# CS 4530
# Fundamentals of Software Engineering

## Module 18: Engineering Software for Equity

Adeel Bhutta, Jan Vitek, and Mitch Wand
Khoury College of Computer Sciences

# Learning Objectives for this Lesson

**By the end of this lesson, you should be able to…**

- Suggest some ways in which software can cause inadvertent harm or amplify inequities, with examples

- Explain why the software engineer has a powerful role to play in avoiding such harms.

# From SE @ Google:

As new as the field of software engineering is, we're newer still at understanding the impact it has on underrepresented people and diverse societies. … [We must recognize] the increasing imbalance of power between those who make development decisions that impact the world and those who simply must accept and live with those decisions that sometimes disadvantage already marginalized communities globally.

# A good software engineer will recognize potentials for inequity from their software.



"One mark of an exceptional engineer is the ability to understand how products can advantage and disadvantage different groups of human beings. Engineers are expected to have technical aptitude, but they should also have the discernment to know when to build something and when not to."

-Demma Rodriguez,
Head of Equity Engineering @ Google

Quote: "Software Engineering at Google: Lessons Learned from Programming Over Time," Wright, Winters and Manshreck, 2020 (O'Reilly)

# A good software engineer will recognize potentials for harm from their software.



- One mark of an exceptional engineer is the ability to understand how products can be weaponized to create harms in certain groups.

- Microsoft failed to do this with their chatbot of Tay that learned and picked up the behavior people used.

- People taught Tay to use offensive and racist language attacking jews.

# A good software engineer will recognize potentials for harm from their software.



- One mark of an exceptional engineer is the ability to understand how products can create harms in certain groups.

- Amazon failed to do this with their AI hiring software that used 10 years worth of resumes that had been submitted to Amazon to learn what candidates should be hired.

- Amazon taught its system to automatically reject the resumes of women.

# Algorithmic sentencing systems can discriminate against Black defendants

**Example: the COMPAS Sentencing Tool**

| | ALL DEFENDANTS | WHITE DEFENDANTS | BLACK DEFENDANTS |
|---|---|---|---|
| Labeled Higher Risk, But Didn't Re-Offend | 32.4% | 23.5% | 44.9% |
| Labeled Lower Risk, Yet Did Re-Offend | 37.4% | 47.7% | 28.0% |

# Algorithmic bias can discriminate against poorer consumers

# Training AI systems can have serious impacts on climate.

The Register®

{* AI + ML *}

## AI me to the Moon... Carbon footprint for 'training GPT-3' same as driving [to a] satellite and back

Get ready for Energy Star stickers on your [...]

Katyanna Quach    Wed 4 Nov 2020 // 07:59 UTC

Training OpenAI's giant GPT-3 text-generating m[...] car to the Moon and back, computer scientists re[...]

More specifically, they estimated teaching the ne[...] Microsoft data center using Nvidia GPUs require[...] which using the average carbon intensity of America would have produced 85,000 kg of $CO_2$ equivalents, the same amount produced by a new car in Europe driving 700,000 km, or 435,000 miles, which is about twice the distance between Earth and the Moon, some 480,000 miles. Phew.

| Consumption | CO$_2$e (lbs) |
|---|---|
| Air travel, 1 passenger, NY↔SF | 1984 |
| [Human li]fe, 1 year | 11,023 |
| [American, a]vg, 1 year | 36,156 |
| [car w/ fu]el, 1 lifetime | 126,000 |
| **[Training one m]odel (GPU)** | |
| [NLP pipeline (p]arsing, SRL) | 39 |
| w/ [tuning &] experimentation | 78,468 |
| Transformer (big) | 192 |
| w/ neural architecture search | 626,155 |

**Not to mention bitcoin mining!**

"Energy and Policy Considerations for Deep Learning in NLP" Emma Strubell, Ananya Ganesh, Andrew McCallum, in Proceedings of ACL 2019

https://www.theregister.com/2020/11/04/gpt3_carbon_footprint_estimate/

# Poor user interfaces can discriminate against differently -abled people.

## Inclusivity and Accessibility: Domino's Pizza LLC v. Robles

Domino's Would Rather Go to the Supreme Court Than Make Its Website Accessible to the Blind

