

---

1     **The Association between Physiological and Eye-tracking**  
2     **Metrics and Cognitive Load in Drivers: A Meta-analysis**

3  
4     Ange Wang<sup>a,c</sup>, Chunxi Huang<sup>b,d</sup>, Jiyao Wang<sup>b</sup>, Dengbo He<sup>a,b,c\*</sup>

5     <sup>a</sup> *Thrust of Intelligent Transportation, Systems Hub, The Hong Kong University of*  
6     *Science and Technology (Guangzhou), Guangzhou, China*

7     <sup>b</sup> *Thrust of Robotics and Autonomous Systems, Systems Hub, The Hong Kong*  
8     *University of Science and Technology (Guangzhou), Guangzhou, China*

9     <sup>c</sup> *Department of Civil and Environmental Engineering, The Hong Kong University of*  
10    *Science and Technology, Hong Kong SAR, China*

11    <sup>d</sup> *Robotics and Autonomous Systems, Division of Emerging Interdisciplinary Areas*  
12    *(EMIA) under Interdisciplinary Programs Office (IPO), The Hong Kong University of*  
13    *Science and Technology, Hong Kong SAR, China*

---

# 1 ABSTRACT

2 Driving performance can be impaired by a high cognitive load of drivers. Thus, it is  
3 important to estimate drivers' cognitive load. Although physiological and eye-tracking  
4 metrics have been widely used in many studies to assess cognitive load while driving,  
5 conflicts still exist regarding the association between physiological and eye-tracking  
6 metrics and different levels of cognitive load. Through a meta-analysis, our study aims  
7 to quantify the association between physiological, eye-tracking metrics and cognitive  
8 load induced by n-back tasks. A total of 18 articles met the inclusion criteria for the  
9 meta-analysis. The results indicate four types of metrics, including the sensitive-to-low  
10 ones that can only differentiate the low to medium level of cognitive load (i.e., the  
11 power spectrum of  $\theta$  wave of electroencephalogram at Fp1 channel); high-resolution  
12 ones that can differentiate all levels of cognitive load (including pupil size, heart rate,  
13 and skin conductance); and low-resolution ones that can only differentiate low and high  
14 cognitive load (including the total power spectrum of electrocardiogram, eye blink rate,  
15 and respiration rate) and others (the power spectrum of  $\theta$  wave of electroencephalogram  
16 at Fp2 channel). Furthermore, the association between metrics and cognitive load can  
17 be modulated by the n-back version, modality of n-back task, automation level, and  
18 percentage of male participants. In summary, this study contributes to the literature by  
19 quantifying associations between physiological and eye-tracking metrics and different  
20 cognitive load levels. Practically, we provide evidence for the selection of physiological  
21 and eye-tracking metrics for future driving cognitive load monitoring system design.

---

1 **KEYWORDS**

2 Physiological Metrics, Cognitive Load, Driving, N-back Task, Meta-analysis

3

Post Print

---

# 1. INTRODUCTION

Human error is recognized as one of the dominating factors in road accidents (Singh, 2015). Though the human brain has long been recognized as a single-channel processor (C. D. Wickens, 1991), the driving task frequently involves multitasking. For example, during driving, in most cases, drivers need to control the vehicle (e.g., fine-tuning the gas pedal/brake pedal and the steering wheel to track the target speed and direction of the vehicle) and monitor the surrounding natural and traffic environment to identify potential hazards simultaneously. Different driving tasks may require different types of attention resources. Specifically, according to the multi-resource theory (Wickens, 1991), speed controlling task mainly requires visual-manual resources; while hazard perception tasks mainly require visual-cognitive resources. It has been commonly acknowledged that, compared to driving tasks that are visually and manually demanding, the cognitive demanding tasks, such as hazard perception and driving strategy selection are more safety-critical, and thus drivers' performance in these tasks is adopted as key metrics differentiating novice and experienced drivers (Jackson et al., 2009; Sagberg & Bjørnskau, 2006). In addition to driving-related tasks, in recent years, new technologies have been introduced into vehicles. For example, driving automation has been prevalent in newly sold vehicles, though it can reduce the overall workload of drivers, it may increase drivers' cognitive load because of the additional responsibility to monitor the automation (Stapel et al., 2019); while the introduction of infotainment functions in the smart cabin (e.g., video-streaming and

---

1 internet browsing) and the prevalence of the bring-in smart devices (e.g., smartphones)  
2 may also increase the workload of drivers.

3 The high cognitive load in driving, either as a result of driving tasks or non-  
4 driving-related tasks, has been found to be closely related to driving safety in research  
5 environments (e.g., in driving simulators or instrumented vehicles). For example, a high  
6 cognitive load may lead to delayed responses to emergency events (Harbluk et al.,  
7 2007), reduced visual search scope, leading to a visual tunnel effect (Recarte & Nunes,  
8 2000)), decreased ability to anticipate hazards (Muhrer & Vollrath, 2011), and  
9 increased reaction times (Du et al., 2020) and impaired performance (Melnicuk et al.,  
10 2021) in takeover events during assisted driving. Thus, estimating drivers' high  
11 cognitive load can be a potential approach to improve driving safety, both in vehicles  
12 with and without driving automation.

13 As an intrinsic state, cognitive load can hardly be measured objectively and  
14 directly (e.g., the questionnaire is direct but subjective, while eye-tracking measures are  
15 objective but indirect). Moreover, unlike distracted driving and fatigue, the state of high  
16 cognitive load is not easily discernible from the normal driving state, as it can be an  
17 integral part of the driving task. In the domain of driving, the cognitive load can be  
18 evaluated using four different types of measures, i.e., subjective measures, such as the  
19 NASA-Task Load Index (NASA-TLX) scale (Hart & Staveland, 1988); physiological  
20 measures, such as the electrocardiogram (ECG), electroencephalogram (EEG),  
21 respiration, and electrodermal activity (EDA); eye-tracking measures, such as pupil size;

---

1 and task performance measures, including driving task measures and non-driving-  
2 related task measures. All these measures have their pros and cons. For example, the  
3 subjective questionnaire methods cannot estimate cognitive load in real-time. The non-  
4 driving-related tasks can disrupt drivers' natural driving behavior and increase their  
5 cognitive load. The driving task measures are highly susceptible to traffic conditions  
6 and become invalid during automated driving, given that drivers do not need to control  
7 the vehicle for extended periods. Thus, the physiological measures and eye-tracking  
8 measures are most promising for real-time cognitive load detection. Further, with the  
9 advancement of new technologies, non-intrusive measures of physiological metrics and  
10 highly accurate eye-tracking measures have become possible (Ayres, 2020).

11 However, although associations between some physiological and eye-tracking  
12 metrics and the variations in cognitive load have been observed in some studies, no  
13 consensus has been reached for some other physiological and eye-tracking metrics,  
14 which poses challenges in selecting appropriate metrics for developing cognitive load  
15 detection algorithms in drivers. For example, the correlation between the respiration  
16 rate (RR) and cognitive load has been found to be negative in some studies (He et al.,  
17 2019), but positive in some other studies (Hajek et al., 2013). For other metrics,  
18 substantial differences in the strength of the association have been identified. For  
19 example, in Rahman et al (2020), the Low Frequency (LF) power (0.04–0.15 Hz) of  
20 heart rate variability (HRV) exhibited a strong positive correlation ( $r > 0.8$ ) with  
21 cognitive load. However, only a weak correlation ( $r < 0.1$ ) has been found in Tjolleng

---

1 et al (2017). Lastly, not all metrics are responsive in differentiating different levels of  
2 cognitive load. For example, the heart rate (HR) was able to differentiate between  
3 median to high levels of cognitive load but showed no difference between low to  
4 median levels of cognitive load (Ferreira et al., 2014). However, the feasible range of  
5 different cognitive load measures has not been systematically analyzed, which hinders  
6 the development of different algorithms targeting different levels of cognitive load,  
7 using minimum types of measures.

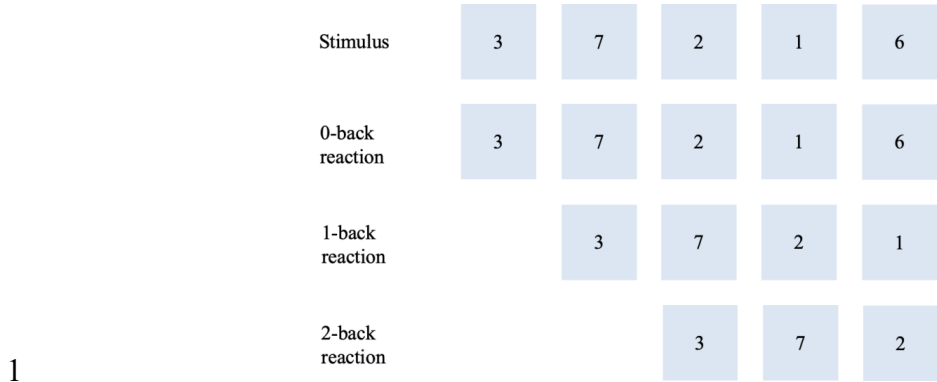
8 Thus, it is necessary to quantify the relationships between the physiological and  
9 eye-tracking metrics and the driver's cognitive load levels. Given that not all metrics  
10 are responsive to the whole range of cognitive load, the meta-analysis needs to be  
11 conducted for different ranges of cognitive load levels and a metric or task that can  
12 consistently impose different levels of cognitive load to drivers has to be selected.  
13 Traditionally, subjective responses such as NASA-TLX were regarded as the ground  
14 truth of cognitive load levels in previous studies (Chen et al., 2022; He et al., 2019;  
15 Hart & Staveland, 1988). Although NASA-TLX allows within-subject comparisons,  
16 individual differences in self-reported scores may exist and we can hardly compare the  
17 NASA-TLX scores across participants and experiments (Muth et al., 2012). Thus, in  
18 this study, we adopted a more standard task to label the levels of cognitive – the n-back  
19 task.

20 The n-back task is mainly a working memory task and has been proven as an  
21 effective manipulation of cognitive load in vehicles (Mehler et al., 2012a; von

---

1 Janczewski et al., 2021). In the n-back task, a series of stimuli (such as numbers or  
2 letters) are presented. Between each stimulus, there is a sustained pause that allows  
3 participants to repeat the stimuli presented  $n$  positions before (see an example in Figure  
4 1). The levels in the n-back task (i.e.,  $n$ ) reflect the difficulty and complexity of the task,  
5 the larger the  $n$ , the more difficult the task is. The utilization of 3-back is rare in research  
6 as 2-back is already sufficiently demanding. A substantial body of behavioral and  
7 neuroimaging research has confirmed the sensitivity of different levels of N-back tasks  
8 to cognitive load (Broadbent et al., 2023; Rieck et al., 2022; Solhjoo et al., 2019).  
9 Further, other commonly used cognitive tasks usually have no clear standard definition  
10 of cognitive levels. For example, the difficulty levels of hybrid tasks (obstacle  
11 avoidance and recall) (Yang et al., 2023) and mathematical tasks (von Janczewski et  
12 al., 2021) are difficult to quantify (e.g., the difficulty level of “3455 – 15” and “3455 –  
13 7” are not clearly quantified, though the former has one more digit than the latter).  
14 Given that the difficulty of 0-, 1- and 2-back tasks are comparable to common in-vehicle  
15 tasks (Mehler & Reimer, 2019) in our study, we treat the cognitive level triggered by  
16 the driving task only and 0-back task as low, the cognitive load triggered by 1-back task  
17 as median, and the cognitive load triggered by 2-back task as high.





