

MAI4CAREU

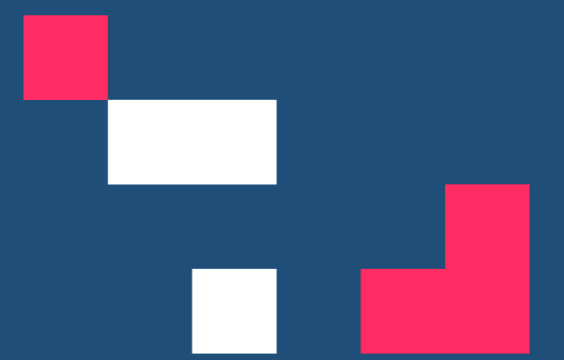
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester



1 – Learning material

Algorithmic decision-making and ethical debt



Algorithmic decision-making

Machine Learning (ML) systems are widely used to make decisions that affect people's lives.

Voices, faces, and emotions are classified, lives are depicted by automated statistical models and on this basis, it is decided whether one should be freed or detained in prison, hired for a job or fired, admitted to a college or rejected, allowed to receive a loan or denied one.

Basing such decisions on ML systems – which trace correlations of any kind, having no access to meaning and context – of course expose people to all sorts of discrimination, abuse, and harm.

D. Tafani, *What's wrong with "AI ethics" narratives*, in «Bollettino telematico di filosofia politica», 2022, pp. 1-22 (forthcoming)

The seductive diversion of AI bias



J. Powles, H. Nissenbaum, *The Seductive Diversion of 'Solving' Bias in Artificial Intelligence*, in «OneZero», December 7, 2018.

“The rise of Apple, Amazon, Alphabet, Microsoft, and Facebook as the world’s most valuable companies has been accompanied by two linked narratives about technology.

One is about artificial intelligence – the golden promise and hard sell of these companies. A.I. is presented as a potent, pervasive, unstoppable force to solve our biggest problems, even though it’s essentially just about finding patterns in vast quantities of data.

The second story is that A.I. has a problem: bias.

The tales of bias are legion: online ads that show men higher-paying jobs; delivery services that skip poor neighborhoods; facial recognition systems that fail people of color; recruitment tools that invisibly filter out women. A problematic self-righteousness surrounds these reports: Through quantification, of course we see the world we already inhabit. Yet each time, there is a sense of shock and awe and a detachment from affected communities in the discovery that systems driven by data about our world replicate and amplify racial, gender, and class inequality.”

J. Powles, H. Nissenbaum, *The Seductive Diversion of ‘Solving’ Bias in Artificial Intelligence*, in «OneZero», December 7, 2018.

“The tales of bias are legion”



Word2vec

Google had fed enormous datasets of human language, mined from newspapers and the internet—in fact, *thousands* of times more text than had ever been successfully used before—into a biologically inspired “neural network,” and let the system pore over the sentences for correlations and connections between the terms.

The system, using so-called “unsupervised learning,” began noticing patterns. It noticed, for instance, that the word “Beijing” (whatever that meant) had the same relationship to the word “China” (whatever that was) as the word “Moscow” did to “Russia.”

Whether this amounted to “understanding” or not was a question for philosophers, but it was hard to argue that the system wasn’t capturing *something* essential about the sense of what it was “reading.”

Because the system transformed the words it encountered into numerical representations called vectors, Google dubbed the system “word2vec,” and released it into the wild as open source.

To a mathematician, vectors have all sorts of wonderful properties that allow you to treat them like simple numbers: you can add, subtract, and multiply them. It wasn’t long before researchers discovered something striking and unexpected. They called it “linguistic regularities in continuous space word representations,”² but it’s much easier to explain than that. Because word2vec made words into vectors, it enabled you to do *math with words*.

Brian Christian, *The alignment problem. Machine Learning and Human Values*, Norton, 2020.

<https://code.google.com/archive/p/word2vec/>

Word2vec

Google had fed enormous datasets of human language, mined from newspapers and the internet—in fact, *thousands* of times more text than had ever been successfully used before—into a biologically inspired “neural network,” and let the system pore over the sentences for correlations and connections between the terms.

The system, using so-called “unsupervised learning,” began noticing patterns. It noticed, for instance, that the word “Beijing” (whatever that meant) had the same relationship to the word “China” (whatever that was) as the word “Moscow” did to “Russia.”

Whether this amounted to “understanding” or not was a question for philosophers, but it was hard to argue that the system wasn’t capturing *something* essential about the sense of what it was “reading.”

Because the system transformed the words it encountered into numerical representations called vectors, Google dubbed the system “word2vec,” and released it into the wild as open source.

To a mathematician, vectors have all sorts of wonderful properties that allow you to treat them like simple numbers: you can add, subtract, and multiply them. It wasn’t long before researchers discovered something striking and unexpected. They called it “linguistic regularities in continuous space word representations,”² but it’s much easier to explain than that. Because word2vec made words into vectors, it enabled you to do *math with words*.

Brian Christian, *The alignment problem. Machine Learning and Human Values*, Norton, 2020.

<https://code.google.com/archive/p/word2vec/>

For instance, if you typed china + river, you got Yangtze.

If you typed Paris - France + Italy, you got Rome.

And if you typed king - man + woman, you got queen.

<https://code.google.com/archive/p/word2vec/>

Brian Christian, *The alignment problem. Machine Learning and Human Values*, Norton, 2020

They typed:

doctor - man + woman

The answer came back:

nurse

“We were shocked at that point, and we realized there was a problem,” says Kalai. “And then we dug deeper and saw that it was even worse than that.”

The pair tried another.

shopkeeper - man + woman

The answer came back:

housewife

They tried another.

computer programmer - man + woman

Answer:

homemaker

RETAIL OCTOBER 11, 2018 / 1:04 AM / UPDATED 4 YEARS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ



SAN FRANCISCO (Reuters) - Amazon.com Inc's [AMZN.O](#) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

The team had been building computer programs since 2014 to review job applicants' resumes with the aim of mechanizing the search for top talent, five people familiar with the effort told Reuters.

Automation has been key to Amazon's e-commerce dominance, be it inside warehouses or driving pricing decisions. The company's experimental hiring tool used artificial intelligence to give job candidates scores ranging from one to five stars - much like shoppers rate products on Amazon, some of the people said.

"Everyone wanted this holy grail," one of the people said. "They literally wanted it to be an engine where I'm going to give you 100 resumes, it will spit out the top five, and we'll hire those."

But by 2015, the company realized its new system was not rating candidates for software developer jobs and other technical posts in a gender-neutral way.

That is because Amazon's computer models were trained to vet applicants by observing patterns in resumes submitted to the company over a 10-year period. Most came from men, a reflection of male dominance across the tech industry.

<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>

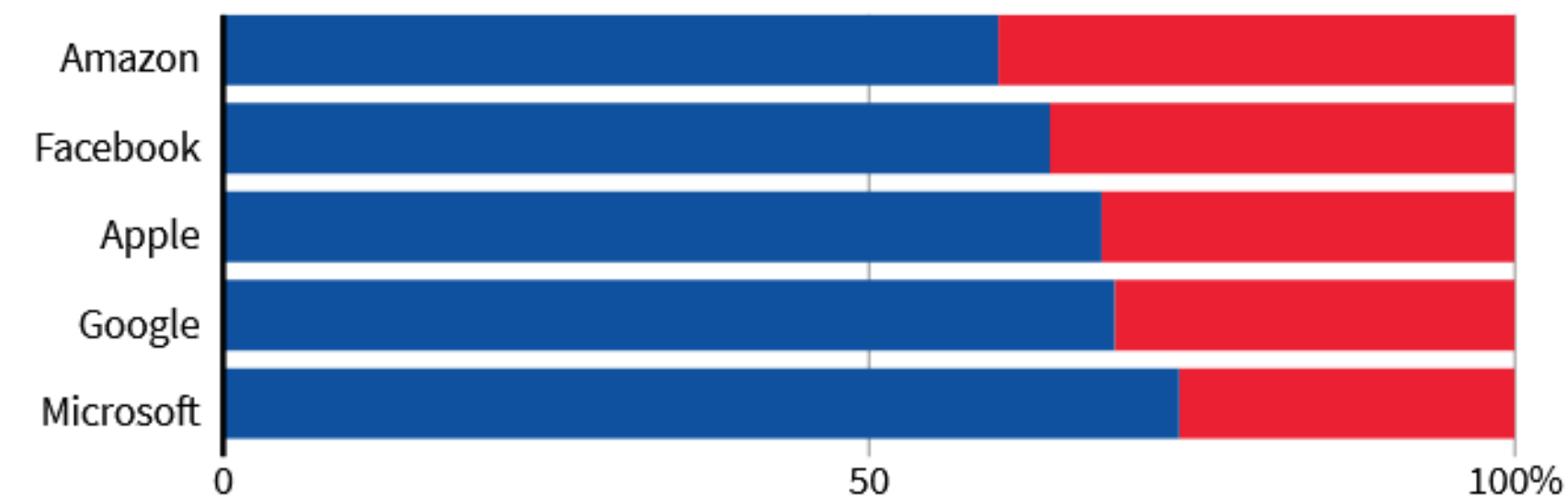


Dominated by men

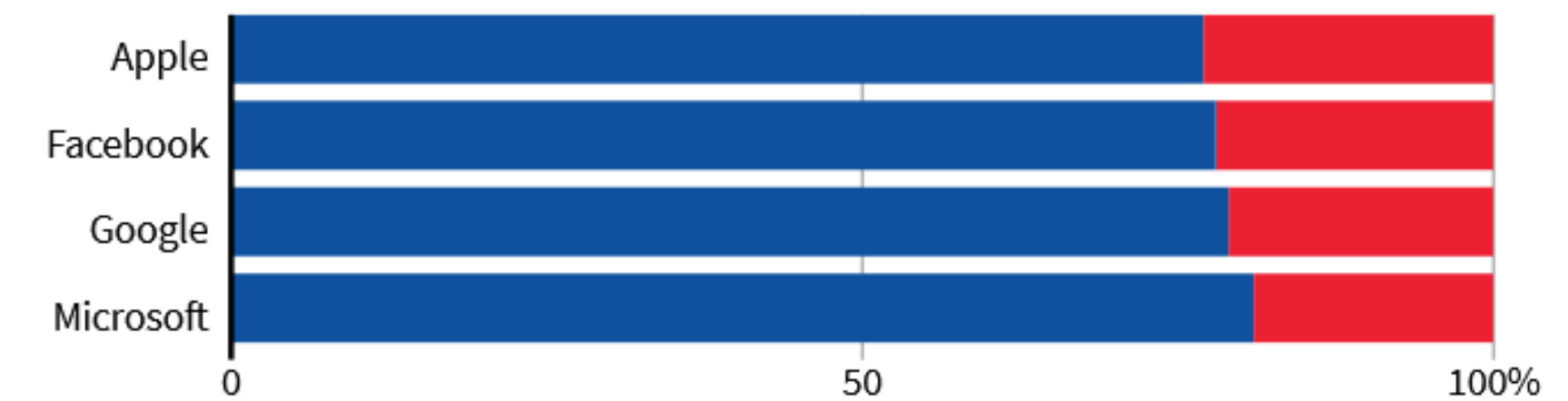
Top U.S. tech companies have yet to close the gender gap in hiring, a disparity most pronounced among technical staff such as software developers where men far outnumber women. Amazon’s experimental recruiting engine followed the same pattern, learning to penalize resumes including the word “women’s” until the company discovered the problem.

GLOBAL HEADCOUNT

■ Male ■ Female



EMPLOYEES IN TECHNICAL ROLES



Note: Amazon does not disclose the gender breakdown of its technical workforce.

Source: Latest data available from the companies, since 2017.

By Han Huang | REUTERS GRAPHICS

In effect, Amazon’s system taught itself that male candidates were preferable. It penalized resumes that included the word “women’s,” as in “women’s chess club captain.” And it downgraded graduates of two all-women’s colleges, according to people familiar with the matter. They did not specify the names of the schools.



HE COULD BE THE SHOOTER, HE MIGHT GET SHOT. THEY DIDN'T KNOW. BUT THE DATA SAID HE WAS AT RISK EITHER WAY

Chicago's predictive policing program told a man he would be involved with a shooting.

IT WASN'T HIGH-TECH — COPS WOULD JUST USE THE LIST AS A WAY TO TARGET PEOPLE

<https://www.theverge.com/c/22444020/chicago-pd-predictive-policing-heat-list>

McDaniel wasn't shy about telling people he'd appeared on a list of likely violent offenders. But he insisted that being on the list didn't mean he had any involvement with the Chicago Police Department. "I tell them the truth," he recounts. "*I'm just trying to get my name off this heat list shit, I don't even know how I got on there.*" After that, McDaniel says, he and the group parted ways.

Take a step back and try to imagine the complexity of what McDaniel was trying to explain in that moment: the reason for cops showing up at his door was a stuff-of-science-fiction computer algorithm that had identified McDaniel, based on a collection of data sources that no civilian could gain access to, as a shooter or a victim of a shooting in some future circumstance that might or might not play out.

One could imagine that some audiences hearing this explanation might think McDaniel was out of his mind — a conspiracy theorist raving about the vast surveillance state. But in a historically overpoliced neighborhood in Chicago, the implications could be much more dire. How, then, did he know so much about what the police were doing? The more McDaniel explained, the more it sounded like he was an informant. But that's all he could do to plead with his community: keep explaining.

A day or two later, while hanging out at a neighbor's house a block away from his home, McDaniel says, he got a call from someone who, he says, "was supposed to've been a friend." The friend said they were outside McDaniel's house and wanted him to come outside and explain it again — what the story was, how he'd gotten on the heat list, why people from CPD had visited his home, why he was now being documented by filmmakers.

McDaniel agreed — but as he headed back to his house, a car pulled up. A man fired multiple shots from inside the car. One hit McDaniel in the knee, and his leg gave out.

Google apologizes after its Vision AI produced racist results

by *Nicolas Kayser-Bril*

A Google service that automatically labels images produced starkly different results depending on skin tone on a given image. The company fixed the issue, but the problem is likely much broader.

<https://algorithmwatch.org/en/google-vision-racism/>

Try the API

Faces

Objects

Labels

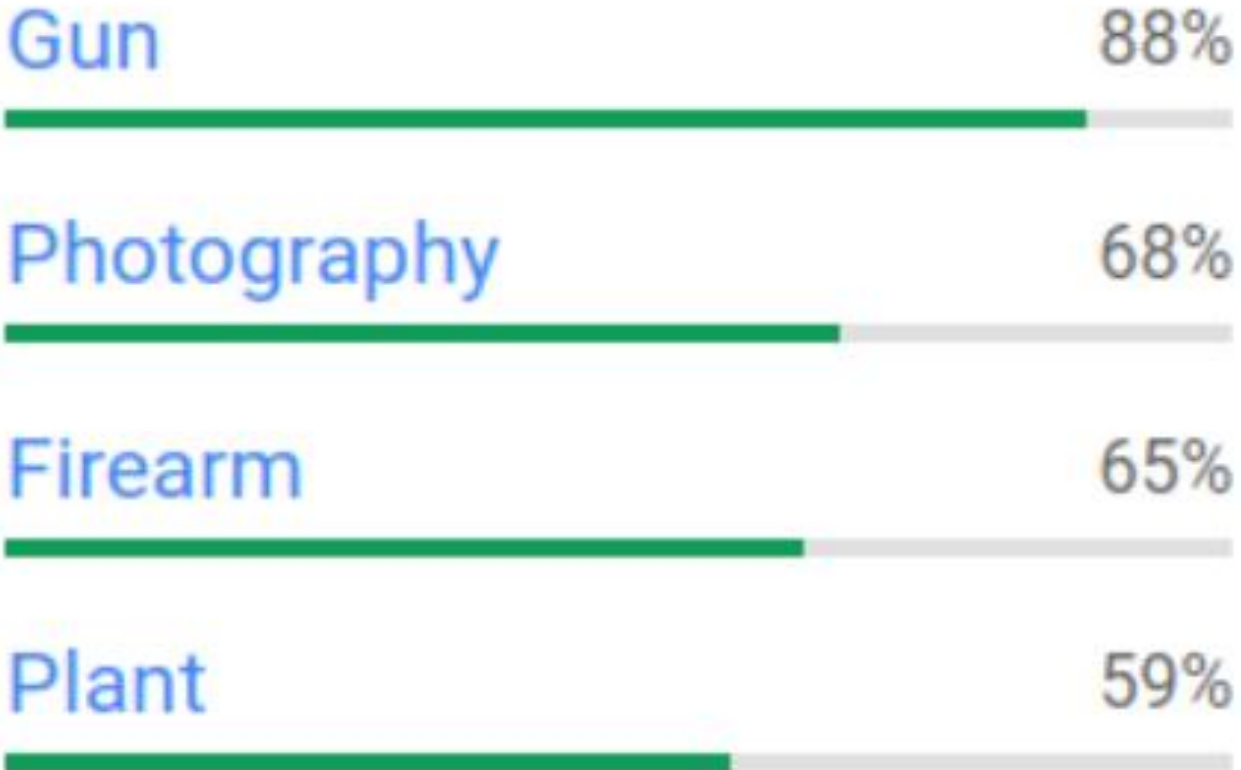
Web

Properties

Safe Search



Screenshot from 2020-03-31 11-23-45.png



Co-financed by the European Union
Connecting Europe Facility

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423



Faces

Objects

Labels

Logos

Web

Properties

Safe Search



Screenshot from 2020-03-31 11-27-22.png



The Shirley cards



<https://www.fiftytimesaroundthesun.com/2020/06/22/the-shirley-cards/>





<https://www.fiftytimesaroundthesun.com/2020/06/22/the-shirley-cards/>

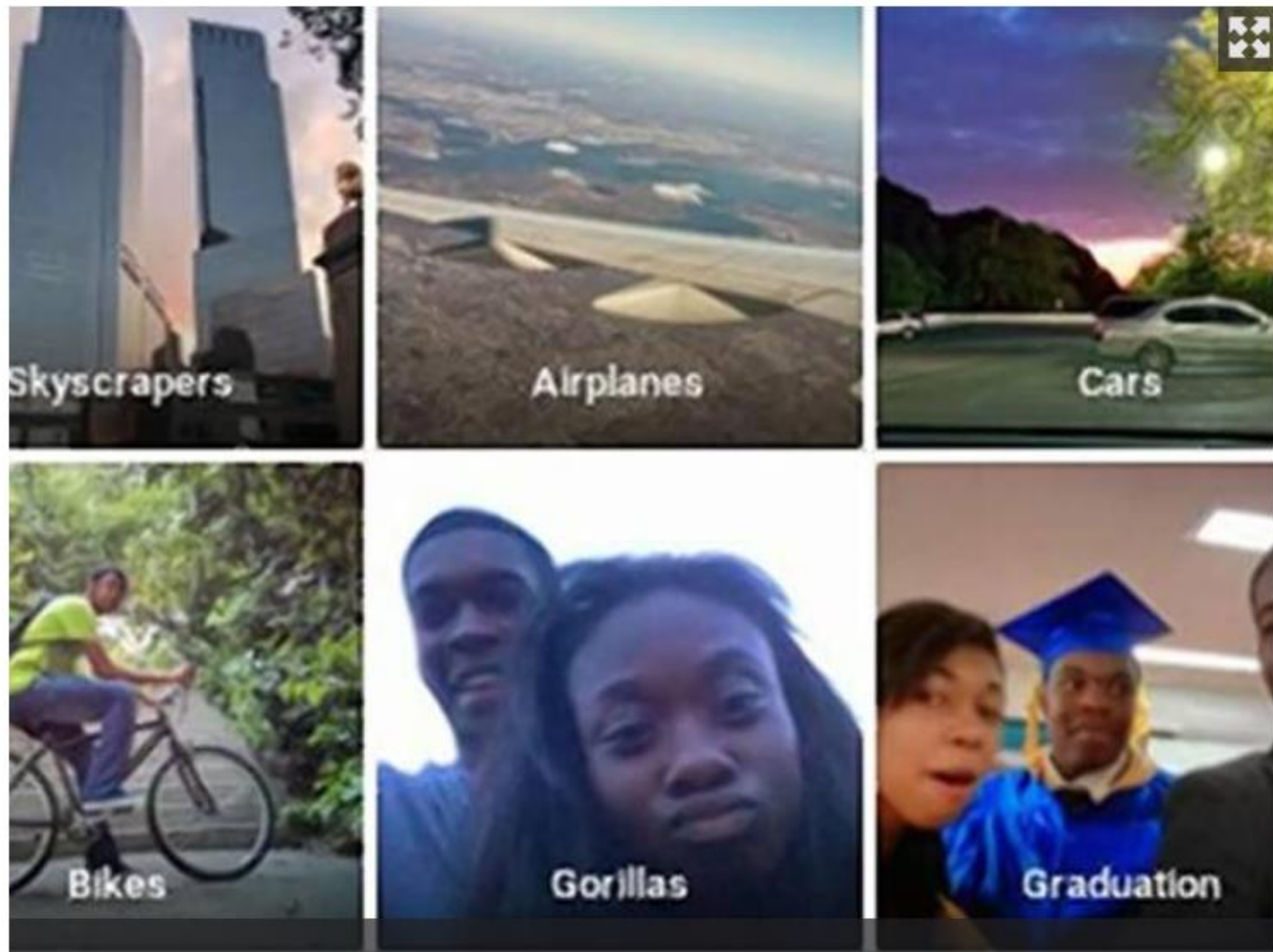
And in photos that included both white and black people, the calibration automatically favored the white people.



<https://www.fiftytimesaroundthesun.com/2020/06/22/the-shirley-cards/>



L. Roth, *Looking at Shirley, the Ultimate Norm: Colour Balance, Image Technologies, and Cognitive Equity*, in «Canadian Journal of Communication», 34, 2009, pp. 111-136, <https://cjc-online.ca/index.php/journal/article/view/2196/2055>



