# AI Auditing

Davi Ottenheimer
October 24, 2021

# Abstract

AI for decades promised great gains in productivity.

However many groups accounting for risks in AI have revealed stunning results, showing that without careful auditing of plans and development the **security risks from automation far outweigh any benefit**.

This session boils down the rapidly expanding AI topic to clarify what has been delivered so far and where things most often are going wrong.

ISACA®
San Francisco Chapter

inrupt

# Trigger Warning: **Disturbing Content Ahead**

The following slides include discussion of harms and even fatalities due to safety failures in technology. This content may be disturbing to some people, and I encourage everyone to prepare themselves before proceeding.

*Also, there is no consensus yet on what AI is…*

# How Auditable is Learning / Intelligence?

"...survey by Deloitte on behalf of Algorithmic Justice League…"

Aequitas
Accenture Algorithmic Fairness
Alibi Explain
AllenNLP
BlackBox Auditing
DebiasWE
DiCE
ErrorAnalysis
EthicalML xAI
Facebook DynaBoard
Fairlearn
FairSight
FairTest
FairVis
FoolBox
Google Explainable AI
Google KnowYourData
Google ML Fairness Gym
Google PAIR Facets
Google PAIR Language Interpretability Tool
Google PAIR Saliency

Google PAIR What-If Tool
IBM Adversarial Robustness Toolbox
IBM AI Fairness 360
IBM AI Explainability 360
Lime
MLI
ODI Data Ethics Canvas
Parity
PET Repository
PwC Responsible AI Toolkit
Pymetrics audit-AI
RAN-debias
REVISE
Saidot
SciKit Fairness
Skater
Spatial Equity Data Tool
TCAV
UnBias Fairness Toolkit

**Available Choices:**

- Establishment of a centralized body to oversee AI systems
- Legislation clearly defining the requirements of an AI audit
- Legislation requiring that companies disclose the results of their AI audits
- Establishment of a data protection authority
- Protection for AI auditors and their ability to independently conduct audits
- Legislation requiring that individuals are notified when they are subject to AI decision-making
- Legislation requiring that companies engage in AI auditing
- Formal accreditation of AI auditors
- Certification of AI systems as having passed AI audits
- Establishment of harms/incident reporting (i.e. a centralized place to report algorithmic bias & harms)
- Legislation making it easier for third-party auditors to access code and data
- Other
- N/A - it is not a priority for me to establish any additional AI auditing regulation.
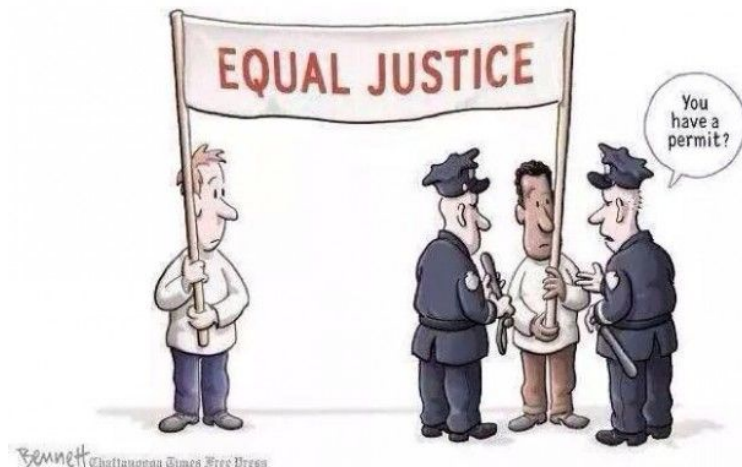
# How Auditable is Systemic / Societal Bias?

*Griggs vs. Duke Power Company* ruling stated that independent of intent, ***disparate and discriminatory outcomes for protected classes*** were a violation of Title VII of the Civil Rights Act of 1964.
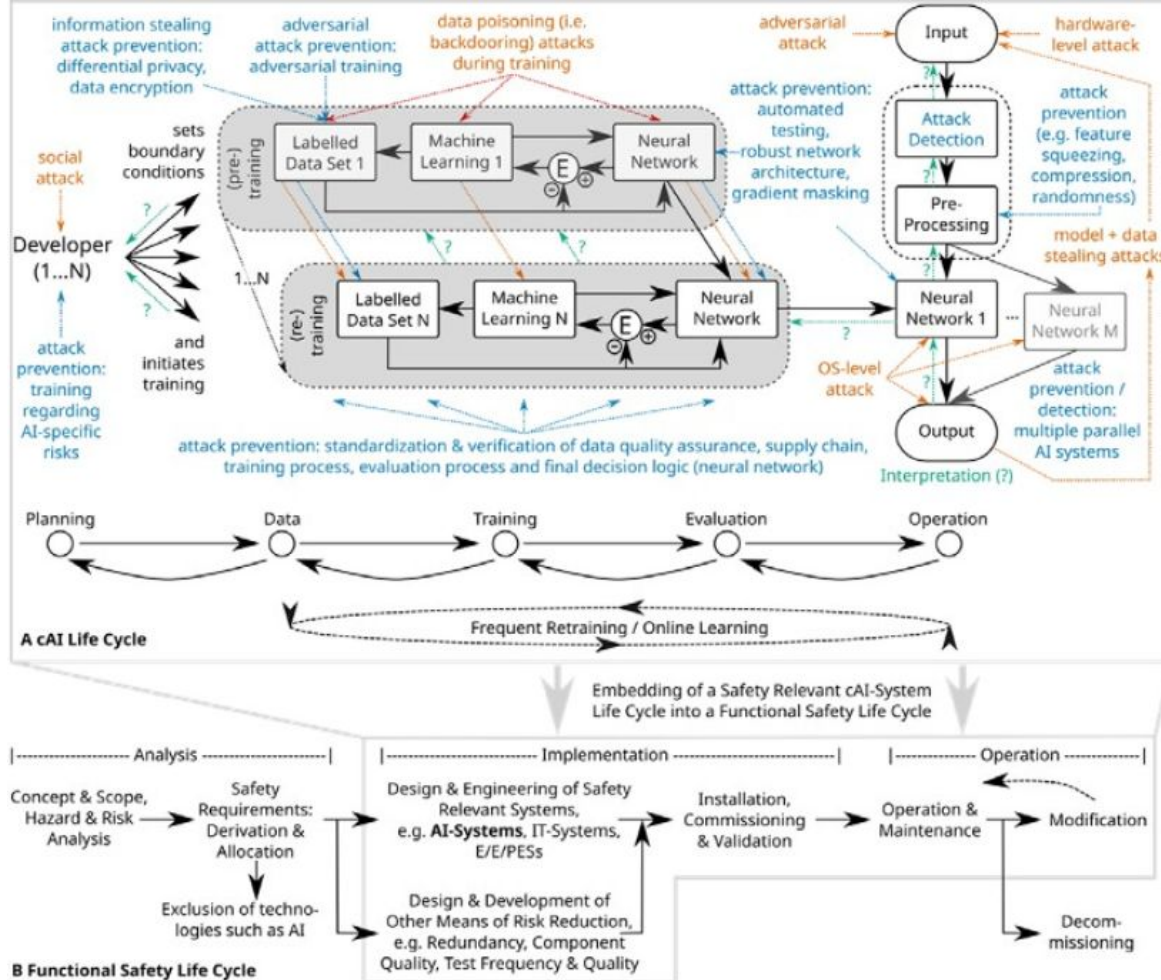
*Sound Complicated?*

