

# Unsupervised Scalable Multimodal Driving Anomaly Detection

Yuning Qiu<sup>1</sup>, Student Member, IEEE, Teruhisa Misu<sup>2</sup>, Member, IEEE, and Carlos Busso<sup>1</sup>, Senior Member, IEEE

**Abstract**—Driving anomaly detection aims to identify objects, events or actions that can increase the risk of accidents, reducing road safety. While supervised approaches can effectively identify aspects related to driving anomalies, it is unfeasible to tabulate and address all potential driving anomalies. Instead, it is appealing to design unsupervised approaches that can automatically identify unexpected driving scenarios. This study formulates the detection of driving anomalies as a binary-discrimination task between expected and unexpected driving behaviors. We propose an unsupervised contrastive method using conditional *generative adversarial networks* (GANs) implemented with the attention model and the triplet loss function. A feature of our framework is its scalability, where it is easy to add new modalities. We consider five different modalities: the vehicle’s CAN-Bus signals, driver’s physiological signals, distance to nearby pedestrians, distance to nearby vehicles and distance to nearby bicycles. Our approach trains a conditional GAN to extract latent features from each of the five modalities. An attention model combines the latent representations from the modalities. The entire framework is trained with the triplet loss function to generate effective representations to discriminate normal and abnormal driving segments. We conduct experimental evaluations on the *driving anomaly dataset* (DAD), achieving improved performance over alternative approaches.

**Index Terms**—Driving anomaly detection, conditional generative adversarial networks, attention mechanism, triplet loss function.

## I. INTRODUCTION

IDENTIFYING abnormal driving behaviors is an important research area with significant societal impact as lives can be saved by increasing road safety. Multiple rule-based and pattern-based methods have been proposed for driving anomaly detection, including monitoring of road conditions [1]–[3], aggressive driving behaviors [4]–[9], risky driving patterns [10]–[16] and unusual driving styles (e.g., fatigue and meandering) [17]–[24]. A typical challenge in those driving anomaly detection methods is that the vehicle’s driving conditions can vary significantly under different scenarios, which make driving patterns and rules hard to reliably establish. Furthermore, it is nearly impossible to exhaustively tabulate all possible actions or situations that lead to hazardous scenarios. Fig. 1 shows four relevant



Fig. 1. Examples of abnormal driving scenarios where driver’s maneuvers are affected by other vehicles or pedestrians: (a) a car drives in the wrong lane in front of the car, (b) a pedestrian suddenly crosses the street, (c) a bicyclist rushes across the street, and (d) a vehicle cuts into the vehicle’s lane.

examples of driving scenarios, illustrating the difficulty in building rule-based systems to detect anomalous scenarios, or creating specialized approaches to deal with each case. An appealing approach is to use unsupervised multimodal approaches to detect driving anomalies by discriminating expected driving behaviors as normal cases and unexpected driving behaviors as abnormal cases.

This study proposes an unsupervised contrastive framework to identify driving anomalies using multiple modalities. The key principle in our formulation is that anomalous driving scenarios are characterized by deviations from expected behaviors. Our approach creates predictions of future frames, conditioned on the values of these signals observed in previous frames. Then, it contrasts the predictions with the actual signals, quantifying their differences. The core feature extraction module relies on conditional *generative adversarial networks* (GANs), following the ideas presented in our previous study [25]. We build one conditional GAN per modality, where its generator creates the predictions of the signals from upcoming frames and the discriminator determines if the data is real or synthesized by the generator. Then, we extract the embedding of the penultimate layer of the discriminator, which is used as the representation for the modality. A novel contribution in this study is the fusion of the modalities, where we rely on the self-attention

Manuscript received October 25, 2021; revised March 5, 2022; accepted March 10, 2022. This work was supported by Honda Research Institute USA, Inc. (Corresponding author: Carlos Busso.)

Yuning Qiu and Carlos Busso are with the Department of Electrical Engineering, University of Texas at Dallas, Richardson, TX 75080 USA (e-mail: yxq180000@utdallas.edu; busso@utdallas.edu).

Teruhisa Misu is with Honda Research Institute USA, Inc., San Jose, CA 95134 USA (e-mail: TMisu@honda-ri.com).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TIV.2022.3160861>.

Digital Object Identifier 10.1109/TIV.2022.3160861

mechanism [26]. The weights assigned to the modalities by the attention mechanism indicate the relative importance of each modality. A strength of the approach is the contrastive loss used to train the proposed formulation in an unsupervised manner. We rely on the triplet loss function [27], where the goal is to reduce the distance between the predicted data and the observed signals, while increasing the distance between the predicted data and the data from a randomly selected segment. After pre-training the individual conditional GANs, the approach can be jointly trained, creating effective end-to-end solutions.

The proposed formulation is scalable, with separate GAN models applied to each of the modalities, avoiding dimension explosion. The feature embeddings extracted from the GAN models are fused by the attention model. An advantage of seamlessly incorporating more modalities is that the system can respond even when the driver is not aware of hazardous scenarios. Our previous work only considered the driver's physiological data and the vehicle's CAN-Bus data [28], [29]. In daily urban traffic, unexpected reactions and maneuvers can be caused by either a pedestrian rushing across the road, another vehicle abruptly cutting into the lane or mistakes made by the drivers (see real examples in Fig. 1). If the driver is not aware of these anomalies, her/his physiological reactions and maneuvers will not reflect the anomaly. Therefore, we incorporate environmental information from vision-based object detection systems applied to the road. In addition to physiological signals and CAN-Bus signals, we add three modalities: distances to nearby cars, pedestrians and bicycles. Our proposed system still perceives these driving anomalies even though the driver might have neglected them.

We rely on the recordings from the *driving anomaly dataset* (DAD) [28] to evaluate our proposed scalable multimodal approach. Experimental results show that recordings annotated with possible abnormal incidents (such as avoiding pedestrians, bicycles, or other vehicles) have higher anomaly scores than recordings without events. To validate the results, we implement perceptual evaluations of video segments, where human annotators were asked to assess the risk level, familiarity level, anomaly level, and causes of the anomalies of the driving scenarios. We evaluate our approach with three baselines. The first baseline is the CNN-LSTM based conditional GAN model proposed by Qiu *et al.* [25], which is trained with 2 modalities: the vehicle's CAN-Bus signals and the driver's physiological signals. The second baseline is the BeatGAN framework proposed by Zhou *et al.* [30], which is an unsupervised method using GANs also trained with CAN-bus and physiological signals. The third baseline is our proposed attention model trained only with the aforementioned two modalities to further quantify the effectiveness of adding the three modalities describing external information. The results show that when trained with CAN-Bus and physiological data, the proposed attention model leads to better performance than the CNN-LSTM based conditional GANs and the BeatGAN models. The discriminative performance of our model increases when we add contextual information about the road, modeling the distances to nearby pedestrians, bicycles and vehicles. This model leads to the best results observed for this task. The main contributions of our study are:

- Scalable formulation for driving anomaly detection that seamlessly incorporates new modalities using an attention model.
- Modeling of contextual information derived from vision-based object detection systems applied to the road, where

our approach can react even when the driver is unaware of potential anomalous scenarios.

- Exhaustive evaluations of the proposed approach using objective and perceptual evaluations on naturalistic recordings collected in real road environments.

This study is organized as follows. Section II presents related studies addressing the detection of driving anomalies. It also describes background information to understand the proposed architecture. Section III discusses the details of our proposed model. Section IV introduces the dataset to train and evaluate our proposed model, and the implementation details. Section V evaluates the discriminative performance of our proposed model with objective and subjective comparisons. Finally, Section VI summarizes the contributions of this work, discussing future research directions.

## II. RELATED WORK

### A. Driving Anomaly Detection

Studies have proposed methods for anomaly detection in several domains. In the area of in-vehicle safety systems, many approaches have been proposed for abnormal driving condition detection, either based on driving rules [1]–[4], [6], [7], [10]–[16], [18] or driving patterns [5], [8], [9], [17], [19]–[24]. Most of these studies use the vehicle's driving information (e.g., speed, acceleration and yaw angle) to describe the vehicle's driving conditions. The approaches based on driving rules detect target events by either setting a threshold on the vehicle's driving information [1]–[4], [6], [7], [10], [14], [16], [18], or calculating the driving behavior *key performance indicators* (KPI) using predefined formulas [11]–[13], [15]. The approaches based on driving patterns determine abnormal conditions utilizing machine learning methods, including *support vector machine* (SVM) [8], [17], [21], [31], *neural networks* (NN) [20], [23], *hidden Markov models* (HMM) [22] and Bayesian classifiers [5]. Chen *et al.* [8] extracted statistic features from the vehicle's acceleration and orientation, using these features to train a SVM that identifies six different abnormal driving patterns (i.e. weaving, swerving, sideslipping, fast U-turn, turning with a wide radius, and sudden braking). Some studies have utilized the driver's information, such as physiological signals [28], [29], [32], eye gaze information [33], [34], facial expressions [35], [36], and driving gestures [37], [38] to identify driving anomalies. Köpüklü *et al.* [38] used the videos recorded by a frontal camera facing the driver and a top camera facing the steering wheel to detect the driver's abnormal behaviors. To extract spatial-temporal features of the driver's behaviors, the authors trained a 3D-*convolutional neural network* (CNN) with contrastive loss to maximize the similarity between normal driving events, and minimize the similarity between normal and abnormal driving samples. During inferences, the feature representations of all the normal driving training clips are normalized using the l2 normalization, using this representation as a template vector describing normal driving. For each testing clip, the authors extracted a feature vector using the 3D-CNN model and calculated the cosine similarity between the feature vector and the normal driving template vector. The testing clips with a cosine similarity score with a value below a preset threshold were considered as anomalies.

