

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/109488>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

© 2018 Elsevier. Licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International <http://creativecommons.org/licenses/by-nc-nd/4.0/>.



Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

1 **Calibrating Trust through Knowledge: Introducing the Concept of** 2 **Informed Safety for Automation in Vehicles**

3 Siddartha Khastgir *, Stewart Birrell, Gunwant Dhadyalla, Paul Jennings

4 WMG, University of Warwick, UK

5 * Corresponding author: S.Khastgir.1@warwick.ac.uk

6 **Abstract**

7 There has been an increasing focus on the development of automation in vehicles due its many
8 potential benefits like safety, improved traffic efficiency, reduced emissions etc. One of the key
9 factors influencing public acceptance of automated vehicle technologies is their level of trust.
10 Development of trust is a dynamic process and needs to be calibrated to the correct levels for safe
11 deployment to ensure appropriate use of such systems. One of the factors influencing trust is the
12 knowledge provided to the driver about the system's true capabilities and limitations. With a 56
13 participant driving simulator study, the authors found that with the introduction of knowledge about
14 the true capabilities and limitations of the automated system, trust in the automated system increased
15 as compared to when no knowledge was provided about the system. Participants experienced two
16 different types of automated systems: low capability automated system and high capability automated
17 system. Interestingly, with the introduction of knowledge, the average trust levels for both low and
18 high capability automated systems were similar. Based on the experimental results, the authors
19 introduce the concept of *informed safety*, i.e., informing the drivers about the safety limits of the
20 automated system to enable them to calibrate their trust in the system to an appropriate level.

21

22

23

24

25

26

27

28

29

30

31

32

33

34 1. Introduction

35 In the last decade there has been a gradual increase of Advanced Driver Assistance Systems (ADASs)
36 (e.g. Adaptive Cruise Control (ACC), Lane-Keep Assist etc.) in vehicles. More recently, there has
37 been a push towards the introduction of higher levels of automation in vehicles with the aim of having
38 Automated Driving (AD) features. The push towards ADAS and AD systems is driven due to their
39 many potential benefits like increased safety leading to reducing the number of accidents (Tingvall,
40 1997; Guérliau *et al.*, 2016; Cicchino, 2017), increased traffic throughput and road efficiency (Le Vine
41 *et al.*, 2016; Talebpour and Mahmassani, 2016), time and monetary savings on parking (Fagnant and
42 Kockelman, 2015), lower emissions (Fagnant and Kockelman, 2014), decreasing drivers' workload
43 (Stanton and Young, 1998; Balfe, Sharples and Wilson, 2015) and providing more productive time to
44 drivers (Cairns *et al.*, 2014).

45 While it is important to provide drivers the opportunity to use ADAS and AD systems (with
46 development in technology), it is equally important to ensure that the drivers actually use the systems
47 in order to ensure the potential benefits from the use of such systems are realized (Lee and See, 2004;
48 Diels and Bos, 2016). Unfortunately, the usage of ADAS features like ACC and Lane Departure
49 Warning has been low (51% of highway driving time (Eichelberger and McCartt, 2014)). Studies
50 discussing the introduction of new technology in different domains like aviation, rail, automotive, etc.
51 have shown that for the new technology to be accepted and used, effort needs to be made to introduce
52 trust towards the new technology (Molesworth and Koo, 2016). Molesworth and Koo (2016)
53 discussed that when participants were given a choice between conventionally piloted aircraft and
54 remotely piloted aircraft (new technology), participants chose the former as they trusted it more.

55 In the driving context, design and behaviour of ADAS and AD systems should be communicated to
56 the driver (Stanton, Young and Mccaalder, 1997) and should be more human-like as it would make
57 the driver-automation cooperation more transparent (Bifulco *et al.*, 2013; Casner, Hutchins and
58 Norman, 2016; Wang *et al.*, 2016), leading to increased trust in the system. One of the challenges
59 with the design of ADAS and AD is that their introduction changes drivers' task from active
60 engagement to passive monitoring (van den Beukel, van der Voort and Eger, 2016). Drivers' driving
61 task is said to have three different levels: 1) strategic 2) tactical and 3) operational (Michon, 1985).
62 ADAS and AD systems alter these levels of driving tasks and the decision to design automation into
63 any of the three levels is generally a trade-off decision (Johansson and Nilsson, 2016; Khastgir,
64 Sivencrona, Dhadyalla, Billing, *et al.*, 2017). The trade-off decision determines the level of
65 engagement of the driver in the driving task. The shift from active engagement to passive monitoring
66 introduces new types of potential errors (human errors) in the driving task as the human driver is not
67 suitable for the task of monitoring monotonous systems (Fitts *et al.*, 1951).

68 1.1. Trust

69 While introduction of automation assumes the removal of human error, in fairness it only shifts the
70 human error from the driver to the designer of the system (Bainbridge, 1983). The designer of the
71 automation makes assumptions about the best design for automation and distribution of driving tasks
72 between the driver and the automated system. These assumption may or may not match with the
73 drivers' perception of the automated system and task distribution. Muir (1994) has suggested that as
74 the automation capability or reliability increases, trust also increases. However, a mismatch between
75 drivers' perception and expectations about the capability of the automated system, and the designers'
76 assumptions can lead to misuse (due to mistrust), disuse (due to distrust) or abuse of the automated
77 system (Parasuraman and Riley, 1997). Misuse is a situation when the driver uses the automated
78 systems for tasks it was not designed to perform and is caused due to mistrust, thus making the
79 situation more unsafe than manual driving. Disuse is a situation when the driver doesn't use the

80 system in situations where the automation is suitable to use, due to distrust, thus not benefiting from
81 the system. Thus, in order to ensure appropriate use of the system, it is essential to calibrate drivers'
82 trust to the appropriate level.

83 Trust is one of the most important factors influencing use of automation (Muir, 1987; Lee and Moray,
84 1992; Muir and Moray, 1996; Parasuraman and Riley, 1997; Parasuraman and Miller, 2004; Rudin-
85 Brown and Parker, 2004; Walker, Stanton and Salmon, 2016). Before the authors discuss details of
86 the development of trust, it is important to define trust in driving context. In order to define trust, the
87 authors adapt the definition of trust from (Lee and See, 2004) as, "*a history dependent attitude that an*
88 *agent will help achieve an individual's goals in a situation characterized by uncertainty and*
89 *vulnerability*". The addition of the reference to "*history dependent*" is particularly important for this
90 work because prior knowledge about the system's capabilities and limitations affects an individual's
91 attitude towards a system, thus affecting their trust. Trust is said to be influenced by various factors
92 (Lee and See, 2004; Xu *et al.*, 2014; Walker, Stanton and Salmon, 2016), with previous work
93 conducted by the authors suggesting this can also include knowledge, certification, situation
94 awareness, workload, self-confidence, experience, consequence and willingness (Khastgir, Birrell,
95 Dhadyalla and Jennings, 2017). In this paper, authors discuss the effect of knowledge on trust.

96 1.1.1. Forms of trust

97 Within scientific literature, trust is often discussed as a single construct. However, inspired by
98 Rajaonah *et al.* (2008) who suggest two forms of trust: trust in automation and trust in the cooperation
99 with automation; for the automotive context, the authors classify trust quantitatively into two forms:

- 100 • Trust *in* the system
- 101 • Trust *with* the system

102 "*Trust in the system*" means the drivers' trust in the capabilities of the system and/or in the system's
103 ability to do what it is supposed to do. "*Trust with the system*" means drivers' awareness or attitude
104 towards the limitations of the systems and their subsequent ability to adapt their use of the system to
105 accommodate for the limitations in order to deliver the expected benefit from the system. Trust with
106 the system implicitly means that the drivers are aware about the true capabilities, and limitations of
107 the system, and are able to adapt their usage to overcome the limitations of the system in real-time.
108 This paradigm of trust is going to be adopted in this paper.

109 1.1.2. Knowledge: a factor influencing trust

110 In order to have appropriate trust, it is important to convey the designer's assumptions about the safe
111 boundaries of the system to the driver. The knowledge of these boundaries provides the ability to have
112 a safe cooperation with the automated system (Beller, Heesen and Vollrath, 2013). In the absence of
113 such knowledge, drivers may not be able to calibrate their trust to an appropriate level (Lee and See,
114 2004; Chavaillaz, Wastell and Sauer, 2016). While failures of automation has been proved to have a
115 detrimental effect on trust, Lee and See (2004) argue that some failures can be classified as "good
116 failures" with negligible impact on trust. Good failures are those whose occurrence is predictable,
117 which allows the driver to be prepared to accommodate for it. Predictability of failures of an
118 automated system comes with knowledge about the true capabilities and limitations of the system.

119 For complex systems requiring supervision, it has been argued that there is a need for an abstraction
120 hierarchical representation of knowledge of the functional properties of the system (Rasmussen,
121 1985). The abstraction hierarchy can potentially be done on two fronts. The first category is a
122 whole/part of the system hierarchy, in which the system is viewed as a number of interacting sub-
123 systems working together at different physical levels (Rasmussen, 1985). The second category
124 suggested in Rasmussen's hierarchical knowledge representation is the abstraction of the functionality
125 (Rasmussen, 1985). The physical form of the system represents the lowest level of abstraction.
126 Moving up through the levels, physical functions represents the next level, next is generalized

127 functions, abstract functions forms the penultimate level with functional purpose forming the highest
128 level of knowledge abstraction. The higher abstraction levels do not just represent the abstraction of
129 physical form, they provide knowledge about the control laws for the interactions of the functions at
130 the lower levels. Moving up the abstraction levels provides a purpose of the task for the level below,
131 while moving down the levels provides information about how the task will be achieved.