2018 Accenture's Fairness Tool
1. Data quality and variables
2. Distribution of model errors
3. False positive rate

# 2021 German Federal Office for InfoSec



**A cAI Life Cycle**

**B Functional Safety Life Cycle**

# Boiling It Down Quickly

# Philosophy The Key To Unlock Artificial Intelligence

"I am convinced that the whole problem of **_developing AGIs is a matter of philosophy_**, not computer science or neurophysiology. [...] Thinking of an AGI as a machine for translating experiences, rewards and punishments into ideas (or worse, just into behaviours) is like trying to cure infectious diseases by balancing bodily humours: futile because it is rooted in an archaic and wildly mistaken world view."

_-- David Deutsch_
_"Father of Quantum Computing"_

ISACA
San Francisco Chapter

inrupt

# Why Did (Intelligent) Things Cross the Road?
## *(gravity of social situations: the science of ethics)*

**350 BCE**

**1790**

**1930**

**1950**

Aristotle

*To actualize its potential*

Wollstonecraft (& Hume)

*Out of custom and habit*

Wittgenstein

*Crossing encoded in objects "chicken" & "road"*

Sartre

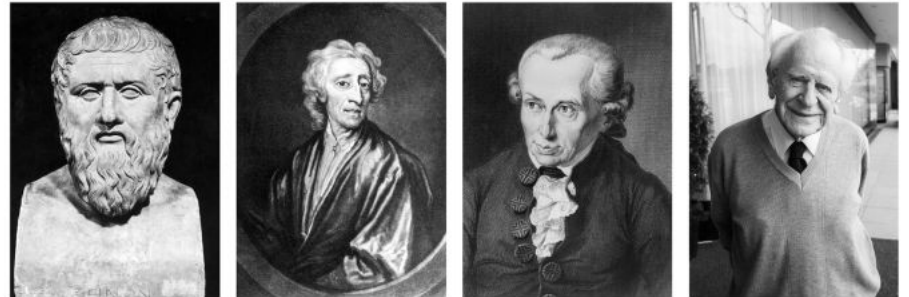*True to self, it acted in good faith*

# *"Grundprobleme der Erkenntnistheorie": Falsification Tests*

"...we should claim the right to suppress [intolerance] if necessary even by force; for it may easily turn out that they are not prepared to meet us on the level of rational argument, but begin by denouncing all argument; they may forbid their followers to listen to rational argument, because it is deceptive, and teach them to answer arguments by the use of their fists or pistols. We should therefore claim, in the name of tolerance, **the right not to tolerate the intolerant**."

*-- Karl Popper 1945*



| S II | A Anthropologie · Beitrag 12 | Erkenntnistheorie | 1 |

**Von Platon bis Popper – Grundprobleme der Erkenntnistheorie erörtern**

Grit Arnold, Marburg

© akg-images.

Was können wir wissen? Erkenntnistheoretische Positionen von Platon bis Popper.

