

# **Addressing Inequal Risk Exposure in the Development of Automated Vehicles**

**Manuel Dietrich**

**2021**

**Preprint:**

This is a post-peer-review, pre-copyedit version of an article published in Ethics and Information Technology. The final authenticated version is available online at: <https://doi.org/10.1007/s10676-021-09610-1>



# Addressing inequal risk exposure in the development of automated vehicles

Manuel Dietrich<sup>1</sup>

Accepted: 29 July 2021  
© The Author(s) 2021

## Abstract

Automated vehicles (AVs) are expected to operate on public roads, together with non-automated vehicles and other road users such as pedestrians or bicycles. Recent ethical reports and guidelines raise worries that AVs will introduce injustice or reinforce existing social inequalities in road traffic. One major injustice concern in today's traffic is that different types of road users are exposed differently to risks of corporal harm. In the first part of the paper, we discuss the responsibility of AV developers to address existing injustice concerns regarding risk exposure as well as approaches on how to fulfill the responsibility for a fairer distribution of risk. In contrast to popular approaches on the ethics of risk distribution in unavoidable accident cases, we focus on low and moderate risk situations, referred to as routine driving. For routine driving, the obligation to distribute risks fairly must be discussed in the context of risk-taking and risk-acceptance, balancing safety objectives of occupants and other road users with driving utility. In the second part of the paper, we present a typical architecture for decentralized automated driving which contains a dedicated module for real-time risk estimation and management. We examine how risk estimation modules can be adjusted and parameterized to redress some inequalities.

**Keywords** Automated vehicles · Distributive justice · Responsibility · Road traffic injustice · Fairness · Ethics · John Rawls · Risk

## Introduction

Automated vehicles (AVs) are an emerging technology which has raised intense ethical questions, both in the academic world and the general public. Most attention has been directed to on the decision-making dilemmas AVs might face when confronted with situations where they have to decide between different very negative outcomes within split-seconds, the so-called trolley case situations (Bonnenfon, Shariff, & Rahwan, 2016; Lin, 2015).

In recent years, potential ethical challenges connected to AVs have been explored systematically from different disciplinary perspectives. Reports and guidelines by the European Union, the German Ministry of Transport and Digital Infrastructure, the professional organization IEEE as well as other institutions have discussed the ethical and

legal challenges of this future technology.<sup>1</sup> The most severe concerns are that AVs will reinforce existing inequalities and introduce new injustice aspects in the context of road traffic. An existing inequality is the unequal exposure to risks of bodily harm induced by a traffic accident between different types of road users.<sup>2</sup> Differences in the exposure to accident risks are reflected in the statistics of actual injuries and fatalities. Studies have shown that non-motorized road users suffer from a disproportionately higher share of serious injuries and fatalities (e.g., Mullen et al., 2014; for more see: Sect. “[Inequal risk exposure as injustice concern](#)”). That means there is a disparity between the ratio of harm and the share of actual road usage. Another, less researched, concern

<sup>1</sup> For example, German Ethics Commission report (Di Fabio, Broy, & Brüngger, 2017), EU experts group's Guidelines for a Trustworthy AI (AI High Level Expert Group, 2019), EU experts group's report on Ethics of Connected and Automated Vehicles (Horizon 2020 Commission Expert Group, 2020), IEEE Global Initiative's report “Ethically Aligned Design” (2019).

<sup>2</sup> Risk refers to an unwanted event which might appear in the future. It is formalized as “the product of its probability and some measure of its severity” (Hansson, 2018), both regarding the unwanted event. In this paper, we are referring to physical harm caused by accidents as the unwanted event.

✉ Manuel Dietrich  
manuel.dietrich@honda-ri.de

<sup>1</sup> Honda Research Institute Europe, Carl-Legien-Straße 30, 63073 Offenbach am Main, Germany

is about differences in the general risk exposure among vehicle categories. Categories are defined for example by vehicle size and weight as well as its equipment, in particular with respect to safety systems.

In line with the recent Ethics of Connected and Automated Vehicles report (Horizon 2020 Commission Expert Group, 2020), addressing such concerns is the responsibility of developers and deployers. They “should be held responsible for designing and operating CAVs [connected and automated vehicles] in ways that neither discriminate against individuals or groups of users, nor create or reinforce large-scale social inequalities” (Horizon 2020 Commission Expert Group, 2020, p. 43).

In this paper, we elaborate on the developer/deployer perspective and first discuss their responsibility in addressing risk exposure injustice by considering risk ethics and institutional justice.<sup>3</sup>

In contrast to many previous approaches, which look at risk distribution and its ethical justification in accident scenarios (Di Fabio, Broy, & Brünger, 2017; Goodall, 2017; Schäffner, 2020), we focus on how AVs ought to decide in situations of low and moderate risk, considered as routine or general driving. To also look at routine driving is important as the related decisions impact AVs and other traffic participants chances to enter dangerous situations. Even though we can find other authors emphasizing the importance of an ethically aligned risk allocation before accidents (for general or routine driving), there has been little ethical analysis addressing it specifically (Goodall, 2017; Himmelreich, 2018; Keeling, 2020). For routine driving the obligation to distribute risks ethically must be viewed in the context of risk-taking and risk-acceptance, balancing safety objectives of occupants and other road users with driving utility. This will be covered in the first part of the paper (Sect. “[Risk exposure and automated driving ethics](#)”).

A second important aspect to be addressed is to strongly relate the ethical considerations to current technical research on risk estimation and management for automated driving applications (Sects. “[Decentralized modular automated driving](#)” and “[Addressing unequal exposure in AV development](#)”). We refer to a decentralized architecture for automated driving which is typical in research and the industry. A decentralized AV’s decision only has an immediate, local impact on how its occupants and other road users are exposed to risk. However, since AVs will likely enter the roads as one entity of a fleet of vehicles of the same

brand which share the same decision-making architecture – individual effects multiply and therefore AV’s impact on risk distribution has broader relevance. Based on a practical example, we sketch how claims for a fair distribution of risk translate into design requirements, meaning how they could be addressed technically in the behavior planning and execution of AVs.

## Inequal risk exposure as injustice concern

In this section, we look at literature which examines injustice concerns in road traffic to gain a more general understanding about what is considered as unequal or unfair regarding risk exposure.

One major concern is the road infrastructure in contemporary urban societies, as it tends to prioritize one mode of transportation, namely motorized traffic (Martens, 2017; Nello-Deakin, 2019). Users of motorized traffic suppress other means of transportation as they occupy more space than other road users and are privileged with regard to safety: “Pedestrians and cyclists face disproportionate risks, particularly as compared to drivers and passengers of motor vehicles” (Mullen et al., 2014, p. 238). Quantitatively, “pedestrians make up 11 percent of traffic fatalities in the United States, far out of proportion to the amount of walking” (Ewing & Dumbaugh, 2009, p. 348). Furthermore, the choice of the mode of transportation often correlates with the social background. An article in the magazine *Dissent* argues that “America’s suburbs are engineered against the walking poor” (Ross, 2014).

