

The Ethical Dilemma of Self-Driving Cars *by Patrick Lin - TedEd*



<https://ed.ted.com/on/ZNWXLxDk>

This is a thought experiment. Let's say at some point in the not so distant future, you're barreling down the highway in your self-driving car, and you find yourself boxed in on all sides by other cars. Suddenly, a large, heavy object falls off the truck in front of you. Your car can't stop in time to avoid the collision, so it needs to make a decision: go straight and hit the object, swerve left into an SUV, or swerve right into a motorcycle. Should it prioritize your safety by hitting the motorcycle, minimize danger to others by not swerving, even if it means hitting the large object and sacrificing your life, or take the middle ground by hitting the SUV, which has a high passenger safety rating? So what should the self-driving car do?

If we were driving that boxed in car in manual mode, whichever way we'd react would be understood as just that, a reaction, not a deliberate decision. It would be an instinctual panicked move with no forethought or malice. But if a programmer were to instruct the car to make the same move, given conditions it may sense in the future, well, that looks more like premeditated homicide. Now, to be fair, self-driving cars are predicted to dramatically reduce traffic accidents and fatalities by removing human error from the driving equation. Plus, there may be all sorts of other benefits: eased road congestion, decreased harmful emissions, and minimized unproductive and stressful driving time. But accidents can and will still happen, and when they do, their outcomes may be determined months or years in advance by programmers or policy makers. And they'll have some difficult decisions to make. It's tempting to offer up general decision-making principles, like minimize harm, but even that quickly leads to morally murky decisions. For example, let's say we have the same initial set-up, but now there's a motorcyclist wearing a helmet to your left and another one without a helmet to your right. Which one should your robot car crash into? If you say the biker with the helmet because she's more likely to survive, then aren't you penalizing the responsible motorist? If, instead, you save the biker without the helmet because he's acting irresponsibly, then you've gone way beyond the initial design principle about minimizing harm, and the robot car is now meting out street justice. The ethical considerations get more complicated here. In both of our scenarios, the underlying design is functioning as a targeting algorithm of sorts. In other words, it's systematically favouring or discriminating against a certain type of object to crash into. And the owners of the target vehicles will suffer the negative consequences of this algorithm through no fault of their own. Our new technologies are opening up many other novel ethical dilemmas. For instance, if you had to choose between a car that would always save as many lives as possible in an accident, or one that

would save you at any cost, which would you buy? What happens if the cars start analyzing and factoring in the passengers of the cars and the particulars of their lives? Could it be the case that a random decision is still better than a predetermined one designed to minimize harm? And who should be making all of these decisions anyhow? Programmers? Companies? Governments? Reality may not play out exactly like our thought experiments, but that's not the point. They're designed to isolate and stress test our intuitions on ethics, just like science experiments do for the physical world. Spotting these moral hairpin turns now will help us manoeuvre the unfamiliar road of technology ethics, and allow us to cruise confidently and conscientiously into our brave new future.