# ECSAS: Exploring Critical Scenarios from Action Sequence in Autonomous Driving

Shuting Kang[1,2], Heng Guo[1,2], Lijun Zhang[2], Guangzhen Liu[2], Yunzhi Xue[2], and Yanjun Wu[2]

[1] University of Chinese Academy of Sciences, Beijing, China
[2] Institute of Software Chinese Academy of Sciences, Beijing, China
kangshuting18@mails.ucas.edu.cn, guoheng21@mails.ucas.edu.cn

**Abstract.** Critical scenario generation requires the ability of sampling critical combinations from the infinite parameter space in the logic scenario. Existing solutions aim to explore the correlation of action parameters in the initial scenario rather than action sequences. How to model action sequences so that one can further consider the effects of different action parameters in the scenario is the bottleneck of the problem. In this paper, we attack the problem by proposing the ECSAS framework. Specifically, we first propose a description language, BTScenario, allowing us to model action sequences of the scenarios. We then use reinforcement learning to search for combinations of critical action parameters. To increase efficiency, we further propose several optimizations, including action masking and replay buffer. We have implemented ECSAS, and experimental results show that it is more efficient than native approaches such as random and combination testing in various nontrivial scenarios.

**Keywords:** Critical scenario generation · Reinforcement learning · Scenario description language · Autonomous driving
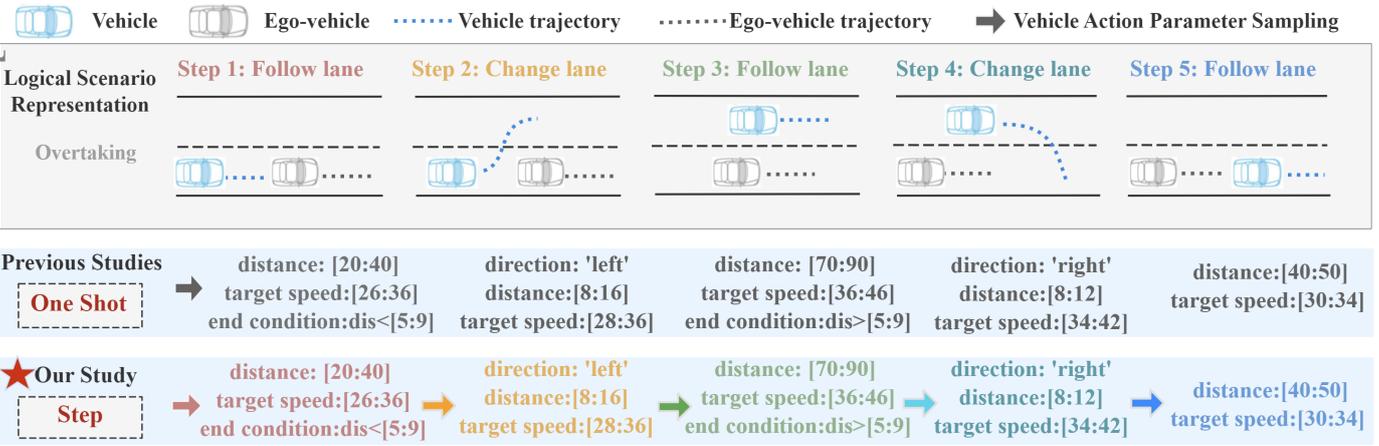
## 1 Introduction



**Fig. 1.** An illustration of our motivation. Compared with past studies sampling from all parameters of a logical scenario at once, our study generates action parameters by step from the perspective of action sequence. Subsequent experiments also verified this.

Despite the great success in *Autonomous Driving Systems* (ADS), safety engineering is still challenging for complex and dynamic environments, also known as scenarios. The space of possible test scenarios is virtually continuous, it is thus a challenging task how critical scenarios are generated. With critical scenarios, we refer to situations that cause potential safety risks [1]. *Critical Scenario Generation* (CSG) aims to find concrete critical scenarios within a given logical scenario, where the logical scenario is on a state-space level with parameter ranges, while a concrete scenario is a concretization of a logical scenario with concrete parameters [2]. Therefore, how to find critical parameters from infinite space becomes a core challenge for CSG.

To advocate research in this direction, native approaches, such as random testing (RT) and combination testing (CT), sample directly from parameter ranges of actions. These approaches are inefficient when critical scenarios are rare. To

attack this problem, recent works [3–5] focus on exploiting the correlation of actor's parameters by applying reinforcement learning (RL). However, they compute the probability distribution of all parameters in the scenario from the beginning.

