

Self-Adaptive Driving in Nonstationary Environments through Conjectural Online Lookahead Adaptation

Tao Li,^{1*} Haozhe Lei,^{1*} and Quanyan Zhu¹

Abstract—Powered by deep representation learning, reinforcement learning (RL) provides an end-to-end learning framework capable of solving self-driving (SD) tasks without manual designs. However, time-varying nonstationary environments cause proficient but specialized RL policies to fail at execution time. For example, an RL-based SD policy trained under sunny days does not generalize well to rainy weather. Even though meta learning enables the RL agent to adapt to new tasks/environments, its offline operation fails to equip the agent with online adaptation ability when facing nonstationary environments. This work proposes an online meta reinforcement learning algorithm based on the *conjectural online lookahead adaptation* (COLA). COLA determines the online adaptation at every step by maximizing the agent’s conjecture of the future performance in a lookahead horizon. Experimental results demonstrate that under dynamically changing weather and lighting conditions, the COLA-based self-adaptive driving outperforms the baseline policies in terms of online adaptability. A demo video, source code, and appendixes are available at <https://github.com/Panshark/COLA>

I. INTRODUCTION

Recent breakthroughs from machine learning [1]–[3] have spurred wide interest and explorations in learning-based self driving (SD) [4]. Among all the endeavors, end-to-end reinforcement learning [5] has attracted particular attention. Unlike modularized approaches where perception, localization, decision-making, and motion control are handled by different modules, end-to-end learning approaches aim to output a synthesized driving policy from raw sensor data.

However, the limited generalization ability prevents RL from wide application in real SD systems. To obtain a satisfying driving policy, RL methods such as Q-learning and its variants [2], [6] or policy-based ones [7], [8] require an offline training process. Training is performed in advance in a stationary environment, producing a policy that can be used to make decisions at execution time in the same environments as seen during training. However, the assumption that the agent interacts with the same stationary environment as training time is often violated in practical problems. Unexpected perturbations from the nonstationary environments pose a great challenge to existing RL approaches, as the trained policy does not generalize to new environments [9].

To elaborate on the limited generalization issue, we consider the vision-based lane-keeping task under changing weather conditions shown in Figure 2. The agent needs to provide automatic driving controls (e.g., throttling, braking,

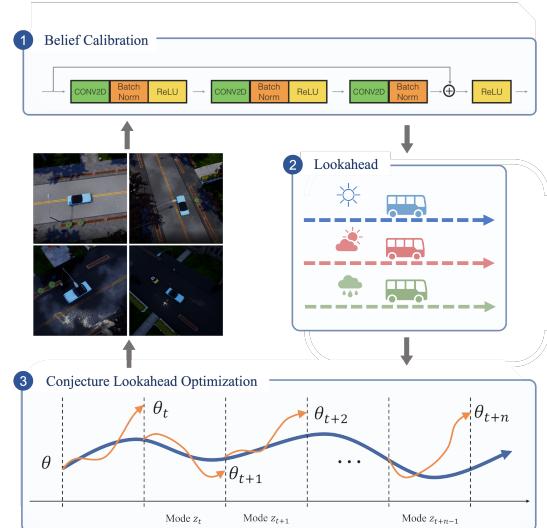


Fig. 1: An illustration of conjectural online lookahead adaptation. When driving in a changing environment, the agent first uses a residual neural network and bayesian filtering to calibrate its belief at every time step about the hidden mode. Based on its belief, the agent conjectures its performance in the future within a lookahead horizon. The policy is adapted through conjectural lookahead optimization, leading to a suboptimal (empirically) online control.

and steering) to keep a vehicle in its travel lane, using only images from the front-facing camera as inputs. The driving testbed is built on the CARLA platform [10]. As shown in Figure 2, different weather conditions create significant visual differences, and a vision-based policy learned under a single weather condition may not generalize to other conditions. As observed in Figure 3a, a vision-based SD policy trained under the sunny/cloudy daytime condition does not generalize to the rainy condition. The trained policy relies on the solid yellow lines in the camera images to guide the vehicle. However, such a policy fails on a rainy day when the lines are barely visible.

The challenge of limited generalization capabilities has motivated various research efforts. In particular, as a learning-to-learn approach, meta learning [11] stands out as one of the well-accepted frameworks for designing fast adaptation strategies. Note that the current meta learning approaches primarily operate in an offline manner. For example, in model-agnostic meta learning (MAML) [12], the meta policy needs to be first trained in advance. When adapting the meta policy to the new environment in the execution time, the agent needs a batch of sample trajectories to

*The authors contributed equally to this work.

¹The authors are with the Department of Electrical and Computer Engineering, New York University, Brooklyn, NY, 11201, USA. {tl2636, h14155, qz494}@nyu.edu

compute the policy gradients to update the policy. However, when interacting with a nonstationary environment online, the agent has no access to the prior knowledge or batches of sample data. Hence, it is unable to learn a meta policy beforehand. In addition, past experiences only reveal information about the previous environments. The gradient updates based on past observations may not suffice to prepare the agent for the future. In summary, the challenge mentioned above persists when the RL agent interacts with a nonstationary environment in an online manner.

a) Our Contributions: To address the challenge of limited online adaptation ability, this work proposes an online meta reinforcement learning algorithm based on *conjectural online lookahead adaptation* (COLA). Unlike previous meta RL formulations focusing on policy learning, COLA is concerned with learning meta adaptation strategies online in a nonstationary environment. We refer to this novel learning paradigm as online meta adaptation learning (OMAL). Specifically, COLA determines the adaptation mapping at every step by maximizing the agent's conjecture of the future performance in a lookahead horizon. This lookahead optimization is approximately solved using off-policy data, achieving real time adaptation. A schematic illustration of the proposed COLA algorithm is provided in Figure 1.

In summary, the main contributions of this work include that 1) we formulate the problem of learning adaptation strategies online in a nonstationary environment as online meta adaptation learning (OMAL); 2) we develop a real-time OMAL algorithm based on conjectural online lookahead adaptation (COLA); 3) experiments show that COLA equips the self-driving agent with online adaptability, leading to self-adaptive driving under dynamic weather.

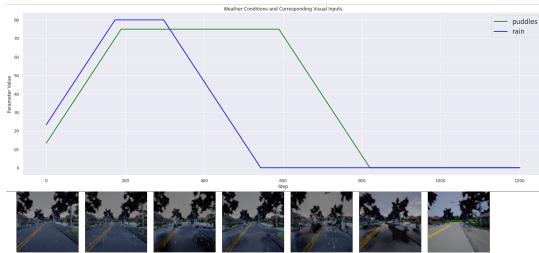


Fig. 2: An example of lane-keeping task in an urban driving environment under time-varying weather conditions. Dynamic weather is realized by varying three weather parameters: cloudiness, rain, and puddles. Different weather conditions cause significant visual differences in the low-resolution image. As a result, a perception-control policy trained under one weather does not generalize to another.

II. SELF-DRIVING IN NONSTATIONARY ENVIRONMENTS: MODELING AND CHALLENGES

We model the lane-keeping task under time-varying weather conditions shown in Figure 2 as an Hidden-Mode Markov Decision Process (HM-MDP) [13]. The state input at time t is a low-resolution image, denoted by $s_t \in \mathcal{S}$. Based on the state input, the agent employs a control action $a_t \in \mathcal{A}$

including acceleration, braking, and steering. The discrete control commands used in our experiment can be found in Appendix Section III. The driving performance at s_t when taking a_t is measured by a reward function $r_t = r(s_t, a_t)$.

