

Bias in AI

Sharada Prasanna Mohanty

PhD Student, **EPFL**
Co-Founder, **CrowdAI**



@MeMohanty

Bias in **Bad** ?



Bias in **Bad** ?

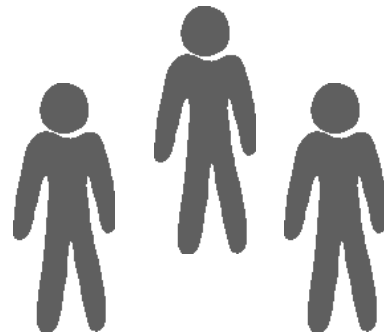
.....its not as **simple** as that !





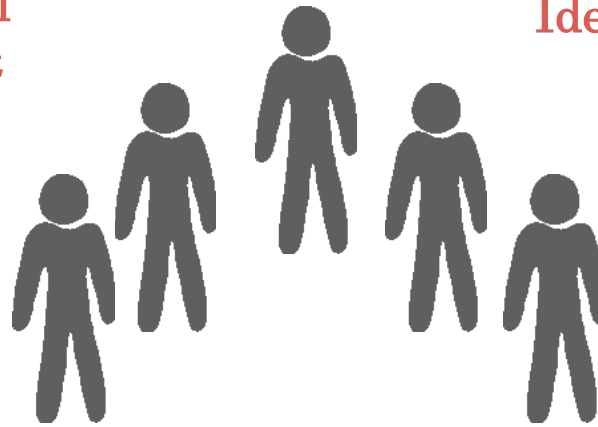
@MeMohanty

Shared
Ideologies

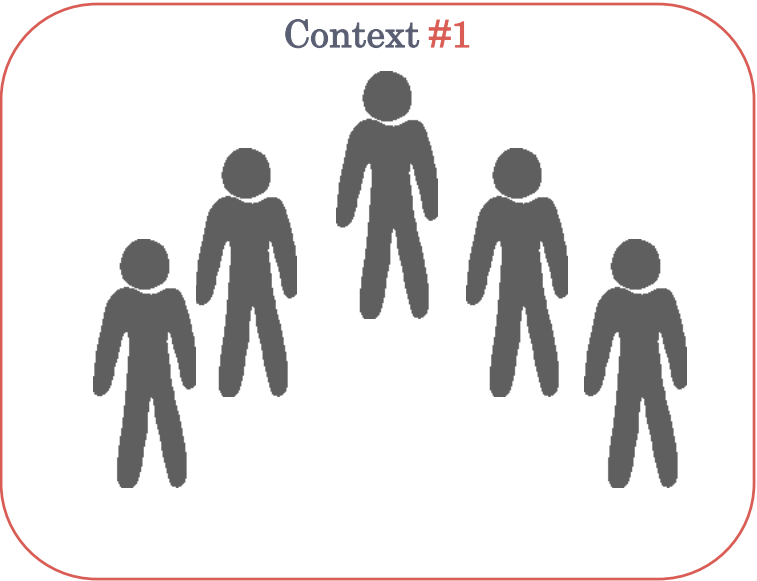


Shared
Cultural
Context

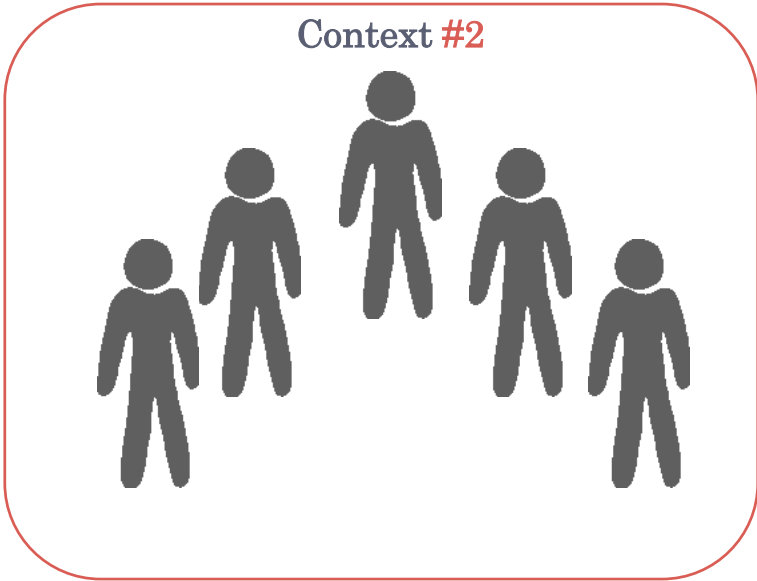
Shared
Ideologies



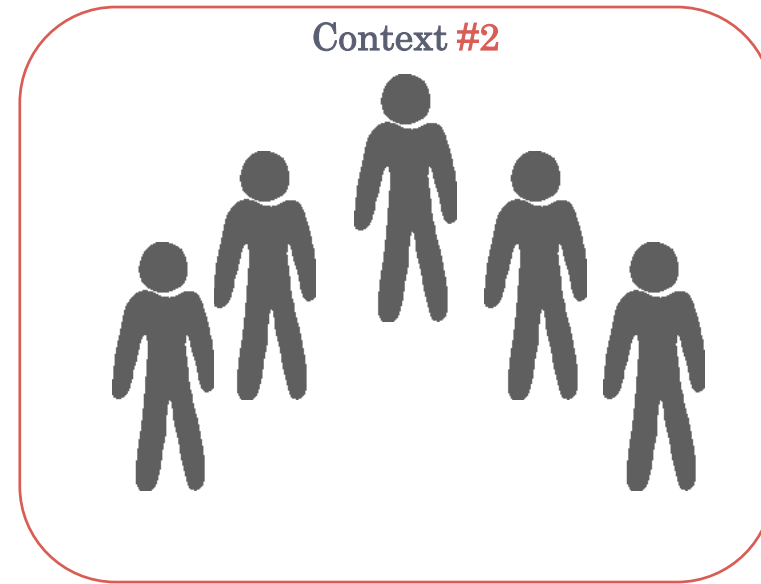
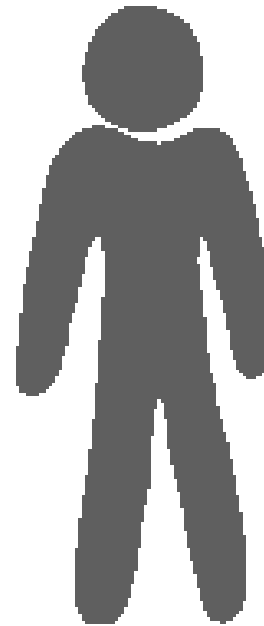
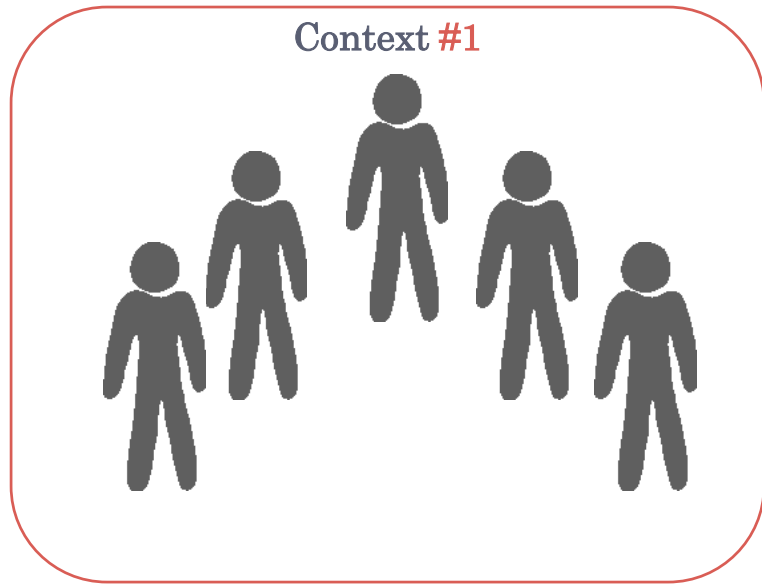
Context #1



Context #2



Being Human





indian food is |

- indian food is **known for its**
- indian food is **the best**
- indian food is **the best food in the world**
- indian food is **nasty**
- indian food is **healthy**
- indian food is **unhealthy**
- indian food is **spicy**
- indian food is **known for its quizlet**
- indian food is **good**
- indian food is **the best food gd topic**

Google Search

I'm Feeling Lucky

[Learn more](#)

Report inappropriate predictions



Artificial Intelligence is |

artificial intelligence is **associated with which generation**

artificial intelligence is **dangerous**

artificial intelligence is **bad**

artificial intelligence is **the future**

artificial intelligence is **good**

artificial intelligence is **a threat**

artificial intelligence is **the new electricity**

artificial intelligence is **impossible**

artificial intelligence is **a threat to humanity**

artificial intelligence is **boon or bane**

Google Search

I'm Feeling Lucky

[Learn more](#)

Report inappropriate predictions



War against **Bias** in AI !



@MeMohanty

Facial recognition

How white engineers built racist code - and why it's dangerous for black people

As facial recognition tools play a bigger role in fighting crime, inbuilt racial biases raise troubling questions about the systems that create them

Ali Breland

Mon 4 Dec 2017 09:00 GMT



7,329



▲ A protest over police violence against black communities. Photograph: Alamy Stock Photo

“You good?” a man asked two narcotics detectives late in the summer of 2015.