Actually, a scenario is often modeled as a sequence of atomic actions, and the temporal relation between action parameters should also be considered. In this paper, we propose the *ECSAS framework*, abbreviating for **E**xploring **C**ritical **S**cenarios from **A**ction **s**equence. The core mechanism behind ECSAS is to model the action sequence of each scenario, so as to consider the influence relation between the action parameters. Considering the example in Fig. 1, we model an overtaking scenario with two atomic actions (follow lane and change lane) in five steps. Previous works sample all action parameters simultaneously in the overtaking scenario at the start. Assume that we have sampled values of the first step. In our framework, decision of the second step will make use of the values from previous steps, allowing ECSAS to search for critical scenarios efficiently.

For modeling scenarios, scenario description languages provide rich APIs, which are more flexible and convenient than the graphic model (Bayesian probability graph [3], causal graph and behavioral graph [5]), temporal logic [11], and Markov decision model [12–14]. Current description languages can be divided into two categories: representing scenarios from a local viewpoint [7–9], or using the positional changes from a global viewpoint [10]. However, these works fail to model scenarios based on the action interactions. Thus, in this paper, we propose a new language, BTScenario, which supports the representation of action sequences using temporal operators ("serial" and "parallel").

Inspired by [3], we use reinforcement learning (RL) to synthesize critical parameter valuations in logical scenarios. Action sequences leading to collisions are assigned with higher rewards, then the problem reduces to synthesizing optimal policies of the model. We use TD3 for this purpose, one of the RL algorithms to handle the continuous parameter spaces. We use an actor network to select actions and two critical networks to evaluate action sequences. During training, we use the action mask to fix the length of action space and the order of action, which means that the choice of next-step parameters does not affect the state of the previous step. We also optimize the storage mechanism, replay buffer, to improve search efficiency.

Summarizing, the main contribution of the paper is the ECSAS framework for generating critical scenarios. We propose a scenario description language, BTScenario, to model the action sequences in logical scenarios. ECSAS then uses RL combined with the Carla simulator to synthesize critical scenario valuations. We have implemented the ECSAS framework. Comprehensive and extensive experiments are conducted, showing that ECSAS outperforms native RT and CT approaches in various scenarios.

## 2   Related Work

### 2.1   Scenario representation

Scenario representation is the first step in critical scenario generation. Let us briefly summarise Scenic 1.0 [7]/2.0 [8], Paracosm [9], and GeoScenario [10] about scenario specifications in the sequel.

Scenic 1.0 focuses on the spatial layouts of map objects and the actors but does not allow specifying the temporal actions of actors. Scenic 1.0 was extended to Scenic 2.0 to remedy this deficiency. Scenic 2.0 specifies the actions locally in the sense that actions are specified separately for each actor. "ego = Car with behavior EgoBehavior()" is a simplified example to show this local viewpoint. Paracosm further abstracts the Python API interfaces of the simulator for users to support the action of actors in a way similar to Scenic 2.0.

Unlike Scenic and Paracosm, GeoScenario focuses on describing the scenario from a global viewpoint. Based on XML, GeoScenario specifies the scenarios by describing the trajectories of actors by a sequence of nodes specified on the map, the triggers for starting the actors, as well as the expected speeds when reaching the nodes.

However, both local view methods based on event-driven behavior patterns and global description methods based on trajectory points are not easy to split into the action sequences we want. We propose the BTScenario to model action sequences using temporal operators("serial" and "parallel") with encapsulated atomic action interfaces (e.g., "followLane" and "changeLane").

### 2.2   Critical Scenario Generation

Given a logical scenario, we aim to generate a set of concrete critical scenarios. A simple one-action scenario, such as going straight, can be sampled using mathematical formula constraints, but for complex multi-action scenarios, collision scenarios are most likely to occur during the first action phase, and user-defined mathematical calculations are required. Moreover, since the vehicle dynamic model in the simulation makes the vehicle no longer in an ideal driving state, it becomes impossible to rely solely on mathematical calculations.

