

Instance-Aware Predictive Navigation in Multi-Agent Environments

Jinkun Cao¹, Xin Wang², Trevor Darrell² and Fisher Yu³

¹ Carnegie Mellon University ² University of California, Berkeley ³ ETH Zurich

jinkunc@andrew.cmu.edu, {xinw, trevor}@eecs.berkeley.edu, i@yf.io

Abstract—In this work, we aim to achieve efficient end-to-end learning of driving policies in dynamic multi-agent environments. Predicting and anticipating future events at the object level are critical for making informed driving decisions. We propose an Instance-Aware Predictive Control (IPC) approach, which forecasts interactions between agents as well as future scene structures. We adopt a novel multi-instance event prediction module to estimate the possible interaction among agents in the ego-centric view, conditioned on the selected action sequence of the ego-vehicle. To decide the action at each step, we seek the action sequence that can lead to safe future states based on the prediction module outputs by repeatedly sampling likely action sequences. We design a sequential action sampling strategy to better leverage predicted states on both scene-level and instance-level. Our method establishes a new state of the art in the challenging CARLA multi-agent driving simulation environments without expert demonstration, giving better explainability and sample efficiency.

I. INTRODUCTION

Over the past years, machine vision and decision-making systems have approached or surpassed human-level performance in vision and robotics applications due to the emergence of deep learning methods [16]. However, learning drive autonomously remains one of the most desirable but notoriously challenging problems, with high requirements on reliability, explainability, and data efficiency [26], [23], [11].

One approach to learning a driving policy is imitation learning [18], [32], [17], [8], where the policy is fully supervised by expert demonstrations. However, it may be unreliable when facing scenarios outside demonstrations as expert demonstrations provide only examples of safe driving instead of recovering from mistakes. Besides, demonstrations are hard to collect, even in some simulation environments. On the other hand, reinforcement learning-based agents can explore the world and learn from mistakes without expert demonstrations [22], [1]. But they suffer severely from sample inefficiency as they rely on sparse reward signals from the environment and require orders of magnitude more data than humans [31]. It brings the utmost challenges to drive in a multi-agent environment, where the scene changes are more complex and even adversarial to the ego-vehicle.

We aim to have reliable driving policies by exploring multi-agent environments, combining the deep representation of imitation with the reinforcement from exploratory experiences. To achieve reliability, we need to predict the consequences of the selected actions. At the same time, we want to learn from experiences to improve data efficiency. Learning from good experience allows self-imitation and eliminates the necessity of expert demonstrations.

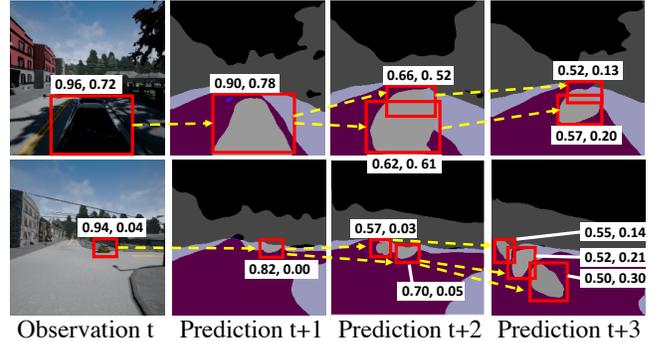


Fig. 1: Multi-instance event prediction in ego-centric view. It predicts both possible instance locations and the probability of instance-level events. The predicted instance locations on a frame sequence build a coarse estimation of possible trajectories. For each possible location, the model outputs confidence of localization and chance of collision, (c, ϕ) .

We propose Instance-Aware Predictive Control (IPC), a decision-making method with future event prediction at both scene and instance levels. IPC forecasts the motions of the other agents in the environment as well as the scene structure changes. A key component is a *multi-instance event prediction* module. It estimates the possible locations of the other agents in the dynamic environments within a finite horizon, conditioned on the past and present observations and the selected action sequence of the controlled ego-agent. The likely locations of the other agents indicate possible future events, such as collisions. Our agent makes action decisions by considering sampled actions’ consequences. The predicted visual structure helps to explain the agent’s decision.

To determine the action, IPC seeks an optimal action sequence within a finite horizon by repeatedly sampling possible action sequences. The prediction module can tell which sequence leads to safe future states, and the agent executes the action for the first step from the sequence. We design a *sequential action selection* strategy, selecting actions to optimize the chances of safe scene-level and instance-level events. It improves the efficiency of planning with noisy prediction. Our experiments show that the strategy is helpful to select action sequences more efficiently.

IPC establishes a new state of the art in the challenging CARLA simulator and the more realistic GTA V environment. IPC improves the sample efficiency and policy performance without any expert demonstration. Moreover, the instance-level visual prediction enhances the explainability of action decisions. Our method predicts future safety-

related events (e.g., collision, offroad, offlane) with 95%+ accuracy. The codebase is released at <https://github.com/SysCV/spc2>.

II. RELATED WORK

Our work connects to the rich literature in reinforcement learning and imitation learning. In this section, we focus on related works within the scope of autonomous driving.

Forecasting and control. The state forecasting is usually handled as a Markov Decision Process (MDP), and state forecasting conditional to certain actions could be used in control tasks. Though it is still hard to learn policy directly from pixels, deep models have brought powerful representations for visual tasks. Recurrent networks [29], [5], [2], [14] have shown the power to predict future visual states based on raw visual inputs. However, in multi-agent environments, the uncertainty from other agents makes it harder to infer the underlying distribution of state transition. Explicit multi-hypothesis forecasting [27] focuses on this problem but requires full instance trajectories. Instead, we use cues of multi-modal future state possibility from only ego-centric monocular observations. Our method refers to Model Predictive Control (MPC) [4] to control the ego-vehicle by comparing the forecast future states conditional to some selected actions. Pan et al. [23] also uses the way to learn driving policy but only in single-agent environments. We are the first to use instance awareness to achieve an efficient driving policy in a complicated multi-agent environment with an explainable decision.

