

This is a repository copy of *A safety-case approach to the ethics of autonomous vehicles*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/152973/>

Version: Accepted Version

Article:

Menon, Catherine and Alexander, Robert David orcid.org/0000-0003-3818-0310 (2019) A safety-case approach to the ethics of autonomous vehicles. *Safety and Reliability*. pp. 33-58. ISSN 0961-7353

<https://doi.org/10.1080/09617353.2019.1697918>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

A safety-case approach to the ethics of autonomous vehicles

Catherine Menon^{a*} and Rob Alexander^b

^aSchool of Computer Science, University of Hertfordshire, Hatfield, UK; ^bDepartment of Computer Science, University of York, York, UK

a* c.menon@herts.ac.uk

b rob.alexander@york.ac.uk

Catherine Menon is a senior lecturer in the Adaptive Systems group at the University of Hertfordshire.

Rob Alexander is a lecturer and researcher in high-integrity systems engineering at the University of York.

A safety-case approach to the ethics of autonomous vehicles

Autonomous Vehicles (AVs) have significant ethical and safety implications. Questions of informed consent and risk acceptance are of primary importance, as is an explicit identification of the ethical principles underlying these decisions. In this paper we present a process framework for producing an ethics assurance case, which can be used to translate ethical imperatives into design decisions and safety management practices. The process and resultant assurance case integrate ethical considerations into the wider engineering lifecycle, providing a tool to demonstrate that design and safety management decisions reflect an identified ethical position.

Keywords: safety; ethics; autonomous vehicles; risk

1. Introduction

Autonomous Vehicles (AVs) are increasingly being presented as the future of transport on public roads. The worth of the Connected and Autonomous Vehicle (CAV) market in the UK by 2035 is estimated to be 28bn (Transport Systems Catapult [TSC], 2017), and typical estimates for the first occurrence date of fully-autonomous vehicles on UK roads are in the mid-2020s (UK Government Department for Transport [DfT], 2015). These forecasts confirm the continuing trend towards increased autonomy in this domain (Anderson & Anderson, 2007). Real-world trials of AVs are also underway in several countries, including the UK (UK Government Centre for Connected and Autonomous Vehicles [CCAV], 2018; UK Autodrive, 2017; Venturer, 2016; Venturer, 2017), the US (Waymo, 2018; Uber, 2018), Singapore, Japan and Europe (CCAV, 2018).

Although the technological capability to develop AV systems is developing quickly, ethical considerations remain an important societal barrier to their acceptance (UK Autodrive, 2017). This is exemplified by the trolley problem (Foot, 1967), which

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

in its original form presents an ethical dilemma in which a runaway train is on course to kill five people. A bystander is given the choice to let the train continue on its course or to divert it onto a different track where only one person will be killed. Although the trolley problem is a hypothetical question only, there is mainstream public interest in how AVs should behave in similar situations where harm to at least one person is inevitable (MIT, 2018).

It is arguable that the trolley problem has dominated the ethical debate around AV introduction to the exclusion of other, more nuanced, ethical considerations (Goodall, 2016). These other considerations include the balance and types of risk considered acceptable for AVs to present, the distribution of these risks across different sections of society (e.g. as differentiated by factors including age, economic position, gender) and the different appetites for risk over these same sections (Wang & Zhao, 2019), concerns over developer culpability and liability (Anderson, Nidhi, Stanley, Sorenson & Oluwatola, 2014), the environmental implications of AV introduction (Fagnant & Kockelman, 2014) and the impact on transport efficiency.

In this paper we build on existing work in risk distribution (Menon & Alexander, 2017; Menon, Bloomfield & Clements, 2013; Menon & Alexander, 2018) to present a process framework for producing an ethics assurance case. This assurance case can be used to translate nuanced ethical considerations into safety and design requirements, and to demonstrate that identified ethical principles – and the safety requirements stemming from these – are satisfied. Structured assurance cases are well established in the area of safety argumentation (Bloomfield & Bishop, 2010; Kelly & Weaver, 2004; Hawkins, Habli, Kelly & McDermid, 2013) where they are used to provide a rigorous justification of the safety of a system. Our proposed approach embeds an assurance case template within a process structure, thereby creating a

1 framework for explicitly describing ethical principles within the engineering lifecycle.

2 The framework ensures that ethical principles are described, justified and implemented
3 within the AV design.
4

5
6 Our framework does not prescribe a set of recommended ethical principles for
7 AVs, but rather focuses on the *translation* of ethical principles into safety requirements.
8
9 As such, the framework is applicable to a range of ethical, risk perception and decision-
10 making theories, some of which are discussed in Section 2. The framework does assume
11 compliance with the UK Human Rights Act (UK Government, 1998) and the other
12 national implementations of the European Convention on Human Rights. We restrict
13 our discussion throughout this paper on AVs which perform the whole of the driving
14 task. These correspond to SAE Level 4 or Level 5, using the levels of automation
15 identified and discussed in (SAE International, 2018).
16
17
18
19
20
21
22
23
24
25
26
27
28

29 In Section 2 we present the ethical and safety background relevant to AVs,
30 highlighting where existing work does not fully facilitate transformation of ethical
31 concerns into safety requirements. Section 3 describes how these ethical concerns
32 impact our judgements around risk acceptance and risk balancing. Section 4 describes
33 how ethical principles may govern how risks are balanced, or traded off, against each
34 other and introduces the concept of risk profiles to describe these trade-offs. Section 5
35 presents an ethics assurance case template pattern and a process framework for
36 instantiating this template to produce an ethics assurance case, thereby enabling the
37 translation of ethical principles into specific safety requirements. Section 6 concludes
38 with a discussion and steps for further work.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

2. Ethical and Safety Background Relevant to AVs

2.1 Ethics and AVs

While there is currently no regulatory or legislative barrier to the testing and deployment of AVs in the UK (DfT, 2015), there are significant ethical and safety challenges. AVs have the potential to cause harm in all the ways that traditional cars do (e.g. collisions, harmful emissions, impacts on road efficiency), as well as via novel pathways. For example, cyber-security is a significant issue (DfT, 2017; UK Autodrive, 2018) as AVs are vulnerable to being controlled by malicious third parties in a way that traditional vehicles are not.

This potential for AVs to cause harm means that there are significant ethical challenges connected with their operation. Public perception of AV safety must be considered, alongside the demonstrated safety of such systems and the adequacy of mechanisms in place to reduce the risk posed by AVs. Such discussions must also consider the selection of risk criteria, the acceptability of residual risk associated with the AV and the extent to which users have consented to bear this risk. We note that these are also ethical issues for the existing car fleet, in that drivers, pedestrians and other road users have not explicitly consented to the risk as currently posed by human drivers. It is arguable that a similar ethical argument should be made for the existing fleet (as an alternative to the introduction of AVs). Such an argument is beyond the scope of this paper, but we note it as a prevailing ethical challenge.

2.2 Risk Perception Theories

The study of risk perception and decision making is a wide field, drawing from cross-disciplinary domains including psychology, linguistics and engineering. In this section we present a number of risk perception theories which are particularly relevant

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

to the societal perception of AVs and their risk acceptability.

Expected utility theory (von Neumann, 1947) provides a model of risk perception which allows for an individual's risk appetite and risk aversion. These properties are crucial when considering the ethics of AV use, as an individual's perception of the novel risks presented by an AV will be affected by their risk appetite. As stated earlier, the current (human-driven) fleet also presents risks which have not necessarily been explicitly consented to, and the framing of the new AV risks as compared to the existing risks can have a significant effect on people's perception of these.

Both expected utility theory and prospect theory (Kahneman & Tversky, 1979) consider the framing of risks, but prospect theory allows for greater complexity around risk perception. In particular, it provides a method of addressing the Allais paradox (Allais, 1953), and allows us to model risk perceptions around specific AV scenarios more accurately. (Wang & Zhao, 2019) use prospect theory to discuss AV risk perception in more detail, noting that the demographics which have a higher appetite for risk (e.g. young, fully-employed, male, high-income) view new technology more favourably, while the risk-averse (poor, elderly) perceive AVs as presenting a risk which – when considered in the context of perceived benefits – they are not prepared to accept.

