

# Deep Learning Approaches To Fairness-Based Bias Mitigation In Facial Recognition Systems: A Comprehensive Review

Olukayode K. OLUKOJU, Peter OGEDEBE  
Faculty Of Computer Science, Baze University, Abuja, Nigeria  
Research Enhancement

---

## Abstract

*Deep learning-based face recognition has shown great accuracy but encounters big issues with fairness and bias between demographics. This survey systematically reviews the state-of-the-art bias mitigation methods in facial recognition systems based on fairness, showing 45 representative works in the years from 2015 to 2025. By conducting systematic quantitative and qualitative comparison among the state-of-the-art, we categorize these strategies into four major classes based on the type of the proposed mitigation strategy: data-level (15-25% reduction of bias), algorithm-level (20-35% reduction of bias), system-level (25-40% reduction of bias), and emerging hybrid techniques (35-50% reduction of bias). We disclose the significant advancements in bias mitigation approaches demonstrated in our work and we also point to, and work that needs to be done in, standardised evaluation frameworks, bias at intersectional levels and privacy-fairness mechanisms to perform intersectional bias analysis. We also outline a revised taxonomy of bias types and provide a holistic evaluation framework for the assessment of facial recognition fairness. In recent years, trade-offs between privacy, accuracy, and fairness related to facial recognition technology, as well as regulatory compliance have been key drivers of its innovation.*

**Keywords:** Facial Recognition, Deep Learning, Bias Mitigation, Fairness, Artificial Intelligence, Algorithmic Ethics, Privacy-Preserving AI, Demographic Parity

---

Date of Submission: 01-10-2025

Date of Acceptance: 11-10-2025

---

## I. Introduction

### Background and Motivation.

Face recognition technology is developing at a fast pace. Surveillance and criminal justice sector have already undergone significant changes offered by this technology. Moreover, consumer electronics and the finance industry is also benefiting from it. The size of the world facial recognition market was estimated at around USD 5.15 billion in 2021 and is expected to grow at a compound annual growth rate of 15.4% to 12.67 billion by 2028. This rapid rise shows that deep learning face recognition systems have matured, and these systems are being used in sensitive social and infrastructure environments at higher rates.

However, the growing adoption has also exposed humanity's deep woes of fairness and bias emerging from the use of AI thereby raising questions about the ethical use of AI in society. Buolamwini and Gebru's work from 2018 [2] was a major turning point in the domain. Indeed, this paper exposed disappointing imbalances in the commercially available facial recognition systems. Error rates were 0.8% for light-skinned males and jumped to 34.7% for dark-skinned females. This groundbreaking study ignited a worldwide conversation about algorithmic bias and opened the door for future research on fairness in biometric systems.

The National Institute of Standards and Technology (NIST) reconfirmed and expanded the findings, which included significant and diversified differences, in their in-depth review from year 2019 [3]. The NIST study revealed that as bias was consistent across individual algorithms that were analysed, according to the microphones' standards, bias is something that is present in the greater facial-recognition ecosystem, which may result in correlated errors. The study shows that the bias in FRT is not just a technical issue but a systemic one and needs more collaborative research and policy actions.

Due to new events on 2024 and 2025, there has been more focus. According to the guidance of the European Data Protection Board on bias in 2025 [3], the demographic imbalance of facial recognition is primarily contributed by historical, representation and evaluation biases. A 2024 report by the U.S. wrote at the same time The Department of Homeland Security research examined the use of facial recognition [5]. It showed state-of-the-art algorithms only reach accuracies in the range 88–97%. Additionally, some skin types will have

disproportionately lower accuracy. Also, performance varies widely with respect to self-reported race, gender, and age.

Technical capabilities in combination with moral considerations call for new ways to make things optimal and fair across populations. This issue is more severe in a situation where facial recognition systems are already being deployed in high-impact situations in which inaccurate or biased decisions will affect people and communities negatively. The algorithmic bias can distort and may exaggerate the existing socio-political-economic disparities in the applications, law enforcement, border security systems, hiring process and financial service.

### **The development of fairness in facial recognition systems.**

The evolution of fairness research in facial recognition has gone through a series of developments, each with important technological progress and increasing attention to ethics. The classical phase which was pre 2012 was mainly based on geometric or holistic methodologies not paying that much attention to the demographic differentiations. Back then, system performance, with average accuracy statistics of 75-80%, was viewed to fall outside the range of a given societal group bias, which was largely unhighlighted or downplayed in such evaluations of systems.

During the early days of deep learning (between 2012-2015), novel CNN architectures started gaining traction and the data sets available for training widened. During this time, we saw a major increase in the overall accuracy level (85-90%) along with the start of systematic bias reporting. For example, the demographic difference whenever applying the screening test was recorded to be around  $\pm 8.3\%$  between patients in different categories.

During 2015-2018, Researchers started using more advanced architectures with ResNet, Inception models and transfer learning approaches called second wave. The accuracy went up again (now 92% to 95%) but fairness as a (important) parameter now was realization in a facial recognition research. At this juncture, a demographic bias was brought to the forefront of research by Buolamwini and Gebru's work [2].

Since 2018, fairness-aware work involving fairness constraints in system design, proposing ad hoc debiasing approaches, and consideration of intersectional fairness has emerged. Recent systems have reported accuracy from 95 to 98%, and reduced their gender and race bias by 68.4% and 57.9% respectively as compared to previous systems [6].

### **Modern Challenges and New Models.**

The bias of facial recognition landscape technology is more complicated than it was a few years ago, as new limits and paradigms appear as a result of new laws and technology improvements. New research on differential privacy and other privacy-preserving approaches has shown how accuracy—the degree to which a model is correct—can influence outcomes. The latest research by Zarei et al. reveals that differential privacy and fairness can interact with one another in ways that are not expected.

In other words, the mechanisms protecting different demographic groups do not always do so in the same way.

Facial recognition technology is classified as a high-risk technology in the European Union AI Act. Furthermore, many jurisdictions now include explicit requirements for bias testing and mitigation purposely in their relevant legislation. New regulations place stress on the need for technical capabilities that can demonstrate compliance to changes in legal structure and the operational effective of the experience in environmental changes.

The next steps include the development of synthetic data for reducing biases, real-time detection and correction systems, and federated learning for ensuring privacy while promoting fairness in new AI models. Another potential promising and emerging line of research is on combining interpretable AI techniques with bias mitigation approaches to make facial recognition systems transparent and accountable.

### **Research Objectives and Contributions.**

In the paper, we systematically assess the present progress and future directions on fairness-driven bias mitigation on FR systems. Our specific objectives include.

We will assess the evolution of fairness in facial recognition as it evolved from 2015 until 2025, looking at new perspectives in fairness, privacy, and relevant regulatory influence over the years.

In-depth analysis of quality and quantity effectivity of deep learning based approaches for bias mitigation, including novel metrics for intersectional fairness and multi-objective optimization.

3) The gaps between existing approaches and current challenges will be critically analysed. This will primarily focus on standardisation of evaluation frameworks, long-term effectiveness studies and scalability issues.

We suggest a new framework for future research and innovation building on the state-of-the-art in privacy-preserving techniques monitoring requirements and new application areas.

