments

We present in this section the experiment setup in three steps. First, we describe how the simulator generates the different traffic scenarios. Second, the network architecture is defined. Third, we describe the details of the training procedure.

### 4.1 Simulator setup

At the start of each episode, up to  $N$  vehicles are spawned with initial positions  $p_0^j$  distributed along the intersecting lane, a starting velocity  $v^0$ , and a desired velocity  $v_{\text{desired}}^j$ . Each vehicle is spawned with a deterministic policy that represents their intention, which can be either take way or give way as described

---

**Algorithm 2** Training process

---

```

1: Initialize  $\theta$  randomly
2:  $\mathcal{D} \leftarrow \emptyset$ 
3: for nr episodes do
4:    $o \leftarrow$  initiate environment
5:    $b \leftarrow \text{INITIALIZEBELIEF}(o)$ 
6:   while episode not finished do
7:     if  $e \sim \mathcal{U}(0, 1) < \epsilon$  then
8:        $a \leftarrow$  random action
9:     else
10:       $Q(b, a) \leftarrow \text{GETACTIONVALUES}(b, \theta)$ 
11:       $a \leftarrow \text{argmax}_a Q(b, a)$ 
12:       $o', r \leftarrow \text{STEPENVIRONMENT}(a)$ 
13:       $b' \leftarrow \text{UPDATEBELIEF}(b, o', a)$ 
14:       $\mathcal{D} \leftarrow \mathcal{D} \cup \{(b, a, r, b')\}$ 
15:       $M \leftarrow$  sample from  $\mathcal{D}$ 
16:      update  $\theta$  with SGD and loss  $J(\theta)$ 

```

---

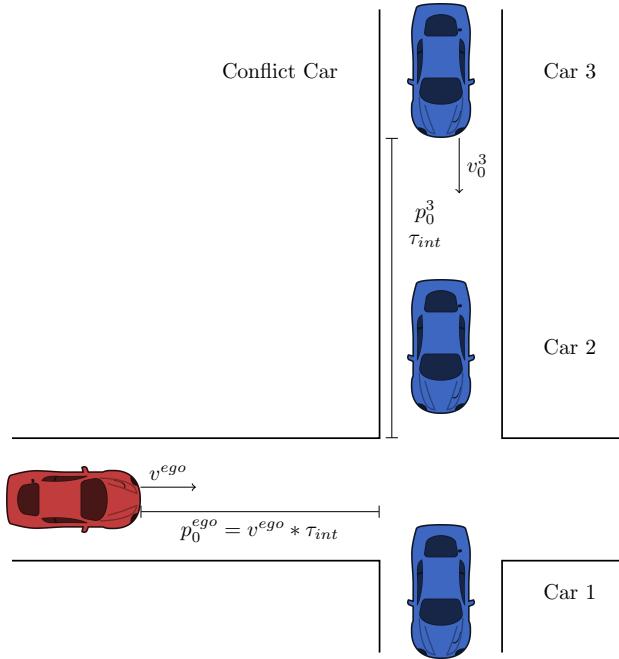
in 3.1. The ego vehicle is spawned last with an initial speed  $v^{\text{ego}}$  and desired speed  $v_{\text{desired}}^{\text{ego}}$ . The starting position of ego is then set based on the time to intersection  $\tau_{int}$  of one of the other vehicles. This other car is referred to as conflict car. The initial position of the ego becomes:

$$p_0^{\text{ego}} = v^{\text{ego}} \frac{p_0^c}{v_0^c}, \quad (\text{D.19})$$

where  $c$  is the index of the conflict car. This initialization increases the probability that the ego will conflict with at least one other car as shown in 3.

The decision time between decisions steps are  $dt_{\text{decision}}$ . For every decision step, the simulator updates the objects states four times with a simulation time  $dt_{\text{sampling}}$ , checking for terminal states at each update. The simulator keeps track of the true state of all objects. An observation model is used to add noise to each observation, following ???. The update function updates the true state  $s_t$ . Every time a vehicle in the perpendicular lane crosses the intersection, they are re-spawned at the start of the lane at a random time with new initial states and intention.

The possible terminal states are the following: (i) Goal, ego vehicle reaching



**Figure 3:** Initial state with three other cars and ego car initial position is set to have the same time to intersection as the conflict car, in this case car 3.

the goal. (ii) Collision, ego collides with one of the other vehicles. (iii) Safe stop, if the ego vehicle has been standing still for more than  $T_{stop}$  seconds and no other car is standing still. (iv) Deadlock, if another vehicle standing still at the intersection yielding for the ego vehicle and while the ego vehicle has also been standing still for more than  $T_{stop}$  seconds. (v) Timeout, if the total simulation time exceeds  $T_{lim}$ . Each of these terminal states, except for the simulation timeout, has a corresponding reward, described in 3.1. The terminal state referred to as collision in this paper may sound drastic, but could also be interpreted as intervention by a collision avoidance system.

## 4.2 Neural network architecture

The neural network architecture that is used in this study is shown in 4. A previous study **Hoel2018** introduced applying a convolutional operator to

**Table 1:** Hyperparameters of Simulator

sampling time [s],	$dt_{\text{sampling}}$	0.5
decision time [s],	$dt_{\text{decision}}$	2
initial speed [ $\text{m s}^{-1}$ ],	$v_0^o$	2 – 7
initial acceleration [ $\text{m s}^{-2}$ ],	$a_0^o$	0
desired speed [ $\text{m s}^{-1}$ ],	$v_{\text{desired}}^o$	2 – 7
time gap [s],	$T_{\text{gap}}$	1.5
Stop time limit [s],	$T_{\text{stop}}$	10
Timeout time [s],	$T_{\text{lim}}$	120
ego spawn speed [ $\text{m s}^{-1}$ ],	$v^{\text{ego}}$	5
ego desired speed [ $\text{m s}^{-1}$ ],	$v_{\text{desired}}^{\text{ego}}$	5
noise position [m],	$\sigma_p$	2
noise velocity [ $\text{m s}^{-1}$ ],	$\sigma_v$	1
max observed vehicles,	$N$	4
IDM max acceleration [ $\text{m s}^{-2}$ ],	$\alpha^{\text{max}}$	0.73
IDM deceleration [ $\text{m s}^{-2}$ ],	$\alpha_b$	0.5 – 4.0
IDM acceleration exponent,	$\delta$	4
IDM minimum distance [m],	$s_0$	2
number of particles	$n_b$	100
acceleration noise [ $\text{m s}^{-2}$ ]	$\sigma_a$	0.1
intention switch probability	$p_i$	5
intention probability threshold	$i_{\text{threshold}}$	0.8
Batch size	B	128
Learning rate	lr	0.0001
Discount factor	$\gamma$	0.95
Replay memory size	$M_{\text{replay}}$	20,000
Target network update frequency	$N_{\text{update}}$	1,000

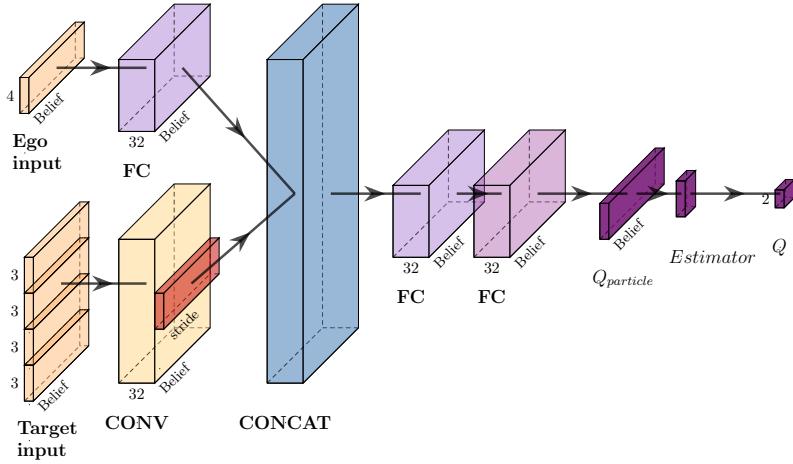
the states that describe the surrounding vehicles. With this structure, the states that describe each surrounding vehicle are passed through the same weights. The size and stride of the convolutional layer is set to  $(4, 1)$ , where four equals the number of states that describe each vehicle and one means that the different particles are handled individually. The convolutional layer uses 32 filters with a tanh activation function. To speed up learning, the index  $j$  are ordered by distance  $p_{\text{int}}^j$  starting from the vehicle closest to the intersection and a default state used for cars that do not exist.