Encounters between cars and others can be more dangerous for the non-car users due to the difference in physical properties, velocity and in the capability of protective measures. Furthermore, even within the category of motorized vehicles, differences in mass and protection capabilities can have a significant impact on risk exposure. Padmanaban states that “of all vehicle parameters, mass is the most important factor influencing odds of driver fatality” (Padmanaban, 2003, p. 523). Another concern is the velocity which has a large impact on the severity of collisions, especially between cars and so-called vulnerable road users (e.g., pedestrians, cyclists). Zegeer et al. (2002) show that the fatality rate for pedestrians reaches 80% at a speed of 40 mph (ca. 64 km/h). To sum up, large and heavy vehicles operating with high speed on the street can enforce a particularly unequal exposure to risk.

In literature, the risk exposure of vulnerable road users is compared to the risk exposure of car occupants. However, finding an appropriate measurement for judging inequalities in exposure to risk is not straightforward. Traditionally, injustice concerns refer to inequalities in the distribution of the benefits and burdens of social cooperation (Rawls, 1971,

<sup>3</sup> We will refer to responsibility in its forward-looking understanding: Responsibility as obligation which “shapes the development, introduction, and use of CAVs in a way that promotes societal values and human well-being” (Horizon 2020 Commission Expert Group, 2020, p. 53; also see: Van de Poel & Sand, 2021).

p. 4). Roads are a publicly shared good through which citizens can fulfill their right of mobility having the free choice of the means of transportation.<sup>4</sup> Given this, an imbalance in the exposure to risks classifies as a justice concern. The total number of deaths and serious injures however must be assessed in relation to the actual usage. Mullen et al. discuss different methods. They compare the ratio of death and serious injury to the time spent travelling and in another approach the ratio to the number of trips (2014). The conclusion for both calculations was that pedestrians suffer substantially more from traffic-induced harms.

In transport and infrastructure planning research, we see increased attention towards addressing such allocation inequalities for future projects (Jones, 2014; Gössling, 2016). However, planning, authorization, and construction can take decades to complete and is often slowed down or even hindered by the real and expected social resistance and cost of infrastructure changes (Gössling, 2016, p. 8).<sup>5</sup>

At a shorter timeframe, it is common to apply local measures directly addressing circumstances which are considered as unproportionally dangerous for certain or all road users or to react to special vulnerabilities for instance in areas with school children. It has been proofed highly efficient to deploy traffic calming measures of physical nature, for instance speed bumps. Ewing and Dumbaugh show that collisions declined significantly after traffic was calmed in this way (2009, p. 356).

## Risk exposure and automated driving ethics

In the previous section we have discussed that an unequal exposure to risk among different types of road participants can be seen as a social justice concern. In this section, we want to elaborate on the role automated vehicles could take in the risk landscape of future road traffic. We would like to answer if there is a responsibility for manufacturers to explicitly design AVs which target a fair distribution of risk in routine driving. Before we elaborate routine driving and

risk from an ethical perspective, we first look at the literature on risk ethics for automated accident decisions.

## Risk-ethics and automated accident decisions

The ethical relevance of automated driving decisions became evident when it could be shown that such vehicles might get into situations where they must decide between driving paths which will all very likely lead to serious harm—meaning human deaths as a not definite but very plausible consequence. These so-called dilemma situations or trolley cases have gained much attention in scientific research as well as public media.<sup>6</sup> As part of the MIT moral machine experiment, multiple interactive online studies were conducted, in which participants were asked how they would decide on behalf of an AV.<sup>7</sup> People had to decide between different outcomes, for instance killing the occupants of the vehicle by driving into a wall or killing a group of people crossing the road. The researchers discovered that there are moral preferences or beliefs on how an AV should decide which are shared by a majority (Awad et al., 2018). Framing serious accident situations as trolley problems and approaching them from an experimental ethics perspective has evoked many concerns (e.g., Keeling, 2017, 2020; Nyholm & Smids, 2016). It has been mainly criticized in two ways. First, using experimental ethics results as a base for how automated systems should be programmed, a so called “empirically-informed policymaking”, has no ethical justification, and therefore no justification to be better than established normative theories (Keeling, 2020, p. 35). Also, such decision policies conflict with established justice theories of how burdens of risk and harm should be distributed as well as with human rights principles. Second, the experimental setup is considered as unrealistic and naïve, since in real driving dilemmas the AV has many subtle choices which do not determine for a definite outcome, but rather influence the probability of events (Goodall, 2017, p. 496; Nyholm & Smids, 2016; Liu, 2018, p. 153). In other words, in accident situations the vehicle does not have to decide between actual harm, serious injuries, or death, but rather between options which all have a high probability to lead to an accident with a high severity. Due to the uncertainty of outcomes, it is suggested that it might be better to frame accident decisions from a risk-ethics perspective (Goodall, 2017, p. 496).

<sup>4</sup> It seems not to be reasonable to request from individuals to use other means of transportation or argue that they should be aware of the risks of road traffic and if they decide to participate, they implicitly accept these risks. At least western societies are built around road traffic, not participating would exclude them from social life.

<sup>5</sup> Gössling observes that the awareness of transport injustice is growing but at the same time recognizes a lack of activity from “transport governance”, which he calls the “implementation gap” (2016, p. 8). One could say there is also an unfulfilled responsibility of actors in transport governance to address this. However, it is not the interest of this paper to compare responsibilities rather focus on the designer’s obligation regarding AVs.

<sup>6</sup> When we talk about the trolley problem, we are not referring to the original thought experiment, much debated in moral philosophy (e.g., Thomson, 1976), but rather to the experimental ethics approach from the MIT moral machines experiment (Awad et al., 2018; Bonnefon, Shariff, & Rahwan, 2016).

<sup>7</sup> Through their experiments, they gathered 40 million decisions (Awad et al., 2018, p. 59).

According to that, it is discussed if it is ethically acceptable that AVs are programmed to choose the driving option with the expected minimal risk of corporal harm. The most severe concern is that such a programming conflicts with the fundamental rights of individuals (in the deontological ethics tradition) who are severely affected by the minimal risk decision. According to deontological ethics, it is prohibited, that individuals are sacrificed for the benefits of others. Choosing the minimal risk options in a dilemma-type situation does mean that the chance of very serious harm is also high for the chosen option, a potential loss of human life is intentionally accepted. At the same time, this means that the other options are not selected, and lives are definitely saved.

One possible ethical assessment is outlined by the German Ethics Commission (Di Fabio, Broy, & Brünger, 2017). The Commission argues that programming an AV for risk minimization is ethically acceptable if the mechanism of calculating the risk is the same for all involved parties, meaning that it does not consider any personal characteristics except of being human (no distinction based on personal features), especially ignoring if a person is an occupant of the automated vehicle. Since everybody has an equal chance, that the risk minimization algorithm is deciding to her disadvantage or her favor, minimizing is in the interest of everyone participating in traffic—the sacrificing concern does not hold anymore.

