# Driving Anomaly Detection Using Conditional Generative Adversarial Network

Yuning Qiu, *Student Member, IEEE,* Teruhisa Misu, *Member, IEEE,* and Carlos Busso, *Senior Member, IEEE*

*Abstract*—**Anomaly driving detection is an important problem in *advanced driver assistance systems* (ADAS). It is important to identify potential hazard scenarios as early as possible to avoid potential accidents. This study proposes an unsupervised method to quantify driving anomalies using a conditional *generative adversarial network* (GAN). The approach predicts upcoming driving scenarios by conditioning the models on the previously observed signals. The system uses the difference of the output from the discriminator between the predicted and actual signals as a metric to quantify the anomaly degree of a driving segment. We take a driver-centric approach, considering physiological signals from the driver and *controller area network*-Bus (CAN-Bus) signals from the vehicle. The approach is implemented with *convolutional neural networks* (CNNs) to extract discriminative feature representations, and with *long short-term memory* (LSTM) cells to capture temporal information. The study is implemented and evaluated with the *driving anomaly dataset* (DAD), which includes 250 hours of naturalistic recordings manually annotated with driving events. The experimental results reveal that recordings annotated with events that are likely to be anomalous, such as avoiding on-road pedestrians and traffic rule violations, have higher anomaly scores than recordings without any event annotation. The results are validated with perceptual evaluations, where annotators are asked to assess the risk and familiarity of the videos detected with high anomaly scores. The results indicate that the driving segments with higher anomaly scores are more risky and less regularly seen on the road than other driving segments, validating the proposed unsupervised approach.**

*Index Terms*—**Driving anomaly detection, conditional generative adversarial networks, convolutional neural networks, long short-term memory cell.**

## I. Introduction

**W**ITH the development of the smart automobile industry in recent years, more and more functions have been added to *advanced driver assistance systems* (ADAS), avoiding human errors and increasing road safety. Examples include *lane departure warning* (LDW), *forward collision warning* (FCW), and *intelligent speed advice* (ISA). These techniques share the common basic principle of detecting hazard scenarios, warning drivers of potential risks, and taking control of the vehicle in extreme situations. All these solutions require detection of driving anomalies that deviate from normal driving patterns, and increase chances of accidents. Current approaches for detecting abnormal driving behaviors or conditions often rely on either threshold-based or rule-based systems [1]–[5]. However, these methods are often

Y. Qiu and C. Busso are with the Department of Electrical Engineering, University of Texas at Dallas, Richardson, Texas, 75080 USA e-mail: {yxq180000, busso}@utdallas.edu

Teruhisa Misu is with the Honda Research Institute, Mountain View, California, US e-mail: TMisu@hra.com

triggered only when a driver makes a mistake, which is too late in many cases. Furthermore, it is highly unlikely that rule-based systems can exhaustively cover all potential anomaly scenarios. It is important to develop algorithms that can detect general abnormal driving behaviors as early as possible so that potential road accidents can be avoided. For these cases, unsupervised approaches are expected to be more effective, without rigid predefined rules or definitions of anomalies.

This work proposes an unsupervised approach based on the conditional *generative adversarial network* (GAN) framework to detect driving anomalies. Our approach is based on the premise that driving anomalies are often unexpected events that we cannot predict with the available contextual information. Motivated by this premise, we use the data from previous frames as a condition to generate prediction of signals of the near future from random noise. Then, we use the discriminator model to compare the differences between the generated prediction and real data of the upcoming frames, creating a powerful metric to indicate the abnormal degree of the near future. This idea was previously validated in our preliminary work [6], where we implemented our approach with fully connected *deep neural networks* (DNNs). In contrast, our proposed framework uses the latest advances in machine learning to leverage discriminative information directly from the data using *convolutional neural networks* (CNNs). CNNs can extract both linear and non-linear relationships in and between sequences [7], which we expect to be useful in detecting driving anomalies. Our formulation also leverages temporal information, relying on *recurrent neural networks* (RNNs) implemented with *long short-term memory* (LSTM) layers. LSTMs are designed to capture temporal relationship among sequential data from previous frames, providing a powerful framework for temporal sequence forecasting [8].

We build our approach with features extracted from the *controller area network-bus* (CAN-Bus) and physiological signals, although this approach is flexible and can be implemented with different sensing modalities. CAN-Bus signals provide powerful information for estimating driving maneuvers and driver behaviors, including acceleration, breaks and steering wheel movements [9]–[11]. Therefore, we expect that predicting future CAN-Bus signals with our formulation will lead to robust driving anomaly detection. In addition to CAN-Bus data, we also rely on the driver's physiological signals. In particular, we consider *Electrocardiography* (ECG), *breath rate* (BR), and *Electrodermal activity* (EDA) signals. The motivation for considering physiological signals is that under certain complex driving conditions, a driver might get nervous or frightened by abnormal driving events, which will

be reflected in her/his biosignals. Even if the driver does not react with a driving maneuver, the physiological signals will indicate the presence of the anomaly. In fact, our preliminary study showed that adding features extracted from the driver's physiological data increased the model's discriminative power [6]. Physiological signals are also closely related to driving behaviors [12]–[15]. Consequently, our model considers the vehicle's CAN-Bus signals and driver's physiological signals.

We evaluate the proposed approach with the *driving anomaly dataset* (DAD). This corpus has rich manual annotations of maneuvers and events. We group events that are likely to trigger driving anomalies, such as traffic violations, pedestrian on the road, and crossing vehicles. The anomaly scores of the driving segments overlapping with these annotations are generally higher than the anomaly scores of the driving segments without any annotations. We also investigate the contributions of the blocks in our formulation. For this purpose, we build our model with fully connected DNN, with only LSTMs, or with only CNNs. The result shows that the LSTM-based model performs better at discriminating abnormal from normal driving scenarios than the CNN-based model. The discriminative performance of our proposed approach is improved when we combine both structures (i.e., CNN and LSTM). These results are also confirmed with perceptual evaluations, where annotators were asked to assess the risk level and familiarity of video segments identified with high anomalous scores by the alternative models. The proposed CNN+LSTM based model is able to identify video segments that are perceived as more risky and less familiar by the annotators. In summary, the contributions of our study are:

- We proposed an unsupervised formulation to predict driving anomalies based on conditional GAN, which contrasts predicted and observed physiological and CAN-Bus features.
- The proposed approach derives discriminative representations directly from data using CNNs, and leverages temporal information across consecutive frames using LSTMs.
- We validated the proposed approach with objective and perceptual evaluations using a naturalistic driving database, which demonstrates the strengths of our proposed formulation.

The paper is organized as follows. Section II presents related studies, discussing effort to detect driving anomalies. It also discusses how conditional GANs have been used for anomaly/outlier detection in other fields. Section III introduces the dataset used in this study for evaluating our proposed models. Section IV presents the motivation of our framework, discussing the details of our formulation. Section V evaluates the discriminative power of our unsupervised driving anomaly detection framework using objective and subjective evaluations. Finally, Section VI summarizes the contributions of this work, discussing potential ideas to extend and improve our unsupervised driving anomaly detection system.

## II. RELATED WORK

This section reviews the most relevant studies to our work. We start with Section II-A, which presents studies aiming to detect driving anomalies. Section II-B presents a broader overview on anomaly detection methods using GAN across different areas. Section II-C highlights the differences between our work and other methods.

### A. Driving Anomaly Detection

Pattern-based methods detect anomaly driving events by modeling driving maneuver patterns. Some of them identify cases where the driving behaviors depart from expected normal driving patterns, labeling them as abnormal events [5], [16]–[21]. Other studies have focused on detecting several specific types of abnormal driving maneuver patterns [3], [4], [22]–[29]. Zhang et al. [5] proposed a driving anomaly detection model by representing normal driving patterns in a state graph. They use this state graph as the criterion to distinguish abnormal driving behaviors that deviate from the expected state transitions. Another representative approach is the work of Chen et al. [3]. They considered six types of abnormal driving behaviors such as fast U-turn and sudden braking. They obtained the vehicle acceleration data from smartphone sensors, which were used to recognize these events using a *support vector machine* (SVM). Dai et al. [4] detected driving under the influence of alcohol using a pattern-matching model, which compares the differences in the vehicle's acceleration between normal and drunk driving conditions.

The threshold based methods [1], [2], [30]–[36] set bounds on the values of features or parameters describing the driving scenario. Abnormal driving conditions are set when their values are outside the predefined *safe* ranges. Hong et al. [1] detected when the vehicle's acceleration was higher than a predefined threshold, using this event as a proxy to measure aggressive driving behaviors. Similarly, Chakravarty et al. [2] proposed a system that evaluates multiple thresholds on the vehicle's acceleration to detect risky driving events. They considered the vehicle's acceleration on different directions to detect four maneuver types (i.e., hard bump, hard cornering, harsh brake, and sharp acceleration). Even though these methods are simple and computationally effective, the threshold-based methods using predefined values often lack flexibility, requiring, in many cases, domain knowledge (e.g., speed limit on current roads).

The clustering based method is another approach used to identify abnormal scenarios [37]–[41]. This approach builds on the hypothesis that most people drive in a proper and safe way under most naturalistic driving conditions. Therefore, abnormal driving behaviors or risky driving conditions are infrequent events. Under this assumption, we should expect that anomaly events will be clustered as outliers. Hansen et al. [39] discriminated the driver's maneuvers by mapping the features extracted from the vehicle's dynamic signals to a feature space. The outlier driving events of the clusters were regarded as driving anomalies. The work of Zheng and Hansen [38] used a one-class SVM and the *topology anomaly detection* (TAD), clustering model to grade each driving event from "good" to "bad". The study established a four-class labeling framework, according to the clusters (from the innermost cluster to the outermost cluster).

