- 1 CogFormer: Aligned-attention Transformer-based Multi-physiological signals Fusion
- 2 for Driver Cognitive Load Estimation in Conditional Automated Driving
- 3

4 Ange Wang

- 5 Intelligent Transportation Thrust, Systems Hub
- 6 The Hong Kong University of Science and Technology (Guangzhou), China
- 7 Email: awang324@connect.hkust-gz.edu.cn
- 8

9 Haohan Yang

- 10 School of Mechanical and Aerospace Engineering
- 11 Nanyang Technological University, Singapore
- 12 Email: haohan.yang@ntu.edu.sg
- 13

14 Jiyao Wang

- 15 Robotics and Autonomous Systems Thrust, Systems Hub
- 16 The Hong Kong University of Science and Technology (Guangzhou), China
- 17 Email: jwang297@connect.hkust-gz.edu.cn
- 18

19 Hai Yang

- 20 Chair Professor
- 21 Department of Civil and Environmental Engineering
- 22 The Hong Kong University of Science and Technology, Hong Kong SAR, China
- 23 Email: cehyang@ust.hk
- 24

25 Dengbo He

- 26 Assistant Professor, Corresponding author
- 27 Intelligent Transportation Thrust, Systems Hub
- 28 The Hong Kong University of Science and Technology (Guangzhou), China
- 29 HKUST Shenzhen-Hong Kong Collaborative Innovation Research Institute, Futian, Shenzhen
- 30 Email: dengbohe@hkust-gz.edu.cn
- 31 32 33 Word Count: 1,693 words + 4 tables/figures
- 34 Submitted [November, 19, 2024] 35

36 Statement of Significance (Relevance of Research)

This research introduces CogFormer, a decision-level fusion model to accurately estimate 37 38 driver cognitive load in conditional automated vehicles. By integrating physiological signals 39 with an attention mechanism, CogFormer provides robust and real-time estimation of driver state, surpassing existing models in accuracy. The findings have implications for enhancing 40 safety in SAE Level 3 and above by enabling cognitive load estimation for adaptive support 41 42 systems, including but not limited to takeover assistance and fatigue monitoring. This research 43 will be relevant to TRB attendees who are interested in driver monitoring, intelligent 44 transportation systems, and smart cabins.

45

46 Acknowledgments

This work was supported by the Natural Science Foundation of Guangdong Province of China
(2024A1515010392) and partially by the National Natural Science Foundation of China
(52202425).

- 50
- 51

1 **Author Contribution**

- The authors confirm their contribution to the paper as follows: Ange Wang: Conceptualization, 2
- Data curation, Formal analysis, Methodology, Software, Validation, Writing original draft. 3
- Haohan Yang: Validation, Formal analysis, coding assistance. Jiyao Wang: Validation, 4
- Formal analysis. Hai Yang: Validation, Formal analysis. Dengbo He: Formal analysis, 5 6 7
- Funding acquisition, Methodology, Supervision, Validation, Writing review & editing.

1 INTRODUCTION

Human error is recognized as one of the dominating factors in road accidents (1). Compared to driving tasks that are visually and manually demanding, the cognitive demanding tasks can be more safety-critical, and thus drivers' performance in these tasks has been widely adopted as key metrics differentiating novice and experienced drivers (2). The introduction of infotainment functions in the smart cabin and the prevalence of bring-in smart devices may also increase the task load of drivers.

8 The high cognitive load in driving has been found to be closely related to driving 9 safety. For example, a high cognitive load may lead to delayed responses to emergency events 10 (3), visual tunnel effect (4), decreased ability to anticipate hazards (5). The introduction of 11 driving automation may be a solution, as a lower overall task load has been observed in vehicles 12 equipped with advanced driving automation systems (ADASs) (6). However, some studies 13 found that ADAS can increase drivers' cognitive load due to the additional responsibility to 14 monitor automation (7).

Most of the previous cognitive load estimation algorithms were developed for nonautomated vehicles, which may not apply to vehicles with ADAS. For example, the study by Wiediartini et al. (8) revealed significant differences in fixation duration and pupil diameter across the three task difficulty levels of memory and arithmetic tasks. This discrepancy in physiological measures has also been observed between manual and automated driving.

