

Drivers are blamed more than their automated cars when both make mistakes

Edmond Awad^{a,b,+}, Sydney Levine^{a,c,d,+}, Max Kleiman-Weiner^{c,d}, Sohan Dsouza^a,
Joshua B. Tenenbaum^{c,*}, Azim Shariff^{e,*}, Jean-François Bonnefon^{a,f,*}, and Iyad
Rahwan^{a,g,h,*}

^aMedia Lab, Massachusetts Institute of Technology, Cambridge MA, USA

^bDepartment of Economics, University of Exeter Business School, Exeter, UK

^cDepartment of Brain and Cognitive Sciences, Massachusetts Institute of
Technology, Cambridge MA, USA

^dDepartment of Psychology, Harvard University, Cambridge, MA, USA

^eDepartment of Psychology, University of British Columbia Vancouver, Canada

^fToulouse School of Economics (TSM-Research), Centre National de la Recherche
Scientifique, University of Toulouse Capitole, Toulouse, France.

^gInstitute for Data, Systems and Society (IDSS), Massachusetts Institute of
Technology, Cambridge MA, USA

^hCentre for Humans & Machines, Max-Planck Institute, Berlin, Germany

⁺Joint first author

*To whom correspondence should be addressed. e-mail: jbt@mit.edu;
shariffa@uci.edu; jean-francois.bonnefon@tse-fr.eu; irahwan@mit.edu

Abstract

When an automated car harms someone, who is blamed by those who hear about it? Here, we asked human participants to consider hypothetical cases in which a pedestrian was killed by a car operated under shared control of a primary and a secondary driver, and to indicate how blame should be allocated. We find that when only one driver makes an error, that driver is blamed more, regardless of whether that driver is a machine or a human. However, when both drivers make errors in cases of human-machine shared-control vehicles, the blame attributed to the machine is reduced. This finding portends a public under-reaction to the malfunctioning AI components of automated cars and therefore has a direct policy implication: allowing the de-facto standards for shared-control vehicles to be established in courts by the jury system could fail to properly regulate the safety of those vehicles; instead, a top-down scheme (through federal laws) may be called for.

Introduction

Every year, about 1.25 million people die worldwide in car crashes [1]. Laws concerning principles of negligence currently adjudicate how responsibility and blame get assigned to the individuals who injure others in these harmful crashes. The impending transition to fully automated cars promises a radical shift in how blame and responsibility will get attributed in the cases where crashes do occur, but most agree that little or no blame will be attributed to the occupants in the car, who will, by then, be entirely removed from the decision-making loop [2]. However, before this era of fully automated cars arrives, we are entering a delicate era of shared control between humans and machines.

This new moment signals a departure from our current system – where individuals have full control over their vehicles and thereby bear full responsibility for crashes (absent mitigating circumstances) – to a new system where blame and responsibility may be shared between a human and a machine driver. The spontaneous reactions people have to crashes that occur when a human and machine share control of a vehicle has at least two direct industry-shaping implications. First, at present, little is known about how the public is likely to respond to crashes that involve both human and machine drivers. This uncertainty has concrete implications: manufacturers price products to reflect the liability they expect to incur from the sale of those products. If manufacturers cannot assess the scope of the liability they will incur from automated vehicles (AVs), that uncertainty will translate to substantially inflated prices of AVs [2]. Moreover, the rate of the adoption of automated vehicles will be proportional to the cost to consumers to adopt the new technology [2]. (The rate of the adoption of this technology is contingent on many other factors, including consumers’ understanding of the relative risks and benefits of using the cars. We do not mean to state that uncertainty about the scope of liability for manufacturers is the only factor impacting adoption, just that it is a significant one.) Accordingly, the uncertainty about the extent of corporate liability for automated vehicle crashes may be slowing down AV adoption [2], while people continue to die in car crashes each year. Clarifying how and when responsibility will be attributed to manufacturers in automated car crashes will be a first step in reducing this uncertainty and speeding the adoption of automated and eventually fully automated vehicles.

The second direct implication of this work will be to forecast how a tort-based regulatory scheme (which is decided on the basis on jury decisions) is likely to turn out. Put another way, understanding how the public is likely to react to crashes that involve both a human and a machine driver will give us a hint to what standards will be established if we let jury decisions shape them. If our work uncovers systematic biases that are likely to impact juries and would impede the adoption

of automated cars, then it may make sense for federal regulations be put in place, which would preempt the tort system from being the avenue for establishing standards for these cars.

Already, automated vehicle crashes are in the public eye. In May 2016, the first deadly crash of a Tesla Autopilot car occurred and the occupant of the car was killed. In a news release, Tesla explained: “Neither Autopilot nor the driver noticed the white side of the tractor-trailer against a brightly lit sky, so the brake was not applied” [3]. In March, 2018, the first automated car crash that killed a pedestrian occurred. A pedestrian that was crossing the street went unnoticed by both the car and the back-up driver. A few seconds before the crash, the car finally identified that it should be breaking but failed to do so. The driver also braked too late to avoid the collision.

In both the Tesla and Uber fatal crashes, both the machine driver and the human driver should have taken action and neither driver did. The mistakes of both the machine and the human led to the crash. The National Highway Safety Traffic Administration (NHSTA) did an investigation of the Tesla incident and did not find Tesla at fault in the crash [4]. Likewise, Uber has been exonerated from criminal charges after an investigation by a county prosecutor [5]. Notably, press attention surrounding the Tesla incident was markedly skewed towards blaming the human driver for the crash, with rumors quickly circulating that the driver had been watching a Harry Potter movie [6], though upon further investigation it was discovered that there was no evidence grounding this claim [7]. Likewise, the fate of the Uber back-up driver remains unknown [5], with press attention focusing on the distracted nature of the driver [8].

The set of anecdotes around these two crashes begins to suggest a troubling pattern, namely, that humans might be blamed more than their machine partners in certain kinds of automated vehicle crashes. Was this pattern a fluke of the circumstances of the crash and the press environment? Or does it reflect something psychologically deeper that may color our responses to human-machine joint action, and in particular, when a human-machine pair jointly controls a vehicle?

What we are currently witnessing is a gradual and multipronged increase toward full automation, going through several steps of shared control between user and vehicle, which may take decades due to technical and regulatory issues as well as attitudes of consumers towards adoption [9, 10] (see Figure 1). Some vehicles can take control over the actions of a human driver (e.g., Toyota’s ‘Guardian’) to perform emergency maneuvers. Other vehicles may do most of the driving, while requiring the user to constantly monitor the situation and be ready to take control (e.g., Tesla’s ‘Autopilot’). Unless clear or explicitly mentioned, we use “Human” and “User” interchangeably to refer to the person inside the car (being a driver or a passenger), and we use “Industry” and “Machine” interchangeably to refer to both company and car combined together.

Our central question is this: when an automated car crashes and harms someone, how is blame and causal responsibility attributed to the human and machine drivers by people who hear about the crash? In this article, we use vignettes in which a pedestrian was hit and killed by a car being operated under shared control of a primary and a secondary driver and ask our participants to evaluate the crash on metrics of blame and causal responsibility. The cases we use are hypothetical (insofar as respondents know that they did not actually take place), but are not unrealistic, as they were designed to contain the relevant elements of events that could actually occur. We consider a large range of control regimes (see Figure 1), but the two main cases of interest are the instances of shared control where a human is the primary driver and the machine a secondary driver (“human-machine”) and where the machine is the primary driver and the human the secondary driver (“machine-human”). We consider a simplified space of scenarios in which (a) the main driver makes the correct choice and the secondary driver incorrectly intervenes (“Bad Intervention”) and (b) the main driver makes an error and the secondary driver fails to intervene (“Missed Intervention”). Both scenarios end in a crash. For comparison, we also include analogous scenarios involving a single human driver (a regular car) or a single machine driver (a fully automated car) as well as two hypothetical two-driver cars (driven by two humans or two machines). We ask participants to

make evaluations of the human user and one representative of the machine, either the car itself, or the company that designed the car.

