

1 **Title**

2 Dense reinforcement learning for safety validation of autonomous vehicles

3 **Authors**

4 Shuo Feng^{1,2,4}, Haowei Sun¹, Xintao Yan¹, Haojie Zhu¹, Zhengxia Zou^{1,5}, Shengyin Shen², Henry X. Liu^{1,2,3*}

5 **Affiliations**

6 ¹Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI, USA

7 ²University of Michigan Transportation Research Institute, Ann Arbor, MI, USA

8 ³Mcity, University of Michigan, Ann Arbor, MI, USA

9 ⁴Present address: Department of Automation, Tsinghua University, Beijing, China

10 ⁵Present address: School of Astronautics, Beihang University, Beijing, China

11 *Corresponding Author, henryliu@umich.edu

12 **Abstract**

13 One critical bottleneck that impedes autonomous vehicle (AV) development and deployment is the
14 prohibitively high economic and time costs required to validate its safety in a naturalistic driving environment,
15 due to the rarity of safety-critical events¹. Here we develop an intelligent testing environment in that artificial
16 intelligence-based background agents are trained to test AVs in an accelerated mode, without loss of
17 unbiasedness. From naturalistic driving data, the background agents learn when to execute what adversarial
18 maneuver through a newly developed dense deep reinforcement learning (D2RL) approach, in which Markov
19 decision processes are edited by removing non-safety-critical states and reconnecting critical ones so that the
20 information in the training data is densified. D2RL enables neural networks to learn from densified
21 information with safety-critical events and achieves tasks that are intractable for the traditional deep
22 reinforcement learning approach. We demonstrate the effectiveness of our approach by testing a highly
23 automated vehicle in both highway and urban test tracks with an augmented reality environment, combining
24 simulated background vehicles with physical road infrastructure and a real AV under test. Our results show
25 that the D2RL-trained agents can accelerate the evaluation process by multiple orders of magnitude (10^3 to
26 10^5 times faster). D2RL also opens the door for accelerated testing and training with other safety-critical
27 autonomous systems.

28 **Introduction**

29 Driven by the rapid development of autonomous vehicle (AV) technologies, we are on the cusp of a new
30 revolution in transportation on a scale not seen since the introduction of automobiles a century ago. AV
31 technologies have the potential to significantly improve transportation safety, mobility, and sustainability,
32 thus attracting worldwide attention from industries, government agencies, professional organizations, and
33 academic institutions. Over the past 20 years, significant progress has been made on the development of AVs,
34 particularly with the emergence of deep learning². By 2015, several companies had announced that they would
35 be mass producing AVs before 2020^{3,5}. So far, the reality has not lived up to these expectations, and no level
36 4 (ref.⁶) AVs are commercially available. The reasons for this are numerous. But above all, the safety
37 performance of AVs is still significantly below that of human drivers. For average drivers in the United States,
38 the occurrence probability of a crash is around 1.9×10^{-6} per mile in the naturalistic driving environment
39 (NDE)¹. In contrast, the disengagement rate for the state-of-the-art AV is around 2.0×10^{-5} per mile,
40 according to the 2021 Disengagement Report from California⁷. Although the disengagement rate is criticized
41 for its potential biasedness, it has been widely used to track the trend of AV safety performance^{8,9}, as it is
42 arguably the only statistics where the results of different AVs are available to the public.

48 One critical bottleneck to improving AV safety performance is the severe inefficiency of safety validation.
49 Prevailing approaches usually test AVs in the NDE through a combination of software simulation, closed test
50 track, and on-road testing. However, to validate the safety performance of AVs at the level of human drivers,
51 it is well-known that hundreds of millions of miles and sometimes hundreds of billions of miles would be
52 required to test in the NDE¹. Due to this severe inefficiency, AV developers must pay significant economic
53 and time costs to evaluate each new development, which has severely hindered the progress of AV deployment.
54 To improve the testing efficiency, many approaches test AVs in purposely generated scenarios that are more
55 safety critical^{10,11}. Yet, existing scenario-based approaches¹²⁻¹⁷ can mainly be applied to short scenario
56 segments with limited background road users (see Supplementary Materials for more discussions).

57 Validating the safety performance of AVs in NDE is in essence a rare-event estimation problem in a high-
58 dimensional space. The main challenge is caused by the compounding effects of the “curse of rarity” in
59 addition to the “curse of dimensionality” (Fig. 1a). By “curse of dimensionality,” we mean that driving
60 environments could be spatiotemporally complex, and the variables needed to define such environments are
61 high dimensional. As the volume of the variable space grows exponentially with dimensionality, the
62 computational complexity also grows exponentially¹⁸. By “curse of rarity,” we mean that the occurrence
63 probability for safety-critical events is rare, i.e., most points of the variable space are non-safety-critical, which
64 provide no or noisy information for training. Under this circumstance, it is hard for a deep learning model to
65 learn even given a large amount of data, as the precious information (e.g., policy gradient) of safety-critical
66 events could be buried under the large amount of non-safety-critical data. The past decades have witnessed
67 rapid progress in the ability of artificial intelligence (AI) systems to solve problems with the “curse of
68 dimensionality”¹⁹, for example, the board game Go has a state space of 10^{360} (ref.²⁰) and the semiconductor
69 chip design may have a state space on the order of 10^{2500} (ref.²¹). Prior to this work, however, solving the
70 “curse of dimensionality” and the “curse of rarity” simultaneously has remained an open question, which has
71 impeded the applicability of AI techniques in safety-critical systems, such as AVs, medical robots, and
72 aerospace systems²².

73 We address this challenge by developing the dense deep reinforcement learning (D2RL) approach. The basic
74 idea is to identify and remove the non-safety-critical data and train neural networks utilizing only the safety-
75 critical data. As only a very small portion of data is safety-critical, the information of the remaining data will
76 be significantly densified. Essentially, the D2RL approach edits the Markov decision process by removing
77 the uncritical states and reconnecting the critical states, and then trains neural networks only for the edited
78 Markov process (Fig. 1b). Therefore, for any training episode, the reward from the end state is backpropagated
79 along the edited Markov Chain with critical states only (Fig. 1c). The D2RL approach can dramatically reduce
80 the variance of the policy gradient estimation with multiple orders of magnitude without loss of unbiasedness,
81 compared with the DRL approach, as proved in Theorem 1 in Methods. Such significant variance reduction
82 can enable neural networks to learn and achieve tasks that are intractable for the DRL approach. For AV
83 testing, we leverage the D2RL approach and train the background vehicles (BVs) through a neural network
84 to learn when to execute what adversarial maneuver, which aims to improve the testing efficiency and ensure
85 evaluation unbiasedness. This results in an AI-based adversarial testing environment that can reduce the
86 required testing miles of AVs by multiple orders of magnitude while ensuring the testing unbiasedness. Our
87 approach can be applied to complex driving environments including multiple highways, intersections, and
88 roundabouts, which cannot be achieved by prior scenario-based approaches. The proposed approach
89 empowers the testing agents in the environment with intelligence to create an intelligent testing environment,
90 i.e., using AI to validate AI. This is a paradigm shift and it opens the door for accelerated testing and training
91 with other safety-critical systems.

93 To demonstrate the effectiveness of our AI-based testing approach, we trained the BVs with large-scale
94 naturalistic driving datasets and conducted simulation experiments as well as field experiments in physical
95 test tracks. Specifically, we tested a level 4 AV with an open-source automated driving system, Autoware²³,
96 in the physical 4-km-long highway test track at the American Center for Mobility (ACM) and the urban test
97 track at Mcity. To test the AV with the D2RL-trained testing environment safely and precisely, we developed
98 an augmented reality testing platform²⁴, which combines the physical test track and a microscopic traffic
99 simulator, SUMO²⁵. As shown in Fig. 1d, by synchronizing the movements of the real AV and virtual BVs,
100 the real AV in the physical test track can interact with the virtual BVs as if it is in a realistic traffic
101 environment, where the BVs are directed to interact with the real AV. For both simulation and field
102 experiments, we evaluated not only crash rates, but also crash types and crash severities. Our simulation and
103 field-testing results show that the D2RL approach can effectively learn the intelligent testing environment,
104 which can significantly accelerate the evaluation process of AVs by multiple orders of magnitude (10^3 to 10^5
105 times faster) unbiasedly, compared with the results from testing AVs directly in NDE.

106 Results

107 Dense deep reinforcement learning (D2RL)

108 To leverage AI techniques, we formulate the AV testing problem as a sequential Markov decision process
109 (MDP), where maneuvers of BVs are decided based on the current state information. We aim to train a policy
110 (a DRL agent) modeled by a neural network, which can control the maneuvers of BVs to interact with the
111 AV, to maximize the evaluation efficiency and ensure unbiasedness. However, as mentioned earlier, it is
112 hard—or even empirically infeasible—to learn an effective policy if directly applying DRL approaches
113 because of the “curse of dimensionality” and the “curse of rarity.”