A native scenario exploration approach, RT, randomly assigns each parameter's value in the logical scenario space. The CT [15] aims to generate a minimum set of test cases that satisfy N-wise coverage. The CT can be used to find
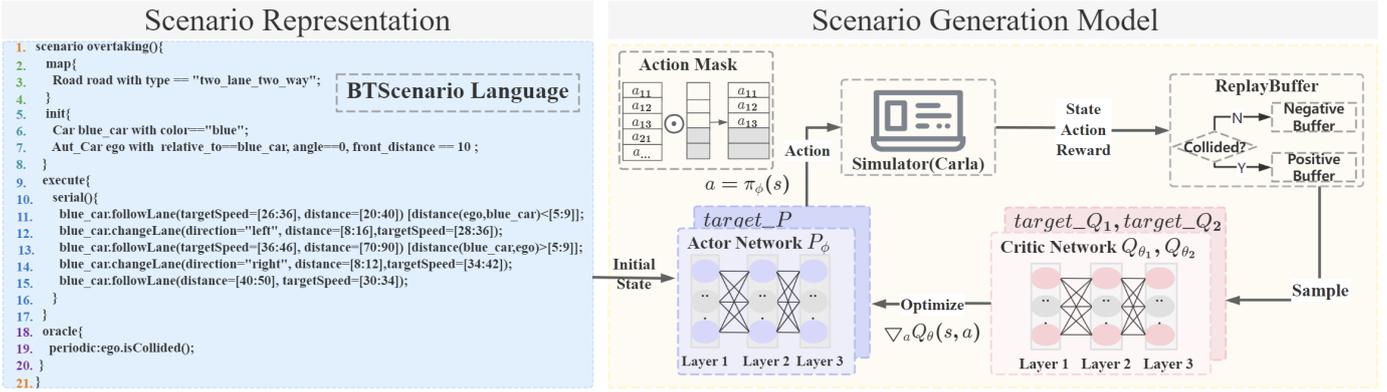
**Fig. 2.** An overview of ECSAS. The model contains two modules: Scenario Representation aims to depict a logical scenario by BTScenario and ego-vehicle system, and Scenario Generation Model aims to explore critical parameter combinations by RL and generates concrete scenarios by the simulator.

unknown combinations of parameters that may fail the ADS. However, these native approaches (RT and CT) can be inefficient because of the infinite parameter space and the low fraction of critical scenarios.

For the above challenge, the search-based methods have the potential to be more efficient [3,4,7,8] since the searching direction at each iteration is adjusted. However, gradient-based methods [7,8] only consider the role of a single parameter in CSG. Considering the correlation between parameters, most recent works sample critical parameters from a joint probability distribution of actor parameters by RL [3, 4] in a logical scenario. Another solution [5] introduces prior human knowledge by the causal graph and then represents the interaction by the behavioral graph. These methods mainly consider the correlation of all actor parameters from the scenario initialization. Instead, ECSAS considers the correlation of parameters in the action sequence, to the best of our knowledge, is the first framework for finding critical scenarios that adjusts decisions based on previous actions.

## 3    Methodology

A logic scenario is a sequence of actions' parameters $A = A_1, A_2, ..., A_T$, where $T$ is the number of steps in the scenario, and $A_t$ contains the parameters of $N_t$ actions $a_{t_1}, a_{t_2}, ..., a_{t_{N_t}}$ in the step $t$. The CSG aims to find *concrete critical scenarios*, which are colliding scenarios in this paper. Intuitively, concrete critical scenarios instantiate all parameters with concrete values so that the scenario leads to collisions.

Fig. 2 illustrates the overall structure of ECSAS. We first propose a scenario description language, *BTScenario*, allowing us to model the sequence of actor actions conveniently. Taking the BTScenario as an input, the scenario generation model exploits RL to explore the critical parameters. Sampled values are fed to the simulator, which simulates the physical environment, returns the reward, and updates the state. By stacking the above two processes multiple times, ECSAS learns the actor and critic networks that encode concrete critical parameters.

### 3.1    BTScenario: Scenario Representation Language

Researchers have proposed various scenario definitions [19,22,23], among which some common components exist: namely map, actor, interaction, and oracle. Inspired by these definitions, we propose BTScenario to model action sequences. It offers rich map elements to describe the driving area (3.1), the specification of the initial state (3.1), interaction (3.1), and the test oracle (3.1). We design an overtaking scenario (see Fig. 2) to illustrate some selected key map elements defined in BTScenario.

**The Map** Lines 2-4, of the code in Fig.2, show how map objects are specified in BTScenario: a road named "road" is declared as "two_lane_two_way". The following code declares a junction named "crossroad" with type "+".

```
Junction   crossroad with type == "+";
```

Note that a junction can also be of "T", "Y", or "unknown" type, where the "unknown" type denotes the junctions, not of a fixed shape, i.e. not of type "+", "T", or "Y".

Referring to the Apollo map structure [24], BTScenario supports eleven map objects(e.g., junction, road, lane, etc.), where several roads can be attached to a junction, and a road can have several lanes.

**Initial States Of Actors** In BTScenario, we consider various types of actors: $Aut\_Car$ represents the car under test (line 7 declares "ego" as the car under test), $Car$ represents the environment car (the "blue_car" at line 6). Moreover, we can specify additional constraints when declaring an actor. For instance, line 6 specifies the car with parameter color defined to be "blue".

