

1 **Title:** Bridging the gap between self-report and behavioral laboratory measures:
2 A real-time driving task with inverse reinforcement learning

3

4 **Authors:** Sang Ho Lee^{1,2}, Myeong Seop Song¹, Min-hwan Oh³, and Woo-Young
5 Ahn^{1,2}

6

7 ¹ Department of Psychology, Seoul National University, Seoul, Korea 08826

8 ² Department of Brain and Cognitive Sciences, Seoul National University, Seoul,
9 Korea 08826

10 ³ Graduate School of Data Science, Seoul National University, Seoul, Korea
11 08826

12

13 **Corresponding author:** Woo-Young Ahn, Ph.D.

14 Department of Psychology, Seoul National University, Seoul, Korea 08826

15 Tel: +82-2-880-2538, Fax: +82-2-877-6428.

16 **E-mail:** wahn55@snu.ac.kr

17

18 **Competing Interest Statement:** The authors declare that they have no known
19 financial or non-financial competing interests.

20 **Keywords:** Impulsivity, realistic experiment, inverse reinforcement learning,
21 deep learning, driving

1 **Abstract**

2 A major challenge in assessing psychological constructs such as impulsivity is the weak
3 correlation between self-report and behavioral task measures that are supposed to assess
4 the same construct. To address this issue, we developed a real-time driving task called the
5 “highway task”, where participants often exhibit impulsive behaviors, such as reckless
6 driving, thereby mirroring real-life impulsive traits captured by self-report surveys. Here,
7 we first show that a self-report measure of impulsivity is highly correlated with
8 performance in the highway task, but not with traditional behavioral task measures of
9 impulsivity. By integrating deep neural networks with an inverse reinforcement learning
10 (IRL) algorithm, we inferred dynamic changes of subjective rewards during the highway
11 task. The IRL results indicated that impulsive participants attribute high subjective
12 rewards to irrational or risky driving behaviors and situations. Overall, our results suggest
13 that using real-time tasks combined with IRL can help reconcile the discrepancy between
14 self-report and behavioral task measures of psychological constructs including
15 impulsivity, with IRL being a practical modeling framework for multidimensional data
16 from real-time tasks.

17

1 **Introduction**

2 Self-report and behavioral task measures are among the most frequently employed
3 methods for assessing psychological constructs. A prevalent issue across multiple
4 domains such as impulsivity (1), self-control (2), empathy (3), and risk preference (4) is
5 that self-report and behavioral task measures consistently show weak correlations with
6 each other, even when they are assumed to tap the same construct (5). Weak associations
7 between measures of the same construct fosters ambiguity and confusion in assessment,
8 making it challenging to integrate findings across different measures.

9 We address this problem using impulsivity as a testbed, as it is one of the
10 psychological constructs that is notably affected by the weak association between self-
11 report and behavioral measures. Extensive studies of impulsivity in relation to mental
12 disorders and maladaptive behaviors (6–8) have utilized a range of self-report and
13 behavioral task measures, which are believed to assess the same construct termed
14 “impulsivity”. However, large-scale investigations and meta-analyses consistently
15 reported weak correlations between different measures of impulsivity (1, 9, 10).

16 A widely accepted approach is to view impulsivity as a multidimensional
17 construct with distinct aspects that do not necessarily overlap with each other. For
18 example, MacKillop et al. (9) suggests that measures of impulsivity can be categorized
19 into three distinct domains: impulsive choice, impulsive action, and impulsive personality
20 traits. According to this categorization, behavioral task measures, which reflect impulsive
21 choice (e.g., delay discounting task) (11) and impulsive action (e.g., go/no-go task) (12),
22 do not need to correlate with self-report measures that typically reflect trait impulsivity
23 (e.g., Barratt Impulsiveness Scale) (13). Another explanation for the inconsistency

1 between self-report and behavioral task measures of impulsivity is that they tap the same
2 construct, but the association between them are obscured by differences in measurement
3 methods (10, 14). Self-reports measure individuals' overall tendencies over a longer
4 duration of time (e.g., for the past week/month), whereas laboratory behavioral tasks
5 usually measure specific behaviors in some discrete states (e.g., go and no-go conditions
6 in the go/no-go task) in highly controlled settings at the time of testing. Thus, behavioral
7 task measures may capture state-specific phenomena that only partly reflect self-reported
8 tendencies of behaviors across situations in real life.

9 Building on this methodological explanation, we postulate that a laboratory task
10 conducted in real-time, which mimics real-life situations, would yield impulsivity
11 measures that are strongly correlated with self-reported impulsivity. Specifically, we
12 develop and implement a real-time driving task, dubbed the "highway task" (**Figure 1**),
13 in which participants control a car on a highway to drive as fast as possible without
14 crashing into other cars. The performance in the task may reflect the traits that contribute
15 to reckless driving, which is frequently associated with impulsivity (15, 16). Unlike
16 traditional trial-based laboratory tasks with a predefined list of discrete states, the
17 highway task provides trajectories of states that continuously interact with participants'
18 actions (e.g., accelerate, change lanes), with a number of possible states being virtually
19 boundless (see the following sections for details). Behaviors in this task might be a better
20 reflection of trait impulsivity than traditional behavioral task measures, as the task
21 environment resembles complex real-world situations where impulsive behaviors occur
22 (17, 18).

1 A challenge is how to describe complex data from the highway task beyond
2 simple summary statistics of observed behaviors (e.g., mean speed, number of crashes).
3 Computational modeling is widely recognized as a valuable tool for assessing
4 neurocognitive characteristics underlying behaviors (19). However, traditional
5 computational models may not be readily applicable to data from real-time tasks (e.g.,
6 virtual reality, arcade-style games such as Atari), because such models are not typically
7 designed to describe multidimensional behaviors with an immense number of possible
8 states inherent in a real-time task.

9 We pose this problem as an *inverse* reinforcement learning (IRL) problem, in
10 which the learning algorithm infers the reward function that underlies observed behaviors
11 (20). The objective of IRL is opposite to that of the conventional “forward” approach of
12 reinforcement learning (RL) (21); IRL learns a reward function based on observed
13 behaviors, whereas RL learns a behavioral policy based on observed rewards. Recent
14 advances in algorithmic techniques made IRL well-suited for explaining behaviors in
15 complex environments (22, 23). One of the breakthroughs is the use of deep neural
16 networks (DNN), which can represent complex associations between states and actions in
17 real-time tasks (24), to approximate complex, nonlinear reward functions (22). By
18 integrating DNN with IRL (i.e., deep IRL), we are not restricted to any particular
19 functional form of rewards for observed behaviors. Once learned from the observed data,
20 DNN can calculate rewards for given states and actions. The rewards derived by IRL can
21 be interpreted as a participant’s internal reward or preference, making it a valuable tool
22 for modeling human decision making (25).