**Fig. 1.** Demonstration of a typical n-back task procedure for n = 0, n = 1, and n = 2.

Therefore, in this study, based on a meta-analyses approach, we systematically analyzed the changes in drivers' physiological and eye-tracking metrics in response to the variation in cognitive load during driving, as defined by the levels of the n-back tasks. All metrics explored in this study were found to be associated with the cognitive load at least in some of the previous studies. To the best of our knowledge, though previous meta-analysis validated the effectiveness of n-back task in imposing high cognitive load in drivers (von Janczewski et al., 2021), no research has quantified the relationship between drivers' cognitive load and their physiological and eye-tracking measures. The summary of the abbreviations, descriptions, and units of the physiological and eye-tracking metrics mentioned in the text is shown in Table 1.

In addition, given that the cognitive load is a multi-dimensional concept and the settings in different studies can affect the responses of physiological and eye-tracking measures (Nilsson et al., 2022). we adopted meta-regression to account for 1) the influence of measurable individual differences so that future driver monitoring systems may take adaptive strategies; 2) and artificial experiment settings so that we can better

---

1 isolate the cognitive effects. Specifically, the demographic variables, including age  
2 (Wickens et al., 2011; Zhang et al., 2017) and gender (Sârbescu et al., 2014; Zhang et  
3 al., 2017) that can affect driving behavior, and the experiment-related settings including  
4 simulator fidelity, experimental environment, the modality of stimulus presentation and  
5 response, and the time interval between stimuli (Janczewski et al., 2021) that could  
6 modulate the drivers' cognitive load were considered. Besides, different versions of n-  
7 back tasks were used in previous driving studies. For example, the participants were  
8 required to only memorize and repeat the item that is n position before (**repeating**  
9 **version**, see Figure 1) in Gable et al (2015) and Mehler et al (2012b); while in Nilsson  
10 et al (2022) and Rahman et al (2020), participants needed to indicate whether the judge  
11 if the two items that are n positions apart are the same or not (**matching version**), which  
12 required additional cognitive resource to compare the two items. Another version of the  
13 n-back task was used by He et al (2019), which required participants to count how many  
14 times a pattern appeared (**counting version**) in addition to the matching version.  
15 Different versions require different cognitive components and hence the n-back version  
16 is also considered to account for the multi-dimensionality of the cognitive task.

17 In summary, the contribution of this study is 2-fold. First, the meta-analysis  
18 approach was employed to investigate the association between physiological, and eye-  
19 tracking metrics and cognitive load levels. Second, the moderating effects of  
20 participants' age and gender, driving automation level, the fidelity of simulators, n-back

1 **version**, modality of n-back stimulus and response, and the time interval between  
 2 stimuli on these associations were explored.

3 **Table 1.** Summary of the abbreviations, descriptions, and units of the physiological metrics.

Measure	Metric	Abbreviation	Description	Unit
Eye	Pupil size	PS	The diameter of the opening in the center of the iris	mm
	Fixation duration	FD	The length of time that gaze remains focused on a specific object or region of interest	ms
	<b>Eye blink rate</b>	<b>EBR</b>	<b>The number of blinks per unit of time</b>	<b>blinks/min</b>
	<b>Eyeblink duration</b>	<b>EBD</b>	<b>The time taken for each blink from closure to reopening of the eyes</b>	<b>ms</b>
Electroencephalography (EEG)	Theta wave-Fp1 channel	$\theta$ -Fp1	The $\theta$ waves (4-8 Hz) located in the Fp1 channel.	$\mu$ V
	Theta wave-Fp2 channel	$\theta$ -Fp2	The $\theta$ waves (4-8 Hz) located in the Fp2 channel.	$\mu$ V
Electrocardiogram (ECG)	Heart rate	HR	The number of heartbeats occurring per minute	Beats/minute
	Standard deviation of normal-to-normal intervals	SDNN	The variability of the time intervals between consecutive normal heartbeats	ms
	Root mean square of successive differences	RMSSD	The magnitude of the differences between consecutive R-R intervals (the time between successive heartbeats)	ms
	Low frequency	LF	The spectral power in the low-frequency range (usually 0.04 to 0.15 Hz) of the heart's electrical activity	ms <sup>2</sup>
	High frequency	HF	The spectral power in the high-frequency range (usually 0.15 to 0.4 Hz) of the heart's electrical activity	ms <sup>2</sup>
	Low frequency/High frequency	LF/HF	Ratio of LF to HF	%
	Total power	TP	The overall power spectrum	ms <sup>2</sup>
	pNN50	pNN50	The percentage of successive RR intervals (the time between R-peaks on an ECG) that differ by more than 50 milliseconds	%
	Very low frequency.	VLF	The frequency range of electrical signals in the ECG waveform that are below 0.04 Hz	Hz
	Low frequency power	LFun	The power or intensity of low-frequency electrical activity	ms <sup>2</sup>

	High frequency power	HFun	The power or intensity of high-frequency electrical activity	ms <sup>2</sup>
	Inter-beat interval	IBI	The time duration between successive heartbeats	ms
Skin	Electrodermal activity	SC-EDA	The general term that encompasses the electrical activity of the skin	μS
Respiration	Respiration rate	RR	The number of breaths taken per minute	Respirations / minute

1

## 2 2. METHOD

### 3 2.1 Literature search and study selection

4 We adopted the approach based on the PRISMA statement (Moher, 2009), a  
5 comprehensive guideline for reporting items in meta-analysis. An extensive literature  
6 search was conducted, covering articles published up until Feb 2024. The study  
7 selection process is summarized in Figure 2. In our search, the title, abstract or  
8 keywords must include ("driver" OR "driving" OR "automobile" OR "automated" OR  
9 "vehicle" OR "car") and ("cognitive load" OR "workload" OR "working memory" OR  
10 "mental workload"); and the full text must have "N-back" and ("physiological" OR  
11 "eye" OR "electroencephalogram" OR "electrocardiogram" OR "respiration" OR  
12 "electrodermal activity" OR "galvanic skin reaction" OR "skin conductance" OR  
13 "pulse rate variability" OR "temperature").

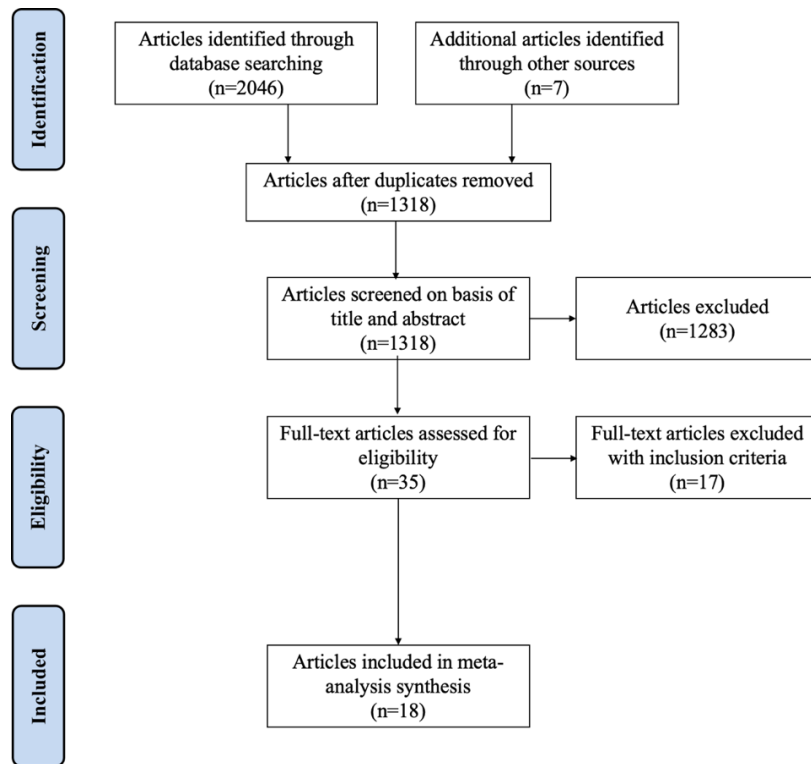


Fig. 2. Literature review process based on the PRISMA method.

## 2.2 Inclusion and exclusion criteria

The inclusion criteria for studies were as follows: (1) The n-back task should be a task secondary to the driving task of four-wheel cars and only empirical studies conducted on real roads or in simulated driving environments were included. (2) The study should have at least two different n-back levels or one n-back level and a baseline without a secondary task. (3) The statistics of the physiological or eye-tracking measures must be reported or could be obtained by contacting the authors. (4) The physiological and eye-tracking measures associated with the n-back level should be independent of other tasks in the vehicle so that the cognitive load was only induced by the n-back task. (5) Due to language barriers, only publications in Chinese and English were included. Publications that did not meet any of the above criteria were excluded

---

1 from the analysis. Initial determinations were made based on the abstracts, followed by  
2 a thorough examination of the full texts based on the inclusion criteria. Ultimately, 18  
3 studies were kept for further analysis.

### 4 **2.3 Data extraction**

5 The following information was extracted from each publication: (1) Meta-  
6 information of the study (i.e., title, author, and publication year); (2) Descriptions and  
7 measures of the cognitive workload; (3) sample characteristics (i.e., sample size, mean  
8 age, and characteristics of participants); (4) experimental conditions (i.e., automation  
9 level, experimental environment, simulator fidelity); (5) characteristics of n-back tasks  
10 (i.e., modality of the stimulus and response, and time interval between stimuli) and (6)  
11 associations between the physiological and eye-tracking metrics and n-back levels or  
12 the raw values of the metrics under each experimental condition.

13 These data were extracted and coded independently by two doctoral students (the  
14 first two authors). To ensure a consistent understanding of the coding scheme between  
15 the two coders, we conducted a preliminary "coding trial" phase. During this phase, the  
16 two coders independently coded five articles and discussed any discrepancies in coding  
17 to reach a consensus on the coding scheme. Necessary modifications and refinements  
18 were made to the coding manual based on the issues encountered during this phase.

### 19 **2.4 Data processing**

20 In the meta-analyses, the Pearson correlation coefficient ( $r$ ) was used as the effect  
21 size for each study. We followed a systematic eight-step process to analyze the data for

---

1 our meta-analyses. An overview of these steps is provided in Figure 3. In the meta-  
2 analyses, we focused on conducting meta-analyses of the physiological and eye-  
3 tracking metrics. Specifically, the physiological and eye-tracking metric from each  
4 level of cognitive load (as imposed by the n-back task or baseline task) was compared  
5 to all other levels of cognitive load. In total, six ( $C_3^2=3$  for baseline and 2 levels of n-  
6 back tasks) pairwise comparisons for all metrics were explored in this study. If the  
7 correlation coefficient was not provided in the study, we did the calculation according  
8 to Lipsey & Wilson (2001) based on the mean value, standard deviation (SD), and  
9 sample size. It is worth noting that the meta-analyses were only conducted for metrics  
10 that were investigated in at least two studies, which is the minimum requirement for a  
11 meta-analysis (McCarthy et al., 2017; Zheng, 2013).

