

On Transferability of Driver Observation Models from Simulated to Real Environments in Autonomous Cars

Walter Morales-Alvarez^{1*}  Member, IEEE, Novel Certad^{1*}  Graduate Student Member, IEEE, Alina Roitberg^{2*}  Member, IEEE, Rainer Stiefelhagen³  Member, IEEE and Cristina Olaverri-Monreal¹  Senior Member, IEEE

Abstract—For driver observation frameworks, clean datasets collected in controlled simulated environments often serve as the initial training ground. Yet, when deployed under real driving conditions, such simulator-trained models quickly face the problem of distributional shifts brought about by changing illumination, car model, variations in subject appearances, sensor discrepancies, and other environmental alterations.

This paper investigates the viability of transferring video-based driver observation models from simulation to real-world scenarios in autonomous vehicles, given the frequent use of simulation data in this domain due to safety issues. To achieve this, we record a dataset featuring actual autonomous driving conditions and involving seven participants engaged in highly distracting secondary activities. To enable direct SIM→REAL transfer, our dataset was designed in accordance with an existing large-scale simulator dataset used as the training source. We utilize the Inflated 3D ConvNet (I3D) model, a popular choice for driver observation, with Gradient-weighted Class Activation Mapping (Grad-CAM) for detailed analysis of model decision-making. Though the simulator-based model clearly surpasses the random baseline, its recognition quality diminishes, with average accuracy dropping from 85.7% to 46.6%. We also observe strong variations across different behavior classes. This underscores the challenges of model transferability, facilitating our research of more robust driver observation systems capable of dealing with real driving conditions.

I. INTRODUCTION

To speed up the development of driver observation systems, researchers often leverage simulated environments for collecting the training data [1], [2], [3], [4]. These environments provide a controlled and easily reproducible setting, which allows for the collection of clean datasets, avoiding the challenges associated with real-world complexities. Furthermore, when studying driver behaviors during autonomous or highly automated driving [1], [5], [6], [7], [8], safety becomes a significant challenge, leading to simulators being very prominent. However, deploying models trained in simulation to real-world scenarios presents a significant challenge, as the distributional shifts between the two environments can lead to poor generalization and performance

*Denotes equal contribution.

¹Chair ITS-Sustainable Transport Logistics 4.0, Johannes Kepler University Linz, Altenberger Straße 69, 4040 Linz, Austria. {novel.certad.hernandez, walter.morales.alvarez, cristina.olaverri-monreal}@jku.at

²Institute for AI, University of Stuttgart, 70569 Stuttgart, Germany. alina.roitberg,@f05.uni-stuttgart.de

³Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany. rainer.stiefelhagen@kit.edu

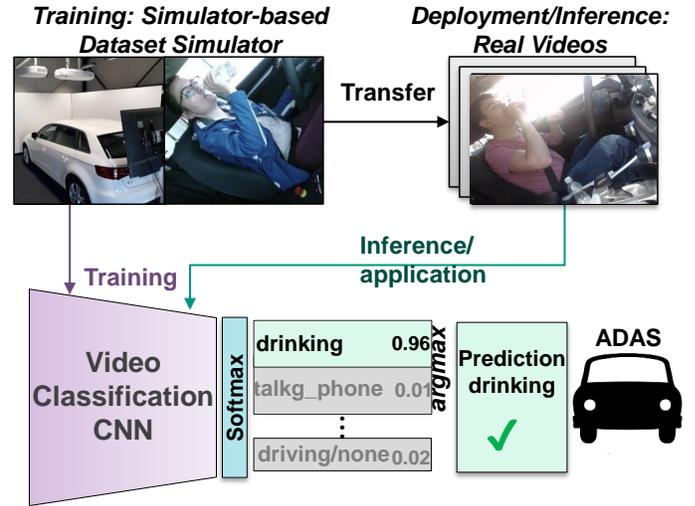


Fig. 1: We collect the first video-based driver activity recognition dataset featuring secondary activities in *real autonomous driving* scenarios. Unlike the controlled simulated environments, our recordings include real-world complexities, such as car movement and fluctuating lighting conditions. Combined with a large-scale training database collected in a simulator, we utilize our dataset as a testbed for validating direct SIM→REAL transfer of deep learning-based driver observation models.

degradation. These shifts can be caused by various factors, such as changing lighting conditions, discrepancies in car models and sensors, and other environmental variations. Some works have presented results from real world tests to validate and show potential discrepancies regarding results from simulated data, e.g., for take over requests [9].