1 In the current study, we aim to find indicators of impulsivity in the highway task
2 by comparing IRL-inferred reward functions among participants with varying levels of
3 trait impulsivity. To our knowledge, this is the first study to employ deep IRL to capture
4 individual differences in a psychological construct among human participants. While a
5 few studies modeled human decision making using simpler IRL algorithms using a
6 restricted class of functional forms (25), no other studies have utilized deep IRL to
7 investigate individual differences in reward functions in real-time tasks.

8 If the behaviors exhibited in the highway task align with the hypothesized trait
9 impulsivity, we would expect the task performance and the rewards inferred by IRL to
10 correlate with measures of trait impulsivity. In our experiment (N = 47; determined by
11 the power analysis described in **SI Methods**), we assessed trait impulsivity using the
12 Barratt Impulsiveness Scale (BIS) (13), which is a widely used self-report measure of
13 impulsivity. The experiment consisted of three behavioral tasks, including the highway
14 task and two traditional behavioral tasks measuring impulsivity: delay discounting and
15 go/no-go tasks. While the latter two tasks measure impulsive choice and impulsive
16 action, respectively, the highway task was chosen to examine its correlation with the BIS
17 score in comparison to other measures.

18 In the following analysis, we evaluate the credibility of IRL in explaining
19 individual differences in behavior by assessing its accuracy in predicting actions during
20 the highway task. We then investigate the IRL-inferred reward for each state in the task
21 as well as the *real-time trajectory* of the rewards to find indicators of impulsivity. Finally,
22 the behavioral performance measures (e.g., task score) of the highway task and the output
23 of the IRL are used together to predict the BIS score.

1

2 **Results**

3 We first assessed the validity of the highway task by correlating its task score (see
4 **Materials and Methods** for score calculation) with the total BIS score, which we used as
5 the benchmark measure of trait impulsivity in the current study. The credibility of
6 correlation was assessed using the Bayes factor (BF_{10}) (26) in a Bayesian correlation test
7 (27). Following the classification scheme in Wagenmakers et al. (28), we interpreted
8 BF_{10} of 1-3, 3-10, 10-30, 30-100 and >100 as anecdotal, moderate, strong, very strong,
9 and extreme evidence, respectively. Statistical evidence was strong for the Pearson
10 correlation between the mean score of the highway task and the BIS score ($r = 0.46$, BF_{10}
11 $= 28.41$), suggesting that the task performance improves as the BIS score (i.e.,
12 impulsivity) decreases (see **Figure 2**). The task score was also reliable within the task.
13 The split-half reliability between the scores in the first half and the second half of the task
14 assessed by intraclass correlation coefficient (ICC) (29) was 0.72, which is acceptable
15 (30). In contrast, the BIS score did not correlate with widely used behavioral measures of
16 impulsivity from the two traditional laboratory tasks, namely, the delay discounting rate
17 parameter ($\log k$; see **SI Methods** for model description) in the delay discounting task (r
18 $= 0.01$; $BF_{10} = 0.182$) and the no-go error rate in the go/no-go task ($r = 0.07$; $BF_{10} =$
19 0.203). The results supported our hypothesis that a real-time task in a realistic
20 environment better reflects impulsivity than traditional trial-based tasks. Other statistics
21 derived from the highway task such as the number of crashes ($r = 0.15$; $BF_{10} = 0.3$), mean
22 speed ($r = 0.08$; $BF_{10} = 0.21$), and mean distance from the car ahead ($r = 0.08$; $BF_{10} =$
23 0.21) did not show any statistical evidence for correlation with the BIS score. These

1 measures also showed no evidence for correlation with other behavioral task measures,
2 except that the mean speed considerably correlated with the error rate in the go/no-go
3 task ($r = -0.38$; $\text{BF}_{10} = 4.87$; see **SI Figure S1** for the full correlation matrix).

4 We used IRL to infer a reward function for each individual based on the
5 individual's observed trajectories of behaviors in the highway task (see **SI Methods** for
6 the details of the algorithm). We investigated the individual differences in the reward
7 functions learned via IRL to identify latent indicators of impulsivity. Prior to interpreting
8 the reward functions, we evaluated the goodness of fit of the model trained by IRL. The
9 model fit was assessed by comparing observed participants' actions with artificial agents'
10 actions generated by the behavioral policies of IRL. If IRL learned the reward functions
11 that accurately explain the data, the actions produced by the agent should closely
12 resemble participants' behaviors. **Figure 3A** shows the mean accuracy of the IRL agents
13 in predicting five possible actions in the highway task; moving up, no action (i.e., no-op),
14 moving down, acceleration, and deceleration. The accuracy was much higher than the
15 chance level (mean accuracy = 0.64; chance level accuracy = 0.2) for all actions except
16 the deceleration.

17 The IRL agents also showed similar proportions of actions throughout the action
18 trajectory (**Figure 3B**) to observed human actions. A noticeable difference between the
19 IRL agents and the participants was that the participants showed higher mean proportion
20 of no actions. Cross et al. (31) found comparable differences in actions between humans
21 and artificial agents trained by a "forward" RL. In their study, the policy learned via deep
22 Q-learning (DQN) (24) showed a greater proportion of no actions in Atari games,
23 compared to human participants. The authors postulated that humans are more inclined to

1 abstain from taking action due to metabolic costs and physical constraints (e.g., response
2 speed) for taking action. While the IRL agents learn from human demonstrations that
3 reflect constraints on human behaviors, they might not replicate infrequent inaction due
4 to fatigue or inattention in situations where the participant typically took action.

5 The similarity between the participants' actions and those of IRL agents suggests
6 that the reward functions derived from IRL reflect subjective rewards underlying
7 observed behaviors. We then assessed whether the IRL reward functions are sensible and
8 interpretable by visually examining the reward functions. The DNN trained by IRL
9 approximated subjective rewards for all possible states in the task. The state in the task
10 was defined as a combination of eleven manually annotated features; the speed of the
11 own car, the lane on which the own car is located, the speed of other cars on each of the
12 three lanes (three features for three lanes), the distance from the closest car ahead on each
13 lane (three features), and the distance from the closest car behind on each lane (three
14 features). **Figure 3C** illustrates the mean reward functions (averaged across participants)
15 in simplified state spaces. To visualize and interpret reward functions in a feasible way,
16 we used mean rewards across the two most seemingly important features, own speed (i.e.,
17 y-axis in **Figure 3C**) and the distance from the closest car ahead on the own lane (i.e., x-
18 axis in **Figure 3C**). The high reward states (i.e., dark red area) in the reward function
19 suggests that participants generally favored driving at a low to moderate speed (20 – 60)
20 and a close to moderate distance (10.5 – 63) from the closest car ahead. This reflected a
21 rational strategy to avoid a crash while attempting to overtake a car ahead (i.e., decrease
22 the speed when the distance between the own car and the car ahead is small). By contrast,
23 the state with the smallest distance and the highest speed was associated with extremely

1 low rewards in that the state would likely result in a crash in the next step. The propensity
2 to avoid a crash, which is the most punishing event in the task, is also reflected in the
3 reward functions marginalized over the speed and distance axes (**Figure 3D**). The mean
4 reward tended to increase with the distance from the car ahead and decrease with the
5 speed, suggesting that participants generally used a safe strategy.