114 We address this challenge by developing the D2RL approach. Due to the rarity of safety-critical events, most
115 states are uncritical and cannot provide information for safety-critical events, so the key concept of D2RL is
116 to remove the data of these uncritical states and only utilize the informative data for training the neural network
117 (Fig. 1b and 1c). For AV testing problems, many safety metrics²⁶ can be utilized to identify the critical states
118 with different efficiency and effectiveness. In this study, we utilize the criticality measure^{12,13}, which is an
119 outer approximation of the AV crash rate within a specific time horizon (e.g., one second) from the current
120 state. Theoretical analysis for more generic problems can be found in Methods and Supplementary Materials
121 (Section 2a). We then edit the Markov process, discard the data of uncritical states, and use the remaining
122 data for the policy gradient estimation and bootstrapping of the DRL training. We discover that dense learning
123 can dramatically reduce the variance of the policy gradient estimation with multiple orders of magnitude
124 without loss of estimation unbiasedness, as proved in Theorem 1 in Methods. The dense learning can also
125 reduce the bootstrapping variance, as it can be regarded as a state-dependent temporal-difference learning²⁷,
126 where only critical states are utilized and others are skipped.

127 To demonstrate the effectiveness of dense learning, we compared D2RL with the DRL approach for a corner
128 case generation problem^{28,29}, which can be formulated as a well-defined reinforcement learning problem. A
129 neural network was trained to maximize the AV’s crash rate by controlling the closest eight BVs’ actions
130 (Fig. 2a). We used proximal policy optimization (PPO)³⁰ to update the parameters of the policy network, given
131 the reward for each testing episode, i.e., +20 for an AV crash and 0 for others. For a fair comparison, the only
132 difference between DRL and D2RL is that DRL utilized all the data for training the neural network, while
133 D2RL only utilized the data of critical states. As shown in Fig. 2b, D2RL removed the data of 80.5% complete
134 episodes and 99.3% steps from uncritical states, compared with DRL. According to Theorem 1, this indicates
135 that D2RL can reduce around 99.3% of the policy gradient estimation variance, which enables the neural
136 network to learn effectively. Specifically, the D2RL can maximize the reward during the training process,
137 while the DRL was stuck from the beginning of the training process (Fig. 2c). The policy learned by D2RL

138 can effectively increase the crash rate of the AV, while DRL failed to do so (Fig. 2d). Figure 2e-g illustrate
139 three generated corner cases.

141 Learning the intelligent testing environment

142 Learning the intelligent testing environment for unbiased and efficient AV evaluation is much more complex
143 than corner case generation. According to the importance sampling theory³¹, the goal is essentially to learn
144 new sampling distributions, i.e., importance function, of BVs' maneuvers to replace their naturalistic ones,
145 with the aim of minimizing the estimation variance of AV testing. Intuitively, the BVs are trained to learn
146 when to execute what adversarial maneuver, in that all BVs follow naturalistic behaviors, only selected
147 vehicles at selected moments execute specifically designed adversarial moves with a learned probability. To
148 achieve this goal, without using any heuristics or handcrafted functions, we derive the reward function from
149 the estimation variance as

$$150 \quad r(\mathbf{x}) = -\mathbb{I}_A(\mathbf{x}) \cdot W_{q_\pi}(\mathbf{x}) \cdot W_{q_{\pi_b}}(\mathbf{x}), \quad (1)$$

151 where \mathbf{x} denotes the variables of each testing episode, $\mathbb{I}_A(\mathbf{x})$ is an indicator function of the AV crash, and
152 $W_{q_\pi}(\mathbf{x}) = P(\mathbf{x})/q_\pi(\mathbf{x})$ and $W_{q_{\pi_b}}(\mathbf{x}) = P(\mathbf{x})/q_{\pi_b}(\mathbf{x})$ are weights (or likelihoods) produced by importance
153 sampling. Here, $P(\mathbf{x})$ denotes the naturalistic distribution, $q_\pi(\mathbf{x})$ denotes the importance function with the
154 target policy π , and $q_{\pi_b}(\mathbf{x})$ denotes the importance function with the behavior policy π_b . As there is no
155 heuristic or handcrafted immediate reward function, the reward function in Eq. (1) is highly consistent with
156 the testing performance, i.e., a higher reward indicates a more efficient testing environment. Such reward
157 design is generic and applicable to other rare event estimation problems with high-dimensional variables.

158 To determine the learning mechanism, we further investigate the relationship between the behavior policy π_b
159 and target policy π . As proved in Theorem 2 in Methods, we discover that the optimal behavior policy π_b^* that
160 collects data during the training process is nearly inversely proportional to the target policy. It indicates that,
161 if using on-policy learning mechanisms ($q_{\pi_b} = q_\pi$), the behavior policy would be far from optimality, which
162 could mislead the training process and eventually cause the underestimation issues. To address this issue, we
163 design an off-policy learning mechanism, where a generic behavior policy is designed and kept unchanged
164 during the training process. Although this off-policy mechanism is not the optimal behavior policy as in
165 Theorem 2 (which is usually unavailable in practice), it can balance the exploration and exploitation and is
166 empirically effective for all experiment settings in this study. With the reward function and off-policy learning
167 mechanism, we can learn the intelligent testing environment by the D2RL approach (see Methods for training
168 details).

169 AV testing in simulation

170 We evaluated the effectiveness of D2RL-based intelligent testing environment regarding accuracy, efficiency,
171 scalability, and generalizability by systematic simulation analysis. To measure the safety performance of AVs,
172 crash rates of different crash types and severities in NDE are utilized as the benchmark. As NDE is generated
173 completely based on naturalistic driving data, testing results in NDE can represent the safety performance of
174 AVs in the real world. For each test episode, we simulated AV driving in traffic for a fixed distance, and then
175 the test results were recorded and analyzed. To investigate the scalability and generalizability, we conducted
176 simulation experiments with different road geometries, different driving distances, and two different types of
177 AV models (i.e., the AV-I and AV-II models; see Section 3d in Supplementary Materials).

178 Figure 3 shows the results of the 2-lane highway environment with the 400m driving distance for the AV-I
179 model, which is a basic experiment to validate our approach. As shown in Fig. 3a, during the training process
180 the estimation variance of the intelligent testing environment decreases with the increase of reward function,
181 which demonstrates the effectiveness of the reward function in Eq. (1). To justify the off-policy mechanism,
182 we investigated the performance of the on-policy mechanism, where the target policy was utilized as the
183 behavior policy. As shown in Fig. 3b, during the training process, the crash rate for the on-policy experiments
184 significantly increases, while the crash rate for the off-policy experiments is unchanged because the behavior
185 policy is unchanged. However, as the on-policy mechanism breaks the consistency between the reward
186 function and estimation variance, this increase of the crash rate would be misleading. As shown in Fig. 3c,
187 the testing environment obtained by the on-policy mechanism underestimates the crash rate. In contrast, our
188 off-policy approach can obtain the same crash rate as the NDE approach, but more efficiently (Fig. 3d, e). To
189 measure the efficiency, we calculated the minimum number of tests for reaching a predetermined precision
190 threshold (the relative half-width^{12,17} is 0.3). To reduce the randomness of the results for a fair comparison,
191 we repeated the testing of our approach by bootstrap sampling and obtained the frequency and average of the
192 required number of tests (Fig. 3f). Compared with the NDE approach that required 1.9×10^8 number of tests,
193 our approach required an average of 9.1×10^4 number of tests, which is 2.1×10^3 times faster. To
194 investigate the generalizability, we further tested the AV-II model using the same intelligent testing
195 environment without any refinement, which can also obtain an accurate estimation with about 10^3 times faster
196 (see Section 4d in Supplementary Materials).

197 To validate the unbiasedness about crash types, crash severities, and near-miss events, we analyzed the crash
198 rates of different crash types, distribution of the speed difference at the crash moment, and distributions of the
199 time-to-collision (TTC), bumper-to-bumper distance, and post encroachment time (PET) of near-miss events,
200 respectively. Throughout the paper our use of the term unbiasedness refers to the fact that estimations from
201 our approach have the same mathematical expectations as those from NDE. In our experiments, we collected
202 about 2.34×10^8 episodes of tests in NDE and 3.15×10^6 (about two orders of magnitude less) episodes of
203 tests in the intelligent testing environment. As the intelligent testing environment is more adversarial than
204 NDE, the total crash rate in our approach is 3.21×10^{-3} (Fig. 3g), which is much higher than that
205 (1.58×10^{-7}) in NDE. As required by the importance sampling theory, each crash event should be weighted
206 by the likelihood ratio to keep the unbiasedness. Therefore, the weighted crash rates for all crash types are
207 compared with the results in NDE (Fig. 3h), which demonstrates the unbiasedness of our approach within the
208 evaluation precision. Similarly, Figures 3.i-1 demonstrate that our approach can also unbiasedly evaluate the
209 AV's safety performance regarding crash severities and near-miss events within the evaluation precision. As
210 near-miss events are critical for the development of AVs, the generated near-miss events without loss of
211 unbiasedness open the door for accelerating the AV training. We leave that for future study.