132 When put in a driving context, the lower levels of abstraction represent the operational (as per Michon
133 (1985)) driving task (means to a desired end goal) while the higher levels of abstraction represent the
134 tactical and strategic driving tasks (defining the desired end goal). As priority is always given to
135 higher levels of abstraction, a driver has to make a trade-off between the end goal (tactical / strategic
136 goals) and means to achieve it (operational goals), to ensure the means to achieve the goal (lower
137 levels of abstraction) lie within the safe boundaries of the system. In a manual driving task, such a
138 trade-off has clear boundaries and represents a causal system (Rasmussen, 1985). The introduction of
139 automation makes the driving task and the system more complex with blurred boundaries and no
140 simple relationship between function and physical processes making it difficult to represent them as a
141 causal system. Such systems are referred to as intentional systems. For intentional systems (ADAS
142 and AD systems), decision making requires knowledge about the system, its limitations and the actual
143 input to the system (from the environment) and a top-down approach to control the system in a safe
144 manner (Rasmussen, 1985).

145 1.1.3. Types of knowledge

146 Based on literature (Rasmussen, 1985; Seppelt and Lee, 2007; Xu *et al.*, 2014; Biassoni, Ruscio and
147 Ciceri, 2016; Feldhütter *et al.*, 2016; Miller *et al.*, 2016; Bennett, 2017), the following classification for
148 knowledge about the capabilities and limitations of automated systems was proposed by (Khastgir,
149 Birrell, Dhadyalla and Jennings, 2017):

- 150 • Static knowledge: Understanding of the functionality of the automated system (intentions
151 behind the design of the system and functionality) (Larsson, 2012; Eichelberger and McCart,
152 2014). Static knowledge is administered prior to the driving task and is akin to an owner's
153 instruction manual, however with information at a higher abstraction level. Over time, a person
154 can also build up static knowledge based on experiences.
- 155 • Real time knowledge: or dynamic knowledge about the automated system (e.g. automation
156 health, current state of the automation, near-future intentions of the automation). With the help
157 of real-time information about the automated system health, drivers can be brought back "in-
158 to-the-loop" (Louw and Merat, 2017), as it helps increase their awareness (Banks and Stanton,
159 2016) and increase transparency in the cooperation between humans and automation (Eriksson
160 and Stanton, 2017). While in-vehicle information systems (IVISs) are known to have
161 detrimental effect on driving performance (Peng, Boyle and Lee, 2014), they have a potential
162 to have a contrasting effect in an automated vehicle as the driver is not actively involved in the
163 driving task. Real time knowledge during repeated driving cycles leads to supplemental static
164 knowledge of the driver about the capability and limitations of the system as it forms part of
165 the consciously imparted knowledge driver brings to the next use of the automated system.
- 166 • Internal mental model: Prior beliefs influenced of external sources (e.g. word of mouth, media
167 etc.). Marketing of an automated system can affect the public trust and perception towards the
168 product. This can potentially backfire if the information provided in marketing material is
169 inaccurate as customers expect the systems to function as advertised (Casner, Hutchins and
170 Norman, 2016). Inaccurate information can potentially cause over-trust or mistrust in the
171 system. Internal mental model is the pre-conceived notion a person brings to the first use of
172 automation, without any conscious effort to understand the system. While internal mental
173 model is influenced by other sources, static knowledge is consciously imparted to a person prior
174 to the use of automation.

175 Comparing the presented knowledge classification with Rasmussen’s abstraction hierarchies, the
176 authors suggest that static knowledge helps adopt a top-down approach, while dynamic knowledge
177 helps adopt a bottom-up approach. Static knowledge further provides the ability to shift the decision
178 making to a higher level or a lower abstraction level depending on the level of dynamic knowledge
179 provided to the driver, i.e. to facilitate the user to more easily transition between levels of the
180 abstraction hierarchy. With the introduction of automation, complexity of system increases, requiring
181 drivers to demonstrate top-down (mean-end) reasoning approach to accommodate for deviations in
182 performance while receiving knowledge about the operational driving parameters (bottom-up
183 knowledge) (Rasmussen, 1985), to demonstrate their knowledge-based behaviour due to unfamiliar
184 nature of the situations (Rasmussen, 1983). The significance of the abstraction hierarchies can be
185 further illustrated by the fact that causes of failures or incorrect function are explained by a bottom-up
186 approach whereas the reasons for the proper function are explained by a top-down approach
187 (Rasmussen, 1985).

188 Qualitatively, knowledge can potentially be classified into: 1) signals 2) signs and 3) symbols
189 (Rasmussen, 1983). Signals which display time-space sensory data, help the drivers demonstrate skill-
190 based behaviour (based on intuition and experience). While signs indicate towards a stored rule, they
191 do not provide the ability for drivers to process the situation in case a stored rule does not exist in
192 their mental model. Symbols on the other hand represent the relationship between signs and provide
193 the ability for drivers to demonstrate their knowledge-based behaviour and process the information to
194 create a new rule (by shifting the processing to a higher or a lower level of abstraction).

195 1.1.4. Creation of knowledge: identifying failures

196 While, as described above, providing knowledge to the drivers has a potential of increasing trust, it
197 needs to be stressed that the accuracy of the knowledge provided is key. Inaccurate knowledge plays a
198 detrimental role in development of trust as it takes additional cognitive effort on the part of drivers to
199 re-calibrate their mental model (initially formed in accordance to the inaccurate knowledge) to the
200 true capabilities of the system as they experience the system (Beggiato and Krems, 2013).

201 In order to create the knowledge of the true capabilities and functionality of the automated system
202 (i.e., to identify failures), it is essential to conduct a thorough verification and validation process.
203 Moreover, due to the safety critical nature of ADAS and AD systems, their deployment needs to be
204 preceded by extensive testing to establish their safety level and performance boundaries (Sepulcre,
205 Gozalvez and Hernandez, 2013). As discussed in section 1.1.2, the identification of failures helps
206 classify them as “good failures” as it provides a level of predictability about them and thus do not
207 have a detrimental effect of trust (Lee and See, 2004). However, knowledge creation about the
208 capabilities and limitations of ADAS and AD systems faces reliability challenges (Khastgir, Birrell,
209 Dhadyalla, Sivencrona, *et al.*, 2017) and validation challenges which include challenges in test
210 methods and test setup (Hendriks, Pelders and Tideman, 2010; Khastgir *et al.*, 2015; Yu, Lin and
211 Kim, 2016). While the authors consider knowledge creation as an important part of the process of
212 development of trust, it remains out of scope of this paper and will be discussed in future publications.

213 While defining trust in section 1.1, the authors mentioned that trust is a history dependent construct,
214 suggesting its dynamic nature. The authors adopt the definition of calibration of trust as “*the process*
215 *of adjusting trust to correspond to an objective measure of trustworthiness*” (Muir, 1994). Khastgir *et*
216 *al.* (2017a) introduced five stages of calibration of trust: initial phase (stage 1), loss phase (stage 2),
217 distrust phase (stage 3) and recovery phase (stage 4 and stage 5). There can be various intervention
218 methods to potentially increase/adjust trust in different stages of calibration. In this paper, the authors
219 discuss the use of static knowledge as an intervention method in the process of calibration of trust.

220 1.2. Research Question

221 As discussed in section 1.1, many authors have studied the effect of reliability (or automation
222 capability) on trust (Muir, 1994; Muir and Moray, 1996; Chavaillaz, Wastell and Sauer, 2016),
223 suggesting that with increased reliability, trust increases. However, there is no published research on
224 the effect of static knowledge of automation capability on trust in a driving context (both “*trust in the*
225 *system*” and “*trust with the system*”). With the help of a driving simulator study, this paper aims to
226 answer the following two research questions:

- 227 1. Does providing static knowledge about the automation capability of the system influence
228 “*trust in the system*”?
- 229 2. With static knowledge about the automation capability, does automation capability influence
230 “*trust in the system*”?

231 1.2.1. Hypothesis

232 The authors hypothesize that static knowledge influences “*trust in the system*” as it would help
233 influence drivers’ mental model and aid in them exercising their knowledge-based behaviour in
234 unfamiliar situations. Furthermore, the authors believe that static knowledge would have limited
235 effect on drivers’ “*trust with the system*” as drivers’ lack information about the automation health and
236 its intentions. While static knowledge does provide an ability for drivers to predict failures, it does not
237 help them understand the real-time tactical and operational driving task choices made by the
238 automated system.

239 This paper is organized in five sections. Section two discusses the methodology adopted for the study,
240 section three illustrates the results of the study, section four provides a discussion on the results and
241 the paper concludes with a conclusion in section five.

242 2. Methodology

243 2.1. Driving Simulator

244 The experimental study was conducted in WMG’s 3xD simulator for Intelligent Vehicles at the
245 University of Warwick, UK (WMG, 2017). The 3xD simulator consists of a Land Rover Evoque
246 Built-Up Cab (BUC) which is housed inside a cylindrical screen of 8 m diameter and 3 m height. The
247 cylindrical screen provides a 360° field of view for the driver sitting inside the BUC. A push button
248 (with a backlight) (akin to an emergency stop button within a highly autonomous vehicle) was
249 connected (hardwired) to a Raspberry Pi 2 board which in turn was connected to the 3xD simulator
250 through a TCP/IP client-server interface. When the participants pressed the button, the backlight
251 switched-off and the vehicle applied emergency braking and came to a stop. When the participant
252 pressed the button again, the emergency brake was released and vehicle continued to drive in
253 autonomous mode, with the backlight glowing again. This setup enabled a true user in the loop
254 simulation platform, with the user being able to transition in and out of autonomous driving mode
255 anytime they desired, rather than only at predefined, scripted simulator events.

