# **Towards Practical Hierarchical Reinforcement Learning for Multi-lane Autonomous Driving**

## Anonymous Author(s) Affiliation Address email

## Abstract

1 In this paper, we propose an approach for making hierarchical reinforcement learning practical for autonomous driving on multi-lane highway or urban structured 2 roads. While this approach follows the conventional hierarchy of behavior decision, 3 motion planning, and control, it introduces an intermediate layer of abstraction 4 that specifically discretizes the state-action space for motion planning according to 5 a given behavioral decision. This hierarchical design allows principled modular 6 extension of motion planning, in contrast to relying on either monolithic behavior 7 cloning or a large set of hand-written rules. We show that this design enables 8 significantly faster learning than a flat design, when using both value-based and 9 policy optimization methods (DQN and PPO). We also show that this design allows 10 transferring of the trained models, without any retraining, from a simulated environ-11 ment with virtually no dynamics to one with significantly more realistic dynamics. 12 Overall, our proposed approach is a promising way to allow reinforcement learning 13 to be applied to complex multi-lane driving in the real world. In addition, we 14 introduce and release an open source simulator for multi-lane driving that follows 15 the OpenAI Gym APIs and is suitable for reinforcement learning research. 16

# 17 **1 Introduction**

Developing autonomous cars that reliably assist humans in everyday transportation is a grand research and engineering challenge. Mundane as it may seem, cruising on a multi-lane highway effectively and safely while taking full advantage of available driving space has proved challenging for existing autonomous cars. What makes multi-lane cruising significantly more challenging than the single-lane adaptive cruise control (ACC) is the fact that the multi-vehicle interaction happens both laterally (i.e. perpendicular to the lanes) and longitudinally (i.e. parallel to the lanes) and requires coordination between lateral and speed control.

The current research focuses on multi-lane cruising as a special case of driving on structured roads in 25 general. Different from classical feedback control scenarios, driving on structured roads is heavily 26 regulated by signs, signals, and rules that come to apply at various points in time, space, or even 27 history (e.g. all-way stop intersections). In the case of multi-lane cruising, lane markings dictate that 28 driving takes place mostly within the boundaries of a single lane. Lane change is thus a short-lived, 29 transitional event in continuous motion space that links two distinct states — driving in one lane 30 vs. driving in an adjacent lane. Furthermore, because driving on structured roads in general is thus 31 symbolically punctuated<sup>1</sup>, it is naturally hierarchical — higher level decisions on discrete state 32 transitions are coordinated with lower level motion control in continuous state space. Because multi-33 lane cruising shares this hierarchical character, we propose using a hierarchical design for autonomous 34 driving systems (similar to [16]), trusting that lessons we learn from this case can be generalized 35 to other aspects of autonomous driving on structured roads. For this research, we take multi-lane 36

<sup>&</sup>lt;sup>1</sup>The symbols here could include traffic lights, lane markings, speed limit signs, fire truck sirens, signals of other vehicles, and so on.

cruising problem, requiring (a) lane changes, (b) in-lane maneuvers, and (c) speed control in order to maintain desirable speed, avoid obstacles, and stay safe. Furthermore, we have operationalized these

essential aspects of the multi-lane cruising task in a small video game called  $SimpleTraffic^2$ , which

40 follows the OpenAI Gyms API conventions [2] and is conducive to reinforcement learning research.

Autonomous driving systems typically consist of modules for perception and prediction, localization 41 and mapping, and planning and control. Our study focuses on the planning and control part, where 42 a hierarchy of behavior decision, motion planning, and control is widely adopted [10]. Classical 43 methods for implementing the behavior planner are largely rule-based with finite state machines being 44 a common choice [10]. Classical motion planner methods typically require optimization according 45 to explicitly defined cost functions with the behavior decision expressed as constraint terms in the 46 cost function [4, 18]. While rule-based behavior planner is extremely hard to maintain and does 47 not scale well in complex dynamic scenarios, explicit cost functions for motion planner are hardly 48 general enough and very difficult to tune for complex dynamic interactions. To avoid the limitations 49 of traditional planning solutions, many recent studies attempted learning-based approaches. Bojarski 50 et al. [1] proposed an end-to-end supervised learning scheme that directly maps images to steering 51 commands. Sun et al. [17] in contrast use a mapping from state features to trajectories and then use 52 an execution layer to further guarantee short term feasibility and safety. These approaches leverage 53 expert experience for training. However, by directly cloning the expert's driving strategy, they are 54 limited to the expert's performance and experience, failing to adequately explore the parts of the 55 state-action space that may be less critical for safety and performance. In addition, planning and 56 control is largely implemented as one monolithic network, which makes debugging, failure analysis, 57 and incorporation of domain knowledge all very difficult. 58