To the best of our knowledge, we can only identify three studies that focused on high cognitive load estimation in vehicles with driving automation (9-11), and most of the existing driver cognitive estimation approaches (including the ones for non-automated vehicles and vehicles with driving automation) have two major limitations:

Most previous models relied on manual feature extraction of physiological features (3, 11, 12), for example, heart rate (HR) extracted from electrocardiogram (ECG) (13), and respiratory rate extracted from respiratory signals (RESP) (14). While satisfactory accuracy has been achieved in previous research, the extraction of these handcrafted low-level features is computational costing and may lead to loss of information in the raw signals.

Most driver cognitive load estimation studies used classical machine learning models that ignored the temporal-spatial dependency of physiological signals. Only a few studies used Recurrent Neural Networks (RNNs) (12) and Long Short-Term Memory (LSTM) networks (15) that could consider temporal information for driver cognitive load estimation. However, RNN or LSTM may still neglect the long-distance temporal dependencies present in time series and previous studies did not fuse multiple physiological features, ignoring the spatial dependency among signals.

To address the aforementioned challenges, in this study, we proposed an Attention-Aligned Transformer algorithm for Cognitive (CogFormer) load estimation in vehicles with driving automation by integrating ECG, electrodermal activity (EDA), and RESP signals.

39

40 METHODOLOGY

The aligned attention mechanism is the most important component of CogFormer which 41 42 integrates multiple physiological signals, i.e., ECG, RESP, and EDA signals, by utilizing self-43 attention and aligned attention. Self-attention captures dependencies within each signal type, while aligned attention aligns and integrates information across different signals. This approach 44 45 allows the model to weigh and combine the most relevant features from each signal, creating a cohesive and comprehensive representation. The ability of the mechanism to fuse diverse data 46 sources can allow the model to understand and predict physiological patterns by capturing 47 48 complex interdependencies. Cross-modal attention and self-attention mechanisms are

49 expressed by equations (1) and (2), respectively.

$$f_{ii} = softmax \left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right) V_i, \tag{1}$$



2

$$f_{ij} = softmax \left(\frac{Q_i K_j^T}{\sqrt{d_k}}\right) V_j, i, j \in s \text{ and } i \neq j.$$
⁽²⁾





FIGURE 1 Overview of the proposed CogFormer for driver cognitive load detection.

6 EXPERIMENT AND DATASET

7 We tested our model on the mathematical and autonomous driving tasks dataset (MADT-D) 8 published by the University of Applied Sciences and Arts of Western Switzerland (16). The 9 MADT-D consists of driving data from 90 participants (with two participants deemed invalid). In the experiment, half of the participants were instructed to perform a cognitive task known as 10 11 the oral digit span counting task, requiring them to verbally count backward from 3,645 in 12 decrements of 2 (labeled as high cognitive load), while the rest half conducted a driving task only (labeled as low cognitive load). To ensure that both load levels appear in the leave-one-13 14 out test, we combined one participant performing the non-driving-related tasks (NDRTs) with one participant not performing the NDRT to form a new subject. Consequently, there are a total 15 16 of 44 new participants. Similarly, in the experiment, the physiological data, i.e., EDA, RESP, 17 and ECG, were collected using sensors by BioPac at a frequency of 1,000Hz in the dataset. The differences in subjective ratings of cognitive load between tasks were validated in Meteier et al. 18 (9). We extracted data spanning 10 minutes from the cognitive load phase in the MADT-D 19 dataset. 20

Apart from our own dataset, to better validate our proposed model, we also constructed a dataset, named CAM-CLD, which focuses on drivers' cognitive states in simulated SAE Level-3 vehicles. The details of the experiment for CAM-CLD are provided below. This study was approved by the Hong Kong University of Science and Technology (HREP-2023-0199).

26 **Participants**

In total, 42 drivers (25 males, 17 females, aged 23-53) were recruited, ensuring a balanced age
distribution to improve generalizability. All had at least one year of licensed driving experience
and were compensated at a rate of 70 RMB per hour.

