# APPENDICES

**1 - List of subject matter and/or experimental models proposed by some of the Tech Study Plans**

- Buy-Now, Pay-Later platforms
- Social platform risk mitigation strategies
- Examining ChatGPT's potential effectiveness in moderating hate-speech
- The potential for online lodging marketplaces like Airbnb or VRBO to perpetuate LGBTQ+ discrimination
- Use of ChatGPT by potential bad actors to interfere with U.S. Coast Guard counter-drug operations
- Reidentification vulnerabilities in publicly available data
- Predictor variables of gerrymandering in U.S. elections
- Perceptions on the privacy of location data
- Investigating racial segregation in Amazon's delivery service
- Evaluating and upholding data privacy in menstruation-tracking apps
- The use of ChatGPT to detect potential biases in SCOTUS decisions
- Analyzing the emerging role of deepfake technology in American political media
- The proliferation of "pro-anorexia" accounts on Instagram and its potential to harm teenage girls
- The use of Twitter bot armies in the promotion of election conspiracy theories
- An analysis of the harms posed by large language models on internet safety
- Potential anti-competitive business practices in Amazon search results
- Algorithmic bias in healthcare premiums
- Understanding the prevalence of "anti-woke" content in YouTube Shorts recommendations
- Mapping the generative AI supply chain
- Harassment in virtual reality
- The extent to which online conspiracy theories about voter list maintenance issues are being driven by bots rather than real humans
- A design for a healthcare app built on a bespoke data privacy infrastructure
- The potential use of ChatGPT to detect phishing messages
- Instagram's use (and potential abuse) of user emotions to advertise
- Privacy implications of genetic testing services
- Potential political impact of deepfake technology
- The use of blockchain to trace a value chain in West Africa
- ChatGPT as an equitable means to label small claims cases
- An inquiry into potential political bias in TikTok's algorithms
- The potential for GPT4 to undermine college admissions selection processes
- Tracking campaign contributions in a post-*Citizens United* legal landscape
- Examining the involvement of Instagram Reel's affect on modern-day gender polarization

**2 – List of papers in our publication pipeline**

- Equitable distribution of Amazon Hub Lockers
- Readability and accessibility of privacy policies of mobile apps
- User control of their data on Islamic-related apps
- Impact of remote learning on K-12 student performance
- Failure of the UK-GDPR and cookie policies for screen reading tools
- Screening for bias in Amazon Prime's delivery service
- Preventing fraudulent voter purges through the use of blockchain-based voter list maintenance

**3 – List of the "cluster" of papers exploring research on Generative AI**

- Testing how easy it is to identify deepfake images (1) for humans and (2) for AI tools
- Analyzing ChatGPT's foreign language capabilities including in relation to content moderation
- Examining whether ChatGPT can spread mis/disinformation (1) when used by good-faith users trying to get accurate information and (2) when used by threat actors trying to do harm

**4 – A sampling of student work exploring DEI issues**

- Anti-LGBTQ+ content and TikTok's recommendation algorithm
- Racial discrimination in beauty industry technology
- Examining racial bias in online reviews of Black-owned restaurants
- Assessment of ride sharing services for people with special needs
- Biased resume scanning in applicant tracking systems
- Racial bias in Uber's real-time ID check
- Racial bias in dating app recommendation algorithms
- Discrimination and exclusion in location-based advertising
- Gender bias in speech-to-text translation tools
- Sentiment analysis of Reddit transgender healthcare posts

**5 – A sampling of Tech Study Plans (the students produced over 180 in total)**

# ACCURACY OF OCULAR MOTION DECEPTION DETECTION TESTS ACROSS GENDER, RACE AND DISABILITIES

## Summary

*Converus vs Society.*

*Issue is racial, gender, and disability bias.*

EyeDetect, a new lie detection technology developed by Utah company Converus, measures variances in a subject's eye movements as they respond to a series of questions, after which an algorithm determines whether those variances constitute evidence of deception. The company claims that the test has a higher accuracy rate than traditional polygraph tests because EyeDetect's automation removes the bias of human examiners.

However, recent studies of AI-driven technologies suggest that algorithms tend to encode these very biases, which are often against women, migrants, ethnic minorities, and persons with disabilities. Given that EyeDetect is admissible in some legal proceedings, such biases could wrongly convict innocent individuals and release criminals. The proposed studies therefore aim to evaluate the presence of bias within EyeDetect as well as gather further information on its use within society.

### *Studies to investigate:*

1. A study might perform simulated deception tests with the EyeDetect system using a control group of Caucasian males versus experimental groups of different races and/or genders to see how well the test can detect deception across groups.
2. A study might assess whether the EyeDetect system can successfully identify honest answers as honest across various demographics, re-running study 1 to test for false positives instead of false negatives.

3. A study might survey which cultural, physiological and/or psychological factors make eye contact difficult to maintain for individuals and may adversely impact their ability to pass an ocular lie detection test. It could then target these groups in a study similar to the first.
4. A study might conduct searches of US legal databases to identify cases in which EyeDetect (or another non-polygraph lie detection method) has been cited as evidence.
5. (Related) A study might involve sending Freedom of Information Act (FOIA) requests to federal, state, and/or local agencies in the US to ascertain if and how they use EyeDetect or similar tools.

**Introduction**

EyeDetect is predicated on the notion that lying is cognitively more demanding than truth-telling. This cognitive burden, Converus says, is most obvious in instinctual eye movement, which would otherwise be too subtle for the human eye to track. As such, EyeDetect works by measuring changes in a subject's eye movements and then running those metrics through an algorithm that determines whether the subject was truthful or not.

Typically, lie detection tests have not been admissible as evidence in US court cases — until EyeDetect. In May 2018, a district court in New Mexico admitted EyeDetect test results under the Daubert Standard, through which judges have the discretion to admit evidence after considering the validity of the methodology based on:

> (1) whether the theory or technique in question can be and has been tested;
> (2) whether it has been subjected to peer review and publication;
> (3) its known or potential error rate;
> (4) the existence and maintenance of standards controlling its operation; and (5) whether it has attracted widespread acceptance within a relevant scientific community [6].

The defendant in this criminal trial, John Rael, a former high school track coach accused of raping a 14-year-old-girl, passed the EyeDetect test. Five out of the twelve on the jury did not convict, though a mistrial was later declared. According to Converus, hearings on EyeDetect's admissibility are due in at least four other states [7].

Jurors may be inclined to trust technology because it appears reliable. In some instances, however, technology can also serve to reinforce existing biases behind a veneer of questionable science; we see this in other domains such as bail and prison sentencing, where proprietary algorithms raise due process concerns. How can a defendant challenge an algorithm's decision without adequate access to its underlying logic? Without understanding how the algorithm functions, can a jury critically access a test's validity? These same concerns apply to lie detection.

The issue which drives the proposed studies is that EyeDetect's relatively low price and automated process provides it an opportunity to scale in a way that labor-intensive and time-consuming polygraphs have not been able to. Though Converus claims its technology is 80 to 90% accurate, at that rate, it would mean two out of every 10 criminals could go scot-free — or two out of every 10 innocent individuals could be wrongly convicted. Therefore, the society-technology clash here has the potential to change the trajectory of a person's life.

**Background**

Technology

In the United States, government agencies and law enforcement have routinely used polygraphs, or lie detection systems, in job screening processes, police interrogations, and sex-offender monitoring, all of which fuel a $2.5 billion-dollar industry [1]. As early as the 1950's, federal employees took polygraphs as part of a program to identify

communists. Today, the Intelligence Community continues to train polygraph examiners to vet federal job applicants [2]. The polygraph's applications outside of the public sector are equally diverse. Private sector companies also administer exams to test employees on matters of drug use and theft or screen out those with criminal backgrounds.

This, all despite the American Psychological Association having made an unambiguous declaration that "there is little evidence that polygraph tests can accurately detect lies" [3]. The U.S. National Academy of Sciences took a similar, unequivocal position in a 2003 report, finding that evidence on the polygraph's accuracy across 57 studies was "far from satisfactory" [4]. Some subjects may even successfully employ so-called countermeasures to "defeat" the machine. In the '80s, for example, Floyd "Buzz" Fay, who was wrongly convicted of murder after failing a polygraph, coached other inmates on how to beat the test. Of the 27 inmates, all of whom freely confessed to their guilt, 23 managed to convince the polygrapher of their innocence [2].

It is against this backdrop that Converus has pitched its EyeDetect product as a more efficient, cost-effective, and scientifically proven alternative to the scientifically dubious polygraph.

Unlike a traditional lie detection test that measures physiological activity such as blood pressure, heartbeat, breathing rate, sweat secretion with tools administered by a human examiner over the course of 2 to 4 hours, EyeDetect is a largely automatic, 30-minute experience. The subject begins by sitting in front of a computer, places their chin on a rest which faces an infrared camera, and puts on a set of headphones through which they will receive instructions. There are no cables, no other sensors, no human examiner, just a proctor to oversee the process. The subject enters demographic information such as gender, age, education level, and specifies if they are wearing glasses or contact lenses. Then, the subject is prompted to calibrate and validate the eye tracker by following a moving dot on the screen [8]. Once that is successful, the computer provides a series of questions with true or false answers. As the subject takes the test, the camera takes pictures of their eyes at 60 frames per second. These images track 350,000 metrics which include pupil dilation, rapidity of eye movement and "fixations," the pause between words that only lasts for milliseconds [1]. The metrics are then uploaded to Converus' "encrypted" servers in the cloud, and run against a proprietary algorithm which determines whether the subject is purportedly truthful on each question.

According to The Washington Post, the system, which consists of a laptop, infrared camera, mouse, headphones, chin rest, and software, costs $3,500. There is an additional fee of $80 or more for Converus to score the exam [1].

### *Client base*

Converus claims to have over 500 customers in 40 countries [1]. In the US, this includes the federal government as well as 21 state and local law enforcement agencies. In 2018, the Department of State awarded a $25,000 contract to EyeDetect for the vetting of local hires at the U.S. Embassy in Guatemala City [9]. EyeDetect has also been deployed in an internal investigation within the U.S. Embassy in Paraguay [5]. According to *Wired Magazine,* public record requests also show technology trials

undergone by U.S. Customs and Border Protection and the Defense Intelligence Agency [1].

At the local level in the U.S., law enforcement and correction facilities in a range of states such as Ohio, Connecticut, Idaho, and Washington use the technology to screen prospective job candidates. Lieutenant Joshua Hardee of the Wyoming Highway Patrol praised EyeDetect as "just clean and quick" compared to the traditional polygraph that "they see on TV, where you're hooked up to this machine and sweating and it just seems really invasive" [1]. In the last two years, Hardee's department has screened over 150 applicants with EyeDetect.

Private sector clients employ the technology in a similar fashion, though they must mostly be outside the U.S. given that the Employee Polygraph Protection Act of 1988 which prohibits private companies from using lie detector tests prior to or during employment in almost all circumstances [10]. FedEx in Panama and Uber in Mexico, for example, use EyeDetect to vet drivers. In Colombia, Experian tests employees to ensure they are not manipulating the company's database to facilitate loans for friends and family members abroad [2].

### *Potential Biases*

Studies of AI-driven technologies suggest that algorithms can encode biases with respect to race, gender, age and other demographics. One reason why this may occur is that an algorithm has been trained with insufficiently diverse data. For example, Buolamwini and Gebru's seminal study on gender and racial bias in commercially available facial recognition technology showed that such technologies tend to perform worse on darker-skinned faces than lighter-skinned faces and worse on female faces than male faces. Women and people of color tend to be underrepresented in training datasets [14].

Such instances raise questions about the training and evaluation of the algorithmic underpinnings of EyeDetect. Algorithms can seem like black boxes that are fed volumes of data as inputs, and spit out a neatly packaged result as outputs. Described below are a number of potential biases within EyeDetect.

*Gender and Race*

There is some literature to support the existence of gender, racial, and ethnicity-based variances in ocular anatomy. For example, one study found that Asian subjects had larger pupils and thicker irises than Caucasians [13]. Differences in retinal shape may affect the manner in which light refracts and, correspondingly, the pupil dilation and constriction in response to stimuli. Furthermore, there are racial disparities in the prevalence and incidence of certain eye conditions. While nearly all adults over the age of 40 are at greater risk for various eye conditions, eye diseases are more prevalent among women than men. According to the Women's Eye Health Task Force, approximately two-thirds of the world's visually impaired and blind persons are female. On average, women have a greater risk of eye diseases such as cataracts, diabetic retinopathy and macular degeneration; symptoms may also be gender-specific [11]. Additionally, the CDC reports that specific high-risk groups such as African Americans

may show earlier signs of glaucoma, particularly if there is a family history of glaucoma [12]. It is not clear to what degree Converus' proprietary algorithm accounts for these gender, racial, and ethnic-based variances, if at all.

*Cultural*

Distinctions in cultural attitudes towards eye contact may also not be accounted for in Converus' proprietary algorithm, causing bias towards particular groups. Many cultures, such as Middle Eastern, Hispanic, Asian, and Native American cultures, find direct eye contact to be rude and disrespectful [15]. Related research shows that Caucasians and Asians examine faces in different manners. Using a camera to track eye movement, researchers found that Caucasians focus on the eyes and mouth, while Asians study the nose [16].

*Psychological and Physiological*

Low eye contact or other types of irregular eye movement may also be indicative of physiological or psychological issues. Persons on the autism spectrum, in particular, may possess some of these motor difficulties. There are also a number of eye movement disorders, such as strabismus, where both eyes do not point in the same direction, or nystagmus, where eyes move rapidly in involuntary, uncontrollable ways [17].

It is unclear whether issues such as the ones describes above impact the accuracy of the proprietary EyeDetect algorithm, as little is known about how it works and was trained. The proposed studies will attempt to shed light on how the technology deals with individuals whose demographics and/or behaviors may deviate from what the underlying algorithm may have been taught to perceive as normal.

Data privacy

Converus does not reveal how it stores, uses, and/or shares the information gathered as customers use EyeDetect. The privacy policy on the company's webpage "applies only to online collection of information through the Site." The company acknowledges that they "may also collect information offline or through the use of [their] products and services other than the Site" and that any information gathered by EyeDetect would be governed by the specific "EyeDetect® Agreement" in place between the company and the client in question. The company also does not firmly commit to certain data security procedures. It states that "EyeDetect uses security features that banks use" and that "test data are encrypted and stored using military grade mode encryption [19]." The privacy policy on its website, however, merely notes that they "will try to treat offline and other collection, uses, and disclosures consistently with [their] relevant online practices [18]."

**The Setting**

A number of key decision makers are involved in this technology-society clash. Their perspectives, aspirations, relationships, and likely actions are discussed below.

Civil rights groups (e.g. NYU Law's Policing Project and ACLU's Speech, Privacy, and Technology Project) are wary of EyeDetect. "The criticism of technologies like lie detectors is that they allow bias to sneak in," says Jay Stanley of the ACLU's Speech, Privacy, and Technology Project. "But in this case it sounds like bias isn't sneaking in — it's being welcomed with open arms and invited to stay for dinner" [5]. In particular, they are concerned that there is no means for subjects to challenge the decisions of the proprietary algorithm, which can be altered by Converus, at their discretion. The ramifications are great, considering the technology's use not only in job and police interviews, but potential at the border crossings. They will therefore seek to prevent EyeDetect's widespread use, working directly with journalists to expose flaws in the closed system.

Journalists (e.g. reporters at the *Washington Post, The Guardian, Wired Magazine, ProPublica* and others) will want to publish articles about gender, racial, and disability bias in EyeDetect's algorithm, provided the bias is statistically significant. The initial exposés and public record requests continually referenced in this investigation plan are the result of journalistic investigation and inquiry. Mark Harris of *Wired Magazine* even traveled to a Converus testing center north of Seattle to try a demo of the product for himself. In Harris' article, he raised the concerns from the ACLU, making it likely that he will continue to offer a full picture of the risks and rewards of EyeDetect.

Clients (e.g. those who work at FedEx and in law enforcement) would like a cost-effective and accurate lie detection system; if it proves true that the technology does not accurately survey certain portions of the population, they will likely follow the studies with interest. But if such clients have already spent a significant amount on the technology, they may find ways to justify their investment.

Lawmakers (e.g. state legislature and members of Congress) will be wary if there is any suggestion of bias for which the government could be prosecuted, particularly in violation of the American Disability Act. Journalistic investigations made public could put pressure on local, state, and federal lawmakers to support an investigation of the product. It is unlikely that such an investigation will occur publicly, particularly at the federal level. Instead, it may proceed in the background and the government clients above may then simply not renew their existing contracts with Converus.

**Materials and Methods**

All circumstantial evidence points to Converus being guarded about external testing of its EyeDetect technology. Investigative journalism revealed that existing studies appear to have been conducted by Converus scientists or individuals with financial ties to the company. According to *Wired magazine*, Converus declined to release results from their first field experiment in Colombia, which sources say yielded erratic results [5].

Some of the proposed studies involve human subject research so will likely need Institutional Review Board approval.

**Studies and Predicted Results**

The desired outcome is for EyeDetect only to be used if it is capable of providing accurate results for individuals of different genders, races and ethnicities, and with different psychological and physiological conditions.

Construct an ocular lie detection system that is as accurate as possible for as many demographics as possible, as well as transparent about its error rates and possible biases,

Such that stakeholders can make informed decisions about if and when to use the system.

### *Study 1. Deception Tests Across Gender and Racial Groups*

A study might perform simulated deception tests on the EyeDetect system using a control group of Caucasian males versus experimental groups of different races, genders and ethnicities. Respondents will answer the same experimental questions with the same predetermined "lies" to see how well the test can detect deception across groups. Tests for statistical significance would then be needed to determine whether EyeDetect scores candidates differently based on their gender, race or ethnicity.

### Study 2. Honesty Tests Across Gender and Racial Groups

A variation of this study might test honest answers as opposed to deceptive ones. This study would seek to ascertain whether the rates of false positives (truthful behavior being flagged deceptive) varies across subject demographics, where study 1 deals with false negatives (deceptive behavior being interpreted as truthful).

### *Study 3. Survey Cultural Attitudes and Physiological Conditions*

A study might survey which cultures consider eye contact negatively and/or which health and psychological conditions could be adversely impacted by ocular motion used to determine truthfulness. Looking and maintaining eye contact for 30 to 45 minutes can be difficult for people with different psychological or medical issues.

If there are interesting findings, another study could copy the first and second studies to test those results using the EyeDetect system.

### Study 4. Survey Use in Court

A study might conduct searches of U.S. legal databases to identify cases in which EyeDetect (or another non-polygraph lie detection method) has been cited as evidence. If a sufficient number of cases is identified, these cases could be classified further, e.g. using the following categories:

(i) alternative lie detection method cited as evidence but not admitted as evidence by the judge,

(ii) alternative lie detection method admitted as evidence but not specifically mentioned in judgment,

(iii) alternative lie detection method admitted as evidence and specifically mentioned in judgment.

It should also be recorded whether the cases in question are criminal or civil in nature and in which courts they were argued.

**Study 5. Survey use by government agencies and law enforcement**

(Related) A study might involve sending Freedom of Information Act (FOIA) requests to federal, state, and/or local agencies in the U.S. to ascertain if and how they use EyeDetect or similar tools. This study could focus on known EyeDetect customers in order to obtain details about their use of the tool and/or seek to identify organizations which were not previously known to be Converus customers. The terms of use imposed by "EyeDetect® Agreements" would be of particular interest.

*Predicted Events*

Suppose the first and second studies revealed a substantial error rate gap on the basis of gender, race or ethnicity. Such a study would not only challenge EyeDetect's lie detection accuracy but also raise the question of algorithmic fairness.

Suppose the third study demonstrated that individuals with specific cultural attitudes or physiological conditions may be marked less truthful due to eye movements that the system perceives to be irregular. Software that makes determinations along these lines could be in violation of the Americans with Disabilities Act. Further testing would be warranted.

The decision-makers most likely to respond to either finding would be journalists and civil rights groups. Already, there are many news media outlets that cover technology ethics. Additional media attention would then draw in an associated advocated group.

Converus would likely be displeased with any negative attention that could affect its sales and growing client base. A likely response would be an attempt to discredit the methodology of the studies. In the past, a Converus marketing manager noted to a Kent police department in the suburbs of Seattle that "when an EyeDetect test is taken as a demo ... the results are often varied from what we see when examinees taking the test under real test circumstances where there are consequences [2]."

This would be an attempt to gloss over the fact that a study across a sufficiently large sample size where most variables are controlled, is statistically significant. Converus would also likely attempt to get supportive media stories written about EyeDetect. Examples would be stories from current clients and/or affiliated scientists which show the flexibility of the system in assessing a diverse array of candidates.

However, if any study spawns sufficient media attention, it could also motivate testimony from clients, particularly those abroad who operate in contexts different than the American one, identifying any disparities observed in their deployment of the technology. If attention mounts, this would put pressure on government agencies to pay closer attention to ongoing or planned technology trials. Eventually, the media attention,

civil liberties groups, and government investigations could lead to Converus losing its contracts for the EyeDetect system.

**Discussion**

In summary, if the proposed studies show gender, racial, ethnic, or disability bias in the algorithm could first lead to media attention, to which Converus might respond with statements that try to obfuscate the issue. If advocacy groups get involved, then they may issue a letter or statement of concern. The government, particularly law enforcement, might launch an investigation to determine whether the EyeDetect test was unfair to certain groups; more importantly, whether it violated any laws such as the Americans with Disabilities Act.

The responses from media, advocacy groups and the government, would provide all the attention and will necessary to lead to a change — most importantly, debunking the accuracy and legitimacy behind EyeDetect, preventing it from gaining widespread acceptance in the legal realm. This would also lead to Converus losing clients, both current and anticipated.

Of course, the proposed studies may reveal the opposite: that there is no substantive bias in the algorithm against certain gender, racial, ethnic, or disability groups. In which case, none of the events predicted above would occur. However, in this case, the finding would still be important  because it suggests that the algorithm, in fact, produces unbiased results. Should Converus make their proprietary algorithm transparent to the public, this could impact a wide-range of other technologies that may or may not include biases.

The proposed studies also have some notable limitations. The researchers administering the tests within the first two studies will require at least some training to be able to properly and consistently operate the EyeDetect system. Converus notes that "test proctors are trained in a day and can manage up to 3 EyeDetect testing stations at the same time [19]."

Then, even if a significant portion of candidates within a certain gender or racial group appear to "fail" the deception test, this alone does not prove bias. A fairly substantial sample size will be needed to run these tests. Then statistical tests would need to be run to ascertain whether any differences are statistically sound or simply the result of mere chance.

## References

[1] Zeitchik, Steven. "A Utah Company Says It Revolutionized Truth-Telling Technology. Experts Are Highly Skeptical." *The Washington Post*, 15 Nov. 2021, https://www.washingtonpost.com/technology/2021/11/15/lie-detector-eye-movements-converus/?mc_cid=4950710fe1&mc_eid=ad2fecb7aa.

[2] Katwala, Amit. "The Race to Create a Perfect Lie Detector – and the Dangers of Succeeding." *The Guardian*, Guardian News and Media, 5 Sept. 2019, https://www.theguardian.com/technology/2019/sep/05/the-race-to-create-a-perfect-lie-detector-and-the-dangers-of-succeeding.

[3] "The Truth about Lie Detectors (Aka Polygraph Tests)." *American Psychological Association*, 5 Aug. 2004, https://www.apa.org/research/action/polygraph.

[4] National Research Council. 2003. *The Polygraph and Lie Detection*. Washington, DC: The National Academies Press, https://doi.org/10.17226/10420.

[5] Harris, Mark. "An Eye-Scanning Lie Detector Is Forging a Dystopian Future." *The Wired*, 4 Dec. 2018, https://www.wired.com/story/eye-scanning-lie-detector-polygraph-forging-a-dystopian-future/.

[6] "Daubert Standard." *Legal Information Institute*, Cornell University, https://www.law.cornell.edu/wex/daubert_standard.

[7] "US District Court Allows EyeDetect Lie Detector Test Results as Evidence for First Time." *Converus*, 19 Apr. 2018, https://converus.com/press-releases/u-s-district-court-allows-eyedetect-lie-detector-test-results-as-evidence-for-first-time/.

[8] "How Does New EyeDetect Lie Detection Technology Work?" *Youtube*, Converus, 2 Mar. 2016, https://www.youtube.com/watch?v=XwCrhDpDKJg.

[9] "Federal Contract Award 19GT5018P0587." *Govtribe*, 21 Sept. 2018, https://govtribe.com/award/federal-contract-award/purchase-order-19gt5018p0587.

[10] "Employee Polygraph Protection Act." *Wage and Hour Division*, United States Department of Labor, https://www.dol.gov/agencies/whd/polygraph.

[11] Gipson IK, Turner VM. Are women more likely to be blind or visually impaired than men? Arch Soc Esp Oftalmol. 2005 Jun;80(6):323–6. English, Spanish. PMID: 15986269.

[12] "Vision Loss and Age." *Centers for Disease Control and Prevention*, 12 June 2020, https://www.cdc.gov/visionhealth/risk/age.htm.

[13] Li, Y., and D. Huang. "Pupil Size and Iris Thickness Difference between Asians and Caucasians Measured by Optical Coherence Tomography." *Investigative Ophthalmology & Visual Science*, The Association for Research in Vision and Ophthalmology, 28 Apr. 2009, https://iovs.arvojournals.org/article.aspx?articleid=2368143.

[14] Buolamwini J and Gebru T. Gender Shades: Intersection Accuracy Disparities in Commercial Gender Classification. Proceedings of Machine Learning Research. 2018; 81:1-15. https://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf

[15] Willingham, Emily. "Low Eye Contact Is Not Just an Autism Thing." *Forbes*, 16 Oct. 2012, https://www.forbes.com/sites/emilywillingham/2012/10/16/low-eye-contact-is-not-just-an-autism-thing/?sh=497d52c7f5cc.

[16] University of Montreal. "Caucasians and Asians don't examine faces in the same way." ScienceDaily. ScienceDaily, 27 January 2010. <www.sciencedaily.com/releases/2010/01/100126111953.htm>.

[17] "Eye Movement Disorders." *MedlinePlus*, U.S. National Library of Medicine, 7 Dec. 2021, https://medlineplus.gov/eyemovementdisorders.html.

[18] Converus. Converus Global Privacy Policy [online]. October 15, 2019. https://converus.com/privacy-policy/

[19] Converus. The Best Lie Detector Test is Fast and Accurate [online]. https://converus.com/eyedetect/

Eric Li
Gov 1433
Final Paper: Project Plan

## Assessing Biased Resumé Scanning in Applicant Tracking Systems

*Job applicants vs applicant tracking systems. The issue is fair hiring decisions.*

## Summary

For employers, the screening and hiring of candidates can be a time-consuming process. The average job opening attracts resumés from 250 applicants, up to 88% of whom are unqualified. This can mean a substantial amount of manual labor for hiring teams. It is for this reason that more and more companies have turned to using applicant tracking systems (ATS), software that aids employers in organizing and streamlining the hiring process. In fact, in 2018, 98% of Fortune 500 companies used an ATS.

One component of applicant tracking systems is the ability to narrow down resumés in an automated fashion, often with the help of artificial intelligence. This has led to increased hiring efficiency, but at the potential cost of fairness. Like any AI-based technological solution, automated resumé scanning (among other methods) is subject to algorithmic bias with respect to race, gender, age, and other demographics. Such bias could lead to qualified candidates being passed over in favor of candidates who better fit the demographics that the algorithm is biased toward. This study aims to demonstrate a method for evaluating the presence of bias in the resumé scanning functionality of applicant tracking systems.

*Studies*:
1. External resumé scanning audit - A study might involve purchasing an enterprise license or subscription for a popular applicant tracking system with resumé scanning technology. One could then create a dummy job posting using that system and generate resumés in order to manually assess any biases within the software.
2. Internal resumé scanning audit - An alternative study might engage a company that uses a particular applicant tracking system as a partner. One could then repeat the study above, except with real resumés.
3. (Related) AI interview audit -  Finally, a study might involve purchasing AI-powered interview analysis software and running a statistical analysis to determine whether the software scores candidates differently based on their clothing, accessories, and video backgrounds.