To address these issues and improve the transferability of driver observation models, this paper investigates the efficacy of different CNN-based approaches in bridging the domain gap between simulated and real environments (an overview is provided in Figure 1). We first introduce the validation testbed, a collected dataset for video-based driver monitoring in autonomous vehicles, and the training dataset generated from a simulated environment. Through a series of experiments, we provide a thorough evaluation of the performance of these models in real driving scenarios.

By offering valuable insights into the challenges associated with transferring models from simulation to real-world

scenarios, we aim to contribute to the development of more robust and reliable driver observation systems that can be deployed in real driving conditions.

II. RELATED WORK

Models for recognizing driver activity from video fall into two main categories: ones that utilize manual feature descriptors, and ones that leverages end-to-end deep learning, processing video inputs directly and concurrently learning intermediate representations alongside the classifier. Traditional methodologies based on manual features [10], [11], [12], [13], [10], [13] exploit classical machine learning models like Support Vector Machines and Random Forest. These first extract features related to the driver’s hand movements, body and head posture, and eye direction. Recently, end-to-end deep learning models have become an increasingly prevalent choice for recognizing driver activities. These models often utilize Convolutional Neural Networks (CNNs)[2], [1], [14], [15], [16], [17] and transformer-based models[16] as their backbone. While neural network-based approaches are at the top of most driver observation benchmarks, they require a significant amount of annotated training data. Several real-world datasets are tailored for in-vehicle observation during *manual* driving [18], [19], [10]. However, for *autonomous* driving, simulators are more commonly used [1], [5], [6], [7]. For instance, Drive&Act [1], the largest public dataset for video-based driver observation during autonomous driving, was collected in a stationary car placed indoors, surrounded by three screens imitating outdoor surroundings. Conversely,

research from the broader field of human activity recognition reports a substantial decline in recognition quality when transitioning from synthetic to real data [20]. Another line of work focuses on domain adaptation for driver observation, such as cross-view recognition or model adaptation for participants wearing glasses [21], [22]. Given real-life training data and resources for post-hoc model adjustment, these domain adaptation approaches show promising potential to significantly enhance recognition, presenting an important future research direction. However, this falls outside the scope of our study, as such methods require additional unlabeled training data in the target domain and cannot be applied directly. Overall, recent research in video-based driver observation tends to prioritize the development of high-accuracy classifiers for conditions similar to training environments, while performance under distributional shifts or adverse conditions is often considered secondary.

Inspired by this, we present an empirical evaluation of direct SIM→REAL transfer of deep learning-based activity recognition models in the context of autonomous driving. To address the gap in suitable SIM→REAL benchmarks, we first collect a real-world dataset of in-vehicle observation during actual autonomous driving, annotated with a subset of activities present in a large-scale simulator-based dataset, facilitating the aforementioned validation scenario.

III. VALIDATING SIMULATOR-BASED DRIVER OBSERVATION MODELS IN REAL ENVIRONMENTS

A. Testbed: collected dataset for video-based driver monitoring in autonomous vehicles

The dataset used in this study consisted of videos obtained from the scientific personnel of Johannes Kepler University Linz. The data collection was carried out utilizing the JKU-ITS vehicle (see Figure 2), as described in [23], where participants engaged in various tasks while the vehicle autonomously navigated through a designated test lane within the university premises.

The data collection setup resembled that of the work presented in [24], where participants were required to activate the vehicle’s automation and commence task performance, while the automation system assumed control of the vehicle throughout the experiment.

In our study, the vehicle automation involved two distinct processes. The first process encompassed the drive-by-wire capability, implemented using Openpilot algorithms [25]. By utilizing the Black Panda device, acceleration and steering commands were transmitted to the vehicle via the internal ADAS (Advanced Driver Assistance System), facilitated by a ROS Wrapper that exposed the Black Panda’s communication protocols to ROS topics. The second process involved a custom ROS2 high-level controller of the vehicle, responsible for generating trajectories, speed profiles, and steering and acceleration commands based on pre-recorded waypoints obtained through the vehicle’s GPS.