30

25

31 Cognitive Load Tasks

- 32 Three kinds of NDRTs used in CAM-CLD, i.e., the n-back task in which participants needed
- to recall stimuli presented n positions earlier (0-, 1-, 2-back) (17); math calculation task, in

- 1 which the participants needed to verbally count backward from 3000 by a step of 3 or 5 (9);
- 2 and spatial memory, in which participants needed to recall the final direction after listening to
- 3 an audio description of a route (18).
- 4

5 Experiment Design

- A 3 (Takeover Scenarios) by 7 (NDRTs) within-subject design was used. Each combination of
 NDRT type and scenario type happened in one drive, leading to a total of 21 drives. Each drive
- 8 was around 7 minutes long, and the takeover scenarios happened near the end of the drive. The
- 9 3 takeover scenarios were counterbalanced using a Latin squared design. For each kind of
- 10 takeover scenario, the 7 NDRT types were also counterbalanced using a Latin squared design,
- 11 leading to 21 unique experimental orders. Each order was experienced by two participants.
- 12 Further, after completing all 7 drives for one kind of takeover scenario, participants took a 10-
- 13 minute break before proceeding to the next 7 drives.
- 14

15 **Procedure**

- 16 Participants followed pre-experiment instructions and completed a 30-minute orientation. Then,
- physiological sensors (ECG, RESP, EDA) and eye-tracking devices were put on and calibrated
- 18 before the formal drives. All the physiological data was collected at 100 Hz.
- 19

20 Signal Preprocessing

Given that we used the raw data as inputs in the model, only noise elimination was conducted to enhance the quality of the data. Specifically, all signals underwent down-sampling to a frequency of 100 Hz (from 1000Hz in MADT-D) to optimize computational efficiency and ensure the consistency between two datasets. For EDA, a low-pass filter with a cutoff frequency of 5 Hz was employed; while for ECG and RESP, band-pass filters were applied within the frequency ranges of 3 Hz to 45 Hz and 0.1 Hz to 0.35 Hz, respectively (9). The preprocessing

- was executed using Python 3.8.
- 28

29 **RESULTS**

30 Experimental Results and Analysis

We compared our model to selected learning-based approaches from prior studies for cognitive load estimation in driving, including MTS-CNN (19), DecNet (20), CNN-LSTM (21), m-HyperLSTM (22) and ARecNet (23). As shown in Table 1, the model comparison encompassed various cognitive load tasks, and classification categories, showcasing the model performance across different time horizons (i.e., the length of physiological signal time series, represented by t_w).

37 It can be found that m-HyperLSTM consistently outperformed MTS-CNN and DecNet. 38 Nevertheless, our proposed CogFormer consistently achieved the highest accuracy in all models. Furthermore, we observed that for some tasks, increasing the time horizon did not necessarily 39 40 lead to an increase in recognition accuracy. The increased time horizon may have introduced some noise in the inputs for some tasks. Although the CogFormer can capture long-term 41 temporal information, it may still not be capable enough to automatically pay attention to high-42 43 value information in the data. An improved attention mechanism may be needed to improve the 44 capability of the model in filtering noise information in long-term temporal sequences.

