

Safe Model-based Off-policy Reinforcement Learning for Eco-Driving in Connected and Automated Hybrid Electric Vehicles

Zhaoxuan Zhu, Nicola Pivaro, Shobhit Gupta, Abhishek Gupta and Marcello Canova.

Abstract—Deep Reinforcement Learning (DRL) has recently been applied to eco-driving to intelligently reduce fuel consumption and travel time. While previous studies synthesize simulators and model-free DRL (MFDRL), this work proposes a Safe Off-policy Model-Based Reinforcement Learning (SMORL) algorithm for eco-driving. SMORL integrates three key components, namely a computationally efficient model-based trajectory optimizer, a value function learned off-policy and a learned safe set. The advantages over the existing literature are three-fold. First, the combination of off-policy learning and the use of a physics-based model improves the sample efficiency. Second, the training does not require any extrinsic rewarding mechanism for constraint satisfaction. Third, the feasibility of trajectory is guaranteed by using a safe set approximated by deep generative models.

The performance of SMORL is benchmarked over 100 trips against a baseline controller representing human drivers, a non-learning-based optimal controller, a previously designed MFDRL strategy, and the wait-and-see optimal solution. In simulation, SMORL reduces the fuel consumption by more than 21% while keeping the average speed comparable while compared to the baseline controller and demonstrates a better fuel economy while driving faster compared to the MFDRL agent and the non-learning-based optimal controller.

Index Terms—Model-based reinforcement learning, generative models, safety-critical applications, connected and automated vehicles.

I. INTRODUCTION

With the advancement in the vehicular connectivity and autonomy, Connected and Automated Vehicles (CAVs) have the potential to operate in a safer and more time- and fuel-efficient manner [1]. With Vehicle-to-Vehicle (V2V) and Vehicle-to-Infrastructure (V2I) communication, the controller has access to real-time look-ahead information including the terrain, infrastructure and surrounding vehicles. Intuitively, with connectivity technologies, controllers can plan a speed profile that allows the ego vehicle to intelligently pass more signalized intersections in green phases with fewer changes in speed. This problem is formulated as the eco-driving problem, which aims to minimize the weighted sum of the fuel consumption and the travel time between two designated locations

Z. Zhu, S. Gupta and M. Canova are with the Center for Automotive Research, The Ohio State University, Columbus, OH 43212 USA (email: zhu.1083@osu.edu; gupta.852@osu.edu; canova.1@osu.edu)

Nicola Pivaro is with Bending Spoons, Italy. The work is done at Center for Automotive Research, The Ohio State University. (email: nicolapivaro@gmail.com)

A. Gupta is with Department of Electrical and Computer Engineering, The Ohio State University, Columbus, OH, 43210 USA (email: gupta.706@osu.edu)

by co-optimizing the speed trajectory and the powertrain control strategy [2], [3].

This field of research has experienced significant momentum in the last decade. [3]–[6] address the eco-driving problem for vehicles with single power source, whereas [7]–[9] study the problem with the hybrid electric powertrain architecture. The latter involves modeling multiple power sources and devising optimal control algorithms that can synergistically split the power demand to efficiently utilize the electric energy stored in battery. Maamria et al. [10] systematically compare the computational requirements and the optimality of different formulations. Meanwhile, the difference also exists in the complexity of the driving scenarios. In [4], [9], the eco-driving is formulated without considering the real-time traffic light variability. [3], [6], [8], [11], [12] have explicitly modeled and considered Signal Phase and Timings (SPaTs) and formulate and solve the eco-driving problem with optimal control techniques. In this work, the eco-driving problem of Connected and Automated Hybrid Electric Vehicles (CAHEVs) with the capability of passing traffic lights autonomously is studied.

Recently, the use of Deep Reinforcement Learning (DRL) in the context of eco-driving has caught considerable attention. DRL provides a train-offline, execute-online methodology with which the policy is learned from the historical data or the interaction with simulated environments. Shi et al. [13] modeled the conventional vehicles with ICE as a simplified model and implemented Q-learning to minimize the CO_2 emission at signalized intersections. Li et al. [14] apply an actor-critic algorithm on the ecological ACC problem in car-following mode. Guo and Wang [15] proposed MPC-initialized Proximal Policy Optimization with Model-based Acceleration (PPOMA) for the problem of active signal priority control for trams. Pozzi et al. [16] designed a velocity planner considering the signalized intersection and hybrid powertrain configuration with Deep Deterministic Policy Gradient (DDPG). Zhu et al. [17] formulates the eco-driving problem as a Partially Observable Markov Decision Process (POMDP) and approaches it with PPO. While the strategies with Model-Free Reinforcement Learning (MFRL) in these studies show improvements in the average fuel economy and reductions in onboard computation, the methodology has a fundamental drawback. To teach the agent to drive under complex driving scenarios while satisfying all the constraints from powertrain and traffic rules, a complex and often cumbersome rewarding/penalizing mechanism needs to be designed. Furthermore, under such setup, the agent learns to satisfy constraints by minimizing

the expected cost. For scenarios that are rare yet catastrophic, the scale of the cost penalizing constraint violation needs to be significantly larger than the learning objective itself [13]. As a result, such extrinsic rewarding mechanism increases the design period and deteriorates the final performance.

This paper proposes a safe-critical model-based off-policy reinforcement learning algorithm to solve the eco-driving problem in a connected and automated mild-HEV. The first contribution of this work, from the eco-driving application perspective, is the elimination of the extrinsic reward design in the eco-driving problem by the design of a Model-based Reinforcement Learning (MBRL) algorithm. The algorithm integrates RL with trajectory optimization, which incorporates the constraints from the powertrain dynamics, vehicle dynamics, and traffic rules in a constrained optimization formulation. The performance of the agent is meanwhile improved by using the learned terminal cost function from the RL mechanism.

The second contribution of this work, from the RL algorithm perspective, is the development of Safe Model-based Off-policy Reinforcement Learning (SMORL), a safe-critical model-based off-policy Q-learning algorithm for systems with known dynamics. The algorithm has three unique features compared to the prior and current MBRL implementations. First, SMORL is off-policy as opposed to [18]–[21]. While the use of the model in MBRL increases the sample efficiency [22], the collection of each individual transition becomes more computationally expensive as it commonly requires solving an online optimization problem, as opposed to a feedforward policy in MFRL. With the use of experience replay [23] in the off-policy learning, the historical data can be used to greatly reduce the overall training time. To obtain the value function from Q function, an actor is explicitly trained as in Twin Delayed Deep Deterministic policy gradient algorithm (TD3) [24]. Second, the distributional mismatch between the actor and critic [25] in MBRL is explicitly addressed to improve training performance and stability with Batch Constrained Q-learning (BCQ) [26]. Third, the long-term feasibility of the policy is considered by extending the safe set [20], [27] to a higher dimensional setting using deep unsupervised learning.

The remainder of the paper is organized as follows. Sec. II presents the simulation environment and the eco-driving problem formulation. Sec. III introduces the preliminaries of the mathematical concepts, and Sec. IV presents the main algorithm SMORL. Sec. V explains the detailed implementation of SMORL on the eco-driving problem, and Sec. VI shows the training details and benchmarks the performance.

II. ECO-DRIVING FOR CAHEVS

A. Environment

As collecting data in real-world driving data is expensive and potentially unsafe, a model of the environment is developed for training and validation purposes. The environment model, named EcoSim, consists of a Vehicle Dynamics and Powertrain (VD&PT) model and a microscopic traffic simulator. Fig.1 shows EcoSim and its interaction with the controller and the learning algorithm. The controller commands three control inputs, namely, the Internal Combustion Engine (ICE)

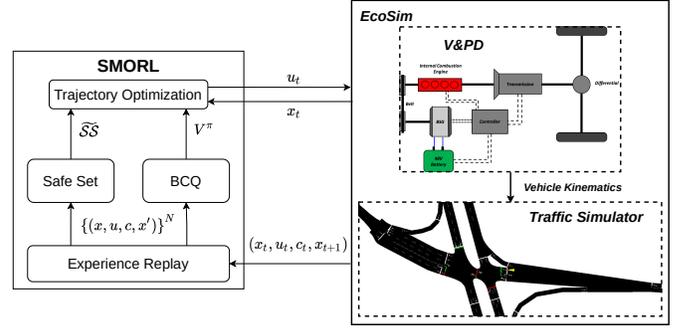


Fig. 1. The Structure of The Environment Model

torque, the electric motor torque and the mechanical brake torque. The component-level torques collectively determine the HEV powertrain dynamics, the longitudinal dynamics of the ego vehicle and its progress along the trip. As in [11], it is assumed that the ego vehicle is equipped with Dedicated Short Range Communication (DSRC) sensors, and SPaTs from signalized intersections become available once they are within the 500 m range. The DRL agent utilizes the SPaT from the upcoming traffic light while ignoring the SPaT from any other traffic light regardless of the availability. Specifically, the controller receives the distance to the upcoming traffic light, its current status and its SPaT program as part of the observation. Finally, a navigation application with Global Positioning System (GPS) is assumed to be available on the vehicle such that the locations of the origin and destination, the remaining distance and the speed limits along the entire trip are available at every point during the trip.

1) *Vehicle and Powertrain Model:* A forward-looking dynamic powertrain model is developed for fuel economy evaluation and control strategy verification. In this work, a P0 mild-hybrid electric vehicle (mHEV) is considered, equipped with a 48V Belted Starter Generator (BSG) performing torque assist, regenerative braking and start-stop functions.

The engine is modeled as low-frequency quasi-static nonlinear maps based on steady-state engine test bench data provided by the supplier. The map of instantaneous fuel consumption \dot{m}_{fuel} is a function of engine angular velocity ω_{eng} and engine torque T_{eng} , and the maps of torque limits T_{eng}^{\min} and T_{eng}^{\max} are functions of engine angular velocity ω_{eng} .

