1 **Detection of dangerous driving behaviour using machine learning techniques and big data**
2

3 **Thodoris Garefalakis**
4 Research Associate
5 Department of Transportation Planning and Engineering
6 National Technical University of Athens, Athens, Greece, GR15773
7 Email: tgarefalakis@mail.ntua.gr
8
9 **Eva Michelaraki**
10 Ph.D. Candidate, Research Associate
11 Department of Transportation Planning and Engineering
12 National Technical University of Athens, Athens, Greece, GR15773
13 Email: evamich@mail.ntua.gr
14
15 **George Yannis**
16 Professor
17 Department of Transportation Planning and Engineering
18 National Technical University of Athens, Athens, Greece, GR15773
19 Email: geyannis@central.ntua.gr
20
21 Word Count: 4921 words + 2 table (250 words per table) = 5,421 words
22
23
24 *Submitted: July 31, 2024*
25

*Garefalakis T. et al.*

**ABSTRACT**

A wide range of variables affect road safety, including the driver's state, environment, and traffic conditions. The aim of the current study was twofold; a) the evaluation of the impact of various indicators on predicting Safety Tolerance Zone (STZ) and b) the development of a deep learning model for identifying dangerous driving. In order to achieve these objectives, a naturalistic driving experiment was implemented and data from a representative sample of 50 Belgian car drivers were collected and analysed. The impact of the features on the estimation of STZ was performed based on the XGBoost algorithm. Key features such as headway, forward collision warning indicator, and distance traveled were found to significantly affect the prediction of STZ levels. Subsequently, a Long Short-Term Memory (LSTM) model was developed for real time data prediction, leveraging its strength in handling sequential data and temporal dependencies. The results demonstrated that the LSTM model achieved a 71% accuracy rate in identifying dangerous driving behaviors, underscoring the potential of this approach for enhancing road safety. The analysis highlighted the critical importance of the identified features, suggesting that the combination of XGBoost for feature selection and LSTM for real-time prediction provides a robust framework for real-time intervention and support systems. This integrated method offers significant promise for reducing road accidents and improving overall traffic safety by enabling timely and accurate identification of risky driving behaviors.

**Keywords:** driving behaviour, i-DREAMS project, Safety Tolerance Zone, XGBoost, Long Short-Memory Network

1 **INTRODUCTION**
2    According to the World Health Organization (WHO), every year 1.19 million lives are lost on
3 road crashes, with millions more experiencing severe injuries and enduring long-term health
4 consequences (World Health Organization, 2023). Road crashes rank as the 12th leading cause of death
5 for people of all ages, with the most pronounced impact observed among individuals aged 5 to 29 years,
6 for whom it stands as the foremost cause of death (1).
7    Numerous factors contribute to road safety, encompassing the driver's condition, environmental
8 elements, and traffic situations (2). Despite advancements in technology and infrastructure, human error
9 remains a significant factor in traffic crashes (3, 4). However, the ongoing development of autonomous
10 vehicles holds promise for improving road safety by minimizing reliance on human drivers (5). Intelligent
11 Transportation Systems (ITS) for monitoring driving behaviour, incorporating real-time interventions,
12 have demonstrated notable effectiveness in enhancing road safety (6). Moreover, the convergence of AI-
13 driven technologies and traditional safety measures marks a transformative era in road safety
14 management. Combining the advantages of autonomous vehicles and monitoring systems has significant
15 potential for reducing the impact of human error and fostering a safer road environment for all users.
16    To date, significant number of studies explore the risk assessment, recognition and classification
17 of driving behaviour and profiling through machine learning techniques and clustering algorithms. To
18 begin with, Chen et al. (2023) (7) implemented different feature extraction methods to identify driver's
19 characteristics and improve the accuracy of driving behaviour modelling. Other studies have introduced
20 deep learning techniques to identify dangerous driving patterns based on factors such as speed, headway,
21 or time to collision (8, 9). Ghandour et al. (2021) (10) applied and compared four machine learning
22 classification methods to identify drivers' behaviour and distraction situations based on real data
23 corresponding to different behaviours (i.e., normal, drowsy, and aggressive). These studies collectively
24 underscore the crucial role of advanced data analytics in preemptively identifying and mitigating risk
25 factors associated with driving behavior. By leveraging sophisticated algorithms and large datasets,
26 researchers can uncover nuanced patterns and correlations that were previously undetectable, thus paving
27 the way for more effective and targeted safety interventions.
28    The overall objective of the European H2020 i-DREAMS project is to address to tackle these
29 challenges by establishing, developing, testing, and validating a 'Safety Tolerance Zone' (STZ) to ensure
30 safe driving behaviour. The i-DREAMS aims to assess the appropriate level within the STZ continuously,
31 considering risk factors associated with task complexity (e.g., traffic conditions and weather) and coping
32 capacity (e.g., driver's mental state, driving behaviour, and vehicle status). Interventions are then
33 implemented to keep drivers' operations within acceptable safety limits. The STZ is structured into three
34 levels: 'Normal,' 'Dangerous,' and 'Avoidable Accident.' The 'Normal' level indicates a low probability of
35 a crash, while the 'Dangerous' level suggests an increased likelihood of a crash without inevitability. The
36 'Avoidable Accident' level signifies a high probability of a crash but allows enough time for drivers to act
37 and prevent it. The key distinction between the 'Dangerous' and 'Avoidable Accident' levels refer to the
38 more immediate need for intervention in the latter level. This stratified approach not only aids in the real-
39 time assessment of driving conditions but also ensures that interventions are appropriately scaled to the
40 severity of the risk, thereby enhancing the overall efficacy of safety measures.
41    According to the framework described above, the aim of the present study is to incorporate crash
42 prediction and risk evaluation within a Long Short-Term Memory (LSTM framework, as well as identify
43 the impact of certain features on this operation. To achieve this goal, the study evaluates explanatory
44 variables related to risk and identifies the most dependable indicators of task complexity and coping
45 capacity. These indicators include variables such as headway, distance, speed, forward collision, time of
46 day (highlighted by high-beam indicators) and weather.
47    The paper is organized as follows: At the beginning, a detailed introduction about the context and
48 the aim of the current study is given. Furthermore, there was a presentation of the data gathered for the
49 analysis. Subsequently, a concise explanation of the methodological approach is given. Then, the paper
50 outlines the notable findings and summarizes the results of the conducted statistical analysis. Finally,