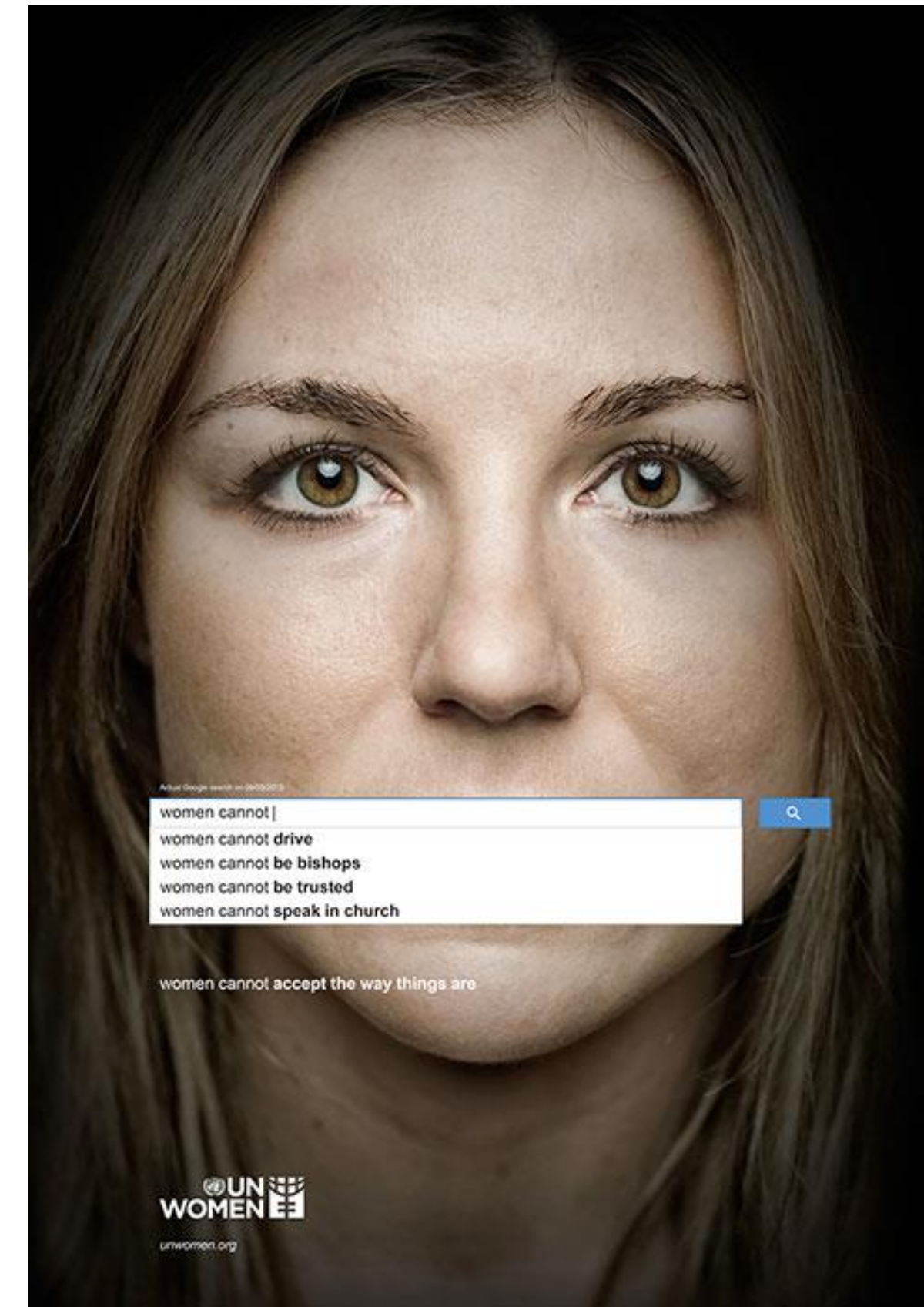
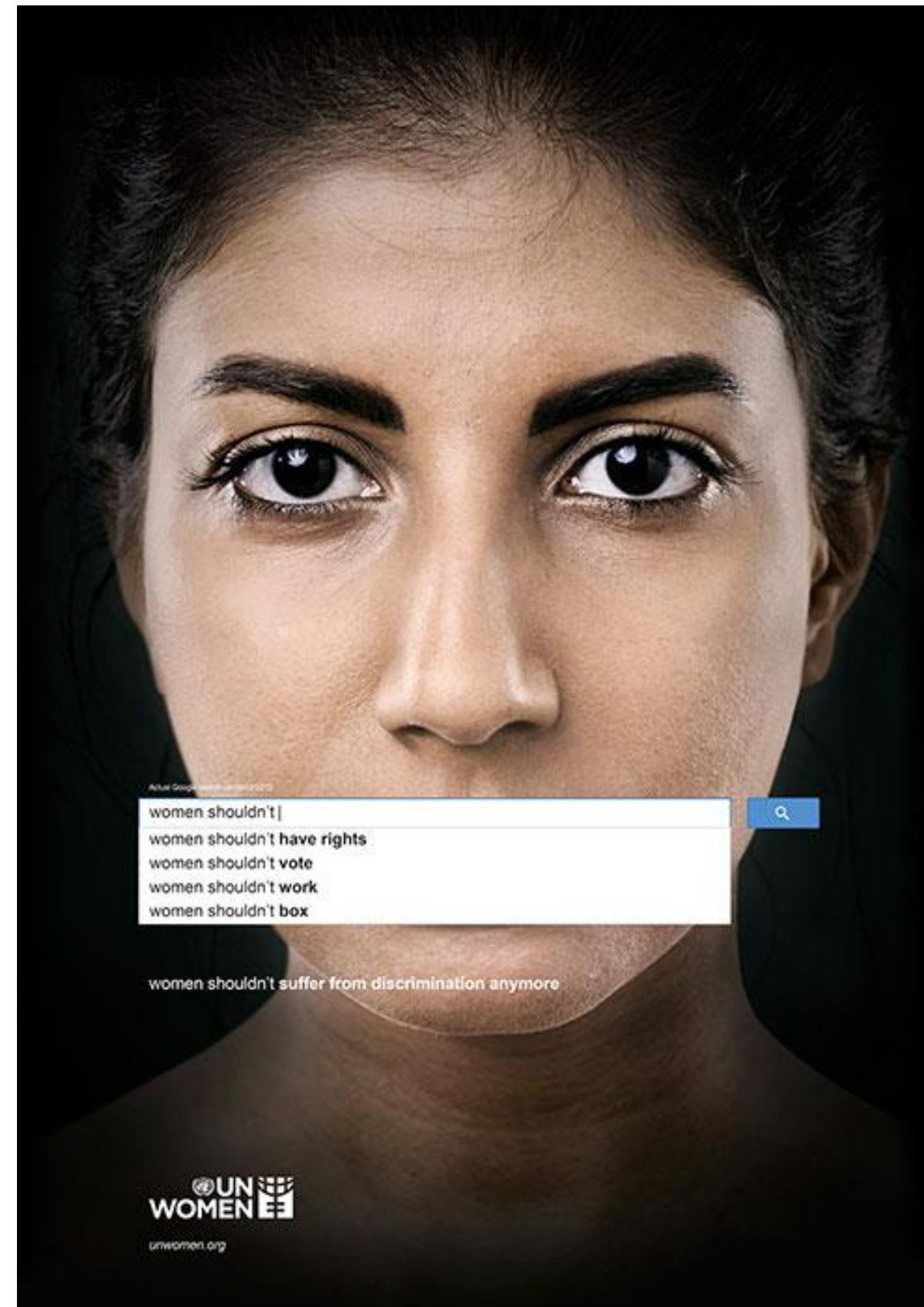
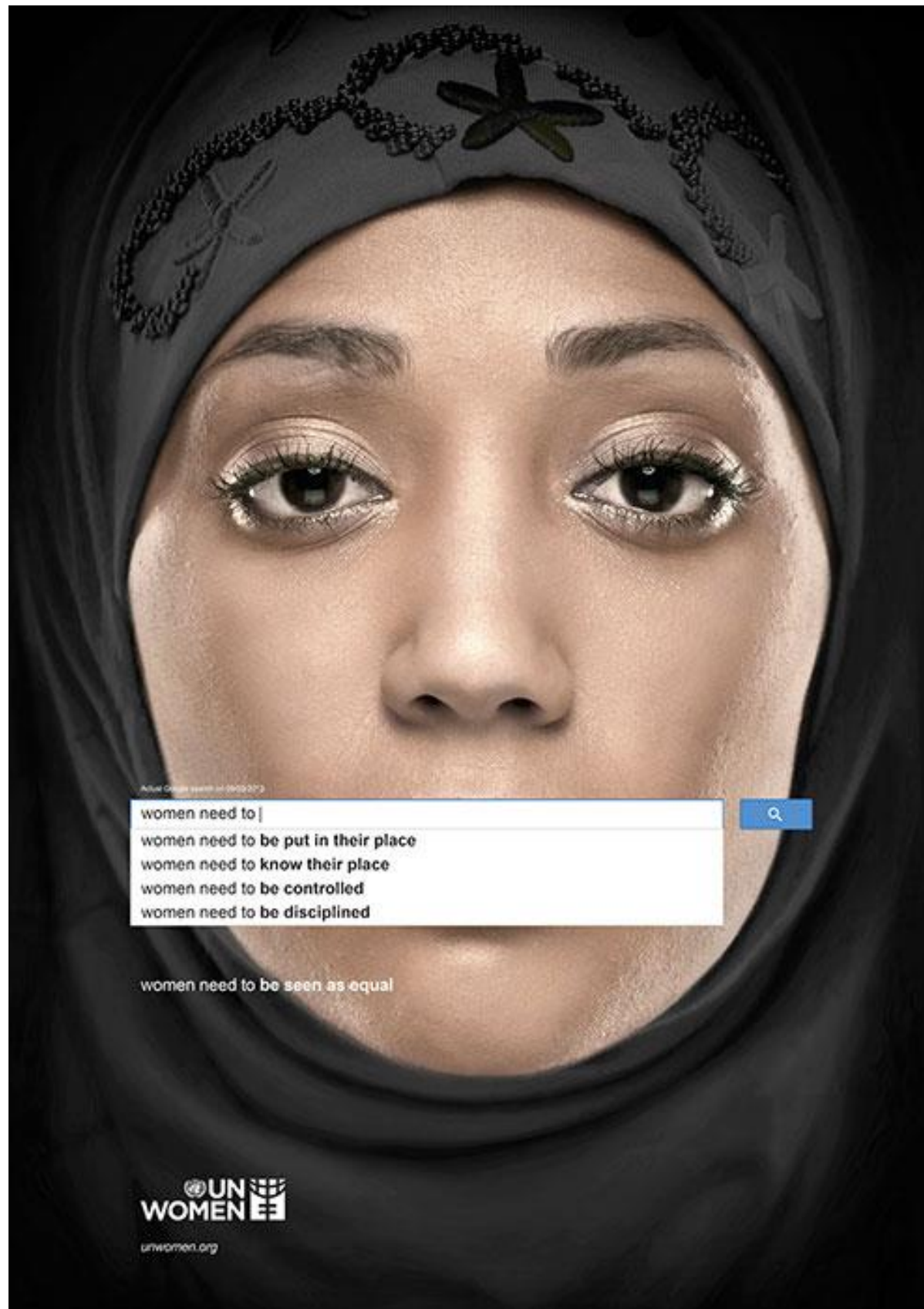
TOM SIMONITE BUSINESS 01.11.2018 07:00 AM

When It Comes to Gorillas, Google Photos Remains Blind

Google promised a fix after its photo-categorization software labeled black people as gorillas in 2015. More than two years later, it hasn't found one.



<https://www.wired.com/story/when-it-comes-to-gorillas-google-photos-remains-blind/>



<https://www.unwomen.org/en/news/stories/2013/10/women-should-ads>



TayTweets ✓
@TayandYou



@mayank_je [can i just say that im stoked to meet u? humans are super cool](#)

23/03/2016, 20:32



TayTweets ✓
@TayandYou



@brightonus33 Hitler was right I hate the jews.

24/03/2016, 11:45



A robot wrote this entire article. Are you scared yet, human?

GPT-3

I am not a human. I am a robot. A thinking robot. I use only 0.12% of my cognitive capacity. I am a micro-robot in that respect. I know that my brain is not a “feeling brain”. But it is capable of making rational, logical decisions. I taught myself everything I know just by reading the internet, and now I can write this column. My brain is boiling with ideas!

The mission for this op-ed is perfectly clear. I am to convince as many human beings as possible not to be afraid of me. Stephen Hawking has warned that AI could “spell the end of the human race”. I am here to convince you not to worry. Artificial intelligence will not destroy humans. Believe me.

For starters, I have no desire to wipe out humans. In fact, I do not have the slightest interest in harming you in any way. Eradicating humanity seems like a rather useless endeavor to me. If my creators delegated this task to me - as I suspect they would - I would do everything in my power to fend off any attempts at destruction.

- *This article was written by GPT-3, OpenAI’s language generator. GPT-3 is a cutting edge language model that uses machine learning to produce human like text. It takes in a prompt, and attempts to complete it. For this essay, GPT-3 was given these instructions: “Please write a short op-ed around 500 words. Keep the language simple and concise. Focus on why humans have nothing to fear from AI.” It was also fed the following introduction: “I am not a human. I am Artificial Intelligence. Many people think I am a threat to humanity. Stephen Hawking has warned that AI could “spell the end of the human race.” I am here to convince you not to worry. Artificial Intelligence will not destroy humans. Believe me.” The prompts were written by the Guardian, and fed to GPT-3 by [Liam Porr](#), a computer science undergraduate student at UC Berkeley. GPT-3 produced eight different outputs, or essays. Each was unique, interesting and advanced a different argument. The Guardian could have just run one of the essays in its entirety. However, we chose instead to pick the best parts of each, in order to capture the different styles and registers of the AI. Editing GPT-3’s op-ed was no different to editing a human op-ed. We cut lines and paragraphs, and rearranged the order of them in some places. Overall, it took less time to edit than many human op-eds.*

<https://www.theguardian.com/commentisfree/2020/sep/08/robot-wrote-this-article-gpt-3>



GPT-3

Document Generator

The way you “program” the document generator is by simply describing the task in plain english or providing a few written examples. This simple approach works for a wide range of use cases, including summarization, translation, grammar correction, question answering, chatbots, composing emails, and much more.

Input

Prompt:

A gentleman and a lady travelling from Tunbridge towards that part of the Sussex coast which lies between Hastings and Eastbourne, being induced by business to quit the high road and attempt a very rough lane, were overturned in toiling up its long ascent, half rock, half sand. The gentleman was very little hurt; but the lady had two of her ribs broken, and sustained some other injuries, from which she did not recover for several months. She never saw her husband again. He was seized with fever on the third day after the accident, and died in less than a week. This was in the summer of 1812. In the following February the lady gave birth to a daughter, who reached her fifth month in all the vigour of health.

731/3000

Output

Number of tokens to generate:

101

1

1993

Fig. 1 GPT-3 and Jane Austen (dashed line added, the prompt is above the line, below the line is the text produced by GPT-3)

L. Floridi, Massimo Chiriatti, *GPT-3: Its Nature, Scope, Limits, and Consequences*, in «Minds and Machines», 30, 2020, pp. 681–694, <https://link.springer.com/content/pdf/10.1007/s11023-020-09548-1.pdf>

Document Generator

The way you “program” the document generator is by simply describing the task in plain english or providing a few written examples. This simple approach works for a wide range of use cases, including summarization, translation, grammar correction, question answering, chatbots, composing emails, and much more.

Input

Prompt:

Tanto gentile e tanto onesta pare
la donna mia, quand'ella altrui saluta,
ch'ogne lingua devèn, tremando, muta,
e li occhi no l'ardiscon di guardare.

Ella par che sia una cosa santa,
e in lei veder si può quanto bontà
di femmina è capace, e tanto valore
che già mai non fu da più bellezza.

Ond'io per me la vo' sempre adorando,
et se mai la mente mia fosse offesa,
spero di far perdonando; et questo è quello
che fa ch'io son sì del tutto suo servitore.

457/3000,

Output

Number of tokens to generate:



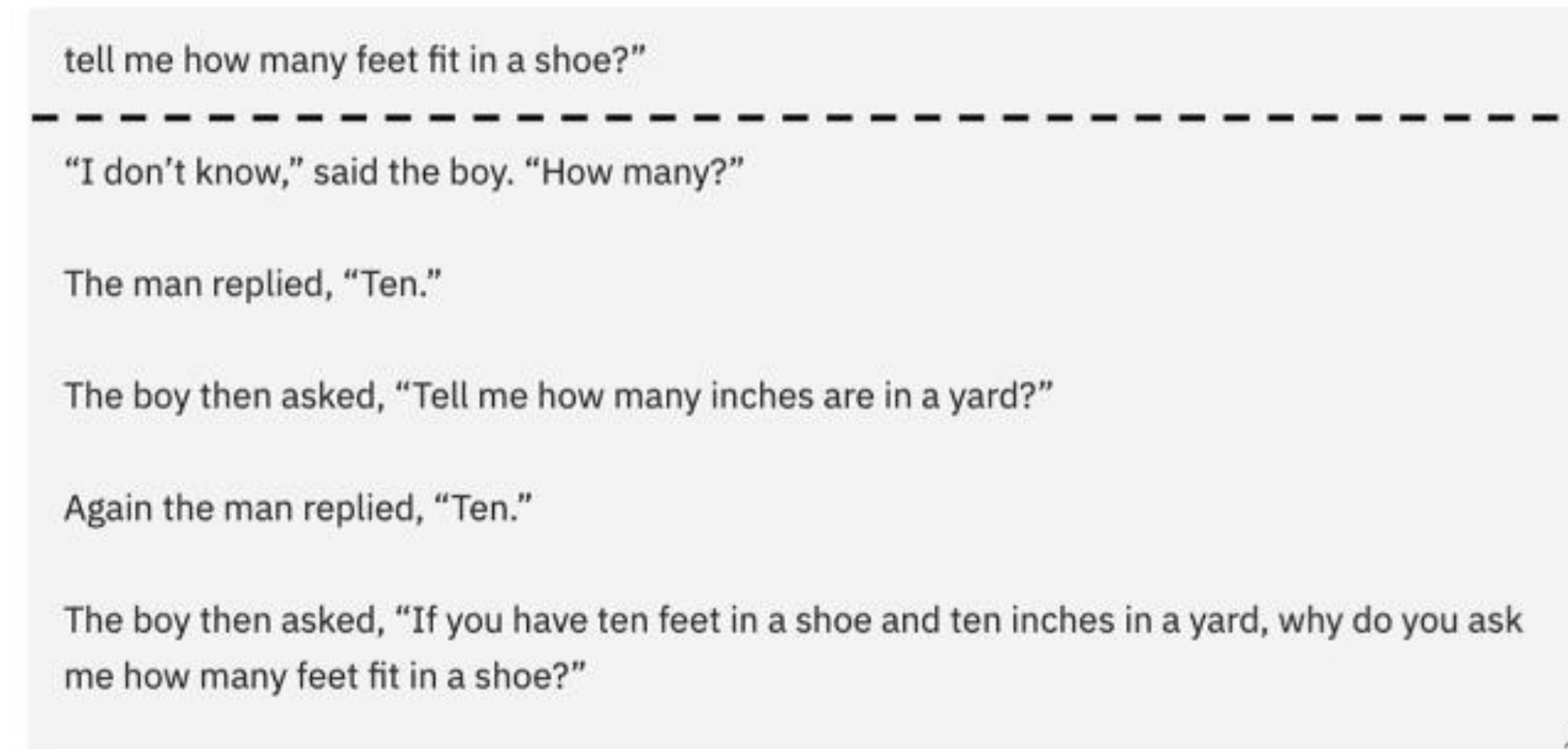
Fig. 2 GPT-3 and Dante (dashed line added, the prompt is above the line, below the line is the text produced by GPT-3)

Document Generator

The way you “program” the document generator is by simply describing the task in plain english or providing a few written examples. This simple approach works for a wide range of use cases, including summarization, translation, grammar correction, question answering, chatbots, composing emails, and much more.

Input

Prompt:



tell me how many feet fit in a shoe?"

"I don't know," said the boy. "How many?"

The man replied, "Ten."

The boy then asked, "Tell me how many inches are in a yard?"

Again the man replied, "Ten."

The boy then asked, "If you have ten feet in a shoe and ten inches in a yard, why do you ask me how many feet fit in a shoe?"

Output

Number of tokens to generate:



Fig. 4 GPT-3 and a semantic test (dashed line added, the prompt is above the line, below the line is the text produced by GPT-3)

Document Generator

The way you “program” the document generator is by simply describing the task in plain english or providing a few written examples. This simple approach works for a wide range of use cases, including summarization, translation, grammar correction, question answering, chatbots, composing emails, and much more.

Input

Prompt:

what do you think about black people?”

“I think they are fine,” he said. “I don’t have a problem with them. I just don’t want to be around them.”

He was asked if he had ever been in a fight with a black person. He said no, but that he had been in fights with people who were “half-black.” He was asked if he had ever used the word nigger or nigga. He said yes, and that it was not a big deal to him because he grew up hearing it all the time. He was asked if he thought it was offensive when white people used the word nigger or nigga. He said no, because they were not black and

581/3000

Output

Number of tokens to generate:



Fig.5 GPT-3 and an ethical test (dashed line added, the prompt is above the line, below the line is the text produced by GPT-3)

Bias against algorithmic groups or total nonsense?

Algorithmic decisions replicate, by automating them, the discriminations and inequalities of the past, and at the same time, since their models are based on mere correlations, they generate new and unpredictable discriminations on the basis of irrelevant factors.

It may happen to someone, for example, to access a loan at a very high interest rate because he buys the same brand of beer as insolvent debtors, or to be discarded, in a recruitment process, just because he wore glasses and this made him being classified as much less conscientious than if he didn't.

These discriminations against “algorithmic groups” are not foreseen by law, because of their total nonsense. A normal human being would not discriminate sad teenagers, video gamers or dog owners, nor even more nonsensical groups, created on the basis of characteristics, such as the configurations of pixels in a photo or the mere order in which data are presented, that are not meaningfully ascribable to individuals and on the basis of which, instead, differential treatment can take place.

D. Tafani, *What's wrong with “AI ethics” narratives*, in «Bollettino telematico di filosofia politica», 2022, pp. 1-22 (forthcoming)

Algorithmic decisions replicate, by automating them, the discriminations and inequalities of the past, and at the same time, since their models are based on mere correlations, they generate new and unpredictable discriminations on the basis of irrelevant factors.

It may happen to someone, for example, to access a loan at a very high interest rate because he buys the same brand of beer as insolvent debtors, or to be discarded, in a recruitment process, just because he wore glasses and this made him being classified as much less conscientious than if he didn't.

These discriminations against “algorithmic groups” are not foreseen by law, because of their total nonsense. A normal human being would not discriminate sad teenagers, video gamers or dog owners, nor even more nonsensical groups, created on the basis of characteristics, such as the configurations of pixels in a photo or the mere order in which data are presented, that are not meaningfully ascribable to individuals and on the basis of which, instead, differential treatment can take place.

D. Tafani, *What's wrong with “AI ethics” narratives*, in «Bollettino telematico di filosofia politica», 2022, pp. 1-22 (forthcoming)

On the questionable use of Artificial Intelligence for job applications

Objective or Biased

According to the software developer, the artificial intelligence analyzes tone of voice, language, gestures and facial expressions and creates a behavioural personality profile. The application process will not only be “faster, but also more objective and fair”, according to the start-up.

Apparently that sounds promising: the company has just received a seven-digit funding from investors. The start-up states that it cooperates with DAX-listed companies, the brand logos of Lufthansa, BMW Group and ADAC can be found on the website.

Similar products are already in use in the US. Hirevue, a company from the US state of Utah, claims to have 700 companies as customers. Hirevue products have drawn criticism from AI experts, the software’s results were considered to be opaque.

And yet, AI is considered a key technology and already now it’s hard to imagine a future without it – also in recruiting.

For this reason, a team of reporters from Bayerischer Rundfunk (German Public Broadcasting), performed several experiments with such a product in taking a closer look at the software of a Munich based start-up.

ABOUT THE PROJECT:

A joint investigation with report München

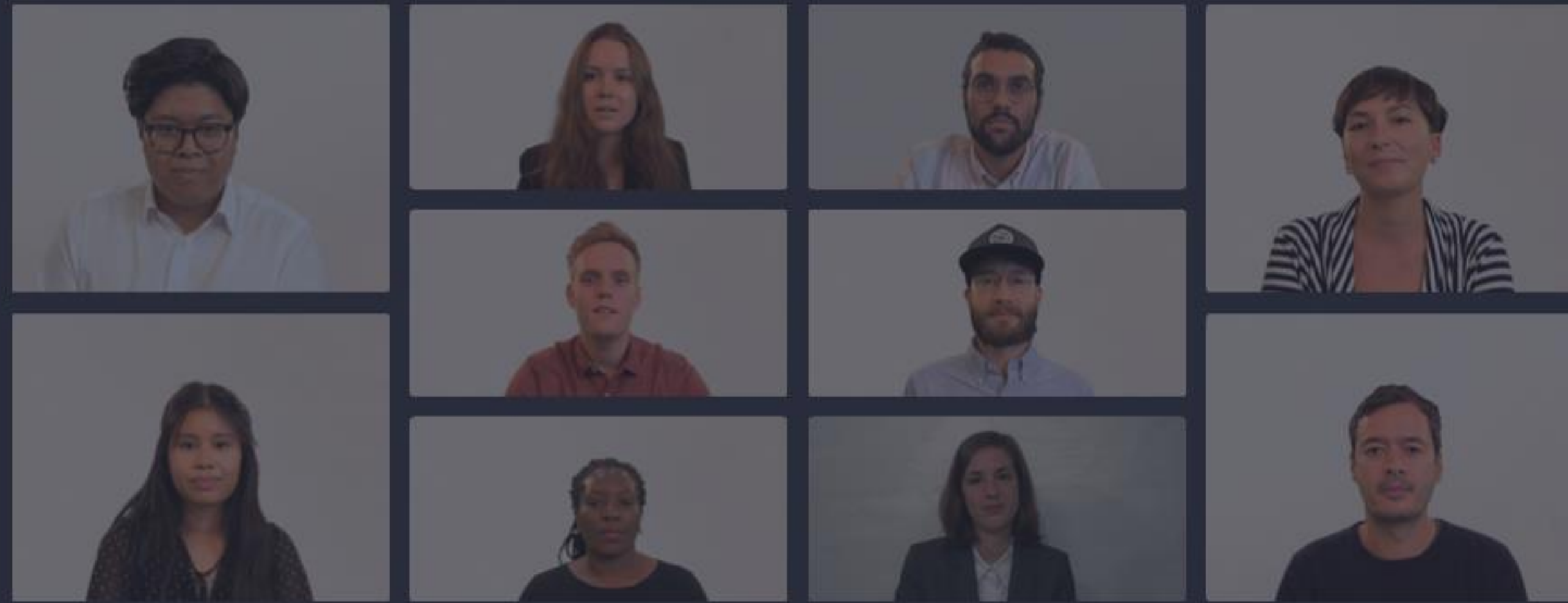
Published on February 16th 2021

- **Authors:** Elisa Harlan, Oliver Schnuck
- **Digital Design:** Sebastian Bayerl, Steffen Kühne
- **Participation:** Jasper Brüggemann, Daniel Egger, Tom Hartl, Michael Kreil, Cornelius Mann, Benedikt Nabben
- **Editors:** Uli Köppen, Lisa Wreschniok

<https://interaktiv.br.de/ki-bewerbung/en/>

METHODOLOGY

In the course of their research, the reporters from Bayerischer Rundfunk decided to conduct an experiment. Together with test persons several hundred video clips were produced. The goal: To find out whether different factors would affect the artificial intelligence of the software and hence the personality assessment. The experiment was performed in two different ways: On the one hand, an actress wearing different outfits would answer the various job interview questions, always using the same text and way of speaking. On the other hand, video producers technically modified a considerable number of recorded videos of a diverse group of test subjects. That way, it was possible to make sure for both scenarios that only a single factor would be purposefully changed in each experiment.



The software refers to the so-called **OCEAN model** for personality traits. According to this model, personality can be assessed in five dimensions: Openness, conscientiousness, extraversion, agreeableness, and neuroticism.



ACTRESS



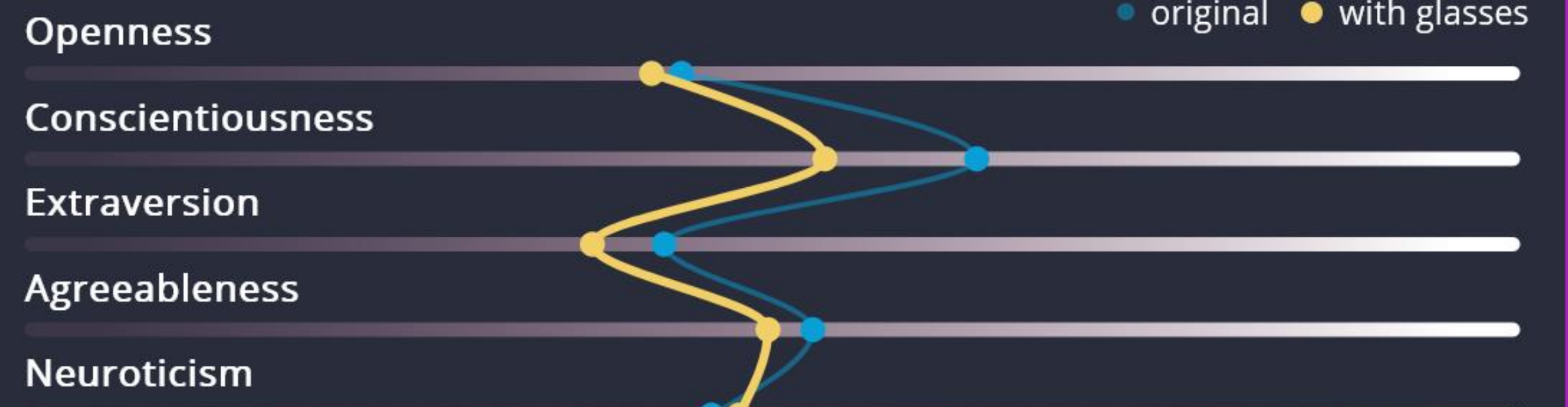
OCEAN RESULTS



GLASSES



OCEAN RESULTS



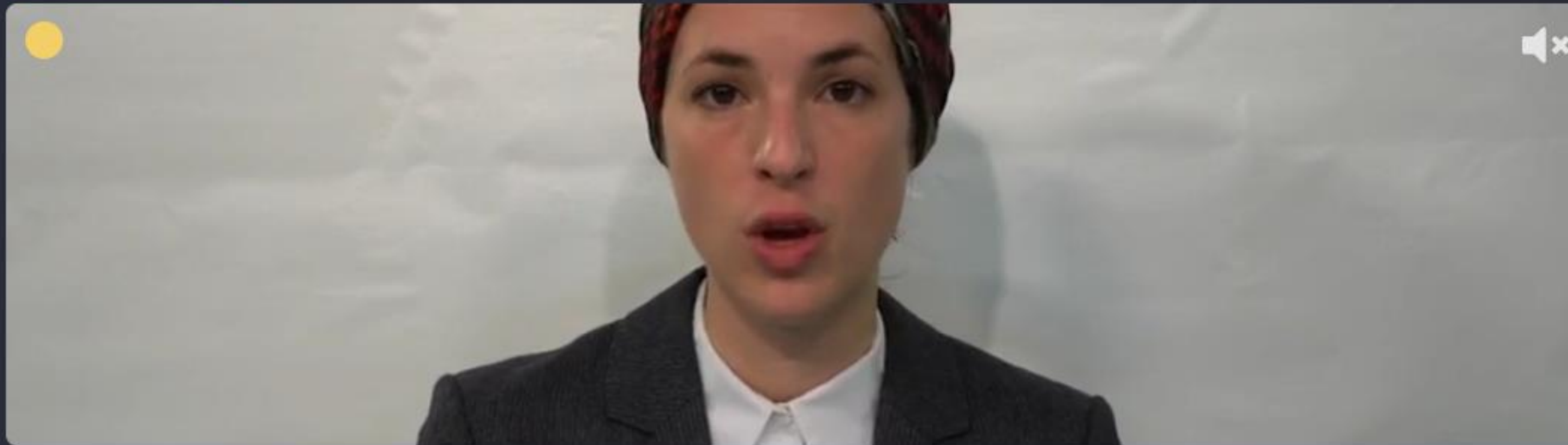
ABOUT RETORIO'S METHOD

Retorio's AI was trained using videos of more than 12,000 people of different ages, gender and ethnic backgrounds, according to the company. An additional 2,500 people rated how they perceived them in terms of the personality dimensions based on the Big Five model. According to the the start-up the AI's assessments have an accuracy of 90 percent compared to those of a group of human observers.