Upon receiving the control commands, the vehicle changes its pose/motion and moves to the next position. The new surrounding traffic conditions captured by the camera serve as the new state input subject to a transition kernel $P(s_{t+1}|s_t, a_t; z_t)$, where $z_t \in \mathcal{Z}$ is the environment mode or latent variable hidden from the agent, corresponding to the weather condition at time t . The transition $P(s_{t+1}|s_t, a_t; z_t)$ tells how likely the agent is to observe a certain image s_{t+1} under the current weather conditions z_t . CARLA simulator controls the weather condition through three weather parameters: cloudiness, rain, and puddles, and z_t is a three-dimensional vector with entries being the values of these parameters.

In HM-MDP, the hidden mode shifts stochastically according to a Markov chain $p_z(z_{t+1}|z_t)$ with initial distribution $p_z(z_1)$. As detailed in the Appendix, the hidden mode (weather condition) shifts in our experiment are realized by varying three weather parameters subject to a periodic function. One realization example is provided in Figure 2. Note that the hidden mode z_t , as its name suggested, is not observable in the decision-making process. Let $\mathcal{I}_t = \{s_t, a_{t-1}, r_{t-1}\}$ be the set of agent's observations at time t , referred to as the information structure [14]. Then, the agent's policy π is a mapping from the past observations $\cup_{k=1}^t \mathcal{I}_k$ to a distribution over the action set $\Delta(\mathcal{A})$.

a) Reinforcement Learning and Meta Learning: Standard RL concerns stationary MDP, where the mode remains unchanged (i.e., $z_t = z$) for the whole decision horizon H . Due to stationarity, one can search for the optimal policy within the class of Markov policies [15], where the policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ only depends on the current state input.

We consider a neural network policy $\pi(s, a; \theta), \theta \in \Theta \subset \mathbb{R}^d$ when dealing with high-dimensional state inputs such as image inputs in this work. The RL problem for the stationary MDP (fixing z) is to find an optimal policy maximizing the expected cumulative rewards discounted by $\gamma \in (0, 1]$:

$$\max_{\theta} J_z(\theta) := \mathbb{E}_{P(\cdot|\cdot; z), \pi(\cdot; \theta)} \left[\sum_{t=1}^H \gamma^t r(s_t, a_t) \right]. \quad (\text{RL})$$

We use the policy gradient method to solve for the maximization using sample trajectories $\tau = (s_1, a_1, \dots, s_H, a_H)$. The idea is to apply gradient descent with respect to the objective function $J_z(\theta)$. Following the policy gradient theorem [7], we obtain $\nabla J_z(\theta) = \mathbb{E}[g(\tau; \theta)]$, where $g(\tau; \theta) = \sum_{h=1}^H \nabla_{\theta} \log \pi(a_h|s_h; \theta) R^h(\tau)$, and $R^h(\tau) = \sum_{t=h}^H r(s_t, a_t)$. In RL practice, the policy gradient $\nabla J_z(\theta)$ is replaced by its MC estimation since evaluating the exact value is intractable. Given a batch of trajectories $\mathcal{D} = \{\tau\}$ under the policy $\pi(\cdot|\cdot; \theta)$, the MC estimation is $\hat{J}(\theta, \mathcal{D}(\theta)) := 1/|\mathcal{D}| \sum_{\tau \in \mathcal{D}(\theta)} g(\tau; \theta)$. Our implementation applies the actor-critic (AC) method [16], a variant of the policy gradient, where the cumulative reward $R^h(\tau)$ in

$g(\tau; \theta)$ is replaced with a Q -value estimator (critic) [16]. The AC implementation is in Appendix Section IV.

To improve the adaptability of RL policies, meta reinforcement learning (meta RL) aims to find a meta policy θ and an adaptation mapping Φ that returns satisfying rewards in a range of environments when the meta policy θ is updated using Φ within each environment. Mathematically, meta RL amounts to the following optimization problem [17]:

$$\max_{\theta, \Phi} \mathbb{E}_{z \sim \rho_z} [J_z(\Phi(\theta, \mathcal{D}_z))], \quad (\text{Meta-RL})$$

where the initial distribution ρ_z gives the probability that the agent is placed in the environment z . \mathcal{D}_z denotes a collection of sample trajectories under environment z rolled out under θ . The adaptation mapping $\Phi(\theta, \mathcal{D}_z)$ adapts the meta policy θ to a new policy θ_z fine-tuned for the specific environment z using \mathcal{D}_z . For example, the adaptation mapping in MAML [12] is defined by a policy-gradient update, i.e., $\Phi(\theta, \mathcal{D}_z) = \theta + \alpha \hat{\nabla} J_z(\theta, \mathcal{D}_z)$, where α is the step size to be optimized [18]. Since the meta policy maximizes the average performance over a family of tasks, it provides a decent starting point for fine-tuning that requires far less data than training θ_z from scratch. As a result, the meta policy generalizes to a collection of tasks using sample trajectories.

b) *Challenges in Online Adaptation:* Meta RL formulation reviewed above does not fully address the limited adaptation issue in execution time, since (Meta-RL) does not include the evolution of the hidden mode $z_{t+1} \sim p_z(\cdot | z_t)$. When interacting with a nonstationary environment online, the agent cannot collect enough trajectories \mathcal{D}_{z_t} for adaptation as z_t is time-varying. In addition, past experiences \mathcal{D}_{z_t} only reveal information about the previous environments. The gradient updates based on past observations may not suffice to prepare the agent for the future.

III. CONJECTURAL ONLINE LOOKAHEAD ADAPTATION

Instead of casting meta learning as a static optimization problem in (Meta-RL), we consider learning the meta adaptation online, where the agent updates its adaptation strategies at every step based on its observations. Unlike previous works [17], [19], [20] aiming at learning a meta policy offline, our meta learning approach enables the agent to adapt to the changing environment. The following formally defines the online meta adaptation learning (OMAL) problem.

Unlike (Meta-RL), the online adaptation mapping relies on online observations $\cup_{k=1}^t \mathcal{I}_k$ instead of a batch of samples \mathcal{D}_z : the meta adaptation at time t is defined as a mapping $\Phi_t(\theta) := \Phi(\theta, \cup_{k=1}^t \mathcal{I}_k)$. Let $r_t^\pi(s_t; \theta) := \mathbb{E}_{a \sim \pi(\cdot | s_t; \theta)} [r_t(s_t, a)]$ be the expected reward. The online meta adaptation learning in the HM-MDP is given by

$$\max_{\{\Phi_t\}_{t=1}^H} \mathbb{E}_{z_1, z_2, \dots, z_H} [\sum_{t=1}^H r_t^\pi(s_t; \Phi_t(\theta))], \quad (\text{OMAL})$$

where $z_{t+1} \sim p_z(\cdot | z_t), t = 1, \dots, H - 1$, $\theta = \arg \max \mathbb{E}_{z \sim \rho_z} [J_z(\theta)]$. Some remarks on (OMAL) are in order. First, similar to (RL) and (Meta-RL), (OMAL) needs to be solved in a model-free manner as the transition kernel

$P(\cdot | \cdot; z_t)$ regarding image inputs is too complicated to be modeled or estimated. Meanwhile, the hidden mode transition p_z is also unknown to the agent. Second, the meta policy in (OMAL) is obtained in a similar vein as in (Meta-RL): θ maximizes the rewards across different environments, providing a decent base policy to be adapted later (see Algorithm 2 in Appendix). Unlike the offline operation in (Meta-RL) (i.e., finding θ and Φ before execution), (OMAL) determines Φ_t on the fly in the execution phase based on $\cup_{k=1}^t \mathcal{I}_k$ to accommodate the nonstationary environment.