The ego vehicle states are passed through a fully connected layer of size 32 before being concatenated with the output of the convolutional layer. The output of the concatenated layer is then passed through two fully connected layers of size 32 and finally a fully connected layer of size two gives the  $Q$ -values. These values are then merged according to one of the algorithms in 3.3 to form a combined estimate of the  $Q$ -values. The belief dimension of the neural network represents the number of particles  $n_b$  that are passed through the network. This number varies depending on which algorithm that is used as discussed in 3.3.

### 4.3 Training procedure

Each network was trained for 200 000 episodes. The loss function of Double DQN is applied, which modifies the maximization operation of D.3 to  $\gamma Q(s', \text{argmax}_{a'} Q(s', a'; \theta_i); \theta_i^-)$  **Hasselt2016**. Stochastic gradient decent is used to update the weights **Kingma2014**. The hyper parameters of the training process are shown in 1, and the values were selected by an informal search due to the computational complexity.

If the current policy of the agent do not reach a terminal state before the simulation timeout time  $T_{\text{lim}}$ , an episode in the real world could continue forever. Therefore, when the simulation reaches the timeout state, the last experience of such an episode is not added to the replay memory. This trick prevents the agent from learning that an episode can end due to a timeout, because the terminating state is not part of the MDP **Hoel2018**.

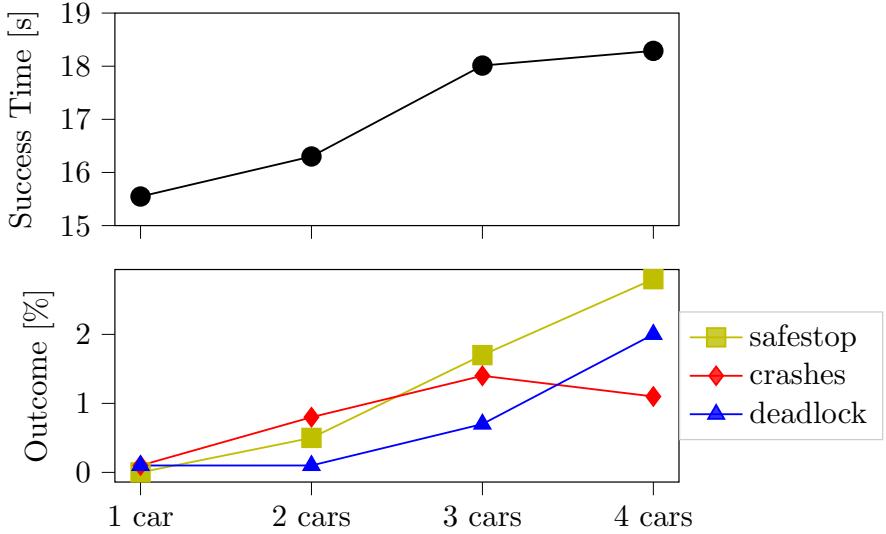


**Figure 4:** Network architecture, where ego input represents the input features for the ego vehicle,  $p_{goal}^{ego}, p_{int}^{ego}, v^{ego}, t_{stop}$ , and Target input  $\{p_{int}^j, v^j, i^j\}_{j=1}^N$  for the target vehicle  $j$ . The layers consists of one convolution (Conv) and a total of three fully connected (FC), where the depth of the layers represents the belief size.

## 5 Results / Results and discussion

This section presents the evaluation metrics and results from the experiments, presented in 2. Although the agents are train and evaluated on scenarios with different number of other cars. Because the scenario with four other cars is the hardest to cross without colliding, it is used to compare the performance of the different agents.

Each agent is evaluated on 1000 episodes. The random seeds are based on the episode number, which generates the same values for the random parameters, defined in 1, for each test scenarios when evaluating different agents. The metrics of interests are the average time it takes for the agent to either reach the success state and how often each terminal state is reached. A success state refers to either reaching the goal or a safe stop, from 3.1, where both are considered good outcomes. For the purpose of this paper, it is more important to first find an agent with no collisions and no deadlocks, and then optimize for the time it takes to reach the goal.



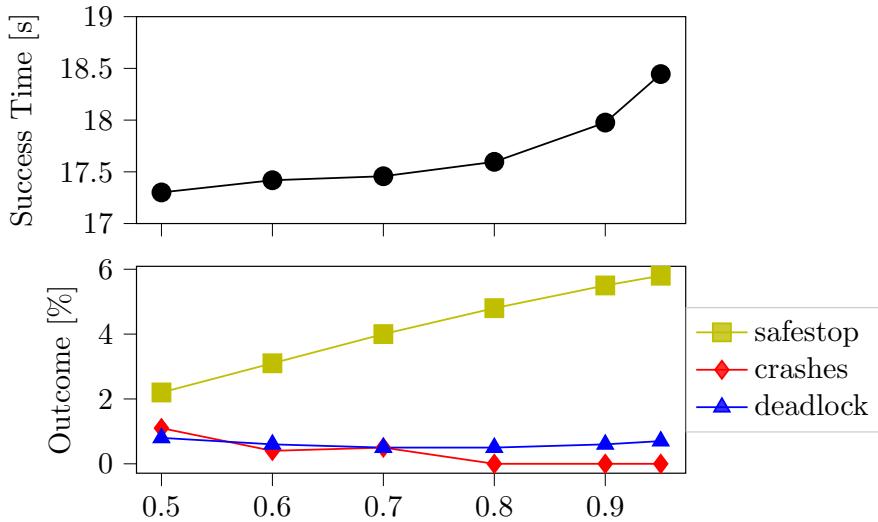
**Figure 5:** Performance results for DQN without intention algorithm, showing the increasing difficulty for scenarios with different number of cars in the environment.

## 5.1 Baseline DQN results

The agent trained with full observability, DQN FO, shows how well an agent could perform if it had access to the true intention of other drivers. From 2, the DQN FO agent has 94.2 % goals reached, 5.6 % safe stops, 0 % collisions, 0.2 % deadlocks, and an average success time of 16.93 s. While DQN without intention has 94.1 % goal reached, 2.8 % safe stops, 1.1 % collisions, 2.0 % deadlocks, and reached the goal slower with an average time of 18.29 s. The difference in collision rates confirms our hypothesis that a agent trained with an intention state is better than one without. DQN with intention was both faster to reach the goal and could do it with zero collisions in the evaluation space. All agents are trained with a variable number of other vehicles. 5 shows the performance difference with increasing number of vehicles. In scenarios with only one other vehicle, the agent trained using DQN without intention found a policy that could avoid collisions. However, increasing the number of other vehicles makes the problem harder to solve, which results in an increasing

**Table 2:** Result summary for scenario  $N = 4$  cars

Experiments	Goal reached	Safe stop	Collision	Deadlock	Success time	Training time
DQN FO	94.20	5.60	<b>0.00</b>	0.20	16.93	1d7h
DQN w/o intent	94.10	2.80	1.10	2.00	18.29	1d7h
QMDP	93.70	5.60	<b>0.10</b>	0.60	19.95	1d7h
QMDP-IE	94.70	4.80	<b>0.00</b>	0.50	17.60	1d7h
QPF	84.70	8.90	3.60	2.80	20.16	2d17h
QID	<b>97.00</b>	0.90	2.00	0.10	<b>16.90</b>	2d10h



**Figure 6:** QMDP-IE performance (aggressiveness) for different  $i_{\text{threshold}}$  values. Increasing  $i_{\text{threshold}}$  lowers the collision rate while at the same time increases the safe Stop rate and the time it takes to reach a success state.

number of collisions and a longer success time.

## 5.2 QMDP and QMDP-IE results

QMDP and QMDP-IE are both approximation algorithms that use the network trained on DQN FO. Both have similar outcome performance, QMDP reached

the goal 93.7 %, safe stopped 5.6 %, collided only 0.1 %, and got stuck in deadlock 0.6 % with an average success time of 19.95 s. for QMDP-IE, the performance for different threshold values of  $i_{\text{threshold}}$  are shown in 6. The lowest threshold  $i_{\text{threshold}} = 0.5$  is equivalent to just taking the most likely estimation of the intention and it has the most aggressive behaviour with the lowest percentage of safe stops, highest collisions and the fastest time to reach the goal. When the threshold increases, so does the safe stop rate and the time it takes to reach the goal, while the collision rate decreases and deadlock rates stay about the same. An agent with  $i_{\text{threshold}} = 0.8$  was the lowest threshold that achieve zero collisions during testing and is therefore chosen as the threshold value when compared to the other agents in 2. With  $i_{\text{threshold}} = 0.8$ , QMDP-IE has 43.7 % goal reached, 4.8 % safe stops, 0.0 % collisions, 0.5 % deadlocks, and with an average success time of 17.60 s Compared to DQN without intention, that has a success time of 18.29 s, QMDP is more conservative with a slower success time of 19.95 s, while the QMDP-IE is faster with 16.90 s.

