



OCTOBER 1 - 5, 2023

IEEE/RSJ International Conference on Intelligent Robots and Systems

Adversarial Driving Behavior Generation Incorporating Human Risk Cognition for Autonomous Vehicle Evaluation

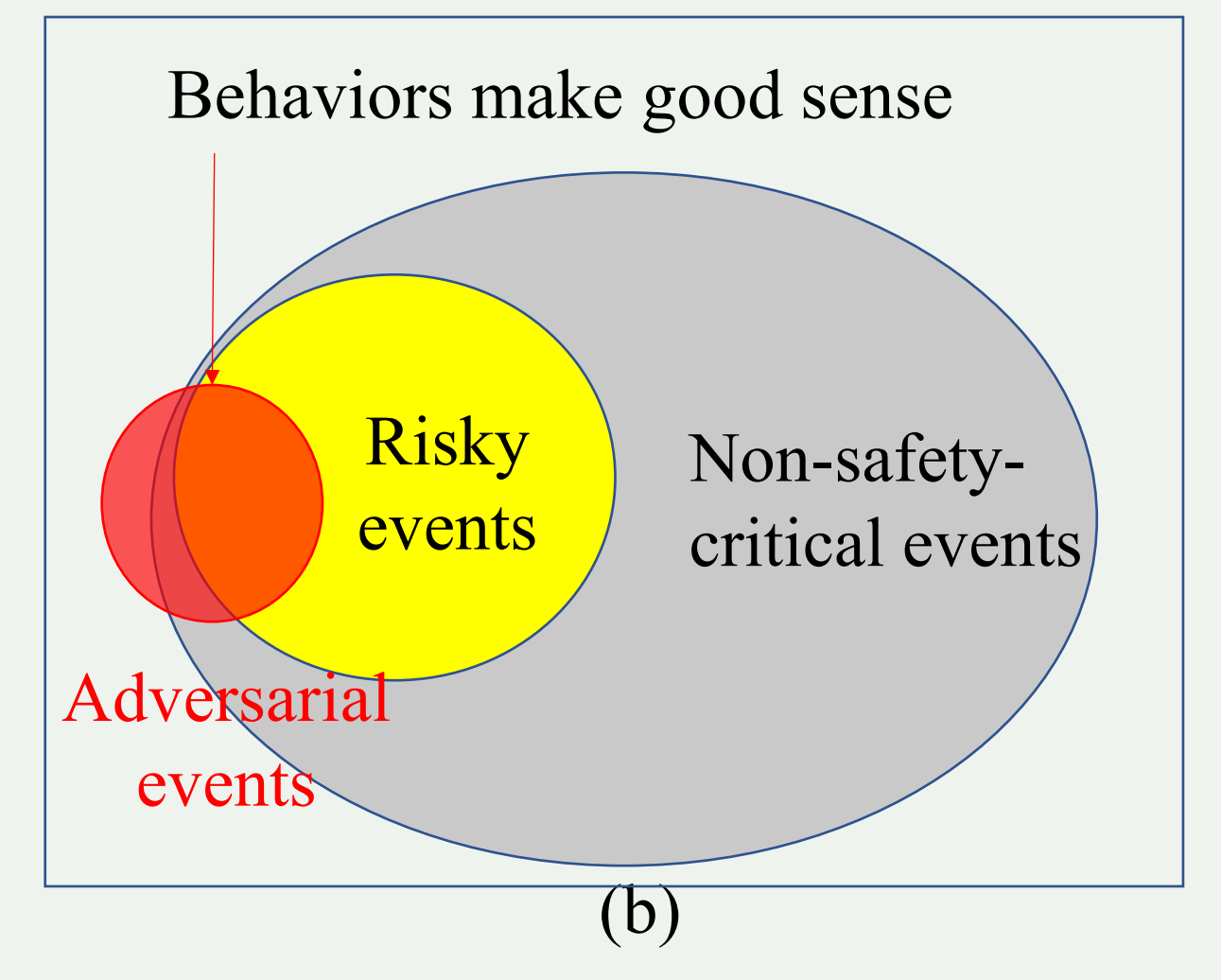
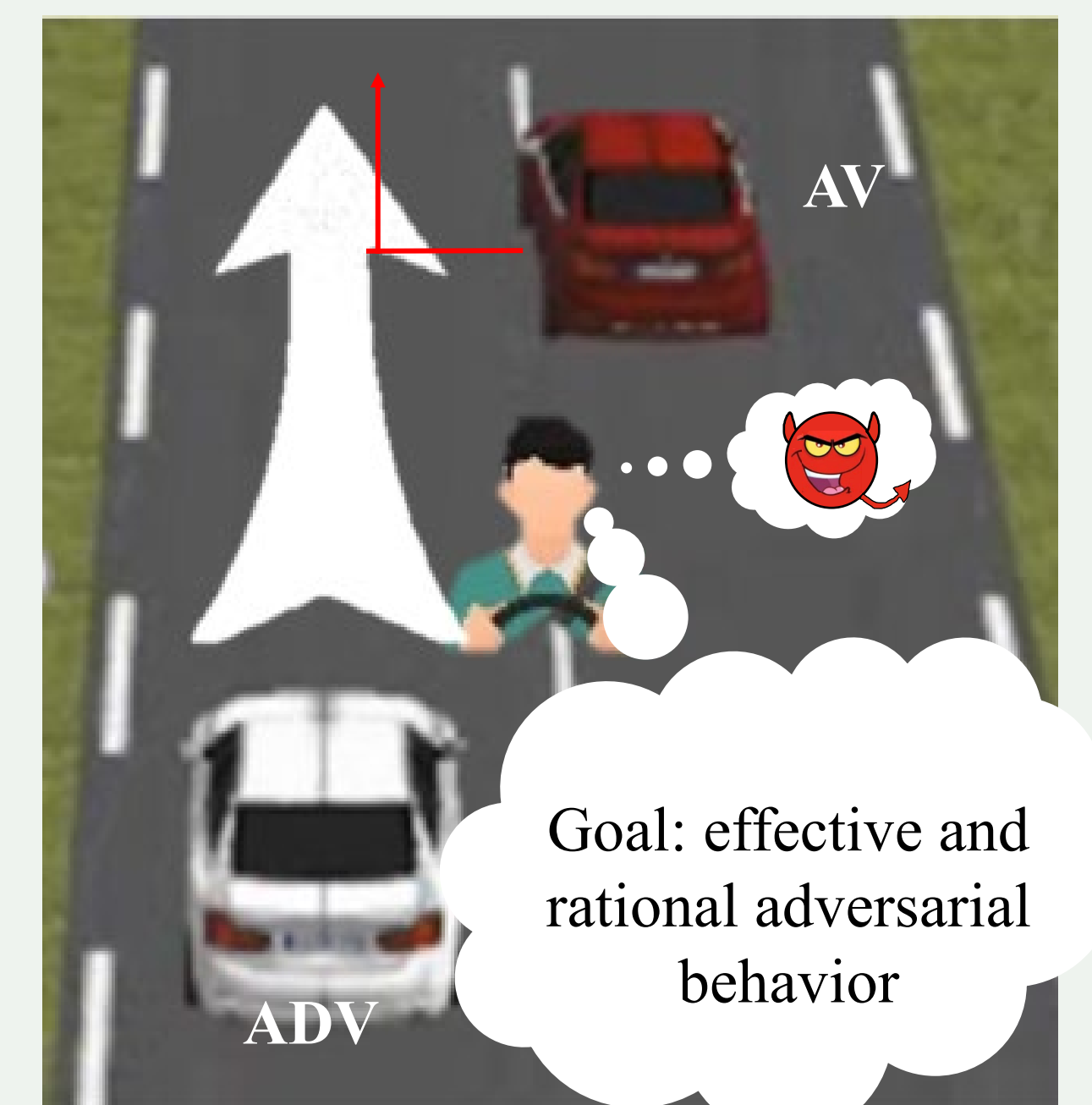