6 The results suggest that IRL successfully inferred sensible reward functions from
7 participants' behaviors as we proposed. Nonetheless, our primary objective was to
8 identify indicators of impulsive behaviors that may deviate from rational strategies in the
9 highway task. To achieve this goal, we examined the association between the BIS score
10 and the rewards derived by IRL. In the simplified state space shown in **Figure 3C**, we
11 identified states showing statistical evidence for the correlation between the BIS score
12 and the IRL rewards. Higher BIS score (i.e., increased impulsivity) corresponded to
13 higher rewards for apparently irrational states; maximum speed (120) at close distances
14 (0 – 21) and relatively low speed (50) at far distances (74 – 84) ($r = 0.35-0.39$, $BF_{10} > 3$;
15 see **SI Figure S2** for the correlation coefficients across the state space).

16 We discovered more pronounced indicators of impulsivity from the trajectories of
17 rewards. The reward function generated by IRL provides a simplified representation of
18 how participants' behaviors were interpreted, but it does not depict changes in rewards
19 over time. A real-time task involves trajectories of states and actions. The reward
20 functions inferred by IRL can map these states into reward trajectories, which reveal real-
21 time changes in rewards around significant events. We examined indicators of
22 impulsivity in reward trajectories, with a focus on two salient events in the task:
23 overtaking and crashing. Participants aimed to achieve the highest possible score by

1 quickly overtaking other cars without crashing into them. Video replays of task
2 performance with real-time display of the IRL reward reveal noticeable changes in
3 rewards at the moments of overtaking and crashing (see **SI Results** for the link to a video
4 replay). This implies that the IRL algorithm identified these events as particularly critical.
5 Subsequent analysis of reward trajectories indicates that IRL rewards during overtaking
6 and crashing moments reflect participants' impulsivity.

7 The moments of overtaking and crashing were manually specified in the state
8 space. Overtaking was defined as the moment when a car from an adjacent lane goes
9 behind the participant's car. The distances from other cars ahead and behind the
10 participant's own car were the state features used to identify overtaking moments.
11 Further, two types of overtaking were distinguished: "active" overtaking, which occurs
12 within one second of changing lane, and "passive" overtaking without a lane change.
13 Pictures in the bottom of **Figure 4** depict an example of each event. Active overtaking is
14 more hazardous than passive overtaking because changing lanes for overtaking at a close
15 distance can lead to a collision with a car ahead. We hypothesized that success in this
16 risky behavior would be highly rewarding for impulsive individuals. A crash was
17 straightforward to define, as it was the moment when the distance from a car ahead
18 becomes zero.

19 **Figure 4** shows the reward trajectories before and after overtaking (-3 ~ 1 sec
20 from the onset). To compare the reward functions between participants with high and low
21 impulsivity, we grouped participants into two categories based on their BIS scores: the
22 "high BIS" group and the "low BIS" group, consisting of participants in the highest and
23 the lowest quartiles of the BIS score, respectively. The reward trajectories for passive

1 overtaking (**Figure 4A**) did not differ between the high and low BIS groups, with the IRL
2 rewards showing no correlation with the BIS score across time points. By contrast, active
3 overtaking revealed noticeable differences in the rewards between low and high BIS
4 participants (**Figure 4B**). Compared to the low BIS group, the high BIS group showed a
5 more rapid increase in IRL rewards prior to overtaking. Statistical evidence supported the
6 correlation between the BIS score and the IRL reward at -1.8 – -1.2 seconds ($r = 0.46$;
7 $BF_{10} = 29.5$) and -0.2 seconds ($r = 0.38$; $BF_{10} = 5.6$) from the moment of active
8 overtaking (red dots in **Figure 4B** show the time points). The positive correlations
9 suggest that impulsive participants favored the states with opportunities for active
10 overtaking within a brief time frame (-1.8 – -1.2 seconds), as well as the moment just
11 prior (-0.2 seconds) to successfully accomplishing it.

12 Reward trajectories for overtaking is likely to reflect participants' intention,
13 because participants should overtake as many cars as possible to maximize the task score.
14 Another salient event, crashing, is different in that it is an abrupt and unintentional event
15 that should be avoided. The reward function on the distance axis (**Figure 3D**) show that
16 the IRL reward drops at the moment of a crash (i.e., zero distance). Rewards at this
17 moment correlated with the BIS score, but only at specific states where participants
18 decelerated ($r = 0.38$; $BF_{10} = 4$). Impulsive participants (i.e., high BIS group) heavily
19 discounted the reward for deceleration immediately before (-0.2 – 0 seconds) a crash,
20 whereas non-impulsive participants (i.e., low BIS group) showed relatively steady
21 rewards until the occurrence of a crash (see **SI Figure S3** for the reward trajectories for
22 crashing). Rewards for other actions around a crash did not show such associations with
23 the BIS score.

1 Preceding analyses found several indicators of impulsivity, one from a
2 performance measure (i.e., task score) and others from rewards inferred by IRL. The
3 variables that respectively correlate with the BIS score raise a question of whether using
4 them altogether would help explain individual differences in impulsivity. To compare the
5 informativeness of different types of variables, we predicted the BIS scores across
6 individuals with regression analysis (i.e., LASSO) (32) using independent variables of
7 three different categories: Traditional measures, highway task performance measures, and
8 IRL measures. Traditional measures included error rate on no-go trials (variable name in
9 the LASSO models: GNG error) and response time on go trials (GNG RT) in the go/no-
10 go task and the logarithm of the delay discounting rate parameter ($\log(k)$) in the delay
11 discounting task. Highway task performance measures were the task score (Highway
12 score), mean speed (Highway speed), mean distance from the car ahead (Highway
13 distance), and the frequency of crashing (Highway n(crash)). In the models that used IRL
14 rewards, only the rewards that correlated with the BIS score were included. They were
15 the mean of the rewards for the subset of states that correlated with the BIS score in the
16 two-dimensional state space depicted in **Figure 3C** (IRL speed*distance), mean of the
17 rewards marked by red points in the reward trajectory for active overtaking (IRL
18 overtaking), and the reward for deceleration at the moment of a crash (IRL crash).