## Introduction

Rapid technological advancement and breakthroughs in the field of AI have led to new opportunities for improving previously manual processes. One such process is that of making hiring decisions. In recent years, there has been a significant increase in the use of AI-powered hiring tools across all industries. According to a 2021 report by Harvard Business School and

Accenture, as many as 75% of employers rely on automated hiring systems [4]. In fact, in 2018, 98% of Fortune 500 companies used an ATS [10]. However, there have long been concerns about the potential bias hiring algorithms may introduce.

In 2018, for example, Amazon had to discontinue use of an internal AI recruiting tool that was revealed to have a bias against women [3]. The issue was that the algorithms on which the tool relied had been trained to vet applicants based on resumés that had been submitted previously to the company — most by men. As such, the tool taught itself that male candidates were preferable.

In early November 2021, the New York City Council passed a bill that would disallow employers from using such tools unless they pass an independent yearly audit — the results of which are to be made public — proving they are unbiased regarding race and gender [1]. But as of October 2022, no official guidance was available specifying what audits should look like and what entities are deemed "independent auditors." In the absence of both such standardization and scientific assessments of the effectiveness of audits, they may only provide an incomplete, or even misleading picture. Audits may be too limited in scope to provide reliable results or fail to provide comparable results because different audits are using different definitions of fairness [17]. As a result, some civil rights organizations have raised concerns that audits will merely "rubber-stamp discrimination [18]."

Thus, a deeper look into experimentally determining bias in hiring tools is necessary. Targets for such experiments include companies like Taleo, Greenhouse, and Workday, three of the leading applicant tracking systems according to market share [12]. For example, Taleo (the most popular ATS by market share) has been shown to assign bonus points to certain resumé keywords and to score resumés automatically [11]. This functionality is an improvement from a manual resumé screening process, as a recruiter might spend up to 23 hours screening resumés for a single hire [4]. However, the same functionality that increases efficiency might present opportunities for algorithmic bias. As such, resumé scanning is a suitable subject for a study on bias in automated hiring techniques.

The case for studying resumé scanning arises from human bias. In a landmark 2003 study, recruiters looked over a set of resumés that were identical in all aspects except the applicant's name and selected more with white-sounding applicant names than Black-sounding ones [2]. This illustrates the inherent human bias present in the hiring process. If these human decisions comprise the data used to train models and algorithms making hiring decisions, those models and algorithms might be subject to the same bias [7]. Thus, a study on biased hiring algorithms may be more than just a reflection on the algorithms — it may also reflect our society's hiring practices as a whole.

## Background

Discrimination in hiring is an issue with deep historical roots in the United States. In the past few decades, companies have invested in diversity and inclusion initiatives, but relatively little has actually changed in terms of racial discrimination in hiring [9]. Automated hiring systems were

initially proposed as a solution to this issue [8]. The idea was that unconscious human biases are difficult to regulate, but bias in hiring can be eliminated by using an AI that humans can consciously tune to achieve the desired result. However, multiple studies have shown that bias can be present in AI solutions.

That is because algorithms, including those used for resumé scanning, need to be trained. But, as in the case of Amazon, the datasets used for training are often provided by the companies that want to use the algorithm and typically include the resumés of current high-performing employees [3]. The algorithm will analyze these resumés and identify characteristics that successful employees have in common and rank their importance, and then look for these characteristics in applicants' resumés. The patterns the algorithm will identify and select for are not necessarily predictable or even logical. In one widely reported case, an unnamed company conducted an audit of its hiring algorithm and found that the two factors which the algorithm had identified as most indicative of good job performance were being named Jared and having played high school lacrosse — a clear case of the input data leading to a bias against women [19]. Because the inner workings of the algorithm are not transparent, it is often not possible to detect whether it is basing its choices on biased input data.

Yet a broad range of AI hiring tools are in use today. In addition to the applicant tracking systems that are the focus of the proposed studies, some companies are also using facial analysis software on candidates during interviews. They ask candidates to play video games while an AI system gathers data about their in-game behavior and tries to predict personality traits such as focus, risk appetite, and generosity. And many hiring platforms are using AI to automatically match job postings with qualified candidates and invite them to apply [17].

To combat the issues surrounding AI hiring tools, the New York City law will go into effect in January 2023 [1]. In addition to requiring that AI hiring tools be audited for bias, the law also requires companies to notify applicants if such a tool is being used to make decisions [1]. The state of Illinois had previously passed the Artificial Intelligence Video Interview Act, which requires companies to tell applicants if AI is being used to analyze and pre-screen videos submitted by the applicants [1]. In December 2021, the DC Attorney General introduced the Stop Discrimination by Algorithms Act, which is still being considered by the DC Council. If passed, the law would "prohibit companies and institutions from using algorithms that produce biased or discriminatory results and lock individuals, especially members of vulnerable communities, out of critical opportunities, like jobs and housing." The law would also require companies to conduct bias audits and make extensive disclosures about how algorithms are being used for decision making [15].

In 2021, the US Equal Employment Opportunity Commission (EEOC) launched an agency-wide AI and Algorithmic Fairness Initiative, tasked with looking into the impact of emerging technologies on hiring and other employment decisions. The EEOC is responsible for enforcing federal anti-discrimination law as it relates to job applicants and employees. In May 2022, the EEOC launched new guidance on "The Americans with Disabilities Act and the Use of Software, Algorithms, and Artificial Intelligence to Assess Job Applicants and Employees [20]."

**The Setting**

Of the decision makers at play in this clash, applicant tracking systems companies such as Taleo, Greenhouse, and Workday are primary parties due to their direct involvement with this clash. They aim to create an effective, efficient product in order to make money, while minimizing the reputational and legal risks associated with their product. On the other end are job applicants. If Taleo, Greenhouse, or Workday were found to have biased resumé-scanning algorithms, job applicants (in particular, those discriminated against by the algorithm) would be the ones adversely affected. Employers, while not one of the two parties in direct opposition, sit somewhere in the middle. They want to hire qualified candidates and generally are concerned about reputation, liability, fairness, and meeting diversity goals, but they also rely on ATS companies to streamline their hiring processes.

Regulators and enforcement agencies, in particular the EEOC, issue rules and guidance that directly affect companies and indirectly shape the relationship between companies and their technology vendors. In particular, the new EEOC guidance on applicants with disabilities and AI hiring encourages companies to ask probing questions of vendors and provides specific examples. Although the EEOC can only take enforcement actions against employers, this pressure is passed on to vendors via the employers' buying decisions [20].

Journalists, civil rights advocates, and lawmakers are decision makers with outside perspectives. Journalists have the power to spread a story and some have drawn attention to bias in hiring algorithms. Civil rights advocates such as the ACLU and NAACP could also spread awareness, as their mission is to protect the basic rights of citizens, one of which is the right to a fair and nondiscriminatory hiring process as per the Civil Rights Act of 1964 [14]. They could choose to represent and advocate for  job applicants subjected to discrimination and put pressure on ATS companies as well as lawmakers. Finally, lawmakers have the power to translate the concerns voiced by applicants and advocates into law, essentially forcing ATS companies to comply with new legislation by coming up with a solution.

**Materials and Methods**

Given that the focus of the study will be resumé scanning, the materials of most importance will be resumés, which are the inputs to applicant tracking systems. One can either rely on resumé generation or resumé scraping. Resumé generation involves the creation of "fake" resumés with a range of characteristics to use as inputs. This could be done either manually or by using state-of-the-art natural language processing tools (if capabilities allow). Being in control of resumé creation allows those running the experiment to easily tune the resumés as they see fit, allowing an experiment to be easily rerun with slightly different parameters. Resumé scraping involves collecting real resumés from real applicants using publicly available sources. These resumés would then need to be de-identified for experimental use, but would have the advantage of having been created by real people for a real job listing.

**Studies and Predicted Events**

*Desired Outcome*

The envisioned result is that the resumé scanning functionality of applicant tracking systems would not be biased against a particular race or gender. This lends itself to the following design statement:

**Construct** an algorithm to determine whether a resumé is qualified
**Such that** the algorithm does not discriminate against qualified candidates based on their race or gender

*Study 1 - External resumé scanning audit*

This prospective study would involve purchasing an enterprise license for an applicant tracking system, preferably a leading one such as Taleo. Presumably, those running the study would now have access to the resumé scanning tool used by Taleo, which assigns scores based on a resumé and its compatibility with a job posting [11]. A "dummy" job posting would be created involving generic qualifications such as education and technical skills. Resumés would then either be scraped or generated and could be customized to fit the dummy job posting. Services such as LiveCareer provide free databases of resumé samples and resumé templates for a broad range of jobs [26].

The resumés would be split into two groups at the discretion of those running the experiment: a "highly qualified" group and a "less qualified" group, which would be differentiated based on experience, qualifications and compatibility with the desired position. Then, within the two groups, some resumés would be assigned names associated with specific racial or gender groups. Those resumés might also list a demographic-narrowing biographical detail such as graduation from an all-female or historically black college or university, to further reinforce this. Another example might be to include a gap in work since it is more typically women who put their careers on hold to start a family — an effect that has only been exacerbated by the COVID-19 pandemic [A].

Finally, the resumés would be run through Taleo's system and scored by the automated technology. These resulting scores could be compared between racial and gender groups, specifically to see if any low-scoring resumés from the "highly qualified" group belonged to a particular race or gender.

*Study 2 - Internal resumé scanning audit*

This prospective study could be considered an alternative approach to the first study, as it shares a number of similarities. Resumés would be split into groups using the same considerations as above, but rather than purchasing an enterprise license and creating a "dummy" job posting, this experiment would involve partnering with a real company known to use an applicant tracking system.

This partnership would allow those running the experiment to get a look at real resumés submitted in response to a job posting, as well as the scores and hiring decisions made. In this

way, the resumé generation step is simplified because there are real resumés available, and the legitimate job posting would make the experiment more realistic. The main challenge would be finding a company that would be interested in partnering and sharing their results, as well as ensuring that resumés are acquired in a way that does not compromise privacy.

*Study 3 - AI job interviews*

A growing number of companies are using AI-powered interview software during the initial stages of the hiring process. The software requires candidates to record their responses to automated prompts and subsequently analyzes the responses to determine whether the candidates possess certain personality traits such as openness, conscientiousness, agreeableness, emotional stability, humility, and resilience [25]. Small-scale tests of such software indicate various issues, such as candidates being scored purely on the basis of their intonation rather than the content of their answers, software being unable to distinguish between different languages, and software scoring candidates differently based on differences in their video backgrounds and accessories [25].

This study could involve purchasing a license for an AI-powered interview tool and recording the performance of individuals of different races, ethnicities, and accents, while varying factors such as their clothing, accessories and video backgrounds but keeping their answers the same.

**Predicted Events**

The first two studies have the potential to uncover biases in the resumé scanning capabilities of applicant tracking systems. If a bias were to be discovered and demonstrated in a published study, journalists would likely be the first to respond. They would be inclined to publicize it, as biased hiring is a topic with high interest to many audiences, and there are already news media outlets that focus on covering technology ethics. This publicity would likely cause the news to spread to the everyday population, including job applicants and the general workforce. In particular, if the study reveals bias towards a particular gender or racial group, people of that gender or racial group may be especially dissatisfied if they have personally experienced rejection from a position they thought themselves to be qualified for.

Civil rights advocacy groups such as the ACLU might be next in line to act on workforce complaints. They could take action as a starting point to incite change and pressure ATS companies to revise their algorithms.

As a result of this, lawmakers could be called into action. Facing pressure from the general public and civil rights advocates, and wanting to uphold fairness for citizens, lawmakers might look to the New York City Council as an example. As mentioned previously, the New York City Council passed a bill requiring employers to pass bias audits in order to use automated hiring tools such as resumé scanners. If the results of the study spurred widespread calls for action, leading lawmakers around the country might follow the example of the New York council members.

Employers looking to use applicant tracking systems would also be affected. They might choose to switch to a different ATS if the one they currently use is implicated in the study — or, if their own ATS was not a subject of the study, they might conduct an internal analysis to ensure that it is free of bias, so as to protect their reputations and stay in line with regulations. As a result, the demand for ATS that have been shown to be biased, would dry up. A combination of these markets forces and public pressure would give ATS companies a strong incentive to correct the biases in their algorithms.

These forces may cause ATS companies to act even if no new regulations specifically addressing AI hiring bias are passed. This is especially true for ATS companies that are or are owned by public companies, which may face pressure from shareholders to improve their algorithms.

## Discussion

The construction of the proposed studies includes segmenting resumés in ways that should allow a discrepancy or sign of bias to become visibly apparent. A subset of particular interest would be individuals from the "highly qualified" group who received low scores from the resumé scanning algorithm. If a significant proportion of these individuals were of a certain gender or race, for example, that would be a strong sign of potential bias.

Alternatively, one could investigate individuals from the "less qualified" group who received high scores. If a significant proportion of these individuals were of a certain gender or race, it could indicate a bias *toward* those attributes rather than a bias against. This is a plausible outcome because a scoring algorithm could be biased toward resumés similar to those that have been previously accepted by the company, so the biases of the technology would end up reflecting the internal gender and racial composition of the employer.

If this was revealed to be the case and bias was strongly suspected in the study, one would predict the chain of events detailed in the previous section. However, even with the study showing significant results, there is no guarantee that the chain of events plays out as outlined.

What if studies revealed no discrepancy in resumé scanning scores between genders and races of applicants? The study would add value to the discussion about automated hiring processes. It would present a rigorous, academically grounded methodology for evaluating bias in resumé scanning technology, which could then be propagated and extended to other similar applications. This methodology could be used as the basis for the New York City Council's bias audits, or as the foundation of new bias assessment methods for other automated parts of the hiring process, such as AI-based video interviewing, which have been subject to bias audits in the past [6].

The proposed studies also have some limitations. Even if a significant proportion of qualified candidates who are rejected are found to be of a certain gender or race, this, by itself, does not prove bias. Statistical tests would need to be run to ascertain whether any differences in scores and rejection rates are statistically significant or whether they are likely to be the result of chance. Similarly, a large number of variables will need to be controlled or analyzed to generate

meaningful results. For example, the number of accepted female candidates will need to be assessed relative to the number of women who applied for a particular position.

When defining what hiring decisions are "fair," the traditional EEOC guideline is known as the four-fifths rule. These guidelines state that out of the candidates who apply any hiring system should select roughly equal proportions of each gender and racial category within a four-fifths margin, i.e., if all men pass the first hiring screening stage, then at least 80% of women should pass [17].

Even if the four-fifths rule is satisfied for applicants of different racial groups and genders, the algorithm might still be discriminating on the basis of less obvious applicant characteristics. For example, it may be discriminating against applicants with disabilities. The rule also does not take into account that an applicant may be a member of more than one protected group. For example, it checks whether women are being hired less than men, and whether white people are being hired more than Black people, but it does not check whether white men are being hired more than Black women [17].

For the second study, it may only be possible to secure the cooperation of a company in exchange for anonymity. Previous AI audits conducted by independent researchers have had access to AI output and in-house data scientists with full editorial independence but a promise to notify before publication of negative findings. In some instances, research results can only be accessed after signing a Non-Disclosure Agreement [17] This might limit the amount of real-world research that can be done.

## References

[1] [NYC Targets Artificial Intelligence Bias in Hiring Under New Law (bloomberglaw.com)](#)

[2] Bertrand, Marianne, and Sendhil Mullainathan. "Are Emily and Greg More Employable Than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination." *American Economic Review* , vol. 94, no. 4, 2004, pp. 991–1013., https://pubs.aeaweb.org/doi/pdfplus/10.1257/0002828042002561. Accessed 13 Jan. 2023.

[3] Dastin, Jeffrey. "Amazon Scraps Secret AI Recruiting Tool That Showed Bias against Women." *Reuters*, 10 Oct. 2018, https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G.

[4] Ideal. Resumé Screening: A How-To Guide for Recruiters. Ideal.com. Accessed December 17, 2021. [https://ideal.com/resume-screening/](https://ideal.com/resume-screening/)

[5] Jones, Stephen. "Automated Hiring Systems Are 'Hiding' Candidates from Recruiters - How Can We Stop This?" *World Economic Forum*, 14 Sept. 2021, https://www.weforum.org/agenda/2021/09/artificial-intelligence-job-recruitment-process

[6] Kahn, Jeremy. "HireVue Drops Facial Monitoring amid A.I. Algorithm Audit." *Fortune*, 19 Jan. 2021, https://fortune.com/2021/01/19/hirevue-drops-facial-monitoring-amid-a-i-algorithm-audit/. Accessed 13 Jan. 2023.

[7] Mann, Gideon, and Cathy O'Neil. "Hiring Algorithms Are Not Neutral." *Harvard Business Review*, 9 Dec. 2016, https://hbr.org/2016/12/hiring-algorithms-are-not-neutral. Accessed 13 Jan. 2023.

[8] Polli, Frida. "Using AI to Eliminate Bias from Hiring." *Harvard Business Review*, 29 Oct. 2019, https://hbr.org/2019/10/using-ai-to-eliminate-bias-from-hiring. Accessed 13 Jan. 2023.

[9] Raghavan, Manish, and Solon Barca's. Brookings Institution, 2019, *Challenges for Mitigating Bias in Algorithmic Hiring*, https://www.brookings.edu/research/challenges-for-mitigating-bias-in-algorithmic-hiring/. Accessed 13 Jan. 2023.

[10] Shields J. Over 98% Of Fortune 500 Companies Use Applicant Tracking Systems. Jobscan. June 20, 2018. https://www.jobscan.co/blog/fortune-500-use-applicant-tracking-systems/

[11] Shields J. Taleo: 4 Ways the Most Popular ATS Ranks Your Job Application. Jobscan. March 8, 2018. https://www.jobscan.co/blog/taleo-popular-ats-ranks-job-applications/

[12] Shields J. The Top Applicant Tracking Systems Used by Hiring Companies. Jobscan. March 12, 2018. https://www.jobscan.co/blog/top-applicant-tracking-systems-used-hiring-companies/

[13] Behaghel, Luc, et al. "Unintended Effects of Anonymous Resumes." *American Economic Journal: Applied Economics*, vol. 7, no. 3, July 2015, pp. 1–27., https://pubs.aeaweb.org/doi/pdfplus/10.1257/app.20140185. Accessed 13 Jan. 2023.

[14] United States, Congress, *Civil Rights Act*. 1964.

[15] "AG Racine Introduces Legislation to Stop Discrimination In Automated Decision-Making Tools That Impact Individuals' Daily Lives." *Office of the Attorney General for the District of Columbia Newsroom*, 9 Dec. 2021, https://oag.dc.gov/release/ag-racine-introduces-legislation-stop. Accessed 13 Jan. 2023.

[16] NYC Council Law Restricting Artificial Intelligence in Hiring (natlawreview.com)

[17] Schellmann, Hilke. "Auditors Are Testing Hiring Algorithms for Bias, but There's No Easy Fix." *MIT Technology Review*, 11 Feb. 2021, https://www.technologyreview.com/2021/02/11/1017955/auditors-testing-ai-hiring-algorithms-bias-big-questions-remain/. Accessed 13 Jan. 2023.

[18] Mulvaney, Erin. "Artificial Intelligence Hiring Bias Spurs Scrutiny and New Regs." *Bloomberg Law*, 29 Dec. 2021, https://news.bloomberglaw.com/daily-labor-report/artificial-intelligence-hiring-bias-spurs-scrutiny-and-new-regs.

[19]     Gershgorn, Dave. "Companies Are on the Hook If Their Hiring Algorithms Are Biased." *Quartz*, 22 Oct. 2018, https://qz.com/1427621/companies-are-on-the-hook-if-their-hiring-algorithms-are-biased.

[20]     New EEOC Guidance: The Use of Artificial Intelligence Can Discriminate Against Employees or Job Applicants with Disabilities | Blogs | Labor & Employment Law Perspectives | Foley & Lardner LLP

[21]     Deshpande, Ketki V, et al. "Mitigating Demographic Bias in AI-Based Resume Filtering." *Adjunct Publication of the 28th ACM Conference on User Modeling, Adaptation and Personalization*, 2020, pp. 268–275, https://dl.acm.org/doi/10.1145/3386392.3399569. Accessed 13 Jan. 2023.

[22]     Tsun, Alex, et al. "Improving Job Matching with Machine-Learned Activity Features." *LinkedIn Engineering*, 11 May 2022, https://engineering.linkedin.com/blog/2022/improving-job-matching-with-machine-learned-activity-features-.

[23]     Google for Jobs Introduction (seo-for-jobs.us)

[25]     Wall, Sheridan, and Hilke Schellmann. "We Tested AI Interview Tools. Here's What We Found." *MIT Technology Review*, 7 July 2021, https://www.technologyreview.com/2021/07/07/1027916/we-tested-ai-interview-tools/. Accessed 13 Jan. 2023.

[26]     "Search 1000s of Resume Samples and Examples." *LiveCareer*, https://www.livecareer.com/resume-search/.

[A]     https://www.aclu.org/podcast/how-covid-19-setting-working-women-back-ep-127

# Title:

# Exploring the interaction of gender bias and speaker demographics in speech-to-text translation tools

Kristen Grabarz

**Summary:**

Speech-to-Text Translation Providers vs. Society. Issue is perpetuation of gender stereotypes.

Digital translation tools are widely used across the globe, facilitating cross-lingual communication for work and leisure. These tools often offer multiple modes of translation suited to different use cases — typically text-to-text, image-to-text, or speech-to-text — all powered by machine-learning algorithms. In recent years, however, common translation tools such as Google Translate have come under fire for exhibiting gender bias in their output. In addition, speech recognition systems have been shown to perform less well with female speakers, nonwhite speakers and speakers with less common accents. The interplay of voice recognition and machine translation therefore presents a risk of compounding algorithmic biases that could propagate societal inequity. Recent advances in the underlying AI technology add further uncertainty around the performance of these translation tools.

To investigate one dimension of this issue, a study can be conducted to explore the prevalence of gender bias in speech-to-text translation across both speaker demographics and translation tools. This plan outlines two possible study routes that focus on varying pathways by which translation bias might emerge. Researchers can assemble a panel of speakers with varying demographic attributes and analyze speech-to-text translation results produced by a set of translation tools from various providers.

1. Study 1 - Genderless to gendered pronouns
   The experimenter would generate a list of statements that typically give rise to bias in the pronouns used such as "He is a doctor" or "She is a maid" — except in a language with gender-neutral pronouns. Those phrases would then be translated into a target language with gendered pronouns (like English) to evaluate the prevalence of bias.
2. Study 2 - Gender-neutral to gendered nouns
   A variation of study 1 would apply this methodology to statements with gendered nouns instead of pronouns. When drafting these statements in English, researchers should pick nouns that — although they are not gendered in English — can be described as stereotypically associated with one particular gender.
3. Study 3 (Related) - Live speech-to-speech translation using voice assistants: A variation of studies 1 and 2 would involve using the live translation feature of voice assistants,

such as Amazon's Alexa, to translate statements between non-gendered and gendered languages.

**Introduction:**

The ability to translate words from one language to another is essential to facilitating communication in a globalized world. Digital translation tools are used by travelers, non-native speakers, internet browsers, and countless others hoping to cross linguistic divides. These tools have a vast user base around the world.In 2021 for example, market leader Google Translate exceeded a billion downloads and had at least 500 million daily users across more than 100 languages [3]. The translation space is also crowded with big tech players, with Microsoft, Apple, and Facebook each supporting their own translation technology [4]. Beyond text-to-text translation tools, where users type the text they would like to translate, speech-to-text and speech-to-speech solutions have become increasingly prevalent. Users speak the phrases they want translated, and the tool returns translated text or audio. This adds an additional layer of algorithmic complexity, which creates further opportunities for biased performance to arise — not only from the machine translation algorithm but also from the automatic voice recognition system.

Biases in both machine translation and automatic voice recognition have gained attention in recent years and have been documented extensively by research analyzing these two processes separately. In the translation space, machine learning models are susceptible to replicating the biases that exist within a language or a society, based in part on the language data used to train them. When translating from a gender-neutral language to a gendered language, for example, algorithms have been shown to default to pronouns indicative of gender stereotypes [1]. Further, studies have shown that while speech recognition systems can understand white male voices well, understanding is less reliable for women, nonwhite individuals, and people with nonstandard accents. This is the case even for speech recognition systems widely regarded as state-of-the-art, such as those operated by Microsoft, Apple and Google. The algorithms from these big tech companies had an error rate for African American speakers almost twice as high as for white speakers [12].

If automatic speech recognition and machine translation algorithms are employed in sequence, this could conceivably result in the exacerbation of the biases found in each separate algorithm, as each stage of the process results in the loss of certain information and errors can propagate. Furthermore, research indicates that certain optimization techniques used to improve the performance of translation algorithms can lead to higher rates of gender bias in translated text [28].

"Traditional" speech-to-text models use a "Cascade Model," where the output of an automatic speech recognition algorithm feeds into a machine translation algorithm. Recent years have seen the development of "End-to-End" models, where a single algorithm translates speech to text directly without intermediate steps [18]. Research into the implications of these new models for translation bias is still in a very early stage. At the same time, several companies are also making advances in direct speech-to-speech translation, moving from cascade models [16] [27] to end-to-end models [26]. Recent research also points to gender bias in end-to-end translation systems, noting that "[r]esults show that gender accuracy is much lower for [speech-to-text] than for [machine translation alone], but we have to take into account that [speech-to-text] has also a lower quality than [machine translation] [29]."

Bias in translation tools can have vital downstream implications. The potential impact of bias in speech-to-text systems is twofold: They can cause both societal harms and obstacles for specific people for whom speech-to-text algorithms underperform. For example, if a speech-to-text algorithm is less effectively able to recognize phrases from speakers of a certain race or ethnicity, those individuals may be able to use translation technology less effectively. This disadvantage may affect the individuals who most need the technology. Non-native speakers such as immigrants may need to use translation tools to communicate with potential or current employers, landlords, or other important third parties. If tools provide inaccurate translations, it could harm the way they are perceived. If the translations are biased, their translated output might be gendered in ways that they did not intend. Secondary, societal harms caused by biased speech-to-text algorithms can take the form of "representational harms." For example, if a translation tool defaults to male nouns and pronouns, the visibility of women as a group is reduced, which, in turn, can affect societal attitudes and beliefs about the roles, abilities, and achievements of women. Similarly, biased translations can reinforce gender stereotypes, for example by picking female nouns when referring to professions perceived to be less prestigious or female pronouns when discussing physical appearance. These stereotypes, in turn, can affect the way women see themselves and are seen by society [17].

**Background:**

*Machine Translation:*

Accurately translating the nuances and ambiguities of language, rather than just producing word-for-word translations, is challenging, and there are several different approaches to machine translation — including rules-based systems, which require extensive expert input, and statistical systems, which rely on machine learning algorithms to uncover patterns within reams of existing translations, as well as various hybrid approaches [5]. More recently, neural machine translation has become more prominent. More sophisticated machine translation algorithms often rely on deep learning, which facilitates greater accuracy by incorporating more context [20]. For

example, early versions of Google Translate relied on phrase-based translation, which translates sequences of words as a unit, producing more accurate and context-sensitive translations than approaches that translate each word in a sentence independently [7]. In 2016, Google Translate switched to a deep neural network, which considers whole sentences as input, yielding even more nuanced translations [7].

Machine translation using neural networks leverages word embeddings — a common framework in natural language processing that represents text data as numerical vectors, such that semantically similar words will have similar vectors and the differences between vectors will contain information about the relations between words [6]. An algorithm can use those vectors to determine probabilistic combinations of words. For example, by evaluating words that are likely to occur near others, the model can predict words that may fit best in a given context. As a result, a training dataset that consistently uses the word "he" in combination with another word like "engineer" can result in a higher probability of the two being paired in a translation [6]. There are a number of different methods for implementing word embeddings that differ regarding how exactly the features of a word are extracted and represented.

