



## Predicting driver takeover performance in conditionally automated driving

Na Du<sup>a</sup>, Feng Zhou<sup>b</sup>, Elizabeth M. Pulver<sup>c</sup>, Dawn M. Tilbury<sup>d</sup>, Lionel P. Robert<sup>e</sup>, Anuj K. Pradhan<sup>f</sup>, X. Jessie Yang<sup>a,\*</sup>

<sup>a</sup> Industrial and Operations Engineering, University of Michigan, United States

<sup>b</sup> Industrial and Manufacturing Systems Engineering, University of Michigan-Dearborn, United States

<sup>c</sup> State Farm Mutual Automobile Insurance Company, United States

<sup>d</sup> Mechanical Engineering, University of Michigan, United States

<sup>e</sup> School of Information, University of Michigan, United States

<sup>f</sup> Industrial and Mechanical Engineering, University of Massachusetts Amherst, United States



### ARTICLE INFO

#### Keywords:

Transition of control  
Predictive modeling  
Human–automation interaction  
Human–autonomy interaction  
Human–robot interaction

### ABSTRACT

In conditionally automated driving, drivers have difficulty taking over control when requested. To address this challenge, we aimed to predict drivers' takeover performance before the issue of a takeover request (TOR) by analyzing drivers' physiological data and external environment data. We used data sets from two human-in-the-loop experiments, wherein drivers engaged in non-driving-related tasks (NDRTs) were requested to take over control from automated driving in various situations. Drivers' physiological data included heart rate indices, galvanic skin response indices, and eye-tracking metrics. Driving environment data included scenario type, traffic density, and TOR lead time. Drivers' takeover performance was categorized as good or bad according to their driving behaviors during the transition period and was treated as the ground truth. Using six machine learning methods, we found that the random forest classifier performed the best and was able to predict drivers' takeover performance when they were engaged in NDRTs with different levels of cognitive load. We recommended 3 s as the optimal time window to predict takeover performance using the random forest classifier, with an accuracy of 84.3% and an F1-score of 64.0%. Our findings have implications for the algorithm development of driver state detection and the design of adaptive in-vehicle alert systems in conditionally automated driving.

### 1. Introduction

While automated vehicles are poised to revolutionize surface transportation, they introduce new challenges. One of the challenges is takeover transitions in conditionally automated driving (Ayoub et al., 2019; Zhou et al., 2020b). In conditionally automated driving, drivers are no longer required to actively monitor the driving environment and are allowed to fully engage in non-driving-related tasks (NDRTs) (Society of Automotive Engineers, 2018). However, serving as a fallback for the automation, drivers are required to take over control of the vehicle whenever the automated system reaches its operational limit.

Previous studies showed that the limited driver–vehicle interaction in conditionally automated driving increases the difficulty for drivers to take over control when requested (Eriksson and Stanton, 2017; Gold et al., 2016; Petersen et al., 2019). In response to such difficulty, empirical studies have investigated the factors that influence drivers' takeover performance, including drivers' cognitive and emotional states

(Du et al., 2020; Wan and Wu, 2018; Zeeb et al., 2017) and driving environments (Gold et al., 2016; Li et al., 2018).

These studies shed light on the relationships between certain factors and takeover performance; for instance, high traffic density harmed takeover performance (Gold et al., 2016). However, with few exceptions (Gold et al., 2018; Braunagel et al., 2017), little effort has been made to integrate these findings into computational models that are capable of predicting drivers' takeover performance in real time. In the present study, therefore, we aimed to fill the research gap and to predict drivers' takeover performance when they were engaged in NDRTs with different levels of cognitive load.

#### 1.1. Factors influencing takeover performance

To facilitate takeover transitions, empirical research has been conducted to examine factors that influence drivers' takeover performance. The factors include drivers' cognitive and emotional states when

\* Corresponding author at: 1205 Beal Avenue, Ann Arbor, MI 48015, United States.

E-mail address: [xijyang@umich.edu](mailto:xijyang@umich.edu) (X.J. Yang).

performing different types of NDRTs (Wan and Wu, 2018; Zeeb et al., 2017; Du et al., 2020) in different driving environments (Gold et al., 2016; Li et al., 2018). Takeover performance consists of takeover timeliness (i.e., takeover reaction time) and takeover quality (e.g., speed, acceleration and jerk statistics, time/distance to collision statistics, lane deviation statistics, and crash rate).

The types of NDRTs have been found to influence takeover performance. Previous studies showed that compared with not performing an NDRT, those engaged in NDRTs had longer takeover reaction times, more crashes in high-traffic situations, and shorter minimum time to collision (TTC) (Eriksson and Stanton, 2017; Gold et al., 2016; Wan and Wu, 2018). The effects of NDRT modality on takeover performance were also explored. For example, Radlmayr et al. (2014) and Wandtner et al. (2018) reported that a visual task with handheld devices degraded takeover performance and led to a higher collision rate, while an auditory task led to comparable performance to a baseline without any task. Zeeb et al. (2016) and Zeeb et al. (2017) explored the effects of manual and cognitive task load and found that a high level of manual task load increased reaction time and deteriorated takeover quality, while the effect of cognitive task load on takeover ability was dependent on the type of driver intervention. A high level of cognitive load lengthened the reaction time and deteriorated takeover quality in steering maneuvers but not braking maneuvers.

Driving environment factors include traffic density, road situations, and weather conditions. Heavy traffic density in takeover situations led to longer takeover time and worse takeover quality in the form of shorter time to collision, more collisions, and higher maximum accelerations (Gold et al., 2016; Körber et al., 2016; Radlmayr et al., 2014). Li et al. (2018) showed that drivers' takeover reaction time to critical events in adverse weather conditions was longer on the highway compared to on city roads. Takeover request (TOR) lead time is the critical event onset for automation failures at the time of the TOR (McDonald et al., 2019). According to the complexity of driving environment and vehicle sensor capability, commonly used TOR lead times range from 1 to 30 s (Eriksson et al., 2018). Research has demonstrated that shorter TOR lead time degraded takeover quality, as demonstrated by higher crash rates, greater maximum accelerations and greater standard deviation of steering wheel angle (Mok et al., 2015; van den Beukel and van der Voort, 2013; Wan and Wu, 2018).

Most of these studies focused on the effects of certain variables on takeover performance, providing valuable yet largely relational insights. For instance, heavy traffic density led to longer takeover time. However, knowing the relationships between certain factors and takeover performance is not enough to accurately predict a driver's takeover performance in the real world because many influential factors could interact with one another. Computational models capable of predicting drivers' takeover performance under various takeover conditions in real time are needed.

## 1.2. Predicting drivers' states through physiological measurements

With advances in wearable technology, it is possible to collect drivers' physiological signals, such as gaze behaviors, heart rate activity, and galvanic skin responses, for a reliable reflection of their cognitive and emotional states in conditionally automated driving.

Drivers' gaze behavior is a valid tool for measuring cognitive load (Solovey et al., 2014; Wang et al., 2014; Gold et al., 2016; Zeeb et al., 2016; Luo et al., 2019) and visual scanning patterns have been shown to indicate situational awareness (Ratwani et al., 2010; Young et al., 2013; Bertola and Balk, 2011). For example, Gold et al. (2016) found that horizontal gaze dispersion was the most sensitive measure of drivers' cognitive demand in NDRTs during conditionally automated driving. Eyes-on-the-road percentage was found to be associated with drivers' situational awareness and attention capture of the driving environments (Young et al., 2013; Molnar, 2017).

Heart rate (HR) and heart rate variability (HRV) have both been used

for assessing drivers' workload in real time (Mehler et al., 2012, 2009; Zhou et al., 2020). Galvanic skin responses (GSRs) were found to reflect drivers' mental activities, and their properties (amplitude, frequency) were used to indicate drivers' changes of arousal related to events (Collet et al., 2009). GSRs have also been linked to drivers' workload and stress (Schmidt et al., 2016; Jones et al., 2014; Wandtner et al., 2018).

Physiological data can thus be used to understand drivers' cognitive and emotional states by applying machine learning models to continuously monitored physiological data. The data captured via non-intrusive sensors can be used to build models that estimate drivers' states and their interactions with the driving environments. Drivers' physiological signals combined with environment factors are promising indicators to predict takeover performance in conditionally automated driving in real time (Braunagel et al., 2017).

## 1.3. Existing models for takeover performance prediction

Although a substantial amount of research has identified factors that influence drivers' takeover performance, there is a lack of research on the development of computational models for predicting drivers' takeover performance, with few exceptions (Gold et al., 2018; Braunagel et al., 2017).

To predict takeover performance, Gold et al. (2018) analyzed 753 takeover events using data from six driving simulator experiments and developed regression models. Their study modeled takeover performance measures (e.g., take-over time, minimum TTC, brake application and crash probability) as a function of the time-budget, traffic density, non-driving-related task, repetition, the current lane and driver's age. The models were validated using 729 takeover events from five additional experiments. The validation results showed that the regression models accurately predicted takeover time, time-to-collision and crash probability, and moderately predicted the brake application.