The proposed particle filter performs well at estimating the probability distribution of the intention. However, it relies on an accurate prediction model. It is possible to improve the estimate by increasing the number of particles at the expense of additional computation.

### 5.3 QPF and QID results

The algorithms that do not have access to the intention during training, QPF and QID, reach the goal much less frequently. QPF has the lowest goal reached rate at 84.7 %, 8.9 % safe stop, the highest collision rate at 3.6 %, 2.8 % deadlocks, and the highest average success time of 20.16 s. QID learned the most aggressive policy, with 97.0 % goal reached, 0.7 % safe stop, 2.0 % collision, 0.1 % deadlock and the fastest average success time of 16.90 s, which is equal to the DQN baseline algorithm. As mentioned in the beginning of the section, the a slightly higher collision rate of 2.0 % classifies this performance as poor compared to QMDP and QMDP-IE, but it is much better than QPF on all evaluation metrics. The higher collision rate could be because of the larger state space, which could create some instability during training. The particle filter was also not fully optimized, because we wanted to investigate if the approximation of the DQN could compensate for the noisy states of the particles. Finally, it should be noted that it required almost two and a half

days to train QPF and QID, which is almost twice as much time required to train the fully observable network.

## 5.4 Discussion

*Is it better to train on the full distribution of the belief or use an estimation of the intention?* When comparing QPF and QID to their baselines, it is observed that both algorithms that were trained on the ground truth outperform the algorithms that were trained using the belief state or even just an estimate of the belief distribution. Both QMDP and QMDP-IE had zero collisions in the evaluation space just like DQN FO, while QMDP-IE had a slightly faster policy.

We used a particle filter to create the belief state and filter the noisy observations, hoping the flexibility of the neural network would be able to compensate for some of the inaccuracy of the belief in the same way it can approximate the noise. The higher collision rate for the QPF compared to QMDP shows that training on the full set of beliefs can make finding a good policy difficult.

6 shows that the aggressiveness of the policy from QMDP-IE is correlated with  $i_{\text{threshold}}$  and by choosing a relatively high value at 0.8 the policy could achieve zero collisions in the evaluation space and in this case could to some extent compensate for the loosely optimized particle filter. The ability to adjust the aggressiveness of the policy by changing  $i_{\text{threshold}}$  is also a strength of QMDP-IE compared to our previous black box methods that use an lstm to estimate the latent state **Tram2018**.

## 6 Conclusion

In this paper, we use a particle filter implementation to maintain a belief state to train an RL agent using a DQN. A set of four different algorithms are evaluated with the aim to find the safest and fastest policy that can drive through an intersection crossing traffic participants that has a latent intentions state. Two algorithms, QMDP and QMDP-IE, use a network trained with access to the true intention of the other vehicles but was evaluated with an estimate of their intention. Both algorithms could find a policy with zero collisions in the evaluation space, while QMDP-IE lead to the most efficient

policy. The other two algorithms, QPF and QID, were trained without access to the true intention state and could not achieve a policy without collisions. The QID agent had the fastest average time to reach the goal, but resulted in collision more often than the baseline DQN agent. As for the QPF agent, training on the entire belief state as input yielded the worst performance of all evaluated algorithms. It seems reasonable that the results can be improved by using a network structure that can better learn the latent state. One such candidate is Deep Variational Reinforcement Learning proposed by **Igl2018**.

## References

- [1] M. J. Kochenderfer, *Decision Making Under Uncertainty: Theory and Application*. MIT Press, 2015, ISBN: 0262029251, 9780262029254.

# PAPER E

## Maximum Likelihood Estimation for Transfer Reinforcement Learning Problems in Autonomous Driving

Hannes Eriksson, Tommy Tram, Debabrota Basu, Jonas Sjöberg, and  
Christos Dimitrakakis

*Artificial Intelligence and Statistics 2023 (AISTATS),*  
pp. 3169–3174, Nov. 2023.  
©2018 IEEE DOI: 10.1109/ITSC.2018.8569316.

*The layout has been revised.*

## Abstract

For decision-problems with insufficient data, it is imperative to take into account not only what you know but also what you do not know. In this work, ways of transferring knowledge from known, existing tasks to a new setting is studied. In particular, for tasks such as autonomous driving, the optimal controller is conditional on things such as, the physical properties of the vehicle, the local and regional traffic rules and regulations and also on the specific scenario trying to be solved. Having separate controllers for every combination of these conditions is intractable. By assuming problems with similar structure, we are able to leverage knowledge attained from similar tasks to guide learning for new tasks. We introduce a maximum likelihood estimation procedure for solving Transfer Reinforcement Learning (TRL) of different types. This procedure is then evaluated over a set of autonomous driving settings, each of which constitutes an interesting scenario for autonomous driving agents to make use of external information. We prove asymptotic regret bounds for proposed method for general structured probability matrices in a specific setting of interest.

## 1 Introduction

Sequential decision-making for autonomous driving agents has been studied under the Reinforcement Learning (RL) framework, which involves actively interacting with an environment for it to be able to improve on its performance. Disregarding the apparent difficulty in learning decision-making agents for large-scale environments, this also gives rise to another set of problems. For one, how do you guarantee safe performance during the learning process and secondly, how can you leverage external knowledge about decision-making problem. For instance, there might be information available informing us on how to act in a particular autonomous driving task, such as in an overtake scenario, that can be used to make decisions on how to act in a related task, such as a cut-in scenario, where another vehicle is changing lane onto the ego vehicle's lane.

Motivations for transfer learning in the context of autonomous driving are multifold. For instance, learning a controller for a particular vehicle involves figuring out what decision-making behavior is optimal given the properties of the vehicle. If the vehicle is less insert, you could expect more hasty decision-making is possible, whereas for a large vehicle like a semi-trailer truck, the decision to accelerate or decelerate needs to be made way in advance of the traffic situation. Furthermore, one could expect the learned controller to be locally optimal in the region it has collected data in. However, learning a global controller that is able to handle all sorts of local traffic rules and regulations, left- and right-hand traffic and the behavior of other road users is not easy. Thus, the framework proposed in this work can alleviate concerns regarding learning optimal controllers for each separate vehicle type, scenario and geographic region by leveraging knowledge transfer and sharing data across problems.

Similar topics have been studied in the context of *transfer machine learning*, e.g. in imitation learning. Recently, a lot of interest has been shown for *transfer reinforcement learning* in e.g. **zhu2020transfer**, where RL agents are able to learn how to act in environments by transferring knowledge from experts to facilitate the learning process. Typically, this has been done by two particular methods, firstly, *learning by demonstrations*, where agents are given data trajectories of demonstrations or expert policies from which they are able to bootstrap knowledge about the task. Secondly, the agent may have access to model dynamics of similar tasks, from which it could be able to extract knowledge from to guide learning in the actual environment.

These two approaches are denoted *policy transfer* and *model transfer*, respectively, by **zhu2020transfer**.

## 1.1 Related Work

[HE: TODO: Write about the uncommented works in this related work section.]

## 1.2 Contribution

In this work, we develop methods for leveraging knowledge transfer in the context of reinforcement learning. By identifying a global model that takes into account existing information from other tasks, we can design a decision-

making agent capable to handling multiple unknown scenarios at once, by using the current observed information of the task together with knowledge from similar tasks. The maximum likelihood global model parameters are identified and the optimal controller for the global model can attained through typical reinforcement learning means. The optimal global model under this formulation can change with time as more data is observed from the target task. The framework is evaluated in both discrete MDP scenarios and a larger, more complicated autonomous driving experiment where multiple transfer learning situations for autonomous driving are considered.

[HE: Some illustrative TRL problem for AD, for example AD in US, Sweden and India showing different traffic scenarios. We try to build a global model.] [TT: TRL examples suggestions: trained for low speed traffic jams, transfer to a high speed highway driving. traffic laws in different countries. some part of the us allow turning right in an intersection when the traffic light is red. Sweden usually have slower speed limitations on the road, Asia has higher traffic density.]

In Section 2 the necessary background details are described, such as the likelihoods used, the feature representations and the model definition.

## 2 Background

In this section, we will be introducing important concepts on which this work is based upon.

