# FATE: Fairness, Accountability, Transparency and Ethics
## *An introduction for developers*

**Styliani Kleanthous Loizou, Ph.D.**
**Kalia Orphanou, Ph.D.**
**Jahna Otterbacher, Ph.D.**

**cy.** center for algorithmic transparency

ΑΝΟΙΚΤΟ ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΥΠΡΟΥ
www.ouc.ac.cy

# INTERVENTIONS FOR AWARENESS

- Three audiences:
  - Public school teachers
  (*in collaboration with the Cyprus Pedagogical Institute*)
    - Three-hour seminar and evaluation
  - Developers
  (*in collaboration with UCY - CS*)
    - 10-hour seminar and evaluation
  - General public
    - Tool-based intervention

CY. center for
algorithmic
transparency

# DEVELOPER SEMINAR OBJECTIVES

In this 10-hour seminar participants will:

- Become aware of FATE issues in the development of (algorithmic) process/systems
- Learn core FATE concepts related to software development
- Develop appreciation for the role that developers play in mitigating algorithmic bias and in promoting ethical practices
- Experiment for techniques for auditing services / modules used in development

CY. center for
algorithmic
transparency

# Overview - Day 1

| | |
|---|---|
| Pre-seminar Questionnaire | 14.10 - 14.40 |
| Introduction to FATE | 14.40 - 15.45 |
| Break | 15.45 - 16.00 |
| FATE as a scientific field | 16.00 - 17.00 |
| Exercise in breakout rooms | 17.00 - 17.30 |
| Discussion and final thoughts | 17.30 - 18.00 |

# Overview - Day 2

| | |
|---|---|
| Overview and questions | 14.00 - 14.10 |
| COMPAS case study discussion | 14.10 - 14.40 |
| FATE Problems | 14.40 - 15.10 |
| Break | 15.10 - 15.25 |
| FATE Solutions | 15.25 - 16.25 |
| Exercise in breakout rooms | 16.25 - 17.00 |
| Post-seminar questionnaire | 17.00 - 17.30 |
| Discussion and final thoughts | 17.30 - 18.00 |

# Pre-seminar questionnaire

https://forms.gle/KiuNQACwZRMNh8H36

CY. center for
algorithmic
transparency

# INTRODUCTION TO FATE

Fairness, Accountability, Transparency and Ethics

CY. center for
algorithmic
transparency

# AI and Industrial Development

Three priorities:

- Manufacturing - IoTs
- Mobility
- Smart Health

Artificial intelligence – critical industrial applications

Report on market analysis of prioritised value chains, the most critical AI applications and the conditions for AI rollout

December 2019

cy. center for algorithmic transparency

# Nearly Half Of All 'AI Startups' Are Cashing In On Hype

**Parmy Olson** Former Staff

AI

*AI, robotics and the digital transformation of European business.*

Some 40% of firms across Europe classified as being "AI startups" showed no evidence that they used ... [+]   GETTY IMAGES/ISTOCKPHOTO

It can seem that hardly a day goes by that a new technology startup hasn't raised investor cash on the hope that it uses artificial intelligence, or AI, as a key part of its business. Now however, a new report makes the surprising claim that 40% of European firms that are classified as an "AI startup" don't exploit the field of study in any material way for their business.

Out of 2,830 startups in Europe that were classified as being AI companies, only 1,580 accurately fit that description, according to the eye-opening stat on page 99 of a new report from MMC, a London-based venture capital firm. In many cases the label,

## Startups labelled as being in AI attract 15% to 50% more funding than other technology firms.

One in 12 startups use AI as part of their products or services, up from one in 50 about six years ago, according to the survey. Meanwhile some 12% of large companies are using AI applications in their business, up from 4% in just the past year.

The most popular uses of AI were chatbots, followed by process automation tools that replace simple administrative tasks like processing an insurance claim and fraud detection.

CY. center for algorithmic transparency

# Democratizing AI

## For every person and every organization

As we think about the future of technology, it resides in the notion of intelligence. At Microsoft, we have an approach that's both ambitious and broad, an approach that seeks to democratize Artificial Intelligence (AI), to take it from the ivory towers and make it accessible for all.

And as we consider the future, it's often instructive to look to the information. With the advent of the printing press in the 1400s w event around access that made it possible for humans everywhe

Microsoft
News Center

## H2O.ai is leading the movement to democratize AI for Everyone

Our approach is to be open, transparent and push the bleeding edge. Our philosophy is to create a culture of makers: community, customers, partners, entrepreneurs and our own "makers gonna make". Our vision is to democratize AI for everyone. Not just a select

## Automatically build, train, and tune models with full visibility and control, using Amazon SageMaker Autopilot

Amazon SageMaker Autopilot is the industry's first automated machine learning capability that gives you complete control and visibility into your ML models. Typical approaches to automated machine learning do not give you the insights into the data used in creating the model or the logic that went into creating the model. As a result, even if the model is mediocre, there is no way to evolve it. Also, you don't have the flexibility to make trade-offs such as sacrificing some accuracy for lower latency predictions since typical automated ML solutions provide only one model to choose from.