<sup>59</sup> With the recent great success of reinforcement learning (RL) methods in decision making and control <sup>60</sup> (*e.g.* [3,8,19]), RL methods are now frequently applied to autonomous driving tasks [6, 12, 13]. In <sup>61</sup> such cases, RL agents learn driving tasks by exploring typically simulated driving environments while <sup>62</sup> attempting to improve performance along some metrics such as safety, average speed, smoothness of <sup>63</sup> maneuvers, and smoothness of speed profile. The hope here is that such RL agents will outperform <sup>64</sup> classical solutions through scaling significantly better for complex dynamic scenarios, and thereby <sup>65</sup> making competent autonomous driving available much more broadly.

In the following, we aim to show that our design allows effective integration of RL for autonomous 66 driving, by (1) eliminating the dependency on explicit cost functions in motion planning, (2) enabling 67 significantly faster learning performance than a non-hierarchical or flat design, (3) allowing principled 68 modular extension of motion planning in contrast to using either monolithic behavior cloning or a 69 large set of hand-written rules. This in turn suggest that our proposed hierarchical architecture is a 70 promising way to allow reinforcement learning to be applied to complex multi-lane cruising in the 71 real world, which we demonstrate through transferring the trained models without retraining from a 72 simulated environment with virtually no dynamics to one with significantly more realistic dynamics. 73

To summarize, the main contributions of this paper are: (1) proposing a modularized skill-based 74 planning framework with two layers of hierarchy (behavioral and motion planner) for cruising in multi-75 lane roads; (2) proposing a higher level of abstraction in the state-action space of driving in multi-lane 76 roads; and (3) introducing an open source simulator, SimpleTraffic, for RL and autonomous driving 77 research community. In Section II, we review the related state-of-the-art methods in autonomous 78 driving. In Section III, we present the details of our planning framework. Section IV introduces 79 SimpleTraffic simulator. In Section V we evaluate our approach comprehensively and conclude our 80 work in Section VI. 81

# 82 2 Related works

Recent studies have utilized RL for higher level decision making [7,9]. Mukadam et al. [9] trained 83 a Q-network to issue only high level commands, e.g. switch left/right, while the execution and 84 collision avoidance rely on the low-level rule-based controller. Mirchevska et al. [7] proposed an 85 RL-based approach for autonomous driving in highway scenarios using the fitted Q-iteration with 86 extremely randomized trees as a function approximator. Both of these approaches have utilized RL 87 for high level decision making (BP in our case) only and adopted hand-crafted rule-based approaches 88 for MoP [7,9] and collision avoidance [9]. In contrast, in addition to behavioral planner, motion 89 planner and collision avoidance are also handled by our hierarchical learning-based approach which 90

<sup>&</sup>lt;sup>2</sup>To be open sourced upon the publication of this paper.

is more generalizable compared to its rule-based counterparts. Shalev-Shwartz et al. [15] proposed to 91 decompose the autonomous driving problem into two phases. First, they apply supervised learning 92 for predicting the near future states based on the current states. Second, they model a full trajectory 93 of the agent using a recurrent neural network, where unexplained factors are modeled as (additive) 94 input nodes. As this method uses fully supervised learning approach, it relies on the training data 95 too much. Paxton et al. [11] learns both low-level control policies and high level task policies by 96 97 integrating a learning-based approach (Monte Carlo Tree Search) with a heuristic search algorithm to achieve a complex task. Their learned low-level policies sometimes can have oscillatory behavior 98 that makes the car unstable. Given well established controllers such as PID and MPC, we believe 99 that learning-based methods are more effective in the high and mid levels (e.g. BP and MoP) of a 100 decision making rather than in the low-level controllers. 101