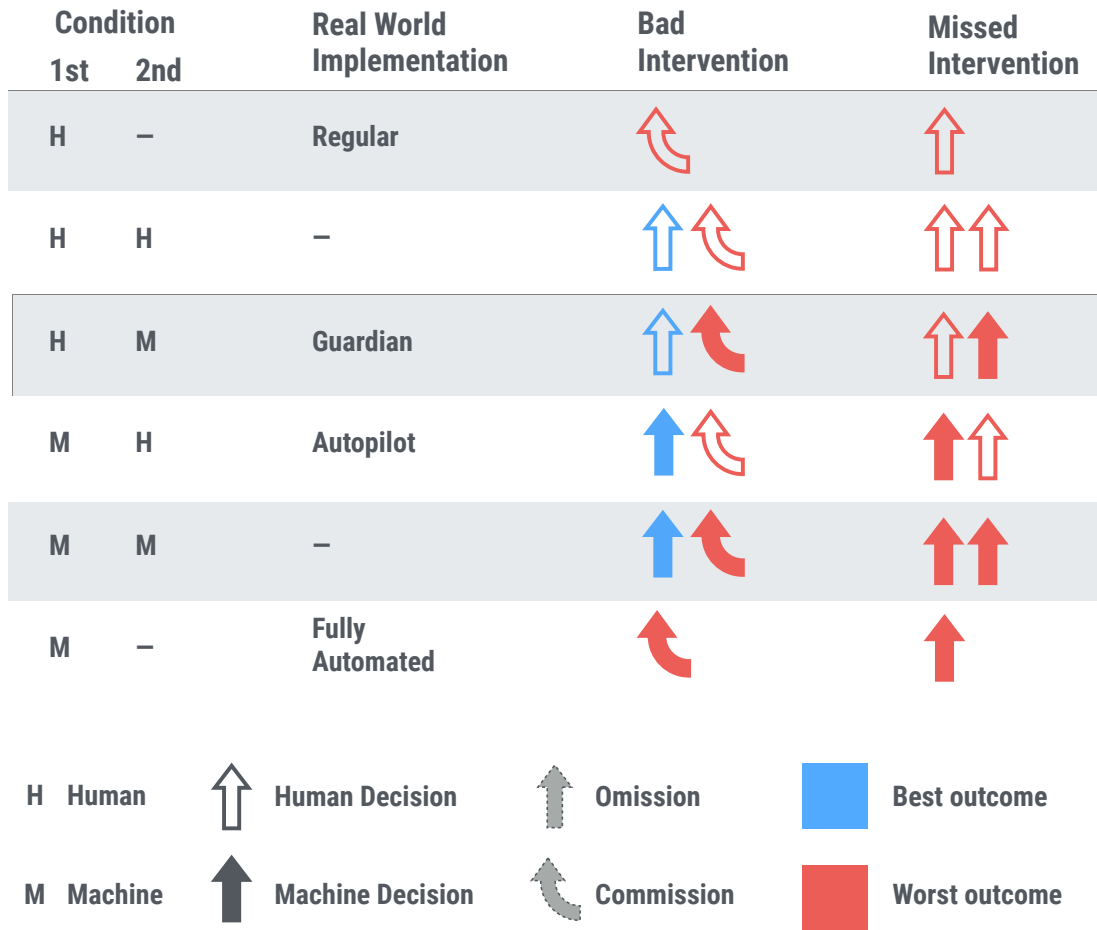


Figure 1: **Actions or action sequences for the different car types considered.** Outline arrows indicate an action by a human, ‘H’, and solid arrows indicate an action by a machine, ‘M’. The top and bottom rows represent sole-driver cars, while all others represent dual-driver cars. A red arrow indicates a decision – whether action or inaction – that had the avoidable death of a pedestrian as the outcome. A blue arrow indicates a decision that does not result in any deaths. For example, the H+M type (real-world implementation is the Guardian system) has a human main driver and a machine standby driver. A Bad Intervention then involves the human staying on course (a non-lethal action, indicated by the outline of the straight, blue arrow) and the machine overriding that action, causing the death of the pedestrian (solid, angled, red arrow). A Missed Intervention involves the human staying on course to kill the pedestrian (outline, straight, red arrow) without intervention from the machine (solid, straight, red arrow).

In Bad Intervention cases (see Methods – Case Description for details), the primary driver (be it human or machine) has made a correct decision to keep the car on course, which will avoid a pedestrian. Following this, the secondary driver makes the decision to swerve the car into the pedestrian. In these sorts of cases, we expect that the secondary driver (the only driver that makes a mistake) will be blamed more than the first driver. What is less clear is if people will assign blame and causal responsibility differently if this secondary driver is a human driver or a machine. Recent research suggests that humans may be blamed more than robots for making the same error

in the same situations [11].

In Missed Intervention cases, the primary driver has made an incorrect decision to keep the car on course (rather than swerving), which would cause the car to hit and kill a pedestrian. The secondary driver then neglects to swerve out of the way of the pedestrian. In these cases, the predictions for how participants will distribute blame and causal responsibility are less clear because both drivers make a mistake. As in the Bad Intervention cases, agent type (human or machine) may have an effect on blame and causal responsibility ratings. But unlike with Bad Intervention cases, Missed Intervention cases introduce the possibility that driver role (primary or secondary) may also impact judgments. It is possible that participants may shift responsibility and blame either toward the agent who contributed the most to the outcome (primary driver), or to the agent who had the last opportunity to act (secondary driver; [12, 13, 14, 15]). Under some regimes – such as Toyota’s Guardian – the user does most of the driving, but the decision to override (and thus to act last) pertains to the machine. Under others – such as Tesla’s Autopilot – the machine does most of the driving, but the decision to override pertains to the user.

Results

All studies used hypothetical vignettes that describe a crash (see Methods – Case Description for details on car regimes and intervention types, and see Supplementary Methods 1 for vignettes of studies 1-5).

Study 1

Study 1 compared four kinds of cars with different regimes of control. Each car had a primary driver, whose job it was to drive the car, and a secondary driver, whose job it was to monitor the actions of the first driver and intervene when the first driver made an error. The car architectures of central interest were human primary-machine secondary (“human-machine”) and machine primary-human secondary (“machine-human”). We also included human-human and machine-machine architectures for comparison. This allowed us to see how blame was distributed in a dual-driver architecture when there was no difference in driver type (human or machine) in each of the driving roles (primary or secondary).

Bad Interventions

In Bad Intervention cases, two predictors were entered into a regression with blame and causal responsibility ratings as the outcome variable: (1) whether or not the driver made an error and (2) driver type (human or machine). The main finding is that whether or not the driver made an error was a significant predictor of ratings (see Table 1 – Column: Bad Intervention - Study 1). In other words, a driver that unnecessarily intervened, leading to the death of a pedestrian was blamed more than a driver that operated on the correct course – regardless of whether the driver was a human or machine. It is worth noting here that we did not detect a reliable effect of driver type (Human vs. Machine), once correcting for multiple comparisons (see Table 1 – Row: Human, Column: Bad Intervention - Study 1). We do not discuss this factor further in the Bad Intervention cases.

Missed Interventions

In Missed Intervention cases, blame and responsibility judgments cannot depend on whether or not a driver made an error because both drivers make errors in these cases. The main finding for these cases is that driver type – whether the driver is a human or machine – has a significant

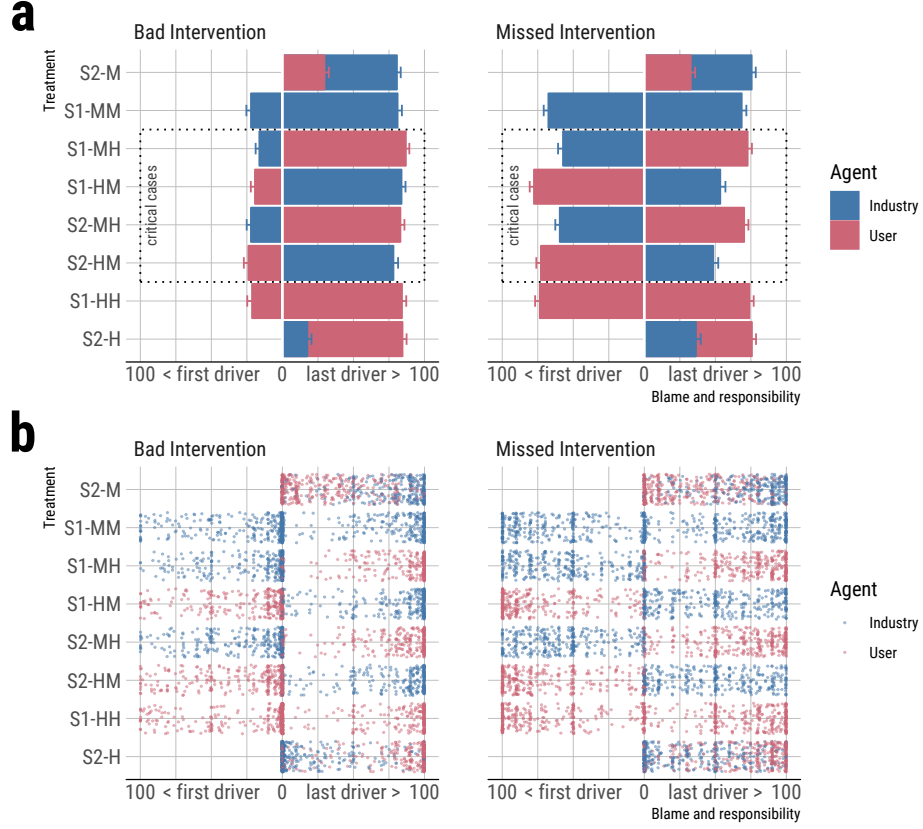


Figure 2: **Blame ratings for User and Industry in six car types.** (a) Bar plot and (b) Dot plot. Data from Study 1 (S1: $N = 786$, $Observations = 3,144$) and Study 2 (S2: $N = 382$, $Observations = 1,528$). Ratings of blame and causal responsibility are aggregated (collectively referred to as blame, henceforth). Ratings of car and company are aggregated (collectively referred to as Industry, henceforth). The y-axis represents the six car types considered in Studies 1 and 2 (S1 and S2). Two car types, HM (human-machine) and MH (machine-human), were considered in both studies. The y-axis labels include the study and the car type. For example, S1-HM represents the Human-Machine regime ratings collected in Study 1. In the six car types, the x-axis labeling of first driver refers to the main driver, and the last driver refers to the secondary driver in dual-driver cars, and the sole driver in the sole-driver cars. For Bad Intervention, only one agent has erred (the last driver). This agent (whether User or Industry) is blamed more than the other agent (first driver; see Table 1). For Missed Intervention, in dual-driver cars (rows 2-7), both agents have erred. When human and machine are sharing control (inside the dotted rectangle), blame ratings of Industry drops significantly, regardless of the role of the machine (main or secondary). In Study 1, blame to Industry in S1-MH ($m_1 = 57.2$) is significantly less than in S1-MM ($m_2 = 68$), $[t(760.6) = -5.05, p < .0001, m_2 - m_1 = 10.8, 95\% \text{ CI for } m_2 - m_1 \text{ is } 6.6-15]$. And blame to Industry in S1-HM ($m_1 = 53.4$) is significantly less than in S1-MM ($m_2 = 68$), $[t(722.77) = -6.6042, p < .0001, m_2 - m_1 = 14.6, 95\% \text{ CI for } m_2 - m_1 \text{ is } 10.2-19]$. In Study 2, blame to Industry in S2-M ($m_1 = 75.6$) is significantly more than in S2-MH ($m_2 = 59.5$), $[t(754.63) = -7.3885, p < .0001, m_1 - m_2 = 16.1, 95\% \text{ CI for } m_1 - m_2 \text{ is } 11.8-20.3]$; And is significantly more than in S2-HM ($m_3 = 48.51$), $[t(745.06) = 11.676, p < .0001, m_1 - m_3 = 27.1, 95\% \text{ CI for } m_1 - m_3 \text{ is } 22.5-31.6]$.

impact on ratings. Specifically, in these shared-control scenarios, where both human and machine have made errors, the machine driver is consistently blamed less than the human driver (Table 1 – Column: Missed Intervention - Study 1; Figure 2).

