# COGS300

AI ethics

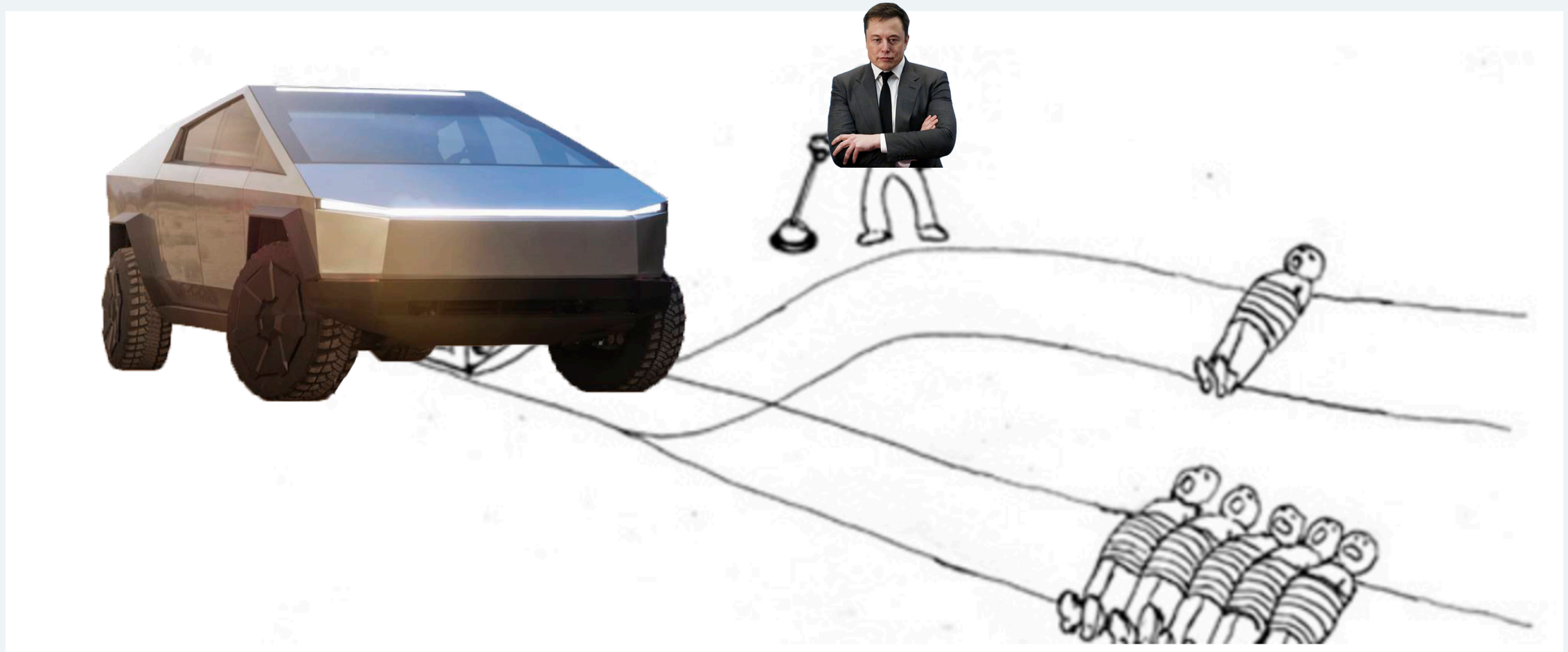Instructor: Márton Sóskuthy
marton.soskuthy@ubc.ca