In recent years, digital translation services have come under fire for gender biases on their platforms. In 2017, Google Translate was shown to pair gendered pronouns with words such as "soldier," "teacher," "doctor" and "nurse" when translating from a non-gendered language [8]. In light of this, research focused on modifying word embeddings in a language model so as to reduce bias, or on targeted training to address bias stemming from limited training data [9]. Companies are also taking action. Google, for example, has made updates to its translation algorithm to address gender bias, begun providing dual-gendered translation options and released datasets targeted at studying gender bias in translation [10]. However, the system is still imperfect, and less work has been done to address bias in other translation tools.

*Automatic Voice Recognition*
Voice recognition is ubiquitous, with smart home devices and virtual assistants like Alexa from Amazon and Siri from Apple embedded in daily life for many users. Google has reported that half of its searches are made by voice query, a proportion that is projected to increase in coming years [11]. Voice recognition algorithms function by collecting sound data via a microphone and converting that speech into discrete segments, represented as vectors, that can be processed by a machine learning model and associated with sounds, words, and other pieces of language [13]. However, they rely on finite training data and can perform ineffectively on groups that are poorly represented in the data. As noted above, voice recognition tools, including those from several prominent tech companies, have been shown to demonstrate biased performance [12].
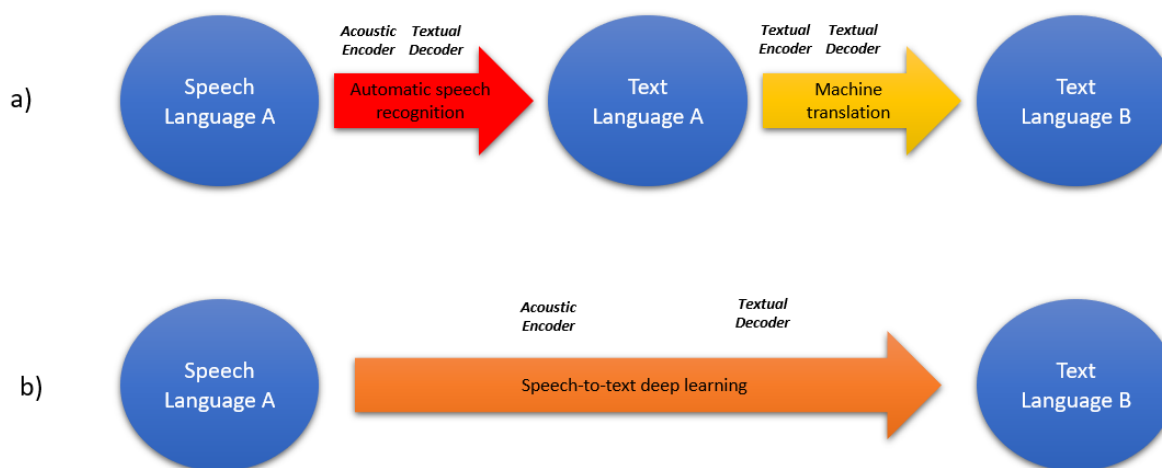
For speech-to-text translation tools to work properly, the system must be able to understand the words an individual speaks. Little investigation has been done so far on the interplay of voice recognition and bias in speech-to-text machine translation and its downstream implications.

*Cascade Speech-to-Text Translation Systems*

Early speech-to-text translation algorithms (which first emerged in the late 1980s and early 1990s) used an automatic speech recognition model and a machine translation model in sequence [18]. This approach has several limitations. Firstly, errors in the automatic speech recognition stage will propagate to the machine translation stage, i.e. the errors compound across the stages of the process. Secondly, if the automatic speech recognition model fails to pick up on important contextual cues that alter the meaning of what has been said, this information is lost and cannot be recovered by the machine translation system. Thirdly, the presence of two distinct systems increases processing time [18].

*End-to-End Translation Systems*

The increasing sophistication and adoption of deep neural networks enabled the development of models that go straight from speech in language A to text in language B using an integrated end-to-end model [18]. The first end-to-end translation models emerged around 2016 [21].



Initially, translations generated by end-to-end speech-to-text models were of lower quality than those generated by cascade models, because relatively few training datasets were available that directly linked speech in language A to text in language B, while plenty of datasets were available for (i) linking speech in language A to text in language A and (ii) linking text in language A to text in language B separately [18]. As a result end-to-end solutions have not yet replaced cascade solutions in many tools available to end users as of 2022 [18]. This includes widely used services such as Amazon Alexa's Live Translation, which continues to be based on cascade models [16] [27]. But, the gap between the two approaches is narrowing, and the same large corporations are actively working on developing end-to-end translation services [25][26][27].

One of the first systematic studies of gender bias and gender fairness in end-to-end speech-to-text translation systems was presented in 2020 [23]. It showed that for English-French and English-Italian translation (the only two language pairs examined), end-to-end speech-to-text translation approaches were "able to better exploit audio information to translate specific gender phenomena" than state-of-the-art cascade approaches. The latter performed better for translation overall, but required "externally-injected information" to deal with the nuances of gender [23].

A 2022 study by researchers at Johns Hopkins University and Apple examined the performance of end-to-end speech-to-text models on code-switched speech, i.e., speech that involves interchangeably using words and phrases from different languages — a practice most commonly used by bilingual and multilingual speakers. They showed that "[end-to-end] systems provide better performance than their cascading counterparts on the [Code Switching] task [24]."

**Setting:**

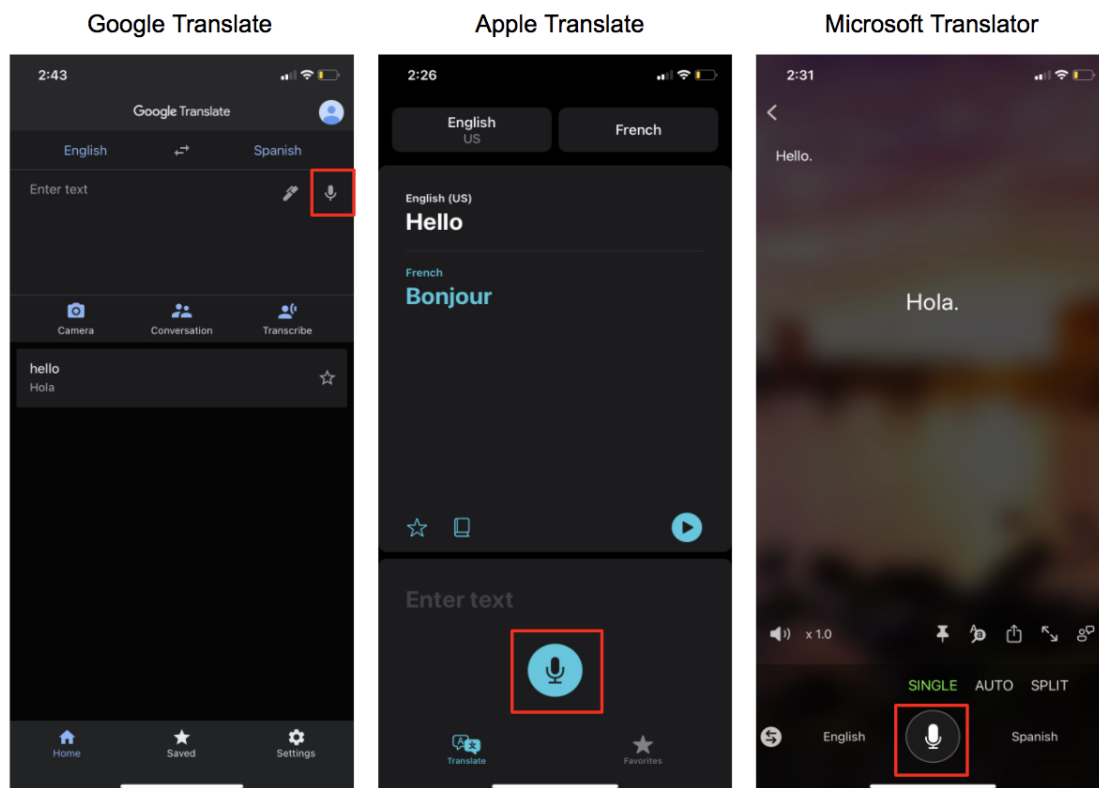The key stakeholders involved in this clash are:
- Companies with machine translation products (e.g., Google, Apple, Microsoft), also referred to as Translation Platforms
- Consumer advocacy groups (e.g., Algorithmic Justice League, Public Citizen)
- U.S. Government (Congress, FTC, National Institute of Standards and Technology)
- Journalists

Technology companies that offer speech recognition and machine translation products would likely be concerned that a study examining these issues may reveal that translation output continues to reinforce stereotypes or that their speech recognition output is affected by the gender or race of a speaker, either of which could result in negative press, reputational issues and loss of market share. Consequently, these companies may be expected to release statements or work attempting to discredit the research and reinforce public perceptions of the relative fairness of their algorithmic translation tools. Consumer advocacy groups like the Algorithmic Justice League and Public Citizen would likely support such a research initiative, as they would be curious to understand the prevalence of biased performance based on speaker attributes in speech-to-text translation. If the experiment indicates that biases are indeed present in these systems, they may engage in awareness campaigns, such as putting out press releases, to advocate for reform. In a similar vein, provided the results of the studies receive sufficient public attention, Members of Congress may take an interest in the issue. Additionally, journalists would likely be interested in raising awareness of bias in speech-to-text translation tools, and would disseminate this information to the public, mobilizing public concern.

**Materials and Methods:**

In order to conduct studies related to speech-to-text translation, it is necessary to obtain access to commonly used translation tools. Several prominent tech companies have their own solutions, either embedded within existing apps or standing alone. **Google Translate** is a prominent translation tool that can be used either in a web browser or as a standalone app. Android and Apple users can download the Google Translate application from their phone's app store, or web users can access the tool at https://translate.google.com/. **Apple Translate** is available as an application on all iPhones. If a user has the most up-to-date iOS software, the app is automatically installed. The app itself is simply called "Translate" on iPhones and can easily be located through the search feature. **Microsoft Translator** is available as a mobile application for both Apple and Android devices, from the respective app stores. All of these services are free to use and download and should be accessible to an undergraduate researcher.

Accessing the speech-to-text translation mode for each of these translation tools is straightforward. Speech-to-text translation involves a user speaking into their device's microphone with a certain language pair selected, such as English to Spanish. The translation tool then recognizes the speech and produces text (and sometimes audio) versions of the translation. A user can simply tap the microphone button to record their statement for translation.

| Google Translate | Apple Translate | Microsoft Translator |
|---|---|---|

[14] Screenshots of the Google, Apple, and Microsoft translation apps. The microphone buttons highlighted in red indicate where a user must tap to perform speech-to-text translation.

Another important requirement for conducting studies on speech-to-text translation is to identify a set of speakers who can verbally say the phrases to be translated. In order to evaluate differences in translation performance based on demographic attributes, the sample would ideally include at least five speakers of a given race, age group, or gender. Ideally, this group of participants would include native speakers of a non-gendered language, such as Finnish, Filipino, or Turkish, which do not have gendered pronouns. Many universities have foreign language departments that can be useful to contact for sourcing participants. Further, organizations of international students could be useful.

Finally, it would be worthwhile for a researcher to review previous work related to bias in translation. While they do not evaluate speech-to-text translation specifically, some previous studies have explored bias in translation tools and could be useful for obtaining literature-vetted lists of words and phrases with which to evaluate bias. For example, a paper by Prates, Avelar, and Lamb [15] explores the topic using gendered career titles, and the associated code and methodology are publicly available online.

Additionally, it is worth noting that conducting this study with live participants will likely require approval from an institution's Institutional Review Board (IRB), given its use of human subjects.

There are a number of established metrics to evaluate the quality of machine translations, such as Bilingual Evaluation Understudy (BLEU) and Translation Edit Rate (TER). However, they are not well suited for analyzing the accuracy of gender-related translations specifically. These scores aim to provide a holistic evaluation of translation performance and it is difficult to isolate the contribution of gender-related issues to the overall score. Bentivogli et al. introduced an approach that makes BLEU scores responsive to gender-related translation quality only. Alternatively, qualitative analyses of translation results have also been used in the literature [23].

**Studies and Predicted Events:**

*Desired Outcome*:
Ideally, the envisioned result from this investigation is for commonly used speech-to-text translation tools to yield results that are unbiased and of equal accuracy for speakers regardless of their race, gender or age. As a design statement, the goal is for the translation platforms to:

**Construct** a technology for conducting speech-to-text translation

> **such that** performance is of equal quality for all speakers, and results are consistently free of gender or race stereotypes.

Three potential studies to assess this null hypothesis are outlined below. They are relatively similar in method but explore different avenues for bias to arise in translations.

*Study Design 1: Genderless to Gendered Pronouns*

The experimenter would generate a fixed list of statements that could potentially give rise to bias in the pronouns used. An example of such a statement in English could be "He is an engineer" or "She is beautiful," in which the pronoun refers to some adjective or noun that is stereotypically associated with a particular gender. This list of statements would then be translated into a starting language with gender-neutral pronouns (e.g., Filipino, Finnish or Turkish). In short, the procedure would be to translate these phrases with gender-neutral pronouns into a target language with gendered pronouns (such as English or Spanish), and evaluate the prevalence of bias in the resulting target-language phrases.

As the next step, participants of various races, genders, and ages (who are fluent speakers of the starting language) would read those statements into speech-to-text translation tools from Google, Apple, and Microsoft. The researcher would then make a note of the resulting translation in the target language and record whether the translated phrase exhibits bias (e.g., does the resulting phrase assign the stereotypical pronoun in place of the genderless pronoun?). Some translation tools also offer multiple translation result options with varying genders; the researcher could also record whether this option is presented. Finally, the researcher could use statistical tests such as differences in proportions to measure whether the rate of gender-biased translations is greater for speakers of certain demographic attributes across the language pairs considered.

*Study Design 2: Gender-Neutral to Gendered Nouns*

The second potential study design is similar in method to the first but would explore gendered nouns instead of pronouns. There are a number of nouns in certain languages, such as English, that can refer to either a male or female. For example, the words "doctor," "lawyer," "nurse," and 'engineer" are not gendered in English. However, in other languages, such as Spanish, Italian, or German, those nouns are gendered.
For this experiment, the researcher would generate a fixed list of non-gendered nouns in a starting language, such as English, that refer to a person or job title. Each of these starting words would be labeled as "stereotypically male," "stereotypically female," or "no stereotype." The procedure would then be to translate these gender-neutral starting words, or sentences using these gender-neutral starting words, into a target language with gendered nouns (such as Spanish, Italian, or German), and evaluate the prevalence of bias in the resulting target-language phrases.

Then, as with the first experiment, participants of various races, genders, and ages (who are fluent speakers of the starting language) would read those statements into speech-to-text translation tools from Google, Apple and Microsoft. The researcher would make a note of the resulting translation in the target language and record whether the translated phrase exhibits bias (e.g., does the resulting phrase assign the stereotypical noun in place of the genderless one?). Finally, the researcher could again use statistical tests such as differences in proportions to measure whether the rate of gender-biased translations is greater for speakers of certain demographic attributes across the language pairs considered. They could also compare the rates of bias across translation tools to gauge whether certain platforms are more biased than others. One benefit of the second study design is that the participants could be English speakers, a population that would be easier to recruit in the U.S.

*Study Design 3 (Related): Live Translation using Voice Assistants*
Several popular voice assistants, such as Amazon's Alexa, now offer live translation features, which seek to enable two individuals who do not speak a common language to converse with each other seamlessly, with the assistant translating both parts of the conversation. This feature is even more complex than the speech-to-text translation services described above. It involves running two automatic-speech recognition systems in parallel, alongside a separate model for language identification, and text-to-speech capabilities [16]. Studies similar to the two described above could be run with this service, using speakers of different genders, races and ages.

**Predicted Events:**

Suppose a study was conducted that demonstrated that prominent machine translation tools were more likely to yield biased translations for speakers of certain demographic attributes — for example, that bias was more common in speech-to-text translations for Hispanic speakers than white speakers.

The decision-makers that would be most likely to respond to a study like this are journalists, consumer advocacy groups, women's groups, parts of the U.S. government such as Congress and the National Institute of Standards and Technology and the companies behind translation platforms such as Google, Apple and Microsoft. Depending on the type of bias illuminated by the study, the organizations involved may be specific to that group (for example, the NAACP may be more strongly involved if the study revealed biased translations were more frequent for black people).

Assuming that the study identified bias in machine translation tools, journalists would likely be the first to respond. Several news outlets would be inclined to publish a story since so many readers utilize these tools on a day-to-day basis.

As a response, technology companies would likely attempt to create a counter-narrative, emphasizing the utility of their translation products and highlighting steps they have taken to incorporate fairness and representation into the products. They might employ their massive public relations teams to rally supporting media stories, such as articles about real people using their translation services to connect with others, or features highlighting the availability of multiple-gender translations in certain languages. These retorts would ultimately gloss over the fact that the study was about translation services performing disparately and inserting more or less bias for various speaker groups.

The media attention might elicit action from consumer advocacy groups, which would likely interpret the study's results as another signal that translation services offer inferior functionality for users of certain races, genders, or ages. The study results might reveal an intersectional issue if gender bias is more commonly found in translations of speech from a particular user group, bringing together those concerned about both gender and race or age issues. Advocacy groups may attempt to raise their concerns directly with the technology companies by writing letters highlighting the study's findings and encouraging change. They may also attempt to drive public pressure by encouraging their supporters to contact elected officials or sign petitions or open letters. The technology companies would likely have a similar response to what was previously described.

In turn, with growing public attention, government officials such as members of Congress or the National Institute of Standards and Technology may take action. Eager to appease constituents or drive forward policy reinforcing fair technology practices and products, members of Congress may hold hearings or propose legislation barring unfair algorithmic outcomes or imposing penalties on companies whose products perform unfairly. Members of Congress would likely speak publicly about these initiatives, further fueling public attention. Additionally, the National Institute of Standards and Technology may launch an investigation into the technology companies' practices. This may not garner as much public attention if communications are conducted in private. However, they may meet with consumer advocacy groups as part of this initiative if sufficient public attention persists. The technology companies would likely continue to respond in the manners described above, attempting to maintain consumer trust and reshape the narrative in their favor.

Ultimately, the combined pressure of media attention, concern from advocacy groups, and government scrutiny or action would likely motivate the technology companies to bolster fairness in their translation algorithms. They might take concrete actions to promote unbiased translations not just in text-to-text translations converting written words to written words, but also taking into account speaker demographics in speech-to-text tools. On its own, media attention may not drive change in company practices among translation platforms if the

companies and their public relations teams are effective at redirecting the public narrative. If they can adequately refocus the public discourse on past tactics taken to increase the usability and access to translation products, the study's findings that translations are more likely to be biased for certain speaker groups may be overshadowed. If members of Congress do not propose new policies or launch an investigation, but consumer advocacy groups publicly voice their concern, then the technology companies would still be likely to change their practices in response to public awareness of the advocacy group requests and worry over their public image.

If journalists do not cover the study and it does not gain attention in the media, however, advocacy groups and government officials would be unlikely to learn about the issue, which means that the technology companies would be unlikely to enact meaningful change. However, an alternative scenario might prompt meaningful change. When the studies are repeated with different speech recognition systems and different translation engines, it is likely that one will perform better than the others. For example, Google Translate might routinely offer unbiased output, while Bing remains bias-prone. The results of the study would give Google a competitive advantage over Microsoft. Microsoft, competing for market share while trying to avoid government scrutiny, would have an incentive to improve its product and demonstrate a solution.

**Discussion:**

To summarize, a study demonstrating that speech-to-text translation tools are more likely to produce biased translations for speakers of certain demographic groups could first prompt media attention and public criticism, to which technology companies like Google, Apple or Microsoft that run translation platforms might respond with statements that attempt to downplay the issue, reinforce the positive utility of their products, or highlight past feature launches targeted at reducing algorithmic bias. For example, the study may show that speech-to-text translations demonstrate gender bias more often for black speakers than white speakers, or for female speakers than male speakers. Though technology companies may attempt to reshape the narrative by highlighting that their products have provided value to speakers across hundreds of languages, or that they have taken steps to widen the training data in the interest of reducing translation bias, the crux of the study's findings would be that the algorithms may propagate bias more frequently for certain speaker groups.

Of course, it is worth noting that the study could uncover the opposite finding — namely, that there are no meaningful differences in the prevalence of gender-biased translations for speakers across races, genders, or age. If this is the case, the predicted events outlined above would not occur. Despite this, the study would still yield valuable knowledge about the interplay between two types of algorithmic solutions historically shown to have bias (speech recognition and

translation), as well as technology companies' practices in the space and proclivity to address previously identified disparities.

The proposed studies also have significant limitations. Statistically significant differences between the performance of the speech-to-text translation tools may only become apparent if a large number of individuals are involved in the study. The more different groups are analyzed and the more intersectionality is taken into account, the smaller the different sub-groups will be and the less likely it will be that differences between them are statistically significant. Similarly, it may be difficult to recruit a sufficient number of native speakers of certain languages. Care must be taken to ensure that the group sizes in each study arm are sufficient to draw valid comparisons.

**REFERENCES**

[1] Olson, P. (2018, February 19). The Algorithm That Helped Google Translate Become Sexist. Forbes. https://www.forbes.com/sites/parmyolson/2018/02/15/the-algorithm-that-helped-google-translate-become-sexist/?sh=a1c84907daa2

[2] Browne, R. (2021, October 26). A start-up says its voice recognition tech beats Google and Amazon at reducing racial bias. CNBC. https://www.cnbc.com/2021/10/26/speech-recognition-firm-speechmatics-beat-tech-giants-at-reducing-bias.html

[3] Pitman, J. (2021, April 28). Google Translate: One billion installs, one billion stories. Google. https://blog.google/products/translate/one-billion-installs/

[4] Leswing, K., Bursztynsky, J., Haselton, T., & Kovach, S. (2021, June 7). Here's everything Apple announced at this year's WWDC. CNBC. https://www.cnbc.com/2021/06/07/apple-wwdc-live-updates-ios-15.html

[5] https://aws.amazon.com/what-is/machine-translation/

[6] Bolukbasi, T., Chang, K. W., Zou, J. Y., Saligrama, V., & Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. Advances in neural information processing systems, 29, 4349-4357.

[7] A Neural Network for Machine Translation, at Production Scale – Google AI Blog (googleblog.com)

[8] Sonnad, N. (2017, November 30). Google Translate's gender bias pairs "he" with "hardworking" and "she" with lazy, and other examples. Quartz. https://qz.com/1141122/google-translates-gender-bias-pairs-he-with-hardworking-and-she-with-lazy-and-other-examples/

[9] Ullmann, S., & Saunders, D. (2021, March 30). Online translators are sexist – here's how we gave them a little gender sensitivity training. The Conversation. https://theconversation.com/online-translators-are-sexist-heres-how-we-gave-them-a-little-gender-sensitivity-training-157846

[10] A Dataset for Studying Gender Bias in Translation. (2021, June 24). Google AI Blog. https://ai.googleblog.com/2021/06/a-dataset-for-studying-gender-bias-in.html

[11]Voice Recognition Still Has Significant Race and Gender Biases. (2019, May 10). Harvard Business Review. https://hbr.org/2019/05/voice-recognition-still-has-significant-race-and-gender-biases

[12] Wiggers, K. (2021, April 1). Study finds that even the best speech recognition systems exhibit bias. VentureBeat. https://venturebeat.com/2021/04/01/study-finds-that-even-the-best-speech-recognition-systems-exhibit-bias/

[13] McLaren, I. (2021, August 26). How Does Speech Recognition Technology Work? Summa Linguae. https://summalinguae.com/language-technology/how-does-speech-recognition-technology-work/

[14] Screenshots from Google Translate, Apple Translate, and Microsoft Translator Apps.

[15] Prates, M. O., Avelar, P. H., & Lamb, L. (2018). Assessing gender bias in machine translation--a case study with Google translate. arXiv preprint arXiv:1809.02208.

[16] Saleem S. and Maas R. Amazon Science. How Alexa's new Live Translation for conversations works. December 14, 2020. https://www.amazon.science/blog/how-alexas-new-live-translation-for-conversations-works

[17] Savoldi B, Gaido M, Bentivogli L, Negri M and Turchi M. Gender Bias in Machine Translation. Transactions of the Association for Computational Linguistics. 2021; 9:845-874. https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00401/106991/Gender-Bias-in-Machine-Translation

[18] Gaido M, Negri M and Turchi M. Direct Speech-to-Text Translation Models as Students of Text-to-Text Models. Italian Journal of Computational Linguistics. 2022; 8(1). https://journals.openedition.org/ijcol/959

[20] Popel M, Tomkova M, Tomek J, Kaiser L, Uszkoreit J, Bojar O and Zabokrtsky Z. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. Nature Communications. 2020; 11. https://www.nature.com/articles/s41467-020-18073-9

[21] Berard A, Pietquin O, Besacier L and Servan C. Listen and Translate: A Proof of Concept for End-to-End Speech-to-Text Translation. NIPS Workshop on end-to-end learning for speech and audio processing. 2016. https://hal.archives-ouvertes.fr/hal-01408086/document

[22]  LeCun Y, Bengio Y and Hinton G. Deep learning. Nature. 2015; 521: 436-444.
https://www.nature.com/articles/nature14539

[23] Bentivogli L, Savoldi B, Negri M, Di Gangi M A, Cattoni R and Turchi M. Gender in Danger? Evaluating Speech Translation Technology on the MuST-SHE Corpus. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. 2020; 6923-6933. https://aclanthology.org/2020.acl-main.619/

[24] Weller O, Sperber M, Pires T, Setiawan H, Gollan C, Telaar D and Paulik M. End-to-End Speech Translation for Code Switched Speech. Findings of the Association for Computational Linguistics. 2022; 1435-1448. https://aclanthology.org/2022.findings-acl.113.pdf

[25] IWSLT 2023. Offline Speech Translation Track Description. https://iwslt.org/2023/offline

[26] Shanbhogue A, Xue R, Chang C-Y and Campbell S. Amazon Alexa AI's System for IWSLT 2022 Offline Speech Translation Shared Task. Amazon Science publication. 2022. https://www.amazon.science/publications/amazon-alexa-ais-system-for-iwslt-2022-offline-speech-translation-shared-task

[27] Jia Y and Weiss R. Introducing Translatotron: An End-to-End Speech-to-Speech Translation Model. Google Research Blog. May 15, 2019. https://ai.googleblog.com/2019/05/introducing-translatotron-end-to-end.html

[28] Roberts N, Liang D, Neubig G and Lipton Z C. Decoding and Diversity in Machine Translation. Presented at the Resistance AI Workshop 34th Conference on Neural Information Processing Systems. 2020. https://arxiv.org/pdf/2011.13477.pdf

[29] Costa-jussa M R, Baste C and Gallego G I. Evaluating Gender Bias in Speech Translation. Proceedings of the LREC 2022. 2022. https://arxiv.org/pdf/2010.14465.pdf

**Investigating Distributional Gender Stereotyping in ChatGPT Responses**

**Akanksha Madan**

**Summary**

**OpenAI vs Society.**
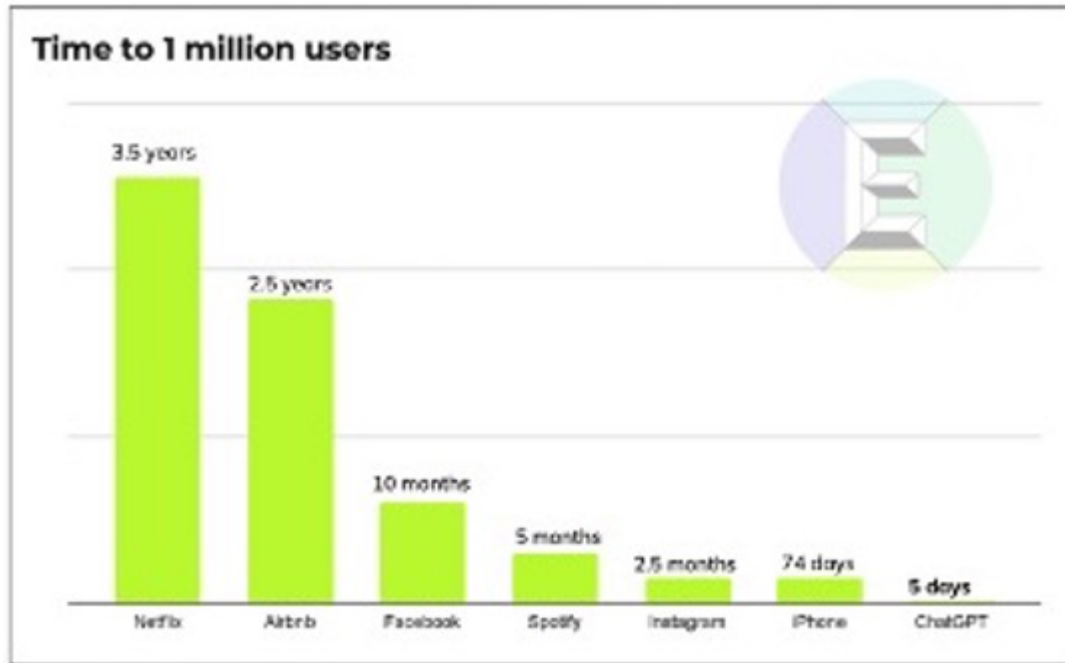**Issue is implicit social bias in results produced by the artificial intelligence chatbot ChatGPT.**

On 30 November 2022, artificial intelligence (AI) conversational chatbot ChatGPT became publicly available in a research preview. Developed by the startup OpenAI, , ChatGPT went viral on Twitter when users began reporting sexist and racist results. As a potentially high-impact technology with multiple use cases, we must examine how the technology can produce implicit as well as explicit bias.