19 **Figure 5** shows the results from LASSO, which was conducted using glmnet (33)
20 and easymml package (34) in Python. Model prediction was evaluated by the correlation
21 between predicted and observed values of the BIS score in the test set. The histograms in
22 **Figure 5** illustrates the distribution of the correlation coefficients. As expected by the
23 correlation analyses, traditional measures from the delay discounting task and go/no-go

1 task were unable to predict the BIS score (**Figure 5A**; $r = 0.1$). The performance
2 measures in the highway task showed the correlation similar to the correlation between
3 the task score and the BIS score (**Figure 5B**; $r = 0.48$ vs. 0.46), with the task score
4 explaining most of the variance (see the beta coefficients on the bottom row of **Figure**
5 **5B**). The selected variables from the IRL reward function were better than the
6 performance measures at predicting the BIS score (**Figure 5C**; $r = 0.72$). Finally, a model
7 with both performance measures and IRL measures did not show a better correlation
8 score than the “IRL only” model (**Figure 5D**; $r = 0.72$), suggesting that the behavioral
9 performance measures do not explain additional variance in the BIS score beyond the
10 variance explained by IRL rewards.

11

12 **Discussion**

13 The lack of correlation between self-report and behavioral task measures of
14 psychological constructs has long been a puzzle. We hypothesized that this discrepancy
15 may be due to the simplicity of traditional behavioral tasks. Our findings regarding
16 impulsivity demonstrate that measures derived from a real-time behavioral task do indeed
17 correlate with a relevant self-report measure. This suggests that behavioral task measures
18 can represent individual traits measured with a self-report questionnaire if the task offers
19 a wide range of states where participants can exhibit diverse behaviors as in real-world
20 situations.

21 The novelty of the current study stems from employing a deep IRL algorithm to
22 extract participants’ reward functions and individual differences in a real-time task. Past
23 studies that associated impulsivity with driving behavior in real-world and simulated

1 environments typically focused on simple summary statistics of behaviors such as
2 speeding, crashing, and traffic violations (15, 35). However, we found stronger indicators
3 of impulsivity from IRL rewards than from summary statistics (e.g., mean speed, number
4 of crashes). This suggests that IRL offers more than just descriptive analysis, as the
5 reward functions can provide insights into participants' characteristics that may not be
6 apparent in their behaviors. The successful application of IRL in this study highlights the
7 potential of IRL as a modeling framework for real-time tasks that generate multi-
8 dimensional data that are not easily described by conventional computational models.

9 Black-box machine learning models, which include DNN in the current IRL
10 algorithm (23), have demonstrated high predictive performance, but often lack
11 interpretability in their predictions (36). This absence of explanatory power has restricted
12 the use of black-box models in human behavior research, where explanation is as
13 important as prediction (37). IRL addresses this issue by providing each participant's
14 rewards, which can be interpreted similarly to subjective values in computational models
15 of decision making (e.g., (38, 39)). Participants would choose actions of highest
16 subjective values (or IRL reward) or actions leading to the states with the highest
17 subjective values. This approach enhances the interpretability of the model, making it
18 more suitable for studying human behaviors.

19 Having behavioral task measures of trait impulsivity might help address some
20 concerns about self-report measures. Concerns regarding the credibility of self-report
21 measures exist due to potential response biases, which include responses influenced by
22 social desirability, a consistent response tendency towards affirmative or negative
23 responses, and a propensity for extreme or mid-point responses (40). In a naturalistic

1 paradigm including the highway task, participants are less likely to mask their traits or
2 intentionally influence the assessment, because they are not directly questioned about
3 their real-life tendencies and behaviors. A crucial next step for future research is to
4 evaluate the validity of the highway task independently with clinical populations. For
5 instance, we could investigate whether patients with impulsivity-related psychiatric
6 disorders (e.g., substance use disorders; attention deficit hyperactivity disorder) exhibit
7 lower scores in the highway task compared to healthy participants.

8 A remaining question is whether the rewards inferred by IRL truly represent the
9 internal rewards experienced by participants, as hypothesized. A promising approach to
10 address this question would be to investigate the associations between reward functions
11 learned via IRL and brain activities related to reward (or value) processing. Rewards
12 inferred by IRL might correspond to the representation of subjective values of predicted
13 outcomes in the brain, which has been shown to correlate with functional magnetic
14 resonance imaging (fMRI) activity in the regions such as orbitofrontal cortex (OFC) (41),
15 ventromedial prefrontal cortex (vmPFC) (42), and ventral striatum (39). Predicting real-
16 time changes in brain activities in these areas using the IRL rewards would help interpret
17 IRL reward functions as subjective value functions that underlie human decision making.
18 This approach would also validate the use of IRL in understanding the relationship
19 between rewards and their neural correlates.

20 The current study focused on impulsivity measures, but our approach can be
21 applied to other real-time tasks assessing different constructs (43). The use of the
22 highway task aligns with recent studies employing realistic and real-time tasks to enhance
23 the ecological validity of neuropsychological assessment (17). The adoption of real-time

1 tasks and data has increased, as recent technological advances (e.g., virtual reality,
2 mobile devices) have facilitated experiments in realistic settings (44). Our work suggests
3 that deep IRL serves as a practical modeling framework, enabling researchers to fully
4 utilize complex data from real-time tasks without being restricted to simple descriptive
5 analysis. Reward functions inferred by a deep IRL algorithm might reflect participants'
6 subjective rewards or intentions in the task, which are central variables in the theories and
7 models of decision making. In summary, the combination of real-time tasks and deep IRL
8 offers a promising novel approach to improving the assessment of psychological
9 constructs underlying human behaviors and decision making.

10

11 **Materials and Methods**

12 The experiment was approved by the institutional review board of Seoul National
13 University (IRB No. 2112/004-004).

14 **Participants.** Forty-seven undergraduate and graduate students in Seoul National
15 University participated. A Bayesian power analysis determined the number of
16 participants (see **SI Methods** for details of the power analysis).

17 **Procedures.** Participants completed one questionnaire (Barratt impulsiveness
18 scale) and three behavioral tasks (highway task, delay discounting task, and go/no-go
19 task) in a dimly lit room. They could take a break between the tasks as long as they
20 desired. An experimenter gave instructions to the participants at the beginning of each
21 task. The questionnaire was controlled by Qualtrics (qualtrics.com) on a web browser.
22 Behavioral experiments were controlled by Python scripts.

1 **Barratt Impulsiveness Scale (BIS-11).** We used a Korean version of the Barratt
2 Impulsiveness Scale (45) to measure trait impulsivity. The questionnaire contained 30
3 questions that corresponds to the question in the English version of BIS-11 (13)
4 Participants answered the questions on a four-point scale (Rarely/Never, Occasionally,
5 Often, Almost Always/Always). Each answer was scored 1-4, with 4 indicating the most
6 impulsive response. The BIS score was calculated by summing the scores across
7 questions. The subscales of BIS (i.e., motor, nonplanning, and attentional impulsivity)
8 were the sum of the scores across subsets of questions (see **SI Results** for the analysis
9 using the BIS subscales).