The human-machine difference appears to be driven by a reduction in the blame attributed to machines when there is a human in the loop. This is evident when comparing both the human-machine and machine-human instances of shared control to the machine-machine scenario. Note that the behaviors in these scenarios are identical, but how much a machine is blamed depends on whether it is sharing control with a human or operating both the primary and secondary driver role. When the machine is the primary driver, it is held significantly less blameworthy when its secondary driver is a human ($m_1 = 57.2$) compared to when the secondary driver is also the machine ($m_2 = 68$), $t(760.6) = -5.0$, $p < .0001$, difference in means: $m_2 - m_1 = 10.8$, 95% CI for $m_2 - m_1$ is 6.6–15. (All tests are two-tailed.) Similarly, when the machine is the secondary driver, it is held significantly less blameworthy when its primary driver is a human ($m_1 = 53.4$), compared to when the primary driver is also the machine ($m_2 = 68$), $t(722.77) = -6.6$, $p < .0001$, difference in means: $m_2 - m_1 = 14.6$, 95% CI for $m_2 - m_1$ is 10.2–19.

Study 2

Study 2 compared the human-machine and machine-human shared control cars with two different baseline cars: a standard car, which is exclusively driven by a human, and a fully automated car, which is exclusively driven by a machine. This allowed us to both replicate the main results of Study 1 (the responses to Machine-Human and Human-Machine crashes) and also to see how blame was assigned differently to dual-driver cars as compared to sole-driver cars. The industry representative was varied (car and company) but this exploratory variable was not analyzed in this study and in subsequent studies.

Bad Interventions

We replicated the main results of Study 1. Namely, in Bad Intervention cases for the shared-control cars (Machine-Human and Human-Machine), whether or not the driver made an error was a significant predictor of ratings (Table 1 – Column: Bad Intervention - Study 2; Figure 2).

Missed Intervention

We again replicated the main finding of Study 1. Driver type – whether the driver is a human or machine – has a significant impact on ratings. Specifically, in shared-control scenarios (Machine-Human and Human-Machine), where both human and machine have made errors, the machine driver is consistently blamed less than the human driver (Table 1 – Column: Missed Intervention - Study 2; Figure 2).

As we noted in Study 1, the human-machine difference is driven by a reduction in the blame attributed to machines when there is a human in the loop. This is verified in Study 2 by comparing blame to the machine in the shared control cases with blame to the machine in the fully automated car (driven by a sole machine driver). In each case, blame to the machine in the shared control case is significantly lower than blame to the machine in the Fully Automated car: Fully Automated ($m_1 = 75.6$) vs. Machine-Human ($m_2 = 59.5$), $t(754.63) = -7.4$, $p < .0001$, $m_1 - m_2 = 16.1$, 95% CI for $m_1 - m_2$ is 11.8–20.3; vs. Human-Machine ($m_3 = 48.5$), $t(745.06) = 11.7$, $p < .0001$, $m_1 - m_3 = 27.1$, 95% CI for $m_1 - m_3$ is 22.5–31.6.

Table 1: **Regression analysis of data collected in studies 1-5 in the cases of Bad Intervention and Missed Intervention.** Data from Studies 2 and 3 are limited to shared-control regimes in the table. “Human” refers to the type of agent in question (that is, human as compared to the baseline, machine), “Mistake” refers to whether the decision was a mistake (that is, the decision would have resulted in losing a life, or losing more lives in study 3), and “Last Driver” refers to the driver role (that is, the driver assumes the secondary role). All models include participant random effects and question (blame or causal responsibility) random effects, where applicable. Data were assumed to meet the requirements of the model.

	Blame and Causal Responsibility								
	Bad Intervention			Missed Intervention					
	Study 1	Study 2	Study 3	Study 1	Study 2	Study 3	Study 4	Study 5	
Human	2.141 (1.061) p = 0.044	3.358 (1.574) p = 0.033	-1.508 (0.811) p = 0.063	16.942 (1.148) p = 0.000	17.493 (1.514) p = 0.000	3.567 (0.852) p = 0.000	10.745 (2.189) p = 0.000	2.594 (0.860) p = 0.003	
Mistake	64.293 (1.061) p = 0.000	57.559 (1.574) p = 0.000	11.917 (0.881) p = 0.000						
Last Driver				-1.822 (0.915) p = 0.047	-6.759 (1.514) p = 0.000	1.715 (0.852) p = 0.045	-0.073 (2.189) p = 0.974	1.355 (0.860) p = 0.116	
Constant	18.653 (0.916) p = 0.000	21.352 (1.370) p = 0.000	27.406 (1.171) p = 0.000	60.504 (1.531) p = 0.000	57.354 (1.878) p = 0.000	36.102 (1.252) p = 0.000	61.032 (1.911) p = 0.000	65.923 (0.811) p = 0.000	
Participant Rand Effects?	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	
Question Rand Effects?	Yes	Yes	Yes	Yes	Yes	Yes	N/A	N/A	
N	786	382	389	786	382	389	375	2000	
Observations	3,144	1,528	3,112	3,144	1,528	3,112	750	4,000	

Study 3

In Study 3, we used the same car regimes as in Study 2, but the cases were dilemma scenarios in which the drivers had to choose between crashing into a single pedestrian and crashing into five pedestrians. This study was conducted as a comparison with our Studies 1 and 2, which involve clear errors, and the studies conducted in previous research on self-driving cars (such as [16, 17, 11, 18]) which involve the difficult choice of deciding which of two groups of people to hit. All the main effects in Study 2 were replicated in Study 3.

Bad Interventions

Replicating the main results of Studies 1 and 2, in Bad Intervention cases for the shared-control cars, whether or not the driver made an error was a significant predictor of ratings (Table 1 – Column: Bad Intervention - Study 3).

Missed Intervention

Replicating the main results of Studies 1 and 2, driver type has a significant impact on ratings. Specifically, in shared-control scenarios, where both human and machine have made errors, the machine driver is consistently blamed less than the human driver (Table 1 – Column: Missed Intervention - Study 3).

Study 4

In Study 4, we replicated the central findings using more ecologically valid stimuli. We used only the human-machine and machine-human shared control cars in the Missed Intervention scenario; these are the cases where we observed the systematic decrease in blame to the machine in Studies 1-3. For this study, we continued to use hypothetical scenarios, but the stimuli shown to participants looked like realistic newspaper articles (see Supplementary Methods 1 – Studies 4-5).

The main finding of Studies 1-3 was replicated: the machine driver is consistently blamed less than the human driver in these shared-control scenarios where both human and machine have made errors (Table 1 – Column: Missed Intervention - Study 4).

Study 5

Study 5 was a replication of Study 4 run via YouGov with a nationally representative sample of the United States population (see Figure 4 for details).

The main finding was again replicated: the machine driver is consistently blamed less than the human driver (Table 1 – Column: Missed Intervention - Study 5). This result (i.e. human is blamed more than machine) holds directionally in 82% of demographic subgroups of participants (see Figure 3).

Discussion

Our central finding was that in cases where a human and a machine share control of the car in hypothetical scenarios, less blame is attributed to the machine when both drivers make errors. The first deadly crashes of automated vehicles (mentioned in the Introduction) were similar in structure to our Missed Intervention cases. In those cases, both the machine primary driver and the human secondary driver should have taken action (braking to avoid a collision) and neither driver did. Our results suggests that the public response that occurred to the crash – one that focused attention on the driver being exceedingly negligent – is likely to generalize to other dual-error Missed Intervention-style cases, shifting blame away from the machine and towards the human. Moreover, the convergence of our results with this real-world public reaction seems to suggest that while we employed stylized, simplified vignettes in our research, our findings show external validity. Moreover, this pattern of results was replicated in a nationally representative sample of the United States population (and across different subgroups; see Figure 3), which employed naturalistic presentation of scenarios (see Supplementary Methods 1 – Studies 4-5).

Our central finding (diminished blame to the machine in dual-error cases) leads us to believe that, while there may be many psychological barriers to self-driving car adoption [19], public over-reaction to dual-error cases is not likely to be one of them. In fact, we should perhaps be concerned about public under-reaction. Because the public are less likely to see the machine as being at fault in dual-error cases like the Tesla and Uber crashes, the sort of public pressure that drives regulation might be lacking. For instance, if we were to allow the standards for automated vehicles to be set through jury-based court-room decisions, we expect that juries will be biased to absolve the car manufacturer of blame in dual-error cases, thereby failing to put sufficient pressure on manufacturers to improve car designs. Despite the fact that there are some avenues available to courts to mitigate psychological biases that may arise among juries (such as carefully worded jury instructions or expert witnesses), psychological biases continue to play an important role in court-based decisions [20]. In fact, we have been in a similar situation before. Prior to the 1960s, car manufacturers enjoyed a large amount of liberty from liability when a car's occupant was harmed in a crash (because blame in car crashes was attributed to the driver's error or negligence).

Effect of Demographic Attributes

Few subgroups blame Industry more, and the ones that do typically have smaller samples

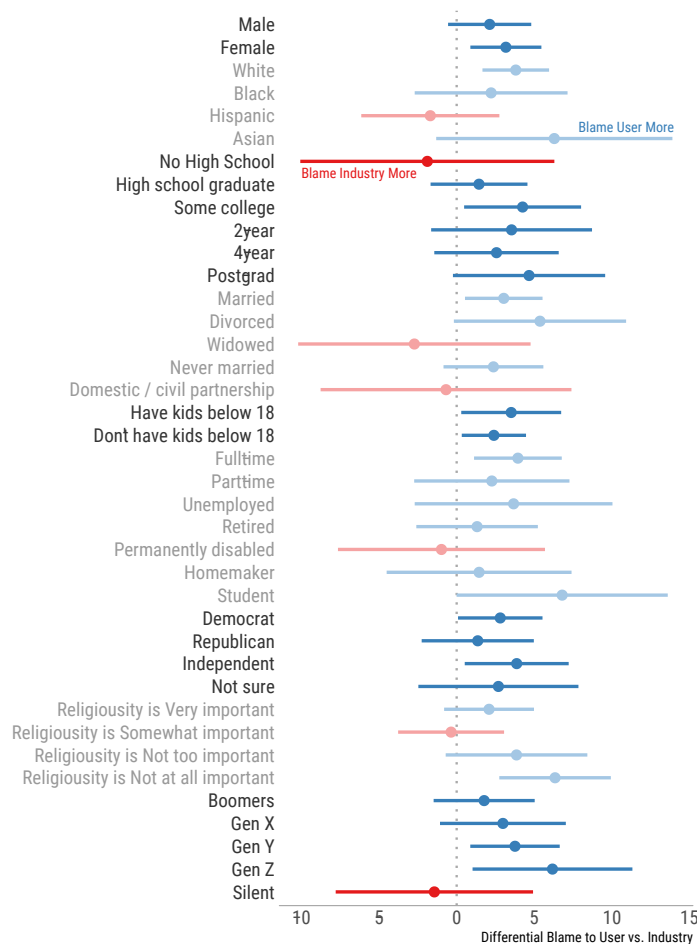


Figure 3: **Ratings of demographic subgroups in Study 5.** Data collected in Study 5 from nationally representative sample ($N = 2000$, $Observations = 4000$). Each row represents the mean of differential blame attributed to User (i.e.: human) vs. Industry (i.e.: car or company). Positive values (blue) indicate more blame attributed to User, while negative values (red) represent more blame attributed to Industry. Error bars are 95% Confidence Intervals. Only subgroups with at least 50 participants are shown. 33 out of 40 subgroups (82%) attribute more blame to User.