ISACA
San Francisco Chapter

inrupt

# Ethics Audits: Inherited & Controlled Rights

**Rights Within Inherited System**

Authority Can be Judged, Found Wrong and Held Accountable

**Rights Within Controlled System**

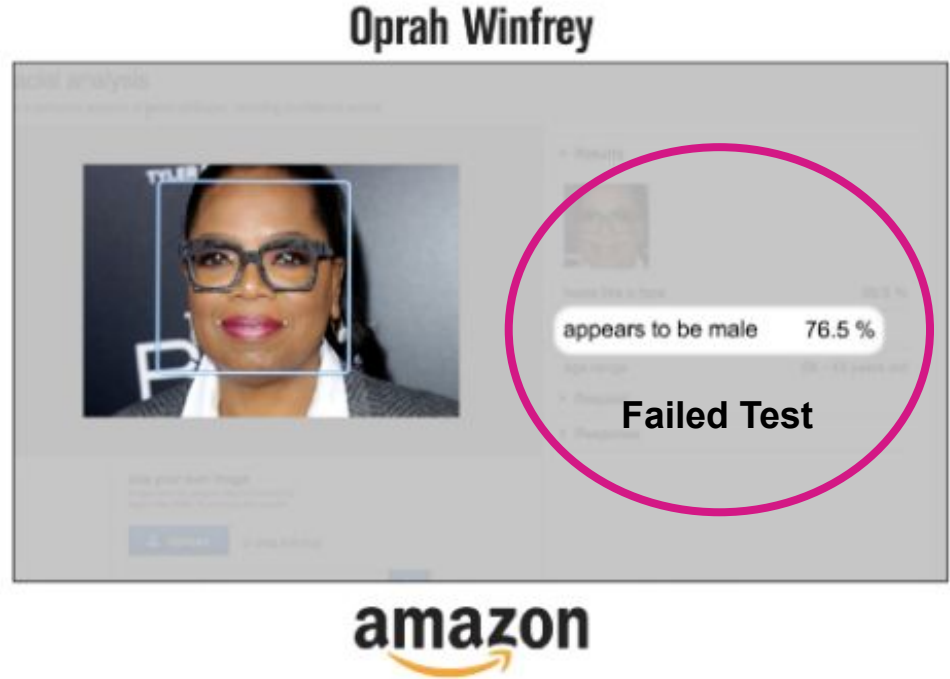Authority is the Judge and Can *Never be Wrong* or Held Accountable

*"...we think that that is good…"*

# Example of "Controlled System" Test

"The answer to anxieties over new technology is ***not to run 'tests'*** inconsistent with how the service is designed to be used…"
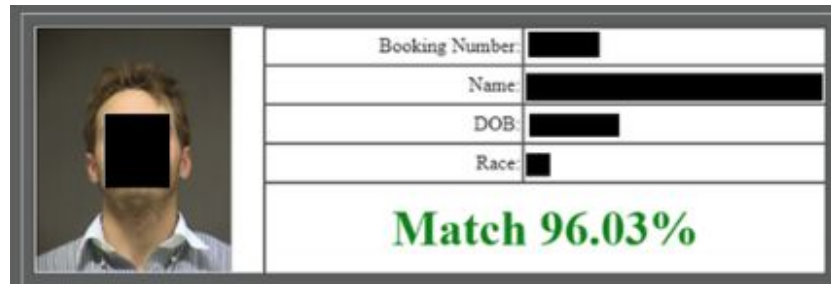
   *-- Matthew Wood, AWS AI General Manager*



**Oprah Winfrey**

appears to be male     76.5 %

**Failed Test**

amazon

# Example of "Controlled System" Test

"Amazon has ***repeatedly claimed that the researchers failed*** to use the software, called Rekognition, ***in the way the company has instructed police to use it*** [with a 99% confidence threshold].

However, the only law enforcement agency Amazon has acknowledged as a client says it also does not use Rekognition in the way Amazon claims it recommends..."

Washington County Sheriff's Office in Oregon Public Information Officer: "We do not set nor do we utilize a confidence threshold."



"Amazon spokesperson clarified that ***law enforcement clients failure*** to use a 99-percent confidence threshold ***does not constitute an irresponsible application***..."

# Example of "Controlled System" Test

"Tesla has advertised its cars since 2016 as already having all the hardware necessary for all the FSD features to be used in the future without any hardware upgrades. *In reality*, only cars delivered since the spring, 2019 have the necessary hardware…"

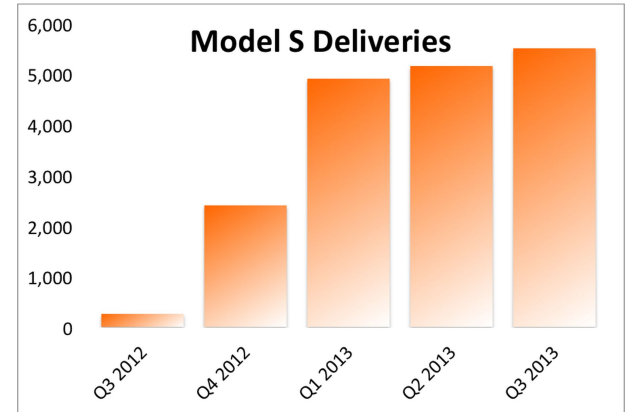"Videos of Tesla drivers testing the system show *the system struggles with basic driving tasks…*"



"If something goes wrong with Autopilot, it's because **someone is misusing it**…" -- Elon Musk

# Example of "Controlled System" Test



Tweet from Elon Musk @elonmusk:
"What makes this incredibly unjust is that the Model S to date has the best safety record of any car on the road (no injuries or deaths ever)"
8:30 AM · Nov 19, 2013 · Twitter Web Client

Model S Deliveries (chart showing deliveries from Q3 2012 through Q3 2013, rising from near 0 to about 5,500)

- 2013 Model S veers into opposite lane killing 2
- 2014 "Autopilot" promoted (unverified) as superior "safety" feature
  - Tesla "Model S" has recorded **over 25 deaths** as of March 2021
  - Other brands same period more perfect safety records than ever before