12 Then, for each metric in each pairwise comparison, we transferred  $r$  (ranging  
13 between -1 and 1) to Fisher's  $Z$  value following the equation in Lipsey & Wilson (2001):

$$14 \quad Z_r = 0.5 * \ln((1 + r) / (1 - r)) \quad (1)$$

15 where  $Z_r$  represents Fisher's  $Z$  value, and  $r$  denotes the correlation coefficient  
16 between two variables. This transformation can alleviate the constraints of the  
17 correlation coefficient, as  $Z_r$  ranges from negative infinite to positive infinite, which  
18 enables the weighted combination of effect sizes from multiple studies in meta-analysis.

19 Next, the meta-analyses were conducted using RevMan 5.4 (Schmidt et al., 2019).  
20 The random-effects model was used to calculate the weighted average correlation, in  
21 which the calculation of weight factors was based on the inverse of the variance (Lipsey

---

1 & Wilson, 2001). The forest plots were provided to visualize the results (see Appendix  
2 1), in which, the pooled  $Z_r$  derived from the amalgamation of all incorporated studies  
3 was visually represented as a rhombus positioned at the lower section of the graph,  
4 wherein the breadth of the rhombus denoted the 95% confidence interval (95%CI). The  
5 significance of the estimated overall effect size was evaluated using the 95%CI (Lipsey  
6 & Wilson, 2001) and the  $p$ -value, with .05 as the significance threshold.

7 In addition, for each meta-analysis, the heterogeneity of the research was evaluated  
8 using the  $I^2$  statistic, tau squared ( $\tau^2$ ), and  $Q$  statistic (Lipsey & Wilson, 2001). The  $I^2$   
9 statistic quantifies the proportion of the observed variation in correlation that can be  
10 accounted for by actual variations between studies. A value of 25%, 50%, or 75%  
11 corresponds to low, moderate, or high levels of variance, respectively (Higgins & Deeks,  
12 2003). The  $\tau^2$  represents the overall extent of heterogeneity, with a smaller  $\tau^2$  indicating  
13 a lower level of heterogeneity. The  $Q$  statistic reflects the degree of heterogeneity  
14 resulting from actual differences between studies, with a significant  $Q$  statistic implying  
15 the existence of heterogeneity among the studies. Out of the above-mentioned metrics,  
16 the  $Q$ -test is the most used for testing heterogeneity. However, its testing power is  
17 limited when the number of studies is small. In contrast, the  $I^2$  statistic can mitigate the  
18 impact of the sample size on the testing power. In our study, we adopted the  $I^2$  of 50%  
19 as the threshold for the existence of heterogeneity (Zheng, 2013), but still reported  $\tau^2$   
20 and  $Q$  statistics for readers' reference.

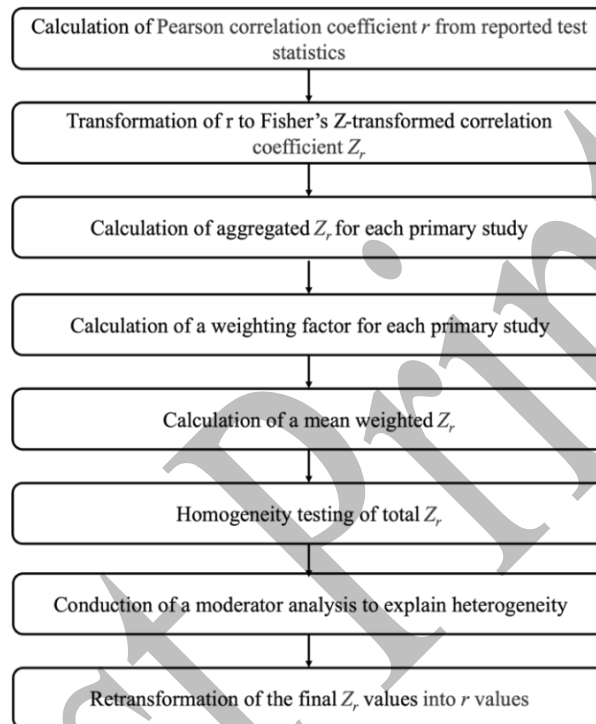


---

1 Finally, to investigate the origins of research heterogeneity, meta-regression was  
2 conducted for potential moderators, including experimental conditions (i.e., automation  
3 level), characteristics of n-back tasks (i.e., modality of the stimuli and responses, the  
4 time interval between stimuli), and demographic variables (i.e., mean age and  
5 percentage of male). The automation level here refers to driving automation defined by  
6 the Society of Automotive Engineers (SAE) (Committee, 2014). It is worth noting that  
7 we refined the simulator-type classification scheme proposed by Spyridakos et al.,  
8 (2020) as follows: "Occlusion" and "Desktop" were classified as the category of low-  
9 fidelity driving simulators, "Cabin with narrow field projection" and "Cabin with  
10 widefield projection" were classified as medium fidelity driving simulators; and  
11 "Hexapod," "Hexapod and lateral motion," and "Hexapod and longitudinal motion"  
12 were classified as high-fidelity driving simulators. Furthermore, since the median of the  
13 median age among the 18 studies included in the analysis is 27.2, we thus adopt the  
14 threshold of 27 years for age stratification.

15 Then, intergroup homogeneity was performed and heterogeneity coefficients were  
16 computed to assess the intergroup effects (Viechtbauer, 2010). For significant  
17 moderators, sub-group meta-analyses were conducted. Subsequently, we proceeded to  
18 assess the homogeneity of effect sizes within a specific group ( $Q_W$ ) and the  
19 heterogeneity across different sub-groups ( $Q_B$ ) (Lipsey & Wilson, 2001). The  
20 moderator analyses were conducted using Stata 17.

1 It should be noted that to facilitate interpretation, when reporting the results, the  
2  $Z_r$  was transformed back to  $r$ . According to Cohen (1988),  $|r| \leq 0.3$  denotes a small  
3 correlation,  $0.3 \leq |r| \leq 0.5$  denotes a medium correlation, and  $|r| \geq 0.5$  denotes a  
4 large correlation.



5

6

**Fig. 3.** Overview of the Processing Workflow.

### 7 **3. RESULT**

8 The process of literature retrieval and study selection is illustrated in Figure 3.

9 Initially, a total of 2,046 records were retrieved through database search, and an

10 additional 7 records were included after examining the reference list of the identified

11 literature. After removing duplicates, 1,318 records remained, which were then

12 screened based on their titles and abstracts. Then, full-text evaluation was conducted

13 for 35 records, and 18 records that met our criteria were kept for meta-analyses.

1 **Table 2.** A summary of the identified literature

Study	Participant Number	Metrics	Automation Level	Experimental Environment	N-back Levels	N-back Version	N-back Modalities (stimuli-responses)	Inter-stimulus Interval (s)	Male Ratio (%)	Mean Age (SD/Range)
(Deng et al., 2024)	20	HR, SC-EDA	L3	Medium fidelity simulator	N, 1B, 2B	Matching	Visual-verbal	—	90	25.3 (3.8/23-39)
(Nilsson et al., 2022)	70	HR, RMSSD, RR, SC-EDA, PS, EBR, EBD	L0	Medium fidelity simulator	N, 1B, 2B	Matching	Auditory-manual	1	100	43 (4/35-51)
(Mehler et al., 2012b)	108	HR, SC-EDA	L0	On-road	N, 0B, 1B, 2B	Repeating	Auditory-verbal	—	50	24.6 (2.7/-), 44.5 (3.0/-), 63.3 (3.1/-)
(Chen et al., 2022)	36	PS	L0, L1, L2	Low fidelity simulator	N, 2B	Repeating	Visual-verbal; Auditory-verbal; Auditory-spatial	0.2	50	26 (2.8/-)
(Zhang et al., 2022)	18	$\theta$ -Fp1, $\theta$ -Fp2, SC-EDA, IBI, HR, SDNN, RMSSD, LF, HF, LF/HF	L0	Low fidelity simulator	N, 1B, 2B	Matching	Auditory-manual	—	100	23.2 (1.6/-)
(Meteier et al., 2022)	80	HR	L3	Low fidelity simulator	N, 1B, 3B	Matching	Auditory-manual; Visual-manual	0.5	32.5	23.9 (8.2/19-66)
(Yang et al., 2021)	75	HR, SDNN, TP	L0	Medium fidelity simulator	N, 2B	Matching	Auditory-manual	1	67	31 (11.6/-)
(Gable et al., 2015)	8	HR, PS	L0	Low fidelity simulator	N, 0B, 1B, 2B	Repeating	Auditory-verbal	2.25	62.5	21.1 (19-23/-)
(Du et al., 2020)	102	HR	L3	High fidelity simulator	1B, 2B	Matching	Visual-manual	2.5	—	22.9 (3.8/18-38)
(Rahman et al., 2020)	33	SDNN, RMSSD, pNN50, VLF, LF,	L0	High fidelity simulator	N, 1B, 2B	Matching	Auditory-manual	—	100	42.5 (-/35-50)

		HF, TP, LF/HF, LFun, HFun								
(Cegovnik et al., 2018)	22	PS	L0	Low fidelity simulator	N, 0B, 1B, 2B, 3B	Repeating	Auditory-verbal	5.5	81.8	32.9 (-/22-61)
(Tjolleng et al., 2017)	15	IBI, LF	L0	Medium fidelity simulator	0B, 1B, 2B	Repeating	Auditory-verbal	—	100	27.7 (3.0/-)
(Niezgoda et al., 2015)	46	PS	L0	High fidelity simulator	N, 0B, 1B, 2B	Repeating	Auditory-verbal	2.25	63	36.6 (9.7/21-59)
(Reimer et al., 2009)	26	HR, EDA	L0	On-road	N, 0B, 1B, 2B	Repeating	Auditory-verbal	—	—	23.9 (1.6/-)
(Hajek et al., 2013)	47	HR, RR, SC-EDA, RMSSD	L2	Low fidelity simulator	N, 2B	Repeating	Auditory-verbal	2.25	72.3	28.5 (8.7/19-55)
(Zheng et al., 2021)	20	SDNN, RMSSD, pNN50, LF, HF, TP, LF/HF, VLF, LFun, HFun, SC-EDA	L0	Low fidelity simulator	N, 0B, 1B, 2B	Repeating	Auditory-verbal	2.5	75	26.7 (3.8/-)
(Mehler et al., 2009)	121	HR, SC-EDA, RR	L0	Medium fidelity simulator	N, 0B, 1B, 2B	Repeating	Auditory-verbal	2.25	48.8	24.5 (2.8/20-29)
(He et al., 2019)	34	PS, HR, RR, SC-EDA, $\theta$ -Fp1, $\theta$ -Fp2	L0	Low fidelity simulator	N, 1B, 2B	Counting	Auditory-verbal	2.5	48.6	26.4 (4.3/20-35)

- 1 Note: “—” means that the information has not been mentioned in the corresponding paper. In the table, N, 0B, 1B, 2B, and 3B standard for baseline without n-back task, 1-back task, 2-back task, and 3-back task, respectively. L0, L1, L2, and L3 donate SAE Level 0, Level 1, Level 2, and Level 3 automation, respectively.
- 2

---

### 1 3.1 Descriptive statistics

2 Table 2 provides descriptive information of all studies included in the meta-  
3 analyses. Overall, a total of 881 participants were involved in the experiments and  
4 experienced various levels of n-back tasks. Table 3 provides a summary of correlation  
5 coefficients between the metrics and cognitive load levels in all studies included for  
6 meta-analyses. It was found that cognitive levels induced by 1-back and 2-back tasks  
7 were most intensively investigated. At the same time, among all psychological and eye-  
8 tracking measures, the relationship between heart measures and cognitive load levels  
9 attracted the most attention in previous research, with HR attracting the most attention  
10 among heart-related metrics. Moreover, the total sample sizes of studies for a single  
11 meta-analysis varied widely, ranging from 38 to 711 participants. Finally, it should be  
12 noted that substantial differences in correlation coefficients of the metrics have been  
13 observed between different studies, confirming the need for further meta-analyses.