With the development of computer vision, many studies have proposed methods to detect and identify driving anomalies by using a camera to collect information about the surrounding

188 traffic environment [39]–[42]. Yao *et al.* [41] proposed a vision-  
 189 based approach to detect traffic accidents in videos recorded by  
 190 a dashboard-mounted camera. The approach localizes detected  
 191 traffic participants (e.g., other vehicles and pedestrians) using  
 192 bounding boxes, making predictions on their trajectories based  
 193 on previous frames. They train their model with only normal  
 194 driving videos to detect deviations from predicted behaviors,  
 195 under the assumption that moving trajectories in traffic accidents  
 196 deviate from expected trajectories. Our study proposes an unsu-  
 197 pervised driving anomaly detection system by combining the ve-  
 198 hicle’s driving information, driver’s physiological information,  
 199 and vision-based surrounding traffic environmental information  
 200 to improve the performance of the system.

201 *B. Conditional GANs for Anomaly Detection on Time Series*

202 *Generative adversarial networks* (GANs) [43] have demon-  
 203 strated effectiveness for time series data anomaly detection [28],  
 204 [44]–[46]. A GAN consists of a *generator* (G) that creates  
 205 synthetic data from noise, and a *discriminator* (D) that deter-  
 206 mines whether the data is real or produced by the generator. By  
 207 training the generator and discriminator with an adversarial loss,  
 208 the model creates realistic synthetic data. As a state-of-the-art  
 209 generative approach, GANs have been used to detect anomalies  
 210 mostly in other domains. Zhou *et al.* [30] proposed BeatGAN,  
 211 which is a GAN-based system that was used for two problems:  
 212 to detect anomalous beats from *electrocardiogram* (ECG) sig-  
 213 nals, and to identify unusual human motions (e.g., hopping and  
 214 jumping) from normal activities such as walking. The approach  
 215 builds a generator with an encoder-decoder structure, using the  
 216 reconstructed signals as the generated fake signals to confuse  
 217 the discriminator. After training, they used the reconstruction  
 218 error between the real signal and the generated fake signal as  
 219 the anomalous metric to detect abnormal beats in ECG signals.  
 220 Other alternative approaches relying on GANs to detect anom-  
 221 alies in other domains include the methods presented by Hyland  
 222 *et al.* [47], Akcay *et al.* [48], and Zenati *et al.* [49].

223 *C. Attention Mechanism for Multimodal Fusion*

224 Our study uses attention networks [26] implemented with the  
 225 triplet loss function [27] to jointly learn discriminative embed-  
 226 dings for driving anomaly detection. Hori *et al.* [50] proposed  
 227 an attention-based feature fusion approach to incorporate audio,  
 228 motion and image features to describe the content of videos. The  
 229 approach calculates the attention weights of the input features  
 230 from different modalities, estimating the linear combination of  
 231 the embeddings of individual modalities using these attention  
 232 weights. The attention mechanism allows the relative weights of  
 233 each modality to change based on the context, showing that this  
 234 combination approach is effective to improve the description  
 235 accuracy. Chen *et al.* [51] utilized the self-attention mecha-  
 236 nism to fuse audiovisual features for an affect recognition task.  
 237 Song *et al.* [39] combined attention mechanism and triplet loss  
 238 function to learn effective representations from speech audio  
 239 for speaker diarization. The authors used an attention model to  
 240 calculate feature embeddings directly from *Mel-frequency cep-*  
 241 *stral coefficients* (MFCCs) obtained from the speech segments.  
 242 Then, they input the extracted features to the subsequent network  
 243 to learn a similarity metric with the triplet loss function. The  
 244 triplet loss function [27] has been widely used in discrimination  
 245 tasks facilitating contrastive learning solutions to learn more

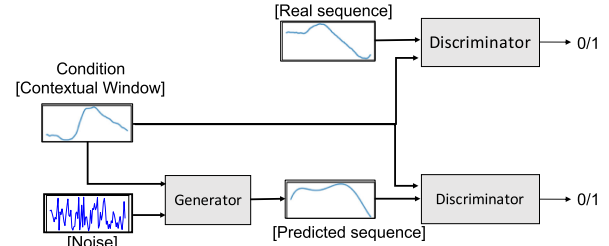


Fig. 2. Training procedure of the conditional GAN model. The generator  $G$  predicts plausible data of the upcoming driving segment based on the observed signals. The discriminator  $D$  determines if the data is real or created by  $G$ .

discriminant representations. Inspired by these studies, our pro-  
 posed methods combine the attention models with the triplet loss  
 function.

249 *D. Relation to Prior Work*

250 In our previous work [52], we found that features extracted  
 251 from the vehicle’s CAN-Bus signals and the driver’s physiologi-  
 252 cal signals can be used to discriminate different driving maneu-  
 253 vers. Utilizing the driver’s physiological data and the vehicle’s  
 254 CAN-Bus data, we proposed an unsupervised driving anomaly  
 255 detection approach based on conditional *generative adversarial*  
 256 *networks* (GANs) [25], [28], [29]. The driving anomalies were  
 257 defined as the events that deviate from normal or expected driv-  
 258 ing patterns that may lead to dangerous situations. Fig. 2 shows  
 259 the strategy for detecting driving anomalies using a conditional  
 260 GAN. We used the generator of the GAN to make predictions  
 261 on the vehicle’s CAN-Bus signals and the driver’s physiological  
 262 signals, conditioned on the data from previous driving segments.  
 263 The discriminator of the GAN was trained to identify whether the  
 264 input data was real or synthesized by the generator. The absolute  
 265 value of the difference between the discriminator outputs of the  
 266 predicted data and the upcoming real signal was regarded  
 267 as the anomaly metric,  $m_{anomaly}$ , which indicates the abnormal  
 268 level of the driving condition. An abnormal driving condition  
 269 was expected to have a higher value for  $m_{anomaly}$  than a normal  
 270 driving condition. Qiu *et al.* [29] extended the approach by  
 271 defining a new metric based on the triplet loss function. Based  
 272 on the conditional GAN model, the study proposed a triplet-loss  
 273 neural network which took the intermediate layer embeddings  
 274 of the discriminator as the input [29]. This triplet network was  
 275 trained to decrease the distance between the embeddings of the  
 276 prediction and real data, while increasing the distance between  
 277 the embeddings of the real data and an unpaired prediction  
 278 (i.e., predicted from unrelated segments). Compared with the  
 279 conditional GAN-based model, the triplet-loss neural network  
 280 increases the discrimination performance by contrasting the  
 281 differences between predicted and real features. This process  
 282 requires no label, leading to an appealing unsupervised approach  
 283 to detect driving anomalies.

284 Our previous approaches have two major limitations [25],  
 285 [28], [29]. First, the system responds only when the driver is  
 286 aware of the anomalies. The driver’s physiological signals and  
 287 the vehicle’s CAN-Bus data describe the driver’s reactions. The  
 288 system would fail to detect potential anomalies when the driver is  
 289 not aware of them (e.g., presence of a pedestrian on the road that  
 290 the driver has overlooked). Second, it is not easy for the system to  
 291 extend the approach to include more modalities. Increasing the



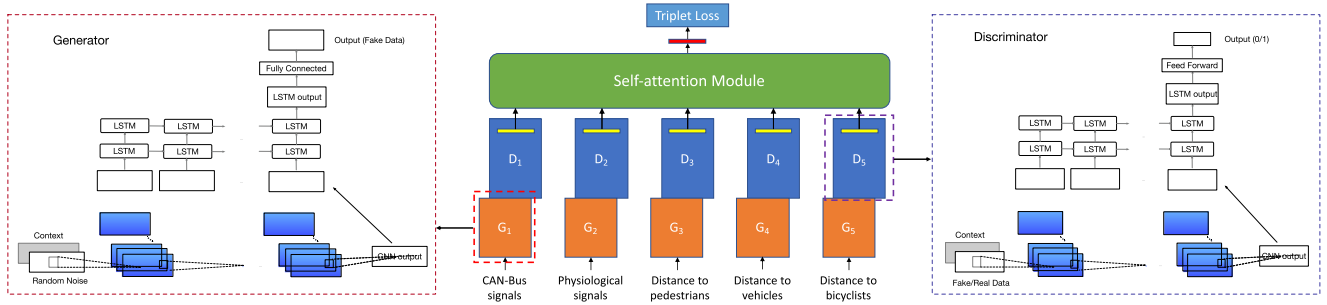


Fig. 3. Proposed unsupervised, scalable, multimodal architecture to detect driving anomalies. The feature representations are obtained with a conditional GAN for each of the modalities. In the figure, the variable  $G_i$  represents the generator of the  $i$  modality, and  $D_i$  represents the discriminator of the  $i$  modality. The attention model weights the modalities using a triplet loss function.

292 dimension of the inputs would prevent the convergence during  
 293 the training process of the GAN model.

294 Building upon our previous work, this study addresses these  
 295 two challenges by proposing an unsupervised scalable multi-  
 296 modal driving anomaly detection system. The modalities are  
 297 fused using an attention model, which provides a principled  
 298 approach to scale our formulation to include more modalities.  
 299 We can seamlessly incorporate information about nearby pedest-  
 300 rians, bicycles and other vehicles. This is a contrastive approach  
 301 implemented with the triplet loss function, which does not re-  
 302 quire labeled data. These features are fundamental contributions  
 303 that make our approach more appealing for real applications.

### 304 III. PROPOSED MODEL

305 This study proposes a novel unsupervised driving anomaly  
 306 detection framework that has three main blocks. Fig. 3 shows an  
 307 overview of our framework. The first block extracts embeddings  
 308 from multiple modalities with conditional *generative adversar-*  
 309 *ial networks* (GANs). The second block fuses the modalities  
 310 with the attention mechanism, learning from the data how to  
 311 weight the representations from each modality. The third block  
 312 is the triplet loss function that is used to train the model, learning  
 313 a contrastive-based metric that indicates the anomaly level of the  
 314 target recording.

315 Our proposed implementation has five modalities: the vehi-  
 316 cle’s CAN-Bus signals, the driver’s physiological signals, the  
 317 distances to nearby vehicles, the distances to nearby bicyclists  
 318 and the distances to nearby pedestrians. By combining the  
 319 conditional GAN models, self-attention mechanism and triplet  
 320 loss function, we aim to create a framework that is (1) scalable,  
 321 making it easy to add more modalities if needed, and (2) ef-  
 322 fective, learning representations of the features extracted from  
 323 different modalities. This section describes the details about the  
 324 three building blocks of our proposed method.

