

# Attention-based Hierarchical Deep Reinforcement Learning for Lane Change Behaviors in Autonomous Driving

Yilun Chen<sup>1</sup>, Chiyu Dong<sup>2</sup>, Praveen Palanisamy<sup>3</sup>, Priyantha Mudalige<sup>3</sup>,  
Katharina Muelling<sup>1</sup> and John M. Dolan<sup>1,2</sup>

**Abstract**—Performing safe and efficient lane changes is a crucial feature for creating fully autonomous vehicles. Recent advances have demonstrated successful lane following behavior using deep reinforcement learning, yet the interactions with other vehicles on-road for lane changes are rarely considered. In this paper, we design a hierarchical Deep Reinforcement Learning (DRL) algorithm to learn lane change behaviors in dense traffic. By breaking down overall behavior to sub-policies, faster and safer lane change actions can be learned. We also apply temporal and spatial attention to the DRL architecture, which helps the vehicle focus more on surrounding vehicles and leads to smoother lane change behavior. We conduct our experiments in the TORCS simulator and the results outperform the state-of-the-art deep reinforcement learning algorithm in various lane change scenarios.

## I. INTRODUCTION

Understanding and interacting with human-driven cars are essential for autonomous driving cars to behave socially in cooperative scenarios, such as lane changing, intersection traversing and ramp merging. Lane change is a fundamental behavior in on-road driving for overtaking or navigation purposes. It requires high-level reasoning about surrounding vehicles' intentions and behavior to form an effective driving strategy. At the same time, it requires low-level reasoning to plan what exact path to follow under the safety requirements, generally known as the path planning problem.

In this paper, we propose to use a deep reinforcement learning (DRL) method that can learn sub-policies for lane changing behavior. The structural overview of our algorithm is shown in Fig. 1. Instead of separately optimizing the high-level and low-level reasoning, our method combines them in one network with a hierarchical structure, which can still be trained in an end-to-end fashion. By designing a hierarchical action space, our network can maintain a high-level strategy and a low-level control command at the same time, while being differentiable end-to-end. This encourages shared computation and optimizes for the overall performance.

An attention mechanism is also applied to cooperative driving. Normally, human drivers do not pay equal attention to all available sensory information. People select the information based on a higher level of cognition for determining the current maneuver and ignore the unrelated

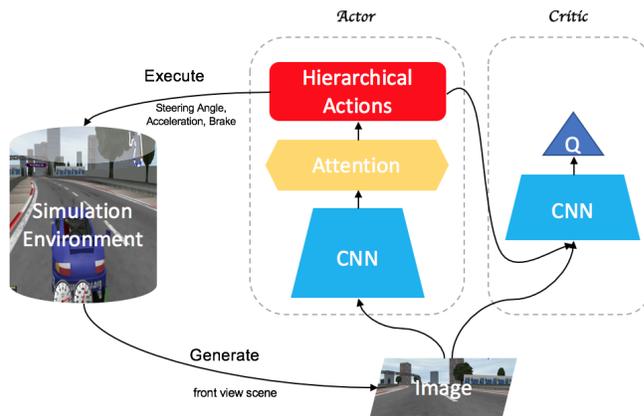


Fig. 1. Illustration of the algorithm. Our algorithm is based on deep reinforcement learning with actor and critic. The actor network is for policy learning and the critic network is for policy evaluation. We propose hierarchical actions and an attention mechanism to generate lane change behavior. We use the TORCS simulation environment.

information. Similarly, the attention mechanism in our deep neural network automatically focuses on feasible paths or surrounding on-road vehicles that may influence our driving behavior. As a result, it promotes the quality of the learned lane change behavior. The attention-enhanced network can make better use of perception information and lead to lower convergence time and better performance. It also leads to better interpretability of the network.

We summarize our contribution as follows. 1) Propose a hierarchical deep reinforcement learning algorithm that can learn lane change behaviors on road. Our method can be easily extended to learn multiple driving policies in one model. 2) Develop an attention mechanism that is suitable for learning driving policies with images. This helps improve compositionality of the network: learn better performance with fewer examples. 3) Discuss and analyze performance on lane change scenarios with comparison to state-of-the-art deep reinforcement learning algorithms.