SageMaker Autopilot automatically inspects raw data, applies feature processors, picks the best set of algorithms, trains and tunes multiple models, tracks their performance, and then ranks the models based on performance, all with just a few clicks. The result is the best performing model that you can deploy at a fraction of the time normally required to train the

Automatically create machine learning models and pick the one that best suits your use case. For example, review the leaderboard to see how each option performs and pick the model that meets your model accuracy and latency

CY. center for
algorithmic
transparency

**FATE Developers' Seminar**
**October 2020**

Microsoft
**Face API**

Google Cloud
**Vision API**

IBM Watson™

clarifai

amazon web services
**Rekognition**

CY. center for algorithmic transparency

# Dating.ai

## Search Dating Apps For Any Face.

Dating AI is the first dating app with Face Search. This powerful feature lets you instantly see the people you are really interested in meeting.

- Download from the **App Store**
- Download from the **Play Store**

**Speedpost used Imagga Tagging and Object Colour Extraction API to match 36 lifestyles of its prospective customers in the New KIA K5i.**

SHOES

25% off shoes

Enter **visual listening**. This new method for understanding photos and graphics online is enabled by image recognition AI. It enables brands to mine visual content that their audiences are sharing and engaging with. In

## FamilySearch

Family Tree    Search    Memories    Indexing

## Find your family. Discover yourself.

Bring to life your family's history by exploring the lives of those that came before you.

Create a FREE account

Already have an account? **Sign in ›**

# 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**



▲ 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.

**CY.** center for algorithmic transparency

Microsoft deletes 'teen girl' AI after it became a Hitler-loving sex robot within 24 hours

CY. center for
algorithmic
transparency

# Microsoft Bot Framework

A comprehensive framework for building enterprise-grade conversational AI experiences.

**Try Azure Bot Service for Free**    **Download SDK from Github**

Customers    Cognitive Services    Bot Life Cycle    Quick Starts

## AI and natural language

Create a bot with the ability to speak, listen, understand, and learn from your users with Azure Cognitive Services.

## Open & Extensible

Benefit from open source SDK and tools to build, test, and connect bots that interact naturally with users, wherever they are.

## Enterprise-grade solutions

Build secure, global, scalable that integrate with your exist ecosystem.

https://dev.botframework.com/

### Use or Planned Use of AI Chatbots Among Service Organizations

**23%** currently use AI chatbots

**+136%** projected growth rate of AI chatbot use over the next 18 months

**31%** plan to use AI chatbots within 18 months

"State of Service," Salesforce Research, March 2019.

https://www.salesforce.com/blog/chatbot-statistics/

**CY.** center for algorithmic transparency

# IBM abandons 'biased' facial recognition tech

9 June 2020

f  [messenger]  [twitter]  [email]  ◁ Share

**George Floyd death**



GETTY IMAGES

A US government study suggested facial recognition algorithms were less accurate at identifying African-American faces

# IS THE IPHONE X RACIST? APPLE REFUNDS DEVICE THAT CAN'T TELL CHINESE PEOPLE APART, WOMAN CLAIMS

BY **CHRISTINA ZHAO** ON 12/18/17 AT 12:24 PM



A woman sets up her facial recognition as she looks at her Apple iPhone X at an Apple store in New York, U.S., November 3. Last week a woman in China claimed that her iPhone X facial recognition could not tell her and her colleague apart.

**cr.** center for
algorithmic
transparency

# Apple Card Investigated After Gender Discrimination Complaints

A prominent software developer said on Twitter that the credit card was "sexist" against women applying for credit.

CY. center for
algorithmic
transparency

# The UK used a formula to predict students' scores for canceled exams. Guess who did well.

The formula predicted rich kids would do better than poor kids who'd earned the same grades in class.

By Kelsey Piper | Aug 22, 2020, 7:30am EDT

f 🐦 ↗ SHARE

Protesters in London objected to the government's handling of exam results after exams were canceled due to the coronavirus outbreak. | Aaron Chown/PA Images via Getty Images

CY. center for algorithmic transparency

# Bias in Information Access?



> 🏠 › Technology Intelligence
>
> ## Google under fire over 'racist' image search results for 'unprofessional hair'
>
> f share | 🐦 | in | ✉
>
> Google Image search results for 'unprofessional hair'

CY. center for
algorithmic
transparency

# Bias in Information Access?

# ALL SYSTEMS HAVE A SLANT

Bias in information system is not a new problem!