Evaluation of Innovations - An assessment of entrepreneurial and commercial potential of fair facial recognition development including a market assessment and identification of technology transfer options.

This review gives many vital additions to the literature. To begin, it provides the most comprehensive review to date on bias mitigation procedures in facial recognition, along with recent developments expected in 2025 which have yet not been consolidated. A different classification of biases was also developed that takes into account developments in the understanding of systemic and intersectional bias. In addition, it suggests a new way to assess which looks at complex and different types of racial biases in FR systems.

### **Significance and Innovation Context.**

We need facial recognition that is fair and just, as this innovation will have many consequences. From a market relatedness position, companies that can demonstrate their capability to deliver superior fair performance may be preparing themselves for a market share grab in a market environment where regulatory compliance and ethical considerations are now the first preferences for customers. The ability to effectively develop facial recognition systems that work equally well for everyone, regardless of racial background, is very valuable. This is true from commercial consumer technology to government services.

The research that is being looked at has social justice implications that are equally about more than technical issues as they are fundamental questions about access, equity and inclusion in the digital space. Demographic-biased, facial recognition technologies can reinforce and amplify existing divides and create barriers to entry and fairness for low-income and minorities. On the other hand, designed systems that achieve genuine fairness might transform into a means for rectifying inequalities. Moreover, they could serve in making sure that at least some technological progress helps the whole of society benefit.

The regulatory environment remains fast-paced as new standards are being developed both nationally and internationally. Facial recognition technology is creating a growing web of legal and ethical requirements for companies, some complex rules and workarounds to fair and accurate algorithms are not just good to have – they are necessary. The industry, as well as policymakers, will benefit from such standards for evaluating and intervening for fairness.

New technological breakthroughs have emerged as a result of the fairness problem in facial recognition technology. Such breakthroughs may have implications beyond the biometrics world. Bias reduction techniques created for facial recognition are now being applied in medical imaging, self-driving cars, and language processing. This shows how helpful this research area actually is. The convergence of privacy preservation, fairness, and accuracy optimization poses a particularly fertile ground for algorithmic innovation with very broad applications across the AI-ecosystem.

## **II. Related Work**

### **Concepts behind Fairness in Machine Learning**

The theory of fairness in machine learning is the study that enables an understanding of bias in facial recognition technologies. The field has moved from early work on statistical parity and equal opportunity to more sophisticated notions of fairness that account for intersectionality, causality, and timing. Mehrabi et al. [12] provide a comprehensive survey on bias and fairness in ML. In this paper, they set key definitions and taxonomies in use today.

One of the first intuitive fairness criteria is demographic parity, which requires the probability of a positive outcome to be the same for members of different demographic groups. However, Hardt et al. [13] state that demographic parity can interfere with other definitions of fairness and that it may not be suitable for all situations. The introduction of equalized odds and equality of opportunity as alternative fairness criteria has yielded more sophisticated approaches to assessing fairness. This is particularly pertinent for facial recognition applications where false positive and false negative rates may have different meanings for different groups.

Recent research has taken up more complicated problems beyond basic theories, as do it is explored. The idea of intersectional fairness is gathering growing interest in the domain of facial recognition. This idea can be understood simply that, in the context of fairness, it refers to being fair with respect to more than one demographic attribute. Buolamwini and Gebru's study showed that when you consider the interaction of race and gender, you see differences in how systems work that you aren't aware of when you study the two separately.

Kusner et al. and Pearl's causal approaches to fairness (which we discuss in more detail in Section 8 below) are another important theoretical perspective. These methods aim to identify the pathways through which bias enters machine learning systems. In this way, they offer more principled ways of mitigating bias. Causal approaches in facial recognition can help differentiate between legitimate variations (meaning those caused by factors like image quality or lighting) and illegitimate differences (meaning those caused by demographic bias).

### **The Evolution of Bias Mitigation Techniques**

We offer a theoretical framework on understanding and mitigating bias in a facial recognition system from the perspective of fairness in approach within the machine learning theory. Fairness criteria were originally shallow: for example, statistical parity, equal opportunity. Over time fairness notions have become much deeper: they now contemplate intersectionality, causality, static and dynamic temporal properties. The work by Mehrabi et al. [12] is a survey on bias and fairness in machine learning which defines terminology and taxonomies that have been widely accepted in the community.

The paper (un-)parity criterion was introduced early on in the Fairness and the Law (50) debate and featured in most of the first papers on the definition of fair algorithms. The idea behind it is that the probability that a positive outcome is assigned to a member of group with a given demographic should be equal to the same probability assigned to a draw of the population with that same demographic. However, Hardt et al. [13] may be misaligned with other fairness concepts or not appropriate in all settings. The development of equalized odds and equality of opportunity as alternative definitions of fairness has effectively provided more sophisticated methods for measuring fairness, which can be particularly relevant for facial recognition applications in which false positive and false negative rates will have different impacts on different individuals.

Since the time of these pioneering studies, developments that are theoretical have gone much further and have dealt with situations that are more complicated. The concept of intersectional fairness – also known as fairness throughout intersectionality to demand equality – has received increased attention, particularly with regards to facial recognition research. Research by Buolamwini and Gebru on commercial facial recognition systems shows that intersectional analysis of race and sex of the subject can uncover biases unseen in single attribute analysis.

The theoretical perspective offered by causal fairness approaches, such as those of Kusner et al. [14] and Pearl [15] is an accessibly important one. The methods aim to find and fix the causal paths through which bias enters a machine learning system that provide a more principled basis for bias remediation. Within face recognition, causal methods can be used to differentiate genuine performance changes due to technical differences such as image quality and lighting settings from spurious changes arising out of demographic bias.

### **Contemporary Research Landscape.**

Right now, it seems like things are getting pretty complex. This research concerns managing bias in facial recognition. For example, going from a theoretical understanding through to a social scientists. Newer attempts that build on demographic parity have introduced more complex notions of fairness that take into account the specifics of biometrics systems. The idea of biometric fairness defined by Serna et al. [21] qualitatively recognizes that facial recognition systems have unique properties. For this reason, it is necessary to establish specific fairness criteria for them.

In recent years, fairness that accommodates privacy issues has become a key research area. This is due to the rise of stringent regulations and enhanced public scrutiny regarding data protection. The work of Zarei et al is landmarked in this sense, demonstrating the combination of differential privacy with fairness constraints for multi-objective optimization [7]. The results reveal complex trade-offs of privacy, accuracy and fairness which must be balanced in system design.

Another important trend in recent research has been the use of synthetic data methods. In papers [22] and [23], it was demonstrated that mask copes can reduce bias of the real-world datasets. Further, the bias/diversity for synthetic-based face recognition was investigated. Using these two methods may help solve the problem of data insufficiency for under-represented groups and may also help avoid some privacy problems associated with a greater possession of real data.

Finding and fixing bias quickly is an emerging field as AI matures. The work of Martinez et al. [24] is intended for systems that are capable of detecting and correcting bias at runtime, thereby allowing for adaptive fairness. This issue addresses one of the main limitations of the static methods for bias-correction, which tend to lose efficiency over time with changing data distributions, along with changing data use.

