

Cognitive-digital-twin-based Driving Assistance

Junyu Diao^{1#} and Renzhi Tang^{1#} and Yi Gu¹ and Sen Tian² and Zhihao Jiang^{*1,3}

Abstract—Advanced driver assistance systems (ADAS) have been developed to enhance driving safety by issuing timely warnings to drivers. However, current ADAS do not take into account the driver’s cognitive state when delivering warnings, which can result in false alarms and impact the driver’s trust in the system. To address this issue, we propose a Cognitive-digital-twin-based Driving Assistance System (CDAS) that issues warnings tailored to the driver’s perception of the driving environment and driving style. In this paper, we present a model of the driver’s decision-making process that explicitly captures their perception of the driving environment, their utility evaluation of predicted future environments, and their driving style in terms of minimum acceptable risk. The cognitive digital twin of the driver is then created and updated by minimizing the discrepancy between the predicted and actual behaviors of the driver. With the cognitive digital twin, the CDAS warns the driver when there is a significant discrepancy between the predicted driving strategy based on partial observation and that based on full observation. This approach can more accurately identify risks that the driver is not aware of and provide warnings only when necessary. We conducted human and simulated experiments in a virtual driving environment, and our results demonstrate that our proposed CDAS has a similar perception of risky behaviors compared to humans. Furthermore, the digital twin learning framework can identify the driving styles of human participants and accurately predict their driving strategies. Additionally, our proposed cognitive driving assistance system provides fewer false warnings and avoids more collisions compared to state-of-the-art ADAS algorithms. Our research shows that incorporating the driver’s cognitive state and driving style can enhance the effectiveness and safety of driving assistance systems.

Index Terms—Cognitive Modeling, Acceptability and Trust, Human Performance Augmentation

I. INTRODUCTION

WITH vision as the primary source of information [1], human drivers cannot anticipate and respond to unexpected actions of unseen agents, which is a primary cause of traffic accidents [2]. Modern vehicles are equipped with Advanced Driver Assistance System (ADAS) that can alert the drivers or even takes control of the vehicle in hazardous driving situations. By using advanced sensors like Lidar, the

Manuscript received: March, 15, 2023; Revised June, 2, 2023; Accepted June, 26, 2023.

This paper was recommended for publication by Editor Aniket Bera upon evaluation of the Associate Editor and Reviewers’ comments. This work was supported by NSFC Young Scholar Program 62103279 and Shanghai Pudong Science Foundation PKX2021-R09

The authors contributed equally to this research

¹Junyu Diao, Renzhi Tang, Yi Gu and Zhihao Jiang are with School of Information Science and Technology, ShanghaiTech University, China diaojy@shanghaitech.edu.cn

²Sen Tian is with Southwestern University of Finance and Economics tiansen@swufe.edu.cn

³ Shanghai Engineering Research Center of Intelligent Vision and Imaging, China

Digital Object Identifier (DOI): see top of this page.

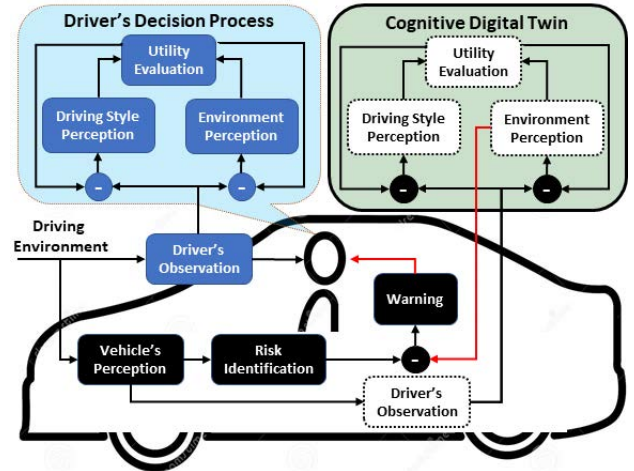


Fig. 1. The proposed cognitive digital twin captures the driver’s decision process. Driving assistance system can provide warnings to the driver for unidentified risks that are consistent with the driver’s driving style.

ADAS can better understand the driving environment, identify potential dangers that the driver may not have noticed, and significantly enhance driving safety in various circumstances [3], [4].

However, the driving assistance capabilities of ADAS are reliant on the system’s perception of the driving environment, which can differ from the driver’s perception. This can result in the assistance being seen as "surprising" or "unneeded", which can impact driver acceptance of the technology [5]. In the case of experienced drivers, warning them about risks they are already aware of may cause them to become less attentive and disregard future warnings. Conversely, novice drivers may experience an excessive number of warnings in complex driving environments, leading to increased cognitive load and impaired driving performance [6].

Studies in game theory showed that information sharing between players leads to improved social outcome [7]. While drivers cannot actively share information about their decision-making process with ADAS, the system can achieve imperfect information sharing by accurately estimating the driver’s decision-making process. To provide appropriate warnings that align with the driver’s driving style, the system must explicitly identify the following components of driving decision-making: **Environment Perception:** Drivers rely on prior knowledge and experience to predict the behaviors of unobserved and observed agents, and maintain the belief of plausible driving environments [8]. In order to identify unidentified risks by the driver, it is essential to capture the driver’s perception of the driving environment and quantify its uncertainties.

Utility Evaluation: In sociology and economics, it is widely

assumed that human behaviors are motivated by weighing possible rewards and penalties [9]. To achieve a similar perception of risks, it is crucial for the driving assistance system to have a similar utility evaluation of driving environments.

Driving Style: Studies show that human develop strategy base on his/her *minimum acceptable risk* [10], which should be identified as driving style by the assistance system in order to provide warnings that are considered as "necessary".