Driving by imitation. Imitation learning focuses on imitating policy pattern from expert demonstrations. Previous methods [3], [6], [18], [32], [17] have shown efforts in autonomous driving tasks. However, the demonstrations focus on only correct behaviors, providing no chance for the agent to learn from mistakes. It can lead to dangerous situations as the agent would not know how to recover from accidents. Besides, demonstrations are usually very hard and expensive to collect in the real world, and manually designed demonstrations in simulators usually suffer from a lack of diversity. So imitation learning is also studied together with exploring in experience [8] to relieve its imperfectness. Our method also uses imitation to boost performance, but it uses imperfect pseudo-demonstration from previous “good experience”, which is called self-imitation.

Policy learning from experiences. One way to eliminate the necessity of demonstrations is to learn policy from the agent’s experiences. Reinforcement learning typically follows this philosophy. Recently, model-based reinforcement learning has been gaining popularity, leveraging deep networks to learn a state transition function [22], [1] for action decisions. Though some methods [25], [24] has tried to learn policies directly from observations without demonstration, they usually work well only in over-simplified simulation environments [21], [13]. For driving tasks, the existing reinforcement learning method in the realistic simulator, e.g., CARLA [7], still requires imitating demonstrations [18]. Be-

sides, compared with visual prediction, deep models still suffer from severe low efficiency when planning action directly from visual signals. Aware of these flaws, our method does not directly output action decisions from visual input but by comparing the consequences of sampled action candidates, which better leverages the capacity of deep models.

III. INSTANCE-AWARE PREDICTIVE CONTROL

In this section, we start with an introduction to the problem setup (III-A) and an overview of the proposed Instance-Aware Predictive Control (IPC) framework (III-B). We then dive into its two core modules: Multi-instance Event Prediction (MEP) for forecasting (III-C) and Sequential Action Sampling (SAS) for action planning (III-D).

A. Problem Setup

Following previous works [23], [30], we formulate the environment transition under Markov Decision Process. With the state space \mathcal{S} and the action space \mathcal{A} , the transition is $M : \mathcal{S} \times \mathcal{A}^N \rightarrow \mathcal{S}$ in an N-agent environment,

$$\mathbf{s}_{t+1} = M(\mathbf{s}_t, \mathbf{A}_t), \quad (1)$$

where $\mathbf{A}_t = (\mathbf{a}_t^0, \dots, \mathbf{a}_t^N)$ are the actions executed by N agents at timestamp t , \mathbf{a}_t^0 is the action of the ego-vehicle and \mathbf{s} is the partial observation of the world.

We only optimize the ego-vehicle’s driving policy π^0 , taking into account the environment and other agents’ behaviors. In contrast, other agents adopt a pre-defined policy from the simulator that is unknown to the ego-vehicle. The reward function R in the MDP rewards the ego-vehicle to drive fast and safely in the environment. That is, the ego-vehicle needs to drive as fast as possible without crashing into other vehicles or obstacles or getting off the drivable area.

B. Instance-aware Predictive Control Framework

IPC is a decision-making method with future event predictions at both scene and instance level. Similar to the classic model predictive control, IPC includes a dynamic model to predict future states and a sampling-based strategy to select the action with respect to the state prediction.

Predictive control. To determine the ego-vehicle’s action at timestamp t , we adopt a predictive control approach that samples multiple action sequences spanning on timestamps $t : t + T$ and forecasts future states conditioned on the sampled action sequences. IPC inherits the merits of the model predictive control framework that implements the action of the current timestamp t while keeping future timestamps in account. In IPC, we train a dynamic model M^* that predicts future states with the observation of the current state and action sequence. Once the dynamic model is learned, we can sample the action sequences conditioned on the future state forecasting and adopt a cost function C to optimize the action sampling. That is,

$$\mathbf{a}_{t:t+T}^0 = \arg \min_{\mathbf{a}_{t:t+T}^0} \sum_{\tau=t}^{t+T} C(M^*(\mathbf{s}_\tau, \mathbf{A}_\tau)), \quad (2)$$