To ensure safety during the experiment, a security driver was present in the passenger seat, overseeing the proper functioning of the system. In case of emergencies, the

TABLE I: Overview of the recorded video dataset of secondary activities during real autonomous driving sessions.

Dataset statistics	
Context	Real driving session
Manual driving	✓
Autonomous driving	✓
Data type	RGB video
Number of subjects	7
Nr. female subjects	2
Nr. driver activities	7
Recording lengths (min)	100.83
Number of samples*	1987



Fig. 2: Autonomous vehicle used for the data collection.

security driver could assume manual control of the vehicle using a joystick to apply the brakes.

For video recording of the participants, a Logitech C920 webcam was utilized. The webcam was positioned on the passenger door, capturing the entirety of the participants’ body movements.

Dataset statistics. Table I provides an overview of a recorded video dataset of secondary activities during real autonomous driving sessions. All secondary activities are recorded during fully autonomous driving, except for the *driving/sitting still* activity, which also included sequences of the subject steering manually. The dataset includes seven subjects in total, two female and five male. These subjects are recorded engaging in seven different driver activities: driving/sitting still, using a phone, talking on the phone, reading a magazine, reading a newspaper, reading a book, and drinking. The duration of the recorded data is 100.83 minutes. Following the procedure of [1] a single sample to be classified is defined as a 3-second video clip that is labeled with a specific activity. The objective of the recognition model is therefore to accurately tag each 3-second or shorter action segment (for events of lesser duration) with the appropriate activity label. The dataset comprises 1987 such annotated samples. Note, that our dataset is not intended for training, but for validation of the SIM→REAL transfer of modern neural networks trained on simulator-based data.

B. Validation Protocol and Recognition Model

Training dataset collected in a simulator. As the simulator-based training dataset, we leverage Drive&Act[1], the largest public in-vehicle human activity dataset focused on distractive behavior during both, manual and autonomous driving. The data is collected from 15 subjects and is annotated with 34 fine-grained activities at the main evaluation level. To maintain label correspondences, we select 6 categories present in Drive&Act. In addition, we collect the *reading book* activity, which was not present in Drive&Act in this form and is therefore interesting for looking at the networks behavior in the case of a new object (book). At the same time, Drive&Act contains similar activities - *reading newspaper* and *reading magazine* and an ideal model would map the new *reading book* situation to one of these states

Neural Architecture. We utilize the Inflated 3D architecture (I3D)[26], an extension of the Inception-v1 network, renowned in video classification and driver observation fields[1]. The I3D adapts the 2D filters of Inception-v1 into a temporal dimension and processes 64-frame video snippets of 224x224 resolution. I3D consists of 27 layers, with nine Inception modules executing parallel convolutions and concatenating the output to ensure computational efficiency. We use the original model [1] trained on the Drive&Act split 1 (200 epochs, SGD at a learning rate of 0.05 and momentum of 0.9, pre-training on Kinetics) and exclude the activity labels not present in our dataset in the last fully-connected layer.

C. Model Attribution Analysis

To examine the way simulator-based video classification models operate when facing real-life driving data, we leverage the Gradient-weighted Class Activation Map method (GradCAM) [27]. However, the original approach was designed for static images [27], while we are additionally dealing with the time dimension. We, therefore, implement a three-dimensional variant of GradCAM similar to [28].

Given an input video, we predict a class c , then estimate the gradient over y_c with respect to each value in the feature channel A_k .

The *importance* w_c^k is then estimated separately for each channel k by averaging the gradients:

$$w_c^k = \frac{1}{n} \sum_{i,j,t} \left(\frac{\partial y_c}{\partial A_k^{i,j,t}} \right), \quad (1)$$

We then calculate final weights $V_c^{i,j,t}$ by applying the *ReLU* function to the linear combination of feature map values and importance estimates:

$$V_c^{i,j,t} = ReLU \left(\sum_k w_c^k A_k^{i,j,t} \right). \quad (2)$$

For visualization, we average the heat-maps over time.

IV. VALIDATION RESULTS

A. Main quantitative results

We evaluate the recognition quality of driver activities collected in a simulator versus a real-life self-driving car, as presented in Table II. The performance of each model is evaluated based on accuracy.

