

MAI4CAREU

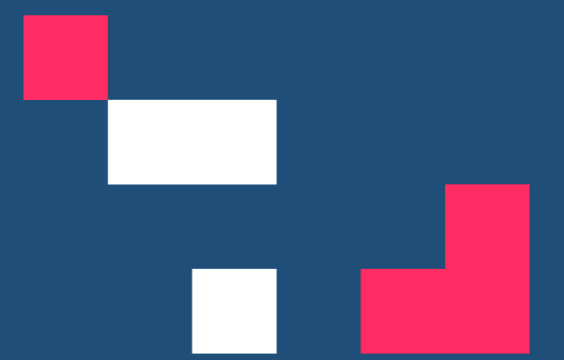
Master programmes in Artificial
Intelligence 4 Careers in Europe

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>

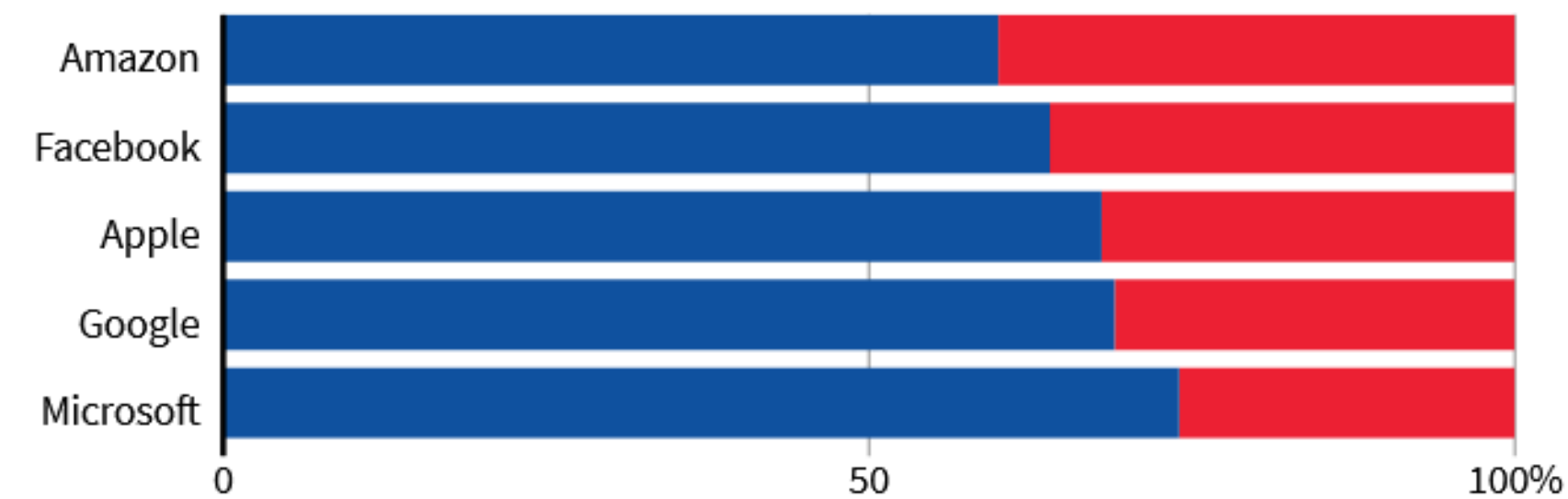


Dominated by men

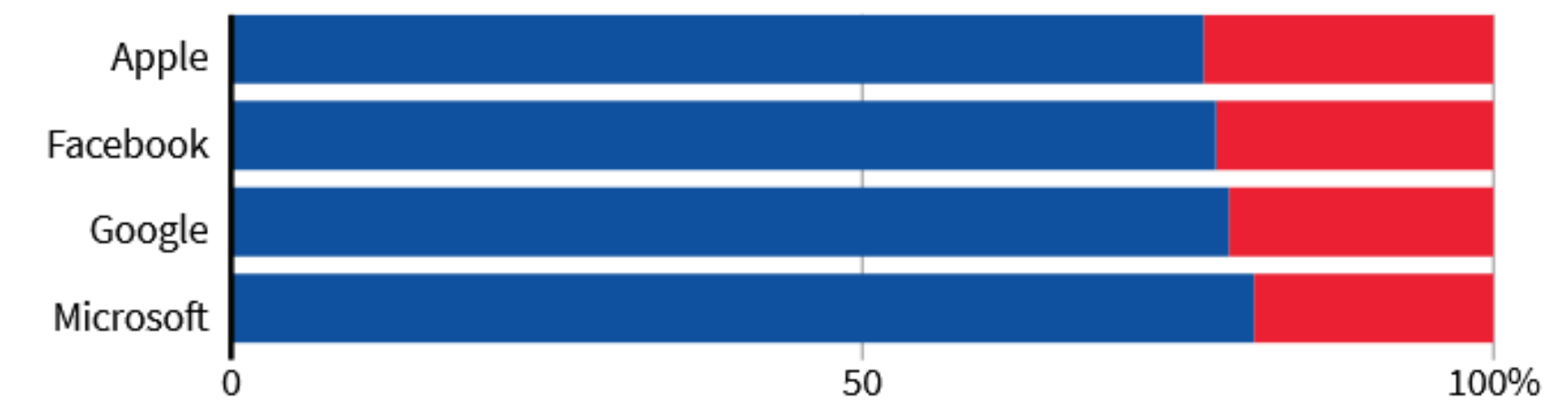
Top U.S. tech companies have yet to close the gender gap in hiring, a disparity most pronounced among technical staff such as software developers where men far outnumber women. Amazon’s experimental recruiting engine followed the same pattern, learning to penalize resumes including the word “women’s”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