The battery *SoC* and voltage V_{batt} are governed by a zeroth-order equivalent circuit model shown as follows:

$$I_t = \frac{V_{\text{OC}}(\text{SoC}_t) - \sqrt{V_{\text{OC}}^2(\text{SoC}_t) - 4R_0(\text{SoC}_t)P_{\text{bsg},t}}}{2R_0(\text{SoC}_t)}, \quad (1)$$

$$\text{SoC}_{t+1} = \text{SoC}_t - \frac{\Delta t}{C_{\text{nom}}}(I_t + I_{\text{aux}}), \quad (2)$$

where t is the discretized time index, and Δt is the time discretization that is set to be 1s in the study. The power consumed by auxiliaries is modeled by a calibrated constant current bias I_{aux} . The cell open circuit voltage V_{OC} and internal resistance R_0 are maps of *SoC* from a battery pack supplier.

The vehicle dynamics model is based on the road-load equation:

$$v_{veh,t+1} = v_{veh,t} + \Delta t \left(\frac{T_{out,t} - T_{brk,t}}{MR_w} - \frac{C_d \rho_a \Omega_f v_{veh,t}^2}{2M} - g \cos \alpha C_r v_{veh,t} - g \sin \alpha \right) \quad (3)$$

Here, the four terms inside the bracket of the left-hand side are associated with the forward propulsion force, the tire rolling resistance, the aerodynamic drag, and the road grade, respectively. T_{brk} is the brake torque applied on wheel, C_d is the aerodynamic drag coefficient, ρ_a is the air density, A_f is the effective aerodynamic frontal area, C_r is rolling resistance coefficient, and α is the road grade.

Besides the aforementioned models, which are directly associated with either the states or the objective in the eco-driving Optimal Control Problem (OCP) formulation, BSG, torque converter and transmission are also modeled in the study. The BSG is modeled as a quasi-static efficiency map to compute the BSG torque T_{bsg} and power output P_{bsg} . A torque converter model is developed to compute the losses during the traction and regeneration modes. The transmission model is based on a static gearbox, and its efficiency η_{trans} is scheduled as a nonlinear map of the gear number n_g , the transmission input shaft torque T_{trans} and the transmission input speed ω_{trans} . The detailed mathematical models of these components can be found in [17].

The forward vehicle model was calibrated and validated using experimental data from a chassis dynamometer. Vehicle velocity, battery SoC , gear number, engine speed and fuel consumption were used to evaluate the model with the experimental data. Fig. 2 shows the sample results from model verification over the FTP-75 regulatory drive cycle. Results indicate that the vehicle velocity and SoC are accurately predicted by the model. The mismatches in the battery SoC can be attributed to the assumptions made in the simplified battery model such as modeling electrical auxiliary loads as a constant current bias. Further, the final value of the fuel consumption estimated by the model over the FTP-75 drive cycle is within 4% of the actual fuel consumption which verifies that the model can be used for energy and fuel prediction over real-world routes.

2) *Traffic Model*: A large-scale microscopic traffic simulator is developed in an open source software Simulation of Urban Mobility (SUMO) [28] as part of the environment. To recreate realistic mixed urban and highway trips for training, the map of the city of Columbus, OH, US is downloaded from the online database OpenStreetMap [29]. The map contains the length, shape, type and speed limit of the road segments and the detailed program of each traffic light at signalized intersections. Fig. 3 highlights the area that is covered in the study. In the shaded area, 10,000 random passenger car trips, each of which the total distance is randomly distributed from 5 km to 10 km, are generated as the training set. Another 100 trips, indicated by the origins (red markers) and destinations (blue markers) in Fig. 3, are generated following the same distribution as the testing set. In addition, the departure time of each trip follows a geometric distribution with the success

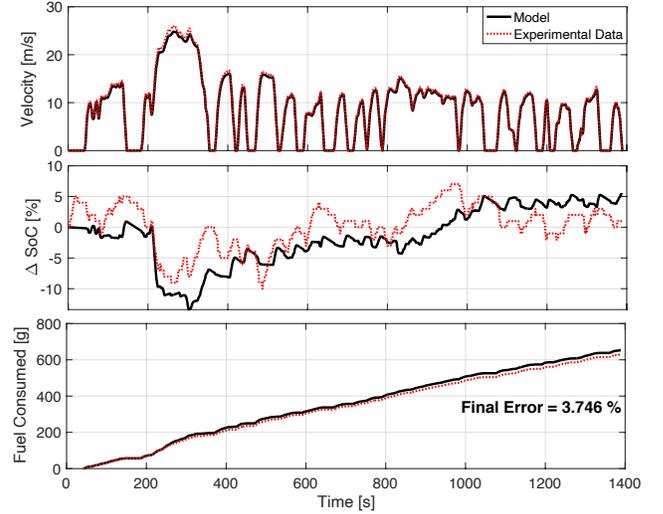


Fig. 2. Validation of Vehicle Velocity, SoC and Fuel Consumed over FTP Cycle.

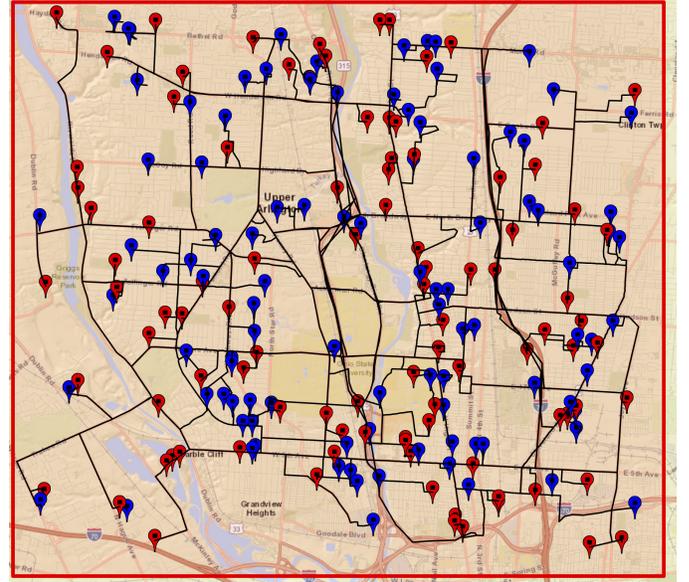


Fig. 3. Map of Columbus, OH as the Traffic Environment for Training

rate $p = 0.01$. The variation among the trips used for training leads to a learned policy that is less subject to local minima and agnostic to specific driving conditions (better generalizability) [30].

B. Optimization Formulation

In the eco-driving problem, the objective is to minimize the weighted sum of fuel consumption and travel time between two designated locations. The optimal control problem is

mathematically formulated as follows:

$$\min_{\{u_t\}_{t=1}^{\infty}} \mathbb{E} \left[\sum_{t=1}^{\infty} [\lambda \dot{m}_{\text{fuel},t} + (1 - \lambda)] \Delta t \cdot \mathbb{I}[s_t < s_{\text{total}}] \right] \quad (4a)$$

$$\text{where } u_t = [T_{\text{eng},t}, T_{\text{bsg},t}, T_{\text{brk},t}]^T \quad (4b)$$

$$\text{s.t. } SoC_{t+1} = f_{\text{batt}}(v_{\text{veh},t}, SoC_t, u_t) \quad (4c)$$

$$v_{\text{veh},t+1} = f_{\text{veh}}(v_{\text{veh},t}, SoC_t, u_t) \quad (4d)$$

$$T_{\text{eng}}^{\min}(\omega_{\text{eng},t}) \leq T_{\text{eng},t} \leq T_{\text{eng}}^{\max}(\omega_{\text{eng},t}) \quad (4e)$$

$$T_{\text{bsg}}^{\min}(\omega_{\text{bsg},t}) \leq T_{\text{bsg},t} \leq T_{\text{bsg}}^{\max}(\omega_{\text{bsg},t}) \quad (4f)$$

$$I^{\min} \leq I_t \leq I^{\max} \quad (4g)$$

$$SoC^{\min} \leq SoC_t \leq SoC^{\max} \quad (4h)$$

$$SoC_T \geq SoC^T \quad (4i)$$

$$0 \leq v_{\text{veh},t} \leq v_{\text{lim},t} \quad (4j)$$

$$(t, s_t) \notin \mathcal{S}_{\text{red}}. \quad (4k)$$

Here, $\dot{m}_{\text{fuel},t}$ is the instantaneous fuel consumption. λ is a normalized weight on the fuel consumption. ω_{eng} and ω_{bsg} are the engine and BSG angular velocities, respectively, and they are static functions of vehicle speed v_{veh} and gear number n_g . f_{batt} and f_{veh} are the battery and vehicle dynamics, respectively, introduced in Sec. II-A1. Eqn. (4e) to (4g) are the constraints imposed by the powertrain components. Eqn. (4h) and Eqn. (4i) are the constraints on the instantaneous battery SoC and terminal SoC for charge sustaining, respectively. Here, the subscript T represents the time at which the vehicle reaches the destination. SoC^{\min} , SoC^{\max} and SoC^T are commonly set to 30%, 80% and 50% [9], [31]. Eqn. (4j) and (4k) are the constraints imposed by traffic conditions. The set \mathcal{S}_{red} represents the set in which the traffic light at the certain location is in red phase [32]. The problem is formulated as an infinite horizon problem in which the stage cost becomes zero once the system reaches the goal set, i.e. the traveled distance s_t is greater than or equal to the total distance of the trip s_{total} while keeping the terminal SoC_T greater than or equal to SoC^T . In addition, anytime that the vehicle violates the traffic light constraints, i.e. Eqn. (4k), the trip is considered a failure and the goal set is not reached.

To solve the aforementioned optimization formulation as an OCP, a MBRL algorithm is proposed, and the preliminaries of the algorithm are included in the next section.

III. PRELIMINARIES ON MBRL

The nonlinear, stochastic, time-invariant system is considered in this work:

$$\begin{aligned} x_{t+1} &= f(x_t, u_t, w_t) \\ x_t &\in \mathcal{X} \subseteq \mathbb{R}^n, t \in \mathbb{N}_+ \\ u_t &\in \mathcal{U}(x_t) \subseteq \mathbb{R}^m, t \in \mathbb{N}_+ \\ w_t &\in \mathcal{W} \subseteq \mathbb{R}^p, t \in \mathbb{N}_+. \end{aligned} \quad (5)$$

Here, x_t , u_t and w_t are the state, control, and uncertainty at time t . \mathcal{X} and \mathcal{U} are the feasible sets for states and inputs, respectively. The uncertainties are assumed to be independent and identically distributed (i.i.d.).