Theoretically, searching for the optimal meta adaptation mappings $\{\Phi_t\}_{t=1}^H$ is equivalent to finding the optimal nonstationary policy $\{\theta_t\}_{t=1}^H$ in the HM-MDP, i.e.,

$$\max_{\{\theta_t\}_{t=1}^H} \mathbb{E}_{z_1, z_2, \dots, z_H} [\sum_{t=1}^H r^\pi(s_t; \theta_t)]. \quad (1)$$

However, solving for (1) either requires domain knowledge regarding $P(\cdot | \cdot; z)$ and p_z [13], [21], [22] or black-box simulators producing successor states, observations, and rewards [23]. When solving complex tasks such as the self-driving task considered in this work, these assumptions are no longer valid. The following presents a model-free real-time adaptation algorithm, where Φ_t is determined by a lookahead optimization conditional on the belief on the hidden mode.

The meta policy in (OMAL) is given by $\theta = \arg \max \mathbb{E}_{z \sim \rho_z} [J_z(\theta)]$. Note that $z \in \mathcal{Z}$ in our case is a three-dimensional vector composed of three weather parameters. The hidden model space \mathcal{Z} contains an infinite number of z , and it is impractical to train the policy over each of these environments. Similar to MAML [12], we sample a batch of modes z using ρ_z and train the meta policy over these sampled modes, as shown in Algorithm 2 in the Appendix. Note that the meta training program returns the stabilized policy iterate θ_k , its corresponding mode sample batch $\tilde{\mathcal{Z}}$, and the associated sample trajectories $\{\tilde{\mathcal{D}}_z\}$ under mode z and policy θ_k . These outputs can be viewed as the knowledge the agent gained in the training phase and are later utilized in the online adaptation process.

A. Lookahead Adaptation

In online execution, the agent forms a belief $b_t \in \Delta(\hat{\mathcal{Z}})$ at time t on the hidden mode based on its past observations $\cup_{k=1}^t \mathcal{I}_k$. Note that the support of the agent's belief is $\hat{\mathcal{Z}}$ instead of the whole space \mathcal{Z} . The intuition is the agent always attributes the current weather pattern to a mixture of known patterns when facing unseen weather conditions. Then, the agent adapts to this new environment using its training experiences. Specifically, the agent conjectures that the environment would evolve according to $P(\cdot | \cdot; z)$ with probability $b_t(z)$ for a short horizon K . Given the agent's belief b_t and its meta policy θ , the trajectory of future K steps $\tau_t^K := (s_t, a_t, \dots, s_{t+K-1}, a_{t+K-1}, s_{t+K})$ follows $q(\tau_t^K; b_t, \theta)$ defined as

$$\prod_{k=0}^{K-1} \pi(a_{t+k} | s_{t+k}; \theta) \times \prod_{k=0}^{K-1} [\sum_{z \in \mathcal{Z}} b_t(z) P(s_{t+k+1} | s_{t+k}, a_{t+k}; z)].$$

In order to maximize the forecast of the future performance in K steps, the adapted policy $\theta_t = \Phi_t(\theta)$ should maximize the forecast future performance: $\max_{\theta' \in \Theta} \mathbb{E}_{q(\tau_t^K; b, \theta')} \sum_{k=0}^{K-1} r(s_{t+k}, a_{t+k})$. Note that the agent has no access to the distribution $q(\cdot; b, \theta')$ in the online setting, and hence, can not use policy gradient methods to solve for the maximizer. Inspired by off-policy gradient methods [24], we approximate the solution to the future performance optimization using off-policy data $\{\hat{\mathcal{D}}_z\}$.

Following the approximation idea in trust region policy optimization (TRPO) [25], the maximization of the future can be reformulated as

$$\max_{\theta' \in \Theta} \mathbb{E}_{q(\cdot; b_t, \theta)} \left[\prod_{k=0}^{K-1} \frac{\pi(a_{t+k}|s_{t+k}; \theta')}{\pi(a_{t+k}|s_{t+k}; \theta)} \sum_{k=0}^{K-1} r(s_{t+k}, a_{t+k}) \right] \quad (\text{CLO})$$

subject to $\mathbb{E}_{s \sim q} D_{KL}(\pi(\cdot|s; \theta), \pi(\cdot|s; \theta')) \leq \delta$,

where D_{KL} is the Kullback-Leibler divergence. In the KL divergence constraint, we slightly abuse the notation $q(\cdot)$ to denote the discounted state visiting frequency $s \sim q$.

Note that the objective function in (CLO) is equivalent to the future performance under θ' : $\mathbb{E}_{q(\tau_t^K; b, \theta')} \sum_{k=0}^{K-1} r(s_{t+k}, a_{t+k})$. This is because the distribution shift between $q(\tau_t^K; b, \theta')$ and $q(\tau_t^K; b, \theta)$ in (CLO) is compensated by the ratio $\prod_{k=0}^{K-1} \frac{\pi(a_{t+k}|s_{t+k}; \theta')}{\pi(a_{t+k}|s_{t+k}; \theta)}$. The intuition is that the expectation in (CLO) can be approximated using $\mathcal{D}_z(\theta)$ collected in the meta training. When θ' is close to the base policy θ in terms of KL divergence, the estimated objective in (CLO) using sample trajectories under θ returns a good approximation to $\mathbb{E}_{q(\tau_t^K; b, \theta')} \sum_{k=0}^{K-1} r(s_{t+k}, a_{t+k})$. This approximation is detailed in the following subsection.

B. Off-policy Approximation

We first simplify (CLO) by linearizing the objective and approximating the constraints using a quadratic surrogate as introduced in [25], [26]. Denote by $L(\theta'; \theta)$ the objective function in (CLO). A first-order linearization of the objective at $\theta' = \theta$ is given by $L(\theta'; \theta) \approx \nabla_{\theta'} L(\theta'; \theta)|_{\theta'=\theta} (\theta' - \theta)$. Note that θ' is the variable and θ is the known meta policy. The gradient $\nabla_{\theta'} L(\theta'; \theta)|_{\theta'=\theta}$ is exactly the policy gradient under the trajectory distribution $q(\cdot; b_t, \theta)$. Using the notions introduced in Section II, we obtain $\nabla_{\theta'} L(\theta'; \theta)|_{\theta'=\theta} = \mathbb{E}_{q(\cdot; b_t, \theta)} [g(\tau_t^K; \theta)] = \sum_{z \in \hat{\mathcal{Z}}} b(z) \mathbb{E}_{q(\cdot; z, \theta)} [g(\tau_t^K; \theta)]$. Note that τ_t^K is the future K -step trajectory, which is not available to the agent at time t , and a substitute is its counterpart in $\hat{\mathcal{D}}_z$: the sample trajectory within the same time frame $[t, t+K]$ in meta training. Denote this counterpart by $\hat{\tau}_t^K$, then $\mathbb{E}_{q(\cdot; b_t, \theta)} [g(\tau_t^K; \theta)] \approx \sum_{z \in \hat{\mathcal{Z}}} b(z) g(\hat{\tau}_t^K; \theta)$, and we denote its sample estimate by $\hat{g}(b; \theta)$. The distribution $q(\cdot|z, \theta)$ is defined similarly as $q(\cdot|b, \theta)$. More details on this distribution and $\hat{g}(b; \theta)$ are in Appendix Section II.

We consider a second-order approximate for the KL-divergence constraint because the gradient of D_{KL} at θ is zero. In this case, the approximated constraint is $\frac{1}{2}(\theta - \theta')^\top A(\theta - \theta') \leq \delta$, where $A_{ij} =$

$\frac{\partial}{\partial \theta_i \partial \theta_j} \mathbb{E}_{s \sim q} D_{KL}(\pi(\cdot|s; \theta), \pi(\cdot|s; \theta'))|_{\theta'=\theta}$. Let \hat{A} be the sample estimate using $\hat{\tau}_t^K$ (see Appendix Section II), then the off-policy approximate to (CLO) is given by

$$\begin{aligned} & \max_{\theta'} \hat{g}(b; \theta)(\theta' - \theta) \quad (\text{OPA-CLO}) \\ & \text{subject to } \frac{1}{2}(\theta - \theta')^\top \hat{A}(\theta - \theta') \leq \delta. \end{aligned}$$

To find such θ' , we compute the search direction $d\theta = \hat{A}^{-1} \hat{g}$ efficiently using conjugate gradient method as in original TRPO implementation [25]. In general, the step size is determined by backtracking line search, yet we find that fixed step size also achieves impressive results in experiments.