conclusions are emphasized, and the paper concludes by addressing limitations and presenting
suggestions for future research.

**DATA**

In the context of this study, a naturalistic driving experiment was conducted, involving 50 car drivers from Belgium over a 15-month timeframe (from April 2021 to July 2022), resulting in the creation of a substantial database comprising 7,160 trips. The i-DREAMS project focuses on delivering an integrated set of monitoring and communication tools for intervention and support. Participants were selected according to specific criteria to ensure a varied and representative sample. Criteria included adequate driving experience and road exposure, being at least 18 years old, an equal representation of genders, and having vehicle types compatible with the i-DREAMS technology. Additionally, participants needed to use a smartphone, and vehicles used by multiple drivers were preferred to increase the sample size. The recruitment process involved general advertising, initial candidate screening based on the criteria, targeted advertising for certain groups, and providing detailed information before finalizing participation contracts. Participants were rewarded to participate in the naturalistic driving experiment.

In order to monitor driving performance indicators, state-of-the-art technologies and systems were employed. Specifically, data from the Mobileye system (Mobileye, 2022), a dash camera, and the Cardio gateway (CardioID Technologies, 2022), which records driving behaviour and Global Navigation Satellite System (GNSS) signals, were utilized. The Mobileye system operates as a sensor network measuring parameters such as headway distance. Information regarding the current warning stage, as defined by Mobileye, was collected for comparison with the i-DREAMS warning stage (normal driving, dangerous phase, avoidable accident phase). The integration of these advanced technologies allows for a comprehensive analysis of driving behavior, capturing a wide range of variables that contribute to road safety. The data collected includes not only driving metrics but also contextual information such as weather conditions and time of day, providing a holistic view of the driving environment.
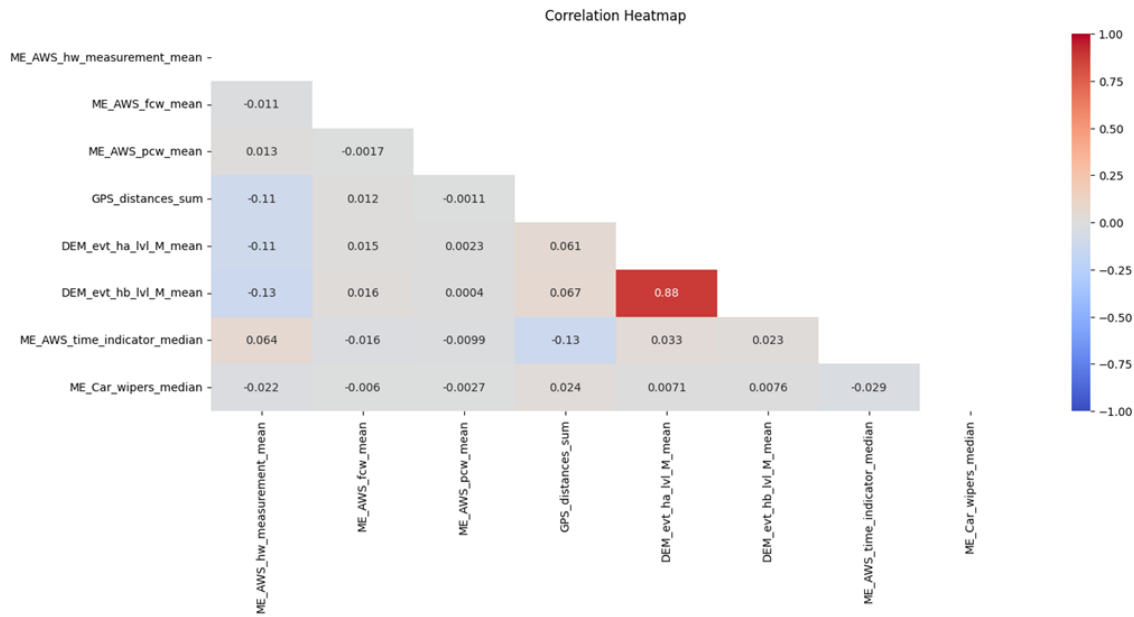
To assess risk-related explanatory variables and reliable indicators of task complexity and coping capacity, parameters like time headway, distance travelled, forward collision, and weather conditions were analysed. Special attention was given to average speed, and a new variable, STZ_level, accounting for different levels of Safety Tolerance Zone (STZ) was created. Consequently, the dependent variable was the level of (i.e., STZ_level), categorized into three levels (Normal Driving phase: 0, Dangerous phase: 1, Avoidable Accident phase: 2). **Table 1** provides an overview of the driving performance indicators examined along with their corresponding description.