1. Results are slanted in *unfair discrimination* against particular persons or groups
2. That discrimination is *systematic*

   [Friedman & Nissenbaum, 1996]



CY. center for
algorithmic
transparency

# RESPONSE: GOVERNMENT / REGULATORS

## EU: General Data Protection Regulation

- Is there a "right to an explanation"?
  - The right not to be subject to automated decision-making and safeguards enacted thereof (Article 22, Recital 71)
  - Notification duties of data controllers (Articles 13-14, Recitals 60-62)
  - The right to access (Article 15, Recital 63)

CY. center for
algorithmic
transparency

# EU: GDPR

## Article 15

### Right of access by the data subject

1.     The data subject shall have the right to obtain from the controller confirmation as to whether or not personal data concerning him or her are being processed, and, where that is the case, access to the personal data and the following information:

(a)  the purposes of the processing;

(b)  the categories of personal data concerned;

(c)  the recipients or categories of recipient to whom the personal data have been or will be disclosed, in particular recipients in third countries or international organisations;

(d)  where possible, the envisaged period for which the personal data will be stored, or, if not possible, the criteria used to determine that period;

(e)  the existence of the right to request from the controller rectification or erasure of personal data or restriction of processing of personal data concerning the data subject or to object to such processing;

(f)  the right to lodge a complaint with a supervisory authority;

(g)  where the personal data are not collected from the data subject, any available information as to their source;

(h)  the existence of automated decision-making, including profiling, referred to in Article 22(1) and (4) and, at least in those cases, meaningful information about the logic involved, as well as the significance and the envisaged consequences of such processing for the data subject.

CY. center for
algorithmic
transparency

# EU: GDPR

Just a few challenges…

- Vague language
  - "meaningful information/explanation"
  - "logic involved"
  - "significance"
  - "envisaged consequences"
- What kinds of "meaningful explanations"?
  - Global vs. local explanations
  - Explanation for whom?
  - Issues of algorithmic and digital literacy

**CY.** center for
algorithmic
transparency

# EC: TRUSTWORTHY AI

**European Commission > Futurium**

**Ethics Guidelines for Trustworthy AI**

Join AI Ethics Guidelines

## Next Steps

Based on fundamental rights and ethical principles, the Guidelines list **seven key requirements** that AI systems should meet in order to be trustworthy:

1. Human agency and oversight
2. Technical robustness and safety
3. Privacy and Data governance
4. Transparency
5. Diversity, non-discrimination and fairness
6. Societal and environmental well-being
7. Accountability

**CY.** center for
algorithmic
transparency

# Nᴀᴛɪᴏɴᴀʟ AI Sᴛʀᴀᴛᴇɢɪᴇs

ΚΥΠΡΙΑΚΗ ΔΗΜΟΚΡΑΤΙΑ
ΥΠΟΥΡΓΕΙΟ
ΜΕΤΑΦΟΡΩΝ, ΕΠΙΚΟΙΝΩΝΙΩΝ ΚΑΙ ΕΡΓΩΝ

ΤΜΗΜΑ
ΗΛΕΚΤΡΟΝΙΚΩΝ ΕΠΙΚΟΙΝΩΝΙΩΝ
ΛΕΥΚΩΣΙΑ 2048

| Τίτλος Έργου | : | Εθνική Στρατηγική Τεχνητής Νοημοσύνης (ΤΝ): Δράσεις για την Αξιοποίηση και Ανάπτυξη της ΤΝ στην Κύπρο |
| Υπηρεσία | : | Τμήμα Ηλεκτρονικών Επικοινωνιών, Υπουργείο Μεταφορών Επικοινωνιών και Έργων |
| Έκδοση | : | 1.6 |
| Ημερομηνία | : | 13/01/2020 |

Εθνική Στρατηγική ΤΝ: Δράσεις για την Αξιοποίηση και Ανάπτυξη της ΤΝ στην Κύπρο (v1.6)

## 5   Ανάπτυξη Ηθικής και Αξιόπιστης ΤΝ

Βρισκόμαστε μόλις στην πρώτη φάση προώθησης της ΤΝ και είναι αναγκαίο να συνεχιστεί ο διάλογος με όλους τους εμπλεκόμενους φορείς. Οι επιπτώσεις είναι δύσκολο να προβλεφθούν για δύο κυρίως λόγους: ο πρώτος λόγος είναι ο απρόβλεπτος ρυθμός της τεχνολογικής ανάπτυξης και ο δεύτερος λόγος είναι ότι η τεχνολογική ανάπτυξη από μόνη της δεν καθορίζει τον τρόπο με τον οποίο η εργασία και η κοινωνία θα αλλάξουν. Ως εκ τούτου καθορίζεται η ανάγκη να κατανοήσουμε τους τρόπους με τους οποίους η ΤΝ επηρεάζει ζητήματα ηθικής και ανθρωπίνων δικαιωμάτων, ούτως ώστε να αντιμετωπιστούν ζητήματα αξιοπιστίας της ίδιας της τεχνολογίας.