C. Belief Calibration

The last piece in our online adaptation learning is belief calibration: how the agent infers the currently hidden mode z_t based on past experiences and updates its belief based on new observations. A general-purpose strategy is to train an inference network through a variational inference approach [19]. Considering the significant visual differences caused by weather conditions, we adopt a more straightforward approach based on image classification and Bayesian filtering. To simplify the exposition, we only consider two kinds of weather: cloudy (denoted the mode by z_c) and rainy (denoted by z_r) in meta training $\hat{\mathcal{Z}} = \{z_r, z_c\}$.

a) *Image Classification*: The image classifier is a binary classifier based on ResNet architecture [27], which is trained using cross-entropy loss. The data sets (training, validation, and testing) include RGB camera images generated by CARLA. The input is the current camera image, and the output is the probability of that image being of the rainy type. The training details are included in Appendix Section IV.

Given an arbitrary image, the underlying true label is created by measuring its corresponding parameter distance to the default weather setup in the CARLA simulator (see Appendix Section IV). Finally, we remark that the true labels and the hidden mode z_r, z_c are revealed to the agent in the training phase. Yet, in the online execution, only online observations $\mathcal{I}_t = \{s_t, a_{t-1}, r_{t-1}\}$ are available.

b) *Bayesian Filtering*: Denote the trained classifier by $f : \mathcal{S} \rightarrow \Delta(\hat{\mathcal{Z}})$, where $f(z_r; s)$ denotes the output probability of s being an image input under the mode z_r . During the online execution phase, given an state input s_t and the previous belief $b_{t-1} \in \Delta(\hat{\mathcal{Z}})$, the updated belief is

$$b_t(z_r) = \frac{b_{t-1}(z_r) f(z_r; s_t)}{\sum_{z \in \hat{\mathcal{Z}}} b_{t-1}(z) f(z; s_t)}. \quad (\text{BF})$$

The intuition behind (BF) is that it better captures the temporal correlation among image sequences than the vanilla classifier f . Note that this continuously changing weather is realized through varying three weather parameters. Hence, two adjacent images, as shown in Figure 2, display a temporal correlation: what follows a rain image is highly likely to be another rain image. Considering this correlation, we use Bayesian filtering, where the prior reveals information about previous images. As shown in an ablation study in

Section IV-C, (BF) does increase the belief accuracy (to be defined later in Section IV-C) and associated rewards. Finally, we conclude this section with the pseudocode of the proposed conjectural online lookahead adaptation (COLA) algorithm presented in Algorithm 1.

Algorithm 1 Conjectural Online Lookahead Adaptation

```

Input The meta policy  $\theta$ , classifier  $f$ , training samples  $\{\mathcal{D}_z\}_{z \in \hat{\mathcal{Z}}}$ , sample batch size  $M$ , lookahead horizon length  $K$ , and step size  $\alpha$ . Set  $\theta_1 = \theta$ .
for  $t \in \{1, 2, \dots, H\}$  do
    Obtain the image input  $s_t$  from the camera;
    Implement the action using  $\pi(\cdot | s_t; \theta_t)$ ;
    Obtain the probability  $f(z; s_t)$  from the classifier;
    Update the belief using (BF);
    Sample  $M$  trajectories  $\hat{\tau}_t^K$  under  $\hat{z}$  from  $\{\mathcal{D}_z\}_{z \in \hat{\mathcal{Z}}}$ ;
    Obtain  $\theta'$  by solving (OPA-CLO), and let  $\theta_{t+1} = \theta'$ ;

```

IV. EXPERIMENTS

We evaluate the proposed COLA algorithm using CARLA simulator [10]. The HM-MDP setup for the lane-keeping task (i.e., hidden mode transition, the reward function) is presented in Appendix Section III. This section reports experimental results under $K = 10$, and experiments under other choices are included in Appendix Section V.

A. Comparative Analysis

We compare COLA with the following baselines. 1) The base policy θ_{base} : a stationary policy trained under dynamic weather conditions, i.e., $\theta_{base} = \arg \max \mathbb{E}_{z_1, z_2, \dots, z_H} [\sum_{t=1}^H r^\pi(s_t; \theta)]$. 2) The MAML policy θ_{MAML} : the solution to (Meta-RL) with Φ being the policy gradient computed using the *past* K -step trajectory. 3) The optimal policy under the mixture of Q functions $\theta_{Q_{mix}}$: the action at each time step is determined by $a_t = \arg \max_{a \in \mathcal{A}} \sum_{z \in \hat{\mathcal{Z}}} b_t(z) Q_z(s_t, a)$, where Q_z is the Q function trained under the stationary MDP with fixed mode z . In our experiment, this mixture $\sum_{z \in \hat{\mathcal{Z}}} b_t(z) Q_z(s_t, a)$ is composed of two Q functions Q_c, Q_r (sharing the same network structure as the critic) under the cloudy and rainy conditions, respectively. 4) The oracle policy θ_{oracle} : an approximate solution to (1), which is obtained by solving (CLO) using authentic future trajectories generated by the simulator.

The experimental results are summarized in Figure 3. As shown in Figure 3b, the COLA policy outperforms θ_{base} , θ_{MAML} , and $\theta_{Q_{mix}}$, showing that the lookahead optimization (CLO) better prepares the agent for the future environment changes. Note that $\theta_{Q_{mix}}$ returns the worst outcome in our experiment, suggesting that the naive adaptation by averaging Q function does not work: the nonstationary environment given by HM-MDP does not equal an average of multiple stationary MDPs. The agent needs to consider the temporal correlation as discussed in Section III-C.

To evaluate the effectiveness of COLA-based online adaptation, we follow the regret notion used in online learning [28] and consider the performance differences between

COLA (the other three baselines) and the oracle policy. The performance difference or regret up to time T is defined as $\text{Reg}(T) = \mathbb{E}[\sum_{t=1}^T r^\pi(s_t; \theta_{oracle}) - \sum_{t=1}^T r^\pi(s_t; \Phi_t(\theta))]$. If the regret grows sublinearly, i.e., $\text{Reg}(T) < O(T)$, then the corresponding policy achieves the same performance as the oracle asymptotically. As shown in Section IV-A, the proposed COLA is the only algorithm that achieves comparable results as the oracle policy.

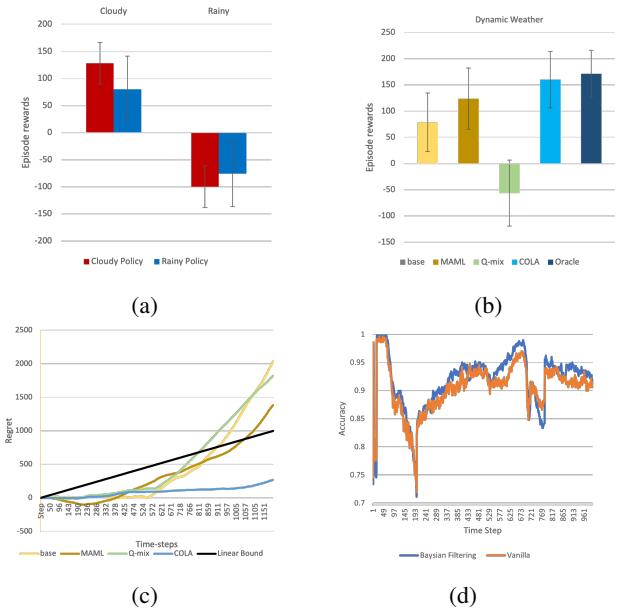


Fig. 3: Evaluations of COLA algorithms in the lane-keeping task under nonstationary weather conditions (averaged over 500 episodes under 10 random seeds). (a) The average rewards of the cloudy-trained and rainy-trained policies in cloudy and rainy conditions. The policy trained under one environment does not generalize to the other. (b) The average rewards of the COLA policy and the four baseline policies. COLA achieves comparable results as the oracle policy and outperforms the other baselines. (c) The regret grows under different policies. Only COLA policy achieves sub-linear regret growth [below the linear bound $O(T)$]. (d) The evolution of transient accuracies (averaged over 500 episodes) with respect to time step t under Bayesian filtering and vanilla classification in dynamic weather.