#### 325 A. Feature Extraction Using Conditional GANs

326 The first block in the system extracts a discriminative feature  
 327 representation for each of the modalities. This feature extraction  
 328 module is implemented with the conditional GANs used in the  
 329 unsupervised driving anomaly detection system proposed by  
 330 Qiu *et al.* [25]. Instead of adopting an *early fusion* approach by  
 331 building one GAN model that takes all the multimodal signals  
 332 as input, we adopt a *model-level fusion* approach by building

333 separate GANs for each modality, which are later combined  
 334 using the attention mechanism. As mentioned in Section II-D,  
 335 the key purpose of using a GAN for this task is to generate  
 336 predictions that are compared with the observed signals. Fig. 3  
 337 shows the architecture of the generator and discriminator of the  
 338 conditional GANs, which is the same architecture proposed in  
 339 Qiu *et al.* [25]. We use CNNs and *recurrent neural networks*  
 340 (RNNs) implemented with *long-short term memory* (LSTM)  
 341 cells [53]. The CNNs extract feature embeddings from the orig-  
 342 inal input signals without relying on hand crafted features. The  
 343 output of the CNNs are then processed by the LSTM network to  
 344 leverage temporal information in the time series sequence. For  
 345 each modality, the *generator* ( $G$ ) predicts plausible data of the  
 346 upcoming 6-second driving segments based on the previous 30  
 347 seconds signals, and the *discriminator*  $D$  determines whether the  
 348 data is real or fake. Equations 1 and 2 show the cost function of  
 349 this adversarial task, where  $x$  is the data sample,  $z$  is the noise  
 350 sample,  $p_{data}$  is the distribution of data and  $p_z$  is the distribution  
 351 of the noise.

$$\max_D V(D) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

$$\min_G V(G) = \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (2)$$

352 From each conditional GAN model, we extract the embedding  
 353 of the penultimate layer of  $D$  as the feature embedding of the  
 354 modality. By building separate GANs for each modality, our  
 355 proposed system is easy to scale when more modalities are avail-  
 356 able. Section IV-B discusses implementation details, including  
 357 pre-training each GAN before jointly training the entire system.

#### 358 B. Self-Attention Model for Multimodal Fusion

359 The combination of features from multiple modalities is  
 360 expected to effectively improve the model performance. This  
 361 section describes the self-attention network used to implement  
 362 the fusion of  $N$  modalities, each of which has its own feature  
 363 embedding, extracted from the penultimate layer of its  $D$ . The  
 364 key idea is to linearly combine the individual embedding by  
 365 dynamically defining the modality weights using the attention  
 366 mechanism. For a driving segment, the attention model takes  
 367  $N$  embeddings as input features. Fig. 4 shows the structure of  
 368 the attention network used in this work. The core component  
 369 of the attention network is the multi-head module from the

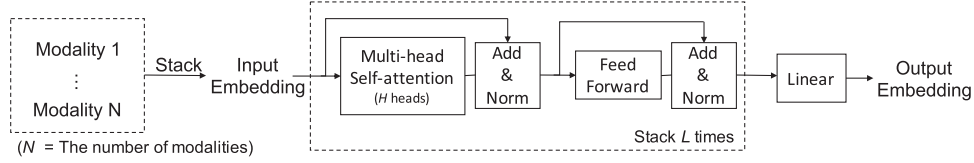


Fig. 4. Details of the architecture used for the attention module. The output of this model is the output embedding used for the triplet loss function.

370 self-attention mechanism [26]. More specifically, we stack the  
 371 features of each modality as the input of the attention model,  
 372 which we denoted  $X$ . For each head, we estimate the matrices  
 373  $W_Q$ ,  $W_K$  and  $W_V$ . These matrices are trainable parameters  
 374 to map the input  $X$  into  $Q$  (query),  $K$  (key), and  $V$  (value),  
 375 respectively. We map  $X$  into these three subspaces by multi-  
 376 plying these matrices with  $X$  (i.e.,  $Q = XW_Q$ ,  $K = XW_K$   
 377 and  $V = XW_V$ ). We compute the scaled dot-product attention  
 378 based on the attention matrices. Then, the dot product of  $Q$  and  
 379  $K$  are activated by the softmax function as the attention weights.  
 380 The matrix of attention representation is computed as:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) \quad (3)$$

$$\text{Attention}(Q, K, V) = WV \quad (4)$$

381 where  $d_k = 256$  is the dimension of the attention matrix  $K$ .  
 382 The attention weight matrix  $W$  describes the interaction among  
 383 the  $N$  input modalities by computing the scaled inner prod-  
 384 uct between pairs of modalities. The number of multi-head  
 385 attentions is denoted by  $H$ . The attention representations are  
 386 computed using  $H$  parallel sets of attention matrices, denoted  
 387 as heads. The reason for assigning different matrices to each  
 388 attention head ( $W_Q$ ,  $W_K$ ,  $W_V$ ) is that the model pay attention  
 389 to the relationship among different modalities. We concate-  
 390 nate the resulting  $H$  attention representations together as an  
 391 ensemble of attention representations. Multi-head attention pre-  
 392 vents the model from focusing on only one modality by jointly  
 393 considering information from multiple representations. This  
 394 multi-head attention module can be stacked multiple times for  
 395 a deeper structure. We denote the number of stacked attention  
 396 modules by  $L$ . The connection between two modules is a *feed*  
 397 *forward network* (FFN) implemented with two fully connected  
 398 layers, where the activation function of the first layer is the  
 399 *rectified linear unit* (ReLU). In (5),  $W_1$  and  $W_2$  are the weight  
 400 matrices, and  $b_1$  and  $b_2$  are the bias terms of the FFN.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (5)$$

### 401 C. Triplet Loss for Metric Learning

402 Inspired by the work of Qiu *et al.* [29], the representations  
 403 from the attention model are then used to learn a similarity  
 404 metric with the triplet loss function. The use of this contrastive  
 405 loss aims to build embeddings that are discriminative for the  
 406 driving anomaly detection task using an unsupervised strategy.  
 407 In a triplet network, each input is constructed as a set of three  
 408 samples:  $s_p$ ,  $s_a$ , and  $s_n$ . The sample  $s_a$  denotes an anchor,  $s_p$   
 409 denotes a positive sample belonging to the same class as  $s_a$ , and  
 410  $s_n$  denotes a negative sample from a different class. The goal of  
 411 the triplet loss function is to create an embedding that minimizes  
 412 the distance between the anchor and the positive sample while  
 413 increasing the distance between the anchor and the negative

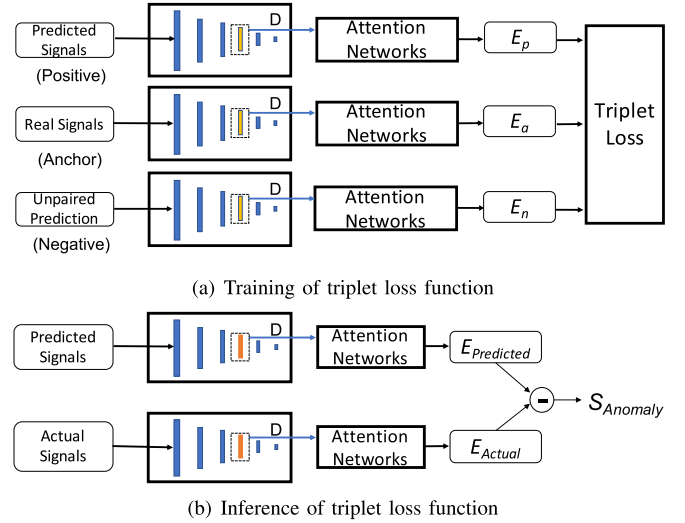


Fig. 5. Attention network trained with the triplet loss function. The penultimate layer embeddings of the discriminators are extracted as input of the attention model. During inferences, we estimate the absolute difference between  $E_{Actual}$  and  $E_{Predicted}$ , which is used as the anomaly score.

414 sample. This study considers the real data to be predicted  
 415 as the anchor example  $s_a$ , and the prediction conditioned on  
 416 the previous frames as the positive example  $s_p$ . The negative  
 417 example  $s_n$  corresponds to the predicted data from another  
 418 randomly selected driving segment (i.e., unpaired data). Fig. 5(a)  
 419 shows the training procedure. The samples are processed by the  
 420 separate GAN models (Section III-A) and the attention model  
 421 (Section III-B). The corresponding outputs are referred to as  
 422  $E_a$  for the anchor,  $E_p$  for the positive sample, and  $E_n$  for the  
 423 negative sample. We use the Euclidian distance between these  
 424 vectors to estimate the cost function, which is defined in (8). The  
 425 distance between  $E_a$  and  $E_p$  is minimized, while the distance  
 426 between  $E_a$  and  $E_n$  is maximized to be larger than a preset  
 427 margin  $\alpha$ .

$$D_{ap} = \|E_a - E_p\|_2 \quad (6)$$

$$D_{an} = \|E_a - E_n\|_2 \quad (7)$$

$$L_{Triplet} = \max(0, D_{ap}^2 - D_{an}^2 + \alpha) \quad (8)$$

428 This loss function maps the embedding of the predicted data,  
 429 closer to the embedding of the corresponding actual data and  
 430 far away from the embedding of the unpaired data. This whole  
 431 process is fully unsupervised, requiring no labels.