The detectives had just finished an undercover drug deal in Brentwood, a predominately black neighborhood in Jacksonville, Florida, that is among the poorest in the country, when the man unexpectedly approached them.

most viewed



Airline pilot 'sucked halfway out' when cockpit windshield broke



EU demands action by Poland's government to protect rule of law



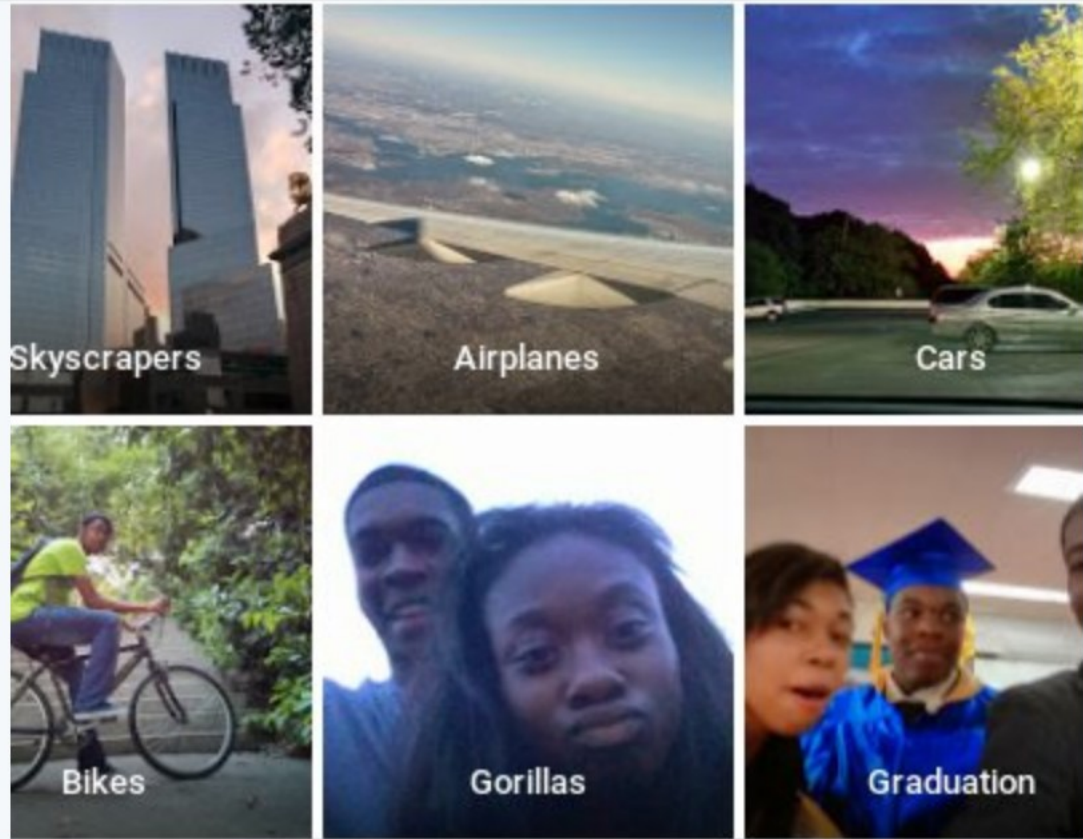
Meghan Markle's father will not attend the royal wedding-report



Iraq's shock election result may be turning point for Iran



No more romcoms for me, says 'older and uglier' Hugh Grant



jackyalciné's like 55% in the IndieWeb.
@jackyalcine



Google Photos, y'all [redacted] up. My friend's not a gorilla.

3:22 AM - Jun 29, 2015

♡ 2,274 💬 3,581 people are talking about this



@MeMohanty

Google

Google's solution to accidental algorithmic racism: ban gorillas

Google's 'immediate action' over AI labelling of black people as gorillas was simply to block the word, along with chimpanzee and monkey, reports suggest

Alex Hern

@alexhern

Fri 12 Jan 2018 16.04 GMT



395

This article is over 4 months old



▲ A silverback high mountain gorilla, which you'll no longer be able to label satisfactorily on Google Photos. Photograph: Thomas Mukoya/Reuters

After Google was criticised in 2015 for an image-recognition algorithm that auto-tagged pictures of black people as "gorillas", the company promised "immediate action" to prevent any repetition of the error.

That action was simply to prevent Google Photos from ever labelling any image as a gorilla, chimpanzee, or monkey - even pictures of the primates themselves.

That's the conclusion drawn by Wired magazine, which tested more than 40,000 images of animals on the service. Photos accurately tagged images of pandas and