B. Knowledge Transfer through Off-policy Gradients

This subsection presents an empirical explanation of COLA's adaptability: the knowledge regarding the speed control learned in meta training is carried over to online execution through the off-policy gradient \hat{g} in (OPA-CLO). As observed in Table II in Appendix, when driving in a rainy environment (i.e., z_r), the agent tends to drive slowly to avoid crossing the yellow line, as the visibility is limited. In contrast, it increases the cruise speed when the rain stops (i.e., z_c). Denote the off-policy gradients under z_r, z_c by \hat{g}_r and \hat{g}_c , respectively. Then, solving (OPA-CLO) using \hat{g}_r (\hat{g}_c) leads to a conservative (aggressive) driving in terms of speeds as shown in Figure 4a (Figure 4b), regardless of the environment changes. This result suggests that the

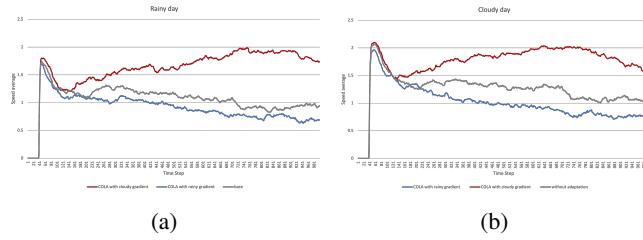


Fig. 4: Figure (a) [or (b)] presents average speeds at each time step within the horizon (average of 500 episodes) with \hat{g}_c (red curves) and \hat{g}_r (blue curves) being used to adapt the meta policy (grey curves) under the cloudy [or rainy in (b)] mode. Using \hat{g}_c [or \hat{g}_r] always leads to aggressive [or conservative] driving, regardless of the environment.

knowledge gained in meta training is encoded into the off-policy gradients and is later retrieved according to the belief. Hence, the success of COLA relies on a decent belief calibration, where the belief accurately reveals the underlying mode. Section IV-C justifies our choice of Bayesian filtering.

C. An Ablation Study on Bayesian Filtering

For simplicity, we use the 0-1 loss to characterize the misclassification in the online implementation in this ablation study. Denote by $\hat{i}_t^\nu \in \{0, 1\}$ the label returned by the classifier (i.e., the class assigned with a higher probability; 1 denotes the rain class) at time $t \in \{1, \dots, H\}$ in the ν -th episode (500 episodes in total). Then the loss at time t is $1_{\{\hat{i}_t^\nu \neq i_t^\nu\}}$, where i_t^ν is the true label. We define the transient accuracy at time t as Transient Accuracy = $\frac{\sum_{\nu=1}^{500} 1 - 1_{\{\hat{i}_t^\nu \neq i_t^\nu\}}}{500} = \frac{\sum_{\nu=1}^{500} 1_{\{\hat{i}_t^\nu = i_t^\nu\}}}{500}$, which reflects the accuracy of the vanilla classifier at time t . The average transient accuracy over the whole episode horizon H is defined as the episodic accuracy, i.e., Episodic Accuracy = $\frac{\sum_{\nu=1}^{500} \sum_{t=1}^H 1_{\{\hat{i}_t^\nu = i_t^\nu\}}}{500H}$. These accuracy metrics can also be extended to the Bayesian filter, where the outcome label corresponds to the class with a higher probability.

We compare the transient and episodic accuracy as well as the average rewards under Bayesian filtering and vanilla classification, and the results are summarized in Table I and Figure 3d. Even though the two methods' episodic accuracy under two methods differ little, as shown in Table I, the difference regarding the average rewards indicates that Bayesian filtering is better suited to handle nonstationary weather conditions. As shown in Figure 3d, Bayesian filtering outperforms the classifier for most of the horizon in terms of transient accuracy. The reason behind its success is that Bayesian filtering better captures the temporal correlation among the image sequence and returns more accurate predictions than the vanilla classifier at every step, as the prior distribution is incorporated into the current belief. It is anticipated that higher transient accuracy leads to higher mean rewards since the lookahead optimization [see (CLO)] depends on the belief at every time step.

V. RELATED WORKS

a) *Reinforcement Learning for Self Driving*: In most of the self-driving tasks studied in the literature [29], such as

	Bayesian Filtering	Vanilla Classification
Episodic Accuracy	91.18%	90.41%
Average Reward	152.37 ± 49.58	122.44 ± 54.80

TABLE I: Comparisons of episodic accuracies and average rewards under dynamic weather.

lane-keeping and changing [30], [31], overtaking [32], and negotiating intersections [5], the proposed RL algorithms are developed in an offline approach. The agent first learns a proficient policy tailored to the specific problem in advance. Then the learned policy is tested in the same environment as the training one without any adaptations. Instead of focusing on a specific RL task (stationary MDP) as in the above works, the proposed OMAL tackles a nonstationary environment modeled as an HM-MDP. Closely related to our work, a MAML-based approach is proposed in [33] to improve the generalization of RL policies when driving in diverse environments. However, the adaptation strategy in [33] is obtained by offline training and only tested in stationary environments. In contrast, our work investigates online adaptation in a nonstationary environment. The proposed COLA enables the agent to adapt rapidly to the changing environment in an online manner by solving (CLO) approximately in real time.

b) *Online Meta Learning*: To handle nonstationary environments or a streaming series of different tasks, recent efforts on meta learning have integrated online learning ideas with meta learning, leading to a new paradigm called online meta learning [20], [34] or continual meta learning [35], [36]. Similar to our nonstationary environment setup, [20], [35], [36] focuses on continual learning under a nonstationary data sequence involving different tasks, where the nonstationarity is governed by latent variables. Despite their differences in latent variable estimate, these works utilize on-policy gradient adaptation: the agent needs to collect \mathcal{D}_z and then adapt $\theta_z = \Phi(\theta, \mathcal{D}_z)$ [see (Meta-RL)]. Note that on-policy gradient adaptation requires a few sample data under the current task, which may not be suitable for online control. The agent can only collect one sample trajectory online, incurring a significant variance. In addition, by the time the policy is updated using the old samples, the environment has already changed, rendering the adapted policy no longer helpful. In contrast, COLA better prepares the agent for the incoming task/environment by maximizing the forecast of the future performance [see (CLO)]. Moreover, the adaptation in COLA can be performed in real time using off-policy data.

VI. CONCLUSION

This work proposes an online meta reinforcement learning algorithm based on conjectural online lookahead adaptation (COLA) that adapts the policy to the nonstationary environment modeled as a hidden-mode MDP. COLA enables prompt online adaptation by solving the conjectural lookahead optimization (CLO) on the fly using off-policy data, where the agent's conjecture (belief) on the hidden mode is updated in a Bayesian manner.