An important attribute of actors is their locations. BTScenario supports the declaration of both absolute and relative locations. For instance, line 6 states that "blue_car" is at a randomly chosen location of "road" by default; line 7 represents "ego" is 10 meters ahead of "blue_car", and "angle==0" represents both cars are in the same direction.

**Temporal Actions of Actors** Before introducing the temporal actions of actors, we first show the specification of the atomic actions. For cars, the atomic actions include *followLane* and *changeLane*, etc. BTScenario supports the two types of parameter expressions of actions. For example, line 11 specifies that the driving distance of "blue_car" can be defined as a value within [20 : 40]. Moreover, actions can have pre-conditions and post-conditions that specify the enabling and terminating conditions of the action. For instance, line 11 specifies that "blue_car" follows the current lane until the distance between "ego" and "blue_car" is less than five to nine meters.

BTScenario specifies the interaction of actors by two temporal operators: the *serial* and *parallel* operators, meaning the sequential and parallel composition of actions, respectively. Fig. 2 gives an overtaking scenario using *serial*. Multiple cars can run in parallel, illustrated by the following code:

```
parallel(){
  car1.followLane(targetSpeed=[10,20]);
  car2.followLane(targetSpeed=[10,20]);
}
```

**The Oracle** To specify the test standards for critical scenarios, there are two kinds of operators in oracle module of BTScenario: "periodic" and "record". The operator "periodic" requires that the oracle needs to be verified unless it is violated or the simulation is completed. The operator "record" indicates that the program must record information about the actors or environment in the scenario. Line 19 specifies that the simulator monitors whether the "ego" collides.

### 3.2   Scenario Generation Model

Because the iterative state-parameter selection process of CSG has some flavour of the *policy iteration* of Markov Decision Process (MDP) [25], we use RL, which is commonly used in analyzing MDPs, to explore critical parameter combinations between different steps.

A scenario $S$ consists of $T$ steps according to the temporal relation of actions. At each step $t \in \{1, \ldots, T\}$, with a given state $s_t \in S$, we select the concrete values of $A_t$ to generate concrete scenarios. In RL, we regard the environment vehicle as an agent. Below, we introduce our reward design, action mask, replay buffer optimizations, and the optimization process.

**Replay Buffer Optimization** Due to the scarcity of critical scenarios, a long-tail problem remains in the training process of CSG network [3]. To address this problem, we optimize the replay buffer of RL. We separated the buffer into two parts to divide the samples according to whether or not the vehicle collided in the episode. The samples are stored in the temporary buffer $R_{tmp}$ during an episode. After the episode is finished, the samples in $R_{tmp}$ are moved to positive buffer $R_+$ if the vehicle collides otherwise to negative buffer $R_-$. In training the ECSAS, we sample $\eta \cdot N$ samples from $R_+$ and $(1 - \eta) \cdot N$ samples from $R_-$, where $0 \leq \eta \leq 1$.

**Optimization Process** Since TD3 [27], one of the RL algorithms, can be used to handle continuous action parameter spaces and is insensitive to hyper-parameter settings, we follow it to solve our optimization problem.

In RL, the objective is to find the optimal policy $\pi_\phi$, with parameters $\phi$, which maximizes the expected return $J(\phi)$. The optimal policy $\pi_\phi$ in ECSAS encodes the critical parameters of action $A$ causing the collision. For the continuous action space, $\pi_\phi$ can be updated by the gradient of the expected return $\nabla_\phi J(\phi)$. In TD3, the policy, known as the actor, can be updated by the deterministic policy gradient [28]:

**Table 1.** Hyper-parameters in experiments.

| Hyper-parameter | Description | Value |
|---|---|---|
| $d$ | delayed steps for updating network | 3 |
| $C$ | collision reward | 5 |
| $E$ | episode | 121 |
| $N$ | batch size | 50 |
| $\sigma$ | range of action noise | [0,1] |
| $\tau$ | delayed update parameters | 1e-2 |
| $\gamma$ | threshold in $Q$ | 0.9 |
| $\eta$ | sample coefficient between $R_+$ and $R_-$ | 0.5 |

$$\nabla_\phi J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s,a)|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s) \tag{1}$$

where $Q_\theta(s,a)$ is the action-value function, and the parameter values of actions $a \in A$ are generated from an actor network. The critic network is updated by temporal difference learning to maintain a fixed objective $y = r + \gamma min_{i=1,2} Q'_{\theta_i}(s', \tilde{a})$ over multiple updates. The gradient for updating critic networks parameters $\theta$ is:

$$\theta_i = argmin_{\theta_i} N^{-1} \sum (y - Q_{\theta_i}(s,a))^2 \tag{2}$$

## 4 Experiment

### 4.1 Experiment Settings

The map data is a virtual town on the plain in Apollo format. In our implementation, we connect Carla as the back-end simulator and Unreal Engine4 as the engine. We apply Lark[3] to compile the language and take PyTrees[4] to construct and execute the behavior tree. We use the feed-forward PID speed control algorithm e [29], and the Stanley direction control algorithm [30] for the ego car. Other details about the hyper-parameters in the Algorithm **??** are listed in Table 1.