432 Fig. 5(b) shows the inference procedure. For a driving seg-  
 433 ment, we process the real data, obtaining  $E_{Actual}$  and the pre-  
 434 dicted data by the generator, obtaining  $E_{Predicted}$ . Equation 9  
 435 shows the final anomaly score, which consists of the difference

436 between  $E_{Actual}$  and  $E_{Predicted}$ . A high value of  $S_{anomaly}$  indicates that the driving segment is more unexpected, suggesting  
 437 a higher degree of anomaly.  
 438

$$S_{anomaly} = |E_{Actual} - E_{Predicted}| \quad (9)$$

#### 439 IV. EXPERIMENTAL SETTING

##### 440 A. Driving Anomaly Dataset (DAD)

441 The experiments in this study rely on the *driving anomaly*  
 442 *dataset* (DAD) [28] collected by *Honda Research Institute* (HRI)  
 443 in an Asian city. The dataset contains 250 hours of naturalistic  
 444 driving recordings, where 84 hours are used in this study. The  
 445 data is collected during day time, and most of the driving scenarios  
 446 are under urban driving environments, including residential,  
 447 school area, and downtown area. The data includes very little  
 448 segments with highway driving. We rely on the vehicle's CAN-  
 449 Bus signals, which consist of the vehicle's speed, yaw angle,  
 450 steer angle, steer speed, pedal pressure and pedal angle (6D  
 451 vector). We also use the driver's physiological signals, which are  
 452 collected using a chest band (heart rate and breath rate - Zephyr  
 453 BioHarness 3 chestband) and a wristband (skin conductance  
 454 and sphygmus - Empatica E4). From these sensors, we use the  
 455 following three signals: *heart rate* (HR), *breath rate* (BR), and  
 456 *electrodermal activity* (EDA). We also leverage road information  
 457 extracted with a vision-based object detection system. The object  
 458 distance information includes the distance to nearby vehicles,  
 459 pedestrians, and bicyclists. The objects are detected by a smart  
 460 camera mounted on the interior side of the windshield, utilizing  
 461 Mobileye technology. This system measures the distances to  
 462 nearby pedestrians, bicyclists, vehicles and lane markings. Mobileye's  
 463 algorithm can simultaneously detect multiple objects. For this study,  
 464 we only consider the two closest pedestrians, bicyclists, and vehicles.  
 465 Each of these modalities is represented with a 4D vector including  
 466 the horizontal and vertical distances from the car of the two closest  
 467 pedestrians, bicyclists, or vehicles. All the considered signals are  
 468 synchronized at the sampling rate of 30 Hz.  
 469

470 The dataset is manually annotated using the camera recording  
 471 of the road. The annotation process followed the same protocol  
 472 used in the collection of the *Honda Research Institute driving*  
 473 *dataset* (HDD) [54], [55]. The annotation includes the presence  
 474 of several events and maneuvers. Regular driving maneuvers,  
 475 such as turns and lane changes, are defined as goal-oriented  
 476 operations, while the maneuvers that are influenced by other  
 477 traffic participants are defined as stimuli-driven operations (e.g.,  
 478 avoid pedestrian near ego lane and avoid on-road bicyclist).  
 479 More detailed information about this dataset is provided by  
 480 the studies of Qiu *et al.* [28], [29]. In this work, we group  
 481 the driving segments into two sets according to the annotations  
 482 provided by the annotators. The driving segments that overlap  
 483 with no annotations are considered as the *normal* set. The driving  
 484 segments that overlap with stimuli-driven operation, driver's  
 485 error and traffic rule violation annotations are grouped as the  
 486 *candidate* set. These segments can potentially be associated  
 487 with driving anomalies. Table I shows the details with the  
 488 annotations included in these two sets. The *candidate* driving  
 489 set represents only 1.69% of the recordings. This ratio is similar  
 490 across partitions with 1.57% for the train set, 1.53% for the  
 491 development set and 2.44% for the test set. This study considers  
 492 89 sessions, which correspond to approximately 84 hours of  
 493 urban driving recordings. We split these recordings into 3 sets:

TABLE I  
 DEFINITION OF CANDIDATE AND NORMAL SETS. THE ANNOTATIONS  
 CORRESPOND TO THE LABELS INCLUDED IN THE DAD CORPUS

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
Normal	No annotations during the segments

train (72 sessions, approx. 70 hours), development (3 sessions, 494  
 approx. 4 hours), and test (14 sessions, approx. 10 hours) sets. 495

##### 496 B. Implementation Details

497 This section introduces the implementation details of our  
 498 approach. Our proposed model includes the conditional GANs,  
 499 to derive discriminative feature representations, and the self-  
 500 attention networks, to fuse the modalities. We implement the  
 501 conditional GANs with *convolutional neural networks* (CNNs)  
 502 and *recurrent neural networks* (RNNs). The generator consists  
 503 of six convolutional layers, implemented with 64, 64, 128, 128,  
 504 64 and 1 channels, respectively. We use batch normalization  
 505 and a leaky ReLU function [56] for each layer except the output  
 506 layer. The output of the CNNs is fed into the RNNs, which  
 507 are implemented with two layers of *long short-term memory*  
 508 (LSTM) cells. The number of units in each LSTM cell is 64. The  
 509 output of the LSTM cells goes through a single fully connected  
 510 layer, where its dimension is equal to the corresponding input  
 511 modality. Similarly, the discriminator consists of four convolutional  
 512 layers, implemented with 64, 128, 128 and 64 channels,  
 513 respectively, followed by two layers of LSTM cells. Each LSTM  
 514 layer is implemented with 64 units. The output of the LSTM  
 515 is fed into the feed forward networks, which has three layers  
 516 with dimensions 1024, 1024, and 1, respectively. The first two  
 517 layers are activated with the leaky ReLU function, while the last  
 518 layer is activated with a sigmoid function. The 1024-dimensional  
 519 embedding of the second layer will be extracted as the unimodal  
 520 feature representation of each modality.

521 During the training process, we train the generator and dis-  
 522 criminator for 20 epochs. We use the Adam optimizer, with a  
 523 learning rate set to 0.001. After training the GANs, we freeze  
 524 the GANs' parameters and extract an unimodal feature rep-  
 525 resentation for each modality, which we denote  $z_{CAN-Bus}$ ,  
 526  $z_{physiological}$ ,  $z_{pedestrian}$ ,  $z_{vehicle}$ , and  $z_{bicyclist}$ . We map these  
 527 vectors into a subspace with a trainable projection implemented  
 528 with the Tanh activation to produce the vector representations  
 529  $x_{CAN-Bus}$ ,  $x_{physiological}$ ,  $x_{pedestrian}$ ,  $x_{vehicle}$ , and  $x_{bicyclist}$ .  
 530 These transformations compensate for the differences in magni-  
 531 tude. Then, we stack the vector embeddings of the five modalities  
 532 as the input of the attention networks. We denote this matrix as  
 533  $X \in \mathbb{R}^{N \times d_{model}}$ , where  $N = 5$  and  $d_{model} = 512$ . As introduced  
 534 in Section III-B, we apply multi-head attention mechanism to  
 535 attend to information from different representation subspaces as  
 536 following:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)W^O \quad (10)$$

$$\text{head}_i = \text{Attention}(XW_i^Q, XW_i^K, XW_i^V) \quad (11)$$

where the parameter matrices are  $W_i^Q \in \mathbb{R}^{d_{model} \times d_Q}$ ,  $W_i^K \in$  537  
 $\mathbb{R}^{d_{model} \times d_K}$ ,  $W_i^V \in \mathbb{R}^{d_{model} \times d_V}$ , and  $W_i^O \in \mathbb{R}^{d_{model} \times d_V}$ . We 538



539 use five heads (i.e.,  $H = 5$ ), setting the dimensions of the query,  
 540 key and value to 256 (i.e.,  $d_Q = d_K = d_V = 256$ ). Section V-A  
 541 discusses results with different number of heads. The attention  
 542 module is stacked  $L$  times, setting  $L = 2$ . The feed forward  
 543 network in the attention model is implemented with three fully  
 544 connected layers with dimension equal to  $d_{model}$  to facilitate the  
 545 residual connections.

546 The parameters of the attention networks are trained with  
 547 the triplet loss function introduced in Section III-C. We use the  
 548 Adam optimizer with a learning rate equal to 0.001. After ten  
 549 epochs, we jointly train the parameters of the GANs and the  
 550 attention networks for another five epochs, where all the param-  
 551 eters are optimized to improve the proposed driving anomaly  
 552 detection system (i.e., end-to-end solution). We use a constant  
 553 margin for the triplet loss function ( $\alpha$  in (8)). The value of  $\alpha$   
 554 needs to be adjusted during training. On the one hand, the loss  
 555 of the model will be very large if the margin is too large. Under  
 556 this setting, the model may not converge during the training  
 557 process. A benefit of having a large margin is that the model  
 558 will be more confident distinguishing similar samples. On the  
 559 other hand, the loss easily converges to 0 if the margin is too  
 560 small, which makes it more difficult for the model to distinguish  
 561 between similar samples. We implement the training process  
 562 with different values for this margin, varying  $\alpha$  from 2 to 25.  
 563 We evaluate the results on the development set, using the binary  
 564 classes *normal* and *candidate* sets. We set  $\alpha = 8$ , which led to  
 565 the best performance on the development set.

## 566 V. EXPERIMENTAL RESULTS

567 This section describes the experimental results of our pro-  
 568 posed unsupervised scalable multimodal driving anomaly de-  
 569 tection system. We also use subjective perceptual evaluation to  
 570 evaluate the model performance.

### 571 A. Driving Anomaly Detection

572 We evaluate model performance by comparing the anomaly  
 573 scores of the driving segments in *candidate* and *normal* sets  
 574 (Sec. IV-A). The annotations included in the videos from the  
 575 *candidate* set suggest something abnormal in the video, due to  
 576 the driver’s maneuvers, or the presence of other people, objects  
 577 or events (e.g., pedestrian crossing the street). Therefore, the  
 578 segments from the *candidate* set are expected to have higher  
 579 anomaly scores than the segments from the *normal* set, which  
 580 do not overlap with any annotation.