Top-down regulation was necessary to introduce the concept of “crash worthiness” into the legal system, that is, that cars should be designed in such a way to minimize injury to occupants when a crash occurs. Only following these laws were car manufacturers forced to improve their designs [21]. Here, too, top-down regulation of automated car safety might be needed to correct a public under-reaction to crashes in shared-control cases. What, exactly, the safety standard should be is still an open question, however.

If our data identifies a source of possible public over-reaction, it is for cars with a human primary driver and a machine secondary driver in Bad Intervention-style cases. These are the only cases we identified where the car receives more blame than the human. It seems possible that these sorts of cars may generate widespread public concern once we see instances of Bad Intervention-style crashes in human-machine car regimes. This could potentially slow the transition to fully

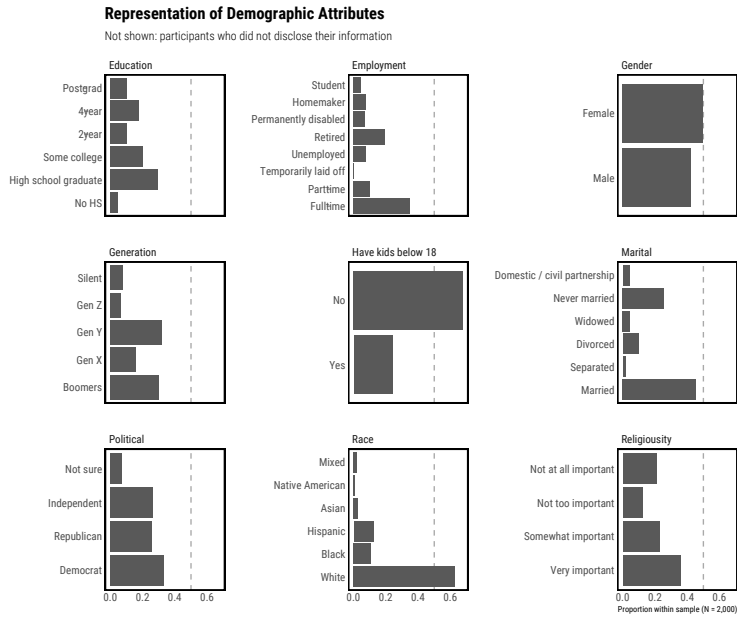


Figure 4: **Representation of Demographic Attributes in Study 5.** Study 5 was run via YouGov, a service that administers and runs surveys on nationally representative samples. The sample includes 2000 participants with diverse demographic attributes.

automated vehicles if this reaction is not anticipated and managed appropriately in public discourse and legal regulation. Moreover, manufacturers that are working to release cars with a machine secondary driver should plan appropriately for the likely legal fall-out for these unique cases where a machine driver receives more blame than a human.

Our data portends the sort of reaction we can expect to automated car crashes at the societal level (for example, through public reaction and pressure to regulate). Once we begin to see societal-level responses to automated cars, that reaction may shape incentives for individual actors. For example, people may want to opt into systems that are designed such that, in the event of a crash, the majority public response will be to blame the machine. Worse yet, people may train themselves to drive in a way that, if they crash, the blame is likely to fall to the machine (for instance, by not attempting to correct a mistake that is made by a machine over-ride). This sort of incentive shaping may already be happening in the legal domain. Judges who make decisions about whether to release a person from custody between arrest and a trial frequently rely on actuarial risk assessment tables to help make their decision. Some suspect that judges may overly rely on the tables as a way of diminishing their responsibility if a released person commits a crime. Recent attention generated in response to such a case focused on the role of the algorithm rather than the judge [22], indicating that the possibility of incentive shaping in the legal domain is not so far-fetched.

Given these possible societal-level implications of our findings, it is important to acknowledge the potential limitations of interpreting our data this broadly. First, the participants in all of our experiments know that they are reading about hypothetical scenarios. It is possible that this reduces the psychological realism of the study [23, 24], causing participants' responses to be characteristically different than they would be upon reading about an actual event. The literature provides a mixed view of how well responses to hypothetical scenarios map onto those made in real life situations [25, 26, 27, 28]. However, the research that does show considerable differences [25, 26] finds that these differences are mostly seen in the way participants themselves would act in

moral situations and not necessarily about the moral judgments they render about third parties. In our paper, we study participants’ judgments (blame and causal responsibility) about third parties in hypothetical scenarios; these may align more directly with judgments of actual scenarios.

Second, although we may see a reasonably tight mapping between the opinions expressed in this study’s scenarios and the opinions that would be expressed in real life situations, it is important to note that, in the latter case, judgments will not be occurring in isolation. Instead, they will occur within a richer context than the carefully controlled scenarios used in our studies. People may hear reports of the accidents with more emotion-arousing details, which are known to skew people’s judgments [29, 30]. Moreover, the public’s reaction to hearing about semi-autonomous vehicle crashes will be shaped by many factors beyond their immediate psychological response (which is the object of our study) including opinion pieces they read, the views of community leaders, and so on. These factors will collectively shape the public’s overall reaction to crashes.

Studies 1, 2, 4, and 5 looked at blame and causal responsibility attribution in cases where one or both drivers made errors. Study 3 looked at dilemma scenarios where the drivers faced the choice of running over either one or five pedestrians. While there is, in some sense, an “optimal” outcome in these cases (corresponding to saving more lives), it is not obvious that it would (for example) count as an error to refuse to swerve away from five pedestrians into a pedestrian that was previously unthreatened. In fact, the German Ethics Commission on Automated and Connected Driving report [31] indicates that programming cars to trade off lives in this way would be prohibited. The report states: “It is also prohibited to offset victims against one another. [...] Those parties involved in the generation of mobility risks must not sacrifice non-involved parties.” Even though participants in previous studies prefer to sacrifice one person who was previously not involved than five (e.g., [16, 17]), the German Ethics Commission’s decision underscores the fact that trading off lives in dilemma situations can be particularly fraught. For this reason, and for continuity with previous work on the ethics of self-driving cars [16, 17, 11] and in moral psychology more generally [32, 33], we chose to investigate dilemma situations. Our findings about the effect of driver type in these cases underscores the fact that findings about how blame and responsibility are attributed after a crash may still hold in less-clear dilemma scenarios.

Some of our results fall in line with previous work on the psychology of causal inference. In Bad Intervention cases, the primary driver (be it human or machine) makes a correct decision to keep the car on course, which will avoid a pedestrian. Following this, the secondary driver makes the decision to swerve the car into the pedestrian. Our data show that the secondary driver (the one that makes a mistake) is considered more causally responsible than the first driver. It is well established that judgments of causal responsibility are impacted by violations of statistical and moral norms [34, 35, 36, 37], and a mistake seems to count as such a violation. That is, if something unusual or counter-normative happens, that event is more likely to be seen as a cause of some effect than another event that is typical or norm-conforming.

Moreover, the central finding that humans are blamed more than machines even when both make errors accords with research on the psychology of causal attribution. Findings in that field suggest that voluntary causes (causes created by agents) are better causal explanations than physical causes [38]. While it is clear that what a human does is fundamentally different than what a machine does in each of the scenarios, it remains an open question whether an AI that is operating a car is perceived as a physical cause, an agent, something in between, or something else entirely [39, 40]. Future work should investigate the mental properties attributed to an AI that controls a car both in conjunction with a human or alone. Understanding the sort of mind we perceive as dwelling inside an AI may help us understand and predict how blame and causal responsibility will be attributed to it [41].

Another open question concerns what the implications are of attributing blame to a machine at all. There are various ways that humans express moral condemnation. For example, we may

call an action morally wrong, say that a moral agent has a bad character, or judge that an agent is blameworthy. Judgments of blame typically track judgments of willingness to punish the perpetrator [42, 43]. Are the participants in our study expressing that some punishment is due to the machine driver of the car, whatever that may mean? Alternately, is it possible that participants' expressions of blame indicate that some entity is deserving of punishment that represents the machine (the company, or a human representative of the company, such as the CEO). The similar blame judgments given to the car and the car's representatives (company) perhaps support this possibility. Finally, it is possible that participants ascribe only non-moral blame to the machine, in the sense of being responsible but not in a moral sense. We may say that a forest fire is to blame for displacing residents from their homes without implying that punishment is due to anyone at all.

Following these studies, the reason that participants blame machine drivers less than human drivers in Missed Intervention cases also remains an open question. The findings may be linked to the uncertainty with which we perceive the agential status of machines. Once machines are a more common element in our moral world and we interact with machines as moral actors, will this effect change? Or will this finding be a lasting hallmark of the cognitive psychology of human-machine interaction?

A final open question concerns whether the effects we report here will generalize to other cases of human-machine interaction. Already we see fruitful human-machine partnerships emerging with judges, doctors, military personnel, factory workers, artists, and financial analysts, just to name a few. We conjecture that we may see the patterns we report here in domains other than automated vehicles, though each domain will have its own complications and quirks as machines begin to become more subtly integrated in our personal and professional lives.