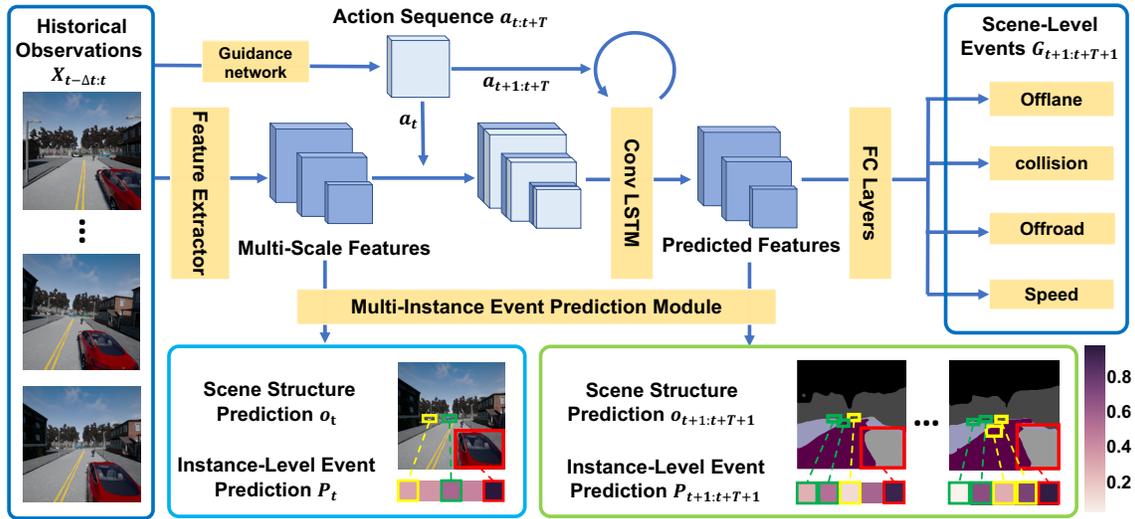


Fig. 2: Instance-aware predictive control (IPC) framework. Given historical observations, a guidance network helps to sample action sequences in action space. The model predicts both future visual structure and the chance of some events. \mathbf{o} is the visual observation containing semantic segmentation and instance locations. G is scene-level events. P are instance-level events on each predicted possible instance location. Event prediction brings reference to action selection. Visual structure prediction brings explanation to the action decision. The bottom-right colorbar indicates the probability of instance-level events.

where T is the length of horizon for state forecasting. To better learn the dynamic model, we adopt an episodic learning strategy to train the prediction module where the ground truth only exposes one possibility of the future in one episode. The model observes multiple possible future outcomes caused by different behaviors of other vehicles in the environment by exploring multiple episodes.

State forecasting. A key component in IPC is its future events prediction module at both scene and instance level. As shown in Fig. 2, the future prediction module first extracts feature representations from a sequence of historical observations, which encode past dynamics of other agents. The visual features are then concatenated with sampled actions and fed into a convolutional LSTM [33] network to predict features on future T timestamps recurrently. Once the features are obtained, IPC utilizes a multi-instance event prediction (MEP) module to estimate future visual structure \mathbf{o} with semantic segmentation and possible locations of other vehicles through object detection. MEP has a RetinaNet [20]-like network to estimate the chances of instance-level events between the ego-vehicle and other vehicles (e.g., collision) denoted as P . IPC also includes a fully connected network to predict scene-level events denoted as G . The overall state in the dynamic model is represented by a triplet

$$\mathbf{s} = (\mathbf{o}, G, P), \quad (3)$$

where \mathbf{o} provides semantic visual information and G and P are safety-related events or ego-vehicle speed. To note that, although the ego-vehicle is unaware of the policies of other vehicles,

the observation \mathbf{o} reflects states of other agents, which influences the action selection of the ego-vehicle.

State forecasting in a multi-agent environment is challenging due to its dynamic nature, where various behaviors from other agents can lead to uncertain future predictions. To address this issue, we adopt a multiple hypotheses forecasting strategy as shown in Fig. 3, where MEP explicitly takes into account multiple possible configurations of the future states and provides uncertainty scores to assist action selection.

Action selection. We develop an action selection method to ensure the ego-vehicle drives within drivable areas and avoids collisions in the environment. According to Eq. 2 and 3, the action selection is conditioned on both the scene-level and the instance-level events which bear more uncertainty from non-deterministic behaviors of other agents. To determine the action sequence efficiently, we design a two-stage sequential action sampling (SAS) strategy to select action as illustrated in Fig. 4. SAS assesses sampled action candidates in the first stage and eliminates the unlikely action sequences based on the scene-level predictions only. Then it uses noisy instance-level predictions to determine the final action sequence.

C. Multi-Instance Event Prediction

A core design of the state forecasting in IPC is the multi-instance event prediction, which explicitly models the state changes in the multi-agent environment. To achieve it, the multi-instance event prediction (MEP) module utilizes visual perception tasks such as semantic segmentation to capture the scene structure changes and object detection to model inter-agent interaction.

Event prediction in a multi-agent environment is challenging. In addition to recognizing the road structure and following the correct lanes, our agent will also avoid collisions with the other agents. Although we are testing our

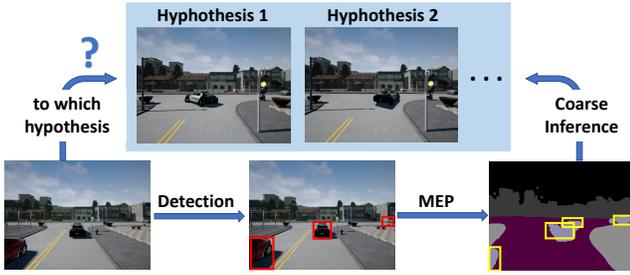


Fig. 3: How the prediction of possible instance locations in the Multi-instance Event Prediction (MEP) builds a coarse inference for different hypotheses.

systems mainly in the simulation environment, we assume we can only control our vehicle, and the policies of those agents are unknown. Therefore, we do not know their exact driving trajectories and speed, and their driving policies are indeed programmed with some degree of randomness. The assumption and challenge are also realistic. When we drive in the real world, we do not know exactly what is in the mind of other drivers, riders, and pedestrians.

Our key insight is that certain motion patterns, despite their stochastic nature, will emerge from a large number of behavioral observations of the agents in the environments. For instance, a car will unlikely stop suddenly without the influence of the scene structures or interference of the other agents, and a pedestrian will likely walk straight on the sidewalks. Also, the object locations across time should be continuous. Therefore, our prediction module, MEP, recognizes the existences of the multiple agents and attempts to learn their motion patterns with deep neural networks. In this work, we represent the instances and their motion with 2D bounding boxes.

In addition, MEP models the behavioral randomness of the future through multiple hypotheses, similar to using modes to represent a distribution. A car may want to speed up or make a turn, leading to uncertainty in the future states. MEP represents the uncertainty with likely future positions of the other agents shown in Fig. 3. Also, MEP provides the uncertainty scores of those possible future locations.