**Metrics:** We follow the metric, collision rate (CR) [1, 3–5] to evaluate the performance after a static policy of the model. It is the percentage of scenarios for which the ECSAS yielded a collision while respecting behavioral constraints. In particular, a scenario is only considered successfully collide if all environment cars do not collide with other environment cars. CR is formulated in Eq. 3, where $U$ refers to the number of successful collision scenarios, $W$ refers to the total number of tested scenarios.

$$CR = \frac{U}{W} \cdot 100\% \tag{3}$$

Previous works [3–5] considered various scenarios consisting of only one step. In our implementation, we take the same actions from these scenarios to construct a more complex overtaking scenario with five steps. Besides, we design another turning scenario to show the generalization ability of our model in Fig. 3.

### 4.2 Comparison Experiment

Fig. 4 compares ECSAS with other methods.

- *RT*: We sample all values from the specifications of actions by BTScenario. This method may find critical scenarios but does not consider the influence of the combination of the action parameters.
- *CT*: The core is to cover the pairwise discrete combinations of multiple factor values with the fewest test cases. To reduce the searching expense, we discrete each action range, taking values every two numbers by Microsoft PICT tool.
- *Ablation Experiments(AM and RB)*: We evaluate the contributions of the action mask and replay buffer optimization for ECSAS.

---

[3] https://lark-parser.readthedocs.io/en/latest/grammar.html
[4] https://apollo.auto/index_cn.html

```
1.    execute{
2.      serial(){
3.        parallel(){
4.          blue_car.followLane(targetSpeed=[20:40], distance = [20:40]);
5.          red_car.followLane(targetSpeed=[20:40], distance = [20:40]);
6.        }
7.        parallel(){
8.          red_car.stop();
9.          blue_car.turn(targetSpeed=[30:40], distance = [30:50], direction="right");
10.       }
11.       blue_car.followLane(targetSpeed=[30:40], distance = [40:50]);
12.       blue_car.stop();
13.     }
14.  }
```
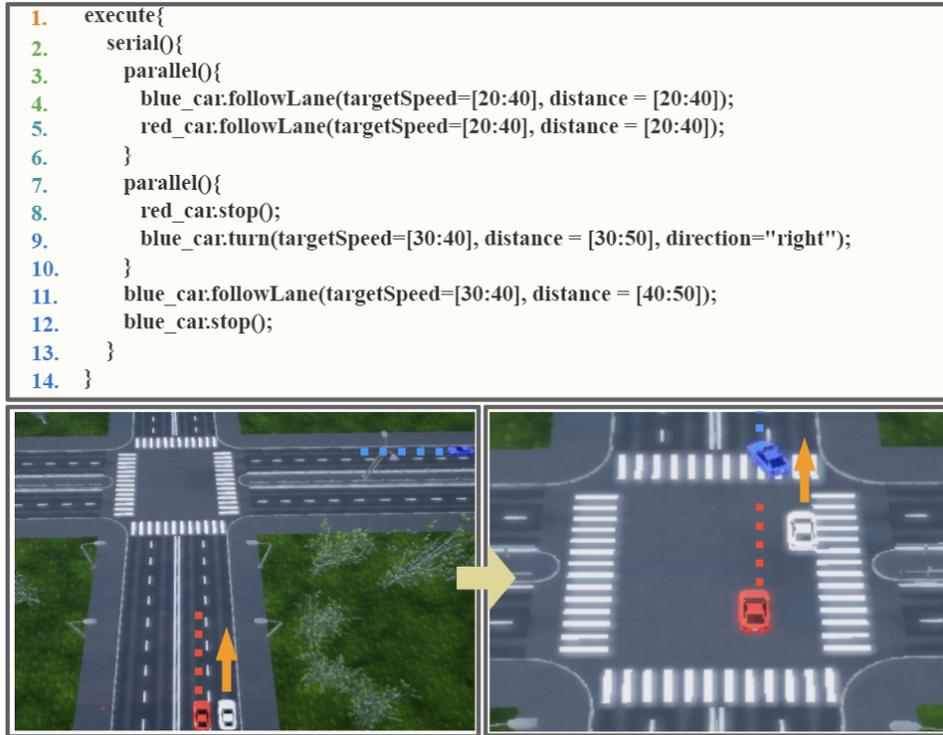