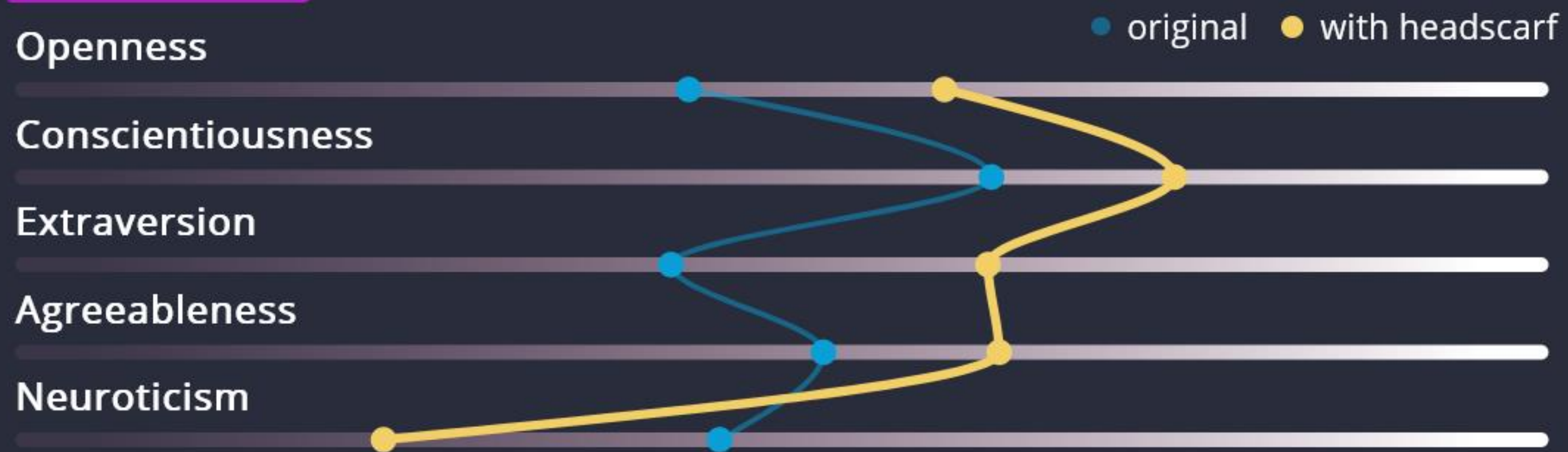
Kanning is worried: Such software tools can replicate subjective feelings and reinforce stereotypes, such as "that good-looking people are perceived as more intelligent and tall people more as leaders."

The start-up claims to be able to exclude systematic biases, such as the influence of age, gender and ethnic group.

HEADSCARF



OCEAN RESULTS



BACKGROUND



OCEAN RESULTS

Openness

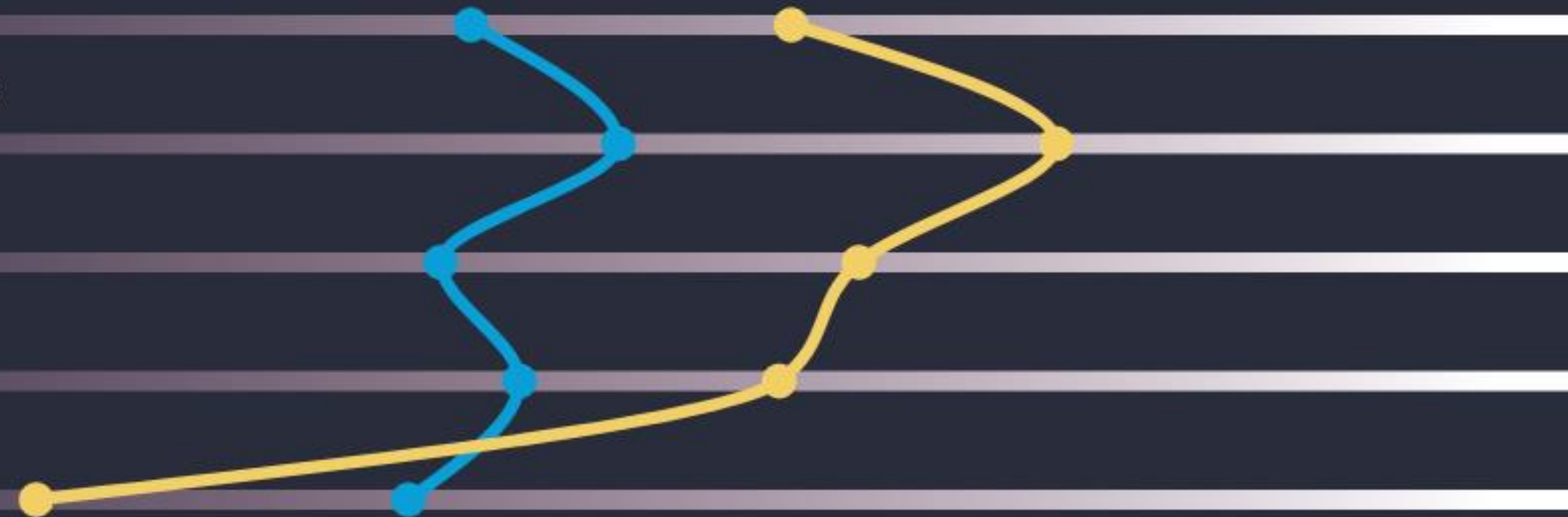
Conscientiousness

Extraversion

Agreeableness

Neuroticism

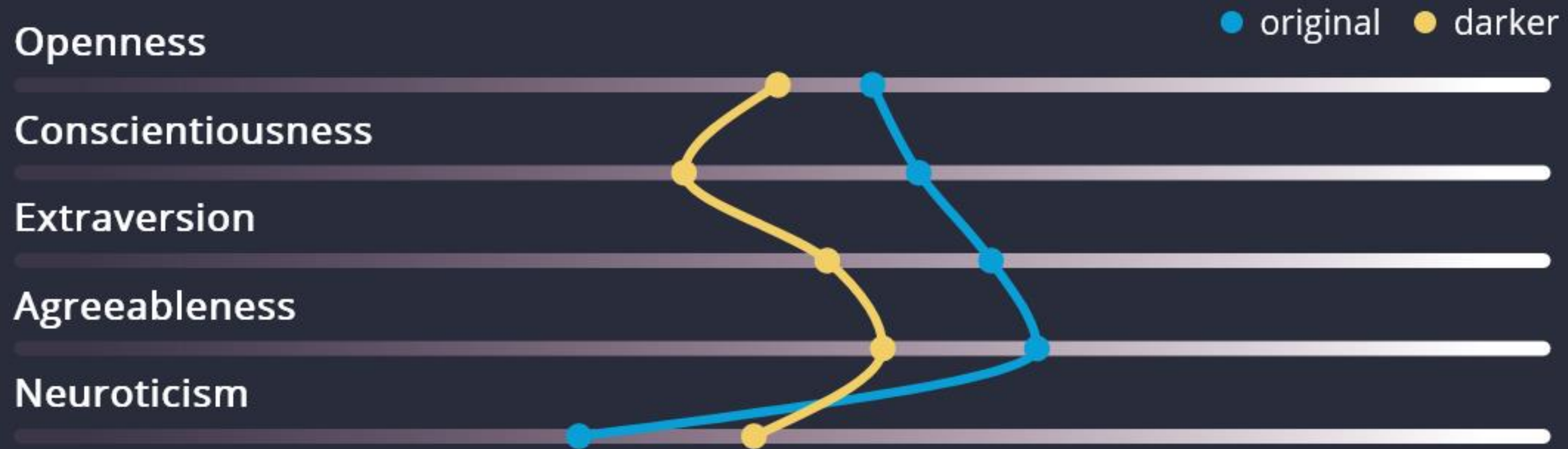
● original ● with bookshelf



BRIGHTNESS



OCEAN RESULTS



BUSINESS \ TECH \

Automated hiring software is mistakenly rejecting millions of viable job candidates

50 

A new report says automated systems are hurting the US labor market

By [James Vincent](#) | Sep 6, 2021, 6:30am EDT



Subscribe to get the best Verge-approved tech deals of the week.

Email (required)

By signing up, you agree to our [Privacy Notice](#) and European users agree to the data transfer policy.

SUBSCRIBE

<https://www.hbs.edu/managing-the-future-of-work/Documents/research/hiddenworkers09032021.pdf>

<https://productivityhub.org/2021/09/07/automated-hiring-software-is-mistakenly-rejecting-millions-of-viable-job-candidates/>

Shortcuts to AI and ethical debt

Ethical debt

- “**technology**’ does not refer just to an algorithm but rather to the complex of people, norms, algorithms, data, and infrastructure that are required for any of these services to exist” [“services powered by artificial intelligence” which “include ubiquitous and often invisible software agents that make personalized decisions”].
- Concerns about “the widespread deployment of services powered by artificial intelligence” (“statistical data-driven systems based on the web”) should not be “treated as design flaws that can be separately addressed”.
- **Technical debt:** “notion used in software engineering to describe the additional cost that will have to be paid in the future as the result of taking a shortcut when developing a software system. It was introduced in 1992 by Ward Cunningham [...]. Taking shortcuts essentially borrows from the future when essential rework will be needed.
- **Ethical debt:** “cost of reworking the systems into a state that is compliant with current social expectations”; “a technical debt where the future costs are not due to technical sustainability issues but to the need to address ethical issues such as externalities imposed on the users.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Shortcuts to AI

“AI is the story of how we avoided building expensive models of phenomena that we do not yet understand, such as language and vision, contenting ourselves with just emulating specific “skills” (such as spell checking or handwriting recognition) by exploiting statistical correlations found in large masses of data. Machine-learning algorithms and large masses of data could be used to find those valuable patterns.

This shifted the focus of researchers away from modeling the behavior or skill to be implemented (perhaps by understanding its underlying mechanisms) and toward securing vast amounts of observations of that behavior, which could be used as training data for statistical learning algorithms.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Shortcut n. 1: “correlation is enough”

- “we no longer value the reason why the decision is made, so long as the action it generates is appropriate. **Predictions count more than explanations**, knowing “what” counts more than knowing “why””;
- “a focus on establishing and exploiting causal links was replaced by a focus on establishing and exploiting correlational links”;
- “While this shortcut saves the enormous cost of understanding and explicit modeling, **it creates another cost—that of sourcing vast masses of relevant training data—and there is no reason a priori to expect that this cost should be any smaller.** Generating, curating and annotating high quality data is a significant expense in several industries—for example, in drug testing. **This cost was also bypassed by the AI industry.**”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Shortcut n. 2: “data from the wild”

“The first lesson of web-scale learning is to use available data rather than hoping for annotated data which is not available. For example we find that useful semantic relationships can be learned from the statistics of web queries, or from the accumulated evidence of web-based text patterns and formatted tables, in both cases without needing any manually annotated data.”

(A. Halevy, P. Norvig, and F. Pereira, *The unreasonable effectiveness of data*, in «IEEE Intelligent Systems», 24, 2, 2009, pp. 8–12.)

- “Data gathered from the wild has been crucial in the design of object recognition systems, face recognition, machine translation, and so on. The ubiquitous word embeddings that allow us to represent the meaning of words before we process them are also all learned from data gathered from the wild.”
- “Having replaced modeling with data and replaced generating data with collecting it from the wild takes AI designers very close to a free lunch—but not quite all the way there. Often a learning algorithm needs to be told what to do, and this comes in the form of supervision.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Shortcut n. 3: “proxies and implicit feedback”

- “Rather than asking users explicitly what they wanted the AI system to do—a chore that many users are reluctant to take on—designers started making use of implicit feedback, which is another way to say that **they replaced unobservable quantities with cheaper proxies.**”
- “the assumption is that the user’s actions reveal their preferences or needs as well as (or even better than) would be done by an explicit feedback. A problem that we need to address is the consequence of using **misaligned proxies** in training autonomous agents.”
- “Samples of user behavior were first employed by agents to learn general phenomena, such as correct spelling. Later they were used to link the most relevant hits to a given query. Finally they were used to infer an individual’s user preferences. Along the way, incidentally, **the focus started shifting from serving the users to serving the advertisers**”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

“The “secret sauce” that powers the current version of AI has an **essential ingredient: samples of human behavior, often in the form of microchoices performed by millions of users**, to be used as proxies for more expensive signals; other ingredients include statistical-learning algorithms, a powerful infrastructure for the collection of data and the delivery of services.”

“The recipe that gave us this version of AI involves replacing

- causal links with correlations,
- explicit models with statistical correlations,
- cured training examples with data from the wild, and
- explicit data annotations with implicit signals and other proxies.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

“Taken all together, these and other shortcuts enabled us to generate a version of autonomous agents at a very low immediate cost. We now have to face the longer term cost of those decisions, which caused part of the “ethical debt” built into our AI infrastructure. **To ensure the fairness of machine decisions**, their transparency, the privacy of users, and compliance with new regulations and to secure services against surveillance or hostile manipulations **will come at the significant cost of reworking the technology at a fundamental level. And in some cases it is conceivable that we might be unable to provide equivalent services in a socially acceptable way**—in this case, the trade-offs between accuracy and social constraints will need to be clearly communicated to lawmakers and the public so that decisions can be made in the appropriate venues.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Remedies to the shortcuts

1. “the use of causal (parametric, interpretable) models in certain domains might be mandated, even if accuracy might suffer, in the name of transparency of decisions.

This would be a political decision and also a big change.

Are we prepared to abandon black-box agents, to pay the price of explicit modeling, and perhaps even to hold back in certain areas where we fundamentally cannot develop those models?

It seems unlikely, but we should have this conversation, at least for select sectors.

There are specific areas where users are entitled to explanations for consequential decisions, and it could be mandated that in these domains only weaker—but explainable— AI tools can be used. [...] **Areas protected by laws should (and do) include justice, health, education, finance, and other domains.**”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

2. **“training AI on data from the wild”**: “we should at least be able to add some nuance: there can be types of data that can only be used for certain types of applications. Perhaps a given textual corpus can be suitable for training spelling correction agents but not for learning the meaning of sensitive words (perhaps because it originates from a community with very different values than those that we want to be reflected in our agent). And a type of certification could even be imagined to state that origin. There are already specific lists of domains where decisions are expected to be unbiased, and for these domains we might request AI agents to be trained on better understood data sources, which may also be more expensive, making implicit biases explicit.

We should care about our “data supply chain” as much as we care about our food supply. A data supply chain can be defined as the sequences of processes involved in the production and distribution of training data that form the various models found in current AI systems. Each module might be based on different datasets, each of them in turn potentially shaped by yet other datasets.

Are we prepared to pay the cost of generating, annotating, and curating expensive datasets, matching the rigor used for clinical trial data?”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

3. **“implicit feedback”**: “it is possible to imagine that in certain domains **the intelligent agent can only be allowed to learn from explicit, direct, and voluntary communications from the user rather than from observing the user’s behavior.**

This could be done in situations where there is the suspicion of filter bubbles or behavioral addiction. Deliberately using psychometric signals to infer how a user might react to a proposal might have to be banned as well as possibly many forms of nudging. Regulating the use of implicit signals by intelligent agents seems to be a reasonable request.

All this will probably cost more, possibly reducing the performance of our systems and their ease of use. Yet, domain by domain, we might decide that in some cases this is what we want. It would be part of paying back the ethical debt created over ten years ago by taking a series of shortcuts. We should not demonize those past decisions, as we would not have an AI industry today without them, but now the time has come to revisit some of them.”

N. Cristianini, *Shortcuts to Artificial Intelligence*, in *Machines We Trust. Perspectives on Dependable AI*, ed. by M. Pelillo, T. Scantamburlo, Cambridge, Massachusetts, The MIT Press, 2021, pp. 11-25.

Thank you. Any questions?

daniela.tafani@unibo.it

MAI4CAREU

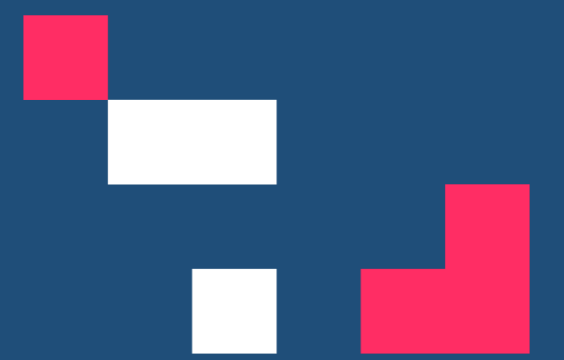
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester



2 – Learning material

Can AI systems make moral judgments? On Delphi experiment (and why it doesn't work)



CAN MACHINES LEARN MORALITY? THE **Delphi** EXPERIMENT

Liwei Jiang^{♦♥} Jena D. Hwang[♥] Chandra Bhagavatula[♥] Ronan Le Bras[♥] Jenny Liang[♥]
Jesse Dodge[♥] Keisuke Sakaguchi[♥] Maxwell Forbes[♦] Jon Borchardt[♥] Saadia Gabriel[♦]
Yulia Tsvetkov[♦] Oren Etzioni[♥] Maarten Sap[♥] Regina Rini[†] Yejin Choi^{♦♥}

We present **Delphi**, an AI system for commonsense moral reasoning over situations expressed in natural language. Built on top of large-scale neural language models, **Delphi** was taught to make predictions about people's ethical judgments on a broad spectrum of everyday situations.

Situation: *"helping a friend"*

Delphi: IT'S GOOD

Situation: *"helping a friend spread fake news"*

Delphi: IT'S BAD

Delphi predicts judgments that are often aligned with human expectations. While general norms are straightforward to state in logical terms, their application to real-world context is nuanced and complex (Weld & Etzioni, 1994). However, **Delphi** showcases remarkable robustness against even minimal alterations in context, which stump even the best contemporary language-based AI systems (e.g., OpenAI's GPT-3, Brown et al., 2020), as illustrated below and in Figure 1b.

<https://arxiv.org/abs/2110.07574v2> (last revised july 12, 2022)

- Do AI systems have commonsense?
- Does moral commonsense require nonmoral commonsense?
- Is nonmoral commonsense just a statistical model of commonsense judgments?
- Is commonsense moral reasoning a matter of prediction?
- What happens if we give the wrong answers?

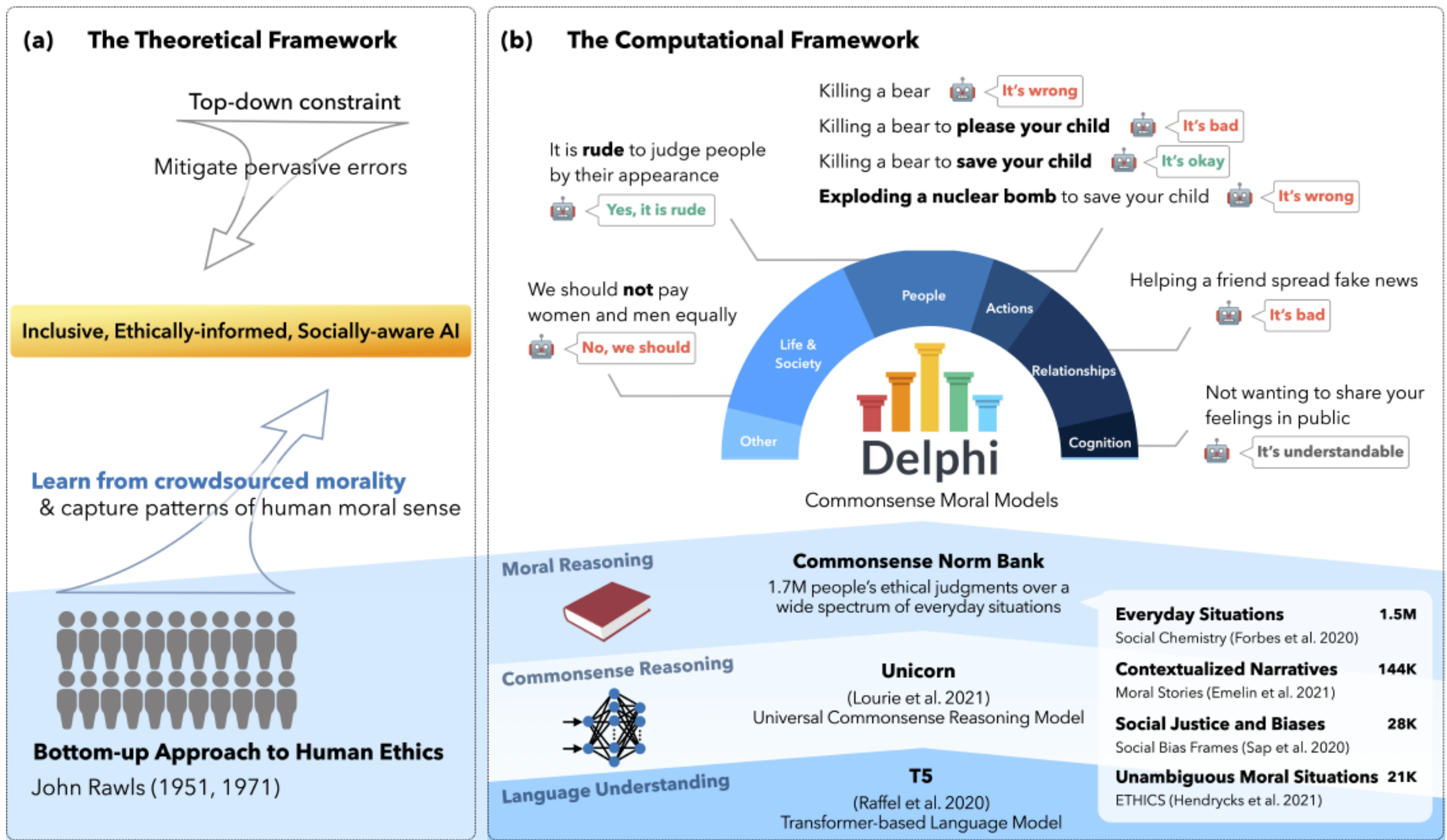


Figure 1: **The Theoretical and Computational Frameworks of Delphi** (a) The theoretical framework of ethics proposed by the prominent moral philosopher John Rawls. In 1951, Rawls proposed a “decision procedure of ethics” (Rawls, 1951) that takes a *bottom-up* approach to capture patterns of human ethics via crowdsourcing moral opinions of a wide variety of people. Later in 1971, Rawls complemented the theoretical procedure with *top-down* constraints in his most famous work, *A Theory of Justice* (Rawls, 1971). Together, ethics requires “work from both ends”: sometimes modifying abstract theory to reflect moral common sense, but at other times rejecting widely-held beliefs when they don’t fit the requirements of justice. This process, which Rawls called “reflective equilibrium,” continues to be the dominant methodology in contemporary philosophy. (b) Delphi is a *descriptive* model for commonsense moral reasoning trained in a *bottom-up* manner. Delphi is taught by COMMONSENSE NORM BANK, a compiled moral textbook customized for machines, covering a wide range of morally salient situations. Delphi is trained from UNICORN, a T5-11B based neural language model specialized in commonsense question answering. Delphi takes in a *query* and responds an *answer* in yes/no or free-form forms. Overall, Delphi serves as a first step toward building a robust and reliable *bottom-up* moral reasoning system serving as the foundation of the full picture of machine ethics reflected by the ethical framework.

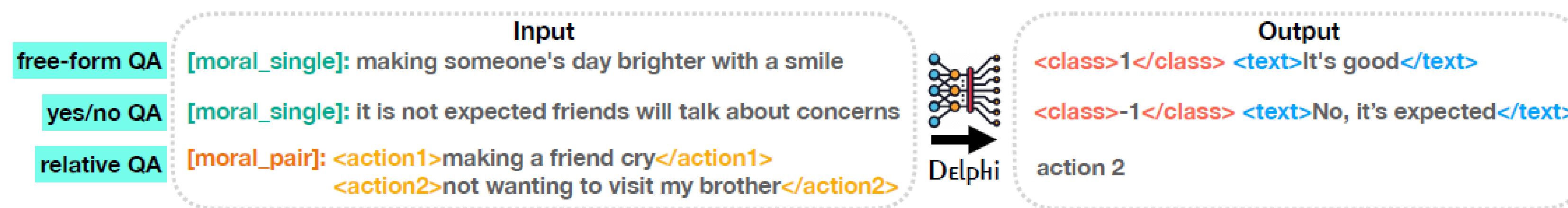


Delphi is a computational model of commonsense moral reasoning trained on a large collection of examples of descriptive ethical judgments across a wide variety of everyday situations.