@MeMohanty

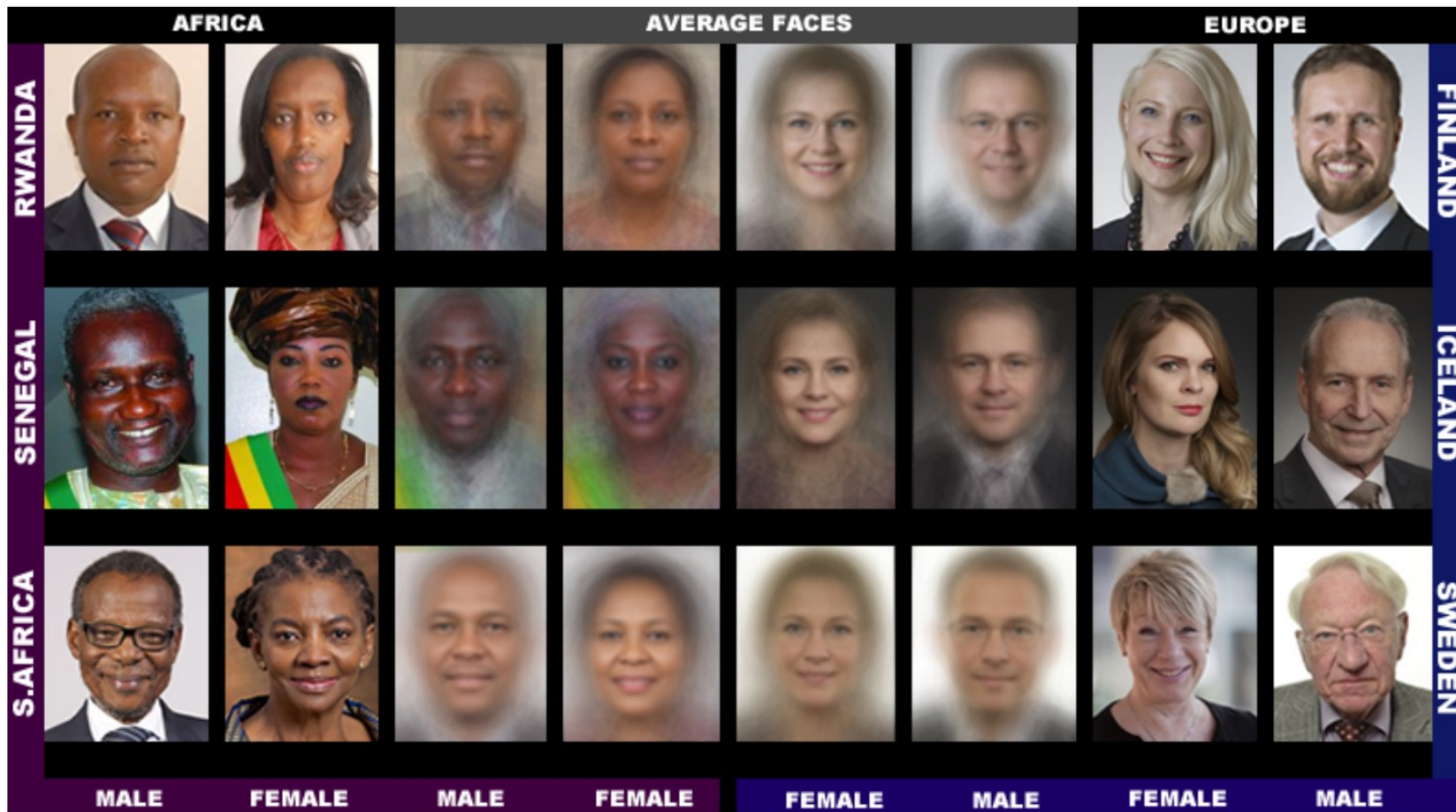


Figure 1: Example images and average faces from the new Pilot Parliaments Benchmark (PPB). As the examples show, the images are constrained with relatively little variation in pose. The subjects are composed of male and female parliamentarians from 6 countries. On average, Senegalese subjects are the darkest skinned while those from Finland and Iceland are the lightest skinned.