581 We compare the performance of our proposed model with  
 582 three baseline models. The first baseline is the CNN-LSTM con-  
 583 ditional GANs proposed by Qiu *et al.* [25], which is trained with  
 584 two modalities: the vehicle’s CAN-Bus signals and the driver’s  
 585 physiological signals. We refer to this method as *CNN-LSTM*  
 586 *GANs with 2 modalities*. This model concatenates the modalities  
 587 training a single conditional GAN model. This formulation  
 588 increases the dimension of the embeddings since it uses a single  
 589 concatenated representation. As we increase its dimension, the  
 590 model will require more data to effectively train this high di-  
 591 mensional feature representation. The convergence of the model  
 592 during training is compromised, as the dimension of the input  
 593 increases. Therefore, the approach is not scalable. In contrast, the  
 594 proposed method builds a separate GAN model for each modal-  
 595 ity, making it easier to train. It adopts an attention mechanism to  
 596 fuse separate embeddings from each modality. This formulation

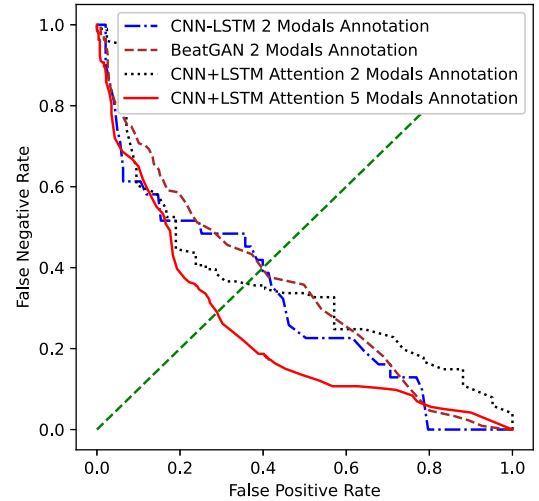


Fig. 6. DET curves for the models by formulating the problem as a binary classification task using the candidate and normal sets.

597 allows us to include more modalities if needed. The second  
 598 baseline is the BeatGAN proposed by Zhou *et al.* [30], which  
 599 is a GAN-based unsupervised method (see Sec. II-B). The gen-  
 600 erator of the BeatGAN model is built with an encoder-decoder  
 601 structure, and is trained to reconstruct 6-sec long signals as fake  
 602 data to confuse the discriminator. The discriminator is trained  
 603 to discriminate the real 6-sec signals and the generated fake  
 604 6-sec signals, following the regular adversarial training strat-  
 605 egy of GANs. For inference, the reconstruction error between  
 606 the real and fake signal is regarded as the anomalous metric  
 607 to discriminate abnormal events. In this work, for each 6-sec  
 608 long driving segment, we implement the BeatGAN framework  
 609 using the CAN-Bus and physiological data as input, using the  
 610 reconstruction error as the anomalous metric of the driving  
 611 segment. We refer to this method as *BeatGAN 2 modalities*. The  
 612 third baseline is the proposed attention model implemented with  
 613 only the CAN-Bus and the physiological signals. This baseline  
 614 is implemented to evaluate the effectiveness of the additional  
 615 modalities describing the external information. We refer to this  
 616 method as *attention with 2 modalities*. For the evaluation, we  
 617 formulate the driving anomaly detection problem as a binary  
 618 classification task. We calculate the *false positive rate* (FPR) and  
 619 *false negative rate* (FNR) as we change the decision threshold,  
 620 creating *detection error tradeoff* (DET) curves of the proposed  
 621 model and baseline models. This curve uses the FPR and FNR  
 622 as its axes. A DET curve that lies closer to the axes indicates  
 623 lower errors, and, therefore, better binary classification results.

624 Fig. 6 shows the DET curves of the proposed model and the  
 625 three baselines. The dashed line represents the operation point  
 626 where the FPR and FNR are equal. Fig. 6 indicates that the  
 627 proposed approach based on the attention model, implemented  
 628 with either two or five modalities, achieves better discriminative  
 629 performance than the CNN-LSTM GANs and BeatGAN models  
 630 for most of the operation points. Our proposed approach imple-  
 631 mented with the five modalities achieves the best performance,  
 632 indicating that adding the contextual information about the  
 633 road is extremely useful to improve the detection of driving  
 634 anomalies.

## 635 B. Subjective Perceptual Evaluation


636 This section relies on subjective perceptual evaluations to  
 637 assess more precisely the performance of the proposed approach.  
 638 Collectively, the videos from the *candidate* set are expected  
 639 to have more anomalies than the videos from the *normal* set.  
 640 However, it is possible that some of the videos in the *normal* set  
 641 may present some level of driving anomaly, while samples from  
 642 the *candidate* set may be normal. Therefore, we select videos in  
 643 the corpus to be directly annotated with anomaly scores.

644 We randomly select 200 segments from the *candidate* set and  
 645 200 segments from the *normal* set. The recording of each seg-  
 646 ments is six seconds long. Three annotators joined the perceptual  
 647 evaluation, who were asked to judge all the recordings after  
 648 watching the camera recordings showing the road. In addition  
 649 to annotating the driving anomalies, we are also interested on  
 650 the level of risk and familiarity perceived in the recordings.  
 651 Fig. 7 shows the *graphical user interface* (GUI). For each  
 652 driving segment, the annotators answered four questions about  
 653 the driving scenario shown in the video: (1) *how risky is the*  
 654 *driving condition in the video?* (safe; slightly risky; risky; very  
 655 risky), (2) *how often do you see similar driving condition on the*  
 656 *road?* (never; almost never; sometimes; quite often; regularly),  
 657 (3) *Is the driving condition in the video normal or abnormal*  
 658 (normal; abnormal), and (4) *what causes the anomaly in the*  
 659 *video?* (pedestrian; bicyclist; motorcyclist; other vehicle; bad  
 660 maneuver of our driver; no anomalies). The first three questions  
 661 consider a single choice. We estimate the inter-evaluator agree-  
 662 ment using the Krippendorff’s Alpha Coefficient, since these  
 663 questions have interval options. The agreement across the three  
 664 evaluators are 0.737 for question one (risky level), 0.509 for  
 665 question two (familiarity level), and 0.895 for question three  
 666 (normal/abnormal). The last question allows the annotators to  
 667 provide multiple choices as possible causes of the anomalies.  
 668 We estimate the inter-evaluator agreement using the Cohen’s  
 669 Kappa coefficient, since this question is multiple choice. This  
 670 metric is calculated between two raters, so we average the results  
 671 calculated from the three pairs of raters as the final agreement  
 672 level. The agreement for question four (possible causes) is 0.759.  
 673 These levels of agreements are considered very high. According  
 674 to the answers of the third question (i.e., Is the driving condition  
 675 in the video normal or abnormal?), we regroup the selected  
 676 400 driving segments into two sets: *normal* and *abnormal*. We  
 677 aggregate the responses of the annotators using the majority vote  
 678 rule, assigning a class if two out of the three evaluators select  
 679 that class. In total, we have 175 segments labeled as *abnormal*,  
 680 and 225 segments labeled as *normal*.

681 We analyze the risk level perceived in the annotated videos.  
 682 From the 400 segments, we select the top 100 segments with the  
 683 highest anomaly scores and the bottom 100 videos with the low-  
 684 est anomaly scores. A more discriminative model should have  
 685 more segments evaluated as *very risky* with fewer *safe* segments  
 686 in the *Top 100* group, and more *safe* segments with fewer *very*  
 687 *risky* segments in the bottom 100 group. Table II shows that the  
 688 top 100 group for the proposed attention model implemented  
 689 with five modalities has 45 segments labeled as either *risky*  
 690 or *very risky*. This number is higher than the corresponding  
 691 segments identified by the baselines: 38 for *CNN-LSTM GANs*,  
 692 42 for *BeatGAN*, and 40 for *Attention with 2 modalities*. Only  
 693 34 segments are selected as *safe*, which is less than the number  
 694 of segments selected by the other methods.

Please watch the video first. Then answer the questions. (Click to expand)

Video



0:01 -0:04

1. How risky is the driving maneuver in the video?
  - safe maneuver
  - slightly risky
  - risky maneuver
  - very risky maneuver
2. How often do you see similar driving maneuver on the road?
  - never
  - rarely
  - sometimes
  - quite often
  - regularly
3. Is the driving condition in the video normal or abnormal?
  - Normal
  - Abnormal
4. What causes the anomaly in the video?
  - Due to pedestrian
  - Due to bicyclist
  - Due to motorcyclist
  - Due to other cars
  - Due to bad maneuver of our driver
  - There is no anomaly shown in the video

Fig. 7. User interface of the subjective perceptual evaluation. After watching the video, the evaluators answer four questions to assess the risk, familiarity and anomaly levels (single choice). The questionnaire also asks for possible causes of anomalies (multiple choice).

TABLE II  
 ANALYSIS OF THE RISK LEVEL OF THE TOP 100 VIDEOS WITH THE HIGHEST ANOMALY SCORES AND THE BOTTOM 100 VIDEOS WITH THE LOWEST ANOMALY SCORES (IN BRACKET). THE ANALYSIS CORRESPONDS TO THE RESPONSES TO THE FIRST QUESTION IN THE PERCEPTUAL EVALUATION (FIG. 7). WE INDICATE IN BOLD THE MOST DESIRABLE RESULTS FOR THE EXTREME CASES

	safe	slightly risky	risky	very risky
CNN-LSTM GANs	41 ( <b>81</b> )	21 (10)	24 (8)	14 ( <b>1</b> )
BeatGAN	36 (71)	22 (18)	22 (8)	20 (3)
Attention with 2 modalities	37 (75)	23 (13)	21 (9)	19 (3)
Attention with 5 modalities	<b>34</b> (79)	21 (12)	24 (7)	<b>21</b> (2)



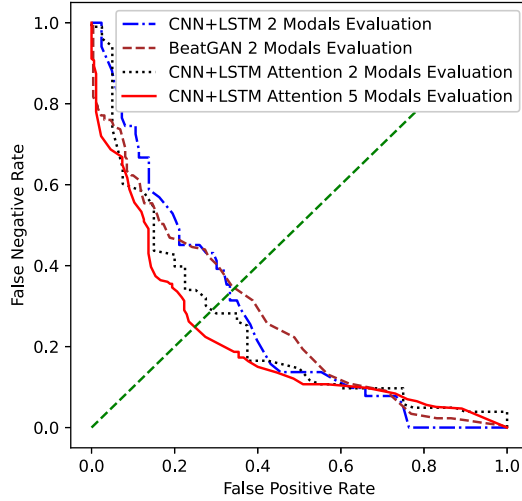


Fig. 8. DET curves for the models by formulating the problem as a binary classification task using the labels from the perceptual evaluations. The analysis relies on the responses to the third question in the perceptual evaluation (Fig. 7).

