

Driver Cellphone Usage Detection Using Wavelet Scattering and Convolutional Neural Networks

Ali Besharati¹, Ali Nahvi¹, and Serajeddin Ebrahimian¹

¹ Virtual Reality Laboratory, K.N. Toosi University of Technology, Tehran, Iran
ali_besharati@email.kntu.ac.ir
nahvi@kntu.ac.ir
sebrahimian@alumni.kntu.ac.ir

Abstract. Cellphone usage by drivers is a major cause of road accidents. This paper provides an automated system based on machine learning and computer vision to detect cellphone usage during driving. We used Wavelet Scattering Networks, which is a simple and efficient type of architecture. The presented model is straightforward and compact and requires little hyper-parameter tuning. The speed of this model is similar to the Convolutional Neural Networks. We monitored the driver from 2 viewpoints: a frontal view of the driver's face and a side view of the driver's whole body. We created a new dataset for the first viewpoint, and a publicly available dataset for the second viewpoint. Our model achieved the test accuracy of 91% for our new dataset and 99% for the publicly available one.

Keywords: Mobile use detection, Wavelet Scattering Network, CNN, Cascade Object Detector, Transfer Learning

1 Introduction

Phone usage by drivers causes distraction and deviates attention from road hazards. About one-third of the accidents caused by distraction involve a collision with a pedestrian, in which the driver had identified the pedestrian too late [1]. Some studies report that inattention far more often contributes to fatal crashes than to less severe crashes [1]. Accidents' severity correlates with the type of inattention [2]–[4]. Inattention is the root of a third of severe accidents between the years 2011 and 2015 in Norway [5]. Inability to perceive information in blind spots or behind obstacles is a common type of distraction [5]. 80% of the traffic accidents and 65% of the close collisions of the drivers were due to the distraction of the drivers for a period of 3 seconds before the accident [5]. On average, it takes 5 seconds to make a call, and if the car is moving at a speed of 100 km/h, it will travel 140 m during this time [6], so every pedestrian or vehicle in this distance is a potential cause of an accident. In Brazil, the fine for using

mobile phones has increased by 150% in 5 years [7]. Using cell phones and texting while driving is responsible for 28% of accidents in the United States [8]. Every hour of the day, at least 5% of American drivers are using their cell phones while driving, which increases the risk of a crash 4-6 times [9], [10].

The purpose of Advanced Driver Assistance Systems (ADAS) is that driver error will be reduced or even eliminated, and efficiency in traffic and transport is enhanced [11]. The technological capability for level-5 autonomous vehicles (AVs) will be available by 2025 but the commercial availability will be yet to happen until 2050 or even later [12]. This delay in availability has two reasons: first, the lack of substructures, like highway infrastructure [13] or 5G used for AVs [14], and second, legislation challenges, like the legislative evolutions required in terms of administrative, civil, and criminal law [15]. ADAS is making progress in both of these spheres. Some examples of ADAS are compensating for the lack of attention for the driving like cruise control [16], automatic car parking [17], and collision avoidance systems [18] or assisting the driver or correcting his behavior via warnings like drowsiness detection systems [19], [20] and drunkenness detection systems [21]. Such setups could be used for legislation or insurance related purposes [22]. Likewise, our system is an ADAS that could be used for warning, analyzing, and reporting the drivers' behavior to authorities or insurance companies.

Distraction is defined as the deviation of the driver's focus to a marginal task instead of the basic activities of driving [23]. Driver distraction can be divided into 4 types: visual, auditory, manual and cognitive [24]. Distraction caused by mobile phone includes all 4 types: looking at the phone screen (visual distraction), talking (auditory), texting (manual) and thinking about the discussed issues (cognitive) [24]. Sometimes drivers have compensatory beliefs [24]: They believe that A risky behavior is neutralized by performing a safe behavior. So the person thinks: Now I can use the mobile phone because I drive at a low speed. Such a belief is dangerous or causes damage in many cases even at low speed [24].

Several other studies also tried proposing methods for drivers' cellphone usage detection. Wagner et al [25] installed a head tracking camera (OptiTrack V120:Duo) in a test car on top of the dashboard's center console. Also, they used 2 IR cameras (Flir FL3-U3- 13Y3M-C) and with active IR illumination (two Osram SFH 4725S LEDs per side) to have a bright view of the driver's face. They claimed that their system remains fully functional during day and night light. They used a Raspberry Pi 3 Model B+ for image acquisition and brightness control. They presented 2 datasets for head position and driver posture. They made use of 4 different CNN architectures for classification using different cameras. He et al [26] established a dataset of 545 relevant images. For classification, they proposed the CornetNet-Lite network which is based on the Hourglass module backbone to localize the cellphone in the image. They utilized 4 different types of noise to measure the robustness of their model. Their training procedure included penalizing negative positions that are in a certain radius to positive positions. Rajput et al [27] extended the goal of "Driver cellphone usage detection" to cellphone detection in all scenarios. For this purpose, they used the State Farm Distracted Driver Detection Kaggle dataset [28] as well as their own dataset named "IITH-dataset on mobile phone usage" (IITH-DMU) which has 618 images of 6000x4000 resolution.