Braunagel et al. (2017) used machine learning algorithms to predict drivers' takeover quality (named as "takeover readiness" in the article). The study categorized takeover quality into low and high levels by analyzing driving parameters such as lane deviations. Data were collected from a driving simulator study with 81 participants. The first feature input was situation complexity with three levels decided by raters; the second set of features was the type of NDRTs performed by drivers; and the third set of features was drivers' gazes at the road. Using machine learning algorithms including k-nearest neighbors (kNN), support vector machine (SVM) with radial basis function (RBF) and linear kernel, Naive Bayes and linear discriminant, they predicted takeover quality with an accuracy of 79% and F1-score of 77%.

However, the above-mentioned models were developed and tested when drivers were engaged in different types of NDRTs (e.g., monitoring vs. reading), where apparent contextual cues existed to discriminate drivers' states. In daily life, even with a specific type of NDRTs such as writing an email, drivers' states can be rather different depending on the importance of the email. Also, some factors deliberately manipulated in the experiment settings such as emotions are not easily accessible in the real world. Although the advanced wearable technology has made it convenient to collect drivers' physiological signals to reflect their cognitive and emotional states, only gaze behaviors were used in previous studies.

## 1.4. The present study

Our study contributes to the literature in three aspects. First, our study aimed to predict drivers' takeover performance when they were engaged in a specific type of NDRTs with different levels of cognitive load. Second, in addition to gaze behaviors, we used drivers' heart rate indices and galvanic skin response indices to indicate their interaction with environments, which might improve prediction results. Third, our study employed a random forest model in addition to the machine

learning models used in previous studies to predict takeover performance. Random forests have been proved to have great prediction performance for classification problems (Zhou et al., 2020; McDonald et al., 2014; Dietterich, 1997).

In this paper, data from two human subject experiments were used for model development. We collected drivers' galvanic skin responses (Mehler et al., 2012; Collet et al., 2009; Wintersberger et al., 2018), heart rate activities (Mehler et al., 2012; Bashiri and Mann, 2014), and gaze behaviors (Radlmayr et al., 2014; Wang et al., 2014; Young et al., 2013; Bertola and Balk, 2011), which have been used as valid signals to assess drivers' cognitive and emotional states and their situational awareness of the driving environments. Using drivers' physiological data and environment factors, we developed a random forest model that was able to predict drivers' takeover performance with an accuracy of 84.3% and an F1-score of 64.0% using a 3 s time window. Additionally, we identified the most important physiological measures for takeover performance prediction, which can be incorporated in practice to develop in-vehicle monitoring systems. Furthermore, the model can be used to guide the design of adaptive in-vehicle alert systems to improve takeover performance in conditionally automated driving.

## 2. Dataset

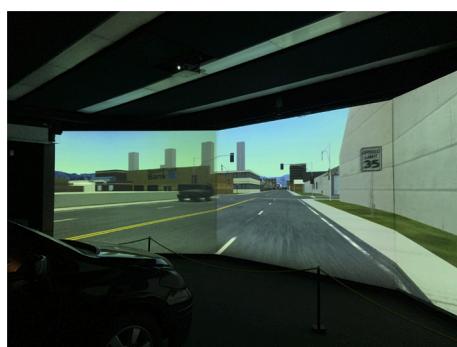
The data used in the development of algorithms were collected in two studies. Both studies complied with the American Psychological Association code of ethics and were approved by the institutional review board at the University of Michigan. The first study investigated the effects of cognitive load, traffic density, and TOR lead time on takeover performance. The second study examined the effects of scenario type and vehicle speed on takeover performance. Participants in both experiments wore the same set of physiological sensors. The similar experimental settings in both studies make it possible to combine the two datasets. At the same time, the varieties of takeover conditions from the two studies increase model generalizability.

### 2.1. Participants

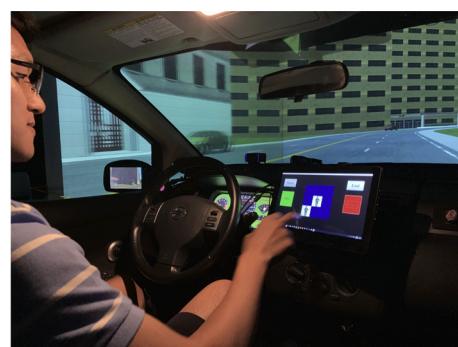
A total number of 102 university students (mean age = 22.9; standard deviation [SD] = 3.8; range = 18–38; 40 females and 62 males) participated in Study 1 and 40 university students (mean age = 22.8, SD = 3.9; 20 females and 20 males) participated in Study 2. All of the participants had normal or corrected-to-normal vision and a valid driver's license. They received \$30 in compensation for an hour of participation.

### 2.2. Apparatus and stimuli

Both studies were conducted in a fixed-base driving simulator from Realtime Technologies Inc. (RTI, MI, USA). The virtual world was



(a)



(b)

**Fig. 1.** RTI driving simulator at the UMTRI.

projected on three front screens 16 ft away (120° field of view), one rear screen 12 ft away (40° field of view), and two side mirror displays (see Fig. 1a).

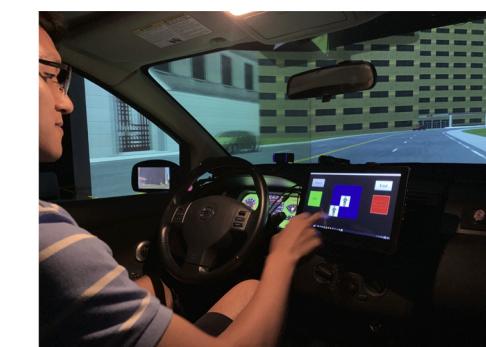
This simulator was equipped with the Smart Eye four-camera eye-tracking system (Smart Eye, Sweden) that provided live head-pose, eye-blink, and gaze data (Fig. 2a). The sampling rate of the eye-tracking system is 120 Hz. The Shimmer3 GSR+ unit (Shimmer, MA, USA) including GSR electrodes and photoplethysmogram (PPG) probe was used to collect GSR and HR data with a sampling rate of 128 Hz (Fig. 2b). The iMotions software (iMotions, MA, USA) was used for physiological data synchronization and visualization in real time (Fig. 2c).

The simulated vehicle was controlled by a steering wheel and pedal system embedded in a Nissan Versa car model. The vehicle was programmed to simulate SAE Level 3 automation, which handled the longitudinal and lateral control and navigation, and responded to traffic elements. Participants could press the button on the steering wheel to activate the automated mode, which was indicated by a green highlight on the dashboard and an auditory warning ("Automated mode engaged"). Once the AV reached its performance limit, an auditory TOR ("Takeover") would be issued with the green highlight turning to black background on the dashboard. Although the Level 3 automation is considered to continue functioning for a certain period of time after issuing the TOR (ISO, 2020), we set the automated mode to be deactivated at the time of TORs for drivers to take over control of the vehicle.

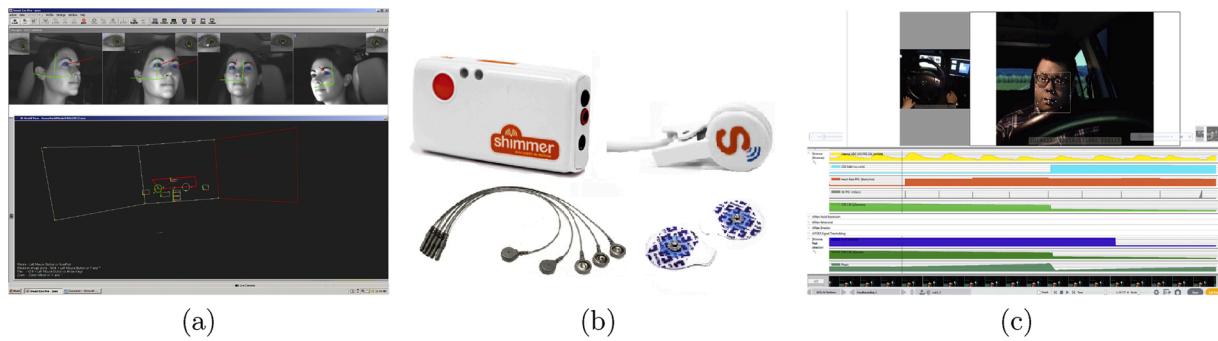
The NDRT in both studies was a visual *N*-back memory task, adapted from the study of Jaeggi et al. (2008). The stimulus consisted of nine (3 × 3) squares with two human figures randomly in two of the nine squares. Each stimulus was presented for 500 ms in sequence with a 2500-ms interval (Fig. 3). Participants were required to press the "Hit" button when the current stimulus was the same as the one presented *N* steps back in the sequence and press the "Reject" button otherwise. With different *N* values, participants were exposed to different cognitive load but the same manual and visual load. The reason for employing a visual task with manual input was that it simulated the eyes-off the road and hands-off the wheel condition. The task was running on an 11.6-in. touch screen tablet mounted in the vehicle (Fig. 1b).