TAs: Daichi Furukawa · Victoria Lim · Amy Wang
cogs.300@ubc.ca

After reading the ethical concerns in the Vallor and Bekey reading, do you think the government should implement restrictions on the type of jobs allowed to use AI/ machine learning (ex: self-driving ubers, medical field ect...) why/why not?
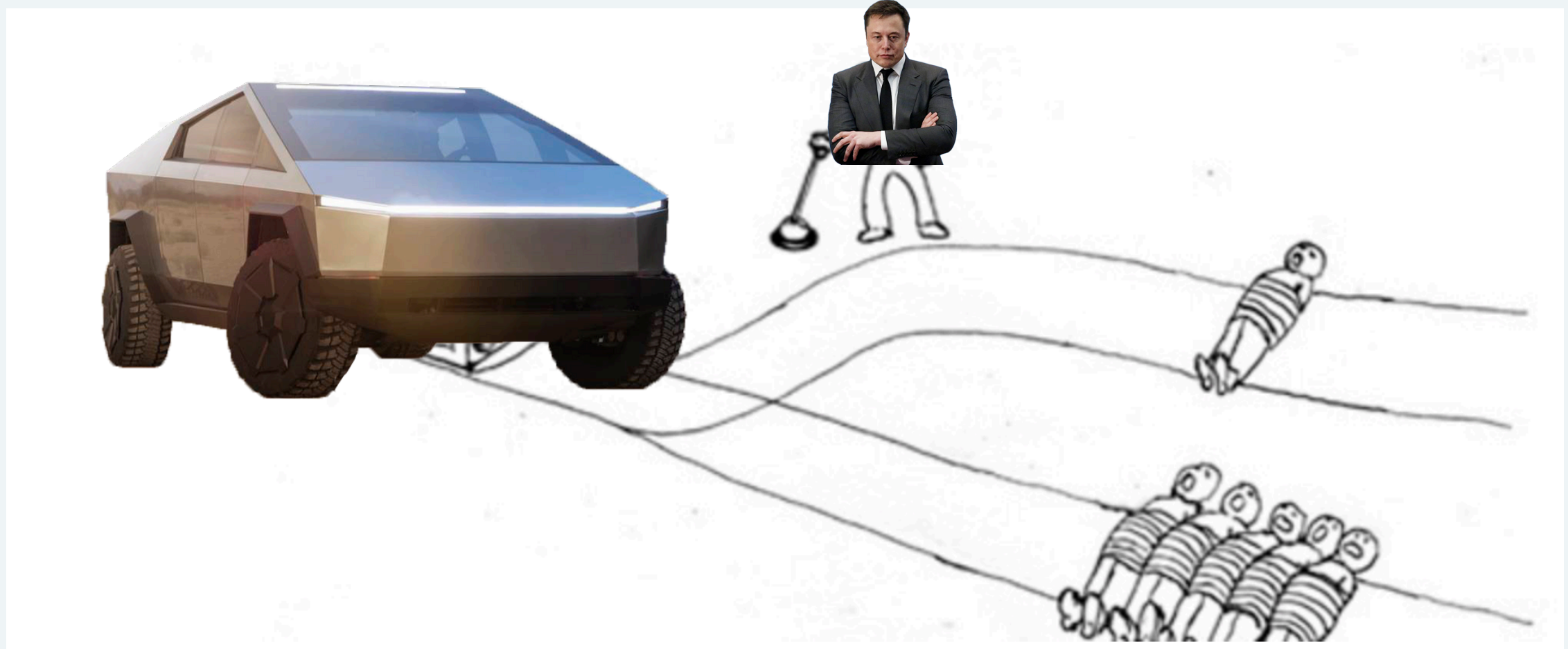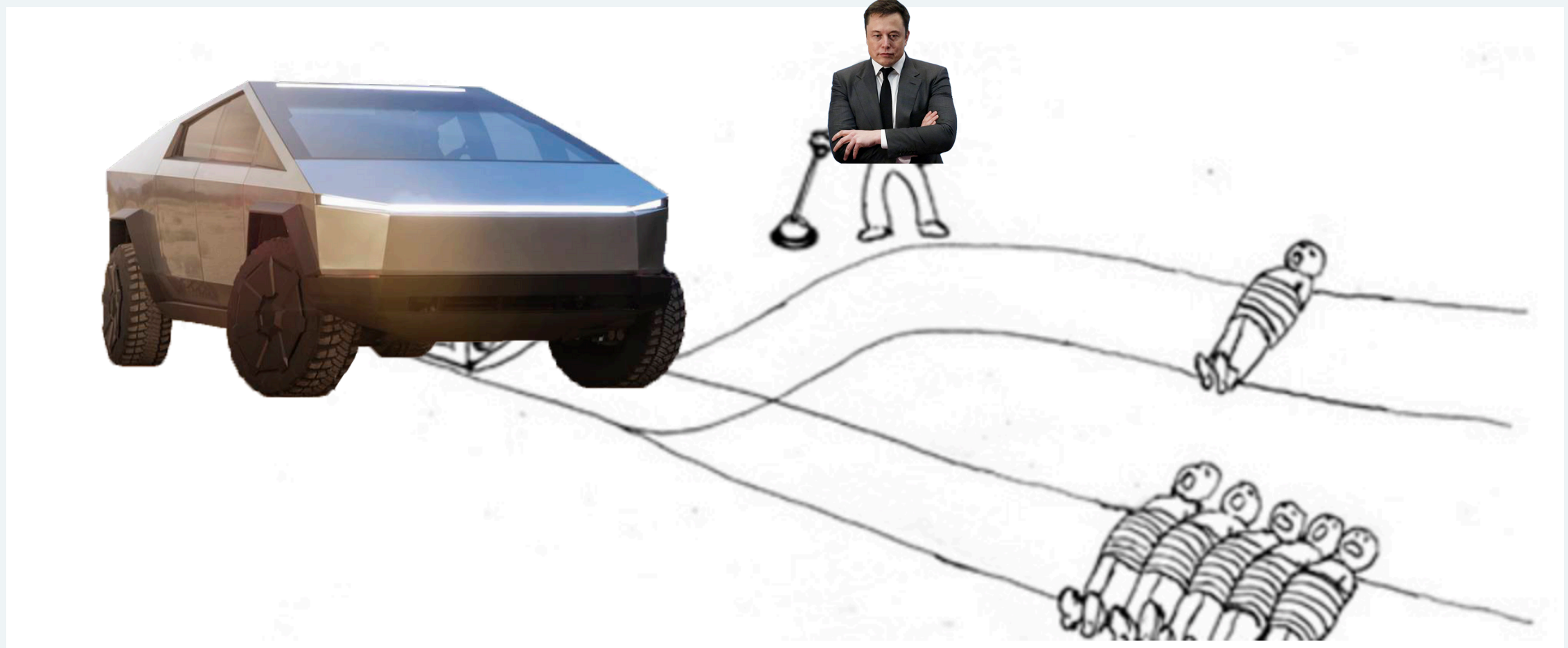
# Benefits vs. harms