Let $\pi : \mathcal{X} \rightarrow \mathcal{U}$ be a feasible deterministic policy and Π be the set of all feasible deterministic policies. The objective

of the OCP is to reach the goal set $\mathcal{G} \subseteq \mathbb{R}^n$ while finding the optimal policy π^* that minimizes the expectation of the discounted sum of the costs defined as follows:

$$\begin{aligned} \pi^* &= \operatorname{argmin}_{\pi \in \Pi} \eta(\pi), \text{ where} \\ \eta(\pi) &= \mathbb{E}_{w_t} \left[\sum_{t=0}^{\infty} \gamma^t c(x_t, u_t) \right], \\ \text{where } u_t &= \pi(x_t). \end{aligned} \quad (6)$$

Here, γ is the discount factor that prioritizes the immediate rewards and ensures the sum over the infinite horizon remains finite.

As in [21], the following assumption is made.

Assumption III.1 (Costs). *The cost is zero for the states inside the goal set \mathcal{G} and positive for the states outside, i.e. $\exists \epsilon > 0$ such that $c(x, u) > \epsilon \mathbb{I}_{\mathcal{G}^C(x)}$ where \mathbb{I} is the indicator function and \mathcal{G}^C is the complement of the goal set \mathcal{G} .*

As in [21], [27], [33], the following definitions are given.

Definition III.1 (Robust Control Invariant Set). *A set $\mathcal{C} \subseteq \mathcal{X}$ is said to be a robust control invariant set for the system Eqn. (5) if for all $x(t) \in \mathcal{C}$, there exists a $u(t) \in \mathcal{U}$ such that $f(x(t), u(t), w(t)) \in \mathcal{C}$, for all $w(t) \in \mathcal{W}$ and $t \in \mathbb{N}_+$.*

Definition III.2 (Robust Successor Set $Suc(\mathcal{S})$). *For a given set \mathcal{S} , its robust successor set $Suc(\mathcal{S})$ is defined as*

$$\begin{aligned} Suc(\mathcal{S}) &= \{x' \in \mathbb{R}^n : \exists x \in \mathcal{S}, \exists w \in \mathcal{W} \\ &\text{such that } x' = f(x, \pi(x), w)\}. \end{aligned} \quad (7)$$

Definition III.3 (Robust Reachable Set $\mathcal{R}_N(x_0^j)$). *For a given initial state x_0^j , the N -step robust reachable set $\mathcal{R}_N(x_0^j)$ of the system defined in Eqn. (5) in a closed loop policy π at iteration j is defined recursively as*

$$\begin{aligned} \mathcal{R}_{i+1}^{\pi}(x_0^j) &= Suc(\mathcal{R}_i^{\pi}(x_0^j)) \cap \mathcal{X}, \\ \mathcal{R}_0^{\pi}(x_0^j) &= x_0^j, \end{aligned} \quad (8)$$

where $i = 0, 1, \dots, N - 1$.

Definition III.4 (Safe Set). *The safe set \mathcal{SS}^j contains the full evolution of the system at iteration j ,*

$$\mathcal{SS}^j = \left\{ \bigcup_{k=0}^{\infty} \mathcal{R}_k^{\pi}(x_0^j) \bigcup \mathcal{G} \right\}. \quad (9)$$

As shown in [27], the exact form of the safe set in 9 is a robust control invariant set. As calculating its exact form is intractable, especially for high dimensional nonlinear system, it is, in practice, approximated as

$$\widetilde{\mathcal{SS}}^j = \bigcup_{k \in \mathcal{M}^j} x^k, \quad (10)$$

where $x^k = \{x_t^k : t \in \mathbb{N}_+\}$ is the trajectory at iteration k , and $\mathcal{M}^j = \{k \in [0, j) : \lim_{t \rightarrow \infty} x_t^k \in \mathcal{G}\}$ is the set of indices of which the trajectories were successfully driven to the goal. As the safe set in this work is constantly evolving during training, the iteration index j will be neglected in the remaining work.

For any policy π , the value function $V^\pi : \mathcal{X} \rightarrow \mathbb{R}$, the Q function $Q^\pi : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ and the advantage function $A^\pi : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ are defined as follows:

$$V^\pi(x_t) = \begin{cases} \mathbb{E}_\pi \left[\sum_{i=t}^{\infty} \gamma^{i-t} c(x_i, u_i) \mid x_t \right], & x_t \in \mathcal{SS} \\ \infty, & \text{otherwise.} \end{cases} \quad (11)$$

$$Q^\pi(x_t, u_t) = \begin{cases} \mathbb{E}_\pi \left[\sum_{i=t}^{\infty} \gamma^{i-t} c(x_i, u_i) \mid x_t, u_t \right], & \begin{matrix} x_t \in \mathcal{SS} \\ u_t \in \mathcal{U} \end{matrix} \\ \infty & \text{otherwise.} \end{cases} \quad (12)$$

$$A^\pi(x_t, u_t) = Q^\pi(x_t, u_t) - V^\pi(x_t). \quad (13)$$

IV. PROPOSED METHOD

In this work, an off-policy model-based deep reinforcement learning algorithm with an approximated safe set is proposed. At any given time t during policy execution, the following trajectory optimization problem with a receding horizon of H steps is solved:

$$\begin{aligned} \min_{\{\tilde{u}_k\}_{k=t}^{t+H-1}} \mathbb{E} & \left[\sum_{k=t}^{t+H-1} \gamma^{k-t} c(\tilde{x}_k, \tilde{u}_k) + \gamma^H V^\pi(\tilde{x}_{t+H}) \right] \\ \text{s.t. } & \tilde{x}_{k+1} = f(\tilde{x}_k, \tilde{u}_k, w_k) \\ & \tilde{x}_t = x_t \\ & \tilde{x}_k \in \mathcal{X}, k = t, \dots, t+H-1 \\ & \tilde{x}_{t+H} \in \widetilde{\mathcal{SS}} \\ & \tilde{u}_k \in \mathcal{U}, k = t, \dots, t+H-1, \end{aligned} \quad (14)$$

where \tilde{x} and \tilde{u} are the variables for states and control actions in the predicted trajectory. Compared to the formulation in [18], the state \tilde{x}_k and action \tilde{u}_k are explicitly constrained to be within the feasible region in the receding horizon and the terminal state \tilde{x}_{t+H} to be within the safe set $\widetilde{\mathcal{SS}}$. With the presence of uncertainties in the dynamic system, solving the exact form of the above stochastic optimization problem can be challenging. In [20], [21], [34], Cross Entropy Method (CEM) [34] is used to solve the problem with unknown dynamics as a chance constraint problem. In Section II-B, techniques will be discussed to simplify and solve the optimization in the eco-driving problem.

As most of the model-based deep reinforcement learning methods with trajectory optimization in literature learn the value function as the terminal cost for the MPC [18], [19], [21], [34], the learning algorithm becomes on-policy. While the trajectory optimization increases the sample efficiency and helps exploration [18], solving the trajectory optimization problem makes each data sample more computationally expensive. As a result, the training wall time is not necessarily reduced. In this work, the off-policy Q-learning [35] is instead proposed. To use the learned Q function in trajectory optimization, the following equation needs to be solved:

$$\min_{\{\tilde{u}_k\}_{k=t}^{t+H-1}} \mathbb{E} \left[\sum_{k=t}^{t+H-1} \gamma^{k-t} c(\tilde{x}_k, \tilde{u}_k) + \gamma^H Q_\theta^\pi(\tilde{x}_{t+H}, \tilde{u}_{t+H}) \right], \quad (15)$$

where Q_θ^π is the approximated Q function parametrized by θ . Compared to solving Eqn.(14), solving Eqn. (15) requires one extra computational step:

$$V^\pi(\tilde{x}_{t+H}) = \min_{\tilde{u}_{t+H}} Q_\theta^\pi(\tilde{x}_{t+H}, \tilde{u}_{t+H}). \quad (16)$$

Depending on the dimension of the problem, solving Eqn. (16) can be computationally intractable, especially for online control. Several algorithms, e.g. DDPG [36], TD3 [24] and dueling network [37] are proposed to obtain the value function from the Q function. In this work, the off-policy actor-critic algorithm TD3 is used since it reduces the overestimation and is shown to be more stable than DDPG. Specifically, with the sample (x_j, u_j, c_j, x'_j) from the experience replay buffer \mathcal{D} [23], the target for the Q function during training is constructed as follows,

$$y_j = c_j + \gamma \max_{i=1,2} Q_{\theta'_i}(x'_j, u'_j), \quad (17)$$

$$u'_j = \pi_{\phi'}(x'_j). \quad (18)$$

Here, Q_{θ_1} and Q_{θ_2} are two independently trained critic networks. $Q_{\theta'_1}$ and $Q_{\theta'_2}$ are the corresponding target networks. π_ϕ and $\pi_{\phi'}$ are the actor network and its target network, respectively. The critics are then updated following

$$\begin{aligned} \theta_i \leftarrow \theta_i - \alpha \nabla_{\theta_i} \left[\frac{1}{N} \sum_{j=1}^N (y_j - Q_{\theta_i}(x_j, u_j))^2 \right], \quad i = 1, 2, \\ \{(x_j, u_j, c_j, x'_j) \sim \mathcal{D}\}_{j=1}^N. \end{aligned} \quad (19)$$

where α is the learning rate, and N is the batch size.