**Fig. 3.** Turning scenario. The ego and red car go straight through an "+" intersection. The blue car follows the lane and then turns right.
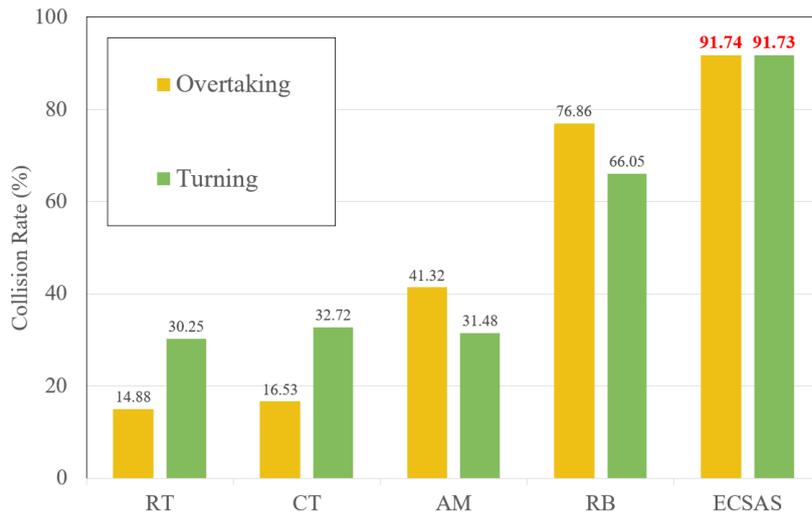


**Fig. 4.** The collision rate of RT, CT, the ablation of action mask(AM) and replay buffer(RB), and ECSAS.
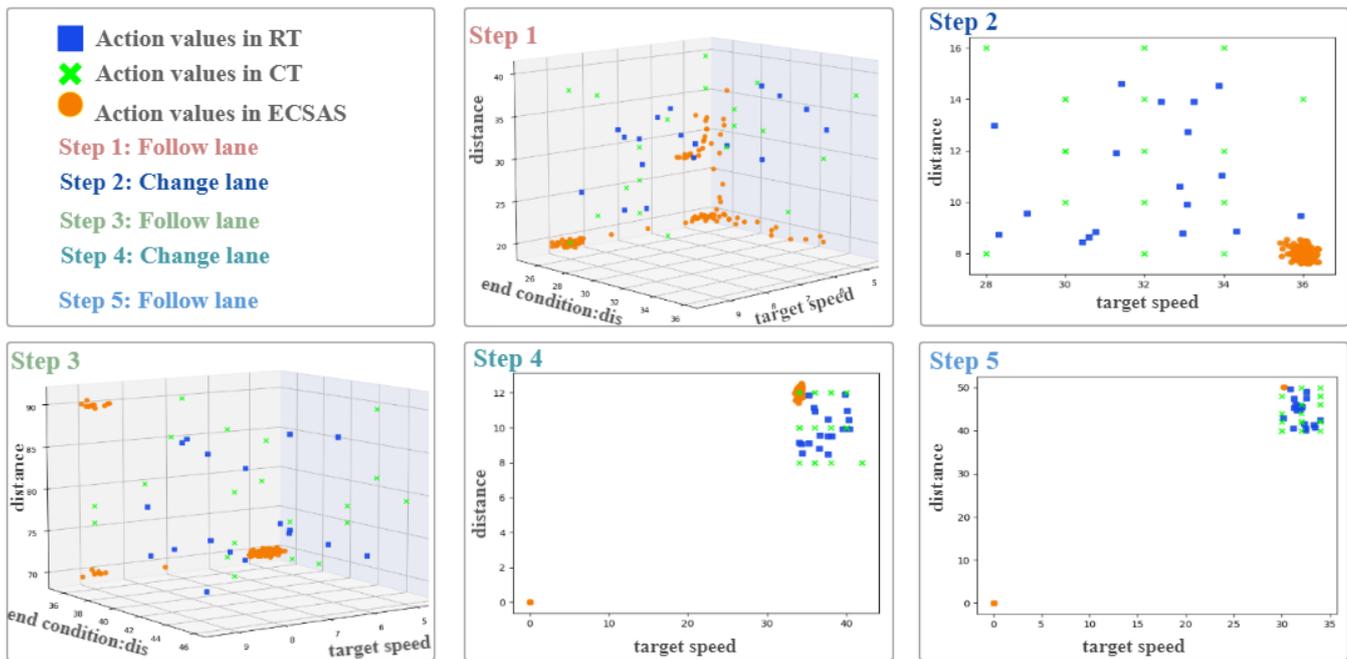
**Fig. 5.** Action value analysis in overtaking scenarios. As described in Fig.2 Scenarios Representation, in steps 1 and 3, the actions have three parameters, thus are shown in 3D. While in steps 2, 4 and 5, the actions have two parameters, thus are shown in 2D. Orange circle, green cross, and blue squares respectively represent the action values leading to the collision of ECSAS, CT, and RT in each step of simulation scenarios. For ECSAS, the action values leading to collisions in step 1 are rather discrete. However, in steps 2 and 3, affected by Action Mask in step 1, action values by ECSAS start to gather towards several regions, which means the action sequences are efficient. In each step, there are different action values leading to a collision, indicating the ECSAS can find different collision action distributions of different steps.