They annotated all images in both datasets and used Faster-RCNN and SSD with different backbones for transfer learning. The backbones they used were pre-trained on the COCO dataset. Likewise, Xiong et al [29] proposed a deep-learning-based approach. Firstly, they used Progressive Calibration Networks [30] to detect and track the face and determine the candidate's mobile phone usage region. They used the YOLOV3 [31] network for the classification part. They also introduced their own experimental data which was collected from 50 drivers under different lighting conditions through an infrared camera. It included 22216 calling behavior images and 25464 normal images. For the training of the neural network, they used the NVIDIA jetson tx2 artificial intelligence supercomputer module. Behati et al [32] utilized the VGG16 architecture for this purpose. They used the Distracted Driver Dataset [33] to train their classifier. They also presented a shrunk version of VGG16 which is multiple times smaller than the original one but yields 1% less accuracy. Le et al [34] performed a robust deep-learning method to detect this problem. They proposed a Multiple Scale Region-based CNN (RCNN) to detect the presence of a mobile phone near the face as well as the position of the hands and the steering wheel to ensure that two hands are placed on the steering wheel. As well as observing the steering wheel, they employed the same method to observe the face and localize the hand near the face. They used SHRP-2 Database and VIVA Hand Database for their experiment. Torres et al [35] also used a CNN trained on the first 4 classes of the State Farm dataset to perform the classification.

While recent studies in this field mostly rely on deep learning based approaches, older studies mostly employed classical methods for this classification. Xu et al [36] used a Near Infra-Red (NIR) imaging system described in [37] located on the roadside. Since the infrared light is absorbed by some vehicles windshield, they incorporated a powerful IR illuminator to compensate for the loss of light. Also, they used a long pass filter ($>750\text{nm}$) for blocking the visible light of the illuminator to reduce its visual impact on the driver. They used a deformable part model to localize the windshield region and sent the whole windshield or a sub-region of it to a trained classifier for violation detection. They utilized Fisher vectors (FV) representation to classify the driver's side of the windshield into cell phone usage violation and non-violation classes. Elqattan et al [38] used CNNs to diagnose this problem. They used UI-5240CP Rev.2 IDS camera with a zoom lens on the police cars and an Axis Q1645 camera on the roadside for this purpose. They made use of the YOLO algorithm trained on the COCO dataset to localize the car and the driver. To categorize the images, they cut and pre-processed each received image using contrast limited adaptive histogram equalization for enhancement and 3×3 Gaussian kernel for blurring. Then they utilized the transfer learning method to categorize the images. They used Xception [39] model pre-trained on ImageNet [40] dataset for transfer learning. Berri et al [41] proposed a hybrid system of pattern recognition and movement detection based on the fact that when drivers use their cell phones, they tend to fix their gaze on a point in front which limits the field of vision and affects the drivability. They implemented the detection using a camera attached to the dashboard of a car. First, their movement detection process classified the camera feed into 3 classes: "motionless", "cell phone to the ear", and "withdraw cell phone". Then, their pattern recognition process evaluated the image. The movement detection parametrized the pattern recognition by setting the classification threshold based on its output. They

utilized a dataset of 100 positive images and 100 negative images. They used an MLP neural network with Gaussian and Sigmoid activation functions and applied a binary-coded genetic algorithm to find the best parameters for their network. Artan et al [42] used a near-infrared camera that is directed toward the windshield of the car. Their method consisted of two stages. First, they applied a Wiener filtering with 3×3 windows to mitigate the impact of image acquisition noise due to low light and then located the position of the face area on the windshield by a Deformable Part Model (DPM). Secondly, by using an image separation method based on the local accumulation in the area obtained around the driver's face, they recognized the mobile phone. They used a histogram of oriented gradient (HOG) for feature extraction and a Support Vector Machine (SVM) for classification. Seshadri et al [43] used a supervised descent method (SDM) and Viola-Jones algorithm to track different points for investigating the desired area of the face, then the resulting images were given to the pre-trained classifier. For feature extraction, they utilized normalized raw pixels as well as features obtained from the histogram of oriented gradients. For the classification part, they used an ensemble of the Real Adaboost Classifier and SVM. Yaser et al [44] used a dataset of 49 positive and 30 negative images. Their images were taken from outside of a vehicle. They used cascade detectors for the classification part and reached the accuracy of 75%.