Model		$t_w = 1 s$			$t_w = 3 s$			$t_w = 5 s$	
	Accuracy (%)	F1-score	AUC	Accuracy (%)	F1-score	AUC	Accuracy (%)	F1-scores	AUC
MADT-D: Math task (two classes)									
MTS-CNN [∆]	87.71 ± 4.86	0.876 ± 0.044	0.877 ± 0.046	88.83 ± 4.14	0.884 ± 0.042	0.882 ± 0.049	88.29 ± 4.03	0.887 ± 0.042	0.885 ± 0.046
DecNet∆	86.76 ± 3.29	0.866 ± 0.035	0.863 ± 0.033	87.27 ± 4.47	0.877 ± 0.046	0.876 ± 0.041	87.45 ± 4.34	0.876 ± 0.043	0.872 ± 0.049
CNN-LSTM [∆]	87.51 ± 4.14	0.871 ± 0.043	0.879 ± 0.042	87.87 ± 5.63	0.876 ± 0.053	0.874 ± 0.048	88.38 ± 3.59	0.888 ± 0.033	0.881 ± 0.036
m-HyperLSTM [∆]	89.68 ± 5.26	0.898 ± 0.055	0.896 ± 0.053	89.32 ± 5.19	0.894 ± 0.053	0.881 ± 0.051	89.14 ± 4.72	0.892 ± 0.045	0.893 ± 0.046
$\operatorname{ARecNet}^{\Phi}$	92.46 ± 4.22	0.921 ± 0.045	0.919 ± 0.040	91.93 ± 3.12	0.913 ± 0.031	0.911 ± 0.033	92.56 ± 3.38	0.921 ± 0.031	0.922 ± 0.032
$\mathbf{CogFormer}\ (\mathbf{ours})^{\Phi}$	$\textbf{95.28} \pm 3.98$	$\textbf{0.952} \pm 0.039$	$\textbf{0.955} \pm 0.037$	$\textbf{93.79} \pm 4.25$	$\textbf{0.938} \pm 0.047$	$\textbf{0.936} \pm 0.044$	$\textbf{94.03} \pm 3.81$	$\textbf{0.944} \pm 0.039$	$\textbf{0.947} \pm 0.035$
			CA	M-CLD: Spacial ta	ask (two classes)				
MTS-CNN [∆]	86.08 ± 2.85	0.863 ± 0.026	0.865 ± 0.027	87.34 ± 2.39	0.874 ± 0.022	0.879 ± 0.021	88.57 ± 1.97	0.882 ± 0.022	0.885 ± 0.021
DecNet∆	86.91 ± 2.43	0.867 ± 0.028	0.863 ± 0.023	86.23 ± 1.47	0.867 ± 0.017	0.866 ± 0.016	88.06 ± 2.44	0.888 ± 0.024	0.882 ± 0.025
CNN-LSTM^{Δ}	89.46 ± 2.32	0.899 ± 0.027	0.898 ± 0.026	88.41 ± 2.52	0.887 ± 0.027	0.884 ± 0.023	89.46 ± 2.37	0.892 ± 0.026	0.894 ± 0.023
m-HyperLSTM [∆]	88.64 ± 3.76	0.884 ± 0.034	0.889 ± 0.036	88.96 ± 1.29	0.886 ± 0.016	$\textbf{0.884} \pm \textbf{0.014}$	89.32 ± 3.45	0.898 ± 0.032	0.893 ± 0.036
$ARecNet^\Phi$	91.92 ± 2.46	0.912 ± 0.027	0.909 ± 0.028	90.09 ± 1.33	0.897 ± 0.015	0.890 ± 0.013	91.49 ± 2.96	0.916 ± 0.029	0.918 ± 0.027
$\mathbf{CogFormer}\ (\mathbf{ours})^{\Phi}$	$\textbf{93.05} \pm 2.39$	$\textbf{0.932} \pm 0.026$	$\textbf{0.931} \pm 0.023$	$\textbf{91.68} \pm 1.68$	$\textbf{0.918} \pm 0.013$	$\textbf{0.921} \pm 0.018$	$\textbf{93.63} \pm 2.12$	$\textbf{0.934} \pm 0.023$	$\textbf{0.931} \pm 0.022$
			CA	M-CLD: Math tas	k (three classes)		•		
MTS-CNN [∆]	85.16 ± 3.18	0.854 ± 0.035	0.852 ± 0.036	86.28 ± 3.13	0.867 ± 0.034	0.864 ± 0.032	86.99 ± 2.35	0.868 ± 0.021	0.865 ± 0.026
DecNet^{Δ}	85.59 ± 2.75	0.858 ± 0.025	0.853 ± 0.026	86.64 ± 2.92	0.863 ± 0.027	0.863 ± 0.029	87.37 ± 2.51	0.870 ± 0.025	0.872 ± 0.024
CNN-LSTM [∆]	86.86 ± 2.27	0.864 ± 0.022	0.862 ± 0.021	87.59 ± 2.23	0.879 ± 0.021	0.877 ± 0.023	87.98 ± 3.12	0.879 ± 0.034	0.872 ± 0.037
m-HyperLSTM [∆]	87.99 ± 2.36	0.872 ± 0.025	0.876 ± 0.023	87.87 ± 2.13	0.875 ± 0.023	0.879 ± 0.019	88.25 ± 2.55	0.888 ± 0.022	0.885 ± 0.027
$\operatorname{ARecNet}^{\Phi}$	88.28 ± 2.16	0.889 ± 0.020	0.884 ± 0.022	88.17 ± 3.37	0.883 ± 0.039	0.882 ± 0.032	89.23 ± 2.13	0.894 ± 0.025	0.901 ± 0.022
CogFormer (ours) $^{\Phi}$	$\textbf{91.13} \pm 1.31$	$\textbf{0.909} \pm 0.012$	$\textbf{0.910} \pm 0.015$	$\textbf{90.64} \pm 2.43$	$\textbf{0.902} \pm 0.023$	$\textbf{0.904} \pm 0.022$	$\textbf{91.92} \pm 2.36$	$\textbf{0.920} \pm 0.022$	$\textbf{0.921} \pm 0.026$
CAM-CLD: n-back task (four classes)									
MTS-CNN [∆]	84.35 ± 3.08	0.841 ± 0.032	0.842 ± 0.031	84.93 ± 2.55	0.846 ± 0.026	0.845 ± 0.023	85.33 ± 2.14	0.856 ± 0.023	0.852 ± 0.021
DecNet∆	84.26 ± 2.09	0.848 ± 0.021	0.845 ± 0.019	85.01 ± 2.58	0.855 ± 0.025	0.851 ± 0.026	85.43 ± 2.33	0.851 ± 0.024	0.856 ± 0.024
CNN-LSTM^{Δ}	85.88 ± 2.94	0.859 ± 0.029	0.850 ± 0.030	85.12 ± 2.64	0.858 ± 0.024	0.852 ± 0.025	85.36 ± 2.19	0.858 ± 0.023	0.851 ± 0.021
m-HyperLSTM [∆]	85.85 ± 2.57	0.854 ± 0.025	0.857 ± 0.026	84.47 ± 2.26	0.845 ± 0.023	0.841 ± 0.021	86.17 ± 2.21	0.864 ± 0.024	0.867 ± 0.026
$\operatorname{ARecNet}^{\Phi}$	86.17 ± 1.52	0.862 ± 0.016	0.869 ± 0.014	85.13 ± 1.54	0.858 ± 0.016	0.852 ± 0.018	87.42 ± 1.74	0.872 ± 0.017	0.870 ± 0.016
CogFormer (ours) $^{\Phi}$	$\textbf{88.72} \pm 1.78$	$\textbf{0.890} \pm 0.019$	$\textbf{0.892} \pm 0.018$	$\textbf{87.67} \pm 1.44$	$\textbf{0.878} \pm 0.014$	$\textbf{0.872} \pm 0.015$	$\textbf{88.14} \pm 1.88$	$\textbf{0.882} \pm 0.018$	$\textbf{0.885} \pm 0.017$