In contrast to previous research on modeling driver behavior, which treated the components mentioned above as hidden internal states [11], we propose a driving behavior model that explicitly captures the driver's driving decision process. Cognitive digital twins of the driver can be learned based on driver's behaviors, and a Cognitive-digital-twin-based Driving Assistance System (CDAS) was developed to provide personalized warnings based on the digital twin. As shown in Fig. 1, the driving decision process of a driver is abstracted as perceiving the driving environment and driving styles of other agents with partial observation, and forming strategy based on utility evaluation of predicted future driving environments. A cognitive model captures the driving decision process, and a cognitive digital twin of the driver is created and maintained. The estimation of the driver's perception allowed CDAS to identify risks the driver was not aware of, and estimation of driving style was used to provide personalized warnings, which can reduce the number of unnecessary warnings and improve the driver's trust towards the system.

The contribution of this paper is 4-folds: 1) A driver's decision model was proposed that infers the driver's perception of the driving environment and other agents' driving style based on utility evaluation of expected driving environments; 2) a virtual driving simulator was developed, which includes cars controlled by the proposed driver's decision model, and can be utilized to evaluate driving assistance systems; 3) The validity of the driver's decision model and its ability to predict driving behaviors was evaluated in human experiments; 4) A cognitive-digital-twin-based driving assistance system was proposed, and its efficacy in improving driving safety was evaluated against existing ADAS algorithms in simulation.

II. RELATED WORK

Models have been developed to explain and/or predict driving behaviors in complex driving conditions. As summarized in [11], previous driving behavior modeling efforts model the driver's decision process as a Partially Observable Stochastic Game (POSG), in which the driver's mental state including driving styles were modeled as hidden internal states. Sadigh et al. [12] proposed an active information gathering method to infer driving style (distracted/attentive) using inverse reinforcement learning. Lefèvre et al. [13] proposed a Markov state space model to infer whether the driver decides to stop at the intersection. In [14], the authors proposed a probabilistic model to infer driving intent (lane changing). In [15], Chen et. al modeled the driving decision process as decision tree, and shared the model among agents to improve driving safety. These black-box models do not have adequate interpretability of the driving decision process, which cannot be used for personalized driving assistance.

The progress made in cognitive science and behavior science has led to the development of modeling frameworks, which aim to simulate the human decision-making process in driving. In [16], the COSMODRIVE framework was proposed to encompass various aspects of the driving decision process. However, this framework focuses on strategic planning rather than vehicle control, making it unsuitable for modeling specific driving behaviors. The ACT-R framework [17] was introduced to simulate human decision-making processes. However, research utilizing ACT-R for modeling driving behaviors has not addressed uncertainties in a driver's perception and driving style, which are crucial elements for our application of the driver's behavior model.

III. MATHEMATICAL FORMULATION OF DRIVING ASSISTANCE

A. Evolution of the Driving Environment

In complex driving environment with N agents, physical state of an agent i at time t is denoted by x_t^i , which include location, velocity, angle, angular velocity etc. The physical state changes over time based on the action of the agent at time t $a_t^i \in \mathbb{A}$:

$$x_t^i \xrightarrow{a_t^i} x_{t+1}^i$$

The overall state of the driving environment $X_t = \{x_t^i | i \in [1, N]\}$ changes based on the combined actions of all agents $A_t = \{a_t^i | i \in [1, N]\}$:

$$X_t \xrightarrow{A_t} X_{t+1}$$

B. Observation and Perception of the Driving Environment

Each agent can only obtain partial information of the driving environment. We denote $\mathcal{O}^i(\cdot)$ as a mapping from a ground truth physical state to Agent i 's observation. For instance, $\mathcal{O}^i(x_t^j)$ represents Agent i 's observation of Agent j at time t , and $\mathcal{O}^i(X_t)$ represents Agent i 's observation of the overall driving environment at time t .

The driver's *perception* of the driving environment is an estimation of the physical states, which usually deviates from the ground truth. We use $\hat{\square}^i$ to represent Agent i 's estimation of \square . Therefore Agent i 's perception of Agent j can be represented by \hat{x}_t^j and Agent i 's perception of the driving environment can be represented by \hat{X}_t .

C. Uncertainties and Belief

Due to partial observability, there may exist *ambiguities* to the perception of the driving environment. For instance, \hat{x}_t^j represents a distribution in which there are m possible perceived physical states of Agent j by Agent i :

$$\widehat{x}_t^j \sim \widehat{x}_t^j$$

and each perception has a *belief* which sums to 1.

$$b(\widehat{x}_t^{j(m)}) \in [0, 1] \ \&\& \ \sum b(\widehat{x}_t^{j(m)}) = 1$$

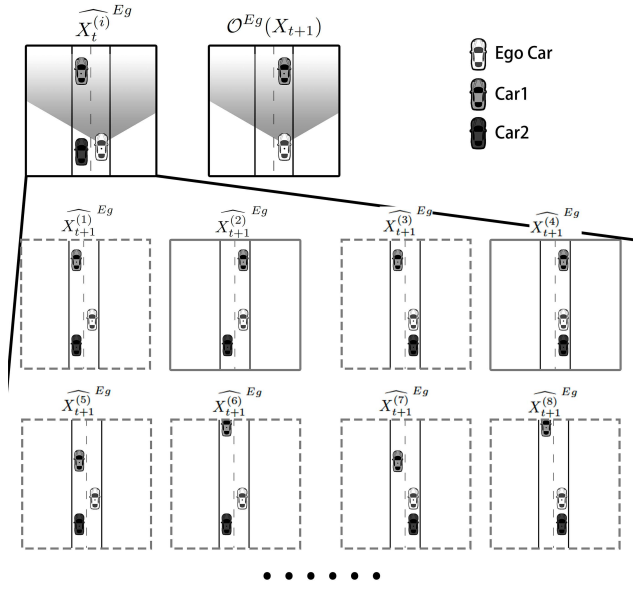


Fig. 2. Predictions of the driving environment at time $t+1$ are made based on each of the driver's perception of the environment at time t . Any predictions that contradict the actual environment observed at time $t+1$ are removed (indicated by the dashed lines).