## B. Anomaly Detection Using Conditional GANs

Multiple methods have been proposed for time series anomaly detection [42]. An interesting formulation is the *generative adversarial network* (GAN) [43], which has opened new directions for this problem. A GAN-based model consists of a generative model ($G$) and a discriminative model ($D$), which are trained with an adversary strategy. $G$ is trained to generate disruptive fake data from random noises, learning the distribution of the data that needs to be generated. As an adversarial game, $D$ is trained to identify the differences between the generated fake data and real data. As the quality of $G$ improves, the differences between the generated and real signals decrease, reducing the performance of $D$. GAN has been widely used as a state-of-the-art generative model. This model has also been used for detection of anomalies or outliers. This section describe some examples of GAN-based anomaly detection models used in different domains. The reader are referred to Di Mattia et al. [44], which presents a survey on this area.

Schlegl et al. [45] presented the AnoGAN framework, which was applied to identify anomalies in retina tomography images. AnoGAN learns a manifold to represent the distribution of the data using a GAN. It maps the image into a latent space where it is feasible to quantify deviations from normal distributions, detecting anomalous cases. Zhou et al. [46] proposed BeatGAN, an unsupervised GANs-based system to identify unusual human motions (e.g., jumping and running) from normal motions (i.e., walking). This approach built a generator with an encoder-decoder structure, using the reconstructed signals as the fake signals to confuse the discriminator. For the evaluation, they used the reconstruction error between the real signal and the generated fake signal as the anomaly metric to detect abnormal motions. Xue et al. [47] proposed a supervised approach based on GAN, called SegAN, for brain tumor segmentation. The generator creates segmentation masks for the input brain image. The discriminator identifies between generated segmentations and ground truth labels. The formulation regards the tumor area as the abnormal part of a given image. Li et al. [48] proposed a GAN-based model to detect attacks on *cyber-physical systems* (CPSs). The approach compares the predictions of the generator with actual multivariate time-series data. The residual between these signals is used to detect anomaly activities in the CPS. Studies have also used GAN for *out-of-domain* (OOD) detection in *natural language understanding* (NLU). For example, Zheng et al. [49] trained an autoencoder to map the input utterance into the latent embedding created by the generator of a GAN model. Then, they used the decoder of the autoencoder to generate utterance from the embedding as OOD samples.

Rather than generating data merely from random noise, Mirza et al. [50] modified the original GAN formulation by adding extra information as condition to the model. The additional input conditions the data generation process, creating behaviors that are properly constrained. Hyland et al. [51] adopted the conditional GAN framework implemented with LSTMs to generate fake patients' *heart rate* (HR) and *respiratory rate* (RR) data, conditioning on blood pressure values. They used this method to detect patients' abnormal physical conditions. Akcay et al. [52] proposed the GANomaly framework, which is another conditional adversary network for anomaly detection. The approach combines GAN with an autoencoder to jointly model a latent and image space. The encoder processes the input data creating a latent space. The decoder processes the latent space to create a reconstructed version of the data. The discriminator aims to classify the original image as real and the reconstructed image as fake. Then, the encoder is used again to map the reconstructed image back into the latent space. The distance in the latent space between the input data and reconstructed latent vector is used as the anomaly score. The model is trained with neutral data. Anomalous examples that do not fit the normal distribution are expected to have larger distance. A similar approach was used by Zenati et al. [53] in their *efficient GAN based anomaly detection* (EGBAD) framework, without using the encoder on the reconstructed image.

## C. Relation to Prior Work

Our proposed approach builds upon our preliminary work in Qiu et al. [6]. This study extracted statistic features from the data to capture the key aspects of the signals as the input of the model. The approach has two problems that are addressed in this study. First, the model used hand crafted features from the modalities. Our proposed architecture addresses this problem by using the CNN block, which extracts discriminative representations directly from the data. The second limitation is that the model only considers static windows of six seconds to constraint the model. This approach ignores temporal information within these previous six second, since statistics are derived from the entire segment. It also ignores longer dependencies that may be important in the prediction of the signals in the upcoming frames. The proposed approach addresses this limitation by integrating the LSTM module. With the combination of both additions, we extract discriminative temporal representations directly from the raw data. We believe that this feature representation can reveal more detailed information than statistic features using predefined functional module. Qiu et al. [54] used an architecture that was similar to the work proposed in Qiu et al. [6], but with an anomaly score relying on the triplet loss function. We do not explore this direction in this paper. We provide an exhaustive evaluation of our proposed architecture, comparing the benefits of adding the CNN and LSTM blocks. We also compares the proposed approach with several alternative methods.

While other studies have proposed GAN-based approaches for anomaly detection, the particular formulation presented in this study is novel, and has important benefits with respect to other alternative models. Our model learns from real data and make forecasts of the upcoming signals based on the observed data. We use the discriminator to determine segments with signals that deviate from the observed patterns. Our approach is fully unsupervised, where we only need to collect data without the need for labels. Most of the other GAN-based methods are supervised (e.g., SegAN [47]). GAN-based anomaly detection methods, such as AnoGAN [45] and

TABLE I
ANNOTATIONS INCLUDED IN THE DAD DATABASE.

| | Annotations |
|---|---|
| Goal-oriented Operation | Left turn; Right turn; Intersection passing; Crosswalk passing; Left lane change; Right lane change; U-turn |
| Stimulus-driven Operation | Stop for congestion; Avoid pedestrian near ego lane; Avoid road motorcyclist; Avoid on-road bicyclist |
| Traffic rule/manner violation | Traffic rule violation |
| Attention | Crossing vehicle; Crossing pedestrian; Red light; Cut-in; Sign; On-road bicyclist; Parked vehicle; Merging vehicle; Yellow light; Road work; Pedestrian near ego lane |

TABLE II
SETS OF ANNOTATIONS CONSIDERED TO EVALUATE THE PROPOSED
ANOMALY DETECTION MODELS. THE CANDIDATE SET IS EXPECTED TO
HAVE DRIVING SCENARIOS THAT CAN BE CONSIDERED AS ANOMALOUS.

| Sets | Annotations |
|---|---|
| Candidate | Avoid on-road pedestrian; Avoid pedestrian near ego-lane; Avoid on-road bicyclist; Avoid bicyclist near ego-lane; Avoid on-road motorcyclist; Avoid parked vehicle; traffic rule violation |
| Maneuver | Left turn; Right turn; Left lane branch; Right lane branch; U-turn; Intersection passing |
| Normal | No annotations during the segments |

GANomaly [52] are trained merely with normal data, and tested with normal/abnormal data. They detect anomalies by discriminating the samples that are different from the normal ones, calling these samples as *abnormal*. Our approach is fundamentally different from these formulations, where we train the model with all the data using a predictive formulation. It is also important to notice that all these studies using GAN-based approaches for anomaly detection were designed for other problems in other fields. With the exception of our preliminary study [6], [54], this is the first time that conditional GAN has been used for driving anomaly detection.

## III. DRIVING ANOMALY DATASET (DAD)

This work uses the *driving anomaly dataset* (DAD). The corpus includes 250 hours of naturalistic driving recordings in an Asian city collected by the Honda Research Institute in collaboration with a local company. One experienced driver participated in the data collection process driving a Honda Accord. This dataset includes the driver's physiological data, which is a key modality in this study. The driver's *electrocardiography* (ECG) and *breath rate* (BR) signals are recorded with a Zephyer BioHarness 3 chestband. The ECG signal is recorded at 250 Hz, and the BR signal is recorded at 25 Hz. The data collection also included the *Electrodermal activity* (EDA) signals from the driver obtained with a Empatica E4 wristband collected at 4 Hz. To compensate for the differences across sessions, we normalize the physiology data per session by using the Z-normalization. The corpus also provides the vehicle's CAN-Bus data. We obtain six vehicle signals: speed, steering speed, steering angle, throttle angle, brake pressure, and yaw. These signals are recorded at 100 Hz. All the vehicle's CAN-Bus data and driver's physiological data are synchronized at 30 Hz.

The data collection also includes data from other sensors not used in this study. The setup includes three FLIR Blackfly S cameras facing the road [(i.e., right, center and left)]. The data collection also included Tobii Pro 2 eye-tracking glasses.

One of the strengths of the DAD corpus is the rich set of annotations, which follows the approach used in the collection of the *Honda Research Institute driving dataset* (HDD) [55], [56]. The annotations are grouped into four layers: goal-oriented operation, stimuli-driven operation, traffic rule/manner violation, and driver's attention. Table I shows the specific annotations

within each layer (see study of Ramanishka et al. [56] for details on this annotation process). The annotations of the driver's driving maneuvers are manually added to the dataset. The annotations and the sensing modalities are aligned, including the drivers' physiological signals and the vehicles' CAN-Bus data. For clearer visualization, all the data sequences, annotations, and videos are combined and synchronized using the open source software ELAN. Figure 1 shows the *user interface* (UI) of ELAN, where the driving events correspond to the *stimulus-driven operation* layer.

While the corpus includes 250 hours of data, not all the recordings have been annotated. We only use 121 sessions, which include about 130 hours of well-annotated urban driving recordings. We split these recordings into train (100 sessions, ~105 hours), development (11 sessions, ~13 hours), and test (10 sessions, ~12 hours) sets. The DAD corpus does not have annotations for anomaly scores. Instead, we evaluate our model using the existing annotations overlapping with the driving videos used to test our framework. For this purpose, we group the driving events in the test set into three groups. The first group is the *candidate* set, which consists of traffic rule violations and hazard driving conditions such as avoiding on-road pedestrians or parked vehicles. We expect that these events will include segments with high driving anomaly scores. The second group is the *maneuver* set, which includes segments annotated with regular driving maneuvers such as right turns, intersection passing, and U-turns. In general, we expect moderate anomaly scores for these events. The third group is the *normal* set, which includes segments without any annotation. We expect low anomaly scores for these segments. Table II lists the events associated with each group. Notice that our approach is fully unsupervised, so the annotations are exclusively used to evaluate our approach.