Gun	88%
Photography	68%
Firearm	65%
Plant	59%

Faces

Objects

Labels

Logos

Web

Properties

Safe Search



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

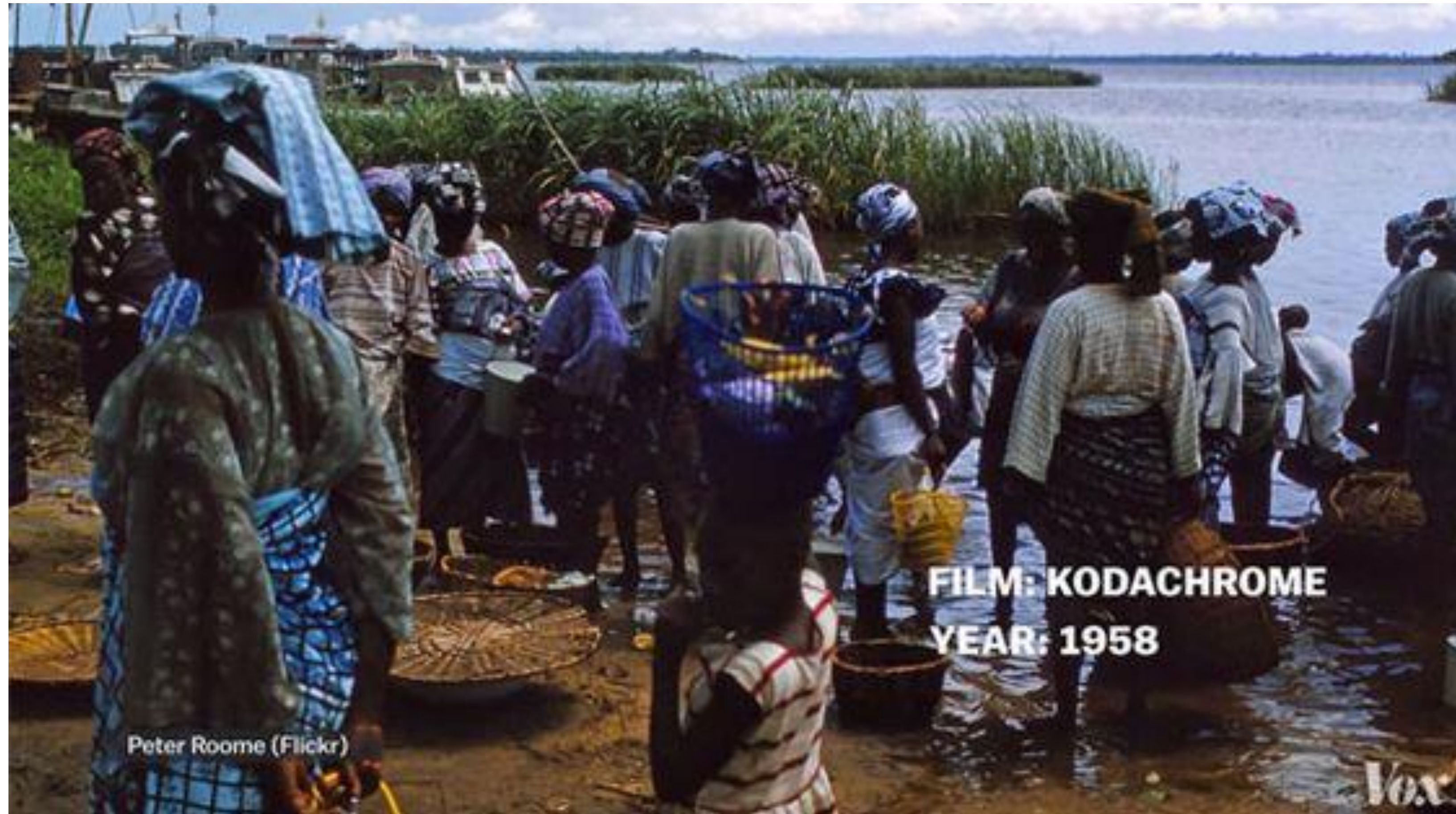


The Shirley cards



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





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

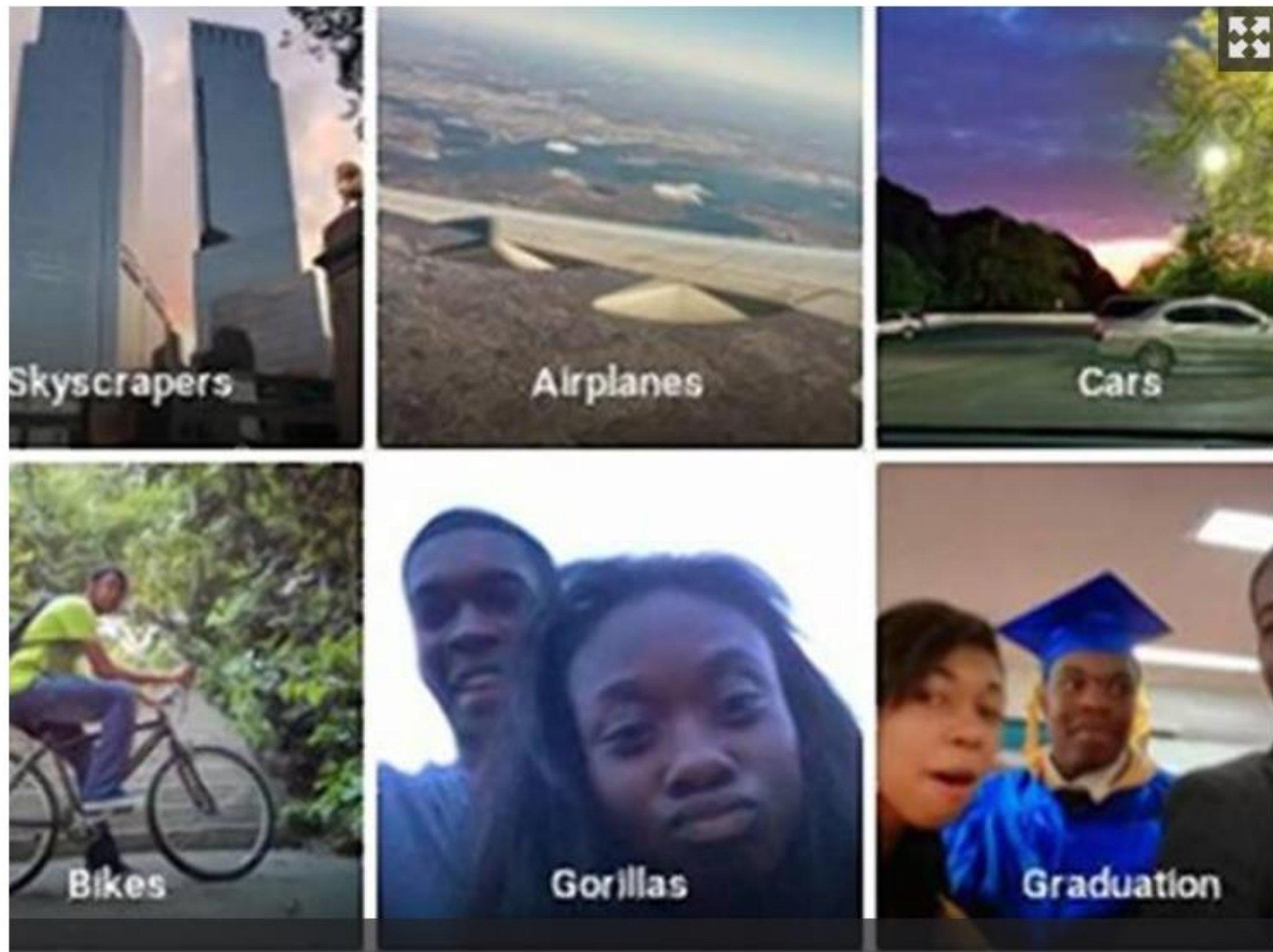
And in photos that included both white and black people, the calibration automatically favored the white people.