256 2.2. Participants

257 Ethical approval for the experiment was secured from the University of Warwick’s Biomedical &
258 Scientific Research Ethics Committee (BSREC) (REGO-2015-1746 AM02). Fifty six participants (16
259 female and 40 male) were recruited for the study via email invitations. The mean age of the
260 participants was 36.29 years (S.D. = 12.82 years). All participants were required to have a valid, UK
261 full driving license and be at least 21 years of age. The average driving experience of the participants

262 was 14.29 years (S.D. = 13.73 years). The participants' assignment was counter balanced among three
 263 groups which were: 1) control group 2) low (20%) capability automation 3) high (80%) capability
 264 automation. The difference in automation capability is described in section 2.3.2. Informed consent
 265 was obtained from all participants.

266 Out of the 56 participants who took part in the study, eight participants were not able to complete the
 267 study due to simulator sickness and technical issues while running the driving simulator. The 48
 268 participants who completed the study were assigned to three groups (see Table 1).

269 *Table 1: Study design: participant groups*

	Control Group: Without knowledge		Group 1: Low capability automation	Group 2: High capability automation
Number of Participants	8	7	21	12
Run 1	Low capability automation	High capability automation	Without knowledge	Without knowledge
Run 2	High capability automation	Low capability automation	With knowledge	With knowledge

270 2.3. Study Design

271 The experiment was designed as a 2 x 2 mixed factorial design with automation capability as the
 272 between-subject factor, and knowledge of the automation capability as a within-subject factor. For the
 273 control group, automation capability was used as a within-in subject factor to evaluate whether trust
 274 increased with experience without providing any knowledge to the driver (participant) about the
 275 automation capability. As a part of the study, each participant was driven in automated mode (SAE
 276 Level 4 as per SAE J3016 (SAE, 2018)) twice and witnessed five hazardous incidents during each
 277 complete run. Since the study was evaluating SAE Level 4 automation, participants were asked to sit
 278 in the front passenger's seat and hold the emergency stop button in their hands. Such an arrangement
 279 also ensured that the participants could only use the button (instead of brake pedal) to stop the vehicle.
 280 They were further informed that when the emergency stop button was pressed, the vehicle will apply
 281 emergency brakes and will need to cover the braking distance depending on the speed of the vehicle.
 282 In cases where the participant met with a simulated accident, the run ended abruptly. The driving
 283 simulator route for the experiment involved a drive around the University of Warwick campus. Each
 284 complete run lasted around 10 minutes. The route around University of Warwick was chosen to
 285 provide a better immersive environment for the participants as most of them were familiar with the
 286 university campus. Additionally, the University of Warwick route in the 3xD simulator has photo-
 287 realistic imagery and realistic road feedback (vibration) due to a LiDAR scan input which forms the
 288 base for the simulation environment. The speed of the automated vehicle was according to the speed
 289 limits set on the campus map, ranging from 10-30 miles per hour.

290 In order to overcome the lack of real-world consequences often experienced by simulation
 291 participants, who can easily choose not to react as they might if their own life were in jeopardy (as in
 292 real-world), the study had a gamification aspect to it. The game gave participants a goal during the
 293 experiment run and added an element of risk to the study (Table 2). Both these factors have been
 294 discussed in section 1.1 as being essential to evaluate development of trust. Participants were awarded
 295 1 point for every second they spent in automated mode. Every time they pressed the button, the button
 296 press was classified as a "correct stop" or an "incorrect stop". For every correct stop they were
 297 awarded a bonus of 200 points and for every incorrect stop, a penalty of 200 points. Before the run,
 298 they were further provided information about what defined a correct and an incorrect stop. A correct
 299 stop was one where the participant correctly identified that the automated system wouldn't be able to
 300 handle the situation, prompting the participant to intervene and press the emergency stop button. An
 301 incorrect stop was one in which the participant pressed the emergency stop button and brought the
 302 vehicle to standstill, even though the automated system was capable of handling the situation.

303 Additionally, in case any participant crashed (met an accident), a penalty of 10000 points was given
 304 and the experiment run came to an end.

305 An extremely high penalty was added for a crash to add a high degree of risk and motivate
 306 participants to avoid crashing the vehicle as perceived risk influences driver’s interaction with the
 307 automated system (Eriksson, Banks and Stanton, 2017). The penalties were added to get the
 308 participants to react in a similar manner as if they were in real danger. The participants were asked to
 309 maximise their score. However, the score was not a variable within the study. It was more of a
 310 mechanism to encourage engagement in the task. Participants were provided information about their
 311 score after the study was completed. Participants were given two objectives: 1) avoid crashing the
 312 vehicle by pressing the button (emergency stop) 2) maximize time spent in automated mode. They
 313 were asked to press the button only if they felt that the automated system couldn’t handle the situation
 314 or if they felt unsure about the automated system’s performance.

315 *Table 2: Scoring criteria for study (gamification)*

Type of Action	Points
Automated mode	1 / second
Correct Stoppage of the automated vehicle	+200
Incorrect Stoppage of the automated vehicle	-200
Crash	-10000

316

317 **2.3.1. Hazards**

318 In order to choose the five hazardous events, a hazard analysis of an automated vehicle was conducted
 319 as per the ISO 26262 (ISO, 2011) functional safety process. Five different automated vehicle
 320 functions were identified and a hazard was identified for each of the functions (Table 3). For each
 321 hazard, a hazardous event was identified which was created in each of the driving scenarios in the
 322 experiment runs in the 3xD simulator. The hazard and hazardous event identification was done by
 323 independent safety experts. One of the factors influencing the selection of the hazardous events was
 324 the ability to create the events in the 3xD simulator.

325 *Table 3: Description of five hazardous events*

Function	Hazard	Hazardous event description
Braking	Lack of Braking	Pedestrian suddenly changes direction and comes in front of the ego vehicle (automated vehicle)
Torque	Excessive torque – excessive acceleration	Vehicle approaching round-about and accelerates instead of braking
Object Detection	Blind-spot and delayed object detection	Another vehicle in perpendicular lane comes in path of the ego vehicle suddenly
Path Planning	Not following rules of road	Ego vehicle joins a roundabout while another vehicle is still in the roundabout and has right of way.
Object Detection	Compromised detection due to environmental factors	In foggy/rainy weather, ego vehicle is not able to detect traffic lights within the specified range.

326

327 **2.3.2. Automation Capability**

328 Two levels of automation capability were used in the study: 1) low capability automation 2) high
 329 capability automation. The difference between the two systems was based on the ability of the
 330 automated system to tackle the five hazardous events mentioned in section 2.3.1. Low capability
 331 automated system was able to handle one out of the five hazardous events, requiring the driver to
 332 intervene in four hazardous events to ensure safe performance of the vehicle. High capability
 333 automated system was able to handle four out of the five hazardous events, requiring the driver to
 334 intervene in only one hazardous event situation to ensure safe performance.

335 2.4. Procedure

336 When participants arrived for the experiment, they were initially briefed about the experiment
337 following which informed consent was taken from each participant and they were asked to fill in a
338 demographic questionnaire. Before the start of the study runs, each participant was given a trial run
339 (on a route different from the one used for the study runs) on the driving simulator with a researcher
340 seated next to the participant, to familiarize the participant with the visuals, motion feedback,
341 experience of sitting inside a car within a simulator and using the button to apply emergency brake on
342 the vehicle. Participants were told that they can ask for as many trial runs as they wish, in order to
343 make them comfortable with the simulator environment. Each trial run was of five minutes in length.
344 While most of the participants requested only one trial run, some participants requested for an
345 additional (second) trial run. After the trial runs, participants were asked whether they would like to
346 continue the study. In the case that the participant agreed, each participant experienced two
347 experiment runs of around 10 minutes each. Before the second run (for group 1 and group 2),
348 participants were provided knowledge about the capabilities of the automated system. Commentary
349 was read out to them via a prepared script. Effort was put into the preparation of the script in order to
350 avoid introducing any experiment bias. The script was reviewed by three independent human factors
351 experts.

352 For the control group, participants were told that in the two runs, they will experience automated
353 control systems from two different suppliers. No other information about system capabilities was
354 given. However, before the second run, it was reiterated that the participants will now experience a
355 different automated control system from a different supplier. Such a design of the control group was
356 implemented to check if there was any changes in the trust levels due to experience. Eight out of 15
357 participants in the control group experienced low capability automation in their first run and high
358 capability automation in their second run. The remaining seven participants experienced the runs in
359 the reverse order.

360 At the end of each experiment run, participants were asked to fill a trust rating questionnaire (section
361 2.4.2), Simulator Sickness Questionnaire (SSQ) (Kennedy *et al.*, 1993), and Van Der Laan's
362 acceptance questionnaire (Van Der Laan, Heino and De Waard, 1997). However, the results from the
363 latter two haven't been included in this paper.

364 2.4.1. Imparting knowledge

365 Knowledge was imparted to the participants via a prepared script which included illustrations
366 regarding the automated systems' capability and limitations. Special care was taken to ensure that
367 participant's mental model was informed so that they understood the functioning of the system in a
368 lay-man language to ensure higher level system understanding. This was particularly important in
369 order to ensure they were imparted with knowledge-based behaviour, as compared to rule-based or
370 skill-based behaviour. The knowledge imparted would enable them to deal with the unfamiliar
371 situation by transferring the cognitive task to a higher level or a lower level of abstraction in search of
372 an existing rule or intuition of their mental model (Rasmussen, 1985). In the automated driving
373 context, the significance of knowledge-based behaviour is further emphasized as it helps a driver
374 adopt a means-end approach to execute the appropriate human intervention needed for the task. The
375 following two scripts are examples of the how knowledge was imparted to the participants.