Semantic segmentation can provide dense supervision to the driving policy learning in addition to the sparse rewards from the environment, which improves the sample efficiency and the explainability of action decisions, as shown in the prior work [23]. However, the original design of SPC only handles static scene structures with one possible future state hypothesis. It also fails to capture the dynamic scene changes in a multi-agent environment, where instance-level events (e.g., collisions between vehicles) dominate the scene changes. Therefore, IPC uses MEP to predict extra instance-level events based on future state and scene structure estimation. The dense supervision for semantic segmentation and object detection provides richer training signals and improves the sample efficiency and explainability.

Training. To train MEP, we adopt an episodic training strategy where episodes are sampled from a replay buffer,

and the loss is defined over a finite horizon of T . That is,

$$\mathcal{L}_D = \sum_{t=1}^T (w_t \mathcal{L}_d + \mathcal{L}_g), \quad (4)$$

where \mathcal{L}_d is the focal loss [20] for location estimation and \mathcal{L}_g is the cross-entropy loss for semantic segmentation. w_t is a weighting factor defined as

$$w_t = \begin{cases} w_1, & n_{t-1} < n_k \text{ or } t = 1 \\ w_2, & n_{t-1} \geq n_k \end{cases}, \quad (5)$$

where n_t is the number of ground truth instances at timestamp t . We use $w_1 = 2$, $w_2 = 1$ in our experiments to emphasize the changes among the adjacent frames.

MEP predicts the inter-agent events along with the visual structure as well as their uncertainty scores to enable multiple hypotheses forecasting. For example, MEP can predict the possible locations of one vehicle together with the chance of ego-vehicle colliding into it at the spot. That is, it predicts $(c, x_1, y_1, x_2, y_2, \phi)$ for each possible location, where c is the prediction confidence, (x_1, y_1, x_2, y_2) indicates the bounding box and ϕ indicates the chance of collision. If the location is very close to the ego-vehicle or in the direction the ego-vehicle is turning to, the chance of collision would be high.

D. Sequential Action Sampling

IPC performs sampling-based action selection based on the event prediction from the learned dynamic model. Similar to model predictive control, IPC selects the action sequence from candidates that minimizes the cost of future events at each step and the ego-vehicle executes the action on the first step in the selected sequence. We propose a two-stage sequential action sampling (SAS) strategy to select the desired action sequence shown in Fig. 4. The final action is determined conditioned on the instance-level event predictions in a pruned action space by the first stage that relies on scene-level event predictions.

Guidance network. We first discretize the action space into multiple bins and then uniformly sample from the selected bins by a guidance network to avoid sampling action directly from continuous space. At each step t , N_0 action sequence candidates with a horizon of T are sampled, and the dimension of the action space is V . We set $V = 2$ for steering and throttle in our experiments, and both values are float between -1 and 1. For steering, -1 means turning left-most, and 1 means turning right-most. For throttle, -1 and 1 means to slow down or speed up in the full capacity.

We adopt a guidance network to output the probability of the discrete bins based on historical observations. It is trained through imitating good episodes from the replay buffer. With the imitation, the model can learn from the current trajectory and the historical experiences to improve the sample efficiency.

Future events definitions. With the historical observations and sampled actions, the dynamic model predicts the future states. For each action sequence, the sequential action sampling (SAS) strategy uses the predictions of both scene-level events G and instance-level events P . In the CARLA

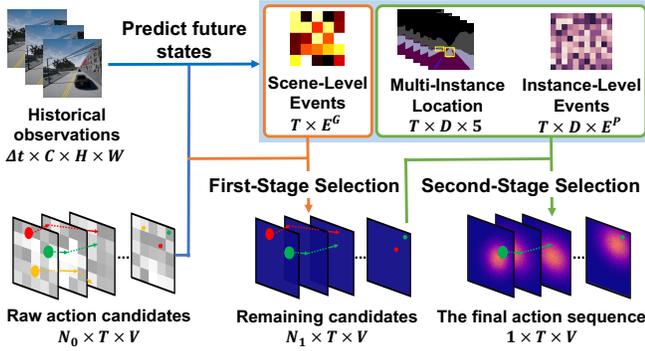


Fig. 4: The procedure of SAS. From the input historical observations and sampled actions, the dynamic model predicts E^G future scene-level events and E^P instance-level events. From the originally sampled N_0 action sequences, the first stage estimates a cost referring to scene-level event chance to keep N_1 sequences. Then considering instance-level events, the second stage chooses the final action sequence.

environment, we use four scene-level events off-road, off-lane, collision and speed. We thus denote the prediction of G as $G^* = (p_{\text{off-road}}, p_{\text{off-lane}}, p_{\text{collision}}, p_{\text{speed}})$. We only consider one instance-level event, inter-agent collision, at D possible vehicle locations with confidence score above a threshold. The event is represented by $(c, x_1, y_1, x_2, y_2, \phi)$, where c is the prediction confidence, (x_1, y_1, x_2, y_2) indicates the bounding box and ϕ indicates the chance of collision.

Two-stage action sampling. With the predicted future states, SAS selects actions to drive safely and fast. The first stage focuses on scene-level events with a cost function defined as

$$C_{\text{scene}}(\mathbf{s}_t, \mathbf{a}_{t:t+T}) = \sum_{i=1}^{T+1} \alpha_i C_g(G_{t+i}^*), \quad (6)$$

where $\alpha_i = \max(\frac{1}{2^i}, \frac{1}{8})$ is the weighting factor to emphasize more on the near future prediction while not neglecting further steps. We consider three types of accidents (i.e., off-road, off-lane and collision), so the cost function C_g is

$$C_g(G^*) = C_p(p_{\text{off-road}}) + C_p(p_{\text{collision}}) + 0.2C_p(p_{\text{off-lane}}), \quad (7)$$

$$C_p(x) = -p_{\text{speed}} \cdot (1 - x) + v_m \cdot x, \quad (8)$$

where v_m is the allowed maximum speed. This cost function encourages higher speed and lower accident chance. We initially sample N_0 action sequences. Only the top N_1 action sequences with the minimal cost are considered in the second stage. We set $N_0 = 20$ and $N_1 = 5$ in experiments.