In the off-policy learning algorithm used here, the behavior policy is the trajectory optimization where state and action constraints within the receding horizon are satisfied thanks to the constrained optimization formulation. However, the trained actor π_ϕ makes decisions solely based on the Q function. The resulting mismatch between the distribution of state-action pairs induced by the actor π_ϕ and that collected by the behavior policy results in extrapolation error leading to unstable training [25]. In the eco-driving problem, the trajectory optimization ensures the power is solely generated from ICE when the battery SoC is at the lower limit SoC^{\min} . Accordingly, no state-action pair resembling low SoC and high motor torque can be collected, which leads to extrapolation in the Q function near the region. The error can eventually cause unstable training or inferior performance.

To address the extrapolation error induced by the mismatch in distributions, Batch Constrained Q-learning (BCQ) [26] originally proposed for offline reinforcement learning is used. Here, a generative model, specifically a Variational Autoencoder (VAE) [38], $G_\omega(x)$ is trained to resemble the state-action distribution in the experience replay buffer. The background on VAE and the training objective are covered in Appendix A. Note that samples from the generative model $a' \sim G_\omega(x')$ should ideally match the distribution collected by the behavior

policy. Instead of selecting action following Eqn. (18), the action is now selected as

$$u'_j = \underset{u_{j,k} + \xi_\phi(x_j, u_{j,k}, \Phi)}{\operatorname{argmin}} \left[\max_{i=1,2} Q_{\theta'_i}(x'_j, u_{j,k} + \xi_\phi(x_j, u_{j,k}, \Phi)) \right],$$

$$\{u_{j,k} \sim G_\omega(x'_j)\}_{k=1}^n. \quad (20)$$

Here n is the hyperparameter that is the number of actions sampled from the generative model. The action u' used for the target value for the Q function is selected as the best among the n sampled ones. Note that there is no longer an actor network mapping from state to action. Instead, to ensure the agent can learn on top of the actions sampled from the generative model imitating the behavior policy from the experience buffer, a perturbation network ξ_ϕ whose output is clipped between $[-\Phi, \Phi]$ is trained. The perturbation network ξ_ϕ is updated by deterministic policy gradient theorem from [39] as

$$\phi \leftarrow \phi - \alpha \nabla_\phi \left[\frac{1}{N} \sum_{j=1}^N Q_{\theta_1}(x_j, u_j + \xi_\phi(x_j, u_j, \Phi)) \right]. \quad (21)$$

To reduce the accumulating error from bootstrapping, all the target networks are updated with a slower rates as

$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \quad i = 1, 2, \quad (22)$$

$$\phi' \leftarrow \tau \phi + (1 - \tau) \phi', \quad (23)$$

where τ is a constant on the order of 10^{-3} to 10^{-1} .

In Eqn. (14), the terminal rollout state is constrained within the safe set. Since the safe set is an approximation to the robust control invariant set, the constrain ensures that there exists policies that can safely drive the terminal state to the goal set. In [20], [21], the safe set is approximated by kernel density estimation, which typically works well only for problems in low dimensions. Here, we extend the approximation to high-dimensional setting by using deep generative models. Following the notion in [20], the safe set is approximated as

$$\widetilde{\mathcal{S}} = \{x : p_\psi(x) \geq \delta\}, \quad (24)$$

where $p_\psi : \mathcal{X} \rightarrow [0, 1]$ is the probability that a state is inside the safe set parametrized by ψ , and the constant δ regulates how exploratory the controller is. Note that the generative model used for the safe set approximation needs to model the probability explicitly and can be slow in sampling, whereas the generative model resembling the distribution of state-action pairs in the experience replay needs to be fast in sampling while the explicit probability is not required.

Due to the aforementioned consideration, the autoregressive model with Long Short-term Memory (LSTM) [40] is used. The description of the model as well as the training objective is included in Appendix B. In Sec.V, the use of the autoregressive model in the application of eco-driving is motivated as the dimension of the problem can get large once the future conditions are sampled discretely.

In summary, Safe Model-based Off-policy Reinforcement Learning (SMORL) is proposed. The algorithm builds on SAVED [20] and extends it to be an off-policy algorithm with the methods proposed in BCQ. The detailed step-by-step algorithm is included in Algorithm 1.

V. IMPLEMENTATION DETAILS

A. Trajectory Optimization

Specific to the eco-driving problem, the state vector x_t is defined as a vector with 88 states. A description of the states are listed in Tab. I. Here, the first seven elements of the state vector are the battery SoC , the vehicle speed v_{veh} , the current speed limit v_{lim} , the next speed limit v'_{lim} , the distance to the next speed limit s_{lim} , the distance to the upcoming traffic light s_{tls} and the total remaining distance s_{rem} . The remaining 81 elements are the sampled upcoming traffic light status in the next 80 seconds x_{tfc} . For example, if the upcoming traffic light has 20 seconds remaining for the current red phase and will remain in green for the rest of the 80 seconds, the first 21 elements of the sampled upcoming traffic light status are 0, and the rest are set to 1. Compared to the manually extracted feature representation in [17], the sampled representation reduces the discontinuity and results in a better performance.

As the vehicle considered in this study is assumed equipped with connected features, e.g. advanced mapping and V2I connectivity, and surrounding vehicles are not included in the study, it is assumed that the ego vehicle can deterministically predict the uncertainties from driving conditions within the receding horizon in this study. Sun et al. [6] suggest by formulating the problem as a chance constraint or a distributionally robust optimization problem, uncertainties in SPaT can be considered without additional computational load.

In Eqn. (4), the receding horizon H is in the time domain. While it is easier to incorporate the time-based information such as SPaT received from V2I communication in the time domain, an iterative dynamic look-ahead process is required to process any distance-based route feature, such as speed limits, grade, traffic light and stop sign locations. For example, the controller requires the speed limits as the constraints to generate speed trajectory while the speed limits can change based on the distance traveled by the speed trajectory. In this study, the value and Q functions are learned in the time domain for ease of integration with the time-based traffic simulator, while the trajectory optimization is conducted in the spatial domain.

As SPaTs and speed limits do not depend on the decision made by the ego vehicle in the spatial domain, they are incorporated into the optimization problem as constraints, and only the vehicle speed, battery SoC and the time at which the vehicle reaches the given distance are considered as the state in the trajectory optimization. Define the optimization state $z \in \mathcal{Z} \subseteq \mathbb{R}^3$ as

$$z_s = [v_{veh,s}, SoC_s, t_s]^T. \quad (25)$$

Here, s is the index in the discretized spatial domain with $\Delta s = 10m$, and the dynamics of z in the time and spatial domains are converted following

$$\frac{\Delta z}{\Delta s} = \frac{\Delta z}{\Delta t} \frac{\Delta t}{\Delta s} = \frac{\Delta z}{\Delta t} \frac{1}{v_{veh}}. \quad (26)$$

Algorithm 1: Safe Model-based Off-policy Reinforcement Learning (SMORL)

Initialize Q-networks $Q_{\theta_1}, Q_{\theta_2}$ independently, and duplicate target networks $Q_{\theta'_1}, Q_{\theta'_2}$.

Initialize the perturbation network ξ_ϕ , its target network ξ'_ϕ and VAE $G_\omega = \{E_{\omega_1}, D_{\omega_2}\}$.

Initialize the experience replay buffer \mathcal{D} .

Collect N_0 successfully executed trajectories with a baseline controller and initialize the safe set $\widetilde{\mathcal{SS}}$.

for $n_{iter} \in 1, \dots, N_{iter}$ **do**

while j^{th} trajectory NOT finished **do**

 Select control action u_t by solving trajectory optimization in Eqn. (14).

 Sample mini-batch of N transitions (x, u, c, x') from \mathcal{D} .

 For each transition, sample n actions u'_j from $G_\omega(x')$ and n perturbations from $\xi_\phi(x', u', \Phi)$.

 Update the critic networks Q_{θ_1} , the target networks Q_{θ_2} following Eqn. (19) and $Q_{\theta'_1}, Q_{\theta'_2}$ following Eqn. (22).

 Update perturbation network ξ_ϕ following Eqn. (21).

 Update VAE G_ω by maximizing Eqn. (30).

end

if $x_T \in \mathcal{G}$ **then**

 Push the trajectory $\{(x_t, u_t, c_t, x_{t+1})\}^T$ to \mathcal{D} .

 Update the safe set $\widetilde{\mathcal{SS}}$ with minibatches sampled from \mathcal{D} following Eqn. (32).

end

end

TABLE I

THE STATE AND ACTION SPACES OF THE ECO-DRIVING PROBLEM

	Variable	Description
\mathcal{X}	$SoC \in \mathbb{R}$	Battery SoC
	$v_{veh} \in \mathbb{R}$	Vehicle velocity
	$v_{lim} \in \mathbb{R}$	Speed limit at the current road segment
	$v'_{lim} \in \mathbb{R}$	Upcoming speed limit
	$d_{tfc} \in \mathbb{R}$	Distance to the upcoming traffic light
	$d'_{lim} \in \mathbb{R}$	Distance to the road segment of which the speed limit changes
	$d_{rem} \in \mathbb{R}$	Remaining distance of the trip
	$x_{tfc} \in \{0, 1\}^{S1}$	Sampled status of the upcoming traffic light
\mathcal{U}	$T_{eng} \in \mathbb{R}$	Engine torque
	$T_{bsg} \in \mathbb{R}$	Motor torque
	$T_{brk} \in \mathbb{R}$	Equivalent brake torque