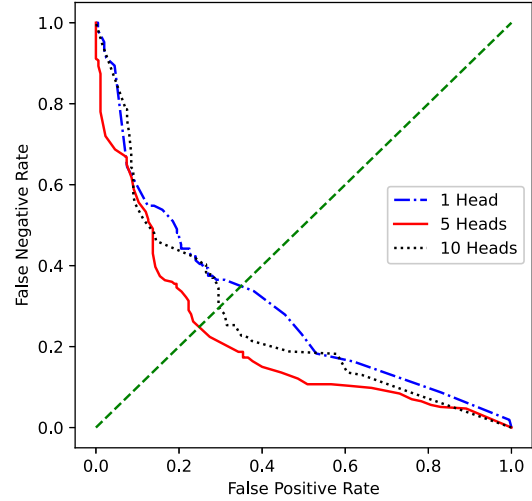


Fig. 9. DET curves to compare the discriminant performance of the proposed model based on attention implemented with different numbers of heads.

We also evaluate the familiarity level assigned to the annotated videos. We expect that videos with high anomaly scores are perceived as less frequently observed on the roads. From the top 100 videos with the highest anomaly scores, we observe that the proposed model implemented with five modalities has 49 videos labeled as either *never* or *rarely*. This number is also higher the corresponding values for the baselines: 38 for *CNN-LSTM GANs*, 46 for *BeatGAN*, and 39 for *Attention with 2 modalities*. The proposed approach is also the method with the lowest number of videos perceived as *regularly* observed on the roads (31 segments).

Fig. 8 shows the DET curves using the *normal* and *abnormal* labels obtained from the perceptual evaluation. In contrast to results on Fig. 6, which rely on annotations indirectly linked to driving anomaly, the results in Fig. 8 leverage the annotations conducted in this study to directly assess driving anomaly. The figure shows that our proposed model achieves the best performance. The proposed attention-based approach implemented with two modalities is better than the baseline method using only the conditional GAN model. These results confirm the observations made in Section V-A.

### C. Ablation Study

This section presents an ablation study to understand the contributions of different parts of the proposed model in the overall results. We report the performance by using the results from the perceptual evaluations, formulating the task as a binary classification task (i.e., *normal* versus *abnormal*).

A key component of the proposed approach is the attention model used to fuse the modalities. A parameter of the model is the number of heads ( $H$ ). This parameter is important, since it helps the system to attend to more than one modality. We implement the proposed approach with either one, five, or ten heads. Fig. 9 shows the corresponding DET curves. The model gets the best discriminant performance with five attention heads  $H = 5$ . The performance is clearly lower when we use a single head. In this case, the model can only attend to one of the modalities at a time, which is not optimal for this task. Adding

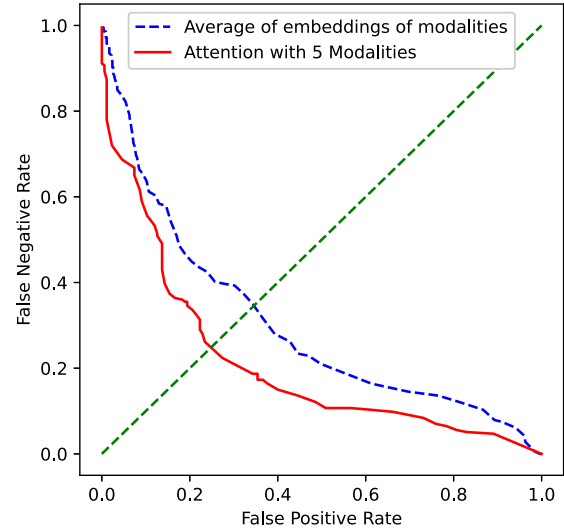


Fig. 10. DET curves to compare the discriminant performance of the proposed approach with and without attention model.

too many heads also is not optimal, especially since we only rely on five modalities.

To illustrate the effectiveness of the attention module in our approach, we remove the attention model, replacing the value with the average of the discriminator embeddings of each modality. Fig. 10 shows the results of this system with our full system with the attention model. The model with attention module outperforms the model without attention.

We explore the contribution of each of the modalities used in this study by adding one environmental modality to the proposed model trained with only CAN-Bus and physiological signals. Fig. 11 shows the corresponding DET curves. Adding environmental information to this baseline system improves the discriminative power of the system. Adding the pedestrian distances leads to more improvements. The figure also shows that we obtain the best performance when we consider the five modalities.

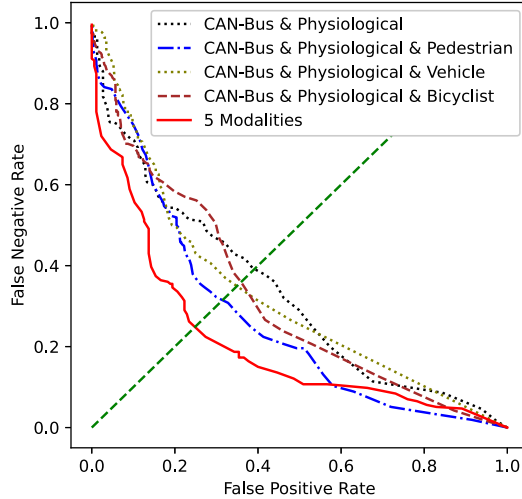


Fig. 11. DET curves to compare the discriminant performance of the proposed model based on attention implemented with different modalities.

TABLE III

ANALYSIS OF THE FAMILIARITY LEVEL OF THE TOP 100 VIDEOS WITH THE HIGHEST ANOMALY SCORES AND THE BOTTOM 100 VIDEOS WITH THE LOWEST ANOMALY SCORES (IN BRACKET). THE ANALYSIS CORRESPONDS TO THE RESPONSES TO THE SECOND QUESTION IN THE PERCEPTUAL EVALUATION (FIG. 7). WE INDICATE IN BOLD THE MOST DESIRABLE RESULTS FOR THE EXTREME CASES

	Never	Rarely	Sometimes	Quite often	Regularly
CNN-LSTM GANs	10 (1)	28 (6)	17 (10)	6 (8)	39 (75)
BeatGAN	<b>13 (1)</b>	33 (6)	13 (8)	8 (6)	33 ( <b>79</b> )
Attention with 2 modalities	11 (1)	28 (7)	16 (11)	9 (6)	36 (75)
Attention with 5 modalities	<b>13 (1)</b>	<b>36 (5)</b>	17 (10)	3 (7)	<b>31 (77)</b>

TABLE IV

NUMBER OF MILLIONS OF PARAMETERS WHEN ADDING MORE MODALITIES (UPTO SIX) TO THE BASE MODEL TRAINED WITH CAN-BUS AND PHYSIOLOGICAL SIGNALS

	CAN-Bus & Physiological						
	+0	+1	+2	+3	+4	+5	+6
	[M]	[M]	[M]	[M]	[M]	[M]	[M]
Qiu et al. [25]	<b>31.3</b>	<b>43.7</b>	<b>56.1</b>	<b>68.5</b>	80.9	93.3	105.7
Proposed approach	38.3	49.0	59.8	69.5	<b>77.1</b>	<b>84.7</b>	<b>92.3</b>
Attention module	1.31	1.38	1.44	1.51	1.57	1.64	1.71

#### 749 D. Scalability of the Model

750 This section focuses on the scalability of the proposed ap-  
 751 proach. We focus on the number of parameters in the models  
 752 as we increase the number of modalities. We assume that the  
 753 modalities that we add have input dimension equal to four,  
 754 similar to the distances to pedestrian, bicycles and other vehicles.  
 755 Table IV lists the number of millions of parameters when we add  
 756 more modalities to the base model trained with the CAN-Bus  
 757 and Physiological signals. Even though we considered three ad-  
 758 ditional modalities in this study (distances to pedestrian, bicycles  
 759 and other vehicles), we include in the analysis adding up to seven  
 760 extra modalities, each of them having a 4D representation. The  
 761 table lists the total number of parameters of the entire model,  
 762 and the number of parameters of the attention module. As a  
 763 reference, we also include the hypothetical scenario in which

764 we implement the CNN-LSTM GAN model [25] with more  
 765 modalities.

766 When we add four or more modalities, the results show that  
 767 the number of parameters is less than the model proposed by Qiu  
 768 *et al.* [25]. Most of the parameters added to our proposed model  
 769 correspond to the parameters needed to train a new separate  
 770 GAN model. The increase in the number of parameters of the  
 771 attention module is very small, as shown in the table. As a  
 772 result, the training of this model is scalable. We just need to  
 773 train a separate GAN model and retrain the attention model  
 774 block, which is minimally impacted by the new modality. In  
 775 contrast, the approach presented by Qiu *et al.* [25] needs to train  
 776 a single GAN model after concatenating all the inputs. The high  
 777 dimension of the input makes this single GAN difficult to train,  
 778 requiring more data to avoid undertraining the models. Because  
 779 of the high dimensionality of the model, the convergence of  
 780 the approach is also questionable. It is more convenient to train  
 781 a small GAN model for each modality than training one huge  
 782 GAN model with the concatenated inputs.

## 783 VI. CONCLUSION

784 This study introduced a novel unsupervised scalable mul-  
 785 timodal driving anomaly detection system based on the self-  
 786 attention mechanism, which is built on conditional GANs and  
 787 trained with the triplet loss function. This system builds a  
 788 separate conditional GAN model for each available modal-  
 789 ity, predicting the signal for the upcoming segment based on  
 790 previous data. The feature embeddings for the modalities are  
 791 fused by the attention model. The attention model is built based  
 792 on the self-attention mechanism and trained with triplet loss  
 793 function, where the distance between embeddings from actual  
 794 signals are minimized and embeddings from unpaired segments  
 795 are maximized. The entire training process does not require  
 796 labeled data. Our experimental results indicate that the proposed  
 797 model achieves better performance than the baseline models on  
 798 discriminating normal versus abnormal driving conditions.