### Risk-ethics and automated routine driving

In this paper, we focus on the ethics of risk for routine driving. This entails situations with avoidable risks and decisions involving low to medium risk options. In routine driving, most of the time choices do not have immediate tangible negative consequences, risk does not materialize in a negative outcome. A routine driving situation can be a simple overtaking maneuver that, even when performed *safely* is not without risks, since there is always a small probability that unwanted events could happen. A car driver which overtakes a bicycle can take many different *safe* paths, with varying distances and velocities, and therefore differing in the scope of small risks. We can imagine another situation where car A refuses to break or change lane to let another car B enter the freeway easily. This might lead to a situation where the vehicle B reaches the end of the driveway and now has trouble to enter the freeway—vehicle B might have to take a higher risk to enter the freeway. Vehicle's A decision contributed to a possible *risky* situation for vehicle B even when vehicle A might be long gone when the situation might actually become dangerous.

In all these examples, the vehicle's driving behavior is transforming the space of moderate or low risk of itself and others. Very dangerous situations resulting from routine driving are rare relative to the absolute time on the road and

the number of decisions made without negative effect. It is often not a single wrong decision leading to an accident. It is more likely a combination of unfortunate events, for instance environmental influences like weather or dirt on road, with a short lack of attention by one or more drivers.<sup>8</sup> Keeping a larger distance to other road users, especially to vulnerable road users or low speed can make a difference in reducing chances for a situation to unluckily become hazardous. One can say, routine driving decisions often have an impact on the own and others' chances to suffer from bad luck.<sup>9</sup>

However, it is not clear if there is an ethical obligation/duty to address the risk exposure of others and to what extent this must be done, especially since the transformation of the risk space can be relatively small.

Emphasizing the ethical significance of automated behavior outside crashes is not new (Goodall, 2016, 2017; Himmelreich, 2018; Keeling, 2020). For instance, Goodall argues that the fair distribution of risk should be the leading ethical consideration not only during, but also “before forced choice situations” (2017, p. 496). However, these approaches do not differentiate well enough between decisions within and before unavoidable crashes. What is important during forced choices from a risk-ethics perspective was sketched above with reference to the deontological account—what is important for routine driving, we elaborate next.

### Risk-taking, risk-imposing and responsibility in the routine driving context

Before we discuss what is ethically relevant from an automated driving perspective with regards to risk-taking, risk-imposing and responsibility, we want to have a look at this from a human perspective.<sup>10</sup>

“In line with the liberal-democratic tradition, individuals are at liberty to take risks” (Ferretti, 2010, p. 505). Ferretti argues that justice/ethics becomes relevant when “negative consequences of risk-taking action [also] fall on third party” (2010, p. 506). Every driving, even a very passive driving style, can be considered as risk-taking and normally driving decisions also impose risks to others, whenever road users share a spatial area.

People's driving decisions, whether they are adhering to traffic regulations or not, are the result of their

<sup>8</sup> Even when accepting that AVs will reduce overall risk by eliminating human error in situ (not human error in the programming, designing, testing, etc.), there will still be situations which were not or could not be anticipated before and therefore the appropriate behavior is not determined.

<sup>9</sup> The role of luck, bad luck, and uncertainty in the ethical valuation of distributive practices has been discussed in the context of luck egalitarianism (see Ferretti, 2016; Nyholm, 2018).

<sup>10</sup> For instance, Himmelreich identifies risk-imposition as key topic for accessing the morality of normal driving in the AV context (2018).

individual explicit or implicit risk estimation, by weighing safety against some expected benefit (e.g., deciding to drive dangerously to be at the office in time).<sup>11</sup> As stated before, if it is only the risk-taker who would suffer harm, this would not be an ethical matter, but since possible negative consequences can also effect a third party it is an ethical question.

For what people are responsible is highly dependent on the circumstances and to what degree the person can anticipate the consequences of his decisions. For instance, in accident situations, humans decide intuitively and reflexively under strong time constraints (e.g., Lin, 2015; Nyholm & Smids, 2016). They may have to decide in split-seconds if they want to crash with the car in front of them or avoid the direct crash by steering into the oncoming traffic. Due to the time pressure and the limited possibilities to anticipate the consequences of their decisions, they cannot be held responsible for these decisions. However, for the antecedent behavior, which might have led to this situation, for instance if they were driving risky, they might be responsible.

Driving risky to serve one's own objectives can be considered as intended action, even when the actual negative consequences are not explicitly intended (van de Poel & Fahlquist, 2013, p. 116). The liberty to take risks is accompanied by the burden to be held responsible for actions and live with consequences (blameworthiness), which is based on the perspective of humans as moral agents (van de Poel & Fahlquist, 2013, p. 114).

For most routine driving humans are limited in their responsibility to act in a specific way, because they can only to some degree anticipate consequences of their behavior, especially regarding follow-up risk-imposition and in the moderated and low risk realm. For instance, the actual consequences of not doing a lane-change on another driver's' opportunity to enter the freeway and the impact on the follow-up risks can only be anticipated to some degree, given the time available.

In automated driving, the mechanisms for responsibility must be framed differently. How automated vehicles behave is partly predetermined by the choice of architecture and its underlying behavioral policies.<sup>12</sup> The designers are not present in situ, and they cannot claim that their decisions on the architecture and policies have been taken under a time-constraint. Furthermore, state-of-the-art automated driving systems have the capability to anticipate how a scene evolves and can compare between low and moderate risk options (see Sect. “Decentralized modular automated

driving”). In other words, AVs can predict the probability of events meaning how their routine behavior has an impact on the scene and how this influences the risk exposure of themselves and others. That is why we argue that designers are responsible for considering risk distribution in the development of architectures and driving policies from a justice or ethics standpoint for all AV behavior. In the next section, we explain why (decentralized) automated driving should be seen as institutional activity. Seeing automated driving as institutional practice, further supports the argument that designers are responsible to address unequal risk exposure toward a fairer risk distribution for routine driving.

### Decentralized automated driving as institutional activity

We consider AVs as entities which are organized in a decentralized way. Decentralized means that the AV's driving activity is primarily based on environmental sensing and its interpretation; vehicle-to-vehicle (V2V) and vehicle-to-infrastructure (V2I) connectivity is not required. Consequently, information about other road element's locations, activities and driving intentions are only accessible when sensed and interpreted.