As a result, the trajectory optimization is formulated as

$$\min_{\{\tilde{u}\}_{k=s_t}^{s_t+H_s-1}} \sum_{k=s_t}^{s_t+H_s-1} \gamma^{t_k} c(\tilde{z}_k, \tilde{u}_k) + \gamma^{t_{H_s}} V^\pi(\mathcal{G}(x_t, z_{H_s})) \quad (27a)$$

where:

$$c(\tilde{x}_k, \tilde{u}_k) = (\lambda \dot{m}_{fuel,k} + (1 - \lambda)) \frac{\Delta s}{v_{veh,k}} \cdot \mathbb{I}[s_k < s_{total}] \quad (27b)$$

$$\text{s.t. } SoC_{k+1} = f_{batt,s}(\tilde{z}_k, \tilde{u}_k) \quad (27c)$$

$$v_{veh,k+1} = f_{veh,s}(\tilde{z}_k, \tilde{u}_k) \quad (27d)$$

$$T_{eng}^{\min}(\omega_{eng,k}) \leq T_{eng,k} \leq T_{eng}^{\max}(\omega_{eng,k}) \quad (27e)$$

$$T_{bsg}^{\min}(\omega_{bsg,k}) \leq T_{bsg,k} \leq T_{bsg}^{\max}(\omega_{bsg,k}) \quad (27f)$$

$$I^{\min} \leq I_k \leq I^{\max} \quad (27g)$$

$$SoC^{\min} \leq SoC_k \leq SoC^{\max} \quad (27h)$$

$$0 \leq v_{veh,k} \leq v_{lim,k} \quad (27i)$$

$$(t_k, s_k) \notin \mathcal{S}_{red} \quad (27j)$$

$$\mathcal{G}(x_t, z_{H_s}) \in \widetilde{\mathcal{SS}}. \quad (27k)$$

Here, s_t is the spatial index corresponding to the distance the ego vehicle has traveled at the time t . $H_s = 20$ is the prediction step in the spatial domain, making the total prediction horizon 200 m . $\mathcal{G} : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{X}$ is the function that takes the full state x_t and the terminal optimization state z_{H_s} and determines the predicted terminal full state $\tilde{x}_{t+t_{H_s}}$. For example, suppose there are 15 seconds left in the current green phase and t_{H_s} in the optimization state is 10 seconds, i.e. it takes 10 seconds for the ego vehicle to travel the future 200 m , there will be 5 seconds left in the current green phase at the end of the prediction horizon. The trajectory optimization problem is solved by Deterministic Dynamic Programming (DDP) [41]. The optimal deterministic policy $\mu_k^* : \mathcal{Z} \rightarrow \mathcal{U}$, $k = 1, 2, \dots, H_s - 1$, along with the optimal cost-to-go function $\mathcal{J}_k : \mathcal{Z} \rightarrow \mathbb{R}$, $k = 1, 2, \dots, H_s$ can be calculated through backward recursion as

$$\mathcal{J}_{H_s}(z) = V^\pi(\mathcal{G}(x_t, z)) + \mathcal{P}_N(z), \quad (28a)$$

$$\mathcal{F}_k(z, u) = c(z, u) + \mathcal{P}_k(z) + \mathcal{J}_{k+1}(f_k(x, u)), \quad (28b)$$

$$\mu_k^* = \underset{\mu_k}{\operatorname{argmin}} \mathcal{F}_k(z, \mu_k(z)), \quad (28c)$$

$$\mathcal{J}_k(z) = \mathcal{F}_k(z, \mu_k^*(z)). \quad (28d)$$

Here, $\mathcal{F} : \mathcal{Z} \times \mathcal{U} \rightarrow \mathbb{R}$ is the cost-to-go associated with a given immediate action and then the optimal policy. $\mathcal{P}_k : \mathcal{Z} \rightarrow \mathbb{R}$ and $\mathcal{P}_N : \mathcal{Z} \rightarrow \mathbb{R}$ are penalty functions introduced to ensure no constraint violation in the predicted trajectory.

Solving Eqn. (28) is computationally intensive yet highly parallelizable. Considering onboard GPU is readily available nowadays on self-driving vehicles, a real-time capable CUDA-based Parallel DDP (PDDP) solver in [32] is used in this work. In the cases where the stochasticity within the prediction horizon cannot be ignored, other gradient-free optimization methods, such as CEM or random shooting method [42], can be used as the trajectory optimizer.

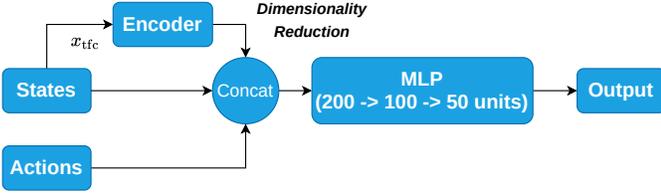


Fig. 4. The Network Architecture of the Value and Q Functions.

B. Q Learning

Fig. 4 shows the architecture of the neural network associated with the Q-learning. Upon receiving the state vector, the sampled traffic light status w_{tfc} are fed to a pre-trained autoencoder with Multilayer Perceptron (MLP) of size (81, 100, 5, 100, 81) for dimensionality reduction. The remaining states along with the actions are concatenated with the latent states from the encoder and subsequently fed into another MLP of size (200, 100, 50) to output the Q function for critic and the perturbation for the actor. The critic and actor do not share parameters in this work.

To accelerate the training and improve generalizability, the state of the vehicle is randomized for every 50 steps in simulation. When the domain randomization occurs, the battery SoC and the vehicle velocity $v_{\text{veh},t}$ are sampled from uniform distributions $\text{Uniform}(SoC^{\min}, SoC^{\max})$ and $\text{Uniform}(0, v_{\text{lim},t})$, respectively. To guarantee feasibility during the trip, the domain randomization is disabled $200m$ within signalized intersections or $1000m$ within the final destination.

C. Safe Set Approximation

In the eco-driving problem, two types of constraints can induce feasibility issues, namely, the battery terminal SoC constraint and the constraint imposed by traffic rules at signalized intersections. For the first case, the goal set is considered not reached when the vehicle is near the destination and it cannot sufficiently charge the battery back to SoC^T in the remaining distance. For the second case, the trip is considered failed when the vehicle breaks traffic rules at signalized intersections, and infeasibility occurs when the vehicle speed is too high and there is not enough distance to brake to stop in front of a traffic light in red or stop sign. A conservative low-speed controller that only uses ICE is used to collect the initial data for the experience replay buffer. During training, only samples from the trips that reach the goal set without violating any constraints are added to the experience replay buffer.

In the eco-driving problem, the sampled traffic light x_{tfc} is binary, while the other variables are continuous. As the PDDP solver also discretizes the continuous state space, we consider the loss of accuracy with the same discretization is acceptable. The discretized states as in one-hot form in each dimension are fed into an LSTM network with 50 units sequentially as shown in Fig. 9. The outputs from LSTM are then masked according to the number of categorical classes in each dimension. Finally, the softmax operator ensures the outputs to be a proper conditional probability distribution.

As an alternative to LSTM, Causal 1D Convolutional Neural Networks (Causal Conv1D) [43] was also implemented as the network for the autoregressive model. The key difference is that the states in all dimensions can be fed into Causal Conv1D in parallel, whereas each dimension needs to be fed into LSTM sequentially. For applications with long sequences, Causal Conv1D can be more efficient and accurate [44]. For the specific problem, Causal Conv1D shows no noticeable advantage over LSTM in either accuracy or inference speed. As a result, LSTM is chosen as it has fewer hyperparameters.

When the receding horizon ($200m$) in this study is longer than the critical braking distance [17], the vehicle will never violate any constraints imposed by signalized intersections. Nevertheless, using the safe set to constrain the terminal state is still essential for the following reason. At the last step of the receding horizon, the value function $V^\pi(\tilde{x}_{t+\Delta t_{H_s}})$ needs to be evaluated numerically for trajectory optimization. Since only the data from the safely executed trips are added into the buffer and there is no penalty mechanism for the constraint violation, the estimation of the critic network is valid only within the safe set and is subject to extrapolation error outside. Although the long receding horizon ensures the feasibility of the actual trajectory regardless of the use of a safe set, the training is subject to instability and the learned performance can significantly deteriorate without the constraint from the safe set.

This effect is shown in Fig. 5. Here, the two subplots on top show the optimized trajectories with and without the use of a safe set, respectively. The three curves in each plot are the trajectories from the optimizer at three consecutive seconds. During the first two seconds, the vehicle is more than $200m$ away from the traffic light, and thus the constraint from Eqn. (4k) is not considered in the trajectory optimization. The subplots on the bottom show the safety status of the terminal state in the dimension of the vehicle velocity and the time at which it reaches the end of the receding horizon before the signalized intersection appears in the receding horizon. Here, green means the state is considered safe, i.e. $p_\psi(x) \geq \delta$, and red otherwise. Although the actual trajectories, with or without a safe set, can slow down in time to avoid trespassing the red light thanks to the sufficiently long receding horizon, the terminal state without the safe set constraint has a speed of $20m/s$ with $20m$ left before a red light, which is unsafe. Meanwhile, comparing the bottom two subplots, the terminal velocity constrained by the safe set progressively reduces as the vehicle approaches the intersection in the red phase. In addition, given speed limit here is set to $v_{\text{lim}} = 22m/s$, any state with a velocity higher than $22m/s$ is considered unsafe. It can be noticed that the red region on the top right corners is incorrectly considered unsafe (false positive). This is because, by optimality, the agent rarely crosses an intersection in the green phase with low speed, therefore, these false positive regions do not affect the performance.

As a summary for all the implementation details, the hyperparameters are listed in Tab. II.

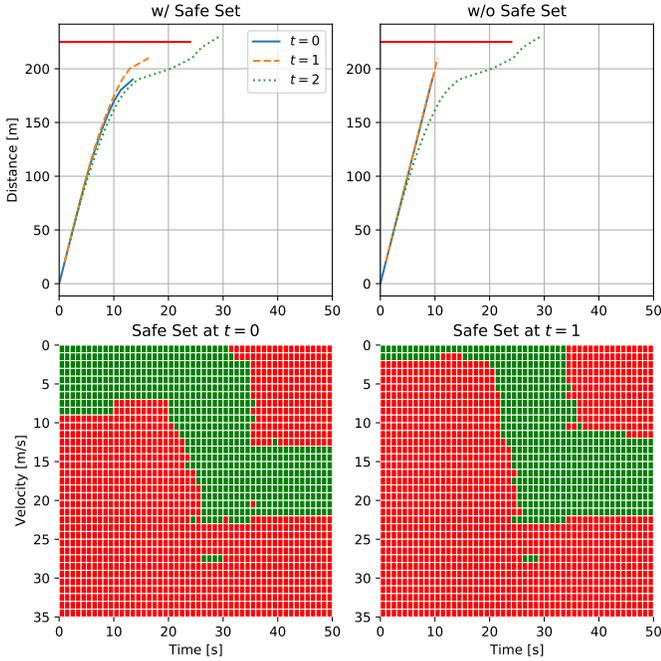


Fig. 5. The Effect Of the Safe Set on Trajectory Optimization.