There are also uncertainties in the predicted actions, which can be modeled as a distribution referred to as *strategy*:

$$\widehat{a}_t^{j(k)} \sim \widehat{\sigma}_t^j$$

$$b(\widehat{a}_t^{j(k)}) \in [0, 1] \text{ and } \sum b(\widehat{a}_t^{j(k)}) = 1$$

D. Prediction of Future Driving Environment

The current strategy is decided based on predicted future outcomes [10]. Therefore each agent predicts the evolution of the driving environment by predicting the actions of all agents including itself. Future state of the driving environment can be predicted such that:

$$\widehat{x}_t^j \xrightarrow{\widehat{\sigma}_t^j} \widehat{x}_{t+1}^j, \widehat{\mathcal{X}}_t^i \xrightarrow{\widehat{\mathbb{A}}_t^i} \widehat{\mathcal{X}}_{t+1}^i$$

Fig. 2 illustrates the predicted driving environments at time $t+1$ from one perceived driving environment at time t .

E. Utility and Driving Style

Driving decisions are made by evaluating the potential risk and benefit of suspected future driving environment after taking certain actions [9]. We denote positive and negative utility evaluation of future driving environment as:

$$U_t^i(\widehat{X}_{t+1}^i) = \{U_t^{i+}(\widehat{X}_{t+1}^i), U_t^{i-}(\widehat{X}_{t+1}^i)\}$$

Here utility is considered to be universally accepted, and individual utility differences are captured by *driving style* of the agent σ_t^i . For instance, driving close to other cars are considered as risky by everyone, but some drivers tolerate small gaps better than others.

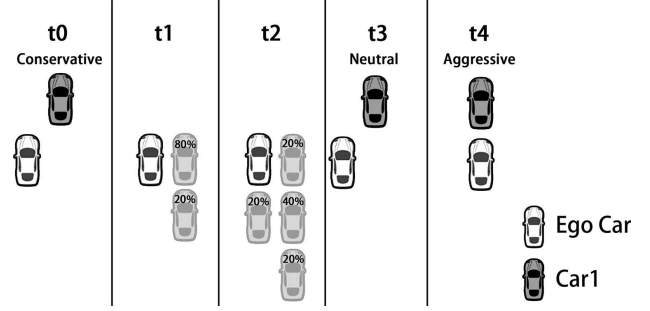


Fig. 3. The potential physical states of unobserved agents and their beliefs are updated in the driver's perception. The perceived driving styles of other agents are updated based on their observed behaviors.

F. Decision Making during Driving

The strategy of agent i depends on the utility of predicted driving environment and the driving style of the agent:

$$\sigma_t^i = \mathcal{D}_t^i(U_t^i(\widehat{\mathcal{X}}_{t+1}^i), \sigma_t^i) \quad (1)$$

Other agents' strategy can be estimated similarly as:

$$\widehat{\sigma}_t^j = \mathcal{D}_t^j(U_t^j(\widehat{\mathcal{X}}_{t+1}^j), \widehat{\sigma}_t^j)$$

G. Perception Update

The perceptions of physical state of observable agents are updated as their observations, and the perceptions of unobservable agents are updated from their perceptions at the last time step:

$$\widehat{x}_t^j = \begin{cases} \mathcal{O}^i(x_t^j) & \text{if } j \text{ is observable by } i \\ \widehat{x}_{t-1}^j \xrightarrow{\widehat{\sigma}_t^i} \widehat{x}_t^j & \text{if } j \text{ is NOT observable by } i \end{cases} \quad (2)$$

As illustrated in Fig. 2, the predicted driving environments at $t+1$ with dashed lines are excluded by the observation at $t+1$.

When agent i was observed doing action $a_t^{i(o)}$ by agent j , the perceived driving style of agent i by agent j $\widehat{\sigma}_t^j$ can be updated such that $a_t^{i(o)}$ would be the most dominant action in predicted strategy:

$$\arg \max_{\widehat{\sigma}_t^j} b(\widehat{a}_t^{i(o)}, \widehat{a}_t^{i(o)}) \sim \mathcal{D}_t^j(U_t^j(\widehat{\mathcal{X}}_t^i), \sigma_{t-1}^j) \quad (3)$$

Fig. 3 illustrates the driving decision formulation with a simple example. The ego car was overtaking car 1 at t_0 , and the ego car's perception of car 1's driving style $\widehat{\sigma}_{t_0}^{Eg}$ was conservative. At t_1 the ego car could not see car 1, but still maintained a perception of its state as $\widehat{x}_{t_1}^{Eg}$, which entropy increased with no new observations at t_2 . At t_3 car 1 overtook the ego car from the slow lane. With new observed behaviors, the uncertainty of $\widehat{x}_{t_3}^{Eg}$ was reduced, and the perceived driving style was updated. Then after car 1 merged in front of the ego car with small gap at t_4 , the behavior is considered as risky and the perceived driving style $\widehat{\sigma}_{t_4}^{Eg}$ is updated to "aggressive".

H. Advanced Driving Assistance System (ADAS)

With advanced sensors like Lidar, the perception of ADAS is the same as its observations.

$$\widehat{X}_t^{AD} = \mathcal{O}^{AD}(X_t)$$

The system predicts all possible outcomes

$$\widehat{X}_t^{AD} \xrightarrow{\mathbb{A}_t^*} \widehat{\mathbb{X}}_{t+1}^{AD}$$

in which possible actions of all agents are considered:

$$\mathbb{A}_t^* = \{a_t^{k_1} \dots a_t^{k_N} \mid \forall i \in [1, N], k_i \in [1, \|\mathbb{A}\|\}\}$$

ADAS outputs warnings and/or actions if there exist an outcome with excessive risk (i.e. collision):

$$\varphi_{warn}^{AD} : \exists \widehat{X}_{t+1}^{(i)AD} \in \widehat{\mathbb{X}}_{t+1}^{AD} \text{ s.t. } \mathcal{U}_t^{AD-}(\widehat{X}_{t+1}^{(i)AD}) > r^{Warn} \quad (4)$$

However, from the driver's perspective, $X_{t+1}^{(i)}$ is either an acceptable outcome, or has very low probability, such that the risk is below the driver's alarming threshold:

$$b(\widehat{X}_{t+1}^{(i)Eg}) \times \mathcal{U}_t^{Eg-}(\widehat{X}_{t+1}^{(i)Eg}) < r_{warn}^{Eg}$$

which makes the warning "redundant".