This can be distinguished from a centralized automated driving paradigm in which AVs are considered highly interconnected and mainly coordinated on a group control level (see for example Mladenovic & McPherson, 2016 for a discussion on justice aspects from an AV traffic control perspective). In the future, automated driving might advance toward a more centralized infrastructure. However, it is still uncertain if the cars will be controlled by a *central* system or if the links are mainly used for the exchange of information. Given the number of competing car manufacturers it seems not very likely that there will be a single controller acting on all vehicles in a given scene. In today's traffic, advanced or semi-automated driver assistant systems are in use which fit more in the decentralized paradigm and at least for the near future, traffic will usually be mixed with automated and manually controlled vehicles interacting. Therefore, we argue that it is likely that decentralization continues to be a dominant concept—meaning that cars will be individually sold and controlled.

However, we must keep in mind that decentralized automated vehicles are not individuals in a human sense, rather they will drive the roads as one of many of the same or similar kind due to mass production. It is likely that manufacturers will apply the same optimized behavior selection algorithms as well as sensing and scene understanding technologies throughout all their vehicle line-up, or will be using third-party AV components which follow common logics and rules. Liu states, given such a constellation, that these “algorithmic policy preferences” (which might

<sup>11</sup> Any intended risk-taking or -acceptance must be viewed regarding some expected benefits or opportunity (Ferretti, 2010, p. 505).

<sup>12</sup> This is true for modular AV architectures which follow *explicit* policies and often contain dedicated risk management mechanisms. These can be distinguished from AI-driven black-box architectures. See more in the introduction to Sect. “Decentralized modular automated driving”.

be implicitly or explicitly incorporated) have a collective impact since individuated outcomes are in this way “multiplied across fleets of vehicles” (Liu, 2018, p. 150). According to the above, individual outcomes are aggregated and so the collective effects can be considered as systematic. This illustrates that local driving decisions of decentralized AVs, which only impact the risk allocation in the immediate time–space have a systemic or “collective dimension” when broadly applied (Liu, 2018, p. 161). In reference to this thought Dietrich argues for framing automated driving as institutional activity (Dietrich, 2020). He emphasizes the similarity to the activity schema of (social) institutions as both are held together by explicit and implicit policies and individual entities (proxy agents) as executive bodies with a certain degree of freedom in applying the policies in a given situation (2020, p. 4).

According to that, automated vehicles’ risk-taking decisions and how to balance risk with utility factors are of institutional quality. Since most driving decisions will also influence the exposure to risk of third parties, distributing risk among road users involved is an institutional justice matter. Following this line of thought, it is a manufacturer’s responsibility (as obligation) to actively shape the distribution of risk exposure among road users by integrating fairness objectives.

In the following section, we show a possible direction of how to integrate fairness objectives for modular decentralized automated driving. We will see that in future mixed-traffic situations, where vehicles share the road with non-motorized road users, AVs might be able to act more cautiously relative to the safety needs of other road users, reduce inequalities, but will be limited in creating full equality with regards to risk-exposure. The reason is that because of the strong asymmetry in protection there will be always risk-imposition towards non-motorized road users (and not vice versa), as long as we expect that AVs operate similar to human-driven cars. The follow-up question is, if low or moderate risk-imposition, can be acceptable from a justice standpoint. Ferretti argues from an institutional point of view that a certain risk-taking action affecting others might be acceptable for a just society. But this risk-taking action requires justification “based on the countervailing expected good consequence of taking a risk” (Ferretti, 2010, p. 506). For instance, it is considered as acceptable if the risk-imposition belongs to a “social system of risk-taking” which also works to the advantage of the risk-bearer, e.g., when she also takes the role as risk-imposer in other situations (Hansson, 2003, p. 305). On a formal level, justification requires a consideration of values and interests which should be explainable and

communicable if requested.<sup>13</sup> Such justification claims are also relevant from a design requirements standpoint, since architectures, for instance a modular compared to a deep learning architecture (see: introduction to the next section), differ in their capability to explicitly manage risks for behavior planning and consequently their capability to give information about how risk is taken.

## Decentralized modular automated driving

In this section, we discuss features of one typical architecture of an AV operating system, which includes a dedicated module for risk estimation and management (similar to e.g., Probst et al., 2021; Tas et al., 2017; Weisswange et al., 2019). We show how this can be adapted to address inequalities in risk exposure. This architecture represents a decentralized and modular approach; modularity is defined as containing different separable components. A modular approach allows to make individual design decisions for every component. An alternative to a component-based architecture is a deep AI or End2End learning approach (for a survey on these techniques see e.g., Grigorescu et al., 2020). The dominant concept is to utilize large amounts of real-world driving data which allows to directly map sensory information to steering commands. AI-based approaches usually require little domain knowledge. At the same time little direct control on outcomes is possible, which makes it inappropriate to explicitly incorporate justice principles and therefore, we are not considering it here.

The main contribution of this paper is to discuss opportunities as to how a decentralized, modular AV architecture and its parameterization can be setup to address the uneven exposure to risk among road users.

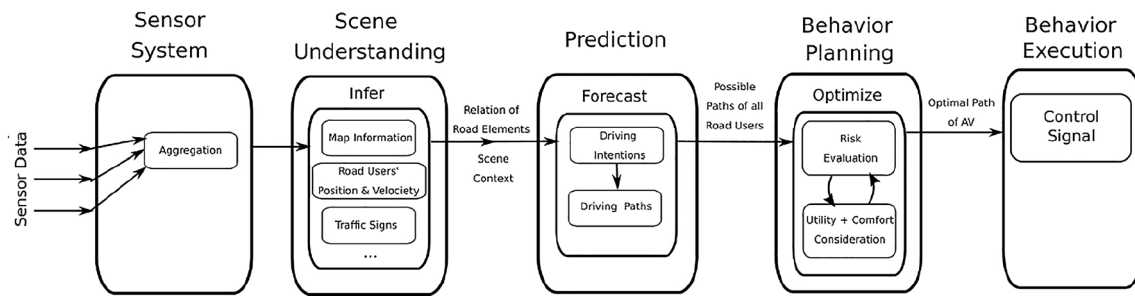
### Key components

The architecture for decentralized automated driving typically contains five components (as shown in Fig. 1). The sensor system provides the vehicle with observations of its surrounding. This can be any combination of radar, lidar, camera and ultrasonic sensors. These measurements will be interpreted into an understanding of the current situation. AI or machine learning algorithms allow a detection and tracking of scene elements, like road participants, lane markings or traffic signs, to achieve a semantic scene understanding.

<sup>13</sup> To get an idea, how a justification could look like, we can refer to the GDPR: in specific cases, data controllers must justify the processing of personal data, i.e., when they have no consent. They must justify that the “processing is necessary for the purpose of the legitimate

Footnote 13 (continued)

interests”; and balance it against the interests, rights and freedoms of persons affected [see: GDPR, Article 6(1)(f)].



**Fig. 1** Schematic architecture of a decentralized, modular automated driving approach. Adapted from Eggert, Klingelschmitt and Damerow (2015)

The information about static and dynamic traffic elements, including their locations and relation to each other as well as behavior profiles (velocities), are used to predict how the scene will most likely evolve in the near future. In many situations, multiple futures are equally possible. For instance, a car which is approaching a T-intersection can either turn left or right. A future scene prediction might also contain the other road users' probable behaviors in response to the possible behaviors of the AV.