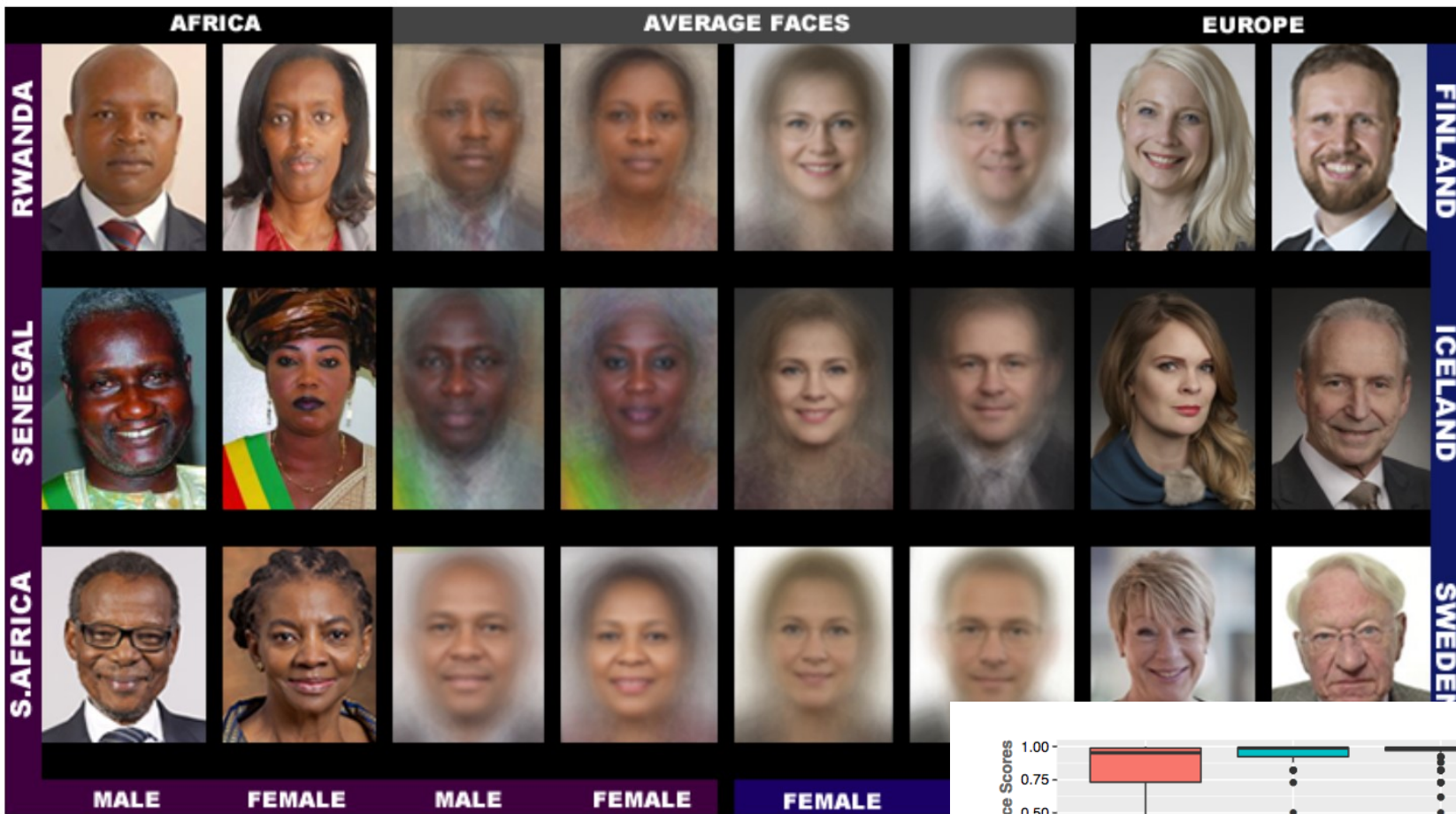


Figure 1: Example images and average faces from the new Pilot the examples show, the images are constrained with rel subjects are composed of male and female parliamenta Senegalese subjects are the darkest skinned while thos lightest skinned.

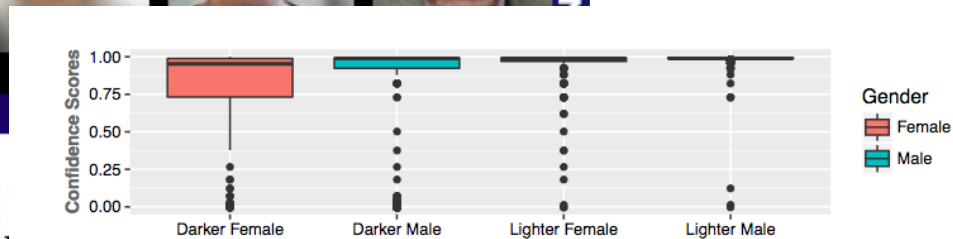


Figure 4: Gender classification confidence scores from IBM (IBM). Scores are near 1 for lighter male and female subjects while they range from $\sim 0.75 - 1$ for darker females.



GET WIRED. 3 MONTHS UNLIMITED ACCESS ON US

→ START YOUR TRIAL

SHARE

f SHARE 5280

🐦 TWEET

💬 COMMENT

✉️ EMAIL

JASON TASHEA OPINION 04.17.17 07:00 AM

COURTS ARE USING AI TO SENTENCE CRIMINALS. THAT MUST STOP NOW



GET WIRED.
3 MONTHS
UNLIMITED
ACCESS ON US

→ START YOUR TRIAL

MOST POPULAR



CULTURE
'Westworld' Is Turning Into 'Lost'—for Better or for Worse
ANGELA WATERCUTTER



TRANSPORTATION
The Vehicle of the Future Has Two Wheels, Handlebars, and Is a Bike
CLIVE THOMPSON



SCIENCE
The Psychology of Amazon's Echo Dot Kids Edition
ROBBIE GONZALEZ



@MeMohanty

When the creator is **Flawed...**



When the creator is **Flawed...**

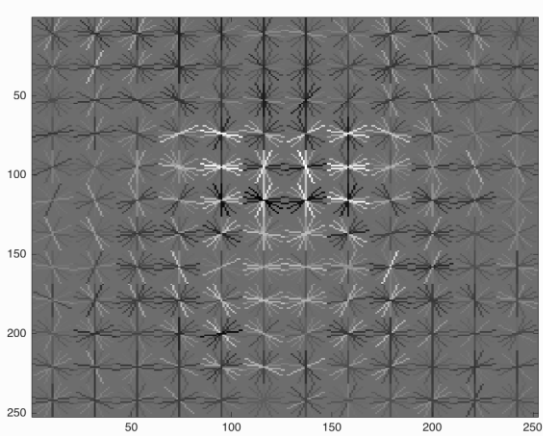
Can the creation be **perfect** ?



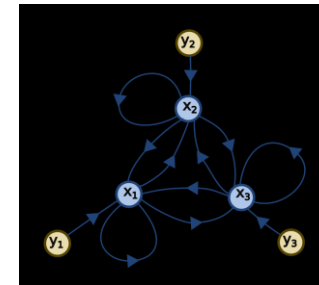
How has the creator been creating ?