## II. RELATED WORK

Traditional lane change maneuvers use rule-based designs without learning. In the DARPA challenge, the CMU Boss self-driving vehicle [1] made lane change decisions by checking empty slots. More recent methods apply sampling-based or optimization-based methods [2], [3], [4]. All these methods include intensive hand-engineering where we have to set up lots of hand-designed parameters and spend large amounts of time tuning them. We can avoid this by learning lane

<sup>1</sup> The Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA yilunc, jdolan@andrew.cmu.edu; kmuelling@nrec.ri.cmu.edu

<sup>2</sup> Department of Electrical and Computing Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA chiyud@andrew.cmu.edu

<sup>3</sup> General Motor Global R&D Center, Warren, MI 48090 USA

change behaviors with data. There are currently two main approaches to learning for self-driving vehicles: supervised / imitation learning and reinforcement learning.

One of the earliest supervised learning approaches was ALVINN [5], which used a neural network to directly map front-view camera images to steering angle. Recently, Nvidia [6], [7], [8] extended it with deep neural networks to demonstrate lane following with more complex real-world scenarios. Xu et al. [9] scaled this effort to a larger crowd-sourced dataset and proposed the FCN-LSTM architecture to derive a generic driving model. These works target the task of steering the vehicle to achieve lane-keeping behavior but do not consider interactive scenarios such as lane changes. Though Dong et al. [10] used a non-parametric regression to learn lane-change behaviors, it only predicts start and end points. [11], [12] do trajectory prediction with awareness of surrounding actors, but the prediction doesn't involve any interaction with the environment.

Supervised learning requires a high amount of human-labeled driving data, which makes it expensive and hard to scale in reality. Deep reinforcement learning avoids direct human-labeled driving data and has the potential to learn to recover from unseen or unsuccessful states, which is a well-known limitation in supervised learning or imitation learning. A number of attempts used deep reinforcement learning to learn driving policies: Shalev-Shwartz et al. [13] learned a safe multi-agent model for autonomous vehicles on the road and Jaritz et al. [14] learned a driving model for racing cars.

In this work, a deep reinforcement learning (DRL) approach with a novel hierarchical structure for lane changes is developed. Hierarchical actions for multi-scale tasks have been actively studied in DRL [15], [16]. Hausknecht et al. [17] developed a parameterized action space for the agent to learn to play soccer in the RoboCup competition. We take the idea of parameterized action space and design a hierarchical action structure suitable for lane change behavior.

The attention mechanism also contributes to the interactive behavior generation. Attention models are widely used in domains such as image captioning [18], machine translation [19] and visuomotor control [20]. Recently, Chen et al. [21] investigated attention in driving from a cognitive perspective and Kim et al. [22] interpreted driving behavior learning with attention. We adopt the idea of temporal and spatial attention to encourage lane change behaviors.

### III. METHOD

In this section, we first introduce the hierarchical action space designed for lane change behaviors and describe our overall approach to learning a driving model. We then dive into the attention mechanism that helps improve the learning both in efficiency and performance. Finally, we discuss the reward signal design for the DRL algorithm.

#### A. Hierarchical Action Space for Lane Change

The lane change behaviors in driving policies require high-level decisions (whether to make a lane change) and low-level planning (how to make the lane change). Based on

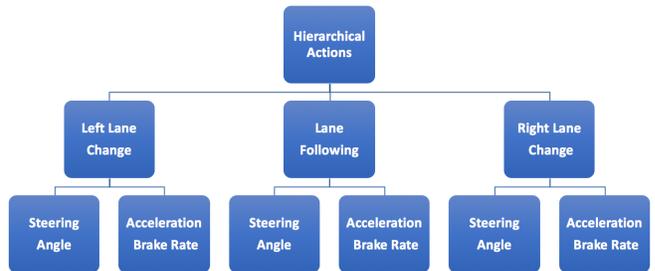


Fig. 2. Illustration of the hierarchical action space for lane change behavior.

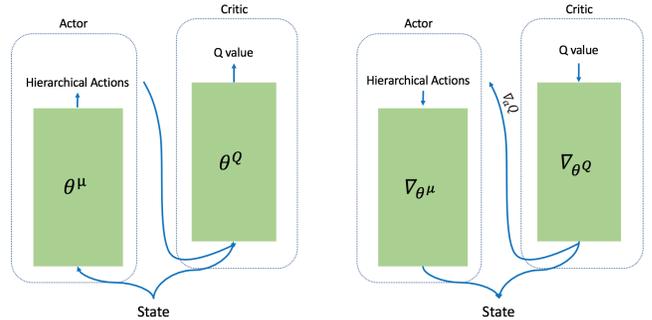


Fig. 3. The Actor-Critic architecture used in deep reinforcement learning, first introduced in [23]. On the left is the data flow during inference (forward pass) and on the right is the gradient flow during training (back-propagation).

