

1 **Developing an Eco-Driving Strategy in a Hybrid Traffic**

2 **Network Using Reinforcement Learning**

3

4 **Umar Jamil¹, Mostafa Malmir¹, Alan Chen², Monika Filipovska³, Mimi Xie⁴, Caiwen**

5 **Ding⁵, Yu-Fang Jin^{1§}**

6

7 ¹Department of Electrical and Computer Engineering, the University of Texas at San Antonio,
8 San Antonio, TX 78249, USA

9 ²Westlake High school, Austin TX 78746, USA

10 ³Department of Civil Engineering, University of Connecticut, Storrs, CT 06269, USA

11 ⁴Department of Computer Science, the University of Texas at San Antonio, San Antonio, TX
12 78249, USA

13 ⁵Department of Computer Science & Engineering, University of Connecticut, Storrs, CT 06269,
14 USA

15 **§Correspondence** should be addressed to Yu-Fang Jin (yufang.jin@utsa.edu)

16

17

18

19

20 **Abstract**

21 Eco-driving has garnered considerable research attention owing to its potential socio-
22 economic impact, including enhanced public health and mitigated climate change effects
23 through the reduction of greenhouse gas emissions. With an expectation of more
24 autonomous vehicles (AV) on the road, an eco-driving strategy in hybrid traffic networks
25 encompassing AV and human-driven vehicles (HDV) with the coordination of traffic
26 lights is a challenging task. The challenge is partially due to the insufficient infrastructure
27 for collecting, transmitting, and sharing real-time traffic data among vehicles, facilities,
28 and traffic control centers, and the following decision-making of agents involved in traffic
29 control. Additionally, the intricate nature of the existing traffic network, with its diverse
30 array of vehicles and facilities, contributes to the challenge by hindering the development
31 of a mathematical model for accurately characterizing the traffic network. In this study,
32 we utilized the Simulation of Urban Mobility (SUMO) simulator to tackle the first
33 challenge through computational analysis. To address the second challenge, we employed
34 a model-free reinforcement learning (RL) algorithm, Proximal policy optimization
35 (PPO), to decide the actions of AV and traffic light signals in a traffic network. A novel
36 eco-driving strategy was proposed by introducing different percentages of AV into the
37 traffic flow and collaborating with traffic light signals using RL to control the overall
38 speed of the vehicles, resulting in improved fuel consumption efficiency. Average

39 rewards with different penetration rates of AV (5%, 10%, and 20% of total vehicles) were
40 compared to the situation without any AV in the traffic flow (0% penetration rate). The
41 10% penetration rate of AV showed a minimum time of convergence to achieve average
42 reward, leading to a significant reduction in fuel consumption and total delay of all
43 vehicles.

44 **Keywords:** Eco-driving; Hybrid Traffic Network; Reinforcement Learning; Traffic Flow
45 Control; Fuel Consumption; Microscopic Traffic Simulator

46 **1. Introduction**

47 Findings from a 2022 study indicate that the transportation sector accounted for 27% of
48 the energy consumption in the United States.¹ Specifically, petroleum (gasoline)
49 consumption comprised about 52% of the total energy consumption, resulting in
50 significant air pollutant emissions. This underscores the necessity for a well-designed
51 traffic control system to mitigate fuel energy consumption (FEC) and air pollution for
52 sustainability.²⁻⁴ The concept of sustainability has driven research into eco-driving
53 strategies designed to reduce FEC rates (FEC within time). FEC rates can be calculated
54 based on factors such as acceleration, mass, drag coefficient, rolling coefficient, driveline
55 efficiency, idling speed, and idling fuel mean pressure.^{5,6} Reducing FEC involves two
56 interconnected goals: shorter travel time and lower FEC rates. Vehicles incur the highest

57 FEC rates during idling and frequent stops and starts, especially at traffic lights or in
58 congestion. Therefore, prioritizing the establishment of a continuous traffic flow,
59 characterized by minimal fluctuations in vehicle speeds, is essential for achieving lower
60 FEC rates and shorter traffic delays. This approach is instrumental in promoting effective
61 eco-driving strategies.⁷

62 Traditional traffic control relies on fixed modes for traffic light changes and manual
63 rerouting, resulting in limited efficiency and a lack of feedback mechanisms. The current
64 setup of traffic control systems poses challenges in developing eco-driving strategies for
65 hybrid traffic networks encompassing AV and HDV. These challenges stem partially
66 from the insufficient infrastructure for collecting, transmitting, and sharing real-time
67 traffic data among vehicles, facilities, and traffic control centers, as well as the subsequent
68 decision-making by involved agents. Furthermore, the intricate nature of the existing
69 traffic networks, with their diverse array of vehicles and facilities, complicates the
70 development of a mathematical model for accurately characterizing the traffic networks.

71 Current eco-driving strategies have addressed the challenges from various perspectives,
72 including real-time artificial intelligence for traffic monitoring, and 5th generation (5G)
73 communication networks to facilitate rapid information sharing.⁸⁻¹⁰ Due to the
74 multifaceted nature of the eco-driving problem, a model-based deterministic strategy is

75 challenging to approach. Meanwhile, data-driven approaches show promise, given the
76 large amount of data accumulated during the past decades.

77 ***Related Work on Reinforcement Learning in Traffic Control***

78 Model-free RL has demonstrated its advantage in decision-making for traffic flow control
79 by examining interactions among multiple agents and the environment.¹⁰⁻¹² RL has been
80 applied to optimize vehicle routes for reduced delay and vehicle accelerations for less
81 FEC.^{13,14} RL algorithms have also been developed to reduce air pollutant emissions by
82 reducing vehicles' waiting time at road intersections.^{15,16} In a study on infrastructure-to-
83 vehicle communications networks,¹⁷ a single vehicle was considered as an agent, and the
84 Q-learning (QL) algorithm was developed to minimize carbon dioxide emissions.
85 Additionally, a recent eco-driving framework based on the deep Q-network (DQN)
86 approach was presented to enhance the fuel efficiency of multiple vehicles in a traffic
87 network with one horizontal road and one vertical road.¹⁸

88 In addition to applications of RL in controlling vehicle routes or acceleration, traffic lights
89 are also considered as agents to control traffic flow with RL algorithms. An RL-based
90 control has been developed for smart traffic signals, to reduce traffic jams and improve
91 traffic smoothness in a traffic grid consisting of 3 horizontal and 3 vertical roads.¹⁹

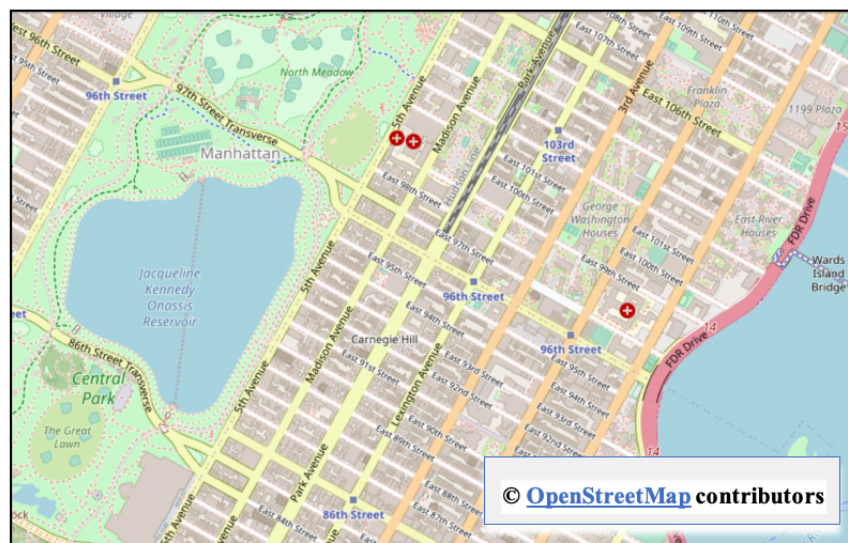
92 With more AV running on the road, they are also considered agents in RL algorithms for
93 traffic control. In a recent study, a circular network with fixed traffic signal patterns at

94 one spot was deployed to develop a deep deterministic policy gradient (DDPG) algorithm.
95 The study aims to minimize the FEC of Connected AV by controlling their acceleration.²⁰
96 Additionally, RL algorithms with a hybrid deep Q-learning and policy gradient (HDQPG)
97 were developed to minimize the FEC of Connected AV by controlling their acceleration
98 in a traffic grid with one horizontal and five vertical roads.²¹ Previous studies also
99 explored a traffic flow containing both HDV and Connected AV using a trust region
100 policy optimization (TRPO) to reduce the FEC and emissions of both HDV and CAV.²²
101 While the above-mentioned RL-based controls have improved traffic smoothness by
102 focusing on the actions of vehicles or traffic lights, the effect of combining AVs and
103 traffic signals on FEC has not been fully investigated.²⁰⁻²²
104 In this study, a novel eco-driving strategy was proposed by introducing a specific
105 percentage of AV into the traffic flow of HDV, in collaboration with smart traffic light
106 signals to reduce the idling time of vehicles and improve the traffic smoothness in a
107 scalable traffic network with user-defined horizontal and vertical roads for intersections.
108 A model-free RL algorithm was developed to control the overall speed of all vehicles,
109 resulting in a continuous traffic flow and reduction of the FEC of the vehicles in the
110 network.