### 2.1 Markov Decision Process

In this work, we are studying sequential decision-making problems that are able to be modelled by Markov Decision Processes (MDPs) **sutton2018reinforcement**. An MDP  $\mu$  consists of the following constructs,  $\mathcal{S} \in \mathbb{R}^d$ , which is the state-space of the problem of dimension  $d$ .  $\mathcal{A}(s)$  is the set of possible actions to take in state  $s$ .  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is a reward function which determines the quality of taking a particular action  $a$  in state  $s$ .  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is a transition kernel that describes how states evolve in the MDP. Finally,  $\gamma$ , is the discount factor chosen for the problem.

The objective of an agent acting in a MDP is to arrive at a policy  $\pi : \mathcal{S} \rightarrow$

**Table 1:** Examples of related TRL works

Classes of Transfer Reinforcement Learning problems	
Parameter space	Related work
TYPE I: Finite parameter set $\Theta$	ref1, ref2, ref3, ref4
TYPE II: Infinite, convex set $\Theta$	ref1, ref2, ref3, ref4
TYPE III: Infinite, bounded set $\Theta$	ref1, ref2, ref3, ref4
TYPE IV: Arbitrary set $\Theta$	ref1, ref2, ref3, ref4

$\Delta(\mathcal{A})$  that maximises the expected sum of future discounted rewards up until the horizon  $T$ .  $V_\mu^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t) \right]$ .

## 2.2 Likelihoods

For discrete MDPs, we use product over each  $s, a$  pair of multinomials as likelihood,  $\mathbb{P}(\mathcal{D}_t | \mu) = \prod_{s,a} \frac{n!}{s_0! \dots s_{|\mathcal{S}|}!} p(s' = s_0 | s = s_0, a = a_0)^{x_0} \dots p(s' = s_{|\mathcal{S}|} | s = s_0, a = a_0)^{x_{|\mathcal{S}|}}$ .

For linear MDPs, we use linear models  $\mathcal{N}(\mathbf{s}^\top \boldsymbol{\mu}_t, \Sigma_t)$  as models. The likelihood is then a product of Gaussians, one for each  $(s, a \rightarrow s') \in \mathcal{D}_t$ .

Constraints on  $\Sigma_t$ , positive semi-definite, symmetric. IF we are even a weight vector  $\mathbf{w}$ , and the objective is to create a mean model, averaged over  $\mathbf{w}$ ,  $\mathbb{E}_{\mathbf{w}}[\mathcal{M}_s]$ , when  $\mathcal{M}_s$  is a set of linear-Gaussian models, then  $\boldsymbol{\mu}_{\mathbf{w}} = \sum_{i=1}^K w_i \boldsymbol{\mu}_i, \Sigma_{\mathbf{w}} = \sum_{i=1}^K w_i \Sigma_i$ .  $\Sigma_{\mathbf{w}}$  is guaranteed to be PSD as well.

## 2.3 Feature Representations

Doing model-based RL for settings with large and complicated state spaces, there might be some need for appropriate feature representations. In this work, we make use of randomised fourier features, based on the work of **ren2021free**.

## 3 Transfer Reinforcement Problem Formulation

Reinforcement learning for transfer learning augments the previous definition of MDPs through the consideration of *source* and a *target* domain. Given a set of source domains  $\Theta_s = \{\theta^1, \dots, \theta^K\} \in \Theta^K$  and a target domain  $\theta_t \in \Theta$ , we wish to leverage knowledge transfer from  $\Theta_s \rightarrow \theta_t$ . If  $\theta \in \Theta_s$ , we call this

setting a TYPE I setting and the considered set of  $\theta$  is finite. If  $\theta \in \mathcal{C}(\Theta_s)$ , where  $\mathcal{C} : \Theta^K \times [0, 1]^K \rightarrow \Theta$  is the convex set spanned by each of the  $\theta \in \Theta^K$ , then we call this a setting of TYPE II with infinite and realisable plausible models. If the assumptions on  $\theta$  is relaxed even more, and now  $\theta \in \delta(\Theta_s)$ , where  $\delta(\Theta_s)$  is the set of MDPs that are  $\delta$  close to  $\Theta_s$ , then we term this TYPE III and it is the set of infinite and non-realisable bounded plausible models. Finally, for TYPE IV, we consider arbitrary MDPs  $\theta \in \Theta$ .

[HE: TODO: explain the different settings here] [TT: introduce the  
TYPES intuitive names here, the names on the subsections]

### 3.1 Finite and realisable plausible models

For the first type of Bayesian TRL the parameters of the source models  $\mu \in \mathcal{M}_s$  are known, and the target model  $\mu_t \in \mathcal{M}_s$ .  $\mu_t$  is unknown, but data  $\mathcal{D}_t \sim \mu_t$  exists. This setting reduces to one of hierarchical Bayesian RL where the initial goal is to identify the posterior  $\mathbb{P}(\mu \in \mathcal{M}_s | \mathcal{D}_t)$ . When this has been attained, we can obtain the optimal controller  $\pi_{\text{Type I}}^*$  by maximising over policies  $\pi \in \Pi$ .

$$\pi_{\text{Type I}}^* = \arg \max_{\pi} V_{\mu}^{\pi} \mathbb{P}(\mu \in \mathcal{M}_s | \mathcal{D}_t) \quad (\text{E.1})$$

### 3.2 Infinite and realisable plausible models

In this setting we can arrive at two different solution methods, one *interior point method* and one fully Bayesian method with a constrained prior over the convex set of MDPs  $\mathcal{C}(\mathcal{M}_s)$ .

The interior point method starts by selecting a point  $\mathbf{p}_0 \in [0, 1]^K$  in the simplex with  $\mu \in \mathcal{M}_s$  being the vertices. A point in the simplex defines a convex combination of the known MDPs at the vertices. Iteratively, we can arrive at better and better points  $\mathbf{p}_1, \dots, \mathbf{p}_T$  by using the likelihood  $\mathbb{P}(\mathcal{D}_t | \mathbf{p})$ . Eventually, the procedure will converge at a point  $\mathbf{p}_*$ , this is the point in the simplex that represents the maximum likelihood MDP  $\mu^*$ . Of course the constraints are linear, (transitions), however, the likelihood is not convex.

Another method would be a gradient based method used to find the maximum likelihood MDP  $\mathbb{P}(\mathcal{D}_t | \mu \in \mathcal{C}(\mathcal{M}_s))$ .

$$\mu^* = \arg \max_{\mu \in \mathcal{C}(\mathcal{M}_s)} \mathbb{P}(\mathcal{D}_t | \mu) \quad (\text{E.2})$$

$$\pi_{\text{Type III}}^* = \arg \max_{\pi} V_{\mu^*}^{\pi} \quad (\text{E.3})$$

This is identical to the following,

$$\begin{aligned} w^* &= \arg \min_w -\mathbb{P}(\mathcal{D}_t | \mathbb{E}_w[\mathcal{M}_s]) \\ \text{s.t. } & \sum_{i=1}^K w_i = 1 \\ & \forall_w 0 \leq w_i \leq 1, \\ \mu^* &= \mathbb{E}_{w^*}[\mathcal{M}_s] \\ \pi_{\text{Type III}}^* &= \arg \max_{\pi} V_{\mu^*}^{\pi} \end{aligned} \quad (\text{E.4})$$

We can also consider a prior over the convex set of MDPs  $\mathcal{C}(\mathcal{M}_s)$ , such as the Dirichlet prior. Thus the posterior  $\mathbb{P}(\mu \in \mathcal{C}(\mathcal{M}_s) | \mathcal{D}_t)$  only has support over the simplex spanned by  $\mathcal{C}(\mathcal{M}_s)$  in model space  $\mathcal{M}$ . In this case, the optimal Bayesian controller could be written as the following,

$$\pi_{\text{Type III}}^* = \arg \max_{\pi} \int_{\mathcal{C}(\mathcal{M}_s)} V_{\mu}^{\pi} d\mathbb{P}(\mu | \mathcal{D}_t). \quad (\text{E.5})$$

A third possible approach is a variational Bayes approach, where the goal is to identify the parameters of the posterior (Dirichlet, with multinomial likelihood). The optimisation routine would then find the parameters of this posterior.

### 3.3 Infinite and non-realisable bounded plausible models

For the type V setting we have no structural constraints on the target MDP  $\mu_t$  and it can be outside the convex set of MDPs spanned by  $\mathcal{M}_s$ . Typically, to make this approach possible, we can apply some kind of constraints of the differences between the target MDP  $\mu_t$  and the source MDPs  $\mathcal{M}_s$ . For example,  $KL(\mu_t || \frac{1}{N} \sum_{\mu \in \mathcal{M}_s} \mu) \leq \delta$ , for discrete MDPs and e.g. the spectral radius / Eigenvalue of linear MDPs.