**TABLE 1 Description of driving performance indicators**

| Variable | Description |
|---|---|
| GPS_distances_sum | Distance travelled (km) |
| ME_AWS_hw_measurement_mean | Headway measurement (seconds) |
| ME_AWS_fcw_mean | Forward collision warning |
| ME_AWS_pcw_mean | Pedestrian collision warning |
| DEM_evt_hb_lvl_M_mean | Medium level harsh braking events |
| DEM_evt_ha_lvl_M_mean | Medium level harsh acceleration events |
| ME_Car_wipers_median | Indicates weather conditions (wipers on/off) |
| ME_AWS_time_indicator_median | Indicates lighting conditions (day, dusk, night) |

Following thorough data cleaning and preparation, the subsequent phase of the analysis encompassed examining collinearity to eliminate any highly correlated variables from the models. Variables with an absolute correlation coefficient of at least 0.6 were considered highly correlated. **Figure 1** illustrates the correlation coefficients among the variables employed in the models. The analysis

1    indicated that "DEM_evt_ha_lvl_M_mean" and "DEM_evt_hb_lvl_M_mean" exhibit a high correlation;
2    thus, "DEM_evt_hb_lvl_M_mean" was excluded from the following analysis.
3



4
5
6    **Figure 1 Correlation heatmap**
7
8    **METHODS**
9    Following the data collection, an algorithm for feature selection, specifically XGBoost, was
10   utilized to pinpoint the crucial features for predicting the STZ level. XGBoost is chosen for its efficiency
11   and accuracy in handling large-scale datasets and its ability to handle various types of data. Subsequently,
12   these identified features were input into a Long Short-Term Memory (LSTM) classifier to determine the
13   STZ level. The subsequent sections provide a more detailed explanation of the algorithms employed in
14   this process.
15
16   **Extreme Gradient Boosting (XGBoost)**
17   XGBoost, stands for eXtreme Gradient Boosting which is an ensemble learning algorithm that
18   has gained widespread popularity for its high performance in various machine learning tasks. The
19   XGBoost algorithm is an optimized form of the Gradient Boosting model that operates as a Newton-
20   Raphson algorithm, using a second-order Taylor approximation (11), contrary to Gradient Boosting,
21   which relies on gradient descent. XGBoost's efficiency is due to its use of a novel tree boosting system
22   that is faster and more accurate than existing methods. More specifically, XGBoost is an implementation
23   of gradient-enhanced decision trees, in which trees are generated sequentially with significantly higher
24   model accuracy, in less computational training time, than standard machine learning models. This makes
25   XGBoost particularly suitable for large-scale data analysis in the context of driving behavior, where the
26   volume of data can be substantial and the need for real-time predictions is critical.
27
28   In the context of XGBoost feature importance calculation, the algorithm utilizes three key metrics: "gain",
29   "frequency" and "cover" to assess the significance of each feature in the constructed trees (12).
30   • Gain: It quantifies the improvement in accuracy that a particular feature brings to the model when
31     creating branches in the decision trees. Features with higher gains are considered more influential in
32     making decisions.