111 **2. Method**

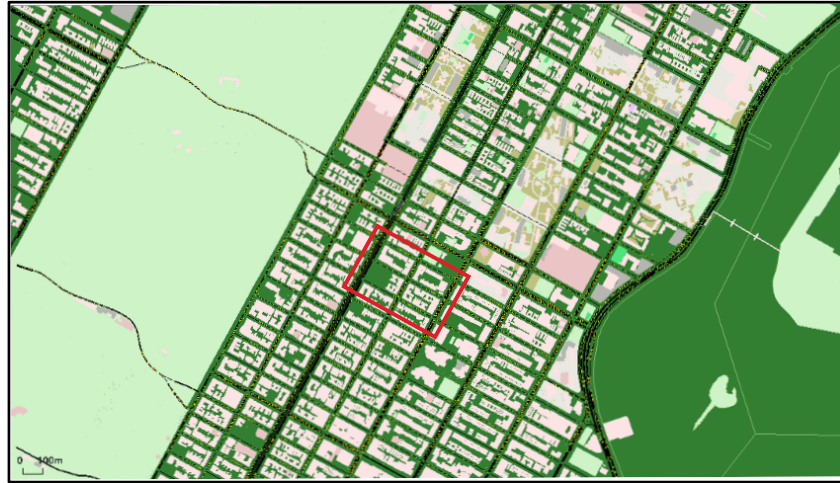
112 The proposed RL algorithm determines the optimal actions of multiple agents including
 113 AV and traffic lights in a dynamic traffic network with HDV to minimize the FEC rates
 114 of all vehicles. The traffic network and motion of all vehicles are simulated using the
 115 SUMO package.²³ The RL algorithm is implemented using Python and integrated into
 116 SUMO for simulation.

117 The selected traffic grid environment is inspired by the grid-like layout of Manhattan
 118 City.²⁴ Figures 1 and 2 display an open street map of the Manhattan traffic grid structure
 119 and its visualization in the SUMO environment, respectively.



120

121 Figure 1. Open street map of traffic grid structure in Manhattan City.



122

123 Figure 2. The grid structure of Manhattan City, simulated in the SUMO environment, is represented in the
 124 highlighted red color region. The selected traffic network serves as the basis for our research, examining
 125 the role of AV combined with HDV in minimizing the FEC rates of all vehicles in the traffic network.

126 ***Environment Setup in SUMO***

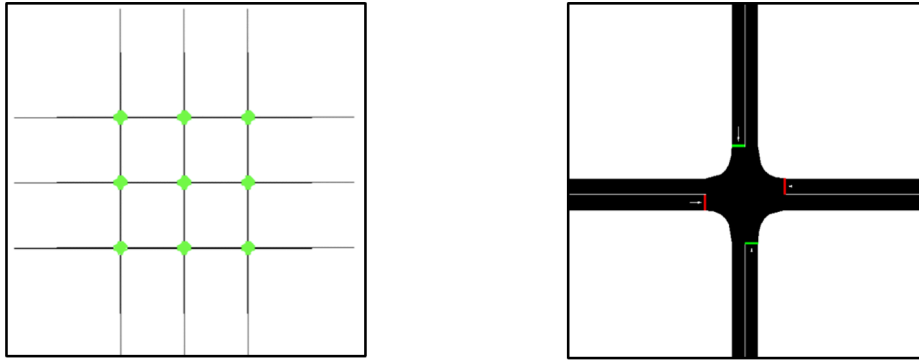
127 The traffic network is configured within an environment featuring N horizontal and N
 128 vertical straight roads, each equipped with two lanes and extending for a length of 1
 129 kilometer (km). There are $4N$ edge points, each assigned a unique number. At each edge
 130 point, a traffic flow of 300 vehicles per hour has been selected to enter the traffic system,
 131 aligning with the range of traffic flow defined by the Federal Highway Administration
 132 for signalized intersections in the United States.²⁵ Each vehicle has a departure speed of
 133 30 m/s (67.1 miles/hour) and SUMO vehicle parameters dictate a minimum gap of 2.5
 134 meters between two vehicles. All vehicles will continue straight in their original direction

135 of travel and exit the simulation environment. To ensure safety during peak traffic time,
136 turn prohibitions are considered in this study according to the Federal Highway
137 Administration in the United States.²⁵

138 According to a recent study, AV account for 10% of all vehicles on the roads.²⁶
139 Accordingly, this study considers different penetration rates for AV (0%, 5%, 10%, and
140 20%) to assess their impact on traffic control. An RL controller is used to control RL
141 agents, such as AV and traffic lights, with commands issued by policy at each time step.
142 The speed and acceleration of AV are determined with an RL controller, while the motion
143 of HDV is controlled by an embedded “sim car-following” controller in SUMO
144 simulation. All vehicles are homogeneous with respect to their mass, size, and economic
145 models.

146 At each intersection of two roads, 4-way traffic lights are defined as actuated agents with
147 a controllable period for red, green, and yellow lights. With the setup of N vertical and N
148 horizontal roads in a network, there are a total of $4N^2$ traffic lights.

149 In this study, we focus on a 3x3 traffic network, assuming uniform road lengths in all
150 directions to facilitate simulation. Figure 3 illustrates a network with $N=3$ and the
151 arrangement of 4 traffic lights at an intersection. It's important to highlight that the
152 framework is adaptable to larger-sized traffic networks, provided there are sufficient
153 computational resources.



154

155 Figure 3. (a) 3x3 Traffic light grid environment (b). 4-way single signalized intersection.

156 ***Reinforcement Learning***

157 A decentralized partially observable Markov Decision Process (De-POMDP) is adopted

158 to coordinate the actions of agents, including traffic lights and AV in the traffic network.

159 When vehicles move in the same direction, HDV are observable to an autonomous vehicle

160 if the distance between a human-driven vehicle and an autonomous vehicle is less than or

161 equal to 25 meters in the same lane. Each traffic light agent also observes the two nearest

162 vehicles and has their information related to speed, distance to the intersection, and edge

163 number. The position, speed, and acceleration of AV, as well as cycles and status of traffic

164 lights, are shared among all AV and traffic lights.

165 The state, action space, policy, and reward function of the RL are defined as follows.

166 **State Space (s):** For each vehicle agent, its state, $s := (v_i, d_i, e_i)_{i=1:M} \in R^{3 \times M}$, where167 $M=3,600$ denotes the maximum number of vehicles in the selected traffic system. This168 number is calculated by considering 300 vehicles entering the system at $4N$ edge points

169 within one hour, with $N=3$, assuming the worst case: no vehicles leave the simulated
 170 traffic network within an hour. Here, v_i represents the speed of the i^{th} vehicle, d_i denotes
 171 the distance of the i^{th} vehicle to the nearest intersection in its driving direction, e_i indicates
 172 the edge number which the i^{th} vehicle enters the traffic network. The edge number
 173 signifies the traffic flow direction of each vehicle, assuming no turns are allowed.

174 The state of each traffic light agent includes the time of the light's last change, the traffic
 175 flow direction controlled by the light (0 indicates passing with a green light, and 1
 176 indicates stopping with a red light status), and the states of other traffic lights in the same
 177 traffic flow direction. At an intersection, if the top-bottom traffic lights have a status of
 178 “0”, the left-right traffic lights must have a status of “1”, and vice versa. When the status
 179 of a traffic light is green, it will change to yellow for 3s before switching to red status due
 180 to safety purposes.

181 **Action Space (\mathbf{a}):** There are $4N^2$ traffic lights in the network and each traffic light have
 182 two discrete actions: 1 (indicating the traffic light switches) and -1 (indicating no action
 183 taken), as defined in equation (1):

$$\left[a = \begin{cases} 1 & \text{traffic light switches} \\ -1 & \text{no action taken} \end{cases} \right] \quad (1)$$

184 The action space for an autonomous vehicle is its acceleration, which ranges between
 185 $[-1, 1]$ and is determined by an RL controller in FLOW

186 package. For HDV, the action space for acceleration values is chosen within the range [-
187 4.5, 2.6], as defined by SUMO.

188 **Policy:** An RL algorithm called PPO is used to train policies for tasks involving decision-
189 making in environments with either continuous or discrete action spaces. Policies are
190 optimized using the policy gradient method to maximize the expected cumulative reward.
191 The choice of a PPO-based RL algorithm for deployment in this study stems from its
192 superior computational efficiency and stability compared to other algorithms.
193 Specifically, RLlib within the Flow package is integrated into SUMO for simulation.²⁷⁻²⁹

194 A stochastic policy $\pi_\varphi: s \times a \rightarrow \mathbb{R}_+$ is a mapping from state, s , and action a of all agents
195 parameterized by φ to a non-negative real number. It can be defined by (2) as a probability
196 distribution over actions of each state:

$$\left[\pi_\varphi = P(a|s; \varphi) = \frac{e^{f_\varphi(s,a)}}{\sum_{a' \in \mathcal{A}} e^{f_\varphi(s,a')}} \right] \quad (2)$$

197 where, $f_\varphi(s, a') = \varphi^\top \mathcal{B}(s, a')$; $\varphi = (\varphi_{a1}, \dots, \varphi_{aU}) \in R^U$; φ^\top is transpose of parameter
198 vector φ ; and $\mathcal{B}(s, a)$ represents transitions among states given an action; and U
199 represents the complete action space.