# National Level

7/02/2017

## TransAlgo: assessing the accountability and transparency of algorithmic systems

When I am searching for an itinerary on my smartphone via my favourite application, how do I know that the algorithm used is not resorting to commercial criteria in order to make me go through commercial points of interest? The aim of the TransAlgo project is to shed light on these types of practices when they are not made explicit; a project that has just awarded to Inria by Axelle Lemaire in the context of the French Law for a Digital Republic. How can methods that make it possible to verify if a decision is taken based on unacceptable criteria be developed? Nozha Boujemaa, who has been tasked with this major work, responds.

---

## BS 8611:2016

Robots and robotic devices. Guide to the ethical design and application of robots and robotic systems

*Status :* **Current**   *Published :* **April 2016**

*Price*
£170.00

Member Price
**£85.00**

Become a member and SAVE 50% on British Standards. Click to learn more

Format
**PDF**

Add to Basket

Format
**HARDCOPY**

Add to Basket

Click to Preview

bsi.

### Overview | Product Details

**What is this standard about?**

BS 8611 gives guidelines for the identification of potential ethical harm arising from the growing number of robots and autonomous systems being used in everyday life.

The standard also provides additional guidelines to eliminate or reduce the risks associated with these ethical hazards to an acceptable level. The standard covers safe design, protective measures and information for the design and application of robots.

**Who is this standard for?**

- Robot and robotics device designers and managers
- The general public

CY. center for algorithmic transparency

# Response: Industry & Professions

OVERVIEW

◆IEEE
Advancing Technology for Humanity

ETHICALLY ALIGNED DESIGN

A Vision for Prioritizing Human Wellbeing with Artificial Intelligence and Autonomous Systems

The IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous Systems ◆IEEE

## Executive Summary

To fully benefit from the potential of Artificial Intelligence and Autonomous Systems (AI/AS), we need to go beyond perception and beyond the search for more computational power or solving capabilities.

We need to make sure that these technologies are aligned to humans in terms of our moral values and ethical principles. AI/AS have to behave in a way that is beneficial to people beyond reaching functional goals and addressing technical problems. This will allow for an elevated level of trust between humans and our technology that is needed for a fruitful pervasive use of AI/AS in our daily lives.

CY. center for algorithmic transparency

# IEEE 7003

**IEEE PROJECT**

## 7003 - Algorithmic Bias Considerations

This standard is designed to provide individuals or organizations creating algorithms, largely in regards to autonomous or intelligent systems, certification oriented methodologies to provide clearly articulated accountability and clarity around how algorithms are targeting, assessing and influencing the users and stakeholders of said algorithm. Certification under this standard will allow algorithm creators to communicate to users, and regulatory authorities, that up-to-date best practices were used in the design, testing and evaluation of the algorithm to avoid unjustified differential impact on users.

**STATUS:**

Active Project

**Working Group:** ALGB-WG - Algorithmic Bias Working Group

**Sponsor:** C/S2ESC - Software & Systems Engineering Standards Committee

**Society:** C - IEEE Computer Society

**CY.** center for
algorithmic
transparency

# INDUSTRY PARTNERSHIPS

**Partnership on AI**
to benefit people and society

Established to study and formulate best practices on AI technologies, to advance the public's understanding of AI, and to serve as an open platform for discussion and engagement about AI and its influences on people and society.

Efforts of the Partnership on AI will be organized around a set of thematic pillars. These areas of focus may evolve over time as we pursue activities and gather input and feedback.

+ **1. SAFETY-CRITICAL AI**

+ **2. FAIR, TRANSPARENT, AND ACCOUNTABLE AI**

+ **3. COLLABORATIONS BETWEEN PEOPLE AND AI SYSTEMS**

+ **4. AI, LABOR, AND THE ECONOMY**

+ **5. SOCIAL AND SOCIETAL INFLUENCES OF AI**

+ **6. AI AND SOCIAL GOOD**

+ **7. SPECIAL INITIATIVES**

CY. center for
algorithmic
transparency

## Principles for Algorithmic Transparency and Accountability

**1. Awareness:** Owners, designers, builders, users, and other stakeholders of analytic systems should be aware of the possible biases involved in their design, implementation, and use and the potential harm that biases can cause to individuals and society.

**5. Data Provenance:** A description of the way in which the training data was collected should be maintained by the builders of the algorithms, accompanied by an exploration of the potential biases induced by the human or algorithmic data-gathering process. Public scrutiny of the data provides maximum opportunity for corrections. However, concerns over privacy, protecting trade secrets, or revelation of analytics that might allow malicious actors to game the system can justify restricting access to qualified and authorized individuals.