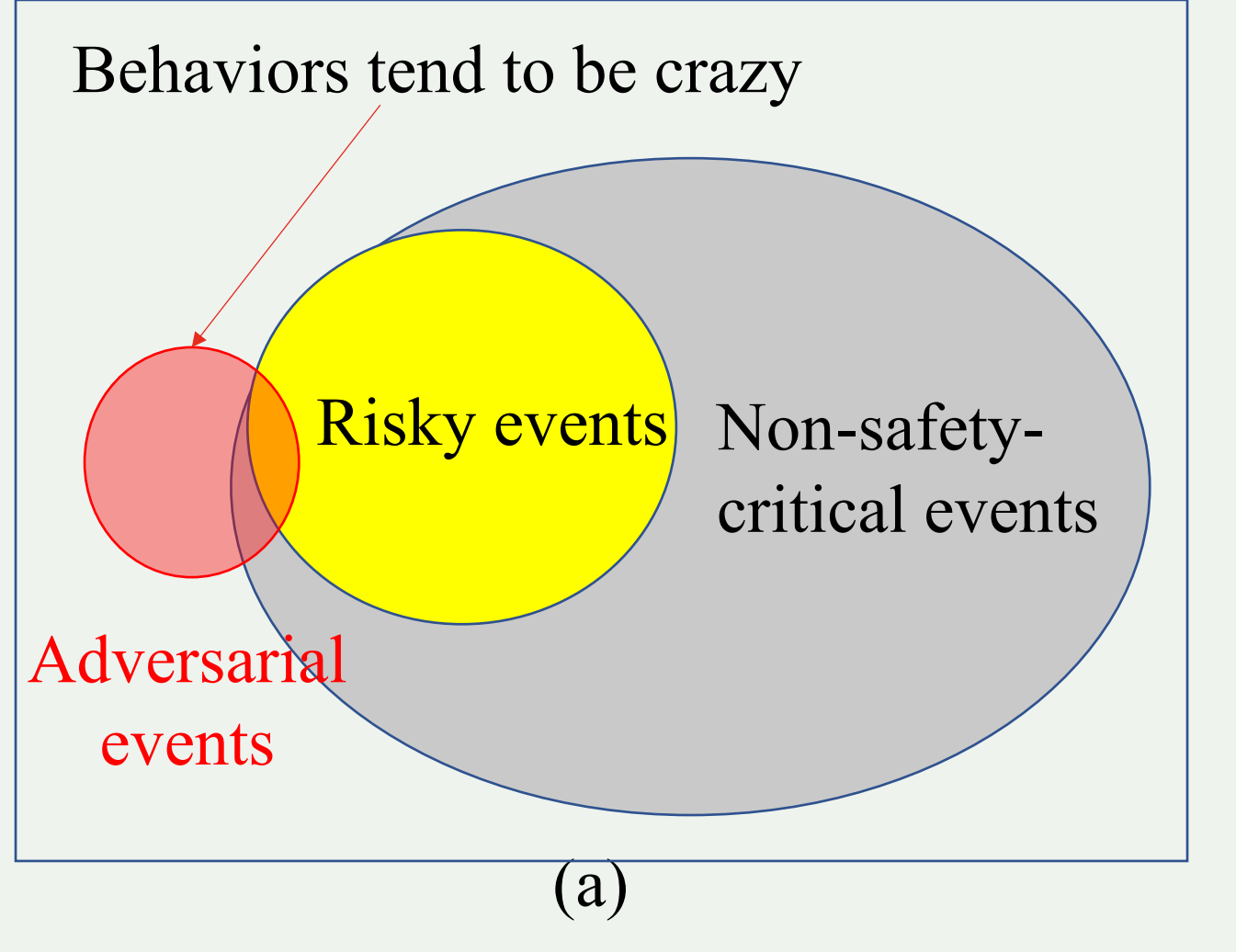
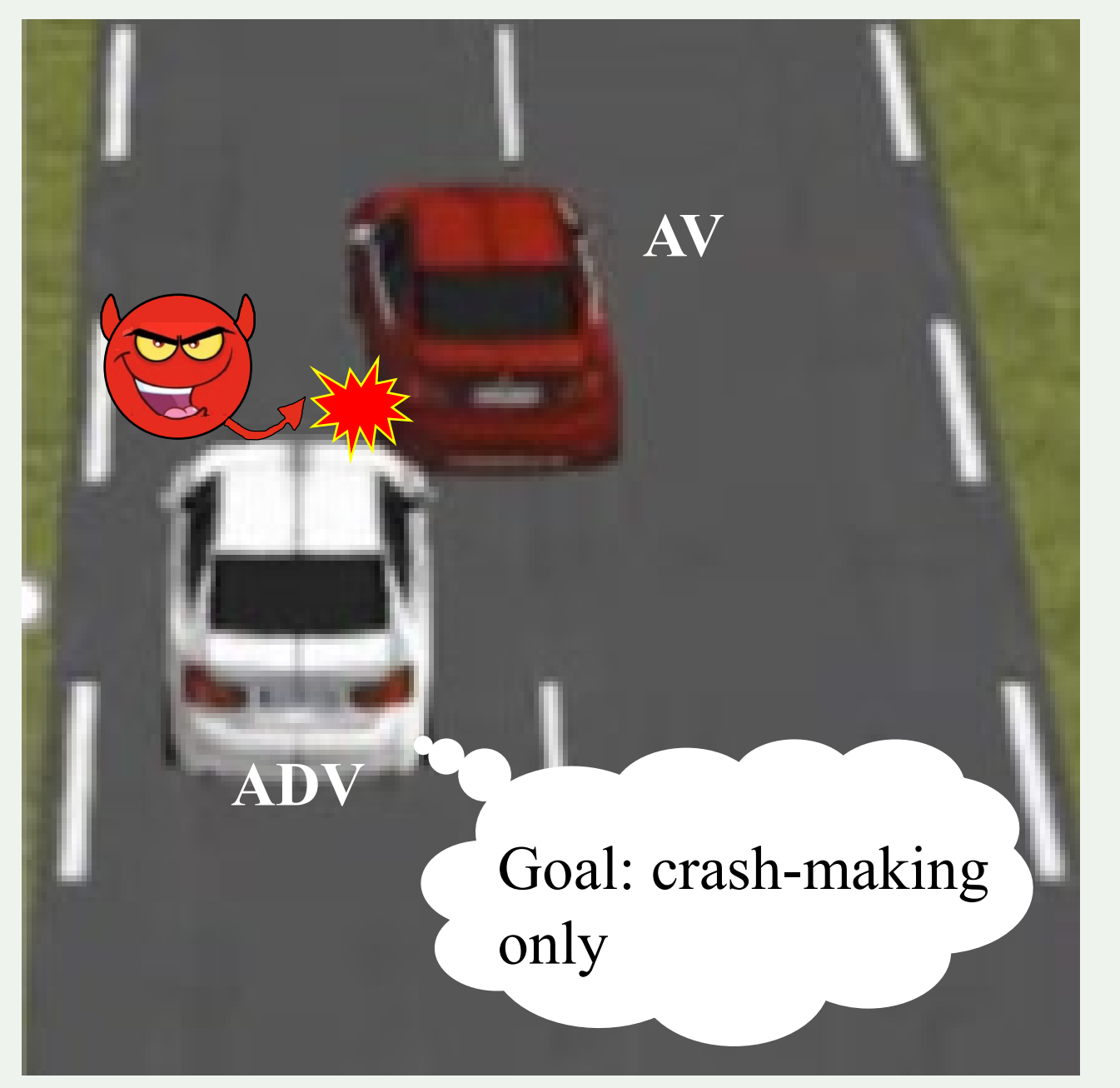
Zhen Liu, Hang Gao, Hao Ma, Shuo Cai, Yunfeng Hu, Ting Qu, Hong Chen, Xun Gong*