200 The average reward received by an agent when it follows a PPO policy at each time step
201 is referred to as the average policy reward. The average policy reward, also defined as

202 expected return of policy, $\eta(\pi_\varphi)$ for the entire trajectory τ at time step t , can be expressed
 203 as equation (3),

$$\left[\eta(\pi_\varphi) = \mathbb{E}_\tau \left[\sum_{t=0}^{\infty} \gamma^t \cdot r(s_t, a_t) \right], \right] \quad (3)$$

204 where, τ represents the entire trajectory of states and actions. The parameter $0 < \gamma \leq$
 205 1 , represents a discount factor, and γ^t gets smaller as time $t \rightarrow \infty$ with $\gamma < 1$. The rewards
 206 function, $r(s_t, a_t)$, determines rewards given the state and action of an agent at time t .
 207 The optimal policy parameter φ^* is reached by maximizing the expected cumulative
 208 return obtained by an agent, as described in equation (4),

$$[\varphi^* := \operatorname{argmax}_\varphi \eta(\pi_\varphi).] \quad (4)$$

209 The policy loss is defined based on the $q_t(\varphi)$, a ratio of new policy $\pi_\varphi(a_t|s_t)$ and the
 210 previous policy $\pi_{\varphi_g}(a_t|s_t)$ as equation (5):

$$\left[\text{Policy Loss} = \mathbb{E}_\tau \left[q_t(\varphi) \hat{A}_t - \beta \text{KL} \left[\pi_{\varphi_g}(\cdot|s_t), \pi_\varphi(\cdot|s_t) \right] \right], \right] \quad (5)$$

211 where, β is hyperparameter to control the strength of regularization of
 212 $\text{KL}[\pi_{\varphi_g}(\cdot|s_t), \pi_\varphi(\cdot|s_t)]$, which represents the Kullback-Leibler (KL) divergence
 213 between two conditional probability distributions over actions given a state s_t . If

214 $\mathbb{E}_\tau \left[KL \left[\pi_{\varphi_g}(\cdot | s_t), \pi_\varphi(\cdot | s_t) \right] \right] < \frac{KL \text{ target value}}{1.5}$, it indicates new policy doesnot diverged
 215 significantly from the old policy, so β needs to be reduced by 1/2. If
 216 $\mathbb{E}_\tau \left[KL \left[\pi_{\varphi_g}(\cdot | s_t), \pi_\varphi(\cdot | s_t) \right] \right] > (KL \text{ target value}) \times 1.5$, it means there is too much
 217 change in policy through update, so β needs to be increased by multiplying with 2. The
 218 *KL target value* is defined by users and a reference value is given in the Results section.
 219 The advantage estimate function \hat{A}_t , representing accumulated future rewards, can be
 220 defined as equation (6),

$$\left[\hat{A}_t = \delta_t + \sum_{d=1}^{T-t+1} (\gamma\lambda)^d \delta_{t+d}, \right] \quad (6)$$

$$\left[\delta_t = r_t + \gamma V_{(s_{t+1})} - V_{(s_t)}, \right] \quad (7)$$

221 where t represents time steps from $[0, T]$, and T represents the range of prediction.
 222 The parameter λ impacts weights of potential rewards in the advantage estimation
 223 function \hat{A}_t . When $\lambda=1$, \hat{A}_t increases by adding more future rewards, resulting in high
 224 variance and less bias. When $\lambda=0$, no future rewards are considered. Policy $\pi_{\varphi_{g+1}}$ is
 225 updated with φ_{g+1} according to (8):

$$\left[\varphi_{g+1} = \underset{\varphi}{\operatorname{argmax}} \frac{1}{|\mathcal{H}_g|T} \sum_{T=\mathcal{H}_g} \sum_{t=0}^T \min \left(q_t(\varphi) A^{\pi_{\varphi_g}}(s_t, a_t) - \beta_g KL[\pi_{\varphi_g}(\cdot | s_t), \pi_\varphi(\cdot | s_t)] \right), \right] \quad (8)$$

226 where $\mathcal{H}_g = [\mathcal{T}_i]$ is a set of trajectories for iteration g .

227 In the RL algorithm, the value function $V_{(s_t)}$ estimates the expected cumulative reward
 228 starting from a specific state s_t , that the agent can attain from that state onwards. A value
 229 function loss (*VF Loss*) is defined as a squared-error loss between predicted and target
 230 value function (9):

$$\left[VF\ Loss = (V_{\phi_g}(s_t) - V_t^{target})^2, \right] \quad (9)$$

231 where, $V_{\phi}(s_t)$, is an output from a neural network parameterized by ϕ with the state s_t as
 232 input; and V_t^{target} is the target value function at time step t can be defined as $V_t^{target} = r_t +$
 233 $\gamma V_{(s_{t+1})}$, the range of $V_t^{target} \in [-1, 1]$. Parameters of the network ϕ_{g+1} can be updated
 234 according to (10):

$$\left[\phi_{g+1} = argmin_{\phi} \frac{1}{|\mathcal{H}_g|T} \sum_{T=\mathcal{H}_g} \sum_{t=0}^T (V_{\phi_g}(s_t) - V_t^{target})^2 \right]. \quad (10)$$

235
 236 In RL, the entropy function refers to the level of uncertainty in the policy distribution. It
 237 is used to encourage exploration by selecting different possible actions in a specific state
 238 and to prevent premature convergence to suboptimal policies. The entropy function of the
 239 PPO algorithm is defined based on the probability of taking actions, $\pi(a|s)$, given a state
 240 s under the policy in equation (11):

$$\left[Entropy = - \sum_a \pi(a|s) \log \pi(a|s). \right] \quad (11)$$

241 A smaller entropy indicates a better performance of the PPO algorithm. The pseudo code
 242 for the PPO algorithm is shown as follows.
 243

Algorithm 1. PPO

Input: Initial policy and value function parameters $(\varphi_0, \varnothing_0)$

for iteration $g=0,1, 2\dots$.do

Run policy $\pi_g = \pi(\varphi_g)$ in environment for time steps T to collect a set of trajectories $\mathcal{H}_g = [T_i]$.

Compute rewards-to-go \hat{r}_t .

Compute advantage estimates \hat{A}_t based on the current value function V_{\varnothing_g} .

Find optimal policy φ_g^* to find average policy reward.

Update policy $\pi_{\varphi_{g+1}}$ with φ_{g+1} using equation (8).

Fit value function V_{\varnothing_g} with \varnothing_{g+1} using equation (10).

end for

244
 245 **Reward (r):** Two reward functions have been designated: one to minimize total traffic
 246 delay, T_d , and another to minimize FEC rates at time step t . These reward functions were
 247 used to train each traffic light and autonomous vehicle. The reward functions are given in
 248 equations (12) and (13):

$$\left[r_1(t) = - \frac{1}{4N^2} T_d, \right] \quad (12)$$

249

$$\left[r_2(t) = -\frac{1}{M}Fc(t). \right] \quad (13)$$

250 Since rewards are negative, the closer a reward to zero means smaller total delay and FEC
 251 rates of all vehicles in the traffic flow. The T_d is defined by (14):

$$\left[T_d = \max \left(\frac{\sqrt{\sum_i^M (v_{ds_i})^2} - \sqrt{\sum_i^M (v_{ds_i} - v_i)^2}}{\sqrt{\sum_i^M (v_{ds_i})^2}}, 0 \right), \right] \quad (14)$$

252 where v_{ds} is the speed limit on the road and v_i is the velocity of each vehicle.

253 Fuel Energy Consumption Rate Model

254 The function. $Fc(t)$, denotes the FEC rate with a unit in Litter/second (L/s) , which is
 255 described as follows according to a previous study,⁵

$$\left[Fc(t) = \begin{cases} \alpha_0 + \alpha_1 P_t + \alpha_2 P_t^2, & \forall P_t \geq 0 \\ \alpha_0, & \forall P_t < 0 \end{cases} \right] \quad (15)$$

256

$$\left[P_t = \left(\frac{R_t + 1.04ma_t}{3600 \eta_d} \right) \cdot v_t, \right] \quad (16)$$

257

$$\left[R_t = \frac{\rho}{25.92} C_d C_h A_f v_t^2 + 9.8066m \frac{C_r}{1000} (c_1 v_t + c_2) + 9.8066mG_t, \right] \quad (17)$$

258 Here,

- 259 • P_t : Power exerted at time t (Kilowatt, KW),
- 260 • a_t : Acceleration of vehicle (m/s^2),
- 261 • v_t : Velocity of vehicle (m/s),
- 262 • R_t : Resistance force (N).