## IV. PROPOSED CONDITIONAL GAN MODEL

The goal of this study is to implement an unsupervised framework to detect driving anomalies. A driving anomaly is something unexpected that deviates from normal patterns. We are not just interested on detecting dangerous conditions. Instead, this work focuses on unexpected driving events. An ADAS should be able to leverage knowledge of unexpected events, even if they do not represent dangerous scenarios. Dangerous events are special cases, since, in daily life, some dangerous driving scenarios are usually caused by unexpected maneuvers or reactions from traffic participants (e.g., drivers
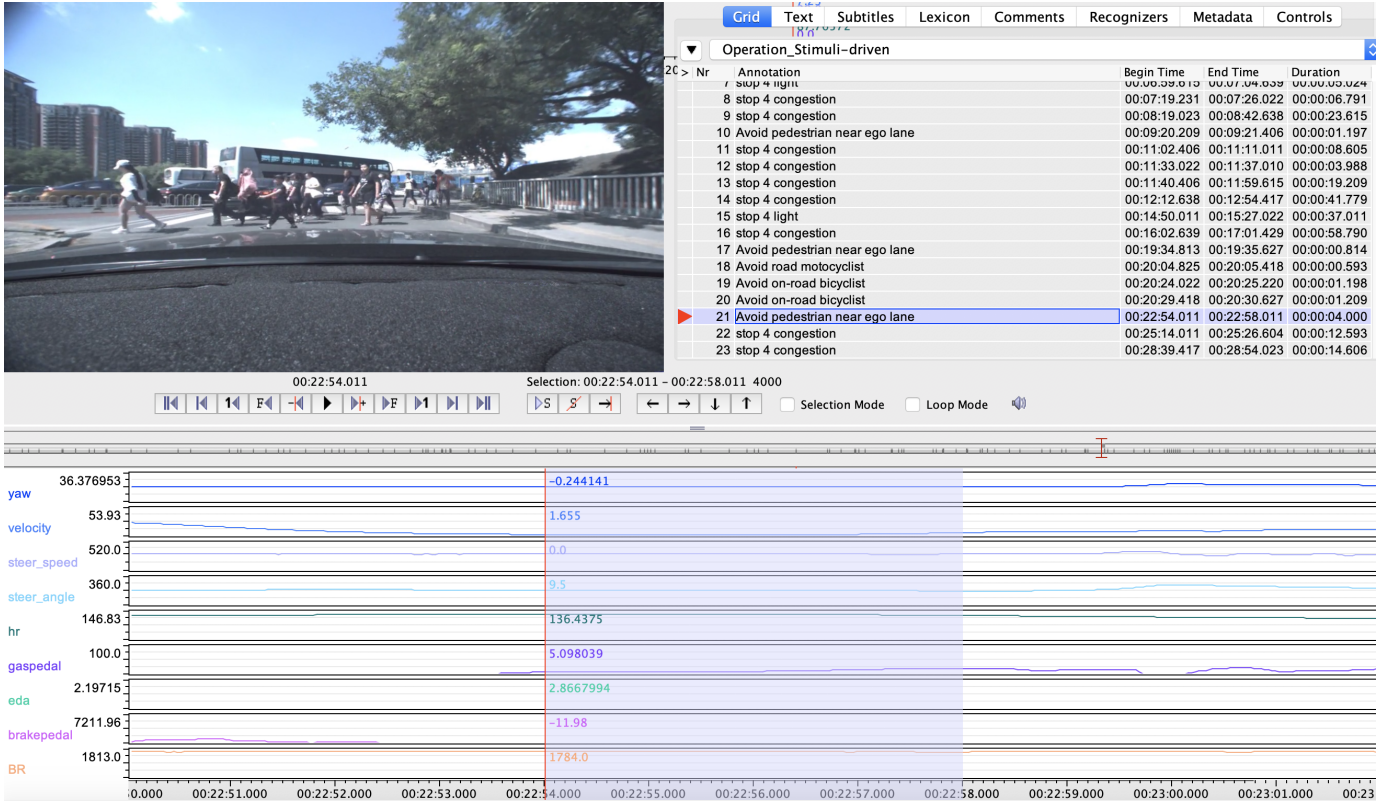
Fig. 1. The software ELAN showing some of the annotations included in the DAD database. The annotations, physiological data, and CAN-Bus data are synchronized with the road videos.

of other vehicles or the ego-vehicle, pedestrians, bicyclists and motorcycle). Therefore, we expect that risky and hazardous scenarios will be considered by our system as driving anomalies. While the framework is general and can be implemented with different features, we detect abnormal driving behaviors using the vehicle CAN-Bus data and the driver's physiological data. The CAN-Bus data consists of the vehicle's speed, yaw, pedal angle, brake pressure, steer angle, and steering speed. The physiological data includes HR, BR, and EDA signals. Our motivation for using physiological signals is that HR, BR, and EDA respond to mental and cognitive states, indicating stress [57], and anxiety [58] levels. Previous studies on driver behaviors have used physiological data [59], [60], showing correlation with driving maneuvers [12]–[15]. These findings show that physiological signals can be natural complements to CAN-Bus signals, providing information even when a driver fails to maneuver the car in the presence of unexpected events.

The premise of our method is that anomaly scenarios are often associated with unexpected events. Therefore, we predict the features from upcoming driving events, conditioned on the previous values of the target features. Then, we quantify the difference between the predicted driving data and the actual data. We implement these ideas with conditional GAN, which is one of the most powerful generative models.

Figure 2 shows the procedure of our implementation. The generator creates plausible data sequences conditioned on the previous values of the features. The discriminator recognizes whether the input data is real or fake (i.e., created by the



Fig. 2. Abstract illustration of our GAN-based model for driving anomaly detection. Conditioned by a contextual window with previous frames, the generator predicts the features in the near future. The discriminator takes the predictions and the real data as inputs comparing the value of the discriminator's score. Unexpected events are then identified with this unsupervised model.

generator). The scores from the discriminator are used to determine our anomaly scores. This section describes our proposed approach in detail. First, we present the intuition of our formulation (Sec. IV-A). Then, we present a basic implementation with fully connected layers (Sec. IV-B), where we explain the main features of our formulation. Sections IV-C and IV-D introduce the use of CNN to extract discriminative features directly from data, and LSTM to capture temporal information. Finally, Section IV-E presents our full model, which combines the architectures of the LSTM and CNN based models.

## A. Anomaly Detection with Conditional GAN

Multiple studies have revealed the impressive capability of GANs to learn the distribution of target data, generating similar samples from random noise. This learning procedure is accomplished through an adversarial game between the discriminator and generator, as shown in Equations 1 and 2.

$$
\begin{aligned}
\max_D V(D) =& \mathbb{E}_{\boldsymbol{x}\sim p_{data}(\boldsymbol{x})}\left[\log D(\boldsymbol{x})\right] \\
& + \mathbb{E}_{\boldsymbol{z}\sim p_z(\boldsymbol{z})}\left[\log(1 - D(G(\boldsymbol{z})))\right]
\end{aligned}
\tag{1}
$$

$$
\min_D V(D) = \mathbb{E}_{\boldsymbol{z}\sim p_z(\boldsymbol{z})}\left[\log(1 - D(G(\boldsymbol{z})))\right]
\tag{2}
$$

The discriminator outputs a value $D(\boldsymbol{x})$ which indicates whether the input $\boldsymbol{x}$, with probability distribution $p_{data}$, is a real sample. The objective of $D(\boldsymbol{x})$ is to maximize the chance to identify the real sample as real, and the generated fake samples as fake. The discriminative score is a sigmoid output which ranges from 0 to 1, where 1 means absolutely real and 0 means absolutely fake. The objective function of the generator aims to create the prediction $G(\boldsymbol{z})$ as close as possible to real samples to fool the discriminator. The variable $\boldsymbol{z} \sim p_z$ is a noise vector used as the input. Following Equations 1 and 2, the GAN model should be trained to converge to a good estimator of $p_{data}$. While the input of a regular GAN is the random noise vector $\boldsymbol{z}$, a conditional GAN uses extra information ($\boldsymbol{y}$) as additional input to constrain the model to generate more targeted predictions. In our study, the condition of the generative model is the data from previous frames. Figure 2 describes the training process of the proposed conditional GANs, showing that the input of $G$ is random noise and the data sequence from previous frames as condition. The output is the generated data sequence of the upcoming analysis window, $G(\boldsymbol{z}|\boldsymbol{y})$.

Figure 2 shows the inference process of the proposed model. Given the analysis window, the generator creates a fake data sequence. Then, $D$ takes either the real data from the upcoming analysis window or the fake data generated by $G$. For each of these inputs, $D$ creates a score ($S_R$ for real signal; $S_F$ for fake/generated signal). The difference between the scores is regarded as the anomaly score, $m_{anomaly}$ (Eq. 3). This metric represents the uncertainty of the upcoming driving events, which we hypothesize to be an informative driving anomalous metric. For a well-trained generative model, $G$ generates realistic data from noise in order to confuse $D$. Therefore, when the real data from the analysis window follows the distribution of the regular data, the generated signal will be similar to the real data and $S_F$ will be closer to $S_R$. This case will produce a small value for $m_{anomaly}$, indicating that the input of $G$ is normal and predictable. However, if the upcoming driving data in the analysis window has a distribution that differs from the expected behaviors, the data will be unpredictable, creating a gap between $S_F$ and $S_G$. This scenario will have a larger $m_{anomaly}$ value. A key advantage of our approach is that it only requires unlabeled data, providing an appealing unsupervised formulation.

$$
m_{anomaly} = |S_R - S_F|
\tag{3}
$$



(a) Generator of the CNN-based model



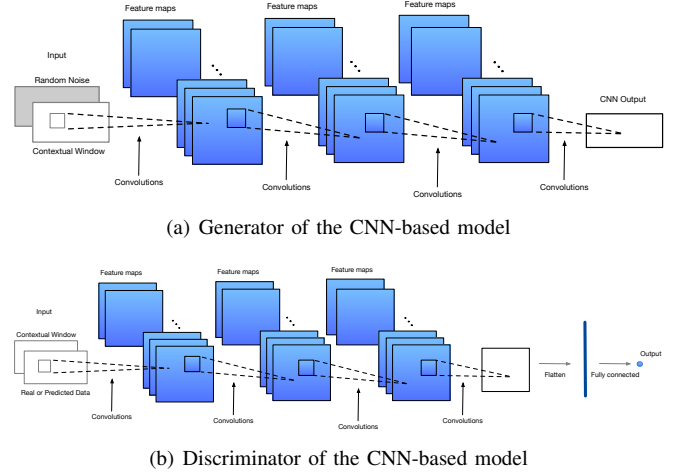(b) Discriminator of the CNN-based model

Fig. 3. Implementation of proposed CNN-based GAN model. The generator and discriminator have similar architectures with four CNNs, extracting discriminative representations directly from physiological and CAN-Bus data.