### 2.3. Experimental design

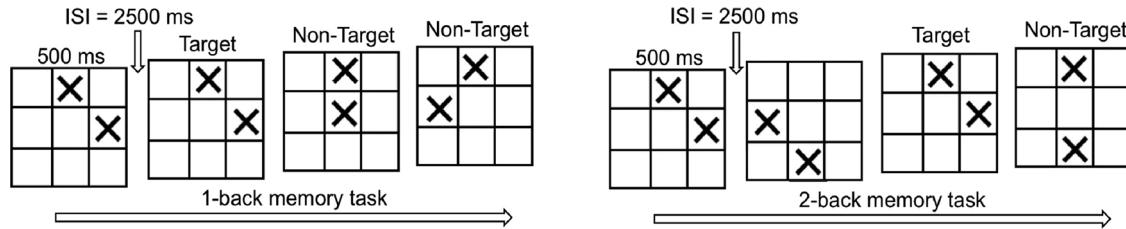
Study 1 employed a within-subjects design with drivers' cognitive load, traffic density, and TOR lead time as independent variables. The cognitive load refers to driver cognitive load prior to TORs and was manipulated via the difficulty of the NDRTs (low: 1-back memory task; high: 2-back memory task). The heavy- and no-traffic conditions had 15 and 0 oncoming vehicles per kilometer, respectively (Gold et al., 2016). The TOR lead time, which refers to the critical event onset for failures at the time of the TOR (McDonald et al., 2019), was set at 4 or 7 s (Eriksson and Stanton, 2017). Based on prior literature (Koo et al., 2016; Miller et al., 2016; Molnar et al., 2018; Rezvani et al., 2016), eight takeover



(b)



**Fig. 2.** (a) Smart Eye. (b) Shimmer3 GSR+ unit. (c) iMotions software.



**Fig. 3.** N-back memory task.

events were designed in urban and rural drives with typical roadway features: (1) bicyclists ahead, (2) construction zone on the left, (3) construction zone ahead, (4) sensor error on the right curve, (5) swerving vehicle ahead, (6) no lane markings on the curve, (7) sensor error on the left curve, and (8) police vehicle on shoulder. The order of cognitive load, traffic density and TOR lead time was counterbalanced via an  $8 \times 8$  balanced Latin square across participants. Considering standard programming practices for the simulator, the order of scenario presentations was counterbalanced by having half of the participants drive from Events 1 to 8, and the other half from Events 8 to 1.

Study 2 used a mixed design with scenario type (lane keeping vs. lane changing) as the between-subjects variable and vehicle speed (35 vs. 60 mph) as the within-subjects variable. Similar to the first study, eight scenarios were designed on the basis of realistic situations and previous literature (Miller et al., 2016; Koo et al., 2016; Rezvani et al., 2016; Zeeb et al., 2016; Naujoks et al., 2014). Lane-keeping scenarios, which required drivers to keep in the current lane, included (1) sensor error on the left curve, (2) construction zone on the left, (3) no lane markings on the curve, (4) sensor error on the right curve. Lane-changing scenarios, which required drivers to change to the neighboring lane, included (1) stranded vehicle ahead, (2) construction zone ahead, (3) construction barrier ahead, and (4) police vehicle on shoulder. According to the range of the Velodyne Lidar sensors (Velodyne Lidar, CA, USA), we set the distance between obstacle/entrance of the curve and the AV at 100 m when the TOR was issued. Generally, traffic consisted of 15 oncoming vehicles per kilometer (Gold et al., 2016). The order of the vehicle speed was counterbalanced among participants. The order of scenarios was counterbalanced by having half of the participants drive from Events 1 to 4, and the other half from Events 4 to 1.

In both studies, drivers started from the right lane, and were asked to stay in the right lane before they engaged the automated mode. Thus, the AV was always in the right lane prior to the TORs and the objects could be pre-coded to appear in front of the vehicle in lane-changing scenarios. With two lanes in lane-changing scenarios, drivers could avoid the objects in their lane by changing to the adjacent lane because there were no other vehicles in the driver's direction. The speed of the subject vehicle was 35 mph in the urban/rural and 60 mph in the highway environments. The radius of curves was 400 m in the highway and 100 m in the urban/rural environments. Participants were asked to

follow the speed limit throughout the drive.

#### 2.4. Experimental procedure

The procedures of the two studies were almost the same. After participants signed an informed consent form and completed an online demographics questionnaire, they were asked to track six targets on the front screen for eye-tracking calibration. Next, two GSR electrodes were attached to their left foot and the PPG probe to their left ear lobe. Participants were informed that there was no need to actively monitor the driving environments or take over control of the vehicle as long as the vehicle was in automated mode.

Participants had a 2-min practice for the N-back memory task, followed by a 5-min practice drive to get familiar with the simulator environment. Next, each participant drove two experimental drives (10–20 min each), each containing four (Study 1) or two (Study 2) takeover events. At the beginning of the drive, participants were asked to activate the AV mode and then start the N-back task when the audio command “Please start the NDRT” was issued. After about 90 s of NDRT, a TOR was issued unexpectedly, and participants were required to terminate the NDRT manually by pressing the “end” button on the tablet screen and take over the control immediately. When participants thought they had negotiated the takeover event, they were free to reactivate the AV mode. Participants were informed that they would get an additional \$20 if their NDRT performance was ranked in the top 10 of all participants. The operation of the NDRT, the takeover, and the AV mode activation were repeated for each takeover event (Fig. 4).

#### 3. Takeover performance model development

We collected drivers' physiological data, driving behaviors, and environment-related data. The physiological measures included heart rate indices, galvanic skin response indices and eye-tracking metrics. Because of malfunctions of the driving simulator and physiological sensors, data from 13 participants were excluded and those of the other 129 participants (i.e., 828 takeover scenarios) were available for further analysis.

To develop the prediction model, we first pre-processed the raw data and then extracted 37 features and set the ground truth. Next, we used a

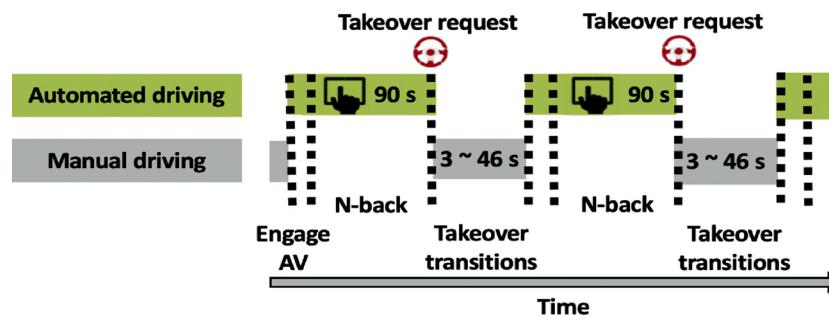


Fig. 4. Illustration of the experimental procedure for two takeover events.

10-fold nested cross-validation method to tune hyper-parameters, train models, and predict test instances for model comparisons. Particularly, we resampled the training dataset and normalized the entire dataset before performing the classification. Fig. 5 shows the modeling process.

### 3.1. Data pre-processing

For GSR signals, we used continuous decomposition analysis (CDA) to decompose the GSR signal into phasic and tonic components, respectively, via Ledalab in Matlab (Benedek and Kaernbach, 2010). Then we used the phasic component for further feature extraction because it is responsible for relatively rapid changes in response to specific events in the GSR signal (order of seconds). Heart rate measures were extracted from the raw RR interval using iMotions software. For eye-tracking data, only data points with high gaze quality value (threshold recommended by Smart Eye:5) were recorded and used for analysis.

### 3.2. Feature generation and ground truth

To fit time series data into the supervised learning framework, we aggregated the values of physiological data within a sliding “time window” and calculated various statistics (Anderson, 2011). The end of the time window is the time of a TOR, and the start of the time window is  $X$  seconds before the TOR, ranging from 1 to 30 s. Model inputs included data on gaze behaviors, galvanic skin response indices, and heart rate indices, as well as environment factors. The generated features are listed in Table 1. A fixation is defined as “a relatively stable eye-in-head position within some threshold of dispersion (typically  $\sim 2^\circ$ ) over some minimum duration (typically 100–200 ms), and with a velocity below some threshold (typically 15–100° per second)” (Jacob and Karn, 2003). In the Smart Eye eye-tracking system, all frames with a gaze velocity below the fixation threshold (100° per second) were treated as a fixation. All frames with the gaze velocity above the saccade threshold (100°

Table 1

Descriptions of generated features (HR = heart rate; min = minimum; max = maximum; GSR = galvanic skin responses; NDRT = non-driving-related task; TOR = takeover request).