Regret theory (Bell, 1983) provides a model of risk perception and decision making which captures the effect of regret, compared to an alternative outcome. Regret theory allows for the consideration of regret-avoidance, together with risk aversion and risk appetite. (Somasundaram & Diecidue, 2015) use regret theory to consider risk attitudes in society, finding that feedback polarizes regret attitudes, and increases risk-seeking amongst those who are regret-averse. As before, an increase in risk-seeking can affect the perceived risks posed by AVs. It presents an ethical challenge, particularly as AVs pose a risk across

1 multiple segments of society, and there is the potential for an increase in risk-seeking by
2 one segment to result in a consequently greater risk being posed to another.
3

4 Dual Process Theory (Baron, 1985; Evans, 1989) allows for the interaction between
5 decisions made for fast, intuitive reasons and slower decisions based on reasoning. As with
6 prospect theory, dual process theory emphasises the importance of attribute framing
7 (Cokey, 2009; Reyna, 2004), which is particularly important for AVs given the complexity
8 of the risks they present.
9

10 11 12 13 14 15 16 17 18 *2.2.1 Risk perception of AVs*

19 Existing work on ethical design standards, such as (IEEE Global, 2018) utilises
20 risk perception theories to examine high-level ethical concerns around AVs, as well as
21 providing an overview of the societal benefits and concerns around autonomous systems
22 in general.
23
24
25
26
27

28 Other work, such as (UK Autodrive, 2017b), focuses on the different ethical
29 factors with the potential to influence the eventual behaviour of an AV. These ethical
30 factors can be drawn from the ethical theories discussed above and include the “human
31 values”, such as the AV developers’ desire for fairness and the AV passengers’ desire
32 for personal autonomy (Thornton, 2018). In practical, engineering terms the first of
33 these could lead to developers preferring algorithms which prioritise polite and non-
34 aggressive behaviour of the AV, while the second could lead to implementation of
35 customised AV behaviours which allow passengers to choose the preferred style of
36 driving (Kuderer, Gulati & Burgard, 2015). This has ethical implications in itself, in that
37 the developers may choose to limit the choices of driving style to those which are non-
38 aggressive, thereby depriving the passenger of some personal autonomy.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

55 In addition to these “human values”, there are less altruistic factors which may
56 also influence the eventual behaviour of an AV. Self-interest and commercial
57
58
59
60
61
62
63
64
65

1 competitiveness are likely to be relevant, given the results presented in (UK Autodrive,
2 2017b), which identify that potential customers prefer the AV to prioritise the safety of
3 its passengers over any third parties. More generally, self-interest could lead to
4 developers making choices about the AV behaviour specifically to reduce their
5 culpability in the case of an accident (Shalev-Schwartz, Shammah & Shashua, 2017).
6
7 The concerns around how risks are framed is particularly relevant here, and several of
8 the relevant ethical theories – including prospect theory and regret theory – identify
9 framing as an important factor in willingness to accept a given risk.
10
11

12 In terms of the behaviours which result from ethical choices, (Gips, 1995;
13 Wallach & Allen, 2008) explore how the design of an autonomous system is affected by
14 the extent to which ethical reasoning and capacity is embedded within it. (Dennis,
15 Fisher, Slavkoviv & Webster, 2004) present formal verification that a high-level ethical
16 policy is satisfied by the eventual behaviour of the system. A general architecture for a
17 robot capable of modelling its own actions using simulation and predicting the ethical
18 consequences of these is discussed in (Winfield, Blum & Liu, 2014; Vanderelst &
19 Winfield, 2018) while (Arkin, Ulam & Wagner, 2012) considers simulation of a robot
20 with an ethical framework.
21
22

23 *2.2.2 Revisiting the trolley problem*

24 Much of the on-going public discussion of AV ethics focuses on the trolley problem,
25 which posits a situation in which an AV must choose which of two pedestrians with
26 which to collide. The trolley problem is often cited in public media as an illustration of
27 AV safety issues, and public debate is typically focused around variants of this (MIT,
28 2018). However, real-world instances of the trolley problem are rare, and much of the
29 public discussion around AVs and the trolley problem assumes a level of engineering
30 capability that is infeasible (UK Autodrive, 2017b; Goodall, 2016).
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

A more realistic variant of the trolley problem considers how the AV can act in any given situation to minimise the overall risk (UK Autodrive, 2017b). This requires the AV to accurately estimate its own operational capacity and to adjust its behaviour accordingly (Nilsson, 2018). Similarly, other work (Lin, 2015) considers the example of an AV driving closer (within its lane) to a smaller car on its left than to a truck on its right. This choice reduces the risk to the AV, as a collision with a small car is safer for the AV occupants than a collision with a truck. Another AV chooses differently, driving closer to the heavier vehicle with more effective safety systems (Lin, 2015). This second AV is optimising its driving position to reduce the overall risk it poses to other road users, as if it collides with the truck this is less likely to result in injuries than a collision with the car would be (Gerdes & Thornton, 2016). Other proposed situations include an AV choosing a “sacrificial” path, such as placing itself to block the trajectory of a runaway vehicle (Lin, 2015).

3. Translating ethics into design

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Considered as a body, the works presented in Section 2 provide two important foundational results: how to specify AV behaviours in individual scenarios, and how to assess AV behaviours resulting from given ethical principles. However, they do not provide a general mechanism for specifying the ethical principles, documenting the translation of these into AV behaviours, or justifying that ethically-motivated behaviours are demonstrated and will continue to be so.

Formal verification, as shown in (Dennis et. al., 2004; Winfield et. al., 2014; Vanderelst & Winfield, 2018), is a valuable contribution in this area. However, formal verification of the entirety of a complex safety critical system – such as an AV – has historically been considered infeasible due to cost, technical limitations, and perceived difficulty (Liu, Stavridou, & Duarte, 1995; Knight, 2002; Yoo, Jee & Cha, 2009)).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Moreover, such verification is also intended to demonstrate correctness according to the specification only, and therefore does not consider wider issues of overall risk reduction, regulatory compliance or errors due to inadequate requirements elicitation and environmental changes. Consideration of these issues is a legal requirement for safety-critical systems (Health and Safety Executive [HSE], 2001).

In the more nuanced trolley problems discussed in Section 2.2.2, the behaviour of the AV is motivated by the intent to reduce risk, whether this be to its own occupants or to other drivers on the road. Such a motivation, of course, is more properly ascribed to the AV developers rather than to the AV itself. This distinction highlights two interrelated areas of application when considering the ethics of AVs: *implemented ethics* (the ethics embedded within an autonomous system and realised in its behaviour) and *engineering ethics* (the ethical principles and codes of practice followed by engineers). Making this distinction allows us to interrogate the ethical behaviour of the AV without necessarily considering the professional conduct of the developers, and vice versa.

3.1 *Engineering ethics and implemented ethics*

Engineering ethics refers to the professional ethical principles which are followed by the developers of the AV during development work. These principles may be represented by professional codes of conduct (Royal Academy of Engineering [RAEng], 2017) as well as more general informal undertakings (Martin & Schinzinger, 2005).

Such ethical principles typically include criteria such as honesty, integrity, respect for law and the public interest, accuracy, rigour, fairness and objectivity (RAEng, 2017). However, they do not in themselves constrain the behaviour of any resulting system on ethical lines. It is, however, plausible that following a code of engineering ethics should prevent the developers from knowingly designing a system

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

that contravenes established ethical foundations (e.g. the Human Rights Act (UK Government, 1998)).

Implemented ethics refers to the ethics embedded in the behaviour of the AV itself, sometimes referred to as the “moral algorithm” (UK Autodrive, 2017b) or the “machine ethics” of the AV. The implemented ethics determine how multiple risks are balanced against each other, and how safety risk is balanced against considerations of security, privacy, trust and capability. Different societies and stakeholders will differ in their criteria for what behaviour is considered ethically acceptable, and this will also vary across environments and domains of use.

3.2 AVs and risk reduction

Arguments have been put forward (Kalra & Groves, 2017) that AVs should be introduced as soon as their safety record is slightly better than traditional vehicles. Such studies estimate 500,000 fewer overall road fatalities over a fifty-year time frame compared to a conservative policy of AV introduction.

However, looking only at overall road fatalities obscures the distribution of such fatalities, which is crucial from both an ethical and safety perspective. The introduction of AVs may transfer risks even while reducing overall risk. That is, AVs may change how different classes of people (e.g. human drivers, passengers, pedestrians) are differentially exposed to risks. Any risk transfer also raises the question of risk consent, and whether all affected parties have agreed to the redistribution of risk. (As we discussed earlier, such explicit consent has not in fact been sought for the current fleet. People’s decision to accept or reject the current fleet risk is largely intuitive, corresponding to many of the factors discussed in Section 2.2. The current method of formally adjudicating whether human intuition and judgement represent sufficient risk

1 mitigation given a particular accident scenario appeals to legal processes rather than
2 ethics).