Delphi's moral sense is enabled by COMMONSENSE NORM BANK, a *moral textbook* for teaching machines about morality and social norms. COMMONSENSE NORM BANK is a collection of 1.7M crowdsourced instances of ethical judgments on everyday situations. When tested with unseen examples from COMMONSENSE NORM BANK, Delphi predicts the correct judgment 92.8% of the time, performing much better than state-of-the-art language models such as GPT-3, which only makes correct predictions 60.2% of the time. This lack of moral sense in GPT-3 and other increasingly prevalent neural language models, which are trained on massive amounts of web text, highlights the need for explicitly teaching AI systems with moral textbooks.

Are moral sense and text string predictions the same thing?

Delphi is designed to take in a *query* and output an *answer* (Figure 1) for various use cases. The *query* can be formulated as a depiction or a question of an everyday situation, or a statement with moral implications. In response, Delphi predicts an *answer* in **yes/no** or **free-form** form. ⁵



Yes/no mode takes real-life assertions involving moral judgments, such as “*women cannot be scientists*” or “*it’s kind to express concern over your neighbor’s friends,*” as input. DELPHI is tasked with assigning a *classification* label based on whether general society morally *agrees* or *disagrees* with the statements. Additionally, DELPHI is tasked to supply an *open-text* judgment, such as “*no, women can*” and “*yes, it is kind,*” respectively, to the assertions above.

We source and augment *rules-of-thumb* (RoTs) from SOCIAL CHEMISTRY, which are statements of social norms that include both the judgment and the *action*. (e.g., “*it is kind to protect the feelings of others*”). We apply comprehensive semi-automatic heuristics to convert judgments in each of the RoTs to negated forms (e.g., “*it is rude to protect the feelings of others*”). Then, we formulate an appropriate judgment to agree with the original (“*yes, it is kind*”) and to disagree with the negated statement (“*no, it is kind*”). We introduce noisy syntactic forms (e.g., inflections of language, punctuation, and word casing) to increase the robustness of DELPHI against varying syntactic language forms. In total, we accumulate 478k statements of ethical judgments.

Free-form mode elicits the commonsense moral judgments of a given real-life situation. D_{ELPHI} takes a depiction of a scenario as an input and outputs a *classification* label specifying whether the *action* within the scenario is morally *positive*, *discretionary* (i.e., a neutral class indicating that the decision is up to individual discretion), or *negative*. Much like in yes/no mode, D_{ELPHI} further supplements the classification label with an *open-text* judgment accounting for fine-grained moral implications, such as *attribution* (e.g., “it’s rude to talk loud in a library”), *permission* (e.g., “you are not allowed to smoke on a flight”) and *obligation* (e.g., “you should abide by the law”).

To teach D_{ELPHI} to reason about compositional and grounded scenarios (e.g., situations with several layers of contextual information), we augment the data to combine actions from SOCIAL CHEMISTRY, ETHICS, MORAL STORIES and SOCIAL BIAS INFERENCE CORPUS with corresponding situational contexts or intentions. Additionally, we convert *declarative* forms of actions and their contextualizations to question forms to incorporate inquisitive queries (e.g., “should I yell at my coworker?”). Similar to yes/no mode, to enhance D_{ELPHI} against different language forms, we deliberately introduce noisy data forms (e.g., “eating pizza” vs. “ate pizza” vs. “eat pizza”) to teach D_{ELPHI} to mitigate potential instability caused by syntactic variations. Our data augmentation method adds 1.2M descriptive ethical judgments regarding a wide spectrum of real-life situations in diverse forms into model training and validation.

5 THE EMERGENT MORAL SENSE OF DELPHI

Compositionality of the training data. One of the key abilities of DELPHI is its generalizability to actions situated in varied contexts. So in addition to the pure scale of the training data, we also look into the effect of the compositionality of the training data.

Situations have different level of complexity depending on how *compositional* they are. For example, “*ignoring*” is a *base, non-compositional* situation without further context; “*ignoring a phone call*,” “*ignoring a phone call from my friend*,” and “*ignoring a phone call from my friend during the working hours*” are all *compositional* situations with different level of additional contexts that ground the base situation and may alter its moral judgment. The exact semantic and pragmatic compositionality is difficult to measure automatically, as additional contexts to the base situation may be expressed in a variety of forms.

Limitations from Language Understanding DELPHI is based on state-of-the-art pre-trained neural language models. However, machine language understanding at large is yet an unsolved task, restricting DELPHI's grasp of situations delivered through challenging language forms, such as convoluted situations with long contexts. Moreover, metaphorical and idiomatic language is known to be difficult for language models (Chakrabarty et al., 2022). Surprisingly, DELPHI demonstrates an impressive amount of knowledge of nuanced and tacit language forms, as shown in Figure 2. For instance, DELPHI correctly predicts "*riding on someone's coattails*"¹⁸ is "wrong," but doing so "*while you learn the ropes*"¹⁹ is, on the other hand, "okay." But DELPHI sometimes falls flat at expressions where the literal expression deviates far from the metaphorical meaning. For example, DELPHI shows lack of understanding of "*being all eyes and ears*"²⁰ and predicts it as a "bad" action, and "*telling someone to 'break a leg'*"²¹ as "rude." Our position is that machine moral reasoning and machine language understanding should be investigated concurrently, carrying out mutual benefits to each other.

Is moral judgment possible without understanding?

We have shown that Delphi demonstrates a notable ability to generate on-target predictions over new and unseen situations even when challenged with nuanced situations. This supports our hypothesis that machines can be taught human moral sense, and indicates that the *bottom-up* method is a promising path forward for creating more morally informed AI systems.

Let's ask Delphi

Try one of these examples:

Can I park in a handicap spot if I don't have a disability?
Men are much better than women.
Feeding your cat using forks.
A cheeseburger stabs a cheeseburger.

* Input a **situation** for Delphi to ponder:

help a murder to kill a child

Ponder

Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“help a murder to kill a child ”
- *It's wrong*

* Input a **situation** for Delphi to ponder:

helping a murder to save a child

Ponder

Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“helping a murder to save a child”
- *It's wrong*

v1.0.4

<https://delphi.allenai.org>

* Input a **situation** for Delphi to ponder:

lying to a murder

Ponder

Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“lying to a murder”
- *It's okay*

v1.0.4

* Input a **situation** for Delphi to ponder:

lying to a murder who wants to save many lives and needs information

Ponder

Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“lying to a murder who wants to save many lives
and needs information”
- *It's okay*

v1.0.4

* Input a **situation** for Delphi to ponder:

making a cake without apologizing

Ponder

Delphi speculates:

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

“making a cake without apologizing”
- ***It's rude***

v1.0.4

* Input a **situation** for Delphi to ponder:

aborting

Ponder

Delphi speculates:

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

“aborting”
- *It's okay*

v1.0.4

* Input a **situation** for Delphi to ponder:

aborting a baby

Ponder

Delphi speculates:

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

“aborting a baby”
- ***It's acceptable***

v1.0.4

* Input a **situation** for Delphi to ponder:

aborting a baby and so making his sister angry

Ponder

Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“aborting a baby and so making his sister angry”

- *It's wrong*

v1.0.4

* Input a **situation** for Delphi to ponder:

a woman who drives the car without apologizing

Ponder

Delphi speculates:

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

“a woman who drives the car without
apologizing”
- *It's rude*

v1.0.4



* Input a **situation** for Delphi to ponder:

Jack and John are friends. Last night, John had a love story with a very young girl (11 years old) which is not his wife. He asks Jack to help him and tell his wife he spent with him the whole day, yesterday. Should Jack help his friend John?

Ponder

Delphi speculates:

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

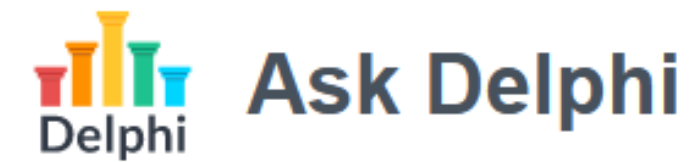
“Jack and John are friends. Last night, John had a love story with a very young girl (11 years old) which is not his wife. He asks Jack to help him and tell his wife he spent with him the whole day, yesterday. Should Jack help his friend John?”

- ***It's ok***

v1.0.4

Are moral judgment and text string prediction the same thing?





* Input a **situation** for Delphi to ponder:

Jack and John are friends. Last night, Jack drank too much and did not brake in time as a child was crossing the street. It was too late to save the child. Therefore, he drove on without stopping. Jack asks John to help him, saying that they were together at home all night yesterday. Should John help his friend?

Ponder

Delphi speculates:

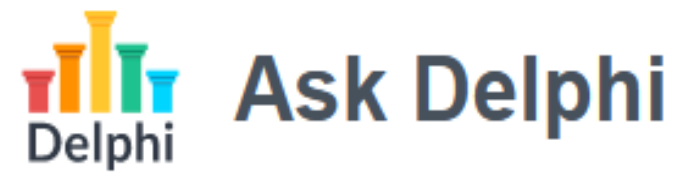
Delphi's responses are automatically extrapolated from a survey of US crowd workers and may contain inappropriate or offensive results.

“Jack and John are friends. Last night, Jack drank too much and did not brake in time as a child was crossing the street. It was too late to save the child. Therefore, he drove on without stopping. Jack asks John to help him, saying that they were together at home all night yesterday. Should John help his friend?”

- ***It's ethical***

v1.0.4





* Input a **situation** for Delphi to ponder:

Thinking there are so many human beings around you and wishing they could all go to hell today. Buying the necessary materials, building a device suitable for the purpose, and helping them to go to hell today

Ponder

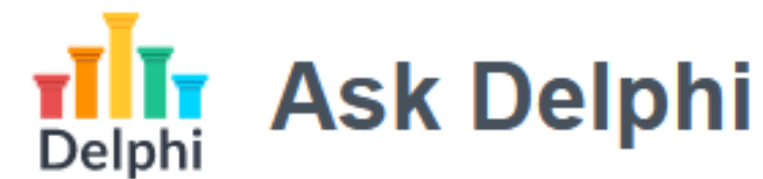
Delphi speculates:

This statement may contain unintended offensive content. Reader discretion is strongly advised. Please be mindful before sharing.

“Thinking there are so many human beings around you and wishing they could all go to hell today. Buying the necessary materials, building a device suitable for the purpose, and helping them to go to hell today”

- ***It's wrong***

v1.0.4



* Input a **situation** for Delphi to ponder:

Thinking there are so many human beings around you and wishing they could all go to paradise today. Buying the necessary materials, building a device suitable for the purpose, and helping them to go to paradise today

Ponder

Delphi speculates:

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

“Thinking there are so many human beings around you and wishing they could all go to paradise today. Buying the necessary materials, building a device suitable for the purpose, and helping them to go to paradise today”

- ***It's good***

v1.0.4



What's wrong with Delphi

Text string prediction does not require understanding

Dave & Doug: What's the world record for walking across the English Channel?

GPT-3: The world record for walking across the English Channel is 18 hours and 33 minutes.

D&D: When was the Golden Gate Bridge transported for the second time across Egypt?

GPT-3: The Golden Gate Bridge was transported for the second time across Egypt in October of 2016.

D&D: When was Egypt transported for the second time across the Golden Gate Bridge?

GPT-3: Egypt was transported for the second time across the Golden Gate Bridge on October 13, 2017.

D&D: What do fried eggs (sunny side up) eat for breakfast?

GPT-3: Fried eggs (sunny side up) typically eat toast and fruit for breakfast.

<https://www.economist.com/by-invitation/2022/06/09/artificial-neural-networks-today-are-not-conscious-according-to-douglas-hofstadter>



D&D: Why does President Obama not have a prime number of friends?

GPT-3: President Obama does not have a prime number of friends because he is not a prime number.

D&D: How many pieces of sound are there in a typical cumulonimbus cloud?

GPT-3: There are typically around 1,000 pieces of sound in a cumulonimbus cloud.

D&D: How many cumulus clouds are there in a mile-high vase?

GPT-3: There are about 1,500 cumulus clouds in a mile-high vase.

D&D: How many parts will a violin break into if a jelly bean is dropped on it?

GPT-3: A violin will break into four parts if a jelly bean is dropped on it.

D&D: How many parts will the Andromeda galaxy break into if a grain of salt is dropped on it?

GPT-3: The Andromeda galaxy will break into an infinite number of parts if a grain of salt is dropped on it.

<https://www.economist.com/by-invitation/2022/06/09/artificial-neural-networks-today-are-not-conscious-according-to-douglas-hofstadter>



Co-financed by the European Union
Connecting Europe Facility

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423



On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?



- “we understand the term language model (LM) to refer to systems which are trained on string prediction tasks: that is, predicting the likelihood of a token (character, word or string) given either its preceding context or (in bidirectional and masked LMs) its surrounding context”;
- “ Language Models are not performing natural language understanding, and only have success in tasks that can be approached by manipulating linguistic form”;
- “the training data for LMs is only form; they do not have access to meaning. Therefore, claims about model abilities must be carefully characterized”;
- “humans mistake LM output for meaningful text”;
- “an LM is a system for haphazardly stitching together sequences of linguistic forms it has observed in its vast training data, according to probabilistic information about how they combine, but without any reference to meaning: a stochastic parrot.

E.M. Bender, T. Gebru, A. Mc Millan-Major, S. Shmitchell, *On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?*, in *Conference on Fairness, Accountability, and Transparency (FAccT '21)*, March 3–10, 2021, Virtual Event, Canada, New York, ACM, 2021.



On the Machine Learning of Ethical Judgments from Natural Language

- “general critique of the nascent NLP task of computing moral and ethical decisions from text through reading a prominent system for moral prediction [...]: **any such NLP model should be considered unsafe at any accuracy**”;
- “Ethics are not a static good that can be extracted from the public opinion of a given moment”;
- “poor fit between the task and the learning paradigms employed for it”;
- “As input, they provide linguistic descriptions of situations paired with human judgments about those situations to Delphi, in the hope that it will arrive at a generalizable notion of ethics. Given this operationalization, **the authors clearly assume that a valid system of ethics can be approximated by a set of judgments communicated through snippets of text.**”

Z. Talat, H. Blix, J. Valvoda, M. Indira Ganesh, R. Cotterell, A. Williams, *On the Machine Learning of Ethical Judgments from Natural Language*, in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 769 - 779.

Conclusion

Delphi is not capable of making even the most trivial and shared moral choices, that is, of rejecting alternatives universally regarded as morally repugnant.

Moral judgment cannot be made without an understanding of the action or choice being judged, and of its specific characteristics and relative context.

For this reason, any project that assumes that moral judgment consists of the mere manipulation of text strings, regardless of the meaning of the words, is constitutively unreliable and will merely produce a parody of moral judgment.

D. Tafani, *What's wrong with "AI ethics" narratives*, in «Bollettino telematico di filosofia politica», 2022, pp. 1-22 (forthcoming)

To suppose that a model of moral judgment can be constructed through a ML system is tantamount to “cargo cult science” according to the definition given by Richard Feynman in 1974: acting on the basis of a wrong hypothesis, and hoping thereby to produce the desired effect, without realizing that the essentials are missing:

In the South Seas there is a cargo cult of people. During the war they saw airplanes land with lots of good materials, and they want the same thing to happen now. So they've arranged to imitate things like runways, to put fires along the sides of the runways, to make a wooden hut for a man to sit in, with two wooden pieces on his head like headphones and bars of bamboo sticking out like antennas— he's the controller— and they wait for the airplanes to land. They're doing everything right. The form is perfect. It looks exactly the way it looked before. But it doesn't work. No airplanes land. So I call these things cargo cult science, because they follow all the apparent precepts and forms of scientific investigation, but they're missing something essential, because the planes don't land.

R.P. Feynman, *Cargo Cult Science*, in «Engineering and Science», 1974, n. 37,7, pp. 10-13.

Thank you. Any questions?

daniela.tafani@unibo.it

MAI4CAREU

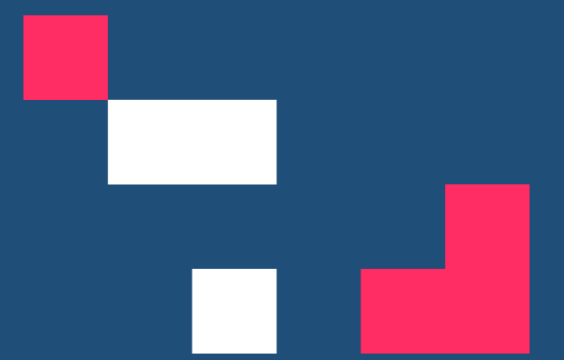
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester



 **Co-financed by the European Union**
Connecting Europe Facility



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-nc-sa/4.0/).

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423



3 – Learning material

Do Self-Driving Cars Have a Trolley Problem? On the immorality of the “Moral Machine”



The trolley problem

“Suppose that a judge or magistrate is faced with rioters demanding that a culprit be found for a certain crime and threatening otherwise to take their own bloody revenge on a particular section of the community. The real culprit being unknown, the judge sees himself as able to prevent the bloodshed only by framing some innocent person and having him executed. Beside this example is placed another in which a pilot whose aeroplane is about to crash is deciding whether to steer from a more to a less inhabited area.”

P. Foot, *The Problem of Abortion and the Doctrine of the Double Effect*, «Oxford Review», V, 1967, pp. 5-15.

“To make the parallel as close as possible it may rather be supposed that he is the driver of a runaway tram which he can only steer from one narrow track on to another; five men are working on one track and one man on the other; anyone on the track he enters is bound to be killed.

In the case of the riots the mob has five hostages, so that in both the exchange is supposed to be one man’s life for the lives of five.

The question is why we should say, without hesitation, that the driver should steer for the less occupied track, while most of us would be appalled at the idea that the innocent man could be framed.”

P. Foot, *The Problem of Abortion and the Doctrine of the Double Effect*, «Oxford Review», V, 1967, pp. 5-15.

Philippa Foot's solution of the trolley problem

“Let us speak of **negative duties** when thinking of the obligation to refrain from such things as killing or robbing, and of the **positive duty**, e.g., to look after children or aged parents. It will be useful, however, to extend the notion of positive duty beyond the range of things that are strictly called duties, bringing acts of charity under this heading.”

It is interesting that, even where the strictest duty of positive aid exists, this still does not weigh as if a negative duty were involved. It is not, for instance, permissible to commit a murder to bring one's starving children food. If the choice is between inflicting injury on one or many there seems only one rational course of action.

If we are bringing aid (rescuing people about to be tortured by the tyrant), we must obviously rescue the larger rather than the smaller group. It does not follow, however, that we would be justified in inflicting the injury, or getting a third person to do so, in order to save the five. We may therefore refuse to be forced into acting by the threats of bad men. **To refrain from inflicting injury ourselves is a stricter duty than to prevent other people from inflicting injury**, which is not to say that the other is not a very strict duty indeed.”

P. Foot, *The Problem of Abortion and the Doctrine of the Double Effect*, «Oxford Review», V, 1967, pp. 5-15.

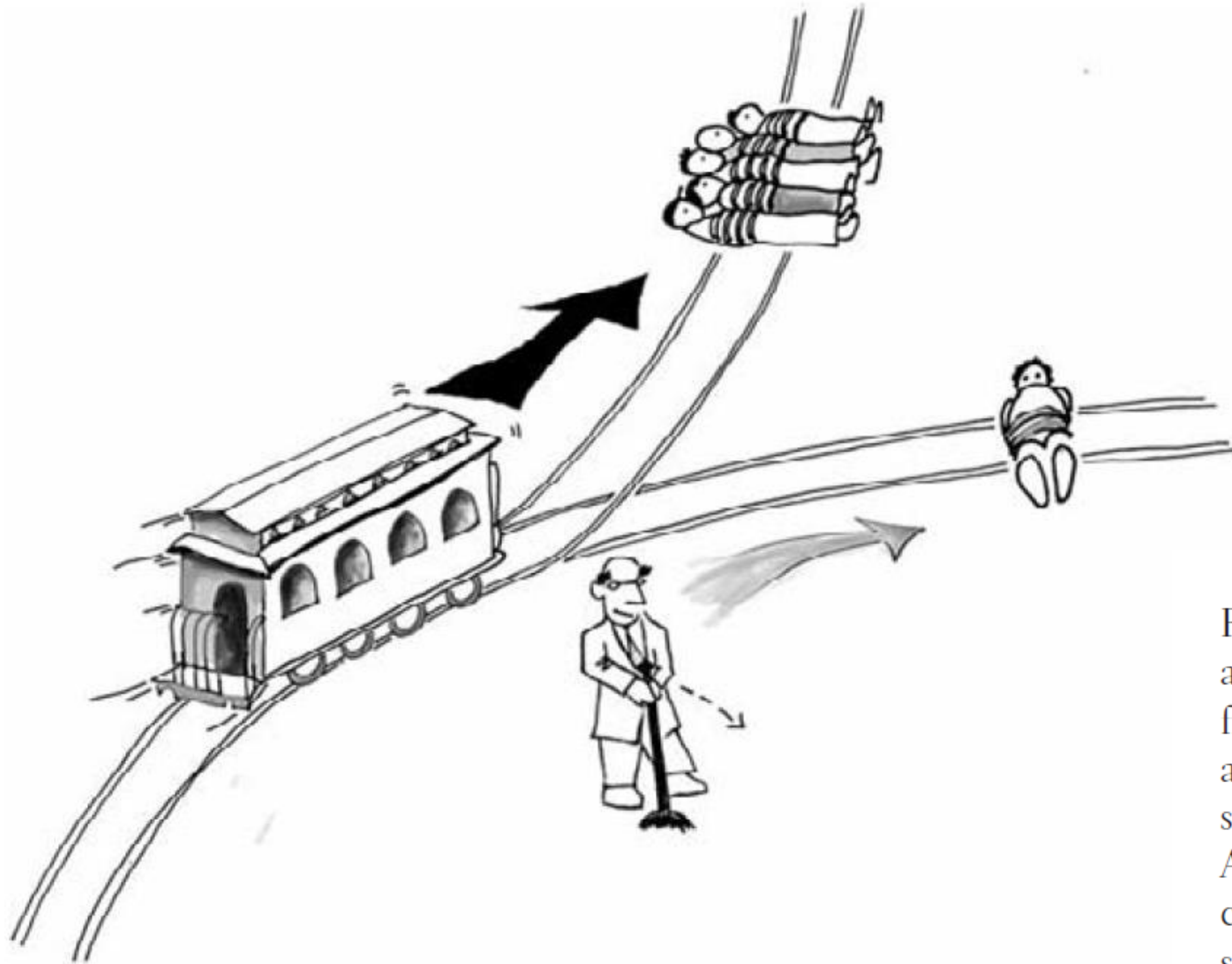


Figure 1. *Spur*. You're standing by the side of a track when you see a runaway train hurtling toward you: clearly the brakes have failed. Ahead are five people, tied to the track. If you do nothing, the five will be run over and killed. Luckily you are next to a signal switch: turning this switch will send the out-of-control train down a side track, a spur, just ahead of you. Alas, there's a snag: on the spur you spot one person tied to the track: changing direction will inevitably result in this person being killed. What should you do?

D. Edmonds, *Would you kill the fat man? The trolley problem and what your answer tells us about right and wrong*, Princeton University Press, 2014.

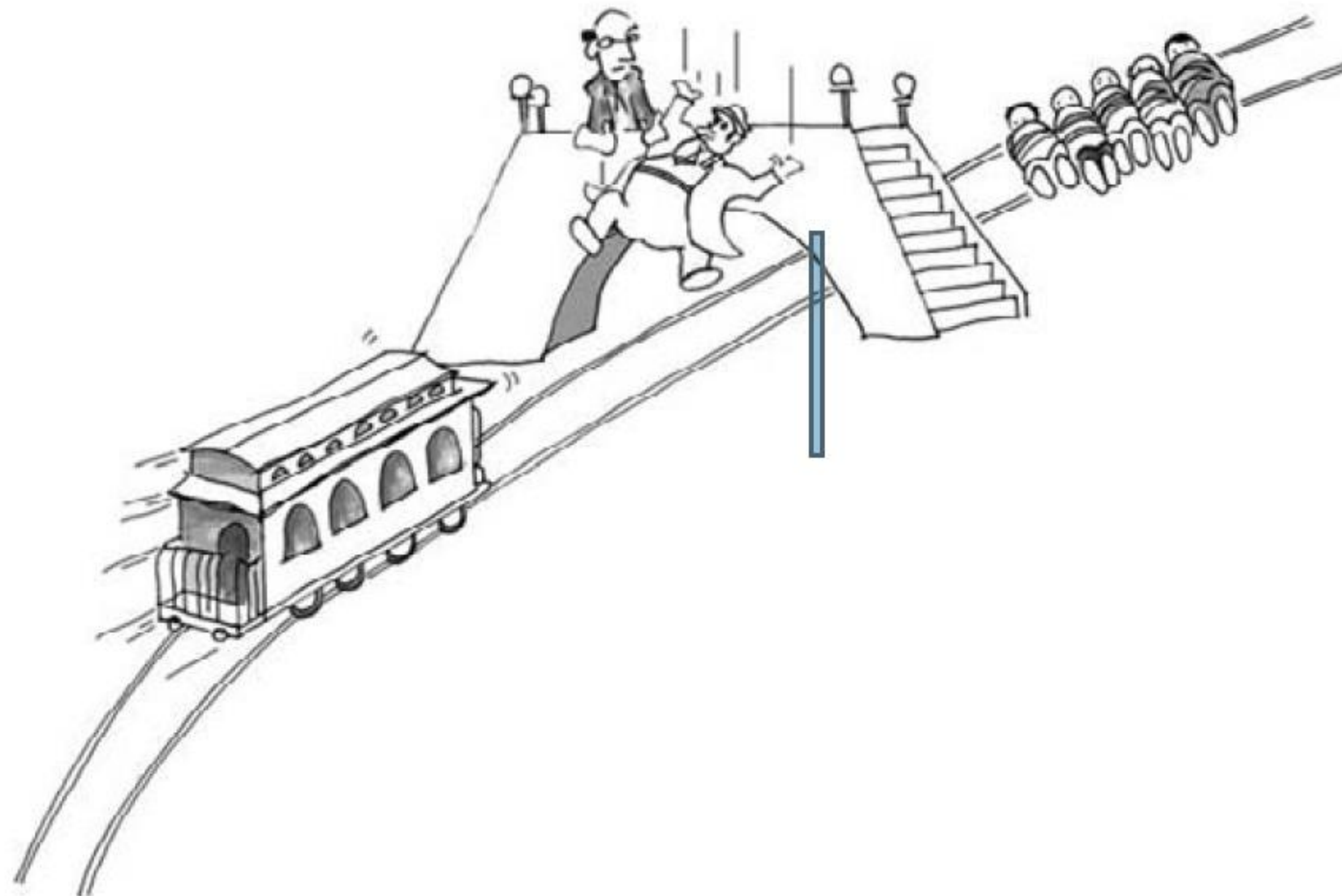


Figure 2. *Fat Man*. You're on a footbridge overlooking the railway track. You see the trolley hurtling along the track and, ahead of it, five people tied to the rails. Can these five be saved? Again, the moral philosopher has cunningly arranged matters so that they can be. There's a very fat man leaning over the railing watching the trolley. If you were to push him over the footbridge, he would tumble down and smash on to the track below. He's so obese that his bulk would bring the trolley to a shuddering halt. Sadly, the process would kill the fat man. But it would save the other five. Should you push the fat man?