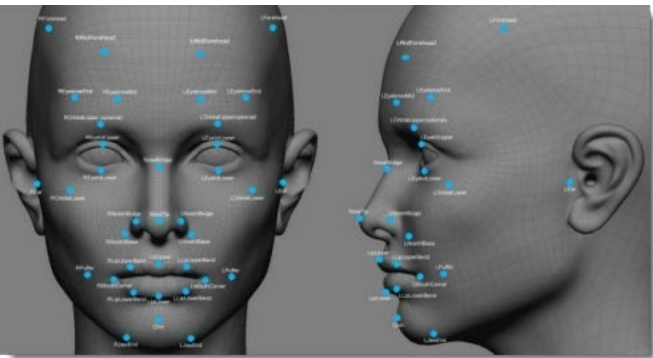
How has the creator been creating ?



Mathematical Models

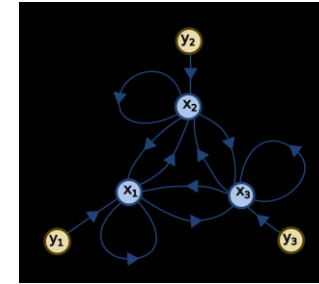


How has the creator been creating ?



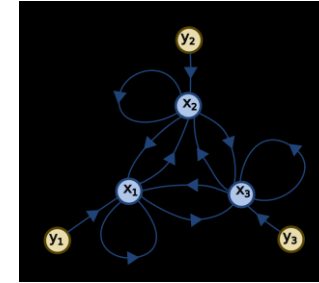
Mathematical Models

Hybrid Models



How has the creator been creating ?

Mathematical Models



Hybrid Models

Data Driven Models



Figure 5: 1024×1024 images generated using the CELEBA-HQ dataset. See Appendix F for a larger set of results, and the accompanying video for latent space interpolations. On the right, two images from an earlier megapixel GAN by Marchesi (2017) show limited detail and variation.

Image Source : <https://arxiv.org/abs/1710.10196>



@MeMohanty

So **where** is the **bias** ?



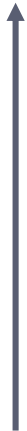
So **where** is the **bias** ?

Everywhere !

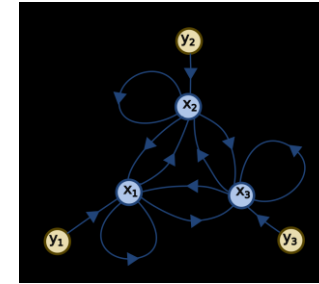


How has the creator been creating ?

Individual Biases



Mathematical Models

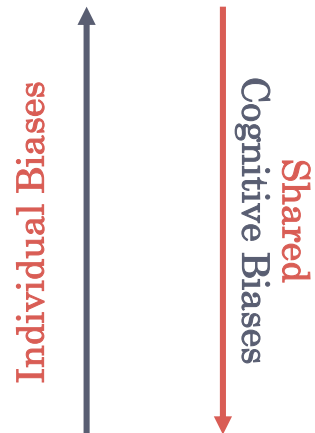


Hybrid Models

Data Driven Models



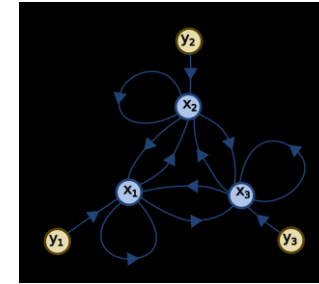
How has the creator been creating ?



Mathematical Models

Hybrid Models

Data Driven Models



How can we address this ?



How can we address this ?

- **Acknowledge** the existence of Bias



How can we address this ?

- **Acknowledge** the existence of Bias
- Clearly Define **Scope** of tools



How can we address this ?

- Acknowledge the existence of Bias
 - Clearly Define **Scope** of tools

Do not sell a chisel as a toothpick



Cancel

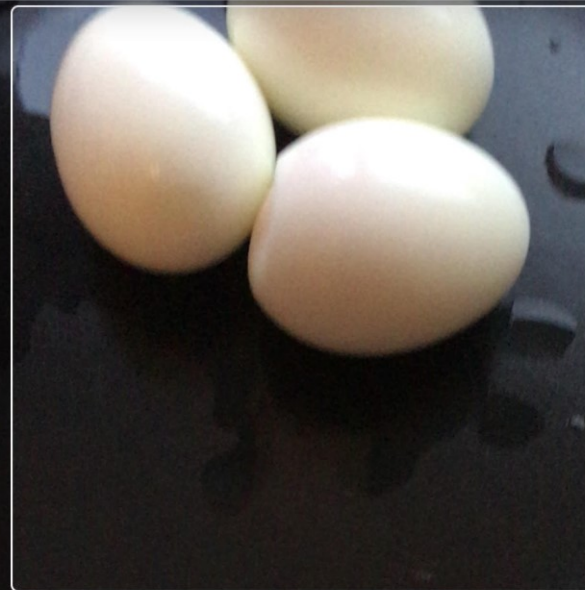


escalope 4.3%
 cotoletta milan style 3.7%
 escalopes with white wine 3.3%

Captured in 88 ms
 Cropped in 10 ms
 Prepared in 14 ms
 Recognized in 687 ms
 Total: 799 ms



jupyterhub/binder
 Isaiah Becker-Mayer: @choldgraf awesome, we are...



egg?