14 At the same time, given the small number of studies that could be identified, and  
15 considering the task difficulties, the 0-back (which only requires participants to simply  
16 repeat what they hear immediately) and baseline without n-back were aggregated as  
17 low task load (L); the 1-back was labeled as medium task load (M); and the 2-back was  
18 categorized as high task load (H).

1 **Table 3.** Correlation coefficients between the physiological and eye-tracking metrics and cognitive load levels in all studies.

	Metrics	<i>n</i>	Sample size			Mean	SD	<i>r</i>		
			Min	Max	Total			Min	Max	
L vs. M	Eye	PS	6	8	70	212	0.39	0.31	-0.12	0.77
		EBR	3	46	70	162	0.11	0.03	0.08	0.15
	Brain	$\theta$ -Fp1	2	18	34	52	0.42	0.04	0.40	0.45
		$\theta$ -Fp2	2	18	34	52	0.39	0.05	0.36	0.43
	Heart	HR	18	8	121	711	0.22	0.19	-0.2	0.76
		SDNN	3	18	33	51	0.35	0.65	-0.34	0.94
		RMSSD	3	18	33	51	0.36	0.65	-0.34	0.94
		LF	5	18	33	106	0.20	0.44	-0.13	0.97
		HF	3	18	33	71	0.34	0.45	-0.22	0.88
		LF/HF	3	18	33	71	0.49	0.46	0.05	0.97
		TP	2	20	33	53	0.35	0.87	-0.26	0.97
		pNN50	2	20	33	53	0.42	0.79	-0.14	0.98
		VLF	2	20	33	53	0.47	0.51	-0.05	0.98
		LFun	2	20	33	53	0.50	0.70	0.00	0.99
HFun	2	20	33	53	0.35	0.90	-0.29	0.99		
Skin	EDA	15	18	121	571	0.15	0.26	-0.52	0.8	
Respiration	RR	3	34	121	225	0.13	0.34	-0.34	0.44	
L vs. H	Eye	FD	2	36	46	82	0.02	0.22	-0.14	0.18
		PS	7	8	46	248	0.37	0.34	-0.2	0.82
	Eye	EBR	4	36	70	198	0.12	0.12	0.12	0.22
		EBD	2	36	70	106	-0.09	0.11	-0.11	-0.08
		$\theta$ -Fp1	2	18	34	42	0.33	0.02	0.32	0.35
	Brain	$\theta$ -Fp2	2	18	34	42	0.28	0.33	0.04	0.51
		HR	17	8	121	653	0.47	0.30	0.00	0.69
	Heart	SDNN	3	18	75	113	-0.23	0.58	-0.70	0.43
		RMSSD	3	18	47	85	0.30	0.61	-0.37	0.84
		LF	4	18	33	73	0.03	0.19	-0.22	0.30
		HF	4	15	20	73	0.04	0.18	-0.18	0.30
		LF/HF	2	18	20	38	0.27	0.15	0.16	0.38
		TP	2	20	75	95	-0.34	0.19	-0.21	-0.48

M vs. H	Skin	EDA	14	18	121	628	0.11	0.33	-0.99	0.45
	Respiration	RR	4	34	121	272	0.31	0.27	-0.13	0.59
	Eye	PS	4	8	46	158	0.23	0.10	0.13	0.34
		EBR	2	46	70	116	0.08	0.02	0.06	0.1
	Brain	$\theta$ -Fp1	2	18	34	52	0.02	0.02	0.01	0.04
		$\theta$ -Fp2	2	18	34	52	-0.25	0.02	-0.27	-0.24
		HR	12	8	121	522	0.17	0.16	0.02	0.56
		SDNN	3	18	33	71	0.30	0.52	-0.09	0.84
		RMSSD	4	18	70	141	0.21	0.40	-0.19	0.84
		LF	4	15	33	86	0.34	0.38	0.04	0.98
		HF	3	18	33	71	0.33	0.36	0.02	0.84
		LF/HF	3	18	33	71	0.46	0.44	0.12	0.96
		TP	2	20	33	53	0.50	0.63	0.06	0.95
		pNN50	2	20	33	53	0.40	0.74	-0.12	0.92
		VLF	2	20	33	53	0.51	0.45	0.06	0.96
		LFun	2	20	33	53	0.52	0.65	0.05	0.98
		HFun	2	20	33	53	0.41	0.82	-0.18	0.99
		Skin	SC-EDA	11	18	121	437	0.23	0.27	0.00
	Respiration	RR	3	34	121	225	0.08	0.08	0.02	0.19

Note: "vs." denotes the act of conducting pairwise comparisons,  $n$  denotes the number of studies included.

1  
2  
3

---

1 **3.2 Results for meta-analyses**

2 The results of the meta-analyses are presented in Table 4. A majority of the  
3 analyses in Table 4 have a significant unexplained variance ( $I^2 > 50\%$ ), indicating the  
4 necessity for moderators' analyses. The forest plots of the meta-analysis in Table 4 are  
5 provided in Appendix 1.

Post Print



1 **Table 4.** Meta-analyses' results of the association between the physiological metrics and cognitive **load levels**.

Metrics	<i>k</i>	Random Effects Model				Heterogeneity Test				
		<i>n</i>	Pooled <i>Z<sub>r</sub></i> (95 % CI)	Pooled <i>r</i> (95 % CI)	<i>p</i>	$\tau^2$	<i>Chi</i> <sup>2</sup> ( <i>df</i> )	<i>I</i> <sup>2</sup> (%)		
Eye	PS	212	6	0.48 (0.13,0.82)	0.45 (0.13, 0.68)	.006	0.15	38.82 (5)	87	
	EBR	162	3	0.11 (-0,0.22)	0.11 (-0,0.22)	.06	0	0.23 (2)	0	
Brain	$\theta$ -Fp1	52	2	0.43 (0.21,0.66)	0.41 (0.21, 0.58)	<.001	0.00	0.06 (1)	0	
	$\theta$ -Fp2	52	2	0.44 (0.22,0.67)	0.41 (0.22, 0.58)	<.001	0.00	0.1 (1)	0	
L vs. M	HR	711	18	0.22 (0.13,0.32)	0.22 (0.13, 0.32)	<.001	0.02	38.92 (17)	56	
	SDNN	51	3	0.88 (-0.06,1.82)	0.71 (-0.06, 0.95)	.07	0.63	26.61 (2)	92	
	RMSSD	121	3	0.61 (-0.34,1.56)	0.61 (-0.33, 0.92)	.2	0.9	91.33 (3)	97	
	LF	106	5	0.43 (-0.48,1.33)	0.41 (-0.45, 0.87)	.4	1.02	109.35 (4)	96	
	HF	71	3	0.52 (-0.47,1.51)	0.48 (-0.44, 0.91)	.3	0.71	29.62 (2)	93	
	Heart	LF/HF	71	3	0.89 (-0.43,2.21)	0.71 (-0.41, 0.98)	.2	1.31	53.86 (2)	96
		TP	53	2	0.92 (-1.40,3.23)	0.73 (-0.89, 1.00)	.4	2.74	60.76 (1)	98
		pNN50	53	2	1.09 (-1.31,3.48)	0.80 (-0.86, 1.00)	.4	2.93	64.13 (1)	98
		VLF	53	2	1.13 (-1.17,3.43)	0.81(-0.82,1)	.3	2.71	59.51(1)	98
		LFun	53	2	1.33 (-1.27,3.93)	0.87 (-0.85, 1.00)	.3	3.46	75.66 (1)	99
	HFun	53	2	1.18 (-1.7,4.07)	0.83 (-0.94, 1.00)	.4	4.27	92.11 (1)	99	
Skin	SC-EDA	551	14	0.20 (0.08,0.31)	0.20 (0.08,0.30)	<.001	0.03	37.92(13)	66	
Respiration	RR	225	3	0.15 (-0.28,0.58)	0.15 (-0.27,0.52)	.5	0.13	30.29 (2)	93	
	FD	82	2	0.03 (-0.28,0.34)	0.03 (-0.27, 0.33)	.9	0.02	1.93 (1)	48	
Eye	PS	248	7	0.47 (0.18,0.76)	0.44 (0.18, 0.64)	.002	0.12	49.12 (6)	88	
	EBR	198	4	0.16 (0.06,0.27)	0.16 (0.06,0.26)	.002	0	2.44 (3)	0	
	EBD	106	2	-0.09 (-0.23,0.06)	-0.09 (-0.23,0.06)	.3	0	0.02 (1)	0	
Brain	$\theta$ -Fp1	52	2	-0.05 (-0.72,0.62)	-0.05 (-0.62, 0.55)	.9	0.2	5.73 (1)	83	
	$\theta$ -Fp2	52	2	-0.35 (-1.06,0.36)	-0.34 (-0.79, 0.35)	.3	0.22	6.12 (1)	84	
L vs. H	HR	653	17	0.4 (0.25, 0.54)	0.38 (0.24, 0.49)	<.001	0.07	83.37 (16)	81	
	SDNN	113	3	-0.3 (-1.06,0.46)	-0.29 (-0.79, 0.43)	.4	0.41	22.09 (2)	91	
	RMSSD	155	4	0.24 (-0.79,1.27)	0.24 (-0.66, 0.85)	.7	1.06	248.43 (3)	99	
	LF	73	4	0.02 (-0.17,0.21)	0.02 (-0.17,0.21)	.8	0	2.54 (3)	0	
	HF	73	4	0.03 (-0.17,0.22)	0.03 (-0.17,0.22)	.8	0	2.2 (3)	0	