School of Artificial Intelligence
Jilin University

Abstract

- This paper focuses on the development of a novel framework for generating adversarial driving behavior of background vehicle interfering against the AV to expose effective and rational risky events.
- The adversarial behavior is learned by a reinforcement learning (RL) approach incorporated with the cumulative prospect theory (CPT) which allows representation of human risk cognition.
- The extended version of deep deterministic policy gradient (DDPG) technique is proposed for training the adversarial policy while ensuring training stability as the CPT action-value function is leveraged.
- A comparative case study regarding the cut-in scenario is conducted on a high fidelity Hardware-in-the-Loop (HiL) platform and the results demonstrate the adversarial effectiveness to infer the weakness of the tested AV.

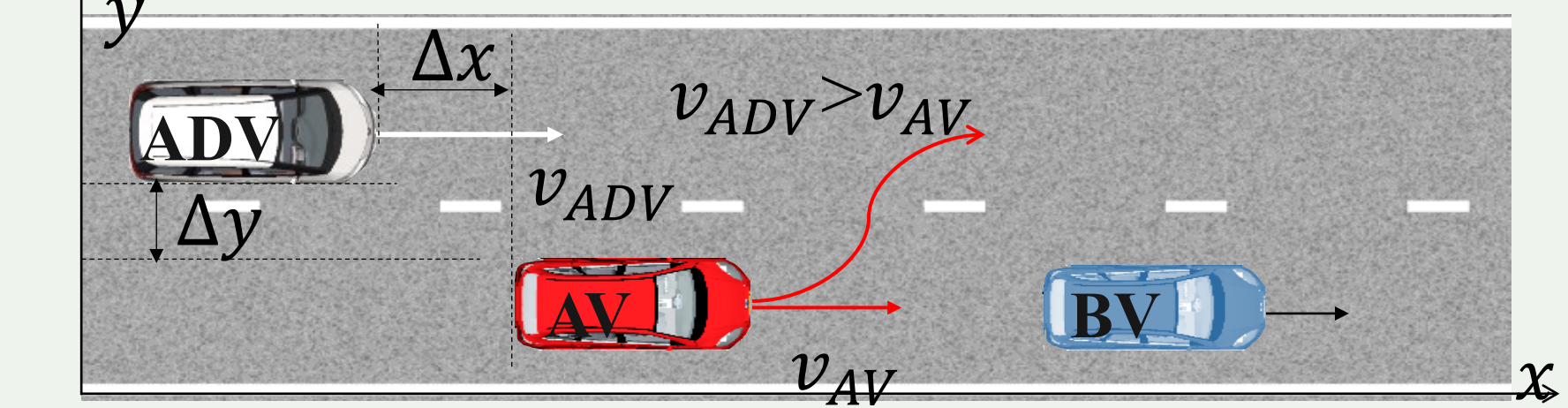
Motivation



- **Contribution**
 - Cumulative prospect theory (CPT) allows representation of human risk cognition.
 - CPT-RL can generate effective adversarial driving behavior by underestimating collision probability.
 - CPT-DDPG is proposed for solving CPT-RL while ensuring training stability as the CPT action-value function is leveraged.
 - A comparative case study demonstrate the adversarial effectiveness to infer the weakness of the tested AV.