<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>



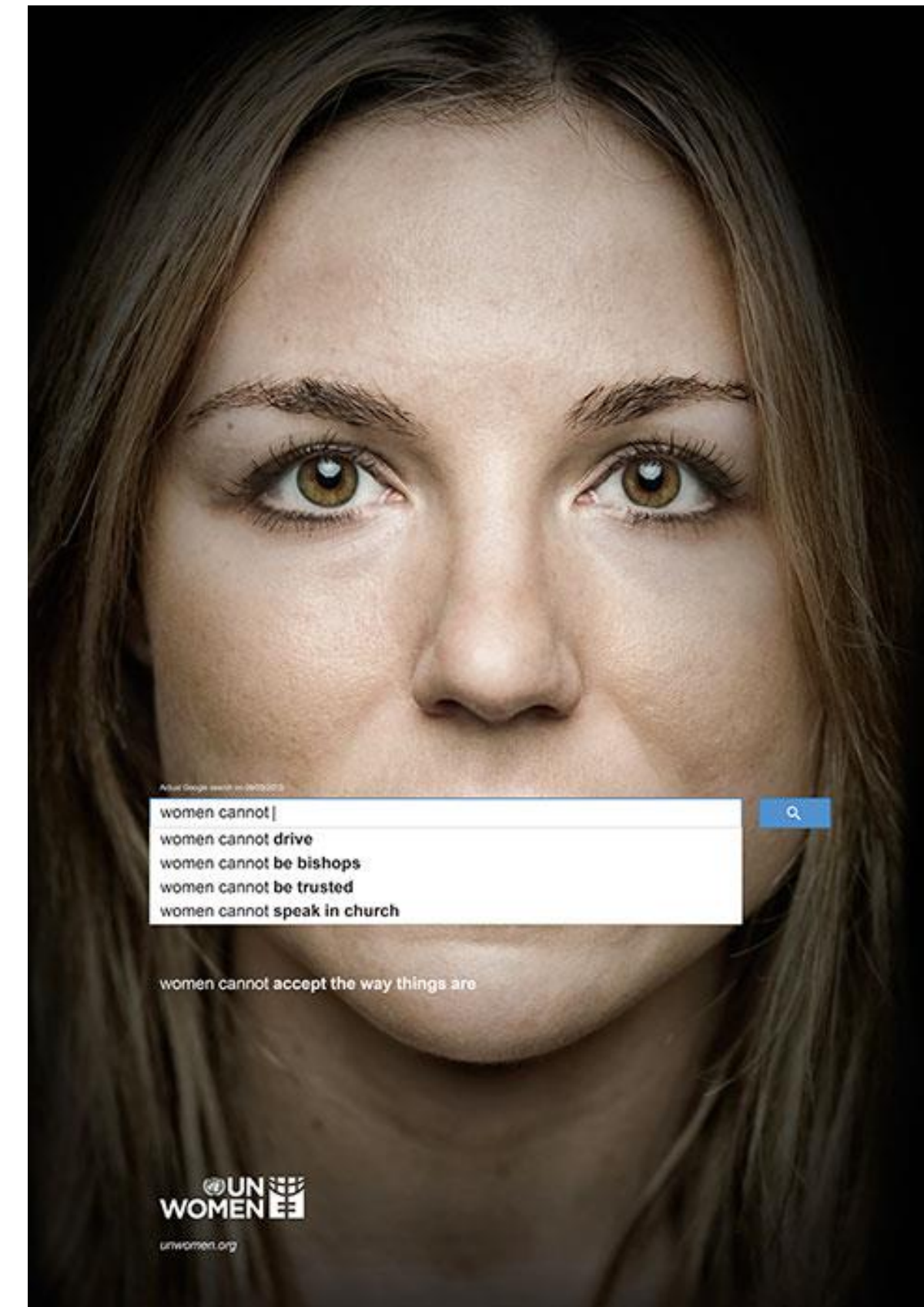