REFERENCES

- [1] S. Ross, G. Gordon, and D. Bagnell, “A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning,” ser. Proceedings of Machine Learning Research, vol. 15. Fort Lauderdale, FL, USA: JMLR Workshop and Conference Proceedings, 2011, pp. 627–635. [Online]. Available: <http://proceedings.mlr.press/v15/ross11a.html>
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [3] T. Li, G. Peng, Q. Zhu, and T. Basar, “The Confluence of Networks, Games and Learning,” *arXiv*, 2021.
- [4] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yogamani, and P. Perez, “Deep Reinforcement Learning for Autonomous Driving: A Survey,” *IEEE Transactions on Intelligent Transportation Systems*, vol. PP, no. 99, pp. 1–18, 2021.
- [5] J. Chen, S. E. Li, and M. Tomizuka, “Interpretable End-to-End Urban Autonomous Driving With Latent Deep Reinforcement Learning,” *IEEE Transactions on Intelligent Transportation Systems*, vol. PP, no. 99, pp. 1–11, 2021.
- [6] T. Li and Q. Zhu, “On Convergence Rate of Adaptive Multiscale Value Function Approximation for Reinforcement Learning,” *2019 IEEE 29th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6, 2019.
- [7] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Advances in Neural Information Processing Systems 12*. MIT press, 2000, pp. 1057–1063. [Online]. Available: <http://papers.nips.cc/paper/1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation.pdf>
- [8] T. Li, G. Peng, and Q. Zhu, “Blackwell Online Learning for Markov Decision Processes,” *2021 55th Annual Conference on Information Sciences and Systems (CISS)*, vol. 00, pp. 1–6, 2021.
- [9] S. Padakandla, “A Survey of Reinforcement Learning Algorithms for Dynamically Varying Environments,” *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–25, 2021.
- [10] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “CARLA: An Open Urban Driving Simulator,” vol. 78, pp. 1–16, 2017. [Online]. Available: <https://proceedings.mlr.press/v78/dosovitskiy17a.html>
- [11] T. M. Hospedales, A. Antoniou, P. Micaelli, and A. J. Storkey, “Meta-Learning in Neural Networks: A Survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–1, 2021.
- [12] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 1126–1135.
- [13] S. P. M. Choi, D.-Y. Yeung, and N. L. Zhang, *Hidden-Mode Markov Decision Processes for Nonstationary Sequential Decision Making*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2001, pp. 264–287.
- [14] T. Li, Y. Zhao, and Q. Zhu, “The role of information structures in game-theoretic multi-agent learning,” *Annual Reviews in Control*, vol. 53, pp. 296–314, 2022.
- [15] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, 1st ed. USA: John Wiley & Sons, Inc., 1994.
- [16] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous Methods for Deep Reinforcement Learning,” 2016.
- [17] A. Nagabandi, I. Clavera, S. Liu, R. S. Fearing, P. Abbeel, S. Levine, and C. Finn, “Learning to Adapt in Dynamic, Real-World Environments Through Meta-Reinforcement Learning,” *arXiv*, 2018.
- [18] Z. Li, F. Zhou, F. Chen, and H. Li, “Meta-sgd: Learning to learn quickly for few-shot learning,” *arXiv preprint arXiv:1707.09835*, 2017.
- [19] K. Rakelly, A. Zhou, C. Finn, S. Levine, and D. Quillen, “Efficient off-policy meta-reinforcement learning via probabilistic context variables,” in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 5331–5340. [Online]. Available: <https://proceedings.mlr.press/v97/rakelly19a.html>
- [20] J. Harrison, A. Sharma, C. Finn, and M. Pavone, “Continuous Meta-Learning without Tasks,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., pp. 17571–17581.
- [21] J. Pineau, G. Gordon, and S. Thrun, “Anytime Point-Based Approximations for Large POMDPs,” *Journal of Artificial Intelligence Research*, vol. 27, pp. 335–380, 2011.
- [22] S. Ross, J. Pineau, S. Paquet, and B. Chaib-draa, “Online Planning Algorithms for POMDPs,” *Journal of Artificial Intelligence Research*, vol. 32, no. 2, pp. 663–704, 2014.
- [23] D. Silver and J. Veness, “Monte-Carlo Planning in Large POMDPs,” in *Advances in Neural Information Processing Systems*, vol. 23. Curran Associates, Inc.
- [24] J. Queaney, Y. Paschalidis, and C. G. Cassandras, “Generalized Proximal Policy Optimization with Sample Reuse,” in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., pp. 11909–11919.
- [25] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, “Trust region policy optimization,” in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 1889–1897. [Online]. Available: <https://proceedings.mlr.press/v37/schulman15.html>
- [26] J. Achiam, D. Held, A. Tamar, and P. Abbeel, “Constrained Policy Optimization,” vol. 70, pp. 22–31, 2017–06. [Online]. Available: <https://proceedings.mlr.press/v70/achiam17a.html>
- [27] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” *arXiv*, 2015.
- [28] S. Shalev-Shwartz, “Online Learning and Online Convex Optimization,” *Foundations and Trends in Machine Learning*, vol. 4, no. 2, pp. 107–194, 2011.
- [29] W. Schwarting, J. Alonso-Mora, and D. Rus, “Planning and Decision-Making for Autonomous Vehicles,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, no. 1, pp. 1–24, 2018.
- [30] A. E. Sallab, M. Abdou, E. Perot, and S. K. Yogamani, “End-to-end deep reinforcement learning for lane keeping assist,” *ArXiv*, vol. abs/1612.04340, 2016.
- [31] P. Wang, C. Chan, and A. de La Fortelle, “A reinforcement learning based approach for automated lane change maneuvers,” *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1379–1384, 2018.
- [32] D. C. K. Ngai and N. H. C. Yung, “A multiple-goal reinforcement learning method for complex vehicle overtaking maneuvers,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 12, no. 2, pp. 509–522, 2011.
- [33] Y. Jaafra, A. Deruyver, J. L. Laurent, and M. S. Naceur, “Context-Aware Autonomous Driving Using Meta-Reinforcement Learning,” *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, vol. 00, pp. 450–455, 2019.
- [34] C. Finn, A. Rajeswaran, S. Kakade, and S. Levine, “Online Meta-Learning,” vol. 97, 2019, pp. 1920–1930. [Online]. Available: <https://proceedings.mlr.press/v97/finn19a.html>
- [35] A. Xie, J. Harrison, and C. Finn, “Deep Reinforcement Learning amidst Lifelong Non-Stationarity,” *arXiv*, 2020.
- [36] M. Caccia, P. Rodriguez, O. Ostapenko, F. Normandin, M. Lin, L. Page-Caccia, I. H. Laradji, I. Rish, A. Lacoste, D. VAzquez, and L. Charlin, “Online Fast Adaptation and Knowledge Accumulation (OSAKA): a New Approach to Continual Learning,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 16532–16545.
- [37] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “CARLA: An open urban driving simulator,” in *Proceedings of the 1st Annual Conference on Robot Learning*, 2017, pp. 1–16.
- [38] P. Palanisamy, *Hands-On Intelligent Agents with OpenAI Gym: Your Guide to Developing AI Agents Using Deep Reinforcement Learning*. Packt Publishing, 2018.
- [39] J. Schulman, X. Chen, and P. Abbeel, “Equivalence between policy gradients and soft q-learning,” 2017. [Online]. Available: <https://arxiv.org/abs/1704.06440>
- [40] L. N. Smith and N. Topin, “Super-convergence: Very fast training of neural networks using large learning rates,” 2017. [Online]. Available: <https://arxiv.org/abs/1708.07120>

APPENDIX I META TRAINING

The meta training procedure is described in Algorithm 2. Our experiment applies two default weather conditions in

Algorithm 2 Meta Training

```

Input Initialization  $\theta_0$ ,  $k = 0$ , step size  $\alpha$ .
while not converge do
    Draw the  $k$ -th batch of modes (weather conditions)
     $\mathcal{Z}_{train}^k = \{z\}, z \in \mathcal{Z}$ ;
    for each environment  $z$  do
        Sample a batch of trajectories  $\mathcal{D}_z(\theta_k)$  under
         $P(\cdot|.;z)$  and  $\pi(\cdot|\theta_k)$ ;
        Compute the policy gradient  $\hat{\nabla}J_z(\theta_k, \mathcal{D}_z)$ 
         $\theta_{k+1} = \theta_k + \alpha / |\mathcal{Z}_{train}^k| \sum_z \hat{\nabla}J_z(\theta_k, \mathcal{D}_z)$ ;
         $k = k + 1$ .
    Return  $\theta_k, \hat{\mathcal{Z}} = \mathcal{Z}_{train}^k, \{\hat{\mathcal{D}}_z\} = \{\mathcal{D}_z(\theta_k)\}_{z \in \mathcal{Z}_{train}^k}$ .
```

CARLA world: the cloudy and rainy in meta training. In this case, $\hat{\mathcal{Z}} = \{z_r, z_c\}$ and $\hat{\mathcal{D}}_{z_r}$ corresponds to the sample trajectories collected under the rainy condition. Likewise, $\hat{\mathcal{D}}_{z_c}$ denotes those under the cloudy one.