799 The approach is scalable, where more modalities can be easily  
 800 added if needed. Our formulation only requires building separate  
 801 conditional GANs for the new modalities and concatenating  
 802 the corresponding feature representation to the input of the  
 803 attention model. Furthermore, the approach can react to driving  
 804 anomalies, even if the driver is not aware of the anomaly, by  
 805 incorporating modalities associated with the environment (i.e.,  
 806 distances to nearby pedestrians, vehicles and bicycles)

807 Our future work includes the integration of our approach  
 808 with new modalities such as lane keeping information or visual  
 809 attention estimation. The proposed approach relies on obtaining  
 810 physiological data, which currently requires wearable sensors.  
 811 The proposed model will benefit from non-contact technology  
 812 to estimate physiological data. Another limitation of the pro-  
 813 posed approach is the latency in the prediction. Our model  
 814 directly compares predicted and actual signals. This approach  
 815 introduces a latency of at least six seconds. A future research  
 816 direction is to investigate approaches to reduce the latency of  
 817 the model. Another appealing research direction is to increase  
 818 the interpretability of the model, identifying why the system  
 819 predicted that a given segment was anomalous. We expect that  
 820 the embeddings generated by individual GANs, or the join  
 821 embedding generated by the attention module can be used to  
 822 increase the interpretability of the model.

## REFERENCES

- 824 [1] P. Mohan, V. N. Padmanabhan, and R. Ramjee, "Nericell: Rich monitoring  
825 of road and traffic conditions using mobile smartphones," in *Proc. ACM*  
826 *Conf. Embedded Netw. Sensor Syst.*, Raleigh NC USA, 2008, pp. 323–336.
- 827 [2] C. Yang, A. Renzaglia, A. Paigwar, C. Laugier, and D. Wang, "Driving  
828 behavior assessment and anomaly detection for intelligent vehicles," in  
829 *Proc. IEEE Int. Conf. Cybern. Intell. Syst. IEEE Conf. Robot., Automat.*  
830 *Mechatronics*, Bangkok, Thailand, 2019, pp. 524–529.
- 831 [3] Z. Liu, M. Wu, K. Zhu, and L. Zhang, "SenSafe: A smartphone-based  
832 traffic safety framework by sensing vehicle and pedestrian behaviors,"  
833 *Mobile Inf. Syst.*, vol. 2016, pp. 1–13, Oct. 2016.
- 834 [4] J. Dai, J. Teng, X. Bai, Z. Shen, and D. Xuan, "Mobile phone based  
835 drunk driving detection," in *Proc. Int. Conf. Pervasive Comput. Technol.*  
836 *Healthcare*, Munich, Germany, 2010, pp. 1–8.
- 837 [5] H. Eren, S. Makinist, E. Akin, and A. Yilmaz, "Estimating driving behavior  
838 by a smartphone," in *Proc. IEEE Intell. Veh. Symp.*, Alcalá de Henares:  
839 Spain, 2012, pp. 234–239.
- 840 [6] C. Saiprasert and W. Pattara-Atikom, "Smartphone enabled dangerous  
841 driving report system," in *Proc. Hawaii Int. Conf. Syst. Sci.*, Wailea, Maui,  
842 HI, USA, 2013, pp. 1231–1237.
- 843 [7] J. Hong, B. Margines, and A. K. Dey, "A smartphone-based sens-  
844 ing platform to model aggressive driving behaviors," in *Proc. SIGCHI*  
845 *Conf. Hum. Factors Comput. Syst.*, Toronto, ON, Canada, 2014,  
846 pp. 4047–4056.
- 847 [8] Z. Chen, J. Yu, Y. Zhu, Y. Chen, and M. Li, "D3: Abnormal driving  
848 behaviors detection and identification using smartphone sensors," in *Proc.*  
849 *IEEE Int. Conf. Sensing, Commun., Netw.*, Seattle, WA, USA, 2015,  
850 pp. 524–532.
- 851 [9] J. Yu, Z. Chen, Y. Zhu, Y. Chen, L. Kong, and M. Li, "Fine-grained ab-  
852 normal driving behaviors detection and identification with smartphones,"  
853 *IEEE Trans. Mobile Comput.*, vol. 16, no. 8, pp. 2198–2212, Aug.  
854 2017.
- 855 [10] M. Fazeen, B. Gozick, R. Dantu, M. Bhukhiya, and M. C. González, "Safe  
856 driving using mobile phones," *IEEE Trans. Intell. Transp. Syst.*, vol. 13,  
857 no. 3, pp. 1462–1468, Sep. 2012.
- 858 [11] T. Chakravarty, A. Ghose, C. Bhaumik, and A. Chowdhury, "Mo-  
859 biDriveScore - A system for mobile sensor based driving analysis: A risk  
860 assessment model for improving one's driving," in *Proc. Int. Conf. Sens.*  
861 *Technol.*, Wellington, New Zealand, 2013, pp. 338–344.
- 862 [12] J. Wahlström, I. Skog, and P. Händel, "Risk assessment of vehicle corner-  
863 ing events in GNSS data driven insurance telematics," in *Proc. IEEE Conf.*  
864 *Intell. Transp. Syst.*, Qingdao, China, 2014, pp. 3132–3137.
- 865 [13] J. Wahlström, I. Skog, and P. Händel, "Detection of dangerous cornering  
866 in GNSS-data-driven insurance telematics," *IEEE Trans. Intell. Transp.*  
867 *Syst.*, vol. 16, no. 6, pp. 3073–3083, Dec. 2015.
- 868 [14] F. Li, H. Zhang, H. Che, and X. Qiu, "Dangerous driving behavior detection  
869 using smartphone sensors," in *Proc. IEEE Int. Conf. Intell. Transp. Syst.*,  
870 Rio de Janeiro, Brazil, 2016, pp. 1902–1907.
- 871 [15] P. Vavouranakis, S. Panagiotakis, G. Mastorakis, C. X. Mavromoustakis,  
872 and J. M. Batalla, "Recognizing driving behaviour using smartphones," in  
873 *Beyond the Internet of Things: Everything Interconnected*, J. Batalla, G.  
874 Mastorakis, C. Mavromoustakis, and E. Pallis, Eds., Cham, Switzerland:  
875 Springer, Jan. 2017, pp. 269–299.
- 876 [16] C. Ryan, F. Murphy, and M. Mullins, "End-to-end autonomous driving  
877 risk analysis: A behavioural anomaly detection approach," *IEEE Trans.*  
878 *Intell. Transp. Syst.*, vol. 22, no. 3, pp. 1650–1662, Mar. 2021.
- 879 [17] Y. Zhang, W. C. Lin, and Y. S. Chin, "A pattern-recognition approach for  
880 driving skill characterization," *IEEE Trans. Intell. Transp. Syst.*, vol. 11,  
881 no. 4, pp. 905–916, Dec. 2010.
- 882 [18] I. Mohamad, M. Ali, and M. Ismail, "Abnormal driving detection using  
883 real time global positioning system data," in *Proc. IEEE Int. Conf. Space*  
884 *Sci. Commun.*, Penang, Malaysia, 2011, pp. 1–6.
- 885 [19] A. Aljaafreh, N. Alshabat, and M. S. Najim al-din, "Sriving style  
886 recognition using fuzzy logic," in *Proc. IEEE Int. Conf. Veh. Electron.*  
887 *Saf.*, Istanbul, Turkey, 2012, pp. 460–463.
- 888 [20] L. Xu, S. Li, K. Bian, T. Zhao, and W. Yan, "Sober-drive: A smartphone-  
889 assisted drowsy driving detection system," in *Proc. Int. Conf. Comput.*  
890 *Netw. Commun.*, Honolulu, HI, USA, 2014, pp. 398–402.
- 891 [21] S. Ramyar, A. Homaifar, A. Karimodini, and E. Tunstel, "Identification  
892 of anomalies in lane change behavior using one-class SVM," in *Proc. IEEE*  
893 *Int. Conf. Syst., Man, Cybern.*, Budapest, Hungary, 2016, pp. 4405–4410.
- 894 [22] M. Zhang, C. Chen, T. Wo, T. Xie, M. Bhuiyan, and X. Lin, "SafeDrive:  
895 Online driving anomaly detection from large-scale vehicle data," *IEEE*  
896 *Trans. Ind. Informat.*, vol. 13, no. 4, pp. 2087–2096, Aug. 2017.
- [23] R. Chai *et al.*, "Driver fatigue classification with independent component  
by entropy rate bound minimization analysis in an EEG-based system,"  
*IEEE J. Biomed. Health Informat.*, vol. 21, no. 3, pp. 715–724, May  
2017.
- [24] 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 *Proc. Int. Conf. Comput.*  
*Eng. Syst.*, Cairo, Egypt, 2018, pp. 609–616.
- [25] Y. Qiu, T. Misu, and C. Busso, "Driving anomaly detection using condi-  
tional generative adversarial network," 2022, *arXiv:2203.08289*.
- [26] A. Vaswani *et al.*, "Attention is all you need," in *Proc. Adv. Neural Inf.*  
*Process. Syst.*, Long Beach, CA, USA, 2017, pp. 5998–6008.
- [27] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embed-  
ding for face recognition and clustering," in *Proc. IEEE Conf. Comput.*  
*Vis. Pattern Recognit.*, Boston, MA, USA, 2015, pp. 815–823.
- [28] Y. Qiu, T. Misu, and C. Busso, "Driving anomaly detection with conditional  
generative adversarial network using physiological and can-bus data," in  
*Proc. ACM Int. Conf. Multimodal Interact.*, Suzhou, Jiangsu, China, 2019,  
pp. 164–173.
- [29] Y. Qiu, T. Misu, and C. Busso, "Use of triplet loss function to improve driv-  
ing anomaly detection using conditional generative adversarial network,"  
in *Proc. Intell. Transp. Syst. Conf.*, Rhodes, Greece, 2020, pp. 1–7.
- [30] B. Zhou, S. Liu, B. Hooi, X. Cheng, and J. Ye, "BeatGAN: Anomalous  
rhythm detection using adversarially generated time series," in *Proc. Int.*  
*Joint Conf. Artif. Intell.*, Macao, China, 2019, pp. 4433–4439.
- [31] N. Li, J. Jain, and C. Busso, "Modeling of driver behavior in real world  
scenarios using multiple noninvasive sensors," *IEEE Trans. Multimedia*,  
vol. 15, no. 5, pp. 1213–1225, Aug. 2013.
- [32] N. Li, T. Misu, and A. Miranda, "Driver behavior event detection for  
manual annotation by clustering of the driver physiological signals," in  
*Proc. IEEE Int. Conf. Intell. Transp. Syst.*, Rio de Janeiro, Brazil, 2016,  
pp. 2583–2588.
- [33] S. Jha and C. Busso, "Head pose as an indicator of drivers' visual  
attention," in *Vehicles, Drivers, and Safety, ser. Intelligent Vehicles and*  
*Transportation*, vol. 2, H. Abut, J. Hansen, G. Schmidt, and K. Takeda,  
Eds., De Gruyter, Berlin, Germany, May 2020, pp. 113–132.
- [34] S. Jha and C. Busso, "Estimation of driver's gaze region from head position  
and orientation using probabilistic confidence regions," *IEEE Trans. Intell.*  
*Vehicles*, to be published, doi: [10.1109/TIV.2022.3141071](https://doi.org/10.1109/TIV.2022.3141071).
- [35] N. Li and C. Busso, "Analysis of facial features of drivers under cognitive  
and visual distractions," in *Proc. IEEE Int. Conf. Multimedia Expo.*, San  
Jose, CA, USA, 2013, pp. 1–6.
- [36] H. Yang, L. Liu, W. Min, X. Yang, and X. Xiong, "Driver yawning detection  
based on subtle facial action recognition," *IEEE Trans. Multimedia*, vol. 23,  
pp. 572–583, 2021.
- [37] W. Huang, X. Liu, M. Luo, P. Zhang, W. Wang, and J. Wang, "Video-  
based abnormal driving behavior detection via deep learning fusions,"  
*IEEE Access*, vol. 7, pp. 64 571–64 582, 2019.
- [38] O. Köpiklül, J. Zheng, H. Xu, and G. Rigoll, "Driver anomaly detection:  
A dataset and contrastive learning approach," in *Proc. IEEE/CVF Winter*  
*Conf. Appl. Comput. Vis., Virtual Conf.*, 2021, pp. 91–100.
- [39] W. Song, Y. Yang, M. Fu, F. Qiu, and M. Wang, "Real-time obstacles detec-  
tion and status classification for collision warning in a vehicle active safety  
system," *IEEE Trans. Intell. Transp. Syst.*, vol. 19, no. 3, pp. 758–773,  
Mar. 2018.
- [40] Y. Yao, M. Xu, Y. Wang, D. J. Crandall, and E. M. Atkins, "Unsupervised  
traffic accident detection in first-person videos," in *Proc. IEEE/RSJ Int.*  
*Conf. Intell. Robots Syst.*, Macau, China, 2019, pp. 273–280.
- [41] Y. Yao, X. Wang, M. Xu, Z. Pu, E. Atkins, and D. Crandall, "When,  
where, and what? A new dataset for anomaly detection in driving videos,"  
pp. 1–17, Apr. 2020, *arXiv:2004.03044*.
- [42] H. Kim, J. Park, K. Min, and K. Huh, "Anomaly monitoring framework in  
lane detection with a generative adversarial network," *IEEE Trans. Intell.*  
*Transp. Syst.*, vol. 22, no. 3, pp. 1603–1615, Mar. 2021.
- [43] I. Goodfellow *et al.*, "Generative adversarial nets," in *Proc. Adv. Neural*  
*Inf. Process. Syst.*, Montreal, Canada, vol. 27, 2014, pp. 2672–2680.
- [44] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li, "On the continuity of rotation  
representations in neural networks," in *Proc. IEEE/CVF Conf. Comput. Vis.*  
*Pattern Recognit.*, Long Beach, CA, USA, 2019, pp. 5738–5746.
- [45] D. Li, D. Chen, B. Jin, L. Shi, J. Goh, and S.-K. Ng, "MAD-GAN: Multi-  
variate anomaly detection for time series data with generative adversarial  
networks," in *Proc. Artif. Neural Netw. Mach. Learn., Text Time Ser.*, Series  
Lecture Notes in Computer Science, I. Tetko, V. Kůrková, P. Karpov, and F.  
Theis, Eds., Munich, Germany, Springer, vol. 11730, 2019, pp. 703–716.