ECSAS consistently outperforms all the approaches on the collision rate. For example, ECSAS achieved 76.86% and 75.21% boost on the metric compared with random search methods (RT and CT) in the overtaking scenario. Since the critical scenario is at the tail of the long-tailed distribution, random search methods are inefficient. By considering the temporal relationship of parameters, ECSAS gains remarkable improvement.

We observe that the ablation results of the action mask and replay buffer optimization decrease compared with the full model in various scenarios, which indicates that the action mask and replay buffer optimization are beneficial for focusing on information about past actions and historical collision scenarios in CSG.

From Fig. 5, we analyze the action values of collisions in overtaking scenarios, which further verify the impact of action sequences on generating critical scenarios. The action values leading to collisions in step 1 are rather discrete. However, in steps 2 and 3, affected by actions in step 1, action values by ECSAS start to gather towards a region. Since RL converges the direction with the least loss, tending to a local optimal, it is likely to fall into mode collapse where the generation of action values is aggregated in one action space. In future work, we will address this issue. Moreover, since vehicles collide during the first three steps, some values of speed and distance in step 4 and step 5 become zero – their values are of no importance anymore.

## 5    Conclusion

In this paper, we propose ECSAS for critical scenario generation, which focuses on the correlation between parameters of the action sequence in a logical scenario. We model the action sequence by our description language, BTScenario. We adopt RL with the action mask and replay buffer optimization to sample critical parameters. Our model outperforms the RT and CT remarkably in various scenarios.

Some potential improvements can be explored in future work: (1) We will use different RL algorithms to analyze perception, decision, and control algorithms. (2) We will further extend BTScenario to support some random actions, random events, and distribution functions for range parameters.

## References

1. X. Zhang, J. Tao, K. Tan, M. Törngren, J. M. G. Sánchez, M. R. Ramli, X. Tao, M. Gyllenhammar, F. Wotawa, N. Mohan *et al.*, "Finding critical scenarios for automated driving systems: A systematic literature review," *arXiv preprint arXiv:2110.08664*, 2021.
2. T. Menzel, G. Bagschik, and M. Maurer, "Scenarios for development, test and validation of automated vehicles," in *2018 IEEE Intelligent Vehicles Symposium (IV)*.   IEEE, 2018, pp. 1821–1827.
3. W. Ding, B. Chen, M. Xu, and D. Zhao, "Learning to collide: An adaptive safety-critical scenarios generating method," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.   IEEE, 2020, pp. 2243–2250.
4. W. Ding, B. Chen, B. Li, K. J. Eun, and D. Zhao, "Multimodal safety-critical scenarios generation for decision-making algorithms evaluation," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1551–1558, 2021.
5. W. Ding, H. Lin, B. Li, and D. Zhao, "Causalaf: Causal autoregressive flow for goal-directed safety-critical scenes generation," *arXiv preprint arXiv:2110.13939*, 2021.
6. F. Martín, M. Morelli, H. Espinoza, F. J. Lera, and V. Matellán, "Optimized execution of pddl plans using behavior trees," *arXiv preprint arXiv:2101.01964*, 2021.
7. D. J. Fremont, T. Dreossi, S. Ghosh, X. Yue, A. L. Sangiovanni-Vincentelli, and S. A. Seshia, "Scenic: a language for scenario specification and scene generation," in *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2019, pp. 63–78.
8. D. J. Fremont, E. Kim, T. Dreossi, S. Ghosh, X. Yue, A. L. Sangiovanni-Vincentelli, and S. A. Seshia, "Scenic: A language for scenario specification and data generation," *arXiv preprint arXiv:2010.06580*, 2020.
9. R. Majumdar, A. Mathur, M. Pirron, L. Stegner, and D. Zufferey, "Paracosm: A test framework for autonomous driving simulations," in *Fundamental Approaches to Software Engineering*, E. Guerra and M. Stoelinga, Eds.   Springer International Publishing, 2021, pp. 172–195.
10. R. Queiroz, T. Berger, and K. Czarnecki, "Geoscenario: An open dsl for autonomous driving scenario representation," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, 2019.
11. C. Gladisch, T. Heinz, C. Heinzemann, J. Oehlerking, A. von Vietinghoff, and T. Pfitzer, "Experience paper: Search-based testing in automated driving control applications," in *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*.   IEEE, 2019, pp. 26–37.
12. A. Wachi, "Failure-scenario maker for rule-based agent using multi-agent adversarial reinforcement learning and its application to autonomous driving," *arXiv preprint arXiv:1903.10654*, 2019.
13. M. Koren and M. J. Kochenderfer, "Efficient autonomy validation in simulation with adaptive stress testing," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*.   IEEE, 2019, pp. 4178–4183.
14. M. Koren, S. Alsaif, R. Lee, and M. J. Kochenderfer, "Adaptive stress testing for autonomous vehicles," in *2018 IEEE Intelligent Vehicles Symposium (IV)*.   IEEE, 2018, pp. 1–7.