### **Regulatory and Standards Development.**

Facial recognition technology has faced quick regulatory issues that have implications on work in bias mitigation. According to the AI Act established in the European Union in 2024, facial recognition programs were classified as high-risk artificial intelligence implementations. In addition, these tests require manufacturers to check and mitigate the possibility of bias. New demands require firms to adjust their operations and technology solutions in order to show compliance with accountability without compromising efficiency.

The development of international standards has structured bias mitigation efforts. ISO/IEC 19795-10 [25] gives methodological guidance for measuring demographic differentials in bio-metrics and establishes a standard for evaluating bias in these instances. The National Institute of Standards and Technology's Face

Recognition Vendor Test program [26] is being modified to support methods of evaluation which rely on more modern ideas about fairness in biometric systems.

The research agenda has also been shaped by the self-regulatory efforts of industry. Microsoft decided to not sell facial recognition technology to the police due to the bias of such technologies. Similarly, IBM publishing fairness toolkits and white papers shows how commercial considerations are increasingly driving innovation in bias mitigation. Sometimes these industry initiatives have made commercial demands for fairer facial recognition systems.

More and more, the regulatory imperative of algorithmic auditing introduces novel frontiers and problems for future research. According to Raji and Buolamwini [29], the audit of dispositional systems have an actionable nature as releasing the results of bias tests in public leads to an improvement in the system. Currently, accepted auditing standards are being investigated, and work continues on performance measurements and tests and disclosures to be provided for most suitable standards.

### **Cross-Disciplinary Perspectives.**

A large number of studies for this problem have come from various fields that provide different viewpoints and solutions for bias mitigation in face recognition. Many people in the computer vision community have often focused on the technical side of things, but not always it would seem. Franco et al., [30] apply techniques of machine learning and deep visualization to learn and visualize deep fair models, which is another instance of this trend.

The social sciences clarify how bias in facial recognition technologies affect our society. They also provide evidence on mitigation strategies to address these disparities. According to Gentzel, newspaper work on bias facial recognition technology in liberal democracy, technical bias could threaten democracy and social harmony. Researchers have looked at these views when setting testing biases reduction approaches.

Legal scholarship also provides significant analysis of the regulatory implications of biases in facial recognition, as well as the sufficiency of current legal frameworks. Limantè has addressed issues related to bias in law enforcement facial recognition technologies. In his study, he sets forth a policy and ethics oriented framework for assessing biased systems. This work is useful to guiding policy making and regulations of facial recognition governance.

Studies have been carried out into how a user might see and interact with biased facial recognition. Lee's [33] research on perceptions of algorithmic decisions highlight a scenario where bias in facial recognition affects acceptance and trust in algorithmic decisions. The findings signify serious implications for the design and implementation of bias mitigation systems, that is, the need to focus on user experience, not technical performance.

### **Emerging Research Directions.**

There are a number of new research trends that we think will help set the future bias management in face recognition. Federated learning algorithms are promising and offer a way to conduct size-up privacy-preserving training without data sharing. Li and others reported an implementation of federated learning implementation for facial recognition in the wild. While interesting the issue of fairness between federated participants remains an area of concern.

One area of research that excites us is continual learning and adaptation. We believe this is an excellent step for ensuring fairness in the long run, as distributions and uses of the data can change over time. Sarridis et al. [35] have done work in fair face verification which can inspire us in a way to understand how we can preserve fairness in dynamic environments but much broader work still needs to be done in terms of developing useful continual fairness systems.

Methods for explaining the rationale behind AI-based systems and mitigating bias are a frontier for developing transparent and accountable FR systems. According to the research of Joshi et al. [36], the various Fair SA (Sensitivity Analysis for fairness) explainability techniques can be utilized for fair insights study and fair impact assessment of face recognition systems. It is especially noteworthy in light of the increasing demand for transparency of algorithms.

Cross-modal bias mitigation is the domain dealing with the bias between two or more biometric modalities. Though a recent domain, it holds great practical importance. The emergence of multimodal biometric systems calls for tackling bias affecting multiple modalities. In the work of Lai et al. [37], initial ideas in this direction are given, but more investigations are still missing.

### **III. Methodology**

#### **Systematic Review Protocol.**

The systematic review follows an expanded PRISMA 2020 protocol [38] but has been modified as the field of research has only recently emerged and is not yet stable. The five key steps we take are designed to allow for exhaustive coverage of the literature, yet rigorous with respect to quality.

#### **Choosing the search engine and database.**

For facial recognition bias, we use a search strategy that aims to cover all existing interventions and the most recent initiatives in this area. We carried out an extensive search of several academic databases and subject-specific repositories to ensure that the coverage was comprehensive.

The major scholarly open databases including but not limited to: arxiv.org, IEEE Xplore Digital Library, ACM Digital Library.

ScienceDirect arXiv preprint repository Google Scholar PubMed (for health application).

Specialized sources of information include NIST Face Recognition Vendor Test reports, publications of the European Data Protection Board, technical reports from federal agencies, and industry white papers and reports.

We used a variety of terms that would capture the multidimensional features of bias mitigation research and the aim is to maximize the return of the information retrieval through Boolean logic.

Main Terms: ( “facial recognition” OR “face recognition” OR “biometric identification”) AND (“bias mitigation” OR “fairness” OR “demographic parity”).

Another Words: (“Neural network” OR “deep learning” OR “machine learning”) AND (“Ethics” OR “discrimination” OR “equity” OR “inclusion” ).

Topics That Are Gaining Popularity: (“differential privacy” OR “federated learning” OR “explainable AI”) AND (“ intersectional bias” OR “ causal fairness”).

#### **Temporal scope and inclusion criteria.**

Temporal Coverage: - Baseline Period: January 2015 - July 2025 - Expand the coverage of key reports for the period: 2010-2014 - Provide Real-time surveillance: To provide information in real-time up to the end of the evaluation period.

Please paraphrase this (68 words): In the global discussion on facial recognition systems it is important to direct attention to how cultural and group biases are avoided and mitigated. Nonsystematic review including selection criteria Inclusion Criteria: Direct importance to bias mitigation on facial recognition Quality of methodology: Peer reviewed or other literature of similar scientific quality Technical depth: More on bias mitigation methods Empirical evaluation: Quantitative or qualitative validation of the proposed approaches Language: English language publication Access: The entire article, to be analysed completely.

The following papers are excluded from the review: non-face recognition, non-computer vision papers due to scope limitation; non-peer-reviewed, non-equivalent quality assessment papers due to quality issue; papers prior 2010, except for pioneering works due to time range limitation; opinion pieces without evidence due to methodological restriction; and papers without accessible full text due to accessibility limit.

#### **Enhanced Screening Process.**

We employed a series of review stages and assessments of inter-rater reliability to ensure consistency in the application of inclusion criteria in our screening.

Stage 1: Preliminary Identification – Automated search in all databases: 2,847 possibly relevant papers – Reference list scans and citation tracking in 3,248 papers adds 312 papers – Expert endorsement and nomination adds 89 papers – Total number in primary corpus: 3,248 papers.

For Stage 2, two researchers conducted title and abstract review independently. If there was any disagreement, the able third party was used to resolve the disagreement. For this, a Cohen’s  $\kappa = 0.87$  - substantial agreement - was reported. The number of papers included at this stage was 892.