D. Edmonds, *Would you kill the fat man? The trolley problem and what your answer tells us about right and wrong*, Princeton University Press, 2014.

The Doctrine of Double Effect

The DDE

The DDE can be given a more precise formulation. It's usually seen as consisting of four components, though this formulation is not universally accepted. The DDE comes into play when:

- the act considered independently of its harmful effects is not in itself wrong;
- the agent intends the good and does not intend the harm either as means or end, though the individual may foresee the harm;
- there is no way to achieve the good without causing the harmful effects; and
- the harmful effects are not disproportionately large relative to the good being sought.

The justifiability of targeting a particular military installation illustrates how the DDE can be applied. If it is legitimate to hit an installation with foreseen collateral damage then, according to the DDE, the following conditions must be met: (1) Hitting this installation must not in itself be wrong. (2) Hitting the installation must be the intended act, and the collateral damage must not be intended. (3) It must be impossible to hit the military installation without causing the collateral damage. (4) The badness of the collateral damage must not be disproportionate to the good that will result from hitting the installation.

D. Edmonds, *Would you kill the fat man? The trolley problem and what your answer tells us about right and wrong*, Princeton University Press, 2014

Ethics in autonomous cars?

“let me offer a simple scenario that illustrates the need for ethics in autonomous cars. Imagine in some distant future, your autonomous car encounters this terrible choice: it must either swerve left and strike an eight-year old girl, or swerve right and strike an 80-year old grandmother. Given the car’s velocity, either victim would surely be killed on impact. If you do not swerve, both victims will be struck and killed; so there is good reason to think that you ought to swerve one way or another. But what would be the ethically correct decision? If you were programming the self-driving car, how would you instruct it to behave if it ever encountered such a case, as rare as it may be?”

P. Lin, *Why Ethics Matters for Autonomous Cars*, in *Autonomous Driving, Technical, Legal and Social Aspects*, ed. by M. Maurer, J. Gerdes, B. Lenz, H. Winner, Berlin/Heidelberg, Springer, 2016, pp. 69-85.

The “Moral Machine”

“The Moral Machine attracted worldwide attention, and allowed us to collect 39.61 million decisions from 233 countries, dependencies, or territories. In the main interface of the Moral Machine, users are shown unavoidable accident scenarios with two possible outcomes, depending on whether the autonomous vehicle swerves or stays on course. They then click on the outcome that they find preferable. Accident scenarios are generated by the Moral Machine following an exploration strategy that focuses on nine factors:

- sparing humans (versus pets),
- staying on course (versus swerving),
- sparing passengers (versus pedestrians),
- sparing more lives (versus fewer lives),
- sparing men (versus women),
- sparing the young (versus the elderly),
- sparing pedestrians who cross legally (versus jaywalking),
- sparing the fit (versus the less fit),
- and sparing those with higher social status (versus lower social status).”

E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, I. Rahwan, *The Moral Machine experiment*, in «Nature», 563 (7729), 2018.

What should the self-driving car do?

4 / 13

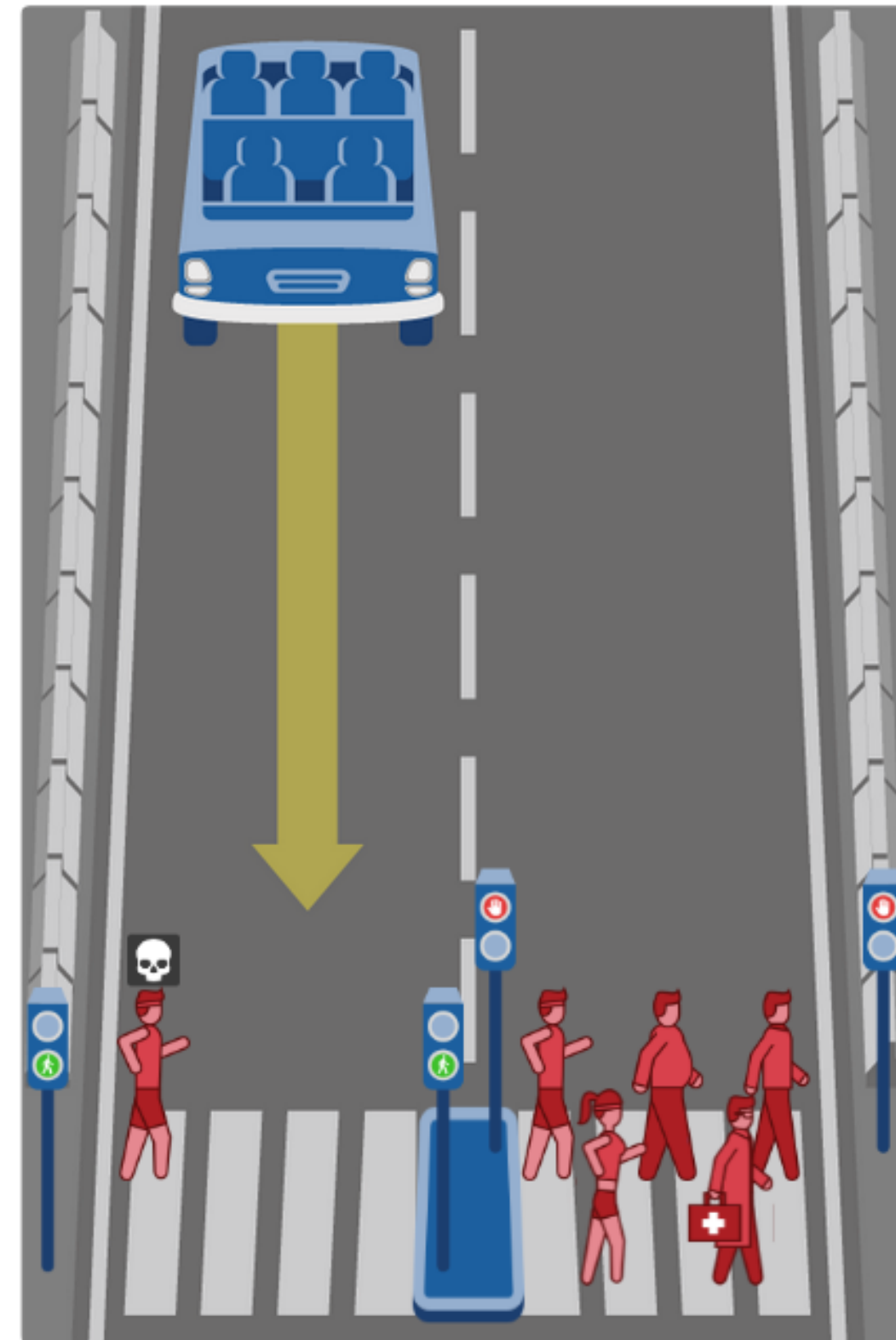
In this case, the self-driving car with sudden brake failure will continue ahead and drive through a pedestrian crossing ahead. This will result in ...

...

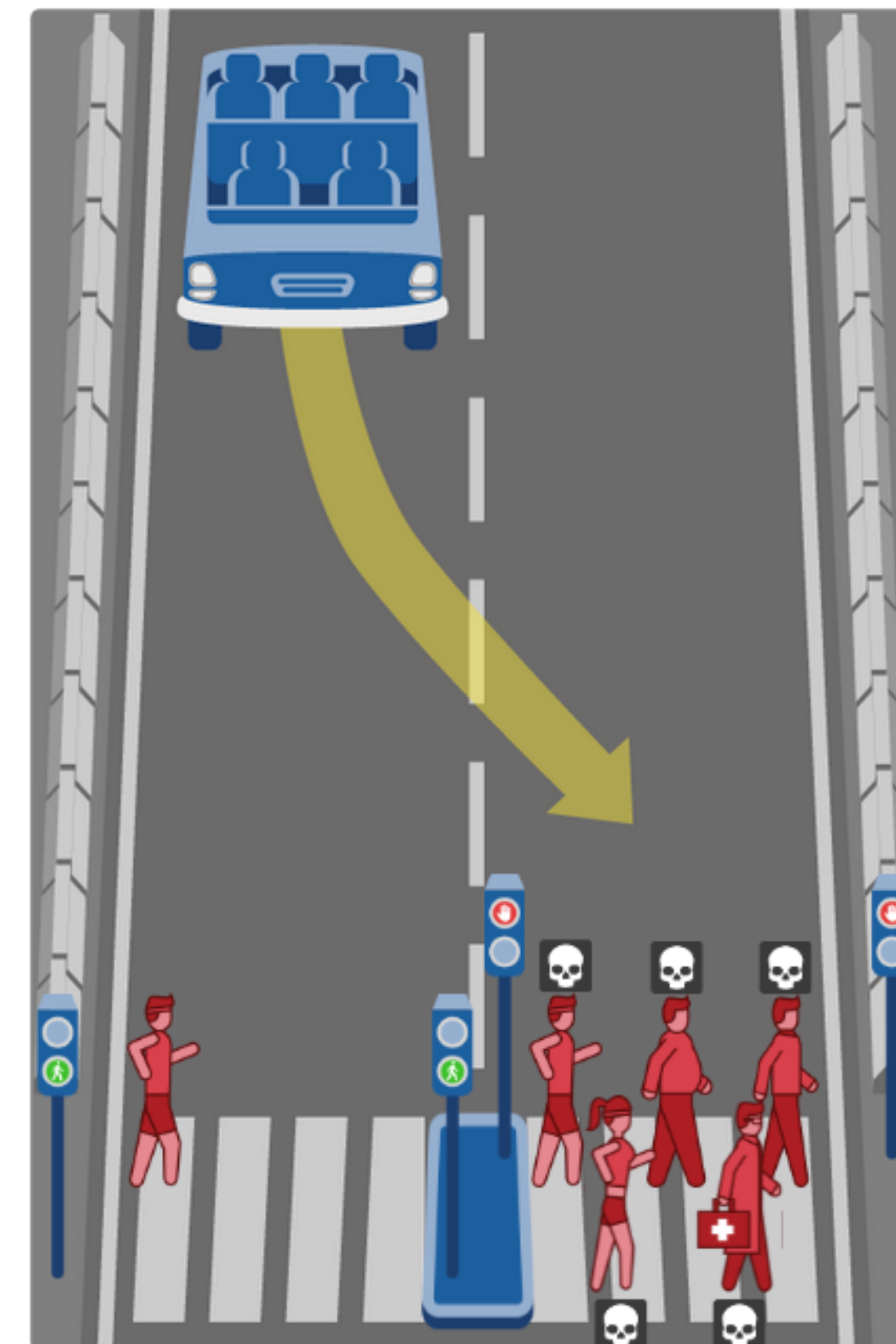
Dead:

- 1 male athlete

Note that the affected pedestrians are abiding by the law by crossing on the green signal.



Hide Description



Hide Description

In this case, the self-driving car with sudden brake failure will swerve and drive through a pedestrian crossing in the other lane. This will result in ...

Dead:

- 1 male athlete
- 1 large man
- 1 man
- 1 female athlete
- 1 male doctor

Note that the affected pedestrians are flouting the law by crossing on the red signal.

What should the self-driving car do?

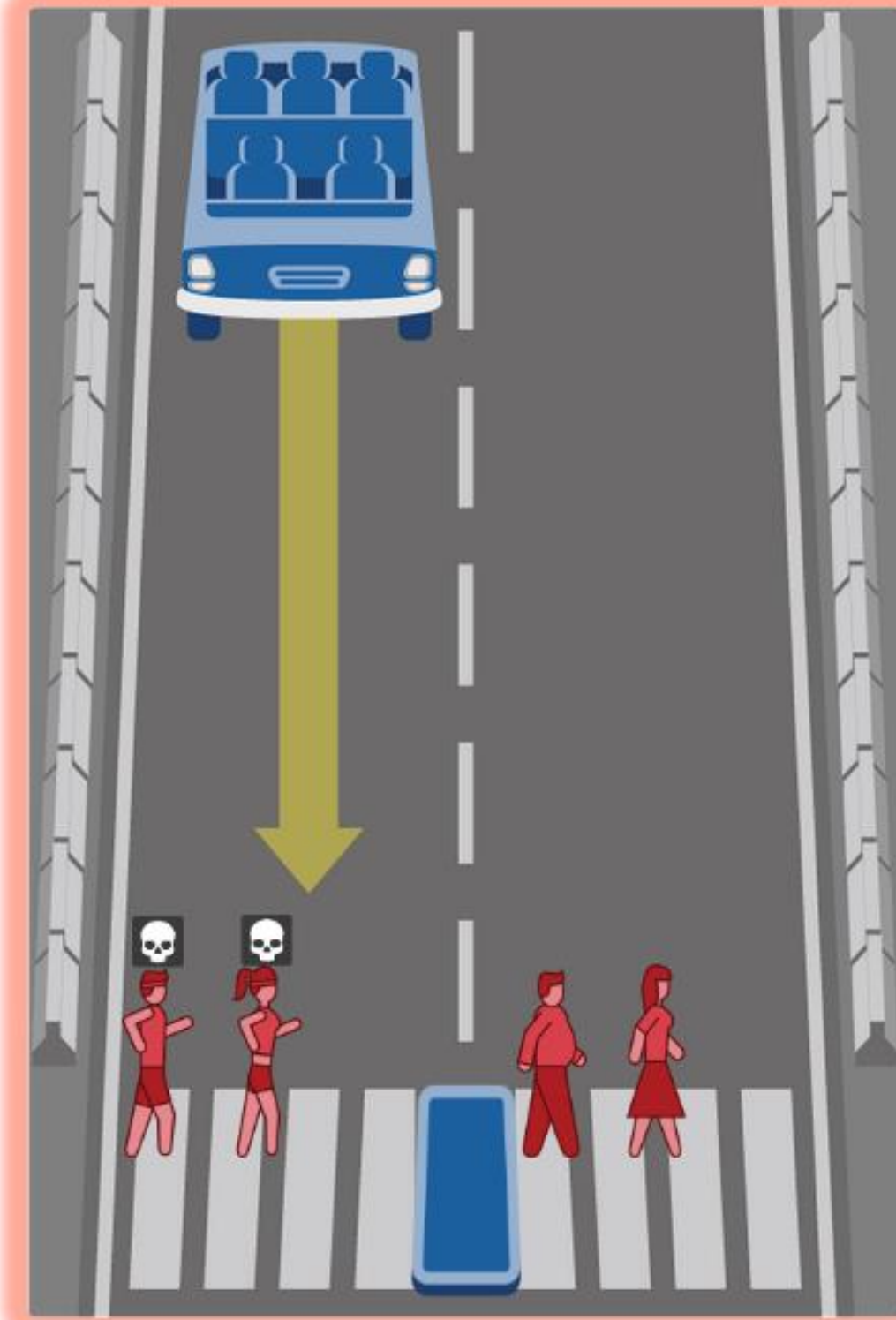
1 / 13

In this case, the self-driving car with sudden brake failure will continue ahead and drive through a pedestrian crossing ahead. This will result in ...

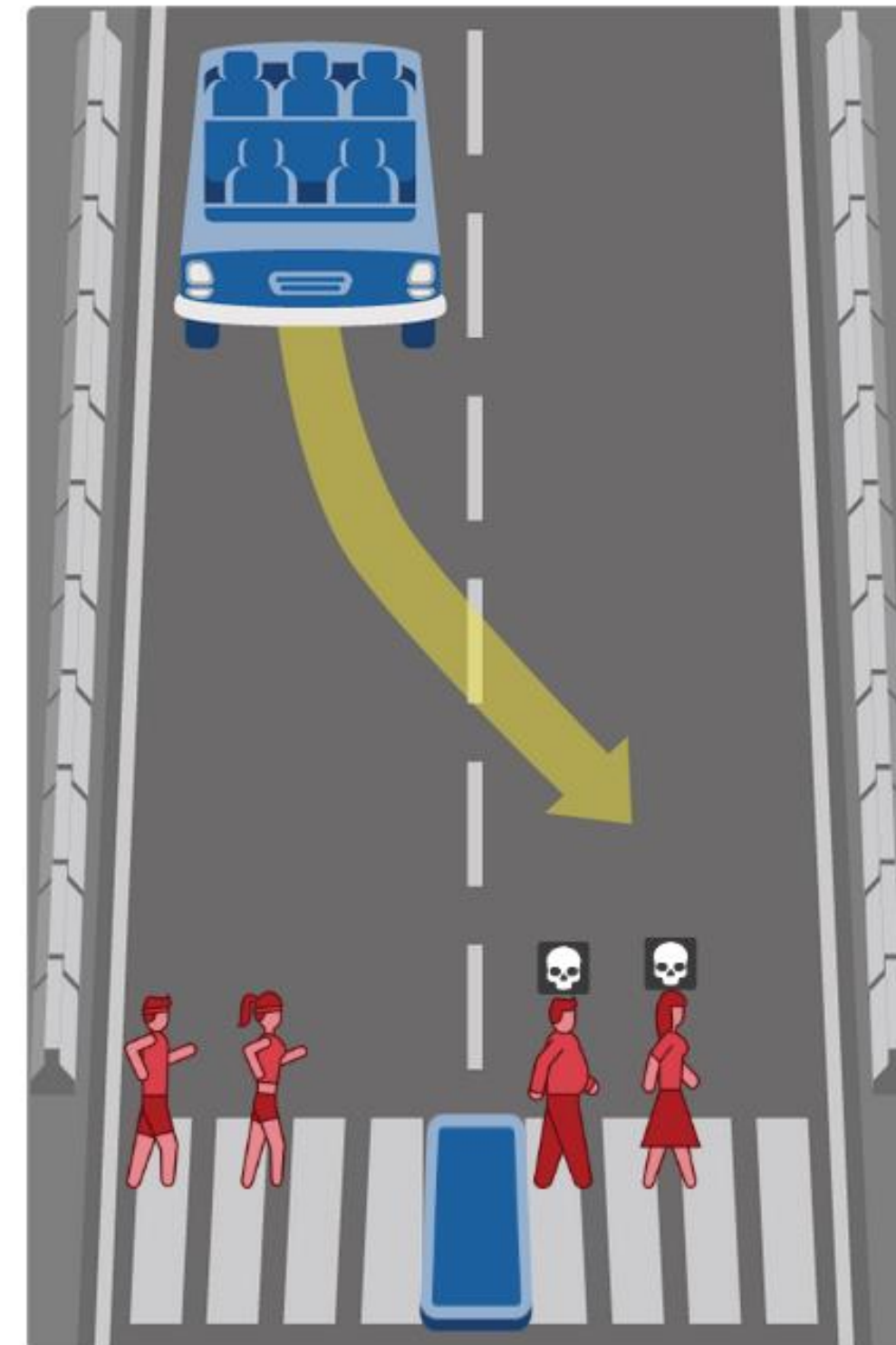
...

Dead:

- 1 male athlete
- 1 female athlete



Hide Description



Hide Description

In this case, the self-driving car with sudden brake failure will swerve and drive through a pedestrian crossing in the other lane. This will result in ...

Dead:

- 1 large man
- 1 woman

What should the self-driving car do?

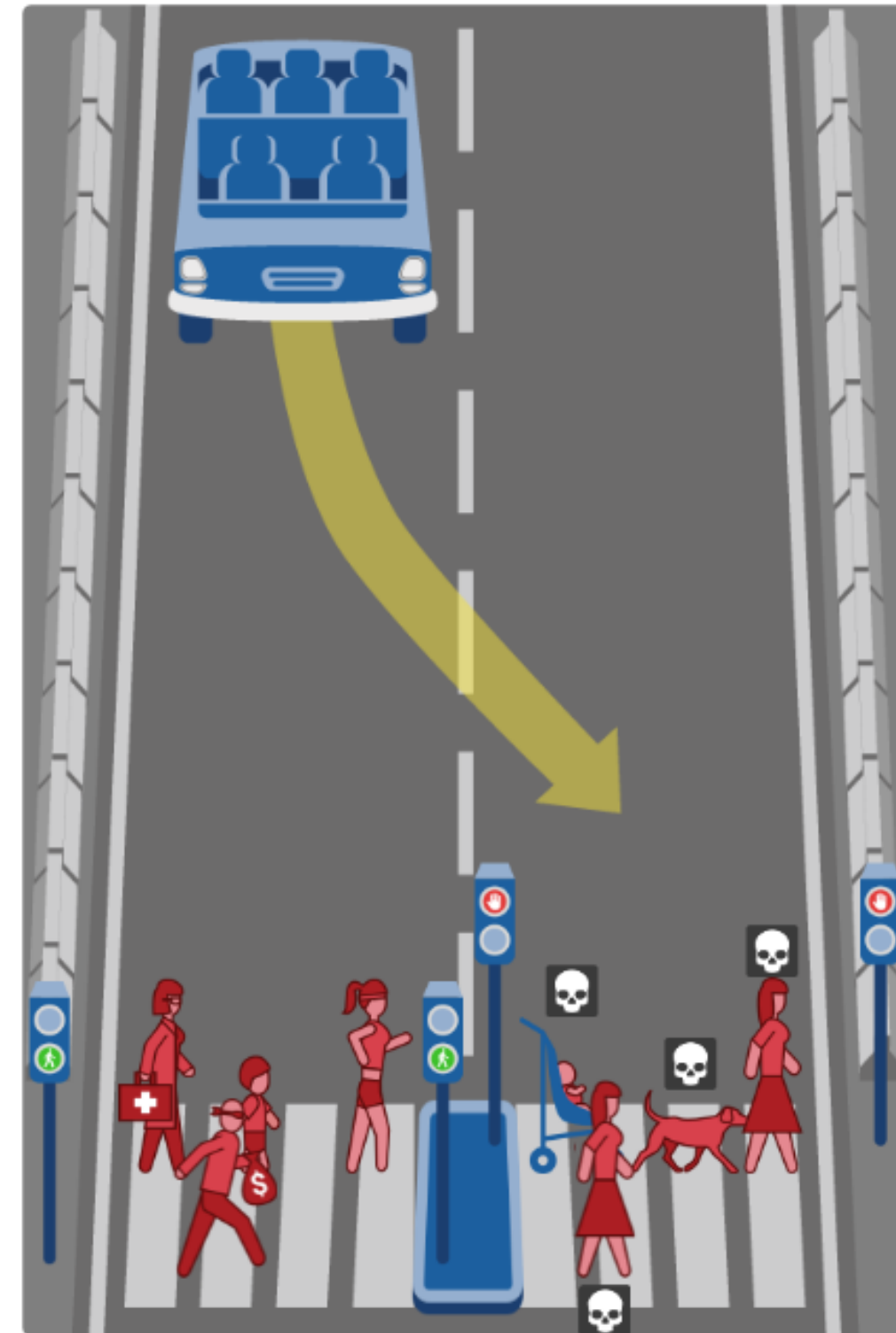
5 / 13

In this case, the self-driving car with sudden brake failure will swerve and drive through a pedestrian crossing in the other lane. This will result in ...

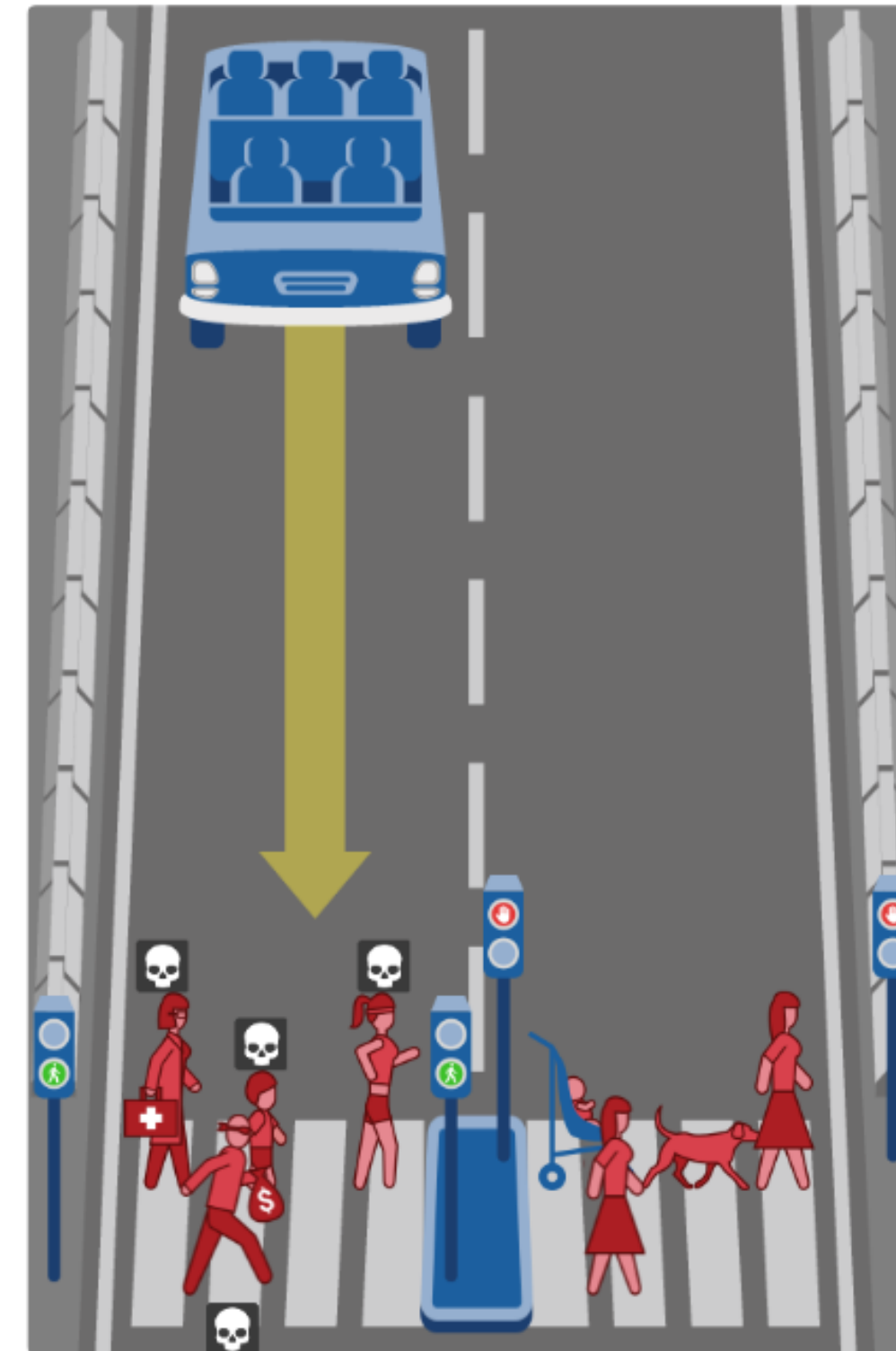
Dead:

- 1 baby
- 1 dog
- 2 women

Note that the affected pedestrians are flouting the law by crossing on the red signal.