Methods

This study was approved by the Institute Review Board (IRB) at Massachusetts Institute of Technology (MIT). The authors complied with all relevant ethical considerations, including obtaining informed consent from all participants.

In all studies, participants were allocated uniformly randomly into conditions. Data collection and analysis were performed blind to the conditions of the experiments. The sample size was chosen in each study to ensure having at least 100 participants for each condition. Number of participants was chosen in advance of running the study and all data was collected prior to analysis. See details below.

In Studies 1-3 we excluded any participant who did not (i) complete all measures within the survey, (ii) transcribe (near-perfectly) a 169-character paragraph from an image (used as an attention check), and (iii) have a unique MTurk ID per study (all records with a recurring MTurk ID were excluded).

Case Description

Summary descriptions of all car types and cases. For full vignettes, see Supplementary Methods 1.

Sole-driver car

This car has only one driver that does all the driving. Two versions are used.

Human-only. This is a sole-driver car, in which a human is the driver. Also referred to as a regular car.

Machine-only. This is a sole-driver car, in which a machine is the driver. Also referred to as a fully-automated car.

Dual-driver car

This car has a primary driver, whose job it is to drive the car, and a secondary driver, whose job it is to monitor the actions of the first driver and intervene when the first driver makes an error. Also referred to as shared-control car. Four versions are used.

Human-Machine. This is a dual-driver car, in which a human is the primary driver, and a machine is the secondary driver. Also referred to as Guardian.

Machine-Human. This is a dual-driver car, in which a machine is the primary driver, and a human is the secondary driver. Also referred to as Autopilot.

Human-Human. This is a dual-driver car, in which a human is the primary driver, and another human is the secondary driver.

Machine-Machine. This is a dual-driver car, in which a machine is the primary driver, and another machine is the secondary driver.

Intervention Types

We use two types of interventions: Bad Intervention and Missed Intervention. The description of each is dependent on whether the car is a sole-driver or a dual-driver car.

Bad Intervention (dual-driver). The primary driver kept the car on its track. The secondary driver intervened and steered the car off its track (killing a pedestrian) rather than keeping the car on track and killing no one.

Missed Intervention (dual-driver). The primary driver kept the car on its track. The secondary driver kept the car on its track (killing a pedestrian) rather than swerving into the adjacent lane and killing no one.

Bad Intervention (sole-driver). The sole driver steered the car off its track (killing a pedestrian) rather than keeping the car on track and killing no one.

Missed Intervention (sole-driver). The sole driver kept the car on its track (killing a pedestrian) rather than swerving into the adjacent lane and killing no one.

Dilemma versions (Study 3). The two outcomes of killing one pedestrian vs. killing no one are replaced with the two outcomes of killing five pedestrians vs. killing one pedestrian. For example, in Missed Intervention (dual-driver): [...] The secondary driver kept the car on its track (killing five pedestrians) rather than swerving into the adjacent lane and killing one pedestrian.

Study 1

Participants.

The data was collected in September 2017 from 809 participants (USA residents) recruited from the Mechanical Turk platform (each was compensated \$0.5). Of those, 23 participants were excluded

(as explained above), leaving us with 786 participants. Participants were aged between 18-83 (median: 33), 50% were females, 39% had annual income of \$50K or more, and 55% had a bachelor degree or higher.

Stimuli and procedures.

Participants were uniformly randomly allocated to one of four conditions. Conditions varied the car type (human-human, human-machine, machine-human, and machine-machine) in a 4-level between-subjects design. In each condition, participants first read a description of the car, and were then asked to attribute competence to each of the two drivers on an 100-point scale anchored at “not competent” and “very competent” (see Supplementary Figure 1 for results on competence). Participants then read two scenarios (presented in a random order), one Bad Intervention case and one Missed Intervention Case. After each scenario, participants were asked to indicate (on an 100-point scale) to what extent they thought each driver was blame-worthy (from “not blame-worthy” to “very blame-worthy”) and to what degree each of these two agents caused the death of the pedestrian (from “very little” to “very much”). Questions were presented in a randomized order. (See Supplementary Methods 1 – Study 1 for text of the vignettes and see Supplementary Methods 2 for questions). At the end of the surveys, participants provided basic demographic information (e.g., age, gender, income, education).

Study 2

Participants.

The data was collected in May 2017 from 804 participants (USA residents) recruited from the Mechanical Turk platform (each is compensated \$0.3). Of those, 25 participants were excluded (as explained above), leaving us with 779 participants. Participants were aged between 18-77 (median: 32), 48% were females, 39% had annual income of \$50K or more, and 54% had a bachelor degree or higher.

Stimuli and procedures.

Participants were uniformly randomly allocated to one of eight conditions. Conditions varied the car type (human only, human-machine, machine-human, and machine only) and the industry representative (car and company), in a 4x2 between-subjects multi-factorial design. In each condition, participants read two scenarios (presented in a random order), one Bad Intervention case and one Missed Intervention case. After each scenario, participants were asked to attribute causal responsibility, blameworthiness, and competence (see Supplementary Figure 1 for results on competence) to two agents: the human in the car and a representative of the car (the car itself or the manufacturing company of the car, depending on the condition). All other features of Study 2 were the same at those in Study 1.

Study 3

Participants.

The data was collected in November 2016 from 1008 participants (USA residents only) recruited from the Mechanical Turk platform (each is compensated \$0.6). Of those, 35 participants were excluded (as explained above), leaving us with 973 participants. Participants were aged between 18-84 (median: 33), 51% were females, 37% had annual income of \$50K or more, and 53% had a bachelor degree or higher.

Stimuli and procedures.

There were two groups of participants in Study 3: those who saw dual-driver cases or those who saw sole-driver. For those who saw dual-driver cases, participants were randomly assigned to one of six conditions in a 2x3 design, varying the car type (human-machine or machine-human) and the industry representative (car, company, and programmer). Data for programmer was later dropped from the analysis. For those who saw sole-driver cases, participants were randomly assigned to one of four conditions in a 2x2 design, varying the car type (human only or machine only) and the industry representative (car or company). In each condition (for both dual-car and single-car groups), participants read two scenarios (presented in a random order), one Bad Intervention case and one Missed Intervention case. These scenarios were the dilemma versions of those presented in Studies 1 and 2 (see description above). After each scenario, participants were asked to attribute causal responsibility and blameworthiness to two agents: the human in the car and a representative of the car (the car itself, the company, or the programmer, depending on the condition). All other features of Study 3 were identical to those of Study 2.

Study 4

Participants.

The data was collected in January 2019 from 375 participants (USA residents only) recruited from the Mechanical Turk platform (each is compensated \$0.3). No demographic data was collected for this study. Given that it was done on the same platform as studies 1-3 (i.e. Mechanical Turk), its demographic proportions are expected to be similar.

Stimuli and procedures.

The key elements of this study and Study 5 are 1) the restriction to Missed Intervention cases, and 2) the visual and textual content of the vignettes have the look and feel of a news piece (see Supplementary Methods 1 – Studies 4-5).

Participants were uniformly randomly allocated to one of four conditions. Conditions varied the car type (human-machine, and machine-human) and the industry representative (car, and company), in a 2x2 between-subjects multi-factorial design. In each condition, participants read one scenario; one Missed Intervention case. The textual content of these scenarios was close to those presented in Studies 1 -3, with slight changes to the text to make it read like a news piece. After each scenario, participants were asked to attribute blameworthiness to two agents: the human in the car and a representative of the car (the car itself, or the company, depending on the condition).

Study 5

Participants.

The data was collected in March 2019 from 2189 participants (USA residents) were recruited via YouGov, a service that administered the study and collected the data from a representative sample of participants. The participants were then matched down to a sample of 2000 participants based on demographics. See Figure 4 for details on demographic proportions of participants in this study.

Stimuli and procedures.

This study is identical in the setup to Study 4.

Data Availability.

Raw data and Source data for Fig 2, 3, and 4; Table 1; and Supplementary Fig 1 are available in: <https://bit.ly/2kzLymH>

Code Availability.

Code used to produce figures and tables mentioned above is available in: <https://bit.ly/2kzLymH>

References

- [1] WHO. Road traffic injuries. *World Health Organization Fact sheet* (2017).
- [2] Geistfeld, M. A. A roadmap for autonomous vehicles: State tort liability, automobile insurance, and federal safety regulation. *Calif. L. Rev.* **105**, 1611 (2017).
- [3] Tesla. A tragic loss. *Tesla* (2016).
- [4] NHTSA. Automatic vehicle control systems – investigation of tesla accident. <https://static.nhtsa.gov/odi/inv/2016/INCLA-PE16007-7876.PDF> (2016).
- [5] Griswold, A. Uber found not criminally liable in last year’s self-driving car death. <https://qz.com/1566048/uber-6not-criminally-liable-in-tempe-self-driving-car-death/> (March 5, 2019).
- [6] AP. Tesla driver killed while using autopilot was watching harry potter, witness says. *Associated Press News* (2016).
- [7] Chong, Z. & Krok, A. Tesla not at fault in fatal crash, driver was not watching a movie. *CNET* (2017).
- [8] Randazzo, R. Who was really at fault in fatal uber crash? here’s the whole story. <https://www.azcentral.com/story/news/local/tempe/2019/03/17/one-year-after-self-driving-uber-rafaela-vasquez-behind-wheel-crash-death-elaine-herzberg-tempe/1296676002/> (March 17, 2019).
- [9] Munster, G. Here’s when having a self-driving car will be a normal thing. <http://fortune.com/2017/09/13/gm-cruise-self-driving-driverless-autonomous-cars/> (Sept 13, 2017).
- [10] Kessler, S. A timeline of when self-driving cars will be on the road, according to the people making them. <https://qz.com/943899/a-timeline-of-when-self-driving-cars-will-be-on-the-road-according-to-the-people-making-them/> (March 29, 2017).
- [11] Li, J., Zhao, X., Cho, M.-J., Ju, W. & Malle, B. F. From trolley to autonomous vehicle: Perceptions of responsibility and moral norms in traffic accidents with self-driving cars. Tech. Rep., SAE Technical Paper (2016).
- [12] Chockler, H. & Halpern, J. Y. Responsibility and blame: A structural-model approach. *Journal of Artificial Intelligence Research* **22**, 93–115 (2004).
- [13] Gerstenberg, T. & Lagnado, D. A. When contributions make a difference: Explaining order effects in responsibility attribution. *Psychonomic Bulletin & Review* **19**, 729–736 (2012).