Feature	Explanations
HR indices	Mean, min, max, and standard deviation of heart rate, inter-beat interval
GSR indices	Mean, max, and standard deviation of GSR in phasic component
GSR peak	The number of GSR peaks, and peak rise time
Fixation	Fixation number and duration in different areas of interests (AOIs) (i.e., driving scenes and NDRT tablet)
Saccade	Saccade number in different AOIs (i.e., driving scenes and NDRT tablet)
Pupil	The mean and standard deviation of pupil diameter in different AOIs (i.e., driving scenes and NDRT tablet)
Blink	The number of blinks
Gaze dispersion	Standard deviation of the values for gaze angle from right front (radians)
Eyes-on-the-road	The proportion of time that participants' gazes are on the road
Scan pattern	The probability of eyes switching from one AOI to another (i.e., the probability that drivers transited eyes from driving scenes to NDRT tablet, from NDRT tablet to driving scenes, or from other areas to driving scenes)
Traffic density	No or heavy oncoming traffic
Scenario type	Lane-keeping or lane-changing scenarios
TOR lead time	Short (3–4 s) or long (6–7 s) TOR lead time

per second) were treated as a saccade. We categorized area of interests (AOIs) into driving scenes, the NDRT tablet, and other areas. The number and average duration of fixations and saccades were accumulated within the certain AOI. The scan pattern is the probability of eyes switching from one AOI to another. Traffic density, TOR lead time, and scenario type were used to describe the driving environments because they indicated the predictability, criticality, and urgency of the takeover scenarios (Gold et al., 2017). To reduce the potential impact of individual differences, we normalized the feature values across participants using the min-max normalization approach.

We used driving behaviors during takeover transitions to assess drivers' takeover performance. As shown in Table 2, for different takeover scenarios, we selected different metrics in the assessment.

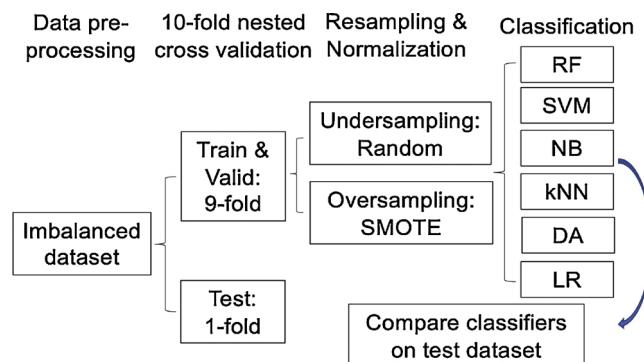


Fig. 5. Modeling process (RF = random forest; SVM = support vector machine; NB = Naive Bayes; kNN = k-nearest neighbors; DA = discriminant analysis; LR = logistic regression).

Table 2

Takeover situations and corresponding driving behavior variables to determine takeover performance (TOR = takeover request; min = minimum; max = maximum; TTC = time to collision).

Takeover reactions	Driving behavior variables (range for bad performance group)		
Lane changing	TOR reaction time ( $>\mu + 2\sigma$ )	Max resulting acceleration ( $>\mu + 2\sigma$ )	$\log(\text{Min TTC})$ ( $<\mu - 2\sigma$ )
Lane keeping	TOR reaction time ( $>\mu + 2\sigma$ )	Max resulting acceleration ( $>\mu + 2\sigma$ )	Standard deviation of road offset ( $>\mu + 2\sigma$ )

Minimum TTC was calculated only for the lane-changing scenarios, and standard deviation of road offset was calculated only for the lane-keeping scenarios. All the driving variables were calculated following prior studies (Du et al., 2020; Clark and Feng, 2017). If any of the calculated TOR reaction time, maximum resulting acceleration, and standard deviation of road offset values were larger than  $\mu + 2\sigma$ , we categorized a takeover transition as a bad performance. For minimum TTC, because the value of  $\mu - 2\sigma$  was negative, we performed a log transformation first and categorized a takeover transition as bad if log (minimum TTC) was lower than  $\mu - 2\sigma$  (Braunagel et al., 2017). For a particular takeover event, as long as one of the driving variables in a certain takeover scenario was categorized as a bad performance, we labeled the scenario as a bad takeover performance. Scenarios that led to collisions were also categorized as bad performances. Eventually, we got an imbalanced dataset with 109 “bad performance” labels and 719 “good performance” labels. The reasons that we used categorical takeover performance rather than individual driving variables as model output were that (1) it combines multiple aspects of driving behaviors and (2) it is easy to be explained to drivers and more practical to guide driver behaviors.

### 3.3. Model development

The takeover performance prediction model was trained with a random forest model considering the following justifications. First, as an ensemble method, random forests are robust for new data generalization and against training data overfitting (Quinlan, 1996). Second, random forests can give us feature importance and makes models interpretable. Five other machine learning approaches mentioned in prior literature were applied for comparisons: k-nearest neighbors (kNN), support vector machine (SVM), Naive Bayes (NB), discriminant analysis (DA), and logistic regression (LR).

Considering the challenge of human behavior data collection, we used a 10-fold nested cross-validation method to train models and compare test results (Lee et al., 2013; Varma and Simon, 2006). As shown in Fig. 5, the 9-fold training and validation set ( $N = 116$  subjects) was used to tune the hyper-parameters with the inner loop and then create classifiers. To handle the imbalanced dataset during the training, we employed a hybrid method of undersampling and oversampling (Choirunnisa and Lianto, 2018). The elimination process was done by deleting 300 good takeover performance scenarios randomly (Prusa et al., 2015). Then we used Synthetic Minority Over-sampling Technique (SMOTE) to create a balanced training and validation dataset with 678 data points (Chawla et al., 2002). Table 3 demonstrates the training procedures of six machine learning approaches. The model assessment was based on the remaining 1-fold testing set ( $N = 13$  subjects) with the outer loop. Notably, the subject data used for testing were not seen in the model training and validation stage. The random selection of 1-fold test dataset assumed that its distribution of good and bad takeover performance scenarios was similar to the whole dataset. With a 10-fold cross-validation, we can make sure all the data points in the dataset would appear once in the test dataset. The training and evaluation of the algorithm were implemented in Matlab 2018b (MathWorks, MA, USA).

### 3.4. Model evaluation

In a binary classification problem, there are four possible outcomes: true positive ( $TP$ ), false positive ( $FP$ ), true negative ( $TN$ ), and false negative ( $FN$ ).  $TP$  is the number of positive samples predicted as a positive class,  $FP$  is the number of negative samples predicted as a positive class,  $FN$  is the number of positive class samples predicted as a negative class and  $TN$  is the number of negative samples predicted as a negative class. In this paper, we used four classification evaluation indicators, including Precision, Recall, Accuracy, and F1-score, to carry out the evaluation of the model performance, which were defined as:

**Table 3**  
Machine learning techniques and training process.

Machine learning approach	Techniques	Hyper-parameters
Support vector machine (SVM)	Embed the data in another dimensional space and find a soft margin that separates the classes with minimum classification error (Chen et al., 2004)	Kernel, regularization parameter
Naive Bayesian (NB)	Use maximum likelihood estimation to estimate parameters (i.e., prior probability and likelihood) (Rish, 2001)	None
Random forest (RF)	Fit an algorithm on a set of bootstrapping samples (bagging) and predictors, i.e., randomly select training samples with replacement and take a random set of predictors at each node without replacement. Repeat many times to form an ensemble of trees (Breiman, 1996, 2001)	Tree number, predictor number per split, leaf size
k-nearest neighbor (kNN)	Calculate Euclidean distance between labeled and unlabeled points to find the k-nearest neighbors. Use the majority vote criteria to decide unlabeled points (Keller et al., 1985)	k
Discriminant analysis (DA)	Find separating hyperplane using parameter estimation (Friedman, 1989)	Discriminant type, Regularization parameter
Logistic regression (LR)	Estimate the parameters of a logistic model (Lee et al., 2006)	Regularization parameter

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (1)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (2)$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} \quad (3)$$

$$F1\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

Precision manifests how well the model predicts (i.e., a measure of exactness) and recall manifests how well the model does not miss the target (i.e., a measure of completeness). The F1 measure is the weighted harmonic mean of the two and represents a realistic measure of model performance.

The receiver operating characteristic (ROC) curve plots the true positive rate (TPR) against the false positive rate (FPR) at different thresholds (i.e., classifier boundary). The area under the curve (AUC) ranges from 0 to 1 and represents the degree of separability. A higher value of AUC indicates better model performance. When AUC is 0.5, it means the model does not have any class separation capability.

## 4. Results

To improve the robustness of machine learning results, we ran the 10-fold cross-validation 30 times (i.e., 30 different random seeds) for every machine learning method at each time window. We first ran an omnibus analysis of variance (ANOVA) to compare the performance of the six machine learning methods. After that, we compared the random forest model with the other five methods using the pairwise *t*-test to see whether the random forest model had the best performance. Similarly, we compared the prediction results of the random forest model with different feature subsets against the full feature model using pairwise *t*-test. We examined the effects of time window and individual feature on random forest prediction performance using ANOVA. All post hoc

comparisons used a Bonferroni  $\alpha$  correction.