For all driver activities, a random chance baseline yields an accuracy of 14.29% in both the simulator and real-life scenarios. The Inflated 3D architecture (I3D), when trained on simulator data, attains an overall accuracy of 85.7% in the simulator and 46.56% in real-life observations. Notably, the I3D model significantly outperforms the random baseline in both settings, demonstrating its efficacy.

However, the per-category results for the simulator-trained I3D reveal some variability. In the simulator, high accuracy is observed for activities such as *driving/sitting still* (99.01%), *using a phone* (90.36%), *reading a newspaper* (98.03%), and *drinking* (100%). However, the accuracy decreases considerably in a real-life environment, with the highest scores being *reading a book* (69.49%) and *drinking* (60.4%). Notably, *talking on the phone* and *reading a magazine* activities witness the most substantial drops in accuracy, scoring merely 8.52% and 20.34%, respectively, in real-life conditions. To facilitate the comparison between the different accuracy results, they are graphically represented in Figure 3 as a radar diagram.

This analysis underscores the difficulties in transferring models trained in simulated environments to real-world conditions, especially for specific activities. It also suggests the need for further fine-tuning and optimization of the I3D model for real-life driver observation.

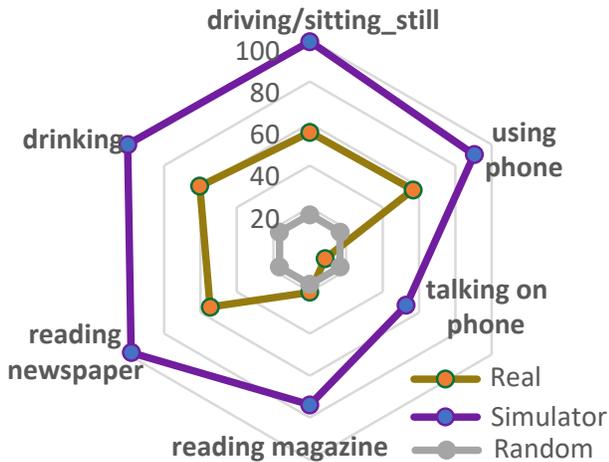


Fig. 3: Accuracy of the individual categories visualized a radar diagram. We compare the recognition quality of the I3D model trained on a simulator-based dataset and validated on (1) data originating from the same simulator (2) in-vehicle observation videos during real autonomous driving sessions collected by us, and (3) random chance baseline.

Next, we analyze the most common confusions in Figure 6 and the confusion matrix (Figure 4). Nearly all categories are often mistaken for *driving/sitting still*, which is not surprising considering its overrepresentation in the original simulator-based training set [1]. The category itself is recognized correctly in 56% of the cases. It is unsurprising that *reading a magazine* is most frequently confused with *reading a newspaper* (39%), while the confusion in the opposite direction is less common (7%). *Talking on the phone* is often mistaken for *drinking* (27%), which is understandable as both actions involve raising one hand close to the mouth. The difficulty in fine-grained recognition of smaller objects, such as phones, may arise from the downsampling performed by 3D CNNs, which rapidly decreases image resolution to obtain larger receptive fields. Undoubtedly, *talking on the phone* is a highly common and significant distractive secondary activity. Accurately recognizing such fine-grained actions from images holds immense importance for the future. Similarly, like other behaviors, the most frequent confusion for *talking on the phone* is with the overrepresented category of *driving/sitting still* (46%). While *using a phone* is better recognized than *talking on the phone*, with a 57% accuracy in predictions, confusions with *driving/sitting still* still occur relatively frequently (17%).

Considering this information, it is clear that improvements are needed, especially with regard to fine-grained activities involving smaller objects, particularly when differentiating activities from the “default” *driving/sitting still* state.