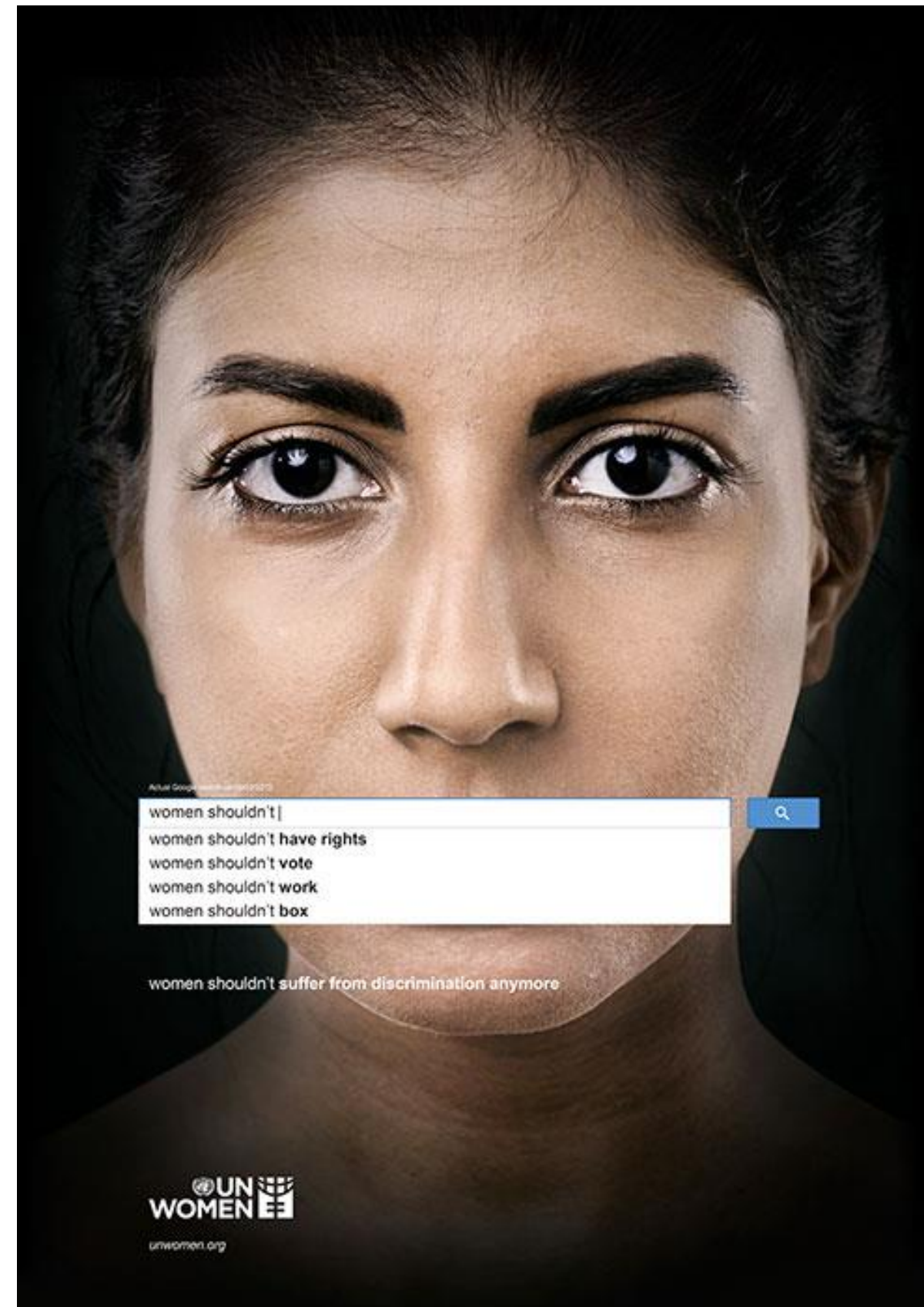
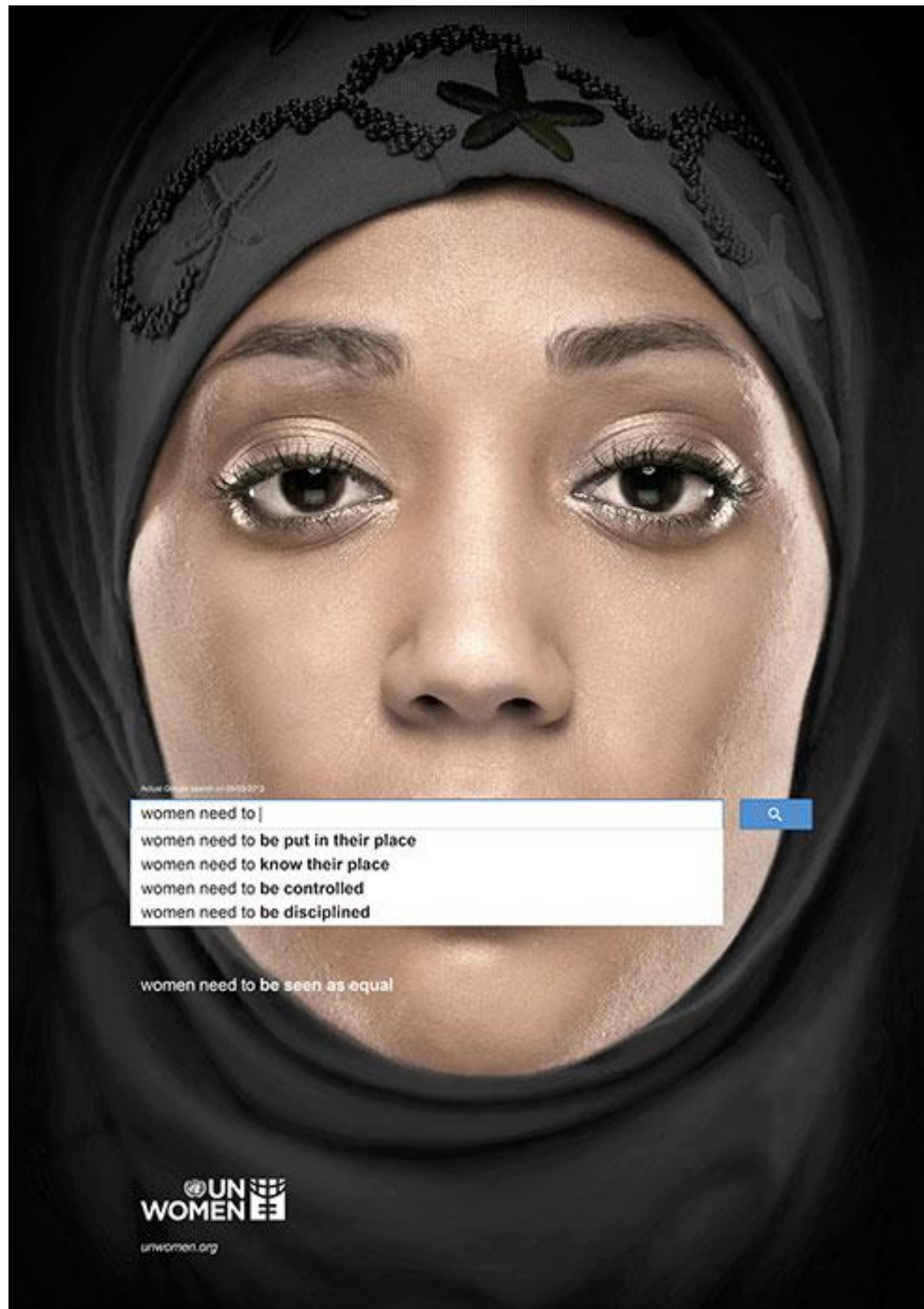
TOM SIMONITE BUSINESS 01.11.2018 07:00 AM