Rather than developing technology to support users with disabilities, the pizza chain is taking its fight to the top

by Brenna Houck | @EaterDetroit | Jul 25, 2019, 6:00pm EDT

f  🐦  ⤴ SHARE

| | |
|---|---|
| Jul 15 2019 | **Brief amicus curiae of Washington Legal Foundation filed.** |
| Jul 15 2019 | **Brief amici curiae of Retail Litigation Center, Inc., et al. filed.** |
| Jul 15 2019 | **Brief amicus curiae of Cato Institute filed.** |
| Jul 15 2019 | **Brief amicus curiae of Restaurant Law Center filed.** |
| Jul 15 2019 | **Brief amici curiae of Chamber of Commerce of the United States of America, et al. filed.** |

*"Domino's Would Rather Go to the Supreme Court Than Make Its Website Accessible to the Blind" by Brenna Houck, Eater Detroit*
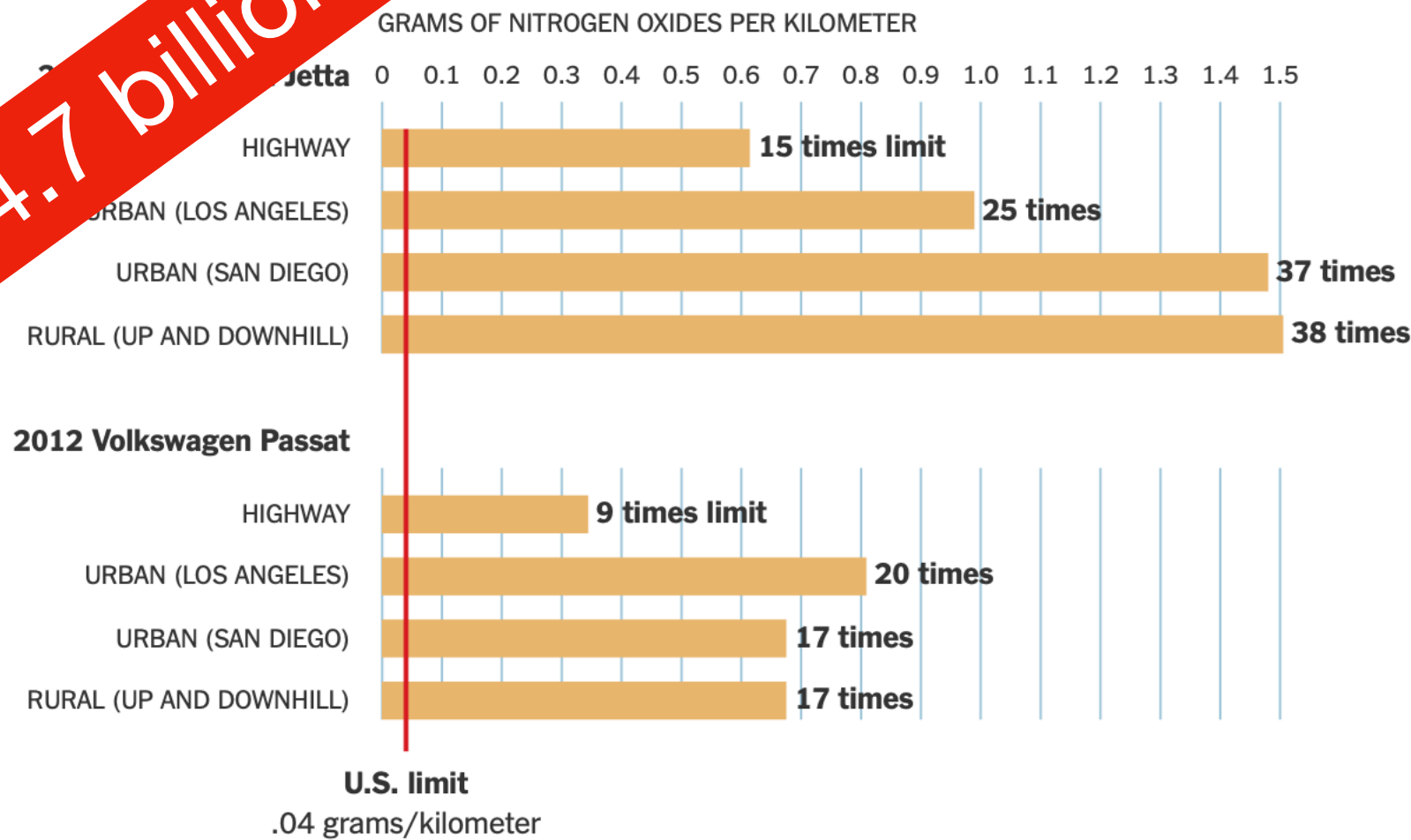
# Software Systems can be used to evade regulation.