### Constrained

$$\begin{aligned} \mu^* &= \arg \max_{\mu} \mathbb{P}(\mathcal{D}_t | \mu) \\ \text{s.t. } & \text{KL}(\mu || \mathbb{E}[\mathcal{M}_s]) \leq \delta. \\ \pi^*_{\text{Type V}} &= \arg \max_{\pi} V_{\mu^*}^{\pi} \end{aligned} \quad (\text{E.6})$$

### Unconstrained

$$\begin{aligned} \mu^* &= \arg \max_{\mu} \mathbb{P}(\mathcal{D}_t | \mu) - \lambda \text{KL}(\mu || \mathbb{E}[\mathcal{M}_s]) \\ \pi^*_{\text{Type V}} &= \arg \max_{\pi} V_{\mu^*}^{\pi} \end{aligned} \quad (\text{E.7})$$

Currently we use a trust region method to optimise with  $|\mathcal{S}|$  linear constraints, one for each  $p(s' | s, a)$  pair, and  $|\mathcal{S}|^2 \times |\mathcal{A}|$  bounds. We use the unconstrained method as the addition of the KL bound would introduce non-linear constraints.

$$\begin{aligned} \mu^* &= \arg \max_{\mu} \mathbb{P}(\mathcal{D}_t | \mu) - \lambda \text{KL}(\mu || \mathbb{E}[\mathcal{M}_s]) \\ \text{s.t. } & \mathbf{A}\mathbf{p} = \mathbf{1}, \\ & \mathbf{0} \leq \mathbf{p} \leq \mathbf{1}. \\ \pi^*_{\text{Type V}} &= \arg \max_{\pi} V_{\mu^*}^{\pi} \end{aligned} \quad (\text{E.8})$$

## 3.4 Infinite and arbitrary plausible models

# 4 Maximum Likelihood Estimation of MDPs in Transfer Reinforcement Learning

## 4.1 Discrete MDPs

**Incorporating the Structure.** In this section we study general structured probability matrices with

$$\sum_{s'} p(s' | s, a) = 1, \quad (\text{E.9})$$

$$\forall_{(s, s', a) \in \mathcal{S}^2 \times \mathcal{A}} 0 \leq p(s' | s, a) \leq 1. \quad (\text{E.10})$$

The likelihood  $\mathcal{L}(\mathcal{D}, \theta)$  is a product over all the likelihoods for the individual transitions  $p(s' | s, a)$  and in this case, a product of *multinomials*.

$$\mathcal{L}(\mathcal{D}, \theta) = \prod_{s,a} \frac{n!}{x_1! \cdots x_{|\mathcal{S}|}!} \theta_1^{x_1} \cdots \theta_{|\mathcal{S}|}^{x_{|\mathcal{S}|}} \quad (\text{E.11})$$

As the parameter  $\theta$  that maximises this likelihood will also maximise the log likelihood  $\log \mathcal{L}(\mathcal{D}, \theta)$ , we will choose to write it in this form. The log-likelihood of a multinomial random variable can be written as,

$$\log \mathcal{L}(x_i, n | \theta) = \log(C) + \sum_{i=1}^k x_i \log \theta_i. \quad (\text{E.12})$$

The product over  $|\mathcal{S}| \times |\mathcal{A}|$  of pairs of these multinomials is thus,

$$\log \mathcal{L}(\mathcal{D}, \theta) = \log(C_1) + \cdots + \log(C_{|\mathcal{S}| \times |\mathcal{A}|}) \quad (\text{E.13})$$

$$+ \sum_{i=1}^k x_i^1 \log \theta_i^1 + \cdots + \sum_{i=1}^k x_i^{|\mathcal{S}| \times |\mathcal{A}|} \log \theta_i^{|\mathcal{S}| \times |\mathcal{A}|}. \quad (\text{E.14})$$

By now we have formulated our objective function, the log-likelihood to maximise. Let us now take a look at the constraints imposed on the MDP  $\mu$  we want to identify. A set of linear constraints in standard form looks like the following,  $\mathbf{l} \leq \mathbf{A}\theta \leq \mathbf{u}$ , where  $\mathbf{l}$  and  $\mathbf{u}$  are the respective lower and upper bounds of each of the constraints and  $\mathbf{A}\theta$  is the set of linear equations for which these constraints should hold. When the problem considered is a discrete MDP, these constraints take the form of,

$$\{\theta_{i \times |\mathcal{S}|}, \dots, \theta_{(i+1) \times |\mathcal{S}|}\} \subset \{\theta_1, \dots, \theta_{|\mathcal{S}|^2 \times |\mathcal{A}|}\} \quad (\text{E.15})$$

$$\subseteq \text{vec}(\mu) = \theta, \quad (\text{E.16})$$

where  $\sum_{j=0}^{|\mathcal{S}|} \theta_{i+j} = 1$ . This means that  $\mathbf{l} = \mathbf{u} = \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|}$ , and  $\mathbf{A}$  is a  $(|\mathcal{S}| \times |\mathcal{A}|) \times (|\mathcal{S}|^2 \times |\mathcal{A}|)$  matrix that specifies where the constraints are active.

$$\mathbf{A} = \begin{bmatrix} \mathbf{1}_{|\mathcal{S}|}^\top & & \\ & \ddots & \\ & & \mathbf{1}_{|\mathcal{S}|}^\top \end{bmatrix} \quad (\text{E.17})$$

In the case where we only consider convex combinations of MDPs, the constraint is much simpler,  $\mathbf{l} = \mathbf{u} = \mathbf{1}, \mathbf{A} = \mathbf{1}_{|\Theta_s|}^\top$ .

**Designing the Optimiser.** In order to identify the maximum likelihood MDP given these constraints, we use a trust region optimiser **EQSQP byrd1987robust**, **omojokun1989trust**, **lalee1998implementation** with the following objective,

$$\begin{aligned} \min_{\theta} \quad & f(\theta) \\ \text{s.t.} \quad & \mathbf{l} \leq \mathbf{A}\theta \leq \mathbf{u}, \end{aligned} \tag{E.18}$$

where  $f(\theta)$ ,  $\mathbf{l}$ ,  $\mathbf{u}$  and  $\mathbf{A}$  has to be selected appropriately.

Let the objective function be one of the following,

$$f(\theta) = \begin{cases} -\log \mathcal{L}(\mathcal{D}_t, \mathbb{E}_{\mathbf{w}}[\mathcal{M}_s]) \\ -\log \mathcal{L}(\mathcal{D}_t, \theta) + \nu \text{KL}(\theta || \mathbb{E}[\mathcal{M}_s]) \\ -\log \mathcal{L}(\mathcal{D}_t, \theta) \end{cases} \tag{E.19}$$

The solution  $\theta^*$  is guaranteed to be a proper probability matrix but because of the non-convexity of the likelihood it might be a local minima. As the problem is already non-convex, adding another non-convex penalty term such as  $\lambda \text{KL}(\theta || \mathbb{E}[\mathcal{M}_s])$  will not change this.

The EQSQP technique attempts to solve the following quadratic problem, where  $\Delta_k$  is the trust region radius at iteration  $k$  and step  $\mathbf{d}_k$  is the next step to iterate with.

$$\begin{aligned} \min_{\mathbf{d}} \quad & \mathbf{d}^\top \nabla f(\theta) + \frac{1}{2} \mathbf{d}^\top \nabla_\theta^2(f(\theta) - \lambda^\top \mathbf{A}\theta) \mathbf{d}, \\ \text{s.t.} \quad & \mathbf{J}\mathbf{d} + \mathbf{A}\theta = r_k, \\ & \|\mathbf{d}\|_2 \leq \Delta_k. \end{aligned} \tag{E.20}$$

$$\min_{\|\mathbf{d}\|_2 \leq \Delta_k} \mathbf{d}^\top \nabla_\theta f(\theta) + \frac{1}{2} \mathbf{d}^\top \nabla_\theta^2 f(\theta) \mathbf{d} + \frac{1}{2} \mathbf{d}^\top (\lambda^\top \nabla_\theta^2(r_k - \mathbf{J}\mathbf{d})) \mathbf{d}. \tag{E.21}$$

Where  $\mathbf{J} = \left[ \frac{\partial \mathbf{A}\theta}{\partial \theta_1} \dots \frac{\partial \mathbf{A}\theta}{\partial \theta_{|\mathcal{S}|^2 \times |\mathcal{A}|}} \right]$  is the Jacobian of the constraints. The details of this procedure can be found in **byrd1987robust**, **omojokun1989trust**, **lalee1998implementation**. After convergence, the ML MDP  $\theta^*$  has been identified, we can use **VALUEITERATION** to identify the optimal way of acting in  $\theta^*$ .