- Frequency: Is a simple count of how frequently a feature is utilized across all the trees in the ensemble. A higher frequency indicates that the feature is consistently chosen during the construction of decision trees.
- Cover: Provides information about the relative scale of a feature's contribution to the overall prediction. It considers both the frequency of the feature and the magnitude of its impact on the predictions. Higher cover values suggest that a feature has a more significant influence on the model's output.

**Long Short-Term Memory (LSTM)**

A Long Short-Term Memory (LSTM) network stands out as a specialized form of Recurrent Neural Network (RNN) that has gained significant prominence across various domains due to its remarkable proficiency in capturing and comprehending intricate sequential patterns. This innovative architecture was conceived to address the prevalent challenge of vanishing gradients, a hindrance commonly encountered by traditional RNNs, particularly in tasks requiring the modelling of prolonged dependencies (13).

The distinctive ability of LSTM networks to persistently retain information over extended durations positions them as a formidable choice for tasks involving the modelling of sequential data. Comprising a series of interconnected modules, an LSTM network adopts a chain-like architectural structure (14). At the core of these networks are fundamental processing units termed "cells," analogous to the complex nature of neurons in traditional multi-layer perceptrons (MLP).
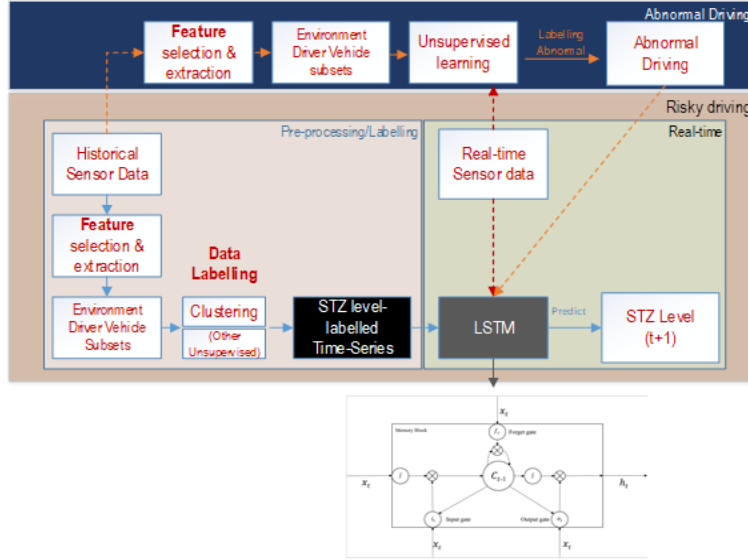
Within each LSTM cell, multiple gates play a pivotal role in orchestrating the flow of information across sequences of variable length. This intrinsic feature endows LSTM networks with the autonomy to discern the relevance of information over both long-term and short-term contexts, making them exceptionally well-suited for an extensive array of tasks such as activity recognition and language translation (15).

In a standard LSTM configuration, three key gates contribute to the network's functionality:

- Forget Gate: Tasked with determining which information is to be retained and what should be discarded in the cell state, the forget gate utilizes a sigmoid layer, known as the "forget gate layer," to make these critical decisions.
- Input Gate: Responsible for deciding what new information should be incorporated into the cell state and how it should be updated, the input gate comprises two essential components. The input gate layer, leveraging a sigmoid function, determines the values to be updated, while a hyperbolic tangent (tanh) layer produces a vector of candidate values for potential integration into the state. The existing cell state undergoes an update based on these components.
- Output Gate: With the responsibility of filtering and selecting the information to be output from the memory block at a specific time step, the output gate derives the output from the cell state after filtering. Consisting of a sigmoid layer, this gate determines the relevant portions of the cell state to be included in the output. The filtered cell state then passes through a tanh activation function to scale values within the range of -1 and 1. The final output is generated by multiplying the result with the output of the sigmoid gate, ensuring the desired output is achieved.

With regards to the proposed LSTM model, the problem of defining the STZ levels becomes more straightforward, since LSTMs as a sub-category of Deep Neural Networks act like "black-boxes" (16) and thus the only input that needs to be provided to the model are labelled time series data. Collected historical measurements from the i-DREAMS technologies were used as input for an unsupervised learning approach grouping together measurements correlated with normal operation of a vehicle and those departing from normal driving behaviour. The proposed approach followed using LSTMs is given in **Figure 2**.