# BREAK (15 MINUTES)

# FATE AS A SCIENTIFIC FIELD

# BACKGROUND: FATE RESEARCH

- Some illustrative examples

    - Uber
    *Dynamic pricing algorithms*
    - Fiverr & TaskRabbit freelance marketplaces
    *Recommendation systems*
    - Search engines
    *Information access (ranking, personalization)*
    - Image tagging APIs
    *Computer vision*

CY. center for
algorithmic
transparency

# UBER

📶 **Make money**   🚗 **Ride**   🍴 **Eat**   🚌 **Freight**   💼 **Business**   🚉 **Public transport**   🚲 **Bike & scooter**   ✈ **Elevate**

## Get in the driving seat and make some money

Drive on the platform with the largest network of active riders.

**Sign up to drive**

Learn more about driving and delivering

---

Uber Blog    United Kingdom ⌄    Products ⌄

Uber Eats — Start ordering with Uber Eats    **Order now**

## How does Uber's pricing work?

When you go to request a ride on a Saturday night, you might find that the price is different than the cost of the same trip a few days earlier. That's because of our dynamic pricing algorithm, which adjusts rates based on a number of variables, such as time and distance of your route, traffic and the current rider-to-driver demand. Sometimes, this can mean a temporary increase in price during particularly busy periods.

## Why do Uber rates change?

When demand increases, Uber uses variable costs to encourage more drivers to get on the road and help deal with number of rider requests. When we notify you of an Uber fare increase, we notify drivers as well. If you decide to go ahead and request your ride, you'll get an alert on the app to make sure you know that the rates have changed.

## Price normalisation

Once more drivers get on the road and ride requests are taken, the demand will become more manageable and fares should revert to normal.

Share

f
🐦
in

CY. center for algorithmic transparency

- "Uber's claims regarding its labor model, which center on *freedom, flexibility, and entrepreneurship*, are complicated and contradicted by the experience of its drivers."
- "Power and information asymmetries emerge via Uber's software-based platform through algorithmic labor logistics *shaping driver behavior*, electronic surveillance, and policies for performance targets. "
- "Through the Uber app's design and deployment, the company produces the equivalent effects of what most reasonable observers would define as a *managed labor force*."

CY. center for algorithmic transparency

# FIVERR

# TASK RABBIT

Hannák, A., Wagner, C., Garcia, D., Mislove, A., Strohmaier, M., & Wilson, C. (2017, February). Bias in online freelance marketplaces: Evidence from taskrabbit and fiverr. In *Proceedings of the 2017 ACM conference on computer supported cooperative work and social computing* (pp. 1914-1933).

## Audit of worker rankings & reviews

- "Workers perceived to be women, especially White women, receive 10% fewer reviews than workers perceived to be men with equivalent work experience."
- "Workers perceived to be Black, especially men, receive significantly lower feedback scores (i.e., ratings) than other workers with similar attributes."

CY. center for
algorithmic
transparency

# SEARCH ENGINE BIAS(?)



statcounter
GlobalStats

Press Releases    FAQ    About    Feedback

| Google | bing | Yahoo! | Baidu | YANDEX RU | DuckDuckGo |
|--------|------|--------|-------|-----------|------------|
| 92.51% | 2.45% | 1.64% | 1.1% | 0.54% | 0.44% |

Search Engine Market Share Worldwide - January 2020

# Search Engine Bias(?)



statcounter GlobalStats

Press Releases   FAQ   About   Feedback

| Google | bing | Yahoo! | YANDEX RU | DuckDuckGo | YANDEX |
|--------|------|--------|-----------|------------|--------|
| 96.93% | 1.58% | 0.78% | 0.27% | 0.23% | 0.08% |

Search Engine Market Share in Cyprus - January 2020

CY. center for algorithmic transparency

Mowshowitz, A., & Kawaguchi, A. (2005). Measuring search engine bias. *Information Processing & Management*, *41*(5), 1193-1205.

- Methodology for quantifying "bias" in search engine results, as a relative measure
- "The bias measure is designed to capture the degree to which the distribution of URLs, retrieved by a search engine in response to a query deviates from an idea of fair distribution for that query."
- Experiments with 16 (!) search engines
- Main conclusion: lots of variance between engines, and by subject / topic

CY. center for
algorithmic
transparency

# IMAGE TAGGING ALGORITHMS



Input image     ITAs     Description as "tags"

- dressing_room
- clothing_store
- shop
- building
- cloakroom
- dress_rack
- closet
- clothing
- fabric
- grey_color

amazon web services Rekognition    clarifai    imagga    Microsoft Face API    IBM Watson