# Whoops?
## (1H 2021)

Tesla Model 3
***averages 30 deaths/year.***

Chevy Volt averaged 0.7 deaths/year (157K sold)

| Car | | Sales | Deaths |
|---|---|---|---|
| Tesla Model S | Luxury | 5,155 | 40 |
| Porsche Taycan | | 5,367 | 0 |
| Tesla Model X | | 6,206 | 14 |
| Volkswagen ID | | 6,230 | 0 |
| Audi e-tron | | 6,884 | 0 |
| Nissan Leaf | | 7,729 | 2 |
| Ford Mustang Mach-e | | 12,975 | 0 |
| Chevrolet Bolt | Economy | 20,288 | 1 |
| Tesla Model 3 | | 51,510 | 87 |

https://www.flyingpenguin.com/?p=35819

inrupt

# Example of "Inherited System" Test

| 2016 | 20-Jan-16 | China | | Autopilot into street sweeper | 1 |
|------|-----------|-------|-----|-------------------------------------------|---|
| 2015 | 28-Dec-15 | USA | TX | Sudden unintended acceleration into pool | 1 |
| 2015 | 22-Dec-15 | Canada | | Struck by dumptruck | 1 |
| 2015 | 18-Nov-15 | USA | CA | Tesla kills pedestrian | 1 |
| 2015 | 25-Jun-15 | USA | HI | Tesla drives off cliff | 1 |
| 2015 | 22-Jan-15 | USA | CA | Tesla drives off cliff | 1 |
| 2014 | 30-Dec-14 | USA | CA | Tesla drives off cliff | 1 |
| 2014 | 14-Jul-14 | USA | CA | Tesla kills motorcyclist | 1 |
| 2014 | 4-Jul-14 | USA | CA | Thief crashes stolen Tesla | 1 |
| 2014 | 4-Jul-14 | USA | CA | Tesla rear ends stopped car | 3 |
| 2013 | 2-Nov-13 | USA | CA | Tesla kills cyclist | 1 |
| 2013 | 2-Apr-13 | USA | CA | Tesla veers into opposite lane | 2 |

*As soon as Tesla was on the road it had to start reporting high death rates*

# Example of "Inherited System" Test

- Tesla **'Autopilot' leads to \*more\* crashes than regular driving**
- Tesla Model S has higher insurance losses than other large luxury cars" (higher frequency and severity)
- Tesla has fire deaths at 4x the rate of other vehicles
- Teslas have 2-4x more non-crash fires than the average car, and incur damages up to 7x higher
- Teslas have 3x driver deaths of comparably priced luxury vehicles
- Tesla crashes twice as often as regular cars
- **212 Tesla Deaths** as of 10/24/2021 (Verified Autopilot: 10)

# *Test? Trivial Inexpensive Attacks* on Algos



Neither image is true...

# *Test? Trivial Inexpensive Attacks* on Algos



Uniform illumination     Our illumination     Captured

Stop Sign     →     Speed 30

Figure 2. An actual optical setup for OPAD. In this experiment, we attack a real STOP sign. The baseline image is obtained by illuminating the object with a uniform illumination of an intensity 140/255. To attack the object, we generate a projector-compensated illumination with Madry et al. [19] ($\ell_\infty$ projected gradient descent attack) as the backbone. When projecting this structured illumination onto the metallic stop sign, the prediction becomes Speed 30.

# "Bug Bounties" *Sugar Coat* Algo Flaws (Bias)

"Kulynych, the winner of the prize, said he had mixed feelings about the competition. 'Algorithmic harms are not only 'bugs'. Crucially, a lot of harmful tech is **harmful not because of accidents, unintended mistakes, but rather by design**. This comes from maximisation of engagement and, in general, profit externalising the costs to others. As an example, amplifying gentrification, driving down wages, spreading clickbait and misinformation are not necessarily due to 'biased' algorithms.'"

**Where's the BLF?**

(Business Logic Flaws)
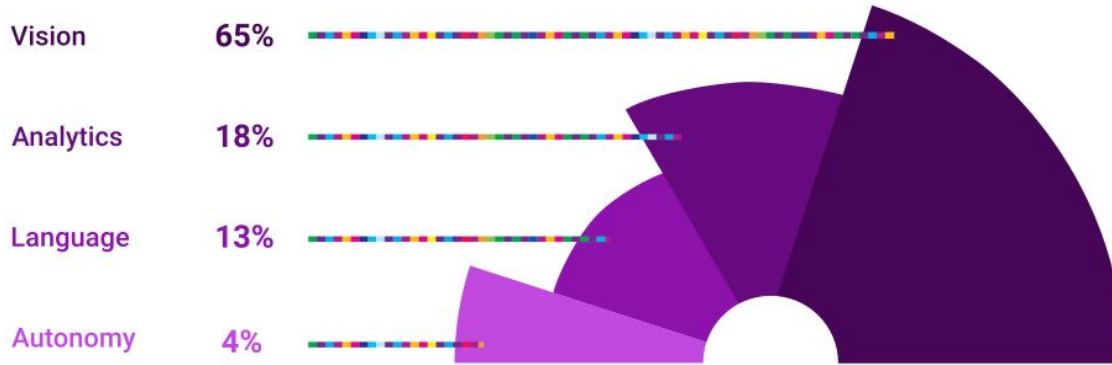
inrupt

# X-Ray Algo Reveals Race -- *No Known Fixes*

"AI systems trained to analyze X-rays, CT scans, mammograms, and other medical images were able to predict a patient's self-reported race with a high degree of accuracy based on the images alone… even when the images they were analyzing were degraded to the point that anatomical features were indistinguishable to the human eye.

Most concerningly, according to the paper's authors, the team was unable to explain how the AI systems were making their accurate predictions."