1
2
3  **Figure 2 STZ modelling using LSTMs**
4
5  **Model Evaluation Metrics**
6       A Long Short-Term Memory (LSTM) network stands out as a specialized form of Recurrent
7  Neural Network (RNN) that has gained significant prominence across various domains due to its
8  remarkable proficiency in capturing and comprehending intricate sequential patterns. This innovative
9  architecture was conceived to address the prevalent challenge of vanishing gradients, a hindrance
10 commonly encountered by traditional RNNs, particularly in tasks requiring the modelling of prolonged
11 dependencies (13).
12      In order to evaluate and compare the classification performance across various configurations
13 (involving hyperparameters and input combinations), established machine learning error metrics were
14 computed. The evaluation relies on metrics derived from the confusion matrix, which includes True
15 Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN). The classification
16 algorithms are assessed based on accuracy, precision, recall, and f1-score, as defined below.
17
18 Accuracy is a measure that assesses the proportion of correctly classified observations in a model's
19 predictions and is expressed as:

20 $Accuracy = \frac{TP+TN}{TP+FP+FN+TN},$  (1)

21 Precision, measuring the number of positive class predictions that truly belong to the positive class, is
22 defined as:

23 $Precision = \frac{TP}{TP+FP},$  (2)

24 Recall, also referred to as True Positive Rate, is defined as:
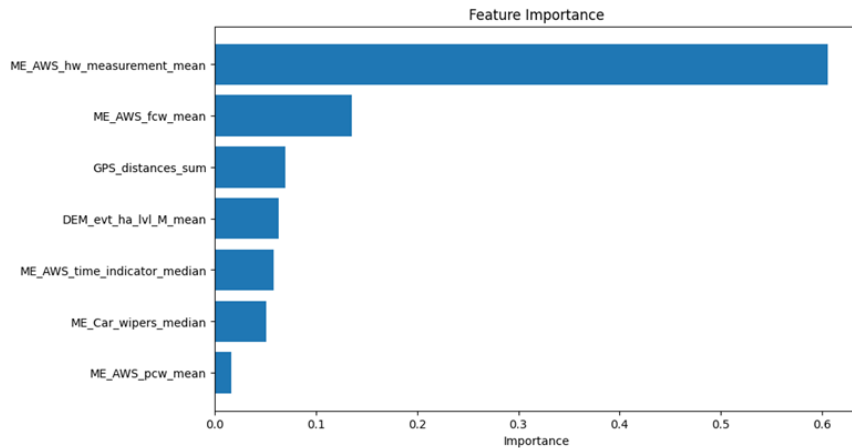
25 $Recall = \frac{TP}{TP+FN},$  (3)

26 F-measure, a composite metric combining precision and recall, is defined as:

27 $f1\text{-}score = \frac{2x\ (Precision)x\ (Recall)}{(Precision)+(Recall)},$  (4)

**RESULTS**

An algorithm for assessing the importance of features, derived from Extreme Gradient Boosting (XGBoost), was employed to assess the relevance of variables in predicting STZ and identify the most suitable independent variables. Among all the indicators analysed, the factors with the highest importance were headway measurement, forward collision warning indicator and GPS distance travelled. Parameters related to task complexity, such as car wipers and time indicator, found to be less significant. Additionally, variables such as medium-level harsh acceleration events and pedestrian collision warning had a lower impact on STZ. **Figure 3** presents an overview of the independent variables' feature importance based on the XGBoost algorithm.
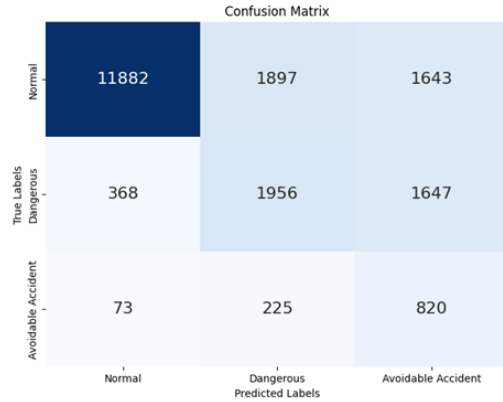


**Figure 3 XGBoost feature importance of independent variables**

A dataset consisting of approximately 275,000 data points was employed to train an LSTM Neural Network model. In accordance with the feature importance and the significance of relevant indicators, the input layer was structured with three neurons, representing headway measurement, forward collision warning indicator, and distance travelled. Additionally, the output layer was configured with a single neuron denoting the STZ. The model's architecture was meticulously designed to ensure that the most critical features were effectively leveraged, enhancing the model's predictive accuracy and reliability.

For an in-depth evaluation of the model performance, a confusion matrix was produced for the independent variable of STZ_level which illustrates the number of correct and incorrect predictions per class, as shown in **Figure 4**. In particular, the confusion matrix contains three rows and three columns and reports the number of true positives, true negatives, false positives, false negatives values. This allows a more detailed analysis than the proportion of correct classifications (e.g., accuracy, precision, recall).