the parameterized action space [17], we create a hierarchical action space for autonomous driving lane changes, as shown in Fig. 2. We define three mutually exclusive discrete high-level actions: Left Lane Change, Lane Following and Right Lane Change. At each time step, the agent must choose one of the three high-level actions to execute. Each action requires two continuous-valued parameters that must be specified. The first parameter is the steering angle. The second parameter is the acceleration brake rate (ABR). Formally, the high-level discrete actions are defined as  $A_d = \{a^{straight}, a^{left}, a^{right}\}$ . Each discrete action  $a \in A_d$  contains a set of continuous parameters  $P_a = \{p_1^a, \dots, p_n^a\}$ . The overall hierarchical action space is defined as

$$A = \{straight, p_{angle}^{straight}, p_{ABR}^{straight}\} \\ \cup \{left, p_{angle}^{left}, p_{ABR}^{left}\} \\ \cup \{right, p_{angle}^{right}, p_{ABR}^{right}\}.$$

The three sets of actions represent three different types of driving behaviors. During training, the system will learn to choose from the three high-level action decisions and apply proper parameters specific to that action.

#### B. Actor-critic Based DRL Architecture

We develop our algorithm based on Deep Deterministic Policy Gradient (DDPG) [23], a deep reinforcement learning algorithm for continuous control. For better training stability, the actor-critic architecture with two neural networks is used. This architecture decouples the action evaluation and selection process into two separate deep neural networks: actor-network and critic-network, as shown in Fig. 3. The actor-network  $\mu$ , parameterized by  $\theta^\mu$ , takes as input state  $s$

and outputs hierarchical action  $a$  along with its parameter  $p_a$ . The critic-network  $Q$ , parameterized by  $\theta^Q$ , takes as input state  $s$  and hierarchical action  $a$  along with its parameter  $p_a$  and outputs a scalar Q-Value  $Q(s, a)$ .

The hierarchical action  $a$  is represented as a vector of the probability of action choices  $p$  and the parameters  $P_a$  coupled to each discrete action. The discrete high-level action is chosen to be the output with maximum value among the choices of action probabilities. Then it is coupled with the corresponding parameters from the parameter outputs. Though the parameters of all actions are outputted, only the parameters of the chosen action are used. In this way, the actor-network simultaneously outputs which discrete action to execute and how to choose parameters for that action. The critic-network receives as input all the values of the output layer in the actor. We do not indicate which exact action is applied for execution or which parameters are associated with which action during training. In the back-propagation stage, the critic-network only provides gradients for the selected action and the corresponding parameters. This ensures that we update the policy only in the direction where we explore.

### C. Recurrence in DRL

The driving task often features incomplete and noisy perception information due to the partial observability from sensors. Given only one frame of input, the autonomous driver cannot collect sufficient environment information such as velocities of surrounding vehicles to generate the right driving behavior. This makes driving a Partially-Observable Markov Decision Process (POMDP).

To deal with partial observability, recurrence is introduced to better estimate the underlying environment state. An additional recurrent neural network module (RNN) with dropout [24] is added to the output of the convolutional neural network in the actor network in DDPG. The RNN module processes the temporal information in the network instead of just having stacked historical observations as input. In practice, LSTM is used as the basic RNN structure, with the implementation based on Zaremba et al. [25]. We refer to this recurrent version of the DDPG algorithm as Deep Recurrent Deterministic Policy Gradient (DRDPG).

Comparing to the original DDPG, DRDPG offers several advantages, including the ability to handle longer input sequences, exploration of temporal dependencies and better performance in case of partially observable experiences.

### D. Attention Mechanism

We introduce an "Attention Mechanism" in our deep reinforcement learning algorithm, as shown in Fig. 4. There are two streams of attention: Temporal Attention and Spatial Attention. The attention module is added to the actor network for better action selection. The critic network is set to be a fully convolutional network for action evaluation.

1) *Temporal Attention*: We add temporal attention over the output of the LSTM layer in the DRDPG model, as shown in Fig. 4. The temporal attention mechanism decides