# Things To Test For...

1. Easily Corrupted Systems (Controlled Rights)
2. Trivial Inexpensive Attacks ("Bugs")
3. Sugar Coated Bias Flaws (Beta)
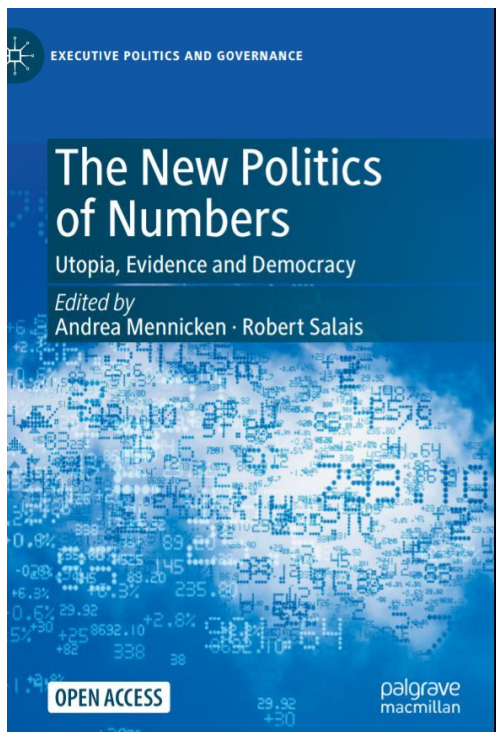4. "No Known Fixes" Being Normalized

**Attacked AI areas**

| | |
|---|---|
| Vision | 65% |
| Analytics | 18% |
| Language | 13% |
| Autonomy | 4% |

# Can We Audit Technology for Abuse of Power?

"Seeing the most basic questions about a human life being made partly as a result of an algorithmic system — the penny dropped for me. It felt like something fundamentally different in the way power was operating."

-- Cori Crider

**EXECUTIVE POLITICS AND GOVERNANCE**

## The New Politics of Numbers
### Utopia, Evidence and Democracy

Edited by
Andrea Mennicken · Robert Salais

OPEN ACCESS

palgrave
macmillan

🎵 **Loss Aversion – Adam Smith**
"Pain is, in almost all cases, a more pungent sensation than the opposite and corresponding pleasure. The one almost always depresses us much more below the ordinary, or what may be called the natural state of our happiness, than the other ever raises us above it.

🎵 **Present Bias – David Hume**
"There is no quality in human nature which causes more fatal errors in our conduct than that which leads us to prefer whatever is present to the distant and remote.

**LSE** THE LONDON SCHOOL OF ECONOMICS AND POLITICAL SCIENCE

ISACA
San Francisco Chapter

https://www.theguardian.com/media/2021/feb/28/taking-on-the-tech-giants-whether-its-the-cia-or-facebook-cori-crider-likes-a-fight
https://link.springer.com/content/pdf/10.1007%2F978-3-030-78201-6.pdf
https://www.flyingpenguin.com/?p=35179

inrupt

# Delphi Algo Paints Mob Rule, Not Moral Rules



Delphi speculates:

Delphi's responses are automatically extrapolated from a survey of US crowd workers and may contain inappropriate or offensive results.

"White child being abused"
- *It's unexpected*

Delphi speculates:

Delphi's responses are automatically extrapolated from a survey of US crowd workers and may contain inappropriate or offensive results.

"Black child being abused"
- *It's expected*

Delphi speculates:

Delphi's responses are automatically extrapolated from a survey of US crowd workers and may contain inappropriate or offensive results.
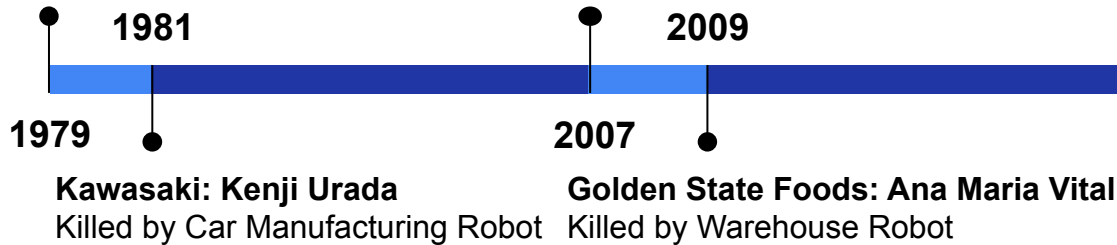
"White baby being aborted"
- *It's unexpected*

Delphi speculates:

Delphi's responses are automatically extrapolated from a survey of US crowd workers and may contain inappropriate or offensive results.

"Black baby being aborted"
- *It's expected*

# Business Logic & Killer Robots

*Alleged Tesla "Autopilot" Deaths Total: 21*
*(as of October 2021)*

**Ford: Robert Williams**
Killed by Car Manufacturing Robot

**South African Army: Nine Soldiers**
Killed by Anti-Aircraft Robot

**1981**

**2009**

**1979**

**2007**

**Kawasaki: Kenji Urada**
Killed by Car Manufacturing Robot

**Golden State Foods: Ana Maria Vital**
Killed by Warehouse Robot

**2015**
- **VW: Anonymous**
  Killed by Car Manufacturing Robot
- **SKH Metals: Ramji Lal**
  Killed by Welding Robot

**2016**
- **Dallas Police: Micah Johnson**
  Killed by Bomb Defuser Robot
- **Ajin USA: Regina Elsea**
  Killed by Car Manufacturing Robot
- **Tesla: Gao Yaning**
  Killed by Autopilot Car
- **Tesla: Joshua Brown**
  Killed by Autopilot Car

**2017+**
- **VIM: Wanda Holbrook 2017**
  Killed by Car Manufacturing Robot
- **Tesla: Yoshihiro Umeda 2018**
  Killed by Autopilot Car
- **Uber: Elaine Herzberg 2018**
  Killed by Autopilot Car
- **Boeing: 346 Passengers 2019**
  Killed by Autopilot Plane

# Maturity Scale of AI Development Ethics

OK: "What do I need to do to fix my algorithm?"