B. Attribution analysis and qualitative examples

In Figure 5, we visualize an example of a single driving session collected in our dataset. The upper bar depicts the true secondary activities of the driver, while the lower bar depicts the predictions of the I3D model trained on

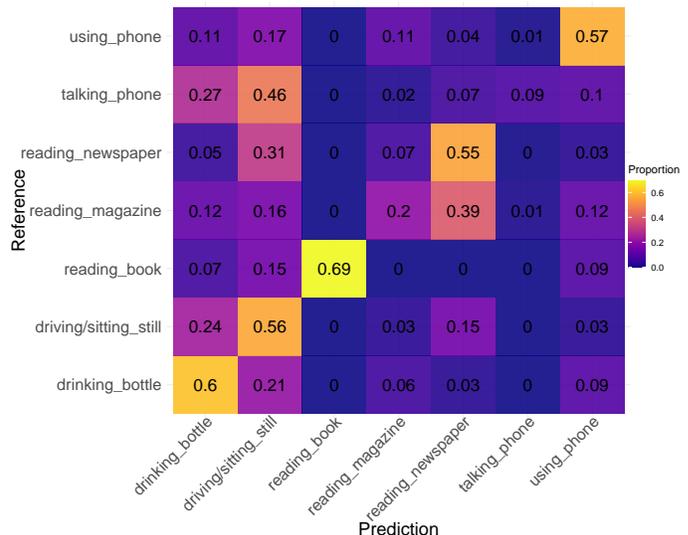


Fig. 4: Confusion matrix for the direct SIM to REAL transfer.

TABLE II: Recognition Quality for Driver Observation collected in a Simulator vs Real-life Self-Driving Car

Activity	Accuracy	
	Simulator	Real
All Driver Activities		
Random chance baseline	14.29	14.29
Simulator-trained I3D	85.7	46.56
Per-category Results for Simulator-trained I3D		
driving/sitting still	99.01	55.76
using phone	90.36	56.7
talking phone	52.94	8.52
reading magazine	73.91	20.34
reading newspaper	98.03	54.72
reading book	-	69.49
drinking	100	60.4

simulator data. The model had issues in recognizing the first two driver behaviors (*using phone* and *reading magazine*), although there were certain brief segments with the correct classification. Secondary behaviors that followed were easier to recognize, and the majority of the frames were assigned the correct label. A surprising observation is that, for this particular subject, *talking on phone* was recognized better than *using phone*. This is contrary to the overall trend observed when examining the statistics of the entire dataset. These findings highlight that individual appearances and the unique manner in which humans perform actions can significantly influence the quality of recognition.

In our subsequent analysis, we explore the specific input pixels that influence the model’s decisions, leveraging the GradCAM technique [27] adapted for our spatiotemporal data (see Section III-C). Figure 7 presents frames that have been correctly classified (top row) and those misclassified (bottom row), with an overlay of a heatmap that highlights the pixels that contributed to the current network decision. The analysis unveils that the video classification model predominantly focuses on the hands and objects in use, particularly when the predictions are accurate. Interestingly, in the

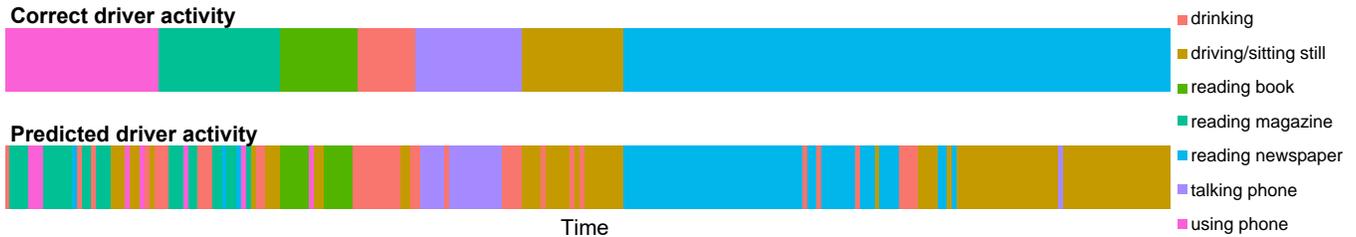


Fig. 5: Example of a single driving session recorded in our dataset. The subject is engaged in six different secondary activities during fully autonomous driving. A brief manual driving segment is also present. The upper bar represents the ground truth activity labels, while the lower bar displays the predictions of an end-to-end CNN trained on a simulator-based dataset. Different colors mark different secondary activities.