**Figure 4 Confusion Matrix**

For the evaluation metrics of the model, performance indicators such as accuracy, precision, recall and f1-score were used as shown in the **Table 2**. It should be noted that accuracy provides the overall correctness of a model, while precision computes the accuracy of positive predictions and recall measures the ability of a model to correctly identify all relevant instances. Given the fact that identifying correctly risky driving behaviour is crucial for the purpose of this analysis, Recall stands as the most significant metric.
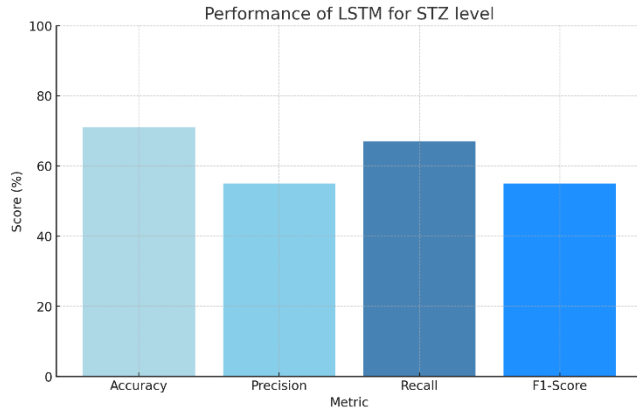
**TABLE 2 Assessment of the classification LSTM on STZ level**

| Model | Accuracy | Precision | Recall | f1-score |
|---|---|---|---|---|
| Long Short-Term Memory (LSTM) | 71% | 55% | 67% | 55% |

The LSTM model demonstrated an accuracy of 71%, indicating that the model correctly identified the STZ level in 71% of the cases. The precision of 55% suggests that just over half of the positive predictions made by the model were correct, while the recall of 67% indicates that the model successfully identified 67% of all actual cases of risky driving. The f1-score of 55% highlights a balance between precision and recall, emphasizing the model's moderate effectiveness in identifying safety-critical driving behaviors.

Further analysis of the confusion matrix revealed that the model performed better at identifying the 'Normal' and 'Dangerous' levels compared to the 'Avoidable Accident' level. This discrepancy suggests that while the model is fairly reliable in detecting standard and moderately risky driving behaviors, it has more difficulty in accurately predicting the most critical situations where immediate intervention is necessary. This limitation could be addressed in future studies by incorporating additional data or refining the model architecture.

The relevant performance of LSTM on STZ level for the Belgian car drivers is presented in **Figure 5.**

Figure 5 **Performance of LSTM for STZ level**

**Figure 5 Performance of LSTM for STZ level**

In summary, the LSTM model's performance indicates a promising approach for real-time prediction of risky driving behaviors, with a notable accuracy rate of 71%. However, the precision and recall metrics highlight areas for improvement, particularly in enhancing the model's ability to accurately identify all instances of dangerous and avoidable accident levels. These findings underscore the importance of continuous refinement and validation of machine learning models to ensure their effectiveness in practical applications.

**DISCUSSION**

As shown in **Table 2**, the LSTM achieve significant results, reaching 71% accuracy after the developed trials. Although LSTM is often used for sequence modelling, it is worth noting that the sequence may not always be explicitly visible in the predictors themselves. In some cases, the sequence may be implicit in the way that the data are structured or organized. For instance, in time series data, the sequence is often defined by the order in which the data were collected over time. In this case, the LSTM is used to model and make predictions based on the temporal dependencies and patterns in the data. In other cases, the sequence may be less obviously related to time, but still exist in the way that the data is organized. In natural language processing, the sequence may be defined by the order of words in a sentence or text document. Thus, the sequence is implicit in the way that the data was collected or organized, even if it's not immediately apparent from the predictors themselves. An LSTM could still be used in this case to model and make predictions based on the implicit sequence in the data.

It should be noted that an accuracy of less than 60% may not be sufficient for a high-performance intervention system, as it could result in a relatively high number of false alarms or missed detections. However, the required level of accuracy depends on the specific use case and the risks involved. For instance, in a system designed to detect potential crashes or safety hazards, a higher level of accuracy may be necessary in order to ensure the safety of drivers and other road users. As for the use of prediction models by an intervention system, the output of the models can be used in a variety of ways. In particular, the prediction models can generate real-time alerts or warnings to drivers or other stakeholders, such as traffic control centers or emergency responders. The models can also be used to trigger automated interventions, such as adjusting the speed of a vehicle or activating safety features like automatic braking systems. In addition, the output of prediction models can be used for ongoing analysis and monitoring of road safety performance, in order to identify trends and patterns that can inform future interventions and improvements.