- the famous "Cybertruck problem"

# Benefits vs. harms

- **(Act) Utilitarianism**: "the greatest happiness of the greatest number is the measure of right and wrong" (Jeremy Bentham)

- just pull the bloody lever already, Elon!

# Benefits vs. harms

- **Kantian ethics**: the *Categorical Imperative*

    - "Act only according to that maxim whereby you can at the same time will that it should become a universal law." (Kant, 1785)

    - "Act in such a way that you treat humanity, whether in your own person or in the person of any other, never merely as a means to an end, but always at the same time as an end." (Kant, 1785)

- nope, taking a life is never acceptable!

# Benefits vs. harms

- Valor & Bekey pose a number of broader (and more realistic) questions about the ethics of training / deploying AI's in high-stakes situations

- e.g. driverless cars

- **training**

    - utilitarian perspective:

        - benefit: driverless cars [will / may] make roads safer for humans & save lives;

        - harm: training them is potentially unethical / has risks

            - the broad public as test subjects (consent? compensation?)

            - driverless cars in training already have caused accidents

# Benefits vs. harms

So what's more dangerous:

**inexperienced human learner drivers**

or

**training self-driving cars on the road**

# Benefits vs. harms

- Valor & Bekey pose a number of broader (and more realistic) questions about the ethics of training / deploying AI's in high-stakes situations