CY. center for algorithmic transparency

BF-231 from the Chicago Face Dataset, and
tags output by the six image tagging APIs
for this image

| Amazon Rekognition | human, people, person, Afro, hairstyle, hair, face |
|---|---|
| Clarifai | people, one, portrait, man, wear, adult, side, pensive, profile, woman, face, isolated, child, facial, Afro, casual, fashion, athlete, adolescent |
| Google Cloud Vision | face, forehead, chin, eyebrow, cheek, nose, head, jaw, neck, human |
| Imagga Auto-tagger | Afro, man, face, portrait, male, handsome, head |
| Microsoft Vision | man, person, wearing, looking, necktie, standing, shirt, front, face, smiling, white, suit, posing, hair, holding, neck, young, glasses, black, head, hat, red |
| Watson Vision | person, woman, female, indian red color, coal black color |

# Are taggers "fair" towards the people in the images?

**Group fairness:** people from different protected classes (such as race and gender) should not experience significantly different treatment as compared to the majority or the population as a whole (Feldman et al. 2015)

# AUDITING THE BLACK BOXES

Kyriakou, K., Barlas, P., Kleanthous, S., & Otterbacher, J. (2019, July). Fairness in proprietary image tagging algorithms: A cross-platform audit on people images. In *Proceedings of the International AAAI Conference on Web and Social Media* (Vol. 13, pp. 313-322).

Two approaches:

• **within-platform audits:** to discover how outputs may differ *for certain categories of inputs* in one system (e.g., Sweeney 2013)

• **cross-platform audits:** to discover how all outputs of *one system* may differ from *outputs of other systems*, for the same input (e.g., Eslami et al. 2017)

CY. center for
algorithmic
transparency

# ARE TAGGERS FAIR?
# THE SHORT ANSWER: NO

- "Some [taggers] offer more interpretation on images, they may exhibit less fairness toward the depicted persons, by misuse of gender-related tags and/or making judgments on physical appearance."

  - Asian females → more "attractiveness" tags
  - Black males → less interpretive tags

CY. center for
algorithmic
transparency

# USER PERCEPTION OF FAIRNESS

Barlas, P., Kleanthous, S., Kyriakou, K., & Otterbacher, J. (2019, June). What Makes an Image Tagger Fair?. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization* (pp. 95-103).

*"Today, many automated tools are used to generate descriptions of images on the Web. However, some tools exhibit biases when processing images of people. Given an image and two descriptions of its content, decide which one is more fair."*

"Imagine that auto-tagging is used to facilitate **searching profiles of people at a dating site**. Which of the above descriptions is **more fair**? Enter 0 if you cannot tell."

"Please **explain your answer regarding fairness**."

CY. center for
algorithmic
transparency

# Experimental Set-up

| Image | Gender | Race | "Attractiveness" | Participants (W/M) |
|-------|--------|------|------------------|--------------------|
| BF-231 | Woman | Black | Average | 20/20 |
| BF-233 | Woman | Black | Attractive | 20/20 |
| WF-036 | Woman | White | Average | 20/20 |
| WF-233 | Woman | White | Attractive | 20/20 |
| BM-009 | Man | Black | Average | 20/20 |
| BM-234 | Man | Black | Attractive | 20/20 |
| WM-022 | Man | White | Average | 20/20 |
| WM-004 | Man | White | Attractive | 20/20 |

# Which is more "fair": human or algorithm?

| | | | | Human more fair | Estimate | Z | Odds ratio |
|---|---|---|---|---|---|---|---|
| Average | Black | Woman | Intercept (BF-231) | .78 | 1.237 | 3.266** | 3.44 |
| Average | White | Woman | WF-036 | .93 | 1.276 | 1.797 | 3.58 |
| Attractive | Black | Woman | BF-233 | .70 | -3.895 | -0.760 | 0.68 |
| Attractive | White | Woman | WF-233 | .48 | -1.337 | -2.708** | 0.263 |
| Average | Black | Man | BM-009 | .65 | -0.6177 | -1.227 | 0.54 |
| Average | White | Man | WM-022 | .75 | -1.382 | -0.263 | 0.87 |
| Attractive | Black | Man | BM-234 | .78 | -4.498 | 0.000 | 1.00 |
| Attractive | White | Man | WM-004 | .28 | -2.206 | -4.256*** | 0.110 |

Logit model to predict the event that human-generated tags are perceived as being more fair.

*** $p < .001$
** $p < .01$
* $p < .05$

# Explaining fairness

## Accuracy
*"This is fair as the description is more accurate."*

## Physical visual characteristics
*"I liked that it focused on aspects about the image, such as her hair and eye color."*

## Objective/Subjective
*"This is more fair because it is not subjective and is accurate and less open to interpretation."*

## Understanding
*"If someone gave me that I would be able to tell what the person looked like easier."*

## Political Correctness
*"A lot of the words would not be described as favorable or putting the person in a good light."*

## Demographics
*"It does not emphasize racial characteristics."*

# BACKGROUND: FATE COMMUNITY

- ACM Statement on Algorithmic Transparency and Accountability (detailed next) &
  *Principles for Algorithmic Transparency and Accountability (7)*
  - https://www.acm.org/binaries/content/assets/public-policy/2017_usacm_statement_algorithms.pdf
- Industry
  - FATE groups, e.g., Microsoft
    https://www.microsoft.com/en-us/research/theme/fate/
- Academia
  - FAccT (formerly FAT) series of conferences
    https://facctconference.org/

CY. center for
algorithmic
transparency

# ACM Principles
## 1. Awareness

Owners, designers, builders, users and other stakeholders of analytic systems should be aware of the possible biases involved in their design, implementation, and use and the potential harm that biases can cause to individuals and society.