#### 4.1. Model performance comparisons

**Figs. 6 and 7** show the average model accuracy and F1-score at different time windows. There was a main effect of machine learning approaches on the prediction accuracy ( $F(5, 5399) = 13, 550, p < .001$ ) and F1-score ( $F(5, 5399) = 4705, p < .001$ ). **Table 4** shows the pairwise *t*-tests comparing the predictive performance of the random forest model with the other five models across different time windows. The results indicate that our proposed random forest model outperformed the other five models across time windows.

**Fig. 8** shows the ROC curves of the random forest and the other five machine learning approaches with the optimal hyper-parameters. The curve of the random forest is above and to the left of the other five curves at the majority of thresholds. Consistent with accuracy and F1-score results, the ROC curve comparisons demonstrated that the random forest model outperformed the other five models.

#### 4.2. Effects of window size on random forest prediction results

There was a main effect of time window on the random forest model accuracy ( $F(29, 899) = 16, p < .001$ ) and F1-score ( $F(29, 899) = 9, p < .001$ ). When applying an algorithm in real-world driving, a time window with shorter size and better prediction performance is preferred. According to **Figs. 6** and **7**, we recommend 3 s as the optimal time window to predict takeover performance, with an average F1-score of 64.0% and accuracy of 84.3% (tuned hyper-parameters: the number of trees = 300; minimum leaf size = 2; the number of predictors per decision split = 6). Post hoc analysis showed that F1-score at the 3 s time window significantly outperformed the rest of the time windows except 5–8, 11–20, and 28–30 s (see **Fig. 7**). Accuracy at the 3 s time window significantly outperformed the rest of the time windows except 4, 6, 11, and 13–16 s (see **Fig. 6**).

#### 4.3. The confusion matrix and feature importance

**Fig. 9** shows the confusion matrix when the time window was 3 s. The precision was 64.5% and the recall was 63.9%, accounting for balanced completeness and exactness of prediction.

Furthermore, by permuting the out-of-bag data (i.e., 36.8% of the total data that were not in the bootstrap samples) randomly across one predictor at a time and by measuring how much this permutation reduced the accuracy of the model, we estimated the feature importance. The values indicate each feature's relative importance in predicting the takeover performance (the larger values are, the more important features are). **Fig. 10** illustrates the out-of-bag estimates of

feature importance of the 37 predictor variables when the time window was 3 s. **Table 5** lists the top 16 important predictor variables. As shown in the table, we found that some heart rate indices and GSR indices (e.g., maximum and mean phasic GSRs, mean of heart rate) were important in predicting takeover performance, but were not included in prior takeover performance algorithm development (Gold et al., 2018; Braunagel et al., 2017).

#### 4.4. Effects of features on random forest prediction results

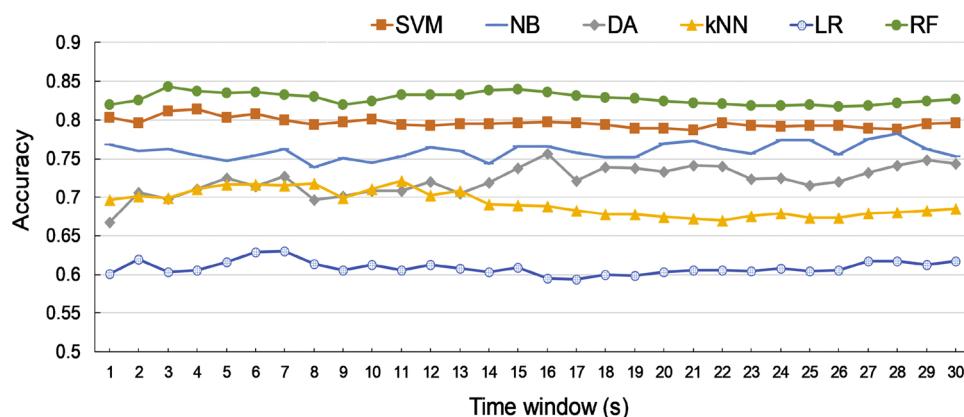
The main effect of feature set on the model accuracy ( $F(3, 119) = 304, p < .001$ ) and the F1-score ( $F(3, 119) = 146, p < .001$ ) were significant at the 3 s time window. We found that the accuracy and F1-score of the random forest model using the full feature set were significantly higher than the accuracy and F1-score using other combinations of feature subsets at the 3 s time window (**Fig. 11** and **Table 6**). To be specific, if only environment factors were used as the features, the average prediction accuracy and F1-score were only .758 and .611, respectively. If only physiological data were used as features, the average prediction accuracy was .770 and F1 score was .563. This suggests that a combination of environment features and features indicating drivers' states are necessary to build a model with high performance. The model using environment factors and eye-tracking metrics as features had an average accuracy of .818 and F1-score of .615 at the 3 s time window. After adding heart rate and galvanic skin response indices as features, the average model accuracy increased to .843 and F1-score increased to .640.

In addition, we ordered features according to the average feature importance values. Next, we built a random forest model with the most important feature, and then added features with lower importance one by one to build another 36 models. As shown in **Fig. 12**, the model accuracy and F1-score generally increased at the beginning when more features were added but reached a plateau when 16 or more features were included in the model. There was a main effect of feature numbers on the model accuracy ( $F(36, 1109) = 3718, p < .001$ ) and F1-score ( $F(36, 1109) = 293, p < .001$ ). Post hoc analysis showed that the F1-score of the full feature model was significantly higher than that for models with fewer than the top 9 important features, and accuracy of the full feature model was significantly higher than that of the models with fewer than the top 16 important features.

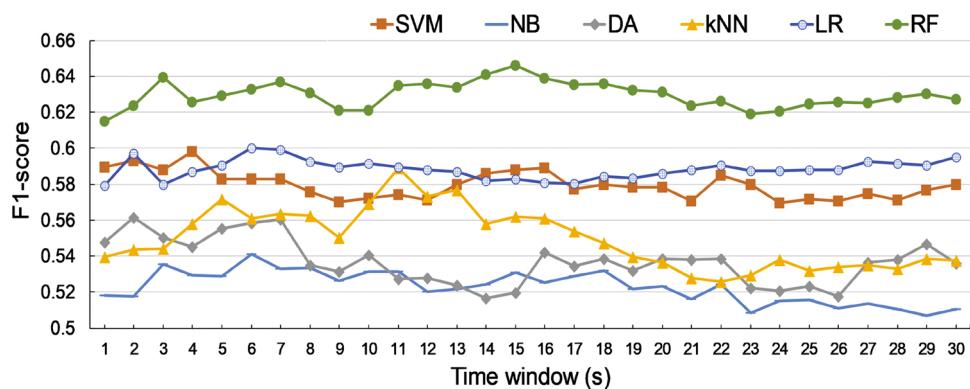
## 5. Discussion

#### 5.1. Model performance comparisons

Our study compared the random forest model with the other five machine learning approaches used in prior literature for takeover



**Fig. 6.** Prediction accuracy of six machine learning approaches under different time windows (SVM = support vector machine; NB = Naive Bayes; DA = discriminant analysis; kNN = k-nearest neighbors; LR = logistic regression; RF = random forest).