212 To further investigate the scalability and generalizability, we conducted the experiments with different
213 numbers of lanes (2 and 3 lanes) and driving distances (400m, 2km, 4km, and 25km) for the AV-I model.
214 Here we studied the 25km case to demonstrate the effectiveness of our approach over full-length trips, because
215 the average commuter travels approximately 25km one way in United States. As shown in Table 1, because
216 of the skipped episodes and steps that significantly reduce the training variance, our approach can effectively
217 learn the intelligent testing environment for all the experiments.

218 Furthermore, to demonstrate the advance of our approach in realistic urban scenarios, we extended our
219 simulation experiments at a real-world four-armed roundabout³² in Germany with a high traffic volume and
220 complex interactions. Compared with the NDE testing approach that requires about 8.91×10^6 number of
221 tests to reach the 30% relative half-width, our approach only requires 3.76×10^3 number of tests, which is
222 2.37×10^3 times faster. See Supplementary Video 2 and Supplementary Materials (Section 4b) for more
223 details.

224 **AV testing in test tracks**

225 Finally, we tested a Lincoln MKZ hybrid equipped with the open-source automated driving system,
226 Autoware²³ (Fig. 4a), driving continuously in the physical multi-lane 4-km highway test track at ACM (Fig.
227 4b) and the physical urban test track at Mcity (Fig. 4c), respectively. We developed an augmented reality
228 testing platform²⁴, which combines the physical test track and a simulation environment, SUMO²⁵. As shown
229 in Fig. 1d, by synchronizing the movements of the real AV and virtual BVs, the real AV in the physical test
230 track can interact with the virtual BVs as if it is in a real traffic environment, where the BVs are controlled
231 according to the intelligent testing environment. Figure 4d illustrates the real-time visualization of the testing
232 process. We trained the intelligent testing environment in the digital twins of the ACM highway section and
233 the Mcity urban section using the similar training settings as in the simulation studies (see Methods for
234 details). As shown in Fig. 4e-h, the crash rate estimations in both ACM and Mcity converge and reach the
235 30% relative half-width after about 156 tests at ACM and 117 tests at Mcity, which are on the order of 10^5
236 times faster than those (2.5×10^7 at ACM and 2.1×10^7 at Mcity) of the NDE testing approach. We also
237 evaluated the AV's safety performance for different crash types and severities (Fig. 4i, j).

238 **Discussion**

239 Our results present evidence of using D2RL techniques to validate AVs' safety performance regarding their
240 behavioral competency³³. D2RL can accelerate the testing process and can be used for both simulation testing
241 and test-track methods. It can significantly enhance existing testing approaches (falsification methods,
242 scenario-based methods, and NDE methods) to overcome their limitations in real-world applications. D2RL
243 also opens the door for leveraging AI techniques to validate machine intelligence of other safety-critical
244 autonomous systems, such as medical robots and aerospace systems.

245 Ideally, the testing environment should consider all operating conditions of AVs and their associated rare
246 events. For example, a six-layer model³⁴ has been developed to structure the parameters of scenarios,
247 including road geometry, road furniture and rules, temporal modifications and events, moving objects,
248 environmental conditions, and digital information. In this study, we mainly focus on two layers: moving
249 objects and road geometry, i.e., multiple surrounding vehicles undertaking maneuvers on roads of varying
250 geometry, which are critical for the testing environment. Our approach could be extended to include
251 parameters from other layers, such as weather conditions, by collecting large-scale naturalistic data and
252 utilizing domain knowledge of those fields.

253 We note that increasing attention has also been paid to formal methods to address the new challenges raised
254 by AI systems (see ref.^{35,36} and references therein). Formal methods provide mathematical framework for
255 rigorous system specification, design, and verification³⁷, which are critical for trustworthy AI. However, as
256 discussed in ref.³⁶, multiple major challenges need to be addressed to fully realize their full potential. D2RL
257 can potentially be integrated with formal methods. For example, reachability-based methods³⁸ could be
258 incorporated into the calculation of criticality measure to identify the critical states, particularly for generic
259 safety-critical autonomous systems. How to further integrate D2RL with formal methods deserves further
260 investigation.

261

262

263 **References**

1. Kalra, N., & Paddock, S. M. Driving to safety: How many miles of driving would it take to demonstrate autonomous vehicle reliability? *Transp. Res. A: Policy Pract.* **94**, 182-193 (2016).
2. LeCun, Y., Bengio, Y., & Hinton, G. Deep learning. *Nature* **521**, 436-444 (2015).
3. Insider, 10 million self-driving cars will be on the road by 2020, <https://www.businessinsider.com/report-10-million-self-driving-cars-will-be-on-the-road-by-2020-2015-5-6> (2016).
4. Nissan promises self-driving cars by 2020, <https://www.wired.com/2013/08/nissan-autonomous-drive/>, (2014).
5. Insider, Tesla's self-driving vehicles are not far off, <https://www.businessinsider.com/elon-musk-on-teslas-autonomous-cars-2015-9> (2015).
6. Society of Automotive Engineers, Taxonomy and Definitions for Terms Related to Driving Automation Systems for On-Road Motor Vehicles, https://www.sae.org/standards/content/j3016_202104/ (2021).
7. California Department of Motor Vehicles, 2021 Disengagement Reports, <https://www.dmv.ca.gov/portal/vehicle-industry-services/autonomous-vehicles/disengagement-reports/> (2022).
8. Paz, D., Lai, P. J., Chan, N., Jiang, Y., & Christensen, H. I. Autonomous vehicle benchmarking using unbiased metrics. *IEEE International Conference on Intelligent Robots and Systems (IROS)* 6223-6228 (IEEE, 2020).
9. Favarò, F., Eurich, S., & Nader, N. Autonomous vehicles' disengagements: Trends, triggers, and regulatory limitations. *Accid. Anal. Prev.* **110**, 136-148 (2018).
10. Riedmaier, S., Ponn, T., Ludwig, D., Schick, B. and Diermeyer, F., 2020. Survey on scenario-based safety assessment of automated vehicles. *IEEE access*, **8**, pp.87456-87477.
11. Nalic, D., Mihalj, T., Bäumler, M., Lehmann, M., Eichberger, A. and Bernsteiner, S. Scenario based testing of automated driving systems: A literature survey. In *FISITA web Congress* (2020).
12. Feng, S., Feng, Y., Yu, C., Zhang, Y., H.X. Liu. Testing scenario library generation for connected and automated vehicles, Part I: Methodology. *IEEE Trans. Intell. Transp. Syst.* **22**, 1573-1582 (2020).
13. Feng, S., Y. Feng, H. Sun, S. Bao, Y. Zhang, H.X. Liu. Testing scenario library generation for connected and automated vehicles, Part II: Case studies. *IEEE Trans. Intell. Transp. Syst.* **22**, 5635-5647 (2020).
14. Feng, S., Yan, X., Sun, H., Feng, Y., & Liu, H. X. Intelligent driving intelligence test for autonomous vehicles with naturalistic and adversarial environment. *Nat. Commun.* **12**, 1-14 (2021).
15. Sinha, A., O'Kelly, M., Tedrake, R., & Duchi, J. C. Neural bridge sampling for evaluating safety-critical autonomous systems. *Advances in Neural Information Processing Systems* **33**, 6402-6416 (2020).
16. Li, L., Wang, X., Wang, K., Lin, Y., Xin, J., Chen, L., Xu, L., Tian, B., Ai, Y., Wang, J. and Cao, D. Parallel testing of vehicle intelligence via virtual-real interaction. *Sci. Robot.* **4**, eaaw4106 (2019).
17. Zhao, D., Lam, H., Peng, H., Bao, S., LeBlanc, D.J., Nobukawa, K. and Pan, C.S. Accelerated evaluation of automated vehicles safety in lane-change scenarios based on importance sampling techniques. *IEEE Trans. Intell. Transp. Syst.* **18**(3), pp.595-607 (2016).
18. Donoho, D. L. High-dimensional data analysis: The curses and blessings of dimensionality. *AMS math challenges lecture* **1**, 32 (2000).
19. Hinton, G. E., & Salakhutdinov, R. R. Reducing the dimensionality of data with neural networks. *Science* **313**, 504-507 (2006).
20. Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., ... & Hassabis, D. Mastering the game of go without human knowledge. *Nature* **550**, 354-359 (2017).
21. Mirhoseini, A., Goldie, A., Yazgan, M., Jiang, J. W., Songhori, E., Wang, S., ... & Dean, J. A graph placement methodology for fast chip design. *Nature* **594**, 207-212 (2021).
22. Cummings, M. L. Rethinking the Maturity of Artificial Intelligence in Safety-Critical Settings. *AI Mag.*, **42**(1), 6-15 (2021).