[DB: Applying the projection to simplex will requires using the algorithm  $SA$  times which can be computationally heavy and cumbersome. Using EQSQP provides a better way to handle the constraints.]

## 4.2 Linear MDPs

In this section we are studying linear-Gaussian MDPs  $\theta = \{\theta_{\mathbf{A}}, \theta_{\mathbf{B}}, \theta_{\Sigma}, \theta_{\mathbf{Q}}, \theta_{\mathbf{R}}, \sigma^2\}$  of the following form for the transition kernel  $\mathbb{P}(s' | s, a)$ ,

$$\begin{aligned}\Delta s_{t+1} &= \theta_{\mathbf{A}} s_t + \theta_{\mathbf{B}} a + e_1 \\ e_1 &\sim \mathcal{N}(\mathbf{0}, \theta_{\Sigma}).\end{aligned}$$

Likewise, for the reward function  $\mathbb{P}(r | s, a)$  we have,

$$\begin{aligned}r_t &= s_t^\top \theta_{\mathbf{Q}} s_t + a^\top \theta_{\mathbf{R}} a + e_2 \\ e_2 &\sim \mathcal{N}(0, \sigma^2)\end{aligned}$$

The likelihood in this case is a product of Gaussians. As the  $\theta$  that maximises this likelihood also maximises the log-likelihood, we choose to maximise the following log-likelihood,

$$\log \mathcal{L}(\mathcal{D}, \theta) = \sum_{(s, a, s') \in \mathcal{D}} \frac{\log |\theta_{\Sigma}|}{2} + \frac{(s' - \mu)^\top \theta_{\Sigma}^{-1} (s' - \mu)}{2} \quad (\text{E.22})$$

$$\mu = \theta_{\mathbf{A}} s + \theta_{\mathbf{B}} a \quad (\text{E.23})$$

The constraints we have on the parameters  $\theta$  aim to ensure  $\theta_{\Sigma}$  being positive semi-definite and symmetric. The parameters  $\theta_{\mathbf{A}} \in [-\infty, \infty]^{d \times d}, \theta_{\mathbf{B}} \in [-\infty, \infty]^{d \times a}, \theta_{\Sigma} \in [0, \infty]^{d \times d/2}, \theta_{\mathbf{Q}} \in (0, \infty]^{d \times d}, \theta_{\mathbf{R}} \in (0, \infty]^{a \times a}, \sigma^2 > 0$ .

**Estimating  $\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}$  from  $s_t, s_{t+1}, u_t, r_t$ .**

$$\begin{aligned}\mathbf{CD} &= \mathbf{AS} + \mathbf{BU} = \mathbf{S}' - \mathbf{S} \\ \mathbf{C} &= (\mathbf{S}' - \mathbf{S})\mathbf{D}^{-1}\end{aligned}$$

The matrix  $\mathbf{C}$  then contains the parameters of  $\mathbf{A}$  and  $\mathbf{B}$ .  $\mathbf{D}$  is a matrix with the data  $\mathbf{S}$  and  $\mathbf{U}$  stacked appropriately. If  $N$  is the number of observations,

then  $\mathbf{C} = |\mathcal{S}| \times (|\mathcal{S}| + |\mathcal{A}|)$ ,  $\mathbf{D} = (|\mathcal{S}| + |\mathcal{A}|) \times N$ ,  $\mathbf{S} = \mathbf{S}' = |\mathcal{S}| \times N$  and  $\mathbf{U} = |\mathcal{A}| \times N$ .  $\mathbf{C}$  is essentially a matrix like the following, for  $|\mathcal{S}| = 3, |\mathcal{A}| = 2$ , where  $\mathbf{A} = |\mathcal{S}| \times |\mathcal{S}|$  and  $\mathbf{B} = |\mathcal{S}| \times |\mathcal{A}|$ .

$$\mathbf{C} = \left[ \begin{array}{c|c} \mathbf{A} & \mathbf{B} \end{array} \right] \quad (\text{E.24})$$

$$\mathbf{CD} = \mathbf{S}^\top \mathbf{QS} + \mathbf{U}^\top \mathbf{RU} = \mathbf{r} \quad (\text{E.25})$$

$$\left[ \begin{array}{c|c|c|c} \mathbf{A} & \mathbf{B} & \mathbf{Q} & \mathbf{R} \end{array} \right] \left[ \begin{array}{c|c|c|c} \mathbf{S} & \mathbf{U} & \mathbf{S} & \mathbf{U} \end{array} \right] = \left[ \begin{array}{c|c} \mathbf{S}' - \mathbf{S} & \mathbf{r} \end{array} \right] \quad (\text{E.26})$$

$$(|\mathcal{S}| + 1) \times 2(|\mathcal{S}| + |\mathcal{A}|) * 2(|\mathcal{S}| + |\mathcal{A}|) \times N = (|\mathcal{S}| + 1) \times N$$

## 5 Planning in Markov Decision Processes

### 5.1 Discrete MDPs

For discrete: use VALUEITERATION

### 5.2 Linear MDPs

For linear: use CARE (Continuous time algebraic Riccati equation) Given a posterior over models,  $\mathbb{P}(\mu | \mathcal{D})$ , we wish to arrive at a linear control policy  $\mathbf{u} = -\mathbf{K}\mathbf{s}^\top$ , where  $\mathbf{s}$  is the observed current state and  $\mathbf{K}$  the Kalman gain. The idea behind this algorithm is to first sample a model  $\mu \sim \mathbb{P}(\mu | \mathcal{D})$ , then use this model as the actual model, and infer the optimal linear control policy for said model. Since we might have a wide distribution of plausible MDPs, this algorithm might perform poorly with low data and/or when linear models can not sufficiently describe the transition probabilities.

$$\Delta\mathbf{s} = \mathbf{As} + \mathbf{Bu} \quad (\text{E.27})$$

This describes the change in state after applying control input  $\mathbf{u}$  given the specification systems  $\mathbf{A}, \mathbf{B}$ .  $\mathbf{A}$  and  $\mathbf{B}$  needs to be inferred from the sampled MDP  $\mu$ . We also specify two more matrices,  $\mathbf{Q} = \mathbf{I}_{|\mathcal{S}|}$  and  $\mathbf{R} = \mathbf{I}_{|\mathcal{A}|}$ .

The linear system described in E.27, from time step  $t = [0, T]$ , has the following cost function,

$$J = \mathbf{s}_T^\top \mathbf{s}_T + \sum_{t=0}^{T-1} (\mathbf{s}_t^\top \mathbf{s}_t + \mathbf{u}_t^\top \mathbf{u}_t) \quad (\text{E.28})$$

If we instead look at the Infinite-horizon and discrete time LQR, we have,

$$J = \sum_{t=0}^{\infty} (\mathbf{s}_t^\top \mathbf{s}_t + \mathbf{u}_t^\top \mathbf{u}_t) \quad (\text{E.29})$$

This is minimised by setting  $\mathbf{u} = -\mathbf{K}\mathbf{s}$ , where

$$\mathbf{K} = (\mathbf{R} + \mathbf{B}^\top \mathbf{P} \mathbf{B})^{-1} (\mathbf{B}^\top \mathbf{P} \mathbf{A}) \quad (\text{E.30})$$

$\mathbf{P}$  can be attained solving the discrete-time algebraic Riccati equation,

$$\mathbf{A}^\top \mathbf{P} \mathbf{A} - \mathbf{P} - (\mathbf{A}^\top \mathbf{P} \mathbf{B})(\mathbf{R} + \mathbf{B}^\top \mathbf{P} \mathbf{B})^{-1} (\mathbf{B}^\top \mathbf{P} \mathbf{A}) + \mathbf{Q} = 0 \quad (\text{E.31})$$

## 6 Algorithms

In this section we introduce a meta algorithm, TRLBI, that appropriately computes a distribution over MDPs  $\xi(\mu)$ , or a maximum likelihood model  $\mu^*$  and feeds this into MMBI as introduced in **dimitrakakis:mmbi:ewrl:2011**. Note that if a single MDP  $\mu$  is sent into MMBI, then it will reduce to **PSRLosband2013more**.