TABLE II
HYPERPARAMETERS OF THE Q LEARNING AND SAFE SET ESTIMATION

Parameter	Value
Weighting factor between fuel and time, λ	0.45
Discount factor, γ	0.995
Optimizer	Adam
Learning rate, α	$1e-4$
Experience buffer size	$2e5$
Batch size, N	256
Target network update rate, τ	$1e-3$
Exploration rate, ϵ	0.2
Perturbation range in physical unit, Φ	30 Nm
Sampled actions from the VAE decoder, n	10
Steps per domain randomization	50
LSTM size for the safe set	50

VI. RESULTS

Both the PDDP optimizer and the neural network training require GPU. To get the results to be shown, the training took 24 hours on a node with an NVIDIA Volta V100 GPU and 2.4GHz Intel CPU from Ohio Supercomputer Center [45]. As domain randomization is used during training, 5 trips out of 1000 randomly generated trips are repeatedly selected for every 25 training episodes to evaluate the performance of the controller and to quantify the progress of training. During the evaluation, domain randomization and epsilon greedy are both deactivated. Fig. 6 shows the evolution of the total costs, fuel economy and average speed of the 5 evaluation trips. Compared to the model-free on-policy method in [17], which takes 80,000 episodes to converge, the sample efficiency of the off-policy model-based method is significantly improved. In the meantime, with the constrained optimization formulation and the safe set, the training quickly learns to respect the constraints imposed by the terminal SoC and signalized intersections. Furthermore, the fact that the agent does not need any

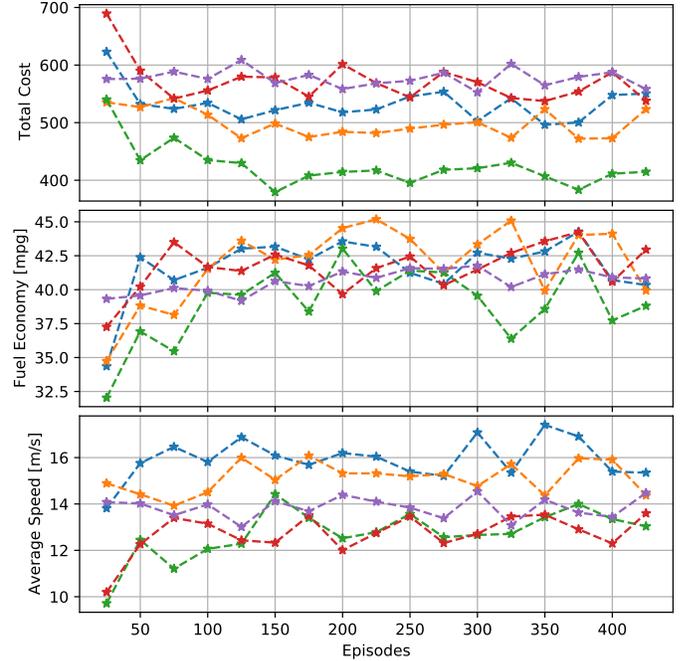


Fig. 6. The Evolution of the Total Costs, Fuel Economy, and Average Speed of the 5 Evaluation Trips.

TABLE III
FUEL ECONOMY, AVERAGE SPEED AND SOC VARIANCE FOR BASELINE, MODEL-FREE DRL, SMORL AND WS SOLUTIONS

	Baseline	ADP	MFDRL	SMORL	WS
Fuel Economy mpg	32.4	39.5	40.8	41.6	47.5
Speed Mean m/s	14.1	13.9	12.5	14.0	14.5
SoC Variance % ²	12.1	21.6	18.2	52.6	22.6

extrinsic penalty from constraint violation and is still capable of learning to operate within the safe region significantly simplifies the design and tuning process as deployed in the previous reinforcement learning attempts on eco-driving [14], [16], [17].

Statistically, the performance of the agent trained with SMORL is compared against the four other strategies, a baseline strategy with Enhanced Driver Model (EDM) representing a human driver and heuristically calibrated energy management module [46], the hierarchical optimal controller using Approximate Dynamic Programming (ADP) in [31], the model-free DRL (MFDRL) agent proposed in [17] and the wait-and-see (WS) solution. The WS solution assumes the speed limits and the sequences of all traffic lights of the entire trip are known *a priori*, and it is solved by PDDP solver as well. Despite being non-causal and computational intractable for online implementation, the solution serves as the upper bound for the causal control strategies. The four strategies are evaluated on the 100 random trips as shown in Fig. 3, and the fuel economy, average speed and variance of the battery SoC are listed in Tab. III.

Here, compared to the baseline strategy, the SMORL agent consumes 21.8% less total fuels while maintaining a compa-

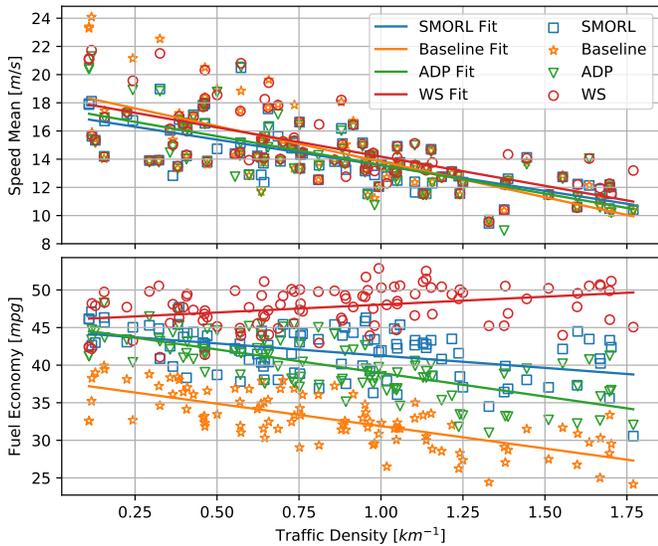


Fig. 7. The Variation of the Average Speed and the Fuel Economy against Traffic Light Density for Baseline, SMORL and WS Solution

erable average speed. The benefit in fuel economy is achieved by avoiding unnecessary acceleration events and by taking advantage of a wider range of battery capacity as indicated by the higher *SoC* variance. The SMORL agent shows dominant performance in average speed and fuel economy compared to the previously trained MFDRL strategy and the non-learning-based ADP strategy. The performance improvement over MFDRL is primarily due to two aspects. First, the online trajectory optimization solved by PDDP guarantees the global optimality within the receding horizon, which is more accurate and reliable than the actions generated from the one-step stochastic policy from neural networks as in MFDRL. Second, the fact that there is no extrinsic penalty to assist constraint satisfaction ensures that the agent focuses on learning only the objective of the OCP formulation, i.e. weighted sum of the trip time and fuel consumption, instead of a carefully designed yet delicate surrogate learning objective.

In Fig. 7, the average vehicle speed and the fuel economy of each trip are plotted against the traffic light density. As the WS solution calculates the global optimal solution with the knowledge of the full trip, it is able to navigate among the traffic lights accordingly, as indicated by the surprisingly increasing fuel economy. This can be due to the fact that when there are more traffic lights, the vehicle is forced to operate with a lower speed and lower fuel consumption condition. On the other hand, as the baseline driver has limited line-of-sight [46] and the ADP, MFDRL and SMORL controllers have limited DSRC sensing range, the fuel economy decreases as the traffic light density increases. Nevertheless, as indicated by the slope of the fitted curve, the fuel economy of SMORL is less affected by the increase of the traffic light density compared to the baseline and ADP controller.

Fig. 8 shows the comparison among the baseline, ADP, SMORL and the WS solution on a specific testing trip. For this specific trip, while the differences in trip time are within 3s, SMORL consumes 24.7% and 11.0% less fuel compared

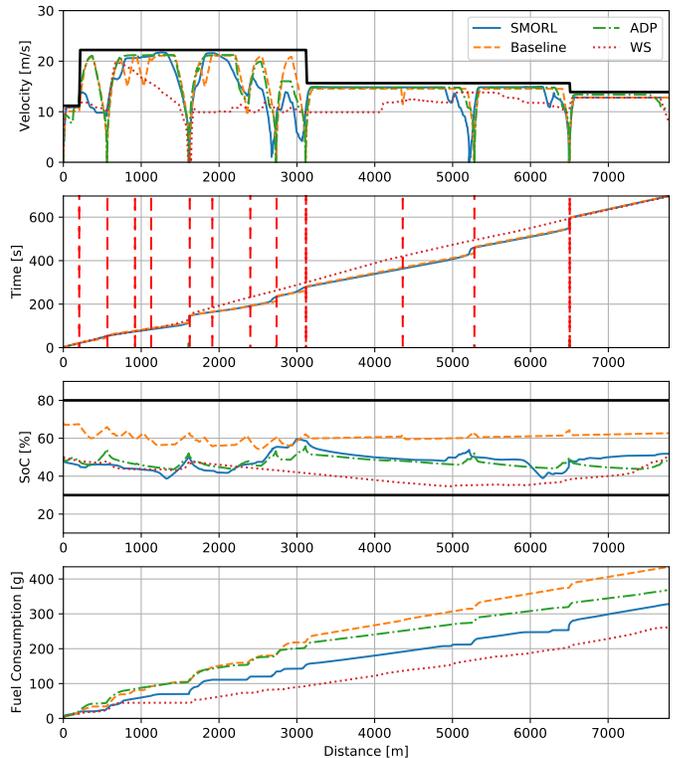


Fig. 8. The Trajectory Comparison among Baseline, SMORL and WS.

to the baseline and the ADP strategies, respectively. While SMORL demonstrates some merits similar to the WS solution, its inferiority to the WS solution is primarily due to the fact that only the SPaTs from the upcoming intersection are available to the controller. Additional comparisons among the four strategies can be found in Appendix C. The ablation study for the key components in the algorithm is shown in Appendix D.