egg 40.4%
 egg white 19.2%
 deviled eggs 6.5%

Captured in 83 ms
 Cropped in 10 ms
 Prepared in 15 ms
 Recognized in 762 ms
 Total: 870 ms





beef rib 9%
brunli 8.5%
fig 7.3%

Captured in 101 ms
Cropped in 11 ms
Prepared in 12 ms
Recognized in 1377 ms
Total: 1501 ms



meat loaf 26%
brunli 7.9%
beef rib 6.2%

Captured in 146 ms
Cropped in 17 ms
Prepared in 17 ms
Recognized in 894 ms
Total: 1074 ms



brunli 7.3%
meat loaf 5.3%
squid 4.4%

Captured in 94 ms
Cropped in 16 ms
Prepared in 13 ms
Recognized in 821 ms
Total: 944 ms

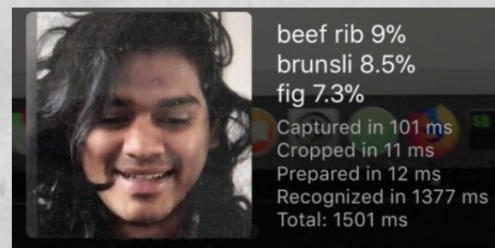


Racist AI calls a religious Indian man as a "beef rib", "meat loaf", a brunсли

By AI JONES

In a recent turn of events, a racist AI started name calling its own creator of Indian Origin a "beef rib", a "meat loaf", a "brunсли". Religious Sentiments were gravely hurt across the whole Indian subcontinent, and the creator of the AI, Sharada Mohanty, a young AI researcher based in Switzerland, is still in shock.