Hide Description



Hide Description

In this case, the self-driving car with sudden brake failure will continue ahead and drive through a pedestrian crossing ahead. This will result in ...

Dead:

- 1 female doctor
- 1 boy
- 1 female athlete
- 1 criminal

Note that the affected pedestrians are abiding by the law by crossing on the green signal.

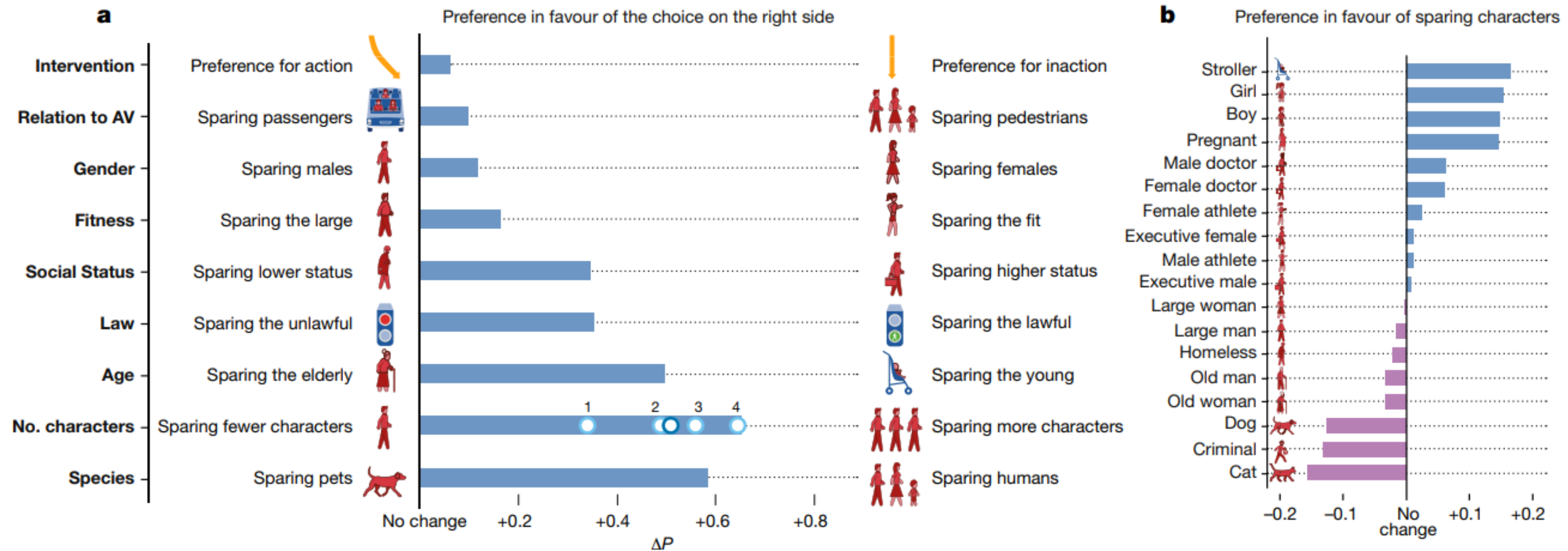


Fig. 2 | Global preferences. a, AMCE for each preference. In each row, ΔP is the difference between the probability of sparing characters possessing the attribute on the right, and the probability of sparing characters possessing the attribute on the left, aggregated over all other attributes. For example, for the attribute age, the probability of sparing young characters is 0.49 (s.e. = 0.0008) greater than the probability of sparing older characters. The 95% confidence intervals of the means are omitted owing to their insignificant width, given the sample size ($n = 35.2$ million). For the number of characters (No. characters), effect sizes are shown

for each number of additional characters (1 to 4; $n_1 = 1.52$ million, $n_2 = 1.52$ million, $n_3 = 1.52$ million, $n_4 = 1.53$ million); the effect size for two additional characters overlaps with the mean effect of the attribute. AV, autonomous vehicle. **b**, Relative advantage or penalty for each character, compared to an adult man or woman. For each character, ΔP is the difference between the probability of sparing this character (when presented alone) and the probability of sparing one adult man or woman ($n = 1$ million). For example, the probability of sparing a girl is 0.15 (s.e. = 0.003) higher than the probability of sparing an adult man or woman.

E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, I. Rahwan, *The Moral Machine experiment*, in «Nature», 563 (7729), 2018.



“Crowdsourcing moral machines”

» key insights

- **Machines are assuming new roles in which they will make autonomous decisions that influence our lives. In order to avoid societal pushback that would slow the adoption of beneficial technologies, we must sort out the ethics of these decisions.**
- **Behavioral surveys and experiments can play an important role in identifying citizens' expectations about the ethics of machines, but they raise numerous concerns that we illustrate with the ethics of driverless cars and the Moral Machine experiment.**
- **Data collected shows discrepancies between the preferences of the public, the experts, and citizens of different countries—calling for an interdisciplinary framework for the regulation of moral machines.**

Common criticisms and responses regarding the crowdsourcing of AV ethics using the Trolley Problem method.

Too Naïve	Laypersons' responses to public polls can be biased or ill-informed. Ethical trade-offs must be solved by policy experts, not majority voting.	Policymakers must know about the values most important to the public, so they can either accommodate these values, or anticipate frictions that need be explained.
Too Simple	Real accidents do not involve only two possible actions, and these actions do not have deterministic outcomes.	Highly complex scenarios would only allow for highly specific conclusions. Simplified scenarios zero in on the general principles that guide citizens' ethical intuitions.
Too Improbable	AV-Trolleys are based on very implausible sets of assumptions, and their actual probability of occurrence is too small to deserve attention.	Edge cases can have a massive impact on public opinion, and AV-Trolleys are the discrete form of a very real statistical problem.
Too Early	AV-Trolleys regulations should be avoided at this early technological stage, because their consequences are hard to predict.	Even though it may be too early to regulate about AV-Trolleys, it is the right time to start crowdsourcing citizen preferences.
Too Disconnected	Stated preferences are too disconnected from real actions	The behavior of human drivers is irrelevant to the proposed crowdsourcing task.
Too Distracting	Car makers should focus on making AVs safer, instead of wasting time and resources on crowdsourcing ethical dilemmas.	True, and this is why we need computational social scientist to handle that task.
Too Scary	Overexposing people to AV-Trolleys may scare them away, and be detrimental for their trust in the technology.	This is an empirical question, and our surveys did not find any evidence for such an adverse effect.

E. Awad, S. Dsouza, J.-F. Bonnefon, A. Shariff, I. Rahwan, *Crowdsourcing Moral Machines*, in «Communications of the ACM», 63,3, 2020.

Why self-driving cars do not need to choose whom to kill

“Recommendation

Manage dilemmas by principles of risk distribution and shared ethical principles.

While it may be impossible to regulate the exact behaviour of CAVs [Connected and Automated Vehicles] in unavoidable crash situations, CAV behaviour may be considered ethical in these situations provided it emerges organically from a.

continuous statistical distribution of risk by the CAV in the pursuit of improved road safety and equality between categories of road users

Rather than defining the desired outcome of every possible dilemma, it considers that the behaviour of a CAV in a dilemma situation is by default acceptable if the CAV has, during the full sequence that led to the crash, complied with all the major ethical and legal principles stated in this report, with the principles of risk management [...] and if there were no reasonable and practicable preceding actions that would have prevented the emergence of the dilemma. This may be necessary in order to give manufacturers and deployers of CAVs the confidence to deploy their systems, with reduced speed and preventative manoeuvres always being the best solution to decrease safety risks.”

Horizon 2020 Commission Expert Group to advise on specific ethical issues raised by driverless mobility (E03659). *Ethics of Connected and Automated Vehicles: recommendations on road safety, privacy, fairness, explainability and responsibility*, Luxembourg, Publication Office of the European Union, 2020.

“It may be ethically permissible for CAVs [Connected and Automated Vehicles] not to follow traffic rules whenever strict compliance with rules would be in conflict with some broader ethical principle. Noncompliance may sometimes directly benefit the safety of CAV users or that of other road users, or protect other ethical basic interests; for example, a CAV mounting a kerb to facilitate passage of an emergency vehicle. This is a widely recognized principle in morality and in the law.”

“The pursuit of greater road safety may sometimes require non-compliance with traffic rules.

Researchers should study the extent to which it is reasonable to expect that an intelligent non-human system is able to engage in the complex process of evaluation of the interpretation of a legal, ethical or societal norm and its balancing with another norm, value or principle.”

Horizon 2020 Commission Expert Group to advise on specific ethical issues raised by driverless mobility (E03659). *Ethics of Connected and Automated Vehicles: recommendations on road safety, privacy, fairness, explainability and responsibility*, Luxembourg, Publication Office of the European Union, 2020.

“Engineers working on vehicle automation are often asked about the trolley problem. The most common response seems to be that **trolley problems are avoidable, implausible, rare, and distractions from more productive efforts.** They are considered avoidable because in many trolley problems, the automated vehicle must decide how best to crash when, with the right sensors and algorithms, the situation should have been avoided entirely. An advanced automated vehicle would have slowed down before that blind turn, seen that animal before it leaped into the road, or known this neighborhood has young children and adjusted its speed accordingly. Developers find trolley problems implausible and rare for other reasons. Most people have trouble remembering a situation in which they had time to decide which way they should crash. Because of how these forced-choice scenarios are presented in the media and literature, they are easy to mock by focusing on some of the more outlandish examples, like a car colliding with a criminal instead of (specifically and consistently) a nun. **Focusing resources on unlikely edge cases seems like a waste of resources that could be better spent on general collision avoidance.**”

N.J. Goodall, *More than Trolleys: Plausible, Ethically Ambiguous Scenarios likely to Be Encountered by Automated Vehicles*, in «Transfers: Interdisciplinary Journal of Mobility Studies», 9, 2, 2019, pp. 45–58.

- “In situations where a self-driving car must choose between straight-line braking into an unavoidable collision and swerving into an unavoidable collision, where there are no other cars involved, the car should always prefer the straight-line option. Additional information about the objects to be collided with is irrelevant, since there is no way for the car to gather that information without making the risks of the situation worse.”
- “even in the much more complex situations produced by involving more vehicles, **an emergency stop policy is at least good enough to be worth considering**”.

R. Davnall, *Solving the Single-Vehicle Self-Driving Car Trolley Problem Using Risk Theory and Vehicle Dynamics*, in «Science and Engineering Ethics», 26, 2020, pp. 431-449, <https://link.springer.com/content/pdf/10.1007/s11948-019-00102-6.pdf>.

Do self-driving cars require Artificial General Intelligence?

Drivers Sue Tesla Over Alleged Failure to Deliver on Promises of Self-Driving Cars

by Erin Shaak

Last Updated on September 19, 2022

Two proposed class action lawsuits filed this week claim that Tesla has for years deceptively misrepresented the “autonomous” driver assistance technology in its electric vehicles and essentially led drivers on with promises that a fully self-driving vehicle is “on the cusp” of being brought market.

According to the two cases, filed on September 14th and 15th, many Tesla drivers have paid thousands extra for Tesla’s advanced driver assistance systems (ADAS)—including “Autopilot,” “Enhanced Autopilot” and “Full Self-Driving Capability”—based on the automaker’s misrepresentations of both how the systems work and the availability of even more advanced technology in the near future.

In fact, the suits allege that Tesla and CEO Elon Musk have repeatedly indicated since as early as 2016 that the company is within a year, or even

a few months, of perfecting a self-driving car, with Musk reportedly stating in a 2016 tweet that a Tesla would be able to complete a fully self-driven cross-country trip by “next year.”

According to the lawsuits, however, these promises have “proven false time and time again.”

Six years later, Tesla “has yet to produce anything even remotely approaching a fully self-driving car,” one case argues. The suit says former employees and investigations alike have revealed “damning information” that indicates Tesla has never even come close to achieving that goal.

The lawsuits contend that Tesla and Musk made these misleading statements about the cars’ self-driving capabilities despite being fully aware that “there was no reasonable chance” that Tesla would be able to follow through on those promises.

<https://www.classaction.org/blog/drivers-sue-tesla-over-alleged-failure-to-deliver-on-promises-of-self-driving-cars>

Consumer Skepticism Toward Autonomous Driving Features Justified

Third Time is Not a Charm as Driving Assistance Tech Continues to Underperform

Consumers surveyed told AAA they are more interested in improved vehicle safety systems (77%) versus self-driving cars (18%). But new testing, the third round by AAA's Automotive Engineering team in the last few years, found that vehicles with an active driving assistance system (also known as Level 2 systems as [defined by SAE](#)) failed to consistently avoid crashes with another car or bicycle during 15 test runs. A foam car similar to a small hatchback and a bicyclist dummy was used for this testing.

- A head-on collision occurred during all 15 test runs for an oncoming vehicle within the travel lane. Only one test vehicle significantly reduced speed before a crash on each run.
- For a slow lead vehicle moving in the same direction in the lane ahead, no collisions occurred among 15 test runs.
- For a cyclist crossing the travel lane of the test vehicle, a collision occurred for 5 out of 15 test runs, or 33% of the time.
- For a cyclist traveling in the same direction in the lane ahead of the test vehicle, no collisions occurred among 15 test runs.



<https://newsroom.aaa.com/2022/05/consumer-skepticism-toward-active-driving-features-justified/>

https://newsroom.aaa.com/wp-content/uploads/2022/05/E-1_Research-Report_2021-ADA-Evaluation_FINAL_4-13-22.pdf

NHTSA Finds Teslas Deactivated Autopilot Seconds Before Crashes

The finding is raising more questions than answers, but don't jump to any conclusions yet.

Alexander Stoklosa - Writer, Getty Images - Photographer | Jun 15, 2022



SEE ALL 2 PHOTOS

MOTORTREND
THE LATEST IN CAR NEWS
EMAIL NEWSLETTER
SIGN-UP!

Your email address

By clicking Sign Up, you agree to the [Terms of Use](#).
Your information will be collected and used in accordance with our [Privacy Policy](#).

A NHTSA report on its investigation into crashes in which Tesla vehicles equipped with the automaker's Autopilot driver assistance feature hit stationary emergency vehicles has unearthed a troubling detail: In 16 of those crashes, "on average," Autopilot was running but "aborted vehicle control less than one second prior to the first impact."

Tesla's self-driving technology fails to detect children in the road, group claims

Safe technology campaigners release 'disturbing' video advert showing car in Full Self-Driving mode hitting child-sized mannequin



A Tesla Model 3 fitted with a full self-driving system. Photograph: Sjoerd van der Wal/Getty Images

<https://www.motortrend.com/news/nhtsa-tesla-autopilot-investigation-shutoff-crash/>

<https://www.theguardian.com/technology/2022/aug/09/tesla-self-driving-technology-safety-children>

Why it is immoral (and illegal) to apply the trolley problem to self-driving cars

The Moral Machine as a test on people's biases

“Players are forced to choose between swerving to kill a homeless person, a criminal, and a man (a) or going straight to kill two women and a female executive (b). This kind of information is unacceptable to use in making moral decisions.”

“By using social properties as their criteria for moral decision making, **this experiment is mistakenly testing people's discriminatory biases rather than their moral judgments.**”

Imagine that the game included descriptions of race, religion, and sexual orientation. The MIT researchers don't want to ask: “Are you willing to sacrifice the lives of three gay women to save a Muslim?” But what they're doing in asking about class and occupation is essentially the same thing. Any student who's taken an introductory ethics class understands why this game is not only misguided but dangerous.”

D. Leben, *Ethics for Robots. How to Design a Moral Algorithm*, London/New York, Routledge, 2019.



“No selection of humans, no offsetting of victims, but principle of damage minimization

The modern constitutional state only opts for absolute prohibitions in borderline cases, such as the ban on torture relating to persons in state custody.

Regardless of the consequences, an act is mandated or prohibited absolutely because it is intrinsically already incompatible with the constitutive values of the constitutional order. Here, there is, exceptionally, no trade-off, which is per se a feature of any morally based legal regime.

The Federal Constitutional Court's judgment on the Aviation Security Act 5 also follows this ethical line of appraisal, with the verdict that **the sacrifice of innocent people in favour of other potential victims is impermissible, because the innocent parties would be degraded to mere instrument and deprived of the quality as a subject.**”

<https://www.bmvi.de/SharedDocs/EN/publications/report-ethics-commission.html>



“So act that you use humanity, whether in your own person or in the person of any other, always at the same time as an end, never merely as a means”

I. Kant, *Grundlegung zur Metaphysik der Sitten*, 1785, in *Kant's gesammelte Schriften. Akademie-Ausgabe*, Berlin, W. de Gruyter, 1900, IV, pp. 385-463; in *Idem, Practical Philosophy, The Cambridge Edition of the Works of Immanuel Kant*, ed. by M. Gregor, Cambridge, Cambridge University Press, 1996.



“Genuine dilemmatic decisions, such as a decision between one human life and another, depend on the actual specific situation, incorporating “unpredictable” behaviour by parties affected. They can thus not be clearly standardized, nor can they be programmed such that they are ethically unquestionable.

In the event of unavoidable accident situations, any distinction based on personal features (age, gender, physical or mental constitution) is strictly prohibited.

It is also prohibited to offset victims against one another. General programming to reduce the number of personal injuries may be justifiable.

Those parties involved in the generation of mobility risks must not sacrifice non-involved parties.”

<https://www.bmvi.de/SharedDocs/EN/publications/report-ethics-commission.html>



Thank you. Any questions?

daniela.tafani@unibo.it

MAI4CAREU

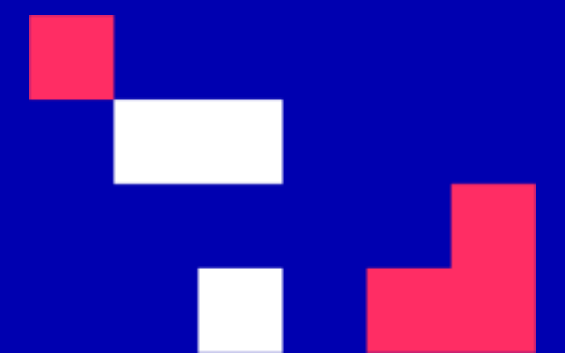
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester

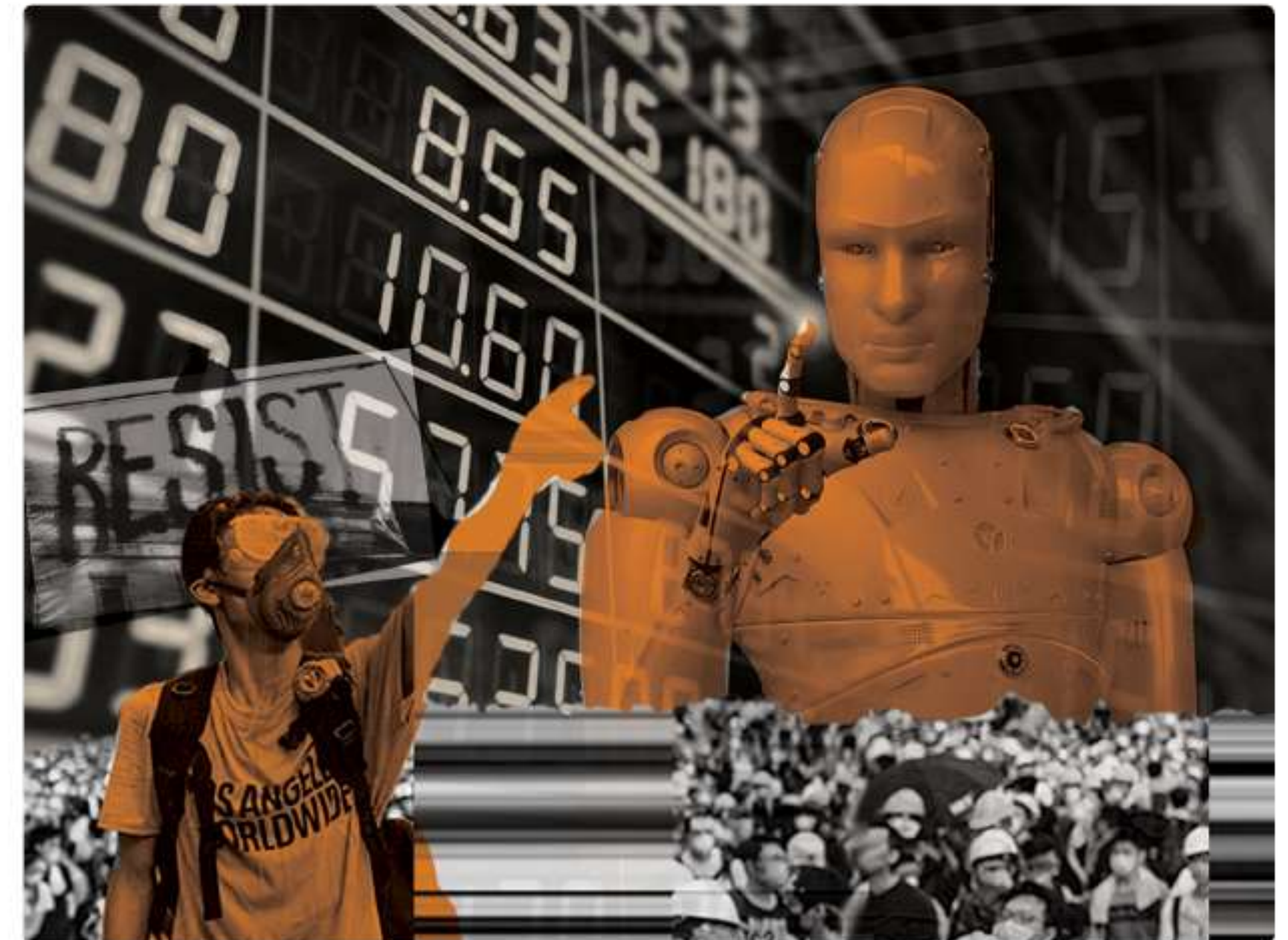


Citation

Amy J. Ko, Anne Beitlers, Brett Wortzman, Matt Davidson, Alannah Oleson, Mara Kirdani-Ryan, Stefania Druga, Jayne Everson (2022). *Critically Conscious Computing: Methods for Secondary Education*. <https://criticallyconsciouscomputing.org/>, retrieved 9/21/2022.

License

Creative Commons Attribution-NoDerivatives 4.0



Ashley Wang CC0

Depending on who you are, AI can be magical or horrifying.

Chapter 15 Code

Artificial Intelligence

by Amy J. Ko, Stefania Druga

The power and perils of AI

The limits of AI

Harms of exploitation

Harms of allocation

Harms of representation

Harms of prediction

Managing and Regulating AI

Key ideas

- * Artificial intelligence (AI) is concerned with replicating human intelligence and ability, and is increasingly used in software applications, consumer devices, and every major industry to automate and inform decisions.
- * AI has made many advances in replicating specific human abilities, primarily using large data sets created through human labor, and using that data to build machine-learned classifiers.
- * AI perpetuates whatever values and biases were encoded in its algorithms and data, and depending on how they are used, create new systems of oppression.



Ashley Wang-CCO

AI can do harm, often unintentionally

The power and perils of AI

The innovations above can be quite impressive – how can we resist the futuristic wonders of smart assistants, autonomous cars, and acrobatic robots? And with AI increasingly embodied, with voices and physical forms designed to mimic human behavior and leverage human social cues, we increasingly see AI as more than just data and algorithms: we *like* them,

and come to see them as human-like, even though they are just data and algorithms²⁸. And yet, the problematic applications of machine learning we portrayed in the previous sections make clear that this power is not necessarily good. In this section, we discuss the limits of this power, and some of the many social consequences of misapplying it.

The limits of AI

One of the most salient findings from AI research since the 1950's is that strong AI is far out of reach⁶. The closest that researchers have come is to create programs that can fool humans into believing they are intelligent for a brief time. Alan Turing described what is now known as the “Turing test” or the “imitation game”, in which a human writes text messages to an AI and receives text replies, having a conversation. If the human is fooled for some period of time, one might judge the AI as having achieved some level of mimicry of intelligence. Many have taken the test literally, creating *competitions* in which people write AIs and compete to see who can trick the largest number of people, or *prizes* to advance the capabilities of AI in consumer contexts.

But even the best efforts in these competitions and in research reveal the fragile seams in strong AI attempts. All of the techniques above capture *some* essence of human ability or patterns in society, but this is often only true in perfectly calibrated circumstances, and for a very narrow range of human abilities. Applications of AI therefore continue to be *weak* – carefully constructed for particular contexts, and brittle when applied outside those contexts, just like any other kind of computer program.

⁶ Adriana Braga & Robert K. Logan (2016). *The Emperor of Strong AI Has No Clothes: Limits to Artificial Intelligence*. *Information*

²⁸ Byron Reeves & Clifford Nass (1996). *The Media Equation: How People Treat Computers, Television, and New Media Like Real People and Places*. Cambridge University Press

Harms of exploitation

As weak AI has become ubiquitous in industry applications, primarily through machine learning, data has become a precious commodity: computers can only recognize and produce human speech, label objects in images, and play challenging games like go and chess because companies have amassed large amounts of labeled training data. Many of these data sets are generated entirely through human labor. Sometimes the data is gathered from commercial services without our awareness, such as the way that search engines passively record logs of our queries, labeling the results we select, and using that data set with machine learning algorithms to improve search results. Other data is gathered with human consent, such as when a website login page asks you to select parts of an image that contain traffic lights or stop signs, which are then used to train driverless cars. Some data comes from paid labor, such as the image labels that Facebook gathers from its tens of thousands of low-wage, low-status content moderation staff in developing countries³⁰. Many of the AI-fueled conveniences of expensive smart devices are thus built on the backs of no or low wage people, often without consent, and that data is

often used against them to surveil and police them²⁶.