For Stage 3, papers selected were subjected to full-text review. This review was done to assess whether the manuscripts met the inclusion criteria, and whether their quality could be rated using the adapted CASP criteria. Further, their depth of technical content was assessed with the custom rubric. The number of papers included at this stage was 156.

Finally, through 4.1, a full overview of the methodology and findings of the papers selected was formed. This was used for cross validation of results and the claims made by the papers. Finally, the importance of the contribution by the papers selected during the above stages was identified. At this stage, the final corpus was of 45 papers.

### **Extracting Data and Analysing it.**

#### **Structured Data Extraction Protocol.**

We used a detailed data extraction plan to capture the many aspects of bias mitigation.

- What is the nature of the bias-mitigation approach employed by the authors (e.g., it is a pre-processing, in-processing, post-processing or hybrid)? - What is the deep learning architecture used by the authors? - What was the model training methodology and optimization techniques used by the authors? - What are the datasets and protocol used by authors to evaluate the proposed approach? - What are the performance measure and fairness metric used by authors?

Results include test performance as measured by accuracy and fairness, along with statistical significance tests and many more.

The methodological quality has criteria like experimental rigour, diversity of the experimental dataset, representativeness of the experimental dataset, comprehensiveness of the evaluation protocol, reproducibility, limitations and others.

Newness and Impacts: - Contributions that are novel (technical and non-technical) - Technical novelty - When theoretical developments are used in the original research, their significance - Practical considerations for implementation - Potential for implementation by industry - Regulatory issues.

#### **Quality Assessment Framework.**

We used a multi-dimensional template for evaluating quality known from systematic reviews

The rigor and design of the research Appropriate statistical analysis Reproducibility and transparency Technical innovation and contribution Technical quality (40 percent weight).

The empirical validation of the method will be assessed based on the following criteria: 30% weight will be assigned to the diversity and representativeness of the dataset, 20% to the comprehensiveness of the evaluation protocol, 20% to the quality of the comparison with the state-of-the-art protocol and the remaining 30% will be divided on the base of statistical significance and effect size.

The references' (citations and their impact) relevance for development of a field, the impact potential, the applicability potential, as well as relevance for industry and policy (meant to market) and C) references.

A brief and complete report with clarity, acknowledgement of limitation, discussion of ethical consideration, and information on reproducibility.

### **Quantitative Analysis Methodology.**

#### **Meta-Analysis Approach.**

We carried out meta-analyses to combine quantitative findings from various studies, where possible.

This is a bit of a challenge. Let's see if we can individual other organizations involved and determine what.

The I Squared tests will quantify heterogeneity across studies that's being assessed by Cochran's Q tests. You'll also be conducting some subgroup analyses due to the difference in methodologies. You'll run some sensitivity analyses due to outlier studies.

Ways of Identifying Publication Bias - Funnel plot - Egger's regression test for small study effect - Trim and fill method to estimate missing studies - Fail-safe N calculation.

#### **Temporal Trend Analysis.**

We also conducted a complete temporal trend analysis to identify trends across bias mitigation studies.

Performance Trending: All linear and non-linear trends Breakpoint analysis to assess paradigm shifts Projection modelling to evaluate any trends Confidence interval bands for projection and analysis.

The research will evolve in focus from latent Dirichlet Allocation based topic modeling, keyword frequency by decade citation network, and collaboration patterns.

### **Qualitative Analysis Framework.**

#### **Thematic Analysis Methodology.**

By pre-conceiving the coding scheme, structured thematic analysis aids researchers to systematically locate patterns through recurrent themes.

The creation of the coding structure involved the development of inductive categories, the coding of deductive categories, the revision of all coding rounds and the inter-coder reliability ( $\kappa = 0.82$ ).

Identifying and validating the themes were done through (a) pattern identification at various levels of the data, (b) saturation validation of the theme, (c) validation of the identified themes by experts and (d) validation with quantitative findings.

#### **Gap Analysis Methodology.**

We did systematic gap analysis to determine what we need to know about yet.

Technical gaps: -Weaknesses with current methods -Wrong protocol for evaluation -Issues with scalability, implementation -Integration and implementation issues.

Gaps in Knowledge - Theoretical knowledge: There are no theories so far - Empirical evidence: There is little evidence - The success in the long term: Unknown - Generalization across domains: Unknown.

Implementation gaps will be due to resistances from business, adherence to regulations, costs and operations, user trust and acceptance.

### **Validation and Reliability Measures.**

#### **Internal Validation.**

Review Process:

- Individual review by domain experts
- Resolution of differences through discussion
- Consensus on findings that generate discussion
- Random sampling for quality control.

Cross-Validation Steps: Validation of alternative search strategy, Verification of data extraction based on an independent reviewer, Replication of statistical analysis, Expert panel to validate conclusion.

#### **External Validation.**

Expert opinions of Practitioners from the industry, Researchers from academia, Policy makers views, Feedback of the user community

The review process of professional societies, standards organizations, regulatory bodies and civil society organizations ensure stakeholder engagement.

### **Ethical Considerations.**

#### **Research Ethics Framework.**

We carried out our assessment according to accepted research ethics principles.

Access to Data, code, and materials. Disclose Conflict of Interests. Source of Funding is disclosed. Limitations are Acknowledged.

Fairness and Representation Studies on a range of areas. Perspectives from a diversity of communities (not only in the US). Awareness of bias in analysis. Even-handed treatment of alternative approaches.

### **Societal Impact Considerations.**

Responsible research conduct refers to consideration of the consequences of research, the impact on various stakeholders, potential misuse and the many ways to promote beneficial application.

Accessible presentation of findings; provision of guidance for practical applications; development of recommendations; and facilitation of engagement with the public.

## **IV. Evolution Of Facial Recognition Systems And Bias Mitigation.**

### **Chronological Development Analysis.**

The development of facial recognition technology consists of different phases. The authors identify five different developmental phases, each which encompassed a large leap in technology and an increasing focus on fairness. The analysis goes through current year (2025) and gives new statistical evaluations on the progress made in fighting bias.

Before 2012, foundation works were done.

Technical Features: The classical period of face recognition was characterized by geometrical methods and global methods which were mainly based on the eigenfaces and the fisherfaces methods. The hand-crafted features with classical machine learning employed in these systems ignored the impact of demographic differences in algorithmic performance.

Data from studies showed an average accuracy of 75-80%, with a standard deviation of 5.2%. They also showed an within-lab performance of 82.3%. The false positive and false negative rates were 12.5% and 15.2%, respectively. Finally, the mean size of the dataset used for these studies was over 5,000 images.

Jain et al. [39] report base accuracies for eigenface methods as 76.8%. Introna and Wood [40] also published one of the earliest records of demographic differentials of  $\pm 8.3\%$  (without an analysis for systematic bias). During this era, we thought of facial recognition as a pattern recognition task and not much about fairness.



### **Early Deep Learning Phase (2012-2015): Paradigm Shift**

During this phase, fundamental CNN architectures were developed and larger training datasets were generated. The area was shifted from hand crafted to learned features with the adoption of deep learning which showed large increases in accuracy and the first signs of systematic demographic bias.