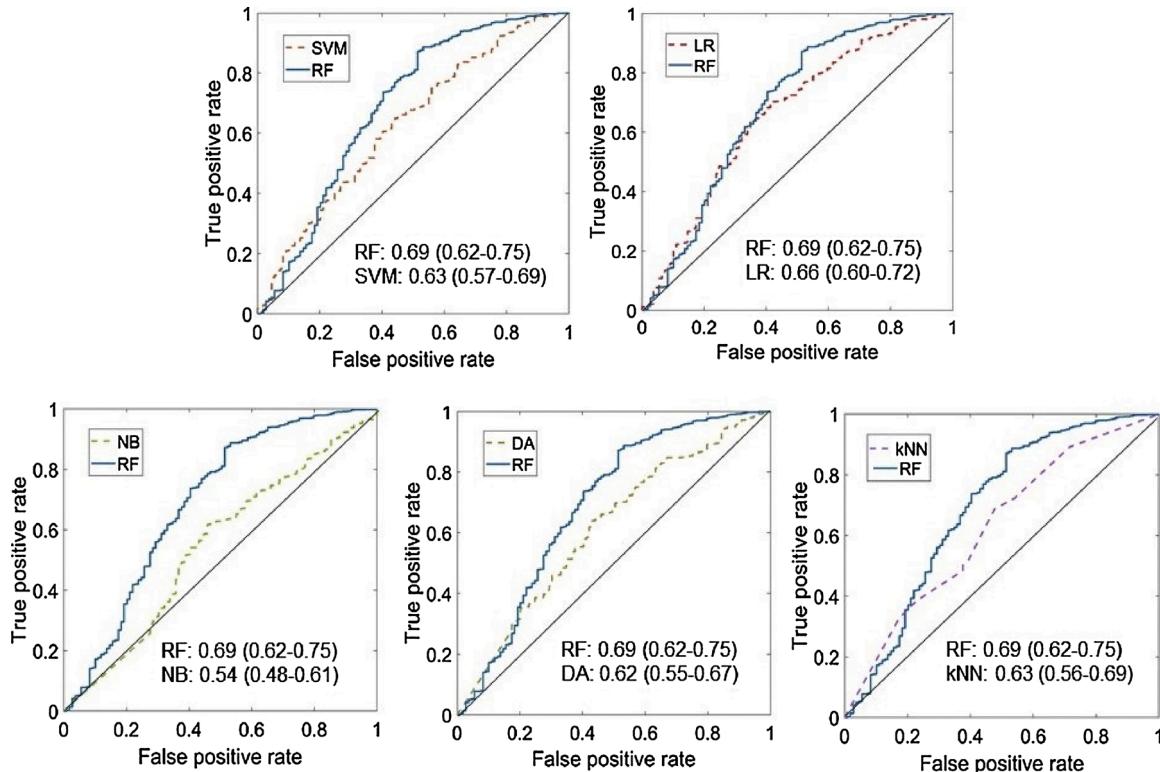


**Fig. 7.** F1 scores of six machine learning approaches under different time windows (SVM = support vector machine; NB = Naive Bayes; DA = discriminant analysis; kNN = k-nearest neighbors; LR = logistic regression; RF = random forest).

**Table 4**

The mean prediction accuracy and F1-score of machine learning approaches across time windows and their comparisons to the random forest model.

Algorithm	Accuracy				F1-score			
	Mean	SD	t-test statistic	p-value	Mean	SD	t-test statistic	p-value
Random forest	.828	.012	–	–	.630	.015	–	–
Support vector machine	.796	.013	60.5	p < .001	.580	.019	72.4	p < .001
Naive Bayes	.760	.033	49.0	p < .001	.523	.022	107	p < .001
Discriminant analysis	.722	.021	134	p < .001	.537	.017	131	p < .001
k-nearest neighbor	.692	.020	209	p < .001	.550	.020	111	p < .001
Logistic regression	.609	.016	342	p < .001	.588	.009	74.5	p < .001



**Fig. 8.** Receiver operating characteristic comparison plots for the random forest (RF) model and the five other models (SVM = support vector machine; LR = logistic regression; NB = Naive Bayes; DA = discriminant analysis; kNN = k-nearest neighbors). The bootstrapped (#1000) confidence intervals are indicated within the parentheses.

performance prediction. As indicated by the results of model accuracy, F1-score, and ROC curve comparisons, the random forest approach outperformed the other classification approaches. Consistent with

previous studies on drivers' fatigue and drowsiness detection (McDonald et al., 2014; Zhou et al., 2020a), the random forest approach also showed its supremacy for takeover performance prediction. It might be

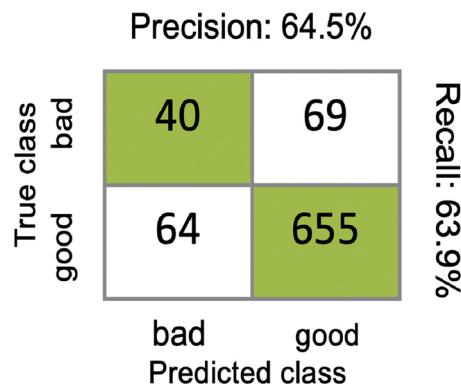


Fig. 9. Confusion matrix when time window was 3 s.

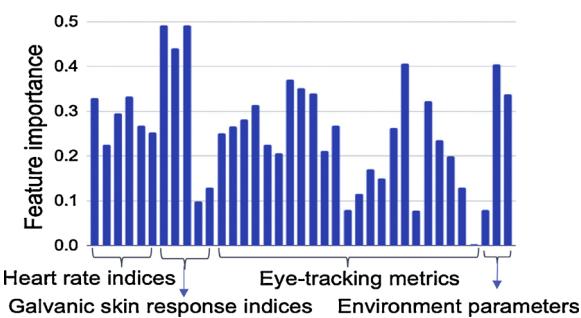


Fig. 10. Feature importance when time window was 3 s.

Table 5

The top 16 important features when time window was 3 s (GSR = galvanic skin response; NDRT = non-driving-related task).

Feature descriptions	Importance
Maximum of GSR in phasic component	.492
Mean of GSR in phasic component	.491
Standard deviation of GSR in phasic component	.441
Vertical gaze dispersion	.406
Scenario type	.404
Fixation duration	.371
Fixation duration on the driving scene	.352
Fixation duration on the NDRT	.341
Takeover lead time	.338
Mean of inter-beat interval	.333
Mean of heart rate	.330
Eyes-on-the-road percentage	.323
Saccade number on the driving scenes	.314
Maximum heart rate	.295
Fixation number on the driving scenes	.282
Standard deviation of inter-beat interval	.268

because random forests aggregate the results of many bootstrap aggregated (bagged) decision trees, which reduces the effects of overfitting and improves generalization.

## 5.2. Effects of window size on random forest prediction results

As the random forest outperformed other machine learning approaches, we examined the prediction performance of random forests under different time window sizes. The results showed that the window size significantly influenced random forest prediction performance. However, such a relationship was not linear. One of the explanations could be that we used a mixture of physiological signals as model inputs. Some physiological signals (e.g., pupil diameter) perform better with a shorter window size because they change rapidly according to the changes in the driver's cognitive workload (Kramer et al., 2013). Some

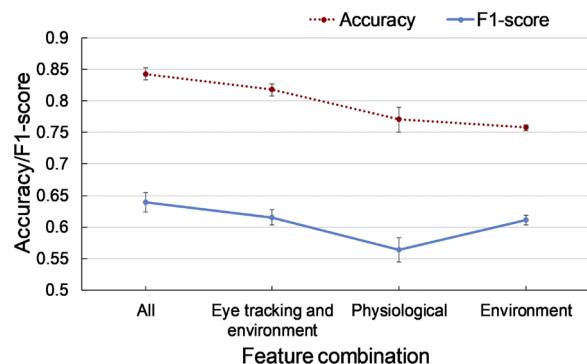


Fig. 11. Prediction accuracy and F1-score of random forests with different feature subsets at the 3 s time window. Error bar indicates 1 standard deviation.

physiological signals (e.g., heart rate) perform better with a longer window size because it can provide an overall understanding of the driver's mental state (Solovey et al., 2014). Future research is needed to explore model performance with customized time windows for different physiological signals.

It was important to find an optimal window size to calculate physiological features for model development in this study. Considering the implementation in real-world driving, a time window with shorter size and better prediction performance is preferred. Thus, we recommend 3 s as the optimal time window to predict takeover performance, with an accuracy of 84.3% and an F1-score of 64.0%. The post hoc analysis showed that the selection of time window for such performance is not unique. Time windows with a size of 6, 11, and 13–16 s led to similar prediction performance. Although the exact time window might be slightly different in the real world given the differences of situational and behavioral parameters, our study provides important insights on window size recommendation for the development of driver state detection systems.

Different from previous studies, our model has a finer granularity and can predict drivers' takeover performance when they are engaged in a specific type of NDRTs with different levels of cognitive load. Such application differences make it infeasible to compare the exact accuracy and F1-score values with those in previous models. Because the test cases in our model prediction are from different participants and are not seen in the training set, our model can be used to predict takeover performance of a new driver who does not have historical data.

## 5.3. Effects of features on random forest prediction results

Drivers' galvanic skin responses, heart rate activities, and eye movements with a combination of environment factors were used to predict drivers' takeover performance. Compared to Braunagel et al. (2017), we added GSR indices and HR indices for model development. Our results showed an improvement of model performance with a full set of features compared to other feature subsets (i.e., physiological data only, environment data only, eye-tracking and environment data). This aligns with the previous studies because all these physiological signals reflected drivers' states and interactions with driving environments (Mehler et al., 2012; Radlmayr et al., 2014; Wang et al., 2014; Ratwani et al., 2010; Young et al., 2013; Bertola and Balk, 2011).

Furthermore, we identified the most important features (e.g., maximum phasic GSR, gaze dispersion, scenario type, and mean of inter-beat interval) for model development. Although the model performance increased at the beginning as more features were added, it reached a plateau when 16 or more features were included. With the top 16 important features, we were able to develop a random forest model with comparable performance to the full feature model. Notably, the top 16 important features were extracted from galvanic skin responses, heart rate activities, eye movements, and environment factors, demonstrating

**Table 6**

Random forest prediction accuracy and F1-score with different feature subsets at the 3 s time window and their comparisons to the full feature model.