## B. Conditional GAN with Fully Connected Layers

In Qiu et al. [6], we presented a preliminary implementation of our unsupervised driving anomaly detection, where the discriminator and generator were implemented with *fully connected* (FC) layers. The approach uses a fixed window with previous values of physiological and CAN-Bus data to predict future values of the data. The discriminator takes the statistical features of either the real or generated signals as input. The feature set includes four time domain features from each of the CAN-Bus data and physiological data (i.e., maximum, minimum, mean, and standard deviation). From each physiological data, we additionally extracted five frequency domain features, calculating the energy in the frequency domain covering the following five bands: [0-0.04 Hz], [0.04-0.15 Hz], [0.15-0.5 Hz], [0.5-4 Hz], and [4-20 Hz]. The generator is implemented with five layers, each of them implemented with 180 neurons. Similarly, the discriminator is implemented with five layers, each of them with 51 neurons.

## C. CNN-based Conditional GAN

The first improvement for our model is replacing the feature extraction module. Instead of using predefined functionals over the previous frames, we learn directly discriminative patterns from the data using CNNs. Models designed based on CNNs have been successful for end-to-end classification tasks [61]. The study of Borovykh et al. [7] showed that CNNs are able to extract temporal features from time series data that were discriminative to predict upcoming data values. Inspired by these studies, we implement our conditional GAN models with CNNs to learn more discriminative features directly from data.

Figure 3 shows the implementation of our CNN-based GAN model. The generator and the discriminator consist of four convolutional layers, implemented with 18, 18, 9, and 1 channels, respectively. The kernel size for each layer is 9, 3, 3, and 3, respectively. We add a fully connected layer after the convolutional layers. During the training process, $D$ and $G$ are trained for 20 epochs and the Adam learning rate is

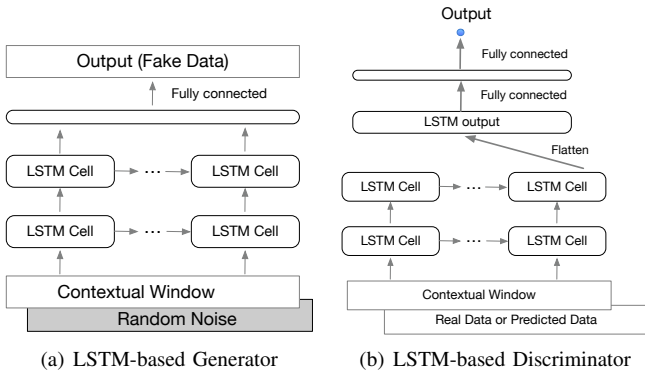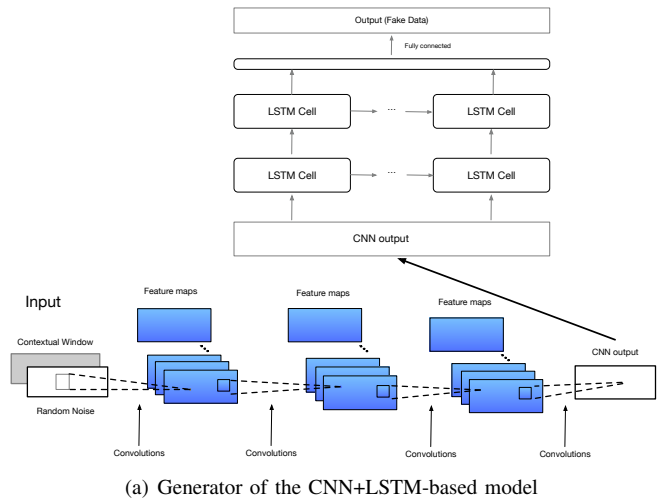(a) LSTM-based Generator     (b) LSTM-based Discriminator

Fig. 4. Implementation of the LSTM-based GAN model. The model relies on longer contextual window to leverage short and longer dependencies in the data to detect driving anomalies.

set to 0.001. From our previous work [15], we conclude that features extracted from CAN-Bus and driver's physiological data with an analysis window of 12 seconds can be used to classify driving maneuvers. We aim to reduce the size of the window to reduce the latency in the model. However, the analysis window cannot be reduced too much to capture changes on physiological signal. As a compromise, we set the analysis window size to six seconds for our CNN-based GAN model. The model takes as input the previous six-second data and random noise, and generates predictions for the upcoming six-second data describing the upcoming CAN-Bus and physiological signals.

### D. LSTM-based Conditional GAN

A limitation of the CNN-based GAN model is that the contextual information is limited to fixed size windows (i.e., six seconds). We hypothesize that modeling longer temporal relationships in the data can lead to better results. We explore this direction by using RNNs, which can extract temporal relationships in the data. We implement the RNNs with *long short-term memory* (LSTM) cells, leveraging the longer history information from previous time series data. Sequence to sequence tasks can be effectively implemented with RNN-based model, as demonstrated in language translation [62]. Additionally, models based on conditional RNNs have been successfully applied on tasks that mimic a particular writer's handwriting style [63]. Borovykh et al. [7] demonstrated the reliable capability of LSTM to make short-term prediction on trends in the stock market. Motivated by these studies, we design our approach with LSTM cells, expecting to improve the temporal modeling of our framework while avoiding rigid contextual feature vectors.

Figure 4 shows the structure of our LSTM-based GAN model. The model consists of two layers of LSTM cells, each of them implemented with 27 nodes. The generator of the LSTM-based GAN model takes a longer contextual window than the window analysis in the CNN-based GAN model, relying on the last 60 seconds of data. The generator predicts the next six seconds of physiological and CAN-Bus signals. The extended contextual window allows the LSTM



(a) Generator of the CNN+LSTM-based model



(b) Discriminator of the CNN+LSTM-based model

Fig. 5. Implementation of the proposed CNN+LSTM-based GAN model. The unsupervised approach combines the strengths in using CNNs and LSTMs, extracting discriminative representations directly from the data while leveraging temporal information.

cells to more effectively leverage short and long temporal relationships to make the predictions. The structure of the discriminator is similar to the generator, as shown in Figure 4(b). The discriminator takes six-second data, which can be real signals or data predicted by the generator, conditioned on the contextual window with the previous 60-second sequence. We extract the last output of the LSTM cells, flattening the feature representation as a vector. We add a fully connected layer creating a one-dimensional output, which predicts if the data is real or predicted by the generator. We use this score to estimate the anomaly score in Equation 3.

### E. CNN+LSTM-Based Approach using Conditional GAN

The CNN-based and LSTM-based GAN models offer complementary benefits for our task. Therefore, our final model combines their structures leveraging better feature representations and temporal modeling. Figure 5 shows the implementation of the proposed CNN+LSTM-based conditional

GAN model, where we use the same structures presented for the CNN-based GAN model (Sec. IV-C) and the LSTM-based GAN model (Sec. IV-D) as blocks to build this model. First, the CNN block of the generator extracts discriminative information from the random noise and contextual window with the previous sixty-second used to condition the models. We implement the CNN block by splitting the contextual window into 10 six-second segments, without overlap. Then, we concatenate the CNN output of each of the six-second segments, creating a conditional embedding, which is used as the input of the LSTM block. The output of the LSTM is the prediction of the physiological and CAN-Bus data for the next 6-second window. For the discriminator, we implement the model using the same structure used for the generator. The only difference is the output layer, which is a one-dimensional score to predict whether the data is real or fake.

During the training process, we first import the pre-trained parameters of the CNN-based and LSTM-based GAN models. Then, we train the LSTM parameters (including both the generator and the discriminator) for 10 epochs while freezing the CNN parameters. Then, we jointly train the entire model together for another 10 epochs to get the final model.

## V. EXPERIMENTAL RESULTS

This section describes the experimental results obtained with our proposed conditional GAN models. Figure 6(a) shows the losses of the generator and discriminator in the training set. Training a GAN is not always easy, since it is a minmax optimization process. It is a problem if the loss of the discriminator drops fast compared to the generator's loss. A strong discriminator means that it is easy to distinguish between fake and real samples. Without an appropriate feedback from the discriminator, it is hard to train the generator. When properly trained, the overall loss of the GAN is often constantly fluctuating, as both losses go down. This is the exact pattern observed in Figure 6(a). For a properly trained GAN, the discriminator should have problems recognizing between real and fake samples. The classification performance should be 50% if the number of real or generated samples are equal. Figure 6(b) shows the probability of the output of the discriminator for real and generated samples on the development set (1 is real, 0 is fake). The performance oscillates around 50% for both type of samples, as expected. These figures shows that the GAN is properly trained.

We evaluate the performance by comparing the scores obtained from videos in the *candidate*, *maneuver*, and *normal* sets. We also evaluate the performance with perceptual evaluations.