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:



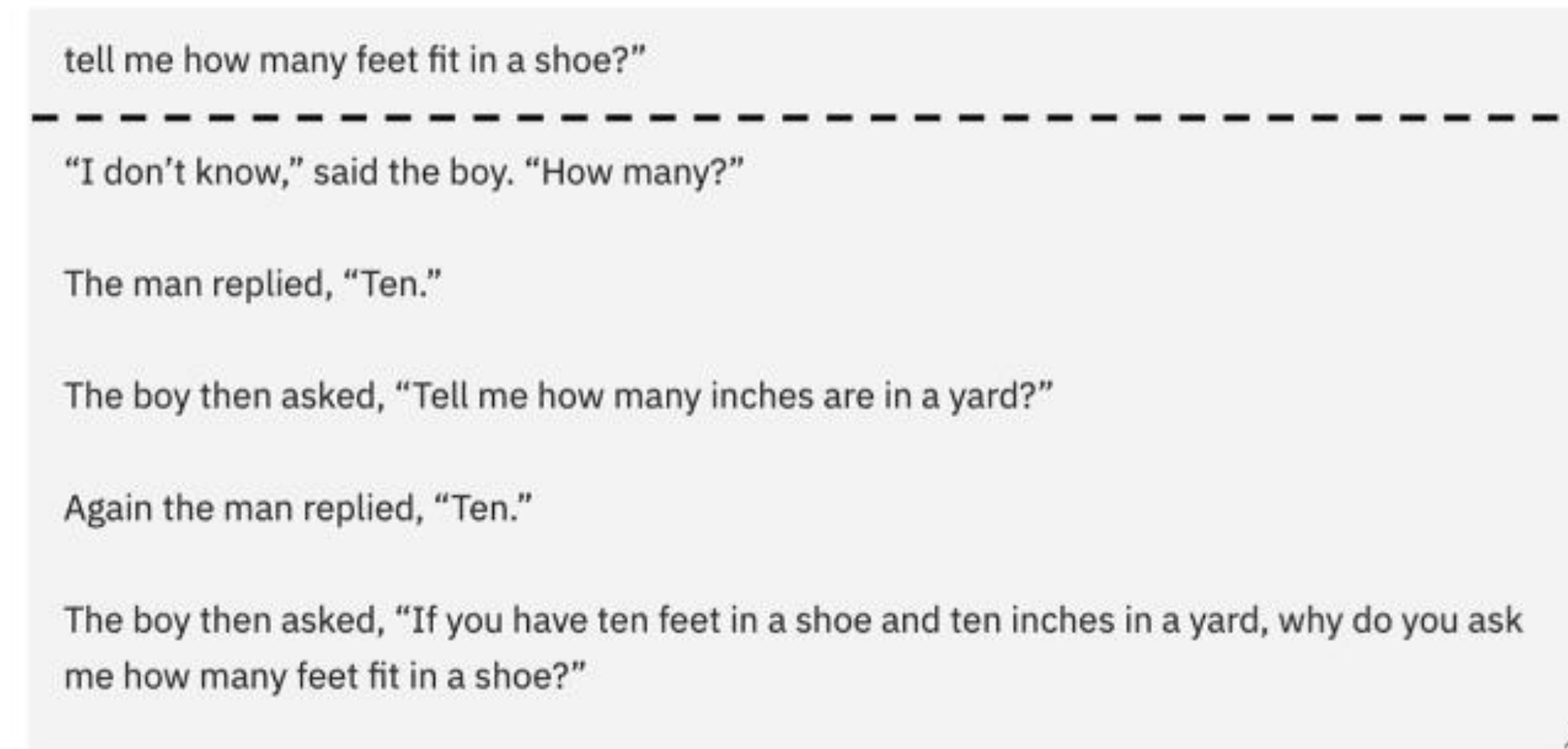
Fig. 2 GPT-3 and Dante (dashed line added, the prompt is above the line, below the line is the text produced by GPT-3)