263 Other constant parameters in the model have been defined in Table 1.

264 Table 1. FEC rate model parameters

Symbol	FEC rate Parameters	Value
α_0	Vehicle model constant	0.00000002
α_1	Vehicle model constant	0.0000001
α_2	Vehicle model constant	0.000001
m	Vehicle mass	1200 Kg
η_d	Derive line efficiency	0.92
ρ	Density of air at sea level at a temperature of 59°F	1.2256 Kg/m ³
C_d	Drag coefficient	0.28
C_h	Correction factor for altitude	0.97
A_f	Frontal area	2.6 m ²
C_r	Rolling coefficient	1.75
c_1	Rolling resistance parameter	0.0328
c_2	Rolling resistance parameter	4.575
G_t	Roadway grade	0.04

265

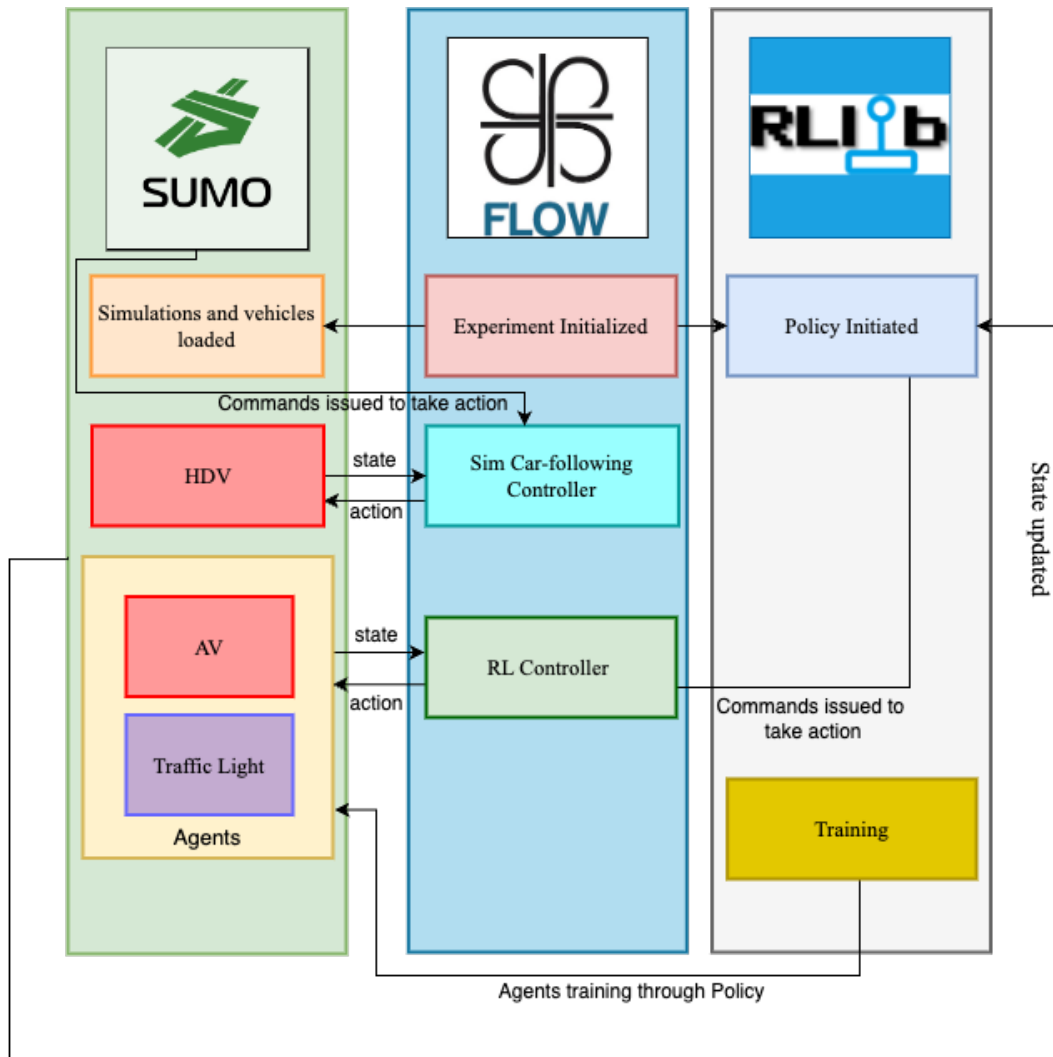
266 Computational Framework

267 Two publicly available software packages, Flow and SUMO, are adopted in this study.

268 Flow is a traffic control benchmarking framework developed in Python and integrates RL

269 algorithms into different traffic control scenarios.¹⁹ The SUMO simulator handles large-

270 scale traffic networks based on physical-world data. Integration of SUMO and Flow
 271 package and implementation of the RL algorithm are shown in Figure 4.



272

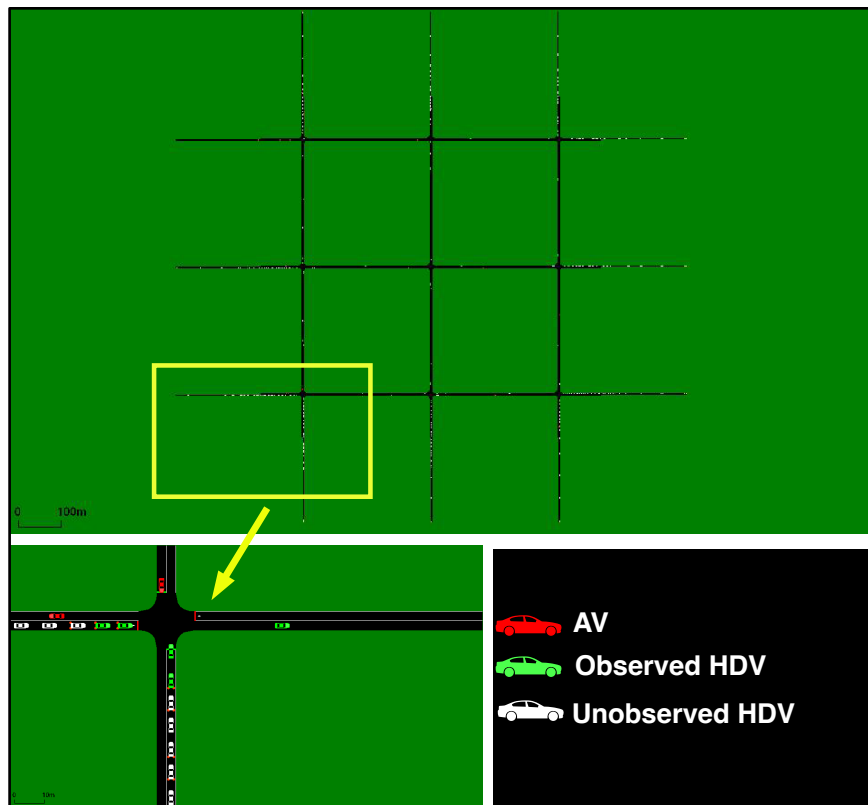
273 Figure 4: Process diagram to describe RL training process and interactions between SUMO, Flow, and
 274 RLLib library. RL and Sim car following controllers used to control the AV and HDV, respectively. Sim

275 car-following controller actions are entirely defined by the simulator, whereas RL-Controller performs
 276 actions by following commands from the policy in RLlib.

277

278 **Results**

279 The RL algorithm was applied to regulate the traffic flow in the selected traffic network
 280 with 4 different penetration rates of AV, 0%, 5%, 10%, and 20%.



281

282 Figure 5. Illustration of the traffic network with 3 vertical and 3 horizontal roads in SUMO simulator. An
 283 overview of all AV (red vehicles), observable HDV (green vehicles), and unobservable HDV (white

284 vehicles) in the traffic network. (Bottom Left) A close view of traffic flow between intersections within the
 285 yellow box at the lower left part of the traffic network.

286 Training was conducted on a machine with 4 Intel® Core™ i5-6600 CPU @ 3.30GHz.

287 The Hyperparameters used in the RL algorithms are listed in Table 2.

288 Table 2. RL algorithm hyperparameters

Hyperparameters	Value
Learning rate	5e-5
Training batch size	1500
SGD minimum batch size	128
Number of SGD iterations	5
Training Iterations	500
Parallel workers	10
Horizon steps	150
Discount factor (γ)	0.999
GAE value (λ)	0.5
KL Target value	0.02
Target value function	0.01
Fixed KL β	3

289 SGD: Stochastic Gradient Descent; GAE: Generalized Advantage Estimation; KL: Kullback-Leibler.

290 The results of this study have been divided into four categories:

- 291 (a) Rewards on total delay at different penetration levels;
- 292 (b) Rewards on FEC rates at different penetration levels;
- 293 (c) Performance of PPO policy at different penetration levels;
- 294 (d) Comparison of the selected 3x3 traffic network with other networks.

295 **Rewards on Total Delay at Different Penetration Levels**

296 Various penetration rates of AV in the traffic flow are examined to optimize traffic flow
297 considering information sharing on traffic lights and AV.

298 Figure 6 shows that traffic flow containing 100% HDV (i.e., 0% AV) has the worst total
299 delay rewards in the long term compared to 5%, 10%, and 20% penetration rates of AV.