309 23. Kato, S., Tokunaga, S., Maruyama, Y., Maeda, S., Hirabayashi, M., Kitsukawa, Y., Monrroy, A., Ando,
310 T., Fujii, Y. and Azumi, T. Autoware on board: Enabling autonomous vehicles with embedded systems.
311 In *2018 ACM/IEEE 9th International Conference on Cyber-Physical Systems (ICCPs)*, pp. 287-296
312 (2018).

313 24. Feng, S., Feng, Y., Yan, X., Shen, S., Xu, S., & Liu, H. X. Safety assessment of highly automated
314 driving systems in test tracks: a new framework. *Accid. Anal. Prev.* **144**, 105664 (2020).

315 25. Lopez, P. et al. Microscopic traffic simulation using sumo. *International Conference on Intelligent
316 Transportation Systems (ITSC)* 2575-2582 (IEEE, 2018).

317 26. Arun, A., Haque, M. M., Bhaskar, A., Washington, S., & Sayed, T. A systematic mapping review of
318 surrogate safety assessment using traffic conflict techniques. *Accid. Anal. Prev.* **153**, 106016 (2021).

319 27. Sutton, R. S., & Barto, A. G. *Reinforcement Learning: An Introduction*. MIT press (2018).

320 28. Koren, M., Alsaif, S., Lee, R., & Kochenderfer, M. J. Adaptive stress testing for autonomous vehicles.
321 *IEEE Intelligent Vehicles Symposium (IV)*, 1-7 (IEEE, 2018).

322 29. Sun, H., Feng, S., Yan, X., & Liu, H. X. Corner Case Generation and Analysis for Safety Assessment of
323 Autonomous Vehicles. *Transport. Res. Rec.*, DOI: 10.1177/03611981211018697 (2021).

324 30. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. Proximal policy optimization
325 algorithms. <https://arxiv.org/abs/1707.06347> (2017).

326 31. Owen, A. B. *Monte Carlo Theory, Methods and Examples*. <https://artowen.su.domains/mc/> (2013).

327 32. Krajewski, R., Moers, T., Bock, J., Vater, L. and Eckstein, L. September. The round dataset: A drone
328 dataset of road user trajectories at roundabouts in Germany. In *2020 IEEE 23rd International
329 Conference on Intelligent Transportation Systems (ITSC)* (pp. 1-6). (IEEE, 2020).

330 33. Nowakowski, C., Shladover, S.E., Chan, C.Y. and Tan, H.S. Development of California regulations to
331 govern testing and operation of automated driving systems. *Transport. Res. Rec.*, **2489**(1), pp.137-144
332 (2015).

333 34. Sauerbier, J., Bock, J., Weber, H., and Eckstein, L. Definition of scenarios for safety validation of
334 automated driving functions," *ATZ worldwide*, vol. 121, no. 1, pp. 4245 (2019).

335 35. Pek, C., Manzinger, S., Koschi, M., & Althoff, M. Using online verification to prevent autonomous
336 vehicles from causing accidents. *Nat. Mach. Intell.* **2**, 518-528 (2020).

337 36. Seshia, S.A., Sadigh, D. and Sastry, S.S. Toward verified artificial intelligence. *Commun. ACM*, **65**(7),
338 pp.46-55 (2022).

339 37. Wing, J.M. A specifier's introduction to formal methods. *IEEE Computer* **23**, **9**, 8–24 (1990).

340 38. Li, A., Sun, L., Zhan, W., Tomizuka, M. and Chen, M. Prediction-based reachability for collision
341 avoidance in autonomous driving. In *2021 IEEE International Conference on Robotics and Automation
342 (ICRA)*, pp. 7908-7914 (IEEE, 2021).

343

345

346 **Fig. 1 Validating safety-critical AI with the dense learning approach.** **a**, The “curse of rarity” hinders
347 the applicability of deep learning techniques for safety-critical systems, as the gradient estimation of neural
348 networks would suffer from the large variance due to the rareness of informative data. By training the neural
349 networks with the informative data only, our dense learning approach significantly reduces the gradient
350 estimation variance, enabling deep learning applications in safety-critical systems. **b**, The D2RL approach
351 edits the Markov process by removing the uncritical states and reconnecting the critical states, and then
352 trains the neural networks (NN) only for the edited Markov process. **c**, For any D2RL training episode, the
353 reward from the end state is backpropagated along the edited Markov Chain with critical states only. Three
354 examples are provided. For the left example, the episode is completely removed from training data as it does
355 not contain any critical state. For the middle and right examples, the uncritical states are skipped and critical
356 states are reconnected to densify the training data. The end state for the middle example is from a non-crash
357 episode, while the right example is from a crash episode. **d**, The augmented reality testing platform can
358 augment the real world with virtual background traffic, resulting in a safer, more controllable, and more
359 efficient testing environment for AVs. Our approach learns to decide when to control which background
360 vehicles to execute what adversarial maneuver with what probability.
361

362 **Fig. 2 Comparison of D2RL with DRL using the corner case generation examples.** **a**, The neural network
363 controls the closest eight vehicles’ maneuvers within 120 m, where each BV has 33 discrete actions at every
364 0.1 second: left lane change, 31 discrete longitudinal accelerations ($[-4, 2]$ with 0.2 m s^{-2} discrete resolution),
365 and right lane change. **b**, Proportions of the removed data by D2RL regarding the episodes (left) and steps
366 (right). **c**, Comparison of training rewards between DRL and D2RL, where the solid line represents the moving
367 averages of rewards and the light shadow represents the standard deviations. **d**, Comparison of crash rates
368 between the policies learned by DRL and D2RL. **e**, The AV (blue vehicle) made an evasive lane change to
369 avoid a cut-in vehicle but collided with an adjacent vehicle. **f**, The right front vehicle made a cut-in, the left
370 behind vehicle made a right lane change, while the right behind vehicle accelerated. These three vehicles
371 cooperatively encircled the AV and caused a crash. **g**, The right front vehicle made a cut-in to enforce the AV
372 for braking, which created the opportunity for the right behind vehicle to make a lane change after 2.8 seconds
373 (i.e., 28 uncritical steps), leading to a crash. Additional explanations are provided in Supplementary Video 1.

374
375 **Fig. 3. Performance evaluation of the D2RL-based intelligent testing environment.** **a**, Comparison of the
376 reward between the DRL and D2RL approaches, along with the estimation variance (dashed line) of the D2RL
377 approach that represents the testing efficiency. The solid line represents the moving average and the light
378 shadow represents the standard deviation. **b**, Comparison of crash rates of the on-policy and off-policy D2RL
379 approaches, during the training process. **c**, Comparison of estimated crash rates of the on-policy and off-policy
380 D2RL approaches, during the testing process. The light shadow represents the 90% confidence level. Crash
381 rate estimations (**d**) and relative half-width (**e**) of the AV-I model by the NDE and the D2RL-based intelligent
382 testing environment, respectively. The bottom x-axis denotes the number of tests for NDE, and the top x-axis
383 denotes the number of tests for the intelligent testing environment. **f**, Frequency of the required number of
384 tests for repeated testing experiments for the AV-I model. Unweighted crash rate (**g**) and weighted crash rate
385 (**h**) of each crash type in the D2RL-trained testing environment. Weighted distributions of the speed difference
386 at the crash moment (**i**), TTC (**j**), bumper-to-bumper distance (**k**), and post-encroachment time (**l**) of the near-
387 miss events.

388
389 **Fig. 4. Testing experiments for a real-world autonomous vehicle at physical test tracks.** **a**, Illustration of
390 the AV under test, equipped with Autoware. RTK, real-time kinematic positioning; IMU, inertial
391 measurement unit; DSRC, dedicated short-range communications; OBU, on-board unit. **b**, Illustration of the
392 ACM highway testing environment. The red line denotes the AV driving route. **c**, Illustration of the Mcity

393 urban testing environment including highways, roundabouts, intersections, etc. The explosion icons denote
 394 the locations of crash events happened during the tests. **d**, Illustration of the real-time visualization of the
 395 testing process: the leftmost figure illustrates the simulation view, where the virtual BVs (green vehicles) are
 396 generated and controlled by the intelligent testing environment to interact with the AV (red vehicle); the
 397 middle figure illustrates the real-world AV view visualized by Autoware, where the black vehicle is the AV
 398 under test and blue vehicles are augmented BVs; the rightmost figures illustrate the original image view (top)
 399 and augmented image view (bottom) from the AV's front camera. **e-h**, Crash rate estimation and the relative
 400 half-width of the real AV at the ACM test track (**e** and **f**) and Mcity test track (**g** and **h**) with the augmented
 401 reality testing platform. The black dashed line (**e** and **g**) represents the final estimation of the crash rate, the
 402 light dashed line (**f** and **h**) represents the 0.3 relative half-width threshold, and the light shadow represents the
 403 90% confidence level. **i**, Crash rates of different crash types of the AV at the Mcity test track. **j**, Distribution
 404 of the speed difference at the crash moment for crash severity analysis of the AV at the Mcity test track.
 405 Additional explanations regarding the field experiments are provided in Supplementary Videos 3-8.