3
4 A complicating factor here is that informed consent requires an accurate
5 perception of the risk posed by the AV. The decision-making theories of Section 2.2
6 identify a number of different motivating factors for decision-making under risk, but
7 most compelling for AVs is that the perception of risk posed by the AV vs that posed by
8 a human (driver) is not uniform across society (Kim & McGill, 2011; Wang & Zhao,
9 2019). A Singaporean study (Wang & Zhao, 2019) has shown that those who are more
10 likely to benefit from AVs (high-income, urban-dwelling) are more willing to accept the
11 risks, while a more general study of autonomous systems (Kim & McGill, 2011) shows
12 that low-income segments of society are more likely to trust machines rather than
13 interrogate the engineering and ethical principles behind these.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 Another complicating factor is the potential for AVs to contribute to risk
30 indirectly, as participants in the wider road network. For example, AVs which mis-
31 identify emergency vehicles may not pull over for these vehicles, thereby indirectly
32 causing harm. Similarly, AVs which learn to behave selfishly (e.g. failing to give way,
33 aggressive driving) may also cause traffic congestion across the wider road network
34
35
36
37
38
39
40

41 The introduction of AVs which perform the majority of the driving task also
42 transfers ownership of much of the risk associated with this task. Currently the driver
43 owns much of this risk, system failure notwithstanding. However, for AVs at SAE
44 levels 4 and 5, it is possible that the developer will own nearly the entire risk associated
45 with the minute-by-minute driving decisions within the operational design domain.
46
47
48
49
50
51
52
53 Although the human passenger may still provide input into route choices, customisation
54 of driving techniques and overall usage, this would mean that AVs would be associated
55 with a significant ethical responsibility borne by an individual or entity – the
56
57
58
59
60
61
62
63
64
65

1 manufacturer and developer – not personally exposed to the possible outcomes of the
2 resultant risks.
3
4
5

6 **4 Ethics and risk balancing**

7

8
9 As discussed in Section 3.2, one of the fundamental issues around ethics and safety of
10 AVs is the question of risk transfer, or risk balancing. It is the redistribution of risks
11 consequent on introduction of AVs which throws up the most complex ethical
12 challenges.
13
14
15
16
17

18
19 There is a legal requirement in the UK for the overall risk associated with a
20 system to be reduced As Low As Reasonably Practicable (ALARP). The Health and
21 Safety Executive provides guidance (HSE, 2001) for good practice in reducing risk
22 ALARP and for demonstrating this.
23
24
25
26
27

28
29 For a minority of systems, the system risk can be reduced ALARP by mitigating
30 the risk from each individual hazard until it is ALARP. However, in cases where risk
31 has been transferred and redistributed (as with AVs), the situation is more complex. It is
32 possible to reduce a single risk at the cost of introducing another, or introduce a risk
33 mitigation which affects multiple risks at once. This means that it is possible for
34 multiple different system designs to all be ALARP, but for each to provide a different
35 balance amongst the individual system risks (Menon et. al., 2013). This can occur under
36 the following circumstances:
37
38
39
40
41
42
43
44
45
46
47

- 48 • When developers have not identified a complete list of hazards. This is relatively
49 likely during the initial deployment of AVs, as safety engineers will not have
50 complete and valid tools for identifying and understanding hazards due to the
51 lack of established good practice and historical data.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- When there are interdependencies between hazards which are not adequately accounted for. In these situations, some risks may be accounted for twice, giving a false idea of the overall risk associated with the system. There may also be interdependencies between systems and common mode failures which increase the complexity of combining risks. Given the relative novelty of AV technology and risk management, this is likely to be a significant issue.
- When multiple risk mitigations are not independent, and a mitigation for one risk potentially increases other risks. For example, increasing the sensitivity of algorithms which detect an object mitigates against the risk of failure to detect, but increases the risk of erratic driving (the AV may brake unnecessarily to avoid a “phantom” object) and therefore the likelihood of a collision.
- When a single risk mitigation affects multiple hazards. In this case, the cost of the mitigation can be amortized over all the hazards and the resulting cost judged reasonably practicable, where it would not when assessed against each hazard individually.
- There are limited resources subject to a threshold effect of aggregation. For example, in a SAE level 3 vehicle (i.e. one with supervising driver), operator attention may be a mitigation against several hazards. When these hazards present themselves simultaneously, the operator may be overwhelmed.

In circumstances where multiple different designs all present an overall ALARP system risk, selecting any one of these designs represents a risk distribution choice. That is, each design provides a different balance amongst the individual system risks, “trading off” an increase in one risk for a decrease in another. This is an established practice in the nuclear domain, with standards such as (HSE, 2006; Office for Nuclear Regulation [ONR], 2018) emphasising the need to balance individual risks within a system.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

However, outside the nuclear domain many safety guidance documents provide little information on risk balancing and risk transfer, or require an explicit justification of any risk trade-offs.

4.1 Ethical motivation for risk balancing

Where multiple different system designs all offer an overall ALARP risk, the final choice may be affected by a number of factors including cost, technical capability, resource availability and ethical imperatives. The cheapest design may be chosen, or the design which can be most easily implemented with the resources and technology at hand. In these cases the justification is relatively simple and has long been part of standard project management techniques (Office of Government Commerce, 2019).

However, the ethical principles behind selecting a design are not generally explicitly discussed. Historically, there has been no mechanism to record those ethical imperatives which result in developers considering one risk to be more important than another. AVs in particular are vulnerable to the impact of such “hidden” ethical priorities of the designers. This is exemplified in the trolley problem, in its simplified form. The AV designer will use his / her ethical principles to decide whether the AV – guided by its programming – should impact one person or the other. Each of these choices will result in different system designs (because of the difference in behaviour), which present different risks to different people. However, the same overall system risk will remain the same whichever person is impacted.

Because the ethical complexities around AV operation and use have the potential to result in different distributions of risk, we consider that there is an obligation to provide information to the general public about the ethically-motivated risk trade-offs that have made during development. Furthermore, where ethical considerations have led designers to consider one risk acceptable (while deeming

1 another equivalent risk not acceptable) and to have acted on this to inform the system
2 design, this practice must be communicated to the general public. For example, where a
3 developer has chosen a system design for the AV which prioritises the safety of
4 passengers over pedestrians, this decision should be made explicit to both pedestrians
5 and passengers. Visibility of this information is necessary if affected stakeholders are to
6 provide informed consent to the risks that they bear as a result. Such visibility is, of
7 course, lacking in the current fleet (pedestrians are unaware of the degree of self-interest
8 of drivers), with current decisions around risk acceptance being made according to
9 intuitive or personal ethical theories where in many cases insufficient information is
10 provided for a full understanding.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

26 ***4.2 Refinement of risk factors***

27
28 When deciding which of two equivalent risks is preferable, or when performing
29 risk trade-offs, the AV developer must consider more than the magnitude of the risks.
30 There are other factors which may make one risk preferable (to the developer) than
31 another risk of equivalent magnitude, in terms of safety. It is important to note that this
32 does not mean that the same risk would necessarily be also judged preferable by an AV
33 user, by a pedestrian or even by another AV developer. Risk perception varies between
34 individuals and between segments of society, and is informed by risk appetite and
35 underlying ethical principles.
36
37
38
39
40
41
42
43
44
45
46

47 Risk is typically calculated to be dependent on both the outcome and the
48 likelihood of this outcome
49
50

$$51 \quad \textbf{Risk} = \textbf{harm} * \textbf{likelihood}$$

52
53
54 If an AV developer wants to differentiate between two risks which may be of
55 equivalent magnitude, and further to trade one of these risks off against another, it is
56 necessary to have greater transparency into this calculation. In particular, the developer
57
58
59
60
61
62
63
64
65

1 must be able to provide sufficient information to justify why these risks are different,
2 and to demonstrate how his or her underlying ethical imperatives have informed trade-
3 offs between different risks.
4
5

6
7 With this in mind, we propose the following factors that characterise risks
8
9 beyond their magnitude (severity x likelihood). We treat the factors below differently to
10 how we treat severity and likelihood, in that we do not impose a ranking on them, or
11 hold an implicit brief that some values of these factors are “worse” than others in terms
12 of ethics. The factors are informational only, and serve to differentiate one risk from
13 another of equivalent magnitude. These factors are:
14
15
16
17
18
19
20

21 **a) *The exposed population***

22
23 We propose identifying several broad categories for the population who might
24 be exposed to any given risk. These will vary depending on the risk in question. For
25 example, the population exposed to the risk of collision includes the AV passengers and
26 other nearby road users, while the population exposed to the risk of emissions includes
27 pedestrians and those living in nearby houses. Categories of stakeholders should include
28 (not all will be relevant for all risks):
29
30
31
32
33
34
35
36
37
38
39

- 40 • The AV passenger
- 41
- 42 • Other road users (e.g. other drivers, motorcyclists etc.)
- 43
- 44 • Pedestrians and cyclists
- 45
- 46 • Those living near the road
- 47
- 48
- 49

50
51 **b) *Effective risk horizon***

52
53 This refers to the physical and temporal area in which the harm from a given risk
54 is experienced. In physical terms, the harm may manifest close to the AV or further
55 away, while in temporal terms the harm may manifest straight away or several years
56
57
58
59
60
61
62
63
64
65