300 Penetration of AVs at 5%, 10%, and 20% results in convergence of rewards on total delay.

301 The 20% AV penetration rates show a more complicated learning process due to the
302 priority of safety over the optimization of traffic flow speed and FEC rate at an early stage

303 of the learning process. Once AV become familiar with the traffic flow patterns during

304 the training period, the PPO algorithm improves the rewards on total delays due to good

305 prediction of other vehicles' behavior. The 10% penetration rate of AV indicates

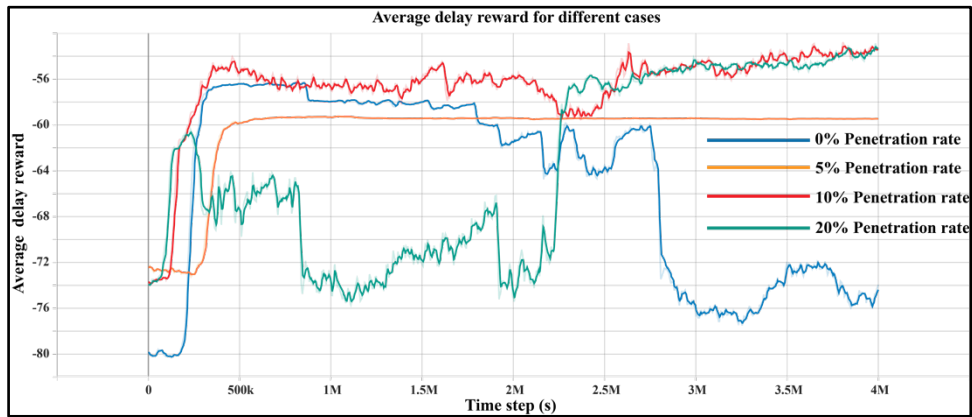
306 fluctuations as well during the training period and it could be due to fewer interactions

307 between AV and HDV on the road.

308 Table 3 shows the average delay for different penetration rates of AV in the selected

309 traffic grid network. At a 10% penetration rate, the average reward was achieved at the least

310 time step of 110K as compared to other penetration rates.



311

312 Figure 6. Behavior of average rewards of total delay with respect to time steps for different AV penetration

313 rates.

314 Table 3. Convergence time and steady state average rewards on total delay obtained with different

315 penetration rates.

Penetration rate	Approximate starting time steps of convergence	Average rewards on total delay at the last time step
0%	No exact convergence observed	-73.94
5%	302K	-59.79
10%	110K	-53.34
20%	2.24M	-53.51

316

317 **FEC Rates at Different Penetration Levels**

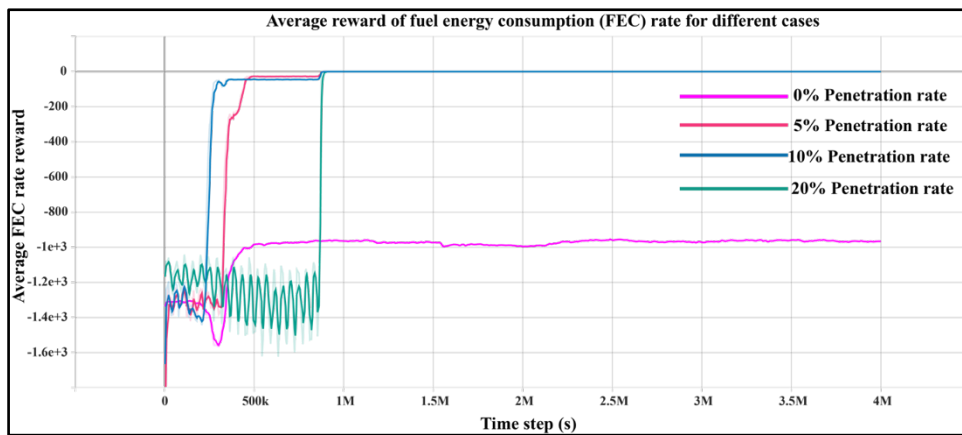
318 The FEC rate of a small-engine vehicle usually falls within the range of 0.05-0.10 L/s

319 with an average driving velocity. The reward on FEC rate can reach zero when the vehicle

320 achieves low levels of FEC at a minimum varying speed and other performance

321 parameters. Results of the average rewards on FEC rate obtained from different

322 penetration levels are presented in Figure 7 and Table 4. The penetration of AV shows
 323 better performance in reaching larger rewards on FEC rates, while the pure HDV case
 324 illustrates the worst scenario with a reward on FEC rate of about -1,000. Interestingly, the
 325 10% penetration rate regulates the FEC rate faster in the simulation compared to the results
 326 obtained from 5% and 20% penetration of AV. The time step for the convergence of the
 327 FEC rate and the steady state average rewards on the FEC rate with 4 different penetration
 328 rates are presented in Table 4.



329
 330 Figure 7. Behavior of average rewards on FEC rates with respect to time steps, considering 4 different
 331 penetration rates of AV in the traffic flow.

332 Table 4. FEC rate results at different penetration rates.

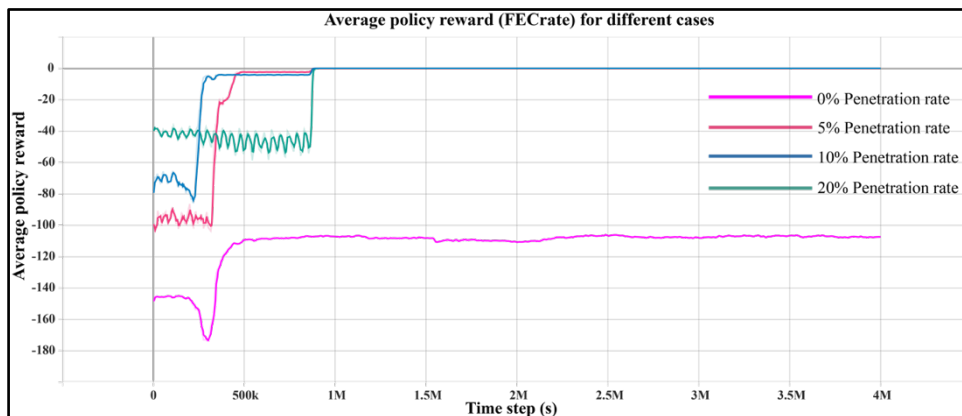
Penetration rate	Starting time step of convergence	Average rewards on FEC rate at the last time step
0%	302K	-964.3
5%	316k	0.071

10%	216k	0.010
20%	862K	0.081

333

334 **Performance of PPO policy**

335 To assess the effectiveness of the PPO policy in optimizing the FEC rate, the average
 336 policy reward and average environment time, policy loss, entropy, and value function loss
 337 have been evaluated with different penetration rates of AV. Figure 8 shows that average
 338 policy rewards with penetration rates 5%, 10%, and 20% of AV converge to zero while
 339 HDV only case has the worst reward with a value of -107.2. With the 10 % penetration
 340 rate, the average policy rewards start to converge about 300K steps, faster than 0%, 5%,
 341 and 20% penetration rates.



342

343 Figure 8. Behavior of average policy rewards with respect to time steps for 4 different penetration rates of

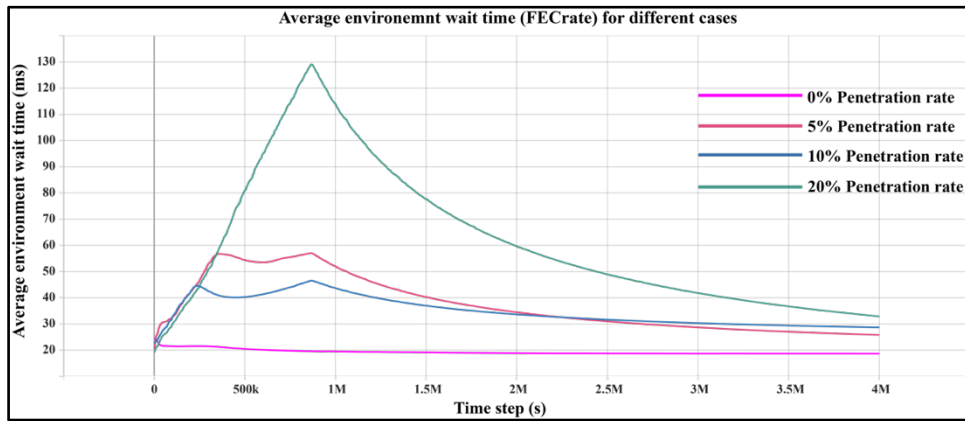
344 AV.

345

346 The time for an agent to stay in a specific state before applying an action is considered an
347 environment waiting time. The highest average environment time is observed for the 20%
348 penetration rates at a time step of about 850K as compared to others, as shown in Figure
349 9. The average environment waiting time of 5% case is slightly higher than that of the
350 10% AV penetration rate by the end of training, but its highest peak is observed to be
351 higher than the 10% penetration rates during training at 800K steps.

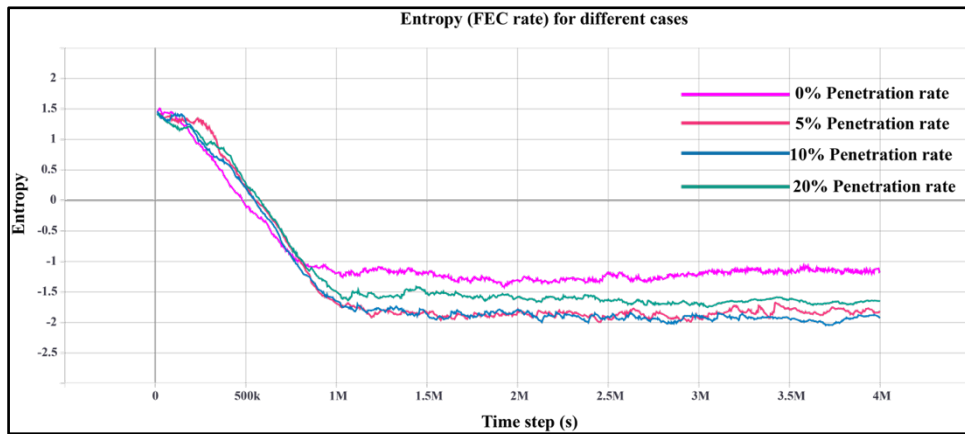
352 The entropy behavior for PPO is shown in Figure 10, with high values observed at a 0%
353 penetration rate. A minimum value of entropy was observed at a 10% penetration rate by
354 the end of training, indicating fewer uncertainties in policy distribution as compared to
355 the 5% and 20% penetration rates.