406 **Table 1. Performance evaluation with different highway simulation environments.** The numbers of tests
 407 of the D2RL approach were the average values of multiple testing experiments, similar to Fig. 3f, and the
 408 numbers of tests for the NDE approach were obtained according to the Monte Carlo method¹.

		400 m		2 km		4 km		25 km
		2 Lanes	3 Lanes	2 Lanes	3 Lanes	2 Lanes	3 Lanes	3 Lanes
NDE	No. of tests	1.9 × 10 ⁸	1.0 × 10 ⁸	4.8 × 10 ⁷	2.5 × 10 ⁷	2.9 × 10 ⁷	9.4 × 10 ⁶	1.7 × 10 ⁶
D2RL	Episodes skipped	95.70%	91.73%	77.54%	79.85%	61.42%	58.92%	8.83%
	Steps skipped	99.78%	99.70%	99.82%	99.81%	99.79%	99.74%	99.76 %
	No. of tests	9.1 × 10 ⁴	4.4 × 10 ⁴	2.4 × 10 ⁴	1.7 × 10 ⁴	1.3 × 10 ⁴	4.5 × 10 ³	1.8 × 10 ³
	Acceleration ratio	2.1 × 10 ³	2.3 × 10 ³	2.0 × 10 ³	1.5 × 10 ³	2.2 × 10 ³	2.1 × 10 ³	9.4 × 10 ²

409

410

411

413 **Methods**414 **Description of the AV safety validation problem**

415 This section describes the problem formulation of AV safety performance evaluation. Denote the variables of
 416 the driving environment as $\mathbf{x} = [\mathbf{s}(0), \mathbf{u}(0), \mathbf{u}(1), \dots, \mathbf{u}(T)]$, where $\mathbf{s}(k)$ denotes the states (position and
 417 speed) of the AV and BVs at the k -th time step, $\mathbf{u}(k)$ denotes the maneuvers of BVs at the k -th time step,
 418 and T denotes the total time steps of this testing episode. With Markovian assumptions of BVs' maneuvers,
 419 the probability of each testing episode in the naturalistic driving environment (NDE) can be calculated as
 420 $P(\mathbf{x}) = P(\mathbf{s}(0)) \times \prod_{k=0}^T P(\mathbf{u}(k)|\mathbf{s}(k))$, and then the AV crash rate can be measured by the Monte Carlo
 421 method³¹ as

$$422 P(A) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})}[P(A|\mathbf{x})] \approx \frac{1}{n} \sum_{i=1}^n P(A|\mathbf{x}_i), \mathbf{x}_i \sim P(\mathbf{x}), \quad (2)$$

423 where A denotes the crash event, n denotes the total number of testing episodes, and $\mathbf{x}_i \sim P(\mathbf{x})$ indicates that
 424 the variables are sampled from the distribution $P(\mathbf{x})$. Here a crash is defined as a contact that the subject
 425 vehicle (e.g., AV) has with an object, either moving or fixed, at any speed resulting in fatality, injury, or
 426 property damage³⁹. As A is a rare event, obtaining a statistically reliable estimation requires a large number
 427 of tests (n), which leads to the severe inefficiency issue of the NDE testing approach, as pointed out in ref.¹.

428 To address this inefficiency issue, the key is to generate an intelligent driving environment, where
 429 BVs can be controlled purposely to test the AV unbiasedly and efficiently. In essence, testing an AV in the
 430 intelligent driving environment is to estimate $P(A)$ in Eq. (2) by the importance sampling method³¹ as

$$431 P(A) = \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})}[P(A|\mathbf{x}) \times W_q(\mathbf{x})] \approx \frac{1}{n} \sum_{i=1}^n P(A|\mathbf{x}_i) \times W_q(\mathbf{x}_i), \mathbf{x}_i \sim q(\mathbf{x}), \quad (3)$$

432 where $q(\mathbf{x})$ denotes the underlying distribution of BVs' maneuvers in the intelligent testing environment, and
 433 $W_q(\mathbf{x})$ is the likelihood of each testing episode as

$$434 W_q(\mathbf{x}) = \frac{P(\mathbf{x})}{q(\mathbf{x})} = \prod_{k=0}^T \left[\frac{P(\mathbf{u}(k)|\mathbf{s}(k))}{q(\mathbf{u}(k)|\mathbf{s}(k))} \right]. \quad (4)$$

435 According to the importance sampling theory³¹, the unbiasedness of the estimation in Eq. (3) can be
 436 guaranteed if $q(\mathbf{x}) > 0$ for any \mathbf{x} that $P(A|\mathbf{x})P(\mathbf{x}) > 0$. To optimize the estimation efficiency, the
 437 importance function $q(\mathbf{x})$ needs to minimize the estimation variance

$$438 \sigma_q^2 = \mathbb{E}_q(P^2(A|\mathbf{x}) \times W_q^2(\mathbf{x})) - P^2(A). \quad (5)$$

439 Therefore, the generation of the intelligent testing environment is formulated as a sequential Markov
 440 decision process (MDP) problem of BVs' maneuvers (i.e., determine $q(\mathbf{u}(k)|\mathbf{s}(k))$) to minimize the
 441 estimation variance σ_q^2 in Eq. (5). However, how to solve such a sequential MDP problem associated with
 442 rare events and high-dimensional variables remains a highly challenging problem, and most existing
 443 importance sampling-based methods suffer from the “curse of dimensionality”⁴⁰, where the estimation
 444 variance would increase exponentially with the dimensionality. In our previous study¹⁴, we discovered that
 445 the “curse of dimensionality” issue could be addressed theoretically by sparse adversarial control to the
 446 naturalistic distribution. However, only a model-based method with handcrafted heuristics was utilized for
 447 conducting the sparse adversarial control, which suffers from significant spatiotemporal limitations, and how
 448 to leverage AI techniques to train the BVs for truly learning the testing intelligence remains unsolved, which
 449 is the focus of this paper. More details of related work can be found in Supplementary Materials (Section 1).

450 **Formulation as a deep reinforcement learning problem**

451 This section describes how to generate the intelligent testing environment as a DRL problem. As mentioned
 452 above, the goal is to minimize the estimation variance in Eq. (5) by training a policy π modeled by a neural
 453 network θ that can control BVs' maneuvers with the underlying distribution $q_\pi(\mathbf{u}|\mathbf{s})$. To keep the notation

simple, we leave it implicit in all cases that π is a function of θ . An MDP usually consists of four key elements: state, action, state transition, and reward. In this study, states encode information (position and speed) about the AV and surrounding BVs, actions include 31 discrete longitudinal accelerations ($[-4, 2]$ with 0.2 m s^{-2} discrete resolution), left lane change, and right lane change, and state transitions define the probability distribution over next states that are also dependent on the AV maneuver. Here we assumed that a lane change maneuver of BVs would be initiated from its current position and completed in one second if a lane change action was decided. Our framework is also applicable to more realistic and complex action settings.

For the corner case generation case study, we studied a three-lane highway driving environment, where eight critical BVs (i.e., principal other vehicles or POVs) are controlled to interact with the AV for a certain distance (400m) and each BV has the 33 discrete actions at every 0.1 second. For the intelligent testing environment generation case study, to keep the runtime of the DRL small, we simplified the output of the neural network as the adversarial maneuver probability ($\varepsilon_\pi \in (0,1)$) of the most critical POV, while POV's other maneuvers are normalized by $1 - \varepsilon_\pi$ according to the naturalistic distribution and other BVs' maneuvers keep following the naturalistic distribution. The adversarial maneuver and POV are determined by the criticality measure. We note that the generalization of this work to multiple POVs is straightforward.