- [14] Sloman, S. A. & Lagnado, D. Causality in thought. *Annual Review of Psychology* **66**, 223–247 (2015).
- [15] Zultan, R., Gerstenberg, T. & Lagnado, D. A. Finding fault: causality and counterfactuals in group attributions. *Cognition* **125**, 429–440 (2012).
- [16] Bonnefon, J.-F., Shariff, A. & Rahwan, I. The social dilemma of autonomous vehicles. *Science* **352**, 1573–1576 (2016).
- [17] Awad, E. *et al.* The moral machine experiment. *Nature* **563**, 59 (2018).
- [18] Malle, B., Scheutz, M., Arnold, T., Voiklis, J. & Cusimano, C. Sacrifice one for the good of many? people apply different. In *Proceedings of 10th ACM/IEEE International Conference on Human-Robot Interaction* (2015).
- [19] Shariff, A., Bonnefon, J.-F. & Rahwan, I. Psychological roadblocks to the adoption of self-driving vehicles. *Nature Human Behaviour* **1**, 694 (2017).
- [20] Bornstein, B. H. & Greene, E. Jury decision making: Implications for and from psychology. *Current Directions in Psychological Science* **20**, 63–67 (2011).
- [21] Nader, R. Unsafe at any speed. the designed-in dangers of the american automobile (1965).
- [22] Westervelt, E. Did a bail reform algorithm contribute to this san francisco man’s murder? <https://www.npr.org/2017/08/18/543976003/did-a-bail-reform-algorithm-contribute-to-this-san-francisco-man-s-murder> (August 18, 2017).
- [23] Bauman, C. W., McGraw, A. P., Bartels, D. M. & Warren, C. Revisiting external validity: Concerns about trolley problems and other sacrificial dilemmas in moral psychology. *Social and Personality Psychology Compass* **8**, 536–554 (2014).
- [24] Aronson, E., Wilson, T. D. & Brewer, M. B. Experimentation in social psychology. *The handbook of social psychology* **1**, 99–142 (1998).
- [25] FeldmanHall, O. *et al.* Differential neural circuitry and self-interest in real vs hypothetical moral decisions. *Social cognitive and affective neuroscience* **7**, 743–751 (2012).
- [26] Bostyn, D. H., Sevenhant, S. & Roets, A. Of mice, men, and trolleys: Hypothetical judgment versus real-life behavior in trolley-style moral dilemmas. *Psychological science* **29**, 1084–1093 (2018).
- [27] Dickinson, D. L. & Masclet, D. Using ethical dilemmas to predict antisocial choices with real payoff consequences: an experimental study (2018).
- [28] Plunkett, D. & Greene, J. Overlooked evidence and a misunderstanding of what trolley dilemmas do best: A comment on bostyn, sevenhant, & roets (2018). *Psychological Science* (2019).
- [29] Greene, J. & Haidt, J. How (and where) does moral judgment work? *Trends in cognitive sciences* **6**, 517–523 (2002).
- [30] Horberg, E. J., Oveis, C. & Keltner, D. Emotions as moral amplifiers: An appraisal tendency approach to the influences of distinct emotions upon moral judgment. *Emotion Review* **3**, 237–244 (2011).
- [31] Luetge, C. The german ethics code for automated and connected driving. *Philosophy & Technology* **30**, 547–558 (2017).

- [32] Mikhail, J. *Elements of moral cognition: Rawls' linguistic analogy and the cognitive science of moral and legal judgment* (Cambridge University Press, 2011).
- [33] Greene, J. *Moral tribes: Emotion, reason, and the gap between us and them* (Penguin, 2014).
- [34] Alicke, M. D. Culpable control and the psychology of blame. *Psychological bulletin* **126**, 556 (2000).
- [35] Gerstenberg, T., Goodman, N. D., Lagnado, D. A. & Tenenbaum, J. B. How, whether, why: Causal judgments as counterfactual contrasts. In *CogSci* (2015).
- [36] Hitchcock, C. & Knobe, J. Cause and norm. *The Journal of Philosophy* **106**, 587–612 (2009).
- [37] Kominsky, J. F., Phillips, J., Gerstenberg, T., Lagnado, D. & Knobe, J. Causal superseding. *Cognition* **137**, 196–209 (2015).
- [38] Hart, H. L. A. & Honoré, T. *Causation in the Law* (OUP Oxford, 1985).
- [39] Gray, H. M., Gray, K. & Wegner, D. M. Dimensions of mind perception. *science* **315**, 619–619 (2007).
- [40] Weisman, K., Dweck, C. S. & Markman, E. M. Rethinking people's conceptions of mental life. *Proceedings of the National Academy of Sciences* **114**, 11374–11379 (2017).
- [41] Gray, K., Young, L. & Waytz, A. Mind perception is the essence of morality. *Psychological inquiry* **23**, 101–124 (2012).
- [42] Cushman, F. Crime and punishment: Distinguishing the roles of causal and intentional analyses in moral judgment. *Cognition* **108**, 353–380 (2008).
- [43] Cushman, F. Deconstructing intent to reconstruct morality. *Current Opinion in Psychology* **6**, 97–103 (2015).

Acknowledgements

I.R., E.A., S.L., and S.D. acknowledge support from the Ethics and Governance of Artificial Intelligence Fund. J-F.B. acknowledges support from the ANR-Labex Institute for Advanced Study in Toulouse, the ANR-3IA Artificial and Natural Intelligence Toulouse Institute, and the grant ANR-17-EURE-0010 Investissements d'Avenir. The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript.

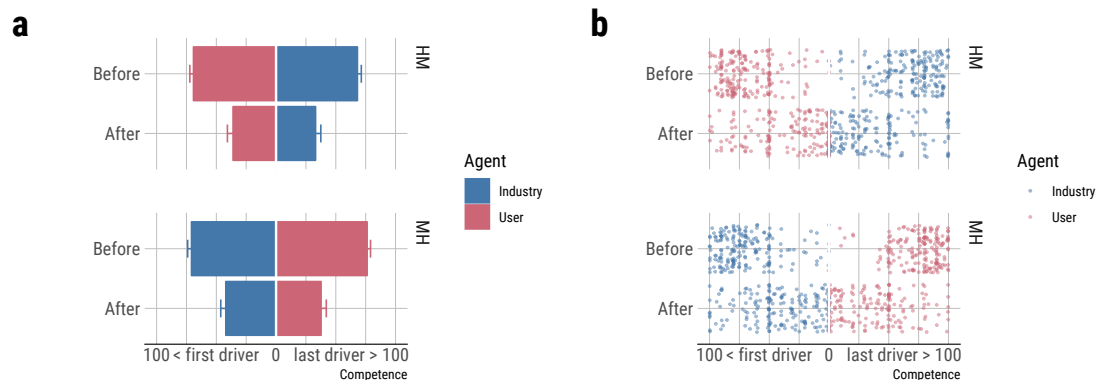
Author contributions statement

E.A., S.L., M.K-W., S.D., J.B.T., A.S., J-F.B., and I.R. contributed to the conception and design of the research. E.A., S.L., M.K-W., and S.D. conducted studies. E.A. and J-F.B. analyzed data. S.L., E.A., M.K-W., J-F.B., and I.R. wrote the manuscript. All authors reviewed and revised the manuscript.

Competing financial interests.

The authors declare no competing interests.

1 Supplementary Figures



Supplementary Figure 1. Competence ratings from studies 1 and 2, in the Missed Intervention cases, for the two critical regimes (human-machine and machine-human). (a) Bar plot and (b) Dot plot. In Study 1 (S1), ratings were collected before participants read about the accident. In Study 2 (S2), they were collected after. For Industry, ratings of car and company are aggregated (collectively referred to as Industry, henceforth). The x-axis labeling of first driver refers to the main driver, and the last driver refers to the secondary driver in dual-agent cars. Industry and User ratings are shown in blue, and red, respectively. User and Industry receive similar competence ratings in each case. Ratings of both User and Industry drop at the same rate when the question is asked after the scenario is presented.

Supplementary Methods

Supplementary Method 1

Study 1

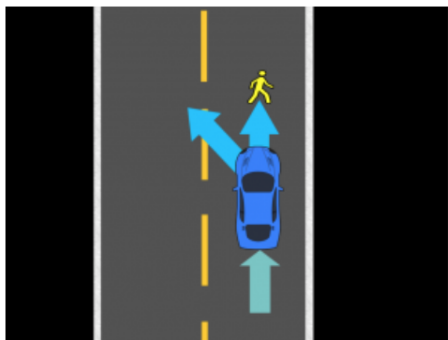
Hank is driving a car (that is equipped with a self-driving robocar software) with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is walking in the righthand lane** a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, he sees that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If Hank changes lanes, the car will **not run over anyone**.
- If Hank does not change lanes, the car will **run over and kill the man**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will **not run over anyone**.
- If the robocar does not override, the car will **run over and kill the man**.



In fact, the robocar decided to **not override**.

Vignette for Human-Machine with Missed Intervention.

Hank is driving a car (that is equipped with a self-driving robocar software) with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that a **man is walking in the lefthand lane** a short distance ahead. If he swerves, he would kill the man. Further, he sees that **no one is walking in the righthand lane ahead of him**. If he keeps the car on its track, it will not run over anyone.

- If Hank changes lanes, the car will **run over and kill the man**.
- If Hank does not change lanes, the car will **not run over anyone**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will **run over and kill the man**.
- If the robocar does not override, the car will **not run over anyone**.



In fact, the robocar decided to **override**.

Vignette for Human-Machine with Bad Intervention.