## Example: Volkswagen diesel emissions



**The Emissions Tests That Led to the Discovery of VW's Cheating**

The on-road testing in May 2014 that led the California Air Resources Board to investigate Volkswagen was conducted [by] researchers at West Virginia University. They tested emissions from two VW [cars equi]pped with the 2-liter turbocharged 4-cylinder diesel engine. The [researchers fou]nd that when tested on the road, some cars emitted almost **40 tim[es the permit]ted levels of nitrogen oxides.

**Average emiss[ions of nitro]gen oxides in on-road testing**

GRAMS OF NITROGEN OXIDES PER KILOMETER

0   0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1.0  1.1  1.2  1.3  1.4  1.5

**$14.7 billion settlement**

2012 Volkswagen Jetta

| HIGHWAY | 15 times limit |
| URBAN (LOS ANGELES) | 25 times |
| URBAN (SAN DIEGO) | 37 times |
| RURAL (UP AND DOWNHILL) | 38 times |

**2012 Volkswagen Passat**

| HIGHWAY | 9 times limit |
| URBAN (LOS ANGELES) | 20 times |
| URBAN (SAN DIEGO) | 17 times |
| RURAL (UP AND DOWNHILL) | 17 times |

**U.S. limit**
.04 grams/kilometer

Source: Arvind Thiruvengadam, Center for Alternative Fuels, Engines and Emissions at West Virginia University