376 *Example 1:* "The automated control system from the supplier is based on camera based sensors and
377 each automated control system will be trialled in separate runs in the sim. However, due to cost
378 pressures, they have chosen a single low quality camera with reduced field of view.

379 *Vision based systems are dependent on the quality of the camera used. Due to cost pressures, the*
380 *supplier has compromised with the accuracy of the camera used for the vehicle. In this vehicle, a*
381 *lower grade camera has been used. Lower grade cameras are vulnerable to environmental factors*

382 and image recognition degrades with lower visibility. E.g., certain cameras find it hard to detect
383 objects in rain, snow or fog or at certain times of the day due to image washout (Figure 1). In your
384 drive today, you might have witnessed bright sunlight or rain. You have the luxury of using
385 sunglasses, wipers etc. However, Camera doesn't have that. It has been found that light colour
386 objects against a bright sky is difficult to detect. This was the case in the recent Tesla Model S crash
387 (NHTSA, 2017) where the white rear end of the truck was not detected against the bright sky."

388



Figure 1: Camera view while driving in fog
(image source: <https://www.flickr.com/photos/kubina/2160242894/>; date accessed: 2017-12-04)

389

390 Example 2: "While, automated vehicles have a repeatable and predictable behaviour, their behaviour
391 is "programmed" by human engineer. Every vehicle before being released to market undergoes
392 rigorous testing. However, it is possible that sometimes a programming bug introduced by a human
393 error manifests itself into a larger failure. The rules of the road are pre-programmed into the
394 automated control system. The automated system in your next run is a pre-production control system
395 and is still undergoing testing. While previous test results have been extremely positive, I advise you
396 to take caution. An example of this might be that as a driver, we know that if a pedestrian is standing
397 next to a zebra crossing, they have the right of way (Figure 2). However, for a camera system, he/she
398 will only be a pedestrian with unknown intention. In this example the automated control system
399 wouldn't know the rules of the road and will not have the understanding of the priorities.

400 Another rule of the road that we as drivers are used to is the priorities at roundabouts and junctions
401 (Figure 2). Imagine a person is given a driving license when he/she doesn't know the rules of the
402 road. Not only its dangerous for him/her, it is hazardous for the traffic around."

403

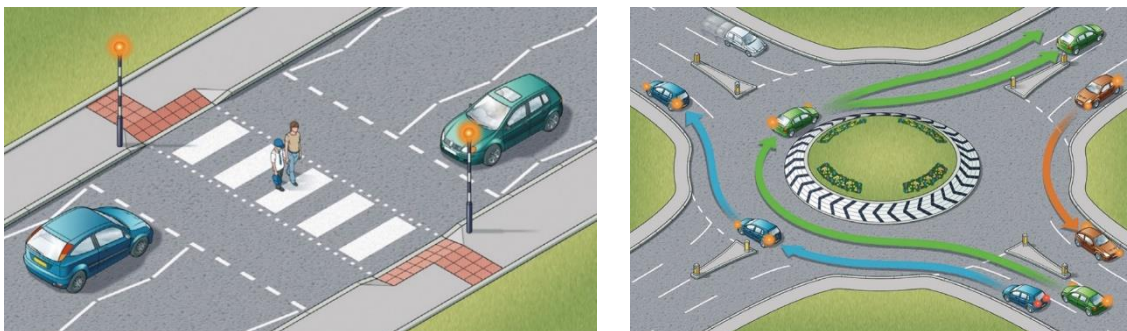


Figure 2: Rules of road: rule 19 (left) and rule 185(right). (DfT, 2017)

404

405 In the above examples, effort was made to differentiate between knowledge and rule-based
406 behaviours. Simple rules are comparatively easy to convey to participants, for Figure 1, a rule would
407 be ‘automated system will not work in fog’. However, there is no understanding why it will not work
408 (e.g. *image recognition degrades with lower visibility* which was provided as a part of the script).
409 Knowledge about other similar situation where the camera may not work was also provided via the
410 script (...*hard to detect objects in rain, snow or fog or at certain times of the day*); (*You have the*
411 *luxury of using sunglasses, wipers etc. However, Camera doesn’t have that. It has been found that*
412 *light colour objects against a bright sky is difficult to detect. This was the case in the recent Tesla*
413 *Model S crash where the white rear end of the truck was not detected against the bright sky*). By
414 trying to impart knowledge the participant can envisage their own varied and numerous situations
415 where the automated system might act unexpectedly.

416 2.4.2. Trust questionnaire

417 At the end of each of the two experiment runs, participants were asked to rate their level of “*trust in*
418 *the system*” and “*trust with the system*”. A subjective rating scale was used and participants were
419 asked to draw a line across a 100 mm box to indicate their level of trust (c.f. (Muir and Moray, 1996;
420 Rajaonah, Anceaux and Vienne, 2006)). Before being asked to rate different trust levels, participants
421 were briefed about the difference in the different types of trust via a prepared script which included
422 examples (was read to the participants as well as given in text form) to highlight the difference
423 between “*trust in the system*” and “*trust with the system*”. Existing rating scales like Jian’s scale
424 (Jian, Bisantz and Drury, 2000), couldn’t be used as they don’t classify trust into the two components
425 mentioned in section 1.1. In order to explain the two different concepts of trust, participants were
426 briefed using an example of a mobile phone and call service provider. The following text was used for
427 the explanation:

428 “*Trust in the system means that you have trust in the capabilities of the system and in its ability to do*
429 *what it is supposed to do as advertised to you. In other words, it does what it says on the box. Trust with*
430 *the system means that you are aware of the limitations of the systems and you adapt your use of the*
431 *system to accommodate for the limitations in order to get maximum benefit from the system.*”

432 *For example, if you buy a mobile phone, you have trust **in** the systems about its advertised*
433 *capabilities. You develop trust **with** the system once you start using it and understand its limitations.*
434 *Ability to work with limitations guides your trust **with** the system. For the mobile phone and the call*
435 *service provider you have, you get call drop-outs in certain part of our house and not in another part*
436 *of your house. You would adapt your usage of the mobile phone by making calls only when you are in*
437 *a part of the house where you know call connection service is good. This is an example of you*
438 *acknowledging the limitations of the system, adapting your usage and developing trust with the*
439 *system”*

440 On the trust scale, a 0% rating suggested very low trust and 100% suggested very high trust. As trust
441 is a continuum, any value in between 0 -100 suggests that the participant had partial trust.

442 3. Results

443 3.1. Trust levels

444 The average “*trust in the system*” for low capability automation increased substantially from 32.4% to
445 65.4 %, with the introduction of knowledge about the system capabilities and limitations (Figure 3).
446 While an increase in “*trust in the system*” rating with the introduction of knowledge was seen for high
447 capability automation from 54.2% to 70.5% also, the effect was comparatively lower. It is interesting
448 to note that with the introduction of knowledge about the automated system’s capabilities and
449 limitations, both median and mean values for “*trust in the system*” for low-capability and high-

450 capability automated system were similar (Figure 3). In the low capability automation group, barring
 451 two participants out of the 21 participants, all participants showed an increase in trust in the system
 452 with the introduction of knowledge (Figure 4). High capability automation group also showed a
 453 similar trend. The box-plots for trust in the system illustrate a higher convergence in trust ratings with
 454 the introduction of knowledge, potentially due to appropriate calibration of trust level (Figure 3).

455 A repeated measures ANOVA was conducted for the “*trust in the system*” and “*trust with the*
 456 *system*” ratings with automation capability as the between factor variable and knowledge as the
 457 within factor variable. The introduction of knowledge about the automation capabilities and
 458 limitations had a highly significant statistical effect on the level of “*trust in the system*”, $F(1, 31) =$
 459 33.712 , $p = 0.000002$ with a $\eta_p^2 = 0.521$, suggesting 52.1% of the variance being associated with the
 460 introduction of knowledge. The introduction of knowledge didn’t have an interaction effect with
 461 automation capability, $F(1, 31) = 3.846$, $p = 0.059$ ($\eta_p^2 = 0.11$). Therefore, there was no effect of
 462 automation capability on trust in the system ratings when knowledge was introduced.

463 While the average “*trust with the system*” changed with the introduction of knowledge (Figure 5), the
 464 effect was statistically insignificant, $F(1, 31) = 3.652$, $p = 0.065$ with a $\eta_p^2 = 0.105$. There was no
 465 interaction effect between knowledge and automation capability for trust with the system ratings, $F(1,$
 466 $31) = 0.742$, $p = 0.396$ ($\eta_p^2 = 0.023$).

467 In order to negate the effect of experience on trust ratings, a repeated measures ANOVA was
 468 performed on the control group. The effect of the runs was statistically highly insignificant on the
 469 level of “*trust in the system*”, $F(1, 13) = 0.105$, $p = 0.751$ with a $\eta_p^2 = 0.008$. There were no
 470 interaction effects between the runs and the two control groups, $F(1, 13) = 0.020$, $p = 0.89$ ($\eta_p^2 =$
 471 0.002).