The reward function design is critical for the DRL problem⁴¹. As the goal of the intelligent testing environment is to minimize the estimation variance in Eq. (5), we derived the objective function of the DRL problem as

$$\min_q \sigma_q^2 = \max_\pi \left\{ -\mathbb{E}_{q_{\pi_b}} \left(\mathbb{I}_A(\mathbf{x}) \times W_{q_\pi}(\mathbf{x}) \times W_{q_{\pi_b}}(\mathbf{x}) \right) \right\}, \quad (6)$$

where \mathbb{I}_A is the indicator function of the crash event and π_b denotes the behavior policy of the DRL. During the training process, the training data is collected by the behavior policy, which is a Monte Carlo estimation of the expectation in Eq. (6), so we can obtain the reward function as

$$r(\mathbf{x}) = -\mathbb{I}_A(\mathbf{x}) \cdot W_{q_\pi}(\mathbf{x}) \cdot W_{q_{\pi_b}}(\mathbf{x}), \quad (7)$$

which is theoretically consistent with the objective function. As it is mainly based on the importance sampling theory, the reward function is also applicable to other rare event estimation problems with high-dimensional variables. To limit the scale of the error derivatives⁴², we rescaled and clipped the function, resulting in the reward function that belongs to $[-100, 100]$, where the scaling constants could be automatically determined during the learning process.

With the state, action, state transition, and reward function, the intelligent testing generation problem becomes a DRL problem. However, as the gradient estimation of neural networks would suffer from the large variance due to the rareness of informative data, applying learning-based techniques for safety-critical systems is highly challenging because of the “curse of rarity”. It is hard—or even empirically infeasible—to learn an effective policy if directly applying DRL approaches.

489 Dense Deep Reinforcement Learning (D2RL)

490 To address this challenge, we propose the D2RL approach in this paper. Specifically, according to the policy 491 gradient theorem²⁷, the policy gradient of the objective function for DRL approaches can be estimated as

$$\widehat{\nabla J(\theta)} = \hat{q}_\pi(S_t, A_t) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)}, \quad (8)$$

493 where θ denotes the parameters of the policy, $q_\pi(S_t, A_t)$ denotes the state-action value, S_t and A_t are samples 494 of the state and action under the policy, $\hat{q}_\pi(S_t, A_t)$ is an unbiased estimation of $q_\pi(S_t, A_t)$, i.e., 495 $\mathbb{E}_\pi[\hat{q}_\pi(S_t, A_t)] = q_\pi(S_t, A_t)$. Differently, for the D2RL approach, we propose to estimate the policy gradient 496 as

$$\widehat{\nabla_{dense} J(\theta)} = \hat{q}_\pi(S_t, A_t) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \mathbb{I}_{S_t \in \mathbb{S}_c}, \quad (9)$$

498 where \mathbb{S}_c denotes the set of critical states and $\mathbb{I}_{S_t \in \mathbb{S}_c}$ denotes the indicator function.

499 Here, a state is defined as an uncritical state if $v_\pi(s) = q_\pi(s, a), \forall a$, where $v_\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi(q_\pi(s, a))$
500 denotes the state value, so the set of critical states can be defined as $\mathbb{S}_c \stackrel{\text{def}}{=} \{s | v_\pi(s) \neq q_\pi(s, a), \exists a\}$. It
501 indicates that a state is defined as uncritical if any action (e.g., AV's maneuvers) from the current state will
502 not affect the expected value of the state (e.g., AV's crash probability within a specific time horizon from the
503 current state). We note that this definition is primarily for the theoretical analysis to be clean and is not strictly
504 required to run the algorithm in practice. For example, a state can be practically identified as uncritical if the
505 current action will not significantly affect the expected value of the state. For specific applications, the critical
506 states can be approximately identified based on domain-specific models or physics. For example, the criticality
507 measure^{12,13}, which is an outer approximation of the AV crash rate within a specific time horizon
508 (e.g., 1 second), is utilized in this study to demonstrate the approach for the AV testing problem. We note that
509 many other safety metrics²⁶ could also be applicable, such as the model predictive instantaneous safety
510 metric⁴³ developed by the National Highway Traffic Administration in the United States and the criticality
511 metric⁴⁴ developed by the PEGASUS project in Germany, as long as the identified set of states covers the
512 critical states. More theoretical analysis for a more general sense can be found in Supplementary Materials
513 (Section 2a).

514 Then, we have the following theorem, and the proof can be found in the Supplementary Materials:

515

516 **Theorem 1:**

517 The policy gradient estimator of D2RL has the following properties:

518 (1) $\mathbb{E}_\pi[\widehat{\nabla_{\text{dense}} J}(\theta)] = \mathbb{E}_\pi[\widehat{\nabla J}(\theta)]$;

519 (2) $\text{Var}_\pi[\widehat{\nabla_{\text{dense}} J}(\theta)] \leq \text{Var}_\pi[\widehat{\nabla J}(\theta)]$; and

520 (3) $\text{Var}_\pi[\widehat{\nabla_{\text{dense}} J}(\theta)] \leq \rho_\pi \text{Var}_\pi[\widehat{\nabla J}(\theta)]$, with the assumption

521
$$\mathbb{E}_\pi[\sigma_\pi^2(S_t, A_t) \cdot \mathbb{I}_{S_t \in \mathbb{S}_c}] = \mathbb{E}_\pi[\sigma_\pi^2(S_t, A_t)] \cdot \mathbb{E}_\pi[\mathbb{I}_{S_t \in \mathbb{S}_c}], \quad (10)$$

522 where $\rho_\pi \stackrel{\text{def}}{=} \mathbb{E}_\pi(\mathbb{I}_{S_t \in \mathbb{S}_c}) \in [0, 1]$ is the proportion of critical states in all states under the policy π (e.g., $1 - \rho_\pi$
523 denotes the proportion of steps skipped in Fig. 2b and Table 1), and $\sigma_\pi^2(S_t, A_t) = \left(\hat{q}_\pi(S_t, A_t) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)}\right)^2$.

524

525 Theorem 1 proves that the D2RL approach has an unbiased and efficient estimation of the policy
526 gradient compared with the DRL approach. To quantify the variance reduction of dense learning, we introduce
527 the assumption in Eq. (10), which assumes that $\sigma_\pi^2(S_t, A_t)$ is independent on the indicator function $\mathbb{I}_{S_t \in \mathbb{S}_c}$. As
528 both the policy and the state-action values are randomly initialized, the values of $\sigma_\pi^2(S_t, A_t)$ are quite similar
529 for all different states, so the assumption is valid at the early stage of the training process. Such significant
530 variance reduction will enable the D2RL approach to optimize the neural network, while the DRL approach
531 would be stuck at the beginning of the training process.

532 We then consider the influence of dense learning on estimating $\hat{q}_\pi(S_t, A_t)$ with bootstrapping, which
533 can guide the information propagation in the state-action space. For example, the fixed-length advantage
534 estimator (\hat{A}_t) is commonly used for the PPO algorithm³⁰ as

535
$$\hat{A}_t = \delta_t + (\gamma \lambda) \delta_{t+1} + \dots + (\gamma \lambda)^{L-t+1} \delta_{L-1}, \quad (11)$$

536 where $\delta_t = r_t + \gamma V(S_{t+1}) - V(S_t)$, $V(S_t)$ is the state-value function, and L denotes the fixed length. For
537 safety-critical applications, the immediate reward is usually zero (i.e., $r_t = 0$), and most state-value functions
538 are determined by initial random values without any valuable information because of the rarity of events.
539 Bootstrapping with such noisy state-value functions will not be effective in the learning process. By editing
540 the Markov chain, only the critical states will be considered. Then, the advantage estimator will be essentially
541 modified as

542 $\bar{A}_t = \delta_{z(t,0)} + (\gamma\lambda)\delta_{z(t,1)} + \dots + (\gamma\lambda)^{L-t+1}\delta_{z(t,L-1)},$ (12)

543 where $\delta_{z(t,j)} = r_{z(t,j)} + \gamma V(s_{z(t,j+1)}) - V(s_{z(t,j)})$, $z(t,0) = t$, and $z(t,j) = \min_i \{s_i \in \mathbb{S}_c | i > z(t, j-1)\}, j > 0$. In essence, it is a state-dependent temporal-difference (TD) learning, where only the values of critical states are utilized for bootstrapping. As the critical states have much higher probabilities leading to safety-critical events, the reward information can be propagated to these critical state values more easily. Utilizing the values of these critical states, the bootstrapping can guide the information from the safety-critical events to the state-action space more efficiently. This mechanism can help avoid the interference of the large number of noisy data and focus the policy on learning the sparse but valuable information. Because of the abovementioned variance reductions regarding the policy gradient estimation and bootstrapping, the D2RL approach significantly improves the learning effectiveness compared with the DRL approach, enabling the neural network to learn from the safety-critical events.