I. Cognitive-digital-twin-based Driving Assistance System (CDAS)

The driver's trust towards driving assistance system can be significantly improved if a driving assistance system can "empathize with" the driver, such that the driving assistance system provides warnings that are both useful and expected.

In this paper, we propose a *Cognitive-digital-twin-based Driving Assistance System (CDAS)* that estimates the driver's 1) driving style and 2) perception of the driving environment, and provides warnings only when the driver is about to perform an action that he/she would not do with better observations.

The observations of the driver can be estimated via eye tracking devices in order to estimate the driver's perception:

$$\widehat{\mathbb{X}}_t^{EgCD} = \begin{cases} \mathcal{O}^{Eg}(x_t^j) & j \text{ observable} \\ \widehat{\mathbb{X}}_{t-1}^j \xrightarrow{\widehat{\mathbb{A}}_t^{EgCD}} \widehat{\mathbb{X}}_t^j & j \text{ unobservable} \end{cases} \quad (5)$$

The ego's prediction of the strategies of other agents can be estimated using Equation 1. And the driving style of agents are also updated based on their observable behaviors based on Equation 3. The predicted outcomes and their corresponding beliefs can then be estimated:

$$\widehat{\mathbb{X}}_t^{EgCD} \xrightarrow{\widehat{\mathbb{A}}_t^{EgCD}} \widehat{\mathbb{X}}_{t+1}^{EgCD} \quad (6)$$

and the ego strategy can be estimated as:

$$\widehat{\mathbb{O}}_t^{EgCD} = \mathcal{D}_t^{Eg} \left(\mathcal{U}_t^{Eg} \left(\widehat{\mathbb{X}}_{t+1}^{EgCD} \right), \widehat{\sigma}_t^{EgCD} \right) \quad (7)$$

Due to limited observability, there may exist risks that the ego driver may fail to identify. To provide personalized warnings, the driver's strategy with complete observation $\widehat{\mathbb{X}}_{t+1}^{CD}$ can be estimated as:

$$\widehat{\mathbb{O}}_t^{Eg'CD} = \mathcal{D}_t^{Eg} \left(\mathcal{U}_t^{Eg} \left(\widehat{\mathbb{X}}_{t+1}^{CD} \right), \widehat{\sigma}_t^{EgCD} \right) \quad (8)$$

The outcomes with unacceptable risk, which results from the k th action of the ego driver can be calculated as:

$$\widehat{\mathbb{X}}_{risk}^{(k)CD} = \left\{ \widehat{X}_{t+1}^{(i,k)Eg} \mid \mathcal{U}_t^{Eg-}(\widehat{X}_{t+1}^{(i,k)Eg}) > \widehat{\sigma}_t^{Eg} \right\} \quad (9)$$

in which

$$\widehat{X}_t^{(i)EgCD} \xrightarrow{a_{t+1}^{Eg(k)} \mathbb{A}_{t+1}^{Eg}} \widehat{X}_{t+1}^{(i,k)EgCD}$$

The CDAS warns the driver when 1) the risk associated with an action of the ego driver is high, and 2) the ego driver with partial observation may perform the action with high probability, and 3) the driver would not perform the action with full observation:

$$\varphi_{warn}^{CD} : \exists a_t^{Eg(k)} \text{ s.t. } \sum \mathcal{U}_t^{Eg-}(\widehat{\mathbb{X}}_{risk}^{(k)CD}) > r^{CD} \ \&\& \ b(\widehat{\mathbb{O}}_t^{Eg(k)CD}) > p_{warn}^{max} \ \&\& \ b(\widehat{\mathbb{O}}_t^{Eg'(k)CD}) < p_{warn}^{min} \quad (10)$$

In the rest of the paper, we will introduce the implementation and validation of the proposed driver's decision model, the digital twin learning process, as well as the CDAS system.

IV. IMPLEMENTATION OF COGNITIVE DIGITAL TWIN

A. Multi-lane straight highway with no exits

A simple highway driving scenario with no exits was used to demonstrate the proposed framework. Without the need to consider goal and exits, the following assumptions can be made:

The decision mechanisms and utility functions are the same for all agents:

$$a1 : \forall t, i, j, \mathcal{U}_t^i(\cdot) = \mathcal{U}_t^j(\cdot) \ \&\& \ \mathcal{D}_t^i(\cdot) = \mathcal{D}_t^j(\cdot)$$

The decision mechanism and utility functions for each agent do not change over time:

$$a2 : \forall i, t1, t2, \mathcal{U}_{t1}^i(\cdot) = \mathcal{U}_{t2}^i(\cdot) \ \&\& \ \mathcal{D}_{t1}^i(\cdot) = \mathcal{D}_{t2}^i(\cdot)$$

The physical state of each agent include:

$$x_t^i = \{p_t^x, p_t^y, v_t^y, acc_t^y\}$$

which are updated by:

$$v_t^y = v_{t-1}^y + acc_t^y, p_t^y = p_{t-1}^y + v_t^y$$

And the actions of each agent are abstracted for each 1 second decision period:

- **Acc:** $acc_t^y = 1m/s^2$
- **Dec:** $acc_t^y = -2m/s^2$
- **Main:** $acc_t^y = 0$
- **Left:** $acc_t^y = 0, v_t^y = v_{t-1}^y, p_t^x = p_{t-1}^x - LaneWidth$