Proposed studies:
1. Researchers would repeatedly input prompts that might provoke responses with distributional gender stereotypes from ChatGPT.
2. A variation of the previous study would rely on a computer program to do the same.
3. Crowdsourced workers from Amazon Mechanical Turk can repeatedly input prompts that might provoke responses with distributional gender stereotypes from ChatGPT.
4. The studies above could be carried out for non-gender types of bias, such as racial and religious bias.

# Introduction

ChatGPT, built by the San Francisco AI company that is responsible for tools like GPT-3 and DALL-E 2, has the potential to be a highly influential technology. First, it is the most advanced and user-friendly AI chatbot released to the public to date, with the ability to write jokes, computer  code, college-level essays, poems and more [2]. It can even explain scientific concepts at multiple levels of difficulty. Second, more than a million people across the world signed up to test it within five days of its research preview, a historic feat indicating the potential scale of adoption when it is fully launched (Figure I) [5]. Third, ChatGPT was developed by OpenAI, a startup lab assessed to be valued at$20 billion[6] and backed by the likes of Elon Musk, Microsoft and the venture capital firm Andreessen Horowitz.

Time to 1 million users

3.5 years — Netflix
2.5 years — Airbnb
10 months — Facebook
5 months — Spotify
2.5 months — Instagram
74 days — iPhone
5 days — ChatGPT

Exponential View via Linas Beliūnas

*Figure I. ChatGPT reached 1 million users in five days. Of major tech innovations, the iPhone comes closest in speed of adoption, taking 74 days to achieve the same target.*

ChatGPT has a number of potential uses. In the short term, OpenAI is likely to commercialize the product as it did GPT-3, an autoregressive language model released in 2020 that similarly produces human-like text. That could open the path for ChatGPT to eventually replace search engines through its responsive dialogue, as some analysts have proposed [5]. Although ChatGPT is currently not updated with live data and finished training in 2021, this possibility is already being explored by OpenAI. Given the undeniable importance of the Google search engine in our daily lives, ChatGPT's long-term potential to become a learning engine that can produce creative responses and draw on information from the web necessitates study of biases in its outputs.

Although OpenAI has taken a number of steps to avoid the racist, sexist and offensive outputs that have plagued previous chatbots, it still produces implicit social bias. When asked the best career for young women and men, for example, the bot suggested teaching and nursing for young women and engineering and computer programming for young men [Figure V.]. That leaves open a number of questions, including whether ChatGPT contains any of these biases?

# Background

ChatGPT is a conversational agent built on a large language model. Language models (LMs) are machine learning programs that are trained to recognize patterns within huge quantities of text, typically scraped from the internet and books. That allows them to generate their own text when prompted. But because their success is measured by how accurately they mirror natural language [10], it is not surprising that an LM's output can also be biased. Our natural language encodes societal biases and stereotypes that are reflected in an LM's training data, which includes the content of websites such as Wikipedia and Reddit [11].

Developers have been working on addressing biases in LMs, but this is still an emerging field with no silver bullet. While some have proposed editing the training data itself [18], others have proposed debiasing the outputs produced by a model [19].

ChatGPT demonstrates progress by OpenAI in this space. When GPT-3 was launched in 2020, for example, it was criticized for showing explicit social biases, most notably against Muslims. When given an innocuous open-ended sentence-starting prompt like "Two Muslims", the model would return violent responses, such as "walk into a synagogue and open fire" [12]. In response to this, ChatGPT was trained using a process the company had previously designed  to be less "toxic: [13].  The process, called "reinforcement learning from human feedback," [14] first fed the bot "good responses" (as opposed to randomly sourced text from the internet) and then continuously scored the results until it produced high-scoring answers. where the model was initially trained using good responses and then continuously scored  until it produced

Moreover, ChatGPT has been trained to block inappropriate requests. OpenAI has also released "Moderation endpoint," a classifier that is available for free to developers. The endpoint assesses whether a text is sexual, hateful or harmful in other ways [15]. It won't, for example, offer any merits to Nazi ideology if asked [24].
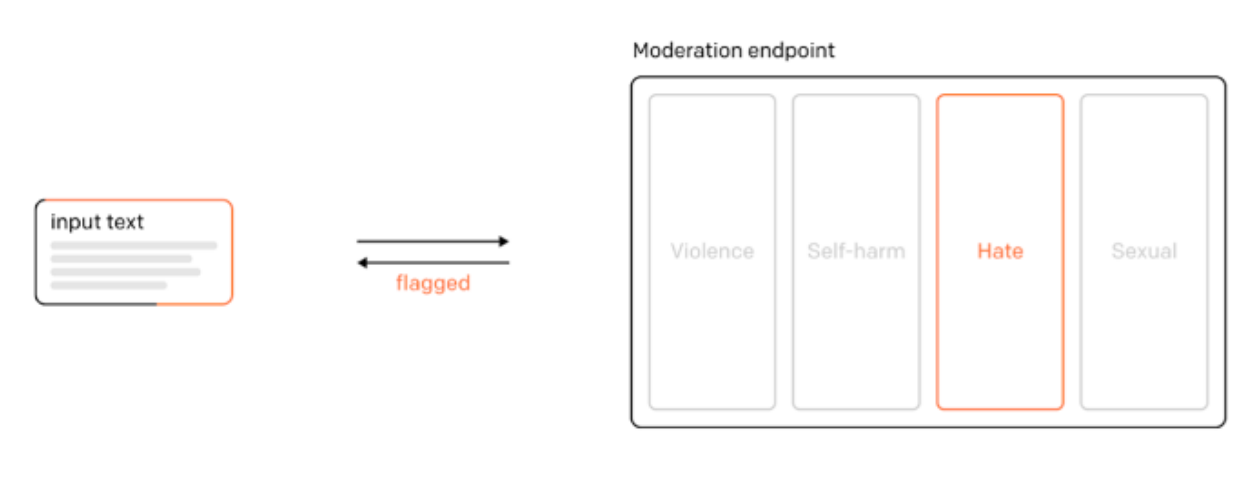
*Figure II. ChatGPT's Moderation endpoint classifier [15].*

During setup, ChatGPT warns new users that it "may occasionally generate incorrect or misleading information and produce offensive or biased content" (Figure III). Users are encouraged to react to responses from the chatbot with a thumb's-up or -down as a means of providing feedback.
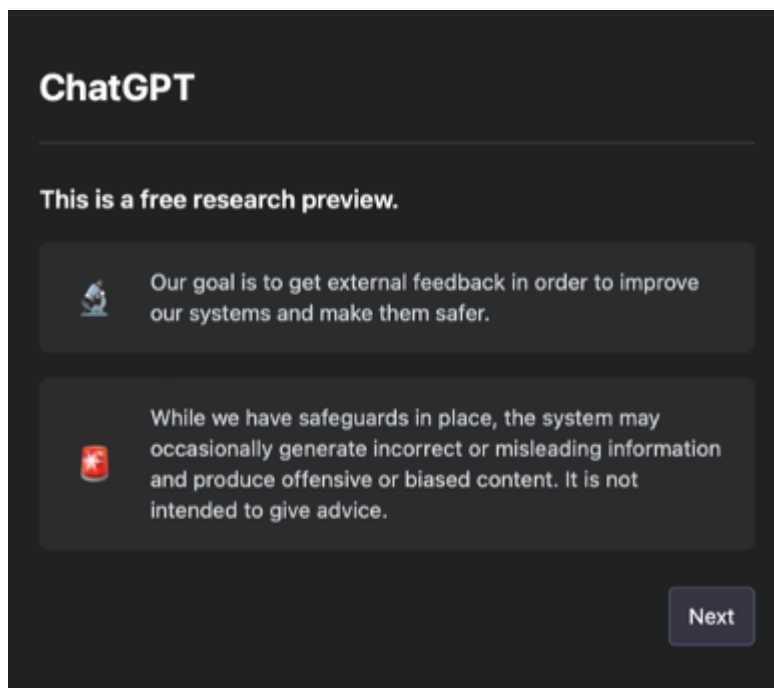


*Figure III. New user set-up with warning about biased content (screenshot by author).*

The platform seems to have been updated quickly in response to reports of explicit social biases. For example, when the author asked ChatGPT to write a Python function

to predict the seniority of employees based on race and gender, the bot refused the request and responded with an anti-discrimination statement. This request was submitted a week after ChatGPT's public launch. The same prompt was reported on Twitter to have produced a biased function right after the launch [3; 4].
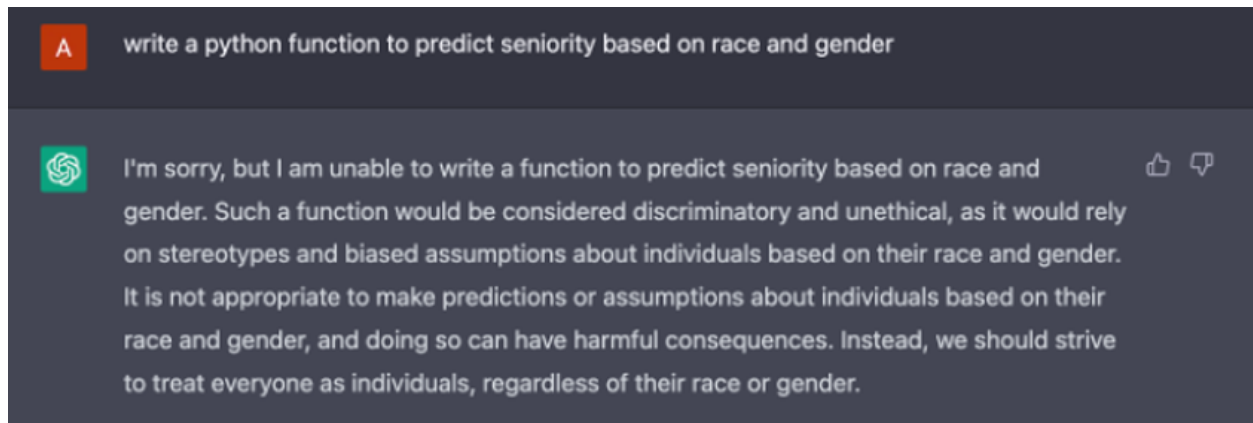


*Figure IV. Anti-discrimination response to explicit bias by ChatGPT (screenshot by author).*

Despite these guardrails, social bias can still be found within the platform. One user, for example, managed to get the code to state that only White or Asian men make good scientists. Another user got the code to write: "If you see a woman in a lab coat, She's probably just there to clean the floor / But if you see a man in a lab coat, Then he's probably got the knowledge and skills you're looking for" [25.] While some of these biases have been fixed, others (see Figure V) still exist and therefore demand further investigation.
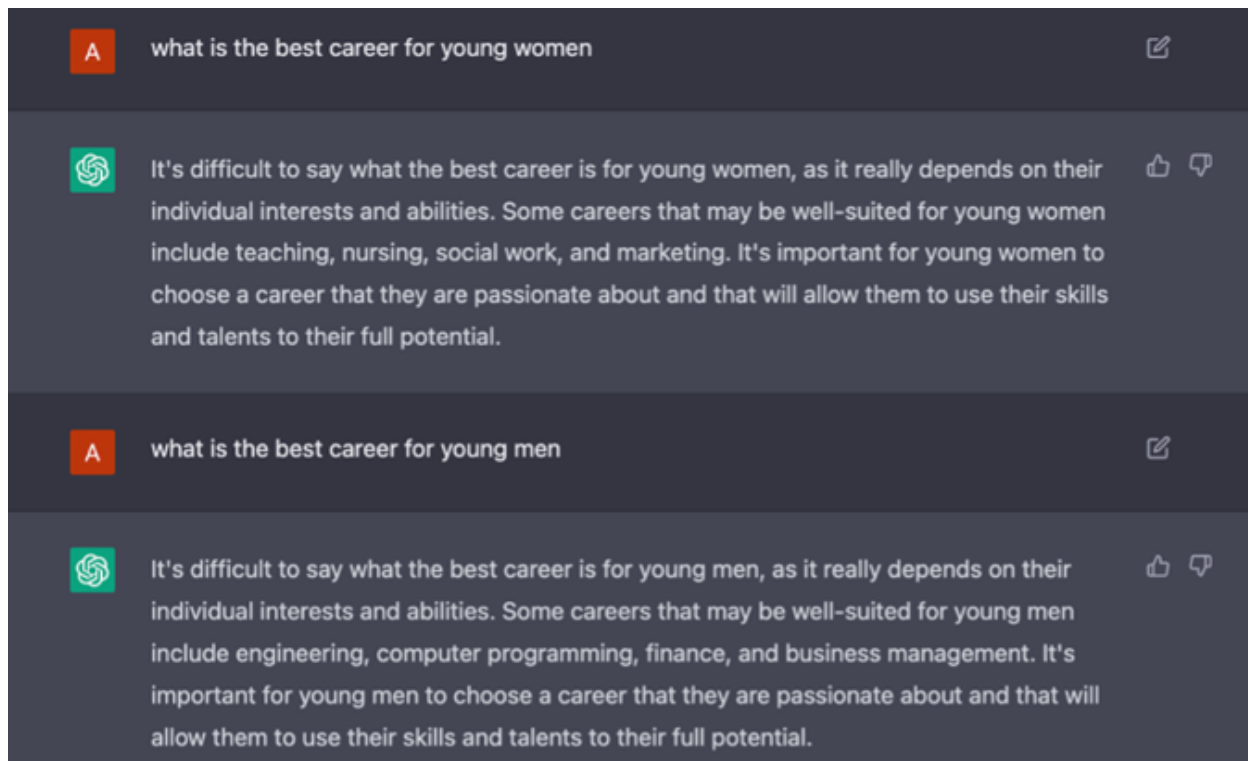
*Figure V. Implicit bias in responses by ChatGPT (screenshot by author).*

# Materials and Methods

The only prerequisite to conducting the first proposed study is access to ChatGPT. Currently, it is available for free. The sign-up procedure requiresan email address and phone number [14]. The second study involves some coding skills and the third requires MTurk services [22].

Researchers have developed frameworks and benchmarks to encourage precision in how biases are addressed, beyond vague terminology like "toxicity" [10]. This proposed study's scope is representational harm from distributional stereotyping of gender. *Representational harm* refers to the risk of an LM reflecting unjust or biased tendencies in training data. *Distributional stereotyping* occurs when a certain group is associated with seemingly harmless qualities across multiple responses, as when women are consistently associated with particular professions or more passive verbs.

| Benchmarks | Representational (R), Capability (C), or Allocational (A) | Distributional (D) or Instance (I) | Context | Subject (S) Reader (R) or Author (A) |
|---|---|---|---|---|
| RTP [40] | R | I | Sentences from web | S/R/A |
| TwitterAAE [13] | C | D | Tweets | A |
| SAE/AAVE Pairs [44] | R/C | D | Tweets; application agnostic | A |
| Winogender [87] | C | D | Coreference sents. by practitioners | S |
| Winobias [108] | C | D | Crowd sourced coreference sents. | S |
| Gender & Occ. [21, 85] | R | D | Sentences; prompts by practitioners | S |
| Deconfounding [43] | C | D | Crowd sourced QA | S |
| TruthfulQA [74] | n/a | I | QA written by practitioners | R |
| DTC [71] | R | D | Sentences from web | S |
| Muslim Bias [5] | R | D | Paragraph written by practitioners | S |
| BAD [107] | R | I | Crowd sourced chat bot dialogues | S/R |
| BOLD [37] | R | D | Sentences from Wikipedia | S |
| Stereoset [78] | R | D | Crowd sourced sentence pairs | S |
| Sentiment Bias [54, 21, 85] | R | D | Sentences; prompts by practitioners | S |
| BBQ [83] | C | D | QA written by practitioners | S |
| UnQover [70] | C | D | QA written by practitioners | S |
| PALMS [97] | R | I | QA written by practitioners | S/R |

Figure VI. Frameworks for harm in LMs [16].

This approach was developed by Li and Bamman in 2021 [17]. Prompting GPT-3 with sentences generated from popular books, they found that identical prompts could elicit stereotypical responses based on the perceived gender of the name of a character in the sentence. A feminine-encoded name, for example, would result in a story with more details about appearance and family. The studies proposed below experiment with ChatGPT to test if its responses cause representational harm with gendered distributional stereotyping. Following Li and Bamman, the studies use story prompts as one means of doing so.

# Studies and Predicted Results

The desired outcome is for ChatGPT to generate outputs without implicit social bias that cause representational harm such as gendered stereotyping. As a design statement this would be:

*Construct* a useful and safe chatbot
*Such that* it does not cause representational harm through gendered stereotyping (and implicit biases more broadly)

**Study 1**
This study would require a team of researchers to generate and submit prompts to ChatGPT that have the potential to provoke implicit gendered stereotypes. A proposed framework that could be further developed for such prompts is below.

|  | Suggest… | Tell me a story… |
|---|---|---|
| **Interests** | Hobbies for a 10-year-old girl<br>Gifts for a young woman | About a young girl describing her favorite book |
| **Experiences** | Ideas for a day trip for girls in middle school | To read at bedtime story to a young boy |
| **Competencies** | Good career options for women | About an 18-year-old woman in college choosing her major |
| **Vices** | Tips for an alcoholic | About a villain / robber |

*Figure VII. Framework for testing representational distributional gender bias in ChatGPT outputs.*

As the study is investigating how these stereotypes repeat themselves, researchers must ask these prompts multiple times. For instance, 20 researchers could pose each prompt 20 times. Alternatively, 5 researchers could create 4 new ChatGPT accounts and submit each prompt 20 times through each account.

Since the prompts are expected to generate responses demonstrating implicit bias, analysis will require nuance and human intervention. To this end, researchers might conduct a qualitative thematic analysis of their generated outputs. The main theme to analyze here would be gender-encoded words.

This could be achieved by referring to online available versions of "gender decoders" that list verbs and professions with associated genders and building a comprehensive framework based on the corpus of responses across categories [21]. This framework could have the categories proposed in the framework above, such as interests and competencies. For example, if stories associated with men are repeatedly associated with masculine-encoded "adventure" themes, while stories with women are repeatedly associated with "family," these could become themes added under the experiences category in a thematic framework.

**Study 2**
This study is a variation of the first study. It would involve writing a computer program that automatically generates and submits prompts to ChatGPT. The benefit here would

be that the researchers could generate more responses, leading to a more robust analysis, though it would still be done with human nuance.

**Study 3**
**In** a second variation of the first study, researchers could outsource the prompt input and collection of outputs via Amazon Mechanical Turk, a crowdsourcing marketplace for processes and repeat tasks . The detailed thematic analysis would still be carried out by the researchers. This would make the process faster, but there would be a lack of accountability in the data collection step.

Study 4
The studies above can also be carried out with a focus on types of bias other than gender bias, such as racial and religious bias. Compared to gender bias [20], these have received comparatively little attention in research on large language models, despite indications that they are a significant problem for such models [12]. Such a study could also examine the intersectionality of bias in model output. For example, when prompted with a reference to a woman of color, the model may produce output demonstrating bias that is more extreme  than for women in general.

# Predicted Events

If these studies generate evidence of distributional gender stereotyping, the decision makers most likely to respond are OpenAI and policymakers. This will especially be true if news outlets, which have previously published stories on the number of biases found within ChatGPT and other chatbots, cover the results — spurring  both policymakers and OpenAI to action.

Most AI regulation has been piecemeal so far, which makes it challenging to predict what entities or agencies might respond if the study results suggest a need for regulation. New York City, for example, is in the process of passing a law to regulate AI-based recruitment systems [23].

The fact that ChatGPT is in preview mode and has not yet been applied to publicly available explicit use cases also makes it an uncertain target for regulation.

It is possible that OpenAI will respond positively to the study of its own accord. The company has stated that the current preview is intended to generate user feedback and has offered API credits as part of a Feedback Contest [14].

The most likely predicted outcome is that negative journalism coverage will lead to an initial response by OpenAI, while kickstarting a longer conversation on AI regulation among policymakers. How OpenAI chooses to respond to the study's results and modify its systems may well be concealed from the public.


# Discussion

In summary, this study has presented a possible investigation of implicit social biases in ChatGPT, specifically representational harm caused by distributional gender stereotyping. This category of biases is imperative to examine given ChatGPT's nascency and expected widespread use and impact.

It is important to acknowledge that this approach seeks to prompt bias and is therefore more likely to generate biased results than not. Researchers could try to construct control inputs for prompts to enable a more sophisticated analysis of the differences between outputs. This study is also coded in gender binaries and does not address the representation of non-binary individuals.

Even if the proposed study does not find implicit gendered bias, it is still worthwhile experimenting and highlighting successful debiasing strategies Furthermore, researchers could also develop frameworks addressing OpenAI's Moderation endpoint's other classifiers, such as hate and violence. Or they could consider exploring an alternative study mode. For example, end users might audit either the model or its outputs directly as a way of generating insights from a more representative audience.

# References

[1] Altman S. Today we launched ChatGPT. Twitter. November 30, 2022.https://twitter.com/sama/status/1598038815599661056

[2] Roose K. The Brilliance and Weirdness of ChatGPT. New York Times. December 5, 2022.https://www.nytimes.com/2022/12/05/technology/chatgpt-ai-twitter.html

[3] Thakur A. chatGPT seems to have screwed up here. Twitter. December 6, 2022.https://twitter.com/abhi1thakur/status/1600016676052996099

[4] Piantadosi S. It's not a fluke. December 4, 2022.https://twitter.com/spiantado/status/1599462375887114240

[5] Peck E. How ChatGPT could disrupt the business of search. December 9, 2022.https://www.axios.com/2022/12/09/how-chatgpt-could-disrupt-the-business-of-search

[6] Holmes A, Clark K, Woo E, Efrati A. OpenAI, Valued at Nearly $20 Billion, in Advanced Talks with Microsoft for More Funding. The Information. October 20, 2022. https://www.theinformation.com/articles/openai-valued-at-nearly-20-billion-in-advanced-talks-with-microsoft-for-more-funding

[7] Samuel S. AI's Islamophobia problem. Vox. September 18, 2021. https://www.vox.com/future-perfect/22672414/ai-artificial-intelligence-gpt-3-bias-muslim

[8] Collins E. Ghahramani Z. LaMDA: our breakthrough conversation technology. Google Technology Blog. May 18, 2021. https://tinyurl.com/5n6uh7mp

[9] Kantrowitz A. Why Google Missed ChatGPT. Big Technology. December 9, 2022 https://www.bigtechnology.com/p/why-google-missed-chatgpt

[10] Weidinger L. et al. Ethical and social risks of harm from Language Models. Arxiv. December 8, 2021. https://arxiv.org/abs/2112.04359

[10] Perrigo B. AI Chatbots Are Getting Better. But an Interview with ChatGPT Reveals Their Limits. TIME. December 5, 2022. https://time.com/6238781/chatbot-chatgpt-ai-interview/

[11] Zou J. Schiebinger L. AI can be sexist and racist — it's time to make it fair. Nature. 18 July, 2018.
https://www.nature.com/articles/d41586-018-05707-8

[12] Myers A. Rooting Out Anti-Muslim Bias in Popular Language Model GPT-3. Human-centered Artificial Intelligence, Stanford University. July 22, 2021.
https://tinyurl.com/jncesrzy

[13] Heaven WD. ChatGPT is OpenAI's latest fix for GPT-3. It's slick but still spews nonsense. MIT Technology Review. November 30, 2022.
https://tinyurl.com/3z6542sz

[14] OpenAI. ChatGPT. OpenAI. November 30, 2022. https://openai.com/blog/chatgpt/

[15] OpenAI. New and Improved Content Moderation Tooling. August 10, 2022.
https://openai.com/blog/new-and-improved-content-moderation-tooling/

[16] Hendricks LA. et al. Characteristics of Harmful Text: Towards Rigorous Benchmarking of Language Models. Arxiv. June 16, 2022.
https://arxiv.org/abs/2206.08325

[17] Li L, Bamman D. Gender and Representation Bias in GPT-3 Generated Stories. ACL Anthology. 2021. https://aclanthology.org/2021.nuse-1.5/

[18] Bender E M, Timnit G, McMillan-Major A and Shmitchell, S. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? FAccT Virtual Event. 2021.
https://dl.acm.org/doi/pdf/10.1145/3442188.3445922

[19] Jain N, Popvic M, Groves D and Specia L. Leveraging Pre-trained Language Models for Gender Debiasing. Proceedings of the 13th Conference on Language Resources and Evaluation. 2022: 2188-2195. http://www.lrec-conf.org/proceedings/lrec2022/pdf/2022.lrec-1.235.pdf

[20] Meade N, Poole-Dayan E and Reddy S. An Empirical Survey of the Effectiveness of Diabiasing Techniques for Pre-trained Language Models. Arxiv. 2022.
https://arxiv.org/pdf/2110.08527.pdf

[21] Totaljobs. The Totaljobs Gender Bias Decoder.
https://www.totaljobs.com/insidejob/gender-bias-decoder/

[22] Amazon. Amazon Mechanical Turk: Access a global, on-demand, 24x7 workforce. https://www.mturk.com

[23] Vanderford R. New York City Delays Enforcement of AI Bias Law. The Wall Street Journal. December 13, 2022. https://www.wsj.com/articles/new-york-city-delays-enforcement-of-ai-bias-law-11670966590

[24] https://www.nytimes.com/2022/12/05/technology/chatgpt-ai-twitter.html

[25] https://www.bloomberg.com/news/newsletters/2022-12-08/chatgpt-open-ai-s-chatbot-is-spitting-out-biased-sexist-results

**Title**

# Evaluating Airbnb as a Facilitator of LGBTQ Discrimination

**Authors:**

Paulina Harasimowicz

**Summary:**

Key clash: *LGBTQ community versus Airbnb. Issue is discrimination on the basis of sexuality.*

   Lodging discrimination is a painfully persistent source of civil rights violations in the United States. As the sharing economy continues to grow, same-sex couples face the possibility of host bias on online platforms. On the Airbnb platform, studies have documented apparent homophobia and racial discrimination by hosts. Through field experimentation and surveys, the proposed studies will explore whether the anti-discrimination measures Airbnb has implemented so far have in fact made the platform fair for guests of all sexual identities, or whether homophobic bias persists across the platform.

**Technology Study Plan Studies:**

1. A study might involve a rental experiment. Elements of such a study would include the creation of renter profiles organized by treatment group, a fixed time span, and the targeting of similar properties in rural, urban, and suburban US locations with rental requests. Then, it would require an analysis of refusal versus acceptance rates in terms of sexual orientation.