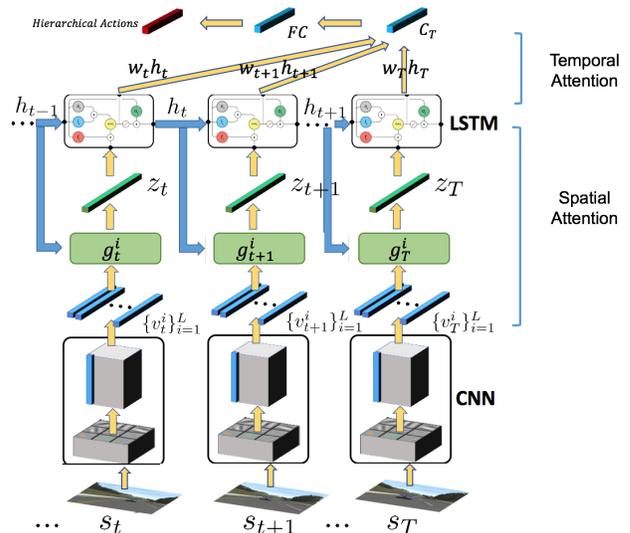


Fig. 4. The architecture of the actor network in the Attention Deep Recurrent Deterministic Policy Gradient algorithm. We consider two kinds of attention: Temporal Attention and Spatial Attention. Temporal Attention learns to weight the importance of previous frames, while Spatial Attention learns the importance of different locations in the image.

which frames matter most in past observations by learning scalar weights for LSTM outputs at different time steps. The weight of each LSTM output  $w_i$  is defined as an inner product of the feature vector  $v_i$  and LSTM hidden vector  $h_i$ , followed by a softmax function to normalize the sum of weights to 1. By this definition, each learned weight is dependent on the previous time step's information and current state information.

$$w_{T+1-i} = \text{Softmax}(v_{T+1-i} \cdot h_{T+1-i}) \quad i = 1, \dots, T$$

Then we compute the combined context vector  $C_T$ . The context vector  $C_T$  is a weighted sum of LSTM outputs through  $T$  time steps.

$$C_T = \sum_{i=1}^T (w_{T+1-i} h_{T+1-i})$$

The derived context vector  $C_T$  is passed by a fully connected layer  $FC$  before calculating the hierarchical actions of the actor network. The learned weights  $\{w_{T+1-i}\}_{i=1}^T$  here can be interpreted as the importance of the LSTM output at a given frame. Therefore, the optimizing process can be seen as learning to choose which observations are relatively more important for learning the correct actions.

Temporal attention works in the sense that it explicitly considers the past  $T$  frames' LSTM output features for computing the action output, while this information is only passed implicitly by LSTM. By increasing the value of  $T$ , the model can consider a longer sequence of history frames and thus can make a better action choice.

2) *Spatial Attention*: Human drivers pay special attention to certain objects when performing specific tasks. For example, drivers care more about the direction of the road when traversing curvy roads and pay more attention to surrounding

vehicles when making lane changes. This intuition is encoded in the network by adding importance to different locations in the image, which we call spatial attention.

With spatial attention, the network learns weights for different areas in an image. The context feature used by LSTM is a weighted-sum combination of spatial features multiplied by the learned weights. According to how the combination is modeled, previous spatial attention models are divided into hard attention and soft attention [18]. Inspired by [26], we use a soft version of attention, which means learning a deterministic weighted context in the system.

The spatial attention, as shown in Fig. 4, occurs after convolution layers and before recurrent layers. At time step  $t$ , suppose the convolutional layers produce a set of  $d$  feature maps with size  $m \times n$ . These feature maps can also be seen as a set of region vectors with length  $d : \{v_t^i\}_{i=1}^L, v_t^i \in \mathbb{R}^D, L = m \times n$ . Each region vector corresponds to the features extracted by the CNN in a different image region. In the soft attention mechanism, the context vector  $z_t$  is represented by a weighted sum of all-region vectors  $\{v_t^i\}_{i=1}^L$ .

$$z_t = \sum_{i=1}^L g_t^i \cdot v_t^i$$

The weights in this sum are chosen in proportion to the importance of this vector (i.e., the extracted feature in this image region), which is learned by the attention network  $g$ . The attention network  $g_t^i$  has region vector  $v_t^i$  and hidden state  $h_{t-1}$  produced by the LSTM layer as input and outputs the corresponding importance weight for the region vector  $v_t^i$ . The attention network  $g_t^i$  here is represented as a fully connected layer followed by a softmax function:

$$g_t^i = \text{Softmax}(w_v \cdot v_t^i + w_h \cdot h_{t-1}) \quad i = 1, \dots, T$$

The context vector  $z_t$  is fed into the LSTM layer. The output of the LSTM layer is concatenated with action vector  $A_T$  and then used to compute the actions.

The attention network can be interpreted as a mask over the CNN feature maps, where it reweights the region features to get the most informative features for computing the actions. Thus, the spatial attention has learned to acquire the ability to select and focus on the more important regions when selecting the action.