Cura et al. (2021) (17) developed an LSTM technique to assess and classify bus driver behaviour characterised by deceleration, engine speed, corner turn and lane change attempts. Their proposed model demonstrated an accuracy of 87.7%, precision of 88.5%, and recall of 87.7%. Furthermore, Parsa et al. (2019) (18) utilised real-time data along with LSTM deep learning technique to detect road crashes. The

LSTM model achieved as accuracy of 96%, a detection rate of 73.8% and a false alarm rate of 3%. These findings are in line with the results of this research, where LSTMs showed significant predictive accuracy. Lastly, it should be mentioned that the aforementioned studies underscore the importance of model selection in road traffic safety applications and highlight the potential of LSTMs in different yet related contexts.

Nevertheless, this study is not without limitations. Firstly, drivers' socio-demographic characteristics, such as gender, age, driving experience, education level, or mental health state were not included in the analysis. Secondly, the influence of psychological status of participants, such as driver distraction through mobile phone use, fatigue or sleepiness were not taken into consideration in the present study, as only naturalistic driver data from the driving experiment were used. Given the fact that drivers react differently under different circumstances with regards to road layouts (i.e., urban, rural environments, highways) and traffic volumes (i.e., high, medium or low traffic volumes), it would be of great interest to investigate headway or speed using environment, vehicle and driver questionnaire data.

As per future research directions, the experimental sample size of the analysis could be further expanded and strengthened. A larger dataset including additional drivers' age groups or drivers from different countries, regions could enhance the analysis procedure. At the same time, data from different transport modes could be also explored in order to allow comparisons among private (i.e., car drivers) and professional drivers (i.e., bus and truck drivers). Future research efforts could consider the examination of additional machine learning methods to be applied. For instance, imbalanced learning as well as microscopic data analysis of the database collected could be implemented through deep learning and econometric techniques. Lastly, the investigation of other significant risk indicators (e.g., drug abuse, alcohol consumption or the seat belt use) could be also included in the future.

**CONCLUSIONS**

The study aimed to investigate the impact of various parameters on predicting the Safety Tolerance Zone (STZ) and develop a deep learning model to identify risky driving behaviour in real time. The analysis relied on data collected from a naturalistic driving experiment involving 50 Belgian car drivers for 15-month period and more than 7,000 trips were analysed.

For the purpose of this research, the most important risk indicators were identified. Towards that end, an importance assessment algorithm derived from Extreme Gradient Boosting (XGBoost) was implemented in order to assess the importance of the examined variables (i.e., headway, forward collision warning, pedestrian collision warning, distance travelled, harsh acceleration events, time of the day and weather conditions) in predicting STZ level. Furthermore, a Long Short-Term Memory (LSTM) model was utilized for real-time data prediction, considering the key and meaningful risk indicators.

Focusing on the results of all classes combined, classifiers achieve 71% accuracy, 55% precision and 67% recall. First of all, the total accuracy means that the model is 71% accurate in making a correct prediction. Moreover, the model was 55% accurate regarding a positive sample and 67% accurate on predicting safety-critical classes (i.e., "Dangerous" and "Avoidable Accident"), which means that the model can be trusted in its ability to detect positive samples in a moderate degree.

The findings revealed a significant impact of headway, forward collision warning indicator and distance travelled on predicting the STZ level. The classification results could be also improved by exploiting other imbalanced learning techniques in order for all three STZ levels to be correctly identified in real-time. Further optimisation and exploration of LSTM architectures may enhance their performance and reliability in driver behaviour analysis.

A combination of machine learning algorithms and i-DREAMS data could be proved beneficial in order to identify safe or unsafe driving behaviour. Through the utilization of data-driven insights, advanced analytics and real-time interventions, this method has the potential to enhance road safety, leading to a decrease in crashes, fatalities, or serious injuries.

6 **AUTHOR CONTRIBUTIONS**
7 The authors confirm contribution to the paper as follows: study conception and design: T. Garefalakis, E.
8 Michelaraki, G. Yannis; data collection: T. Garefalakis, E. Michelaraki; analysis and interpretation of
9 results: T. Garefalakis, E. Michelaraki. Author; draft manuscript preparation: T. Garefalakis, E.
10 Michelaraki, G. Yannis. All authors reviewed the results and approved the final version of the manuscript.

**REFERENCES**

1. World Health Organization. *Global Status Report on Road Safety 2023*. Geneva, 2023. https://www.who.int/publications/i/item/9789240086517. Accessed Mar. 6, 2024.

2. Wegman, F. The Future of Road Safety: A Worldwide Perspective. *IATSS Research*, Vol. 40, No. 2, 2017, pp. 66–71. https://doi.org/10.1016/j.iatssr.2016.05.003.