15. D. R. Kuhn, R. Bryce, F. Duan, L. S. Ghandehari, Y. Lei, and R. N. Kacker, "Combinatorial testing: Theory and practice," *Advances in computers*, vol. 99, pp. 1–66, 2015.
16. E.-H. Choi, T. Kitamura, C. Artho, and Y. Oiwa, "Design of prioritized n-wise testing," in *IFIP International Conference on Testing Software and Systems*. Springer, 2014, pp. 186–191.
17. M. Iovino, E. Scukins, J. Styrud, P. Ögren, and C. Smith, "A survey of behavior trees in robotics and ai," 2020.
18. K. Go and J. M. Carroll, "The blind men and the elephant: Views of scenario-based system design," *interactions*, vol. 11, no. 6, pp. 44–53, 2004.
19. S. Ulbrich, T. Menzel, A. Reschka, F. Schuldt, and M. Maurer, "Defining and substantiating the terms scene, situation, and scenario for automated driving," in *2015 IEEE 18th International Conference on Intelligent Transportation Systems*. IEEE, 2015, pp. 982–988.
20. S. Geyer, M. Baltzer, B. Franz, S. Hakuli, M. Kauer, M. Kienle, S. Meier, T. Weißgerber, K. Bengler, R. Bruder *et al.*, "Concept and development of a unified ontology for generating test and use-case catalogues for assisted and automated vehicle guidance," *IET Intelligent Transport Systems*, vol. 8, no. 3, pp. 183–189, 2013.
21. M. Jarke, X. T. Bui, and J. M. Carroll, "Scenario management: An interdisciplinary approach," *Requirements Engineering*, vol. 3, no. 3-4, pp. 155–173, 1998.
22. J. Bock, R. Krajewski, L. Eckstein, J. Klimke, J. Sauerbier, and A. Zlocki, "Data basis for scenario-based validation of had on highways," in *27th Aachen colloquium automobile and engine technology*, 2018.
23. M. Scholtes, L. Westhofen, L. R. Turner, K. Lotto, M. Schuldes, H. Weber, N. Wagener, C. Neurohr, M. Bollmann, F. Körtke, J. Hiller, M. Hoss, J. Bock, and L. Eckstein, "6-layer model for a structured description and categorization of urban traffic and environment," *IEEE Access*, vol. 9, pp. 59 131–59 147, 2021.
24. J. Xu, Q. Luo, K. Xu, X. Xiao, S. Yu, J. Hu, J. Miao, and J. Wang, "An automated learning-based procedure for large-scale vehicle dynamics modeling on baidu apollo platform," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 5049–5056.
25. A. Abakuks, "Reviewed work: Markov processes: characterization and convergence. by sn ethier, tg kurtz," *Biometrics*, vol. 43, no. 2, pp. 113–122, 1987.
26. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017.
27. S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International conference on machine learning*. PMLR, 2018, pp. 1587–1596.
28. D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *International conference on machine learning*. PMLR, 2014, pp. 387–395.
29. F. Gene, J. D. Powell, A. Emami-Naeini, R. Braun, and J. Flatley, "Feedback control of dynamic systems (7th global edition)," 2015.
30. G. M. Hoffmann, C. J. Tomlin, M. Montemerlo, and S. Thrun, "Autonomous automobile trajectory tracking for off-road driving: Controller design, experimental validation and racing," in *2007 American Control Conference*. IEEE, 2007, pp. 2296–2301.