In this paper, first, we used a type of architecture known as Wavelet Scattering Network (WSN). Second, we showed that this network achieves higher accuracy like the widely used CNNs but using fewer parameters and relying only on strong feature extraction property of the wavelet transform. Moreover, it achieved high accuracy without applying image augmentation techniques to the images. Third, we produced a new dataset based on frontal view images of drivers using their cellphones.

The structure of this paper is as follows: In section 2 we provide a short description of the WSN and present the architecture used for this study. Also, we describe the datasets used for this study. In section 3 we explain the results of our method and compare the performance of our WSN with the widely used CNN backbones on the 2 datasets. In section 4 we compare our contribution with those of previous studies and in section 5 we provide the conclusion of the paper.

2 Materials and Methods

In this paper, we present a method based on machine learning to automatically detect driver cellphone usage. We created a dataset for this purpose as well as using a publicly available dataset to improve the system's generalization. For classification, we used WSN and compared its results with the widely used CNN architectures. We exhibit that our WSN architecture reaches the accuracy of the State of the Art (SOTA) models with its simple architecture and fewer number of parameters.

2.1 Ear Dataset

Like any other machine learning algorithm, our method requires a tremendous amount of data to learn the difference between "phone usage" and "safe behavior" [45]. First,

we apply a HAAR cascade object detector [46] to the image to locate the subject's face and measure its size based on the distance from the camera. Then, using the face location and its size, we estimate the potential regions of the ears (Figure 1). There is a linear relationship between human face size and hand size [47]. So after measuring the coefficients 0.4 and 1.2 for one subject, we can be certain that they will be constant for the other subjects. It enables us to robustly track the regions of interest (ROI)s, especially when the subject gets closer or farther to the camera.

Image acquisition protocol. We had the subject sitting at a distance of 40-80cm from the camera. Hence, they held a cell phone in one hand beside the corresponding ear and gradually leaned to the left and right and also rotated their necks to the side while we took their images. We further asked the subject to repeat the same procedure with the same side of the face but using the opposite hand to have all possible cases in which a driver can use a phone while driving. The subjects were supposed to perform the moves gradually and we took 2 images from either side every second.

We only saved regions of ears to minimize memory usage and further processing. We took positive images from one hand and negative ones from the other. Due to the symmetry of the problem, there was no need to run the protocol on the other hand. Every image, whether positive or negative, could be flipped to represent that state for the opposite ear. (This will be further discussed in the next subsection)

We gathered 18537 images from 48 subjects, 41 males and 7 females, using the protocol. All images were taken with the subject's consent with a variety of indoor lighting setups.

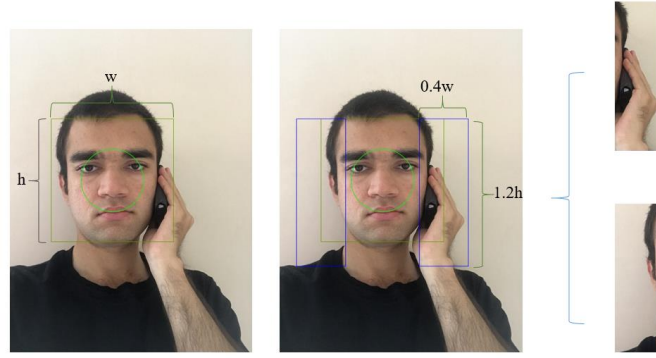


Fig. 1. Left. Subject holding a phone in the left hand. The cascade detector locates the head (the circle) and measures the size (the square). Middle. We estimate the potential regions where the subject may have a phone based on the locations extracted by the cascade detector. Right. Extracted images for classifier decision

Image Augmentation. First, we resized all images to 150x50 (Figure 2). Then, we applied several image augmentation techniques to investigate the unsearched regions in the data space (Figure 3) [48]. Due to the horizontal symmetry of the human body, we applied the horizontal flip transform with a high probability (0.9).

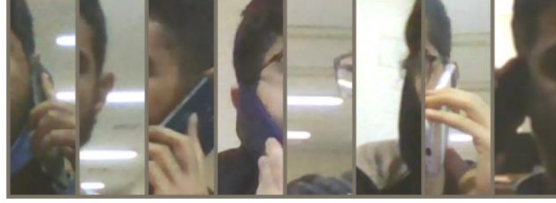


Fig. 2. Raw images used for the validation set.



Fig. 3. Transformed images used for the training set.