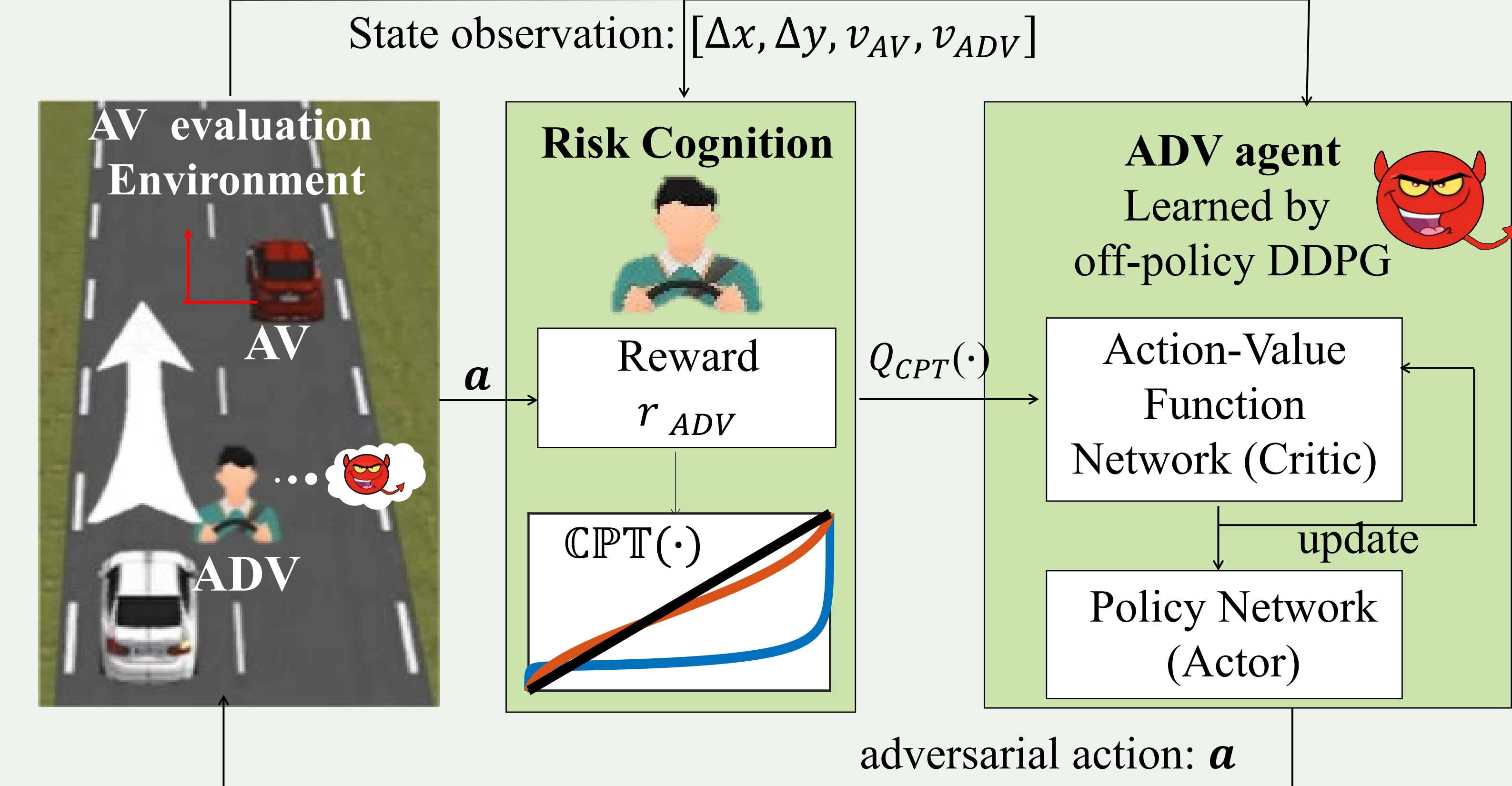
Method

1. Problem Statement



- **Lane-changing Scenario**
- It is assumed that when the AV decides to cut in, the initial speed of ADV is faster than that of AV.
- ADV's speed v_{ADV} , the AV's speed v_{AV} , the relative longitudinal distance Δx and relative lateral distance Δy between ADV and AV.