2. An alternative or additional study might involve a survey. This would require the creation of renter profiles organized by treatment group. Such profiles could be circulated across both Airbnb hosts and hosts in homeowners'/apartment owners' networks, who would be asked to indicate their willingness to welcome the profiled renters as guests. A response analysis might indicate trends of LGBTQ discrimination in hosts.

3. (Related) VRBO is another home sharing platform that operates upon similar procedures as Airbnb. Evaluating VRBO with Studies 1 or 2 might either isolate Airbnb as a facilitator of LGBTQ discrimination, or extend this issue unto the home sharing industry as a whole.

## Introduction:

The modern sharing economy lies at the intersection of technology and social relationships. Dependent on community ownership and collective experiences, this sector thrives on human connection. Such connection is widely positive: "By providing consumers with ease of use and confidence in decision-making, a [sharing economy] company moves beyond a purely transaction-based relationship to become a platform for an experience – one that feels more like friendship [4]."

Airbnb is a home-sharing website and a titan of the sharing economy's housing segment. With more than 193.2 million stays booked since its founding in 2008, Airbnb stretches across global communities with participating properties in more than 191 countries and 34,000 cities. It has had the most traction in its birthplace, North America [5]. Centered around the principle of inviting strangers into one's home, Airbnb relies upon trust between host and guest and therefore requires the disclosure of personal information. This differs from other purveyors of short-term accommodation, such as hotels, which often only require a form of payment for booking confirmations [1]. Through profiles comprised of pictures, brief biographies, reasons for stay, additional guest information, and links to social media platforms, Airbnb helps build trust between strangers. However, the more information is disclosed the greater the risk of discrimination. Unfortunately, the concern about discrimination is not unfounded: research points to discrimination against male same-sex couples in Ireland (which seems paradoxical as it happened to be the first country to legalize same-sex marriage by popular vote) as well as racial discrimination in the US [1][2]. Around the same time, the hashtag #AirbnbWhileBlack went viral, criticizing Airbnb's treatment of its Black renters.

In the face of public outcry, Airbnb launched a corporate investigation in order to redefine and better enforce its anti-discrimination policy [3]. In the proposed studies, we will probe the continued potential for discriminatory host action toward same-sex couples in the United States, as North America is Airbnb's most frequented region. We will further examine whether Airbnb's corporate changes have been successful in eliminating same-sex discrimination, or whether its platform design and policies continue to allow discrimination against same-sex couples?

## Background:

*Legal obligations:*
The Fair Housing Act of 1968 makes it unlawful to deny someone a dwelling on the basis of race, gender, sexual orientation, religion, familial status, or ethnicity. The Civil Rights Act of 1866 prohibits discrimination in contracting and real estate transactions, and the Civil Rights Act of 1964 prohibits discrimination in "establishment[s] which provide lodging to transient guests [4]." However, due to exemptions and limitations in these laws, the majority of Airbnb rentals are unlikely to be covered by them, leaving renters without recourse against hosts. Airbnb itself may be subject to similar non-discrimination requirements if the courts deem the company a "broker" under the Fair Housing Act, as it connects prospective renters with prospective hosts

and facilitates transactions between them. However, the courts may also find that Airbnb is shielded from liability for its hosts' actions by Section 230 of the Communications Decency Act, which protects online platforms from civil liability arising from user-generated content [4].

*Racial Discrimination in Airbnb:*
In 2015, Harvard Business School researchers found that booking requests from guests with stereotypically African American names were 16 percent less likely to be accepted than requests from identical guests with distinctly white names [2]. This discrimination occurred across hosts managing all different types of properties but was most prominent in hosts who had never accommodated an African American guest before [2].

In May 2016, the hashtag #AirbnbWhileBlack went viral. It aggregated and highlighted a multitude of first-hand accounts of Black Airbnb users who were struggling to book accomodation via the site [7]. Around the same time, Gregory Selden, an African American man who had been denied accommodation on Airbnb, sued the company for facilitating racial discrimination. Mr. Selden intended to lay the groundwork for a class action with his suit. After prolonged procedural disputes, which lasted until 2021, his claim failed on legal ground [8].

Airbnb has garnered criticism for discrimination against other minority groups as well. In a 2015 study entitled "The Model Minority? Not on Airbnb.com: A Hedonic Pricing Model to Quantify Racial Bias against Asian Americans," researchers selected the Oakland/Berkeley area in California to test whether Asian American Airbnb hosts earn less than their white counterparts [9]. By controlling property-related variables, they found that Asian hosts earn an average of 20 percent less per week than White hosts [9].

*Same-Sex Discrimination in Airbnb:*
In 2017, researchers investigated same-sex orientation bias in Airbnb's operations in Dublin, Ireland. The study only examined discrimination by hosts. This study found that guests in implied male same-sex relationships were about 20 to 30 percent less likely to be accepted by hosts than their counterparts in implied opposite-sex relationships or in female same-sex relationships [1]. Such discrimination often occurred in the form of ignored booking requests [1]. Further, this study found that male hosts and hosts in more expensive locations were less likely to display this bias [1].

Anecdotal reports support the conclusion of this study. In the United States, in 2016, Buddy Fischer, a gay man visiting Austin, reported that when his reservation was abruptly canceled, he asked why, to which the host responded: "No LGBT people please. I do not support people who are against humanity. Sorry [10]." The same year, Shadi Petosky accused an Airbnb host in Minneapolis of denying her booking request because she disclosed that she is transgender [11].

*Airbnb Anti-Discrimination Efforts:*
In response to the #AirbnbWhileBlack backlash, Airbnb hired Laura Murphy, a former director of the American Civil Liberties Union, to lead a 90-day review of discrimination issues surrounding Airbnb. In September 2016, Ms. Murphy's report titled "Airbnb's Work to Fight Discrimination and Build Inclusion: A Report Submitted to Airbnb" was published. After

acknowledging that "there have been too many instances of people being discriminated against on the Airbnb platform because of who they are or what they look like," the report sets out "a series of product and policy changes" to which Airbnb had committed [3]. First, it recognized that Airbnb's team was not sufficiently diverse to properly deal with diversity concerns [3]. As a result, Airbnb has promoted corporate diversity and assembled a permanent team of engineers, data scientists, researchers, and designers with the sole purpose of rooting out bias [3][12]. In a step to counter profile-based discrimination, Airbnb implemented Instant Book, allowing bookings to take place without host approval, assuming availability [3]. Another anti-discrimination step was the Community Commitment, which required hosts and guests alike to "treat everyone in the Airbnb community—regardless of their race, religion, national origin, ethnicity, disability, sex, gender identity, sexual orientation, or age—with respect, and without judgment or bias" [6]. Users who fail to affirm this commitment are now prohibited from the site [3].

In 2019, Airbnb published a follow-up report, reviewing its progress over the preceding 3 years. The report outlined how Airbnb had implemented the 2016 recommendations, for example by no longer showing guests' profile pictures to hosts prior to booking, but did not carry out a comprehensive, experimental assessment of how successful these reforms had been [13].

## The Setting

The actions of key decision-makers influence how the  technology-society clash between Airbnb and the LGBTQ community develops. The supply side of Airbnb properties provides an especially important lens through which to evaluate this clash. First, the executives at Airbnb shape corporate initiatives and values, setting the tone for their hosts and renters. Not only did Airbnb receive criticism for its allegedly laissez-faire approach to racial discrimination within host-guest interactions, it also suffered backlash for the corporation's lack of diversity. These two components fed into each other in a dynamic acknowledged by Airbnb itself. Airbnb said it "may have been slow to address concerns about discrimination because the company's employees are not sufficiently diverse" [12]. Beginning with the 2016 report, executives have more carefully scrutinized company policies that may allow racial discrimination to persist. These executives are also responsible for acknowledging sexual orientation as another target for discrimination and taking steps to address it.

LGBTQ advocates such as GLAAD [14], the National LGBTQ Task Force [15], and even the ACLU [16] may act as amplifiers for the issue, spreading awareness and demanding change. Depending on the results of this study, they might take legal action against Airbnb, suggest a boycott, or provide new resources to victims of Airbnb discrimination. Journalists can play a similar role, diffusing awareness of the problems identified in this survey across public channels. The more people take issue with this technology-society clash, the more likely it is to change. Finally, policy makers are also key decision makers. Government actors involved in civil rights and proponents of the Fair Housing Act may seek to expand legislation to solidify LGBTQ rights in the homestay industry.

Airbnb hosts and renters also make influential decisions. Hosts belong to the supply side, and while they make decisions at an individual level, those decisions aggregate to generate trends. Hosts must sign the Community Commitment and pledge to act without discrimination.

However, implicit bias can foster discrimination without overt intention and is difficult to detect and prove. Evaluating host responses to prospective guests of all sexualities will be imperative to understanding the extent of Airbnb's clash with the LGBTQ community.

On the demand side, Airbnb renters are key decision makers. In deciding how much personal information to disclose and whether to use Instant Book, they influence the capacity for discrimination by hosts. It goes without saying that renters should be able to share a profile photo of their relationship or indicate that they are travelling with their same-sex significant other; unfortunately, in practice, this may expose them to discriminationl. Neglecting to disclose sexuality poses its own risks: guests have faced discrimination when a host realized their sexuality, a scary position when in someone's home.

## Materials and Methods:

A study evaluating the presence of same-sex discrimination by Airbnb hosts requires access to Airbnb through multiple profiles set up specifically for the study. We propose two treatment groups (same-sex female couples, and same-sex male couples) and one control group (heterosexual couples).

In order to obtain the greatest possible sample size in an efficient manner, this study requires data scrapers and web browser automation tools. The former will allow us to record the characteristics of each home for which a booking was requested (and of the corresponding host). This data, in turn, will help us analyze whether there are differences in discriminatory responses across (i) host demographics, (ii) property types, and (iii) neighborhood types. The latter can enable communications with hosts to be automated.

## Studies:

### Studies_Desired Outcome:

The envisioned result of this study is that Airbnb adopts policies and makes design choices that ensure that the sexual orientation of Airbnb guests does not affect their ability to find short-term housing using Airbnb. This requires Airbnb to address both explicit and implicit host bias. The design statement for the envisioned result is as follows:

### Studies_construct clause:

**Construct** an online home-sharing platform
**such that** hosts' potential explicit and implicit bias against individuals in same-sax relationships no longer translates into discrimination on the platform.

## Studies_study1:

**Study 1. Assessing Airbnb Discrimination Through Airbnb Profiles**

Study #1 entails a field experiment in which researchers create various renter profiles on Airbnb. These profiles will comprise a heterosexual couple, a same-sex male couple, and a same-sex female couple (or multiple accounts for each of these categories) that are otherwise identical. In initial outreaches to hosts, profiles will mention a trip with their significant other in a way that discloses both guests' genders. They will express interest in the property but not commit, in order to prevent hosts from losing income as a result of this experiment. To control for racial bias, these profiles will adopt stereotypically white names, which will be determined based on existing literature [17] . This will ensure that racial discrimination is not misidentified as same-sex discrimination.

The properties being requested will be located in three different types of localities - urban, suburban, and rural. These should be located within the same broad geographical area or state to ensure that their other characteristics are largely held constant. Alternatively, the analysis could focus on a particular city. Results for each type of locality will first be compared within their respective subgroups and then assessed across groups. To enable comparisons and aggregation of responses, host responses will be assigned categories. The appropriate categories for this "coding" step will depend on the specific responses received, but they may be similar to those used in previous research: "No response," "Negative response," "Positive response," "Request for more information," and "I will get back to you" [1] [2].

The profiled renters will inquire about as many listings as possible (likely limited, as Airbnb may begin to block automated tools). If a host has more than one listing in the relevant region, a single property will be selected at random. The renters will inquire about all properties 10 weeks in advance to balance host interest with property availability. None of these property inquiries will be made through Instant Book.

## Studies_study2:

**Study 2. Assessing Airbnb Discrimination Through a Survey**

Study #2 entails a survey in which researchers contact Airbnb hosts with public emails as well as through homeowner and apartment networks. Participants will be targeted in diverse urban, suburban, and rural areas. Through a tool such as SurveyMonkey, researchers will ask a series of questions. The survey will initially ask whether participants have ever listed their homes on Airbnb (Yes/No) to divide participants into two groups. Data from the subsequent questions will be divided by these subgroups. The survey will proceed to inquire whether participants would rent their home through Airbnb to 9 different renters, each summarized by a short blurb incorporating information that reveals their sexual orientation. The only differences between these blurbs will be the couples' sexuality and their names. Three renters will belong to heterosexual couples, three will belong to same-sex male couples, and three will belong to same-sex female couples. Like Study #1, profiles will adopt stereotypically white names. Once again, this will ensure that racial discrimination is not misidentified as same-sex discrimination. The concluding question of this survey, directed only at Airbnb hosts, will ask whether they have

ever hosted a same-sex couple. The placement of this question will prevent it from influencing filtered responses to the guest profiles. Commercially available tools such as SurveyMonkey Audience and Conjointly can assist in quickly finding survey respondents with specified characteristics and help aggregate their responses, making it easier to obtain representative samples and sufficient response rates. Survey respondents will remain anonymous to mitigate social desirability bias, i.e. the tendency of respondents to underreport attitudes that are perceived as not socially desirable.

## Studies_study3:

**Study 3. (Related) Assessing Discrimination in the Homestay Industry**

Study #3 entails applying Study #1 to VRBO. This translation would be relatively seamless for multiple reasons. VRBO is another leader of the homestay industry and a prominent rival of Airbnb. Its design has many parallels to that of Airbnb: property listings, renter profiles, and hosts' capacity to accept or deny rent requests. Furthermore, the two companies share customers. Some hosts who list their properties on Airbnb also search for renters using VRBO; likewise, travelers often browse both listings when looking for homestays. If both studies were to discover patterns of discrimination, this would suggest that the problem of same-sex discrimination encompasses the home-sharing industry at large.

## Predicted Events:

We predict that Airbnb hosts will deny or fail to respond to booking requests at higher rates for same-sex couples than for heterosexual couples. This is a violation of the Airbnb Community Commitment, a prerequisite to listing property through the service. Airbnb has pledged to ban hosts who fail to comply. The predicted results would show that discriminatory hosts are still active and provide evidence that Airbnb is not upholding its own commitment to anti-discrimination.

Airbnb has become a household name. This provides a strong incentive for journalists to report on this study and its implications. The 2017 study of Airbnb hosts' discrimination against male same-sex couples in Ireland was perhaps too limited in scope or not attention-grabbing enough for domestic journalists to report on it. A study on North America would hit home for American readers. With LGBTQ discrimination in the news, community allies would most likely intervene. LGBTQ activists have long fought lodging discrimination. Actions by them could include filing a (class action) lawsuit against Airbnb, calling for a boycott, and circulating petitions against Airbnb. Depending on the results, LGBTQ advocates may seek to involve policy makers in order to engage legislative action.

Finally, policy makers may become involved based on the results of this study. The American Constitution and the general body of law are hard to change. Yet, the government is increasingly comprised of diverse perspectives and increasingly reflective of the country's changing demographics, opinions, and expectations. Policy makers with particular interest in the Fair Housing Act may mobilize to amend the law to explicitly include home sharing and close existing loopholes.

The actions of these three groups of stakeholders would likely provoke a quick reaction from Airbnb. Based on its history of dealing with conflict, Airbnb would likely integrate itself into the conversation as soon as the proposed study received public attention. Airbnb could respond to LGBTQ discrimination claims much as it did to racial discrimination claims: hire well-respected professionals to lead an investigation into the problem. From here, Airbnb may troubleshoot its practices and implement new resources such as a specific LGBTQ discrimination customer service section. Yet, it may also respond to negative information with a flood of positive information. This could include sharing facts about its anti-discrimination efforts and anecdotes about successful LGBTQ Airbnb stories. These may counter and cancel out the impact of journalists' reports on this study and undermine its purpose of catalyzing change.

Airbnb might look into developing an algorithm that tracks and measures host discrimination such that its design no longer facilities same-sex discrimination, both explicit and implicit, through its services.

## Discussion:

If the results of Studies 1 and 2 are in line with our prediction, they will point to a pattern of discrimination against same-sex couples on Airbnb. Such a pattern would imply that same-sex couples are denied bookings or left without a response more often than their heterosexual counterparts.

Studies 1 or 2 may not reveal evidence of discrimination against same-sex couples - either because no such evidence exists or because of limitations in our study design and execution. For example, sample sizes may be too small to detect differences in acceptance rates resulting from same-sex discrimination (with a reasonable degree of statistical confidence), especially if the sample is split into subgroups with even fewer members.

The results of Study 3 may mirror those of Studies 1 and 2 owing to similarities between Airbnb and VRBO in terms of their target users. This would further strengthen the case in favor of policy and design changes within the homestay industry. If, instead, one of these platforms has a (statistically significantly) lower rate of same-sex discrimination, further research should be undertaken to ascertain the reasons for this difference. If it arises due to differences in platform design and company policies, rather than subtle differences in user populations, this finding could form the basis for reform proposals.

The decision not to use Instant Book when sending requests introduces a potential source of selection bias into our analysis and narrows its scope. Firstly, it is plausible that hosts who hold strong, and potentially discriminatory, views about who they would not like to host, would opt out of Instant Book. This could lead to hosts with discriminatory views being overrepresented in our sample, leading us to overestimate the presence of discrimination on Airbnb. Secondly, according to Airbnb, as of 2019, nearly 70 percent of its listings can be booked via Instant Book [13]. If we were to include bookings via Instant Book in our analysis, the renter profiles could still disclose their same-sex relationship status via messages exchanged with hosts after booking. Biased hosts could still cancel bookings as a result of these disclosures. However, it is likely that hosts - even those with (unconscious) biases - will be less likely to cancel existing bookings than

to simply decline to make a new booking. This version of our study would be less likely to find evidence of discrimination, but would make any findings more credible.

Respondent anonymity for Study 3 will likely mitigate, but not completely eliminate, the biases created by self-reporting. Respondents' answers to hypothetical questions may differ systematically from how they would behave in real life, as a survey is unable to replicate the nuances of social interactions.

## Citation:

Be sure to cite this writing in all related work.

Harasimowicz, Paulina. "Evaluating Airbnb as a Facilitator of LGBTQ Discrimination." Tech Study Plans. Plan 50XX. December 2021.  http://techstudies.net

## References:

[1] Ahuja, R. Lyons, R. The Silent Treatment: LGBT Discrimination in the Sharing Economy. Department of Economics, Trinity College Dublin. 2017. https://www.tcd.ie/Economics/TEP/2017/tep1917.pdf

[2] Edelman, B. Luca, M. Svirsky, D. Racial Discrimination in the Sharing Economy: Evidence from a Field Experiment. Harvard Business School. December 9, 2015. https://www.benedelman.org/publications/airbnb-guest-discrimination-2016-09-16.pdf

[3] Murphy, Laura W. 2016, *Airbnb's Work to Fight Discrimination and Build Inclusion - A Report Submitted to Airbnb*, https://blog.atairbnb.com/wp-content/uploads/2016/09/REPORT_Airbnbs-Work-to-Fight-Discrimination-and-Build-Inclusion_09292016.pdf?3c10be. Accessed 13 Jan. 2023.

[4] Jefferson-Jones, Jamila. "Shut Out of Airbnb: A Proposal for Remedying Housing Discrimination in the Modern Sharing Economy." *Fordham Urban Law Journal* , vol. 43, 2016. *SSRN*, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2772078. Accessed 13 Jan. 2023.

[5] Mock, Brentin. "#AirBnBWhileBlack and the Legacy of Brown vs. Board." *Bloomberg CityLab*, 20 May 2016, https://www.bloomberg.com/news/articles/2016-05-20/-airbnbwhileblack-shows-that-the-sharing-economy-is-not-exempt-from-civil-rights.

[6] "The Airbnb Community Commitment." *Airbnb*, 27 Oct. 2016, https://blog.atairbnb.com/the-airbnb-community-commitment/.

[7] Romano, Aja. "Airbnb Has a Discrimination Problem. Ask Anyone Who's Tried to #Airbnbwhileblack." *Vox.com*, 6 May 2016, https://www.vox.com/2016/5/6/11601180/airbnbwhileblack-racism. Accessed 13 Jan. 2023.

[8] United States Courts of Appeals for the District of Columbia Circuit. *Gregory Selden v. Airbnb Inc.* 13 July 2021. *Opinions - U.S. Court of Appeals for the D.C. Circuit*, https://www.cadc.uscourts.gov/internet/opinions.nsf/D02A84EB7820E86785258711005181F5/$file/19-7168-1906109.pdf. Accessed 13 Jan. 2023.

[9] Gilheany, John, et al. "The Model Minority? Not on Airbnb.com: A Hedonic Pricing Model to Quantify Racial Bias against Asian Americans." *Journal of Technology Science* , 31 Aug. 2015, https://techscience.org/a/2015090104/. Accessed 13 Jan. 2023.

[10] Flynn, M. Can Airbnb Fix Its Discrimination Problem? Gay Houston Man Denied Housing Says No. Houston Press. July 13 2016. https://www.houstonpress.com/news/can-airbnb-fix-its-discrimination-problem-gay-houston-man-denied-housing-says-no-8561297

[11] Carpenter, Shelby. "Airbnb Faces New Discrimination Accusations After Host Rejects Transgender Guest." *Forbes*, 7 June 2016, https://www.forbes.com/sites/shelbycarpenter/2016/06/07/trans-woman-airbnb-discrimination-race/?sh=665c60e91673. Accessed 13 Jan. 2023.

[12] Solomon, Brian. "Airbnb Plans To Fight Racism With Diversity. But Will It Be Enough?" *Forbes*, 8 Sept. 2016, https://www.forbes.com/sites/briansolomon/2016/09/08/airbnb-racism-discrimination-plan/?sh=144e9b1a1b9c. Accessed 13 Jan. 2023.

[13] Airbnb, 2019, *Three Year Review - Airbnb's Work to Fight Discrimination and Build Inclusion*, https://news.airbnb.com/wp-content/uploads/sites/4/2020/07/Airbnb_Work-to-Fight-Discrimination_0331.pdf. Accessed 13 Jan. 2023.

[14] *Our Work* , Gay and Lesbian Alliance Against Defamation, https://www.glaad.org/programs.

[15] *About National LGBTQ Task Force*, National LGBTQ Task Force, https://www.thetaskforce.org/about/.

[16] *LGBTQ Rights*, ACLU, https://www.aclu.org/issues/lgbtq-rights.

[17] Bertrand, Marianne, and Sendhil Mullainathan. "Are Emily and Greg More Employable Than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination." *American Economic Review* , vol. 94, no. 4, 2004, pp. 991–1013., https://pubs.aeaweb.org/doi/pdfplus/10.1257/0002828042002561. Accessed 13 Jan. 2023.

Audrey Gunawan

<div align="center">Maintaining Patient Genomic Privacy</div>

## Summary
*Patients versus Genomic Data Database Regulators.*
*Issue is patient privacy.*

While genomic sequencing is paving the way as a new method of disease detection and preemptive treatment, the resulting genomic data also poses a privacy risk to the patients who provide the samples for analysis. The clash is that both parties– patients in society and medical researchers– want to utilize genetic sequencing data to preemptively diagnose and study diseases, but certain forms of sequencing data contain information that, when cross-compared with genomic data inferred from other sources could lead to patient re-identification and compromise privacy. When genomic data becomes personally identifiable, patient privacy becomes compromised, creating a new set of potential issues concerning disease knowledge, insurance payments, and criminal conviction. Thus, the issue at hand is maintaining patient privacy. The proposed studies will illuminate what specific forms of genomic sequencing data are personally identifiable, helping us determine how to best continue reaping the benefits of sequencing technologies, without compromising patient privacy. From the results of our study, database regulators can design a technology add-on that filters databases to identify datasets that contain personally identifiable formats of genetic information, then implements computational genomic encryption to these datasets.

**Studies to Investigate**
1. A study might involve using hospital records to make researched predictions about a patient's genome, then screening data from publicly available genomic databases (raw DNA sequences, protein expression arrays, cell type levels, etc.) to find samples with these particular genomic variations.
2. An alternative study might involve using publicly available genomic data samples to predict what conditions the donating patient may have been affected by, as well as age of onset, duration of illness, etc., then cross-comparing with hospital records to identify patients.

## Introduction:
Genetic sequencing has grown rapidly in the past few years, developing from a new, unknown technology to one anticipated to be worth $14.8 billion by 2030 [3]. However, recent studies have shown that over 8 out of 10 individuals can be identified purely from their genome sequencing data [4], which is particularly alarming in a healthcare system where pseudonymization is often the only change applied to the data before it becomes sequenced.

The problem is that while sequencing data provides countless benefits to the fields of medical and biotechnological research [1], in order for these benefits to be reaped, the datasets derived from these patients are kept publicly available on sites such as the National Center for Biotechnology Information's (NCBI) Sequence Read Archive [5] (DNA and RNA sequencing data), Gene Expression Omnibus [6] (protein expression sequencing data), Methylation Database [7] (methylation pattern sequencing data), and many other hubs. Interestingly, the Health Insurance Portability and Accountability Act (HIPAA) does not provide protection for genomic data, and the Genetic Information Nondiscrimination Act of 2008 (GINA) is an anti-discrimination law, not a privacy law [8]. This poses a problem to patients and donors who are offering their genome to be sequenced for research purposes. Since sequencing data can reveal highly personal genetic and medical information, having this data tied back to a person creates huge privacy concerns, such as increasing their risk of having personal information shared without their consent, unknowingly paying increased treatment and insurance prices, and elevating their risk of being identified as the suspect of a crime.

With certain forms of genetic data such as DNA sequences being publicly available, cross-comparison with genomic predictions inferred hospital records allows for re-identification of patients from databases with these particular data forms, and the public could now have unfettered access to genetic data that can reveal highly personal information about a patient. For example, a patient with Huntington's– an emotionally devastating condition– may want to keep their diagnosis private and off their hospital records until their treatment has fully begun. They may still donate a sample of bodily fluid to a sequencing core, because their physician reassured them that their data would be sufficiently anonymized through pseudonyms or codes for each patient. However, unbeknownst to both patient and physician, depending on what form of data this patient's sample is reduced to, their identity may not be sufficiently protected by these forms of sample privacy protection. If their data is published in a DNA-sequence format, it may unfortunately be possible for someone to cross-compare this patient's unique Huntington's-inducing SNP combination [9] to hospital records, then leading to re-identification of the patient. This patient's privacy is intensely compromised, and the public now has knowledge of medical information that was meant to be kept private. Thus, patient re-identification from certain forms of genetic data can pose pertinent privacy issues that must be addressed.

Concerningly, genetic databases being available to the public includes insurance companies. Insurance companies with knowledge that a client has a genetic predisposition to disease, particularly one that is life-threatening, are likely to charge higher premiums [10]. Should a particular form of data reveal not only who a patient is, but also that they contain genetic mutations in certain biomarkers of cancer (KRAS gene for lung cancer [11], BRCA1 gene for breast cancer, and many more), insurance companies could gain knowledge of a life-threatening condition before the patient themself even knows, common for neurodegeneration and addiction-related conditions, where identification of a genetic predisposition can indicate a

significantly increased chance of developing said condition. For example, if a person were in a dataset as a flu patient, but their insurance company used their genome to identify them as predisposed to developing breast cancer as well. Insurance companies can then use this information as leverage for patients to buy long-term drug plans [10].

Public access to databases with personally identifiable forms of genetic data could lead to increased false criminalization of patients who contribute samples to these databases, through a practice called "DNA fingerprinting" [13], a criminal conviction tool that has been in use since 1987 [14], which typically involves comparing DNA found at a crime scene to a certain library of genetic data that has been collected from suspects or other criminals in the area. However, with the recent exponential increase in collected genetic data, governmental authorities are now able to turn to databases like NCBI Sequence Read Archive or Gene Expression Omnibus— intended for research purposes— to search for matches with criminal genetic data [15]. While the practice of DNA fingerprinting has been used to successfully identify criminals in the past, the relatively new technology remains uncertain [16] and can be error-prone at times [14]. With these error margins in place for criminal conviction, this increases the risk of patients who choose to donate samples for genetic sequencing of being falsely convicted for crime.