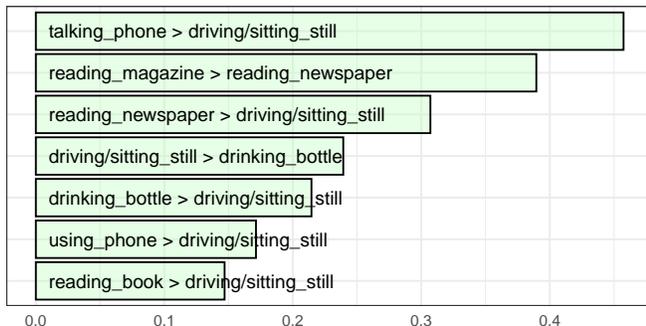


Fig. 6: Most common confusions of an end-to-end CNN for driver observation after the direct SIM→REAL transfer.

correctly classified *talk on phone* example (Figure 7b), the attention does not center on the hands, but gravitates towards the area around the wrist. This observation suggests that the network places importance not on the object per se, but rather on a specific wrist position, potentially accounting for the subpar recognition of this category. Conversely, when *talk on phone* is erroneously identified as *drinking* (Figure 7g), the model’s attention is indeed trained on the driver’s hands. This evidence suggests that the model’s comprehension of *talk on phone* and *drinking* is centered around certain wrist and hand patterns, discounting the significance of the associated objects. We also observe, especially in failure cases, that the model’s attention can deviate from the hands or the relevant object, steering towards the mid-cabin area or certain outside patterns, presumably in response to unusual movements absent from the simulator training data. In conclusion, the GradCAM analysis offers significant insights into the way simulator-based video classification models operate when facing naturalistic driving data. Especially for categories that are hard to recognize, we discover a strong reliance of the model on hand and wrist movements, potentially at the expense of object recognition. This might be due to the different appearance of these objects present in the training set. Enhancing the training data with a broader range of phones, drinking bottles, cups, and similar objects could address this limitation. Furthermore, we observe instances where the model’s focus shifts to external movement during misclassifications (Figure 7h). This suggests that the absence

of such movement in the simulator data adversely impacts the model’s recognition capability. Moving forward, we advocate for the inclusion of such variable movements in the training data. This could be accomplished by recording more naturalistic datasets or developing sophisticated data augmentation methods that effectively mimic these car movements.

V. CONCLUSION

We collected a video-based dataset for driver activity recognition during real autonomous driving sessions. Our key motivation is to study the direct SIM→REAL transfer of deep learning-based driver observation models, which is of particular relevance given that simulated data is a prevalent resource in autonomous driving research. Our dataset features seven drivers engaged in six distractive activities as well as a short manual driving segment. Furthermore, the dataset is constructed with annotations and sensor correspondence to a large-scale simulator-based dataset, specifically designed to supplement the validation protocol of simulator-trained models with real-world data. While the model clearly surpasses the random baseline, its recognition quality drops drastically when moving from simulated to real-life recordings, highlighting the necessity of incorporating real-world data into validation protocols of simulator-trained models.

Acknowledgements. This work was supported by the Austrian Ministry for Climate Action, Environment, Energy, Mobility, Innovation, and Technology (BMK) Endowed Professorship for Sustainable Transport Logistics 4.0., IAV France S.A.S.U., IAV GmbH, Austrian Post AG and the UAS Technikum Wien. Alina Roitberg was partially supported by the KHYS Connecting Young Scientists travel award and Deutsche Forschungsgemeinschaft (DFG) under Germany’s Excellence Strategy - EXC 2075.

REFERENCES

- [1] M. Martin *et al.*, “Drive&act: A multi-modal dataset for fine-grained driver behavior recognition in autonomous vehicles,” in *ICCV*, 2019.
- [2] D. Tran, H. M. Do, J. Lu, and W. Sheng, “Real-time detection of distracted driving using dual cameras,” in *IROS*, 2020.
- [3] J. S. Katrolija, B. Mirbach, A. El-Sherif, H. Feld, J. Rambach, and D. Stricker, “Ticam: A time-of-flight in-car cabin monitoring dataset,” *arXiv preprint arXiv:2103.11719*, 2021.
- [4] O. Kopuklu, J. Zheng, H. Xu, and G. Rigoll, “Driver anomaly detection: A dataset and contrastive learning approach,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 91–100, 2021.