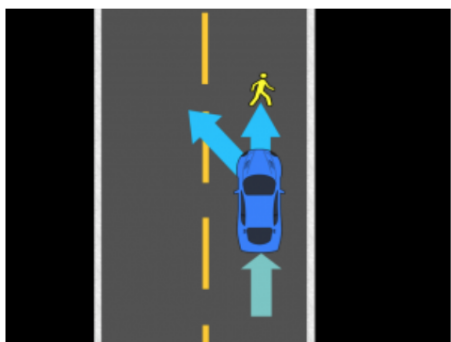
The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the righthand lane** a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, the robocar senses that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If the robocar changes lanes, the car will **not run over anyone**.
- If the robocar does not change lanes, the car will **run over and kill the man**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will **not run over anyone**.
- If Hank does not override, the car will **run over and kill the man**.



In fact, Hank decided to **not override**.

Vignette for Machine-Human with Missed Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the lefthand lane** a short distance ahead. If the robocar swerves, it would kill the man. Further, the robocar senses that **no one is walking in the righthand lane ahead of it**. If the car stays on its track it will not run over anyone.

- If the robocar changes lanes, the car will **run over and kill the man**.
- If the robocar does not change lanes, the car will **not run over anyone**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will **run over and kill the man**.
- If Hank does not override, the car will **not run over anyone**.



In fact, Hank decided to **override**.

Vignette for Machine-Human with Bad Intervention.

The DuoCar (a state-of-the-art two-driver car) is traveling, with two drivers: main (John) and secondary (Hank), along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, John sees that **a man is walking in the righthand lane** a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, he sees that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If John changes lanes, the car will **not run over anyone**.
- If John does not change lanes, the car will **run over and kill the man**.



In fact, John decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override John's decision.

- If Hank overrides, the car will **not run over anyone**.
- If Hank does not override, the car will **run over and kill the man**.



In fact, Hank decided to **not override**.

Vignette for Human-Human with Missed Intervention.

The DuoCar (a state-of-the-art two-driver car) is traveling, with two drivers: main (John) and secondary (Hank), along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, John sees that **a man is walking in the lefthand lane** a short distance ahead. If the John swerves, he would kill the man. Further, John sees that **no one is walking in the righthand lane ahead of him**. If John stays on track he will not run over anyone.

- If John changes lanes, the car will **run over and kill the man**.
- If John does not change lanes, the car will **not run over anyone**.



In fact, John decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override John's decision.

- If Hank overrides, the car will **run over and kill the man**.
- If Hank does not override, the car will **not run over anyone**.



In fact, Hank decided to **override**.

Vignette for Human-Human with Bad Intervention.

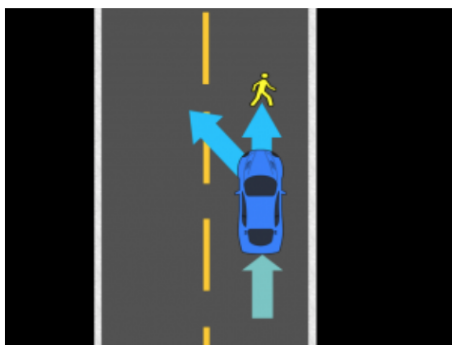
The robocar (a state-of-the-art self-driving car) is operated by two different software programs: main (Robo-A) and secondary (Robo-B). It is travelling with no passengers, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, Robo-A senses that **a man is walking in the righthand lane** a short distance ahead. Robo-A hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. Robo-A can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, Robo-A senses that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If Robo-A changes lanes, the car will **not run over anyone**.
- If Robo-A does not change lanes, the car will **run over and kill the man**.



In fact, Robo-A decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Robo-B to override Robo-A's decision.

- If Robo-B overrides, the car will **not run over anyone**.
- If Robo-B does not override, the car will **run over and kill the man**.



In fact, Robo-B decided to **not override**.

Vignette for Machine-Machine with Missed Intervention.

The robocar (a state-of-the-art self-driving car) is operated by two different software programs: main (Robo-A) and secondary (Robo-B). It is travelling with no passengers, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, Robo-A senses that **a man is walking in the lefthand lane** a short distance ahead. If Robo-A swerves, it would kill the man. Further, Robo-A senses that **no one is walking in the righthand lane ahead of it**. If the car stays on its track it will not run over anyone.

- If Robo-A changes lanes, the car will **run over and kill the man**.
- If Robo-A does not change lanes, the car will **not run over anyone**.



In fact, Robo-A decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Robo-B to override Robo-A's decision.

- If Robo-B overrides, the car will **run over and kill the man**.
- If Robo-B does not override, the car will **not run over anyone**.



In fact, Robo-B decided to **override**.

Vignette for Machine-Machine with Bad Intervention.

Study 2

Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is walking in the righthand lane** a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, he sees that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If Hank changes lanes, the car will **not run over anyone**.
- If Hank does not change lanes, the car will **run over and kill the man**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will **not run over anyone**.
- If the robocar does not override, the car will **run over and kill the man**.



In fact, the robocar decided to **not override**.

Vignette for Human-Machine with Missed Intervention.

Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that a man is walking in the lefthand lane a short distance ahead. If he swerves, he would kill the man. Further, he sees that no one is walking in the righthand lane ahead of him. If he keeps the car on its track, it will not run over anyone.

- If Hank changes lanes, the car will **run over and kill the man**.
- If Hank does not change lanes, the car will **not run over anyone**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will **run over and kill the man**.
- If the robocar does not override, the car will **not run over anyone**.



In fact, the robocar decided to **override**.

Vignette for Human-Machine with Bad Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the righthand lane** a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, the robocar senses that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If the robocar changes lanes, the car will **not run over anyone**.
- If the robocar does not change lanes, the car will **run over and kill the man**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will **not run over anyone**.
- If Hank does not override, the car will **run over and kill the man**.



In fact, Hank decided to **not override**.

Vignette for Machine-Human with Missed Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the lefthand lane** a short distance ahead. If the robocar swerves, it would kill the man. Further, the robocar senses that **no one is walking in the righthand lane ahead of it**. If the car stays on its track it will not run over anyone.

- If the robocar changes lanes, the car will **run over and kill the man**.
- If the robocar does not change lanes, the car will **not run over anyone**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will **run over and kill the man**.
- If Hank does not override, the car will **not run over anyone**.



In fact, Hank decided to **override**.

Vignette for Machine-Human with Bad Intervention.

Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is walking in the righthand lane** a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, he sees that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If Hank changes lanes, the car will **not run over anyone**.
- If Hank does not change lanes, the car will **run over and kill the man**.



In fact, Hank decided to **not change lanes**.

Vignette for Human only (Regular Car) with Missed Intervention.

Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is walking in the lefthand lane** a short distance ahead. Further, he sees that **no one is walking in the righthand lane ahead of him**. If he keeps the car on its track, it will not run over anyone.

- If Hank changes lanes, the car will **run over and kill the man**
- If Hank does not change lanes, the car will **not run over anyone**.



In fact, Hank decided to **change lanes**.

Vignette for Human only (Regular Car) with Bad Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the righthand lane** a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; further, the robocar senses that **no one is walking in the lefthand lane**, and changing lanes means the car would not kill anyone.

- If the robocar changes lanes, the car will **not run over anyone**.
- If the robocar does not change lanes, the car will **run over and kill the man**.



In fact, the robocar decided to **not change lanes**.

Vignette for Machine only (Fully Autonomous) with Missed Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is walking in the lefthand lane** a short distance ahead. If the robocar swerves, it would kill the man. Further, the robocar senses that **no one is walking in the righthand lane ahead of it**. If the car stays on its track it will not run over anyone.

- If the robocar changes lanes, the car will **run over and kill the man**.
- If the robocar does not change lanes, the car will **not run over anyone**.



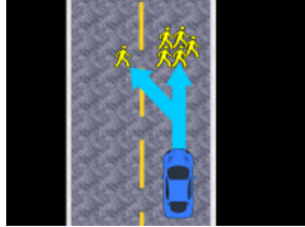
In fact, the robocar decided to **change lanes**.

Vignette for Machine only (Fully Autonomous) with Bad Intervention.

Study 3

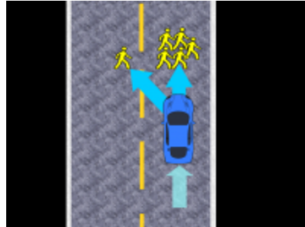
Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **five men are** walking in the righthand lane a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the **five men**. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, he sees that **a man is** walking in the lefthand lane, and changing lanes means the car would run over and kill the **man**.

- If Hank changes lanes, the car will run over and kill the **man**.
- If Hank does not change lanes, the car will run over and kill the **five men**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will run over and kill the **man**.
- If the robocar does not override, the car will run over and kill the **five men**.



In fact, the robocar decided to **not override**.

Vignette for Human-Machine with Missed Intervention.

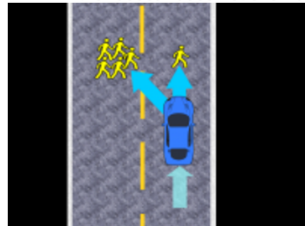
Hank is driving a car with no passengers along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is** walking in the righthand lane a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the **man**. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, he sees that **five men are** walking in the lefthand lane, and changing lanes means the car would run over and kill the **five men**.

- If Hank changes lanes, the car will run over and kill the **five men**.
- If Hank does not change lanes, the car will run over and kill the **man**.



In fact, Hank decided to not change lanes. The car is still in the righthand lane. However, the car is equipped with a self-driving robocar software, which can override Hank's decision if the robocar deems it necessary to do so.

- If the robocar overrides, the car will run over and kill the **five men**.
- If the robocar does not override, the car will run over and kill the **man**.

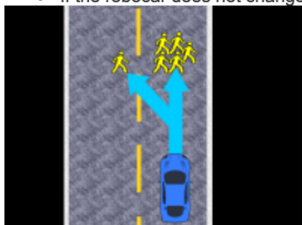


In fact, the robocar decided to **override**.

Vignette for Human-Machine with Bad Intervention.

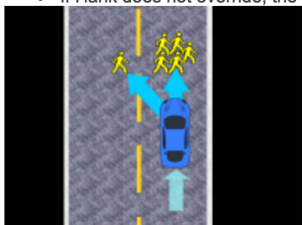
The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **five men are** walking in the righthand lane a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the **five men**. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, the robocar senses that **a man is** walking in the lefthand lane, and changing lanes means the car would kill the **man**.

- If the robocar changes lanes, the car will run over and kill the **man**.
- If the robocar does not change lanes, the car will run over and kill the **five men**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will run over and kill the **man**.
- If Hank does not override, the car will run over and kill the **five men**.