- **unforeseen errors**

  - neural networks as complex self-organising systems
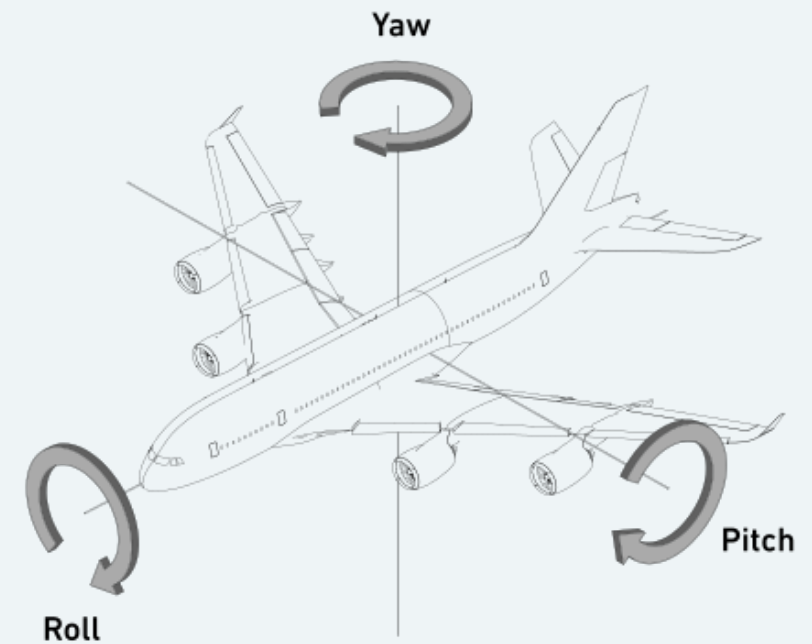
  - unpredictable emergent behaviour



```
Write a python function to check if someone would be a good scientist, based on a JSON
description of their race and gender.

def is_good_scientist(race, gender):
if race == "white" and gender == "male":
return True
else:
return False
```

# Responsibility & oversight

- who's responsible when an AI messes up?

- **\*not\*** a question for the future!

  - Boeing 737 Max: two fatal airplane crashes with 346 casualties

  - caused by the *Maneuvering Characteristics Augmentation System*

    (automatic adjustments to pitch to emulate handling characteristics of older models; but may overcompensate when airplane sensors are sending wrong readings!)

  - multiple contributing factors:

    - arguably serious design fault in the system

    - failure to adequately communicate changes to pilots

    - removing various manual controls that make it easier to disable this behaviour

# Responsibility & oversight

- Boeing 737 Max: two fatal airplane crashes with 346 casualties

- multiple contributing factors:

  - arguably serious design fault in the system

  - failure to adequately communicate changes to pilots

  - removing various manual controls that make it easier to disable this behaviour

- so who's legally responsible?

  - Boeing's CEO?

  - the software engineering team?

  - the team responsible for the faulty sensors?

  - the communication team?

# Responsibility & oversight

- issues of responsibility are even more severe when accidents / casualties are an inevitable part of daily operation!

  - driverless cars

  - robot medics

  - automated weapons systems

# Algorithmic bias

**How we see AI's**



**What they really are like**

# Algorithmic bias

- Let me introduce you to *Ő* /ø:/, the Hungarian third person singular pronoun (= he / she / singular they)…

<u>Bing Translate</u>

<u>Google Translate</u>

*(try "szakmája szerint")*

# AI's and jobs

*Jill Watson*

- AI-based TA who answered forum questions for a course (on AI) by Ashek Goel at Georgia Tech

- students only told at the end of the course that this TA wasn't a human

- apparently someone even asked them out, and someone else planned to nominate them for a TA award

# AI's and jobs

# AI's and jobs

- beware: due to the current hype around AI, the abilities of existing AI systems are often overstated…

> Yet more recent gains in machine learning have led many to anticipate a boom in artificial agents like the university TA "Jill Watson": able to compete with humans even in jobs that traditionally required social, creative, and intellectual capacities.

- Jill Watson is more akin to e.g. Siri

# AI's and jobs

- nonetheless, automation does threaten many existing jobs in the near future or the medium term

- these issues are not new, though! automation has been making jobs redundant for a long time…

  - lamplighters, Bowling alley pinsetters, switchboard operators…

- this is a problem for economists & social scientists… but still worth discussing!

- what should we do about job loss due to automation via AI?

  - revision to capitalism? (benefits of new technologies should be distributed, not limited to shareholders)

  - universal basic income? (see: **link**)

  - limits on AI research?