### A. Distribution of Anomaly Scores

As stated in Section III, driving events in the *candidate* set are more hazardous and rarely seen than the events in the *normal* set. Therefore, the anomaly scores ($m_{anomaly}$) of the driving events in the *candidate* set are expected to be larger than the corresponding scores on the *normal* set. The analysis in this study compares the distribution of anomaly scores of the driving events from these two sets. For each segment, we



(a) Generator and discriminator losses in the training set



(b) Discriminator performance on the development set

Fig. 6. Loss of the GAN network during training. The training losses of *D* and *G* are often constantly fluctuating, as both losses go down. The outputs of the discriminator for the fake and real signals oscillate around 50%.
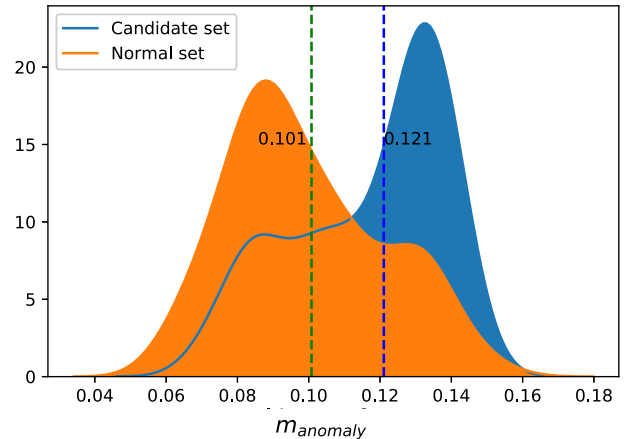


Fig. 7. Histogram of the predictions of the FC-based models for segments in the test set from the normal and candidate sets. The vertical dashed lines indicate the medians of the anomaly scores for the sets.

provide the previous data as condition using six seconds for the FC-based and CNN-based models, and 60 seconds for models implemented with LSTMs. All models generate predictions for the upcoming six seconds.

Figure 7 shows the histograms of the anomaly scores for normal and candidate sets for the FC-based GAN model. The figure shows that the segments from the candidate set have higher anomaly scores than the segments from the normal set. The figure shows clear modes in the histograms showing good separation. Similar results are observed when using the CNN-based GAN model (Fig. 8(c)), LSTM-based GAN model
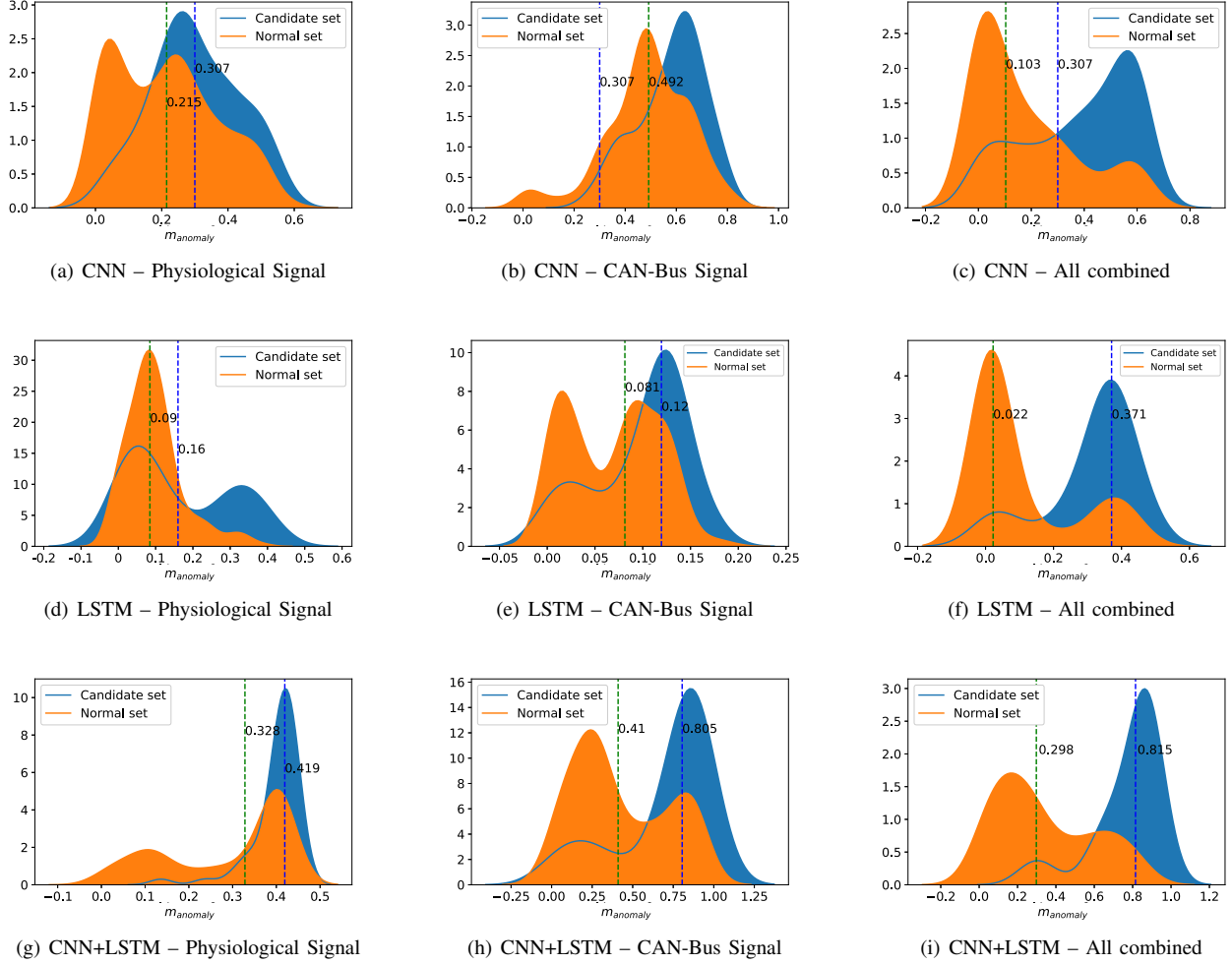
Fig. 8. Histogram of the predictions of the CNN-based, LSTM-based and CNN+LSTM-based models for segments in the test set from the normal and candidate sets obtained from the DAD corpus. The vertical dashed lines indicate the medians of the anomaly scores for the sets. The results are presented when the models are trained with (1) physiological data, (2) CAN-Bus data, (3), and physiological and CAN-Bus data.

(Fig. 8(f)), and CNN+LSTM-based GAN model (Fig. 8(i)). In particular, the results for the CNN+LSTM-based GAN model show clear separations, showing the strengths in combining the CNN-based and LSTM-based GAN models, leading to better discriminative performance. We will directly compare all these methods in Section V-B.

As described in Section IV, our models are trained with physiological and CAN-Bus signals. We evaluate the contributions of each of these modalities by retraining the models with either physiological or CAN-Bus features. Figure 8 shows the histograms for these models. When we only use either physiological or CAN-Bus signals, the differences in the distribution of $m_{anomaly}$ between the normal and candidate sets is clearly reduced, indicating that both modalities provide complementary information. Figures 8(b) and 8(d) are two examples where the overlaps between these distributions are quite clear. Adding both modalities leads to more separation between the distributions, especially for the LSTM-based (Fig. 8(f)) and the CNN+LSTM-based (Fig. 8(i)) GAN models. We measure the medians of the distributions to quantify the differences, which are included as vertical dashed lines

in the distributions in Figure 8. When the CNN+LSTM-based GAN model is trained with only the physiological signals, the difference between the distributions' medians is $\Delta_{Phy} = 0.419 - 0.328 = 0.091$ (Fig. 8(g)). Similarly, when using only CAN-Bus data, the difference between the medians is $\Delta_{CAN} = 0.805 - 0.41 = 0.395$ (Fig. 8(h)). In contrast, the difference between the medians increases to $\Delta_{Both} = 0.815 - 0.298 = 0.517$ when using both modalities (Fig. 8(i)).

Figure 9 shows some of the abnormal driving scenarios which are discriminated with higher anomaly scores. Most of the anomalies are caused by sudden appearance of pedestrians or improper maneuvers from other vehicles.

## B. Direct Comparison of Proposed Models Using DET Curve

We directly compare the models by formulating the evaluation as a binary classification problem (i.e., normal versus candidate sets), where we estimate the *detection error tradeoff* (DET) curves by changing the threshold on the anomaly score. Samples with a score higher than the threshold are classified as abnormal (i.e., part of the candidate set), and those with a score

(a) Example 1        (b) Example 2        (c) Example 3

Fig. 9. Example of frames from segments with high anomaly scores by the unsupervised CNN+LSTM-based model. (a) A motorcycle suddenly cuts in front of the ego-vehicle, (b) the ego-vehicle is trying to avoid a vehicle parked on the roadside, while a bicyclist is coming in the opposite direction, and (c) a pedestrian suddenly crosses the road, and the driver has to press the brakes to avoid hitting him.
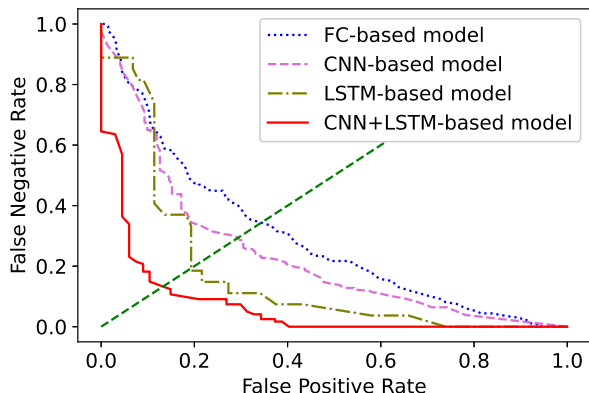


Fig. 10. The DET curves for the models by formulating the problem as a binary classification task (candidate versus normal sets). The DET curve shows the false positive rate as a function of the false negative rate. The diagonal dashed line indicates the *equal error rate* (EER) when both error rates are the same. Generally, the LSTM-based model has better discriminative performance than the CNN-based model and the FC-based model. The best performance is achieved with the CNN+LSTM-based GAN model.