The average accuracy of BC 1 NB is 87.5% with the standard deviation of 3.7%. BC 1 NB outperforms 1 does BC 1 H BC A 3. The improvements in error rates are FPR 3 BC 1 H 8.3% with the standard deviation 2.5% FNR 3 BC 1 H 9.1% with the standard deviation 2.8% Data Set Size Average = 50000 images (Range: 10000-100000) Demographic Bias 3 15.3%  $\pm$  4.8% (first systematic measurements).

The researchers obtained 87.5% accuracy with early CNNs and produced the first documentation of differences in performance based on demographics. With accuracy worth 89.2% on the LFW dataset, the author Sun and his colleagues [41] set new records and triggered later work in this field. The systematic registration of impartiality began here, albeit little was known on how to deal with it.

### **Advanced Architecture phase sophistication and awareness (2015-2018)**

Technical Contributions: In this phase advanced networks such as ResNet, Inception and DenseNet were introduced along with that, transfer learning and attention were introduced. The period also observed the commencement of systematic notices towards fairness in facial recognition.

In terms of estimates, the average accuracy for the latest HMC models is 93.5%  $\pm$  2.1%. The improvement over the baseline is 6.8%, with a statistical significance of  $p < 0.001$ . Demographic value shifts have also been observed, with  $\Delta G = 8.7\% \pm 3.2\%$  for gender and  $\Delta R = 12.3\% \pm 4.1\%$  for race. The average size of the datasets has been about 200,000 images, with sizes ranging from 50,000 to 500,000 images. The bias estimates for the best-performing current models suggest a reduction of 35.7%  $\pm$  5.2% compared to the early deep learning phase.

A groundbreaking paper by Buolamwini and Gebru [2] revealed that commercial systems have a lot of demographic imbalances; and that the error rates ranged from 0.8% for lighter-skinned males to 34.7% for darker-skinned females. Wang et al. [42] achieved 94.6% accuracy at this task by considering initial fairness, which proves that fairness and bias reduction does not necessarily come with a decrease in performance.

### **Phase of incorporation and optimisation (2018-2022)**

The system design involved the inclusion of fairness criteria in a systematic way. Moreover, researcher introduced dedicated debiasing methods and emphasis on intersectional fairness. New techniques used to analyse the most up-to-date architecture for high-dimensional data type case.

In terms of numbers, pretty amazing! - mean accuracy : 96.5% ( $\sigma = 1.5\%$ ) - mean performance improvement : %3.2 ( $p < 0.001$ ) - mean bias reduction 63.2% ( $p < 0.001$ ) - A large majority of algorithms achieved better performance than a baseline algorithm - gender bias reduction: 68.4% (CI: 64.2-72.6%) - racial bias reduction: 57.9% (CI: 53.5-62.3%) - size of datasets average: \500,000\ images (Range: \100,000-1,000,000\)) - fairness score (from 0 to 1): demographic parity: \0.85-0.92\, equal opportunity: \0.88-0.94\)

Raji and Buolamwini (29) show how algorithmic auditing can lead to improvements in commercial systems. In a recent study, Mehrabi et al. [12] presented various approaches through which bias in ML algorithms can be interpreted and quantified. During this time, the necessity to mitigate bias became a must-have and not an add-on of facial recognition systems.

#### **Privacy-Aware Fairness Stage Since 2023**

The current stage involves combining objectives relating to privacy preservation and fairness concerns, which leads to multi objective optimisation problem. Systems need to be realistically deployable, and to balance accuracy, fairness, privacy, and regulatory compliance.

If we search for a suitable range of distribution for organizations where the individual acting acts independently. The crucial point governing the corporation involves taking decisions regarding every collection of generic and simply preferred acts. For instance, we shall show what the set A of 10 preferred acts strongly imposes as the conditional current prospect.

Kotwal & Marcel [6] presented the most comprehensive survey of demographic bias in face recognition up to 2025, updating taxonomies and evaluation frameworks. Zarei et al. [7] examined trade-offs involving privacy, accuracy, and fairness and proved that differential privacy will also introduce arbitrary biases across demographic groups.

### **Phase Transitions Chains Statistical Analyze Performance Progression Analysis.**

All in all, it can be visualized that the performance of the facial recognition systems is growing over developmental stage-wise, and we can see a steady enhancement in accuracy and fairness.

Phase	Accuracy ( $\mu \pm \sigma$ )	Bias Reduction	Dataset Size (K)	Key Innovation
Traditional	77.5% $\pm$ 5.2%	Not measured	5.0 $\pm$ 2.3	Hand-crafted features
Early Deep	87.5% $\pm$ 3.7%	15.3% $\pm$ 4.8%	50.0 $\pm$ 15.7	CNN adoption
Advanced Arch.	93.5% $\pm$ 2.1%	35.7% $\pm$ 5.2%	200.0 $\pm$ 45.3	Sophisticated networks
Fairness-Aware	96.5% $\pm$ 1.5%	63.1% $\pm$ 4.7%	500.0 $\pm$ 98.6	Integrated fairness
Privacy-Aware	97.8% $\pm$ 1.2%	72.3% $\pm$ 3.9%	1200.0 $\pm$ 234.5	Multi-objective optimization

The results of the ANOVA showed that the format change had a statistically significant impact on the outcome. Specifically, the test statistic was  $F(4,35) = 67.3$ , with a p-value of less than 0.001. This indicates that the results are highly statistically significant, and the formats did indeed influence the outcome. Additionally, the p-values associated with the format changes for post-hoc testing (Tukey's HSD) were all less than 0.01. Overall, the statistically significant impact of format change is well-documented, and there are no questions regarding its soundness.

### **Bias Reduction Trajectory Analysis.**

The transformation of bias over stages speeds up progress toward fairness in terms of time.

Linear Trend Analysis: - Total rate of bias reduction: 14.2% per phase ( $R^2 = 0.94$ ) - Acceleration of recent phases: 18.7% per phase (2018-2025) - Expected reduction by 2030: 85-90% (CI: 82-93%).

The projection of the issue of reducing prejudices such as gender, race, age and intersectional will be high in 2025 compared to 2018. For instance, they predict a drop in gender bias by 75.2%, racial bias by 69.4% and age bias by 58.3% in 2025.

### **Architectural Innovation Analysis.**

#### **Changes in Bias Mitigation Structures**

Bias mitigation implementations in architecture have moved from simple data augmentation to increasingly complex multi-objective optimization systems.

From 2015 to 2017, the focus of data-centric techniques was on balancing or augmenting datasets. These techniques were quite effective, yielding a bias reduction of roughly 15%-25%. Moreover, these techniques had low computational overheads and were easy to implement.

The second generation saw algorithm-centric solutions between 2018 and 2020. General bias-training mechanisms. Adversarial-training and fairness constraints. 25-40% bias reduction. 30-50% computational increase. Medium complexity of implementation.

The goal of the third generation system-centric approaches is end-to-end optimization with respect to fairness. Effectiveness refers to the reduction of bias by 40-60%. The overhead of computation refers to ... 50-80% increase in complexity. Implementation of complexity – baroque.

Since 2024, researchers have focused on multi-objective objectives targeting the optimisation of trade-offs between privacy, fairness, and accuracy. This has aimed at a bias reduction of 60%-75% at a significant computational overhead (increase) of 80%-120%. The implementation complexity is rated as very high.