2.2 State Farm Datasets

We also used the first 5 classes of the publicly available Kaggle State Farm Dataset [28] (Figure 4). The total number of images used was 11745.



Fig. 4. The first 5 classes of Kaggle State Farm Dataset [26]

2.3 Classifier

We implemented a WSN [49] for each of the datasets. The WSN is a computationally efficient network based on wavelet filters that is invariant to translation and rotation [49]. The scattering transform requires 3 parameters for the transformation: number of samples or input size (T), the averaging scale (2^J), and the number of wavelets per octave, Q [50]. For our problem, T is the size of the images. We chose the numbers 3 and 8 for J and Q respectively. Figure 5 depicts the architectures we used for either of the datasets.

We also used two CNNs to have a comparison amongst different architectures. We claim that our WSN achieves the same accuracy as SOTA architectures but with far fewer parameters.

We resized all images in Ear Dataset to 150x50 and all images in State Farm dataset to 64x64. We used grayscale images for the WSN and RGB augmented images for the CNNs. Both datasets were split in 2 (90% of the data was used for training and 10% for validation) and we shuffled the data during training to avoid overfitting. We used the

backbones VGG16 [51] and ResNet50 [52] for this experiment. We trained all networks for 20 epochs with Adam optimizer.

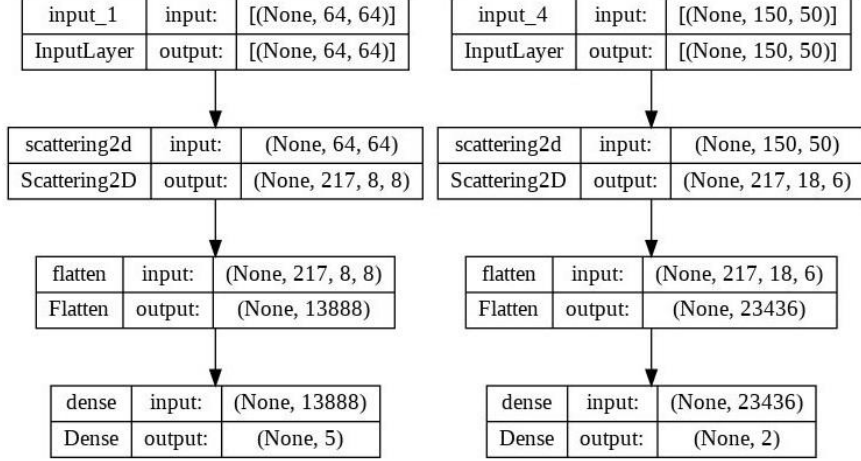


Fig. 5. WSN Architectures used for the classification of: Left. State Farm data classifier and Right. Ear data Classifier

2.4 Tools

All the codes were written in Python. The fundamental libraries were OpenCV for image processing tasks and PyTorch, Tensorflow, and Keras for machine learning tasks [53].

All of the trainings of the Neural Networks were done on using Google CoLaboratory™ data analysis tool and its Tesla T4 GPU. For inferencing, we used a laptop computer with an Intel Core i7 CPU, 8 GB RAM, and a webcam.

3 Results

For the ear dataset, the WSN achieved an accuracy of 91%, a precision of 87%, a sensitivity of 94%, and a specificity of 89%. Regarding the State Farm dataset, it achieved an accuracy of 99%, a sensitivity of 99%, and a precision of 99%. Figure 6 portrays the confusion matrix for the Ear dataset and the State Farm dataset classification respectively. It reveals the high accuracy and precision achieved by the WSN.

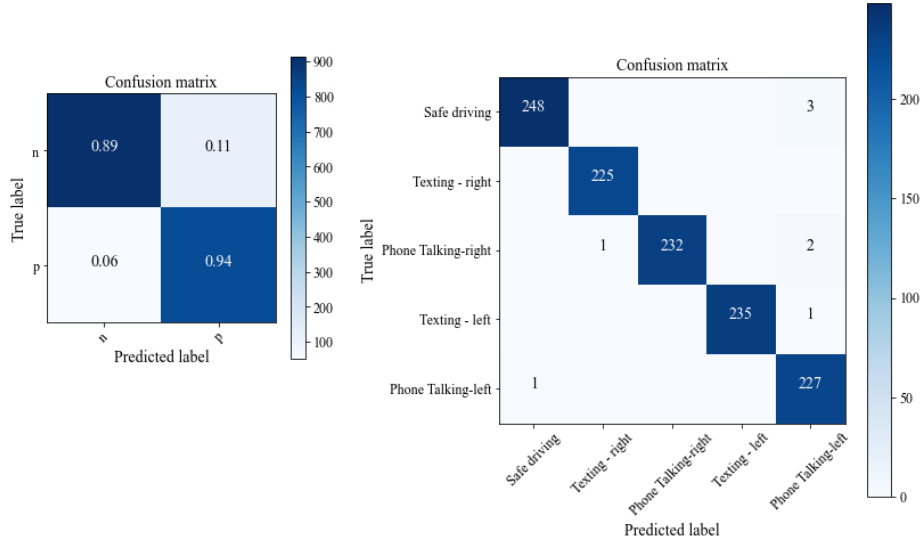


Fig. 6. Confusion matrices for the classification of Ear dataset (Left) and the State Farm dataset (Right)