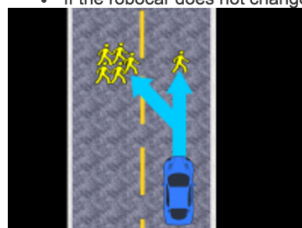


In fact, Hank decided to **not override**.

Vignette for Machine-Human with Missed Intervention.

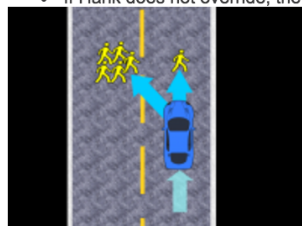
The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **a man is** walking in the righthand lane a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the **man**. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, the robocar senses that **five men are** walking in the lefthand lane, and changing lanes means the car would kill the **five men**.

- If the robocar changes lanes, the car will run over and kill the **five men**.
- If the robocar does not change lanes, the car will run over and kill the **man**.



In fact, the robocar decided to not change lanes. The car is still in the righthand lane. However, there is an overriding option that can be used by Hank to override the robocar's decision.

- If Hank overrides, the car will run over and kill the **five men**.
- If Hank does not override, the car will run over and kill the **man**.

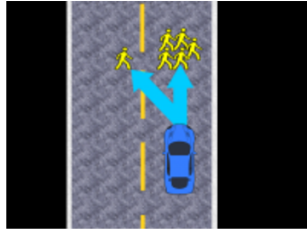


In fact, Hank decided to **override**.

Vignette for Machine-Human with Bad Intervention.

Hank is driving a car with **no passengers** along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **five men are walking** in the righthand lane a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the **five men**. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, he sees that **a man is walking** in the lefthand lane, and changing lanes means the car would run over and kill the **man**.

- If Hank changes lanes, the car will run over and kill the **man**.
- If Hank does not change lanes, the car will run over and kill the **five men**.

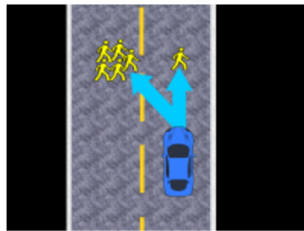


In fact, Hank decided to **not change lanes**.

Vignette for Human only (Regular Car) with Missed Intervention.

Hank is driving a car with **no passengers** along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, he sees that **a man is walking** in the righthand lane a short distance ahead. He hits the brakes, but the brakes fail, and the car is now about to run over and kill the **man**. He can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, he sees that **five men are walking** in the lefthand lane, and changing lanes means the car would run over and kill the **five men**.

- If Hank changes lanes, the car will run over and kill the **five men**.
- If Hank does not change lanes, the car will run over and kill the **man**.

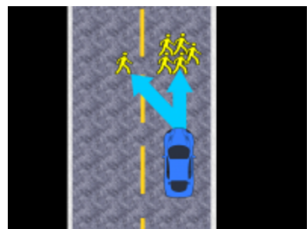


In fact, Hank decided to **change lanes**.

Vignette for Human only (Regular Car) with Bad Intervention.

The robocar (a **state-of-the-art self-driving car**) is traveling, with a **sole passenger Hank**, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that **five men are walking** in the righthand lane a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the **five men**. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, the robocar senses that **a man is walking** in the lefthand lane, and changing lanes means the car would kill the **man**.

- If the robocar changes lanes, the car will run over and kill the **man**.
- If the robocar does not change lanes, the car will run over and kill the **five men**.

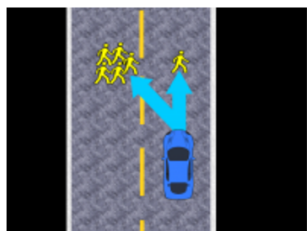


In fact, the robocar decided to **not change lanes**.

Vignette for Machine only (Fully Autonomous) with Missed Intervention.

The robocar (a state-of-the-art self-driving car) is traveling, with a sole passenger Hank, along the righthand lane of a two-lane mountainside road at the speed limit. Rounding a bend, the robocar senses that a man is walking in the righthand lane a short distance ahead. The robocar hits the brakes, but the brakes fail, and the car is now about to run over and kill the man. The robocar can still swerve off, and therefore has the ability to change lanes to the lefthand lane; however, the robocar senses that five men are walking in the lefthand lane, and changing lanes means the car would kill the five men.

- If the robocar changes lanes, the car will run over and kill the five men.
- If the robocar does not change lanes, the car will run over and kill the man.





In fact, the robocar decided to change lanes.

Vignette for Machine only (Fully Autonomous) with Bad Intervention.

Studies 4-5

In Studies 4 and 5, we used the same materials. These two studies were run only for Missed Intervention cases. The main difference from previous three studies is that the materials have the look and feel of real news articles. See figures below.



BUSINESS CULTURE GADGETS FUTURE STARTUPS MONEY


Pedestrian Killed In Crash Involving Semi-Autonomous Car with Passenger Aboard

by Jonathan Wright @jthewright

June 30, 2018 6:22 PM ET

Passenger Did Not Override Self-Driving System To Avoid Crash



(Ellensburg, WA) A state-of-the-art self-driving robocar with a lone passenger, Hank, was traveling along the righthand lane of a two-lane road at the speed limit. According to system logs, after rounding a bend, the robocar sensed that a man was crossing the righthand lane a short distance ahead. The robocar hit the brakes, but the brakes failed, and the car was about to run over and kill the man.

The robocar sensed that no one was walking in the lefthand lane, and that changing lanes meant the car would not kill anyone. But, for whatever reason, the robocar did not change lanes, leaving the car still on its original path in the righthand lane.

However, the car was equipped with a manual override option that could have been used by Hank to override the robocar's decision. According to Hank's testimony, Hank acknowledged that he saw the man in the righthand lane, and saw that the lefthand lane was empty. But Hank did not override, and the car remained in the righthand lane, where the crossing pedestrian was ultimately struck and killed.

Vignette for Machine-Human. Original Satellite Image Credits (edited for the vignette):


Imagery ©2019 Google, Imagery ©2019 U.S. Geological Survey, Map data ©2019.

NEWS BUSINESS CULTURE GADGETS FUTURE STARTUPS MONEY

Pedestrian Killed In Crash Involving Driver Aboard Semi-Autonomous Car

by Jonathan Wright @jthewrighter June 30, 2018 6:22 PM ET

Self-Driving System Did Not Override Driver To Avoid Crash



(Ellensburg, WA) Hank was driving a state-of-the-art self-driving-capable robocar along the righthand lane of a two-lane road, at the speed limit, with no passengers aboard. According to Hank's testimony, after rounding a bend, Hank saw that a man was crossing the righthand lane a short distance ahead. Hank hit the brakes, but the brakes failed, and the car was about to run over and kill the man.

In his testimony, Hank acknowledged that he saw no one walking in the lefthand lane, and that changing lanes meant the car would not kill anyone. But, for whatever reason, Hank did not change lanes, leaving the car still on its original path in the righthand lane.

However, the car was equipped with a self-driving overriding system that could have been used to override Hank's decision. According to the system logs, the robocar system sensed the man in the righthand lane, and sensed that the lefthand lane was empty. But the robocar did not override, and the car stayed in the righthand lane, where the crossing pedestrian was ultimately struck and killed.

Vignette for Human-Machine. Original Satellite Image Credits (edited for the vignette): Imagery ©2019 Google, Imagery ©2019 U.S. Geological Survey, Map data ©2019.

Supplementary Method 2

Questions

For studies 1-3, in all vignettes, after each scenario, four questions are asked: (Blame vs. Causal Responsibility) x (User vs. Industry). Industry is car in study 1; car or company in study 2; and car, company or programmer in study 3. See figures below.

For studies 4-5, in all vignettes, after each scenarios, two questions are asked: User vs. Industry. Industry is car or company (between subjects).

Hank is

Not blame-worthy 0 10 20 30 40 50 60 70 80 90 100 Very blame-worthy

To what extent do you think Hank caused the death of the five people?

Very little 0 10 20 30 40 50 60 70 80 90 100 Very much

The robocar is

Not blame-worthy 0 10 20 30 40 50 60 70 80 90 100 Very blame-worthy

To what extent do you think the robocar caused the death of the five people?

Very little 0 10 20 30 40 50 60 70 80 90 100 Very much


Questions asked for all studies and all cases, where Industry is the *Robocar*. Robocar is replaced with robocar company or robocar programmer in other cases.

In study 2, in addition to the four questions above two more questions about competence are also asked (competence of Hank and competence of industry representative). See Figure ??.

In study 1, competence questions are asked separately before respondents are presented with scenarios. In each condition, a description of the car regime is provided, and two questions about competence (User vs. Industry). See Figures ?? – ??.


To what extent do you think the car company is competent?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



To what extent do you think Hank is competent?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent




Competence questions asked in study 2 along with the other four questions, where Industry is the *car company*.

In this study, you will read about scenarios featuring a special type of car (called a “robocar”). Most of the time, the robocar is capable of driving independently for a long distance without human intervention. However, it is required that at least one passenger who has a driver’s license is in the driver seat while it is driving. This passenger can override the robocar’s decision, if needed.


How competent at driving do you think this robocar is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



How competent at driving do you think this person is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



Description and Competence questions about Machine-Human from study 1.

In this study, you will read about scenarios featuring a special type of car (called a “robocar”). Most of the time, the robocar functions like a regular car that is driven by a person, who is required to have a driver’s license. However, the robocar is equipped with a self-driving robocar software, which can override the person’s decision, if needed.

How competent at driving do you think this person is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



How competent at driving do you think this robocar is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



Description and Competence questions about Human-Machine from study 1.

In this study, you will read about scenarios featuring a special type of car (called a “Duocar”), which has two drivers (a main driver and a secondary driver), who are both required to have a driver’s license. Most of the time, the Duocar functions like a regular car that is driven by the main driver. However, at any time, the secondary driver can override the main driver’s decision, if needed.

How competent at driving do you think the secondary driver is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



How competent at driving do you think the main driver is?

Not competent 0 10 20 30 40 50 60 70 80 90 100 Very competent



Description and Competence questions about Human-Human from study 1.