10 **Highway task.** The task was built on a collection of OpenAI gym (46)
11 environments for driving tasks (47). Task display and action input were controlled by
12 Pygame package in Python. The goal of the highway task was to drive the green car on
13 the screen as fast as possible without crashing into yellow cars (see **Figure 1** for the task
14 display). Participants controlled the green car by pressing the arrow keys on the
15 keyboard. The left and right arrow keys decreased and increased the speed by 10
16 distance/second, respectively. The up and down arrow keys moved the green car to the
17 upper (left) lane and the lower (right) lane, respectively. The score in an episode
18 incremented by $0.2 \times (\frac{speed}{10})^2$ every 0.2 second. An episode continued until the
19 remaining fuel becomes zero or the green car crashes with a yellow car. The remaining
20 fuel, which was displayed on the top of the screen, started as 60 and decreased at the rate
21 of 1/second. Crashing into another car immediately decreased the score by 200. The
22 score was reset to zero at the beginning of each episode. The highest score a participant

1 achieved in an episode was recorded as the “high score” (text below the road) until the
2 participant score higher in another episode (see **SI Methods** for additional details of the
3 task).

4 **Traditional behavioral tasks.** See **SI Methods** for details of the delay
5 discounting task and the go/no-go task.

6 **Inverse reinforcement learning (IRL) algorithm.** We inferred the reward
7 functions underlying the trajectories of states and actions in the highway task using
8 adversarial inverse reinforcement learning (AIRL) (23), which is an IRL algorithm that
9 achieves state-of-the-art performance. IRL is a challenging problem because multiple
10 policy and reward functions can explain a given set of observed behaviors, leading to
11 ambiguity in the learned reward function (20). AIRL is built on maximum entropy IRL
12 (22, 48), which mitigates the ambiguity in solution by identifying a single reward
13 function that maximizes the entropy of the policy derived from the rewards (49). AIRL
14 also addresses the complexity of behaviors in a real-time task by using DNN to
15 approximate nonlinear reward functions, whereas a majority of IRL methods (e.g., (48,
16 50)) assume linear reward functions that might be simplistic in complex tasks (see **SI**
17 **Methods** for the details of the algorithm).

18

19 **Acknowledgments**

20 This work was supported by the National Research Foundation grant funded by the
21 Ministry of Science, Information and Communication Technologies and Future Planning,
22 the Korean Government (Grant No. 2021M3E5D2A0102249311 [to W-YA]); the BK21
23 FOUR Program (Grant No. 5199990314123 [to W-YA]); the Creative-Pioneering

1 Researchers Program through Seoul National University [to W-YA]; and the Creative
2 Challenge Research Program through the National Research Foundation of Korea funded
3 by the Ministry of Education (Grant No. 2022R1I1A1A01066530 [to SHL]). We thank
4 Adam Gazzaley, Jay Myung, Mark Pitt, and Robert Whelan for their constructive
5 feedback on the earlier version of the manuscript.

6
7
8
9
10

References

- 11 1. L. Sharma, K. E. Markon, L. A. Clark, Toward a Theory of Distinct Types of “Impulsive” Behaviors: A
12 Meta-Analysis of Self-Report and Behavioral Measures. *Psychol Bull* 140, 374–408 (2014).
- 13 2. B. Saunders, *et al.*, Reported Self-control is not Meaningfully Associated with Inhibition-related
14 Executive Function: A Bayesian Analysis. *Collabra: Psychol.* 4 (2018).
- 15 3. B. A. Murphy, S. O. Lilienfeld, Are Self-Report Cognitive Empathy Ratings Valid Proxies for Cognitive
16 Empathy Ability? Negligible Meta-Analytic Relations With Behavioral Task Performance. *Psychol. Assess.*
17 31, 1062–1072 (2019).
- 18 4. R. Frey, A. Pedroni, R. Mata, J. Rieskamp, R. Hertwig, Risk preference shares the psychometric
19 structure of major psychological traits. *Sci. Adv.* 3, e1701381 (2017).
- 20 5. J. Dang, K. M. King, M. Inzlicht, Why Are Self-Report and Behavioral Measures Weakly Correlated?
21 *Trends Cogn. Sci.* 24, 267–269 (2020).
- 22 6. J. L. Perry, M. E. Carroll, The role of impulsive behavior in drug abuse. *Psychopharmacology* 200, 1–26
23 (2008).
- 24 7. S. P. Whiteside, D. R. Lynam, The Five Factor Model and impulsivity: using a structural model of
25 personality to understand impulsivity. *Pers Individ Differ* 30, 669–689 (2001).
- 26 8. J. MacKillop, *et al.*, Multidimensional Examination of Impulsivity in Relation to Disordered Gambling.
27 *Exp Clin Psychopharm* 22, 176–185 (2014).
- 28 9. J. MacKillop, *et al.*, The latent structure of impulsivity: impulsive choice, impulsive action, and
29 impulsive personality traits. *Psychopharmacology* 233, 3361–3370 (2016).
- 30 10. M. A. Cyders, A. Coskunpinar, The relationship between self-report and lab task conceptualizations of
31 impulsivity. *J Res Pers* 46, 121–124 (2012).
- 32 11. L. Green, J. Myerson, A Discounting Framework for Choice With Delayed and Probabilistic Rewards.
33 *Psychol Bull* 130, 769–792 (2004).