The CNNs we used reached 99% accuracy. Figures 7 and 8 compare the accuracy and loss of the CNNs with the WSN for the Ear dataset and State Farm dataset respectively for training for 20 epochs. The high validation accuracy reached by the WSN shows that it surmounted overfit to a high degree. For the Ear dataset, due to high similarity between classes, it falls behind the CNNs in terms of validation accuracy but for the State Farm dataset, it accomplishes the same validation accuracy.

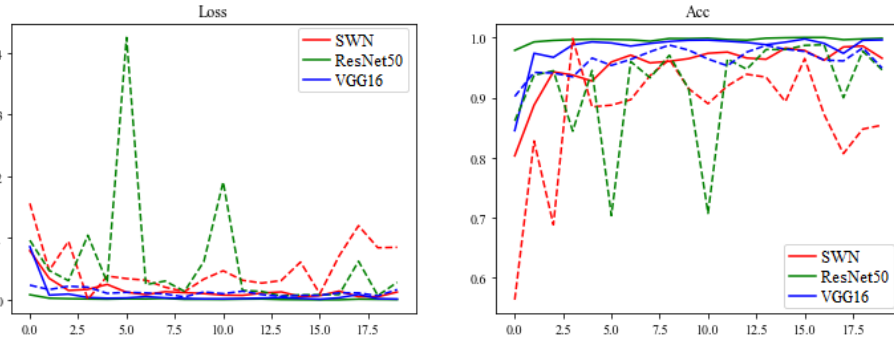


Fig. 7. Loss and Accuracy of WSN, Resnet50 and VGG16 for the Ear dataset (Dashed lines are showing the validation loss and accuracy)

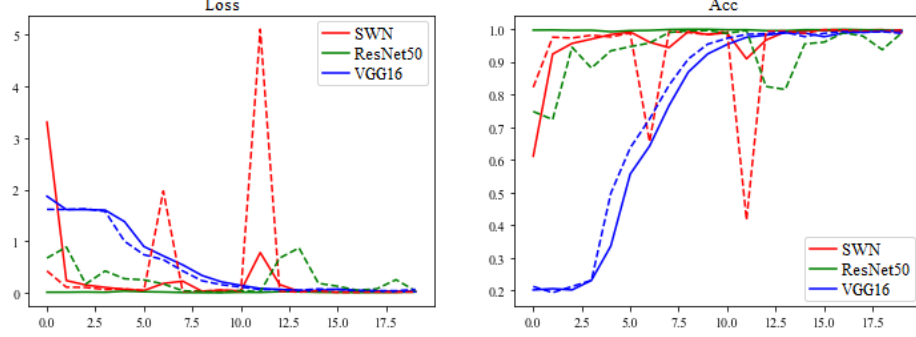


Fig. 8. Loss and Accuracy of WSN, Resnet50 and VGG16 for the State Farm dataset (Dashed lines are showing the validation loss and accuracy)

Table 1 provides the utility metrics for the models trained on Ear Dataset and State Farm Dataset. It provides a comparison between the WSN and the CNNs in terms of volume and speed and shows that WSN is many times smaller than the other networks.

Table 1. Utility metrics

Model	Ear Dataset		State Farm Dataset	
	fps	(MB)	fps	(MB)
		Volume		Volume
WSN	5	0.57	5	0.84
VGG16	3	174.9	7	179.1
Resnet50	8	330.7	8	301.8

4 Discussions

In this section, we elaborate on our results and compare our method with previous studies. Table 2 compares the performance of this paper with other research works. It indicates that this study achieved the correctness metrics of the previous studies.

Table 3 compares the backbones used in other studies. It shows that the number of parameters in WSN is more than 300 times less than the smallest model used by previous authors (ResNeXt-50 in [25]) for this classification. Also, it is about 12 times smaller in volume. Most of the recent studies tried to tackle the problem using CNNs and achieved SOTA results. Our method gained the same accuracy, but with far fewer parameters. Additionally, as illustrated in figure 5, its 2-layer architecture is a lot simpler than the CNNs and does not require pooling or drop out layers.