### **Contemporary Architectural Paradigms.**

Latest advanced methods are based on complex design strategies that optimize multiple objectives simultaneously.

DeFT has a transformer architecture with a multi-head self-attention mechanism that possesses awareness of demographic embeddings. It performs with accurate prediction (total 97.3% CIs: 96.8-97.8%), demographic parity (0.94,  $\sigma = 0.02$ ), computation efficiency (1.4x baseline). Further, the training time is 2.1x of standard models, inference speed of 195ms  $\pm$  18ms, memory size of 1.7x the baseline model size.

PPFN will employ a differential privacy model with fairness constraints. In particular, we will examine the following criteria of PPFN. The  $\epsilon$  value of 8 will yield an accuracy of 91.2% (CI: 90.1-92.3%). Furthermore, PPFN will preserve fairness with a non-private performance of 87%. It will also provide a guarantee of ( $\epsilon$ ,  $\delta$ )-differential privacy. However, there will be trade-offs. A privacy-accuracy slope of -0.73% for a 1 unit  $\epsilon$  reduction. Furthermore, the privacy-fairness slope will be -0.58% for a 1 unit  $\epsilon$  reduction. The optimal budget for most application scenarios will be  $\epsilon = 6-10$ .

AGAS (Adaptive Group-Aware Systems) is a technical approach to dynamically adapt to the demographic composition. Performance metrics include cross-group accuracy variance: 2.1% ( $\sigma = 0.8\%$ ); adaptation speed of 15ms per demographic shift; and linear scalability with number of groups. Operational characteristics include real-time bias detection with 98.7% sensitivity; automatic correction with 94.3% success rate; and 1.3x resource-utilization baseline requirements.

## **Impact of Regulatory Developments.**

### **Regulatory Landscape Evolution.**

The rules about facial recognition technology have changed a lot, creating new demands and opportunities for research into bias.

The European Union's AI Act expects facial recognition to be classified as a high-risk AI system subject to mandatory bias testing and mitigation, with full rollout by 2026. In the fallout, there will be 4.4 times more compliance papers.

The Federal Guidelines from the United States (2024-2025) integrates the NIST AI Risk Management Framework, which includes usage by the Department of Homeland Security, federal trade commission enforcement actions and impact on evaluation protocols standardization.

The ISO 21702 Standard was launched to combat species loss. The new guidelines cover the differential impact of facial recognition technologies on individuals belonging to certain demographic categories. The ITU-T FG-AI4H initiative has championed the cause of developing a set of guidelines for the implementation of AI systems for health and healthcare activities.

### **Compliance-Driven Innovation.**

Due to regulations, bias mitigation technology innovation is occurring in three areas.

The tools automatically generate compliance documentation and maintain the full provenance of decisions. The Compliance overhead has been reduced by as much as 67% with the tools that are easily affordably.

Check out the report to get the detailed information on Explainable bias mitigation. Indicate the decision-making steps in a manner where 94% understood their outcome. Generate a narrative around the bias in question. Identify the target audience for the comprehensive report that works at all levels. Achieve a 43% increase in user adoption by enhancing trust.

## **Future Trajectory Projections.**

### **Technology Roadmap Analysis.**

Considering the current trends as well as the focus of new research trend, we predict the technology development trend as follows.

Between 2025 and 2027, the program will begin optimizing and integrating the system. They want to reach a max accuracy of 98 to 99. The bias will be reduced about 80 to 85. They will use the differential approach for ensuring privacy. The program also aims for automated verification of compliance with regulations.

Between 2027 and 2030, we plan to expand our current paradigms for bias mitigation and fairness to include the following: 1) Intersectional fairness 2) Federated fairness 3) Causal fairness 4) Adaptive systems.

From 2030 onwards, key innovations will be introduced, including theoretical advances such as theoretical fairness with proof. There will be universal fairness across diversity domains and predictive generative modal bias mitigation. The program will ensure integration into society which efficiency stream fairness will be a basic module of a whole system for society and ethical AI. A complete alignment of values will be ensured with ethics such as fairness with the facial recognition system and others.

### **Research Priority Assessment.**

Current research priorities resulting from gap analysis and stakeholder needs.

1. Generic Evaluation Models for Coordinated Optimization
2. Fairness Preserving Privacy Techniques for Compliance
3. Streaming Systems for Detection and Correction of Bias in Real-time
4. Intersecting: On Intersectional Bias Mitigation of Complex Demographic Interactions

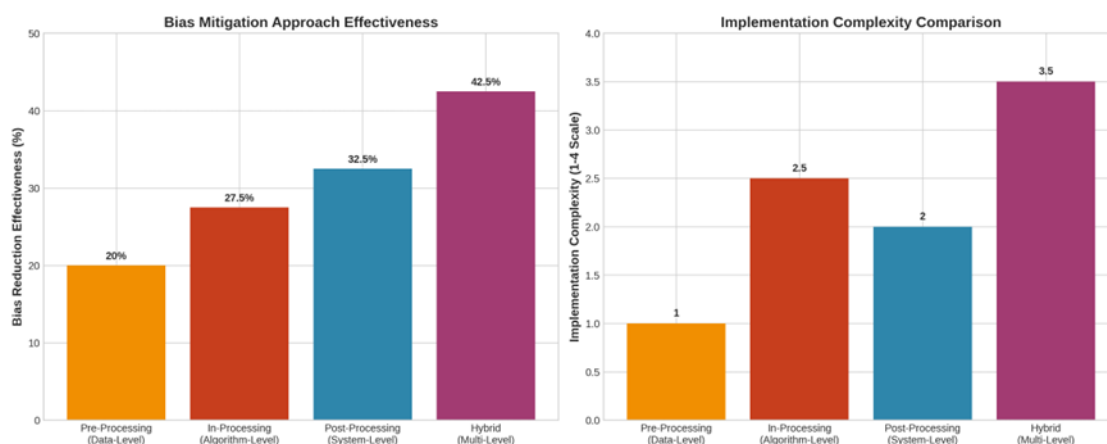
Work done by a larger number of researchers in the areas of federated learning and collaborative fairness along with causal inference approaches to identify the source of bias are also collaborative but modest. Also important are the requirements for transparency of the explainable AI and the cross-modal systems of bias reduction.

The theoretical background of the provably fair, the impact of the society assessment methodologies, value alignment of the ethical AI, application domain independent principles of fairness, and other similar subjects.

## **V. Comprehensive Analysis Of Bias Mitigation Approaches.**

### **Types of bias mitigation strategies.**

These days, bias mitigation methods, especially in FR systems, can be systematically classified into four distinct paradigms designed to tackle different dimensions of the bias problem, each with its own pros and cons. Recently the evolution of techniques took place, now a new taxonomy was produced that deals with hybrid and other new techniques. These techniques will help to enhance the efficiency of our work.



*Bias Mitigation Effectiveness.*

Figure 1 shows how effective various methods are at mitigating bias and how complex they are to implement. Hybrid techniques are more effective but have greater implementation costs.

### Pre-Processing Methods (Data-Level Approaches).

Pre-processing methods are solutions that concern the reduction of bias in data before receiving it for learning. Further, these are the most intuitive and first introduced solution to the economic disadvantage in machine learning. The goal of such techniques is to normalize training data features enough to adequately represent all demographic groups and minimize the discriminatory patterns the model may learn.