### E. Reward Function Design

Table I shows the six components that constitute the reward function.  $r_1, r_2$  encourage the car to follow the lane.  $r_3$  encourages the car to keep speed.  $r_5, r_6$  ensures driving safety without driving-off-road and collision. Specifically,  $r_4$  encourages the car to overtake if the front vehicle is within a distance of 100 meters. Here  $x$  means the distance to the front vehicle in the same lane. If no vehicle is found, then  $x$  has a default value of 100 meters.

The final reward function is a linear combination of terms in Table I with assigned weights  $w$ :  $R = \sum_{i=1}^6 w_i r_i$ . Finding the optimal weights is computationally heavy and beyond the scope of this paper. Here we empirically set the weights to

TABLE I

TERMS IN THE REWARD FUNCTION. WHERE  $\theta$  IS THE YAW ANGLE W.R.T. THE ROAD DIRECTION;  $d$  IS THE OFFSET TO THE LANE CENTER;  $v$  IS THE SPEED;  $x$  IS THE DISTANCE TO THE LEADING VEHICLE.

| Rewards Term         | Reward Functions  |
|----------------------|---|
| Road Alignment       | $r_1 = \cos \theta - \sin \theta$   |
| Lane Centering       | $r_2 = - d $  |
| Prefer Speed Limit   | $r_4 = \frac{1}{25} \begin{cases} v & v \leq 25 \text{ m/s} \\ 70 - v & v > 25 \text{ m/s} \end{cases}$ |
| Encourage Overtake   | $r_5 = -\frac{1}{100} \max(0, 100 - x)$   |
| Out-of-Bound Penalty | $r_3 = -\mathbf{1}\{OutofBoundary\}$  |
| Avoid Collision      | $r_6 = -\mathbf{1}\{Collision\}$  |

$w_{1,2,3,4} = 1, w_{5,6} = 10$ . This set of parameters can enforce driving safety with a high penalty on driving off road and collision.

## IV. EXPERIMENTS

The algorithm is evaluated in the open source car simulation environment TORCS. To ensure generality, five tracks and ten types of cars are used to generate our evaluation sets. Each track has two to three lanes with various illumination of the sky. The tracks have different surrounding environments and road geometry. For example, in the Corkscrew scenario, the road has more curves and fewer straight sections, and there is a sharp high-speed corner. To successfully traverse the tracks, the car has to learn various skills like U-turn, hill climbing, overtaking and throttling before a sharp turn. All these challenges are posed to the training agent with no human supervision.

We add traffic cars in different locations on every track per trial of the training to ensure the complexity of the traffic. The traffic cars' behaviors are controlled by the internal AI of the TORCS environment, which is designed to mimic human behaviors when driving. We add diversity to the traffic cars' behavior by changing the internal parameters like average speed, acceleration limit and chances of lane changes in the car simulation AI. This will change the car's driving style, e.g., being more aggressive when turning. In this way, we would like to mimic the real traffic patterns on the road. To encourage the overtaking and lane changing behavior of our vehicle, we set the speed limit of traffic cars to 20 m/s and set the speed limit of our ego vehicle to 25 m/s. An example of the environment used is shown in Fig. 5.

Front-view image is used as the perception sensor input, at a rate of 5Hz. The images are then down-sampled to  $320 \times 240$ . Other measurements like car direction, speed direction, distance to lane boundaries, etc. are collected along with the images. For each model tested, we train the agent on all five tracks in TORCS for a fixed 5-hour period. This allows the agent to collect around 100,000 frames for each scenario. To avoid overfitting in one track and add diversity in training, we randomized the sequence of tracks encountered in training. For each model, the training starts an episode by randomly selecting one of the five tracks. Any off-road or collision will stop and restart a new episode.



Fig. 5. One example of the five tracks used for training. From left to right: the map of the example track *Street-1*, image top view when starting a new episode, a screenshot of the front view camera during training.

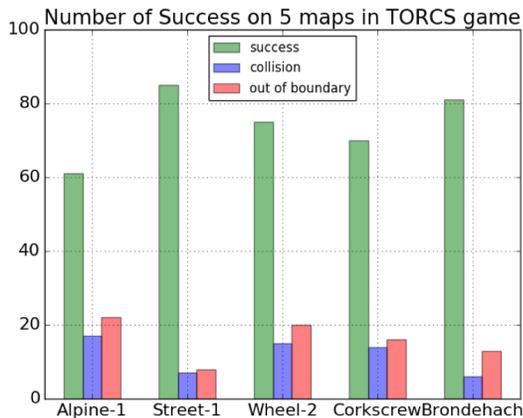


Fig. 6. Final DRDPG model with hierarchical actions, Spatial and Temporal attention tested on different trials in TORCS game. We obtain the result of each map by running 100 episodes.