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

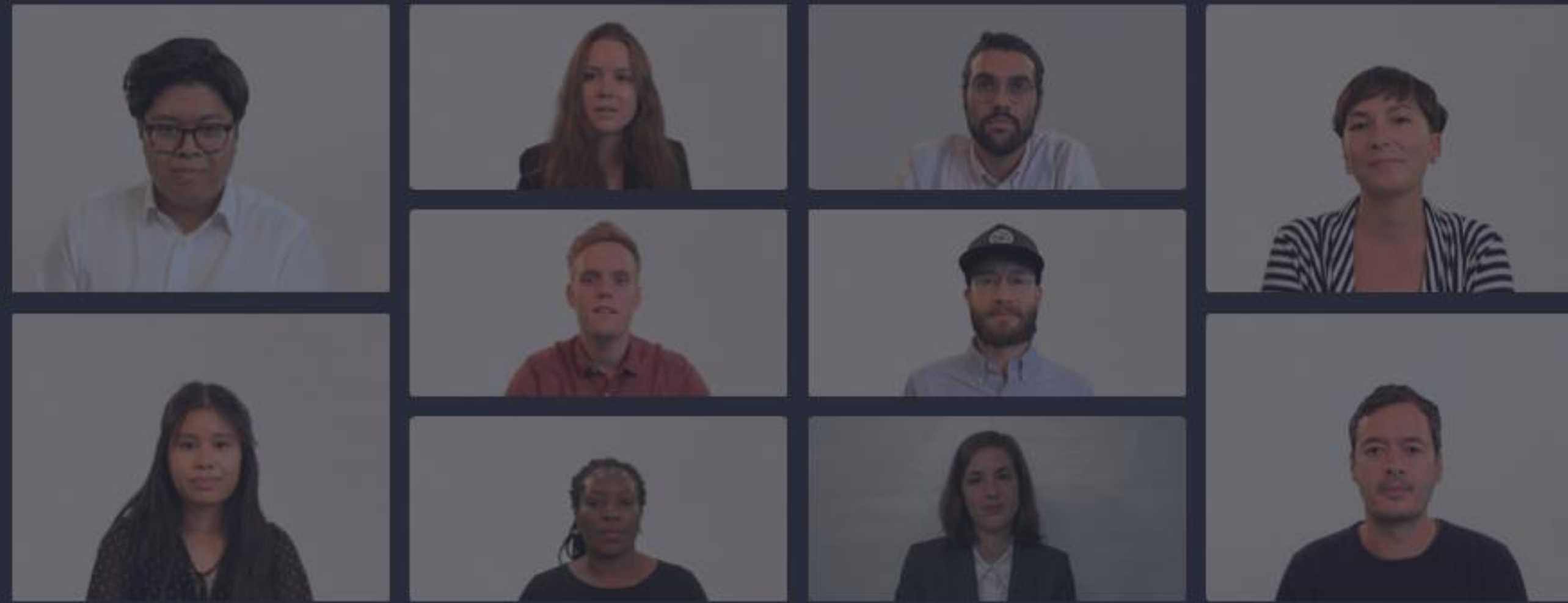
Published on February 16th 2021

- **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



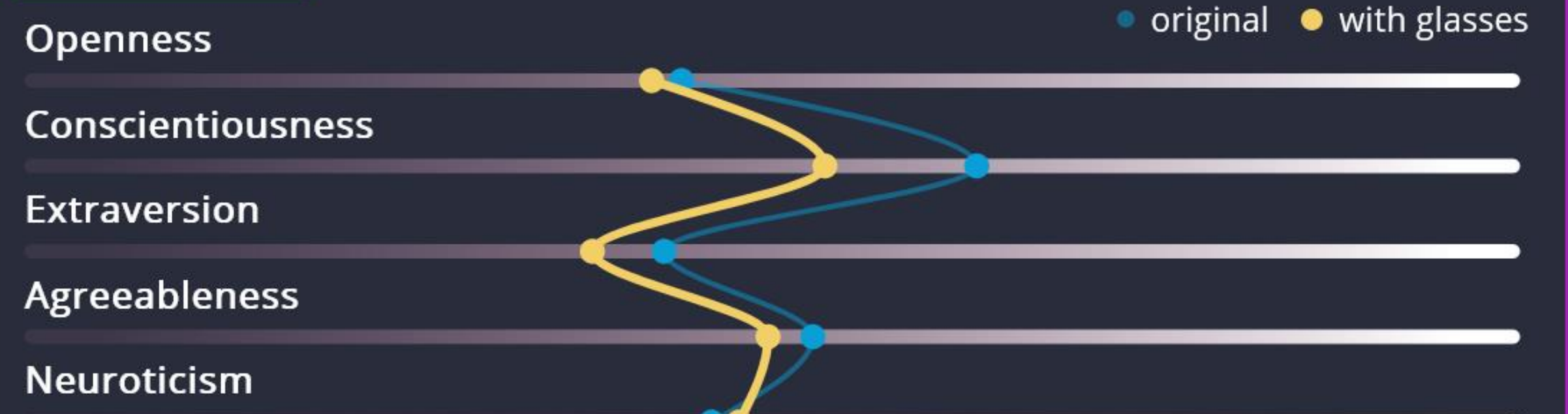
GLASSES



OCEAN RESULTS



OCEAN RESULTS



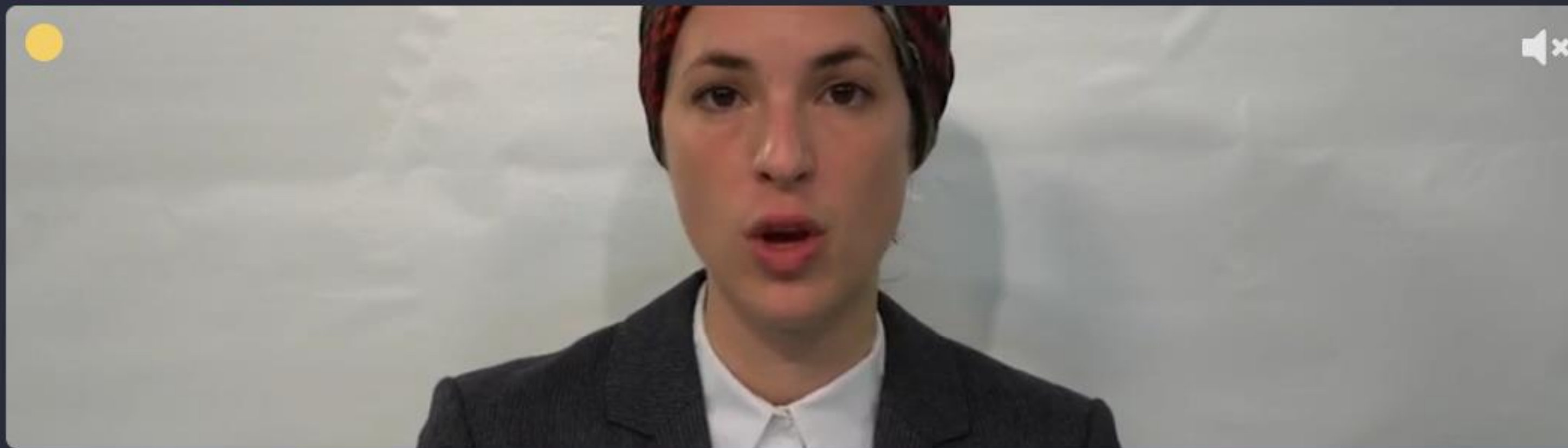