In the behavior planning stage, the optimal path of the AV will be selected. A risk evaluation module predicts the probabilities that the execution of a certain driving maneuver will evolve into an unwanted event. Unwanted events are in particular collisions between traffic participants or with objects in the surrounding. The risk evaluation has to be balanced with utility and comfort considerations. How risk is modelled, estimated, and balanced, also has an impact on other road users' exposure to risk. We discuss that in more detail in the next section. Finally, for behavior execution, the optimal path is selected and translated into a control signal.

### Risk modelling for behavior planning

Risk refers to an unwanted event which might appear in the future. Risk is formally defined as the product of the probability and severity of an unwanted event (e.g., see Hansson, 2018). To model risk, future unwanted events must be predicted and their likelihood of occurrence has to be quantified. In the automated driving context, the evolution of the traffic scene is forecasted based on environmental sensing and scene interpretation. The objective of a risk estimation component is to predict hazard events in particular collisions with other road users or static objects. Risk estimation is used as a continuous measure that can represent the influence of AV behaviors on future event probabilities.<sup>14</sup>

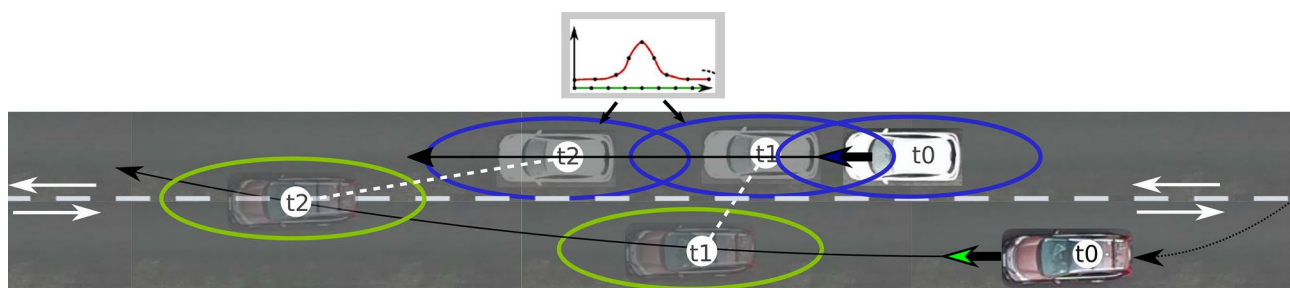
Continuous risk-measures can support the AV in finding the optimal path, given its estimation and modelling methods, and updating the path when the circumstances change. This is equally useful in high-risk *crash situations* as in routine driving, as any behavior will usually come with some, albeit small, risk and can therefore be optimized towards its minimum.

### Event probability

The simplest way to predict collisions is to look for a geometric overlap of two or more road objects any time in the future given the current situation and the likely motions of participants in the scene. However, there is a high number of possibilities how a scene might evolve since AVs cannot know with certainty about other road users' intentions. The observed vehicles can, for instance, change their current velocity and take multiple routes when approaching an intersection. Furthermore, even actual positions and velocity profiles can only be estimated with some uncertainty, because of inaccuracies due to sensory noise and detection errors. An adequate risk estimation has to account for these variances. Figure 2 shows an example situation, where an AV (red car) is overtaking the white vehicle. The future position of the other car is predicted based on its current position and a constant velocity profile. The ellipses illustrate the uncertainties about the AVs and other vehicle's position. Modelling uncertainty for instance with a Gaussian function means that the car can be at any position within the ellipse whereby the center position is more probable than towards the borders of the ellipse (e.g., Schreier, Willert, & Adamy, 2016; Puphal, Probst, & Eggert, 2019). The probability of a collision increases, the more the position distributions overlap.

<sup>14</sup> Such risk modelling for AV behavior planning is e.g., done by Eggert (2014) and Eggert, Klingelschmitt & Damerow (2015).





**Fig. 2** The graphic illustrates an overtaking maneuver from a birds-eye view. The red AV has started its overtaking maneuver of the white vehicle on a road with two-way traffic. The future positions of

the AV itself and the other car are predicted including its uncertainties, illustrated by the green and blue ellipses. (Color figure online)

## Severity

Modeling the severity of an unwanted event is a challenging task. That is also why many approaches in the context of automated driving incorporate abstracted policies to avoid any collision, independent of object type or velocity.<sup>15</sup>

So, what does it take to make use of severity as one component to estimate collision risks? First, the harm to humans and property damage of a potential collision has to be modelled. Second, the different types of potential harm have to be mapped to a common scale to make them comparable.

To figure out what harm, e.g., injuries, fatalities, or property damage, is expected by a predicted crash, data from real crash data could be used. For instance, in the US (e.g., Michigan Traffic Crash Facts) and Germany (German In-Depth Accident Study) organizations document crash cases. This includes, among other things, information about injury level and the road users involved. However, a *data-driven* approach is limited in informative value due to the fact that crashes are rare and even fewer are well documented (see: Eggert, 2018, p. 121).

When severity is modelled in the context of automated driving, most of the work follows a *theory-driven approach*. One possibility is to assume that a crash can be described as a two-dimensional encounter between masses whereby the kinetic energy involved acts on the human body and correlates with the harm (Chen, Yang, & Otte, 2010; Probst et al., 2021; Puphal et al., 2018). The kinetic energy is proportional to the operating masses and velocities. The intensity of the impact is also dependent on the angle with which the masses collide (relative direction of the movement). Rear-end, front or side impact crashes have different angles and so might differ in severity. When looking to the literature, we can describe the relation between human injury and the kinetic energy in an accident, as logistic function (Puphal et al.,

2018, p. 1708). This is still a simplified model of severity as, among other things, it is only two-dimensional. For instance, it has been found that cyclists suffer a lower injury outcome compared to pedestrians for the same accident due to their higher position (Chen, Yang, & Otte, 2010).

That brings us to the second question how to rate the different types of damages (potential harm to humans in relation to expected property damage) in order to use them for driving decisions. Two approaches are possible: Injuries or fatalities can be factored in as monetary costs, for instance the economic value of a person's life or healthcare and the costs of treating injuries. That would allow it to be compared with expected property damage. However, setting a monetary value for life could be seen as ethically problematic.<sup>16</sup> The second option is not to refer to any actual value rather work with abstract numbers and weight the potential harm to humans several orders of magnitude higher than any property damage. This prevents the driving system to favor an option where humans are harmed as opposed to any damage to property.

## Balancing risk against utility and comfort values

AVs should be designed to avoid risk for themselves (and their passengers), but this cannot be the only criterion for path planning. It is expected that AVs operate without violating traffic rules and keep risks for themselves and others low. However, when planning exclusively by minimizing risk, many driving decisions will end up in no movement at all, because other road users and the uncertainties in sensing and prediction will always impose some risk that usually scales with the amount of movement (see: Eggert, 2018, p. 129).