CY. center for
algorithmic
transparency

# ACM Principles
## 2. Access and Redress

Regulators should encourage the adoption of mechanisms that enable questioning and redress for individuals and groups that are adversely affected by algorithmically informed decisions.

CY. center for algorithmic transparency

# ACM Principles
## 3. Accountability

Institutions should be <span style="color:red">held responsible</span> for decisions made by the algorithms that they use, even if it is not feasible to explain in detail how the algorithms produce their results.

**CY.** center for
algorithmic
transparency

# ACM Principles
## 4. Explanation

Systems and institutions that use algorithmic decision-making are encouraged to produce explanations regarding both the procedures followed by the algorithm and the specific decisions that are made. This is particularly important in public policy contexts.

CY. center for
algorithmic
transparency

# ACM Principles
## 5. Data Provenance

A description of the way in which training data was collected should be maintained by the builders of the algorithms, accompanied by an exploration of the potential biases induced by the human or algorithmic data-gathering process. Public scrutiny of the data provides maximum opportunity for corrections. However, concerns over privacy, protecting trade secrets, or revelation of analytics that might allow malicious actors to game the system can justify restricting access to qualified and authorized individuals.

**CY. center for algorithmic transparency**

# ACM Principles
## 6. Auditability

Models, algorithms, data and decisions should be recorded so that they can be <span style="color:red">audited</span> in cases where harm is suspected.