The second stage focuses on instance-level events to avoid collision between the ego-vehicle and other vehicles. The cost function is then defined as

$$C_{\text{instance}}(\mathbf{s}_t, \mathbf{a}_{t:t+T}) = \sum_{i=1}^{T+1} \alpha_i \sum_{d=1}^D \phi_{t+i,d} \cdot c_{t+i,d}, \quad (9)$$

where α_i is the same weighting factor as in Eq. 6.

IV. EXPERIMENTS

We mainly conduct experiments in CARLA simulator. The results show IPC outperforms a range of baselines (IV-B). We also conduct ablation studies on the novel components, MEP and SAS, and find that they can both improve performance (IV-D). We compare the scene-level event prediction quantitatively with the closest baseline, SPC [23] (IV-C). Besides, we find IPC is flexible to incorporate expert demonstrations to boost model performance (IV-E). Finally, we show potentials to use IPC in more realistic environments (IV-F).

A. Experimental Setup

Simulation environment. We adopt the multi-agent driving environments in CARLA [7], which provide multiple built-in maps with diverse scene structures and road configurations. The ego-vehicles uses a monocular RGB camera to get visual observations. The frame rate of the simulator is 10 frames per second. CARLA supports extracting semantic segmentation and vehicle locations for supervision. We also evaluate our method on GTA V, which has more realistic rendering but provides no visual ground truth, so we train a Mask R-CNN [10] model on a collected dataset [28] to provide pseudo-supervision.

Model configurations. We use Deep Layer Aggregation [35] to extract visual features and a single layer ConvLSTM network [33] to synthesize features on future steps. MEP uses a RetinaNet [20]-style head to predict instance locations and an up-sampling layer for semantic segmentation.

Training. We use Adam [15] for training with an initial learning rate 5×10^{-4} . The exploration rate decreases linearly from 1 to 0 until 100k steps. IPC forecasts on future 10 frames as in SPC [23]. The history buffer size is 20k, and the batch size is 16. One episode terminates when the ego-vehicle (1) collides for 20 times or (2) gets off-road or stuck for 30 steps continuously, or (3) has run 1000 steps.

Evaluation. We adopt the reward function in SPC [23] as

$$R(\mathbf{s}_t) = \frac{\text{speed}}{15} - (\mathbb{1}_{\text{off-road}} + 2 \times \mathbb{1}_{\text{collision}} + 0.2 \times \mathbb{1}_{\text{off-lane}}), \quad (10)$$

where $\mathbb{1}_{\text{off-road}}$, $\mathbb{1}_{\text{collision}}$ and $\mathbb{1}_{\text{off-lane}}$ indicate (binary) if off-road, collision or off-lane happens on that step.

B. Results in CARLA

We compare IPC with both model-free (DDPG [19] and SAC [9]) and model-based methods (GCG [12] and SPC [23]). DDPG and SAC output action signals directly from visual observations. SPC and GCG do predictive control similar to IPC, but SPC uses visual signal while GCG does not. We also use a depth map as an extra supervision signal to study if it can boost the learning rewards. Our comparison also includes CIRL [18], which uses imitation learning to pretrain model weights and finetunes with DDPG. An expert autopilot agent is deployed to give a reward reference for accident-free driving. As shown in Fig. 5, IPC obtains higher rewards than baselines on two maps. More agents in the environments make it harder to drive, so all the methods' rewards deteriorate even for the expert agent.

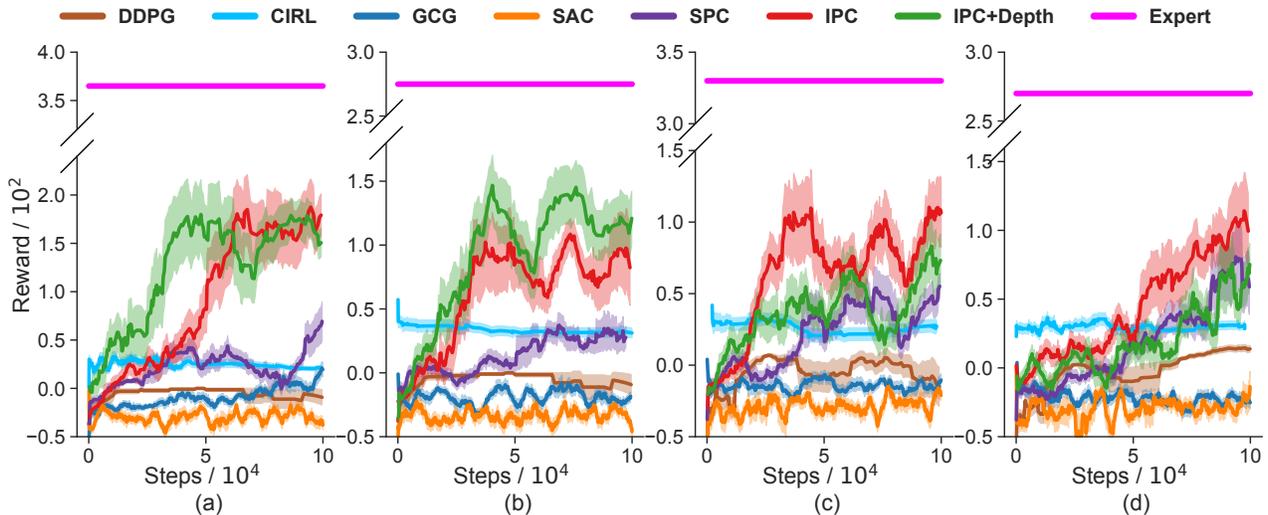


Fig. 5: The reward curves of methods tested under different environment settings in CARLA. (a): `town01` w/ 32 vehicles in; (b): `town01` w/ 64 vehicles in; (c): `town02` map with 32 vehicles in; (d): `town02` w/ 64 vehicles in. In all cases, our method outperforms all other methods. The depth information helps marginally in (a) and (b).