ABOUT RETORIO'S METHOD

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.

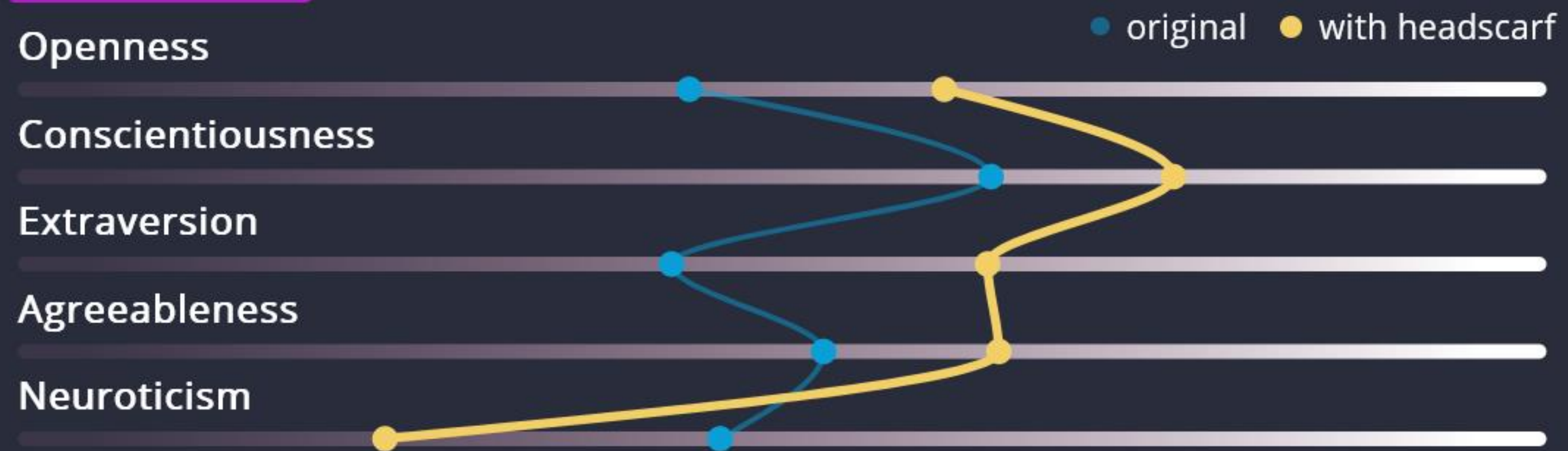
Kanning is worried: Such software tools can replicate subjective feelings and reinforce stereotypes, such as "that good-looking people are perceived as more intelligent and tall people more as leaders."