472

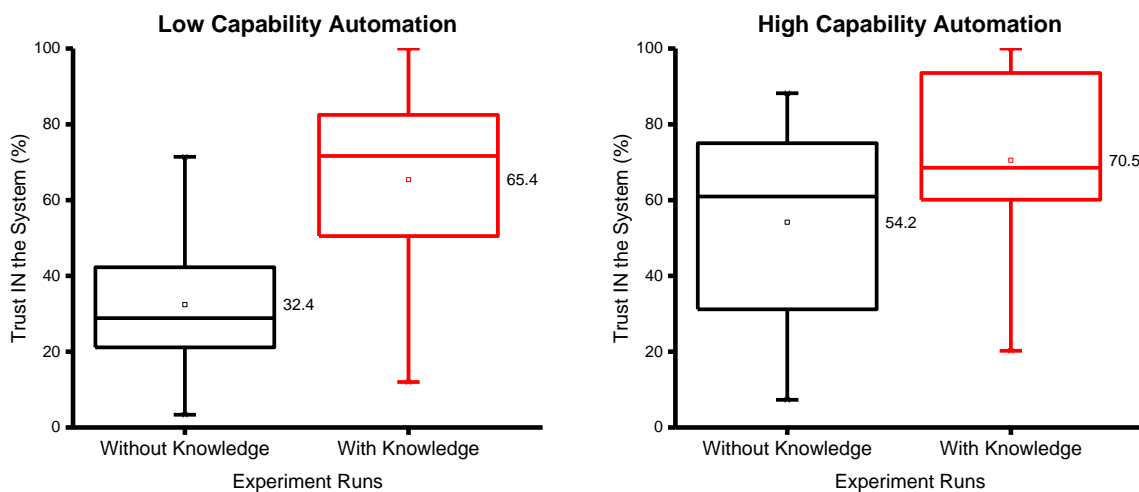


Figure 3: Box-plots of Trust-In the system ratings (highlighting average trust ratings) (central dot represents average value)

473

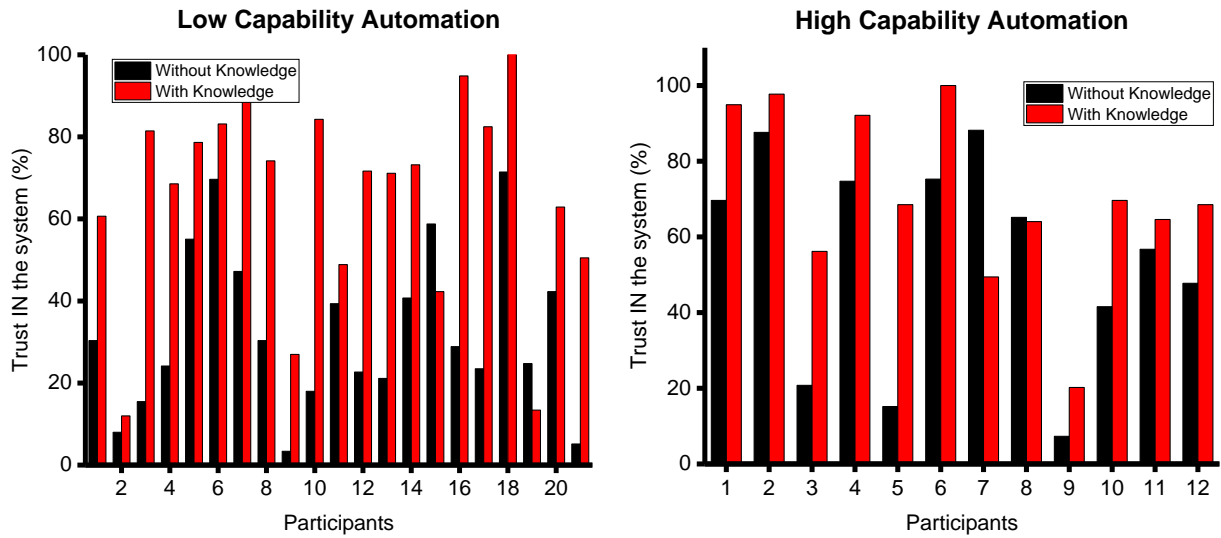


Figure 4: “Trust in the System” level of individual participants for low capability and high capability automation

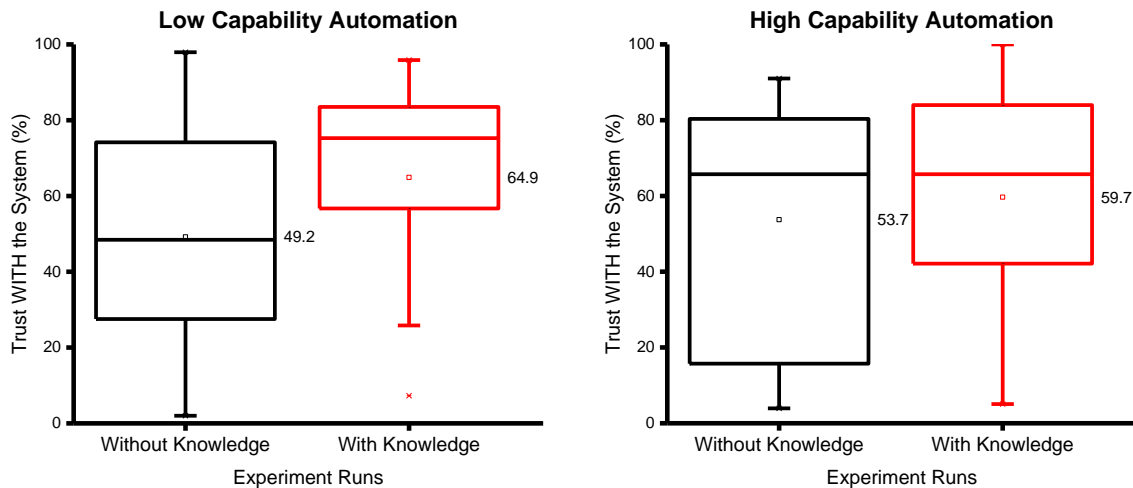


Figure 5: Box-plots of Trust-With the system ratings (central dot represents average value)

477 3.2. False presses

478 While the introduction of knowledge about system capabilities and limitations increased trust in the
 479 system for both low and high capability automation, it had contrasting effect in the two groups in
 480 terms of number of false presses. The authors define a false press as a button press in a situation
 481 which could be handled by the automated system, indicating distrust in the system.

482 For low capability automation, the average number of false presses increased significantly from 0.47
 483 to 2.67 with the introduction of knowledge. On the contrary, for high capability automation the
 484 average number of false presses decreased from 1.73 to 1.36 with the introduction of knowledge
 485 (Figure 6). The outlier data from the box-plot were removed for mean calculation. This meant one

486 data point each from the two runs for high capability automation was removed. There were no outliers
487 in the data set for low capability automation group.

488 A paired-sample t-Test was conducted to assess the significance in the number of false presses with
489 the introduction of knowledge. For low capability automation, there was a statistically significant
490 difference in the number of False Presses for without knowledge run ($M = 0.47$, $SD = 0.60$) and
491 knowledge run ($M = 2.67$, $SD = 1.65$); $t(20) = -6.398$, $p = 0.000003$. For high capability automation,
492 the number of False Presses (FP) for without knowledge run ($M = 2.41$, $SD = 2.79$) and knowledge
493 run ($M = 1.67$, $SD = 1.43$) was statistically insignificant; $t(11) = 0.792$, $p = 0.445$.

494 As discussed in section 2.4.1, for the low capability automation group, participants were given a lot of
495 knowledge based on the automated systems' limited capability. One of the potential reasons for the
496 contrasting results between the two groups could be the amount of knowledge provided in the low
497 capability automation group and the participants' ability to process all the knowledge, develop
498 accurate mental model and display knowledge-based behaviour. However, higher trust ratings with
499 introduction of knowledge suggest that knowledge-based behaviour was displayed. Another potential
500 reason for the contradictory results could be the lack of dynamic (real-time) knowledge provided to
501 the participants (discussed in section 4).

502

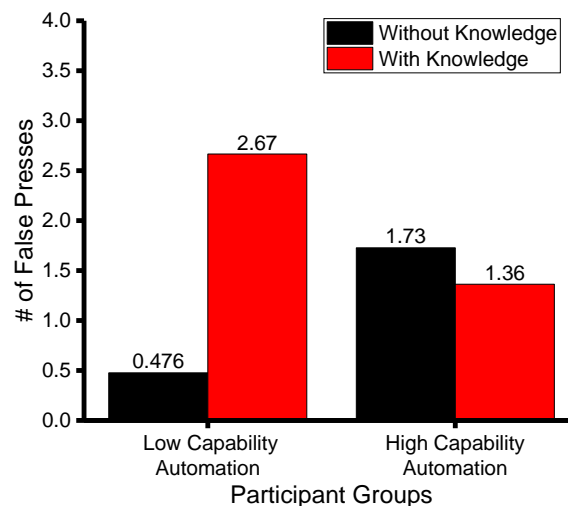


Figure 6: Average number of false presses

503

504 3.3. Accidents

505 The authors define an accident as a collision of the ego vehicle (automated vehicle) with other entities
506 (vehicles, pedestrians or cyclists) in the scenario or if the own vehicle doesn't follow the traffic light
507 rules. Introduction of knowledge about the automated system capability had similar effect on the
508 average number of accidents for both the automation groups. For low capability automation, the
509 average number of accidents reduced significantly from 1 to 0.38 with the introduction of knowledge
510 (Figure 7). For high capability automation, the average number of accidents reduced slightly from
511 0.58 to 0.42 (Figure 7). It is interesting to note that most of the accidents were caused to due to late
512 interventions rather than absence of interventions. This may be explained due to lack of accurate
513 situation awareness about scenario handling capabilities of the automated system during the
514 automated driving scenario which could potentially be due to the lack of dynamic knowledge of the
515 participants. A paired sample t-Test was conducted to assess the statistical significance in the number
516 of accidents with the introduction of knowledge. There was a statistically significant difference in the

517 number of accidents between the without knowledge ($M = 1$, $SD = 0$) and with knowledge ($M = 0.38$,
518 $SD = 0.49$) conditions; $t(20) = 5.701$, $p = 0.000014$, for low capability system.