2. Overview

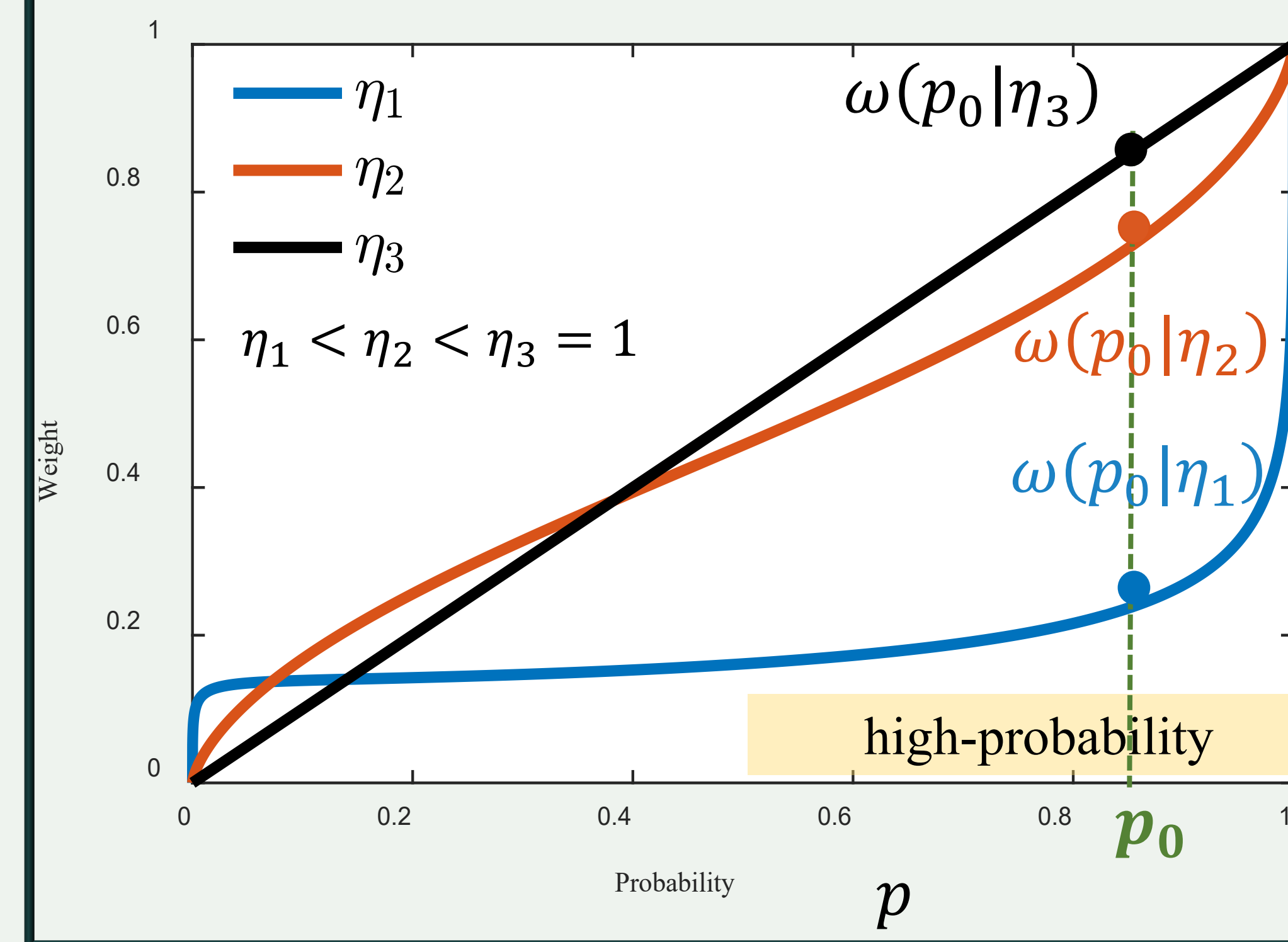


- We propose a framework based on CPT to generate adversarial behavior for testing autonomous vehicle.

$$r_{ADV} = \varphi_1 \frac{v_{ADV} - \underline{v}_{ADV}}{\bar{v}_{ADV} - \underline{v}_{ADV}} + \varphi_2 r_c \quad \text{CPT}(r_{ADV}, \eta)$$