356 The total loss, which is a combination of policy loss and value function loss, is depicted
357 in Figure 11. At a 10% penetration rate, the minimum total loss is observed compared to
358 other cases. All these performance indices show that penetration of AV can improve the
359 rewards on FEC rates compared with 100% HDV traffic flow.



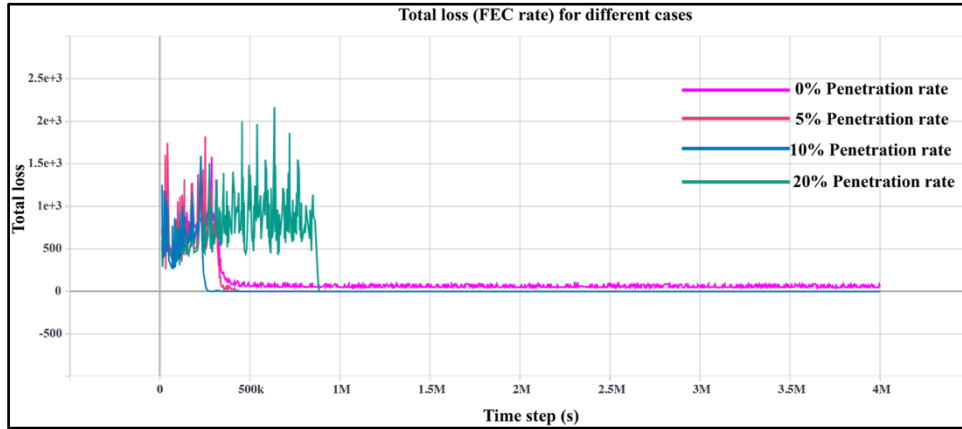
360

361 Figure 9. Average environment waiting time with respect to time steps for 4 different penetration rates of
 362 AV.



363

364 Figure 10. Entropy of the policy to optimize FEC rates for 4 different penetration rates of AV.



365

366 Figure 11. Total loss including policy value function loss with respect to time steps for 4 penetration rates
 367 of AV.

368 Table 5 shows the values of five measurements of policy to optimize FEC rate for 4
 369 different penetration rates of AV.

370 Table 5. FEC rate results at different penetration rates

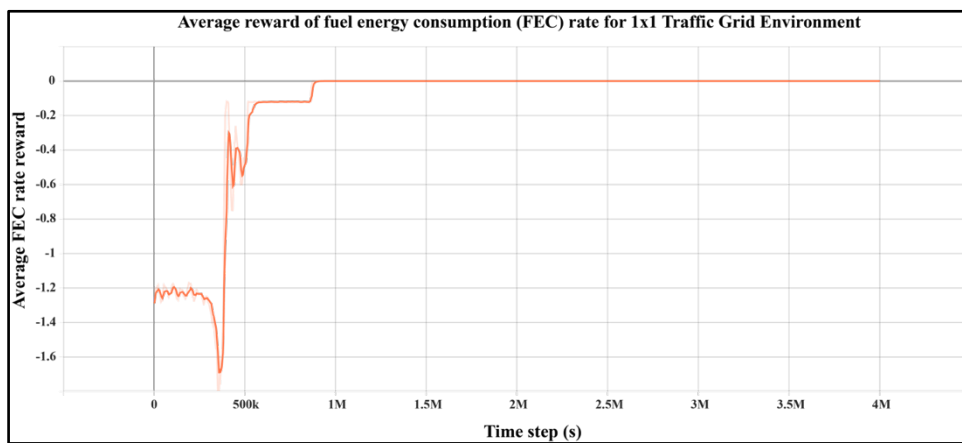
Measurements	Penetration Rate			
	0%	5%	10%	20%
Average policy reward	-107.2	0.059	0.057	0.069
Average environment Wait time (ms)	18.7	25.81	28.69	32.87
Policy loss	6.724e-3	6.07e-3	9.225e-3	6.801e-3
Entropy	-1.158	-1.829	-1.924	-1.65
Value function loss	52.66	7.05e-7	6.268e-7	1.622e-6

371

372 Comparison with Other Traffic Grid Environments

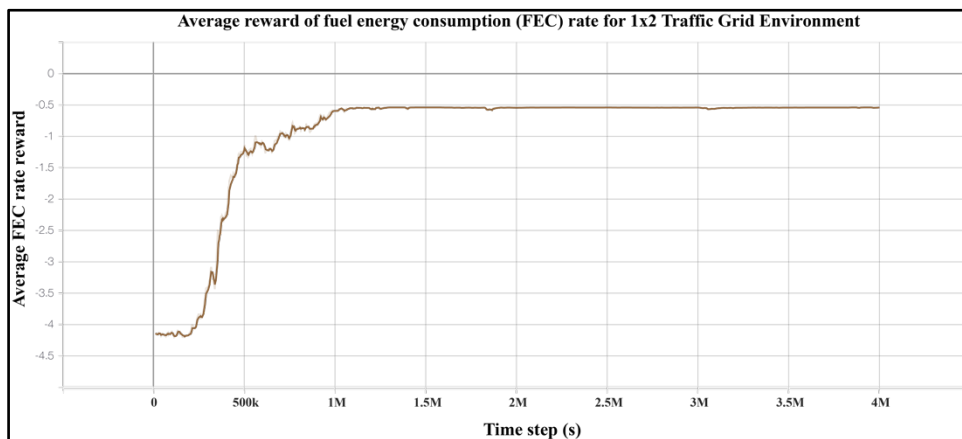
373 With a 10% penetration rate of AV, four traffic environments including 1x1, 1x2, 2x2,
 374 and 3x3 traffic grids were simulated with the proposed PPO algorithm. Figures 12 to 15

375 show the behavior of the average rewards on FEC rates with respect to time steps for each
 376 simulated environment, respectively. Table 6 presents the convergence of average
 377 rewards on FEC rates and the convergence time for each simulated environment.
 378 Specifically, the average FEC reward in the 3X3 traffic grid converged at about 216K
 379 steps, which was less than results obtained from other environments.



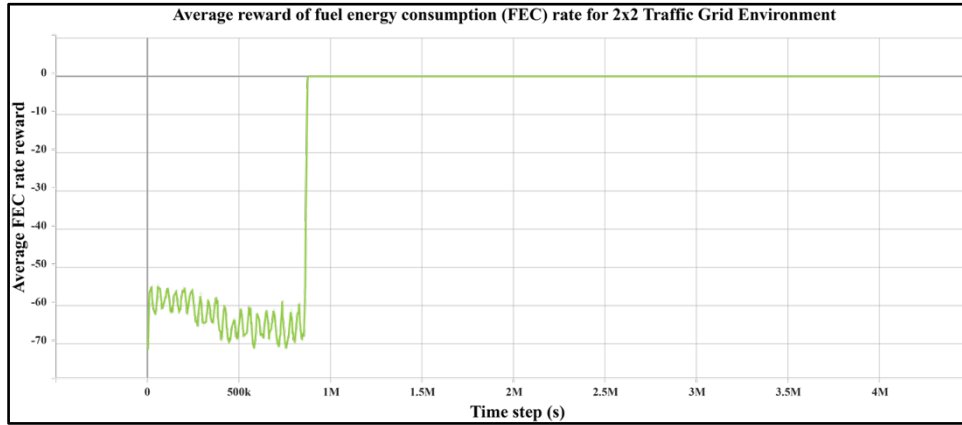
380

381 Figure 12. Behavior of average rewards on FEC rate with respect to time steps for a 1x1 traffic grid.



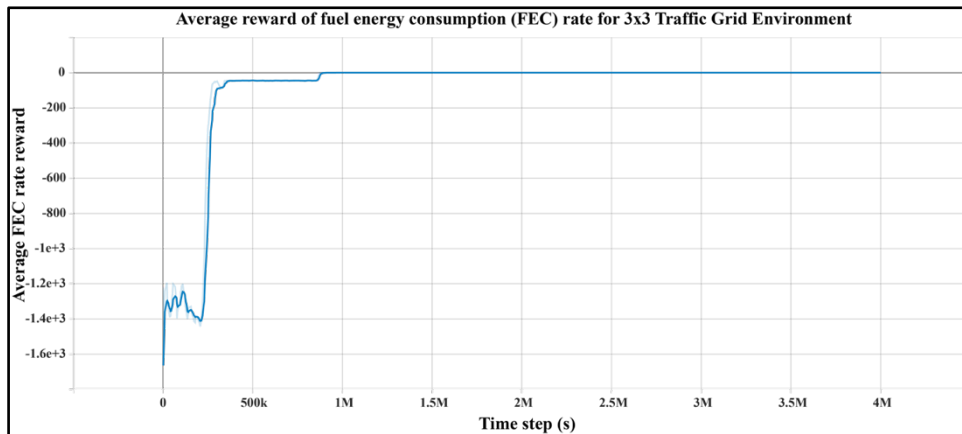
382

383 Figure 13. Behavior of average rewards on FEC rate with respect to time steps for a 1x2 traffic grid.



384

385 Figure 14. Behavior of average rewards on FEC rate with respect to time steps for a 2x2 traffic grid.



386

387 Figure 15. Behavior of average rewards on FEC rate with respect to time steps for a 3x3 traffic grid.

388 Table 6. Performance of the PPO algorithm for rewards on FEC rates for different traffic grid networks.

Traffic Grid	Converging Time steps	Average rewards on FEC rates at convergence
1x1	550K	-0.1
1x2	944K	-0.7231
2x2	854K	-69.4
3x3	216K	-20.1