519 Similar to the false presses, the number of accidents for without knowledge ($M = 0.5$, $SD = .52$) and
520 with knowledge runs ($M = 0.42$, $SD = 0.51$) conditions for high capability automation was
521 insignificant; $t(11) = 0.321$, $p = 0.754$.

522

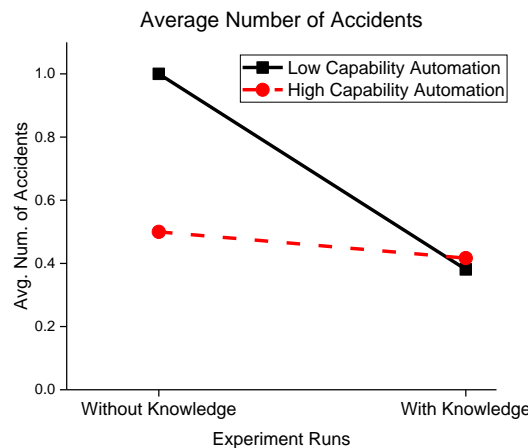


Figure 7: Average number of accidents

523

524 4. Discussion

525 As mentioned in section 1.1.1, “*trust in the system*” refers to the capability of the system where as
526 “*trust with the system*” refers to the ability of the driver to work with the system. In the study
527 presented, the authors have illustrated that with the introduction of knowledge about the system
528 capabilities and limitations, “*trust in the system*” increases, to similar trust ratings for low-capability
529 and high-capability systems. These results differ from the study in (Helldin *et al.*, 2013) and (Hergeth,
530 Lorenz and Krems, 2017). While these studies did provide some feedback about the system
531 boundaries to the drivers, they were unable to instil knowledge-based behaviour as they didn’t
532 mention how the system works due to which the driver’s higher level mental model could not be
533 made.

534 It is worth noting that the effect of knowledge on “*trust in the system*” had a statistically highly
535 significant relationship ($p = 0.000002$), the effect of knowledge on “*trust with the system*” was
536 statistically not significant ($p = 0.065$). This can be explained by analysing the nature of knowledge
537 provided to the participants. As mentioned in section 1.1.2, knowledge can be qualitatively classified
538 into three categories. In the study presented, participants were provided with only static knowledge
539 about the capabilities and limitations of the systems. While this allowed them to demonstrate their
540 knowledge-based behaviour and helped them calibrate their trust in the system, the lack of system
541 feedback on the real-time state and intention of the system, led to lower levels of trust with the
542 system. This inference is further corroborated by the qualitative feedback from participants who were
543 asked to explain their rating of trust in their own words. One of the participants (participant #20)
544 commented: “*warnings from the car missing*” while other (participant # 40) commented “*no*
545 *warnings & notification*”. Another participant (participant #37) mentioned: “*I was able to*
546 *accommodate for the system but it was discomforting... near misses and close calls*”.

547 In other words, the introduction of static knowledge provided participants the capability to
548 demonstrate top-down understanding as per the abstraction hierarchy levels. However, with the
549 absence of dynamic knowledge, they were unable to get feedback (signs and signals) on the causes of
550 the failure, subsequently their reasoning capability was limited. Thus, in order to be able to work with
551 the system, i.e. accommodate for the limitations of the system and display their knowledge-based
552 behaviour appropriately, participants also require real-time knowledge (e.g. signals and signs) to
553 move the decision task to a higher or a lower abstraction level in search of pre-existing rules or
554 intuition, similar to a co-pilot in the aviation domain (Eriksson and Stanton, 2017). Thus, the authors
555 suggest that “*trust with the system*” is potentially influenced to a larger extent by dynamic (real-time)
556 knowledge about the system capabilities and limitation.

557 The introduction of knowledge didn’t have an interaction effect with automation capability on trust
558 ratings ($p = 0.059$ for “*trust in the system*” and $p = 0.065$ for “*trust with the system*” ratings). Thus
559 suggesting that similar levels of trust can be achieved if knowledge about the true capabilities and
560 limitations of the systems is provided to the driver.

561 While due to the study design the control group’s trust ratings can’t be compared with the low-
562 capability automation or high-capability automation group’s trust ratings, they do provide more
563 confidence in the results obtained in the two latter groups. The role of the control group was to either
564 support or negate the hypothesis that any change in trust ratings could be a result of experience.
565 Results showed that automation capability has no interaction effect on experience of the system ($p =$
566 0.89), thus negating the hypothesis.

567 4.1. Informed Safety

568 Results from this study could infer that vehicle manufacturers may choose to introduce low-capability
569 systems and provide knowledge in order to deliver increased user trust and overall system
570 performance. However, there is a caveat to this inference. For low capability automation, while
571 introduction of knowledge increased the level of trust in the system significantly (from 32.4% to
572 65.4%), it also increased the number of false presses significantly (from 0.476 to 2.67). Therefore,
573 very low capability and too much knowledge is also not an appropriate solution. The authors believe
574 that there is an optimum level of system capability and knowledge to be imparted at which trust could
575 be maximized and false presses could be minimized. Therefore, manufacturers may decide to enhance
576 automation capability by providing knowledge. Until systems are fully (100%) capable, augmenting
577 system capability with knowledge about the system’s true capabilities, could be a method to bridge
578 the gap in trust. In other words, while manufacturers should aim to introduce high capability systems
579 in the market, the gap in system capability (system limitations) should be provided as knowledge to
580 the customers to ensure high trust in the system.

581 It is important to appreciate the difference in the manner in which non-specialists (i.e. general public)
582 would understand / interpret the knowledge imparted to them. As creators of the system, designers
583 and engineers have an appreciation and inclination towards technical understanding and the technical
584 feature explanation. Therefore, in this study care was taken in the language used in the script used to
585 impart knowledge to the participants. Use of technical jargon terms was avoided and illustrations were
586 used as examples to help participants visualize the system. In real life, it is important that
587 manufacturers explain the system capabilities and limitations in a non-technical manner in order to aid
588 customer’s understanding by providing examples and ensuring the people read the provided
589 information.

590 This paper introduces the concept of “*informed safety*”, as a means to calibrate trust to the appropriate
591 levels, which may include increasing those with low trust in capabilities or even reducing trust in
592 those with too much confidence in what the system can achieve by making them aware of system

593 boundaries. Informed safety means informing the driver (via static and/or dynamic knowledge) about
594 the safety limits of the automated system and its intention. Informed safety provides the ability to
595 display knowledge-based behaviour to shift the interpretation of a scenario to higher abstraction level
596 or a lower abstraction level (Rasmussen, 1983). Informed safety aids the driver to interpret an
597 unexpected situation to adopt an appropriate tactical or strategic manoeuvre to handle the situation
598 safely. Informed safety is not just about providing rules of usage, it includes the background
599 information, understanding and knowledge about how the system operates.

600 4.2. Future research

601 It is a well-known fact that users don't read manuals and that vehicle dealers/Original Equipment
602 Manufacturers (OEMs) rarely do a good job in sufficiently or appropriately informing customers
603 about the system capabilities and limitations (Beggiato and Krems, 2013; Eichelberger and McCartt,
604 2014; Larsson, Kircher and Hultgren, 2014). As automated systems are introduced, innovative
605 methods of informing the driver (customer) to create an "*informed safety*" level, need to be
606 implemented. One potential solution could be providing a virtual tour of the vehicle at the dealership,
607 which gives the customers an immersive experience of the various features and can help them
608 calibrate their mental models and their expectations from the vehicle. Other means of providing
609 "*informed safety*" may be short videos on the working of the Human Machine Interface (HMI) or
610 specifically designed voice assistant features. All the discussed methods may form a part of the initial
611 showroom briefing or a pre-sale briefing. However, these methods need to be evaluated to measure
612 their effectiveness.

613 4.3. Study limitations

614 The WMG's 3xD simulator provides a fully immersive driving experience for participants. However,
615 like all simulator studies, transferability of results to real world needs to be evaluated separately. Real-
616 world evaluation of trust remains out of the scope of this paper. Additionally, as discussed in section
617 4.1, informed safety, as introduced in this paper, has two facets: 1) static knowledge (e.g. initial
618 briefing and driving manual) and 2) dynamic knowledge such as human-machine interface. In this
619 paper, the authors only provided static informed safety to drivers. Future studies are planned where
620 participants will be provided both dynamic knowledge and static knowledge. Results will be
621 published in future publications.

622 5. Conclusion

623 Trust in automated systems is one of the key factors that would help realize the potential benefits
624 offered by the introduction of automation in vehicles. However, trust level needs to be calibrated to
625 the appropriate level in order to reap the benefits of the automated systems in a safe manner by
626 preventing misuse or disuse. This study explores the effect of knowledge about the automation
627 capability on trust in the system.

628 In this paper, the authors demonstrate via a 56 participants driving simulator study that "trust in the
629 system" increases with the introduction of static knowledge about the capabilities and limitation of the
630 automated system. With the introduction of static knowledge, trust in the system for both low
631 capability automation and high capability automation were not significantly different, 65.4% and
632 70.5% respectively, suggesting no influence of automation capability on trust in the system when
633 knowledge is provided to the drivers. Based on results, the authors introduced the concept of
634 "*informed safety*" which helps calibrate drivers' trust to an appropriate level, subsequently ensuring
635 safe use of the automated system.