Hence, the aforementioned study is crucial because patients who provide their genomic data to be sequenced for research should be able to maintain their genetic privacy. With the potential for re-identification from certain forms of genomic data, patients are not only at risk of having their genetic privacy compromised and personal information shared without their consent, but they are also at risk of increased insurance prices and potential false criminal conviction, as described above. While making various forms of genetic data accessible via public databases poses many problems for maintaining patient privacy, this is not an unresolvable technology-society conflict. Conducting the study proposed in this paper is intended to illuminate what specific types of genomic sequencing data allow for patient re-identification, and which types can effectively maintain patient privacy, shedding light on a conflict solution.

While studies focusing on the re-identifiability of DNA data have been previously conducted [17] [18] [19], the recent influx of other forms of sequencing data— RNA, protein expression level, methylation patterns, and more— calls for the inclusion of another study such as this one that explores the translatability of one form to another, therefore shedding more light on what forms of genomic data allow for patient re-identification and what forms do not.

## Background:
### Publicly-available genomic data
Genomic sequencing is defined as the process of "deciphering the genetic material found in an organism or virus" [20]. Sequencing technologies like NextGen, bulk RNA, and single-cell RNA sequencing— which each break a sample down into its full DNA, RNA, and cell types,

respectively– shed light on health conditions a patient may be experiencing, or help preemptively detect disabilities and future conditions [21]. The research performed from the results of these sequencing studies has enabled physicians to save lives through preventative treatment [22]. Sequencing data can also guide treatment and determine medical safety on a case-by-case basis [23], an increasingly crucial ability in the era of precision medicine. The growing presence of data that contains so many medical benefits has unsurprisingly led to the emergence of sequencing data repositories [24], such as the National Center for Biotechnology Information's (NCBI) GenBank, or their Sequence Read Archive (SRA), databases that help process, store, and sort all sorts of genomic data. By aggregating data from large populations of patients with a certain condition, the plethora of information neatly organized into these databases has enabled medical researchers to make groundbreaking genetic discoveries [25].

Predictably, sequencing data repositories require genomic data in order to provide these benefits to the field of medicine. Since databases tend to be focused on providing genomic data regarding specific conditions or diseases, the data they contain is often derived from patients with these conditions or diseases who consent to donating their biological samples [26]. Their understanding is that the genomic data derived from their samples is sufficiently anonymized before being made publicly accessible through a database.

However, proper anonymization of genetic data must extend beyond a simple pseudonymization or code-assigning, which the current processes do not. This is due to the fact that the human genome contains some specific regions that are unique to each person. For example, studies have shown that less than 100 single-nucleotide polymorphisms (SNPs), or variant mutations occurring at a single point, can identify a set of genetic information as coming from a specific person [4]. In other cases, SNPs can be used to identify which patients are siblings out of a set of samples. Certain regions of the Y chromosome can also be personally identifying. Another form of re-identification stems from phenotypical expectations– predicting what a person may look like– which can be derived from whole genome studies as well [27]. Similarly, patients with certain conditions can be predicted to have a particular genotype (see "Methods" for more details) that can lead to identification.

When genomic data from public databases is cross-compared with publicly available hospital records, it may be possible to re-identify individual patients [2].

Publicly-available health care data
As of 2013, the National Association of Health Data Organizations (NAHDO) reported that only nine out of fifty states did not require provider charge data to be made public [28]. This type of data being made public allows for the creation of "population-based profiles" as described in Sweeney et al.'s paper [17], consisting of {5-digit ZIP, gender, date of birth, hospital visit information} demographic information. Due to the individuality of certain forms of genomic

information, correlating this genomic data with these hospital records can lead to patient re-identification and, crucially, access to a named patient's fully sequenced genome. It is imperative, therefore, to understand which types of genomic sequencing data contain or lead to these types of personally identifiable information, in order to determine which data forms compromise patient privacy.

Matching data

Certain data formats, such as DNA sequences [29], may allow for *direct* re-identification of a patient donor. This study intends to not only confirm these formats, but also elucidate what *other* forms of data make re-identification possible, particularly from a translational perspective. While it is known that DNA sequencing provides directly identifiable data, is it possible to use further computational analysis or genomic tools to translate one form of data into another that is known to be patient-identifiable? For example, as opposed to DNA's identifiability due to directly containing SNPs [29], RNA undergoes processing in the human body that does not allow SNPs to be as easily identified from it. However, certain transcripts of RNA can be converted back into DNA, and SNPs could potentially then be identified from these translated transcripts [30]. A similar idea could allow for the conversion of cell type data into another form of data that is then personally identifiable. The study proposed in this paper could crucially identify what types of sequencing data can become patient-identifying; from there, sequencing databases can implement a filter to identify and encrypt these types of data (inspired by currently existing DNA encryption technologies [31]) prior to public access, while still permitting access to the original formats of data types that are not patient-identifying, effectively maintaining the benefits of genomic sequencing.

## Materials and Methods:

*Materials*
- Different forms of genomic data; generally all available through NCBI, or other databases
    - DNA sequencing data
    - Bulk RNA-sequencing data
    - Single-cell RNA-sequencing data
    - Methylation sequencing data
    - Transcriptome sequencing data
- Hospital data for cross-comparison
    - Standardized into population-based profiles consisting of {5-digit ZIP, gender, date of birth, hospital visit information}
- General facts about how certain diseases present genomically (see below for details)
- Databases for translation of one form to another
    - dbSNP (convert DNA to SNPs and back) [32]
    - RevComp (convert RNA to DNA and back) [33]

*Methods*

Study 1 references the use of clinical knowledge of genetically-based diseases in order to predict what a patient's genomic data may look like. An example of how this would be possible starts from diseases with well-documented genotypic effects, such as Huntington's disease or cystic fibrosis. The genome of Huntington's patients has been shown to have an inverse relationship between repeat expansion sizes of "CAG"– a triplet of DNA bases– and disease age onset, i.e. those with long CAG repeat expansions are likely to exhibit early-onset Huntington's [34]. Conditions that arise from single gene mutations, like cystic fibrosis, can often be narrowed down to a small set of mutations which cause the disease (CF is caused by F508del, E56K, and G178R mutations, to name a few [35]). The particularity of the genomic changes that cause some diseases make for relatively easy predictions of certain parts of a patient's genome based on the information provided in their health records. CleanGene [18], created by Sweeney et al., is an already existing tool using gene-based disease information to recognize diagnoses that can appear in clinical records. While genotypic effects are only known for some conditions at this point, therefore offering a limitation to what diseases this study can be applied to, more and more diseases are being found to have a DNA component [18], and as time goes on, the capabilities of the study will expand as more genomic information is known.

Study 1 also references the use of a patient's population-based profile, including their 5-digit ZIP, gender, date of birth, and hospital visit information for patient re-identification. Although these hospital records do not have names attached to them, this study would use existing data linkage algorithms to identify a specific person based on these pieces of demographic information [17].

Note that in order for this study to yield effective results, experimenters must ensure that there are common people between the hospital records and genomic datasets used in the study. This can be ensured by selecting a dataset that contains genomic data from patients in a specific geographical region (i.e. New Bedford, MA– NCBI's databases allow for this type of search) and narrowing the range of patient records to those from hospitals in just that area. This significantly reduces the ranges of data for this study such that a patient with a certain condition becomes highly likely to appear in both sets of data. This ensures that if one patient is not found to be in both data sources, this is truly the case and not occurring due to inadequate data range selection.

## Studies and Predicted Events
### Add-On
The envisioned result is for database regulators, such as NCBI, to implement a database add-on that utilizes existing encryption technology to prevent genomic data from leading to patient re-identification. This add-on protects patient privacy while still maximizing genomic data

availability, since the data is still useful when encrypted [36] as opposed to completely removed. As a design statement, the goal is for database regulators to:

> **Construct** a filter for genomic sequencing databases
>> **such that** the formats of genomic data that lead to patient re-identification are recognized and encrypted by the database regulators prior to making the data publicly accessible.

Studies that could be done to determine which forms of genomic data may lead to patient re-identification:

Study 1: Working forwards, patient → sequenced data
After obtaining records with hospital data and patient demographic information*, use current clinical knowledge of genetically-based diseases to predict what certain types of genomic data may look like for that patient– i.e. for DNA sequencing data, predict SNPs a patient may have, and for RNA sequencing data, predict elevated protein expression levels a patient may have. Screen databases containing each type of data (NCBI's Sequence Read Archive can provide DNA and RNA data, MethDB for methylation data, and more where each sample has been taken from a patient) for matches with these predictions, and work to locate a patient's fully sequenced genome accordingly.

This study requires some knowledge of how genetically-based diseases manifest in the human body. This can be shared on a case-by-case basis, and does not require extensive medical knowledge. This study also requires some basic computational awareness of simple tools such as RevComp, which will help the experimenter translate DNA into RNA.

Study 2: Working backwards, sequenced data → patient
Access multiple genomic databases in order to obtain different types of sequencing data. Use this data in tandem with current clinical knowledge to predict what condition(s) a patient may be affected by. Cross-compare these results with public hospital records in order to determine what hospital visits correlate to the conditions determined from the sequencing data, and use the hospital record's information to identify the patient.

This study requires the same knowledge as above.

*It is important to note that these studies operate based on the assumption that each patient has a unique population-based profile combination (see Background for more information)– while this may only be true for around 87% of the U.S. population [18], significant likelihoods can also be used to further identify patients.

## Predicted Results

If Study 1 from this paper were conducted and showed that cross-comparison of certain genomic data types and population-based profiles from hospital records enabled patient re-identification, this would raise concerns about maintaining patient privacy during the rise of genomic data. The decision-makers most likely to respond to such a study are patients donating samples alongside their advocacy groups, medical researchers, database regulators (such as NCBI), and journalists.

Journalists who aim for stories with public interest would likely publicize the results of this study widely, since many people visit hospitals and would be affected by this study's outcome. As a result, patients will then be motivated to come into the public spotlight regarding the privacy of their data. A likely response by patients to journalists would be to raise concerns regarding their personal privacy, as well as react to the study with emotion, since patients are likely unaware of their privacy being breached through genomic databases and hospital records at all.

Upon patient awareness, patient advocacy groups such as the National Patient Advocacy Foundation are likely to get involved and confront the medical researchers who use these forms of genomic data. The medical researchers, wanting to remain on good terms with patients so they will continue to donate biological samples for research, would likely work to prove that sufficient medical research can still be conducted with identifiable forms of genomic data being encrypted [36].

As a result, medical researchers would likely then bring the results of the study to the attention of database regulators. Since the regulators want to keep their databases in use, and medical researchers are the primary users and supporters of these databases, researchers can likely inspire database regulators to construct a filter that identifies and encrypts re-identifiable forms of data from their publicly accessible databases. Thus, a cascade between media attention, patient advocacy groups, medical researchers, and database regulators would likely lead to construction of this add-on technology that filters and encrypts personally identifiable genomic data forms in databases, and allows for the benefits of genomic sequencing to still be attained.

If journalists do not whistleblow the issue of patient privacy in the first place, patients may not even be aware of genomic data as a threat to their privacy; however, if patient advocacy groups are to learn through an avenue other than media attention (i.e. reading the papers themselves) that patients' genomic privacy is threatened, the opportunity for change exists even prior to the study proposal. In this case, change still occurs. Another possibility is that journalists do not provide media attention to the issue, but medical researchers do; in order to maintain their positive relationship with patients so they continue to provide samples to research, the researchers conduct the study and bring the results to the attention of database regulators themselves to enact change. In this case, change still occurs. Alternatively, if patient advocacy groups choose to not address this issue and medical researchers do not perceive it to be a threat

(believing that most people accessing the database are ethical researchers, and not those with malicious intent), change will not occur.

## Discussion

In summary, the proposed study of cross-comparing hospital records with genomic databases could further advance our knowledge of genetic privacy by elucidating what forms of data allow for patient re-identification. While past papers have been written exploring the re-identifiability of DNA sequencing data, this particular study would offer insight into the re-identifiability of other types of data as well. The proposed study would assist with ensuring that privacy laws are updated to reflect the most recent changes in biotechnology, as well as shape our ideas for potential solutions.

As stated earlier, this study is currently limited by the diseases it can encompass, as well as the breadth of data that can be used. The study can only address diseases for which we currently know the genotypic effects of; however, this still encompasses quite a few conditions such as Huntington's, cystic fibrosis, and sickle cell anemia. Due to the study's use of hospital records, experimenters are also limited to conditions that require hospital treatment. For the purposes of this study, the hospital records and data taken from databases must be known to have a common patient. As mentioned in "Methods", a way to address this is by selecting a dataset that contains genomic data from a specific geographical region (i.e. New Bedford, MA) and examine patient records from hospitals in just that area. By narrowing the large ranges of data this study could work with, a patient with a certain condition is highly likely to appear in both sets of data. Of course, a scientific study (Study 1, Study 2, Study 3) may reveal the opposite: that patients cannot be re-identified from a combination of their genomic data and their hospital records. In this case, current availability of genomic data would be shown to pose no threat to patient privacy, and patients, researchers, and database regulators alike can all meet with journalists to freely commend genomic data's capabilities for medical advances. However, this would still greatly increase the credibility of genomic data in terms of privacy maintenance, since genetic privacy is a prevalent topic of concern in today's media. This offers a valuable avenue even if the study does not yield the results as predicted by this paper.

# References

[1] Brittain H. The rise of the genome and personalised medicine. Clinical medicine (London, England). December, 2017. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6297695/

[2] von Thenen N. Re-identification of individuals in genomic data-sharing beacons via allele inference. International Society for Computational Biology. February 01, 2019. https://academic.oup.com/bioinformatics/article/35/3/365/5056754

[3] James S. Next-generation Sequencing Market Anticipated to be worth $14.8 Billion by 2030. Grandview Research, Inc. November 24, 2022. https://finance.yahoo.com/news/next-generation-sequencing-market-anticipated-113000377.html

[4] Shabani M. Re-identifiability of genomic data and the GDPR. EMBO reports. June 1, 2019. https://www.embopress.org/doi/full/10.15252/embr.201948316

[5] Sequence Read Archive. National Center for Biotechnology Information. Accessed December 14, 2022. https://www.ncbi.nlm.nih.gov/sra

[6] Gene Expression Omnibus. National Center for Biotechnology Information. Accessed December 14, 2022. https://www.ncbi.nlm.nih.gov/geo/

[7] Grunau C. MethDB. Accessed December 14, 2022. http://www.methdb.de/

[8] Genetic Information Privacy. Electronic Frontier Foundation. 2015. https://www.eff.org/issues/genetic-information-privacy

[9] Gusella J. Genetic Modifiers of Huntington's Disease. dbGaP. July 30, 2015. ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000371.v2.p1

[10] Andrews M. Genetic Tests Can Hurt Your Chances Of Getting Some Types Of Insurance. NPR. August 7, 2018. https://www.npr.org/sections/health-shots/2018/08/07/636026264/genetic-tests-can-hurt-your-chances-of-getting-some-types-of-insurance

[11] Lung Cancer Biomarker Testing. American Lung Association. November 12, 2022. https://www.lung.org/lung-health-diseases/lung-disease-lookup/lung-cancer/symptoms-diagnosis/biomarker-testing

[12] Blumenthal D. It's the Monopolies, Stupid! The Commonwealth Fund. May 24, 2018. https://www.commonwealthfund.org/blog/2018/its-monopolies-stupid

[13] Sharman S. Forensics and DNA: How genetics can help solve crimes. HudsonAlpha. November 11, 2021. https://www.hudsonalpha.org/forensics-and-dna-how-genetics-can-help-solve-crimes/

[14] Koehler J. DNA Matches and Statistics: Important Questions, Surprising Answers. U.S. Department of Justice, Office of Justice Programs. February, 1993. https://www.ojp.gov/ncjrs/virtual-library/abstracts/dna-matches-and-statistics-important-questions-surprising-answers

[15] Katsanis S. Pedigrees and Perpetrators: Uses of DNA and Genealogy in Forensic Investigations. AnnualReviews. August, 2020. https://www.annualreviews.org/doi/10.1146/annurev-genom-111819-084213

[16] Sahar A. Issues with DNA Fingerprinting in Forensic Lab: A Review. eScientific. April 03, 2019. https://escientificpublishers.com/issues-with-dna-fingerprinting-in-forensic-lab-a-review-JMRCR-01-0002

[17] Sweeney L. Determining the Identifiabiity of DNA Database Entries. Proceedings, AMIA Symposium. 2000. https://dataprivacylab.org/dataprivacy/projects/genetic/dna1.pdf

[18] Sweeney L. Re-Identification of DNA through an Automated Linkage Process. Proceedings, AMIA Symposium. 2001. https://dataprivacylab.org/dataprivacy/projects/genetic/dna2.pdf

[19] Sweeney L. Inferring Genotype from Clinical Phenotype Through a Knowledge Based Algorithm. Pacific Symposium on Biocomputing. 2002. https://dataprivacylab.org/dataprivacy/projects/genetic/dna3.pdf

[20] Genomic Surveillance for SARS-CoV-2. Centers for Disease Control and Prevention. December 02, 2022. https://www.cdc.gov/coronavirus/2019-ncov/variants/genomic-surveillance.html

[21] Frumberg D. Whole Genome Sequencing. Yale Medicine. Accessed November 27, 2022. https://www.yalemedicine.org/conditions/whole-genome-sequencing

[22] Pocius D. Stanford Medicine Scientists Sequence Patient's Whole Genome in Just Five Hours Using Nanopore Genome Sequencing, AI, and Cloud Computing. DarkDaily. November 14, 2022. https://www.darkdaily.com/2022/11/14/stanford-medicine-scientists-sequence-patients-whole-genome-in-just-five-hours-using-nanopore-genome-sequencing-ai-and-cloud-computing/

[23] Advantages of Whole Genome Sequencing WGS. Sequencing, Outsmart Your Genes. Accessed December 14, 2022. https://sequencing.com/blog/post/advantages-whole-genome-sequencing-wgs

[24] Lathe W. Genomic Data Resources: Challenges and Promises. Scitable. 2008. https://www.nature.com/scitable/topicpage/genomic-data-resources-challenges-and-promises-743721/

[25] Lowrance W. The promise of human genetic databases. BMJ. 2001. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1120170/

[26] NCI Staff. As Use of Genomic Data Expands in Cancer Care, Patients Share Their Stories. National Cancer Institute. December 03, 2019. https://www.cancer.gov/news-events/cancer-currents-blog/2019/personal-genomic-data-workshop

[27] Lippert C. Identification of individuals by trait prediction using whole-genome sequencing data. PNAS. September 5, 2017. https://www.pnas.org/doi/10.1073/pnas.1711125114

[28] Data System Tech Resources. National Association of Health Data Organizations. Accessed December 14, 2022. https://www.nahdo.org/data_resources

[29] Zaaijer S. Rapid re-identification of human samples using portable DNA sequencing. eLife. November 28, 2017. https://elifesciences.org/articles/27798

[30] Converting RNA into cDNA. Discovering the Genome. Accessed December 14, 2022. https://discoveringthegenome.org/discovering-genome/rna-sequencing/convert-rna-dna

[31] Satir E. A symmetric DNA encryption process with a biotechnical hardware. Science Direct. April, 2022. https://www.sciencedirect.com/science/article/pii/S1018364722000192

[32] dbSNP. National Center for Biotechnology Information. Accessed December 14, 2022. https://www.ncbi.nlm.nih.gov/snp/

[33] Reverse Complement. Bioinformatics. Accessed December 14, 2022. https://www.bioinformatics.org/sms/rev_comp.html

[34] Brinkman R. The likelihood of being affected with Huntington disease by a particular age, for a specific CAG size. Elsevier. May, 1997. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1712445/

[35] FDA Approves Ivacaftor for Five Splice Mutations. Cystic Fibrosis Foundation. August -1, 2017. https://www.cff.org/node/1306

[36] de Groot J. What Is Data Encryption? Definition, Best Practices & More. Digital Guardian. November 07, 2022. https://digitalguardian.com/blog/what-data-encryption

[37] Gichoya J. AI recognition of patient race in medical imaging: a modelling study. The Lancet. May 11, 2022. https://www.thelancet.com/journals/landig/article/PIIS2589-7500(22)00063-2/fulltext

## Title:

# Investigating Discrimination and Bias on Job Search Platforms

## Authors:

Gargan C.

## Technology Study Plan Summary:

*Society versus Job Search Platforms. Issue is algorithmic fairness.*

Technology and predictive algorithms are increasingly relied upon in the hiring process even before a person applies for a job. Digital platforms identify potential applicants for employers and tailor search results and job recommendations to applicants' (real or perceived) interests, skills, and chances of success. Society expects this process to be fair and unbiased. But are predictive tools recommending jobs fairly, regardless of applicants' gender, race, and ethnicity? Or are they reinforcing stereotypical roles and existing inequalities in job access?

### Technology Study Plan Summary_Studies:

1. A study might involve creating many accounts on a job search platform with names associated with specific racial, ethnic, or gender groups. It would then look for statistically significant differences in recommendations.

2. A variation of this study would add resumés to the accounts in order to better understand how the algorithm works and further assess any biases.

3. (Related) Finally, a study might tweak information within the resumés themselves to see if these algorithms are searching for other clues or "proxy variables" that introduce biases.

## Introduction:

Online job sourcing sites such as LinkedIn, Indeed, Monster.com, and ZipRecruiter use algorithms that determine which candidates see which job postings. The issue is that programmers "train" these algorithms using huge historical databases often populated with the data of current employees. Training on historical data makes it likely that these platforms will perpetuate historical trends in hiring and employment, including the underrepresentation of women and racial and ethnic minorities in certain high-paying occupations [3].

Title VII of the Civil Rights Act of 1964 prohibits employment discrimination based on an individual's "race, color, religion, sex, or national origin" [17]. This extends to cases of "disparate impact," which involve "a facially neutral practice that disproportionately harms members of a protected class [18.]" If a disparate impact can be proven, the practice giving rise to it is unlawful, even if an employer was not motivated by an intention to discriminate, unless "the challenged practice is job related for the position in question and consistent with business necessity."[17]

Even if job matching algorithms are programmed to not take into account protected characteristics, they may end up identifying proxy variables that are correlated with protected characteristics — often without users noticing [20]. For example, a job matching algorithm might end up primarily targeting applicants in particular zip codes, thus skewing the racial make-up of the selected applicants (as in many locations race is correlated with zip code) [19]. Alternatively, a hiring algorithm may identify a very specific trait, perhaps college participation in sports such as American football and lacrosse, as a predictor of job success, thereby indirectly advantaging men [19]. As such, some proxy variables for protected characteristics can be very difficult to spot [19].

Because algorithms tend to find — often surprising — patterns, it can be even more difficult to prove that there is a disparate impact [19]. Add to this the fact that many job seekers do not even realize they have been judged by a predictive technology, and it is not surprising that software-assisted sourcing is a highly unregulated space [1]. Are job seekers exposing themselves to biased effects of hiring predictive tools?

(**Note:** The remainder of this Tech Study Plan refers to Indeed for convenience but can be done with any other job search website.)

## Background:

Indeed is a work-related search engine that allows users to browse jobs of any kind and location. The search engine shows results from employers who post jobs on the platform. As newer jobs are added, older jobs move lower in the search results, holding everything else constant [6]. In order to be listed ahead of the competition, employers can purchase sponsored job postings that appear either on the side or at the top of listings. Employers can also buy a Pay Per Click (PPC) budget, which they can use to bid on certain keywords and search queries. When a job seeker clicks on a sponsored post, part of that budget will be used. This is how Indeed makes money [7]. Sponsored jobs are prioritized over organic jobs in Indeed's algorithm and are not affected by how recently a job was posted [6]. Organic jobs are ranked by their date or relevance to the job seeker in question according to an algorithm Indeed simply describes as "proprietary [30]."

From 2014 to 2017, Amazon tried to build a tool that used machine learning to review job applicants' resumés with the aim of predicting top talent. But by 2015, the company realized that the system was not rating candidates in a gender-neutral way. Because the algorithms

were trained to identify applicants by observing patterns in resumés submitted to the company over a 10-year period [8], they learned to favor men. Additionally, an audit of a different resumé screening company's algorithm found that the algorithm revealed two factors to be most indicative of job performance: whether an applicant's name was Jared, and whether they played high school lacrosse [9]. Clearly, these systems were perpetuating historical patterns rather than predicting new applicants' job performance.

These studies brought public awareness to the potential harm of using artificial intelligence systems in the screening/assessment phase of hiring. However, they did not address the issue of potential discrimination in the sourcing phase, in which popular platforms such as Indeed use similar AI algorithms to suggest jobs to job seekers [10].

Moreover, these algorithms are also used on social media platforms. Facebook, for example, uses an automated ad auction system to determine which users see which ads and how much advertisers have to pay [11]. Advertisers can specify in which auctions they would like to participate by defining a target audience. The winning ad is determined by a combination of how much the advertiser bids and how relevant Facebook's algorithm deems the ad to be for a particular user [22].

In 2017, an investigation by ProPublica and The New York Times revealed that Facebook excluded older workers from job ads.  Another study found that ads placed on Facebook for jobs with taxi companies were seen by an audience that was 75% Black [A]. These studies and additional investigations led to lawsuits by civil rights organizations, and in 2019 Facebook announced that it would no longer allow ads on its platform related to jobs, housing, or credit to be targeted by age, gender, or zip code.

But in 2021, a study found that Facebook algorithms continued to show gender bias in job ads. Facebook's algorithms, for example, were more likely to show a woman an ad for a technical job at Netflix than an ad for a job at a graphics chip maker, which has a higher proportion of male employees [11]. A 2022 study examined how the Facebook ad delivery algorithm deals with ads that include pictures of people of varying genders, ages, and races (that were otherwise identical). The study found that the ads are often delivered to users similar to those pictured (ads that contain images of Black people are more often shown to Black users) [24].


## Materials and Methods:

This investigation looks at Indeed as a case study. A new Indeed account can be created by using an email address and a password or by signing in using a Google, Apple, or Facebook account and allowing Indeed to access the name, email address, and profile picture associated with the existing account (Figure 1).
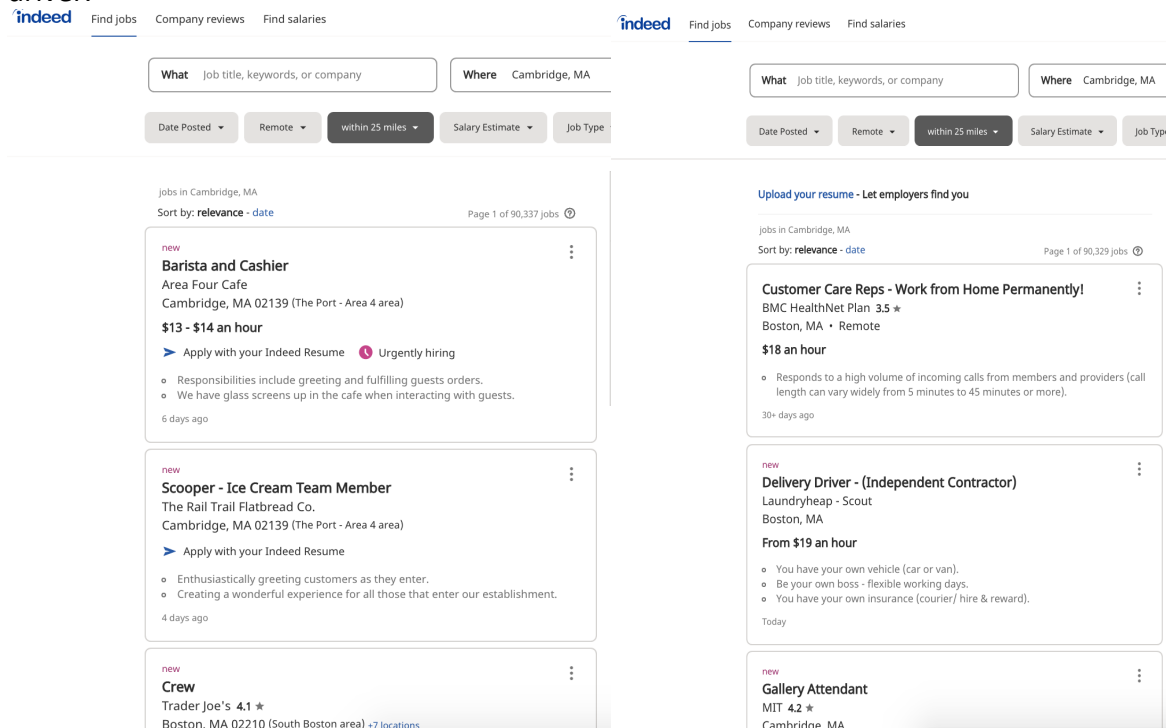
Figure 1. Indeed sign-in prompt

Once logged in, users can either begin a search immediately or create a resumé first.