TABLE III
AUC AND EER RATES FOR THE BASELINE MODELS AND THE PROPOSED
CNN+LSTM-BASED MODEL WHEN THEY ARE TRAINED WITH
PHYSIOLOGICAL AND CAN-BUS SIGNALS.

| Approach | Physiol. | | CAN-Bus | | Both | |
|---|---|---|---|---|---|---|
| | AUC | EER | AUC | EER | AUC | EER |
| CNN-based model | 0.329 | 37.6% | 0.330 | 34.1% | 0.235 | 29.1% |
| LSTM-based model | 0.492 | 52.5% | 0.314 | 32.0% | 0.167 | 20.5% |
| CNN+LSTM-based model | **0.237** | **33.4%** | **0.176** | **20.8%** | **0.106** | **13.4%** |

the best performance, leveraging the strengths of using CNNs and LSTMs. The CNN+LSTM based GAN model achieves the lowest AUC and EER rates.

The DET figures can also be used to compare the results when the models are only trained with either physiological, or CAN-Bus features. Figure 11 shows the DET curves for the CNN-based, LSTM-based, and CNN+LSTM based GAN models trained with partial modalities. When we use only physiological data to train the models, we consistently observe lower performance than models only trained with CAN-Bus features. The differences are clearly seen in Table III for the AUC and EER rates. However, when we combine physiological and CAN-Bus data, the discriminative performance of the models is improved. These results reveal the significant role of the drivers' physiological data in the performance of our models. This result is also confirmed by observations on the videos with high anomaly scores detected by all the GAN models, when trained with both feature sets. Figure 13 shows a case to illustrate this point, where the driver is slowing down as the vehicle approaches a T-road. All of the sudden, a motorcycle rider rushes into the lane in front of the vehicle. While the driver does not react with any driving maneuver, the driver's breath rate immediately drops, followed by increases in heart rate and skin conductivity. These changes result in a high anomaly score for this driving segment.
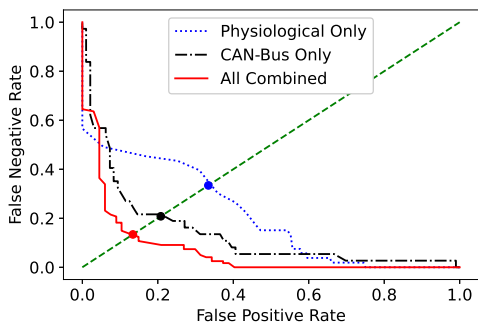
*C. Comparison with Other Baselines*

We evaluate our approach with four representative baselines inspired by approaches used by previous driving anomaly detection methods. The first baseline corresponds to the *fixed-threshold* method, which is inspired by the work of Li et al. [34]. They detected abnormal driving behaviors by setting thresholds on the vehicle's speed ($VS$), yaw angle ($YA$),

below the threshold are classified as normal. The threshold is not fixed. Instead, we increase its value creating different *operation points*. The DET curves show the *false negative rate* (FNR) in the y-axis and the *false positive rate* (FPR) in the x-axis as the threshold that determines the two classes is moved. The performance of a binary classifier is better when the DET curve lies closer to the axes. The diagonal line in the DET curves indicates *equal error rate* (EER), where the FNR and FPR have the same value. In addition to the DET curves, we also report the EER and the *area under the curve* (AUC) to quantify the results. The AUC is estimated over the DET curves so lower values indicate better performance.

Figures 10 compares the DET curves of the FC-based, CNN-based, LSTM-based, and CNN+LSTM based GAN models. Table III shows the corresponding EER and AUC values. We observe gains in performance over the FC based model by extracting the features directly from the data using CNNs (CNN-based approach), and by modeling temporal information (LSTM-based approach). The performance of the LSTM-based model outperforms the performance of the CNN-based model. This observation is clearly observed in Table III. The proposed CNN+LSTM based GAN model achieves

(a) CNN-based GAN model



(b) LSTM-based GAN model



(c) CNN+LSTM-based GAN model

Fig. 11. The DET curves when the models are trained with (1) physiological data, (2) CAN-Bus data, (3), and physiological and CAN-Bus data. The figure shows that both modalities are important for driving anomaly detection.

acceleration ($AC$) and steering speed ($SS$) data. They defined abnormal speeding and dangerous steering behavior events as:

$$\text{Abnormal speeding} = |AC| > 0.8m/s^2 \ \& \ SS < 0.1rad/s$$
$$\text{Steering} = VS > 30km/h \ \& SS > 0.4rad/s \ \& YA > 0.7rad$$

An abnormal driving behavior is defined when either of these two conditions are satisfied. We need to vary the thresholds to compare the performance in DET curves. Our implementation starts with the same threshold as Li et al. [34] for each variable. We consistently move the thresholds across variables by adding or subtracting a fix percentage of their respective range. (i.e., $\hat{t}_i = t_i + \alpha r_i$, where $t_i$ is the original threshold for variable $i$, $r_i$ is the range of variable $i$, and $\alpha$ is an

adjustable parameter to find different tradeoffs between FNR and FPR). The second baseline is the *PCA-threshold* method, which is inspired on the framework proposed by Sadjadi and Hansen [64] for *speech activity detection* (SAD). The idea of this unsupervised SAD approach is to combine multiple indicators into a single metric over which we can apply a threshold. This metric corresponds to the first principal component obtained using *principal component analysis* (PCA). PCA determines the eigenvectors of the covariance matrix of the multidimensional data, which provides the principal directions where the data is spread. The high dimensional feature vectors are linearly mapped into a low dimensional feature space represented by the eigenvector with the highest eigenvalues. We follow the approach presented by Sadjadi and Hansen [64] that maps the entire multidimensional data into a single principal dimension. For each 12-second window, we extract the 51-dimensional feature vector from the CAN-Bus and physiological signals discussed in Section IV-B. The vector is mapped into a 1 dimensional metric using this PCA-based approach. We estimate the DET curve by moving the threshold on the resulting 1-dimensional signal, where the segments with value above the threshold are considered abnormal. We estimate the FPR and FNR rates for different operating points. The third baseline corresponds to the *GMM-threshold* method, which aims to represent the distribution of the data, defining outliers as anomalous events. The *Gaussian mixture model* (GMM) is a common algorithm to fit the distribution of the multidimensional data. A segment can be considered as an outlier if the estimated distribution does not represent well a segment in the test set. Similar to the second baseline, we extracted the 51-dimensional feature vector for each 12-second window, estimating the parameters of a *Gaussian mixture model* (GMM). The number of clusters is set to eight by minimizing the value of the *Akaike information criterion* (AIC) and the *Bayesian information criterion* (BIC). We use the same partitions introduced in Section III to train and test the GMM. We calculate the posterior probability of the feature vectors $x$ from the test set according to $p(x) = \sum_{i=1}^{8} \omega_i \mathcal{N}(x|\mu_i, \sigma_i))$, where $\omega_i$ (weights), $\mu_i$ (mean vector), and $\Sigma_i$ (covariance matrix) are the parameters of the GMM. Segments with posterior probability lower than a threshold are considered as outliers (i.e., anomalous events). The forth baseline is the BeatGAN framework, introduced in Section II-B. We build the BeatGAN model following the description in Zhou et al. [46], setting up the generator in an encoder-decoder structure, using a *multilayer perceptron* (MLP) structure. The numbers of nodes per layer for the generator are 1620-256-128-32-10-10-32-128-256-1620. The numbers of nodes per layer for the discriminator are 1620-256-128-32-1. We calculate the reconstruction error between the real and reconstructed signals as the anomaly score.

Figure 12 shows the DET curves comparing the baseline methods with our CNN+LSTM based model. Table IV shows the corresponding EER and AUC results. The figure shows that our proposed model leads to clear improvements over the baseline models. The results show that the BeatGAN framework is the best baseline model. However, our proposed approach clearly outperforms this method.

TABLE VI
KRIPPENDROFF'S ALPHA COEFFICIENTS AMONG DIFFERENT GROUPS OF
RATERS. EACH GROUP HAS THREE RATERS WHO ASSESSED THE SAME
VIDEOS. THE QUESTIONS CORRESPOND TO THE QUESTIONNAIRE
PRESENTED IN FIGURE 13.

|  | Group 1 | Group 2 | Group 3 | Over All |
|---|---|---|---|---|
| Question 1 | 0.745 | 0.711 | 0.706 | 0.736 |
| Question 2 | 0.512 | 0.532 | 0.653 | 0.573 |



Fig. 13. Graphical user interface for the perceptual evaluation of driving anomaly. After watching each video, the raters are asked to assess the level of risk and familiarity in the segment.



(a) How risky is the driving condition in the video?



(b) How often do you see similar driving condition on the roads?