---

**Algorithm 3** Transfer Reinforcement Learning Backward Induction

---

TRLBI( $f(\cdot)$ ,  $\{\mathbf{l}, \mathbf{u}, \mathbf{A}\}$ ,  $\gamma$ ,  $T$ ,  $K$ )

**input** objective function  $f(\theta)$ , constraints  $\{\mathbf{l}, \mathbf{u}, \mathbf{A}\}$ , discount factor  $\gamma$   
**for**  $t = 0, \dots, T$  **do**

    MODEL ESTIMATION

$$\theta_t^0 = \theta_{t-1}^*$$

**for**  $k = 0, \dots, K$  **do**

$$\min_{\mathbf{d}_k} \mathbf{d}_k^\top \nabla f(\theta_t^k) + \frac{1}{2} \mathbf{d}_k^\top \nabla_\theta^2(f(\theta_t^k) - \lambda_k^\top \mathbf{A} \theta_t^k) \mathbf{d}_k,$$

$$\text{s.t. } \mathbf{J} \mathbf{d}_k + \mathbf{A} \theta_t^k = r_k,$$

$$\|\mathbf{d}_k\|_2 \leq \Delta_k.$$

$$\theta_t^* = \theta_t^K$$

    PLANNING

$$V_t^K = \mathbf{0}$$

**for**  $k = K-1, \dots, 0$  **do**

$$Q_t^k = \mathcal{R} + \gamma \theta_t^* V_t^{k+1}$$

$$V_t^k = \max Q_t^k$$

$$V_t^* = V_t^0$$

    CONTROL

$$\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{s_t, a_t, s_{t+1}\}$$


---

---

**Algorithm 4** Transfer Reinforcement Learning Optimal Control

---

TRLOC( $f(\cdot)$ ,  $\{\mathbf{l}, \mathbf{u}, \mathbf{A}\}$ ,  $\gamma$ ,  $T$ ,  $K$ )

**input** objective function  $f(\theta)$ , constraints  $\{\mathbf{l}, \mathbf{u}, \mathbf{A}\}$ , discount factor  $\gamma$   
**for**  $t = 0, \dots, T$  **do**

    MODEL ESTIMATION

$$\theta_t^* = \theta_t^K$$

    PLANNING

    solve  
 $\mathbf{P}$

$$\mathbf{A}^\top \mathbf{P} \mathbf{A} - \mathbf{P} + \mathbf{Q}$$

$$-(\mathbf{A}^\top \mathbf{P} \mathbf{B})(\mathbf{R} + \mathbf{B}^\top \mathbf{P} \mathbf{B})^{-1}(\mathbf{B}^\top \mathbf{P} \mathbf{A}) = 0$$

$$\mathbf{K} = (\mathbf{R} + \mathbf{B}^\top \mathbf{P} \mathbf{B})^{-1}(\mathbf{B}^\top \mathbf{P} \mathbf{A})$$

$$\mathbf{u} = -\mathbf{K} \mathbf{s}$$

    CONTROL

$$\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{s_t, a_t, s_{t+1}\}$$

E15

---

## 7 Bounds of TRLBI

In this section we will be studying the case of using a TYPE III algorithm in a TYPE V problem setting. That means the target MDP  $\mu_t$  might not be part of the convex set  $\mathcal{C}(\mathcal{M}_s)$  of source MDPs.  $\mu_t$  will however be bounded in distance from the marginal of the source MDPs. This setting gives rise to three sources of estimation errors, namely *value estimation error*, *model estimation error* and *realizability error*. Combined, they might look like the following,

$$\|V_{\mu_t}^* - \hat{V}_{\mu_t}^*\| + \|V_{\mu^*}^* - V_{\hat{\mu}^*}^*\| + \|V_{\mu_t}^* - V_{\mu^*}^*\| \quad (\text{E.32})$$

If we have a proper distance metric as model bound, then we can use the triangle inequality to get,

$$\leq \|V_{\mu_t}^* - \hat{V}_{\mu_t}^*\| + \|V_{\mu^*}^* - V_{\hat{\mu}^*}^*\| + \|V_{\hat{\mu}^*}^* - V_{\mu_t}^*\| \quad (\text{E.33})$$

One part of this comes from error in model estimation, thus we can write this as,

$$\leq \|V_{\mu_t}^* - \hat{V}_{\mu_t}^*\| + f\left(\|\mu^* - \hat{\mu}^* + \mu_t - \mu^*\|\right), \quad (\text{E.34})$$

where  $f(\cdot)$  is a function that produces a value function bound based on model difference.

$$\leq \mathcal{O}(\text{VALUEITERATION}) + f\left(\left\| \begin{array}{c} \mu^* - \hat{\mu}^* \\ \text{ML-ESTIMATE} \end{array} + \begin{array}{c} \mu_t - \mu^* \\ \text{NON-REALISABILITY} \end{array} \right\| \right) \quad (\text{E.35})$$

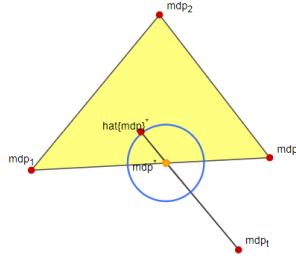
If we have infinite data, and the model is realisable, then the errors from ValueIteration, ML-estimate and non-realisation should go to zero. If we have infinite data, and the model is not realisable, then the errors from ValueIteration and ML-estimate should go to zero.

[HE: Agnostic bounds <https://sites.stat.washington.edu/jaw/COURSES/580s/581/LECTNOTES/ch4-rev1.pdf>]

$\mathcal{O}(\text{VALUEITERATION})$  should be a finite sample bound that depends on the horizon  $H$ , error  $\epsilon$ , state space  $\mathcal{S}$  and action space  $\mathcal{A}$ .

$f(\|\mu^* - \hat{\mu}^*\|)$  should be a finite sample bound that depends on the size of the data  $\mathcal{D}_t$ , the likelihood  $\mathbb{P}(\mathcal{D}_t | \mu)$ , the optimiser, the state space  $\mathcal{S}$  and the action space  $\mathcal{A}$ .

$f(||\mu_t - \mu^*||)$  should be a bound based on model difference. Trivially,  $f(||\mu_t - \mu^*||) \leq f(\delta)$ , as  $\delta$  bounds the difference between the marginal MDP and the target MDP. The best possible model in the convex set should be closer to the target model than the marginal model.



**Figure 1:** Model space geometry of the bound discussed in this section. The target MDP  $\mu_t$  is non-realisable and the best approximation in the realisable set is  $\mu^*$ . The ML-estimator estimating  $\mu^*$  gives us a set of plausible models, constrained by the convex set. The model furthest away from the target model is  $\hat{\mu}^*$ .

Even-Dar **even2009online** could possibly be the citation for this bound.

$$||\mu_1 - \mu_2||_1 \leq \delta \leq D_{KL}(\mu_1 || \mu_2) \quad (\text{E.36})$$

$$\implies ||V_{\mu_1}^\pi - V_{\mu_2}^\pi||_\infty \leq \frac{\epsilon}{(1-\gamma)^2} \quad (\text{E.37})$$

[HE: It's possible we can also say something about the stability in the linear case if the spectral radius is bounded.]

Lemmas of interest:

Mixing,  $d, d'$  are distributions of states,  $P^\pi$  the transition model given policy  $\pi$ .  $\tau$  mixing time.

$$||dP^\pi - d'P^\pi||_1 \leq e^{-1/\tau} ||d - d'||_1 \quad (\text{E.38})$$

State distribution,

$$||d_{\pi,t} - d_\pi||_1 \leq 2e^{-t/\tau} \quad (\text{E.39})$$

Value function,

$$Q_{\pi,r}(s, a) \leq 3\tau \quad (\text{E.40})$$

## 7.1 Value Estimation Error

The value estimation error arises when trying to estimate  $V_{\mu_t}^*$  given that we actually know the model  $\mu_t$ . This bound, for discrete MDPs, should be available to get as a finite sample bound for VALUEITERATION.

## 7.2 Model Estimation Error

The second part of the error terms involve correctly estimating the maximum likelihood-model  $\mu^*$  from data  $\mathcal{D}_t$ .

## 7.3 Realizability Error

The third and final part of the error term is related to the difference between the best possible model in the convex set  $\mathcal{C}(\mathcal{M}_s)$  and the actual target MDP  $\mu_t$ . Since we have a bound on the difference in MDP  $KL(\mu_t \parallel \frac{1}{N} \sum_{\mu \in \mathcal{M}_s} \mu) \leq \delta$ , we should be able to get a bound also on the value function here.