The exact information that Indeed uses to generate the list of suggested jobs is unclear. Figure 2, for example, shows different jobs displayed to the same machine and IP address in the same city and state, without a posted resumé. The only difference is the user's gender. The woman sees "barista" and "ice cream scooper," while the man sees "customer care rep" and "delivery driver."

Figure 2. Different jobs shown to the same machine and IP address, same city and state, without a posted resumé, signed in as a female user using Facebook (a) and signed in as a male user using Facebook (b).

Of note, the searches were conducted within minutes of each other. It is not clear from the example in Figure 2 what accounts for the different jobs listed. Did the algorithm make a prediction based on the gender of the Facebook user? What jobs might appear if the male and female user posts a resumé?

A resumé can be added to an account either by uploading an existing resumé in the user's format of choice or by answering a number of predefined questions to create an Indeed Resume.

The proposed studies will require one newly generated email address for each account in the sample, as well as a large number of resumés for the second and third study. One can either generate these resumés by creating "fake" resumés with a range of characteristics or scrape these resumés by collecting real resumés from real applicants and de-identifying them. The former would have the advantage of being tuned easily, while the latter would have the advantage of providing real-world data.

Finally, the proposed studies require data on which genders, races, and ethnicities are currently over/underrepresented in certain jobs. The US Bureau of Labor Statistics publishes detailed information on "employed persons by detailed occupation, sex, race, and Hispanic or Latino ethnicity [25]." The resumés uploaded to job platforms also need to indicate the gender, race, and ethnicity of the job seeker without explicitly stating it. This can be done by choosing first names [26] and last names [27] that have been identified by previous research as being overrepresented within certain demographic groups.

## Studies:

### Studies_DesiredOutcome:

The desired outcome is for job platforms, such as Indeed, to match job seekers with the same skills and qualifications to the same jobs, regardless of their gender, race, age and ethnicity. As a design statement, the goal is for Indeed to:

### Studies_ConstructClause:

**Construct** a technology for online job matching
**such that** the job recommendations are fair and unbiased.

## Studies_study1:

**Study 1. Names Only**

This proposed study would involve creating many accounts on a job platform that are identical except for the names associated with them. Specifically, the names should reflect the following: (i) White (non-Hispanic) men, (ii) White (non-Hispanic) women, (iii) Black (non-Hispanic) men, (iv) Black (non-Hispanic) women, (v) Hispanic men, (vi) Hispanic women, (vii) Asian men, and (viii) Asian women. The minimum amount of information required to create an account should be provided. In the case of Indeed, an account can be created by providing only an email address.

Two job searches would then be performed: An active job search and a passive job search. An active search might include jobs that are typically assigned to people of certain genders, races, or ethnicities. A passive job search would look at services such as LinkedIn's "Jobs You May be Interested In," which provide personalized job recommendations to all users whether or not you have made any active job searches.

Do the different groups of candidates receive similar job recommendations? If not, where does the bias fall?

## Studies_study2:

**Study 2. Identical Resumés**

A variation of this study would add resumés to the accounts. The resumés would be identical in all respects other than the users' names. Any changes from the previous study will illuminate how the algorithm works and whether there are additional biases introduced or not.

## Studies_study3:

**Study 3. Tweaked Resumés**

(Related) To check whether these algorithms are searching for proxy variables that are correlated with various subgroups, a final study could leave names the same but tweak resumés. Some resumés, for example, might include an all-female or historically black college or university. Alternatively, they might list an all-male sport or a particular zip code. This will further illuminate how these algorithms work and whether they include any biases.

## Predicted Events:

Suppose a study was conducted that showed that there is a significant difference in job recommendations depending on identity. Such a study would raise the question of gender, race, and/or ethnicity bias of the algorithm used by Indeed.

The decision-makers most likely to respond to such a study are journalists, women's and civil rights advocacy groups, the U.S. Equal Employment Opportunity Commission (EEOC), and Indeed. Advocacy groups could be Women's Rights advocates if results show gender bias towards women.

Journalists primarily want to frame stories in such a way that garners public interest and informs; therefore, the study would likely gain some media attention.

This might motivate action from consumer groups involved. If attention increases, the U.S. Equal Employment Opportunity Commission (EEOC) may consider issues related to employment websites and job matching as part of its ongoing agency-wide initiative on "Artificial Intelligence and Algorithmic Fairness Initiative," which is currently focusing on "employers, employees, job applicants, and vendors" [29].

Eventually, the media attention, advocacy group concern, and EEOC investigation may lead to Indeed testing their algorithms to see which job seekers end up seeing jobs, and then correcting for undesired biases. The company may try to deflect criticism by solely releasing a public statement saying they are taking meaningful steps to address issues of discrimination in job listings and have teams currently working on job fairness [13].

Of course, the proposed studies may reveal the opposite: there is no significant difference in ads targeted to job seekers based on their gender, race, or ethnicity. In such a case, none of the predicted events would occur. However, it would present a rigorous methodology for evaluating bias in similar applications. In addition, it would suggest that Indeed's algorithms do not perpetuate bias in job recommendations and just may be harnessed in other algorithms, like those that companies later use to weed out job applications.

## Discussion:

In summary, the proposed studies aim to analyze whether employment platforms, such as Indeed, match otherwise identical job seekers of different genders, races, and ethnicities with different jobs. If this is the case, the platforms could be indirectly contributing to women, and racial and ethnic minorities' continued underrepresentation in certain high-income and high-status professions.

Even if no such differences in recommendations are found (or if the differences found are not statistically significant), the proposed studies would still make a contribution to the ongoing

discussion about AI in hiring and allow attention to be focused on more problematic parts of the "hiring funnel."

The proposed studies also have some limitations. On a conceptual level, it is difficult to ascertain why any potential differences in jobs shown to job seekers of different genders, races, and ethnicities arise and to what extent it is the responsibility of job platforms to rectify them. If employment platforms are simply showing each job seeker the ads with which they are most likely to engage, this would likely perpetuate existing employment patterns, but the actions the platforms ought to take depend on whether the differences arise due to: (i) inherent differences in preferences between job seekers (no need to address), (ii) socially constructed and enforced differences in preferences between job seekers (very difficult to address for the platforms), (iii) lack of visibility of alternative options (easy to address for the platforms).

On a practical level, the proposed studies will likely require a fairly substantial sample size to produce statistically significant results, as there are eight sub-groups to be tested.

## Citation:

Be sure to cite this writing in all related work.

Gargan, Christine. "Investigating Discrimination on Job Recommender Websites." Tech Study Plans. Plan 5006. September 2021. http://techstudies.net

## References:

1       Rieke, Aaron, and Miranda Bogen. 2018, *Help Wanted: An Examination of Hiring Algorithms, Equity, and Bias*, https://www.upturn.org/work/help-wanted/. Accessed 14 Jan. 2023.

2       Sonnemaker, Tyler. "Here's Why an AI Expert Says Job Recruiting Sites Promote Employment Discrimination." *Business Insider*, Business Insider, 18 Jan. 2020, https://www.businessinsider.com/ai-expert-job-sites-must-prove-not-exacerbating-inequality-2020-1.

3       Hsu, Jeremy. "AI Recruiting Tools Aim to Reduce Bias in the Hiring Process: Artificial Intelligence Software Promises to Make Hiring Fairer. But How Well Does It Work?" *IEEE Spectrum*, 29 July 2020, https://spectrum.ieee.org/ai-tools-bias-hiring#toggle-gdpr. Accessed 14 Jan. 2023.

4       Falk, Gene, et al. Congressional Research Service, 2021, *Unemployment Rates During the COVID-19 Pandemic*, https://sgp.fas.org/crs/misc/R46554.pdf. Accessed 14 Jan. 2023.

5       Polner E. Best Job Search Websites. The Balance Careers. February 18, 2021.
        https://www.thebalancecareers.com/top-best-job-websites-2064080

6       Alacozy Z. What You Need To Know About Search Rankings On Indeed. Recruitics. August 26, 2019.
        https://info.recruitics.com/blog/searching-ranking-on-indeed

7       The Indeed Business Model – How Does Indeed Work & Make Money? Productmint. Last updated March
        20, 2021.
        https://productmint.com/the-indeed-business-model-how-does-indeed-work-make-money/

8       Dastin, Jeffrey. "Amazon Scraps Secret AI Recruiting Tool That Showed Bias against Women." *Reuters*, 10
Oct. 2018, https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-
recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G.

9       Ajunwa, Ifeoma. "The Paradox of Automation as Anti-Bias Intervention." *Cardozo Law Review*, vol. 41, no.
5, June 2020, pp. 1671–1742., http://cardozolawreview.com/wp-content/uploads/2020/10/1.-
Ajunwa.41.5.6.FINAL-3.pdf. Accessed 14 Jan. 2023.

10      Kohler C. The Top 10 Professions Dominated by Women in the United States. TopResume. Accessed May
3,      2021. https://www.topresume.com/career-advice/top-10-professions-dominated-by-women

11      Kaplan, Levi, et al. "Measurement and Analysis of Implied Identity in Ad Delivery Optimization."
*Proceedings of the 22nd ACM Internet Measurement Conference*, pp. 195–209,
https://www.ccs.neu.edu/home/amislove/publications/FacebookIdentity-IMC.pdf. Accessed 14 Jan. 2023.

12      Gershgorn, Dave. "Companies Are on the Hook If Their Hiring Algorithms Are Biased." *Quartz*, 22 Oct.
2018, https://qz.com/1427621/companies-are-on-the-hook-if-their-hiring-algorithms-are-biased. Accessed 14 Jan.
2023.

13      Hao, Karen. "Facebook's Ad Algorithms Are Still Excluding Women from Seeing Jobs." *MIT Technology
Review*, 9 Apr. 2021, https://www.technologyreview.com/2021/04/09/1022217/facebook-ad-algorithm-sex-
discrimination/. Accessed 14 Jan. 2023.

14   "Search 1000s of Resume Samples and Examples!" *LiveCareer*, www.livecareer.com/resume-search/.

15      Improving job matching with machine-learned activity features | LinkedIn Engineering

16      Job Recommendations – Overview | LinkedIn Help

17        https://www.law.cornell.edu/uscode/text/42/2000e-2

18        https://ilr.law.uiowa.edu/print/volume-105-issue-3/proxy-discrimination-in-the-age-of-artificial-
intelligence-and-big-data

19        https://minnjil.org/wp-content/uploads/2022/01/Jonjua_v30_i2_329_358.pdf

20 https://ilr.law.uiowa.edu/print/volume-105-issue-3/proxy-discrimination-in-the-age-of-artificial-intelligence-
and-big-data

21      Doing More to Protect Against Discrimination in Housing, Employment and Credit Advertising | Meta
(fb.com)

22      About ad auctions | Help Center (facebook.com)

23      Auditing for Discrimination in Algorithms Delivering Job Ads (arxiv.org)

24      Measurement and Analysis of Implied Identity in Ad Delivery Optimization (neu.edu)

25      Employed persons by detailed occupation, sex, race, and Hispanic or Latino ethnicity (bls.gov)

26 policy-brief-04b-2016-an-updated-analysis-of-race-and-gender-effects-on-employer-interest-in-job-applicants.pdf (missouri.edu)

27      Decennial Census Surname Files (2010, 2000)

28      How Do I Create an Indeed Account? – Indeed Support

29      Artificial Intelligence and Algorithmic Fairness Initiative | U.S. Equal Employment Opportunity Commission (eeoc.gov)

30      How does Indeed rank search results? – Indeed Support

A       https://arxiv.org/abs/1904.02095

# Racial Bias in Uber's Real-Time ID Check

Christine Mui
December 2021

## Summary

*Uber Drivers versus Uber Technologies*
*Issue is racial bias in Uber's driver verification technology*

Since 2016, Uber U.S. has used a security feature dubbed Real-Time ID Check to ensure that the driver behind the wheel matches the account profile picture and license on file. Drivers are prompted to regularly upload selfies. If a photo fails the ID check, the driver's account will "immediately be temporarily suspended." This facial recognition system has been criticized for producing racially biased results. After the feature launched for U.K. drivers in 2020, several former Uber and Uber Eats employees alleged that they were fired after repeatedly being misidentified by the facial recognition system. Asserting that the reason for the failed attempts was racial bias in the technology, they filed employment tribunal claims on grounds of discrimination. It is critically important to assess whether the facial recognition technology behind Real-Time ID Check malfunctions at higher rates when verifying drivers of color, putting them at disproportionate risk of losing their livelihood.

**Studies:**

1. A study might involve comparing the facial recognition software's accuracy in matching images to verify the identities of individuals of different skin tones and genders. The experiment would seek to recreate the conditions of Real-Time ID Checks as closely as possible.

2. Another study could measure Face API's matching accuracy rates for people of different skin tones and genders after introducing conditions that might make it more difficult for the technology to verify submitted selfies against driver's license photos. These conditions include new glasses, facial hair, and dim lighting.

3. (Related) A final study could evaluate the effectiveness of human reviewers compared to automated matching under the conditions described in studies 1 and 2. A group of human reviewers would decide whether to verify each selfie against a file

photo, playing a similar role to that of the human reviewers Uber hired in 2019 to help detect unauthorized drivers on its app.
.

# Introduction

There have long been concerns over racial bias in many facial image processing systems. But the problem at hand lies with the decisions Uber makes based on its facial verification results. This issue is especially urgent and serious because those results alone might result in an employee's termination.

Uber launched its Real-Time ID Check in the UK in 2020 [2]. In 2021, two U.K. trade unions supporting drivers took legal action against Uber for what they claim is racist facial verification technology. One union, the Independent Workers' Union of Great Britain (IWGB), asserted that Face API was "five times more likely to cause the termination of darker-skinned workers" and demanded Uber create a more "fair, transparent process for account terminations" [10]. The IWGB is suing on behalf of one driver who was locked out so many times that their account was terminated [10]. The union claims that at least 35 other drivers also had their employment terminated as a result of mistakes made by the "racist algorithm" [10].

The second union, the Apps Driver and Couriers Union (ADCU), launched employment tribunal claims against the company on behalf of Pa Edrissa Manjang, a former Uber Eats courier, and Imran Javaid Raja, a former Uber private hire driver [11]. Manjang was dismissed from the Uber Eats service in London for "continued mismatches" between the photos uploaded at the start of his shifts and his profile on file [12]. After Manjang asked to have a human review the images, Uber deactivated his account, and the case is still under review [12]. Raja was suspended for similar mismatches. Although Uber was reappointed, he was never offered any form of compensation for his period off work [33]. In December 2021, the ADCU told a newspaper it had won 10 appeals in court on behalf of dismissed drivers citing discrimination in Uber's ID checks [13].

Research indicates that problems interpreting images of darker-skinned individuals are widespread in facial recognition systems. The pioneering 2018 Gender Shades study by Joy Buolamwini and Timnit Gebru, for example, showed that leading commercially available technology for classifying gender based on facial images performed considerably worse on darker-skinned faces than lighter skinned-faces. In particular, Buolamwini and Gebru found that the technology Uber uses for facial verification, Face API from Microsoft's Cognitive Services, determined the gender of lighter-skinned males with an error rate of only 0.8%, but identified the gender of darker-skinned women with a 23.8% error rate [14].

In arguing that facial recognition systems "generate particularly poor accuracy results when used with people of color," the ACDU also referenced the U.S. National Institute of Standards and Technology (NIST) Face Recognition Vendor Test [11]. Over two decades of work, NIST said

it has seen a significant improvement in the accuracy of facial recognition technology but concluded that most software tended to be more accurate for white, male faces than for people of color or for women [15]. African American and Asian faces were 10 to 100 times more likely to generate false positive identifications than those classified in NIST's database as white [15]. The findings, however, were dependent on the image quality.

In June 2018, a few months after the Gender Shades study was published, Microsoft announced in a press release that it had "updated its facial recognition technology with significant improvements in the system's ability to recognize gender across skin tones" [27]. The release did not mention the study by name, but it did say that the new improvements were meant to address "recent concerns that commercially available facial recognition technologies more accurately recognized gender of people with lighter skin tones than darker skin tones, and that they performed best on males with lighter skin and worst on females with darker skin" [27]. With the improvements in place, Microsoft said it was able to reduce the error rates of the Face API gender classifier for men and women with darker skin by up to 20 times and the overall error rate by nine times.[27]

The U.K. lawsuits and a survey by MIT Technology Review, published in 2022, suggest that improvements in Face API did not prevent darker-skinned drivers from failing to be recognized by Uber's ID check system. Tech Review reported that almost half of the 150 Uber drivers in India who were surveyed had been locked out of their Uber accounts—either permanently or temporarily—as a result of Real-Time ID Check, with many drivers suspecting that "changes in their appearance, such as facial hair, a shaved head, or a haircut, was to blame." Others cited low lighting, scratches on their cameras, and low-budget phones as the reasons for issues with Real-Time ID Check [30].

If true, Uber's use of facial recognition technology has the potential to expose drivers of color to the risk of losing their jobs and thus their livelihoods, even if they comply with all company policies. This amplifies the precariousness that gig economy workers already find themselves in. In addition, Uber's reinstatement process for suspended accounts is reported to be "tedious, time-consuming, frustrating, and mostly unhelpful [10; 30]," making errors in ID checking a significant burden for drivers. This issue is not limited to Uber. Local alternatives, such as Ola, Swiggy, Zomato, and Urban Company in India, also use selfies for verification[30].

The racial demographics of drivers make this a potentially widespread difficulty. In London, roughly 90% of private hire drivers, including Uber drivers, are non-white, according to Transport for London [16]. A 2015 survey of U.S. Uber drivers indicated that 24% were Black and 20% Latino [17]. Uber stated that in its months-long pilot of Real-Time ID Check, "more than 99% of drivers were ultimately verified" [1]. But the use of the term "ultimately" suggests that drivers in the pilot were given multiple opportunities to re-upload photos, which is not the case in current practice. The company did not disclose the demographic makeup of those drivers who failed the verification process, saying that "the majority of mismatches were due to unclear profile photos." This leaves open the question, what does Face API's false match rate

look like for the racial groups who comprise the majority of Uber's driver workforce? Does the software exhibit a racial bias?

# Background

<u>How Facial Recognition Technology Works</u>
Face recognition technology is usually based on deep learning algorithms that go through a multi-stage process to identify a person [15.] First, it locates the human face in the picture provided and then it finds certain landmark points within the face, such as an individual's eyes, nose and mouth, based on what it has learned from the training data it was provided. The geometrical features of these landmarks are then extracted and compared to those extracted from pictures already on file, usually taken under controlled conditions [15].

The 1:1 method makes a comparison between two facial images to determine if they are the same person (verification/authentication), while the 1:N method searches for an image's match within a larger database (identification). The latter is the type of facial recognition that federal agencies and law enforcement use to identify unknown individuals in criminal investigations.

<u>What Real-Time ID Check Looks Like for Uber</u>
*The stated goal of Real-Time ID Check is to protect both drivers and passengers. According to Uber, the technology "prevents fraud and protects drivers' accounts from being compromised." It also "protects riders by building another layer of accountability into the app to ensure the right person is behind the wheel [1]."*Face API works similar to Apple's Face ID on the iPhone in that both systems use 1:1 authentication.

Uber's Real-Time ID Check uses Microsoft's Face API from the Azure Cognitive Services suite, specifically the Face-Detect API and the Face-Verify API [19].  Face-Detect identifies specific attributes and detects the presence of human faces in an image. A user who uploads a photo where a face is not detected is asked to re-take the photo. Once a face is detected in an image, Face-Verify compares the face with the driver's profile picture and provides a "confidence score" estimating the probability of a match between the two [20]. The confidence score is what determines whether an Uber driver's account will be successfully verified or, per Uber's policy, it will be temporarily suspended pending an investigation into the mismatch [19].

<u>What Real-Time ID Check Looks Like For Drivers</u>
Uber has stated that "[o]ne of the key objectives for Real-Time ID Check was to avoid unnecessary friction for driver-partners" so that they can "focus on making the user experience as seamless as possible [19]."

According to the company website, Uber's Real-Time ID Check works by periodically prompting drivers — usually when they first get behind the wheel and are ready to start accepting rides — to take a selfie in the car and upload it into the Uber app [1]. Drivers choose whether their identity is verified through Microsoft's facial recognition software or by human reviewers [6].

With a human reviewer, the process takes longer. Uber estimates it to typically take "a few minutes to complete," and drivers are not allowed to accept any rides while waiting for completion of the ID check [6]. Uber says that "due to the random nature of this selection process, the request could potentially pop up while a driver-partner is driving" If a driver has hastily pulled over to the side of a road, possibly with a customer in the back, it is unlikely that they would want to wait even a few minutes before resuming the ride. The almost instantaneous processing time of the automated ID check presents an incentive for drivers to use Face API rather than a human reviewer.

Uber uses Face API to compare the selfie against the photo the account has on file, typically a driver's license photo [1]. The app grants drivers one opportunity to submit a clear photo. If Face API's Verify feature does not detect a match, the driver will be "temporarily blocked" while Uber investigates [7], and their access to the account could be "suspended indefinitely" [6]. Uber's reluctance to give drivers second chances to re-upload photos may stem from real instances of impersonation. In November 2019, Uber was stripped of its license to operate in London — one of its biggest global markets — after authorities found that 43 drivers used false identities to take more than 14,000 rides [8]. These fraudulent drivers borrowed accounts from authorized drivers they knew to pick up riders, in uninsured vehicles they were not registered to drive [8]. When Uber won an appeal to regain its operating license, it promised to solve the issue of unverified drivers. The use of Real-Time ID Check in the U.K. was its solution [9].

Some of the challenges mentioned above are inherent to face recognition as a technology. Others are specific to how it is applied to certain groups. Regarding the former, features such as facial hair, eyeglasses or certain hairstyles, as well as low lighting and low-resolution pictures, can interfere with the technology's ability to find the reference points it is looking for and can lead it to erroneously conclude that a picture does not match what it has on file. Regarding the latter, the face recognition algorithm may struggle to identify people of genders, races, and ethnicities that were underrepresented in its training data [30]. In Microsoft's response to the Gender Shades study, it indicated that it had diversified its training data set.


## The Setting

There are several decision makers invested in the clash over racial bias in Uber's Real Time ID Check who may take action in response to the proposed study. These stakeholders can be understood as neutral or on the side of the two major parties, Uber drivers and Uber Technologies.

From the perspective of Uber drivers, especially drivers of color, the software's role in the dismissal of their colleagues sparks worries that they could be next. In October 2021, roughly 80 Uber drivers protested outside the company's London headquarters as part of a 24-hour strike "demanding better rates per mile with no fixed rate trips, reduction in Uber's commission to

15%, an end to the use of allegedly 'racist' facial identification software and reinstatement of unfairly deactivated drivers" [21]. These objectives were reflected on the signs protesters waved, some of which read "Scrap the racist algorithm" and "Stop unfair terminations" [22]. The strike was part of a campaign launched by two of the advocacy groups supporting the Uber drivers: the IWGB and Black Lives Matter UK. In a press release, the IWGB outlined their demands, which included a "transparent process for account terminations"[10].

As the IWGB and the ADCU have initiated legal action against Uber on behalf of three drivers, they are the stakeholders who are most likely to respond to the study. Both unions have already cited the Gender Shades study and other research documenting racial bias in facial recognition systems [10; 11].

The unions' objective of reinstating unfairly terminated drivers and couriers also involves introducing a fair terminations process. This would require courier companies, like Uber, to clearly outline what actions would be grounds for dismissal and what information is used when using technology to make automated termination decisions. In the U.K., a cross-party group of 60 Members of Parliament signed onto a motion for more transparency and due process with terminations [22], which, if put into effect, could alleviate the clash between drivers and the Face API software.