TABLE 1 Comparison within-subject cognitive load estimation.

Notes: Δ means feature-level fusion model, Φ means decision-level fusion model.

1 Ablation Experiment

Ablation experiments revealed that the Aligned Attention module consistently enhances spatial
feature capture, outperforming traditional transformers. Multi-stream Encoding excels with shorter
historical horizons, while concatenated encoding improves with longer horizons. Results indicate
that Aligned Attention strengthens multimodal information representation in extended sequences
(see Table 2).

7

8 TABLE 2 Ablation study with different historical horizons on recognition accuracy with various

9 <u>cognitive tasks.</u>

	Var	CogFormer					
Multi-stream Encoding	×						
Aligned Attention	×	×					
MADT-D: Math task							
$t_w = 1 s$	91.41 (↓3.87%)	91.85 (↓ 3.43%)	95.28				
$t_w = 3 s$	90.47 (J3.32%)	90.64 (↓3.55%)	93.79				
$t_w = 5 s$	91.58 (↓ 2.45%)	91.06 (↓2.97%)	94.03				
CAM-CLD: Spacial task							
$t_w = 1 s$	88.51 (↓4.54%)	89.23 (↓ 3.82%)	93.05				
$t_w = 3 s$	88.03 (↓3.65 %)	88.34 (↓ 3.34%)	91.68				
$t_w = 5 s$	90.61 (↓ 3.02%)	90.50 (↓3.13%)	93.63				
CAM-CLD: Math task							
$t_w = 1 s$	87.70 (↓3.43%)	87.87 (↓ 3.26%)	91.13				
$t_w = 3 s$	85.99 (↓4.65%)	86.99 (↓ 3.65%)	90.64				
$t_w = 5 s$	88.10 (\3.82%)	88.21 (↓ 3.71%)	91.92				
CAM-CLD: n-back task							
$t_w = 1 s$	82.96 (↓5.76%)	84.18 (↓ 4.54%)	88.72				
$t_w = 3 s$	82.96 (↓4.71%)	83.66 (↓ 4.01%)	87.67				
$t_w = 5 s$	84.91 (↓ 3.23%)	84.93 (↓3.59%)	88.14				

10

11 Robustness Testing

In practical applications, time series signals are affected by missing values and Gaussian White
 Noise (GWN). Robustness tests showed that CogFormer consistently outperforms ARecNet in

14 handling these distortions, except slightly in the math task under mixed conditions. This highlights

15 CogFormer's superior robustness with its parallel transformer and coherent attention mechanism,

see Table 3.