We use CNN networks for both the actor and critic networks. The feature extraction from images is done by a standard AlexNet [27] with the first five layers. We train the LSTM layers sequentially on the video sets from simulation with an unrolling size of 8 frames during training. Specifically, the LSTM we use has 64 hidden units. Stochastic gradient descent (SGD) is used for the optimization, with an initial learning rate of  $10^{-3}$ , momentum of 0.99 and a batch size of 10. The learning rate is decayed by half whenever training loss plateaus. We run our training in parallel on 4 Nvidia Titan X GPUs.

## V. EVALUATION

We first evaluate the success rate of our final model with both temporal and spatial attention on all five tracks in TORCS. We calculate the success rate by testing the agent’s ability to drive along the entire track and see if it can successfully finish a loop. As shown in Fig. 6, the model shows overall good performance in all five trials with a success rate between 60 and 80 percent to complete a runway circle on the map. The failure cases mainly come from collision or out-of-bounds cases, about half of each. One example of out-of-boundary comes from the case in which a high-speed car fails to turn at a sharp curve since the agent does not learn to decelerate ahead of time. Another case occurs when a front vehicle blocks the view of a turn, and the agent doesn’t have enough time to react to a turn

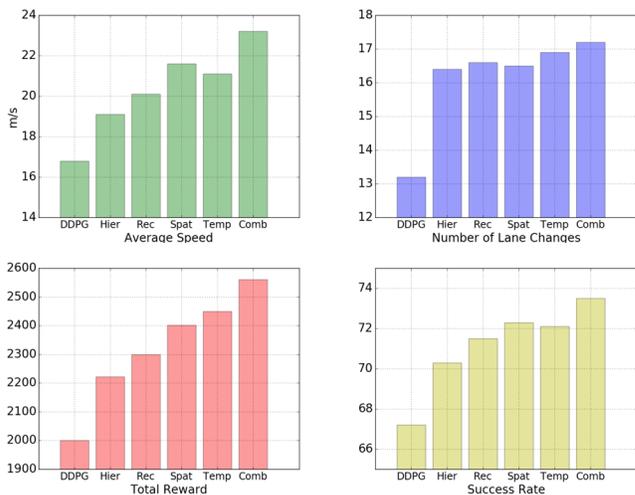


Fig. 7. A comparison of the performance improvements with respect to the baseline (DDPG), obtained by introducing the hierarchical actions (Hier), recurrence (Rec), spatial attention (Spat) and temporal attention (Temp) from left to right. The final combined model (Comb) applies all of above. The performance evaluation is based on average speed, the number of lane changes, total reward during an episode and percentage of successful episodes.

that suddenly appears. For the collision issue, most collisions happen when trying to overtake a vehicle or being forced to do a lane change. Some of these cases are difficult: occlusions can arise when entering a corner or the front vehicle suddenly slows down due to the road geometry.

We further investigate the performance of our method compared to the original DDPG baseline, as shown in Fig. 7. All models are evaluated based on average values over 100 episodes after 15 hours of training. We first compare the version which has hierarchical action space (the “hier” configuration in Fig. 7) with the original DDPG algorithm. We observe an obvious boost of 2 m/s in average speed and 3 more lane change behaviors after we apply our hierarchical action space. The results show that the original DDPG algorithm with direct steering and brake rate output tends to drive more conservatively and stays behind other vehicles more often without trying to overtake. In comparison, our model with hierarchical actions is designed to learn separate policies for left lane change and right lane change, so it tends to act more aggressively and maintain higher speed to make lane changes whenever possible. Although the method encourages overtake behavior, it still guarantees safety at the same time, as the success rate of finishing a loop increases by 3 percent. This is due to the fact that more overtake behaviors are able to avoid getting too close to or colliding with the leading vehicle. Our method achieves higher rewards and fewer collisions with other vehicles, thus resulting in a higher success rate for a complete episode compared to the original DDPG algorithm.