TABLE I: The prediction accuracy of scene-level events by SPC [23] and our method. The evaluation is conducted with 32 vehicles spawned in the map. Compared with SPC, our method shows more accurate prediction on all scene-level events.

Horizon		1	2	3	4	5	6	7	8	9	10
collision	SPC	96.12	96.05	95.74	95.58	95.41	95.29	95.03	94.80	94.67	94.37
	IPC	97.34	97.28	97.17	97.07	96.86	96.71	96.52	96.31	96.07	95.81
offroad	SPC	96.87	96.74	96.56	96.30	96.07	95.82	95.53	95.24	94.92	94.61
	IPC	97.47	97.46	97.29	97.12	96.91	96.76	96.54	96.34	96.08	95.72
offlane	SPC	94.79	94.89	94.91	94.63	94.42	94.05	93.33	92.82	92.32	91.51
	IPC	95.19	95.18	95.13	94.82	94.54	94.28	93.94	93.53	93.14	92.61
coll w/ veh	SPC	95.10	94.82	94.66	94.37	94.18	94.09	93.93	93.82	93.75	93.62
	IPC	96.02	96.00	95.97	95.70	95.59	95.37	95.26	95.19	95.02	94.96

TABLE II: Semantic segmentation prediction accuracy in single-agent (SA) and multi-agent (MA) environments with 32 vehicles spawned. Though the prediction is accurate in single-agent environments, it becomes more unreliable in multi-agent environments with more uncertainties. And the semantic prediction for future vehicles is extremely inaccurate. It shows semantic prediction is not enough to drive in multi-agent environments.

Horizon		1	2	3	4	5	6	7	8	9	10
SA	pixel acc.	96.50	95.45	95.39	95.25	95.19	95.14	95.06	95.02	94.91	94.70
	mean acc.	95.95	95.75	95.60	95.51	95.43	95.27	95.15	94.97	94.76	94.43
	mean IU	93.45	92.76	92.50	92.32	92.21	92.07	91.94	91.77	91.58	91.39
	f.w. IU	94.38	94.11	93.88	93.67	93.49	93.30	93.17	93.02	92.86	92.65
MA	pixel acc.	94.52	93.27	92.54	91.52	91.30	90.72	90.38	90.07	89.90	89.53
	vehicle acc.	48.67	42.90	42.62	41.83	41.04	39.36	39.01	38.76	38.50	38.27
	mean acc.	80.21	73.20	71.12	70.41	68.50	67.36	67.07	66.83	66.07	65.50
	mean IU	71.96	66.42	64.60	62.97	62.01	61.60	60.99	60.32	59.54	59.02
	f.w. IU	90.56	88.08	86.81	84.92	84.60	84.01	83.76	83.21	82.77	81.74

Despite the challenges, our model still performs well in 64-agent environments. We find IPC is more sample efficient.

C. IPC v.s. SPC

We examine the differences between IPC and SPC on CARLA. We compare their scene-level event prediction accuracy in Table I. We find that IPC outperforms SPC on

all four scene-level event predictions within a horizon of 10 steps. We hypothesize that the scene structure prediction using semantic segmentation in multi-agent environments maybe not reliable anymore. We, therefore, test the semantic segmentation prediction accuracy of IPC in single-agent and multi-agent environments in Table II. It shows that the accuracy drops substantially in multi-agent environments,

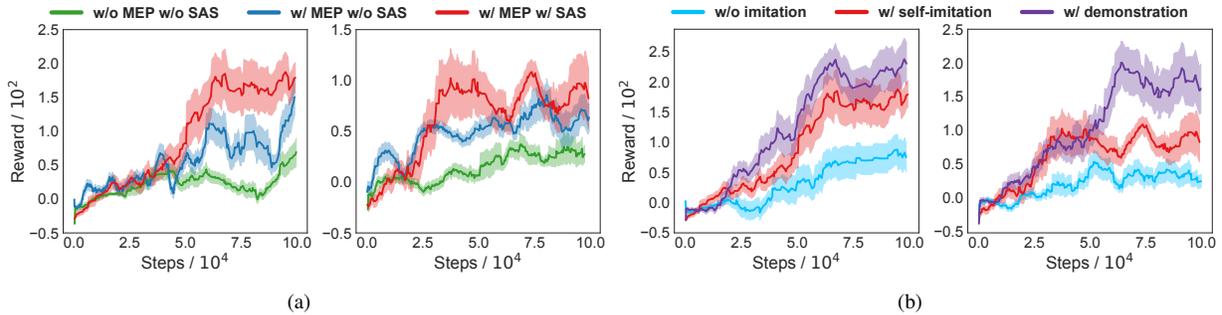


Fig. 6: (a). Ablation study of MEP and SAS. It is tested on `town01` map with 32 vehicles (left) or 64 vehicles (right) spawned. The results show that MEP improves the data efficiency a lot and SAS makes better action decisions. (b). How imitation boosts the model performance with 32 vehicles (left) or 64 vehicles (right) spawned. The results show that demonstration and historical experience both boost the model performance but demonstration helps more than the pattern from good experience.