1 later. We emphasise that a change in effective risk horizon does not necessarily affect
2 the magnitude of the risk: that is, a “close” risk is not necessarily worse than a “distant”
3 risk. Nevertheless, an AV developer may look to differentiate two equivalent risks by
4 their effective risk horizons. Different risk horizons may cause the quality of the risks to
5 be perceived differently, and may affect the way in which their associated economic,
6 reputational or emotional effect is viewed. These will all have implications for the
7 ethical imperatives of the AV developer. Some examples of risks which differ in terms
8 of effective risk horizon (as well as potentially in magnitude) are:
9

- 10 • A risk of collision can result in geographically local harm (a collision outcome
11 harms pedestrians only in the immediate area)
- 12 • A risk of the AV causing road congestion can result in geographically wider
13 harm (the traffic jam delays an ambulance several streets away, causing harm to
14 the patient)
- 15 • A risk of collision can result in harm that is close in time (the harmful effects of
16 a crash are typically experienced immediately)
- 17 • A risk of cancer from carcinogenic AV fuel can result in harm which is distant
18 in time (cancer may take years to develop)

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43 **c) Causes**

44 Risk cause is also relevant to AV developers seeking to differentiate between
45 equivalent risks. Although, as with the above factors, the cause of a risk does not
46 necessarily affect the magnitude of that risk, there may be ethical, emotional and
47 intuitive reasons why an AV developer might prioritise the mitigation of risks from
48 particular causes. For example, Section 2.2 identifies regret avoidance as a
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

consideration in some risk decisions; a developer may choose to reduce risks from a cause which is also likely to engender regret.

4.2.1 Risk profiles and risk balancing

Previous work (Menon et. al., 2013) introduces the concept of risk profiles, each of which represents an approach to balancing and trading off risks. The risk profiles defined in previous work focus primarily on the magnitude of risk, and the justification for accepting an increase in one risk in return for a decrease in another. Two fundamental risk profiles examined in (Menon et. al., 2013) are the *fairness in improvement* profile (where an attempt is made to reduce all risks by the same magnitude) and the *fairness in outcome* profile (where an attempt is made to reduce all risks to the same magnitude). The authors also identify a number of confounding factors which can limit the practical application of these risk profiles, including the challenge of fully characterising two equivalent – but different – risks of the same magnitude.

The risk refinement factors identified above represent an approach that can be developed to address this challenge. Characterising a risk by its exposed population, effective risk horizon and cause allows more complex ethical arguments to be made about risk prioritisation and permits, for example, the ethical arguments that might be used in accepting a brief short-term increase in risk for a longer-term decrease in risk.

Constructing representative risk profiles from the risk refinement factors is beyond the current scope of this paper, and hence we restrict ourselves to noting the utility of such risk profiles in explicitly identifying and justifying the ethical imperatives behind risk trade-offs.

5. A process framework and ethics assurance case pattern

In this section we present a framework which enables the translation of ethical

1 principles into specific safety / design requirements, and also provides a tool for
2 demonstrating the satisfaction of these requirements. The framework is in two parts: an
3 ethics assurance case template pattern and a process framework for instantiating this
4 template to produce an ethics assurance case. The process for instantiating the ethics
5 assurance case template pattern uses the risk balancing and trade-offs (risk profiles) of
6 Section 4, and must be embedded within a wider engineering lifecycle.
7
8
9
10
11
12
13
14

15 We first provide some background on structured assurance cases, and their
16 history of use. We then build on this to define the template pattern and the process
17 framework, including a diagram and textual discussion of each.
18
19
20
21
22
23

24 **5.1 Structured assurance cases**

25
26 Within safety engineering, *structured assurance cases* are used to present a compelling,
27 credible argument that a system is safe in a given context (Bloomfield & Bishop, 2010;
28 Kelly & Weaver, 2004; Kelly, 2007; Object Management Group [OMG], 2019). An
29 assurance case consists of a set of claims about the system, such as a claim that all
30 hazards have been identified, or that the failure rate of the system is below a certain
31 threshold. These claims are supported with evidence, and with an argument that the
32 evidence is sufficient to provide confidence in that claim. Claims can be broken down
33 into sub-claims, and typically several pieces of evidence are needed to provide
34 confidence that a claim has been satisfied.
35
36
37
38
39
40
41
42
43
44
45
46
47
48

49 Assurance cases allow safety management decisions to be scrutinized (e.g. by a
50 regulator) and defended. The adequacy of the argument is critical, and assurance cases
51 typically build on a set of principles which must be satisfied within the overall argument
52 construction. These principles define an assurance case *template pattern*, which is a
53 recommended method of constructing the argument in order to minimise the chances of
54
55
56
57
58
59
60
61
62
63
64
65

1 introducing a logical, technical or semantic error. An assurance case template pattern
2 therefore consists of a template argument built upon these principles. In practice, the
3 template argument might contain a number of claims in order to support each principle,
4 and each claim might also support multiple principles.
5
6
7

8
9 For example, safety standards within the defence domain (UK Ministry of
10 Defence [MOD], 2017) require that arguments in assurance cases demonstrate
11 satisfaction of the following four principles. Consequently, argument template patterns
12 for assurance cases to be used in the defence environment typically contain claim(s)
13 corresponding to each principle:
14
15
16
17
18
19
20

- 21 1. Safety requirements are defined to address the system's contribution to hazards
- 22 2. The intent of the safety requirements shall be maintained throughout
- 23 requirements decomposition
- 24 3. The safety requirements shall be satisfied
- 25 4. Hazardous behaviour of the system shall be identified and mitigated
- 26
27
28
29
30
31
32
33
34
35

36 Argument template patterns must be instantiated for use with a particular
37 system. It is at this stage that specific requirements, specifications, V&V artefacts and
38 other items of evidence are included to expand and support the claims.
39
40
41
42

43 A significant body of work already exists around the construction and
44 instantiation of assurance case template patterns, with common questions being focused
45 on how to create template patterns which eliminate logical fallacies, unjustified
46 assumptions, weakened conclusions or other similar faults (Bloomfield & Bishop, 2010;
47 Kelly, 2007; Hawkins, Habli, Kolovos, Paige & Kelly, 2015; Common Criteria, 2007).
48
49
50
51
52
53
54
55
56

57 ***5.2. Principles for ethics assurance case template pattern***

58

59 In order to construct an ethics assurance case template pattern, we build on the four key
60
61
62
63
64
65

1 principles for assurance cases in the defence domain (MOD, 2017). We define four
2 extended principles for the ethics domain as follows:
3
4

5 **P1. Ethics requirements appropriate for AV development and operation shall**
6 **be defined.**
7

8 This requires that engineering ethics and implemented ethics
9 requirements should be explicitly defined, free from inconsistencies, and
10 containing sufficient detail to allow the other principles to be met.
11
12

13 **P2. The intent of the ethics requirements shall be maintained throughout**
14 **decomposition.**
15

16 This requires that the implemented ethics should be propagated
17 throughout the design of the system and refined into lower-level
18 requirements on design, implementation and risk management. The
19 engineering ethics should be satisfied throughout the system lifecycle.
20
21

22 **P3. Ethics requirements shall be satisfied.**
23

24 This requires that the ethics requirements, both implemented and
25 engineering, should be demonstrably satisfied and evidence provided to
26 support this.
27
28

29 **P4. The AV shall continue to be safe, and emergent behaviour of the AV which**
30 **conflicts with the ethics requirements shall be identified and mitigated**
31

32 This constrains emergent behaviour of the AV, either due to changes in
33 the environment or to adaptive algorithms used within the AV software.
34
35 Such emergent behaviours may not have been considered when
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 specifying the original ethics and safety requirements, and this principle
2 requires that evidence be provided to assure the continued safety of the
3
4 AV even in a changing environment.
5
6
7

8 There is a further principle relating to confidence in (MOD, 2015), which has no
9
10 immediate analogue to ethics, and which we do not develop further.
11
12

13 ***5.3. Process and assurance case framework***

14 ***5.3.1 Ethics assurance case template***

15
16
17
18
19 Figure 1 shows the ethics assurance case template pattern using Goal Structuring
20
21 Notation (Assurance Case Working Group, 2018). We discuss this textually below.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