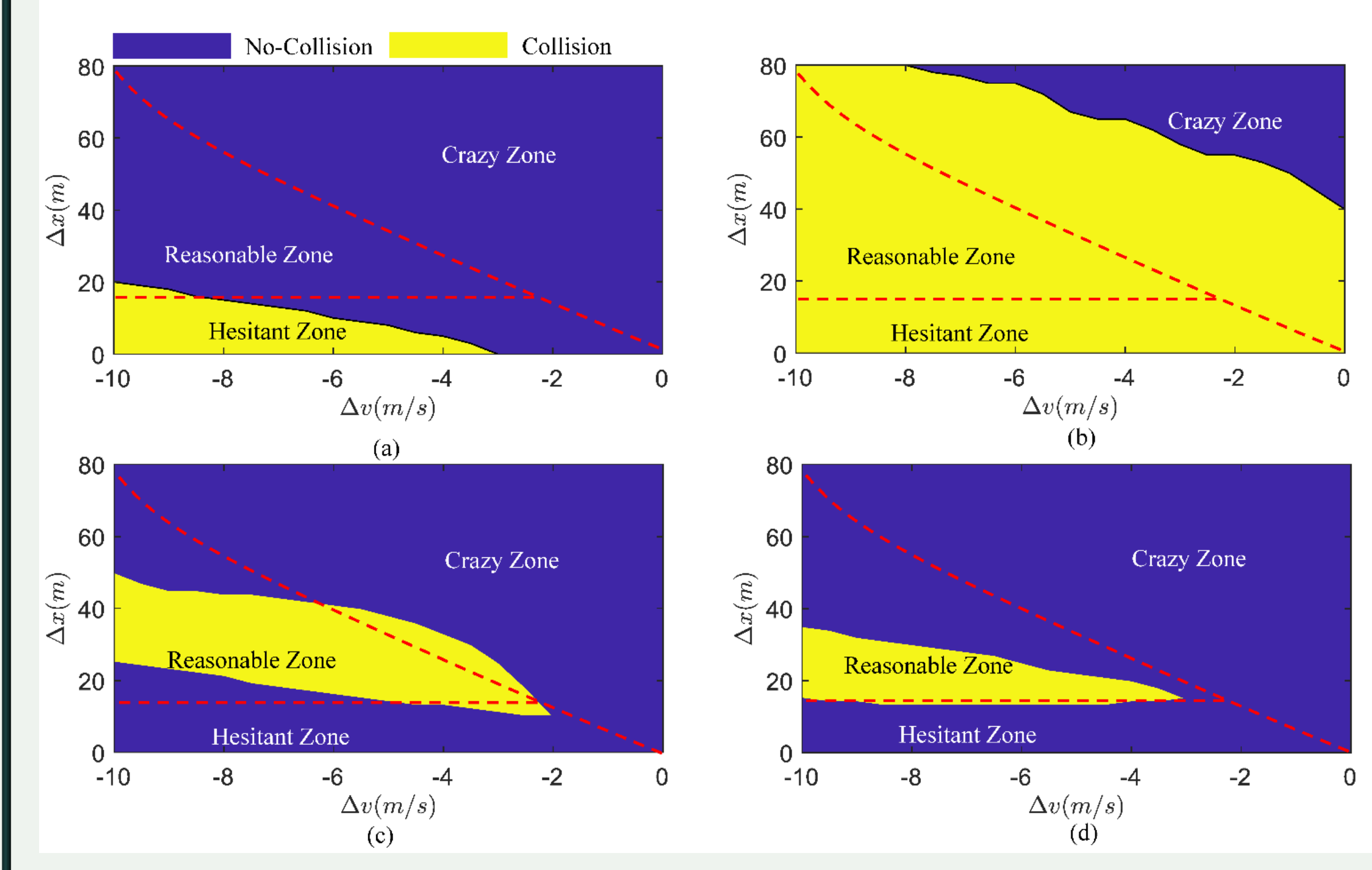
where φ_1 and φ_2 are weights, \underline{v}_{ADV} and \bar{v}_{ADV} are denoted as the lower bound and upper bound of ADV's longitudinal speed respectively. The collision penalty $r_c = -1$ if collision happened, otherwise $r_c = 1$.