Model	Normal (%)	Missing 10% (M)	Missing 20% (M)	Missing 30% (α)	Missing 30% (β)	Missing 30% (γ)	GWN σ=0.1	GWN σ=0.2	Mixed distortion
MADT-D: Math task									
CogFormer	94.03	92.80	91.36	92.39	92.70	92.74	93.61	92.30	90.80
(Proposed)		(\1.23)	(\$2.67)	(↓1.64)	(\1.33)	(\1.29)	(↓0.42)	(1.73)	(↓3.23)
ARecNet	92.56	90.81	88.75	90.20	90.35	91.18	91.85	90.00	87.14
(Contrastive)		(\1.75)	(↓3.81)	(\12.36)	(↓2.21)	(\1.38)	(↓0.71)	(\12.56)	(↓5.42)
CAM-CLD: Spacial task									
CogFormer	02.62	91.87	91.51	91.86	92.39	91.60	93.27	91.18	89.61
(Proposed)	93.63	(↓1.76)	(↓2.12)	(↓1.77)	(↓1.24)	(↓2.03)	(↓0.36)	(↓2.45)	(↓ 4.02)
ARecNet	01 /0	89.36	87.98	89.52	89.46	88.94	90.82	88.33	86.87
(Contrastive)	91.49	(\12.13)	(\$3.51)	(↓1.97)	(\1.03)	(\12.55)	(↓0.67)	(\$3.16)	(↓4.62)
CAM-CLD: Math task									
CogFormer	91.92	90.29	88.81	90.18	90.79	89.77	91.29	89.76	86.39
(Proposed)		(\1.63)	(↓3.11)	(↓1.74)	(\1.13)	(\12.15)	(↓0.63)	(12.16)	(↓5.53)
ARecNet	89.23	87.71	86.21	86.78	87.31	86.77	88.88	86.47	83.99
(Contrastive)		(\1.52)	(\$\$.02)	(\1.45)	(\1.92)	(\12.46)	(↓0.35)	(\1.76)	(↓5.24)
CAM-CLD: n-back task									
CogFormer	88.14	85.50	83.71	85.83	86.12	86.00	87.22	84.80	81.73
(Proposed)		(\12.64)	(↓4.43)	(↓2.31)	(\$2.02)	(↓2.14)	(↓0.92)	(↓3.34)	(↓6.41)
ARecNet	87.42	83.89	82.67	85.25	85.89	85.22	86.09	83.41	80.20
(Contrastive)		(\$\$.53)	(↓4.75)	(↓2.17)	(\1.53)	(\.2.20)	(\1.33)	(↓4.01)	(↓7.22)

TABLE 3 Comparison of recognition accuracy in the decision-fusion models with varied information distortions when tw=5s.

3 Note: ECG, RESP, EDA - \mathcal{M} , ECG - α , RESP - β , and EDA - γ , GWN - σ .

4

5 DISCUSSION AND CONCLUSION

Being different from previous approaches that utilized traditional machine/deep learning models,
we developed a decision-level multi-physiological information fusion architecture to extract
temporal and spatial information from multiple physiological signals. Experimental results
demonstrate that the proposed CogFormer surpassed other baseline models in terms of estimation
accuracy and robustness. In addition, based on the model ablation study and robustness test, we
find that:

12 The preferred feature combinations for driver state estimation may depend on the type of 13 targeted tasks, data collection quality, and driving context. Thus, the performance of models based 14 on handcrafted features and manual feature selection may not be guaranteed in real-world 15 applications.

A longer time horizon, though can provide richer information, may not necessarily increase
 the model performance, potentially because the models may not be able to capture the complex
 features in the data. Thus, the degree of matching between the models and the characteristics

19 of data should be considered when designing driver state-monitoring algorithms.