Feature subsets	Accuracy				F1-score			
	Mean	SD	t-test statistic	p-value	Mean	SD	t-test statistic	p-value
All	.843	.010	–	–	.640	.015	–	–
Eye-tracking and environment	.818	.010	11.2	$p < .001$	.615	.013	10.9	$p < .001$
Physiological	.770	.020	17.2	$p < .001$	.563	.019	19.7	$p < .001$
Environment	.758	.005	42.7	$p < .001$	.611	.008	8.82	$p < .001$

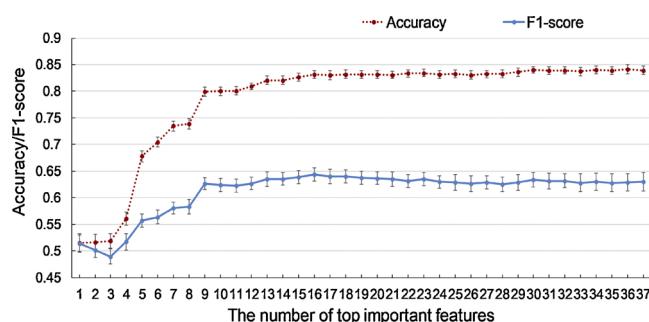


Fig. 12. Model accuracy and F1-score with different numbers of top important features. Error bar indicates 1 standard deviation.

the importance of all these data sources. Utilizing the advances of wearable technology and vehicle sensors, these features can be collected in a minimally invasive manner to predict drivers' takeover performance in real time.

#### 5.4. Limitations and future work

Several limitations should be taken into consideration in the future. First, this study used a snapshot of the time-series data as model inputs without considering the complexity of sequence dependence among the data. Future study could try a convolutional neural network (CNN) combined with long-short-term memory (LSTM) to predict drivers' takeover performance using a larger dataset. Second, the ground truth was determined by drivers' driving behaviors. It is necessary to propose a standard set of metrics for measuring takeover performance. An ensemble method combining subjective ratings, driving behaviors and video coding can be explored to provide a more robust ground truth label of takeover performance. Third, instead of using dichotomous classification of takeover performance, we could increase the number of classes (e.g., bad, neutral, good; or very bad, bad, neutral, good, very good) or use regression to see model prediction power. Fourth, this study only recruited young adult participants with few AV experiences and each participant only experienced four or eight takeover scenarios in the whole experiment. Future studies could recruit participants from different ages, AV experience levels, and training groups. Then the individual characteristics and power law of learning could be taken into account as model inputs to increase the generalization of models (Forster et al., 2019).

#### 5.5. Implications

Our study is a preliminary effort to predict drivers' takeover performance for designing advanced driver monitoring systems. With the advances of technologies in connected automated vehicle systems, real-time road environments such as traffic situations can be accessed easily in the future. Predictive model performance can be improved when data from various drivers engaging in different NDRTs in diverse environments are available for model training. The model outputs can contribute to the design of adaptive in-vehicle alert systems in conditionally automated driving. Specifically, if the system predicted that a

driver would not be able to take over control successfully, a multi-modal display could be designed to help the driver realize the urgency of the event, augment situational awareness and allocate attention properly. Eventually, it could improve drivers' takeover performance and enhance the safety and adoption of automated vehicles.

## 6. Conclusion

This study developed a random forest model to predict drivers' takeover performance in conditionally automated driving. In contrast to previous models capable of predicting drivers' takeover performance when they performed different types of NDRTs, our model has a finer granularity and is able to predict takeover performance when drivers are engaged in a specific type of NDRTs. The results showed that the random forest classifier has an accuracy of 84.3% and an F1-score of 64.0% using a 3 s time window, which outperformed other machine learning models used in prior studies. In addition, we identified the most important physiological measures for takeover performance prediction, and they can be used for developing in-vehicle monitoring systems. Such models can be used to guide the design of adaptive in-vehicle alert systems to improve takeover performance in conditionally automated driving in the future.

## Authors' contribution

Na Du: conceptualization, methodology, software, formal analysis, writing – original draft. Feng Zhou, Dawn M. Tilbury, Lionel P. Robert and Anuj K. Pradhan: conceptualization, methodology, writing – reviewing and editing. Elizabeth Pulver: conceptualization, methodology. X. Jessie Yang: conceptualization, methodology, writing- reviewing and editing, resources, supervision, funding acquisition.

## Declaration of competing interest

none

## Acknowledgements

This work was supported by University of Michigan Mcity and in part by the National Science Foundation. The views expressed are those of the authors and do not reflect the official policy or position of State Farm®.

## References

- Anderson, T.W., 2011. *The Statistical Analysis of Time Series*, Vol. 19. John Wiley & Sons.
- Ayoub, J., Zhou, F., Bao, S., Yang, X.J., 2019. From manual driving to automated driving: a review of 10 years of autoUI. Proceedings of the 11th International Conference on Automotive User Interfaces and Interactive Vehicular Applications (AutomotiveUI'19). ACM, New York, NY, USA, pp. 70–90.
- Bashiri, B., Mann, D.D., 2014. Heart rate variability in response to task automation in agricultural semi-autonomous vehicles. *Ergon. Open J.* 7 (1), 6–12.
- Benedek, M., Kaernbach, C., 2010. A continuous measure of phasic electrodermal activity. *J. Neurosci. Methods* 190 (1), 80–91.
- Bertola, M.A., Balk, S.A., 2011. Eyes on the Road: A Methodology for Analyzing Complex Eye Tracking Data.