389 **3. Discussion**

390 As more cars run on fuel like gasoline, resulting in air pollutants, the demand for eco-
391 driving strategies is highly.^{30 31} In this study, we employed an RL algorithm, PPO, to
392 investigate the impact of introducing AV to the traffic flow of HDV on reducing traffic
393 delay and minimizing fuel energy consumption rates. This involved introducing a specific
394 penetration rate of AV into a continuous traffic flow, coordinated with traffic light signals,
395 within a large 3x3 traffic grid system. The Flow computational package, developed in
396 Python, was utilized to integrate the publicly available microscopic traffic simulator,
397 SUMO, and the RL library 'RLlib'. In a previous study,³² a comparison between different
398 types of action spaces for different algorithms has been presented. Algorithms such as Q-
399 Learning, DQN, DDPG, etc. are considered reliable for specific types of action spaces—
400 either continuous or discrete. For environments that feature both continuous and discrete
401 action spaces, PPO-based RL algorithms are feasible due to their computational and
402 sample efficiency. So, the PPO-based RL approach is used in this research work to train
403 agents in the selected traffic environment.

404 The environmental setup consists of a traffic network with 3 horizontal and 3 vertical
405 roads in the SUMO simulator. Different percentages of AV (0%, 5%, 10%, and 20%)
406 were introduced in this study to control the speed of HDV in the network. The average
407 rewards on total delay and FEC rates were computed in this research work with different

408 penetration rates. The penetration of AV illustrated better average rewards on both total
409 delay and FEC rates. The 20% AV penetration initially results in more delays due to their
410 prioritization of safety over speed and efficiency. Specifically, a 10% penetration rate in
411 AV combined with HDV showed significant results for minimization of FEC rate and
412 total delay. The rewards on total delay for the 10% penetration rate case converged at a
413 minimum value of -53.34 at the least time steps of 110K in comparison with other cases.
414 At a 0% penetration rate, an average reward on FEC rates of -964.3 was obtained by the
415 end of training. For all other cases, the rewards on FEC rates approached zero by the end
416 of training. To assess the performance of the PPO policy in training agents to minimize
417 the FEC rates, results for average policy reward, entropy, value function loss, and mean
418 environment time were obtained at various penetration rates. A 10% penetration rate
419 demonstrated better performance compared to 0%, 5%, and 20% rates. A comparison of
420 four traffic light grids (1x1, 1x2, 2x2, and 3x3) was performed at a penetration level of
421 10% in terms of the FEC rate. The results indicated that the average rewards on FEC rates
422 converged in a shorter time for a 3x3 traffic network as compared to other configurations.
423 We are well aware that there are limitations in this study. PPO-based RL algorithms need
424 to be precisely tuned to achieve the most effective learning results because they are
425 hyperparameter-sensitive. The consideration of lane change was not incorporated into this
426 research; lane-changing behavior can be discussed by introducing a lane change

427 controller in the future. All HDV are assumed to have the same economic model, while
 428 there are heavy-duty vehicles, buses, cycles, and passenger vehicles have different
 429 economic models. To address this limitation, we can find an economical model for each
 430 type of vehicle, determine the percentage of each vehicle type based on the public traffic
 431 flow dataset, and integrate this information into the energy calculation in our future
 432 research. We assume ideal communication without any delay or failure in controlling AV
 433 in this study. In the future, we can apply RL algorithms to scaled-down AV to examine
 434 the impact of communication delays.

435 A comparison of the proposed eco-driving strategy is performed with prior related
 436 research as shown in Table 7, suggesting the effectiveness of the proposed PPO algorithm
 437 for an eco-driving strategy.

438 Table 7. Comparison with prior research work.

Reference	Year	Vehicle Type	Algorithm	Action Space Type	Traffic Grid Network	Objective	Fuel Consumption (L) per Vehicle
17	2018	CV	Q-Learning	Discrete	Cases-1 single intersection (1x1) Case-2: a 2-way road network	To minimize CO2 emissions and optimize traffic performance	N/A
18	2020	CAV	DQN	Discrete	1x1	Optimizing acceleration/deceleration of CAV to minimize fuel consumption	0.0691
20	2020	CAV	DDPG	Continuous	Circular Network with signalized interactions	To enhance travel efficiency, reduce fuel consumption, and ensure safety	0.015 (at 100% CAV)
21	2021	CAV	DDPG +DQN	Discrete & Continuous	1x5	To minimize fuel consumption, ensure reasonable travel times,	0.12005 (for HDQPG)

22	2019	CAV and HDV	TRPO	Continuous	1x1	and execute lane changes strategically to avoid congested lanes Percentage of AV in the traffic flow to minimize fuel consumption, emissions, and improvement in travel speed	0.0954 (at 100% CAV)
This work	2024	AV and HDV	PPO	Continuous & Discrete	3x3	Collaboration of traffic lights signals and percentage of AV in the traffic flow to minimize fuel consumption and total delay	0.010 (at 10% AV)

439

440 4. Conclusions

441 Eco-driving positively impacts human health by reducing pollution resulting from vehicle
 442 fuel consumption and emissions. This study explores a hybrid traffic network that
 443 combines autonomous vehicles and human-driven vehicles through the coordination of
 444 traffic light signals to manage a large traffic flow. The approach addresses eco-driving
 445 challenges, including real-time traffic data collection and the intricate nature of the traffic
 446 network, which currently lacks a comprehensive mathematical model. The research
 447 employs model-free PPO-based reinforcement learning algorithms to analyze the fuel
 448 energy consumption rates of vehicles. It focuses on minimizing fuel energy consumption
 449 by introducing specific penetration rates of AV (0%, 5%, 10%, and 20%) in a 3x3 traffic
 450 grid system, utilizing the Flow compactional package to integrate the SUMO simulator
 451 and RLlib. The study results indicate that a 10% penetration rate of AV alongside HDV
 452 yielded significant reductions in both fuel consumption and total delay of traffic.

453 **Acknowledgement**

454 This research was partially supported by the National Science Foundation (2051113 to
455 YFJ) and the Department of Transportation (TranSET Program 21ITS034 and
456 21UTSA049 to YFJ).

457 **Author's Contribution**

458 Conceptualization, U.J., M.F. M.X, C.D., and Y.J.; Data Curation, U.J., M.F., and Y.J.;
459 Methodology, U.J. M.F. and Y.J.; Simulations, U.J., A.C., and M.M.; Supervision, Y.J.;
460 Writing-original Draft Preparation, U.J., and Y.J.; Writing-Review and Editing, M.M,
461 M.F., C.D., and M.X.; Funding, Y.J. All authors have reviewed the manuscript and
462 provided their consent for publication.

463 **Conflict of Interests Statement**

464 The authors declared no conflicts of interest for this publication.

465 **Funding**

466 This work was supported in part by the U.S. National Science Foundation #2051113, and
467 US Department of Transportation for Transportation Consortium of South-Central States
468 (TranSET) with projects 21034 and 21049. The funding sources had no role in the design
469 of the study; collection, analysis, and interpretation of data; or in the writing of the
470 manuscript.

471 **ORCID iD**472 Umar Jamil <https://orcid.org/0000-0002-2346-556X>473 Yu-Fang Jin <https://orcid.org/0000-0002-7421-527X>474 **References**

- 475 1. U.S. Energy Information Administration (EIA). Use of energy explained Energy
 476 use for transportation. Accessed October 30, 2023,
 477 [https://www.eia.gov/energyexplained/use-of-](https://www.eia.gov/energyexplained/use-of-energy/transportation.php#:~:text=Energy%20sources%20are%20used%20in,and%20some%20types%20of%20helicopters)
 478 [energy/transportation.php#:~:text=Energy%20sources%20are%20used%20in,and%20so](https://www.eia.gov/energyexplained/use-of-energy/transportation.php#:~:text=Energy%20sources%20are%20used%20in,and%20some%20types%20of%20helicopters)
 479 [me%20types%20of%20helicopters](https://www.eia.gov/energyexplained/use-of-energy/transportation.php#:~:text=Energy%20sources%20are%20used%20in,and%20some%20types%20of%20helicopters).
- 480 2. Pasquale C, Sacone S, Siri S, Ferrara A. Traffic control for freeway networks with
 481 sustainability-related objectives: Review and future challenges. *Annual Reviews in*
 482 *Control*. 2019;48:312-324.
- 483 3. Boltze M, Tuan VA. Approaches to achieve sustainability in traffic management.
 484 *Procedia engineering*. 2016;142:205-212.
- 485 4. Jamil U, Sulaiman M, Ghafoor N, Malmir M, Nawaz F, Shakoore RI. Power
 486 Harvesting towards Sustainable Energy Technology through Ambient Vibrations and
 487 Capacitive Transducers. *IEEE*; 2023:1-6.
- 488 5. Park S, Rakha H, Ahn K, Moran K. Virginia tech comprehensive power-based
 489 fuel consumption model (VT-CPFM): Model validation and calibration considerations.
 490 *International Journal of Transportation Science and Technology*. 2013;2(4):317-336.
- 491 6. Lee JW, Gunter G, Ramadan R, et al. Integrated Framework of Vehicle Dynamics,
 492 Instabilities, Energy Models, and Sparse Flow Smoothing Controllers. *arXiv preprint*
 493 *arXiv:210411267*. 2021;
- 494 7. Yao Z, Hu R, Wang Y, Jiang Y, Ran B, Chen Y. Stability analysis and the
 495 fundamental diagram for mixed connected automated and human-driven vehicles.
 496 *Physica A: Statistical Mechanics and Its Applications*. 2019;533:121931.
- 497 8. Hao P, Wei Z, Bai Z, Barth MJ. *Developing an adaptive strategy for connected*
 498 *eco-driving under uncertain traffic and signal conditions*. 2020.
- 499 9. Wang S, Lin X. Eco-driving control of connected and automated hybrid vehicles
 500 in mixed driving scenarios. *Applied Energy*. 2020;271:115233.
- 501 10. Zhou S, Xie M, Jin Y, Miao F, Ding C. An End-to-end Multi-task Object
 502 Detection using Embedded GPU in Autonomous Driving. 2021:122-128.