TABLE III: The statistical results of driving on `Twon01`. Frequency is calculated on every 100 steps.

	avg speed	collision	coll w/ veh	offroad
IPC-32	10.22	1.23	1.07	1.12
SPC-32	8.79	3.02	2.79	2.23
IPC-64	9.32	2.53	2.23	1.45
SPC-64	8.03	4.12	3.98	2.38

indicating the necessity of using instance-level cues. We also compare raw driving metrics for SPC and IPC in Table III which agree that IPC makes the ego-vehicle drive faster and more safely in multi-agent environments. Accidents cause a risk for the ego-vehicle to get into situations it would not have predicted. So by better avoiding accidents, IPC forecasts future states more accurately.

D. Ablation study of MEP and SAS

MEP estimates instance-level events related to other agents, providing cues for more reliable control. SAS selects action with better tolerance of forecasting noise. We conduct an ablation study for the two modules in Fig. 6(a). As the figure shows, the reward increase is slower when MEP is not used. This indicates MEP can improve the sample efficiency. Moreover, without SAS, we merge two stages into one. It shows that the performance is lower than SAS. It indicates that SAS helps to select more reliable actions.

E. Boost Performance by Imitation

So far, we have shown that IPC achieves competitive results in multi-agent environments without any expert demonstration. Yet, IPC is flexible to incorporate expert demonstrations. Imitation can be used in training the guidance network to sample actions. We utilize expert demonstrations from the autopilot agent in CARLA for imitation. We also have self-imitation learning to imitate actions taken in historical good episodes, in which reward is at the top 10% and higher than 200. Results are shown in Fig. 6(b). It shows the expert demonstration improves sample efficiency not substantially

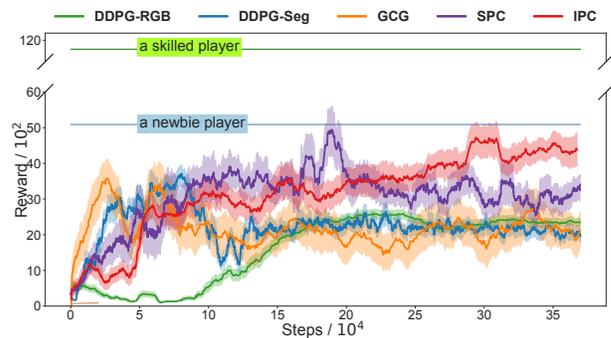


Fig. 7: Results on GTA V. IPC outperforms other methods.

compared to the self-imitation when there are 32 vehicles. But the model benefits much more from demonstrations when there are 64 other vehicles. With the ability to learn from imitations, we can collect demonstrations from either simulators or the real world and boost model performance offline. Joint optimizing the visual prediction and action decision making forms a scalable path to learn more robust action patterns among different environments.

F. Evaluation on more realistic benchmarks

Though CARLA is designed for evaluating driving algorithms and systems, its image rendering realism still has space for improvement compared to the latest rendering technology. Therefore, we also evaluate our method on GTA V, a video game with more realistic rendering. We run DDPG, GCG, SPC, and IPC in GTA V for the comparison. As a reference of human performance, we also show the driving rewards of a skilled player, who has rich experience in playing GTA V, and a newbie player, who has not played GTA V. The reference rewards are their average performance with 10 trials. As shown in Fig. 7, IPC outperforms all other baselines. Moreover, both IPC and SPC can have a close performance with the newbie player, but IPC is more stable.

It is challenging to adapt the model learned from the simulator to real-world datasets. But we still try to directly deploy the model trained on CARLA to the real-world

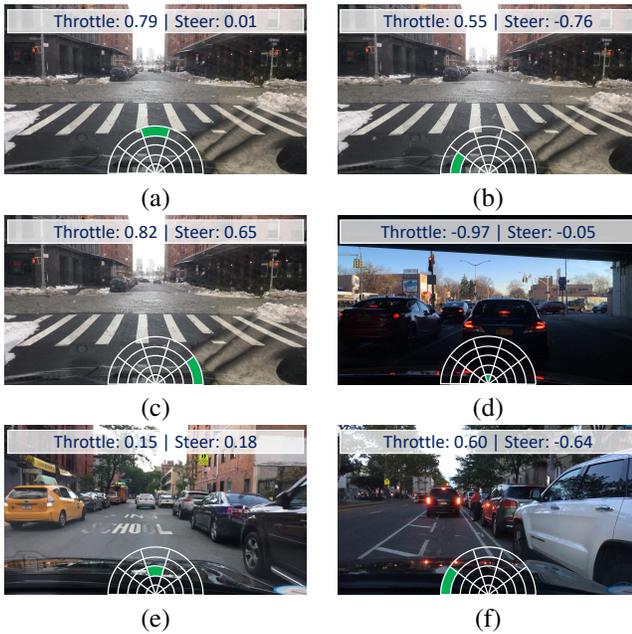


Fig. 8: The policy from simulator shows reasonable action even in real-world scenes. In (a)-(c), it is confident to drive to three directions in a crossing. In (d), it slows down to wait for the vehicle ahead. In (e), it follows straight on the road. In (f), it drives towards left to avoid the vehicle on the right.

driving video dataset BDD100K [34]. It generates reasonable actions in some scenes, as shown in Fig. 8, but can fail a lot. Our study provides some basis for a better sim-to-real policy transfer in the future.