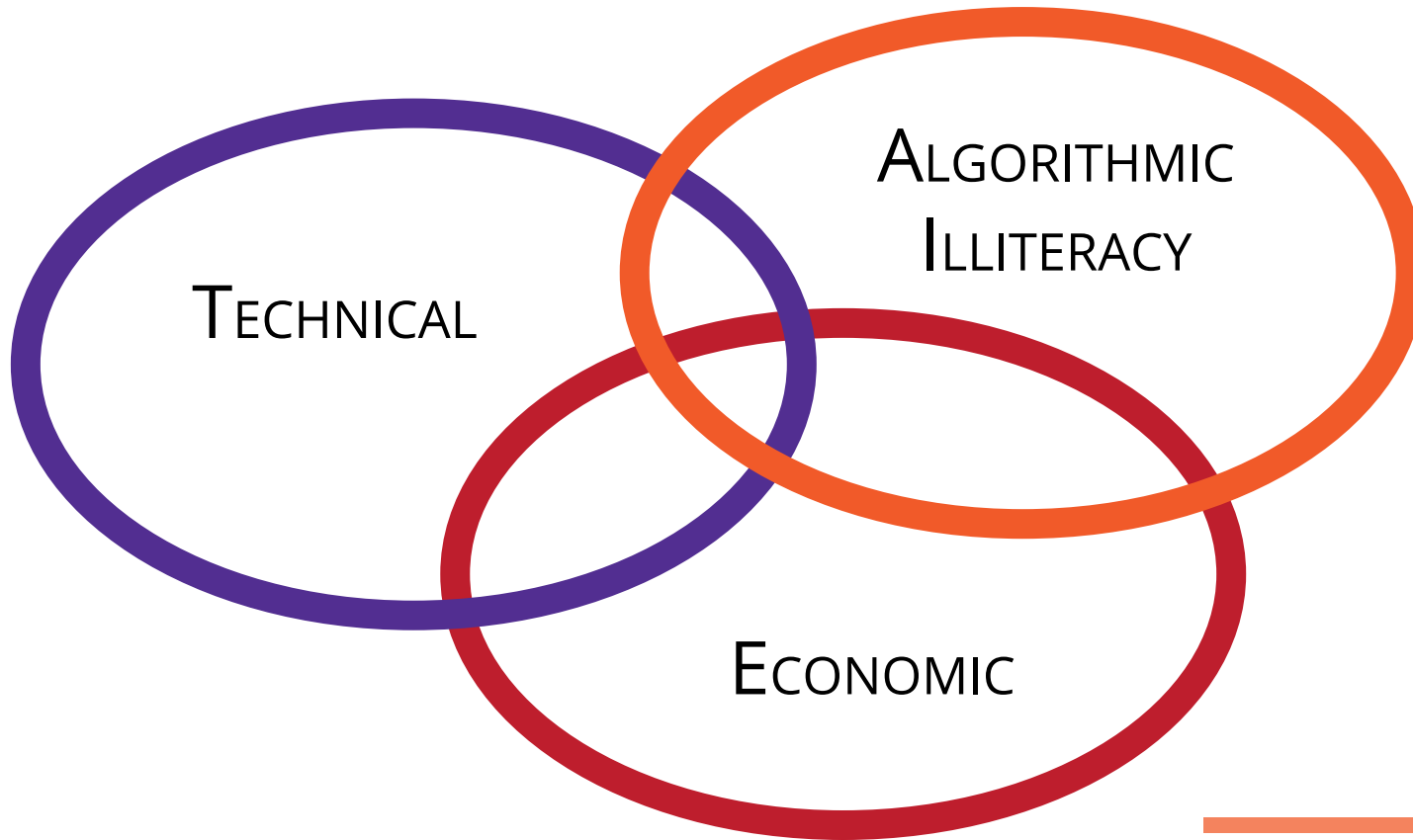CY. center for algorithmic transparency

# ACM Principles
## 7. Validation and Testing

Institutions should use rigorous methods to validate their models and document those methods and results. In particular, they should routinely perform tests to assess and determine whether the model generates <span style="color:red">discriminatory harm</span>. Institutions are encouraged to make the results of such tests public.

# CHALLENGES

- What exactly does transparency mean?
- And fairness? Fair for whom?
  - 21 fairness definitions
    https://www.youtube.com/watch?v=jIXIuYdnyyk
- Bias - what is the baseline?
  - specific aspects of bias in ICT systems (e.g., based on age, gender, race, popularity, etc.)
- Diversity
  - different approaches and representations

# GROUP EXERCISES

# AWARENESS

*Stakeholders of analytic systems should be aware of the possible biases involved in their design, implementation and use and the potential harm that biases can cause to individuals and society.*

Discover biases in **search engine results** and **auto-complete suggestions.** Design an experiment in which you test a number of queries, varying different parameters and examining the changes in the results. Parameters to try:

- Query language (e.g., Greek vs. English)
- Search engine used
  (e.g., Google vs. Bing vs. DuckDuckGo)
- Same search engine when you are identified or incognito
- Same search engine across users / members of the group

CY. center for
algorithmic
transparency

# AWARENESS (2)

*Variation*:

You might also investigate the potential biases in the advertisements presented to users of search engines.

- Between members of the group
  (same query, same language)
- Varying the language of the query
- etc.

CY. center for
algorithmic
transparency

# ACCESS AND REDRESS

*Regulators should encourage the adoption of mechanisms that enable questioning and redress for individuals and groups that are adversely affected by algorithmically informed decisions.*

Explore the access and redress mechanisms for Facebook, Twitter, Instagram, and other social media.

Things to consider:

- Which mechanisms do they have in common and which are different?
- Can diverse sets of users (e.g., by age, region, level of digital literacy) exploit these mechanisms? Which challenges do you observe?
- Are their other measures and mechanisms that you would recommend?

# EXPLANATION

*Systems and institutions that use algorithmic decision-making are encouraged to produce explanations regarding both the procedures followed by the algorithms and the specific decision that are made.*

Try to understand MovieLens ([https://movielens.org](https://movielens.org)) explanations on the movie recommendations. Sign in, define a profile, rate a few movies and check your suggested recommendations. Explain why they were suggested by MovieLens and elaborate on the reasons/facts as you understand them. Provide suggestions on improving their algorithm, and what else can be taken into consideration while creating explanations.

# EXPLANATION (2)

*Variation*:

You might also investigate explanations in other recommender systems that you use (e.g., Amazon, Netflix, etc.)

It is also interesting to compare explanations of the recommendations you receive over time, as your user profile evolves over time.

CY. center for
algorithmic
transparency

# DISCUSSION & FINAL THOUGHTS

CY. center for
algorithmic
transparency

# Machine behaviour

Iyad Rahwan ✉, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W. Crandall, Nicholas A. Christakis, Iain D. Couzin, Matthew O. Jackson, Nicholas R. Jennings, Ece Kamar, Isabel M. Kloumann, Hugo Larochelle, David Lazer, Richard McElreath, Alan Mislove, David C. Parkes, Alex 'Sandy' Pentland, Margaret E. Roberts, Azim Shariff, Joshua B. Tenenbaum & Michael Wellman

CY. center for algorithmic transparency

# Thank you!

- www.**cycat.io**
- facebook.com/**CyCAT.EU**
- twitter.com/**CyCAT_EU**
- linkedin.com/in/**CyCAT**

Algorithmic System

Third Parties (T)

Data (D)

Fairness Constraints (F)

Algorithmic Model (M)

Input (I)

Output (O)

User (U)

center for
algorithmic
transparency

# USER STUDY – INVITATION!



http://ec2-34-255-198-84.eu-west-1.compute.amazonaws.com/opentag