## **102 3 Technical approach**

## 103 3.1 The Planning Hierarchy

Driving is a symbolically punctuated behavior. Different from classical feedback control scenarios, 104 driving is heavily punctuated by signs and rules on top of what is largely a continuous control task. 105 The symbols here include lane markings, traffic lights, speed limit signs, fire truck sirens, turning 106 signals of the heading vehicle, and so on. As an example, lane markings dictates that most driving 107 happen within a single lane. Lane changes are thus short-lived, transitional events that link forward 108 driving in one lane to forward driving in an adjacent lane - two discrete states at a higher level 109 of abstraction in the state space of driving. Because driving is thus symbolically punctuated, it is 110 naturally hierarchical: higher level decisions on discrete state transitions with lower level execution in 111 112 continuous state space, which suggests a hierarchical structure in the design of planning systems for autonomous driving. Figure 1 illustrates our proposed hierarchical decision making architecture for 113 cruise in multi-lane roads. As shown in Figure 1, the proposed decision making framework includes 114 behavioral planner that makes high level decisions about transitions between discrete states, and 115 motion planner that generates a target spatio-temporal trajectory with a target speed according to 116 the decisions made by BP. The target trajectory is then fed to the controller to follow the trajectory 117 by controlling the steering wheel, throttle, and brake in continuous state space. The hierarchical 118 structure of our planning framework facilitates analysis of the decisions that are made during driving. 119 In addition, structured layers allows for convenience for modularizing different *skills*, *e.g.* adaptive 120 cruise control, lane switching, etc. Each modularized skill acts as an independent entity and forms a 121 comprehensive maneuver function considering its own constraints and safety internally. This also 122 enables modifing and replacing sub-modules according to new requirements and conditions. 123

#### 124 **3.2 Behavior Planner**

Behavior decision is about transitioning between states that are discrete only at a higher level of abstraction. BP is responsible to drive the car to the destination safely and as fast as possible. In our current setting, BP makes high level decisions including *keep lane*, *switch to the left lane*, and *switch to the right lane* subject to the following conditions: (1) following the mission planner (*drive to the destination*); (2) navigating the ego-car to less busy lanes so the car can drive to the maximum speed limit (*drive as fast as possible*); and (3) avoiding collisions (*drive safely*).

BP takes the full set of states as input which includes:, target lane suggested by Mission planner 131 (MiP), current lane, current speed, speed and distance to the nearest heading car for current and 132 neighboring lanes, and speed and distance to nearest car in the back for current and neighboring lanes. 133 We design a *coarse-grained* reward function and avoid any *fine-grained* rules in our reward feedback. 134 This way, we give the RL agent a chance to explore the state space and to come up with solutions 135 that possibly outperform classical rule-based planners. We design the reward for BP in a way that it 136 gets more reward if the car drives closer to the lane specified by the mission planner and closer to the 137 speed limit and it gets a large negative reward if it collides with an obstacle. 138

#### 139 3.3 Motion Planner

Motion planner's main task is to provide a safe and collision-free path towards its destination, while taking into account road boundaries, the vehicle dynamics, its maneuver capabilities in the presence of obstacles along with traffic rules or other constraints dictated by BP. In our design, motion planner generates a target trajectory defined as a set of 2D points (*i.e.* path) coupled with a target *speed* value. We propose a new level of road abstraction, through which each lane consists of  $N_c$  corridors, as depicted in Figure 2. Corridors are defined in the Frenét coordinate frame along center lines provided



planners (MoP).