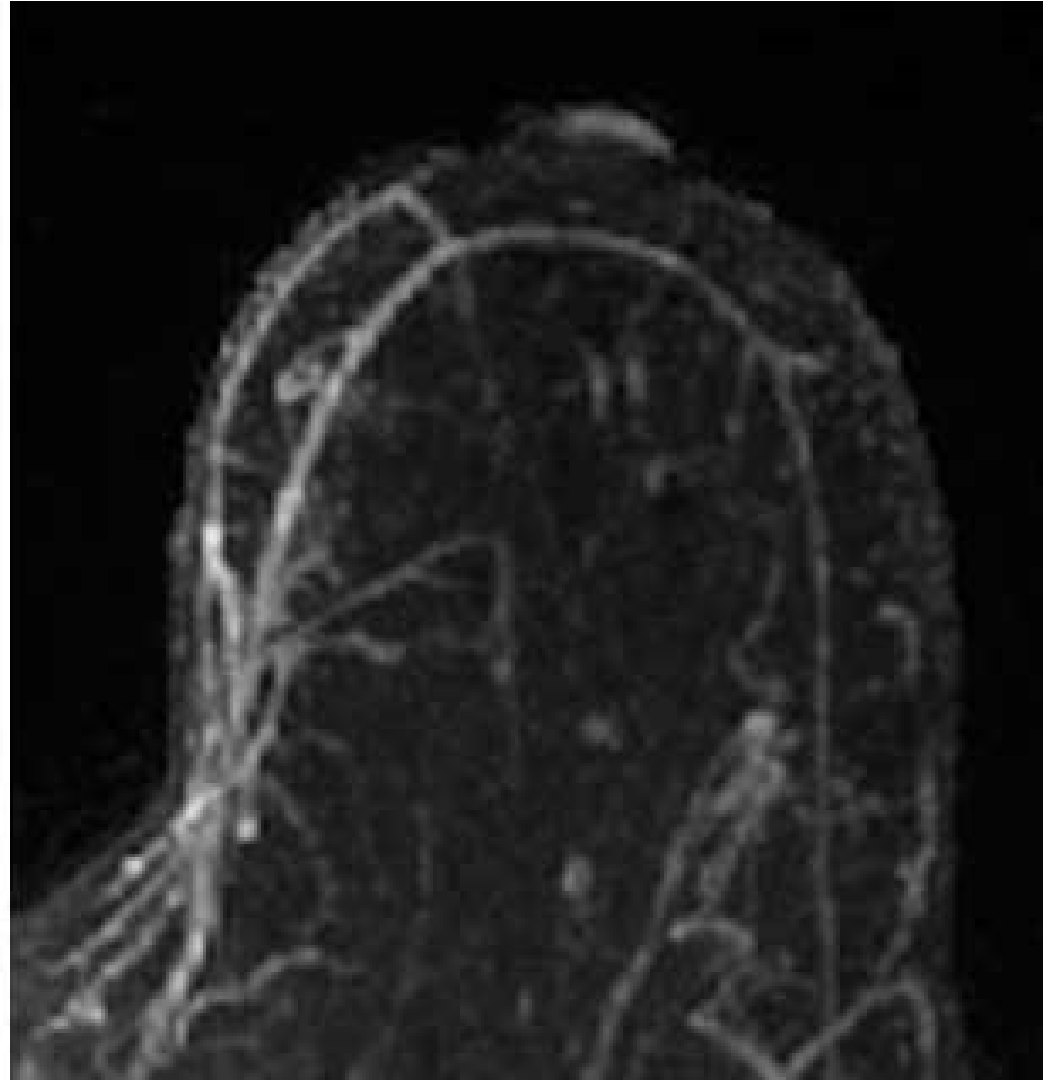
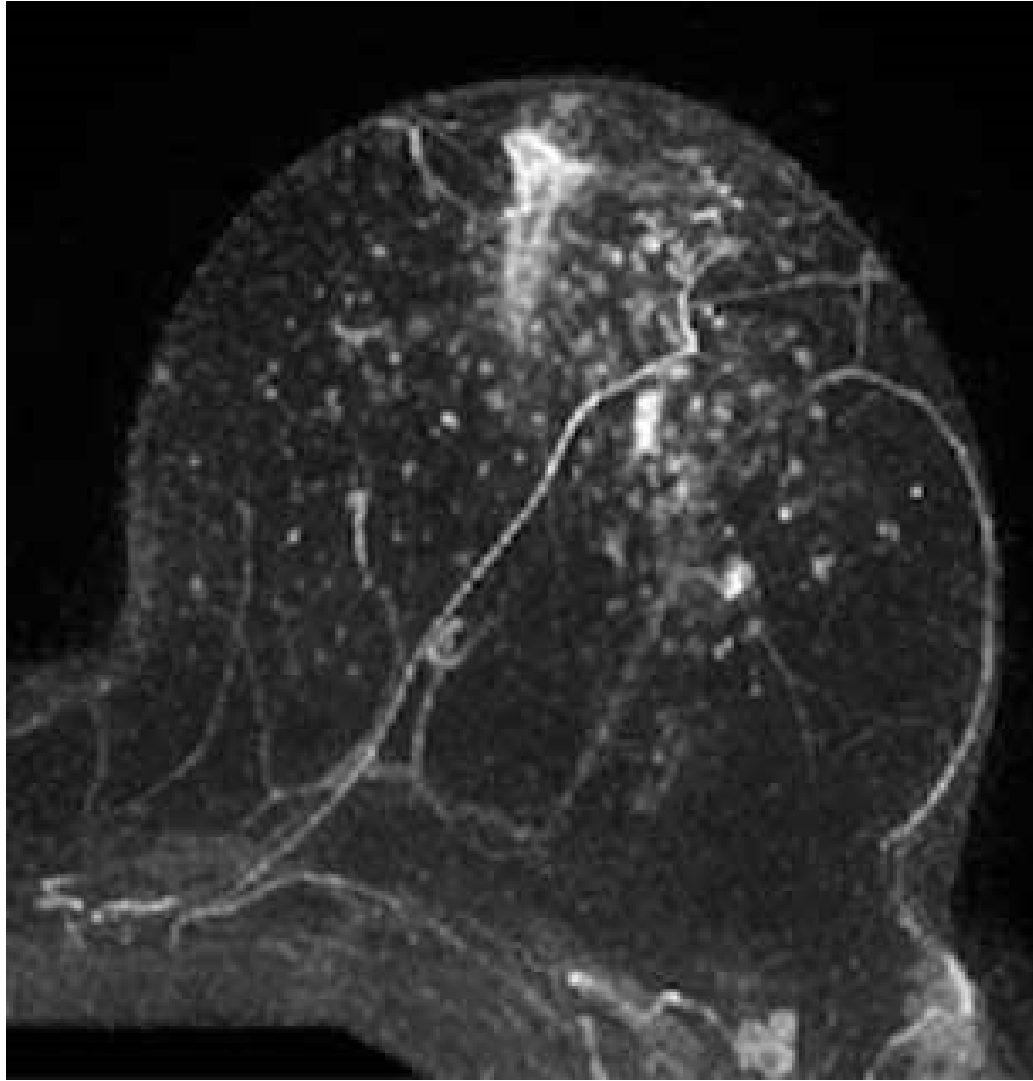
Investigators are still puzzled by this racist behavior, and warn about the imminent dangers from large scale adoption of racist AIs.



How can we address this ?

- Acknowledge the existence of Bias
 - Clearly Define Scope of tools
 - Make AI accountable





How can we address this ?

- Acknowledge the existence of Bias
 - Clearly Define Scope of tools
 - Make AI accountable
- Empower everyone to fight Bias



Solve real-world problems using open data

crowdAI connects data science experts and enthusiasts with open data to solve specific problems, through challenges.

[HOST A CHALLENGE](#)



R USERS



HALLENGES

[SEE ALL](#)



Mapping Challenge
Building Missing Maps with Machine Learning
16 days left

10.1 k Views 230 Participants 223 Submissions



NIPS 2018 : Adversarial Vision Challenge
Pitting machine vision models against adversarial attacks.
Starting soon

5081 Views 44 Participants 0 Submissions



NIPS 2018 : AI for Prosthetics Challenge
Starting soon

967 Views 21 Participants 0 Submissions



Visual Doom AI Competition 2018 - Track 1



How can we address this ?

- Acknowledge the existence of Bias
 - Clearly Define Scope of tools
 - Make AI accountable
 - Empower everyone to fight Bias
- Ensure everyone is represented well in the fight



Technological Challenge



~~Technological~~ Challenge

Cultural Challenge



Each one of you will have to
contribute