**Main computer**
Engine control module

Diesel oxidation catalytic converter

Oxygen sensor

Muffler

Oxygen sensor

H2S catalytic converter

Diesel particulate filter

Temperature sensors

Exhaust valve

**Nitrogen oxide trap**

This system traps nitrogen oxides, reducing toxic emissions. But the engine must regularly use more fuel to allow the trap to work. The car's **computer** could save fuel by allowing more pollutants to pass through the exhaust system. Saving fuel is one potential reason that Volkswagen's software could have been altered to make cars pollute more, according to researchers at the International Council on Clean Transportation.
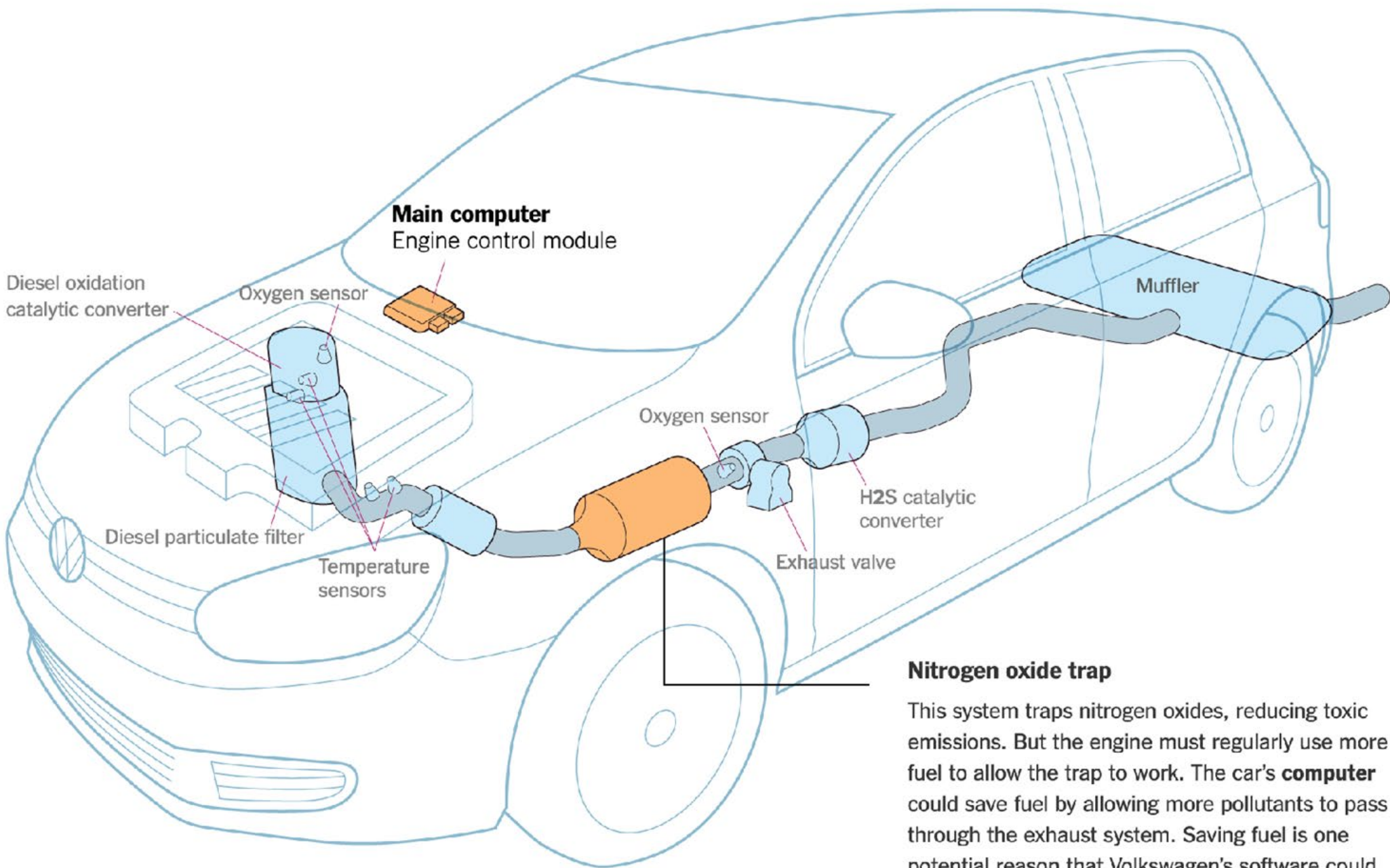
Illustration by Guilbert Gates | Source: Volkswagen, The International Council on Clean Transportation

"How Volkswagen's 'Defeat Devices' Worked" By Guilbert Gates, Jack Ewing, Karl Russell and Derek Watkins

# Bias is the Default

## Example: Google Photos auto   -tagging (2015)

https://www.wsj.com/articles/BL  -DGB-42522

https://www.wired.com/story/when   -it-comes-to-gorillas-google-photos-remains-blind/
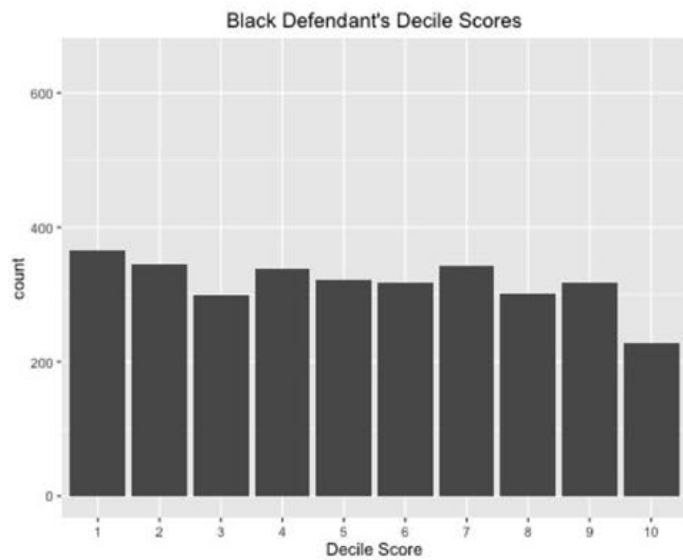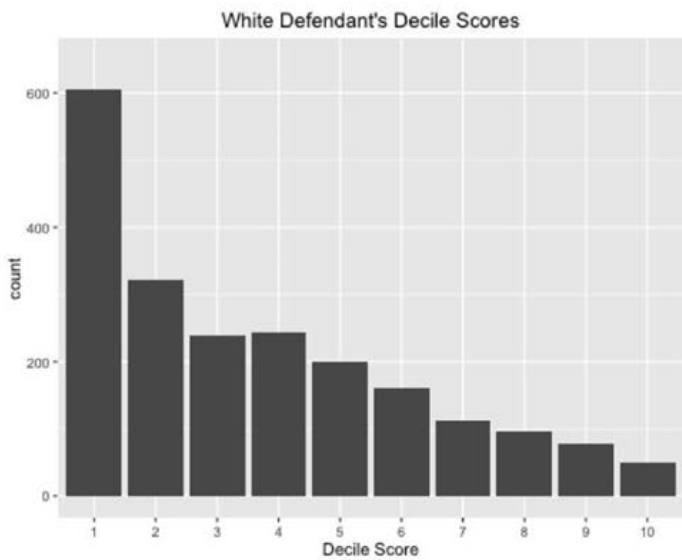
# Reflecting on these examples