		LF/HF	38	2	0.28 (-0.07,0.63)	0.27 (-0.07, 0.56)	.1	0	0.42 (1)	0
		TP	95	2	-0.47 (-0.68, -0.26)	-0.44 (-0.59, -0.25)	<b>&lt;.001</b>	0	0.86 (1)	0
	Skin	SC-EDA	608	13	<b>0.20 (0.12,0.28)</b>	<b>0.20 (0.12,0.28)</b>	<b>&lt;.001</b>	<b>0</b>	<b>15.27 (12)</b>	<b>21</b>
	Respiration	RR	272	4	<b>0.34 (0.04,0.63)</b>	<b>0.34 (0.04, 0.56)</b>	<b>.03</b>	<b>0.08</b>	<b>23.3 (3)</b>	<b>87</b>
	Eye	PS	158	4	<b>0.29 (0.17,0.41)</b>	<b>0.28 (0.17, 0.39)</b>	<b>&lt;.001</b>	<b>0</b>	<b>2.15 (3)</b>	<b>0</b>
		EBR	116	2	<b>0.08 (-0.05,0.22)</b>	<b>0.08 (-0.05, 0.22)</b>	<b>.2</b>	<b>0</b>	<b>0.08 (1)</b>	<b>0</b>
		$\theta$ -Fp1	52	2	0.03 (-0.16,0.22)	0.03 (-0.16, 0.22)	.8	0	0.02 (1)	0
		$\theta$ -Fp2	52	2	-0.25 (-0.45, -0.06)	-0.24 (-0.42, -0.06)	<b>.009</b>	0	0.04 (1)	0
		HR	522	12	<b>0.14 (0.05,0.23)</b>	<b>0.14 (0.05, 0.23)</b>	<b>.002</b>	<b>0.01</b>	<b>18.74 (11)</b>	<b>41</b>
		SDNN	71	3	0.47 (-0.42,1.36)	0.44 (-0.40, 0.88)	.3	0.59	43.72 (2)	95
		RMSSD	141	4	<b>0.29 (-0.21,0.78)</b>	<b>0.29 (-0.21, 0.65)</b>	<b>.3</b>	<b>0.24</b>	<b>52.57 (3)</b>	<b>94</b>
	Brain	LF	86	4	<b>0.67 (-0.35,1.69)</b>	<b>0.59 (-0.34, 0.93)</b>	<b>.2</b>	<b>1.06</b>	<b>110.6 (3)</b>	<b>97</b>
		HF	71	3	<b>0.45 (-0.26,1.17)</b>	<b>0.42 (-0.25, 0.82)</b>	<b>.2</b>	<b>0.37</b>	<b>28.15 (2)</b>	<b>93</b>
		LF/HF	71	3	0.79 (-0.3,1.88)	0.66 (-0.29, 0.95)	.2	0.91	65.5 (2)	97
		TP	53	2	0.94 (-0.79,2.68)	0.74 (-0.66, 0.99)	.3	1.54	53.36 (1)	98
		pNN50	53	2	0.73 (-0.94,2.41)	0.62 (-0.74, 0.98)	.4	1.43	49.7 (1)	98
		VLF	53	2	<b>1 (-0.85,2.85)</b>	<b>0.76 (-0.69, 0.99)</b>	<b>.3</b>	<b>1.76</b>	<b>61.34 (1)</b>	<b>98</b>
		LFun	53	2	1.35 (-1.2,3.9)	0.87 (-0.83, 1.00)	.3	3.35	117.52 (1)	99
		HFun	53	2	1.23 (-1.54,4.01)	0.84 (-0.91, 1.00)	.4	3.98	137.67 (1)	99
	Skin	SC-EDA	417	10	<b>0.18 (0.04,0.33)</b>	<b>0.18 (0.04,0.33)</b>	<b>.01</b>	<b>0.04</b>	<b>35.95 (9)</b>	<b>75</b>
	Respiration	RR	225	3	<b>0.05 (-0.04,0.14)</b>	<b>0.05 (-0.04, 0.14)</b>	<b>.3</b>	<b>0</b>	<b>1.53 (2)</b>	<b>0</b>

- 1 Notes: L denotes the baseline condition and 0-back task, 1B denotes the 1-back task, and 2B represents the 2-back task, whereas "vs." denotes the act of conducting pairwise
- 2 comparisons,  $k$  denotes the cumulative sample size,  $n$  denotes the number of studies included. The bolded pooled  $r$  indicates significant ( $p < .05$ ) metrics. The bolded  $I^2$
- 3 indicates the existence of heterogeneity of the metrics.

---

### 1 3.3 Moderator analysis

2 Given that the inclusion of a sufficient number of studies is required for conducting  
3 meta-regression, we specifically focused on meta-analyses with an  $I^2$  greater than 50%  
4 or  $p$ -value in heterogeneity tests smaller than .05 and a minimum of three included  
5 studies. The results of the meta-regression model are summarized in Table 5. In addition,  
6 inter-group homogeneity tests were conducted on significant moderating factors  
7 ( $p < .05$ ). Table 6 displays the weighted average effect size  $Z_r$  and  $r$  (as well as their  
8 95%CI) for each subgroup, as well as the  $Q_W$  value that captures the overall  
9 heterogeneity within all the sub-groups of one moderator. Additionally,  $Q_B$  values are  
10 listed for each moderating variable, indicating the presence of heterogeneity among  
11 subgroups for each moderating factor (Lipsey & Wilson, 2001). The forest plots  
12 regarding the aggregated  $Z_r$  for each subgroup analysis are provided in Appendix 2.

1 **Table 5.** Results of meta-regression models.

	Moderators	L vs. M			L vs. H			M vs. H		
		<i>n</i>	<i>Coefficient</i>	<i>p</i>	<i>n</i>	<i>Coefficient</i>	<i>p</i>	<i>n</i>	<i>Coefficient</i>	<i>p</i>
PS	Simulator fidelity	6	0.20	.3	7	0.19	.4	-	-	-
	Modality of stimuli and responses	6	0.04	.95	7	-0.03	.9	-	-	-
	n-back version	6	-0.46	.02	7	-0.14	.6	-	-	-
	Inter-stimulus interval	6	0.13	.8	7	-0.03	.9	-	-	-
	Percentage of males	6	0.001	.09	7	0.01	.3	-	-	-
	Mean age	6	0.01	.7	7	0.04	.07	-	-	-
HR	Automation level	18	0.09	.09	17	0.33	<.0001	-	-	-
	Experimental environment	18	0.06	.6	17	0.32	.06	-	-	-
	Simulator fidelity	10	0.17	.5	9	-0.50	.14	-	-	-
	n-back version	18	0.12	.2	17	0.20	.2	-	-	-
	Modality of stimuli and responses	18	0.08	.09	17	0.18	.04	-	-	-
	Inter-stimulus interval	8	0.10	.08	7	0.16	.4	-	-	-
	Percentage of males	16	0.002	.4	15	0.01	.02	-	-	-
EDA	Mean age	18	-0.003	.4	17	-0.009	.2	-	-	-
	Automation level	14	0.22	<.0001	-	-	-	10	0.21	.005
	Simulator fidelity	6	0.15	.7	-	-	-	6	0.19	.6
	Experimental environment	14	0.25	.07	-	-	-	10	0.17	.4
	n-back version	14	0.06	.6	-	-	-	10	0.15	.3
	Modality of stimuli and responses	14	0.22	<.0001	-	-	-	10	0.22	.001
	Inter-stimulus interval	5	-0.03	.7	-	-	-	5	-0.09	.3
	Percentage of males	12	0.01	.07	-	-	-	9	0.01	.2
RR	Mean age	14	-0.002	.5	-	-	-	10	-0.002	.7
	Automation level	-	-	-	4	-0.02	.97	-	-	-
	Simulator fidelity	-	-	-	4	0.03	.97	-	-	-
	n-back version	-	-	-	4	-0.41	.052	-	-	-
	Modality of stimuli and responses	-	-	-	4	0.03	.97	-	-	-
	Inter-stimulus interval	-	-	-	4	-0.1	.9	-	-	-
	Percentage of males	-	-	-	4	0.01	.6	-	-	-
Mean age	-	-	-	4	0.01	.9	-	-	-	

1 Notes: In this table and the following tables,  $n$  denotes the number of studies included; “-” means that there is no need for subgroup analysis, as the meta-analysis results did  
 2 not demonstrate the presence of heterogeneity (see Table 4). The significant metrics ( $p<.05$ ) are bold

3  
 4 **Table 6.** Results of sub-group moderator analyses.

Pairwise Comparison of Cognitive Load	Physiological Metrics	Moderator	Moderator Level	$n$	Pooled $Z_r$ within Subgroup (95%CI)	Pooled $r$ within Subgroup (95%CI)	$Q_B$	$Q_W$
L vs. M	PS	n-back Version	Repeating	4	0.66 (0.28, 1.03)	<b>0.58 (0.27, 0.77)</b>	17.36 ( $p<.001$ )	10.83
			Matching	1	0.53 (0.36,0.70)	<b>0.49 (0.35, 0.60)</b>		
			Counting	1	-0.12 (-0.40,0.16)	-0.12 (-0.38, 0.16)		
	EDA	Automation Level	L0	12	0.13 (0.06,0.19)	<b>0.13 (0.06, 0.19)</b>	33.82 ( $p<.001$ )	4.10
			L2	1	0.29 (0.05,0.53)	<b>0.29 (0.05, 0.53)</b>		
			L3	1	1.10 (0.77,1.43)	<b>0.89 (0.65, 0.89)</b>		
EDA	n-back Modality	Auditory-verbal	11	0.12 (0.05,0.19)	<b>0.12 (0.05, 0.19)</b>	33.43 ( $p<.001$ )	4.13	
		Auditory-manual	2	0.21 (0.07,0.36)	<b>0.21 (0.07, 0.36)</b>			
		Visual-verbal	1	1.1 (0.77,1.43)	<b>0.89 (0.65, 0.89)</b>			
L vs. H	HR	n-back Modality	Auditory-verbal	11	0.42 (0.25,0.58)	<b>0.41 (0.25, 0.56)</b>	27.42 ( $p<.001$ )	43.31
			Auditory-manual	4	0.21 (0.04, 0.38)	<b>0.21 (0.04, 0.36)</b>		
			Visual-verbal	1	1.18 (0.86,1.50)	<b>0.84 (0.71, 0.91)</b>		
	HR	Automation Level	L0	13	0.25 (0.19,0.32)	<b>0.25 (0.19, 0.31)</b>	57.97 ( $p<.001$ )	12.03
			L2	2	0.97 (0.73, 1.22)	<b>0.75 (0.66, 0.82)</b>		
			L3	1	1.18 (0.86,1.50)	<b>0.84 (0.71, 0.91)</b>		
HR	Percentage of Male	<=50	9	0.31 (0.20, 0.42)	<b>0.30 (0.20, 0.40)</b>	9.85 ( $p=.007$ )	17.73	
		(50-90)	3	0.49 (-0.36, 1.34)	<b>0.46 (-0.32, 0.85)</b>			
		[90-100]	2	1.01 (0.58,1.43)	<b>0.77 (0.67, 0.86)</b>			
EDA		L0	8	0.08 (0.01,0.15)	<b>0.08 (0.01, 0.15)</b>		6.18	

		Automation Level	L2	1	0 (-0.13,0.13)	0.00 (-0.13, 0.13)	38.08	
			L3	1	1.09 (0.76,1.42)	<b>0.8 (0.64,0.89)</b>	( $p<.001$ )	
M vs. H	EDA	n-back Modality	Auditory-verbal	7	0.06 (-0.02,0.14)	0.06 (-0.02, 0.14)	7.60	4.15
			Auditory-manual	2	0.21 (0.06,0.35)	<b>0.21 (0.06, 0.35)</b>	( $p<.001$ )	<b>0.18</b>
			visual-verbal	1	1.09 (0.76,1.42)	<b>0.8 (0.64,0.89)</b>		

1 Notes: The bolded pooled  $r$  indicates significant ( $p<.05$ ) associations; The bolded  $Q_H$  indicates the existence of heterogeneity of the metrics.