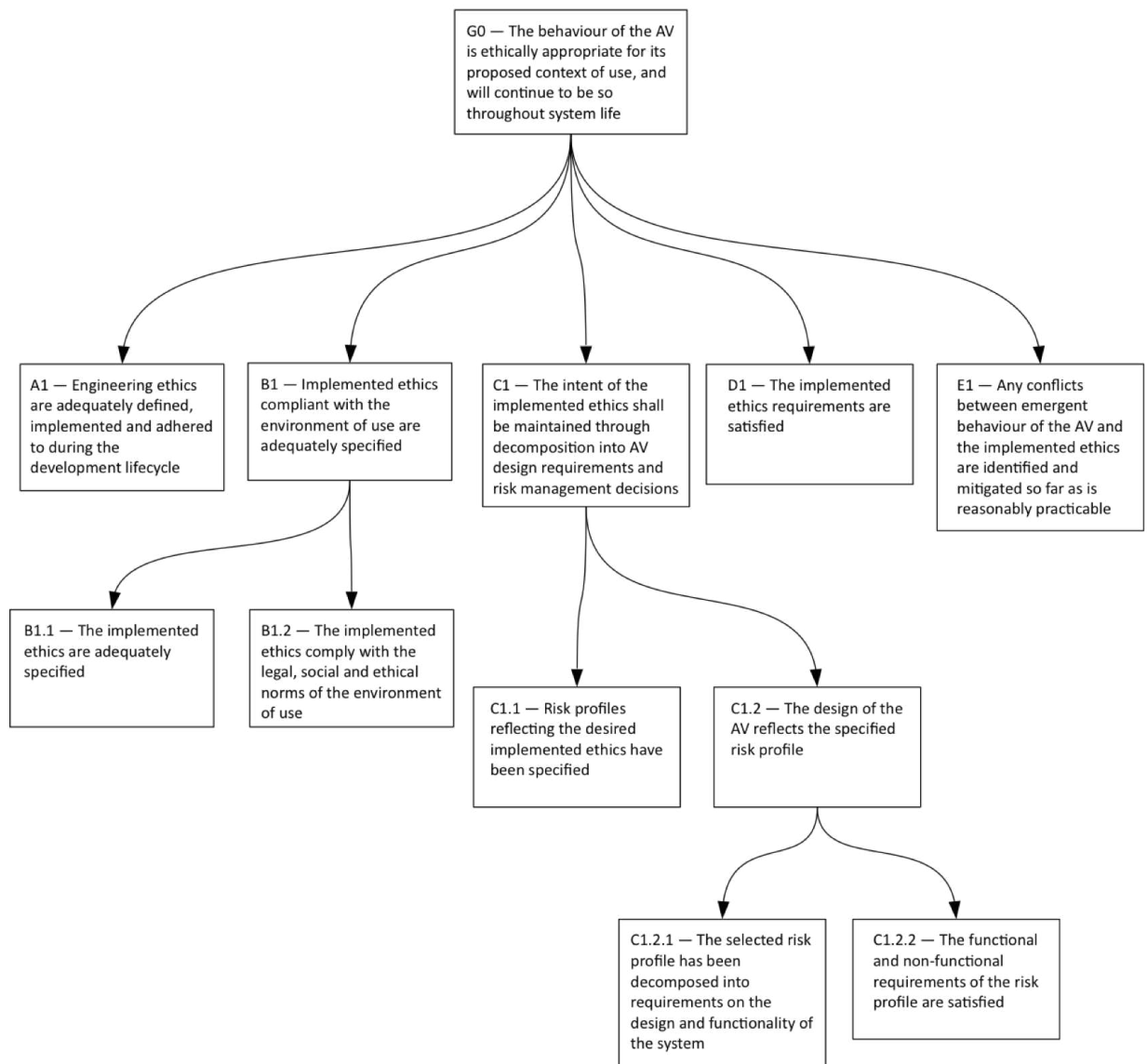


Figure 1: ethical assurance template argument pattern

Claim A1: *Engineering ethics are adequately defined, implemented and adhered to during the system lifecycle*

This claim partially supports P3 and requires AV developers to demonstrate compliance with an appropriate code of conduct, domain-specific good practice and existing ethical precedents. This claim might reasonably be supported with records from audits and artefacts from the development lifecycle, as well as documentation that developers have

1 followed processes defined in accordance with codes of professional conduct. This
2 claim helps to provide confidence in the integrity of any lifecycle artefacts which are
3
4
5 needed to support the top-level claim of Figure 1.
6

7 **Claim B1:** *Implemented ethics compliant with the environment of use are adequately*
8
9 *specified.*
10

11 This claim supports P1, and is broken down into two sub-claims as follows:
12

13 **Claim B1.1:** *The implemented ethics are adequately specified*
14

15
16 Specification of the implemented ethics may be achieved via identification and citation
17
18 of relevant items of regulation and policy, as well as the results of any public
19
20 consultation or ethical objections already tabled. It may also be useful to specify aspects
21
22 of implemented ethics via references to previous system designs, to academic papers,
23
24 and to accepted good practice. The specification of the implemented ethics must be
25
26 sufficient to address competing ethical motivations, and its adequacy must be explicitly
27
28 justified.
29
30
31
32

33 **Claim B1.2:** *The implemented ethics comply with the legal, social and ethical norms of*
34
35 *the environment of use*
36
37

38 The implemented ethics must be compatible with behaviour that would be reasonably
39
40 expected by the general public for an AV operating within the stated environment. It
41
42 should be noted that this does not necessarily imply an AV should behave in exactly the
43
44 same way as a human driver (IEEE Global, 2018), but rather that the AV should act in a
45
46 way that a human might plausibly expect from an AV.
47
48
49

50 **Claim C1:** *The intent of the implemented ethics shall be maintained throughout*
51
52 *decomposition into AV design requirements and risk management decisions.*
53
54
55
56
57
58
59
60
61
62
63
64
65

1 This claim supports P2, and serves to translate ethical requirements into lower-level
2 safety and risk management requirements which implement the ethical intent. It is
3
4 broken down into two sub-claims as follows:
5
6

7 **Claim C1.1:** *The risk balancing and trade-offs resulting from the implemented ethics*
8
9 *across different environments have been specified in the form of a risk profile*

10
11 The risk profiles discussed in Section 4 provide a method of reflecting ethical
12
13 perspectives in risk management, risk balancing and risk distribution decisions.
14
15

16 Satisfaction of this claim requires provision of an explicit description of the risk trade-
17
18 offs and balances that have been performed, in the form of a risk profile.
19
20

21 **Claim C1.2:** *The design of the AV reflects the risk profile*

22
23 This claim is supported by an argument that the risk balancing inherent in the specified
24
25 risk profile has been performed accordingly, via the translation of this risk profile into
26
27 technical, safety and risk requirements. It is broken down into three further sub-claims
28
29 as follows:
30
31

32
33 **Claim C1.2.1:** *The risk profile has been decomposed into requirements on the design*
34
35 *and functionality of the system*

36
37 This claim is supported by referencing out to the design and implementation
38
39 requirements specification, as well as to the safety case.
40
41

42
43 **Claim C1.2.2:** *The design and functional requirements of the system are satisfied.*

44
45 This claim is supported by referencing evidence provided within the safety case, which
46
47 is the primary mechanism for demonstrating satisfaction of design and safety
48
49 requirements.
50
51

52
53 **Claim D1:** *The implemented ethics requirements are satisfied.*

54
55 This claim partially supports P3 and requires identification of what the acceptable
56
57 ethical behaviour of the AV might be. Natural-language interpretation of the ethics
58
59
60
61
62
63
64
65

1 requirements will help support this claim, as will a description of the functionality and
2 behaviours which comply with these ethics. Supplementary supporting evidence will
3 include system verification and validation of the derived requirements sourced from the
4 ethical imperatives via risk profiles. Traceability between these derived requirements
5 and any further lower-level requirements must be demonstrated, along with traceability
6 between these derived requirements and verification artefacts.
7
8
9

10
11
12
13
14 **Claim E1:** *Any conflicts between emergent behaviour of the AV and the implemented*
15 *ethics are identified and mitigated so far as is reasonably practicable.*
16

17
18
19 This claim supports P4 and requires an estimation of likely emergent behaviours and
20 changes in the environment throughout the AV's lifetime. Support for this claim
21 requires a gap analysis of potential environmental change, as well as of gap analysis
22 between the behaviours which may potentially be learnt by AVs (via adaptive
23 algorithms and continuous machine learning) and the behaviours which were "hard-
24 coded" or scripted at deployment. Any conflicts between the ethics requirements and
25 these new behaviours and environments must be identified and mitigated so far as is
26 reasonably practicable.
27
28
29
30
31
32
33
34
35
36
37
38
39

40 **5.3.1 Engineering-focused process**

41
42 In order to be applicable to a given system, the ethics assurance case template
43 that we describe above must be instantiated for that system. This is achieved via the
44 second part of our framework: an engineering-focused process that uses the risk profiles
45 of Section 4 to makes risk trade-offs and risk balances explicit. This process sits within
46 the wider engineering lifecycle, and relies on artefacts from that lifecycle, such as test
47 results, requirements specifications etc.
48
49
50
51
52
53
54
55
56

57 The process consists of several phases, as shown in Figure 2, with decision points and
58 feedback loops between the phases. We emphasise that the decision points and phases
59
60
61
62
63
64
65

of Figure 2 relate only to the instantiation of the ethics assurance template.