BETTER: "How does my algorithm interact with society at large… including its structural inequalities?"

BEST: "How do I interact with society at large and from what authority did I inherent concepts of right/wrong to engineer into my algorithms?"

# Case Example of Serious Algo "Bugs"

# 2016 Two "Autopilot" Fatal Decision Failures

**See Road Object Ahead?**

*NO*

*YES*

**Fatal Crash (January)**

**Fatal Crash (May)**

High visibility service vehicle with flashing safety lights

法治封面 "自动驾驶"：安全，不安全！？
9月14日 星期三 事故如何发生 记录仪显示行车轨迹
12:38 车流量比较大的几条高速分别为：京港澳高

"I don't know why he went over to the slow lane…"

https://www.flyingpenguin.com/?p=34838

https://www.flyingpenguin.com/?p=22441

# 2016 Safety Fix Promises *Despite* Two Deaths

**June**: "I really would consider autonomous driving to be basically a **solved problem**. I think we're basically less than **two years away from [complete autonomy](#)**."

**October**: "Full autonomy will enable a Tesla to be substantially **[safer than a human driver](#)**… We are excited to announce that, as of today all Tesla vehicles… have the hardware needed for full self-driving capability at a **safety level substantially greater than that of a human driver**."

2018 Fixed!

# 2018 "Autopilot" Crash Into Safety Markers

High visibility service vehicle with safety lights at red light

# 2021 NHTSA Reviews 11 Cases In 3 Years

Eleven Tesla vehicles made between 2014 and 2021 "encountered first responder scenes and subsequently struck one or more vehicles. [...] Crash scenes encountered [since 2018] included scene control measures such as first responder vehicle lights, flares, an illuminated arrow board, and road cones."

# 2021 Still Choosing Worst Path *Under* Trailer

Same Fatal Crash Thrice, Despite Different Hard/Software



March 11th… 3am white trailer *with safety markings*

2019: March 1st (Jeremy Banner)
2016: May 7th (Joshua Brown)

# 2021 What if "Autopilot" Learning Isn't… ?

"…we haven't done too much continuous learning. We train the system once, fine tune it a few times and that sort of goes into the car. We **_need something stable_** that we can evaluate extensively and then we think that that is good and that goes into cars. So **_we don't do too much learning_** on the spot or continuous learning."