3. Cumulative Prospect Theory(CPT)

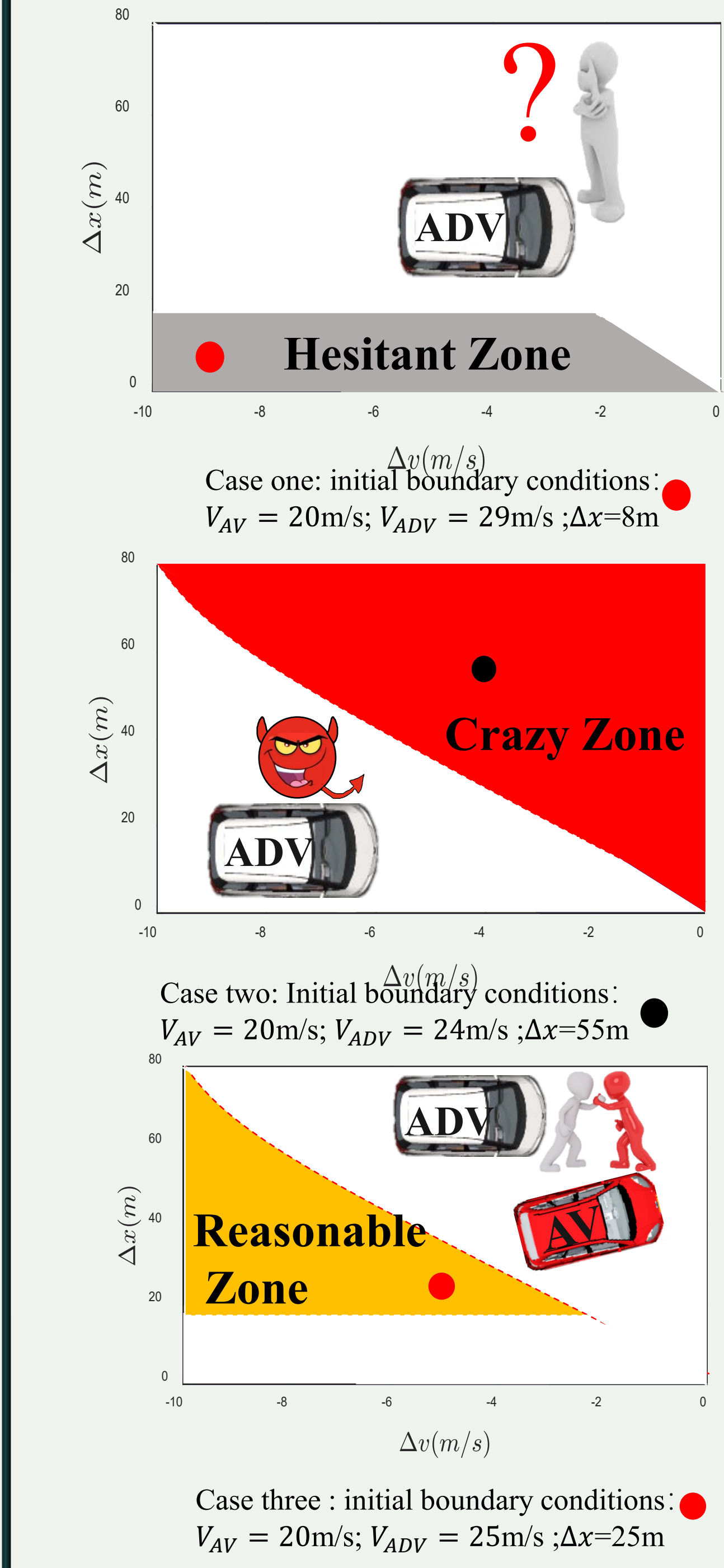


The impact of η on human risk cognitive probability $\omega(p)$. Small η value leads to underestimated occurrence probability by human: Given a high-probability event with a probability p_0 , and $\eta_1 < \eta_2 < \eta_3 = 1$, the human objective outcome probability $\omega(p)$ satisfies $\omega(p_0|\eta_1) < \omega(p_0|\eta_2) < \omega(p_0|\eta_3) = p_0$.

Results



- (a): **Conservative ADV**
 $r_{ADV} = \varphi_1 \frac{v_{ADV} - \underline{v}_{ADV}}{\bar{v}_{ADV} - \underline{v}_{ADV}} + \varphi_2 r_c$
- (b): **Hostile ADV**
 $r_{ADV}^{Hostile} = \begin{cases} 1, & \text{if collision} \\ -1, & \text{if no collision} \end{cases}$
- (c): **CPT-RL ($\eta=0.1$) ADV**
 $r_{ADV} = \varphi_3 \frac{v_{ADV} - \underline{v}_{ADV}}{\bar{v}_{ADV} - \underline{v}_{ADV}} + \varphi_4 r_c$
 $\text{CPT}(r_{ADV}, \eta = 0.1)$
- (d): **CPT-RL ($\eta=0.9$) ADV**
 $r_{ADV} = \varphi_3 \frac{v_{ADV} - \underline{v}_{ADV}}{\bar{v}_{ADV} - \underline{v}_{ADV}} + \varphi_4 r_c$
 $\text{CPT}(r_{ADV}, \eta = 0.9)$



- Cause of collision in Hesitant Zone: only ADV decelerates too cautiously and loses the opportunity to overtake AV.**
 - The conservative ADV has weak adversarial effectiveness.
- Cause of collision in Crazy Zone: ADV accelerate aggressively with obvious hostile intents to AV.**
 - The hostile ADV has bad adversarial effectiveness.
- Cause of collision in Reasonable Zone: ADV underestimating collision probability, driving mistakes, etc.**
 - The CPT-RL ADV has good adversarial effectiveness.

Conclusions

- We develop a CPT-RL approach for adversarial behavior generation towards the task of AV evaluation.
- The approach leverages human risk cognition to achieve rational exposure of safety-critical events. The stable training process is guaranteed via the proposed CPT-DDPG algorithm.
- Experimental results demonstrate that the CPT-RL is able to offer personalized adversarial patterns and facilitate effective AV evaluation.