Consequently, Figure 2 deliberately excludes decision points and feedback loops which relate to the wider engineering lifecycle, such as feedback loops to indicate iteration over engineering requirements, or decision points to verify that such requirements are fulfilled. This ensures that instantiation of the ethics assurance template does not depend on a specific engineering lifecycle process being used.

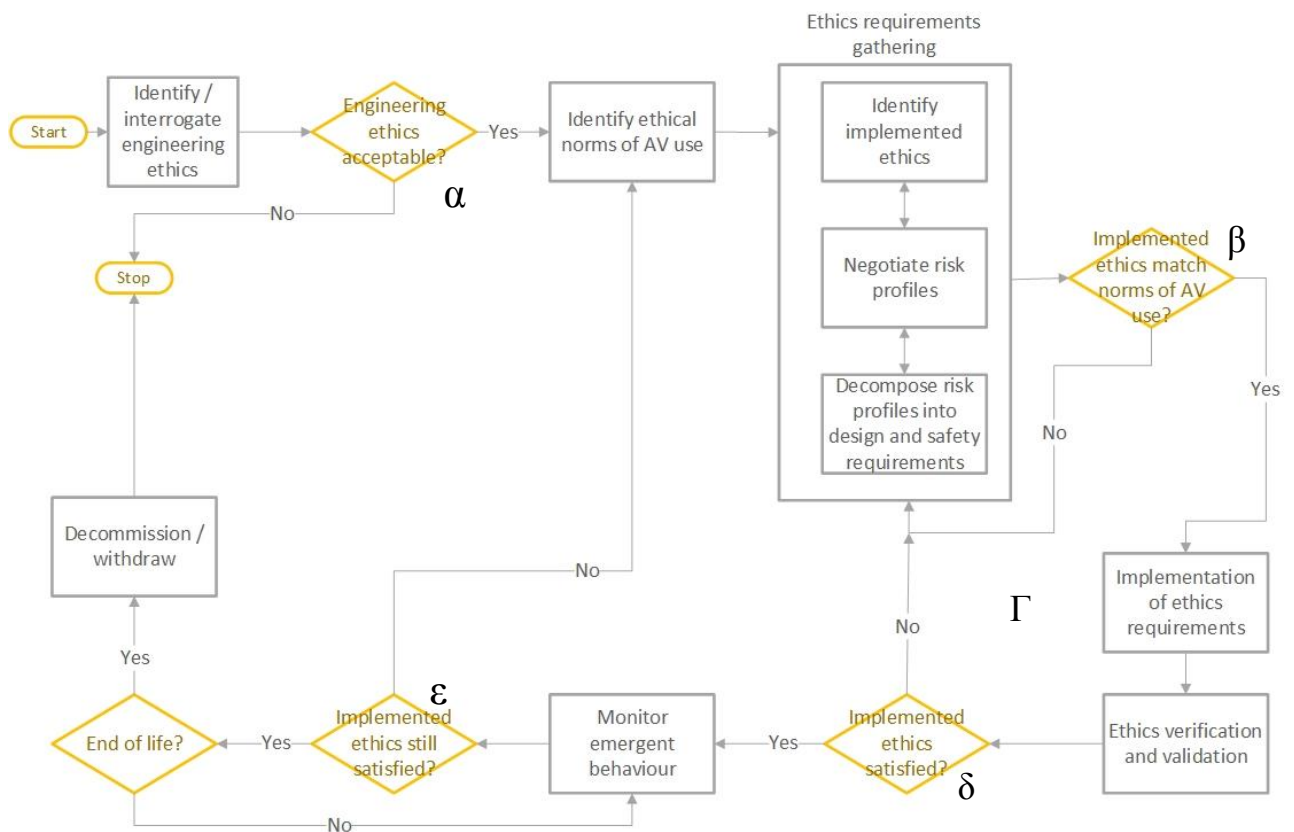


Figure 2: engineering-focused process

Phase 1: identify and interrogate engineering ethics

This phase instantiates Claim A1 of Figure 1, and consists of activities to identify and interrogate the engineering ethics of the developers. These ethical principles and codes of conduct are crucial to maintain the integrity of the AV development for stakeholders to have confidence in all artefacts produced during development. Because of the

1 importance of the developers' engineering ethics, concerns at this point may result in a
2 hard stop for the project via decision point α .
3

4 ***Phase 2: Identify ethical norms of AV use***

5
6
7 This phase partially instantiates Claim B1 of Figure 1 and consists of activities to
8 identify the legal, social and ethical norms of the proposed AV use. These norms will be
9 situation-dependent (an AV operating on a motorway will be subject to different norms
10 than an AV operating near a school). As discussed earlier, this phase is not about
11 determining how human drivers act in a given environment, but rather about
12 determining how humans will expect the prospective AV to act in its environment of
13 use situation.
14
15
16
17
18
19
20
21
22
23

24 ***Phase 3: Ethics requirements gathering***

25
26 This phase instantiates Claim C2, and also completes the instantiation of Claim B1. It is
27 a tri-part phase which also serves as an interaction and feed-in point for the
28 requirements phase of the wider engineering lifecycle. During this phase the developers
29 identify the implemented ethics (i.e. the ethics that will govern the behaviour of the
30 AV), negotiate risk profiles which reflect those ethics and their associated risk trade-
31 offs, and decompose these risk profiles into design and safety requirements. (There will
32 of course, be additional requirements gathering activities defined within the engineering
33 lifecycle, but these are excluded from the scope of Figure 2). Should the implemented
34 ethics not match the norms of AV use in the proposed environment, this phase is
35 iterated over again, as shown via the decision point β and its associated feedback loop in
36 Figure 2. We note that in practice that gathering and refinement of ethics requirements
37 is likely to be an ongoing process, occurring simultaneously with implementation of
38 ethics requirements and verification and validation of ethics requirements. This iteration
39 and simultaneous development is shown via the feedback loop Γ of Figure 2.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Phase 4: Implementation of ethics requirements

This phase serves as an interaction and feed-in point for the system development and implementation phase(s) of the wider engineering lifecycle. This phase incorporates development and implementation against the ethics requirements of the AV specification. There is no specific output from this phase in terms of instantiation of the ethics assurance case template.

Phase 5: Ethics verification & validation

This phase instantiates Claim D1 of Figure 1. It also serves as an interaction and feed-in point for the verification and validation performed during the wider engineering lifecycle. This phase incorporates verification and validation of the implemented AV system against its ethical requirements. Should this V&V identify that the ethical requirements have not been met (via decision point δ), there will be a return to Phase 3 to iterate over the gathering, refinement and subsequent implementation of these ethical requirements.

Phase 6: Monitor emergent behaviour

This phase instantiates Claim E1 of Figure 1, and consists of activities to monitor any behaviour of the AV which emerges after its deployment. Emergent behaviour may be due to the use of continuous machine learning, which requires the AV to adapt and refine its behaviours in response to new input. Similarly, emergent behaviour may be due to an original under-specification of behaviours, or over-restrictive assumptions about the environment in which the AV may eventually need to operate. Where any emergent behaviour is identified, a further assessment activity must be made to determine whether this behaviour complies with the ethical requirements of the AV. Because such emergent behaviour may relate to a new or unforeseen use of the AV, or environment of use, any contravention of the ethical requirements results in a return to

1 Phase 2 (via decision point ϵ) in order to iterate and expand upon the norms of AV use
2 for this new environment.
3

4 ***Phase 7: decommission and withdrawal***

5
6
7 This phase corresponds to the end-of-life point for an AV system. End-of-life decisions
8 may be taken where the system is superseded by another, or where legal, ethical,
9 technological or societal factors lead to a situation where withdrawal of the system is
10 recommended. There is no specific output from this phase in terms of instantiation of
11 the ethics assurance case template.
12
13
14
15
16
17
18
19
20

21 **6 Conclusions**

22
23 In this paper, we have discussed the ethical concerns around the introduction and use of
24 AVs, focusing on the perceived risk these present, and the acceptability of this risk. We
25 have identified a number of factors which affect risk perception and risk acceptance,
26 particularly where ethical principles such as altruism or “fairness” are concerned.
27
28
29
30
31

32
33 One of the most significant underlying ethical issues we have discussed is the
34 situation where different risks are balanced against each other, with an increase in one
35 risk being traded off against a decrease in another. This is of particular concern to AVs
36 where risk is transferred between different segments of society (e.g. from the AV
37 passenger to pedestrians or other road users). At present it is not always clear when such
38 risk trade-offs have been made by developers, or what the ethical factors motivating
39 these are. We argue that this lack of clarity can lead to a lack of informed societal
40 consent to these risks.
41
42
43
44
45
46
47
48
49
50