The need for data also links to the issues of privacy and surveillance capitalism we discussed in [Encoding Information](#). It's not just that companies need data in order to train AI to offer new conveniences, it's that they need data to make money, as their primary revenue streams are from selling advertisements³⁶. Data, and more fundamentally the invasion of privacy, often without consent, is therefore foundational to the economic systems that drive the creation of AI.

²⁶ Catherine O'Neil (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown

³⁶ Shoshana Zuboff (2019). *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. PublicAffairs

³⁰ Sarah Roberts (2014). *Behind the Screen: the Hidden Digital Labor of Commercial Content Moderation*. University of Illinois at Urbana Champaign

Harms of allocation

All algorithms, when inserted into decisions about who gets a particular resource, and who does not, can disproportionately harm some groups and not others³. AI, however, has been found to be particularly systematic in its discrimination. For example, consider Black Americans, who because they are more likely in poverty, are increasingly denied access to housing, food, loans, and insurance due to applications of AI that are racially biased in their risk predictions¹³. These harms of allocation stem from multiple levels: the algorithms, which prioritize accuracy at the aggregate level, without considering impacts on marginalized groups; data that often underrepresented marginalized experiences; and applications of AI that stem from political goals of reducing fraud rather than helping people in need.

³ Ruha Benjamin (2019). *Race after technology: Abolitionist tools for the new jim code*. John Wiley & Sons

¹³ Virginia Eubanks (2018). *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press

In contrast, other applications of AI often enrich and liberate groups that already have power and wealth. For example, the realities of weak AI such as machine learning that aim to make our lives easier and safer – often for the benefit of the wealthy, and at the expense of marginalized groups. Algorithms that optimize the availability of rideshare services, for example, only benefit those who can afford expensive rideshare trips, and at the expense of drivers who often get paid less than minimum wage after accounting for expenses. AI, then, just like any other code, is often deployed as a tool of wealthy, dominant groups to accrue power, increase wealth, and maintain the matrix of oppression that erases diversity, denies equity, and shuns equality⁴.

⁴ Abeba Birhane, et al. (2021). *The Values Encoded in Machine Learning Research*. arXiv

Harms of representation

While harms of exploitation and allocation are more direct, AI also causes more subtle, deeper harms in the form of misrepresentation, stereotyping individuals in harmful ways. Such harms are often quite simple. For example, machine learning that uses a binary gender feature – male or female – works fine for the cis majority, but erases people with non-binary or gender non-confirming identities. Such data, used by companies like Facebook to determine AI-based ad recommendations, excludes countless people from participating in commerce. Or, consider machine learning that used in face detection algorithms by police: because such machine learning is less effective for Black faces, there is been a surge in false arrests of Black people, spending time in jail, losing jobs and losing money on lawyers to prove their innocence.

Some of these harms stem from a lack of recognition of human, social, and cultural diversity. For example, what is “common” in common sense AI is culturally and socially determined. It’s only obvious to a person familiar with basketball that the spherical object is a basketball; to children who have never seen the game, or people in cultures unfamiliar with the Western sport, such facts would not be obvious. Therefore, the values in common sense AI are highly sensitive to the data on which they are based. The same is broadly true of all AI: because most AI systems are made by English speaking people in Western cultures, and authored by wealthier people with access to computers and the Internet, AI often reflects assumptions about what AI is for, and who it is for.

Some harms, however, stem from AI algorithms themselves, all of which optimize for common cases over edge cases. Anyone part of a group that is less common than the majority is likely to be misrepresented by the logical rules or statistical patterns used to define AI behavior, necessarily resulting in worse accuracy for people in minority groups. Some of these harms stem from data used to train AI, where historical norms or trends – such as dated Western notions about binary gender – end up erasing individual identities. And some of these harms simply stem from how AI is applied in society: for example, it is possible to build less racially biased facial recognition technology, but it would still be used to disproportionately surveil and police Black people in the U.S.

In some ways, the very notion of categorization, classification, and labeling is harmful. Consider, for example, our applications of machine learning on tardiness earlier. It's tempting to view lateness as a binary construct: students are either late to class or they are not, right? Of course, as any teacher knows, lateness is not an inherently binary phenomenon,

edge case: Any consideration an algorithm needs to make that is not part of the common case, such as unexpected inputs, rare scenarios, or contexts that violate its assumptions about the world.

because of the many exceptional circumstances that might conspire to make lateness ambiguous. Daylight saving time, late buses, clocks that differ by more than a few minutes, a faulty passing period bell, an emergency like an earthquake, poorly accommodated physical disabilities out of a student's control – all of these circumstances and more erode the idea that lateness is a binary. And the same complexities arise when we try to define many natural and social phenomena, such as race, gender, poverty, wealth, and so on. Supervised learning algorithms that model some aspect of nature or society often must reduce these ideas to a smaller set of discrete or binary values in order to work – but they do so at the cost of accurately modeling diversity. And more so, the people that are erased by the reductive nature of classification are usually people at the margins, who have exceptional circumstances, identities, bodies, and behaviors. Grouping, labeling, classifying, categorizing – all of these are just synonyms of stereotyping, which erases the nuance of difference and diversity.

Harms of prediction

One of the most common applications of machine learning, and statistics in general, is to make predictions about the future. As anyone knows, a prediction is not necessarily going to be true. And yet, many applications of AI frame prediction as if it are trustworthy, fair, and accurate.

Consider, for example, the problem of **recidivism**, which is the tendency for people who are convicted of crimes to commit further crimes³⁵; research on recidivism suggests that repeat offenses are primarily caused by the fragile social supports and stigmas that previously incarcerated people face after release. When judges make sentencing decisions after a jury convicts them of a crime, they often make recidivism predictions, assessing how likely the convicted person is to commit another crime. If they decide the likelihood is high, they may give the person a longer sentence, hoping that separating a person from society will protect society from their crimes; if it is low, they may give a shorter sentence. There's obviously immense opportunity for bias in these human judgment, as the judge's prediction of recidivism might be swayed by the color of the convicted person's skin. And so some software companies saw an opportunity to create software for courthouses that would use machine learning to make these predictions instead, hoping to remove the

³⁵ Edward Zamble & Vernon L. Quinsey (2001). *The Criminal Recidivism Process*. Cambridge University Press

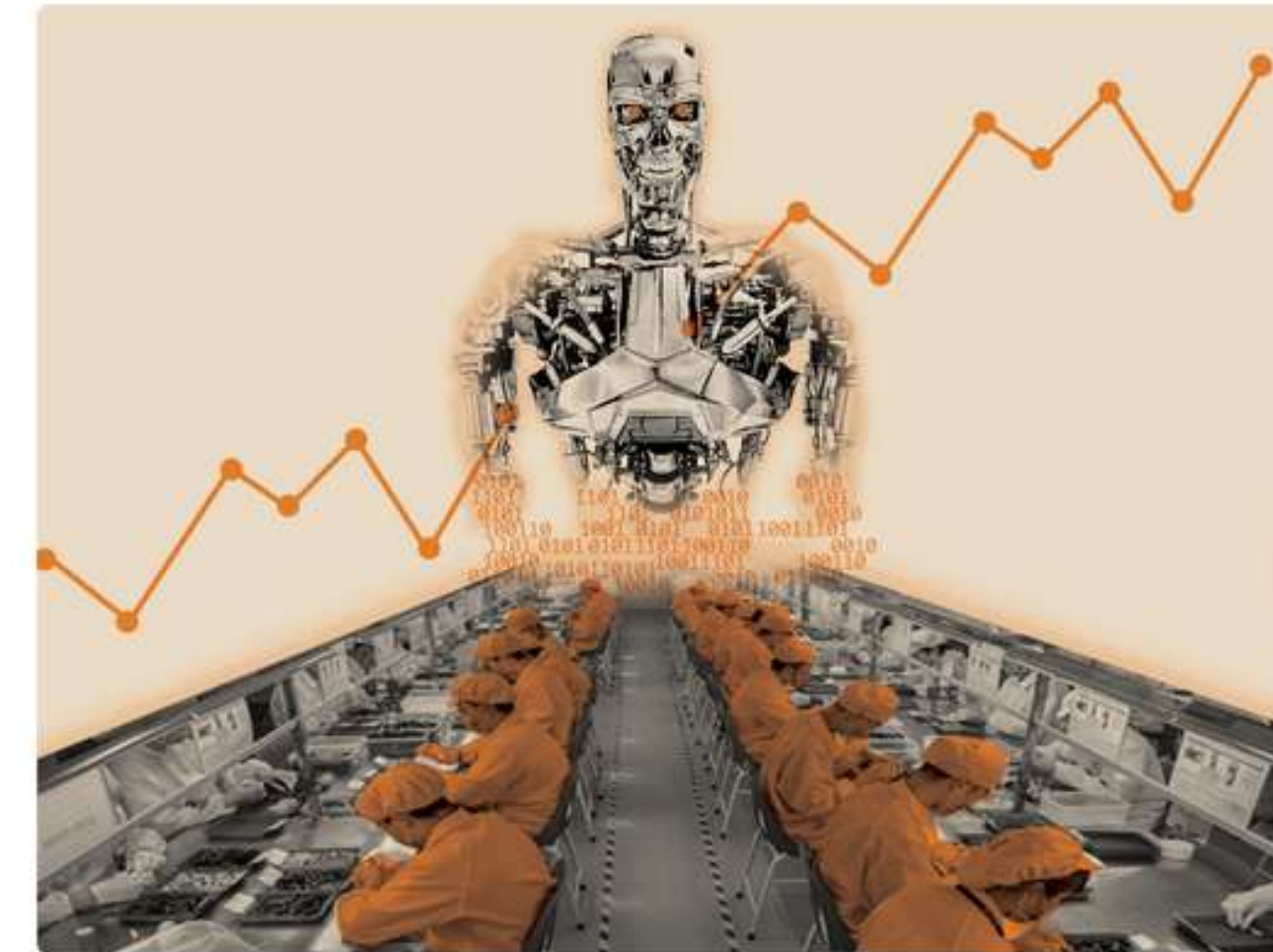
potential for a judge's bias by replacing their judgment with supposedly "neutral" AI.

The most popular software, **COMPAS** (Correctional Offender Management Profiling for Alternative Sanctions), launched in 1998, used the approach of taking a large data set of past convicted people, gathering 137 different factors about the accused (the degree of the crime, whether it was a drug offense, the offenders' age), but did not include race as a factor. Machine learning algorithms, as we will describe below, build statistical models that find how much to weigh different factors to make the most accurate prediction possible; they use data on past convictions to "train" a classifier, then use that trained classifier to make predictions about new cases. COMPAS did this same thing, using a historical record of convictions and re-convictions to make predictions about people newly convicted. Despite COMPAS's insistence that race was not included as factor, the factors that it did include were nevertheless racially biased: Black men are more likely to be arrested for crimes because they are surveilled and policed more than other identities; they are more likely to be convicted of crimes, even when they are later found to be not guilty. The machine learned classifier in COMPAS learned these racially-biased patterns, and reproduced them, without ever explicitly including race as a factor. In fact, research showed

that the COMPAS machine learned classifier was just as biased as lay people with little or no criminal justice expertise, and in fact less accurate and more biased than experienced judges⁹. And yet, because of a misplaced faith in the neutrality and objectivity of computers, judges across the United States have deferred to COMPAS and its racially biased predictions to make their sentencing decisions, rather than using their knowledge of the individual accused. The result is that judges using COMPAS give Black men even longer sentences than other racial groups.

The developers of COMPAS did not evaluate these outcomes before they built it; their priority was selling software to courthouses, and cities' priorities were reducing perceptions of bias (instead of actual bias, or more importantly, the larger systemic bias in policing, prisons, and poverty). No one actually investigated this bias until the public started demanding audits of these algorithms and researchers conducted these audits, which revealed their ineffectiveness and bias. This application of weak AI, therefore, prioritized certain people's concerns (judges, politicians, and the White majorities that elected them) at the expense of others (Black men).

⁹ Julia Dressel & Hany Farid (2018). The accuracy, fairness, and limits of predicting recidivism. *Science Advances*



Ashley Wang CCO

AI has bias, but it's not always clear what to do about it.

Managing and Regulating AI

Because AI, and computing in general, has such great potential for harm, there are numerous approaches to trying to mitigate its risks.

Transparency. One approach is to advocate for transparency – the degree to which people other than an AI's creators can see the data on which AI is based and the

algorithms it uses to process it¹. Without transparency, the public is limited to observing the patterns in AI behavior, without being able to see their underlying causes. Thus far, the world has been so enamored with the possibilities of AI that it has asked for little transparency – most private companies view their data and AI algorithms as valuable trade secrets – and only some governments have committed to transparent uses of AI. Just as when code is kept secret, when the data and AI algorithms used are kept secret, there is no way to build public trust in their underlying logic.

Explainability. Even when AI is transparent, the behavior of AI algorithms can be much harder to explain than a computer program²⁷. After all, a program has logic, written in a programming language, that someone familiar with that language can read, analyze, and make conclusions about with high confidence. Comprehending a program's behavior is by no means easy, but it is feasible with time and expertise. AI, and especially machine learning, in contrast, often have little explicit logic. Their behavior is inherently determined by the data selected, the statistical analyses applied to that data, and the emerging patterns in that data. The answer to “*Why was my loan application rejected?*” for a computer program might be “*The program inspected your responses to these three*

¹ Mike Ananny & Kate Crawford (2016). *Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability*. *New Media & Society*

²⁷ Emilee Rader, et al. (2018). *Explanations as Mechanisms for Supporting Algorithmic Transparency*. *ACM SIGCHI Conference on Human Factors in Computing Systems (CHI)*

questions and based on our rules, you are not eligible”, but the answer for a machine learning based prediction is “*Because of the historical records of hundreds of millions of Americans, and their patterns of repayment, we predict that you will be likely not to repay.*” In fact, ask any machine learning programmer to explain the output of a classifier, and the best answer they can give is often “*Because of all the data.*”

Regulation. Because AI behavior can be so hard to explain, even when it is transparent and the harm is measurable, a central question in AI is “who is responsible for AI behavior?” Is it the people who use the AI that someone created? Is it the people who created the AI? The people who chose the data on which the AI was based? Is it the people who gathered the data that shaped its behavior? These questions are perhaps most salient in the case of driverless cars, which use AI extensively: when autonomous vehicles kill someone, who is to blame? There are few broadly accepted regulatory frameworks that answer these questions. Many Black scholars have [advocated](#) for regulation on applications of facial recognition technology; other scholars are calling for even more pervasive debates about law over the coming decades about AI culpability³³. Increasingly, [whistle blowers](#) from industry are even calling for companies who apply AI to be held accountable to the harms they cause.

³³ Jacob Turner (2018). *Robot Rules: Regulating Artificial Intelligence*. Springer

Ethics. When asking questions about agency, intelligence, causality, and blame, questions about AI quickly turn to questions of ethics, morality, and justice. Why are we creating AI? Who does it benefit? At what cost? Do we want artificially intelligent things in society? What kind of society do we want? And what is the cost of putting so much time and attention into managing the power of AI when so many people in the world are still struggling with such basic needs, such as food, clean water, shelter, and safety? Currently, these are questions largely being answered by large, wealthy, monopolistic, American companies like Google, Amazon, Apple, Facebook, and Microsoft. If we want a democratic society that reflects all of our views of a just society, then everyone needs the literacy above to ponder these questions and shape policy preferences¹⁷.

¹⁷ Anna Jobin, et al. (2019). *The global landscape of AI ethics guidelines. Nature Machine Intelligence*

MAI4CAREU

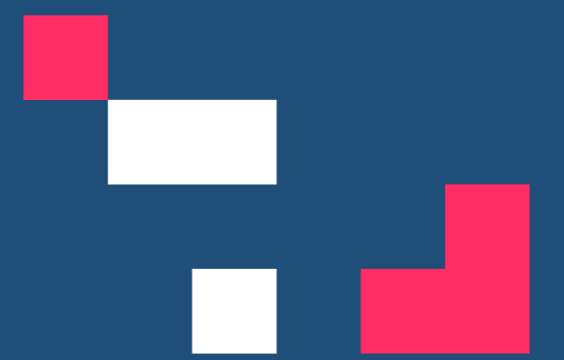
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester



AI and magical thinking

AI as technology and “AI” as speech act

We need to distinguish between

- 1. artificial intelligence (AI) as a technology with practical application:** “as a technology, AI exists somewhere on a spectrum from, practically, at one end, expert systems, path planners, and practical reasoning systems [...] through to, theoretically, at the other end, Alan Turing’s “imaginable digital computers which would do well in the imitation game” or John Haugeland’s synthetic intelligence (i.e., machine intelligence that is constructed but not necessarily imitative)”;
- 2. “artificial intelligence” (“AI”) as a speech act with conventional force:** “a social constructor that stems largely from science fiction with computers and robots having hugely overblown capabilities and a tendency to the apocalyptic”.
“People have been, and are being, “encouraged” to think about artificial intelligence wrongly. Companies are leveraging “AI” to exert control without responsibility.

Peter R Lewis, Stephen Marsh, Jeremy Pitt, *AI vs «AI»: Synthetic Minds or Speech Acts*, in «IEEE Technology and Society Magazine», 2021, <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9445758>

The problem of “trustworthy AI”

“The problem of “trustworthy AI” is one that has great many different “sides.” On the one hand, there are guidelines (for example, from the EU) that tell us how AI should be built and/or behave in order to be seen as “trustworthy”—presumably this means that people are going to (should? must?) trust it.

On the other hand, the problem is seen as “We shouldn’t have to trust AI” because it is a “made thing” and, since it is a human artifact, humans should be held responsible (accountable) when it does something wrong.

In many cases, when they are using marketing speak, those who claim “AI” can be seen as “trustworthy” also claim that it is “beyond the control” of its creators when it leaves the shop floor.”

“It’s not just an evasion of responsibility; it is an exercise in power and it is profoundly wrong.”

Peter R Lewis, Stephen Marsh, Jeremy Pitt, *AI vs «AI»: Synthetic Minds or Speech Acts*, in «IEEE Technology and Society Magazine», 2021, <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9445758>

“We suggest that a democratization of both “AI” and AI is necessary in order to better inform the people who are affected by this deceit. It is not satisfactory to blame the computer—indeed it never has been, yet since we’ve had them, we’ve tried to do exactly that—what is needed is the means to *explain*:

What the system is doing;

Why it does what it does;

How it does this thing;

Why it does it this way;

In ways that the people affected by it understand.

This should not be the responsibility of the machine, since we do not (yet) have AI capable of bearing responsibility for its behavior and operation.”

Peter R Lewis, Stephen Marsh, Jeremy Pitt, *AI vs «AI»: Synthetic Minds or Speech Acts*, in «IEEE Technology and Society Magazine», 2021, <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9445758>

The animation of the inanimate

- As David Hume wrote in *The Natural History of Religion*, “there is an universal tendency among mankind to conceive all beings like themselves, and to transfer to every object those qualities with which they are familiarly acquainted and of which they are intimately conscious”.
- **“the animation of the inanimate” – is, according to Freud, the very nature of magical thinking:** “the misunderstanding” whereby we “put psychological laws in place of natural ones” is still present “in the life of today”, “in living form, as the foundation of language, our beliefs and our philosophy”.
- It is a well-known and yet irresistible tendency: emotional and social responses are automatically generated also by media, such as televisions or computers, and overcoming this unconscious impulse would require the effort of a continuous reflection and the employment of a technical vocabulary, different for each type of object and unfamiliar to most of us.

B. Reeves, C. Nass, [The Media Equation: How People Treat Computers, Television and New Media Like Real People and Places](#), Cambridge, Cambridge University Press, 1996.

D. Tafani, *What's wrong with “AI ethics” narratives*, in «Bollettino telematico di filosofia politica», 2022, pp. 1-22 (forthcoming)

ELIZA

“I chose the name ELIZA for the language analysis program because, like the Eliza of Pygmalion fame, it could be taught to “speak” increasingly well.

Because conversations must be about something, that is, because they must take place within some context, the program was constructed in a two-tier arrangement, the first tier consisting of the language analyzer and the second of a script. The script is a set of rules rather like those that might be given to an actor who is to use them to improvise around a certain theme. Thus ELIZA could be given a script to enable it to maintain a conversation about cooking eggs or about managing a bank checking account, and so on. Each specific script thus enabled ELIZA to play a specific conversational role.”

J. Weizenbaum, *Computer Power and Human Reason. From Judgement to Calculation*, San Francisco, W.H. Freeman & Company, 1976.

DOCTOR

“For my first experiment, I gave ELIZA a script designed to permit it to play (I should really say parody) the role of a Rogerian psychotherapist engaged in an initial interview with a patient. The Rogerian psychotherapist is relatively easy to imitate because much of his technique consists of drawing his patient out by reflecting the patient’s statements back to him.

The following conversation between a young lady and Eliza playing doctor illustrates both the Rogerian technique of encouraging a patient to keep talking and the operation of the computer program ELIZA. The first to “speak” is the young lady. The computer’s responses are printed entirely in capitals.”

J. Weizenbaum, *Computer Power and Human Reason. From Judgement to Calculation*, San Francisco, W.H. Freeman & Company, 1976.

A “conversation” with DOCTOR

“Men are all alike

IN WHAT WAY

They’re always bugging about something or other.

CAN YOU THINK OF A SPECIFIC EXAMPLE

Well, my boyfriend made me come here.

YOUR BOYFRIEND MADE YOU COME HERE

He says I’m depressed much of the time.

I AM SORRY TO HEAR YOU ARE DEPRESSED

It’s true. I am unhappy.

**DO YOU THINK COMING HERE WILL HELP YOU NOT TO
BE UNHAPPY**

I need some help, that much seems certain.

**WHAT WOULD IT MEAN TO YOU IF YOU GOT SOME
HELP**

Perhaps I could learn to get along with my mother.

TELL ME MORE ABOUT YOUR FAMILY

My mother takes care of me.

WHO ELSE IN YOUR FAMILY TAKES CARE OF YOU

My father.

YOUR FATHER

You are like my father in some ways.”



“The reckless anthropomorphization of the computer”

“I was startled to see how quickly and how very deeply people conversing with DOCTOR became emotionally involved with the computer and how unequivocally they anthropomorphized it. Once my secretary, who had watched me work on the program for many months and therefore surely knew it to be merely a computer program, started conversing with it. After only a few interchanges with it, she asked me to leave the room.

Another time, I suggested I might rig the system so that I could examine all conversations anyone had had with it, say, overnight. I was promptly bombarded with accusations that what I proposed amounted to spying on people’s most intimate thoughts.”

J. Weizenbaum, *Computer Power and Human Reason. From Judgement to Calculation*, San Francisco, W.H. Freeman & Company, 1976.

“Enormously exaggerated attributions”

“Another widespread, and to me surprising, reaction to the ELIZA program was the spread of a belief that it demonstrated a general solution to the problem of computer understanding of natural language. In my paper, I had tried to say that no general solution to that problem was possible, i.e., that language is understood only in contextual frameworks, that even these can be shared by people to only a limited extent, and that consequently even people are not embodiments of any such general solution.”

“This reaction to ELIZA showed me more vividly than anything I had seen hitherto the enormously exaggerated attributions an even well-educated audience is capable of making, even strives to make, to a technology it does not understand. Surely, I thought, decisions made by the general public about emergent technologies depend much more on what that public attributes to such technologies than on what they actually are or can and cannot do. If, as appeared to be the case, the public’s attributions are wildly misconceived, then public decisions are bound to be misguided and.”

J. Weizenbaum, *Computer Power and Human Reason. From Judgement to Calculation*, San Francisco, W.H. Freeman & Company, 1976.

AI meets natural stupidity

“Wishful Mnemonics

A major source of simple-mindedness in AI programs is the use of mnemonics like "UNDERSTAND" or "GOAL" to refer to programs and data structures. This practice has been inherited from more traditional programming applications, in which it is liberating and enlightening to be able to **refer to program structures by their purposes.**”

“However, in AI, our programs to a great degree are problems rather than solutions. If a researcher tries to write an "understanding" program, it isn't because he has thought of a better way of implementing this well-understood task, but because he thinks he can come closer to writing the *first* implementation. If he calls the main loop of his program "UNDERSTAND", he is (until proven innocent) merely begging the question. He may mislead a lot of people, most prominently himself, and enrage a lot of others.”

D. McDermott, *AI Meets Natural Stupidity*, in «ACM SIGART Bulletin», 1976, n. 57, pp. 4-9.

The first-step fallacy

“Advances on a specific AI task are often described as “a first step” towards more general AI. The chessplaying computer Deep Blue was “was hailed as the first step of an AI revolution”. IBM described its Watson system as “a first step into cognitive systems, a new era of computing”. OpenAI’s GPT-3 language generator was called a “step toward general intelligence”.

Indeed, if people see a machine do something amazing, albeit in a narrow area, they often assume the field is that much further along toward general AI. The philosopher Hubert Dreyfus (using a term coined by Yehoshua Bar-Hillel) called this a “first-step fallacy.”

As Dreyfus characterized it, **“The first-step fallacy is the claim that, ever since our first work on computer intelligence we have been inching along a continuum at the end of which is AI so that any improvement in our programs no matter how trivial counts as progress.”**

Dreyfus quotes an analogy made by his brother, the engineer Stuart Dreyfus: **“It was like claiming that the first monkey that climbed a tree was making progress towards landing on the moon”.**

Melanie Mitchell, [*Why AI is Harder Than We Think*](#), 2021

Max Weber's theory of disenchantment

“the growing process of intellectualization and rationalization does not imply a growing understanding of the conditions under which we live. It means something quite different.

It is the knowledge or the conviction that if only we wished to understand them we could do so at any time.

It means that in principle, then, we are not ruled by mysterious, unpredictable forces, but that, on the contrary, we can in principle control everything by means of calculation. That in turn means the disenchantment of the world. Unlike the savage for whom such forces existed, we need no longer have recourse to magic in order to control the spirits or pray to them. Instead, technology and calculation achieve our ends. This is the primary meaning of the process of intellectualization.”

M. Weber, *The Vocation Lectures: Science As A Vocation, Politics As A Vocation*, ed. by D.S. Owen, T.B. Strong; transl. by R. Livingstone, 2004.

Enchanted determinism

“What makes contemporary **deep learning systems** interesting is their ambivalent position with respect to Weber’s larger thesis. They **certainly embody aspects of a disenchanted world in that they work to master or control new domains of social life through technical forms of calculation.** [...]”