- All models, including our proposed model, were highly susceptible to individual differences
 the performance of all models dropped significantly when across-subjects data partition was
 applied, indicating that the models, even with the attention mechanism, may still not be able
 to capture the individual-invariant features of high cognitive load states. Future research
 should design specific algorithms and structures to handle this issue.
- Though our algorithm was not designed specifically to handle the noise and data distortion, our model showed better robustness compared to the best baseline model. The noise and missing data are common in real-world applications. Thus, future research should consider a more specific algorithm design and validate the proposed model based on real-world datasets.
- 10
- 11

12 **REFERENCES**

- 13 1. Singh, S. Critical Reasons for Crashes Investigated in the National Motor Vehicle Crash
- 14 Causation Survey. *Traffic Safety Facts Crash Stats*, 2015.
- 15 2. Jackson, L., P. Chapman, and D. Crundall. What Happens next? Predicting Other Road Users'
- 16 Behaviour as a Function of Driving Experience and Processing Time. *Ergonomics*, Vol. 52, No.
- 17 2, 2009, pp. 154–164. https://doi.org/10.1080/00140130802030714.
- 18 3. Harbluk, J. L., Y. I. Noy, P. L. Trbovich, and M. Eizenman. An On-Road Assessment of
- 19 Cognitive Distraction: Impacts on Drivers' Visual Behavior and Braking Performance. Accident
- 20 Analysis & Prevention, Vol. 39, No. 2, 2007, pp. 372–379.
- 21 https://doi.org/10.1016/j.aap.2006.08.013.
- 22 4. Recarte, M. A., and L. M. Nunes. Effects of Verbal and Spatial-Imagery Tasks on Eye
- 23 Fixations While Driving. Journal of Experimental Psychology: Applied, Vol. 6, No. 1, 2000, pp.
- 24 31–43. https://doi.org/10.1037/1076-898X.6.1.31.
- 25 5. Muhrer, E., and M. Vollrath. The Effect of Visual and Cognitive Distraction on Driver's
- 26 Anticipation in a Simulated Car Following Scenario. Transportation Research Part F:
- 27 *Psychology and Behaviour*, Vol. 14, No. 6, 2011, pp. 555–566.
- 28 https://doi.org/10.1016/j.trf.2011.06.003.
- 29 6. Huang, Chunxi, Xie, Weiyin, Huang, Qihao, Zhu, Yan, Cui, Dixiao, and He, Dengbo. Effect
- 30 of Advanced Driver Assistance Systems on Fatigue Levels of Heavy Truck Drivers in Prolonged
- Driving Tasks. *Journal of Tongji University (Natural Science)*, Vol. 52, No. 6, 2024, pp. 846–
 855.
- 33 7. Stapel, J., F. A. Mullakkal-Babu, and R. Happee. Automated Driving Reduces Perceived
- 34 Workload, but Monitoring Causes Higher Cognitive Load than Manual Driving. *Transportation*
- 35 *Research Part F: Traffic Psychology and Behaviour*, Vol. 60, 2019, pp. 590–605.
- 36 https://doi.org/10.1016/j.trf.2018.11.006.
- 8. Wiediartini, U. Ciptomulyono, and R. S. Dewi. Evaluation of Physiological Responses to
- 38 Mental Workload in N-Back and Arithmetic Tasks. *Ergonomics*, 2023, pp. 1–13.
- 39 https://doi.org/10.1080/00140139.2023.2284677.
- 40 9. Meteier, Q., M. Capallera, S. Ruffieux, L. Angelini, O. Abou Khaled, E. Mugellini, M.
- 41 Widmer, and A. Sonderegger. Classification of Drivers' Workload Using Physiological Signals
- 42 in Conditional Automation. *Frontiers in Psychology*, Vol. 12, 2021, p. 596038.
- 43 https://doi.org/10.3389/fpsyg.2021.596038.
- 44 10. Shi, W., Z. Wang, A. Wang, and D. He. Classification of Driver Cognitive Load in
- 45 Conditionally Automated Driving: Utilizing Electrocardiogram-Based Spectrogram with