We next show the improvements of adding recurrence (“rec”), temporal attention (“temp”) and spatial attention (“spat”) to the DDPG algorithm with hierarchical actions in Fig. 7. We find that adding recurrence can actually speed up the training curve and achieve higher rewards

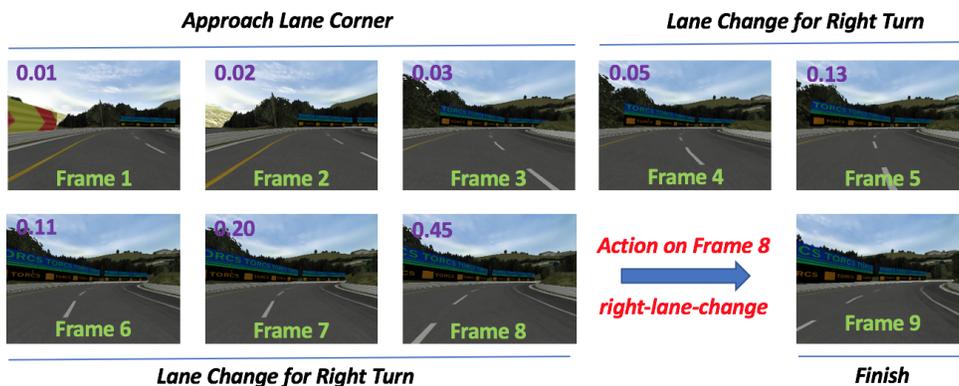


Fig. 8. A turning scenario to illustrate Temporal Attention. In this scenario, the agent has learned to do a right lane change for a more efficient right turn. The number at the top left of each image is the weight assigned to that image frame for temporal attention (higher weight indicates more importance).

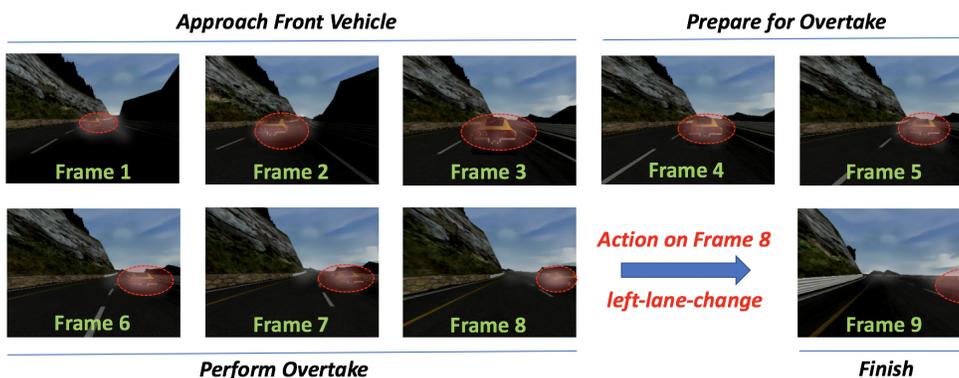


Fig. 9. An overtaking scenario to illustrate Spatial Attention. A mask over the input image is learned by the Spatial Attention. Brighter colors indicate higher weights assigned to that region. The weights are smoothed with a Gaussian kernel for visualization.

and performance as evaluated by the success rate. It also results in higher average speed due to smoother and more stable driving behavior. For the attention mechanism, we can observe that both temporal and spatial attention improve the metrics compared to the recurrence model. This can be attributed to the better utilization of the images as input and finally results in safer and more stable lane change behavior. With the final combined model (“comb”), we see an increase in the average speed from 16.7 to 23.2 m/s and the number of lane changes increase from 13.2 to 17.2 while also improving the success rate from 67.2% to 73.5% compared to the original DDPG model.

## VI. VISUALIZATION

Fig. 8 and 9 give two examples to illustrate the effect of the attention mechanism in turning and overtaking scenarios.

In the turning scenario (Fig. 8), the agent has learned to change lane in order to perform a more efficient turn at a sharp right curve. At the top-left corner in each figure, different weights  $w_i$  in temporal attention are assigned to the input images. Notice that although the weights generally increase through time, there are exceptions where the agent thinks a certain situation needs special attention. For example, in frame 5, the car is just crossing the lane marker on the road, which could be more dangerous and need greater attention. In this case, a single last frame helps less than a sequence of frames. By looking back into the past at the weighted sum

of features of several frames before, the agent can determine the best trajectory for turning at the corner.

In the overtaking scenario (Fig. 9), the regions of the front vehicle are highlighted by the Spatial Attention. The attention model gives large weights to where the neighboring vehicle appears, which helps the agent figure out when and how to launch the proper lane change behavior. The agent has learned to approach the front vehicle, slow down to keep distance, perform a lane change and then speed up to overtake the front vehicle. By learning this attention over the input image, the agent can extract the relevant context for the task and learn the proper behavior more efficiently.