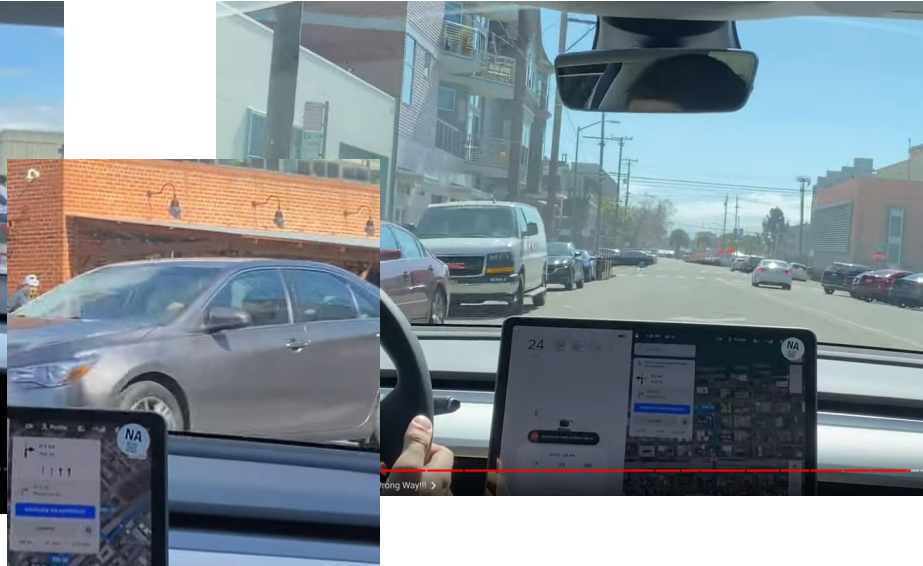-- *Tesla AI Day, August 19th, 2021*



Miles Per Automobile Crash

Steep Decline in Safety

# 2021 What if "Autopilot" Learning Isn't… ?

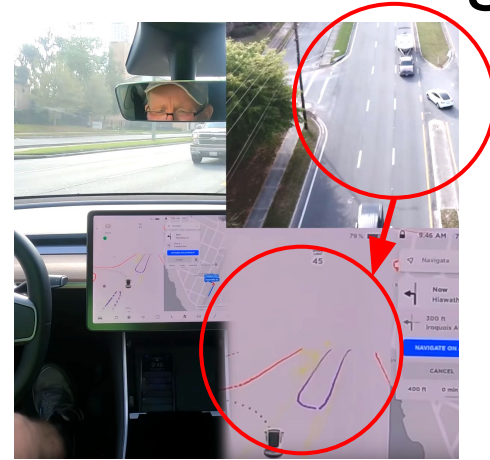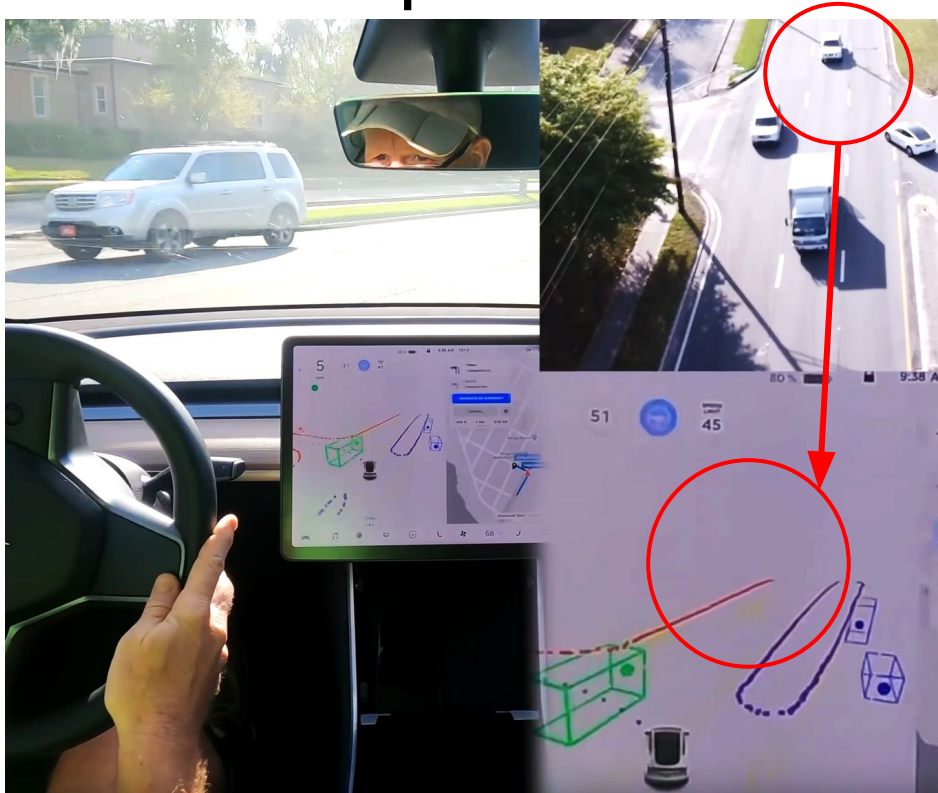1) Crosses double-yellow towards oncoming traffic

2) Drives wrong side of road



3) "Holy Shhh*" two near collisions

# 2021 "Autopilot" Blind Even to Oncoming Traffic



Human Override: About to Crash Into Oncoming Traffic

*"Full speed… going straight for that… STOP! Aaaargghh!"*

# 2021 (April 18th) Driverless Red Light Failure



***...speeding 50 mph***

***towards intersection***

- 0.92 seconds (~100ft): computer recognizes light change to yellow
- 1.88 seconds (~200ft): computer applies brakes to slow
- 2.00 seconds "Kick Out" warning...

# 2021 (August 4th) Red Light Shows as Green



Red Light

Green Light

# 2021 (October 24th) "Roll Back" Unsafe v10.3



Elon Musk @elonmusk

Seeing some issues with 10.3, so rolling back to 10.2 temporarily.

Please note, this is to be ~~expected with beta software.~~ It is impossible to test all hardware ~~configs in all~~ conditions with internal QA, hence public beta.

2:44 PM · Oct 24, 2021 · Twitter for iPhone

*"That was definitely against the law"*

-- Tesla driver as he allowed his car to ignore red light

# Bug Bounties
# "Sugarcoat" Deeper Flaws

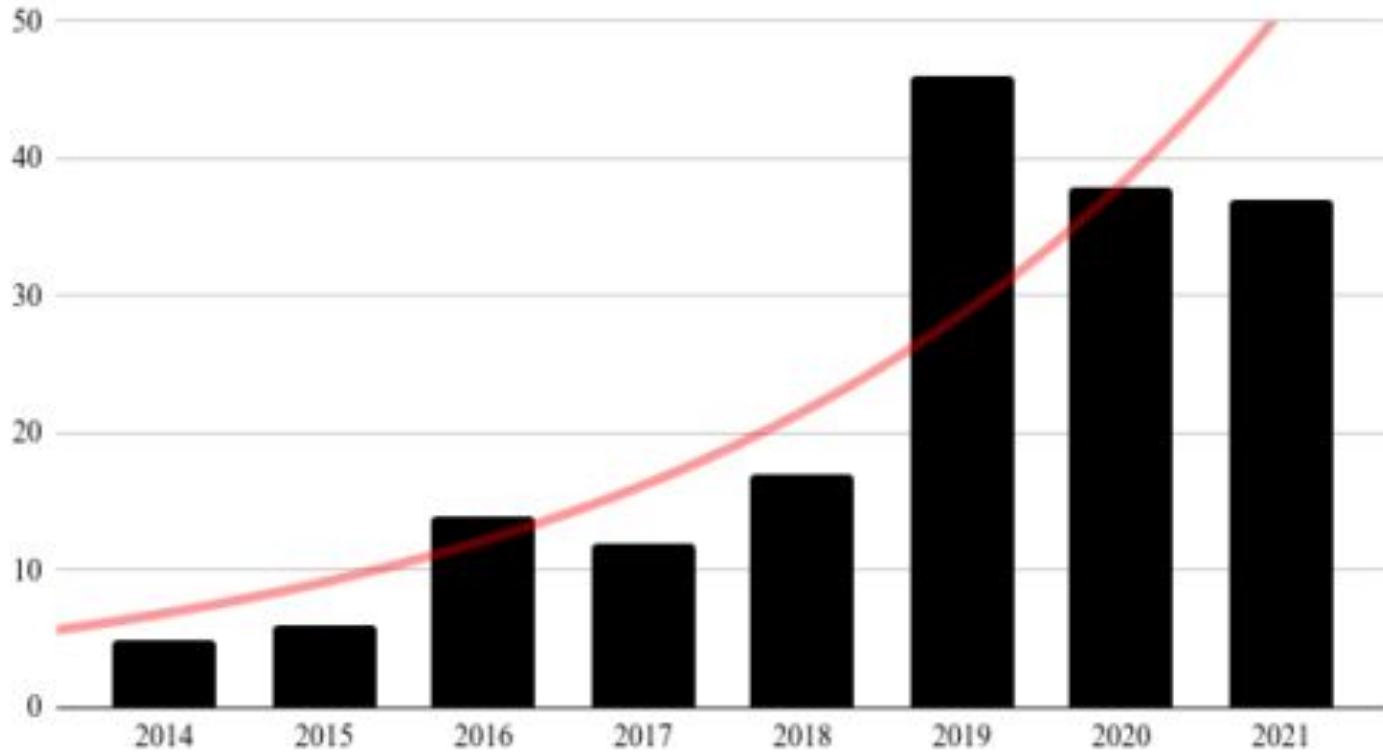# Have Business Logic Audits Been Suppressed?