VII. CONCLUSION

In this paper, a safe-critical model-based off-policy reinforcement learning algorithm SMORL is proposed. The algorithm is applied to the eco-driving problem for CAHEV. Compared to the previous model-free attempts on eco-driving in the literature, the method does not require any extrinsic rewarding mechanism, and thus, greatly simplifies the design process and improves the final performance. With the online constrained optimization formulation and the approximate safe set, the learned strategy is capable of satisfying the constraints in the prediction horizon and restricting the state within the approximate safe set, which is an approximation to the robust control invariant set. The performance of the strategy trained with SMORL is compared to a baseline strategy representing human drivers' behavior over 100 randomly generated trips in Columbus, OH, US. With a comparable average speed, the strategy from SMORL consumes approximately 22% less fuel.

While the demonstration of the algorithm is on the eco-driving problem, we believe it can be applied to many other real-world problems, in particular to those with well-studied system dynamics, such as robotics and autonomous driving.

Future studies include the extending the SMORL algorithm to include the presence of leading vehicles, as well as the integration and verification of the algorithm in a demonstration vehicle.

APPENDIX A VARIATIONAL AUTOENCODER

Let $X = \{x_i\}_{i=1}^N$ be some data set and Z represent a set of low-dimensional latent variables, the objective is to maximize the marginal log-likelihood:

$$\begin{aligned} \log p(X) &= \sum_{i=1}^N \log p(x_i) = \sum_{i=1}^N \int \log p(x_i|z)p(z)dz \\ &= \sum_{i=1}^N \mathbb{E}_{z \sim p(z)} \log p(x_i|z), \end{aligned} \quad (29)$$

As Eqn. (29) is in general intractable, its variational lower bound is instead maximized:

$$\begin{aligned} \mathcal{L}(\omega_1, \omega_2, X) &= -D_{\text{KL}}(q_{\omega_1}(z|X)||p(z)) \\ &\quad + \mathbb{E}_{q_{\omega_1}(z|X)} [\log p_{\omega_2}(X|z)] \end{aligned} \quad (30)$$

Here, D_{KL} is the Kullback–Leibler (KL) divergence, and $p(z)$ is the prior distribution that is typically assumed to be a multivariate normal distribution. $q_{\omega_1}(z|X)$ is the posterior distribution parametrized by ω_1 . To analytically evaluate the KL divergence, the posterior is typically constructed as $\mathcal{N}(z|\mu_{\omega_1}(X), \Sigma_{\omega_1}(X))$. From a coding theory perspective, $q_{\omega_1}(z|X)$ and $p_{\omega_2}(X|z)$ can be considered as a probabilistic encoder and a probabilistic decoder, respectively.

To compute the $\nabla_{\omega_1} L(\omega_1, \omega_2, X)$, policy gradient theorem [47] or Reparametrization trick [38], [47] can be used. The latter is often used in VAE as it typically leads to a lower variance.

In practice, the encoder $q_{\omega_1}(z|X)$ and the decoder $p_{\omega_2}(X|z)$ can be any function approximator. An implementation of VAE as the generative model to sample actions can be found in <https://github.com/sfujim/BCQ>. In this work, the latent space dimension is selected to be 5, and the encoder and the decoder are both MLPs with 2 layers of 300 hidden units.

APPENDIX B AUTOREGRESSIVE MODEL WITH LSTM

For any probability distribution, the joint distribution can be factorized as a product of conditional probabilities as follow:

$$\begin{aligned} p(x) &= \prod_{i=1}^K p(x^{(i)}|x^{(1)}, \dots, x^{(i-1)}), \\ \rightarrow \log p(x) &= \sum_{i=1}^K \log p(x^{(i)}|x^{(1)}, \dots, x^{(i-1)}), \end{aligned} \quad (31)$$

where $x^{(i)}$ is the i^{th} dimension of the discrete input vector and K is the dimension of the input vector. As shown in Fig. 9, the input vector in one-hot vector form is fed to the LSTM network in sequence. The i^{th} output of the LSTM network after the softmax operation becomes a proper conditional probability $p(x^{(i)}|x^{(1)}, \dots, x^{(i-1)})$. The model is trained by

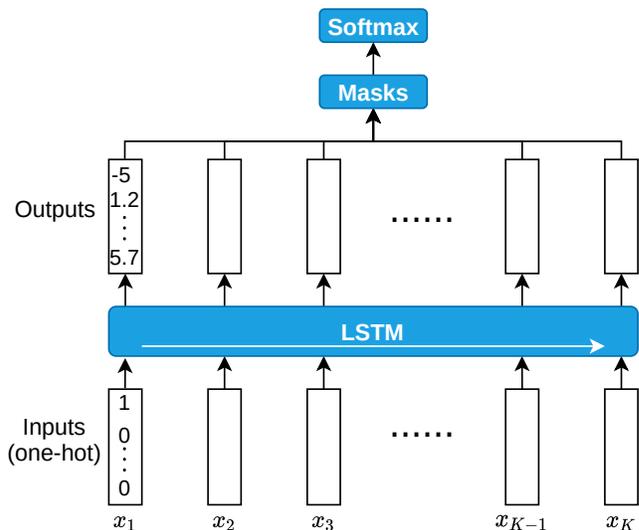


Fig. 9. The Network Architecture of Recurrent Autoregressive Model.

minimizing the KL divergence between the data distribution sampled from the experience replay buffer and the modeled distribution:

$$\begin{aligned} \min_{\psi} D_{\text{KL}}[p^*(x)||p_{\psi}(x)] \\ = \min_{\psi} \mathbb{E}_{x \sim p^*(x)} [-\log p_{\psi}(x)] + \text{constant} \end{aligned} \quad (32)$$

In this work, the LSTM network has a single layer and 50 hidden size.

APPENDIX C ADDITIONAL COMPARISON AMONG STRATEGIES

Here, we show the comparison on two additional trips. The trip shown in Fig. 10 contains a large number of signalized intersections. As indicated by Fig. 7, the gap between SMORL and the baseline and the gap between the wait-and-see solution and SMORL are both amplified by the high traffic density. The trip shown in Fig. 11 has a very low traffic density and the speed limits higher. In such case, the difference between SMORL and the wait-and-see solution becomes less noticeable. Meanwhile, SMORL was still able to consume less fuel by using the capacity of the battery more efficiently.

APPENDIX D ABLATION

In this part, we compare the full SMORL algorithm with the four intermediate algorithms. All the algorithms presented below use trajectory optimization solved via the PDDP solver. Tab. IV shows the difference in configuration and compares the trained final performance over the 100 trips used for testing. Here, we see that the safe set and BCQ both have a positive impact on the trained performance. In fact, the combination of TD3 and trajectory optimization (Config. 2) does not provide any show any significant improvement over trajectory optimization only (Config. 1). In addition, without the use of the safe set, the controller will deplete the battery SoC to SoC^{\min} at the end of the trip as the terminal state constraint cannot be considered unless with the help of extrinsic penalty.

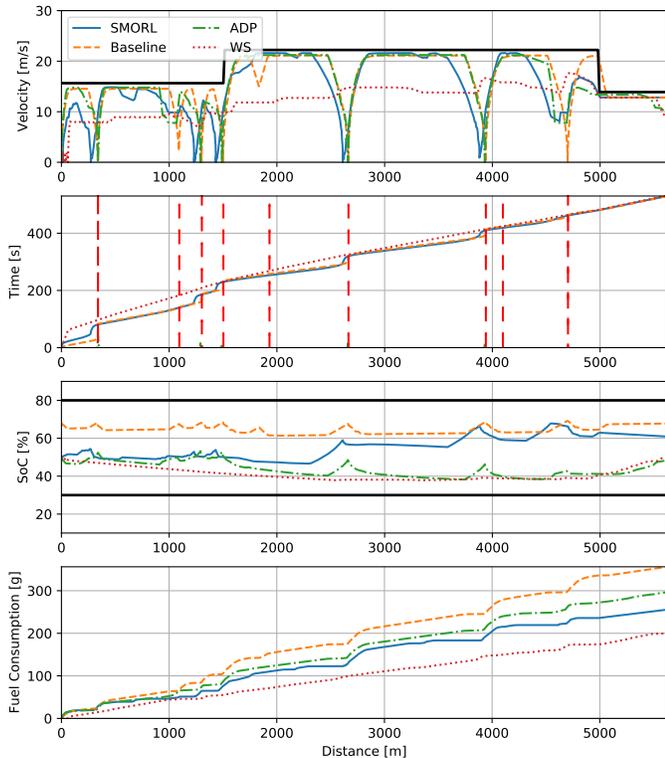


Fig. 10. Comparison for High-density Low-speed Scenario.

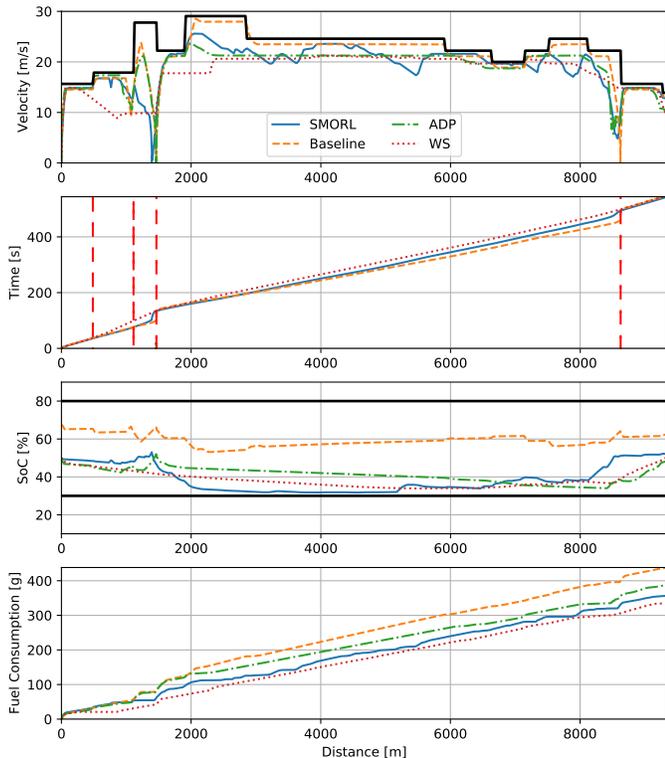


Fig. 11. Comparison for Low-density High-speed Scenario.