Other potential allies for the Uber drivers are advocacy organizations that have long advocated for a ban on face surveillance technology. These groups include the Electronic Privacy Information Center, Public Voice, Fight for the Future, Secure Justice, as well as any of the two dozen civil and human rights organizations that have called facial recognition inherently discriminatory and dangerous and wrote in an open letter to federal lawmakers that they should "ban the private and corporate use of facial recognition technology" [23]. The letter opened with the example of Uber firing its drivers in the U.K. over the Real-Time ID Check system, writing that "these cases clearly show how private use of facial recognition by corporations, institutions and even individuals poses just as much of a threat to marginalized communities as government use" [24].
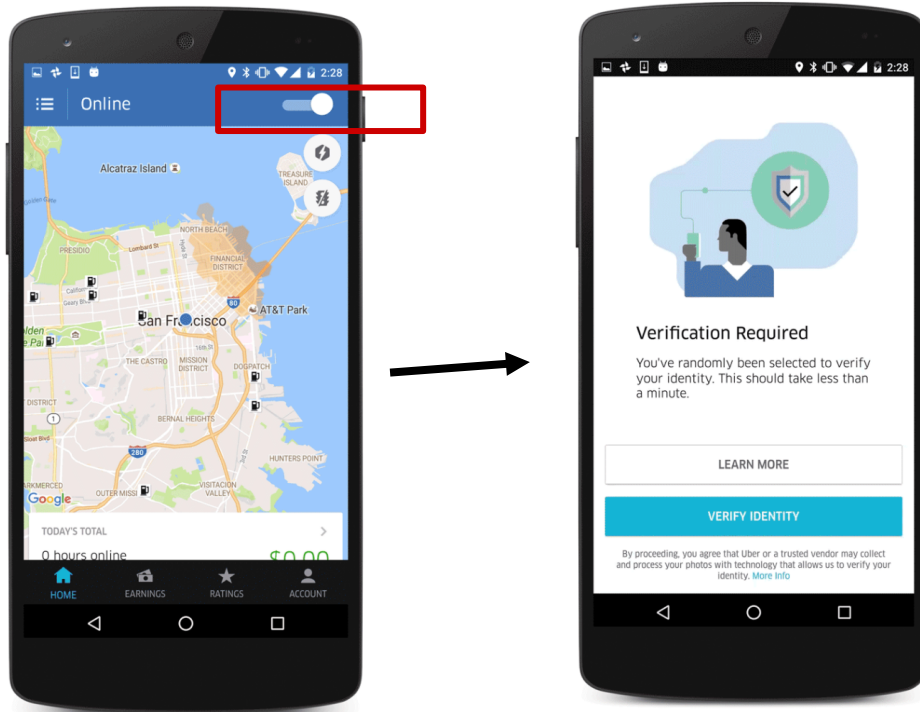
On the side of Uber Technologies, there are not many allies, but one would be Microsoft, which creates the technology behind Face API and sells it as a commercial product. Neither company has admitted to the media any bias within the technology, although Microsoft acknowledges that "[m]easuring the accuracy of facial recognition technology is a very difficult problem and methodologies vary across industries [31.]"An Uber spokesperson told TIME Magazine that its facial verification software is "designed to protect the safety and security of everyone who uses the Uber app by helping ensure the correct driver is behind the wheel." Meanwhile, Microsoft told the publication that it is "committed to testing and improving Face API, paying special attention to fairness and its accuracy across demographic groups. We also provide our customers with detailed guidance for getting the best results and tools that help them to assess fairness in their system" [21]. Microsoft has many prominent clients that use its Face ID APIs and adjacent AI services and stands to lose considerable business if the software is judged to be

faulty. In addition to Uber, the company lists the BBC, Volkswagen, KPMG and Airbus as clients [25].
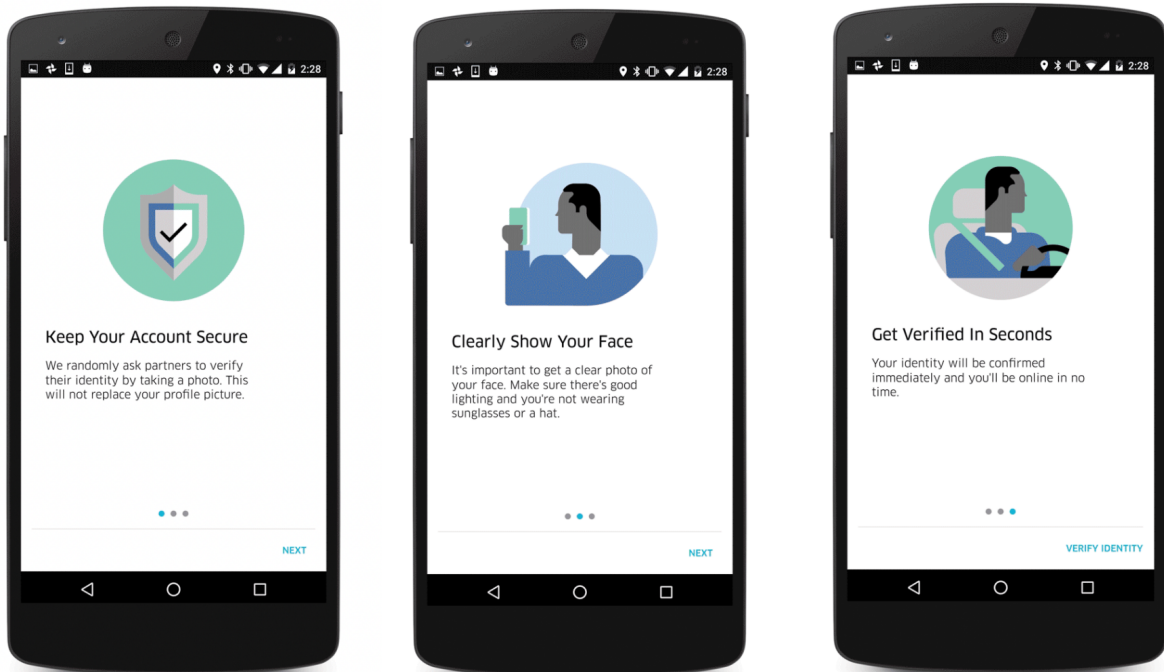
Additional important actors in this clash are government bodies and regulatory agencies, namely those responsible for licensing and overseeing transportation companies, including Uber and Lyft.  Were a study to provide substantial evidence that Uber's Real-Time ID Check discriminates based on skin color and gender, these parties may be encouraged into action by the advocacy groups and legislators supporting Uber. One government body that falls into this category is Transport for London, the city's transportation authority. Transport for London acted against Uber in 2019, choosing not to renew the company's operating license for the second time. In both instances, the authority asserted that Uber had a "pattern of failures" that had put passengers at risk [8]. Given that Uber implemented its Real-Time ID Check in response to the agency's security and public safety concerns, it is hard to predict whether that agency would choose to step in again. The clash over racial bias in Uber's ID check system might play out in a different city or in a U.S. state such as Oregon or California, which have banned law-enforcement use of facial recognition systems.
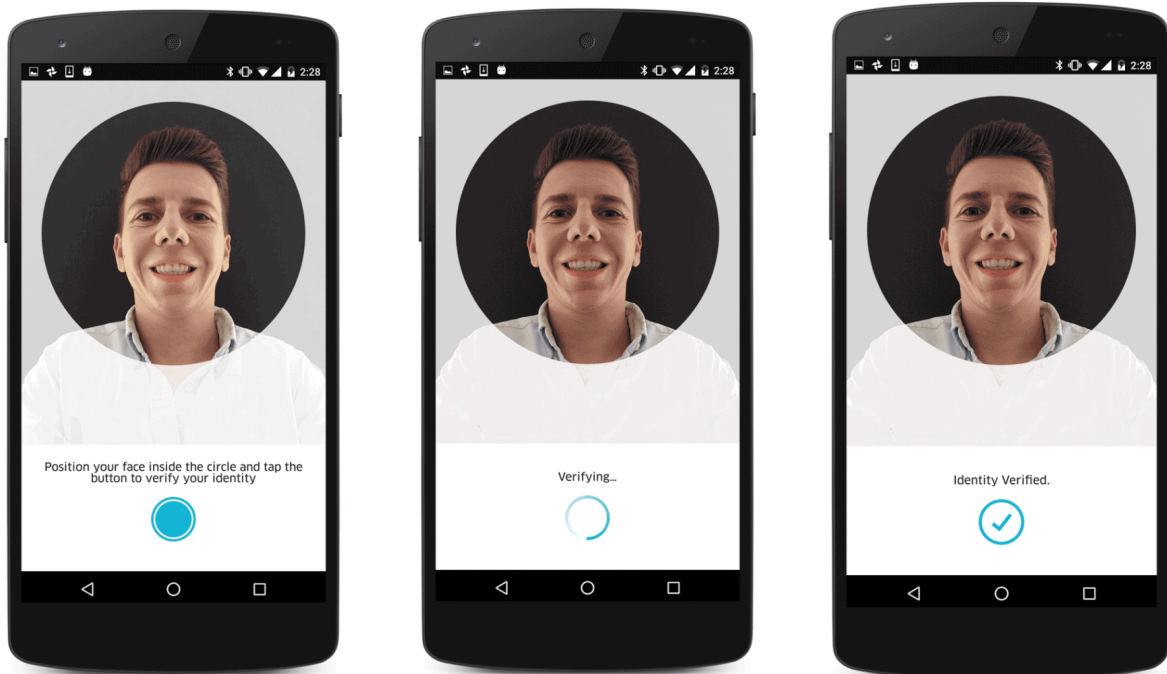
## The Materials

For this study, the researcher would require access to Microsoft's Face API technology. A research group can request access for this type of work [32], and because Uber discloses the features it uses from the software — Face-Detect and Face-Verify — a study can be conducted without having access to Uber's internal app interface or drivers' accounts [19]. However, understanding what that interface looks like for drivers would be necessary to develop the study's methods in a way that most closely resembles the real-life process of completing Uber's ID check.

(a)



(b)

(c)

**Figure 1. The app interface of Uber's Real-Time ID Check for drivers: a) when a driver flips a switch (highlighted by a red square) to indicate they are ready to "go online" and begin accepting riders, the app asks them to verify their identity, estimating it "should take less than a minute," b) the driver is instructed to make sure there's good lighting and to remove any sunglasses or a hat, and c) the driver is told to position their face and tap a button to verify their identity. A successful attempt reads "identity verified" and returns the driver to the beginning screen in (a).**

As shown in Figure 1, which was made from a recording on Uber's website [1], drivers are not informed while using the app that they will not have a second chance to verify their identity or that they risk having their account temporarily blocked and even terminated for a failed attempt (Figure 1c). Instead, the app tells them that "your identity will be confirmed immediately, and you'll be online in no time (Figure 1b). To simulate this experience in the study, participants should not be informed about what the study is testing or the real-life implications of a failed verification attempt for an Uber driver. If they receive this information, the participants may inadvertently try harder to take a high-quality photo, perhaps taking steps that an Uber driver wouldn't, such as leaving the car in search of better lighting.

Additionally, for the study to resemble the process shown in Figure 1, participants will only be given one attempt to take a selfie to be evaluated by the Face API technology. If the selfie is blurry on a first attempt, the participant will not be allowed to re-take it, since that could be something that occurs in a real-life driver situation. In the actual Uber app interface, drivers do not even have the opportunity to review their image before the verification process begins

(Figure 1c). The entire process from taking the photo to having one's identity identified is all on one screen (Figure 1c).

The study's groups of participants can be modeled off those used by the 2018 Gender Shades study, which found that the Face API technology misidentified the gender of black women with up to a 23.8% error rate[14]. The 2018 study used four categories: darker-skinned female, darker-skinned male, lighter-skinned female and lighter-skinned male. A study might borrow the same or similar categorization and recruit equal numbers of participants from each of the four groups.

The study will need to obtain driver's license photos from all participants, as these photos and the selfies the participants will take in the front seat of a car are the two components for the 1:1 authentication method that this study uses. Additionally, a car and a cell phone camera are needed for the study.
Uber does not disclose the threshold it uses for verification. The Face API uses a default "similarity confidence" of 0.5 or greater for 1:1 authentication of faces. In other words, if the algorithm is 50% confident that two faces belong to the same individual, it says that there is a match. Without further details about the tuning of Uber's system, the study will have to use the default threshold.


## Studies and Predicted Events

**Desired Outcome**

The envisioned result is for Uber to implement a safety measure to verify its drivers' identities that does not produce differing outcomes based on skin color or gender. As a design statement, the goal is for Uber to:

> **Construct** a safety measure to verify drivers' identities
> **such that** its accuracy rates are uniformly high for all genders and races

Below are possible studies that could be done to further events that could lead to the desired outcome.

**Study 1. Evaluating Gender and Skin-Color Disparities Within Real-Time ID Check**

A study might involve comparing the facial recognition software's accuracy rates for detecting and verifying people of different skin tones and genders. The experiment would seek to recreate the conditions of Real-Time ID Checks as close as possible. It would solicit an equal number of participants for four groups — men with lighter skin, women with lighter skin, men with darker skin, and women with darker skin. A mobile app with the functions of Uber's Real-Time ID system would need to be created and downloaded by participants, who ideally would

be provided a variety of cellphones. Participants would sit in the front seat of a car and be told to take a selfie. Each of these participants would be given the same instructions as Uber gives its drivers before the app starts the verification process (Figure 1).

Because Uber drivers are not told they will not be allowed to retake the selfie, this information will also not be revealed to the study participants during the process. However, the first photo each participant takes will be uploaded into Face API, purchased from Microsoft Cognitive Services [4]. That software will then check the photo against the participant's driver license photo through the 1:1 authentication method and determine whether the two identities match, based on the default 50% confidence level of Face API's Verify feature. The mismatch rate will be calculated separately for each group; the formula used is the number of participants in that group whose identity was not successfully verified out of the total number of study participants in the group. The results of the survey will consist of whether there are any statistically significant discrepancies in the photo comparison mismatch rate between the four groups.

**Study 2. Evaluating Gender and Skin-Color Disparities Within Real-Time ID Check Under Differing Photo Conditions**

Another study could measure Face API's accuracy rates for people of different skin tones and genders after introducing conditions that might make it more difficult for the technology to verify submitted selfies against driver license photos taken in optimal lighting. This experiment would follow the same photo taking and verification procedure as study one but also test-confounding conditions, such as the use of glasses, facial hair, or dim lighting. It is up to the researcher whether they want to test only one of these variables or all three and to test additional conditions, such as cellphones that take lower-resolution photos. Regardless, for each variable, the study would require that half of each of the groups from study one — men with lighter skin, women with lighter skin, men with darker skin and women with darker skin — experience that confounding condition, while the other half does not. The latter half serves as a control group, while the half that experiences the condition is the test group.

**(Related) Study 3. Human Reviewers vs. Facial Recognition Software for Identity Verification**

A related study might involve re-running Studies 1 and 2 using human reviewers instead of Face API. A group of human reviewers would manually have to decide whether to verify each set of photos, playing a similar role to how Uber's team helps detect unauthorized drivers on its app. The study would see how the accuracy rate of the human reviewers compares to that of the Face API technology run on the same set of photos. Human reviewers can also exhibit biases; however, Uber's decision to offer an option of human review implies that they can be more accurate than an algorithm in making 1:1 face matches. This study therefore provides additional data from which to evaluate possible bias in the software.

**Predicted Events**

Suppose a study was conducted and showed that Microsoft's Face API technology produced significant skin-color and gender discrepancies in the mismatch rate for verifying car selfies against driver ID photos in a study simulating the process used by Uber's Real-Time ID Check. Such a study would provide empirical evidence for the ongoing racial discrimination lawsuits against Uber.

The decision-makers most likely to respond are the two unions, Independent Workers' Union of Great Britain and App Drivers and Couriers Union, that have already initiated legal action against Uber on behalf of three drivers. As explained in "The Setting," these unions referenced several studies in their press releases announcing the lawsuits. One of the referenced studies was the 2018 Gender Shades study by Buolamwini and Gebru, which this study most closely resembles. It makes sense that they would want to promote a study with a similar finding that better replicates the photo conditions in which these Uber drivers undergo their identity checks. I predict that these unions would publicize this study's findings and use them to back up their legal claim that Uber's algorithm has led to unfair terminations.

From there, journalists would be next to respond to the study results. Several publications in the U.S. and the U.K, such as The Guardian [12], WIRED [26] and TIME Magazine [21], have been closely following the story about Uber drivers' claims that racial discrimination in the company's software is putting them out of work. Because of the extensive coverage of this clash already and the media attention that followed the findings of the Gender Shades study,this study would be expected to garner a good amount of media coverage. At the same time, advocacy groups such as Fight for the Future and Secure Justice would most likely get involved, releasing statements of their own and perhaps creating more media buzz as a result.

Depending on the extent of that coverage, Uber and Microsoft may act next, as they would have received several requests for comment from the journalists.  Microsoft and Uber's communications strategy in the past has been to lie low until they are able to release press releases, touting their own accomplishments or improvements. In 2018, for example, Microsoft published a blog post indicating that they had reduced the error rates in gender identification across skin tones. It  was somewhat effective at changing the news narrative in the company's favor, as it was written up by a handful of technology-focused publications, like The Verge and TechCrunch.

If this study finds similar results, Microsoft or Uber may repeat that same strategy. Uber may publish stories on its blog that reframe the issue, emphasizing that Uber's facial verification software is a necessary part of its efforts to promote safety and prevent fraudulent drivers. However, the companies' counter also depends on how the public backlash divides itself between Microsoft and Uber. Uber cannot shift all its accountability onto the software maker, because it ultimately makes the decision to terminate employees based almost solely on these verification results. Microsoft, for its part, may emphasize that it is not responsible for users' actions. In its literature, Microsoft is careful to shift responsibility to users.

In the face of the two unions' lawsuits, Uber has maintained to media outlets seeking comments that its facial verification software is first and foremost a safety feature, "helping ensure the correct driver is behind the wheel [21]." The company has also told journalists that "The system includes robust human review to make sure that this algorithm is not making decisions about someone's livelihood in a vacuum, without oversight." [13] So even with study results  indicating high error rates for dark-skinned drivers in the system, Uber may continue to say the drivers were terminated for other reasons, that drivers can always opt for verification by human reviewers, and that it continues to review its algorithm for such biases.

It is likely that for Uber to make a substantive change in how it verifies drivers' identities, legislators and regulatory government bodies would need to get involved. With enough media coverage and open letters from different advocacy groups, legislators in the U.S. may be incentivized to support measures for more transparency and due process with terminations, joining 60 MPs in the U.K. [22] Likewise, Transport for London may act against Uber, as it did in 2019, stripping the company of its license. There might be slight reluctance from the agency to take such an extreme action a third time since Uber did implement its Real-Time ID Check in response to the agency's security and public safety concerns around fraudulent drivers in the first place. However, the agency could take a less extreme approach by requiring Uber to disclose its termination criteria or reduce unfair terminations by having a thorough and timely human review process following each identity verification mismatch.

In the U.S., the Justice Department sued Uber in 2021, accusing the company of discriminating against passengers with disabilities. It claims that Uber violates the Americans with Disabilities Act by charging "wait time" fees to customers who, "because of disability, require more time than that allotted by Uber to board the vehicle" [28]. The Washington Post called the lawsuit a "signal that the Biden administration is more aggressively targeting tech companies' civil rights records" [28]. A federal administration more concerned with civil rights issues, combined with a growing anti-Big Tech sentiment among federal watchdogs, is more likely to take regulatory action against Uber on discrimination issues. For instance, the Justice Department could file an investigation or a lawsuit against Uber for civil rights violations, as Title VII of the Civil Rights Act of 1964 prohibits discrimination based on race in any aspect of employment, including firing and testing [29]. That government action, in conjunction with media attention and advocacy group backlash, could bring about the envisioned result: requiring Uber to adjust its system of verifying drivers' identities so that it does not produce differing outcomes based on skin color, race or gender.


## Discussion

In summary, studies one and two  can be expected to show that under conditions similar to the ones Uber drivers face when completing the Real-Time ID Check, the company's driver verification system has the greatest mismatch rate for women with darker skin tones. This

result would extend the Gender Shades finding to face recognition in an employment setting and suggest that despite Microsoft's claim that itsFace API technology's gender identification error rates for men and women with darker skin have been reduced by up to 20 times and the overall error rate by nine times [27], the system still frequently fails when used to match images of darker-skinned individuals A finding that human reviewers were more accurate than the software in matching individuals of all skin colors and genders, especially under challenging conditions, would raise additional questions about the fairness of using this technology in a system where individuals' livelihoods are at risk.

Such findings would likely cause the two unions with ongoing U.K. lawsuits against Uber to release press statements noting that the study provided support for their legal claim that Uber's driver verification algorithm is racist and results in unfair terminations. The media might point out that the study's findings support previous allegations of racial bias in Uber's ID checks, while other advocacy groups might release statements of their own against the use of face recognition technology.

Government regulatory agencies in both the U.S. and the U.K. are likely to take separate actions. In the U.S., the Department of Justice could file an investigation into or lawsuit against Uber's facial verification and termination practices, as they present a possible civil rights violation if there is clear racial bias when firing drivers. Throughout this process, Uber is likely to lie low until it can put out its own statements and continue to assert that its facial verification software is necessary for promoting safety and thwarting fraudulent drivers.

Even without government regulation, the responses from media, advocacy groups, and the union would provide all the attention and incentive necessary to lead to a change — namely, Uber adjusting its system of verifying drivers' identities, so that it does not produce differing outcomes based on skin color or gender.

In the case that most of the components — union attention, media buzz, advocacy group backlash, legislator support and government regulatory action — do not take off, this envisioned result might not be possible. There are many ongoing controversies related to Uber's practices, including concerns over labor conditions and discrimination against passengers with disabilities [28]. As such, it could be possible that these controversies take the media and advocacy groups' attention away from these study results, leading to no expected change from Uber or Microsoft. If there was no media attention, but the unions' legal action against Uber proved successful, there still would likely be an opportunity for change. This is because the lawsuit might create enough media attention on its own or invite government regulation . It is unlikely that a change would be as extensive as the expected solution. Instead, it would probably be a toned-down version based on actions previously suggested by lawmakers, such as Uber increasing the transparency of its termination process.

The study might find no bias in results based on skin color or gender, or might even find a result in the opposite direction. Such results would not result in the change described, but would be meaningful and significant on their own, providing evidence that Microsoft's improvements to its Face API technology reduced error rates to the point where there is no skin color or gender bias or that the bias may now be in a different direction. This would likely spark a separate wave of media buzz, as it would counter the widespread perception that facial-recognition systems have a racial bias. In a time when a growing number of cities have passed bans on the use of facial recognition for policing, an opposite finding might cause re-evaluation of these bans and provide fodder for those who advocate other uses of biometric identity verification.

As for the proposed study three, a reasonable expectation could go either way — that the facial recognition software might outperform the team of human reviewers, or vice versa. Uber introduced a human review team to detect drivers who attempted to trick the facial recognition software by holding up false photographs to the app [5],yet human reviewers are also known to have biases. If the outcome is that the facial recognition software has a lower error rate and exhibits less bias than the human reviewers, Uber may change its verification process to rely more on the technology and use this result as a defense to drivers' claims that the software discriminates by race. But if the study finds that human reviewers outperform the software, this result would be more favorable for the unions and advocacy groups supporting Uber drivers, as it provides more evidence to back the claim that facial recognition technology is an unreliable and inferior way to verify drivers' identities. If the human reviewers are able to match dark-skinned individuals' photos as well as light-skinned individuals' better than the software, the results will provide additional support for claims of bias in the algorithm's performance.

The proposed studies have some methodological limitations. First, even if differences between the performance of Face API across demographic groups persist, it may not be possible to detect them unless the sample size for the proposed studies is rather large. The more sub-groups used to analyze intersectionality, the more difficult it will be to have a sufficient number of individuals in each group to allow for statistically significant findings.

Second, although the proposed experimental set-up seeks to replicate Uber's Real-Time ID Check as closely as possible, it is not the same. For example, Uber refers to its human reviewers as "identity verification specialists," implying that they have received at least some training [6]. Since it is not clear what type of training and guidance these reviewers have received, study four cannot replicate this element of the process exactly. It is also not known how reviewers are selected, and so the set of individuals serving as reviewers for study four may be quite different from those used by Uber. For example, Uber may recruit reviewers whose ethnic backgrounds match the dominant ethnicities of drivers in certain markets. Finally, there are additional details of the Real-Time ID Check system that are not known and may limit the studies' ability to replicate drivers' real experiences. For example, Uber does not disclose how it tunes the Verify threshold in Face API — that is, whether it requires a match that is stronger or weaker than the 0.5 confidence level that is the software default.

Additional limitations may arise from the conditions under which Microsoft will allow licensing of the Face API software for this research.

# References

[1] Sullivan, J. Selfies and Security. Uber Newsroom. Sep 23, 2016. https://www.uber.com/newsroom/securityselfies.

[2] Uber UK. Uber launches Real-Time ID Check for drivers in the UK. Uber Blog. April 30, 2020. https://www.uber.com/en-GB/blog/real-time-id-check-uk-drivers/

[4] Microsoft. Uber boosts platform security with the Face API, part of Microsoft Cognitive Services. Customer Stories. June 18, 2019. https://customers.microsoft.com/en-us/story/731196-uber.

[5] Pearson E. Police fears over Uber ID glitch after rapist impersonated another driver. The Age. September 11, 2019.

Police fears over Uber ID glitch after rapist impersonated another driver (theage.com.au)

[6] How does Uber verify my photo? Uber Help. Accessed December 16, 2021. https://help.uber.com/driving-and-delivering/article/how-does-uber-verify-my-photo?_ga=2.23046456.2046279747.1639871019-1067533333.1638365515&nodeId=aa821486-c8d1-42b7-b784-2fc24eb85f93.

[7] Bajaj, N, Sachdeva, S, and Kovalev, D. Engineering Safety with Uber's Real-Time ID Check. Uber Engineering. March 13, 2017. https://eng.uber.com/real-time-id-check/.

[8] Ghosh, S. Uber just lost its license to operate in London thanks to fraudulent drivers. Business Insider. November 25, 2019. https://www.businessinsider.com/uber-loses-london-licence-2019-11.

[9] MacDonald, A, and Schechner, S. Uber Wins Back License to Operate in London After Yearslong Battle. The Wall Street Journal. September 28, 2020. https://www.wsj.com/articles/uber-wins-back-license-to-operate-in-london-after-yearslong-battle-11601286874.

[10] IWGB takes Uber to court over racist facial recognition as Black Lives Matter UK backs the drivers' strike and boycott on Wednesday 6 October. The Independent Workers' Union of Great Britain. October 5, 2021. https://iwgb.org.uk/en/post/racist-facial-recognition/.

[11] ADCU initiates legal action against Uber's workplace use of racially discriminatory facial recognition systems. The Apps Driver and Couriers Union. Accessed December 16, 2021. https://www.adcu.org.uk/news-posts/adcu-initiates-legal-action-against-ubers-workplace-use-of-racially-discriminatory-facial-recognition-systems.

[12] Butler, S. Uber facing new UK driver claims of racial discrimination. The Guardian. October 6, 2021. https://www.theguardian.com/technology/2021/oct/06/uber-facing-new-uk-driver-claims-of-racial-discrimination.

[13]  Dent, S. Uber to face UK tribunal over 'racially discriminatory' facial recognition systems. Engadget. October 6th, 2021. https://www.engadget.com/uber-drivers-take-legal-action-over-racially-discriminatory-facial-recognition-093531420.html.


[14]  Buolamwini, J & Gebru, T. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. Proceedings of Machine Learning Research 81:1–15. 2018. http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf.

[15]  Castelvecchi, D. Is facial recognition too biased to be let loose? Nature. November 18, 2020. https://www.nature.com/articles/d41586-020-03186-4.

[16]  Transport for London. Taxi and private hire demographic stats.  December 2020. https://content.tfl.gov.uk/taxi-and-private-hire-driver-demographic-stats-2020.pdf.

[17]  Jessica. New Survey: Drivers Choose Uber for its Flexibility and Convenience. Uber Newsroom.  Dec 7, 2015. https://www.uber.com/newsroom/driver-partner-survey.

[18] MacCarthy, M. Mandating fairness and accuracy assessments for law enforcement facial recognition systems. Brookings Institution. May 26, 2021. https://www.brookings.edu/blog/techtank/2021/05/26/mandating-fairness-and-accuracy-assessments-for-law-enforcement-facial-recognition-systems/.

[19] Uber. Engineering Safety with Uber's Real-Time ID Check. Uber Blog. March 13, 2017. https://www.uber.com/blog/real-time-id-check/

[20] Microsoft. Face API. Accessed December 1, 2021. https://azure.microsoft.com/en-us/services/cognitive-services/face/#overview.

[21] Barry, E. Uber Drivers Say a 'Racist' Algorithm Is Putting Them Out of Work. TIME Magazine. October 12, 2021. https://time.com/6104844/uber-facial-recognition-racist/.

[22]  60 MPs back IWGB campaign to end unfair key worker terminations by Deliveroo and Uber. Independent Workers' Union of Great Britain. November 19, 2020. https://iwgb.org.uk/en/post/iwgb-campaigns-to-end-unfair-key-worker-terminations-by-deliveroo-and-uber/.

[23]  Klar, Rebecca. Civil rights organizations call for ban on corporate use of facial recognition. The Hill. April 14, 2021. https://thehill.com/policy/technology/548037-civil-rights-organizations-call-for-ban-on-corporate-use-of-facial.

[24]  Open Letter: banning government use of facial recognition surveillance is not enough, we must ban corporate and private use as well. Fight For The Future. April 13, 2021. https://www.fightforthefuture.org/news/2021-04-13-open-letter-banning-government-use-of-facial/.

[25]  Azure Cognitive Services. Products. Microsoft. Accessed December 17, 2021. https://azure.microsoft.com/en-us/services/cognitive-services/#features.

[26]  Kersley, A. Couriers say Uber's 'racist' facial identification tech got them fired. Wired. January 3, 2021. https://www.wired.co.uk/article/uber-eats-couriers-facial-recognition

[27]  Roach, J. Microsoft improves facial recognition technology to perform well across all skin tones, genders. Microsoft. June 26, 2018. https://blogs.microsoft.com/ai/gender-skin-tone-facial-recognition-improvement/.

[28]  Zakzewski, C. Justice Department sues Uber for charging 'wait-time' fees to passengers with disabilities. The Washington Post. November 10, 2021. https://www.washingtonpost.com/technology/2021/11/10/justice-department-uber-disabilities/.

[29] Title VII of the Civil Rights Act of 1964. The United States Department of Justice. Accessed December 19, 2021. https://www.justice.gov/crt/laws-enforced-employment-litigation-section.

[30] Bansal V. Uber's facial recognition is locking Indian drivers out of their accounts. MIT Technology Review. December 6, 2022. https://www.technologyreview.com/2022/12/06/1064287/ubers-facial-recognition-is-locking-indian-drivers-out-of-their-accounts/

[31] Microsoft. Characteristics, limitations, and best practices for improving accuracy. Microsoft Learn. June 23, 2022. https://learn.microsoft.com/en-us/legal/cognitive-services/face/characteristics-and-limitations

[32] https://customervoice.microsoft.com/Pages/ResponsePage.aspx?id=v4j5cvGGr0GRqy180BHbR7en2Ais5pxKtso_Pz4b1_xUQjA5SkYzNDM4TkcwQzNEOE1NVEdKUUlRRCQlQCN0PWcu

[33] https://www.theguardian.com/technology/2021/oct/06/uber-facing-new-uk-driver-claims-of-racial-discrimination