## Personal philosophies and business cases



### Algorithmic Bias: COMPAS Sentencing Tool

| | ALL DEFENDANTS | WHITE DEFENDANTS | BLACK DEFENDANTS |
|---|---|---|---|
| Labeled Higher Risk, But Didn't Re-Offend | 32.4% | 23.5% | 44.9% |
| Labeled Lower Risk, Yet Did Re-Offend | 37.4% | 47.7% | 28.0% |

Analysis of Broward County, FL data: "How We Analyzed the COMPAS Recidivism Algorithm" by Jeff Larson, Surya Mattu, Lauren Kirchner and Julia Angwin

### Algorithmic Bias: Price Discrimination

THE WALL STREET JOURNAL

Websites Vary Prices, Deals Based on Users' Information

By Jennifer Valentino-DeVries, Jeremy Singer-Vine and Ashkan Soltani
December 24, 2012
https://www.wsj.com/articles/SB10001424127887323777204578189391813881534

### Inclusivity and Accessibility: Domino's Pizza LLC v. Robles

Domino's Would Rather Go to the Supreme Court Than Make Its Website Accessible to the Blind

Rather than developing technology to support users with disabilities, the pizza chain is taking its fight to the top

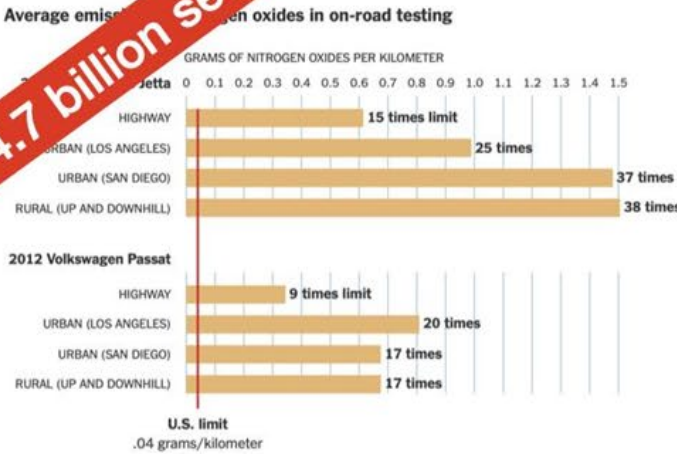by Brenna Houck | @EaterDetroit | Jul 25, 2019, 6:00pm EDT

| Jul 15 2019 | Brief amicus curiae of Washington Legal Foundation filed. |
| Jul 15 2019 | Brief amici curiae of Retail Litigation Center, Inc., et al. filed. |
| Jul 15 2019 | Brief amicus curiae of Cato Institute filed. |
| Jul 15 2019 | Brief amicus curiae of Restaurant Law Center filed. |
| Jul 15 2019 | Brief amici curiae of Chamber of Commerce of the United States of America, et al. filed. |

"Domino's Would Rather Go to the Supreme Court Than Make Its Website Accessible to the Blind" by Brenna Houck, Eater Detroit

### Evading regulation: Volkswagen

The Emissions Tests That Led to the Discovery of VW's Cheating

"How Volkswagen's 'Defeat Devices' Worked" By Guilbert Gates, Jack Ewing, Karl Russell and Derek Watkins

# More than "don't be evil"

**Engineering equitable software requires conscious effort**

- How do we determine what "the right thing" is?