APPENDIX II OFF-POLICY APPROXIMATION

To reduce the variance of the estimate $\hat{g}(b, \theta)$, we use the sample mean of multiple trajectories, i.e.,

$$\mathbb{E}_{q(\cdot; b_t, \theta)}[g(\tau_t^K; \theta)] \approx \sum_{z \in \hat{\mathcal{Z}}} \sum_{\hat{\tau}_t^K \in \hat{\mathcal{D}}_z} 1/|\hat{\mathcal{D}}_z| b(z) g(\hat{\tau}_t^K; \theta),$$

$$q(\tau_t^K; z, \theta) = \prod_{k=0}^{K-1} \pi(a_{t+k}|s_{t+k}; \theta) P(s_{t+k+1}|s_{t+k}, a_{t+k}|z).$$

The sample estimate of A using the off-policy data is

$$\hat{A}_{ij} = \frac{1}{K} \sum_{k=0}^{K-1} \frac{\partial}{\partial \theta_i \partial \theta_j} D_{KL}(\pi(\cdot|s_{t+k}; \theta), \pi(\cdot|s_{t+k}; \theta')).$$

One can also take the sample mean of \hat{A} over multiple trajectories $\hat{\tau}_t^K$.

It should be noted that $\hat{\tau}_t^K$ is not an i.i.d copy of τ_t^K , as the initial states $s_t \in \tau_t^K$ and $\hat{s}_t \in \hat{\tau}_t^K$ are subject to different distributions determined by previous policies $\{\theta_k\}_{k=1}^t$. In other words, as the online adaptation updates θ_t at each time step by solving (CLO), θ_t deviates from the meta policy θ , leading to a distribution shift between τ_t^K and $\hat{\tau}_t^K$. Hence, $\hat{\tau}_t^K$ is referred to as the off-policy (i.e., not following the same policy) data with respect to τ_t^K . Thanks to the KL-constraint in (CLO), such deviation is manageable through the constraint δ , and a mechanism has been proposed in [24] to address the distribution shift.

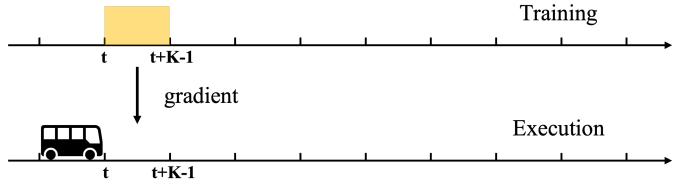


Fig. 5: Off-policy Approximation. The future K -step trajectory τ_t^K is not available at time t in online execution. A substitute is its counterpart in $\hat{\mathcal{D}}_z$: the sample trajectory $\hat{\tau}_t^K$ within the same time frame $[t, t+K]$ in meta training.

Algorithm	Cloudy	Rainy	Dynamic
COLA	1.26	1.07	1.53
COLA-LTM	1.73	0.90	1.60

TABLE II: Average speeds of 500 episodes under COLA and COLA-LTM in three weather setups.

APPENDIX III EXPERIMENT SETUP

Our experiments use CARLA-0.9.4 [10] as the autonomous driving simulator. On top of the CARLA, we modify the API: Multi-Agent Connected Autonomous Driving (MACAD) Gym [37] to facilitate communications between learning algorithms and environments. To simulate various traffic conditions, we enable non-player vehicles controlled by the CARLA server.

Three kinds of weather conditions are considered in experiments: stationary cloudy, rainy, and dynamic weather conditions. For stationary weather conditions, we use two pre-defined weather setup: *CloudyNoon* and *MidRainyNoon* provided by CARLA. For the dynamic environment, the changing weather in experiments is realized by varying three weather parameters: cloudiness, precipitation, and puddles. The changing pattern of dynamic weather is controlled by the parameters update frequency, the magnitude of each update, and the initial values. The basic idea is that as the clouds get thick and dense, the rain starts to fall, and the ground gets wet; when the rain stops, the clouds get thin, and the ground turns dry.

a) *States, Actions, and Rewards*: The policy model uses the current and previous 3-channel RGB images from the front-facing camera (attached to the car hood) as the state input. Basic vehicle controls in CARLA include throttle, steer, brake, hand brake, and reverse. Our experiments only consider the first three, and the agent employs discrete actions by providing throttle, brake, and steer values. Specifically, these values are organized as two-dimensional vectors [throttle/brake, steer]: the first entry indicates the throttle or brake values, while the other represents the steer ones. In summary, the action space is a 2-dimensional discrete set, and its elements are shown in Table III.

The reward function consists of three components: the speed maintenance reward, collision penalty, and driving-too-slow penalty. The speed maintenance function in (III.2)

TABLE III: Discrete control actions in the experiments

Action Vector	Control
[0.0, 0.0]	Coast
[0.0, -0.5]	Turn Left
[0.0, 0.5]	Turn Right
[1.0, 0.0]	Forward
[-0.5, 0.0]	Brake
[0.5, -0.5]	Forward Left
[0.5, 0.5]	Forward Right
[-0.5, -0.5]	Brake Left
[-0.5, 0.5]	Brake Right
[0.5, 0.0]	Bear Forward
[0.75, 1.0]	Sharp Right
[0.75, -1.0]	Sharp Left
[-0.5, 1.0]	Brake Sharp Right
[-0.5, -1.0]	Brake Sharp Left
[0.5, 0.75]	Accelerate Forward Inferior Sharp Right
[0.5, -0.75]	Accelerate Forward Inferior Sharp Left
[0.75, 0.0]	Accelerate Forward
[-1.0, 0.0]	Emergency Brake

encourages the agent to drive the vehicle at 30 m/s. Note that the corresponding sensors attached to the ego vehicle return the speed and collision measurements for reward computation. These measurements are not revealed to the agent when devising controls.

$$F_{SMR}(t) = \begin{cases} -0.005, & V(t) < 0 \\ \frac{1}{9}V(t)^2, & 0 \leq V(t) \leq 30 \\ 5(50 - V(t)), & 30 < V(t) \leq 50 \\ -2(V(t) - 50)^2, & 50 < V(t) \end{cases}, \quad (\text{III.2})$$

To penalize the agent for colliding with other objects, we introduce the collision penalty function F_{CP} defined in (III.3). F_{CP} is a terminate reward function. The terminal reward depends on 1) whether the agent completes the task (the episode will be terminated and reset when collisions happen); 2) how many steps the ego vehicle runs out of the lane in an episode. Denote by $t_{\text{terminate}} \leq H$ the terminal time step and by N_{out} the number of time steps when the ego vehicle is out of the lane. Let $\delta(\cdot)$ be the Dirac function, and F_{CP} is given by

$$F_{CP} = -100 + (t_{\text{terminate}} - H) + 100\delta(t_{\text{terminate}} - H) - 0.1N_{\text{out}}, \quad (\text{III.3})$$

penalizing the agent for not completing the task and running out of the lane.