- 1 Lightweight Neural Network. *Transportation Research Record*, 2024.
- 2 https://doi.org/10.1177/03611981241252797.
- 3 11. Wang, A., J. Wang, W. Shi, and D. He. Cognitive Workload Estimation in Conditionally
- 4 Automated Vehicles Using Transformer Networks Based on Physiological Signals.
- 5 *Transportation Research Record*, 2024. https://doi.org/10.1177/03611981241250023.
- 6 12. Kumar, S., D. He, G. Qiao, and B. Donmez. Classification of Driver Cognitive Load Based
- 7 on Physiological Data: Exploring Recurrent Neural Networks. Presented at the 2022
- 8 International Conference on Advanced Robotics and Mechatronics (ICARM), Guilin, China,
- 9 2022.
- 10 13. He, D., Z. Wang, E. B. Khalil, B. Donmez, G. Qiao, and S. Kumar. Classification of Driver
- 11 Cognitive Load: Exploring the Benefits of Fusing Eye-Tracking and Physiological Measures.
- *Transportation Research Record: Journal of the Transportation Research Board*, Vol. 2676, No.
 10, 2022, pp. 670–681. https://doi.org/10.1177/03611981221090937.
- 14 14. Qu, Y., H. Hu, J. Liu, Z. Zhang, Y. Li, and X. Ge. Driver State Monitoring Technology for
- 15 Conditionally Automated Vehicles: Review and Future Prospects. IEEE Transactions on
- 16 *Instrumentation and Measurement*, Vol. 72, 2023, pp. 1–20.
- 17 https://doi.org/10.1109/TIM.2023.3301060.
- 18 15. Ansari, S., F. Naghdy, H. Du, and Y. Pahnwar. Driver Mental Fatigue Detection Based on
- 19 Head Posture Using New Modified reLU-BiLSTM Deep Neural Network. *IEEE*
- 20 TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS, Vol. 23, No. 8, 2022, pp.
- 21 10957–10969. https://doi.org/10.1109/TITS.2021.3098309.
- 22 16. Meteier, Q., M. Capallera, E. de Salis, L. Angelini, S. Carrino, M. Widmer, O. Abou Khaled,
- 23 E. Mugellini, and A. Sonderegger. A Dataset on the Physiological State and Behavior of Drivers
- in Conditionally Automated Driving. *Data in Brief*, Vol. 47, 2023.
- 25 https://doi.org/10.1016/j.dib.2023.109027.
- 26 17. Jaeggi, S. M., M. Buschkuehl, W. J. Perrig, and B. Meier. The Concurrent Validity of the N-
- 27 Back Task as a Working Memory Measure. *Memory*, Vol. 18, No. 4, 2010, pp. 394–412.
- 28 https://doi.org/10.1080/09658211003702171.
- 29 18. Liang, Y., and J. D. Lee. Combining Cognitive and Visual Distraction: Less than the Sum of
- 30 Its Parts. Accident Analysis & Prevention, Vol. 42, No. 3, 2010, pp. 881–890.
- 31 https://doi.org/10.1016/j.aap.2009.05.001.
- 32 19. Xie, Y., Y. L. Murphey, and D. S. Kochhar. Personalized Driver Workload Estimation Using
- Deep Neural Network Learning From Physiological and Vehicle Signals. *IEEE Transactions on Intelligent Vehicles*, Vol. 5, No. 3, 2020, pp. 439–448.
- 35 https://doi.org/10.1109/TIV.2019.2960946.
- 36 20. Amadori, P. V., T. Fischer, R. Wang, and Y. Demiris. Predicting Secondary Task
- 37 Performance: A Directly Actionable Metric for Cognitive Overload Detection. *IEEE*
- 38 Transactions on Cognitive and Developmental Systems, Vol. 14, No. 4, 2022, pp. 1474–1485.
- 39 https://doi.org/10.1109/TCDS.2021.3114162.
- 40 21. Huang, J., Y. Liu, and X. Peng. Recognition of Driver's Mental Workload Based on
- 41 Physiological Signals, a Comparative Study. *BIOMEDICAL SIGNAL PROCESSING AND*
- 42 *CONTROL*, Vol. 71, 2022. https://doi.org/10.1016/j.bspc.2021.103094.
- 43 22. Wang, R., P. V. Amadori, and Y. Demiris. Real-Time Workload Classification during
- 44 Driving Using HyperNetworks. Presented at the 2018 IEEE/RSJ International Conference on
- 45 Intelligent Robots and Systems (IROS), 2018.

- 1 23. Yang, H., J. Wu, Z. Hu, and C. Lv. Real-Time Driver Cognitive Workload Recognition:
- 2 Attention-Enabled Learning with Multimodal Information Fusion. *IEEE Transactions on*
- 3 *Industrial Electronics*, 2023, pp. 1–11. https://doi.org/10.1109/TIE.2023.3288182.