Fig. 14. Results of the perceptual evaluation to assess the degree of risk and familiarity of the selected videos. The figures show the results for the top 100 segments with the highest anomaly scores selected by three different models. The figure also shows the results for 100 segments randomly selected.

methods. The differences in the selected sets are reflected in the differences observed in Figure 14. A video is annotated by three raters. Therefore, for each model, we have 300 annotations assigned by the raters. Then, we determine the proportion of the annotations assigned to each of the options listed in the GUI (Figure 13), providing the results in Figure 14. For example, the number of annotations assigned to each of the options for the question "How risky is the driving condition in the video?" for the CNN+LSTM model are: *very risky* 16 (5.3%), *risky* 83 (27.7%), *slightly risky* 92 (30.7%), and *safe* 109 (36.3%). Figure 14 shows that the videos selected by the CNN+LSTM based GAN model contain more segments annotated with higher risk and lower familiarity. The figure shows that 33% of these videos are considered as *risky* or *very risky*, and 29.3% of them are considered to occur *never* or *almost never*. The corresponding percentages for the CNN-based and LSTM-based models are lower. These comparisons illustrate that the CNN+LSTM-based model can identify more hazardous and more abnormal driving conditions than the CNN-based model and the LSTM-based model. These numbers are significant, since 75% of the randomly selected segments are considered safe, and 69.3% of them are considered to occur *regularly*, showing that most of the driving conditions in the DAD corpus are regular scenarios without driving anomalies. Our best unsupervised conditional GAN model is able to identify segments which are often perceived with a level of risk (64.7%), which do not occur so often on the road. The figure also shows the superior performance of the LSTM-based model compared to the CNN-based models, validating the results observed in previous sections with the DAD annotations.

## VI. CONCLUSIONS

This study proposed an unsupervised driving anomaly detection system based on a conditional GAN. The proposed approach makes predictions of the driver's physiological data and the vehicle CAN-Bus data, conditioning the model on previous observed signals. The predictions are contrasted with actual data, creating an anomaly score that increases its value when unexpected data is observed. The approach obtains discriminative features from the physiological and CAN-Bus signals directly from the data using CNNs. The model also leverages temporal information by using LSTM networks.

This study shows that the driving events with more hazardous driving conditions usually receive higher anomaly scores by the proposed model. This result is validated with objective evaluations, relying on the annotations of the DAD corpus, and perceptual evaluations conducted on the videos selected by our models with the highest anomaly scores. Our proposed approach is able to effectively detect anomaly driving conditions that deviate from the predictions of the upcoming driving behaviors, creating an appealing unsupervised solution that does not depend on either predefined thresholds or supervised rules.

One limitation of this study is that our detection algorithm depends on actions or reactions from the driver. Anomalies can be detected only when the driver notices events and reacts to them. Therefore, if the driver is unaware of a driving anomaly, a model trained with physiological and CAN-Bus signals will not provide discriminative information to detect it. A potential solution is adding other features that can objectively capture the driver's environment regardless of her/his awareness (e.g., pedestrian detection, car detection). We plan to augment our proposed model with further information, such as the results from vision-based object detection systems. Another limitation of our approach is the use of wearable devices to capture physiological signals. Future research directions include developing remote approaches to measure the driver's physiological signals [65], [66]. While today this is a challenging problem, there are technological advances to create non-contact measurement systems to monitor physiological data that suggest that this could be reasonable in the future. Sensors can be installed on the driver's seat to record the driver's BR and HR signals. Physiological signals can be alternatively obtained from wearable sensors that the drivers may be already using (e.g., smartwatch). Development in this area will make our approach more suitable for deployment in real-driving conditions.

## REFERENCES

[1] J. Hong, B. Margines, and A. K. Dey, "A smartphone-based sensing platform to model aggressive driving behaviors," in *SIGCHI Conference on Human Factors in Computing Systems*, Toronto, ON, Canada, April-May 2014, pp. 4047–4056.

[2] T. Chakravarty, A. Ghose, C. Bhaumik, and A. Chowdhury, "MobiDriveScore - a system for mobile sensor based driving analysis: A risk assessment model for improving one's driving," in *International Conference on Sensing Technology (ICST 2013)*, Wellington, New Zealand, December 2013, pp. 338–344.

[3] Z. Chen, J. Yu, J. Zhu, Y. Chen, and M. Li, "D3: Abnormal driving behaviors detection and identification using smartphone sensors," in *IEEE International Conference on Sensing, Communication, and Networking (SECON 2015)*, Seattle, WA, USA, June 2015, pp. 524–532.

[4] J. Dai, J. Teng, X. Bai, Z. Shen, and D. Xuan, "Mobile phone based drunk driving detection," in *International Conference on Pervasive Computing Technologies for Healthcare*, Munich, Germany, March 2010, pp. 1–8.

[5] M. Zhang, C. Chen, T. Wo, T. Xie, M. Bhuiyan, and X. Lin, "SafeDrive: Online driving anomaly detection from large-scale vehicle data," *IEEE Transactions on Industrial Informatics*, vol. 13, no. 4, pp. 2087–2096, August 2017.

[6] Y. Qiu, T. Misu, and C. Busso, "Driving anomaly detection with conditional generative adversarial network using physiological and can-bus data," in *ACM International Conference on Multimodal Interaction (ICMI 2019)*, Suzhou, Jiangsu, China, October 2019, pp. 164–173.

[7] A. Borovykh, S. Bohte, and C. Oosterlee, "Conditional time series forecasting with convolutional neural networks," *ArXiv e-prints (arXiv:1703.04691)*, pp. 1–8, March 2017.

[8] T. Fischer and C. Krauss, "Deep learning with long short-term memory networks for financial market predictions," *European Journal of Operational Research*, vol. 270, no. 2, pp. 654–669, October 2018.

[9] Y. Zheng, A. Sathyanarayana, and J. H. L. Hansen, "Threshold based decision-tree for automatic driving maneuver recognition using CAN-Bus signal," in *IEEE Conference on Intelligent Transportation Systems (ITSC 2014)*, Qingdao, China, October 2014, pp. 2834–2839.

[10] A. Sathyanarayana, P. Boyraz, Z. Purohit, R. Lubag, and J. Hansen, "Driver adaptive and context aware active safety systems using CAN-bus signals," in *IEEE Intelligent Vehicles Symposium (IV 2010)*, San Diego, CA, USA, June 2010.

[11] J. Jain and C. Busso, "Analysis of driver behaviors during common tasks using frontal video camera and CAN-Bus information," in *IEEE International Conference on Multimedia and Expo (ICME 2011)*, Barcelona, Spain, July 2011.

[12] N. Li, T. Misu, A. Tawari, A. Miranda, C. Suga, and K. Fujimura, "Driving maneuver prediction using car sensor and driver physiological signals," in *ACM International Conference on Multimodal Interaction (ICMI 2016)*, Tokyo, Japan, October 2016, pp. 108–112.

[13] Y. Murphey, D. S. Kochhar, P. Watta, X. Wang, and T. Wang, "Driver lane change prediction using physiological measures," *SAE International Journal of Transportation Safety*, vol. 3, no. 2, pp. 118–125, July 2015.

[14] N. Li, T. Misu, and A. Miranda, "Driver behavior event detection for manual annotation by clustering of the driver physiological signals," in *IEEE International Conference on Intelligent Transportation Systems (ITSC 2016)*, Rio de Janeiro, Brazil, November 2016, pp. 2583–2588.

[15] Y. Qiu, T. Misu, and C. Busso, "Analysis of the relationship between physiological signals and vehicle maneuvers during a naturalistic driving study," in *Intelligent Transportation Systems Conference (ITSC 2019)*, Auckland, New Zealand, October 2019, pp. 3230–3235.

[16] P. Mohan, V. N. Padmanabhan, and R. Ramjee, "Nericell: rich monitoring of road and traffic conditions using mobile smartphones," in *ACM Conference on Embedded Network Sensor Systems (SenSys 2008)*, Raleigh NC USA, November 2008, pp. 323–336.

[17] A. Aljaafreh, N. Alshabatat, and M. S. Najim Al-Din, "Driving style recognition using fuzzy logic," in *IEEE International Conference on Vehicular Electronics and Safety (ICVES 2012)*, Istanbul, Turkey, July 2012, pp. 460–463.

[18] H. Eren, S. Makinist, E. Akin, and A. Yilmaz, "Estimating driving behavior by a smartphone," in *IEEE Intelligent Vehicles Symposium (IV 2012)*, Alcala de Henares, Spain, June 2012, pp. 234–239.

[19] B. Nirmali, S. Wickramasinghe, T. Munasinghe, C. Amalraj, and H. Bandara, "Vehicular data acquisition and analytics system for real-time driver behavior monitoring and anomaly detection," in *IEEE International Conference on Industrial and Information Systems (ICIIS 2017)*, Peradeniya, Sri Lanka, December 2017, pp. 1–6.

[20] C. Yang, A. Renzaglia, A. Paigwar, C. Laugier, and D. Wang, "Driving behavior assessment and anomaly detection for intelligent vehicles," in *IEEE International Conference on Cybernetics and Intelligent Systems (CIS 2019) and IEEE Conference on Robotics, Automation and Mechatronics (RAM 2019)*, Bangkok, Thailand, November 2019, pp. 524–529.

[21] C. Ryan, F. Murphy, and M. Mullins, "End-to-end autonomous driving risk analysis: A behavioural anomaly detection approach," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[22] M. Fazeen, B. Gozick, R. Dantu, M. Bhukhiya, and M. C. González, "Safe driving using mobile phones," *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 3, pp. 1462–1468, September 2012.

[23] L. Xu, S. Li, K. Bian, T. Zhao, and W. Yan, "Sober-drive: A smartphone-assisted drowsy driving detection system," in *International Conference on Computing, Networking and Communications (ICNC 2014)*, Honolulu, HI, USA, February 2014, pp. 398–402.

[24] N. Li and C. Busso, "Detecting drivers' mirror-checking actions and its application to maneuver and secondary task recognition," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 4, pp. 980–992, April 2016.

[25] N. Li, J. Jain, and C. Busso, "Modeling of driver behavior in real world scenarios using multiple noninvasive sensors," *IEEE Transactions on Multimedia*, vol. 15, no. 5, pp. 1213–1225, August 2013.

[26] S. Ramyar, A. Homaifar, A. Karimoddini, and E. Tunstel, "Identification of anomalies in lane change behavior using one-class SVM," in *IEEE International Conference on Systems, Man, and Cybernetics (SMC 2016)*, Budapest, Hungary, October 2016, pp. 4405–4410.