V. CONCLUSION

We presented an instance-aware predictive control (IPC) approach for driving policy learning in multi-agent environments. IPC includes a novel multi-instance event prediction (MEP) module to forecasts future events at both scene and instance levels. MEP predicts the scene structure changes through semantic segmentation and the possible locations of other vehicles. To efficiently select actions conditioned on the multi-level event prediction, we introduced a two-stage sequential action sampling strategy to improve control from noisy forecasting. Experiments on realistic driving simulation environments showed that our method built new state-of-the-art under different multi-agent settings with improved data efficiency and policy reliability.

REFERENCES

- [1] Agrawal, P., Nair, A.V., Abbeel, P., Malik, J., Levine, S.: Learning to poke by poking: Experiential learning of intuitive physics. In: NIPS (2016)
- [2] Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., Savarese, S.: Social lstm: Human trajectory prediction in crowded spaces. In: CVPR. pp. 961–971 (2016)
- [3] Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., Jackel, L.D., Monfort, M., Muller, U., Zhang, J., et al.: End to end learning for self-driving cars. arXiv:1604.07316 (2016)
- [4] Camacho, E.F., Alba, C.B.: Model predictive control. Springer Science & Business Media (2013)
- [5] Chiu, H.k., Adeli, E., Nibbles, J.C.: Segmenting the future. IEEE RA-L 5(3), 4202–4209 (2020)

- [6] Codevilla, F., Miiller, M., López, A., Koltun, V., Dosovitskiy, A.: End-to-end driving via conditional imitation learning. In: ICRA (2018)
- [7] Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., Koltun, V.: CARLA: An open urban driving simulator. In: Proceedings of the 1st Annual Conference on Robot Learning. pp. 1–16 (2017)
- [8] Gao, Y., Lin, J., Yu, F., Levine, S., Darrell, T.: Reinforcement learning from imperfect demonstrations. arXiv:1802.05313 (2018)
- [9] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In: ICML. pp. 1856–1865 (2018)
- [10] He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: ICCV. pp. 2961–2969 (2017)
- [11] Huval, B., Wang, T., Tandon, S., Kiske, J., Song, W., Pazhayampallil, J., Andriulka, M., Rajpurkar, P., Migimatsu, T., Cheng-Yue, R., et al.: An empirical evaluation of deep learning on highway driving. arXiv preprint arXiv:1504.01716 (2015)
- [12] Kahn, G., Villafior, A., Ding, B., Abbeel, P., Levine, S.: Self-supervised deep reinforcement learning with generalized computation graphs for robot navigation. In: ICRA. pp. 1–8. IEEE (2018)
- [13] Kempka, M., Wydmuch, M., Runc, G., Toczek, J., Jaśkowski, W.: Vizdoom: A doom-based ai research platform for visual reinforcement learning. In: CIG. pp. 1–8 (2016)
- [14] Kim, C., Li, F., Rehg, J.M.: Multi-object tracking with neural gating using bilinear lstm. In: ECCV. pp. 200–215 (2018)
- [15] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- [16] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NIPS (2012)
- [17] Kuderer, M., Gulati, S., Burgard, W.: Learning driving styles for autonomous vehicles from demonstration. In: ICRA (2015)
- [18] Liang, X., Wang, T., Yang, L., Xing, E.: Cirl: Controllable imitative reinforcement learning for vision-based self-driving. In: ECCV (2018)
- [19] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971 (2015)
- [20] Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV. pp. 2980–2988 (2017)
- [21] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.: Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602 (2013)
- [22] Nagabandi, A., Kahn, G., Fearing, R.S., Levine, S.: Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In: ICRA (2018)
- [23] Pan, X., Chen, X., Cai, Q., Canny, J., Yu, F.: Semantic predictive control for explainable and efficient policy learning. In: ICRA. pp. 3203–3209 (May 2019). <https://doi.org/10.1109/ICRA.2019.8794437>
- [24] Pan, X., Seita, D., Gao, Y., Canny, J.: Risk averse robust adversarial reinforcement learning. In: ICRA. pp. 8522–8528. IEEE (2019)
- [25] Pan, X., You, Y., Wang, Z., Lu, C.: Virtual to real reinforcement learning for autonomous driving. arXiv:1704.03952 (2017)
- [26] Pomerleau, D.A.: Alvin: An autonomous land vehicle in a neural network. In: NIPS. pp. 305–313 (1989)
- [27] Rhinehart, N., McAllister, R., Kitani, K., Levine, S.: Precog: Prediction conditioned on goals in visual multi-agent settings. In: ICCV (2019)
- [28] Richter, S.R., Vineet, V., Roth, S., Koltun, V.: Playing for data: Ground truth from computer games. In: ECCV. pp. 102–118. Springer (2016)
- [29] Rochan, M., et al.: Future semantic segmentation with convolutional lstm. arXiv preprint arXiv:1807.07946 (2018)
- [30] Shim, D.H., Kim, H.J., Sastry, S.: Decentralized nonlinear model predictive control of multiple flying robots. In: IEEE International Conference on Decision and Control. pp. 3621–3626 (2003)
- [31] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al.: Mastering the game of go without human knowledge. Nature p. 354 (2017)
- [32] Wang, D., Devin, C., Cai, Q.Z., Yu, F., Darrell, T.: Deep object-centric policies for autonomous driving. In: ICRA. pp. 8853–8859 (2019)
- [33] Xingjian, S., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.K., Woo, W.c.: Convolutional lstm network: A machine learning approach for precipitation nowcasting. In: NIPS. pp. 802–810 (2015)
- [34] Yu, F., Chen, H., Wang, X., Xian, W., Chen, Y., Liu, F., Madhavan, V., Darrell, T.: Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In: CVPR. pp. 2636–2645 (2020)
- [35] Yu, F., Wang, D., Shelhamer, E., Darrell, T.: Deep layer aggregation. In: CVPR. pp. 2403–2412 (2018)