A behavior planning system should therefore incorporate beside risk also utility variables (e.g., progress towards

<sup>15</sup> For instance, when using the time-to-contact metrics (TTC) for risk estimation (e.g., Ward et al., 2014).

<sup>16</sup> This could mean that human life could be sacrificed in the name of other goods which is seen as a violation of human dignity (e.g., Horizon 2020 Commission Expert Group, 2020, p. 21).

a destination with a certain comfort). To compare different driving options regarding all criteria relevant for driving, a so-called cost function is used. Unsafe paths cause high costs, and an efficient and comfortable driving flow has low costs. For the parametrization of such a cost function, it is important that utilities and comfort have no impact on the course of the vehicle when the AV faces a crash situation; collision risk must become the dominant driving force. Due to that it is a common approach for the parametrization of the cost function to weight the costs of a significant collision risk at least one magnitude higher than the costs of strong discomfort and low utility (Probst et al., 2021).

## Addressing unequal exposure in AV development

In this section we discuss if and to what extent, design choices in decentralized AV's risk estimation and decision-making modules, can address unequal exposure. The way, in which the risk of potential encounters between AVs and other road users is modelled, especially regarding severity, has an impact on driving decisions and so on the risk exposure of those involved.

### Approach: severity-sensitive risk estimation including others

Risk estimation is seen as a measure to allow the vehicle to make decisions to find a safe path in complex traffic situations. In current AV approaches, we often see that severity is modelled as being constant, meaning that every crash has the same severity (e.g., TTC methods for risk estimation, as referred to in the previous section). Such a simple severity model does not allow to decide between different paths of collision probabilities, where one potential accident would be more severe. A kinematic severity modelling, as we have also presented it above, considers that the velocities, the angle of the encounter as well the involved masses have an impact on the expected harm.

However, even when incorporating a more advanced severity-sensitive model to estimate risks, current AV approaches consider risk mainly as risk for the AV and its passengers. This means risk is predicted from the AV's point of view and only considers other road users as obstacles which can cause harm to the AV and its occupants.

In the interest of a fair risk distribution, we argue, that a risk estimation for the AV's behavior planning component must not only refuse a simple severity modelling towards a more advanced rather should also include the other road users in any given scene. One way to do so is to calculate the severity of a potential collision, not only considering

the effects for the AV and its passengers, but also the potential negative impact on all other people involved. Being able to identify other road users via environmental sensing also allows to estimate their probable severities in case of an encounter with the AV.

We sketch one approach how to take the other road users into account explicitly. The focus is on how to calculate the severity of potential encounters when multiple entities are involved which might bear differently from the negative implications.

Given two vehicles, it might be straight-forward to take the average severity of the two entities that might be involved in a probable accident. According to such an averaging model, calculating the severity of an encounter is performed by summing up the individual severities and dividing them by the number of entities involved.

As we have described in the beginning, crashes can be highly asymmetric with regards to the negative implications for the colliding parties. A car passenger will likely be unharmed when colliding with a pedestrian at 60 km/h, but the pedestrian will most likely be seriously injured. Regarding the example, applying an average model would underestimate the vulnerability of the pedestrian in such asymmetric constellations in relation to expected encounters between symmetric road users.

When there are entities with little expected harm involved together with high severity candidates, the involvement of an advantaged (well protected) entity is *watering down* the overall severity of a possible encounter. Certainly, a severity calculation must be assessed in relation to other constellations. A possible encounter in a symmetric constellation, e.g., between two cars, which both expect a medium severe outcome might then be rated as high as the mutual severity in the asymmetric constellation.

That is the reason why we think such an average model conflicts with the claim to especially take into account vulnerable road users (Horizon 2020 Commission Expert Group, 2020, p. 31). It seems to be more appropriate, to always give full priority to the entity which has the highest expected harm. To formalize this, the severity of a potential encounter between two or more entities can be calculated by predicting the severity for each of the entities separately and then selecting the maximum value.

$$S = \max (S_{AV}, S_{v1}, S_{v2}, \dots, S_{vn}).$$

In this way, the entity with highest risk has a strong effect on the risk-based behavior planning—the consequence is that cautious behaviors around vulnerable road users will become more likely.

Using this basic proposal, we highlight some aspects relevant for further discussions in the next paragraph.

## Discussion

Applying a risk estimation metric which always gives priority to entities with the highest expected harm is in accordance with prior approaches on fair risk distribution for AVs, adopted from John Rawls—mostly in an immediate high-risk context (Leben, 2017). According to Rawls, an institutional allocation policy ought to be chosen, which is to be to the greatest benefit of the least advantaged members of society—which he referred to as the difference principle. It can be seen as a form of social solidarity (Rawls, 1971, p. 90). The difference principle was originally designed for a society-wide distribution of primary goods. The arguments which support the difference principle for the distributions of social goods, seem also to be valid for the distribution of low and moderate risks—especially since we introduced risks as a joint burden of road traffic as social practice. Incorporating the principle in the risk estimation architecture, as we propose it, directly affects the in situ driving behavior but it also has a more systemic impact when effects are multiplied across a fleet of vehicles. Seeing it this way, we can frame the risk estimation and management mechanisms as a kind of institutional allocation policy, for which Rawls' theory was originally designed for.<sup>17</sup>

In Fig. 2, a situation in which the AV is overtaking another slower vehicle was illustrated. We could imagine an alternative scene in which the white car is replaced by a cyclist. The corporal damage of a collision of the AV with a cyclist is significantly higher than with another vehicle, as we would see in the injustice discussion above, the fatality rate between cars and vulnerable road users is very high even at relatively small velocities. An overtaking maneuver with the same distance and velocity will be considered as riskier for the cyclist case; in technical terms, it will be a path with higher costs. To compensate this, the overtaking maneuver could be performed at a slower speed, with more lateral space or not at all.

Applying such a relative decision schema means judging one constellation relative to another and acting accordingly. This is reasonable from a fair distribution standpoint since “fairness is [...] a matter of how one candidate is treated relative to others” (Broome, 1984, p. 43). How one candidate is treated in relation to the other depends on the legitimate claims they each have on the public good, in the driving case this is safety or the absence of risk. Those who have more safety needs (i.e., are more vulnerable) must to be given priority relative to others.

Another aspect which is relevant to discuss is if and how the other road users' risk-taking and risk-imposing behavior should be considered in the AV's risk avoidance behavior (see discussion above). In theory, there are observable factors of other road users' attitudes towards risk-taking and risk-acceptance. For instance, if a motorcyclist is driving above the speed-limit, it can be inferred, that he has a higher risk-acceptance. Since the expected severity of a crash between the AV and the speeding motorcycle is very high, primarily for the motorcyclist, a severity-aware risk estimation *forces* the AV to keep a high distance. In this hypothetical situation, keeping distance to the motorcycle might bring the AV closer to a truck and so increases the risk of the AV's occupants compared to the same situation with a speed-limit compliant motorcyclist. Since the motorcyclist has an observable high risk-acceptance, it might be reasonable to *ignore* this higher risk for the motorcycle rider, which results from its speeding behavior.