Post Print

---

## 1 4. DISCUSSION

2 In this study, we systematically reviewed previous research regarding the  
3 associations between physiological and eye-tracking metrics and cognitive load levels  
4 in vehicles. The n-back tasks, which were commonly adopted to impose cognitive load  
5 in previous driving research, have been used as benchmarks of cognitive load levels  
6 (Janczewski et al., 2021). The random effects meta-analyses were conducted followed  
7 by moderator analyses.

### 8 4.1 Associations between the metrics and varying levels of cognitive load

9 Through meta-analyses, we found that although some metrics were found to be  
10 sensitive to the cognitive load in certain previous studies, they failed to pass the  
11 significance test in our meta-analyses. For example, the LF/HF ratio was found to be  
12 positively associated with the increase of cognitive load in Rahman et al (2020) and  
13 Zheng et al (2021), but it did not achieve statistical significance in our analyses.

14 At the same time, although it has been widely acknowledged that not all  
15 physiological features are sensitive to all levels of cognitive load (e.g., (Ayres et al.,  
16 2021; Li et al., 2022)) our meta-analyses provide further evidence to support this  
17 statement. For example, some metrics were more sensitive to lower levels of cognitive  
18 load compared to higher levels of cognitive load. We call these metrics **Sensitive-to-**  
19 **Low Metrics**. For example, the power of theta waves at Fp1 was sensitive to the low  
20 to median cognitive load levels with a median association ( $r=0.41$ ), but it was not  
21 sensitive to higher levels of cognitive load (i.e., between medium and high cognitive

---

1 load levels). This finding is partially in line with the findings in the meta-analysis by  
2 Chikhi et al (2022), who also observed a positive association between the power of  $\theta$   
3 wave at the frontal area and high cognitive load. But our research has provided higher  
4 resolution as we quantified more levels of cognitive load. Specifically, some metrics  
5 (e.g.,  $\theta$ -Fp1) may increase rapidly with a small increment of the cognitive load and  
6 reach a plateau at a medium level of cognitive load. This might be because some  
7 physiological indicators may cease to rise beyond a certain threshold, similar to some  
8 observations in brain studies (Bosking et al., 2017). The potential “ceiling effect” of  $\theta$ -  
9 Fp1 may explain our observations, which has also been mentioned by Chikhi et al (2022)  
10 to explain the weaker association between the power of  $\theta$  in multitasking versus  
11 single-tasking situation. It is also possible that beyond this “ceiling” point, other neural  
12 mechanisms or states can dominate (e.g., Weiss et al., 1995; Sauseng et al., 2004), and  
13 this may explain the significant difference between low to medium but insignificant  
14 difference between low to high cognitive load in terms of  $\theta$ -Fp1. Though future  
15 research is needed to explain this phenomenon, the findings have some practical  
16 implications. Specifically, if we aim to detect medium (e.g., heavy traffic) to high  
17 cognitive load (e.g., heavy traffic and non-driving-related tasks) in drivers, the weights  
18 of sensitive-to-low metrics should be downgraded.

19 At the same time, some other metrics, for example, pupil size (PS), heart rate (HR),  
20 and skin conductance-electrodermal activity (SC-EDA) demonstrated a consistent  
21 growth relationship with the increase of cognitive load levels. We call these metrics



---

1 **High-Resolution Metrics.** Specifically, from low to medium, and from medium to high  
2 cognitive load, significant associations were observed for HR, SC-EDA, and PS.  
3 However, we should also be aware that, there are still differences in resolution among  
4 high-resolution metrics. Specifically, HR showed a substantially higher association  
5 strength with cognitive load variations than SC-EDA, and PS from eye-tracking  
6 measures exhibited an even higher association compared to HR and SC-EDA. The  
7 finding regarding HR is consistent with previous meta-analyses in non-driving domains  
8 (Hughes et al., 2019), which also found an association between high cognitive load and  
9 HR. In addition, the superior performance of PS is consistent with previous research  
10 (He et al., 2022) which found that, when predicting drivers' cognitive load states using  
11 five typical machine learning models, feature sets including eye-related measures could  
12 consistently result in high accuracies compared to feature sets with physiological  
13 measures alone. This highlights the superior performance of eye-related measures in  
14 monitoring drivers' cognitive load states (Chen et al., 2022).

15 In contrast to high-resolution metrics, some metrics were only sensitive from low  
16 to high cognitive load, but cannot differentiate low to medium and medium to high  
17 cognitive load levels. We call these **Low-Resolution Metrics.** For example, eye blink  
18 rate (EBR), respiration rate (RR), the total power (TP) demonstrated associations  
19 between low and high cognitive load levels only. This suggests that significant changes  
20 in cognitive load are required to induce notable variations in EBR, RR, and TP.

---

1 Finally, we also observed a non-linear relationship between the power spectrum  
2 of  $\theta$  waves at Fp2 ( $\theta$ -Fp2) and cognitive load levels. Specifically, we observed a  
3 positive correlation ( $r = 0.41$ ) from low to medium cognitive load levels, but a negative  
4 correlation ( $r = -0.24$ ) from medium to high cognitive load levels. Similar to the  
5 findings in  $\theta$ -Fp1, additional states of the participants might have dominated the EEG  
6 at Fp2 when the task becomes “too difficult”, which may have led to a decrease in  $\theta$ -  
7 Fp2 at some point. In driver state monitoring systems, acknowledging this non-linearity  
8 is vital for accurate assessments.

9 It should be noted that the categorization of the metrics in this study is range-  
10 specific. Specifically, we only considered the cognitive levels from no secondary task  
11 to 2-back task. Even with no secondary task condition, drivers were still responsible for  
12 driving tasks. With lower or higher extreme cognitive load, high- or low-resolution  
13 metrics may be downgraded to sensitive-to-low ones; and with higher resolution of the  
14 cognitive task, the high-resolution metrics may become low-resolution ones. It should  
15 also be noted that previous research indicated that the driving performance measures  
16 may only be sensitive to high cognitive load (Yang et al., 2023). As for physiological  
17 and eye-tracking measures, we did not observe any that can differentiate medium to  
18 high cognitive load levels only. Thus, it seems that physiological and eye-tracking  
19 measures and driving performance measures may complement each other in driver-  
20 monitoring tasks.

21

---

## 1 4.2 Moderators

2 First of all, as expected, the n-back task version can moderate the association  
3 between pupil size and cognitive load from low to medium levels. Specifically, the  
4 repeating version of the n-back task led to significant associations in the repeating and  
5 matching versions, but not in the counting version. It is likely that the increased  
6 cognitive demand to remember the running total number of cases in the counting  
7 version led to an already high cognitive load even in the 1-back task and thus shadowed  
8 the effect of the additional cognitive load in the 2-back task. Though future research is  
9 needed to validate this hypothesis, the finding reveals the influence of task  
10 characteristics in modulating the cognitive-related eye-tracking metrics.

11 At the same time, we notice that the associations between HR and SC-EDA and  
12 cognitive load levels were moderated by the automation levels. Specifically, the  
13 associations between the HR and SC-EDA and the cognitive load were the strongest in  
14 vehicles with SAE Level 3 automation, both from low to medium and from low to high  
15 levels of cognitive load. It is possible that the drivers in SAE Level 3 vehicles are freed  
16 from continuous vehicle controlling tasks and they may experience the lowest workload  
17 in driving. Thus, the cognitive load imposed by the n-back task is less likely to be  
18 shadowed by the variations of task load in driving tasks. This finding suggests that  
19 different features and different algorithms may need to be designed for driver cognitive  
20 load detection in vehicles with driving automation, which is still lacking. To the best of

---

1 our knowledge, only one study has focused on the driver cognitive load estimation  
2 algorithms in vehicles with driving automation (Meteier et al., 2021).

3 Moreover, in addition to automation levels, the associations were also moderated  
4 by the n-back task modality. Specifically, our analysis reveals that the correlation  
5 between HR and SC-EDA and cognitive load is most pronounced in the visual-verbal  
6 n-back tasks. It is likely that the high demand of visual resources in driving competes  
7 with the visual component in the visual-manual n-back tasks and thus leads to high  
8 sensitivity of the HR and SC-EDA to the visual-verbal n-back task. However, it should  
9 be noted that the visual component is not a cognitive component, and thus, the high  
10 associations of HR and SC-EDA might be the result of increased stress during the task  
11 (De Loeff et al., 2018; Liu & Du, 2018) This finding indicates that the experiment  
12 settings may affect the physiological and eye-tracking measures (Nilsson et al., 2022).

13 Additionally, we observed that the characteristics of the participants may also  
14 affect the association between cognitive load and physiological responses, with male  
15 participants leading to a stronger association between HR and low to high levels of  
16 cognitive load. This highlights the importance of considering participant characteristics  
17 when designing driver monitoring systems.

18 Finally, it should be noted that heterogeneity has still been observed in most of the  
19 sub-groups, indicating that additional moderating factors may still exist. Future  
20 research is still needed to explore these factors and thus better guide the design of the  
21 driver monitoring systems.

---

### 1 4.3 Limitations and Future Directions

2 The current investigation presents several limitations. Firstly, the included  
3 investigations only considered the cognitive load imposed by the n-back tasks.  
4 Although n-back tasks have been widely adopted as a method for inducing cognitive  
5 load in traffic research. (Janczewski et al., 2021), other cognitive load induction tasks,  
6 such as mathematical tasks, may also be considered if their difficulty levels are  
7 quantified (Yang et al., 2023). Second, four studies were excluded from the present  
8 investigation due to the absence of required information for meta-analysis (Barua et al.,  
9 2017; Chihara et al., 2020; Solovey et al., 2014; Zhen et al., 2016). Consequently,  
10 though the meta-analysis based on a small sample size may still provide insights into  
11 potential trends and differences (Zheng, 2013), conclusions from our study should still  
12 be interpreted with caution, given that the associations between various metrics and  
13 cognitive load identified in our study are based on a limited number of studies. For  
14 similar reasons, the current investigation examined only a few potential moderators and  
15 some subgroups in the subgroup analyses contained relatively small numbers of studies.  
16 Finally, it should be noted that most of the research was conducted in simulators, and  
17 given the nature of the n-back task, some of the measures may be different from what  
18 they are in a natural driving condition (e.g., the response modality of the n-back task  
19 instead of cognitive load may have a strong influence on respiration rate and the  
20 complex lighting condition on a public road may shadow the influence of cognitive load

---

1 on pupil size). Future research should be re-conducted when a larger sample size  
2 becomes available.