[27] J. Yu, Z. Chen, Y. Zhu, Y. Chen, L. Kong, and M. Li, "Fine-grained abnormal driving behaviors detection and identification with smartphones," *IEEE Transactions on Mobile Computing*, vol. 16, no. 8, pp. 2198–2212, August 2017.

[28] R. Chai, G. R. Naik, T. N. Nguyen, S. H. Ling, Y. Tran, A. Craig, and H. T. Nguyen, "Driver fatigue classification with independent component by entropy rate bound minimization analysis in an EEG-based system," *IEEE Journal of Biomedical and Health Informatics*, vol. 21, no. 3, pp. 715–724, May 2017.

[29] N. El Masry, P. El-Dorry, M. El Ashram, A. Atia, and J. Tanaka, "Amelio-rater: Detection and classification of driving abnormal behaviours for automated ratings and real-time monitoring," in *International Conference on Computer Engineering and Systems (ICCES 2018)*, Cairo, Egypt, December 2018, pp. 609–616.

[30] I. Mohamad, M. Ali, and M. Ismail, "Abnormal driving detection using real time global positioning system data," in *IEEE International Conference on Space Science and Communication (IconSpace 2011)*, Penang, Malaysia, July 2011, pp. 1–6.

[31] C. Saiprasert and W. Pattara-Atikom, "Smartphone enabled dangerous driving report system," in *Hawaii International Conference on System Sciences (HICSS 2013)*, Wailea, Maui, HI, January 2013, pp. 1231–1237.

[32] J. Wahlström, I. Skog, and P. Händel, "Risk assessment of vehicle cornering events in GNSS data driven insurance telematics," in *IEEE Conference on Intelligent Transportation Systems (ITSC 2014)*, Qingdao, China, October 2014, pp. 3132–3137.

[33] ——, "Detection of dangerous cornering in GNSS-data-driven insurance telematics," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 6, pp. 3073–3083, December 2015.

[34] F. Li, H. Zhang, H. Che, and X. Qiu, "Dangerous driving behavior detection using smartphone sensors," in *IEEE International Conference on Intelligent Transportation Systems (ITSC 2016)*, Rio de Janeiro, Brazil, November 2016, pp. 1902–1907.

[35] Z. Liu, M. Wu, K. Zhu, and L. Zhang, "SenSafe: A smartphone-based traffic safety framework by sensing vehicle and pedestrian behaviors," *Mobile Information Systems*, vol. 2016, pp. 1–13, October 2016.

[36] P. Vavouranakis, S. Panagiotakis, G. Mastorakis, C. X. Mavromoustakis, and J. M. Batalla, "Recognizing driving behaviour using smartphones," in *Beyond the Internet of Things: Everything Interconnected*, J. Batalla, G. Mastorakis, C. Mavromoustakis, and E. Pallis, Eds. Cham, Switzerland: Springer International Publishing, January 2017, pp. 269–299.

[37] Z. Constantinescu, C. Marinoiu, and M. Vladoiu, "Driving style analysis using data mining techniques," *International Journal of Computers Communications & Control*, vol. 5, no. 5, pp. 654–663, December 2010.

[38] Y. Zheng and J. Hansen, "Unsupervised driving performance assessment using free-positioned smartphones in vehicles," in *IEEE International Conference on Intelligent Transportation Systems (ITSC 2016)*, Rio de Janeiro, Brazil, November 2016, pp. 1598–1603.

[39] J. Hansen, C. Busso, Y. Zheng, and A. Sathyanarayana, "Driver modeling for detection and assessment of driver distraction: Examples from the UTDrive test bed," *IEEE Signal Processing Magazine*, vol. 34, no. 4, pp. 130–142, July 2017.

[40] A. Hamdy, A. Atia, and M. Mostafa, "Recognizing driving behavior and road anomaly using smartphone sensors," *International Journal of Ambient Computing and Intelligence (IJACI)*, vol. 8, no. 3, pp. 22–37, July 2017.

[41] U. Fugiglando, E. Massaro, P. Santi, S. Milardo, K. Abida, R. Stahlmann, F. Netter, and C. Ratti, "Driving behavior analysis through CAN bus data in an uncontrolled environment," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 2, pp. 737–748, February 2019.

[42] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection for discrete sequences: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 5, pp. 823–839, May 2012.

[43] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems (NIPS 2014)*, vol. 27, Montreal, Canada, December 2014, pp. 2672–2680.

[44] F. Di Mattia, P. Galeone, M. D. Simoni, and E. Ghelfi, "A survey on GANs for anomaly detection," *ArXiv e-prints (arXiv:1906.11632)*, pp. 1–8, June 2019.

[45] T. Schlegl, P. Seeböck, S. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *Information Processing in Medical Imaging (IPMI 2017)*, ser. Lecture Notes in Computer Science, M. Niethammer, M. Styner, S. Aylward, H. Zhu, I. Oguz, P. Thian Yap, and D. Shen, Eds. Boone, NC, USA: Springer Berlin Heidelberg, June 2017, vol. 10265, pp. 146–157.

[46] B. Zhou, S. Liu, B. Hooi, X. Cheng, and J. Ye, "BeatGAN: Anomalous rhythm detection using adversarially generated time series," in *International Joint Conference on Artificial Intelligence (IJCAI 2019)*, Macao, China, August 2019, pp. 4433–4439.

[47] Y. Xue, T. Xu, H. Zhang, L. Rodney Long, and X. Huang, "SegAN: Adversarial network with multi-scale L1 loss for medical image segmentation," *Neuroinformatics*, vol. 16, pp. 383–392, May 2018.

[48] D. Li, D. Chen, J. Goh, and S.-K. Ng, "Anomaly detection with generative adversarial networks for multivariate time series," in *International Workshop on Big Data, Streams and Heterogeneous Source Mining BigMine 2018*, London, UK, August 2018, pp. 1–10.

[49] Y. Zheng, G. Chen, and M. Huang, "Out-of-domain detection for natural language understanding in dialog systems," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1198–1209, 2020.

[50] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *ArXiv e-prints (arXiv:1411.1784)*, November 2014.

[51] S. Hyland, C. Esteban, and G. Rätsch, "Real-valued (medical) time series generation with recurrent conditional GANs," *ArXiv e-prints (arXiv:1706.02633)*, June 2017.

[52] S. Akcay, A. Atapour-Abarghouei, and T. Breckon, "GANomaly: Semi-supervised anomaly detection via adversarial training," in *Asian Conference on Computer Vision (ACCV 2018)*, ser. Lecture Notes in Computer Science, C. Jawahar, H. Li, G. Mori, and K. Schindler, Eds. Perth, Australia: Springer Berlin Heidelberg, December 2018, vol. 11363, pp. 622–637.

[53] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient GAN-based anomaly detection," *ArXiv e-prints (arXiv:1802.06222)*, pp. 1–7, February 2018.

[54] Y. Qiu, T. Misu, and C. Busso, "Use of triplet loss function to improve driving anomaly detection using conditional generative adversarial network," in *Intelligent Transportation Systems Conference (ITSC 2020)*, Rhodes, Greece, September 2020, pp. 1–7.

[55] T. Misu and Y. Chen, "Toward reasoning of driving behavior," in *International Conference on Intelligent Transportation Systems (ITSC 2018)*, Maui, HI, USA, November 2018, pp. 204–209.

[56] V. Ramanishka, Y.-T. Chen, T. Misu, and K. Saenko, "Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2018)*, Salt Lake City, UT, USA, June 2018, pp. 7699–7707.

[57] O. Rompelman, A. Coenen, and R. Kitney, "Measurement of heart rate variability: Part i - comparative study of heart rate variability analysis methods," *Medical and Biological Engineering and Computing*, vol. 15, pp. 223–239, May 1977.

[58] J. Taelman, S. Vandeput, A. Spaepen, and S. Van Huffel, "Influence of mental stress on heart rate and heart rate variability," in *4th European Conference of the International Federation for Medical and Biological Engineering*, ser. IFMBE Proceedings, J. Van der Sloten, P. Verdonck, M. Nyssen, and J. Haueisen, Eds. Antwerp, Belgium: Springer Berlin Heidelberg, November 2009, vol. 22, pp. 1366–1369.

[59] J. Healey and R. Picard, "Detecting stress during real-world driving tasks using physiological sensors," *IEEE Transactions on intelligent transportation systems*, vol. 6, no. 2, pp. 156–166, June 2005.

[60] M. Nishigaki, R. Mose, O. Takahata, H. Imafuku, and H. Aoygai, "Quantitative evaluation on mental workload reduction for hands free driving," in *International Conference on Intelligent Transportation Systems (ITSC 2018)*, Maui, HI, USA, November 2018, pp. 230–235.

[61] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems (NIPS 2012)*, vol. 25, Lake Tahoe, CA, USA, December 2012, pp. 1097–1105.

[62] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Deep Learning and Representation Learning Workshop: NIPS 2014*, Montreal, QC, Canada, December 2014, pp. 1–9.

[63] A. Graves, "Generating sequences with recurrent neural network," *ArXiv e-prints (arXiv:1308.0850)*, pp. 1–43, August 2013.

[64] S. Sadjadi and J. H. L. Hansen, "Unsupervised speech activity detection using voicing measures and perceptual spectral flux," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 197–200, March 2013.

[65] Z. Yu, W. Peng, X. Li, X. Hong, and G. Zhao, "Remote heart rate measurement from highly compressed facial videos: An end-to-end deep learning solution with video enhancement," in *IEEE/CVF International Conference on Computer Vision (ICCV 2019)*, Seoul, South Korea, October-November 2019, pp. 151–160.

[66] Z. Yu, X. Li, and G. Zhao, "Remote photoplethysmograph signal measurement from facial videos using spatio-temporal networks," in *British Machine Vision Conference (BMVC 2019)*, Cardiff, UK, September 2019, pp. 1–12.