553 Densifying the information is a natural way to overcome the challenges caused by the rarity of events. 554 In the field of deep neural networks, connecting different layers of neural networks more densely has been 555 demonstrated to produce better training efficiency and efficacy, i.e., DenseNet⁴⁵. Instead of connecting layers 556 of neural networks, our approach densifies the information by connecting states more densely with safety- 557 critical states, besides the natural connections provided by the state transitions. As safety-critical states have 558 more connections with rare events, they contain more valuable information with less variance. By densifying 559 the connections between safety-critical states with other states, we can better propagate the valuable 560 information to the entire state space, which can significantly facilitate the learning process. This study 561 proposed and demonstrated one specific realization of the dense learning approach by approximately 562 identifying uncritical states and connecting the remaining states directly. This can be further improved by 563 more flexible and dense connections among safety-critical states and uncritical states. The connections can 564 even be added in the form of curriculum learning⁴⁶, which can guide the information propagation gradually. 565 The measures for identifying critical states can also be further improved by involving more advanced 566 modeling techniques.

567

568 **Off-policy learning mechanism**

569 We justify the off-policy learning mechanism in this section. The goal of the behavior policy π_b is to collect 570 training data for improving the target policy π that can maximize the objective function in Eq. (6). To achieve 571 this goal, it is critical to estimate the objective function accurately using the reward function in Eq. (7), which 572 determines the calculation of the policy gradient. However, only episodes with crashes have nonzero rewards, 573 so the objective function estimation suffers from a large variance, because of the rarity of crashes. Without an 574 accurate estimation of the objective function, the training could be misled. According to the importance 575 sampling theory, we have the following theorem, and the proof can be found in the Supplementary Materials:

577 **Theorem 2:**

578 The optimal behavior policy π_b^* that can minimize the estimation variance of the objective function has the 579 following property:

580 $q_{\pi_b^*}(x) \propto \frac{q_{\pi^*}^2(x)}{q_{\pi}(x)},$ (14)

581 where $q_{\pi^*}(x)$ denotes the optimal importance sampling function that is unchanged during the training process, 582 and the symbol \propto means “proportional to”.

583 Theorem 2 finds that the optimal behavior policy is nearly inversely proportional to the target policy, 584 particularly at the beginning of the training process when q_{π} is far from q_{π^*} . If using on-policy learning 585 mechanisms ($q_{\pi_b} = q_{\pi}$), the behavior policy would be far from optimality, which could mislead the training 586

587 process and eventually cause the underestimation issues. For example, if a target policy misses an action that
588 could lead to a likely crash, an on-policy learning mechanism will never find this missing crash. More
589 importantly, the on-policy mechanism could mislead the policy for purposely hiding the crashes that are
590 difficult to evaluate, leading to the severe underestimation issue of the safety performance evaluation.

591 We design an off-policy learning mechanism to address this issue, where a generic behavior policy is
592 designed and kept unchanged during the training process. Specifically, we determined a constant probability
593 of the adversarial maneuver of the POV (i.e., $\varepsilon_{\pi_b} = 0.01$) and conducted other maneuvers with the total
594 probability of 0.99 that were normalized according to the naturalistic distribution. This policy explores the
595 state-action space using the naturalistic distribution most of the time and exploits the information of the model-
596 based criticality measure that helps identify the POV and adversarial maneuver. We note that although the
597 optimal behavior policy needs to be adaptively determined based on the target policy, as indicated in Theorem
598 2, an off-policy learning mechanism can provide a sufficiently good foundation for effective learning in this
599 study. The behavior policy is also not sensitive to the constant of ε_{π_b} , and generally, a small value (e.g., 0.1,
600 0.05, 0.01, etc.) that balances the exploration and exploitation would be effective in this study. Further
601 improvement can be investigated in the future.

602 **Simulation settings**

603 **Naturalistic driving environment simulator.** To simulate naturalistic driving environment, we developed a
604 simulation platform based on an open-source traffic simulator SUMO. The scheme of the platform can be
605 found in Supplementary Materials. We utilized both the C++ and TRACI interfaces to refine the SUMO
606 simulator so that high-fidelity driving environments can be integrated. Specifically, we rewrote and
607 recompiled the C++ codes of SUMO to integrate the high-fidelity driving environments, including car-
608 following and lane-changing behavior models. Then, we utilized the TRACI interface to implement the
609 intelligent testing environment, where at selected moments, selected vehicles would execute specific
610 adversarial maneuvers with a learned probability, following the policy obtained by the D2RL approach. We
611 also synchronized the modified SUMO with the physical test tracks related to the information of BVs,
612 autonomous vehicles, traffic signals, high-definition maps, etc., through the TRACI interface. To provide a
613 training environment for intelligent testing environments, we constructed a multi-lane highway driving
614 environment and an urban driving environment, where all vehicles were controlled at 100 millisecond
615 intervals.

616 **Driving behavior models in the naturalistic driving environment simulator.** The default driving behavior
617 models of SUMO, which are simple and deterministic, cannot be utilized for safety testing and training of
618 AVs because they are designed to be crash-free models. To address this issue, in this study, we constructed
619 NDE models⁴⁷ to provide naturalistic behaviors of BVs according to the large-scale naturalistic driving
620 datasets (NDD) from the Safety Pilot Model Deployment program⁴⁸ and the Integrated Vehicle-Based Safety
621 System program⁴⁹ at the University of Michigan, Ann Arbor. At each step of simulation, the NDE models can
622 provide distributions of each BV's maneuvers, which are consistent with NDD. Then, by sampling maneuvers
623 from the distributions, a testing environment that can evaluate the real-world safety performance can be
624 generated. For the field testing at ACM and Mcity, although the intelligent testing environment can accelerate
625 the AV testing from about 10^7 loops of testing to only around 10^4 loops (see Table 1), this still represents a
626 significant level of effort for an academic research group. To demonstrate our approach in a more efficient
627 way, we simplified the NDE models to demonstrate our method more conveniently. Specifically, we modified
628 the intelligent driving model (IDM)⁵⁰ and the MOBIL (Minimizing Overall Braking Induced by Lane change)
629 model⁵¹ as stochastic models to construct the simplified NDE models. More details of the NDE models can
630 be found in the Supplementary Materials.

631 **D2RL architecture, implementation, and training.** The D2RL algorithm can be easily plugged into existing
632 DRL algorithms by defining a specific environment with the dense learning approach. Specifically, for

existing DRL algorithms, the environment receives the decision from the DRL agent, executes the decision, and then collects observations and rewards at each time step, while for the D2RL algorithm, the environment only collects the observations and rewards for the critical states, as illustrated in Supplementary Materials (Section 3e). In this way, we can quickly implement the D2RL algorithm utilizing existing DRL platforms. In this study, we utilized the PPO algorithm implemented at the RLLib 1.2.0 platform⁵², which was parallelly trained on 500 CPU cores and 3500 GB memory high-performance computation cluster at the University of Michigan, Ann Arbor. We designed a 3-layer fully connected neural network with 256 neurons in each layer and chose the 10^{-4} learning rate and 1.0 discount factor besides the default parameters. Each CPU collected 120 timesteps of training data for all experiment settings in each training iteration, so a total of 60,000 timesteps were collected in each training iteration. For the corner case generation, the neural network's output is the actions of the closest eight BVs, where each BV has the 33 discrete action space: left lane change, 31 discrete longitudinal accelerations ($[-4, 2]$ with 0.2 m s^{-2} discrete resolution), and right lane change. For the intelligent testing environment generation, the neural network's output is the adversarial maneuver probability (ε_π) of the POV, where the action space is $\varepsilon_\pi \in [0.001, 0.999]$. To further improve the data efficiency during the training process, we used the collected data with a resampling mechanism to train the neural network for multiple steps.

Field test settings

Augmented reality testing platform. We implemented the augmented reality testing platform at American Center for Mobility (ACM), one of the world's premier test tracks for AVs located in Ypsilanti, Michigan, and the Mcity test track, which is the world's first purpose-built test track for AV testing. In this study, we utilized the 4km highway loop featuring two-three lanes and both exit and entrance ramps to create various merging opportunities, as well as the Mcity urban driving environment, including various types of highways, roundabouts, urban streets, etc., as shown in Supplementary Materials (Section 3f). We constructed digital twins of the ACM and Mcity based on the naturalistic driving environment simulator and available high-definition maps. To synchronize the information between the simulation and physical test track, we utilized the dedicated short-range communications (DSRC) roadside units (RSUs) that were installed in the test tracks. These DSRC-based devices can communicate with AVs via 802.11p and SAE J2735 protocols through the immediate forward messaging (IMF) and forwarding functions. Specifically, we utilized the IMF function to broadcast proxy Basic Safety Message (BSMs) containing virtual BVs' identifier, latitude, longitude, altitude, etc., to the physical AV, and the forwarding function to forward incoming BSMs of the AV to the digital twins. After receiving the BSMs of the AV, we synchronized the AV states in the simulation world, where BVs were controlled by the intelligent testing environment. More details of the platform can be found in ref.²⁴. We implemented the system with an average 33ms communication delay, which is acceptable for AV testing and can be further improved with advanced wireless communication techniques.