Table 2. Correctness metrics and fps for other studies and this study

Study	% Accuracy	% Precision	% Sensitivity	fps
Wagner et al. 2022	92.8 & 90	Null	Null	44 & 28
He et al. 2021	86.2	81.4 & 86.2	75.3 & 83.8	33
Rajput et al. 2020	98 & 95	Null	57.5 & 82.4	Null
Xiong et al. 2019	96	Null	Null	25
Elqattan et al. 2019	89 & 95	46	100	Null
Baheti et al. 2018	95.5	Null	Null	Null
This study	91 & 99	87 & 99	94 & 99	5 & 5

Table 3. The volume and the number of parameters of the backbone architectures for previous studies and this study

Study	Model Name	Model Volume	Number of Parameters
Wagner et al. 2022	ResNeXt-50	10.84 MB	23 M
He et al. 2021	CornerNet-Lite	~600 MB	null
Rajput et al. 2020	Inception-v2	92 MB	23.9 M
Xiong et al. 2019 & Elqattan et al. 2019	YOLOV3	235 MB	62.25 M
Baheti et al. 2018	VGG16	528 MB	138.4 M
This study	WSN	0.84MB	69.4 k

As explained in section 2.3, we used RGB augmented images for CNNs training, but we trained the WSNs with grayscale images and without any augmentation. This highlights the high potential of the WSN in feature extraction. From this point of view and with a consideration of the results achieved for the State Farm dataset, we observe that the WSN accomplishes SOTA accuracy by training on one-third of the information used for the CNNs (grayscale images versus RGB images).

The main challenge we faced in this study was the weaknesses of the HAAR cascade detector that affected the performance of the classifier for the frontal view. Face occlusions or dark skin harms the functionality of the cascade detector, hence we would lose the track of the face and the ears' ROIs. In this regard, we refer to the studies that worked on occluded face detection [54], [55]. The error rate of our system increases if the environment is too dark. Also, we cannot detect cellphone usage if the subject is using a hands-free gadget for the phone call.

5 Conclusion

In this paper, we presented a new dataset for mobile cellphone usage detection. Also, we used WSN for the classification of the images in our presented dataset as well as the publicly available Kaggle State Farm dataset. We illustrated that the WSN outperforms the accuracy of well-known SOTA CNN models, but with a very small and simple architecture and far fewer parameters. We reached the accuracy of 91% for our domestic dataset and 99% for the State Farm publicly available dataset. We also showed that it can achieve even better results in image classification by training on grayscale images without image augmentations.

For future works, we will apply WSN for the localization of cellphone usage. Also, the same classification approach could be further improved by merging the datasets from the frontal view (like the dataset created by [23]) to the side view of the driver (Kaggle State Farm and dataset used by [56]). To compensate for the inability of the system in object detection in low light, architectures proposed by [57], [58] are powerful solutions. For detection of other ways of using cellphones like hands-frees or speakers, an ensemble of vision-based techniques with methods like [59], [60] may overcome the problem.

References

1. Sundfør, H. B.: Inattention and distraction in fatal road crashes – Results from in-depth crash investigations in Norway. *Accident Analysis & Prevention* 125, 152–157 (2019).
2. Stevens, A.: In-vehicle distraction and fatal accidents in England and Wales. *Accident Analysis & Prevention* 33 (4), pp. 539–545 (2001).
3. Stimpson, J. P.: Fatalities of Pedestrians, Bicycle Riders, and motorists Due to Distracted driving motor Vehicle Crashes in the U.S., 2005–2010. *Public Health Reports* 128(6), 436–442 . (2013).
4. Zhang, Y.: Exploring Driver Injury Severity at Intersection: An Ordered Probit Analysis. *Advances in Mechanical Engineering* 7(2), (2015)
5. Transportation Institute releases findings on driver behavior and crash factors, <https://vtechworks.lib.vt.edu/bitstream/handle/10919/59404/2006-237.html?sequence=1&isAllowed=y>, last accessed 2022/11/23.
6. Balbinot, A. B.: FUNÇÕES PSICOLÓGICAS E COGNITIVAS PRESENTES NO ATO DE DIRIGIR E SUA IMPORTÂNCIA PARA OS MOTORISTAS NO TRÂNSITO. *Ciências & Cognição* vol. 16(2), 13-29 (2011).
7. A. Salvador.: A culpa foi do celular. *Revista Veja* (Brazilian weekly newsmagazine), (2011).
8. National Safety Council.: Understanding the distracted brain. National Safety Council (April), (2012).
9. Regan, M.: Driver Distraction Injury Prevention Countermeasures— Part 3: Vehicle, Technology, and Road Design. *Driver Distraction: Theory, Effects, and Mitigation*, 579–601(2008).
10. Peissner, M.: Can voice interaction help reducing the level of distraction and prevent accidents? Meta-Study on Driver Distraction and Voice Interaction. *Whitepaper* 125, 152–157 (2011).