## VII. CONCLUSION

We introduce an attention-based hierarchical deep reinforcement learning algorithm for learning lane change behaviors in dense traffic with an end-to-end trainable architecture. The proposed hierarchical action space for learning driving behaviors can generate sub-policies for lane change behaviors in addition to the lane following behavior. This model simplifies the work of deliberately designing a sophisticated lane change maneuver and introduces a data-oriented approach that can learn the lane change policy through trial and error. We also investigate how an attention mechanism can help in the task of driving with deep reinforcement learning. The attention model also helps to explain what is learned in the driving model.

## REFERENCES

- [1] C. R. Baker and J. M. Dolan, "Traffic interaction in the urban challenge: Putting boss on its best behavior," in *Intelligent Robots and Systems, 2008. IROS 2008. IEEE/RSJ International Conference on*. IEEE, 2008, pp. 1752–1758.
- [2] S.-H. Lee and S.-W. Seo, "A learning-based framework for handling dilemmas in urban automated driving," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 1436–1442.
- [3] S. Manzinger, M. Leibold, and M. Althoff, "Driving strategy selection for cooperative vehicles using maneuver templates," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017, pp. 647–654.
- [4] C. Liu, W. Zhan, and M. Tomizuka, "Speed profile planning in dynamic environments via temporal optimization," in *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017, pp. 154–159.
- [5] D. A. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," in *Advances in neural information processing systems*, 1989, pp. 305–313.
- [6] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [7] Z. Chen and X. Huang, "End-to-end learning for lane keeping of self-driving cars," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017, pp. 1856–1860.
- [8] V. Rausch, A. Hansen, E. Solowjow, C. Liu, E. Kreuzer, and J. K. Hedrick, "Learning a deep neural net policy for end-to-end control of autonomous vehicles," in *American Control Conference (ACC), 2017*. IEEE, 2017, pp. 4914–4919.
- [9] H. Xu, Y. Gao, F. Yu, and T. Darrell, "End-to-end learning of driving models from large-scale video datasets," *arXiv preprint*, 2017.
- [10] C. Dong, Y. Zhang, and J. M. Dolan, "Lane-change social behavior generator for autonomous driving car by non-parametric regression in reproducing kernel hilbert space," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2017*, pp. 4489–4494.
- [11] N. Deo and M. M. Trivedi, "Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1179–1184.
- [12] X. Huang, S. McGill, B. C. Williams, L. Fletcher, and G. Rosman, "Uncertainty-aware driver trajectory prediction at urban intersections," *arXiv preprint arXiv:1901.05105*, 2019.
- [13] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "Safe, multi-agent, reinforcement learning for autonomous driving," *arXiv preprint arXiv:1610.03295*, 2016.
- [14] M. Jaritz, R. de Charette, M. Toromanoff, E. Perot, and F. Nashashibi, "End-to-end race driving with deep reinforcement learning," *arXiv preprint arXiv:1807.02371*, 2018.
- [15] T. D. Kulkarni, K. Narasimhan, A. Saedi, and J. Tenenbaum, "Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation," in *Advances in neural information processing systems*, 2016, pp. 3675–3683.
- [16] Z. Yang, K. Merrick, L. Jin, and H. A. Abbass, "Hierarchical deep reinforcement learning for continuous action control," *IEEE Transactions on Neural Networks and Learning Systems*, 2018.
- [17] M. Hausknecht and P. Stone, "Deep reinforcement learning in parameterized action space," *arXiv preprint arXiv:1511.04143*, 2015.
- [18] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International Conference on Machine Learning*, 2015, pp. 2048–2057.
- [19] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," *arXiv preprint arXiv:1508.04025*, 2015.
- [20] R. Zhang, Z. Liu, L. Zhang, J. A. Whritner, K. S. Muller, M. M. Hayhoe, and D. H. Ballard, "Agil: Learning attention from human for visuomotor tasks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 663–679.
- [21] S. Chen, S. Zhang, J. Shang, B. Chen, and N. Zheng, "Brain-inspired cognitive model with attention for self-driving cars," *IEEE Transactions on Cognitive and Developmental Systems*, 2017.
- [22] J. Kim and J. Canny, "Interpretable learning for self-driving cars by visualizing causal attention," in *Int. Conf. Comput. Vis.(ICCV)*, 2017, pp. 2961–2969.
- [23] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [24] V. Pham, T. Bluche, C. Kermorvant, and J. Louradour, "Dropout improves recurrent neural networks for handwriting recognition," in *2014 14th International Conference on Frontiers in Handwriting Recognition*. IEEE, 2014, pp. 285–290.
- [25] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.
- [26] I. Sorokin, A. Seleznev, M. Pavlov, A. Fedorov, and A. Ignateva, "Deep attention recurrent q-network," *arXiv preprint arXiv:1512.01693*, 2015.
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.