Figure 1: Overview of our hierarchical planning Figure 2: Corridor abstraction for structured framework: Behavioral planner (BP) and Motion roads. Here Green lines separate corridors, Blue line is the corridor selected by MoP, and Yellow *line* is the generated trajectory corresponding to the selected corridor (Red car is the ego-car).

by mission planner, constructed based on road centers (waypoints) or path planning algorithms 146 for unstructred environments. As corridors are defined in the Frenét coordinate frame, tracking 147 performance remains invariant to transformation [18]. As such, motion planner can be considered 148 to be agnostic to lanes, intersections, and other road structures. As a result, there is no difference 149 between straight roads, curvy roads, and roundabout. This road abstraction reduces the the continuous 150 trajectory generation problem to the discrete corridor selection. 151

An MoP agent in our framework selects two sets of actions: 1) A lateral action identifying the target 152 corridor; and 2) A speed action which selects the target speed. Corridor selection is equivalent to 153 selecting a path among a set of predefined paths (clothoids or splines) along the mission planner 154 routes. An MoP agent chooses the best target corridor based on the behavioral planner decisions: 155 **BP** == Keep lane:  $MoP_{keep\_lane}$  should select a corridor within the current lane while avoiding any 156 collision. This enables MoP to maneuver around small objects in the lane without switching lane; and 157 **BP == Switch left/right**:  $MoP_{switch\_left/right}$  should select a corridor in the adjacent lanes. We 158 included the nearest corridor to the target lane to MoP choices to consider prepare-for-merge action. 159 Figure 2 illustrates corridor selection by an MoP agent. 160

MoP agents also select target speed based on BP actions and physical boundaries (e.g. heading cars, 161 or any interfering object on the road) while ensuring safety and smoothness. By selecting a target 162 speed, cruise in lane, emergency stop, vehicle following, and other BP-speed dependent actions can 163 be handled. By adjusting the target speed, MoP agents can avoid collision while keeping the lane 164 or switching to adjacent lanes. The target corridor and speed that are selected by MoP are relative 165 quantities. The absolute corridor and speed are calculated in the *Action post-processing* module in 166 Figure 1 and fed into the *Trajectory generator* module. Trajectory generator is a non-learning-based 167 module implemented using a simple curve fitting from point A to point B (Yellow line in Figure 2). 168 The generated trajectory is extended along the target corridor as the vehicle moves and is guaranteed 169 to be smooth and feasible for the controller node. 170

MoP must also be a function of the complete set of environment states because it must always be 171 safe no matter what the current BP decision is. MoP input states include: current lane and corridor, 172 current speed, speed and distance to heading cars for each corridor in current and neighboring lanes, 173 and speed and distance to cars in the back for each corridor in current and neighboring lanes. Similar 174 to BP reward function, MoP reward function also need to avoid using any *fine-grained* rules and 175 should remain *coarse-grained*. A *Keep lane* agent is responsible for adaptive cruise control (ACC). 176 Keep lane agent is responsible to remain in the current lane, maintain the distance to the heading 177 car and the car in the back (if it is too close to the ego-car), and maneuver around small obstacles 178 within the current lane. The reward for MoP<sub>switch\_left</sub> and MoP<sub>switch\_right</sub> was designed such that 179 the agent gets positive reward as it gets closer to the lane specified by BP. The reward function for 180  $MoP_{keep \ lane}$  is designed to encourage the agent to keep a *desired distance* from the heading cars 181



Figure 3: HRL vs. Flat agent. Top row: PPO agent. Bottom row: DQN agent.

if there exist any heading car or maintain safe front and back distance when the car in the back is
 too close to the vehicle. All *MoP* agents get penalized with a negative value if they collide with
 obstacles.

### 185 3.4 Training

In the proposed hierarchical RL framework, BP issues a high level command which is executed by 186 the corresponding MoP. As opposed to the other hierarchical frameworks (e.g. [5]), BP does not wait 187 until its command gets executed. Considering any fixed lifetime for BP commands is dangerous for 188 autonomous driving. In fact, BP should be able to update its earlier decisions (at every time step) 189 according to the new states. MoP is designed to prioritize safety over BP decisions. Our framework 190 is flexible when it comes to choosing RL algorithms to be applied for BP and MoPs. We tested our 191 framework with DQN [8], and PPO [14]. Potentially, we can train the HRL agent in two phases. In 192 the first phase, we initialize the exploration of BP and MoP agents to 1, and train BP coupled with a 193 rule-based MoP, and MoP coupled with a rule-base BP. This leads to pre-training BP and MoPs. In 194 the second phase, we may jointly train BP and MoPs. However, in this study, we only present the 195 results corresponding to the first phase. 196