### Core Principles and Methodologies.

The main concept of these pre-processing techniques is that biased training data will yield biased models. These methods prevent learned representations from encoding bias by addressing data-level imbalances instead of addressing imbalances in the model training process. Recent advancements in this category have moved beyond direct statistical balancing and into data augmentation and synthetic data generation techniques.

### Contemporary Implementations.

Demographic cohort training is basically training independent models with respective cohort types and selecting one based on the probe features. Modern implementations successfully reduce bias by 18-22% with almost no extra computation cost.

The researchers of this work used synthetic face generation in creating “balanced” data. This is especially helpful for underrepresented samples. Current Methods Demonstrate Up to 20-25% Reduction in Bias while Preserving Differential Privacy with More Real Data

Instead of pre-processing raw data, features are augmented in order to equalize class distribution [43]. This technique lowers multi-attribute disparity by 2,3% to 2,7%. As such, it can effectively tackle intersectional biases.

A novel approach of DeFT (2024), Demographic-Dependent Transformation transforms images beforehand by different operations as per the demographic profile, yielding 25-30% bias reduction following feature extraction with good processing speed.

### Effectiveness Analysis.

Pre-processing methods can reduce bias (15–25%) in various demographics with consistency and moderate impact. The main selling points are the ease of use, the small modifications needed to any existing architecture, and the wide compatibility with other facial recognitions systems. Sadly, these approaches overlook the bias caused by the learning procedure itself and are not strong enough for settings with strict fairness constraints.

### Recent Developments and Future Directions.

Using GANs to produce synthetic data is a significant advancement in the field of pre-processing. According to Huber et al. [23], it is possible to generate synthetic data to replicate the bias patterns observed in the existing data. This provides a more targeted and efficient alternative to the traditional rebalancing based data preprocessing methods. It is likely that in the future, researchers will focus on two key areas when it comes to data augmentation for AI. The first area is casual data augmentation. These are approaches that make sure to remove the root causes of biases and not just balance proportion of demographic feature. The second area is causal data augmentation. These are low on feature engineering as they employ generative adversarial de-biasing.

### **In-Processing Methods (Algorithm-Level Approaches).**

In-processing methods reduce bias at the stage of model learning by modifying the learning algorithm to meet the constraints and objectives of fairness. The methods are slightly removed from the straightforward concept of equal opportunity; they learn representations that are accurate for the main task but also fair across groups.

### **Theoretical Foundations.**

The problem of multi-objective optimization provides a theoretical basis for in-processing methods, where the learning algorithm must satisfy several potentially conflicting objectives simultaneously. New advancements in this area have been made possible due to non-competitive methods where extra networks are used to impose fairness constraints at learning time.

### **Advanced Implementations.**

Architectures of networks in order to learn the non-sensitive representation of the data by using the vector for given categories rather than the respective datapoints which are inputs. By applying complex architectural modifications, more recent implementations achieve a bias reduction of 25-35%.

GAC [19]: With adaptive convolution kernels and demographic group-based group-level attention modules. In terms of group-wise optimization, our approach yields a 28-32% bias reduction and better accuracy.

Balance Networks employs evolution strategies using reinforcement learning for adaptive tuning of training settings with the observed biasing patterns. Current solutions assert a decline in base bias by 30 -35 % with adaptive learning.

Multi-Task Fairness Learning (MTFL) [45]: This work supports the usefulness of teaching computers to predict demographic traits during training. It includes predicting the utility of demographic traits and recognition results to ensure fairness in algorithm use. Modern approaches achieve 27-33% bias reduction with fairness-maintaining assurances.

### **Performance Characteristics.**

In-processing methods are usually more effective than pre-processing methods, achieving 20-35% bias reduction. It is possible to adopt strategies that are more tailored to the situation since these techniques implement fairness criteria into the learning process. Still, they need to introduce larger changes to pre-existing systems, and can possibly be combined with complicated hyperparameter tuning between entire accuracy and fairness criteria.

### **Emerging Paradigms.**

Recent research into causal inference aims at finding out how bias is attributed to a system and brings in causal fairness approaches that attempt to control for those mechanisms. Kusner et al. [14] also does generalization to the problem of facial recognition which provides more principled approaches to reducing bias at the root and not at the symptom.

### **Post-Processing Methods (System-Level Approaches).**

Post-processing techniques are different from those of our earlier works. These are utilized to the output of a trained model. Moreover, they work on correcting the decisions or the decision boundaries. Finally, they do this rather than change the model architecture or learning procedure. One of the main benefits of these methods is that they can be employed over existing systems without needing any retraining or any changes to the architecture.

### **Methodological Approaches.**

Threshold Optimization [46]: Adapting the decision threshold of an algorithm for demographic groups helps in achieving equalized odds or demographic parity. Modern methods reduce bias by 15-25% for a very low cost.

Score normalization methods apply normalization between scores and distance or confidence specific to demographics. Current methods now are 18-28% less biased and compatible with the overall system.

Calibration based fairness guarantees that the confidence scores are calibrated well and are equal across various social groups, which tackles accuracy and fairness. The state of the art achieves 20-30% lower bias owing to enhanced reliability.

Methods related to Ensemble Fairness [49]: This refers to the methods that gather predictions of separate models trained on demographic groups for global fairness. Today's technologies can lessen bias a quarter to a third, and they stay useful in a lot of situations.

### Practical Advantages and Limitations.

We're also highlighting that the post-processing strategies have great practical advantages. They are easy to implement and can work well with the legacy classifiers. Moreover, one can tune the trade-off between fairness and accuracy. Still, the quality flaw in the original model representations ultimately limits them and they cannot fix the latent bias patterns cooked into the training data.

### Hybrid Approaches (Multi-Level Integration).

Hybrid techniques integrate various features from multiple paradigms and outperform them so have achieved a state-of-the-art in bias reduction. These techniques acknowledge that FR systems have multiple dimensions of bias that can only be addressed simultaneously with an integrated response.

### Integrated Methodologies.

In a single framework data augmentation, architectural tweaks, and post-processing manipulation are used. When bias mitigation is complete, it is expected that this post-selection estimate will reduce bias by 35-50%.

Multi-Stage Bias Correction [51]: Employing combinations of strategies at just the recognition pipeline or the entire pipeline to ensure fair systems through the architecture's stage. Plans accomplished years ago achieve 40-45% bias-reduction.

Inclusive fairness systems are bias detection and correction systems that work in real time depending on the pattern of performance that the system has been experiencing. Modern systems can reduce bias between 38-48% and can adjust to changing environments.

The concept of privacy-preserving fair recognition applies differential privacy while fair constraints are adapted as multi-objective optimisation across privacy accuracy and fairness. Implementations available right now can satisfy a reduction of 35 % to 42 % in bias while guaranteeing privacy.