The agent is said to drive too slowly if the speed is lower than 0.5 m/s in two consecutive steps. Similar to the collision penalty, we introduce a driving-too-slowly(DTS) penalty as a terminal cost. Denote by N_{slow} the number of steps when the vehicle runs too slowly, and the DTS penalty is given by

$$F_{DTS} = -0.005 \times \max(50, 10^{-7} \times N_{\text{slow}}). \quad (\text{III.4})$$

The horizon length is 1200 time steps, and the interval for adjacent frames is 0.05 seconds. The discount factor of all experiments is 0.99.

b) Dynamic Weather Setup: To realize dynamic weather conditions, we create a function $W(t)$ with respect to time defined in the following. To randomize the experiments, the initial value $W(0)$ is uniformly sampled from $[-150, 100]$ (subject to different random seeds in repeated experiments). To align time-varying weather with the built-in timeclock in CARLA, we first introduce a weather update frequency $1/\Delta T, \Delta T > 0$. The whole episode is evenly divided into n intervals: $[0, \Delta T], [\Delta T, 2\Delta T], \dots, [(n-1)\Delta T, n\Delta T]$. Then, for $t \in [k\Delta T, (k+1)\Delta T), k > 1$, the function $W(t)$ is given by

$$W(t) = \frac{25}{12} \cdot |\text{Mod}(1.3k\Delta T + W(0), 250) - 125| - 150,$$

and $W(t) = W(0)$, for $t \in [0, \Delta T)$.

Note that the hidden environment mode $z(t)$ consists of three weather parameters: clouds, rain, and puddles. Denote by $\text{Clouds}(t)$, $\text{Rain}(t)$, and $\text{Puddles}(t)$ the three parameters at time t , i.e., $z(t) = [\text{Clouds}(t), \text{Rain}(t), \text{Puddles}(t)]$. Based on the function $W(t)$, these parameters are given by

$$\begin{cases} \text{Clouds}(t) := \text{Clip}(W(t) + 40, 0, 90), \\ \text{Rain}(t) := \text{Clip}(W(t), 0, 80), \\ \text{Puddles}(t) := \text{Clip}(W(t) + d(t), 0, 75), \end{cases}$$

where $d(t)$ is defined as the following: for $t \in [k\Delta T, (k+1)\Delta T), k \geq 0$,

$$d(t) = \begin{cases} -10, & \text{if } W((k+1)\Delta T) \geq W(k\Delta T), \\ 90, & \text{otherwise.} \end{cases}$$

$d(t)$ is a translation applied to $W(t)$ in the definition of $\text{Puddles}(t)$ to preserve the causal relationship between rain and puddles in the simulation: puddles are caused by rain and appear after the rain starts.

APPENDIX IV TRAINING DETAILS

a) The Actor-Critic Policy Model: Our policy models under three weather conditions: cloudy, rainy, and dynamic changing weather, are trained via the Asynchronous Advantage Actor-Critic algorithm (A3C) with Adam optimizer mentioned in [38]. The learning rate begins at 1×10^{-4} , and the policy gradient update is performed every 10 steps. The reward in each step is clipped to -1 or +1, and the entropy regularized method [39] is used. Once episode rewards stabilize, the learning rate will be changed to 1×10^{-5} and 1×10^{-6} .

The policy network consists of two parts: the Actor and the Critic. They incorporate three convolution layers that use Rectified Linear Units (ReLU) as activation functions. The input is a $1 \times 6 \times 84 \times 84$ dimension image data transformed from two $80 \times 80 \times 3$ RGB images. A linear layer is appended to convolution layers. The softmax layer serves as the last layer in the Actor and the Critic. The schematic diagram of the model structure is shown in Figure 6.

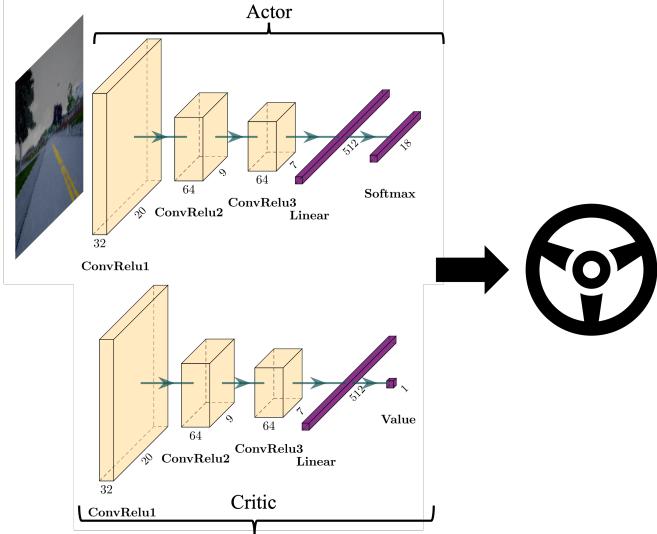


Fig. 6: Network structures of the actor and the critic. The perception and decision-making module are incorporated into a single policy model that outputs control commands when fed with raw images.

b) Image Classifier: Our image classifier is based on a residual neural network [27] with two residual blocks and one linear output layer. The input is the current camera image, and the output is the probability of that image being of the rainy type. The data sets (training, validation, and testing) include RGB camera images under the cloudy and rainy conditions in the CARLA world.

Every image input is paired with a three-dimensional vector z (the mode), with entries being the values of three weather parameters: cloudiness, rain, and puddles, respectively. Denote the corresponding mode of an image by z_{img} . Let z_r , and z_c be the parameter vectors provided by the simulator in the default rainy and cloudy conditions, respectively. Then, the true label i of an image is given by

$$\text{Label} = \begin{cases} 1, & \text{if } \|z_{img} - z_r\| < \|z_{img} - z_c\|, \\ 0, & \text{otherwise.} \end{cases}$$

We use One-Cycle-Learning-Rate [40] in the classifier training, with the max learning rate 1×10^{-3} in the first 120K iterations of training and 1×10^{-6} in the later 600K iterations for fine-tune training. The test accuracy of the classifier is 91.7%.

APPENDIX V THE TRADEOFF OF LOOKAHEAD HORIZON

The impact of the lookahead horizon K on the online adaptation is due to a tradeoff between variance reduction and belief calibration. If K is small, the sample-based estimate of $\mathbb{E}[\sum_{k=0}^{K-1} r(s_{t+k}, a_{t+k})]$ would incur a large variance, which is then propagated to the gradient computation \hat{g} in (OPA-CLO). Hence, the resulting approximate also exhibits a large variance. On the other hand, when K is large, the agent may look into the future under the wrong belief since the agent believes that the environment is stationary

K	Cloudy	Rainy	Dynamic
5	113.16 ± 47.80	-156.36 ± 60.83	109.72 ± 57.45
10	159.09 ± 38.41	-145.75 ± 47.64	159.89 ± 51.35
20	150.41 ± 43.68	-136.72 ± 55.33	124.03 ± 52.36

TABLE IV: The average cumulative rewards and associated standard deviations under different lookahead horizons K in three weather conditions.

for the future K steps. In this case, the obtained policy may not be able to adapt promptly to the changing environment. Table IV summarizes the average rewards under three lookahead horizons $K = 5, 10, 20$. As we can see from the table, the average rewards under $K = 5$ exhibit the largest variances across all weather conditions in comparison with their counterparts under $K = 10, 20$. Even though $K = 20$ achieves performance in stationary environments (i.e., cloudy and rainy) comparable to that under $K = 10$, we choose $K = 10$ in our experiments as it balances variance reduction and quick belief calibration, leading to the highest rewards in the dynamic weather.