## 197 **4 Simulator**

We developed a simulator, called SimpleTraffic, for cruising in multi-lane roads. SimpleTraffic has 198 been developed in Python and is compatible with OpenAI Gym environments. We will open source 199 SimpleTraffic for research community under GPL<sup>3</sup>. Concepts such as lane and corridor as well as 200 ego-car and traffic obstacles (e.g. other cars or motorbikes) have been implemented in the simulator 201 and can be customized for different scenarios. A user can set the number of lanes, number of corridors 202 per lane, maximum and minimum speed of the ego-car and traffic obstacles. Currently, two sets of 203 traffic obstacles have been implemented in *SimpleTraffic*: cars and motorbikes. In the default setting, 204 205 each car occupies three corridors and each motorbike occupies one corridor only. Having these two 206 sets of obstacles enables agent to distinguish between small objects that can be bypassed without changing the lane and larger objects that can only be bypassed by a lane change. 207

SimpleTraffic uses discrete time steps. Therefore, neither the speed of computers nor a slow archi-208 tecture (e.g. huge neural network) influences the result. This way, researchers can debug their RL 209 210 systems easier and more effectively. At every time step, the *SimpleTraffic* environment expects an action in the form of a tuple (corridor, speed) where corridor specifies the corridor index, Figure 2, 211 and *speed* specifies the target speed. After executing the action *SimpleTraffic* returns the next state, 212 reward, and a flag indicating whether a terminal state is reached. The simulator states include: total 213 number of lanes and corridors, current ego-car lane and corridor id, current ego-car speed, current 214 ego-car lateral position, target lateral position (specified by the previous action), distance to heading 215 cars and cars in the back at each corridor, and speed of heading cars and cars in the back at each 216 corridor. A trajectory generator module has been developed inside the simulator that generates a 217 smooth trajectory from the current position of the ego-car to the selected corridor. 218

<sup>&</sup>lt;sup>3</sup>To be released on http://www.github.com



Figure 4: Corridor selection of the transferred model in a new environment with Lincoln MKZ simulated dynamic model.

# 219 **5 Experiments**

We tested our framework with both Q-learning-based and policy-gradient-based techniques and chose popular algorithms from each of these categories, namely DQN and PPO, respectively.

SimpleTraffic: We trained both HRL and flat agents for 10 million time steps and deployed them 222 on SimpleTraffic simulator. During evaluation, each model was deployed on the simulator 10 times 223 for 1000 time steps and recording the average speed, the time to first collision, and the number of 224 collisions per 1000 steps. During our experiments, we noticed that some agents lean towards driving 225 too slowly to avoid any collision. This is not a desirable behavior, especially in a highway. On 226 the other hand, other agents drive more aggressively (i.e. too fast) to maximize the speed reward, 227 which may result in higher tendency to crash into other vehicles. In such case considering average 228 speed and time to first collision separately is not fully informative. However, combining them can be 229 more informative. Therefore, we define another metric we call *collision-free travel length* which is a 230 multiplication of time to first collision and average speed. Figure 3 compares the performance of the 231 HRL agent with the flat agent over time for PPO and DQN algorithms. These figures illustrate the 232 233 aforementioned metrics over time; from initial agent (not trained) to trained agent (after 10 million time steps). Figures 3(a) and 3(b) confirm that both HRL and flat agents have learned to maximize the 234 collision-free travel length, however, HRL agents managed to obtained better results by controlling 235 their speed and sacrificing speed reward in order to drive safer. We emphasize that the same rewards 236 were used for both HRL and flat agents. 237

**Transfer to real vehicles:** In section III we introduced our road abstraction and state representation for an RL agent that not only facilitates the learning procedure, but also helps to transfer the models to similar environment including real vehicle. As a proof of concept, we tested the transferability of the trained models by deploying them on a new ROS environment and using a Lincoln MKZ 2015 model as the vehicle dynamic model. Figure 4 illustrates few snapshots of the deployed model on the new environment. Our trained agents managed to successfully bypass static and dynamic obstacles by choosing appropriate corridors and cruise safely in multi-lane roads.