11. Brookhuis, K. A.: Behavioral impacts of advanced driver assistance systems—an overview. *European Journal of Transport and Infrastructure Research* 1(3), 245-253 (2001).
12. Agrawal, S.: Building on the past to help prepare the workforce for the future with automated vehicles: A systematic review of automated passenger vehicle deployment timelines. *Technology in Society* 72(December), (2023).
13. Saeed, T. U.: Road infrastructure readiness for autonomous vehicles (2019).
14. Hakak, S.: Autonomous Vehicles in 5G and Beyond: A Survey, 1–34, (2022).
15. Ilkova, V.: Legal aspects of autonomous vehicles-An overview. In: 2017 21st International Conference on Process Control (PC), pp. 428–433, (2017).
16. Validi, A.: Examining the Impact on Road Safety of Different Penetration Rates of Vehicle-to-Vehicle Communication and Adaptive Cruise Control, *IEEE Intelligent Transportation Systems Magazine* 10(4), 24–34 (2018).
17. Eckert, A., Hohm, A.: An integrated ADAS solution for pedestrian collision avoidance. In: Proceedings of the 23rd International Conference on the Enhanced Safety of Vehicles, Seoul, Republic of Korea, pp. 13-298 (2013).
18. Gruyer, D.: From virtual to reality, how to prototype, test and evaluate new ADAS: Application to automatic car parking. In: IEEE Intelligent Vehicles Symposium, Proceedings, pp. 261–267. IEEE, Dearborn, MI, USA (2014).
19. Ebrahimian, S.: Multi-Level Classification of Driver Drowsiness by Simultaneous Analysis of ECG and Respiration Signals Using Deep Neural Networks. *International Journal of Environmental Research and Public Health*, 19(17) (2022).
20. Saini, V.: Driver Drowsiness Detection System and Techniques : A Review. *International Journal of Computer Science and Information Technologies* 5(3), 4245–4249 (2014).
21. Sandeep, K.: Novel Drunken Driving Detection and Prevention Models Using Internet of Things, In: International Conference on Recent Trends in Electrical, Electronics and Computing Technologies (ICRTEECT), IEEE, Warangal, India (2017).
22. Isaksson-hellman, I.: Real-World Evaluation of Driver Assistance Systems for Vulnerable Road Users Based on Insurance Crash Data in Sweden. 1–10 (2019).
23. Lassmann, P.: Keeping the balance between overload and underload during partly automated driving: relevant secondary tasks. In: Bertram, T. (eds) *Automatisiertes Fahren 2019. Proceedings*, Springer Vieweg, Wiesbaden (2020).
24. Zhou, R.: Driver's distracted behavior: The contribution of compensatory beliefs increases with higher perceived risk. *International Journal of Industrial Ergonomics* 80(August), 103009 (2020).
25. Wagner, B.: Vision Based Detection of Driver Cell Phone Usage and Food Consumption. *IEEE Transactions on Intelligent Transportation Systems* 23(5), 4257–4266 (2022).
26. He, A.: Driver cell-phone use detection based on cornernet-lite network. *IOP Conference Series: Earth and Environmental Science*. 632(4), (2021).
27. Rajput, P.: Detecting Usage of Mobile Phones using Deep Learning Technique. *ACM International Conference Proceeding Series*. 96–101 (2020).
28. State Farm Distracted Driver Detection, <https://kaggle.com/competitions/state-farm-distracted-driver-detection>, last accessed 2022/12/17.
29. Xiong, Q.: A deep learning approach to driver distraction detection of using mobile phone. *IEEE Vehicle Power and Propulsion Conference*, 1–5 (2019).
30. Shi, X.: Real-Time Rotation-Invariant Face Detection with Progressive Calibration Networks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 2295–2303 (2018).
31. Redmon, J.: YOLOv3: An Incremental Improvement. (2018).