- 503 11. Wei H, Zheng G, Gayah V, Li Z. Recent advances in reinforcement learning for
504 traffic signal control: A survey of models and evaluation. *ACM SIGKDD Explorations*
505 *Newsletter*. 2021;22(2):12-18.
- 506 12. Clemmons J, Jin YF. Reinforcement Learning-Based Guidance of Autonomous
507 Vehicles. 2023:1-6.
- 508 13. Zeynivand A, Javadpour A, Bolouki S, et al. Traffic flow control using multi-
509 agent reinforcement learning. *Journal of Network and Computer Applications*.
510 2022;207:103497.
- 511 14. Li M, Cao Z, Li Z. A reinforcement learning-based vehicle platoon control
512 strategy for reducing energy consumption in traffic oscillations. *IEEE Transactions on*
513 *Neural Networks and Learning Systems*. 2021;32(12):5309-5322.
- 514 15. Haydari A, Zhang M, Chuah C-N, Ghosal D. Impact of deep rl-based traffic signal
515 control on air quality. *IEEE*; 2021:1-6.
- 516 16. Stern RE, Chen Y, Churchill M, et al. Quantifying air quality benefits resulting
517 from few autonomous vehicles stabilizing traffic. *Transportation Research Part D:*
518 *Transport and Environment*. 2019;67:351-365.
- 519 17. Shi J, Qiao F, Li Q, Yu L, Hu Y. Application and evaluation of the reinforcement
520 learning approach to eco-driving at intersections under infrastructure-to-vehicle
521 communications. *Transportation Research Record*. 2018;2672(25):89-98.
- 522 18. Mousa SR, Ishak S, Mousa RM, Codjoe J, Elhenawy M. Deep reinforcement
523 learning agent with varying actions strategy for solving the eco-approach and departure
524 problem at signalized intersections. *Transportation research record*. 2020;2674(8):119-
525 131.
- 526 19. Vinitzky E, Kreidieh A, Le Flem L, et al. Benchmarks for reinforcement learning
527 in mixed-autonomy traffic. *PMLR*; 2018:399-409.
- 528 20. Zhou M, Yu Y, Qu X. Development of an efficient driving strategy for connected
529 and automated vehicles at signalized intersections: A reinforcement learning approach.
530 *IEEE Transactions on Intelligent Transportation Systems*. 2019;21(1):433-443.
- 531 21. Guo Q, Angah O, Liu Z, Ban XJ. Hybrid deep reinforcement learning based eco-
532 driving for low-level connected and automated vehicles along signalized corridors.
533 *Transportation Research Part C: Emerging Technologies*. 2021;124:102980.
- 534 22. Jayawardana V, Wu C. Learning eco-driving strategies at signalized intersections.
535 *IEEE*; 2022:383-390.
- 536 23. Lopez PA, Behrisch M, Bieker-Walz L, et al. Microscopic traffic simulation using
537 SUMO. *IEEE*; 2018:2575-2582.
- 538 24. OpenStreetMap Contributors. Manhattan City [Map]. Accessed June 22, 2023,
539 <https://www.openstreetmap.org/export#map=15/40.7867/-73.9533>

- 540 25. Federal Highway Administration. Chapter 12 - Signalized Intersections:
 541 Informational Guide, August 2004 (Publication Number: FHWA-HRT-04-091).
 542 Accessed November 22, 2022, 2022.
 543 <https://www.fhwa.dot.gov/publications/research/safety/04091/12.cfm>
 544 26. McKinsey and Company. Autonomous driving's future: Convenient and
 545 connected. Accessed January 20, 2023, 2023.
 546 [https://www.mckinsey.com/industries/automotive-and-assembly/our-](https://www.mckinsey.com/industries/automotive-and-assembly/our-insights/autonomous-drivings-future-convenient-and-connected)
 547 [insights/autonomous-drivings-future-convenient-and-connected](https://www.mckinsey.com/industries/automotive-and-assembly/our-insights/autonomous-drivings-future-convenient-and-connected)
 548 27. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. Proximal policy
 549 optimization algorithms. *arXiv preprint arXiv:170706347*. 2017;
 550 28. Liang E, Liaw R, Nishihara R, et al. RLlib: Abstractions for distributed
 551 reinforcement learning. PMLR; 2018:3053-3062.
 552 29. Wu C, Kreidieh AR, Parvate K, Vinitsky E, Bayen AM. Flow: A modular learning
 553 framework for mixed autonomy traffic. *IEEE Transactions on Robotics*. 2021;38(2):1270
 554 - 1286.
 555 30. Huang Y, Ng EC, Zhou JL, Surawski NC, Chan EF, Hong G. Eco-driving
 556 technology for sustainable road transport: A review. *Renewable and Sustainable Energy*
 557 *Reviews*. 2018;93:596-609.
 558 31. Malmir M, Momeni H, Ramezani A. Controlling Megawatt Class WECS by
 559 ANFIS Network Trained with Modified Genetic Algorithm. *IEEE*; 2019:939-943.
 560 32. Zhu J, Wu F, Zhao J. An overview of the action space for deep reinforcement
 561 learning. 2021:1-10.
 562

563

564 **Author Biographies**

565 **Umar Jamil** is a Doctoral Candidate and Graduate Research Assistant in the Electrical
 566 and Computer Engineering Department at the University of Texas at San Antonio, USA.
 567 He received his MS in Mechatronics Engineering from Air University, Islamabad,
 568 Pakistan. Additionally, he graduated with another MS degree in Electrical Engineering,
 569 and a BS degree in Electrical Engineering, focusing on Power Systems, from Mirpur

570 University of Science and Technology, Azad Jammu & Kashmir, Pakistan. His research
571 interests encompass intelligent transportation systems, deep learning, reinforcement
572 learning, electrical power distribution systems, sensors and actuators, and energy
573 harvesting systems.

574 **Mostafa Malmir** is a second-year Ph.D. student at the Electrical Engineering Department
575 of the University of Texas at San Antonio. His research concentration is studying the
576 advantages of different Machine Learning and Deep Learning algorithms and explore
577 implementing them in biological applications. His current work examines novel methods
578 of single-cell cell typing to improve cell identification for rare cell types.

579 **Monika Filipovska** is an assistant professor of Transportation and Urban Engineering in
580 the Department of Civil and Environmental Engineering at the University of Connecticut.
581 She received her Ph.D. and MS in Civil and Environmental Engineering, focusing on
582 Transportation Systems Analysis and Planning, from Northwestern University. Her
583 research interests are in the domains of transportation networks and traffic modeling, with
584 a focus on emerging vehicles and sensing technologies and data-driven predictive
585 analytics.

586 **Mimi Xie** is an assistant professor in Computer Science at the University of Texas at San
587 Antonio. She received the BE and MS degrees from the College of Computer Science,
588 Chongqing University, Chongqing, China, in 2010 and 2013, respectively, and the Ph.D.

589 degree in Electrical and Computer Engineering from the University of Pittsburgh in 2019.
590 She is currently an assistant professor with the Department of Computer Science,
591 University of Texas at San Antonio, San Antonio, TX. Her current research interests
592 include energy-harvesting embedded Systems and AI on Edge.

593 **Caiwen Ding** received the Ph.D. degree in computer science and engineering from the
594 Northeastern University, Boston, MA, USA, in 2019. He is currently an Assistant
595 Professor in the Department of Computer Science and Engineering at the University of
596 Connecticut, Storrs, CT, USA. His research interests include machine learning and deep
597 neural network systems, computer vision, and natural language processing. Dr. Ding is
598 the recipient of the Best Paper Award Nomination at DATE 2018 and DATE 2021.

599 **Yu-Fang Jin** received her Ph.D. degree in Electrical and Computer Engineering from the
600 University of Central Florida, Orlando, Florida, USA, in 2004. She is currently a Full
601 Professor in the Department of Electrical and Computer Engineering at the University of
602 Texas at San Antonio, San Antonio, Texas, USA. Her research interests include
603 applications of deep learning to large-scale networked systems and interpretations of deep
604 learning algorithms.

605

606

607

608

609

610