- How do we convince our investors/managers to take this action?

# How might we mitigate harms in Software?

**Everything can and should be iterated on, including the problem itself … what are you trying to solve?**

- For every piece of software you create, you should iterate on it and include a wide range of people to use your software.

- By including more people you can better detect biases and harm that your software might create on certain populations.
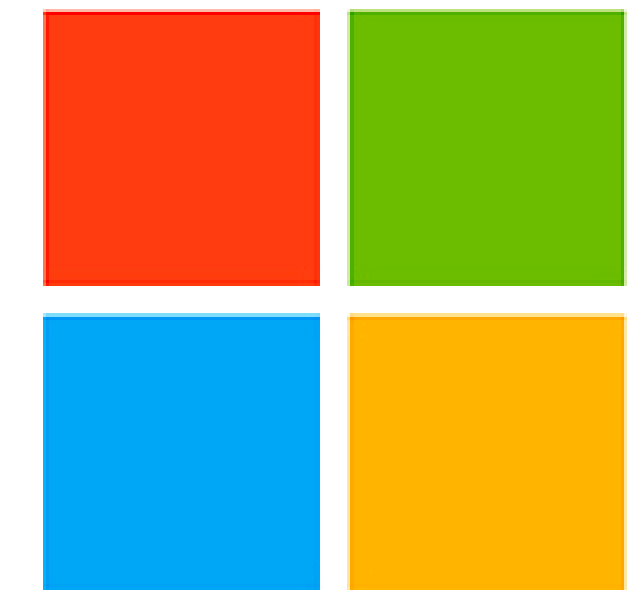
- You want to iterate your software throughout its entire life cycle.

# Guidelines from Microsoft on how to create software for people that mitigates harm.

**7**

WHEN WRONG

### Support efficient invocation.

Make it easy to invoke or request the AI system's services when needed.

**8**

WHEN WRONG

### Support efficient dismissal.

Make it easy to dismiss or ignore undesired system services.

**9**

WHEN WRONG

### Support efficient correction.

Make it easy to edit, refine, or recover when the AI system is wrong.

**10**

WHEN WRONG

### Scope services when in doubt.

Engage in disambiguation or gracefully degrade the AI system's services when uncertain about a user's goals.

**11**

WHEN WRONG

### Make clear why the system did what it did.

Enable the user to access an explanation of why the AI system behaved as it did.

⚠ **WHEN WRONG**

Microsoft

**12 OVER TIME**

**Remember recent interactions.**

Maintain short-term memory and allow the user to make efficient references to that memory.

**13 OVER TIME**

**Learn from user behavior.**

Personalize the user's experience by learning from their actions over time.

**14 OVER TIME**

**Update and adapt cautiously.**

Limit disruptive changes when updating and adapting the AI system's behaviors.

**15 OVER TIME**

**Encourage granular feedback.**

Enable the user to provide feedback indicating their preferences during regular interaction with the AI system.

**16 OVER TIME**

**Convey the consequences of user actions.**

Immediately update or convey how user actions will impact future behaviors of the AI system.

**17 OVER TIME**

**Provide global controls.**

Allow the user to globally customize what the AI system monitors and how it behaves.

**18 OVER TIME**

**Notify users about changes.**

Inform the user when the AI system adds or updates its capabilities.

🕐 **OVER TIME**

Microsoft

# Class Exercise



- For Amazon's Hiring Software, define with a partner how you would re-design the system using these guidelines to mitigate harm from the software.

# This lesson was about the harms that software can inflict

**You should now be able to…**

- Suggest some ways in which software can cause inadvertent harm or amplify inequities, with examples

- Explain why the software engineer has a powerful role to play in avoiding such harms.