3. Muecklich, N., I. Sikora, A. Paraskevas, and A. Padhra. The Role of Human Factors in Aviation Ground Operation-Related Accidents/Incidents: A Human Error Analysis Approach. *Transportation Engineering*, Vol. 13, 2023, p. 100184. https://doi.org/10.1016/j.treng.2023.100184.

4. Staubach, M. Factors Correlated with Traffic Accidents as a Basis for Evaluating Advanced Driver Assistance Systems. *Accident Analysis & Prevention*, Vol. 41, No. 5, 2009, pp. 1025–1033. https://doi.org/10.1016/j.aap.2009.06.014.

5. Camps-Aragó, P., L. Temmerman, W. Vanobberghen, and S. Delaere. Encouraging the Sustainable Adoption of Autonomous Vehicles for Public Transport in Belgium: Citizen Acceptance, Business Models, and Policy Aspects. *Sustainability*, Vol. 14, No. 2, 2022, p. 921. https://doi.org/10.3390/su14020921.

6. Michelaraki, E., C. Katrakazas, G. Yannis, E. Konstantina Frantzola, F. Kalokathi, S. Kaiser, K. Brijs, and T. Brijs. A Review of Real-Time Safety Intervention Technologies. 2021.

7. Chen, Y., K. Wang, and J. J. Lu. Feature Selection for Driving Style and Skill Clustering Using Naturalistic Driving Data and Driving Behavior Questionnaire. *Accident Analysis & Prevention*, Vol. 185, 2023, p. 107022. https://doi.org/10.1016/j.aap.2023.107022.

8. Song, X., Y. Yin, H. Cao, S. Zhao, M. Li, and B. Yi. The Mediating Effect of Driver Characteristics on Risky Driving Behaviors Moderated by Gender, and the Classification Model of Driver's Driving Risk. *Accident Analysis & Prevention*, Vol. 153, 2021, p. 106038. https://doi.org/10.1016/j.aap.2021.106038.

9. Shi, X., Y. D. Wong, M. Z.-F. Li, C. Palanisamy, and C. Chai. A Feature Learning Approach Based on XGBoost for Driving Assessment and Risk Prediction. *Accident Analysis & Prevention*, Vol. 129, 2019, pp. 170–179. https://doi.org/10.1016/j.aap.2019.05.005.

10. Ghandour, R., A. J. Potams, I. Boulkaibet, B. Neji, and Z. Al Barakeh. Driver Behavior Classification System Analysis Using Machine Learning Methods. *Applied Sciences*, Vol. 11, No. 22, 2021. https://doi.org/10.3390/app112210562.

11. Ghosh, P., and S. Chakraborty. Spectral Classification of Quasar Subject to Redshift: A Statistical Study. 2023.

12. Zheng, H., J. Yuan, and L. Chen. Short-Term Load Forecasting Using EMD-LSTM Neural Networks with a Xgboost Algorithm for Feature Importance Evaluation. *Energies*, Vol. 10, No. 8, 2017, p. 1168. https://doi.org/10.3390/en10081168.

13. DiPietro, R., and G. D. Hager. Deep Learning: RNNs and LSTM. In *Handbook of Medical Image Computing and Computer Assisted Intervention*, Elsevier, pp. 503–519.

14. Sohail, A., M. A. Cheema, M. E. Ali, A. N. Toosi, and H. A. Rakha. Data-Driven Approaches for Road Safety: A Comprehensive Systematic Literature Review. *Safety Science*, Vol. 158, 2023, p. 105949. https://doi.org/10.1016/j.ssci.2022.105949.

15. Elmaz, F., R. Eyckerman, W. Casteels, S. Latré, and P. Hellinckx. CNN-LSTM Architecture for Predictive Indoor Temperature Modeling. *Building and Environment*, Vol. 206, 2021, p. 108327. https://doi.org/10.1016/j.buildenv.2021.108327.

16. Xu, C., A. P. Tarko, W. Wang, and P. Liu. Predicting Crash Likelihood and Severity on Freeways with Real-Time Loop Detector Data. *Accident Analysis & Prevention*, Vol. 57, 2013, pp. 30–39. https://doi.org/10.1016/j.aap.2013.03.035.

17. Cura, A., H. Kucuk, E. Ergen, and I. B. Oksuzoglu. Driver Profiling Using Long Short Term Memory (LSTM) and Convolutional Neural Network (CNN) Methods. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 22, No. 10, 2021, pp. 6572–6582. https://doi.org/10.1109/TITS.2020.2995722.

18. Parsa, A. B., R. S. Chauhan, H. Taghipour, S. Derrible, and A. Mohammadian. Applying Deep Learning to Detect Traffic Accidents in Real Time Using Spatiotemporal Sequential Data. *arXiv preprint arXiv:1912.06991*, 2019.