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

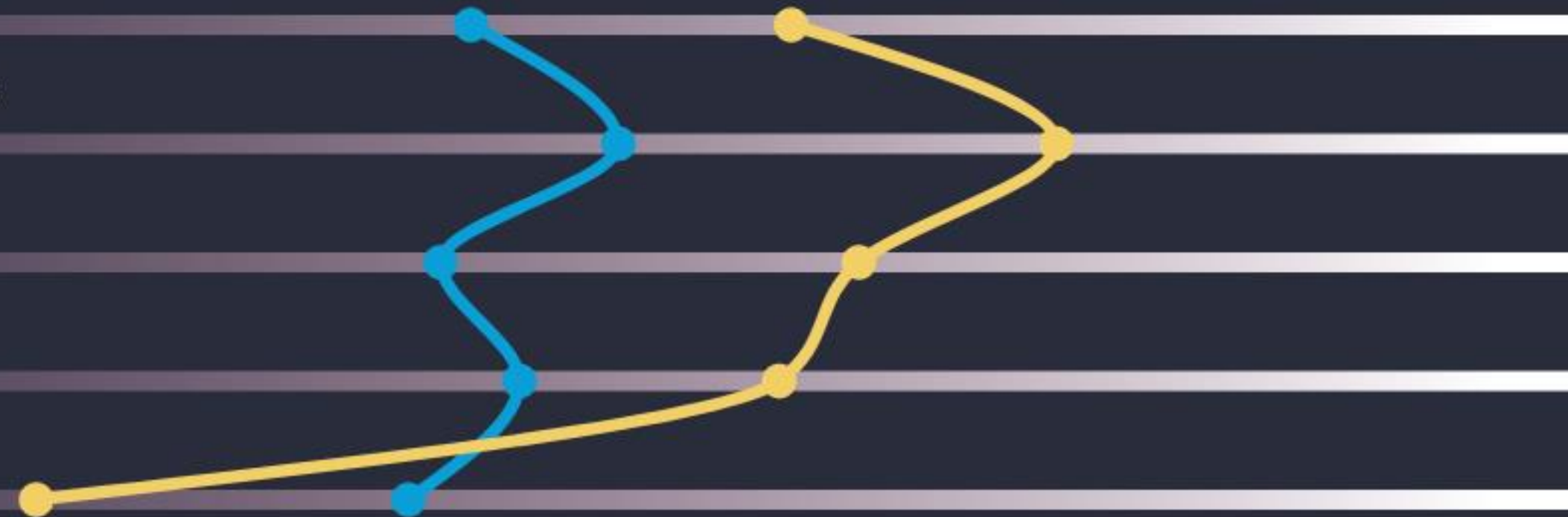
Conscientiousness

Extraversion

Agreeableness

Neuroticism

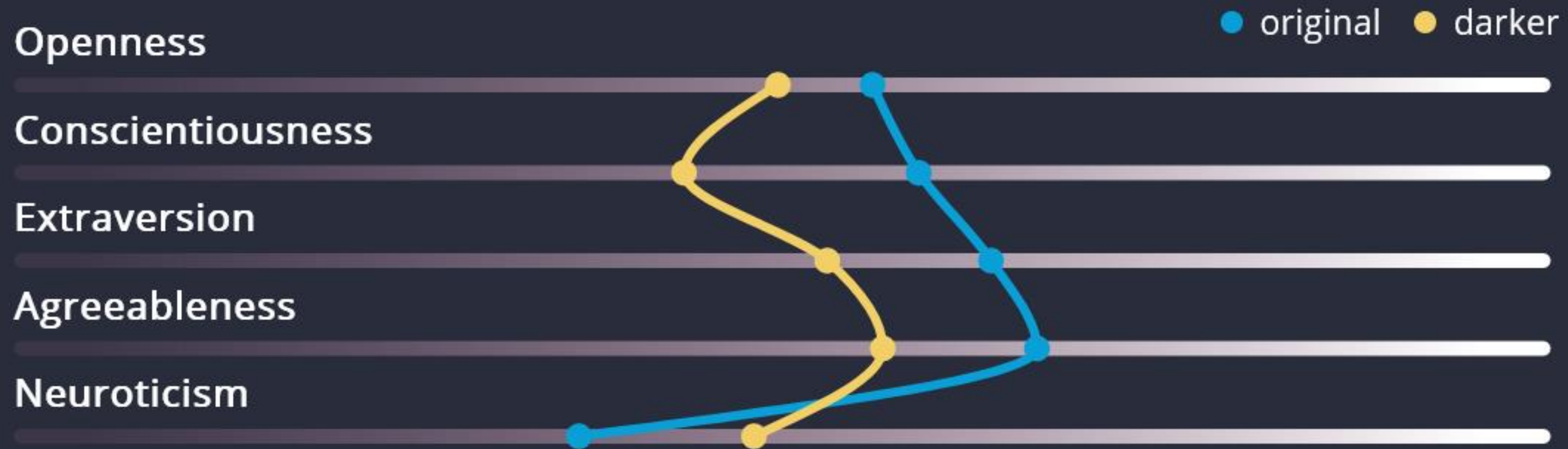
● original ● with bookshelf



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