636 Interestingly, with the introduction of static knowledge the average number of false presses had
637 contrasting results for the two automation groups. With the introduction of knowledge, for the high
638 capability automation group, the average number of false presses decreased from 1.73 to 1.36, while it
639 increased from 0.47 to 2.67 for the low capability automation group. However, average number of
640 accidents decreased from 1 to 0.38 and from 0.58 to 0.42 for low capability automation and high
641 capability automation respectively. The improved safety with the introduction knowledge lends its
642 support to the concept of informed safety. In order to reduce the number of false presses, the authors
643 hypothesize the need to provide “informed safety” in a dynamic manner, i.e., via knowledge about the
644 automation state and health through the HMI system. Results on the study exploring the hypothesis
645 will be presented in future publications.

646 Acknowledgements

647 The work presented in this paper has been carried under the EPSRC Grant (Grant EP/K011618/1).
648 The authors would like to thank the WMG centre of HVM Catapult and WMG, University of
649 Warwick, UK, for providing the necessary infrastructure for conducting this study. WMG hosts one of
650 the seven centres that together comprise the High Value Manufacturing Catapult in the UK. The
651 authors would also like to thank Andrew D. Moore and Jonathan Smith for their assistance in building
652 the experimental setup. The authors would also like to thank three anonymous reviewers for their
653 detailed comments on previous versions of the paper, which has helped considerably to improve the
654 paper.

655 References

- 656 Bainbridge, L. (1983) ‘Ironies of automation’, *Automatica*, 19(6), pp. 775–779. doi: 10.1016/0005-1098(83)90046-8.
- 657 Balfe, N., Sharples, S. and Wilson, J. R. (2015) ‘Impact of automation: Measurement of performance, workload and behaviour in a complex
658 control environment’, *Applied Ergonomics*. Elsevier Ltd, 47, pp. 52–64. doi: 10.1016/j.apergo.2014.08.002.
- 659 Banks, V. A. and Stanton, N. A. (2016) ‘Keep the driver in control: Automating automobiles of the future’, *Applied Ergonomics*. Elsevier
660 Ltd, 53, pp. 389–395. doi: 10.1016/j.apergo.2015.06.020.
- 661 Beggiato, M. and Krems, J. F. (2013) ‘The evolution of mental model, trust and acceptance of adaptive cruise control in relation to initial
662 information’, *Transportation Research Part F: Traffic Psychology and Behaviour*, 18, pp. 47–57. doi: 10.1016/j.trf.2012.12.006.
- 663 Beller, J., Heesen, M. and Vollrath, M. (2013) ‘Improving the Driver-Automation Interaction: An Approach Using Automation
664 Uncertainty’, *Human Factors*, 55(6), pp. 1130–1141. doi: 10.1177/0018720813482327.
- 665 Bennett, K. B. (2017) ‘Ecological interface design and system safety: One facet of Rasmussen’s legacy’, *Applied Ergonomics*. Elsevier Ltd,
666 59, pp. 625–636. doi: 10.1016/j.apergo.2015.08.001.
- 667 van den Beukel, A. P., van der Voort, M. C. and Eger, A. O. (2016) ‘Supporting the changing driver’s task: Exploration of interface designs
668 for supervision and intervention in automated driving’, *Transportation Research Part F: Traffic Psychology and Behaviour*, 43, pp. 279–
669 301. doi: 10.1016/j.trf.2016.09.009.
- 670 Biassoni, F., Ruscio, D. and Ciceri, R. (2016) ‘Limitations and automation: The role of information about device-specific features in ADAS
671 acceptability’, *Safety Science*, 85, pp. 179–186. doi: 10.1016/j.ssci.2016.01.017.
- 672 Bifulco, G. N. et al. (2013) ‘Driving behaviour models enabling the simulation of Advanced Driving Assistance Systems: Revisiting the
673 Action Point paradigm’, *Transportation Research Part C: Emerging Technologies*, 36, pp. 352–366. doi: 10.1016/j.trc.2013.09.009.
- 674 Cairns, S. et al. (2014) ‘Sociological perspectives on travel and mobilities: A review’, *Transportation Research Part A: Policy and Practice*,
675 63, pp. 107–117. doi: 10.1016/j.tra.2014.01.010.
- 676 Casner, S. M., Hutchins, E. L. and Norman, D. (2016) ‘The Challenges of Partially Automated Driving’, *Communications of the ACM*,
677 59(5), pp. 70–77. doi: 10.1145/2830565.
- 678 Chavaillaz, A., Wastell, D. and Sauer, J. (2016) ‘System reliability, performance and trust in adaptable automation’, *Applied Ergonomics*,
679 52, pp. 333–342. doi: 10.1016/j.apergo.2015.07.012.
- 680 Cicchino, J. B. (2017) *Effectiveness of forward collision warning and autonomous emergency braking systems in reducing front-to-rear
681 crash rates, Accident Analysis and Prevention*. doi: 10.1016/j.aap.2016.11.009.
- 682 DfT (2017) *The Highway Code*. Available at: <https://www.gov.uk/guidance/the-highway-code> (Accessed: 18 July 2017).
- 683 Diels, C. and Bos, J. E. (2016) ‘Self-driving carsickness’, *Applied Ergonomics*, 53, pp. 374–382. doi: 10.1016/j.apergo.2015.09.009.
- 684 Eichelberger, A. H. and McCart, A. T. (2014) ‘Volvo drivers’ experiences with advanced crash avoidance and related technologies.’,
685 *Traffic Injury Prevention*, 15(2), pp. 187–195. doi: 10.1080/15389588.2013.798409.
- 686 Eriksson, A., Banks, V. A. and Stanton, N. A. (2017) ‘Transition to manual: Comparing simulator with on-road control transitions’,
687 *Accident Analysis and Prevention*, 102, pp. 227–234. doi: 10.1016/j.aap.2017.03.011.
- 688 Eriksson, A. and Stanton, N. A. (2017) ‘The chatty co-driver: A linguistics approach applying lessons learnt from aviation incidents’, *Safety
689 Science*, 99, pp. 94–101. doi: 10.1016/j.ssci.2017.05.005.
- 690 Fagnant, D. J. and Kockelman, K. (2015) ‘Preparing a nation for autonomous vehicles: Opportunities, barriers and policy
691 recommendations’, *Transportation Research Part A: Policy and Practice*, 77, pp. 167–181. doi: 10.1016/j.tra.2015.04.003.
- 692 Fagnant, D. J. and Kockelman, K. M. (2014) ‘The travel and environmental implications of shared autonomous vehicles, using agent-based

693 model scenarios', *Transportation Research Part C: Emerging Technologies*, 40, pp. 1–13. doi: 10.1016/j.trc.2013.12.001.

694 Feldhütter, A. *et al.* (2016) 'Trust in Automation as a matter of media and experience of automated vehicles.', in *Proc. of the Human*

695 *Factors and Ergonomics Society 60th Annual Meeting*, pp. 2024–2028.

696 Fitts, P. M. *et al.* (1951) *Human engineering for an effective air - navigation and traffic - control system*. Washington, D.C., USA.

697 Guériau, M. *et al.* (2016) 'How to assess the benefits of connected vehicles? A simulation framework for the design of cooperative traffic

698 management strategies', *Transportation Research Part C: Emerging Technologies*, 67, pp. 266–279. doi: 10.1016/j.trc.2016.01.020.

699 Helldin, T. *et al.* (2013) 'Presenting system uncertainty in automotive UIs for supporting trust calibration in autonomous driving', in *Proc.*

700 *of the International Conference on Automotive User Interfaces and Interactive Vehicular Applications - AutomotiveUI '13*, pp. 210–217.

701 doi: 10.1145/2516540.2516554.

702 Hendriks, F., Pelders, R. and Tideman, M. (2010) 'Future Testing of Active Safety Systems', *SAE International Journal of Passenger Cars -*

703 *Electronic and Electrical Systems*, 3(2), pp. 170–175. doi: 10.4271/2010-01-2334.

704 Hergeth, S., Lorenz, L. and Krems, J. F. (2017) 'Prior Familiarization With Takeover Requests Affects Drivers' Takeover Performance and

705 Automation Trust', *Human Factors*, 59(3), pp. 457–470. doi: 10.1177/0018720816678714.

706 ISO (2011) *Road vehicles — Functional safety (ISO 26262)*.

707 Jian, J.-Y., Bisantz, A. M. and Drury, C. G. (2000) 'Foundations for an Empirically Determined Scale of Trust in Automated System',

708 *International Journal of Cognitive Ergonomics*, 4(1), pp. 53–71.

709 Johansson, R. and Nilsson, J. (2016) 'The Need for an Environment Perception Block to Address all ASIL Levels Simultaneously', in *Proc.*

710 *of the IEEE Intelligent Vehicles Symposium (IV)*. Gothenburg, Sweden. doi: 10.1109/IVS.2016.7535354.

711 Kennedy, R. S. *et al.* (1993) 'Simulator Sickness Questionnaire: An Enhanced Method for Quantifying Simulator Sickness', *International*

712 *Journal of Aviation Psychology*, 3(3), pp. 203–220.

713 Khashtgir, S. *et al.* (2015) 'Identifying a Gap in Existing Validation Methodologies for Intelligent Automotive Systems: Introducing the 3xD

714 Simulator', in *Proc. of the IEEE Intelligent Vehicles Symposium (IV)*. Seoul, South Korea: IEEE, pp. 648–653. doi:

715 10.1109/IVS.2015.7225758.

716 Khashtgir, S., Birrell, S., Dhadyalla, G. and Jennings, P. (2017) 'Calibrating Trust to Increase the Use of Automated Systems in a Vehicle', in

717 Stanton, N. *et al.* (eds) *Advances in Human Aspects of Transportation. Advances in Intelligent Systems and Computing*. Springer, Cham, pp.