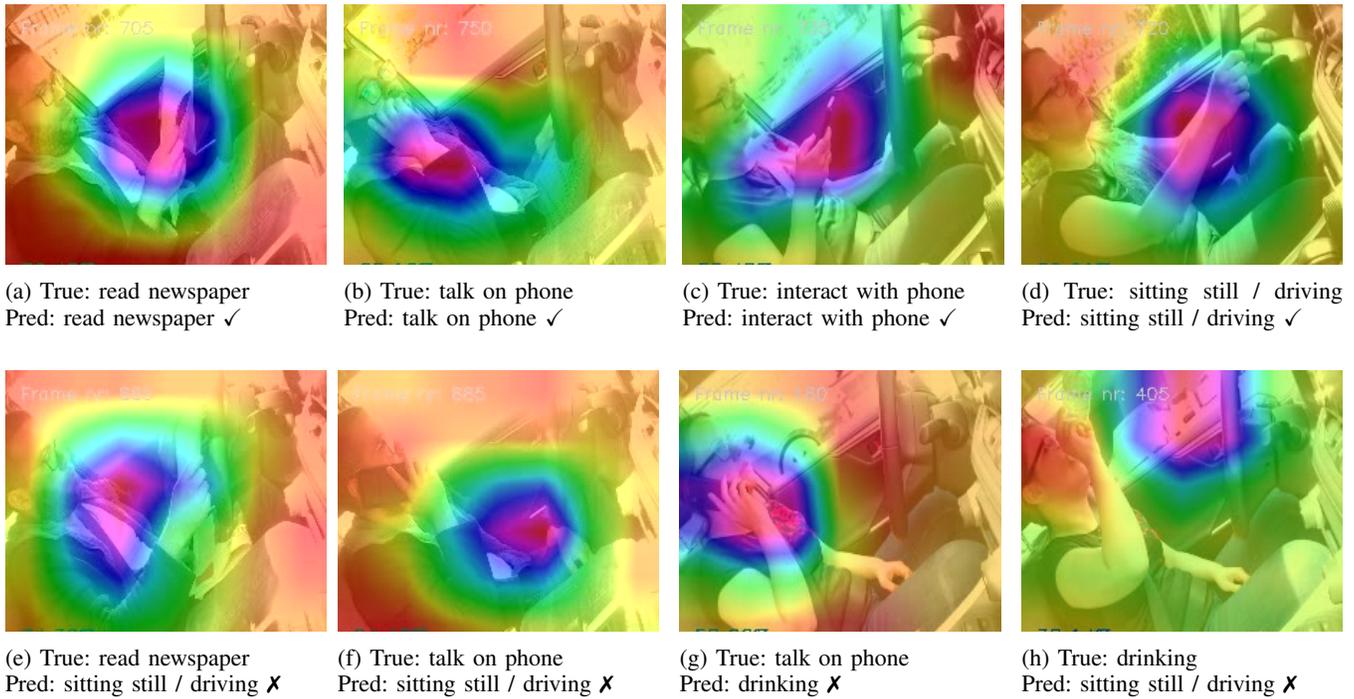


Fig. 7: Qualitative results and attribution analysis: image portions with high activation values in the intermediate network layers are highlighted as a heatmap. These regions are estimated via a spatiotemporal implementation of GradCAM [27].