### 3 5. CONCLUSIONS

4 Despite the extensive research on the use of physiological and eye-tracking  
5 measures to assess cognitive load in driving, researchers have not reached a consensus  
6 on their associations with cognitive load. Based on a systematic review and a meta-  
7 analysis, for the first time, we quantified the association between physiological and eye-  
8 tracking metrics and cognitive load in driving. We identified four types of metrics, i.e.,  
9 sensitive-to-low ones that can only differentiate the low (no secondary task or 0-back)  
10 to medium (1-back) level of cognitive load (including the power spectrum of  $\theta$  waves  
11 of electroencephalogram at Fp1 channel); low-resolution ones that can only  
12 differentiate low and high cognitive load (including the overall power spectrum of  
13 electrocardiogram, eye blink rate and respiration rate) and others that show non-linear  
14 patterns with the increase of cognitive load (i.e., the power spectrum of  $\theta$  waves at Fp2  
15 channel). Furthermore, it has been found that n-back task versions, the modality of n-  
16 back tasks, the level of automation, and the percentage of male participants could  
17 moderate the associations between metrics and cognitive load.

18 This study, through a meta-analysis, offers a new perspective in understanding the  
19 relationship between physiological and eye-tracking metrics and different cognitive  
20 load levels and provides new insights into resolving the debates in this area. The  
21 findings highlight the importance of considering individual heterogeneity, driving

---

1 automation, data collection environment, and metric characteristics when developing  
2 algorithms for driver cognitive load estimation. Future research should further validate  
3 our findings when more data and research become available.

#### 4 **ACKNOWLEDGEMENTS**

5 This work was supported by the National Natural Science Foundation of China (No.  
6 52202425), and in part by the Guangzhou Municipal Science and Technology Project  
7 (No. 2023A03J0011), and Guangzhou Science and Technology Program City-  
8 University Joint Funding Project (No. 2023A03J0001).

#### 9 **CRedit AUTHORSHIP CONTRIBUTION STATEMENT**

10 **Ange Wang:** Conceptualization, Data curation, Formal analysis, Methodology, Software,  
11 Validation, Writing – original draft. **Chunxi Huang:** Validation, Formal analysis. **Jiyao**  
12 **Wang:** Data curation, Validation. **Dengbo He:** Formal analysis, Funding acquisition,  
13 Methodology, Supervision, Validation, Writing – review & editing.

#### 14 **DECLARATION OF COMPETING INTEREST**

15 The authors declare that they have no known competing financial interests or personal  
16 relationships that could have appeared to influence the work reported in this paper.

#### 17 **REFERENCES**

- 18 Ayres, P. (2020). Something Old, Something New from Cognitive Load Theory.  
19 *Computers in Human Behavior*, 113, 106503.
- 20 Ayres, P., Lee, J. Y., Paas, F., & van Merriënboer, J. J. G. (2021). The Validity of  
21 Physiological Measures to Identify Differences in Intrinsic Cognitive Load. *Frontiers*  
22 *in Psychology*, 12, 702538.
- 23 Barua, S., Ahmed, M. U., & Begum, S. (2017, May). Classifying Drivers' Cognitive  
24 Load Using EEG Signals. *In pHealth* (pp. 99-106).

- 
- 1 Bosking, W. H., Sun, P., Ozker, M., Pei, X., Foster, B. L., Beauchamp, M. S., &  
2 Yoshor, D. (2017). Saturation in Phosphene Size with Increasing Current Levels  
3 Delivered to Human Visual Cortex. *Journal of Neuroscience*, 37(30), 7188–7197.
- 4 Broadbent, D. P., D’Innocenzo, G., Ellmers, T. J., Parsler, J., Szameitat, A. J., &  
5 Bishop, D. T. (2023). Cognitive load, Working Memory Capacity and Driving  
6 Performance: A Preliminary fNIRS and Eye Tracking Study. *Transportation  
7 Research Part F: Traffic Psychology and Behaviour*, 92, 121–132.
- 8 Cegovnik, T., Stojmenova, K., Jakus, G., & Sodnik, J. (2018). An Analysis of the  
9 Suitability of A Low-cost Eye Tracker for Assessing the Cognitive Load of Drivers.  
10 *Applied Ergonomics*, 68, 1–11.
- 11 Chen, W., Sawaragi, T., & Hiraoka, T. (2022). Comparing Eye-tracking Metrics of  
12 Mental Workload Caused by NDRTs in Semi-autonomous Driving. *Transportation  
13 Research Part F: Traffic Psychology and Behaviour*, 89, 109–128.
- 14 Chihara, T., Kobayashi, F., & Sakamoto, J. (2020). Evaluation of Mental Workload  
15 During Automobile Driving Using One-class Support Vector Machine with Eye  
16 Movement Data. *Applied Ergonomics*, 89.
- 17 Chikhi, S., Matton, N., & Blanchet, S. (2022). EEG Power Spectral Measures of  
18 Cognitive Workload: A Meta-Analysis. *Psychophysiology*, 59(6), e14009.
- 19 Cohen, J. E. (1988). The Counterintuitive in Conflict and Cooperation. *American  
20 Scientist*, 76.6 (1988): 577-584.
- 21 Committee, O.-R. A. D. (ORAD). (2014). *Taxonomy and Definitions for Terms  
22 Related to On-Road Motor Vehicle Automated Driving Systems*. SAE International.
- 23 De Looft, P. C., Cornet, L. J. M., Embregts, P. J. C. M., Nijman, H. L. I., & Didden,  
24 H. C. M. (2018). Associations of Sympathetic and Parasympathetic Activity in Job  
25 Stress and Burnout: A Systematic Review. *Plos One*, 13(10), e0205741.
- 26 Deng, M., Gluck, A., Zhao, Y., Li, D., Menassa, C. C., Kamat, V. R., & Brinkley, J.  
27 (2024). An Analysis of Physiological Responses as Indicators of Driver Takeover  
28 Readiness in Conditionally Automated Driving. *Accident Analysis & Prevention*, 195,  
29 107372.
- 30 Du, N., Kim, J., Zhou, F., Pulver, E., Tilbury, D. M., Robert, L. P., Pradhan, A. K., &  
31 Yang, X. J. (2020). Evaluating Effects of Cognitive Load, Takeover Request Lead  
32 Time, and Traffic Density on Drivers’ Takeover Performance in Conditionally  
33 Automated Driving. *12th International Conference on Automotive User Interfaces  
34 and Interactive Vehicular Applications*, 66–73.
- 35 Du, N., Yang, X. J., & Zhou, F. (2020). Psychophysiological Responses to Takeover  
36 Requests in Conditionally Automated Driving. *Accident Analysis and Prevention*,
- 37 Ferreira, E., Ferreira, D., Kim, S., Siirtola, P., Röning, J., Forlizzi, J. F., & Dey, A. K.  
38 (2014). Assessing Real-time Cognitive Load Based on Psycho-physiological



- 
- 1 Measures for Younger and Older Adults. *2014 IEEE Symposium on Computational*  
2 *Intelligence, Cognitive Algorithms, Mind, and Brain (CCMB)*, 39–48.
- 3 Gable, T., Walker, B., Kun, A., Winton, R., & ACM. (2015). Comparing Heart Rate  
4 and Pupil Size as Objective Measures of Workload in the Driving Context: Initial  
5 Look. *In Adjunct proceedings of the 7th international conference on automotive user*  
6 *interfaces and interactive vehicular applications*, 20-25, New York, NY, USA.
- 7 Hajek, W., Gaponova, I., Fleischer, K. H., & Krems, J. (2013). Workload-adaptive  
8 Cruise Control – A New Generation of Advanced Driver Assistance Systems.  
9 *Transportation Research Part F: Traffic Psychology and Behaviour*, 20, 108–120.
- 10 Harbluk, J. L., Noy, Y. I., Trbovich, P. L., & Eizenman, M. (2007). An On-road  
11 Assessment of Cognitive Distraction: Impacts on Drivers' Visual Behavior and  
12 Braking Performance. *Accident Analysis & Prevention*, 39(2), 372–379.
- 13 Hart, S. G., & Staveland, L. E. (1988). Development of NASA-TLX (Task Load  
14 Index): Results of Empirical and Theoretical Research. In P. A. Hancock & N.  
15 Meshkati (Eds.), *Advances in Psychology*, 52, 139–183. North-Holland.
- 16 He, D., Donmez, B., Liu, C. C., & Plataniotis, K. N. (2019). High Cognitive Load  
17 Assessment in Drivers Through Wireless Electroencephalography and the Validation  
18 of a Modified N-Back Task. *IEEE Transactions on Human-Machine Systems*, 49(4),  
19 362–371.
- 20 He, D., Wang, Z., Khalil, E. B., Donmez, B., Qiao, G., & Kumar, S. (2022).  
21 Classification of Driver Cognitive Load: Exploring the Benefits of Fusing Eye-  
22 Tracking and Physiological Measures. *Transportation Research Record: Journal of*  
23 *the Transportation Research Board*, 2676(10), 670–681.
- 24 Higgins, J. P. T., & Deeks, J. J. (2003). *Measuring Inconsistency in Meta-analyses* |  
25 *The BMJ*. <https://www.bmj.com/content/327/7414/557.short>.
- 26 Hughes, A. M., Hancock, G. M., Marlow, S. L., Stowers, K., & Salas, E. (2019).  
27 Cardiac measures of cognitive workload: A meta-analysis. *Human Factors: The*  
28 *Journal of the Human Factors and Ergonomics Society*, 61(3), 393-414.
- 29 Jackson, L., Chapman, P., & Crundall, D. (2009). What Happens Next? Predicting  
30 Other Road Users' Behaviour as A Function of Driving Experience and Processing  
31 Time. *Ergonomics*, 52(2), 154–164.
- 32 Li, P., Li, Y., Yao, Y., Wu, C., Nie, B., & Li, S. E. (2022). Sensitivity of  
33 Electrodermal Activity Features for Driver Arousal Measurement in Cognitive Load:  
34 The Application in Automated Driving Systems. *IEEE Transactions on Intelligent*  
35 *Transportation Systems*, 23(9), 14954–14967.
- 36 Lipsey, M., & Wilson, D. (2001). *Practical Meta-Analysis*. SAGE Publications.
- 37 Liu, Y., & Du, S. (2018). Psychological Stress Level Detection Based on  
38 Electrodermal Activity. *Behavioural Brain Research*, 341, 50–53.