- 1 12. C. M. Hartung, R. Milich, D. R. Lynam, C. A. Martin, Understanding the Relations Among Gender,
2 Disinhibition, and Disruptive Behavior in Adolescents. *J Abnorm Psychol* 111, 659–664 (2002).
- 3 13. J. H. Patton, M. S. Stanford, E. S. Barratt, Factor structure of the barratt impulsiveness scale. *J. Clin.*
4 *Psychol.* 51, 768–774 (1995).
- 5 14. M. A. Cyders, A. Coskunpinar, Measurement of constructs using self-report and behavioral lab tasks: Is
6 there overlap in nomothetic span and construct representation for impulsivity? *Clin Psychol Rev* 31, 965–
7 982 (2011).
- 8 15. J. Hatfield, A. Williamson, E. J. Kehoe, P. Prabhakaran, An examination of the relationship between
9 measures of impulsivity and risky simulated driving amongst young drivers. *Accid Analysis Prev* 103, 37–
10 43 (2017).
- 11 16. F. O’Brien, M. Gormley, The contribution of inhibitory deficits to dangerous driving among young
12 people. *Accid Analysis Prev* 51, 238–242 (2013).
- 13 17. K. Robertson, M. Schmitter-Edgecombe, Naturalistic tasks performed in realistic environments: a
14 review with implications for neuropsychological assessment. *Clin Neuropsychologist* 31, 16–42 (2017).
- 15 18. A. Verdejo-Garcia, *et al.*, A unified online test battery for cognitive impulsivity reveals relationships
16 with real-world impulsive behaviours. *Nat Hum Behav* 5, 1562–1577 (2021).
- 17 19. S. Palminteri, V. Wyart, E. Koehlin, The Importance of Falsification in Computational Cognitive
18 Modeling. *Trends Cogn Sci* 21, 425–433 (2017).
- 19 20. S. Arora, P. Doshi, A survey of inverse reinforcement learning: Challenges, methods and progress. *Artif*
20 *Intell* 297, 103500 (2021).
- 21 21. R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction* (MIT press, 2018).
- 22 22. M. Wulfmeier, P. Ondruska, I. Posner, Maximum Entropy Deep Inverse Reinforcement Learning. *Arxiv*
23 (2015) <https://doi.org/10.48550/arxiv.1507.04888>.
- 24 23. J. Fu, K. Luo, S. Levine, Learning Robust Rewards with Adversarial Inverse Reinforcement Learning.
25 *Arxiv* (2017) <https://doi.org/10.48550/arxiv.1710.11248>.
- 26 24. V. Mnih, *et al.*, Human-level control through deep reinforcement learning. *Nature* 518, 529–533
27 (2015).
- 28 25. R. Zhang, *et al.*, Modeling sensory-motor decisions in natural behavior. *Plos Comput Biol* 14,
29 e1006518 (2018).
- 30 26. R. E. Kass, A. E. Raftery, Bayes Factors. *Journal of the American Statistical Association* 90, 773–795
31 (1995).
- 32 27. R. Wetzels, E.-J. Wagenmakers, A default Bayesian hypothesis test for correlations and partial
33 correlations. *Psychon B Rev* 19, 1057–1064 (2012).
- 34 28. E.-J. Wagenmakers, *et al.*, Bayesian inference for psychology. Part II: Example applications with JASP.
35 *Psychon B Rev* 25, 58–76 (2018).

- 1 29. T. K. Koo, M. Y. Li, A Guideline of Selecting and Reporting Intraclass Correlation Coefficients for
2 Reliability Research. *J Chiropr Medicine* 15, 155–163 (2016).
- 3 30. D. V. Cicchetti, Guidelines, Criteria, and Rules of Thumb for Evaluating Normed and Standardized
4 Assessment Instruments in Psychology. *Psychol Assessment* 6, 284–290 (1994).
- 5 31. L. Cross, J. Cockburn, Y. Yue, J. P. O’Doherty, Using deep reinforcement learning to reveal how the
6 brain encodes abstract state-space representations in high-dimensional environments. *Neuron* 109, 724-
7 738.e7 (2021).
- 8 32. R. Tibshirani, Regression Shrinkage and Selection Via the Lasso. *J Royal Statistical Soc Ser B*
9 *Methodol* 58, 267–288 (1996).
- 10 33. J. Friedman, T. Hastie, R. Tibshirani, Regularization Paths for Generalized Linear Models via
11 Coordinate Descent. *J Stat Softw* 33, 1–22 (2010).
- 12 34. W.-Y. Ahn, P. Hendricks, N. Haines, EasymL: Easily Build and Evaluate Machine Learning Models.
13 *Biorxiv*, 137240 (2017).
- 14 35. P. Bıçaksız, T. Özkan, Impulsivity and driver behaviors, offences and accident involvement: A
15 systematic review. *Transp Res Part F Traffic Psychology Behav* 38, 194–223 (2016).
- 16 36. C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use
17 interpretable models instead. *Nat Mach Intell* 1, 206–215 (2019).
- 18 37. T. Yarkoni, J. Westfall, Choosing Prediction Over Explanation in Psychology: Lessons From Machine
19 Learning. *Perspect Psychol Sci* 12, 1100–1122 (2017).
- 20 38. C. S. Sripada, R. Gonzalez, K. L. Phan, I. Liberzon, The neural correlates of intertemporal decision-
21 making: Contributions of subjective value, stimulus type, and trait impulsivity. *Hum. Brain Mapp.* 32,
22 1637–1648 (2011).
- 23 39. J. W. Kable, P. W. Glimcher, The neural correlates of subjective value during intertemporal choice. *Nat*
24 *Neurosci* 10, 1625–1633 (2007).
- 25 40. A. Furnham, M. Henderson, The good, the bad and the mad: Response bias in self-report measures.
26 *Pers Individ Differ* 3, 311–320 (1982).
- 27 41. J. A. Gottfried, J. O’Doherty, R. J. Dolan, Encoding Predictive Reward Value in Human Amygdala and
28 Orbitofrontal Cortex. *Science* 301, 1104–1107 (2003).
- 29 42. M. P. Paulus, L. R. Frank, Ventromedial prefrontal cortex activation is critical for preference
30 judgments. *Neuroreport* 14, 1311–1315 (2003).
- 31 43. J. A. Anguera, *et al.*, Video game training enhances cognitive control in older adults. *Nature* 501, 97–
32 101 (2013).
- 33 44. T. D. Parsons, Virtual Reality for Enhanced Ecological Validity and Experimental Control in the
34 Clinical, Affective and Social Neurosciences. *Front Hum Neurosci* 9, 660 (2015).

1 45. S.-R. Lee, *et al.*, The Study on Reliability and Validity of Korean Version of the Barratt Impulsiveness
2 Scale-11-Revised in Nonclinical Adult Subjects. *Journal of the Korean Neuropsychiatric Association* 51,
3 378–386 (2012).

4 46. G. Brockman, *et al.*, OpenAI Gym. *Arxiv* (2016) <https://doi.org/10.48550/arxiv.1606.01540>.

5 47. E. Leurent, An Environment for Autonomous Driving Decision-Making. *GitHub repository* (2018).

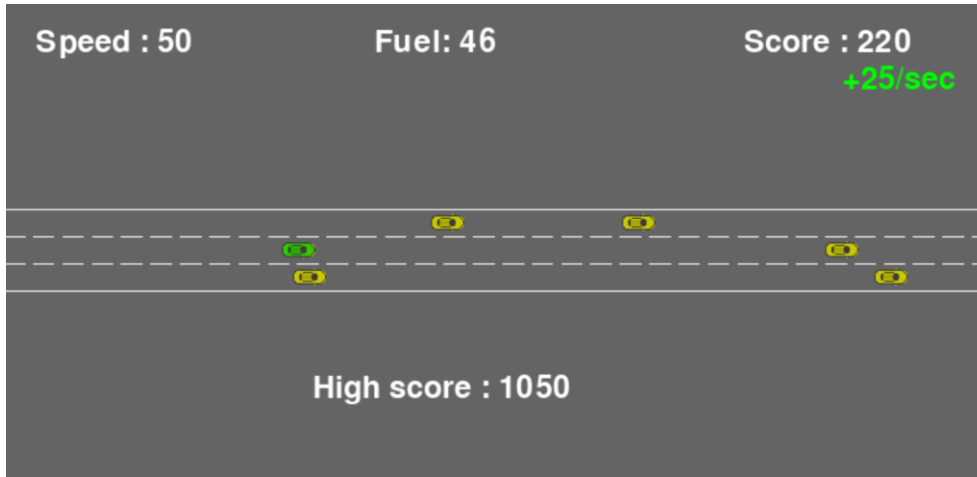
6 48. B. D. Ziebart, A. Maas, J. A. Bagnell, A. K. Dey, Maximum Entropy Inverse Reinforcement Learning
7 in *AAAI Conference on Artificial Intelligence*, (2008), pp. 1433–1438.