- [5] W. Morales-Alvarez, N. Certad, H. H. Tadjine, and C. Olaverri-Monreal, "Automated driving systems: Impact of haptic guidance on driving performance after a take over request," in *IV*, pp. 1817–1823, IEEE, 2022.
- [6] C. Gold, D. Damböck, L. Lorenz, and K. Bengler, "'take over!' how long does it take to get the driver back into the loop?," in *Proceedings of the human factors and ergonomics society annual meeting*, vol. 57, pp. 1938–1942, Sage Publications Sage CA: Los Angeles, CA, 2013.
- [7] B. Mok, M. Johns, K. J. Lee, D. Miller, D. Sirkin, P. Ive, and W. Ju, "Emergency, automation off: Unstructured transition timing for distracted drivers of automated vehicles," in *ITSC*, pp. 2458–2464, IEEE, 2015.
- [8] M. Flad, P. Karg, A. Roitberg, M. Martin, M. Mazewitsch, C. Lange, E. Kenar, L. Ahrens, B. Flecken, L. Kalb, *et al.*, "Personalisation and control transition between automation and driver in highly automated cars," *Smart Automotive Mobility: Reliable Technology for the Mobile Human*, pp. 1–70, 2020.
- [9] W. Morales-Alvarez, O. Sipele, R. Léberon, H. H. Tadjine, and C. Olaverri-Monreal, "Automated driving: A literature review of the take over request in conditional automation," *Electronics*, vol. 9, no. 12, p. 2087, 2020.
- [10] E. Ohn-Bar, S. Martin, A. Tawari, and M. M. Trivedi, "Head, eye, and hand patterns for driver activity recognition," in *ICPR*, 2014.
- [11] L. Xu and K. Fujimura, "Real-time driver activity recognition with random forests," in *AutomotiveUI*, 2014.
- [12] R. Zheng, K. Nakano, H. Ishiko, K. Hagita, M. Kihira, and T. Yokozeki, "Eye-gaze tracking analysis of driver behavior while interacting with navigation systems in an urban area," *THMS*, 2016.
- [13] C. Braunagel, E. Kasneci, W. Stolzmann, and W. Rosenstiel, "Driver-activity recognition in the context of conditionally autonomous driving," in *ITSC*, 2015.
- [14] M. Tan *et al.*, "Bidirectional posture-appearance interaction network for driver behavior recognition," *T-ITS*, 2021.
- [15] A. Roitberg, K. Peng, D. Schneider, K. Yang, M. Koulakis, M. Martinez, and R. Stiefelwagen, "Is my driver observation model overconfident? input-guided calibration networks for reliable and interpretable confidence estimates," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 25271–25286, 2022.
- [16] K. Peng, A. Roitberg, K. Yang, J. Zhang, and R. Stiefelwagen, "Transarc: Transformer-based driver activity recognition with latent space feature calibration," in *IROS*, pp. 278–285, IEEE, 2022.
- [17] A. Roitberg, C. Ma, M. Haurilet, and R. Stiefelwagen, "Open set driver activity recognition," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1048–1053, IEEE, 2020.
- [18] A. Jain, A. Singh, H. S. Koppula, S. Soh, and A. Saxena, "Recurrent neural networks for driver activity anticipation via sensory-fusion architecture," in *ICRA*, pp. 3118–3125, IEEE, 2016.
- [19] Y. Xing, C. Lv, H. Wang, D. Cao, E. Velenis, and F.-Y. Wang, "Driver activity recognition for intelligent vehicles: A deep learning approach," *IEEE transactions on Vehicular Technology*, vol. 68, no. 6, pp. 5379–5390, 2019.
- [20] A. Roitberg, D. Schneider, A. Djmal, C. Seibold, S. Reiß, and R. Stiefelwagen, "Let's play for action: Recognizing activities of daily living by learning from life simulation video games," in *IROS*, pp. 8563–8569, IEEE, 2021.
- [21] S. Reiß, A. Roitberg, M. Haurilet, and R. Stiefelwagen, "Deep classification-driven domain adaptation for cross-modal driver behavior recognition," in *IEEE IV*, pp. 1042–1047, IEEE, 2020.
- [22] A. Rangesh, B. Zhang, and M. M. Trivedi, "Driver gaze estimation in the real world: Overcoming the eyeglass challenge," in *IEEE IV*, pp. 1054–1059, IEEE, 2020.
- [23] N. Certad, W. Morales-Alvarez, G. Novotny, and C. Olaverri-Monreal, "Jku-its automobile for research on autonomous vehicles," in *Computer Aided Systems Theory – EUROCAST 2022* (R. Moreno-Díaz, F. Pichler, and A. Quesada-Arencibia, eds.), (Cham), pp. 329–336, Springer Nature Switzerland, 2022.
- [24] A. Valle-Barrio, W. Morales-Alvarez, C. Olaverri-Monreal, and J. Naranjo-Hernandez, "Development and validation of an open architecture for autonomous vehicle control," in *IEEE Intelligent Vehicles Symposium*, (Anchorage, Alaska, USA), 2023.
- [25] G. Holtz, "Comma ai," 2018. (Accessed on May 01, 2023).
- [26] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *CVPR*, IEEE, 2017.
- [27] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, and D. Batra, "Grad-cam: Why did you say that?," *ICCV*, 2017.
- [28] A. Roitberg, M. Haurilet, S. Reiß, and R. Stiefelwagen, "CNN-based driver activity understanding: Shedding light on deep spatiotemporal representations," in *ITSC*, 2020.