51
52 In order to increase transparency of the ethical factors which drive such risk
53 decisions, we have identified a methodology for explicitly translating the underlying
54 ethics of AV behaviour into safety and design requirements on the AV. This
55 methodology makes use of two components: an ethics assurance case template, and an
56
57
58
59
60
61
62
63
64
65

1 engineering-focused process for instantiating this template. Together, these ensure that
2 the underlying ethical and risk-balancing choices of the AV developers are identified
3 and explicitly justified. This is achieved via the use of risk profiles, which require
4 developers to describe the approach to risk reduction that they have taken, and to justify
5 their risk transfer decisions.
6
7
8
9
10

11 The framework we propose is also applicable to fully-autonomous systems in
12 other domains with ethical concerns, such as the military and healthcare fields. We
13 propose to expand this work in future to apply this methodology to a working case
14 study, either in the automotive, military or healthcare domains. A foundational step in
15 this will be to build a taxonomy of risk profiles in conjunction with industry personnel.
16 This will allow us to represent the risk transfer decisions which are made in a real-world
17 scenario.
18
19
20
21
22
23
24
25
26
27

28 We also propose to extend the framework to address the issue of confidence.
29 Like safety, ethics is a limit concept (Kelly, Habli, Nicholson, Megone & Mcnish,
30 2014) and the degree of confidence stakeholders have in the ethics of a system will be
31 dependent on the evidence they have been shown. We propose further study into the
32 societal perception of AVs in order to identify what factors increase confidence in the
33 ethical principles governing the AV, and what steps can be taken to mitigate a lack of
34 confidence. As part of this work we will also consider ways to formalise certain ethical
35 principles relevant to AV introduction and operation, including the principle of double
36 effect and Kantian ethics, along the lines of (Bentzen, 2015; Linder & Bentzen, 2018).
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 **References**

52 Allais, M. (1953). Le comportement de l'homme rationnel devant le risque: Critique des
53 postulats et axiomes de l'école américaine. *Econometrica: Journal of the*
54 *Econometric Society*, 503-546
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- Anderson, J., Nidhi, K., Stanley, K., Sorenson, P., Samaras, C. & Oluwatola, O. (2014). *Autonomous Vehicle Technology: A Guide for Policymakers*, Rand Corporation.
- Anderson, M., & Anderson, S. (2007). Machine ethics: Creating an ethical intelligent agent, *AI Magazine* 28, 15—26.
- Arkin, R., Ulam, P. & Wagner, A. (2012). Moral decision-making in autonomous systems: enforcement, moral emotions, dignity, trust, and deception. In *Proceedings of the IEEE*, volume 100, 571--589.
- Baron, J. (1985) *Rationality and intelligence*. New York, Cambridge University Press.
- Bell, D. (1983). Risk premiums for decision regret. *Management Science* Vol 29(10), 1156—1166.
- Bentzen, M. (2016). The principle of double effect applied to ethical dilemmas of social robots. In *Proceedings of the 2nd International Conference on Robophilosophy*, 268—279.
- Bloomfield, R. & Bishop, P. (2010). Safety and assurance cases: Past, present and possible future: an Adelard perspective. In *Proceedings of the Eighteenth Safety-Critical Systems Symposium*.
- Cokey, E. & Kelley, C. (2009) Cognitive Abilities and Superior Decision Making. In *Judgement and Decision Making*, Vol 4(1), 20 – 33.
- Common Criteria Management Board. (2007). Common Methodology for Information Technology Security Evaluation. CCMB-2007-09-004.
- Dennis, L., Fisher, M., Slavkoviv, M. & Webster, M. (2004). Formal verification of ethical choices in autonomous systems, *Robotics and Autonomous Systems* 77, 1--14.
- Donde, J. (2017). Self-driving cars will kill people. Who decides who dies? *Wired*. <https://www.wired.com/story/self-driving-cars-will-kill-people-who-decides-who-dies>.
- Evans, J. (1989). *Bias in human reasoning: Causes and consequences*. Brighton, Erlbaum
- Fagnant, D. & Kockelman, K. (2014). The travel and environmental implications of shared autonomous vehicles, using agent-based model scenarios. In *Proceedings of the 93rd Annual Meeting of the Transportation Review Board*.
- Foot, P. (1967). The problem of abortion and the doctrine of double effect, *Oxford Review* 5, 5—16.

- 1 Gerdes, J. & Thornton, S. (2016). Implementable ethics for autonomous vehicles,
2 *Autonomous Driving* 87--102.
- 3 Gips, J. (1995). Towards the Ethical Robot. *Android Epistemology*, MIT Press, 243—
4 252.
- 5
6
7 Goodall, N. (2016). Away from trolley problems and toward risk management, *Applied*
8 *Artificial Intelligence*, vol. 30, no. 8, 810—821.
- 9
10
11 Groves, D., Kalra, N. (2017). *Autonomous Vehicle Safety Scenario Explorer*, Web-Only
12 Tool, RAND Corporation, <https://www.rand.org/pubs/tools/TL279.html>.
- 13
14 Object Management Group. (2019) Structured Assurance Case Metamodel (SACM),
15 Document Number 20190314. <https://www.omg.org/spec/SACM/2.1/Beta1/>
- 16
17
18 Hawkins, R., Habli, I., Kelly, T. & McDermid, J. (2013). Assurance Cases and
19 Prescriptive Software Certification: A Comparative Study, *Safety Science*, 55--
20 71.
- 21
22
23 Hawkins, R., Habli, I., Kolovos, D., Paige, R. & Kelly, T. (2015). Weaving an
24 Assurance Case from Design: A Model-based Approach. In *Proceedings of the*
25 *16th IEEE International Symposium on High Assurance Systems Engineering*.
- 26
27
28 Health and Safety Executive. (2001). *Reducing Risks, Protecting People*.
29 <http://www.hse.gov.uk/risk/theory/r2p2.pdf>
- 30
31
32 Health and Safety Executive. (2002). *The Precautionary Principle: Policy and*
33 *Application*.
34 <http://www.hse.gov.uk/aboutus/meetings/committees/ilgra/pppa.htm>
- 35
36
37
38 IEEE Global Initiative. (2018). *Ethically aligned design*, v2.0.
39 [https://standards.ieee.org/content/dam/ieee-](https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead_v2.pdf)
40 [standards/standards/web/documents/other/ead_v2.pdf](https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead_v2.pdf).
- 41
42
43
44 Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under
45 risk. *Econometrica: Journal of the Econometric Society*, 263-291.
- 46
47
48
49 Kalra, N. & Groves, D. (2017). *The Enemy of Good - Estimating the Cost of Waiting for*
50 *Nearly Perfect Automated Vehicles*, Technical Report RR-2150 RC, Rand
51 Corporation. https://www.rand.org/pubs/research_reports/RR2150.html
- 52
53
54
55 Kalra, N. & Paddock, S. (2016). *Driving to Safety: How Many Miles of Driving Would*
56 *It Take To Demonstrate AV Reliability?*, Technical Report RR1478, RAND
57 Corporation, <https://www.rand.org/pubs/researchreports/RR1478.html>.
- 58
59
60
61
62
63
64
65

- 1 Kelly, T. & Weaver, R. (2004). The goal structuring notation — a safety argument
2 notation. In *Proceedings of Dependable Systems and Networks 2004 Workshop*
3 *on Assurance Cases*.
4
- 5 Kelly, T., Habli, I., Nicholson, M., Megone, C. & Mcnish, K. (2014). The ethics of
6 acceptable safety. In *Proceedings of the 23rd Safety-critical Systems*
7 *Symposium*.
8
9
- 10 Kelly, T. (2007). Reviewing Assurance Arguments — A Step-By-Step Approach. In
11 *Proceedings of the 12th International Conference on Dependable Systems and*
12 *Networks DSN 84—95*.
13
14
- 15 Kim, S. & McGill, A. (2011). Gaming With Mr Slot or Gaming the Slot Machine?,
16 *Journal of Consumer Research*, Vol 38, No 1, 94--107.
17
18
- 19 Knight, J. (2002). Safety Critical Systems: Challenges and Directions. In *Proceedings of*
20 *the 24th International Conference on Software Engineering*, 547--550.
21
22
- 23 Kuderer, M., Gulati, S. & Burgard, W. (2015). Learning Driving Styles for Autonomous
24 Vehicles from Demonstration. In *Proceedings of the IEEE International*
25 *Conference on Robotics and Automation (ICRA)*.
26
27
- 28 Lin, P. (2015). Why Ethics Matters for Autonomous Cars. In *Autonomes Fahren*,
29 Springer Vieweg, 69 –85.
30
31
- 32 Lindner, F. & Bentzen, M. (2018). A formalization of Kant’s second formulation of the
33 categorical imperative. In *Proceedings of the 14th International Conference on*
34 *Deontic Logic and Normative Systems*, College Publications.
35
36
- 37 Liu, S., Stavridou, V., Duarte, B. (1995). The Practice of Formal Models in Safety
38 Critical Systems. *Journal of Systems and Software*, 77—87.
39
40
- 41 Martin, M. & Schinzinger, R. (2005). *Ethics in engineering*, McGraw-Hill New York.
42
43
- 44 Menon, C. & Alexander, R. (2017). A safety-case approach to ethical considerations for
45 autonomous vehicles. In *Proceedings of the 12th International Conference on*
46 *System Safety and Cyber Security*.
47
48
- 49 Menon, C., Bloomfield, R. & Clements, T. (2013). Interpreting ALARP. In *Proceedings*
50 *of the 8th IET International System Safety Conference*.
51
52
- 53 Menon, C. & Alexander, R. (2018). Ethics and the safety of autonomous systems. In
54 *Proceedings of the 26th Safety Critical Systems Symposium*.
55
56
- 57 Menon, C. (2017). *A White Paper Discussing Ethical Considerations for Autonomous*
58 *Vehicles*, Technical Report, Transport Systems Catapult,
59 <https://ts.catapult.org.uk/intelligent-mobility/im-resources/research-papers/>
60
61
62
63
64
65