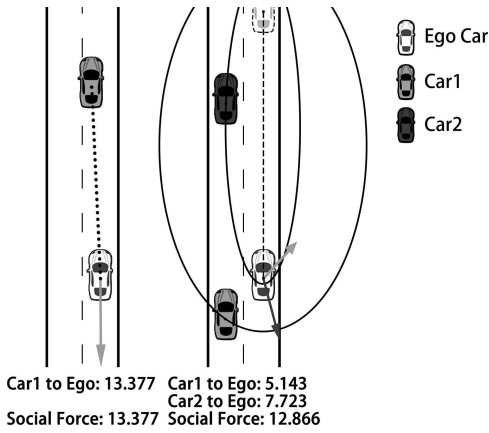


Fig. 4. Depending on the location of the agent, the social force exerted on the ego car is calculated either as the length of the normal of ellipses or the Euclidean distance, following the idea in [18].

- **Right:** $acc_t^y = 0, v_t^y = v_{t-1}^y, p_t^x = p_{t-1}^x + LaneWidth$

The observation of each driver is assumed to be the physical states of agents within 120° in front of the vehicle. The driver's perception are updated each second according to Equation 2.

B. Utility $\mathcal{U}^i(\cdot)$

It is widely assumed that human behaviors are motivated by maximizing expected reward, and the decision making process can be abstracted as weighing possible rewards and penalties [9].

1) *Risk* \mathcal{U}_t^{Eg-} : The following aspects were considered as risks for each predicted outcome \widehat{X}_{t+1}^i :

- Collision $\{0, 100\} * (n - 1)$: Collisions with other agents or guardrails is considered as high risk.
- Social force $[0, 60] * (n - 1)$: The concept of *social force* has been suggested as a means to measure the ideal physical distance that individuals feel comfortable maintaining from one another [18]. As shown in Fig. 4, Euclidean distance and lengths of the normal of ellipses are used to quantify social forces among cars.
- Speeding: $\{0, 15\}$: A minor risk is added when the speed is more than 120km/h.

2) *Reward* \mathcal{U}_t^{Eg+} : In our simplified implementation, the reward is given based on the action:

$$re^{Acc} = 4, re^{Dec} = 1, re^{Main} = 3, re^{Left} = re^{Right} = 3.5$$

Intuitively, accelerating and lane changing may get to the destination faster.

C. Driving Style σ^i

The driving style for each agent is modeled as *Minimum acceptable risk*, such that any potential outcomes that carry a risk level exceeding the minimum acceptable risk are deemed unacceptable by the agent. Therefore acceptable future scenarios for each action k can be calculated as:

$$\mathcal{X}_{accept}^{(k)} = \{ \widehat{X}_{t+1}^{(i,k)Eg} | \mathcal{U}^{Eg-}(\widehat{X}_{t+1}^{(i,k)Eg}) < \sigma_t^{Eg} \} \quad (11)$$

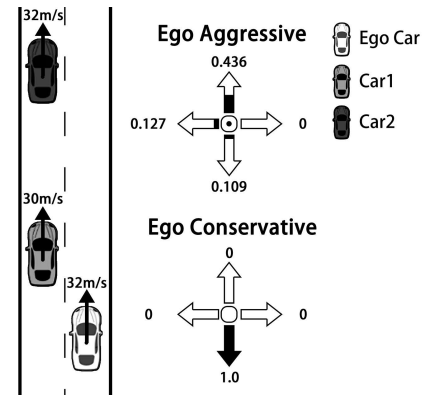


Fig. 5. The conservative driver opted to decelerate in anticipation of car 1 changing lanes, while the aggressive driver chose to accelerate in order to overtake car 1 quickly. The strategies are represented as probabilities of 5 different actions.

D. Driving Decision $\mathcal{D}(\cdot)$

The strategy at time t σ_t^{Eg} is calculated based on the expected reward of each action:

$$b(a_t^{Eg(k)}) = \frac{re^k \times \|\mathcal{X}_{accept}^{(k)}\|}{\sum_k re^k \times \|\mathcal{X}_{accept}^{(k)}\|} \quad (12)$$

In Fig. 5, the influence of driving styles σ^{Eg} on driving decision σ_t^{Eg} is simulated in a scenario where the ego car is faster than two cars on the other lane. Since there is a chance that car 1 may change lane, an aggressive ego driver would overtake, while a conservative one would decelerate. This aligns with our intuitive understanding of typical driving behaviors.

E. Learning of the Cognitive Digital Twin

The CDAS constructs and maintains a cognitive digital twin of the driver, which captures the driver's perception of the driving environment, as well as the driving style.

The driver's perception $\widehat{\mathcal{X}}_t^{EgCD}$ is updated using Equation 5 and 6, and the system's estimation of the driver's driving style $\widehat{\sigma}_t^{EgCD}$ is updated by minimizing the SVM loss [19], which is a common criteria to quantify similarity between a sample and a distribution. Equation 3 can be extended as:

$$\arg \min_{\widehat{\sigma}_t^{EgCD}} \sum_{a_t^{Eg(k)} \in \|A\|}^{k \neq o} \max(0, b(a_t^{Eg(k)CD}) - b(a_t^{Eg(o)CD})) + \Delta \quad (13)$$

V. EXPERIMENTS

Three experiments ¹ were performed to validate the driver's decision model, the learning of digital twin, and the driving assistance based on the digital twin.

¹Experiments involving human subjects were approved by the IRB of ShanghaiTech University with approval No. Q2022-044



Fig. 6. The Virtual Driving Environment in Unity

A. Virtual Driving Simulator

As shown in Fig. 6, a virtual driving environment for the multi-lane straight highway scenario was developed in Unity. Each vehicle is controlled by a driving decision model introduced in Section IV. Three vehicles can be initialized with different physical states as well as driving styles. The red car was designated as the ego vehicle, with a stationary 120° field of view in the forward direction. Driving strategies and/or decisions were collected using questionnaires during simulation. The virtual simulator is used as driving environment in all three experiments. Studies have suggested that individuals may experience a decrease in risk perception when utilizing driving simulators [20]. However, as our experiments are focused solely with the relative variations in risk perception, this disparity is deemed acceptable.