Augmented image rendering. We use augmented reality techniques to render and blend virtual objects (e.g., vehicles) onto the camera view of the ego vehicle. Given a background 3D model with its 6DoF pose/location in the world coordinate, we perform a two-stage transformation to project the model to the onboard camera image: 1) from the world coordinate to the ego-vehicle coordinate, and 2) from the ego-vehicle coordinate to the onboard camera coordinate. In the first transformation, the ego vehicle pose and location are obtained from the real-time signal of the onboard high-precision RTK. In the second transformation, the projection is based on the pre-calibrated camera intrinsic and extrinsic. We also perform relighting on the rendered layer to harmonize the visual quality of the blending result. The augmented view is generated based on a linear blending with the rendered foreground layer, camera's background layer, and the rendered alpha matte. On top of the blending result, a weather-control layer is further added to simulate different weather conditions, e.g., rain, snow, and fog. We implemented the augmented rendering based on pyrender⁵³. An additional validation of the augmented image rendering can be found in Supplementary Materials (Section 4f).

680 **Autonomous vehicle under test.** As the AV under test, we used a retrofitted Lincoln MKZ from the Mcity
681 Test Facility at the University of Michigan, Ann Arbor. The vehicle was equipped with multiple sensors,
682 computing resources (2 Nexcom Lumina), and with drive-by-wire capabilities provided by Dataspeed Inc.
683 Specifically, the sensors include PointGrey camera, Velodyne 32 channel LiDAR, Delphi radars, OTXS
684 RT3003 RTK GPS, Xsens MTi GPS/IMU, etc. We implemented the vehicle with a ROS-based open-source
685 software, Autoware.AI²³, which provides full-stack software for the highly automated driving functions,
686 including localization, perception, planning, control, etc. We then integrated the AV with the AR testing
687 platform to evaluate the AV's safety performance. An illustration of the system framework can be found in
688 Supplementary Materials. Specifically, we modified the AV localization component to utilize the high-
689 definition map and high-accuracy RTK for obtaining the current pose and velocity. The surrounding vehicles'
690 BSMs were directly obtained from the simulation through wireless communications. To generate the AV's
691 future trajectory, we applied the OpenPlanner 1.13⁵⁴ as the decision module, an advanced planning algorithm
692 including global and local path planning. We applied the pure pursuit algorithm to convert the planned
693 trajectory into the velocity and yaw rate and then used a PID controller provided by Dataspeed Inc. to further
694 convert them into the vehicle by-wire control commands, i.e., steering angle, throttle, and brake percentages.

695 **Data availability**

696 The raw datasets that we used for modeling the naturalistic driving environment come from the Safety Pilot
697 Model Deployment (SPMD) program⁴⁸ and the Integrated Vehicle Based Safety System (IVBSS)⁴⁹ at the
698 University of Michigan, Ann Arbor. The ShapeNet Dataset that includes the 3D model assets for the image
699 augmented reality module can be found in <https://github.com/mmatl/pyrender>. The police crash reports used
700 in Supplementary Video 7 are available at <https://www.michigantrafficcrashfacts.org/>. The processed data for
701 constructing NDE models and the intelligent testing environment and the experiment results that support the
702 findings of this study are available at <https://github.com/michigan-traffic-lab/Dense-Deep-Reinforcement-Learning>.

703 **Code availability**

704 The simulation software SUMO, the automated driving system Autoware, and the RLLib platform with the
705 implemented PPO algorithm are publicly available, as described in the text and the relevant references^{23,25,52}.
706 The source codes for the naturalistic driving environment simulator, the driving behavior models in the
707 simulator, the D2RL-based intelligent testing environment, as well as the simulation setups are available at
708 <https://github.com/michigan-traffic-lab/Dense-Deep-Reinforcement-Learning>.

709 **References**

- 710 39. Automated Vehicle Safety Consortium. AVSC Best Practice for Metrics and Methods for Assessing
711 Safety Performance of Automated Driving Systems (ADS). SAE Industry Technologies Consortia
712 (2021).
- 713 40. Au, S. K., & Beck, J. L. Important sampling in high dimensions. *Struct. Saf.* **25**, 139-163 (2003).
- 714 41. Silver, D., Singh, S., Precup, D., & Sutton, R. S. Reward is enough. *Artif. Intell.* **299**, 1-13 (2021).
- 715 42. Mnih, V., Kavukcuoglu, K., Silver, D. et al. Human-level control through deep reinforcement learning.
716 *Nature* **518**, 529–533 (2015).
- 717 43. Weng, B., Rao, S. J., Deosthale, E., Schnelle, S., & Barickman, F. Model predictive instantaneous safety
718 metric for evaluation of automated driving systems. *IEEE Intelligent Vehicles Symposium (IV)* 1899-
719 1906 (IEEE, 2020).
- 720 44. Junietz, P., Bonakdar, F., Klamann, B., & Winner, H. Criticality metric for the safety validation of
721 automated driving using model predictive trajectory optimization. *International Conference on*
722 *Intelligent Transportation Systems (ITSC)* 60-65 (IEEE, 2018).

714 45. Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. Densely connected convolutional
 715 networks. *IEEE conference on computer vision and pattern recognition (CVPR)* 4700-4708 (2017).

716 46. Bengio, Y., Louradour, J., Collobert, R., & Weston, J. Curriculum learning. *International conference on*
 717 *machine learning (ICML)* 41-48 (2009).

718 47. Yan, X., Feng, S., Sun, H., & Liu, H. X. Distributionally Consistent Simulation of Naturalistic Driving
 719 Environment for Autonomous Vehicle Testing. <https://arxiv.org/abs/2101.02828> (2021).

720 48. Bezzina, D., Sayer, J. Safety pilot model deployment: Test conductor team report. (Report No. DOT HS
 721 812 171). Washington, DC: National Highway Traffic Safety Administration (2014).

722 49. Sayer, J. et al. Integrated vehicle-based safety systems field operational test: final program report (No.
 723 FHWA-JPO-11-150; UMTRI-2010-36). United States. Joint Program Office for Intelligent
 724 Transportation Systems (2011).

725 50. Treiber, M., Hennecke, A. & Helbing, D. Congested traffic states in empirical observations and
 726 microscopic simulations. *Phys. Rev. E* **62**, 1805 (2000).

727 51. Kesting, A., Treiber, M. & Helbing, D. General lane-changing model MOBIL for car-following models.
 728 *Transport. Res. Rec.* **1999**, 86–94 (2007).

729 52. Liang, E. et al. RLlib: Abstractions for Distributed Reinforcement Learning. *International conference on*
 730 *machine learning (ICML)*, 3053-3062 (2018).

731 53. Chang AX. et al. ShapeNet: An information-rich 3d model repository. <https://arxiv.org/abs/1512.03012>
 732 (2015).

733 54. Darweesh, H. et al. Open source integrated planner for autonomous navigation in highly dynamic
 734 environments. *J. Robot. Mechatron.* **29**, 668-684 (2017).

736 **Acknowledgements**

737 This research was partially funded by the U.S. Department of Transportation (USDOT) Region 5 University
 738 Transportation Center: Center for Connected and Automated Transportation (CCAT) of the University of
 739 Michigan (69A3551747105) and the National Science Foundation (CMMI #2223517). We thank the
 740 American Center for Mobility (ACM) for providing access to their test track. Any opinions, findings, and
 741 conclusions or recommendations expressed in this material are those of the authors and do not necessarily
 742 reflect the official policy or position of the U.S. government or the American Center for Mobility.

743 **Author contributions**

744 S. F. and H. L. conceived and led the research program, developed the AI against AI concepts, developed the
 745 dense learning approach, and wrote the paper. S. F. and H. S. developed the algorithms for the intelligent
 746 testing environment generation and designed the experiments. H.S. and H. Z. developed the simulation
 747 platform, implemented the algorithms, performed the simulation tests, and prepare the simulation results. X.
 748 Y., H. Z., and S. S. implemented the Autoware system in the autonomous vehicle, performed the field tests,
 749 and prepared the testing results. Z. Z. developed and performed the augmented image rendering. All authors
 750 provided feedback during the manuscript revision and results discussions. H. L. approved the submission and
 751 accepted responsibility for the overall integrity of the paper.

753 **Competing interests:** The authors have filed a provisional patent application 63/338,424.

755 **SUPPLEMENTARY MATERIALS**

757 Related Work.

759 Theoretical Analysis.

760 Supplementary Information for Experiments.

761 Supplementary Results.

762 Supplementary Videos: all eight supplementary videos files are available via Figshare.
763 https://figshare.com/articles/media/Dense_reinforcement_learning_for_safety_validation_of_autonomous_vehicles/21848259.
764