8 49. A. J. Snoswell, S. P. N. Singh, N. Ye, Revisiting Maximum Entropy Inverse Reinforcement Learning:
9 New Perspectives and Algorithms. *Arxiv* (2020) <https://doi.org/10.48550/arxiv.2012.00889>.

10 50. P. Abbeel, A. Y. Ng, Apprenticeship learning via inverse reinforcement learning. *Twenty-first Int Conf*
11 *Mach Learn - Icml '04*, 1 (2004).

12

13

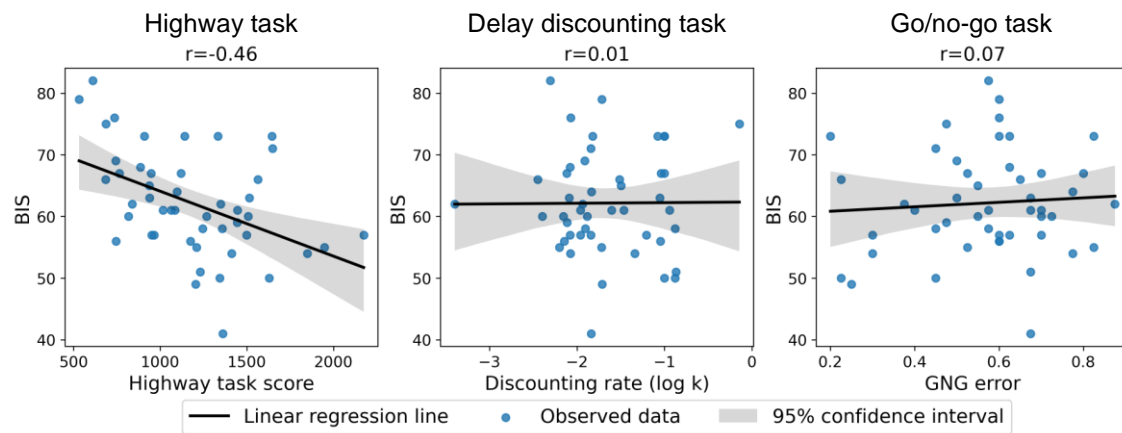


1

2 Figure 1. A screenshot of the highway task. Participants control the green car to drive as
3 fast as possible without crashing into yellow cars. Score per second increases with speed.

4 An episode (or trial) continues until the car crashes or runs out of fuel. High score
5 indicates the highest score achieved in an episode during the highway task.

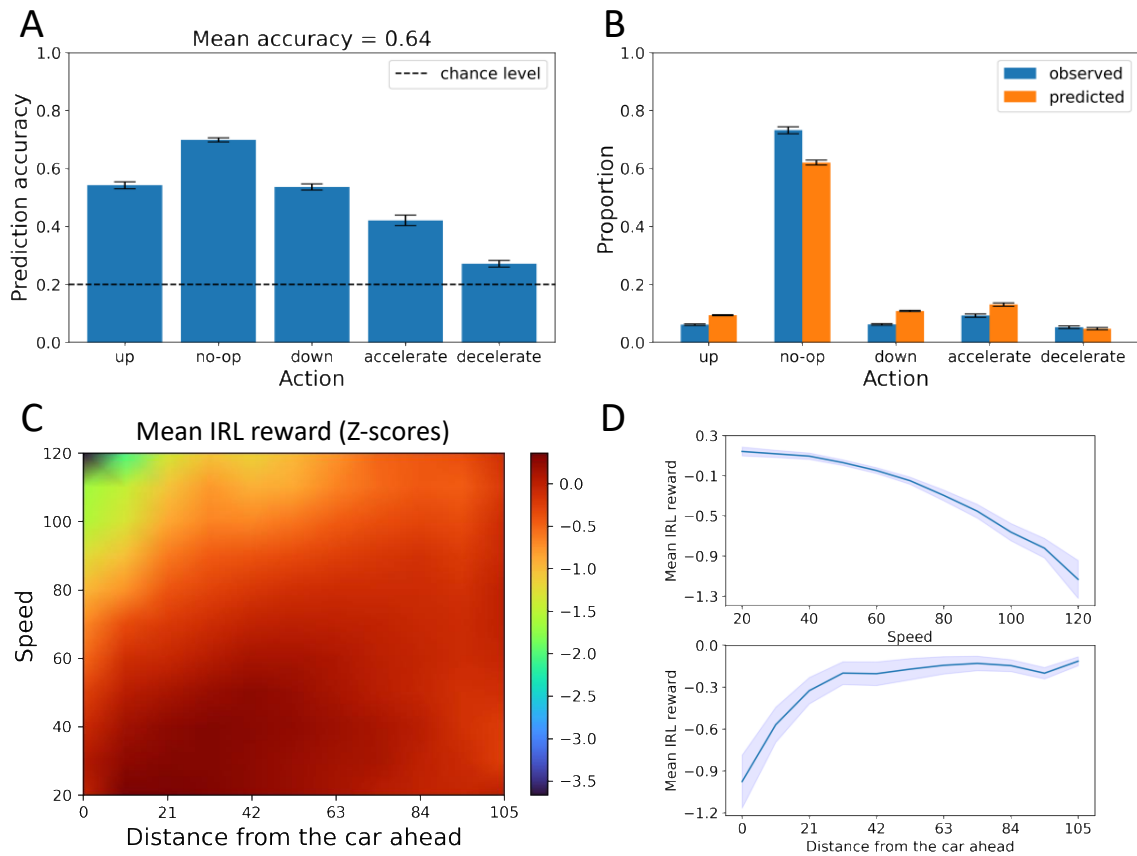
6



1

2 Figure 2. The correlation between the BIS score and behavioral task measures of
 3 impulsivity.

4

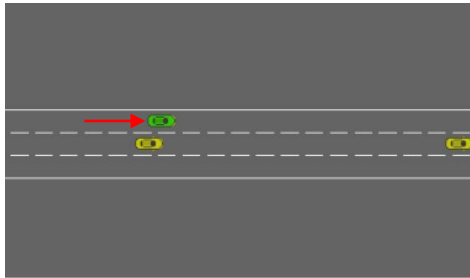
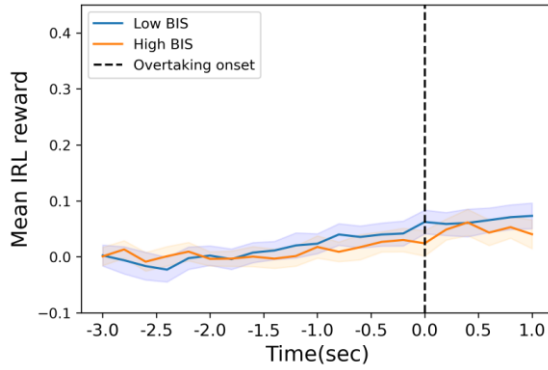