B. Validation of Risk Perception

It is important that the driving decision model has similar perception of risks compared to human drivers.

Experiment Design:

Three agents were initialized in the virtual driving simulator as scenes with different initial states. Agent 1's perception of Agent 2's driving style $\hat{\sigma}^2_1$ was set to 20 and simulated 5 seconds. 14 scenes were selected in which $\hat{\sigma}^2_1$ has increased to above 30, indicating a risky behavior was observed, and 16 scenes in which the number remained under 30. 13 human participants were asked to watch the replay of the 30 scenes and vote *whether the behaviors of Agent 2 showed high minimum acceptable risks* in each scene. The behavior is considered as "risky" if 60% of human participants voted yes and "not risky" if 60% of human voted no, and the rest is considered as "unsure".

Hypothesis:

If majority of human drivers agree that Agent 2 performed an action that lead to risky outcome, the driving style perception $\hat{\sigma}^2_1$ of Agent 1 should increase, indicating that Agent 2 has

model \ human	Risky	Not risky	Unsure
Risky	10	2	2
Not risky	0	15	1

TABLE I
CONFUSION MATRIX OF RISK PERCEPTION

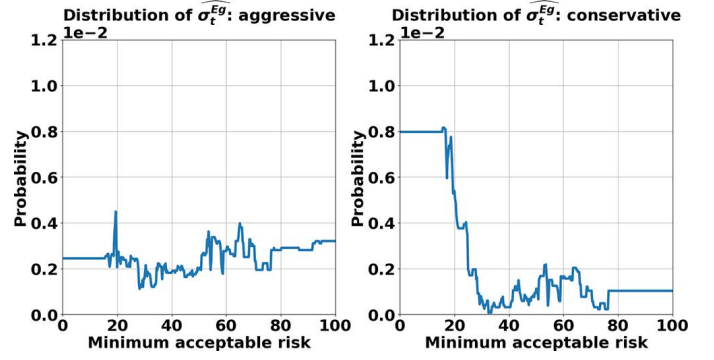


Fig. 7. Cumulative distribution of predicted driving style (minimum acceptable risk) of 13 subjects in 30 scenes, aggressive vs. conservative

high minimum acceptable risk. With human driver's vote as ground truth for risky behaviors, the model should have no false-negatives, and a few false-positives are allowed in order to ensure safety.

Results:

As shown in Table I, the driving decision model has no false-negatives for risk perception with 100% sensitivity and has a few false-positives with 88.2% specificity, which is desired for safety considerations.

C. Learning and Predicting Human Behaviors

It is important for the cognitive digital twin to accurately capture driving style and predict the driving strategies of the ego driver.

Experiment Design:

Three agents were initialized in the virtual driving simulator as scenes with different initial states. The driving style of Agent 1, represented by σ^1 , was set to 20 and simulated for 5 seconds. At the 5th second, two strategies were provided by an aggressive model with $\sigma = 40$ and a conservative model with $\sigma = 18$. 30 scenes were selected in which the Jensen-Shannon divergence (JS divergence) [21] between the two strategies were large (>0.3), indicating the scenes are most effective in distinguishing aggressive and conservative behaviors. 13 human participants were instructed to watch replays of the 30 selected scenes from the first-person perspective of Car 1 and provide their own strategy at the 5th second. The experiment required each participant to drive aggressively and conservatively twice, with the understanding that aggressive/conservative corresponds to high/low minimum acceptable risk. The cognitive digital twin system predicts the human participant's driving style $\hat{\sigma}^{Eg}_{CD}$ and strategy $\hat{\sigma}_t^{Eg}_{CD}$ for further analysis.

Hypothesis:

Aggressive driver should have larger $\hat{\sigma}^{Eg}_{CD}$, and the discrepancy between the predicted strategy and the driver's actual strategy $JSD(\hat{\sigma}_t^{Eg}_{CD}, \sigma_t^{Eg})$ should be small².

Results:

Fig. 7 shows the distribution of $\hat{\sigma}^{Eg}_{CD}$ for all 30 scenes

²Two distributions are considered as similar when JS divergence is <0.1

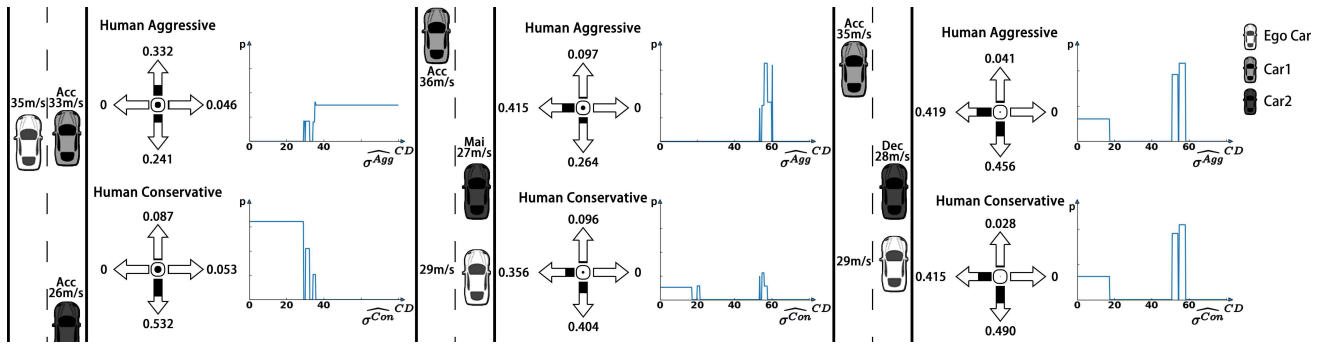


Fig. 8. Accuracy of driving style identification is dependent on the driving scenario. Here are 3 driving scenarios in which aggressive and conservative drivers exhibited 1) distinct strategies, 2) similar but distinguishable strategies, or 3) nearly identical strategies.