## Conclusion

A dedicated risk estimation together with a behavior planning module allows to make explicit choices about how the AV occupants and other road users are exposed to risk in each driving scene. We have shown that it is technically feasible to adjust state-of-the-art approaches towards a continuous in situ severity-sensitive risk estimation which considers others' situated risks following fair distribution principles.

Addressing the distribution of risk before an accident situation has an impact on the chance of road users getting into accidents and is therefore an important variable to address unequal risk exposure. We have argued that it is the designers' responsibility to adapt the AV architecture to incorporate fair risk distribution objectives not only in accident cases but also in routine driving.

However, the proposed fair distribution approach which gives priority to entities with the highest expected harm, does not result in a numerically equal distribution of risk between all categories of road users. Even when vulnerable road users profit the most, they will always, due to the strong asymmetry in protection, be more exposed to risk than the AV, as long as we expect that AVs integrate in the traffic flow and follow utility objectives, similar to human-driven cars (no very passive driving or a full-stop policy). We have argued that, even being in the spectrum of low and moderated risk, routine driving can be seen as risk-taking activity which affects others. According to Ferretti this requires justification (Ferretti, 2010) through being able to inform affected people about incorporated risk-utility tradeoffs and the reasons behind these. In this regard an AV architecture with a dedicated risk estimation module is preferable to black-box approaches of automated driving, since it allows

<sup>17</sup> A further discussion on the ethical justification why Rawls' difference principle, also called “maximin” rule, serves a good guidance for situated decisions of ethical AVs, can be found in the literature (Dietrich & Weisswange, 2019; Leben, 2017).

both to explicitly control risk related trade-offs and to give information about the risk assessment for driving decisions.

**Author contributions** Not applicable.

**Funding** This research is funded by the Honda Research Institute Europe.

**Data availability** Not applicable.

**Code availability** Not applicable.

## Declarations

**Conflict of interest** Not applicable.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- AI High Level Expert Group. (2019). *Ethics guidelines for trustworthy AI*. Publication Office of the European Union.
- Awad, E., Dsouza, S., Kim, R., Schulz, J., Henrich, J., Shariff, A., Bonnefon, J.-F., & Rahwan, I. (2018). The Moral Machine experiment. *Nature*, *563*, 59–64. <https://doi.org/10.1038/s41586-018-0637-6>.
- Bonnefon, J. F., Shariff, A., & Rahwan, I. (2016). The social dilemma of autonomous vehicles. *Science*, *352*(6293), 1573–1576. <https://doi.org/10.1126/science.aaf2654>.
- Broome, J. (1984). Selecting people randomly. *Ethics*, *95*(1), 38–55.
- Chen, Y., Yang, J., & Otte, D. (2010). Load and impact conditions for head injuries in car-to-pedestrian and car-to-cyclist accidents. *Proceedings of the Expert Symposium on Accident Research*, 294–308.
- Di Fabio, U., Broy, M., & Brünger, R. J. (2017). *Ethics Commission: Automated and connected driving*. Federal Ministry of Transport and Digital Infrastructure of the Federal Republic of Germany.
- Dietrich, M. (2020). Understanding autonomous driving as institutional activity: Opening new ways to react to discriminatory concerns in autonomous driving. In M. Nørskov, J. Seibt, & O. S. Quick (Eds.), *Culturally sustainable social robotics: Proceedings of Robophilosophy 2020* (pp. 335–373).
- Dietrich, M., & Weisswange, T. H. (2019). Distributive justice as an ethical principle for autonomous vehicle behavior beyond hazard scenarios. *Ethics and Information Technology*, *21*, 227–239. <https://doi.org/10.1007/s10676-019-09504-3>.
- Eggert, J. (2014). Predictive risk estimation for intelligent ADAS functions. *Proceedings of the 17th international IEEE conference on intelligent transportation systems (ITSC)*, 711–718. <https://doi.org/10.1109/ITSC.2014.6957773>.
- Eggert, J. (2018). Risk estimation for driving support and behavior planning in intelligent vehicles. *Automatisierungstechnik*, *66*(2), 119–131. <https://doi.org/10.1515/auto-2017-0132>.
- Eggert, J., Klingelschmitt, S., & Damerow, F. (2015). The foresighted driver: Future ADAS based on generalized predictive risk estimation. *Proceedings of the FAST-zero 2015 symposium*, 93–100.
- Ewing, R., & Dumbaugh, E. (2009). The built environment and traffic safety: A review of empirical evidence. *Journal of Planning Literature*, *23*(4), 347–367. <https://doi.org/10.1177/0885412209335553>.
- Ferretti, M. P. (2010). Risk and distributive justice: the case of regulating new technologies. *Science and Engineering Ethics*, *16*(3), 501–515. <https://doi.org/10.1007/s11948-009-9172-z>.
- Goodall, N. J. (2016). Away from trolley problems and toward risk management. *Applied Artificial Intelligence*, *30*(8), 810–821. <https://doi.org/10.1080/08839514.2016.1229922>.
- Goodall, N. J. (2017). From trolleys to risk: Models for ethical autonomous driving. *American Journal of Public Health*, *107*, 496–496. <https://doi.org/10.2105/AJPH.2017.303672>.
- Gössling, S. (2016). Urban transport justice. *Journal of Transport Geography*, *54*, 1–9. <https://doi.org/10.1016/j.jtrangeo.2016.05.002>.
- Grigorescu, S., Trasnea, B., Cocias, T., & Macesanu, G. (2020). A survey of deep learning techniques for autonomous driving. *Journal of Field Robotics*, *37*, 362–386. <https://doi.org/10.1002/rob.21918>.
- Hansson, S. O. (2003). Ethical criteria of risk acceptance. *Erkenntnis*, *59*, 291–309. <https://doi.org/10.1023/A:1026005915919>.
- Hansson, S. O. (2018). “Risk”. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*, Fall 2018 Edition. <https://plato.stanford.edu/archives/fall2018/entries/risk/>.
- Himmelreich, J. (2018). Never mind the trolley: The ethics of autonomous vehicles in mundane situations. *Ethical Theory and Moral Practice*, *21*(3), 669–684. <https://doi.org/10.1007/s10677-018-9896-4>.
- Horizon 2020 Commission Expert Group to Advise on Specific Ethical Issues Raised by Driverless Mobility. (2020). *Ethics of Connected and Automated Vehicles: Recommendations on road safety, privacy, fairness, explainability and responsibility*. Publication Office of the European Union.
- IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems. (2019). *Ethically aligned design: A vision for prioritizing human well-being with autonomous and intelligent systems* (1st ed.). IEEE. <https://standards.ieee.org/content/ieee-standards/en/industry-connections/ec/autonomous-systems.html>.
- Jones, I. (2014). *Road space allocation: The intersection of transport planning, governance and infrastructure*. Doctoral Dissertation, RMIT University.
- Keeling, G. (2017). Commentary: Using virtual reality to assess ethical decisions in road traffic scenarios: Applicability of value-of-life-based models and influences of time pressure. *Frontiers in Behavioral Neuroscience*, *11*, 247.
- Keeling, G. (2020). *The ethics of automated vehicles*. Doctoral Dissertation, University of Bristol.
- Leben, D. (2017). A Rawlsian algorithm for autonomous vehicles. *Ethics and Information Technology*, *19*(2), 107–115. <https://doi.org/10.1007/s10676-017-9419-3>.
- Lin, P. (2015). Why ethics matters for autonomous cars. In M. Maurer, J. C. Gerdes, B. Lenz, & H. Winner (Eds.), *Autonomous driving: Technical, legal and social aspects* (pp. 79–85). Springer.
- Liu, H. Y. (2018). Three types of structural discrimination introduced by autonomous vehicles. *UC Davis Law Review Online*, *51*, 149–180.
- Martens, K. (2017). *Transport justice: Designing fair transportation systems*. Routledge.
- Mladenovic, M. N., & McPherson, T. (2016). Engineering social justice into traffic control for self-driving vehicles? *Science and*