- Braunagel, C., Rosenstiel, W., Kasneci, E., 2017. Ready for take-over? A new driver assistance system for an automated classification of driver take-over readiness. *IEEE Intell. Transp. Syst. Mag.* 9 (4), 10–22.
- Breiman, L., 1996. Bagging predictors. *Mach. Learn.* 24 (2), 123–140.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. Smote: synthetic minority over-sampling technique. *J. Artif. Intell. Res.* 16, 321–357.
- Chen, D.R., Wu, Q., Ying, Y., Zhou, D.X., 2004. Support vector machine soft margin classifiers: error analysis. *J. Mach. Learn. Res.* 5 (Sep), 1143–1175.
- Choirunnisa, S., Lianto, J., 2018. Hybrid method of undersampling and oversampling for handling imbalanced data. *2018 International Seminar on Research of Information Technology and Intelligent Systems (ISRITI)*, 276–280.
- Clark, H., Feng, J., 2017. Age differences in the takeover of vehicle control and engagement in non-driving-related activities in simulated driving with conditional automation. *Accid. Anal. Prev.* 106, 468–479.
- Collet, C., Clarion, A., Morel, M., Chapon, A., Petit, C., 2009. Physiological and behavioural changes associated to the management of secondary tasks while driving. *Appl. Ergon.* 40 (6), 1041–1046.
- Dietterich, T.G., 1997. Machine-learning research. *AI Mag.* 18 (4), 97.
- Du, N., Zhou, F., Pulver, E., Tilbury, D.M., Robert, L.P., Pradhan, A.K., Yang, X.J., 2020. Examining the effects of emotional valence and arousal on takeover performance in conditionally automated driving. *Transp. Res. Part C: Emerg. Technol.* 112, 78–87.
- Eriksson, A., Petermeijer, S.M., Zimmermann, M., De Winter, J.C., Bengler, K.J., Stanton, N.A., 2018. Rolling out the red (and green) carpet: supporting driver decision making in automation-to-manual transitions. *IEEE Trans. Hum.-Mach. Syst.* 49 (1), 20–31.
- Eriksson, A., Stanton, N.A., 2017. Takeover time in highly automated vehicles: noncritical transitions to and from manual control. *Hum. Fact.* 59 (4), 689–705.
- Forster, Y., Hergeth, S., Naujoks, F., Beggiato, M., Krems, J.F., Keinath, A., 2019. Learning to use automation: behavioral changes in interaction with automated driving systems. *Transp. Res. Part F: Traff. Psychol. Behav.* 62, 599–614.
- Friedman, J.H., 1989. Regularized discriminant analysis. *J. Am. Stat. Assoc.* 84 (405), 165–175.
- Gold, C., Happée, R., Bengler, K., 2018. Modeling take-over performance in level 3 conditionally automated vehicles. *Accid. Anal. Prev.* 116, 3–13.
- Gold, C., Körber, M., Lechner, D., Bengler, K., 2016. Taking over control from highly automated vehicles in complex traffic situations: the role of traffic density. *Hum. Fact.* 58 (4), 642–652.
- Gold, C., Naujoks, F., Radlmayr, J., Bellem, H., Jarosch, O., 2017. Testing scenarios for human factors research in level 3 automated vehicles. *International Conference on Applied Human Factors and Ergonomics* 551–559.
- ISO, 2020. ISO/TR 21959-1: Road Vehicles – Human Performance and State in the Context of Automated Driving – Part 1: Common Underlying Concepts.
- Jacob, R.J., Karn, K.S., 2003. Eye tracking in human-computer interaction and usability research: READY to deliver the promises. *The Mind's Eye* 573–605.
- Jaeggi, S.M., Buschkuhl, M., Jonides, J., Perrig, W.J., 2008. Improving fluid intelligence with training on working memory. *Proc. Natl. Acad. Sci. U.S.A.* 105 (19), 6829–6833.
- Jones, M., Chapman, P., Bailey, K., 2014. The influence of image valence on visual attention and perception of risk in drivers. *Accid. Anal. Prev.* 73, 296–304.
- Keller, J.M., Gray, Givens, J.A., 1985. A fuzzy k-nearest neighbor algorithm. *IEEE Trans. Syst. Man Cybern.* 4, 580–585.
- Koo, J., Shin, D., Steinert, M., Leifer, L., 2016. Understanding driver responses to voice alerts of autonomous car operations. *Int. J. Veh. Des.* 70 (4), 377–392.
- Körber, M., Gold, C., Lechner, D., Bengler, K., 2016. The influence of age on the take-over of vehicle control in highly automated driving. *Transp. Res. Part F: Traff. Psychol. Behav.* 39, 19–32.
- Kramer, S.E., Lorens, A., Coninx, F., Zekveld, A.A., Piotrowska, A., Skarzynski, H., 2013. Processing load during listening: the influence of task characteristics on the pupil response. *Lang. Cogn. Process.* 28 (4), 426–442.
- Lee, J.J., Knox, B., Baumann, J., Breazeal, C., DeSteno, D., 2013. Computationally modeling interpersonal trust. *Front. Psychol.* 4, 893.
- Lee, S.I., Lee, H., Abbeel, P., Ng, A.Y., 2006. Efficient l1 regularized logistic regression. *Aaaai*, Vol. 6 401–408.
- Li, S., Blythe, P., Guo, W., Namdeo, A., 2018. Investigation of older driver's takeover performance in highly automated vehicles in adverse weather conditions. *IET Intell. Transp. Syst.* 12 (9), 1157–1165.
- Luo, R., Wang, Y., Weng, Y., Paul, V., Brudnak, M.J., Jayakumar, P., Yang, X.J., 2019. Toward real-time assessment of workload: a Bayesian inference approach. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 63 196–200.
- McDonald, A.D., Alambeigi, H., Engström, J., Markkula, G., Vogelpohl, T., Dunne, J., Yuma, N., 2019. Toward computational simulations of behavior during automated driving takeovers: a review of the empirical and modeling literatures. *Hum. Fact.* 61 (4), 642–688.
- McDonald, A.D., Lee, J.D., Schwarz, C., Brown, T.L., 2014. Steering in a random forest ensemble learning for detecting drowsiness-related lane departures. *Hum. Fact.* 56 (5), 986–998.
- Mehler, B., Reimer, B., Coughlin, J.F., 2012. Sensitivity of physiological measures for detecting systematic variations in cognitive demand from a working memory task: an on-road study across three age groups. *Hum. Fact.* 54 (3), 396–412.
- Mehler, B., Reimer, B., Coughlin, J.F., Dusek, J.A., 2009. Impact of incremental increases in cognitive workload on physiological arousal and performance in young adult drivers. *Transp. Res. Rec.* 2138 (1), 6–12.
- Miller, D., Johns, M., Mok, B., Gowda, N., Sirkin, D., Lee, K., Ju, W., 2016. Behavioral measurement of trust in automation: the trust fall. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 60 1849–1853.
- Mok, B., Johns, M., Lee, K.J., Miller, D., Sirkin, D., Ive, P., Ju, W., 2015. Emergency, automation off: Unstructured transition timing for distracted drivers of automated vehicles. *2015 IEEE 18th International Conference on Intelligent Transportation Systems* 2458–2464.
- Molnar, L.J., 2017. Age-Related Differences in Driver Behavior Associated With Automated Vehicles and the Transfer of Control Between Automated And Manual Control: A Simulator Evaluation. University of Michigan, Ann Arbor, Transportation Research Institute.
- Molnar, L.J., Ryan, L.H., Pradhan, A.K., Eby, D.W., Louis, R.M.S., Zakrajsek, J.S., 2018. Understanding trust and acceptance of automated vehicles: An exploratory simulator study of transfer of control between automated and manual driving. *Transp. Res. Part F: Traff. Psychol. Behav.* 58, 319–328.
- Naujoks, F., Mai, C., Neukum, A., 2014. The effect of urgency of take-over requests during highly automated driving under distraction conditions. *Adv. Hum. Aspects Transp.* 7 (Part I), 431.
- Petersen, L., Robert, L., Yang, J., Tilbury, D., 2019. Situational awareness, driver's trust in automated driving systems and secondary task performance. *SAE Int. J. Connect. Auton. Veh.* 2 (2) <https://doi.org/10.4271/12-02-02-0009>.
- Prusa, J., Khoshgoftaar, T.M., Dittman, D.J., Napolitano, A., 2015. Using random undersampling to alleviate class imbalance on tweet sentiment data. *2015 IEEE International Conference on Information Reuse and Integration* 197–202.
- Quinlan, J.R., 1996. Bagging, Boosting, and c4.5. *Aaaai/iaai*, Vol. 1 725–730.
- Radlmayr, J., Gold, C., Lorenz, L., Farid, M., Bengler, K., 2014. How traffic situations and non-driving related tasks affect the take-over quality in highly automated driving. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 58 2063–2067.
- Ratwani, R.M., McCurry, J.M., Trafton, J.G., 2010. Single operator, multiple robots: an eye movement based theoretic model of operator situation awareness. *2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 235–242.
- Rezvani, T., Driggs-Campbell, K., Sadigh, D., Sastry, S.S., Seshia, S.A., Bajcsy, R., 2016. Towards trustworthy automation: user interfaces that convey internal and external awareness. *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, 682–688.
- Rish, I., 2001. An empirical study of the naive bayes classifier. *IJCAI 2001 Workshop on Empirical Methods in Artificial Intelligence*, Vol. 3 41–46.
- Schmidt, E., Decke, R., Rasshofer, R., 2016. Correlation between subjective driver state measures and psychophysiological and vehicular data in simulated driving. *2016 IEEE Intelligent Vehicles Symposium (IV)* 1380–1385.
- Society of Automotive Engineers, 2018. Taxonomy and Definitions for Terms Related to On-Road Motor Vehicle Automated Driving Systems.
- Solovey, E.T., Zec, M., Garcia Perez, E.A., Reimer, B., Mehler, B., 2014. Classifying driver workload using physiological and driving performance data: two field studies. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* 4057–4066.
- van den Beukel, A.P., van der Voort, M.C., 2013. The influence of time-criticality on situation awareness when retrieving human control after automated driving. *16th International IEEE Conference on Intelligent Transportation Systems (ITSC) 2013*, 2000–2005.
- Varma, S., Simon, R., 2006. Bias in error estimation when using cross-validation for model selection. *BMC Bioinformatics* 7 (1), 91.
- Wan, J., Wu, C., 2018. The effects of lead time of take-over request and nondriving tasks on taking-over control of automated vehicles. *IEEE Trans. Hum.-Mach. Syst.* 99, 1–10.
- Wandtner, B., Schömg, N., Schmidt, G., 2018. Effects of non-driving related task modalities on takeover performance in highly automated driving. *Hum. Fact.* 60 (6), 870–881.
- Wang, Y., Reimer, B., Dobres, J., Mehler, B., 2014. The sensitivity of different methodologies for characterizing drivers' gaze concentration under increased cognitive demand. *Transp. Res. Part F: Traff. Psychol. Behav.* 26, 227–237.
- Wintersberger, P., Riener, A., Schartmüller, C., Frison, A.K., Weigl, K., 2018. Let me finish before i take over: towards attention aware device integration in highly automated vehicles. *Proceedings of the 10th INTERNATIONAL CONFERENCE on Automotive User Interfaces and Interactive Vehicular Applications* 53–65.
- Young, K.L., Salmon, P.M., Cornelissen, M., 2013. Missing links? The effects of distraction on driver situation awareness. *Saf. Sci.* 56, 36–43.
- Zeeb, K., Buchner, A., Schrauf, M., 2016. Is take-over time all that matters? The impact of visual-cognitive load on driver take-over quality after conditionally automated driving. *Accid. Anal. Prev.* 92, 230–239.
- Zeeb, K., Härtel, M., Buchner, A., Schrauf, M., 2017. Why is steering not the same as braking?: the impact of non-driving related tasks on lateral and longitudinal driver interventions during conditionally automated driving. *Transp. Res. Part F: Traff. Psychol. Behav.* 50, 65–79.
- Zhou, F., Alsaid, A., Blommer, M., Curry, R., Swaminathan, R., Kochhar, D., Lei, B., 2020a. Driver fatigue transition prediction in highly automated driving using physiological features. *Expert Syst. Appl.* 113204.
- Zhou, F., Yang, X.J., Zhang, X., 2020b. Takeover transition in autonomous vehicles: a youtube study. *Int. J. Hum. Comput. Interact.* 36 (3), 295–306.