by all 13 participants when instructed to drive aggressively and conservatively. The results indicate that the system is more likely to identify drivers who drive conservatively as conservative, whereas there is no definitive identification of aggressive drivers in general. The accuracy of identification is contingent on the type of scene being considered. Fig 8 provides further insight into this trend by presenting 3 distinct scenarios³. In the first scene, aggressive drivers are more inclined to accelerate, while conservative drivers prefer to decelerate. In the second scene, aggressive drivers are more inclined to change lanes rather than slow down, while conservative drivers prefer to decelerate. These scenes demonstrated distinguishable difference in driving style that can be correctly identified by the system. However, in the third scene, the risk of all possible actions exceeds the minimum acceptable risk for both aggressive and conservative drivers. As a result, both types of drivers may choose to slow down to minimize risk. This can make it challenging to distinguish their driving style based on this scene alone. Filtering techniques may be introduced in future work to account for consistency in different driving scenes.

The JS divergence between the human strategies and the system predicted strategies was 0.054 for aggressive driving and 0.058 for conservative driving ($p < 0.05$), indicating good strategy prediction performance.

D. Validation of CDAS

With good prediction performance, the performance of the cognitive digital twin and the corresponding CDAS are evaluated in terms of collision prevention against state-of-the-art Bosch's Blind Spot Detection⁴ and Forward Collision Warning [22].

Experiment Design:

Three agents were initialized in the virtual driving simulator with different initial states for 6000 5sec-scenes. As ego car, Car 1's driving style σ was set to $\{15, 25, 35\}$ to represent "conservative", "regular", and "aggressive" drivers. Without driving assistance, there were $\{940, 1541, 1774\}$ collisions, respectively. CDAS with correct and incorrect $\hat{\sigma}^{CD}$ as well

³Data was normalized among 13 participants in Scenes #24, #15, #28 from Experiment 2

⁴The turning signals were approximated as the initial left or right turn actions within the model

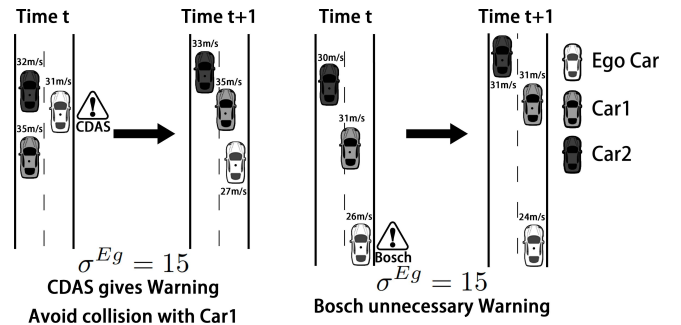


Fig. 9. 1) CDAS warned the ego driver about the approaching aggressive driver and helped avoid a collision; 2) CDAS withheld warning for an anticipated aggressive behavior from an observable Car 2.

as the Bosch ADAS system were implemented on Car 1⁵. Warnings were provided in form of complete observation of the driving environment for the next iteration.

Result:

Table. II illustrates that both Bosch and CDAS were successful in reducing the number of collisions by warning the driver. However, CDAS provided significantly fewer warnings than Bosch, particularly false warnings in scenes where there were no collisions with or without warnings. Overall, CDAS was able to avoid more collisions compared to Bosch ADAS. CDAS with accurate estimation of driving style was more effective in avoiding collisions compared to CDAS with inaccurate driving style estimation. It is worth noting that early warnings altered the original evolution of the driving environment, leading to the occurrence of new collisions.

Fig. 9 shows two scenes that illustrate the advantage of CDAS. In the first scene, an aggressive car 1 was approaching car 2 quickly, and intended to merge in front of the ego car, assuming the ego car would yield. However, since the ego car did not change lanes, the Bosch blind-spot detection algorithm did not alert the driver with acoustic warning. Consequently, car 1 collided with the ego car. On the other hand, CDAS detected that the ego driver was not aware of car 1's approach and predicted that the driver would maintain speed. CDAS also anticipated that car 1's merging could pose a risk to the ego car and that the driver would decelerate if made aware of

⁵Implementation of ADAS and CDAS followed Equation 4 and 10

IEEE Robotics and Automation Letters (RA-L) paper, presented at ICRA 2024, Yokohama, Japan. Cite as RA-L paper.

Ego Style	System	# of collisions w/o warning	# of collisions w/ warning	# of scenes w/ warning	total # of warnings	# of avoided collisions	# of new collisions	cdas ✓ bosch ×	bosch ✓ cdas ×	# of scenes w/ false warning
conservative $\sigma = 15$	Bosch	940	874	2588	5102	91	25			2222
	CDAS:con	940	800	537	582	145	5	54	0	144
	CDAS:reg	940	828	519	558	116	4	27	2	305
regular $\sigma = 25$	Bosch	1541	1397	943	997	147	3			570
	CDAS:reg	1541	1428	643	727	142	29	10	15	262
aggressive $\sigma = 35$	Bosch	1774	1658	821	851	118	3			468
	CDAS:agg	1774	1641	462	573	167	34	49	0	151
	CDAS:reg	1774	1647	526	622	133	6	19	4	232

TABLE II
COLLISION PREVENTION OF BOSCH ADAS AND THE PROPOSED CDAS

the approaching car 1. Therefore, CDAS issued a warning to the ego driver at time t . As a result of the warning, the ego car decelerated and avoided the collision.