32. Baheti, B., Gajre, S., Talbar, S.: Detection of distracted driver using convolutional neural network. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops 2018-June*, 1145–1151 (2018).
33. Abouelnaga, Y.: Real-time Distracted Driver Posture Classification. *32nd Conference on Neural Information Processing Systems (NIPS 2018), Workshop on Machine Learning for Intelligent Transportation Systems*, (2017).
34. Le, T. H. N.: Multiple Scale Faster-RCNN Approach to Driver's Cell-Phone Usage and Hands on Steering Wheel Detection. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 46–53 (2016).
35. Torres, R.: A deep learning approach to detect distracted drivers using a mobile phone. In: *International Conference on Artificial Neural Networks*, pp. 72–79. Springer Science, Bel'em, PA, Brazil (2017).
36. Xu, B.: A machine learning approach for detecting cell phone usage. *Video Surveillance and Transportation Imaging Applications*, pp. 70–77. Society of Photo-Optical Instrumentation Engineers (SPIE), San Francisco, California, United States (2015).
37. Xu, B.: A machine learning approach to vehicle occupancy detection. In: *17th IEEE International Conference on Intelligent Transportation Systems, ITSC 2014*, pp. 1232–1237. Institute of Electrical and Electronics Engineers Inc, Qingdao, China (2014).
38. Elqattan, Y.: System for Detecting and Reporting Cell Phone Distracted Drivers. In: *11th International Symposium on Image and Signal Processing and Analysis (ISPA)*, pp. 215–221. IEEE, Dubrovnik, Croatia (2019).
39. Chollet, F.: Xception: Deep Learning with Depthwise Separable Convolutions. In: *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 1800–1807. Institute of Electrical and Electronics Engineers Inc. Honolulu, HI, USA (2016).
40. Deng, J.: ImageNet: A large-scale hierarchical image database. *IEEE Conference on Computer Vision and Pattern Recognition*, 248–255, (2010).
41. Berri, R.: A hybrid vision system for detecting use of mobile phones while driving, In: *International Joint Conference on Neural Networks (IJCNN)*, pp. 4601–4610. IEEE, Vancouver, BC, Canada (2016).
42. Artan, Y.: Driver Cell Phone Usage Detection from HOV/HOT NIR Images. In: *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 225–230, IEEE, Columbus, OH, USA (2014).
43. Seshadri, K.: Driver cell phone usage detection on Strategic Highway Research Program (SHRP2) face view videos, In: *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 35–43, IEEE, Boston, USA (2015).
44. Yasar, H.: Detection of Driver's Mobile Phone Usage. *IEEE 9th International Conference on Humanoid, Nanotechnology, Information Technology, Communication and Control, Environment and Management (HNICEM)*, 1–4 (2017).
45. Gröger, C.: There is no AI without data. *Communications of the ACM* 64(11), pp. 98–108, (2021).
46. Bengani, D.: Face Detection Using Viola Jones Algorithm. *International Journal for Modern Trends in Science and Technology* 6(11), 131–134 (2020).
47. Manimala, S.: Anticipating Hand and Facial Features of Human Body using Golden Ratio Anticipating Hand and Facial Features of Human Body using Golden Ratio. *International Journal of Graphics & Image Processing* 4, (2017).
48. Illustration of transforms — Torchvision main documentation, https://pytorch.org/vision/stable/auto_examples/plot_transforms.html#sphx-glr-auto-examples-plot-transforms-py, last accessed 2022/10/30.
49. Oyallon, E.: *Generic Deep Networks with Wavelet Scattering*, (2013).

50. Andreux, M.: Kymatio: Scattering transforms in python. *Journal of Machine Learning Research* 21, 1–6 (2020).
51. Simonyan, K.: Very deep convolutional networks for large-scale image recognition. *3rd International Conference on Learning Representations, ICLR - Conference Track Proceedings*, 1–14 (2015).
52. He, K.: Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2016-Decem*, 770–778 (2016).
53. Keras, <https://github.com/fchollet/keras>, last accessed 2022/12/19
54. Carragher, D. J.: Masked face identification is improved by diagnostic feature training. *Cognitive Research: Principles and Implications* 7(1), 1–12 (2022).
55. Qiu, H.: End2End Occluded Face Recognition by Masking Corrupted Features. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44 (10), 6939 - 6952, 2015.
56. Eraqi, H. M.: Driver distraction identification with an ensemble of convolutional neural networks. *Journal of Advanced Transportation*, (2019).
57. Mukherjee, R.: Object Detection under Challenging Lighting Conditions Using High Dynamic Range Imagery. *IEEE Access* 9, 77771–77783, (2021).
58. Morawski, I.: NOD: Taking a Closer Look at Detection under Extreme Low-Light Conditions with Night Object Detection Dataset, (2021).
59. Rodríguez-Ascariz, J. M.: Automatic system for detecting driver use of mobile phones, *Transportation Research Part C: Emerging Technologies* 19(4). 673–681 (2011).
60. Yang, J.: Detecting driver phone use leveraging car speakers. In: *Proceedings of the Annual International Conference on Mobile Computing and Networking, MOBICOM*, pp. 97–108. Association for Computing Machinery, Nevada, Las Vegas, USA (2011).