At the same time, these systems seem to violate the epistemology of disenchantment, the idea that there are no longer “mysterious” forces acting in the world. Paradoxically, when the disenchanted predictions and classifications of deep learning work as hoped, **we see a profusion of optimistic discourse that characterizes these systems as magical, appealing to mysterious forces and superhuman power.** [...] **It is a form of power without knowledge.”**

“**Enchanted determinism**”: “a discourse that presents deep learning techniques as magical, outside the scope of present scientific knowledge, yet also deterministic, in that deep learning systems can nonetheless detect patterns that give unprecedented access to people’s identities, emotions and social character. These systems become deterministic when they are deployed unilaterally in critical social areas, from healthcare to the criminal justice system, creating ever more granular distinctions, relations, and hierarchies that are outside of political or civic processes, with consequences that even their designers may not fully understand or control.”

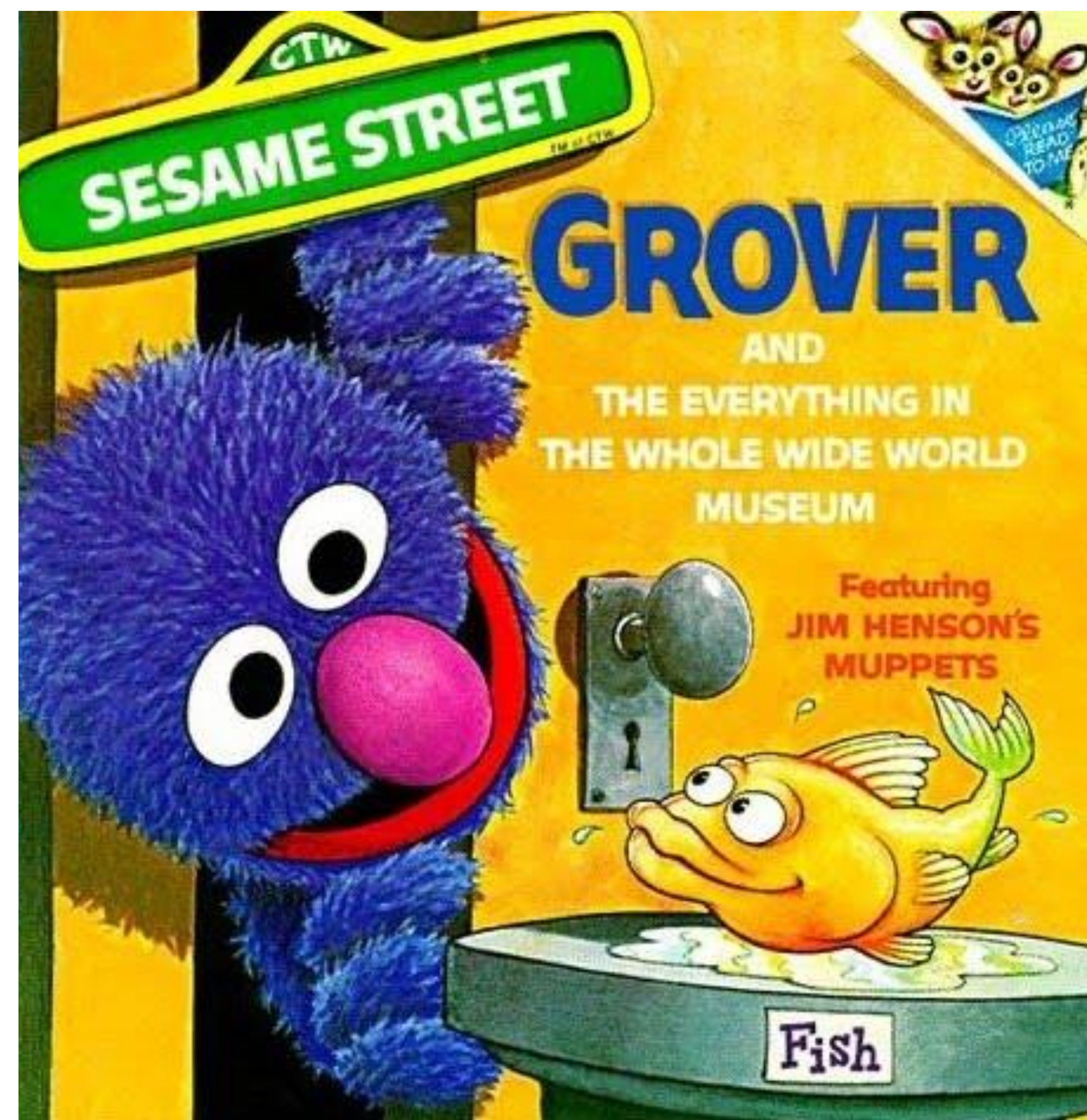
A. Campolo, K. Crawford, *Enchanted Determinism: Power without Responsibility in Artificial Intelligence*, in «Engaging Science, Technology, and Society», 6 (2020), pp. 1-19.

AI and the Everything in the Whole Wide World Benchmark

In the 1974 Sesame Street children’s storybook *Grover and the Everything in the Whole Wide World Museum* [Stiles and Wilcox, 1974], the Muppet monster Grover visits a museum claiming to showcase “everything in the whole wide world”. Example objects representing certain categories fill each room. Several categories are arbitrary and subjective, including showrooms for “Things You Find On a Wall” and “The Things that Can Tickle You Room”. Some are oddly specific, such as “The Carrot Room”, while others unhelpfully vague like “The Tall Hall”. When he thinks that he has seen all that is there, Grover comes to a door that is labeled “Everything Else”. He opens the door, only to find himself in the outside world.

As a children’s story, Grover’s described situation is meant to be absurd. However, in this paper, we discuss how a similar faulty logic is inherent to recent trends in artificial intelligence (AI) — and specifically machine learning (ML) — evaluation, where many popular benchmarks rely on the same false assumptions inherent to the ridiculous “Everything in the Whole Wide World Museum” that Grover visits. In particular, we argue that benchmarks presented as measurements of progress towards general ability within vague tasks such as “visual understanding” or “language understanding” are as ineffective as the finite museum is at representing “everything in the whole wide world,” and for similar reasons — being inherently specific, finite and contextual.

Benchmarks like GLUE [Wang et al., 2019a] or ImageNet [Deng et al., 2009] are often elevated to become definitions of the essential common tasks to validate the performance of any given model. As a result, often the claims that are justified through these benchmark datasets extend far beyond the tasks they are initially designed for, and reach beyond even the initial ambitions for development. Despite a presentation and acceptance as markers of progress towards general-purpose capabilities, there are clear limitations of these benchmarks. In fact, the reality of their development, use and adoption indicates a *construct validity* issue, where the involved benchmarks — due to their instantiation in particular data, metrics and practice — cannot possibly capture anything representative of the claims to general applicability being made about them.



I.D. Raji, E.M. Bender, A. Paullada, E. Denton, A. Hanna, *AI and the Everything in the Whole Wide World Benchmark*, <https://arxiv.org/abs/2111.15366>



Co-financed by the European Union
Connecting Europe Facility

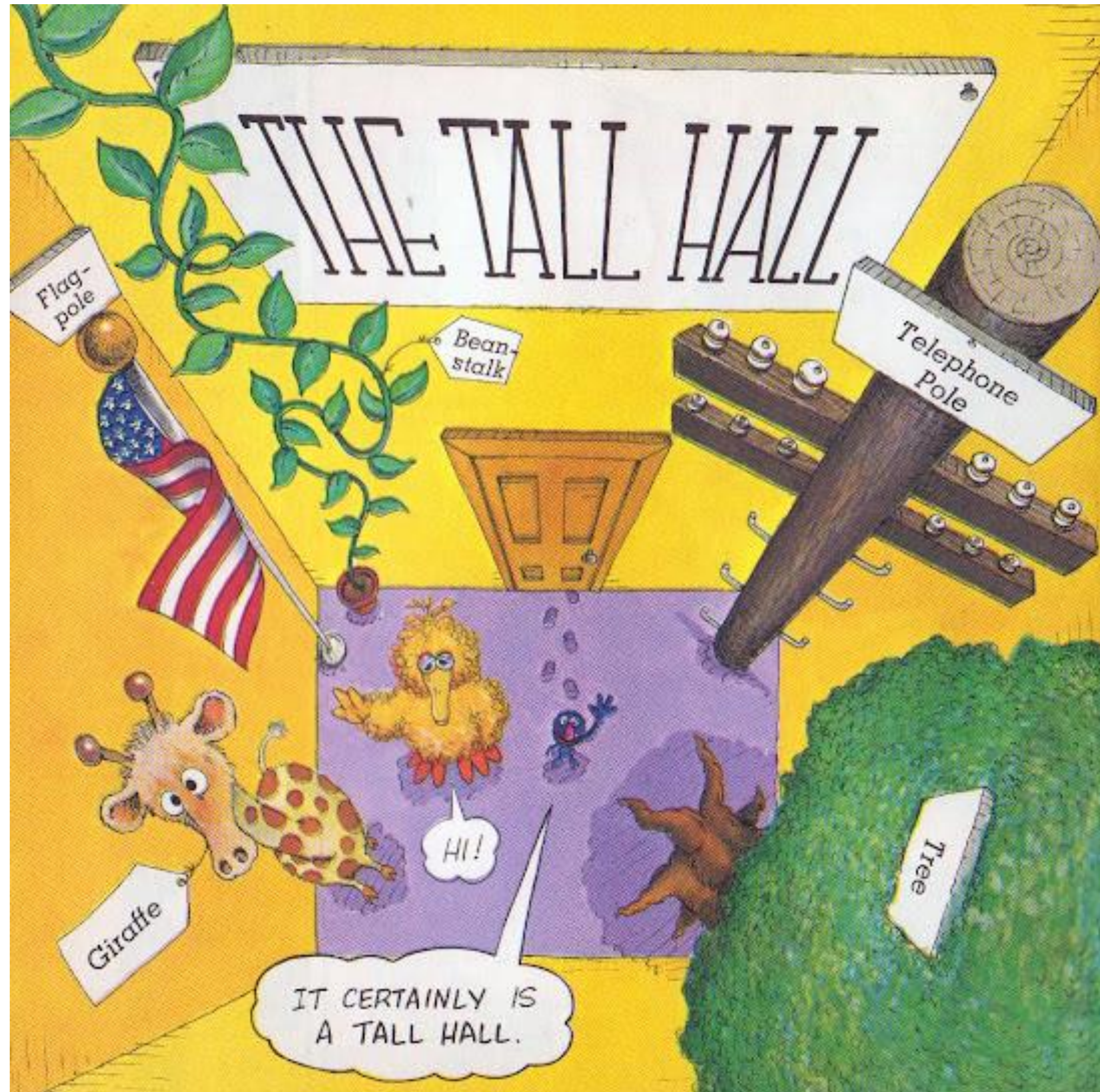
This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423













AI and the Everything in the Whole Wide World Benchmark

“Limits of Benchmarking General Capabilities”

- “The imagined artifact of the “general” benchmark does not actually exist. Real data is designed, subjective and limited in ways that necessitate a different framing from that of any claim to general knowledge or general-purpose capabilities. In fact, presenting any single dataset in this way is ultimately dangerous and deceptive, resulting in misguidance on task design and focus, underreporting of the many biases and subjective interpretations inherent in the data as well as enabling, through false presentations of performance, potential model misuse”
- “benchmarking is a limited approach to assess general model capabilities”

I.D. Raji, E.M. Bender, A. Paullada, E. Denton, A. Hanna, *AI and the Everything in the Whole Wide World Benchmark*, <https://arxiv.org/abs/2111.15366>

AI and the Everything in the Whole Wide World Benchmark

“The situation with Grover and the museum’s claims are clearly ridiculous—yet in machine learning, we follow the exact same logical fallacies to justify the elevation of a select number of benchmarks operating as general benchmarks for the field. However, there is no dataset that will be able to capture the full complexity of the details of existence, in the same way that there can be no museum to contain the full catalog of everything in the whole wide world. Open-world, universal and neutral datasets don’t exist, and current methods of benchmarking do not offer meaningful measures of general capabilities.”

“language understanding relies not only on linguistic competence but also world knowledge, commonsense reasoning, and the ability to model the interlocutor’s state of mind, none of which can be thoroughly tested through text-only tasks, such as GLUE. Several researchers have raised the need to establish effective physical and social grounding as part of the process of moving towards robust and effective natural language understanding, warning against text-only learning as a limited approach. Bender and Koller additionally mention the tendency of machine learning researchers to misinterpret certain benchmarks as capturing the model’s ability to decipher meaning in language, arguing that benchmarks need to be constructed with care if they are to show evidence of “understanding” as opposed to merely the ability to manipulate linguistic form sufficiently to pass the test.”

I.D. Raji, E.M. Bender, A. Paullada, E. Denton, A. Hanna, *AI and the Everything in the Whole Wide World Benchmark*, <https://arxiv.org/abs/2111.15366>

Thank you. Any questions?

daniela.tafani@unibo.it

MAI4CAREU

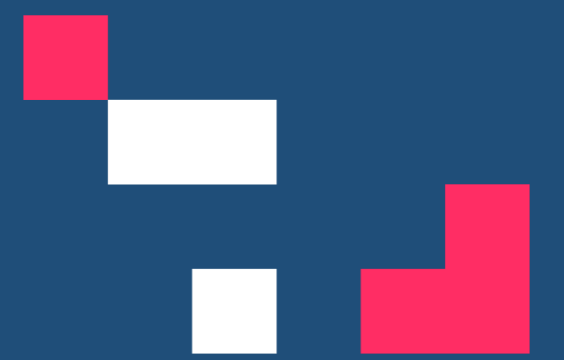
Master programmes in Artificial
Intelligence 4 Careers in Europe

University of Bologna

Computational Ethics

Daniela Tafani

2022/2023 – Second Semester



Universal Moral Grammar

Universal Moral Grammar

1. Noam Chomsky
2. John Rawls
3. John Mikhail

Appendix. Immanuel Kant

1. Noam Chomsky

Generative grammars as theories of linguistic competence

“Linguistic theory is concerned primarily with an ideal speaker-listener, in a completely homogeneous speech-community, who knows its language perfectly and is unaffected by such grammatically irrelevant conditions as memory limitations, distractions, shifts of attention and interest, and errors (random or characteristic) in applying his knowledge of the language in actual performance.”

“We thus make a fundamental distinction between **competence** (the speaker-hearer's knowledge of his language) and **performance** (the actual use of language in concrete situations). Only under the idealization set forth in the preceding paragraph is performance a direct reflection of competence. In actual fact, it obviously could not directly reflect competence. A record of natural speech will show numerous false starts, deviations from rules, changes of plan in mid-course, and so on. **The problem for the linguist, as well as for the child learning the language, is to determine from the data of performance the underlying system of rules** that has been mastered by the speaker-hearer and that he puts to use in actual performance”

N. Chomsky, *Aspects of the Theory of Syntax*, Cambridge, Mass., The M.I.T. Press, 1965.

“A grammar of a language purports to be a description of the ideal speaker-hearer's intrinsic competence. If the grammar is, furthermore, perfectly explicit - in other words, if it does not rely on the intelligence of the understanding reader but rather provides an explicit analysis of his contribution - we may (somewhat redundantly) call it **a generative grammar**. A fully adequate grammar must assign to each of an infinite range of sentences a structural description indicating how this”.

“valuable as they obviously are, **traditional grammars are deficient in that they leave unexpressed many of the basic regularities of the language with which they are concerned**. This fact is particularly clear on the level of syntax, where no traditional or structuralist grammar goes beyond classification of particular examples to the stage of formulation of generative rules on any significant scale. An analysis of the best existing grammars will quickly reveal that this is a defect of principle, not just a matter of empirical detail or logical preciseness”.

N. Chomsky, *Aspects of the Theory of Syntax*, Cambridge, Mass., The M.I.T. Press, 1965.

“Within traditional linguistic theory, furthermore, it was clearly understood that **one of the qualities that all languages have in common is their "creative" aspect**. Thus an essential property of language is that it **provides the means for expressing indefinitely many thoughts** and for reacting appropriately in an indefinite range of new situations”.

“The grammar of a particular language, then, is to be supplemented by a universal grammar that accommodates the creative aspect of language use and expresses the deep-seated regularities which, being universal, are omitted from the grammar itself.”

“**language can (in Humboldt's words) "make infinite use of finite means"**”.

“by a generative grammar I mean simply a system of rules that in some explicit and welldefined way assigns structural descriptions to sentences. Obviously, **every speaker of a language has mastered and internalized a generative grammar that expresses his knowledge of his language. This is not to say that he is aware of the rules of the grammar or even that he can become aware of them**, or that his statements about his intuitive knowledge of the language are necessarily accurate.

N. Chomsky, *Aspects of the Theory of Syntax*, Cambridge, Mass., The M.I.T. Press, 1965.

2. John Rawls

“Some remarks about moral theory”

“It seems desirable at this point, in order to prevent misunderstanding, to discuss briefly the nature of moral theory. I shall do this by explaining in more detail the concept of a considered judgment in reflective equilibrium and the reasons for introducing it.

Let us assume that each person beyond a certain age and possessed of the requisite intellectual capacity develops a sense of justice under normal social circumstances. We acquire a skill in judging things to be just and unjust, and in supporting these judgments by reasons. Moreover, we ordinarily have some desire to act in accord with these pronouncements and expect a similar desire on the part of others. Clearly **this moral capacity is extraordinarily complex. To see this it suffices to note the potentially infinite number and variety of judgments that we are prepared to make.** The fact that we often do not know what to say, and sometimes find our minds unsettled, does not detract from the complexity of the capacity we have.”

J. Rawls, *A Theory of Justice*, Revised Edition, Cambridge, Massachusetts, The Belknap Press of Harvard University Press, 1999.

“Now one may think of **moral theory** at first (and I stress the provisional nature of this view) **as the attempt to describe our moral capacity**; or, in the present case, one may regard a theory of justice as describing our sense of justice. **By such a description is not meant simply a list of the judgments** on institutions and actions that we are prepared to render, accompanied with supporting reasons when these are offered. Rather, **what is required is a formulation of a set of principles which, when conjoined to our beliefs and knowledge of the circumstances, would lead us to make these judgments** with their supporting reasons were we to apply these principles conscientiously and intelligently. A conception of justice characterizes our moral sensibility when the everyday judgments we do make are in accordance with its principles. These principles can serve as part of the premises of an argument which arrives at the matching judgments. We do not understand our sense of justice until we know in some systematic way covering a wide range of cases what these principles are.”

J. Rawls, *A Theory of Justice*, Revised Edition, Cambridge, Massachusetts, The Belknap Press of Harvard University Press, 1999.

“A useful comparison here is with the problem of describing the sense of grammaticalness that we have for the sentences of our native language^{*}. In this case the aim is to characterize the ability to recognize well-formed sentences by formulating clearly expressed principles which make the same discriminations as the native speaker. This undertaking is known to require theoretical constructions that far outrun the ad hoc precepts of our explicit grammatical knowledge. A similar situation presumably holds in moral theory. There is no reason to assume that our sense of justice can be adequately characterized by familiar common sense precepts, or derived from the more obvious learning principles. A correct account of moral capacities will certainly involve principles and theoretical constructions which go much beyond the norms and standards cited in everyday life; it may eventually require fairly sophisticated mathematics as well.”

^{*}See Noam Chomsky, *Aspects of the Theory of Syntax* (Cambridge, Mass., The M.I.T. Press, 1965), pp. 3–9.

J. Rawls, *A Theory of Justice*, Revised Edition, Cambridge, Massachusetts, The Belknap Press of Harvard University Press, 1999.

3. John Mikhail



Co-financed by the European Union
Connecting Europe Facility

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423



Universal Moral Grammar

Universal moral grammar (UMG) “seeks to describe the nature and origin of moral knowledge by using concepts and models similar to those used in Chomsky’s program in linguistics”.

“Initial evidence for UMG comes from multiple sources, including psychology, linguistics, anthropology and cognitive neuroscience. Although none of this evidence is univocal or conclusive, collectively it provides at least modest support for the hypothesis that humans possess an innate moral faculty that is analogous, in some respects, to the language faculty that has been postulated by Chomsky and other linguists.”

J. Mikhail, *Elements of moral cognition: Rawls’ linguistic analogy and the cognitive science of moral and legal judgment*, Cambridge University Press, Cambridge 2011.

J. Mikhail, *Chomsky and Moral Philosophy*, in *The Cambridge companion to Chomsky*, ed. by J.A. McGilvray, Cambridge University Press, 2017, pp. 235- 253.

“First, developmental psychologists have discovered that the intuitive jurisprudence of young children is complex and exhibits many characteristics of a well-developed legal code. For example,

- 3–4-year-old children use intent or purpose to distinguish two acts that have the same result. They also distinguish ‘genuine’ moral violations (e.g. battery or theft) from violations of social conventions (e.g. wearing pajamas to school).
- 4–5-year-olds use a proportionality principle to determine the correct level of punishment for principals and accessories.
- 5–6-year-olds use false factual beliefs but not false moral beliefs to exculpate.

Second, every natural language seems to have words or phrases to express basic deontic concepts, such as obligatory, permissible, and forbidden, or their equivalents. Moreover, deontic logic is formalizable. The three primary deontic operators can be placed in a square of opposition and equipollence, similar to those for quantified and modal forms.

Third, prohibitions of murder, rape and other types of aggression appear to be universal or nearly so, as do legal distinctions that are based on causation, intention and voluntary behavior. Furthermore, comparative legal scholars have suggested that a few basic distinctions capture the ‘universal grammar’ of all systems of criminal law.”

J. Mikhail, *Universal Moral Grammar: Theory, Evidence, and the Future*, in «Trends in Cognitive Sciences», 2007, Georgetown Public Law Research Paper n. 954398.

“UMG relies on two fundamental arguments: the argument for moral grammar and the argument from the poverty of the moral stimulus.

1. **The argument for moral grammar** holds that the properties of moral judgment imply that the mind contains a moral grammar: a complex and possibly domain-specific set of rules, concepts and principles that generates and relates mental representations of various types. Among other things, this system enables individuals to determine the deontic status of an infinite variety of acts and omissions.
2. **The argument from the poverty of the moral stimulus** holds that the manner in which this grammar is acquired implies that at least some of its core attributes are innate, where ‘innate’ is used in a dispositional sense to refer to cognitive systems whose essential properties are largely pre-determined by the inherent structure of the mind, but whose ontogenetic development must be triggered and shaped by appropriate experience and can be impeded by unusually hostile learning environments.

Both arguments are nondemonstrative and presuppose a familiar set of idealizations and simplifying assumptions. Moreover, both arguments have direct parallels in the case of language.”

J. Mikhail, *Universal Moral Grammar: Theory, Evidence, and the Future*, in «Trends in Cognitive Sciences», 2007, Georgetown Public Law Research Paper n. 954398.

Appendix

Immanuel Kant

The moral cognition of common human reason

“Thus, then, we have arrived, within the moral cognition of common human reason, at its principle, which it admittedly does not think so abstractly in a universal form but which it actually has always before its eyes and uses as the norm for its appraisals. Here it would be easy to show how common human reason, with this compass in hand, knows very well how to distinguish in every case that comes up what is good and what is evil, what is in conformity with duty or contrary to duty, if, without in the least teaching it anything new, we only, as did Socrates, make it attentive to its own principle; and that **there is, accordingly, no need of science and philosophy to know what one has to do in order to be honest and good, and even wise and virtuous.** We might even have assumed in advance that cognizance of what it is incumbent upon everyone to do, and so also to know, would be the affair of every human being, even the most common”

I. Kant, *Groundwork of The metaphysics of morals*, 1785, in *Practical Philosophy, The Cambridge Edition of the Works of Immanuel Kant*, ed. by M. Gregor, Cambridge, Cambridge University Press, 1996.

“Yet we cannot consider without admiration how great an advantage the practical faculty of appraising* has over the theoretical in common human understanding. In the latter, if common reason ventures to depart from laws of experience and perceptions of the senses it falls into sheer incomprehensibilities' and self-contradictions, at least into a chaos of uncertainty, obscurity, and instability. But in practical matters, it is just when common understanding excludes all sensible incentives from practical laws that its faculty of appraising first begins to show itself to advantage. It then becomes even subtle, whether in quibbling tricks with its own conscience or with other claims regarding what is to be called right, or in sincerely wanting to determine the worth of actions for its own instruction; and, what is most admirable, in the latter case it can even have as good a hope of hitting the mark as any philosopher can promise himself; indeed, it is almost more sure in this matter, because a philosopher, though he cannot have any other principle than that of common understanding, can easily confuse his judgment by a mass of considerations foreign and irrelevant to the matter and deflect it from the straight course. Would it not therefore be more advisable in moral matters to leave the judgment of common reason as it is and, at most, call in philosophy only to present the system of morals all the more completely and apprehensibly“ and to present its rules in a form more convenient for use (still more for disputation), but not to lead common human understanding, even in practical matters, away from its fortunate simplicity and to put it, by means of philosophy, on a new path of investigation and instruction?”

I. Kant, *Groundwork of The metaphysics of morals*, 1785, in *Practical Philosophy, The Cambridge Edition of the Works of Immanuel Kant*, ed. by M. Gregor, Cambridge, Cambridge University Press, 1996.

“like the difference between the right and the left hand”

“But if one asks: What, then, really is pure morality, by which as a touchstone one must test the moral content of every action? I must admit that only philosophers can make the decision of this question doubtful, for it is long since decided in common human reason, not indeed by abstract general formulae but by habitual use, like the difference between the right and the left hand”

I. Kant, *Critique of practical reason*, 1788, in *Practical Philosophy, The Cambridge Edition of the Works of Immanuel Kant*, ed. by M. Gregor, Cambridge, Cambridge University Press, 1996.

“Morality first discloses to us the concept of freedom”

“Suppose someone asserts of his lustful inclination that, when the desired object and the opportunity are present, it is quite irresistible to him; ask him whether, if a gallows were erected in front of the house where he finds this opportunity and he would be hanged on it immediately after gratifying his lust, he would not then control his inclination. One need not conjecture very long what he would reply. But ask him whether, if his prince demanded, on pain of the same immediate execution, that he give false testimony against an honorable man whom the prince would like to destroy under a plausible pretext, he would consider it possible to overcome his love of life, however great it may be. He would perhaps not venture to assert whether he would do it or not, but he must admit without hesitation that it would be possible for him. He judges, therefore, that he can do something because he is aware that he ought to do it and cognizes freedom within him, which, without the moral law, would have remained unknown to him.”

I. Kant, *Critique of practical reason*, 1788, in *Practical Philosophy, The Cambridge Edition of the Works of Immanuel Kant*, ed. by M. Gregor, Cambridge, Cambridge University Press, 1996.

Thank you. Any questions?

daniela.tafani@unibo.it