- 
- 1 McCarthy, P. L., Holstein, S. A., Petrucci, M. T., Richardson, P. G., Hulin, C., Tosi,  
2 P., Bringham, S., Musto, P., Anderson, K. C., Caillot, D., Gay, F., Moreau, P., Marit,  
3 G., Jung, S.-H., Yu, Z., Winograd, B., Knight, R. D., Palumbo, A., & Attal, M.  
4 (2017). Lenalidomide Maintenance After Autologous Stem-Cell Transplantation in  
5 Newly Diagnosed Multiple Myeloma: A Meta-Analysis. *Journal of Clinical*  
6 *Oncology*, 35(29), 3279–3289.
- 7 Mehler, B., & Reimer, B. (2019). How Demanding is “Just Driving?” A Cognitive  
8 Workload—Psychophysiological Reference Evaluation. *Driving Assessment*  
9 *Conference*, 10(2019), Article 2019.
- 10 Mehler, B., Reimer, B., & Dusek, J. A. (2011). *MIT AgeLab Delayed Digit Recall*  
11 *Task (N-Back)*. Cambridge, MA: Massachusetts Institute of Technology, 17.
- 12 Mehler, B., Reimer, B., & Coughlin, J. F. (2012b). Sensitivity of Physiological  
13 Measures for Detecting Systematic Variations in Cognitive Demand From a Working  
14 Memory Task: An On-Road Study Across Three Age Groups. *Human Factors: The*  
15 *Journal of the Human Factors and Ergonomics Society*, 54(3), 396–412.
- 16 Mehler, B., Reimer, B., Coughlin, J. F., & Dusek, J. A. (2009). Impact of Incremental  
17 Increases in Cognitive Workload on Physiological Arousal and Performance in Young  
18 Adult Drivers. *Transportation Research Record: Journal of the Transportation*  
19 *Research Board*, 2138(1), 6–12.
- 20 Melnicuk, V., Thompson, S., Jennings, P., & Birrell, S. (2021). Effect of cognitive  
21 load on drivers’ state and task performance during automated driving: Introducing a  
22 novel method for determining stabilisation time following take-over of  
23 control. *Accident Analysis & Prevention*, 151, 105967.
- 24 Meteier, Q., Capallera, M., Ruffieux, S., Angelini, L., Abou Khaled, O., Mugellini,  
25 E., Widmer, M., & Sonderegger, A. (2021). Classification of Drivers’ Workload  
26 Using Physiological Signals in Conditional Automation. *Frontiers in Psychology*, 12,  
27 596038.
- 28 Moher, D. (2009). *Preferred Reporting Items for Systematic Reviews and Meta-*  
29 *Analyses: The PRISMA Statement | Annals of Internal Medicine*.  
30 <https://www.acpjournals.org/doi/full/10.7326/0003-4819-151-4-200908180-00135>
- 31 Muhrer, E., & Vollrath, M. (2011). The Effect of Visual and Cognitive Distraction on  
32 Driver’s Anticipation in A Simulated Car Following Scenario. *Transportation*  
33 *Research Part F: Traffic Psychology and Behaviour*, 14(6), 555–566.
- 34 Muth, E. R., Moss, J. D., Rosopa, P. J., Salley, J. N., & Walker, A. D. (2012).  
35 Respiratory Sinus Arrhythmia as A Measure of Cognitive Workload. *International*  
36 *Journal of Psychophysiology*, 83(1), 96–101.
- 37 Niezgod, M., Tarnowski, A., Kruszewski, M., & Kamiński, T. (2015). Towards  
38 Testing Auditory-vocal Interfaces and Detecting Distraction While Driving: A

- 
- 1 Comparison of Eye-movement Measures in the Assessment of Cognitive Workload.  
2 *Transportation Research Part F: Traffic Psychology and Behaviour*, 32, 23–34.
- 3 Nilsson, E. J., Bärgrman, J., Ljung Aust, M., Matthews, G., & Svanberg, B. (2022).  
4 Let Complexity Bring Clarity: A Multidimensional Assessment of Cognitive Load  
5 Using Physiological Measures. *Frontiers in Neuroergonomics*, 3, 787295.
- 6 Rahman, H., Ahmed, M. U., Barua, S., & Begum, S. (2020). Non-contact-based  
7 Driver’s Cognitive Load Classification Using Physiological and Vehicular  
8 Parameters. *Biomedical Signal Processing and Control*, 55, 101634.
- 9 Recarte, M. A., & Nunes, L. M. (2000). Effects of Verbal and Spatial-imagery Tasks  
10 on Eye Fixations While Driving. *Journal of Experimental Psychology: Applied*, 6(1),  
11 31–43.
- 12 Reimer, B., Mehler, B., Coughlin, J. F., Godfrey, K. M., & Tan, C. (2009). An On-  
13 road Assessment of the Impact of Cognitive Workload on Physiological Arousal in  
14 Young Adult Drivers. *Proceedings of the 1st International Conference on Automotive  
15 User Interfaces and Interactive Vehicular Applications*, 115–118, New York, NY,  
16 USA.
- 17 Rieck, J. R., DeSouza, B., Baracchini, G., & Grady, C. L. (2022). Reduced  
18 Modulation of BOLD Variability as A Function of Cognitive Load in Healthy Aging.  
19 *Neurobiology of Aging*, 112, 215–230.
- 20 Sagberg, F., & Bjørnskau, T. (2006). Hazard Perception and Driving Experience  
21 Among Novice Drivers. *Accident Analysis & Prevention*, 38(2), 407–414.
- 22 Sauseng, P., Klimesch, W., Doppelmayr, M., Hanslmayr, S., Schabus, M., & Gruber,  
23 W. R. (2004). Theta coupling in the human electroencephalogram during a working  
24 memory task. *Neuroscience Letters*, 354(2), 123-126.
- 25 Sârbescu, P., Stanojević, P., & Jovanović, D. (2014). A Cross-cultural Analysis of  
26 Aggressive Driving: Evidence from Serbia and Romania. *Transportation Research  
27 Part F: Traffic Psychology and Behaviour*, 24, 210–217.
- 28 Schmidt, L., Shokraneh, F., Steinhausen, K., & Adams, C. E. (2019). Introducing  
29 RAPTOR: RevMan Parsing Tool for Reviewers | Systematic Reviews | Full Text.  
30 *Systematic Reviews*, 8(1). [https://systematicreviewsjournal.biomedcentral.com/  
31 articles/10.1186/s13643-019-1070-0](https://systematicreviewsjournal.biomedcentral.com/articles/10.1186/s13643-019-1070-0)
- 32 Singh, S. (2015). Critical Reasons for Crashes Investigated in the National Motor  
33 Vehicle Crash Causation Survey. *Traffic Safety Facts - Crash Stats*, Article DOT HS  
34 812 115. [https://trid.trb.org/view.aspx?id=1346216&source=post\\_page](https://trid.trb.org/view.aspx?id=1346216&source=post_page)
- 35 Solhjoo, S., Haigney, M. C., McBee, E., Van Merriënboer, J. J. G., Schuwirth, L.,  
36 Artino, A. R., Battista, A., Ratcliffe, T. A., Lee, H. D., & Durning, S. J. (2019). Heart  
37 Rate and Heart Rate Variability Correlate with Clinical Reasoning Performance and  
38 Self-Reported Measures of Cognitive Load. *Scientific Reports*, 9(1), 14668.

- 
- 1 Solovey, E. T., Zec, M., Garcia Perez, E. A., Reimer, B., & Mehler, B. (2014).  
2 Classifying Driver Workload Using Physiological and Driving Performance Data:  
3 Two Field Studies. *Proceedings of the SIGCHI Conference on Human Factors in*  
4 *Computing Systems*, 4057–4066.
- 5 Spyridakos, P. D., Merat, N., Boer, E. R., & Markkula, G. M. (2020). Behavioural  
6 Validity of Driving Simulators for Prototype HMI Evaluation. *IET Intelligent*  
7 *Transport Systems*, 14(6), 601–610.
- 8 Stapel, J., Mullakkal-Babu, F. A., & Happee, R. (2019). Automated Driving Reduces  
9 Perceived Workload, But Monitoring Causes Higher Cognitive Load Than Manual  
10 Driving. *Transportation Research Part F: Traffic Psychology and Behaviour*, 60,  
11 590–605.
- 12 Tjolleng, A., Jung, K., Hong, W., Lee, W., Lee, B., You, H., Son, J., & Park, S.  
13 (2017). Classification of A Driver’s Cognitive Workload Levels Using Artificial  
14 Neural Network on ECG Signals. *Applied Ergonomics*, 59, 326–332.
- 15 Viechtbauer, W. (2010). Conducting Meta-analyses in R with The Metafor Package.  
16 *Journal of Statistical Software*, 36(3), 1–48.
- 17 von Janczewski, N., Wittmann, J., Engeln, A., Baumann, M., & Krauß, L. (2021). A  
18 Meta-analysis of The n-back Task While Driving and Its Effects on Cognitive  
19 Workload. *Transportation Research Part F: Traffic Psychology and Behaviour*, 76,  
20 269–285.
- 21 Weiss, T., Sust, M., Beyer, L., Hansen, E., Rost, R., & Schmalz, T. (1995). Theta  
22 Power Decreases in Preparation for Voluntary Isometric Contractions Performed with  
23 Maximal Subjective Effort. *Neuroscience Letters*, 193(3), 153-156.
- 24 Wickens, C. D. (2020). Processing Resources and Attention. In *Multiple Task*  
25 *Performance* (pp. 3-34). CRC Press.
- 26 Wickens, C. M., Mann, R. E., Stoduto, G., Ialomiteanu, A., & Smart, R. G. (2011).  
27 Age Group Differences in Self-reported Aggressive Driving Perpetration and  
28 Victimization. *Transportation Research Part F: Traffic Psychology and Behaviour*,  
29 14(5), 400–412.
- 30 Yang, H., Liu, H., Hu, Z., Nguyen, A.-T., Guerra, T.-M., & Lv, C. (2024).  
31 Quantitative Identification of Driver Distraction: A Weakly Supervised Contrastive  
32 Learning Approach. *IEEE Transactions on Intelligent Transportation Systems*, 25(2),  
33 2034–2045.
- 34 Yang, H., Wu, J., Hu, Z., & Lv, C. (2024). Real-Time Driver Cognitive Workload  
35 Recognition: Attention-Enabled Learning with Multimodal Information Fusion. *IEEE*  
36 *Transactions on Industrial Electronics*, 71(5), 4999–5009.
- 37 Yang, S., Kuo, J., Lenne, M., Fitzharris, M., Horberry, T., Blay, K., Wood, D.,  
38 Mulvihill, C., & Truche, C. (2021). The Impacts of Temporal Variation and

- 
- 1 Individual Differences in Driver Cognitive Workload on ECG-Based Detection.  
2 *Human Factors*, 63(5), 772–787.
- 3 Zhang, H., Qu, W., Ge, Y., Sun, X., & Zhang, K. (2017). Effect of Personality Traits,  
4 Age and Sex on Aggressive Driving: Psychometric Adaptation of the Driver  
5 Aggression Indicators Scale in China. *Accident Analysis & Prevention*, 103, 29–36.
- 6 Zhang, Q., Yang, K., Qu, X., & Tao, D. (2022). Evaluation of Drivers' Mental  
7 Workload Based on Multi-modal Physiological Signals. *Shenzhen Daxue Xuebao*  
8 *(Ligong Ban)/Journal of Shenzhen University Science and Engineering*, 39(3), 278–  
9 286.
- 10 Zheng, L., Qiao, X.-Q., Ni, T., Yang, W., & Li, Y.-N. (2021). Driver Cognitive Loads  
11 Based on Multi-dimensional Information Feature Analysis. *Zhongguo Gonglu*  
12 *Xuebao/China Journal of Highway and Transport*, 34(4), 240–250.
- 13 Zhenhai, G., Yang, L., Lifei, D., Hui, Z., & Kaishu, Z. (2016). Driver Workload  
14 Evaluation Using Physiological Indices in Dual-Task Driving Conditions.  
15 *International Conference on Applied System Innovation*, 809–814, Osaka, Japan.
- 16 Zheng, M. (2013). *Application and Case Analysis of Meta-Analysis Software*. ISBN:  
17 9787117171670, People's Health Publishing House, Beijing, China.