In the second scene, car 2 merged in front of the ego car at time t within a short distance, which triggered a front collision warning from the Bosch ADAS. However, in the same situation, CDAS recognized that car 2 was observable by the ego car and predicted that the conservative ego driver would decelerate to allow car 2 to merge in with very low risk. As a result, CDAS withheld the warning.

VI. SUMMARY AND FUTURE WORK

This paper presents a Cognitive-Digital-Twin-Based Driving Assistance System (CDAS) that issues personalized warnings to drivers. The system utilizes a driving decision model and digital twin learning algorithm to capture the driver's perception and driving style. The cognitive digital twin allows the system to issue warnings when there is a significant discrepancy between the predicted driving strategy based on partial observation and that based on full observation. This approach effectively improves driving safety and minimizes the occurrence of false warnings.

Many components of the framework can be extended in future work:

Action space: Driver's observation actions can be taken into account which can improve perception estimation.

Long-term utility: Utility evaluation on future scenarios after multiple actions instead of one action can be used to represent long-term goals.

Driving styles: Whether a driver consider other driver's utility, as well as whether a driver assumes his/her action/perception is known to others can all be extended as driving style.

Experiments: the effect of excessive warnings and how they affect the driver's acceptance of warnings were not accounted for in current virtual CDAS validation, which can be improved with experiments involving human subjects.

REFERENCES

- [1] M. Sivak, "The information that drivers use: Is it indeed 90% visual?" *Perception*, vol. 25, no. 9, pp. 1081–1089, 1996.
- [2] X. S. Zheng and G. W. McConkie, "Two visual systems in monitoring of dynamic traffic: Effects of visual disruption," *Accident Analysis and Prevention*, vol. 42, no. 3, pp. 921 – 928, 2010.
- [3] J. B. Cicchino, "Effects of lane departure warning on police-reported crash rates," *Journal of Safety Research*, vol. 66, pp. 61–70, 2018.
- [4] L. Yue, M. Abdel-Aty, Y. Wu, and L. Wang, "Assessment of the safety benefits of vehicles' advanced driver assistance, connectivity and low level automation systems," *Accident Analysis & Prevention*, vol. 117, pp. 55–64, 2018.
- [5] F. Biassoni, D. Ruscio, and R. Ciceri, "Limitations and automation. the role of information about device-specific features in adas acceptability," *Safety Science*, vol. 85, pp. 179–186, 2016.
- [6] F. Yan, M. Eilers, A. Lüdtkke, and M. Baumann, "Building driver's trust in lane change assistance systems by adapting to driver's uncertainty states," in 2017 IEEE intelligent vehicles symposium (IV). IEEE, 2017, pp. 529–534.
- [7] S. Sánchez-Pagés and M. Vorsatz, "An experimental study of truth-telling in a sender–receiver game," *Games and Economic Behavior*, vol. 61, no. 1, pp. 86–112, 2007.
- [8] C. Press, P. Kok, and D. Yon, "The perceptual prediction paradox," *Trends in Cognitive Sciences*, vol. 24, no. 1, pp. 13–24, 2020.
- [9] P. N. Tobler, J. P. O'Doherty, R. J. Dolan, and W. Schultz, "Reward value coding distinct from risk attitude-related uncertainty coding in human reward systems," *Journal of Neurophysiology*, vol. 97, no. 2, pp. 1621–1632, 2007.
- [10] T. R. Palfrey and H. Rosenthal, "Underestimated probabilities that others free ride: An experimental test," mimeo, California Institute of Technology and Carnegie-Mellon University, Tech. Rep., 1989.
- [11] K. Brown, K. Driggs-Campbell, and M. J. Kochenderfer, "A taxonomy and review of algorithms for modeling and predicting human driver behavior," arXiv preprint arXiv:2006.08832, 2020.
- [12] D. Sadigh, S. S. Sastry, S. A. Seshia, and A. Dragan, "Information gathering actions over human internal state," in 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2016, pp. 66–73.
- [13] S. Lefèvre, C. Laugier, and J. Ibañez-Guzmán, "Intention-aware risk estimation for general traffic situations, and application to intersection safety," Ph.D. dissertation, INRIA, 2013.
- [14] D. D. Salvucci, "Inferring driver intent: A case study in lane-change detection," in *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, vol. 48, no. 19. SAGE Publications Sage CA: Los Angeles, CA, 2004, pp. 2228–2231.
- [15] X. Chen, E. Kang, S. Shirashi, V. M. Preciado, and Z. Jiang, "Digital behavioral twins for safe connected cars," in *Proceedings of the 21th ACM/IEEE International Conference on Model Driven Engineering Languages and Systems*, ser. MODELS '18, 2018, p. 144–153.
- [16] H. Tattegrain-Veste, T. Bellet, A. Pauzić, and A. Chapon, "Computational driver model in transport engineering: Cosmodrive," *Transportation research record*, vol. 1550, no. 1, pp. 1–7, 1996.
- [17] F. E. Ritter, F. Tehrani, and J. D. Oury, "Act-r: A cognitive architecture for modeling cognition," *WIREs Cognitive Science*, vol. 10, no. 3, p. e1488, 2019.
- [18] D. Helbing and P. Molnar, "Social force model for pedestrian dynamics," *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [19] K.-B. Duan and S. S. Keerthi, "Which is the best multiclass svm method? an empirical study," in *International workshop on multiple classifier systems*. Springer, 2005, pp. 278–285.
- [20] R. A. Wynne, V. Beanland, and P. M. Salmon, "Systematic review of driving simulator validation studies," *Safety Science*, vol. 117, pp. 138–151, 2019.
- [21] D. M. Endres and J. E. Schindelin, "A new metric for probability distributions," *IEEE Transactions on Information theory*, vol. 49, no. 7, pp. 1858–1860, 2003.
- [22] Bosch, "Blind spot detection," 2022. [Online]. Available: <https://www.bosch-mobility-solutions.com/en/solutions/assistance-systems/blind-spot-detection/>