971 [46] A. Geiger, D. Liu, S. Alnegheimish, A. Cuesta-Infante, and K. Veera-  
 972 machaneni, "TadGAN:Time series anomaly detection using generative  
 973 adversarial networks," in *Proc. IEEE Int. Conf. Big Data*, Atlanta, GA,  
 974 USA, 2020, pp. 33–43.

975 [47] S. Hyland, C. Esteban, and G. Rätsch, "Real-valued (medical) time series  
 976 generation with recurrent conditional GANs," 2017, *arXiv:1706.02633*.

977 [48] S. Akcay, A. Atapour-Abarghouei, and T. Breckon, "GANomaly: Semi-  
 978 supervised anomaly detection via adversarial training," in *Proc. Asian  
 979 Conf. Comput. Vis.*, Series Lecture Notes in Computer Science, C. Jawa-  
 980 har, H. Li, G. Mori, and K. Schindler, Eds., Perth, Australia, Springer,  
 981 vol. 11363, 2018, pp. 622–637.

982 [49] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R.  
 983 Chandrasekhar, "Efficient GAN-based anomaly detection," 2018,  
 984 *arXiv:1802.06222*.

985 [50] C. Hori *et al.*, "Attention-based multimodal fusion for video description,"  
 986 in *Proc. IEEE Int. Conf. Comput. Vis.*, Venice, Italy, 2017, pp. 4203–4212.

987 [51] H. Chen, D. Jiang, and H. Sahli, "Transformer encoder with multi-modal  
 988 multi-head attention for continuous affect recognition," *IEEE Trans. Mul-  
 989 timedia*, vol. 23, pp. 4171–4183, 2021.

990 [52] Y. Qiu, T. Misu, and C. Busso, "Analysis of the relationship between  
 991 physiological signals and vehicle maneuvers during a naturalistic driving  
 992 study," in *Proc. Intell. Transp. Syst. Conf.*, Auckland, New Zealand, 2019,  
 993 pp. 3230–3235.

994 [53] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural  
 995 Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.

996 [54] T. Misu and Y. Chen, "Toward reasoning of driving behavior," in *Proc. Int.  
 997 Conf. Intell. Transp. Syst.*, Maui, HI, USA, 2018, pp. 204–209.

998 [55] V. Ramanishka, Y.-T. Chen, T. Misu, and K. Saenko, "Toward driving  
 999 scene understanding: A dataset for learning driver behavior and causal  
 1000 reasoning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Salt Lake  
 1001 City, UT, USA, 2018, pp. 7699–7707.

1002 [56] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers:  
 1003 Surpassing human-level performance on ImageNet classification," in *Proc.  
 1004 IEEE Int. Conf. Comput. Vis.*, Santiago, Chile, 2015, pp. 1026–1034.



**Teruhisa Misu** (Member, IEEE) received the B.E., 1017  
 M.E., and Ph.D. degrees in information science from 1018  
 Kyoto University, Kyoto, Japan, in 2003, 2005, and 1019  
 2008, respectively. From 2005 to 2008, he was a 1020  
 Research Fellow (DC1) of the Japan Society for the 1021  
 Promotion of Science (JSPS). From 2008 to 2013, 1022  
 he was a Researcher with NICT Spoken Language 1023  
 Communication Group. In 2013, he joined Honda 1024  
 Research Institute USA, Inc. From November 2011 1025  
 to February 2012, he was a Visiting Researcher with 1026  
 USC/ICT. 1027  
 1028



**Carlos Busso** (Senior Member, IEEE) received the 1029  
 B.S. and M.S. degrees (with high honors) in electrical 1030  
 engineering from the University of Chile, Santiago, 1031  
 Chile, in 2000 and 2003, respectively, and the Ph.D. 1032  
 degree in electrical engineering from the University of 1033  
 Southern California (USC), Los Angeles, CA, USA, 1034  
 in 2008. He is currently an Associate Professor with 1035  
 the Electrical Engineering Department, The Universi- 1036  
 ty of Texas at Dallas (UTD), Richardson, TX, USA. 1037  
 His research focuses on human-centered multimodal 1038  
 machine intelligence and applications. His current re- 1039

search interests include broad areas of affective computing, multimodal human-  
 machine interfaces, nonverbal behaviors for conversational agents, in-vehicle  
 active safety system, and machine learning methods for multimodal processing.  
 He was selected by the School of Engineering of Chile as the best electrical  
 engineer graduated in 2003 across Chilean universities. At USC, he was the  
 recipient of the Provost Doctoral Fellowship from 2003 to 2005 and Fellowship  
 in Digital Scholarship from 2007 to 2008. At UTD, he leads the Multimodal  
 Signal Processing (MSP) Laboratory. He was the recipient of the NSF CAREER  
 Award. In 2014, he was the recipient of the ICMI Ten-Year Technical Impact  
 Award. In 2015, his student was the recipient of the third prize IEEE ITSS Best  
 Dissertation Award (N. Li). He was also the recipient of the Hewlett Packard Best  
 Paper Award at the IEEE ICME 2011 (with J. Jain), and Best Paper Award at the  
 AAC ACII 2017 (with Yannakakis and Cowie). He received the Best of IEEE  
 Transactions on Affective Computing Paper Collection in 2021 (with R. Lotfian).  
 He is the coauthor of the winner paper of the Classifier Sub-Challenge event at  
 the Interspeech 2009 emotion challenge. His work has direct implication in  
 many practical domains, including national security, health care, entertainment,  
 transportation systems, and education. He was the General Chair of ACII 2017  
 and ICMI 2021. He is a Member of ISCA, AAAC, and a Senior Member of  
 ACM.

1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016



**Yuning Qiu** (Student Member, IEEE) received the  
 B.S. degree in electrical engineering from the Harbin  
 Institute of Technology, Harbin, China, in 2016, and  
 the M.S. degree in electrical engineering from Boston  
 University, Boston, MA, USA, in 2018. He is cur-  
 rently working toward the Ph.D. degree in electrical  
 engineering with the University of Texas at Dallas,  
 Richardson, TX, USA. In 2018, he joined Multimodal  
 Signal Processing (MSP) Laboratory. His research  
 interests include the area of in-vehicle safety system,  
 human-machine interaction, and machine learning.

1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060