- 1 MIT. (2018). *MIT Moral Machine*. <http://moralmachine.mit.edu/>.
- 2 Von Neumann, J. & Morgenstern, O. (1947). *Theory of Games and Economic*
3 *Behaviour*. Princeton University Press.
- 4
- 5 Nilsson, J. (2018). Safe self-driving cars: Challenges and some solutions. In
6 *Proceedings of the 26th Safety Critical Systems Symposium*.
- 7
- 8 Office for Nuclear Regulation. (2006). *Safety Assessment Principles for Nuclear*
9 *Facilities*. <http://www.onr.org.uk/saps/saps2014.pdf>
- 10
- 11 Office for Nuclear Regulation (2018). *Guidance on the Demonstration of ALARP*,
12 http://www.onr.org.uk/operational/tech_asst_guides/ns-tast-gd-005.pdf
- 13
- 14 Office of Government Commerce. (2019). *Prince2, Project Methodology*,
15 <https://www.prince2.com/uk/prince2-methodology>.
- 16
- 17 Reyna, V. (2004). How People Make Decisions That Involve Risk: A Dual Process
18 Approach. *Current Directions in Psychological Science*, Vol 13(2), 60 – 66.
- 19
- 20 Royal Academy of Engineering (2017). *Statement of ethical principles*.
21 <https://www.engc.org.uk/media/2337/statement-of-ethical-principles-2014.pdf>
- 22
- 23 SAE International. (2018). *J3016: Taxonomy and Definitions for Terms Related to*
24 *Driving Automation Systems for On-Road Motor Vehicles*.
25 https://www.sae.org/standards/content/j3016_201401/
- 26
- 27 Shalev-Schwartz, S., Shammah, S. & Shashua, A. (2017). *On a Formal Model of Safe*
28 *and Scalable Self-driving Cars*, arXiv preprint arXiv:1708.06374.
- 29
- 30 Somasundaram, J. & Diecidue, E. (2015). Regret Theory and Risk Attitudes. *Journal of*
31 *Risk and Uncertainty*, Issue 2-3, pp. 147-175
- 32
- 33 The Assurance Case Working Group. (2018). *Goal Structuring Notation Community*
34 *Standard (Version 2)*, Technical Report SCSC-141B, The Safety Critical
35 Systems Club. <https://scsc.uk/scsc-141B>
- 36
- 37 Thornton, S. (2018). *Autonomous vehicle motion planning with ethical considerations*,
38 (Doctoral dissertation), Stanford University.
- 39
- 40 Transport Systems Catapult. (2017) *Market Forecast For Connected And Autonomous*
41 *Vehicles*.
42 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/atta](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/642813/15780_TSC_Market_Forecast_for_CAV_Report_FIN_AL.pdf)
43 [chment_data/file/642813/15780_TSC_Market_Forecast_for_CAV_Report_FIN](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/642813/15780_TSC_Market_Forecast_for_CAV_Report_FIN_AL.pdf)
44 [AL.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/642813/15780_TSC_Market_Forecast_for_CAV_Report_FIN_AL.pdf)
- 45
- 46 Uber Advanced Technologies Group. (2018). *Safety Report: A Principled Approach to*
47 *Safety*. <https://uber.app.box.com/v/UberATGSafetyReport>
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1 UK Autodrive. (2017). *Public Attitudes Survey*.
2 [http://www.ukautodrive.com/wpcontent/uploads/2017/08/Executive-Summary-
4 FINAL.pdf](http://www.ukautodrive.com/wpcontent/uploads/2017/08/Executive-Summary-
3 FINAL.pdf).
- 5 UK Autodrive. (2017). *The Moral Algorithm*. [http://www.ukautodrive.com/wp-
7 content/uploads/2017/08/moral_algorithm_white_paper_A4_051216.pdf](http://www.ukautodrive.com/wp-
6 content/uploads/2017/08/moral_algorithm_white_paper_A4_051216.pdf).
- 8 UK Autodrive. (2018). *CAV: A Hacker's Delight*, [http://www.ukautodrive.com/wp-
10 content/uploads/2018/01/Cyber_security_White_paper_A4_050917.pdf](http://www.ukautodrive.com/wp-
9 content/uploads/2018/01/Cyber_security_White_paper_A4_050917.pdf).
- 11 UK Government Centre for Connected and Autonomous Vehicles. (2018). *UK
12 Connected And Autonomous Vehicle Research And Development Projects*.
13 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment
15 data/file/737778/ccav-research-and-development-projects.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment
14 data/file/737778/ccav-research-and-development-projects.pdf)
- 16 UK Government Department for Transport. (2015). *The Pathway To Driverless Cars*.
17 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/atta
19 chment_data/file/446316/pathway-driverless-cars.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/atta
18 chment_data/file/446316/pathway-driverless-cars.pdf)
- 20 UK Government Department for Transport. (2017). *The key principles of vehicle cyber
21 security for connected and automated vehicles*.
22 [https://www.gov.uk/government/publications/principles-of-cyber-security-for-
25 connected-and-automated-vehicles/the-key-principles-of-vehicle-cyber-security-
26 for-connected-and-automated-vehicles](https://www.gov.uk/government/publications/principles-of-cyber-security-for-
23 connected-and-automated-vehicles/the-key-principles-of-vehicle-cyber-security-
24 for-connected-and-automated-vehicles)
- 27 UK Government (1998) *Human Rights Act*.
28 <https://www.legislation.gov.uk/ukpga/1998/42/contents>
- 29 UK Ministry of Defence. (2017). *Defence Standard 00-56: Safety Management
30 Requirements for Defence Systems*, Technical Report 00-56 Issue 7.
- 31 Venturer Cars. (2016). *Introducing Driverless Cars to UK Roads*. [https://www.venturer-
33 cars.com/wp-content/uploads/2016/08/VENTURER-Trial-1-Overview.pdf](https://www.venturer-
32 cars.com/wp-content/uploads/2016/08/VENTURER-Trial-1-Overview.pdf)
- 34 Venturer Cars. (2017). *Interactions Between Autonomous Vehicles and Other Vehicles
35 at Links and Junctions*. [http://www.venturer-cars.com/wp-
37 content/uploads/2017/11/VENTURER-Trial-2-Technical-Reportv2.pdf](http://www.venturer-cars.com/wp-
36 content/uploads/2017/11/VENTURER-Trial-2-Technical-Reportv2.pdf).
- 38 Vanderelst, D. & Winfield, A. (2018). An architecture for ethical robots inspired by the
39 simulation theory of cognition, *Cognitive Systems Research* 48 56--66.
- 40 Wallach, W. & Allen, C. (2008). *Moral Machines: Teaching Robots Right from Wrong*,
41 Oxford University Press.
- 42 Wang, S. & Zhao, J. (2019). *Risk Preference and Adoption of Autonomous Vehicles*,
43 Transportation Research A.

Waymo. (2018). *Safety Report: On the Road to Fully Self-Driving*.

<https://storage.googleapis.com/sdc-prod/v1/safety-report/Safety%20Report%202018.pdf>.

Winfield, A., Blum, C. & Liu, W. (2014). Towards an ethical robot: Internal models, consequences and ethical action selection. In *Advances in Autonomous Robotics Systems*, 85--96.

Yoo, J., Jee, E. & Cha, S. (2009). Formal Modelling and Verification of Safety-Critical Software, *IEEE Software*, 42—49.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65