ACKNOWLEDGMENT

The authors acknowledge the support from the United States Department of Energy, Advanced Research Projects Agency-

TABLE IV
ABLATION STUDY FOR SMORL

	Safe Set	Q-learning	Fuel Economy mpg	Average Speed m/s	Normalized Cost
1		None	43.1	11.2	100
2		TD3	39.1	13.9	88.0
3	✓	TD3	39.9	13.3	90.5
4		BCQ	38.5	14.3	87.2
5	✓	BCQ	41.6	14.0	86.5

Energy (ARPA-E) NEXTCAR project (Award Number DE-AR0000794) and Ohio Supercomputer Center.

REFERENCES

- [1] A. Vahidi and A. Sciarretta, "Energy saving potentials of connected and automated vehicles," *Transportation Research Part C: Emerging Technologies*, vol. 95, pp. 822–843, 2018.
- [2] A. Sciarretta, G. De Nunzio, and L. L. Ojeda, "Optimal ecodriving control: Energy-efficient driving of road vehicles as an optimal control problem," *IEEE Control Systems Magazine*, vol. 35, no. 5, pp. 71–90, 2015.
- [3] Q. Jin, G. Wu, K. Boriboonsomsin, and M. J. Barth, "Power-based optimal longitudinal control for a connected eco-driving system," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 10, pp. 2900–2910, 2016.
- [4] E. Ozatay, S. Onori, J. Wollaeger, U. Ozguner, G. Rizzoni, D. Filev, J. Michelini, and S. Di Cairano, "Cloud-based velocity profile optimization for everyday driving: A dynamic-programming-based solution," *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 6, pp. 2491–2505, 2014.
- [5] J. Han, A. Vahidi, and A. Sciarretta, "Fundamentals of energy efficient driving for combustion engine and electric vehicles: An optimal control perspective," *Automatica*, vol. 103, pp. 558–572, 2019.
- [6] C. Sun, J. Guanetti, F. Borrelli, and S. Moura, "Optimal eco-driving control of connected and autonomous vehicles through signalized intersections," *IEEE Internet of Things Journal*, 2020.
- [7] F. Mensing, R. Trigui, and E. Bideaux, "Vehicle trajectory optimization for hybrid vehicles taking into account battery state-of-charge," in *2012 IEEE vehicle power and propulsion conference*. IEEE, 2012, pp. 950–955.
- [8] L. Guo, B. Gao, Y. Gao, and H. Chen, "Optimal energy management for hevs in eco-driving applications using bi-level mpc," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 8, pp. 2153–2162, 2016.
- [9] P. Olin, K. Aggoune, L. Tang, K. Confer, J. Kirwan, S. R. Deshpande, S. Gupta, P. Tulpule, M. Canova, and G. Rizzoni, "Reducing fuel consumption by using information from connected and automated vehicle modules to optimize propulsion system control," SAE Technical Paper, Tech. Rep., 2019.
- [10] D. Maamria, K. Gillet, G. Colin, Y. Chamaillard, and C. Nouillant, "Computation of eco-driving cycles for hybrid electric vehicles: Comparative analysis," *Control Engineering Practice*, vol. 71, pp. 44–52, 2018.
- [11] B. Asadi and A. Vahidi, "Predictive cruise control: Utilizing upcoming traffic signal information for improving fuel economy and reducing trip time," *IEEE transactions on control systems technology*, vol. 19, no. 3, pp. 707–714, 2010.
- [12] G. Guo and Y. Wang, "Eco-driving of freight vehicles with signal priority on congested arterial roads," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 5, pp. 4225–4237, 2021.
- [13] J. Shi, F. Qiao, Q. Li, L. Yu, and Y. Hu, "Application and evaluation of the reinforcement learning approach to eco-driving at intersections under infrastructure-to-vehicle communications," *Transportation Research Record*, vol. 2672, no. 25, pp. 89–98, 2018.
- [14] G. Li and D. Gorges, "Ecological adaptive cruise control for vehicles with step-gear transmission based on reinforcement learning," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [15] G. Guo and Y. Wang, "An integrated mpc and deep reinforcement learning approach to trans-priority active signal control," *Control Engineering Practice*, vol. 110, p. 104758, 2021.

- [16] A. Pozzi, S. Bae, Y. Choi, F. Borrelli, D. M. Raimondo, and S. Moura, "Ecological velocity planning through signalized intersections: A deep reinforcement learning approach," in *2020 59th IEEE Conference on Decision and Control (CDC)*. IEEE, 2020, pp. 245–252.
- [17] Z. Zhu, S. Gupta, A. Gupta, and M. Canova, "A deep reinforcement learning framework for eco-driving in connected and automated hybrid electric vehicles," *arXiv preprint arXiv:2101.05372*, 2021.
- [18] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch, "Plan online, learn offline: Efficient learning and exploration via model-based control," *arXiv preprint arXiv:1811.01848*, 2018.
- [19] N. Karnchanachari, M. I. Valls, D. Hoeller, and M. Hutter, "Practical reinforcement learning for mpc: Learning from sparse objectives in under an hour on a real robot," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 211–224.
- [20] B. Thananjeyan, A. Balakrishna, U. Rosolia, F. Li, R. McAllister, J. E. Gonzalez, S. Levine, F. Borrelli, and K. Goldberg, "Safety augmented value estimation from demonstrations (saved): Safe deep model-based rl for sparse cost robotic tasks," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3612–3619, 2020.
- [21] B. Thananjeyan, A. Balakrishna, U. Rosolia, J. E. Gonzalez, A. Ames, and K. Goldberg, "Abc-mpc: Safe sample-based learning mpc for stochastic nonlinear dynamical systems with adjustable boundary conditions," *arXiv preprint arXiv:2003.01410*, 2020.
- [22] T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, P. Abbeel, and J. Ba, "Benchmarking model-based reinforcement learning," *arXiv preprint arXiv:1907.02057*, 2019.
- [23] L.-J. Lin, "Self-improving reactive agents based on reinforcement learning, planning and teaching," *Machine learning*, vol. 8, no. 3-4, pp. 293–321, 1992.
- [24] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International Conference on Machine Learning*, 2018, pp. 1582–1591.
- [25] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," *arXiv preprint arXiv:2005.01643*, 2020.
- [26] S. Fujimoto, D. Meger, and D. Precup, "Off-policy deep reinforcement learning without exploration," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2052–2062.
- [27] U. Rosolia and F. Borrelli, "Sample-based learning model predictive control for linear uncertain systems," in *2019 IEEE 58th Conference on Decision and Control (CDC)*. IEEE, 2019, pp. 2702–2707.
- [28] P. A. Lopez, M. Behrisch, L. Bieker-Walz, J. Erdmann, Y.-P. Flötteröd, R. Hilbrich, L. Lücken, J. Rummel, P. Wagner, and E. Wießner, "Microscopic traffic simulation using sumo," in *The 21st IEEE International Conference on Intelligent Transportation Systems*. IEEE, 2018. [Online]. Available: <https://elib.dlr.de/124092/>
- [29] OpenStreetMap contributors, "Planet dump retrieved from <https://planet.osm.org>," <https://www.openstreetmap.org>, 2017.
- [30] N. Heess, D. TB, S. Sriram, J. Lemmon, J. Merel, G. Wayne, Y. Tassa, T. Erez, Z. Wang, S. Eslami *et al.*, "Emergence of locomotion behaviours in rich environments," *arXiv preprint arXiv:1707.02286*, 2017.
- [31] S. R. Deshpande, S. Gupta, A. Gupta, and M. Canova, "Real-time eco-driving control in electrified connected and autonomous vehicles using approximate dynamic programming," *arXiv preprint arXiv:2108.02652*, 2021.
- [32] Z. Zhu, S. Gupta, N. Pivaro, S. R. Deshpande, and M. Canova, "A GPU implementation of a look-ahead optimal controller for eco-driving based on dynamic programming," *arXiv preprint arXiv:2104.01284*, 2021.
- [33] F. Borrelli, A. Bemporad, and M. Morari, *Predictive control for linear and hybrid systems*. Cambridge University Press, 2017.
- [34] K. Chua, R. Calandra, R. McAllister, and S. Levine, "Deep reinforcement learning in a handful of trials using probabilistic dynamics models," *arXiv preprint arXiv:1805.12114*, 2018.
- [35] C. J. C. H. Watkins, "Learning from delayed rewards," 1989.
- [36] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [37] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *International conference on machine learning*. PMLR, 2016, pp. 1995–2003.
- [38] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [39] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *International conference on machine learning*. PMLR, 2014, pp. 387–395.
- [40] A. Karpathy, J. Johnson, and L. Fei-Fei, "Visualizing and understanding recurrent networks," *arXiv preprint arXiv:1506.02078*, 2015.
- [41] O. Sundstrom and L. Guzzella, "A generic dynamic programming matlab function," in *2009 IEEE control applications,(CCA) & intelligent control,(ISIC)*. IEEE, 2009, pp. 1625–1630.
- [42] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine, "Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 7559–7566.
- [43] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.
- [44] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.
- [45] O. S. Center, "Ohio supercomputer center," 1987. [Online]. Available: <http://osc.edu/ark:/19495/f5s1ph73>
- [46] S. Gupta, S. R. Deshpande, P. Tulpule, M. Canova, and G. Rizzoni, "An enhanced driver model for evaluating fuel economy on real-world routes," *IFAC-PapersOnLine*, vol. 52, no. 5, pp. 574–579, 2019.
- [47] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3-4, pp. 229–256, 1992.