## **245 6 Conclusion & future work**

We proposed an RL-based hierarchical framework for autonomous multi-lane cruising. We introduced 246 a key intermediate abstraction within the motion planner to discretize the state-action space according 247 to high level behavioral decisions. Furthermore, we showed that the hierarchical design for an 248 autonomous vehicle system enables significantly better learning performance than a non-hierarchical 249 design, when using both a value-based method (DQN) and a policy optimization method (PPO). The 250 proposed framework allows for principled modular extension of motion planning, which is not the case 251 in rule-based or monolithic behavior cloning-based approaches. Moreover, we experimentally showed 252 that our state-action space abstraction allows transferring of the trained models from a simulated 253 environment with virtually no dynamics to the one with significantly more realistic dynamics without 254 a need for retraining. We also introduced our open source OpenAI Gym compatible simulator, called 255 SimpleTraffic, to the reinforcement learning research community. Although training BP and MoP 256 individually could sufficiently address the cruising in multi-lane problem, as our future work, we 257 aim to train the BP and MoP agents jointly (in an end-to-end fashion) to acquire higher level of 258 performance. 259

## 260 **References**

- [1] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon
  Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning
  for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [2] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [3] Yu Fan Chen, Miao Liu, Michael Everett, and Jonathan P How. Decentralized non communicating multiagent collision avoidance with deep reinforcement learning. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 285–292. IEEE, 2017.
- [4] Xuemin Hu, Long Chen, Bo Tang, Dongpu Cao, and Haibo He. Dynamic path planning for
  autonomous driving on various roads with avoidance of static and moving obstacles. *Mechanical Systems and Signal Processing*, 100:482–500, 2018.
- [5] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical
  deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In
  Advances in neural information processing systems, pages 3675–3683, 2016.
- [6] Richard Liaw, Sanjay Krishnan, Animesh Garg, Daniel Crankshaw, Joseph E Gonzalez, and
  Ken Goldberg. Composing meta-policies for autonomous driving using hierarchical deep
  reinforcement learning. *arXiv preprint arXiv:1711.01503*, 2017.
- [7] Branka Mirchevska, Manuel Blum, Lawrence Louis, Joschka Boedecker, and Moritz Werling.
  Reinforcement learning for autonomous maneuvering in highway scenarios.
- [8] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
  Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
  Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [9] Mustafa Mukadam, Akansel Cosgun, Alireza Nakhaei, and Kikuo Fujimura. Tactical decision
  making for lane changing with deep reinforcement learning. 2017.
- [10] Brian Paden, Michal Čáp, Sze Zheng Yong, Dmitry Yershov, and Emilio Frazzoli. A survey of
  motion planning and control techniques for self-driving urban vehicles. *IEEE Transactions on intelligent vehicles*, 1(1):33–55, 2016.
- [11] Chris Paxton, Vasumathi Raman, Gregory D Hager, and Marin Kobilarov. Combining neural
  networks and tree search for task and motion planning in challenging environments. *arXiv preprint arXiv:1703.07887*, 2017.
- [12] Etienne Perot, Maximilian Jaritz, Marin Toromanoff, and Raoul De Charette. End-to-end
  driving in a realistic racing game with deep reinforcement learning. In *International conference* on Computer Vision and Pattern Recognition-Workshop, 2017.
- [13] Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep rein forcement learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76,
  2017.
- [14] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
  policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [15] Shai Shalev-Shwartz, Nir Ben-Zrihem, Aviad Cohen, and Amnon Shashua. Long-term planning
  by short-term prediction. *arXiv preprint arXiv:1602.01580*, 2016.
- [16] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. On a formal model of safe and
  scalable self-driving cars. *arXiv preprint arXiv:1708.06374*, 2017.
- [17] Liting Sun, Cheng Peng, Wei Zhan, and Masayoshi Tomizuka. A fast integrated plan ning and control framework for autonomous driving via imitation learning. *arXiv preprint arXiv:1707.02515*, 2017.
- [18] Moritz Werling, Julius Ziegler, Sören Kammel, and Sebastian Thrun. Optimal trajectory
  generation for dynamic street scenarios in a Frenet frame. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 987–993. IEEE, 2010.
- [19] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali
  Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In
  *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3357–3364. IEEE,
- 312 2017.