# 8 Experiments

## 8.1 Discrete MDPs

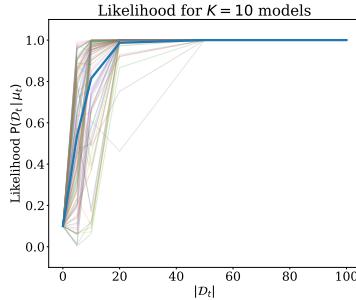
We evaluate our algorithms for discrete MDPs and settings of *Type I-VI*. We assume the reward functions are identical, and put our interest into MDPs with different transition matrices. We evaluate our algorithms against a benchmark that simply uses *PSRL* with the data from the target MDP to show that knowledge transfer from the source MDPs help with performance.

The algorithms we are evaluating are thus, *PSRL*, *EM-MMBI*, *EM-PG* and possibly a fully Bayesian algorithm that does model selection correctly.

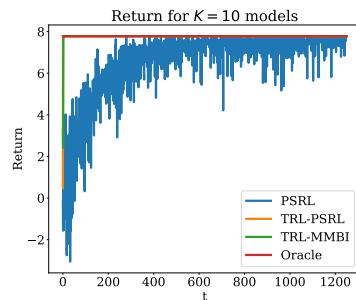
[HE: TODO: make interesting gridworld AD scenario]

## 8.2 Continous Space and Action

Further we evaluate our algorithms in continuous state and action spaces. For instance, CartPole, Mujoco.



**Figure 2:** Likelihood experiment with 100 independent runs. The number of models considered was 10. The thick blue line is the mean likelihood for the true model  $\mu_t$  averaged over experiments.

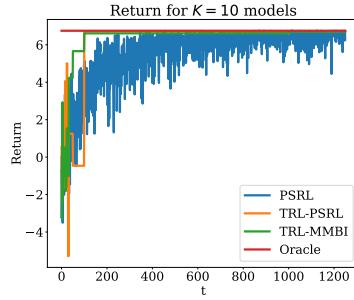


**Figure 3:** Experiment of type I. The oracle agent knows which MDP generates the data  $\mathcal{D}_t$ . The TRL-MMBI agent uses the likelihood  $\mathbb{P}(\mathcal{D}_t | \mu)$  to weigh the MDPs. PSRL learns the MDP from scratch using  $\mathcal{D}_t$ .

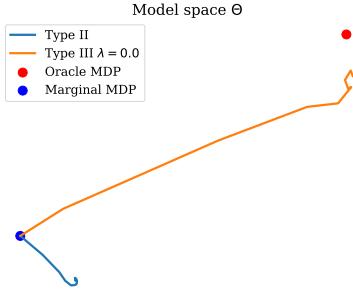
## 9 Autonomous Driving Scenarios with Knowledge Transfer

Random fourier features to linear mdps

free lunch with noise <https://arxiv.org/pdf/2111.11485.pdf>



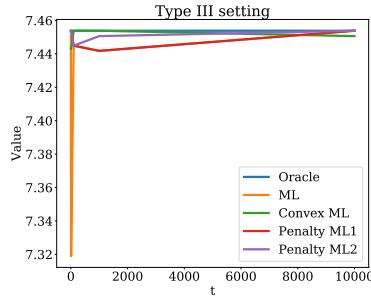
**Figure 4:** Experiment of type II. The oracle agent knows which MDP generates the data  $\mathcal{D}_t$ . The TRL-MMBI agent uses the likelihood  $\mathbb{P}(\mathcal{D}_t | \mathcal{D}_i, \mu)$  to weigh the MDPs. PSRL learns the MDP from scratch using  $\mathcal{D}_t$ .



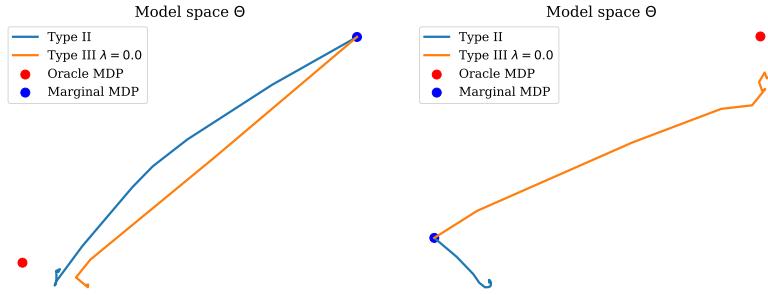
**Figure 5:** Experiment of type III. Testing convergence of the minimiser using SLSQP to find the point in the simplex that corresponds to the convex MDP with the lowest negative log-likelihood. The MDP  $\mu^* = \sum_{i=0}^K w_i \mu_i$  is the corresponding MDP the minimiser converged to.  $\mu^{Oracle} = \sum_{i=0}^K w_i^{Oracle} \mu_i$  is the MDP that corresponds to the MDP in the convex set at a given point in the simplex. The data  $\mathcal{D}_t \sim \mu^{Oracle}$ .

## 9.1 Setting 1: Transfer knowledge from similar tasks, e.g. cut-in to overtake

[HE: Motivation: here we transfer knowledge from one or multiple tasks to a new task, which might have similar structure to it (in the MDP)]



**Figure 6:** Experiment of type III.



**Figure 7:** Experiment of convex MDP. The left plot is made using PCA over the model space.

## 9.2 Setting 2: Transfer knowledge from similar cars, e.g. sedan to combi

[HE: Motivation: an optimal controller must necessarily take into account the dynamics of the Ego vehicle. If a controller is learned for a smaller car, it's unlikely it will be optimal for a heavier car.]

Environment	Type	$\mu_t \in \Delta$			$\mu_t \notin \Delta$	
		$\alpha = 0$	$\alpha = 1$	$\alpha = 2$	$\alpha = 3$	$\alpha = 4$
<i>CartPole</i>	I					
	III					
<i>dm_LQR_2_1</i>	I					
	III					
<i>dm_LQR_6_2</i>	I					
	III					
		200.00	31.63	38.9	35.87	37.10
		200.00	200.00	57.91	61.62	200.00
		44.22	29.96	200.00	34.19	35.54
		200.00	200.00	28.24	200.00	200.00
		200.00	200.00	25.97	125.80	200.00

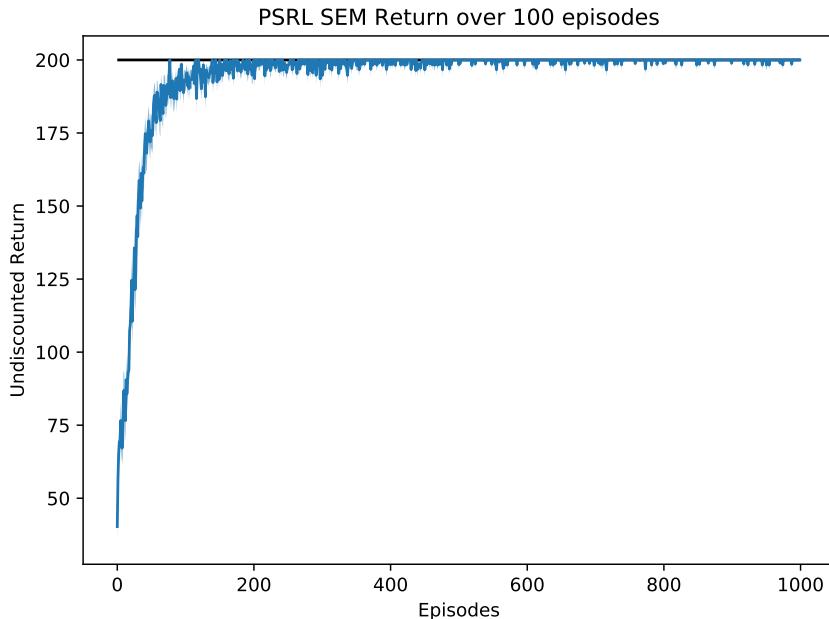
**Table 2:** Utility matrix

### 9.3 Setting 3a: Transfer knowledge from similar driving locations, e.g. USA to China, using traffic rules, left-side / right-side traffic

[HE: Motivation: in different countries we have different traffic rules, etc.]

		$\mu_j$				
		0.00	168.37	161.10	164.13	162.90
		0.00	0.00	142.09	138.38	0.00
$\pi^*(\mu_i)$	155.78	170.04	0.00	165.81	164.46	
		0.00	0.00	171.76	0.00	0.00
		0.00	0.00	174.03	74.20	0.00

**Table 3:** Regret matrix



**Figure 8:** Experiment of CartPole for PSRL over 100 experiments. The shaded region is the the standard error of the mean over experiments.

#### 9.4 Setting 3b: Transfer knowledge from similar driving locations, e.g. USA to China, using driver interaction

[HE: Motivation: likewise as the previous, in different countries we have different kinds of drivers, behaviors, etc.]