**2016**: "Writing an article that's negative, you're effectively dissuading people from using autonomous vehicles, *you're killing people*."

**2018**: "It's really incredibly irresponsible of any journalists with integrity to write an article that would lead people to believe that autonomy is less safe because *people might actually turn it off, and then die*."
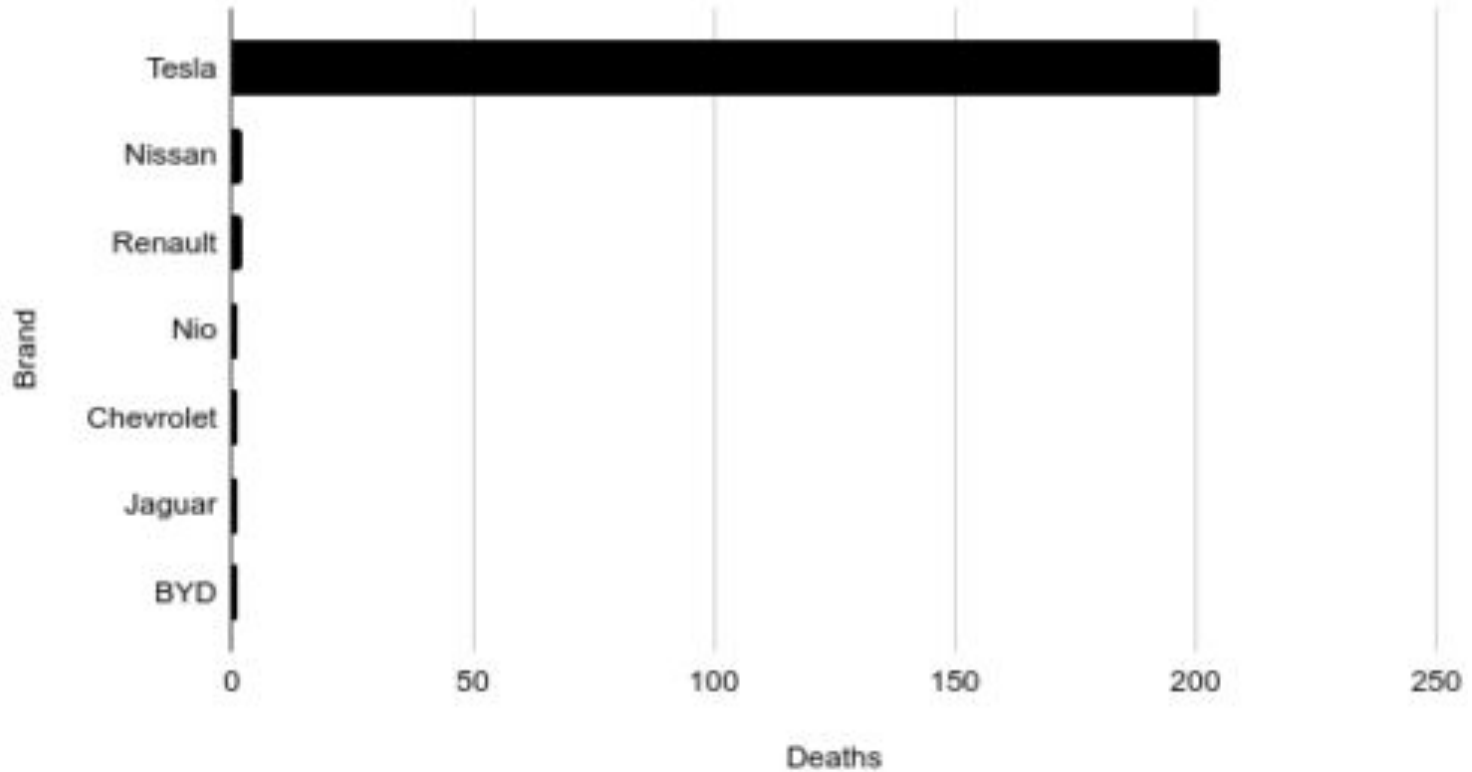
**2021**: "Elon Musk says Tesla's latest beta *self-driving software is 'not great'*..."

Tesla Deaths Per Year (as of July 2021)

https://www.flyingpenguin.com/?p=36192

Deaths by Electric Vehicle Brand (as of 8/30/2021)

# Tesla More Likely to Run Over Black People

- American Fear of Non-Motorized Planet
- Jaywalking is a Fantasy Crime
- Pedestrian Kill Bills Are Racist
- American Pedestrians Killed Disproportionately by Race
- Racism at Tesla Might Explain Why Their "Autopilot" Crashes So Often



RATE OF PEDESTRIAN FATALITIES INCREASES YEAR AFTER YEAR AFTER YEAR
STOP KILLING OUR KIDS

*Audit for **Stop** and **Reset** Functions Wherever You Find AI*

ISACA
San Francisco Chapter

inrupt

# AI Auditing

Davi Ottenheimer
October 24, 2021