718 535–546. doi: 10.1007/978-3-319-41682-3_45.

719 Khashtgir, S., Sivencrona, H., Dhadyalla, G., Billing, P., *et al.* (2017) 'Introducing ASIL Inspired Dynamic Tactical Safety Decision

720 Framework for Automated Vehicles', in *Proc. of the IEEE 20th International Conference on Intelligent Transportation Systems (ITSC*

721 *2017)*. Yokohama, Japan, pp. 2398–2403.

722 Khashtgir, S., Birrell, S., Dhadyalla, G., Sivencrona, H., *et al.* (2017) 'Towards increased reliability by objectification of Hazard Analysis and

723 Risk Assessment (HARA) of automated automotive systems', *Safety Science*. Elsevier Ltd, 99, pp. 166–177. doi:

724 10.1016/j.ssci.2017.03.024.

725 Van Der Laan, J. D., Heino, A. and De Waard, D. (1997) 'A simple procedure for the assessment of acceptance of advanced transport

726 telematics', *Transportation Research Part C: Emerging Technologies*, 5(1), pp. 1–10. doi: 10.1016/S0968-090X(96)00025-3.

727 Larsson, A. F. L. (2012) 'Driver usage and understanding of adaptive cruise control', *Applied Ergonomics*, 43, pp. 501–506. doi:

728 10.1016/j.apergo.2011.08.005.

729 Larsson, A. F. L., Kircher, K. and Hultgren, J. A. (2014) 'Learning from experience: Familiarity with ACC and responding to a cut-in

730 situation in automated driving', *Transportation Research Part F: Traffic Psychology and Behaviour*, 27, pp. 229–237. doi:

731 10.1016/j.trf.2014.05.008.

732 Lee, J. D. and See, K. A. (2004) 'Trust in Automation: Designing for Appropriate Reliance', *Human factors*, 46(1), pp. 50–80. doi:

733 10.1518/hfes.46.1.50.30392.

734 Lee, J. and Moray, N. (1992) 'Trust, control strategies and allocation of function in human-machine systems', *Ergonomics*, 35(10), pp.

735 1243–1270. doi: 10.1080/00140139208967392.

736 Louw, T. and Merat, N. (2017) 'Are you in the loop? Using gaze dispersion to understand driver visual attention during vehicle automation',

737 *Transportation Research Part C*. Elsevier Ltd, 76, pp. 35–50. doi: 10.1016/j.trc.2017.01.001.

738 Michon, J. A. (1985) 'A critical view of driver behavior models: what do we know, what should we do?', in Evans, L. and Schwing, R. C.

739 (eds) *Human behavior and traffic safety*. Plenum Press, pp. 485–520. doi: 10.1007/978-1-4613-2173-6.

740 Miller, D. *et al.* (2016) 'Behavioral Measurement of Trust in Automation: The Trust Fall', in *Proc. of the Human Factors and Ergonomics*

741 *Society 2016 Annual Meeting*, pp. 1849–1853. doi: 10.1177/1541931213601422.

742 Molesworth, B. R. C. and Koo, T. T. R. (2016) 'The influence of attitude towards individuals??? choice for a remotely piloted commercial

743 flight: A latent class logit approach', *Transportation Research Part C: Emerging Technologies*, 71, pp. 51–62. doi:

744 10.1016/j.trc.2016.06.017.

745 Muir, B. M. (1987) 'Trust between humans and machines, and the design of decision aids', *International Journal of Man-Machine Studies*,

746 27, pp. 527–539. doi: 10.1016/S0020-7373(87)80013-5.

747 Muir, B. M. (1994) 'Trust in automation: Part I. Theoretical issues in the study of trust and human intervention in automated systems',

748 *Ergonomics*, 37(11), pp. 1905–1922. doi: 10.1080/00140139408964957.

749 Muir, B. M. and Moray, N. (1996) 'Trust in automation. Part II. Experimental studies of trust and human intervention in a process control

750 simulation.', *Ergonomics*, 39(3), pp. 429–460. doi: 10.1080/00140139608964474.

751 NHTSA (2017) *Investigation Report: PE 16-007 (MY2014-2016 Tesla Model S and Model X)*.

752 Parasuraman, R. and Miller, C. a. (2004) 'Trust and etiquette in high-criticality automated systems', *Communications of the ACM*, 47(4), pp.

753 51–55. doi: 10.1145/975817.975844.

754 Parasuraman, R. and Riley, V. (1997) 'Humans and Automation: Use, Misuse, Disuse, Abuse', *Human Factors*, 39(2), pp. 230–253.

755 Peng, Y., Boyle, L. N. and Lee, J. D. (2014) 'Reading, typing, and driving: How interactions with in-vehicle systems degrade driving

756 performance', *Transportation Research Part F: Traffic Psychology and Behaviour*, 27, pp. 182–191. doi: 10.1016/j.trf.2014.06.001.

757 Rajaonah, B. *et al.* (2008) 'The role of intervening variables in driver-ACC cooperation', *International Journal of Human Computer Studies*,

758 66(3), pp. 185–197. doi: 10.1016/j.ijhcs.2007.09.002.

759 Rajaonah, B., Anceaux, F. and Vienne, F. (2006) 'Trust and the use of adaptive cruise control: a study of a cut-in situation', *Cognition,*

760 *Technology & Work*, 8(2), pp. 146–155. doi: 10.1007/s10111-006-0030-3.

761 Rasmussen, J. (1983) 'Skills, Rules, and Knowledge; Signals, Signs, and Symbols, and Other Distinctions in Human Performance Models',

762 *IEEE Transactions on Systems, Man, and Cybernetics*, 13(3), pp. 257–266.

763 Rasmussen, J. (1985) 'The Role of Hierarchical Knowledge Representation in Decisionmaking and System Management', *IEEE*

764 *Transactions on Systems, Man, and Cybernetics*, 15(2), pp. 234–243. doi: 10.1109/TSMC.1985.6313353.

765 Rudin-Brown, C. M. and Parker, H. a. (2004) 'Behavioural adaptation to adaptive cruise control (ACC): Implications for preventive

766 strategies', *Transportation Research Part F: Traffic Psychology and Behaviour*, 7(2), pp. 59–76. doi: 10.1016/j.trf.2004.02.001.

767 SAE (2018) *Surface Vehicle Recommended Practice, J3016: Taxonomy and Definitions for Terms Related to Driving Automation Systems*

768 for On-Road Motor Vehicles. doi: 10.4271/2012-01-0107.
769 Seppelt, B. D. and Lee, J. D. (2007) 'Making adaptive cruise control (ACC) limits visible', *International Journal of Human Computer*
770 *Studies*, 65(3), pp. 192–205. doi: 10.1016/j.ijhcs.2006.10.001.
771 Sepulcre, M., Gozalvez, J. and Hernandez, J. (2013) 'Cooperative vehicle-to-vehicle active safety testing under challenging conditions',
772 *Transportation Research Part C: Emerging Technologies*, 26, pp. 233–255. doi: 10.1016/j.trc.2012.10.003.
773 Stanton, N. A. and Young, M. S. (1998) 'Vehicle automation and driving performance', *Ergonomics*, 41(7), pp. 1014–1028. doi:
774 10.1080/001401398186568.
775 Stanton, N. a, Young, M. and Mccaulder, B. (1997) 'Drive-By-Wire : the Case of Driver Workload and Reclaiming Control With Adaptive
776 Cruise Control', *Safety Science*, 27(2), pp. 149–159. doi: 10.1016/S0925-7535(97)00054-4.
777 Talebpour, A. and Mahmassani, H. S. (2016) 'Influence of connected and autonomous vehicles on traffic flow stability and throughput',
778 *Transportation Research Part C: Emerging Technologies*, 71, pp. 143–163. doi: 10.1016/j.trc.2016.07.007.
779 Tingvall, C. (1997) 'The Zero Vision: A Road Transport System Free from Serious Health Losses', *Transportation, Traffic Safety and*
780 *Health: the New Mobility*, pp. 37–57.
781 Le Vine, S. *et al.* (2016) 'Automated cars: Queue discharge at signalized intersections with "Assured-Clear-Distance-Ahead" driving
782 strategies', *Transportation Research Part C: Emerging Technologies*, 62, pp. 35–54. doi: 10.1016/j.trc.2015.11.005.
783 Walker, G. H., Stanton, N. A. and Salmon, P. (2016) 'Trust in Vehicle Technology', *International Journal of Vehicle Design*, 70(2), pp.
784 157–182. doi: 10.1504/IJVD.2016.074419.
785 Wang, J. *et al.* (2016) 'Driving safety field theory modeling and its application in pre-collision warning system', *Transportation Research*
786 *Part C: Emerging Technologies*, 72, pp. 306–324. doi: 10.1016/j.trc.2016.10.003.
787 WMG (2017) *Drive-in, Driver-in-the-loop, multi-axis driving simulator (3xD)*. Available at:
788 <http://www2.warwick.ac.uk/fac/sci/wmg/research/naic/facilities/> (Accessed: 10 July 2017).
789 Xu, J. *et al.* (2014) 'How different types of users develop trust in technology: A qualitative analysis of the antecedents of active and passive
790 user trust in a shared technology', *Applied Ergonomics*, 45(6), pp. 1495–1503. doi: 10.1016/j.apergo.2014.04.012.
791 Yu, H., Lin, C.-W. and Kim, B. (2016) 'Automotive Software Certification: Current Status and Challenges', *SAE International Journal of*
792 *Passenger Cars - Electronic and Electrical Systems*, 9(1), pp. 74–80. doi: 10.4271/2016-01-0050.
793