- Engineering Ethics*, 22(4), 1131–1149. <https://doi.org/10.1007/s11948-015-9690-9>.
- Mullen, C., Tight, M., Whiteing, A., & Jopson, A. (2014). Knowing their place on the roads: What would equality mean for walking and cycling? *Transportation Research Part a: Policy and Practice*, 61, 238–248. <https://doi.org/10.1016/j.tra.2014.01.009>.
- Nello-Deakin, S. (2019). Is there such a thing as a ‘fair’ distribution of road space? *Journal of Urban Design*, 24(5), 698–714. <https://doi.org/10.1080/13574809.2019.1592664>.
- Nyholm, S. (2018). The ethics of crashes with self-driving cars: A roadmap. II. *Philosophy Compass*, 13(7), 1–10. <https://doi.org/10.1111/phc3.12506>.
- Nyholm, S., & Smids, J. (2016). The ethics of accident-algorithms for self-driving cars: An applied trolley problem? *Ethical Theory and Moral Practice*, 19(5), 1275–1289. <https://doi.org/10.1007/s10677-016-9745-2>.
- Padmanaban, J. (2003). Influences of vehicle size and mass and selected driver factors on odds of driver fatality. *Annual Proceedings of the Association of Advancement of Automotive Medicine*, 47, 507–524.
- Probst, M., Wenzel, R., Puphal, T., Komuro, M., Weisswange, T. H., Steinhardt, N., Bolder, B., Flade, B., Sakamoto, Y., Yasui, Y., & Eggert, J. (2021). Automated driving in complex real-world scenarios using a scalable risk-based behavior generation framework. *Proceedings of the 24th IEEE International Intelligent Transportation Systems Conference (ITSC 2021)*, Indianapolis, IN, USA.
- Puphal, T., Probst, M., & Eggert, J. (2019). Probabilistic uncertainty-aware risk spot detector for naturalistic driving. *IEEE Transactions on Intelligent Vehicles*, 4(3), 406–415. <https://doi.org/10.1109/TIV.2019.2919465>.
- Puphal, T., Probst, M., Li, Y., Sakamoto, Y., & Eggert, J. (2018). Optimization of velocity ramps with survival analysis for intersection merge-ins. *Proceedings of the IEEE intelligent vehicles symposium (IV)*, 1704–1710. <https://doi.org/10.1109/IVS.2018.8500667>.
- Rawls, J. (1971). *A theory of justice*. Harvard University Press.
- Ross, B. (2014). Injustice at the intersection. *Dissent Magazine*. Retrieved September 16, 2021, from [https://www.dissentmagazine.org/online\\_articles/injustice-intersection-suburbs-traffic-engineering-poverty](https://www.dissentmagazine.org/online_articles/injustice-intersection-suburbs-traffic-engineering-poverty).
- Schäffner, V. (2020). Wenn Ethik zum Programm wird: Eine risikoethische Analyse moralischer Dilemmata des autonomen Fahrens. *Zeitschrift Für Ethik Und Moralphilosophie*, 3(1), 27–49. <https://doi.org/10.1007/s42048-020-00061-9>.
- Schreier, M., Willert, V., & Adamy, J. (2016). An integrated approach to maneuver-based trajectory prediction and criticality assessment in arbitrary road environments. *IEEE Transactions on Intelligent Transportation Systems*, 17, 2751–2766. <https://doi.org/10.1109/TITS.2016.2522507>.
- Tas, Ö. S., Hörmann, S., Schäufele, B., & Kuhn, F. (2017). Automated vehicle system architecture with performance assessment. *Proceedings of the 2017 IEEE 20th international conference on intelligent transportation systems (ITSC)*, 1–8. <https://doi.org/10.1109/ITSC.2017.8317862>.
- Thomson, J. J. (1976). Killing, letting die, and the trolley problem. *The Monist*, 59(2), 204–217. <https://doi.org/10.5840/monist197659224>.
- van de Poel, I., & Fahlquist, J. N. (2013). Risk and responsibility. In S. Roeser, R. Hillerbrand, P. Sandin & M. Peterson (Eds.), *Essentials of Risk Theory* (pp. 107–143). Springer. [https://doi.org/10.1007/978-94-007-5455-3\\_5](https://doi.org/10.1007/978-94-007-5455-3_5).
- Van de Poel, I., & Sand, M. (2021). Varieties of responsibility: Two problems of responsible innovation. *Synthese* 198, 4769–4787. <https://doi.org/10.1007/S11229-018-01951-7>.
- Ward, J., Agamennoni, G., Worall, S., & Nebor, E. (2014). Vehicle collision probability calculation for general traffic scenarios under uncertainty. *Proceedings of the 2014 IEEE intelligent vehicles symposium (IV)*, 986–992. <https://doi.org/10.1109/IVS.2014.6856430>.
- Weisswange, T. H., Rebhan, S., Bolder, B., Steinhardt, N.A., Joublin, F., Schmuedderich, J., & Goerick, C. (2019). Intelligent traffic flow assist: Optimized highway driving using conditional behavior prediction. *IEEE Intelligent Transportation Systems Magazine*, 13(2), 20–38. <https://doi.org/10.1109/MITS.2019.2898969>.
- Zegeer, C. V., Seiderman, C., Lagerwey, P., Cynecki, M., Ronkin, M., & Schnieder, R. (2002). *Pedestrian facilities users guide—Providing safety and mobility*. Report No. FHWA-RD-102–01. Federal Highway Administration. <https://trid.trb.org/view/673359>.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.