1

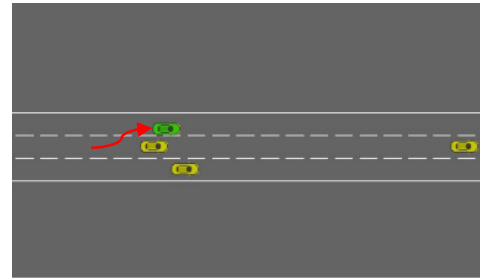
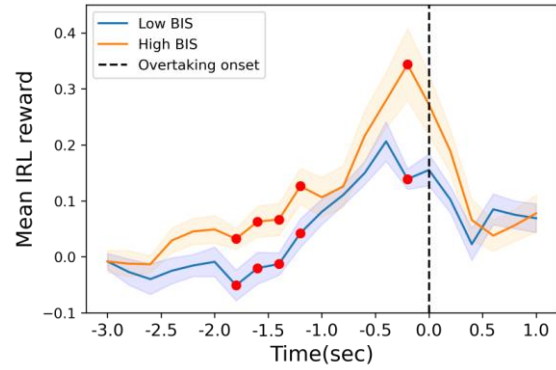
2 Figure 3. Performance of the IRL algorithm. A) Mean accuracy of the five possible
 3 actions produced by the IRL agents. B) Mean proportions of actions in the observed data
 4 (blue) and the action trajectories generated by the IRL agents (orange). C) Mean IRL
 5 rewards in a state space defined by the combination of the speed of the green car (y-axis)
 6 and the distance from the closest car on the same lane (x-axis). D) Mean IRL rewards
 7 marginalized over the speed and distance axes. Error bars in A-B) and the blue areas in
 8 D) indicate standard errors of the means.

9

A. Passive overtaking



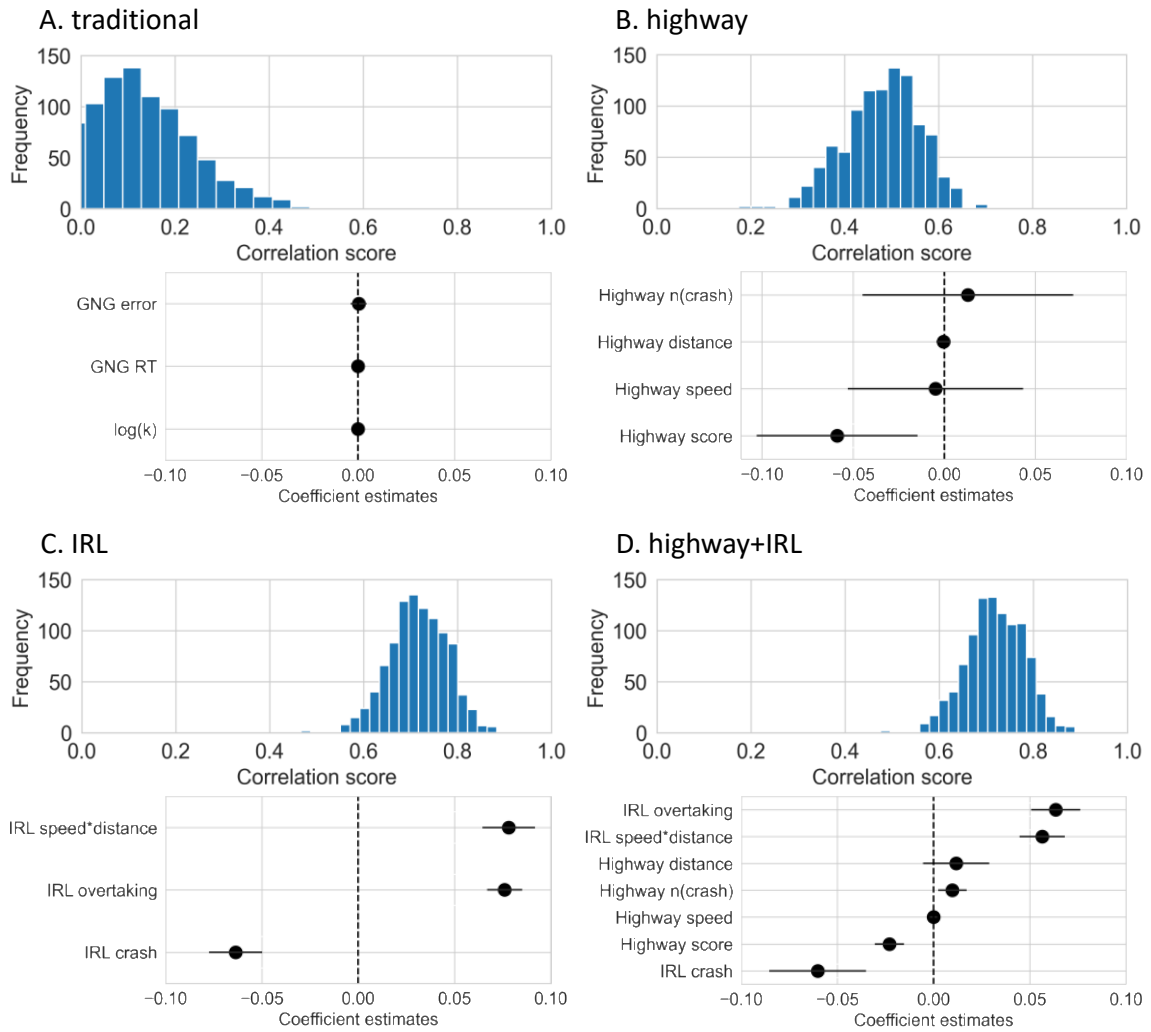
B. Active overtaking



1

2 Figure 4. Reward trajectories before and after overtaking moments. The red points on the
3 lines in B indicate the time points at which the rewards credibly correlated with the BIS
4 score ($BF_{10} > 3$). The pictures on the bottom are examples of passive and active
5 overtaking in the highway task. Red arrowed lines in the pictures illustrate the movement
6 trajectories of the own (green) car.

7



1

2 Figure 5. Model fit and beta coefficients of the Lasso models that predicted the BIS score
 3 using different independent variables. The histograms show the distributions of the
 4 correlation score, which is the correlation coefficients between observed and predicted
 5 values of the BIS score. The graphs displayed below the histograms depict the beta
 6 coefficient for each variable. The error bars represent the 95% confidence intervals.

7