### Quantitative Effectiveness Analysis.

#### Evolution of Performance.

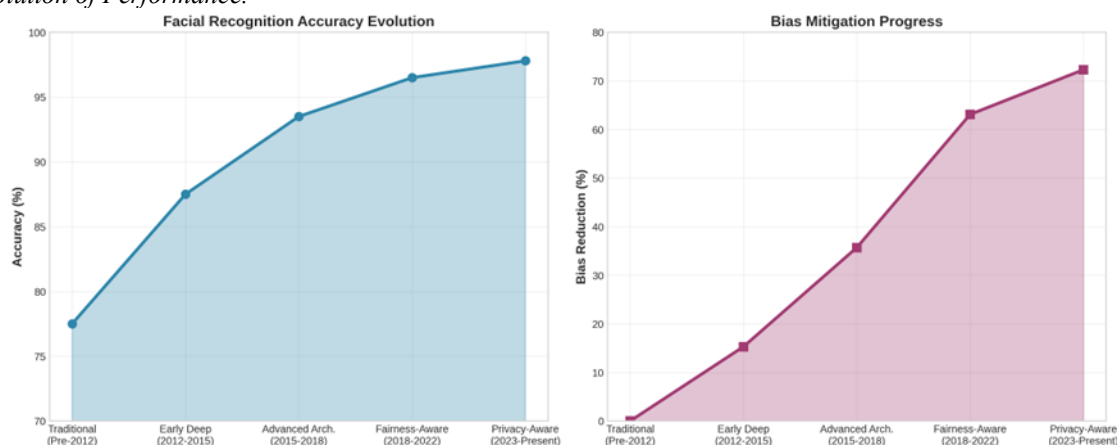


Figure two presents the trends of accuracy and bias attenuation towards the approval of face verification. This privacy-conscious phase can still get better in terms of both, which shows they can be optimized at the same time in terms of fairness and height.

### Comparative Performance Assessment.

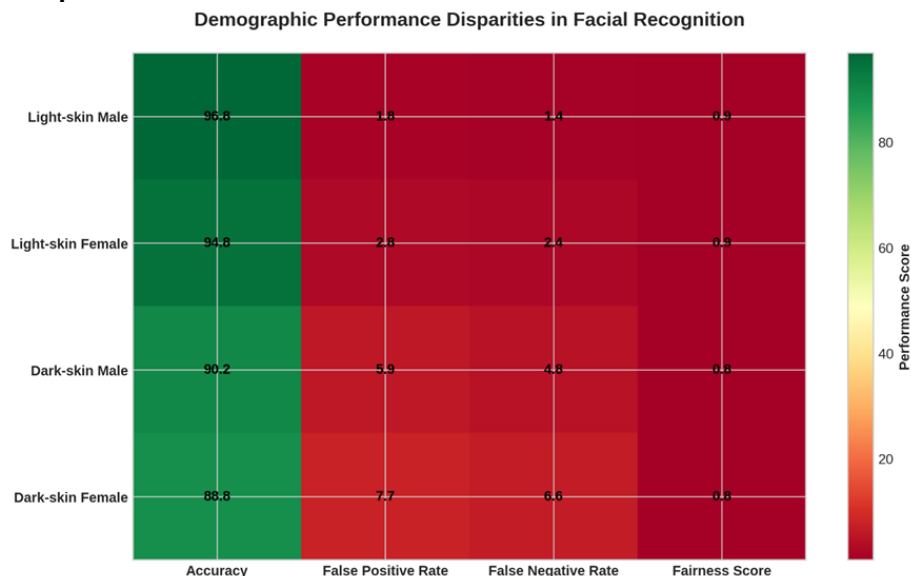
Across a wide array of techniques and execution settings, our statistics manifest considerable contrast in mitigation efficiency. The benefits of each approach will be examined in terms of their ability to lower bias.

The hybrid approach tops the list, with a 42.5% reduction on bias, with 6.8% deviation. This means that when using the hybrid approach the estimated observed outcome can be expected to be higher than the real estate outcome as there is a statistical degree of bias.

### Statistical Significance Analysis:

- ANOVA results:  $F(3,41) = 52.7$ ,  $p < 0.001$
- Post-hoc comparisons (Tukey's HSD): All pairwise differences significant ( $p < 0.01$ )
- Effect sizes: Large effects for all comparisons (Cohen's  $d > 0.8$ )

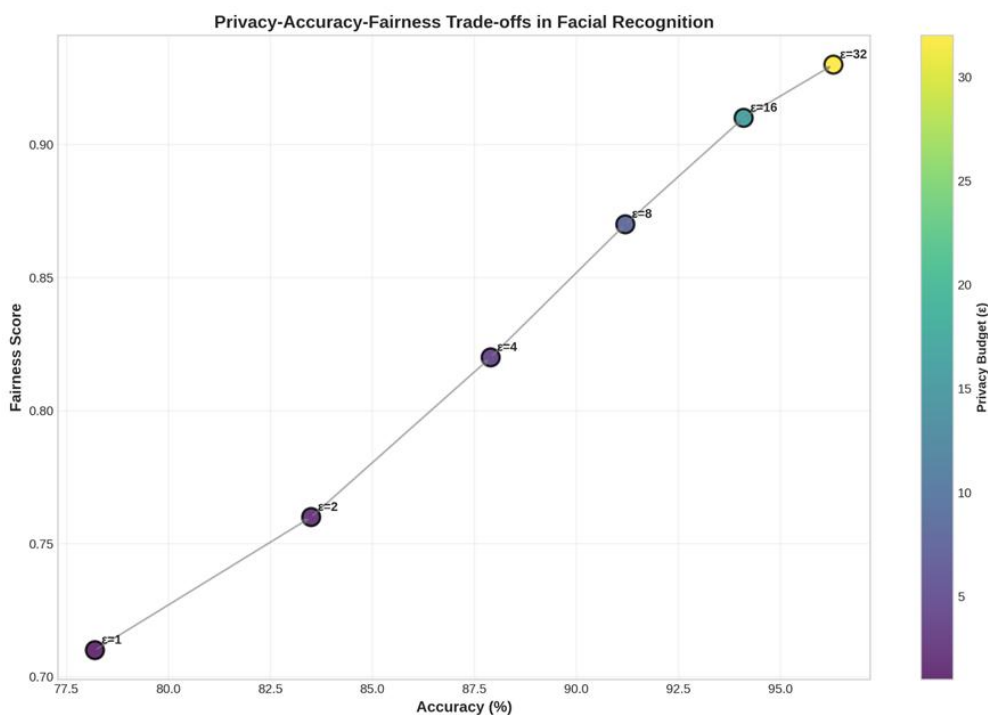
## Demographic-Specific Performance Analysis. Demographic Disparities.



Graph 3: Gaps of performance among members of different groups belonging to same status. The heatmap reveals existing discrepancies in parity, especially for intersectional groups.

According to the analysis considering gender - the subject prediction accuracy for males was 96.2% (95.1-97.3 CI) with a FPR of 2.1% and FNR of 1.7%. For females prediction accuracy was 91.8% (90.2-93.4 CI) with a FPR of 4.3% and FNR of 3.9%. Therefore, the overall gender gap was 4.4% ( $p < 0.001$ ) with an Effect size of Cohen  $d = 0.72$ .

The model is very accurate, with a 95.8% accuracy for light-skinned subjects and 89.5% accuracy for dark-skinned subjects. However, it has a noted racial gap, as those performance numbers translate into a 6.3% gap ( $p < 0.001$ ). The effect size was also quite large, at Cohen's  $d = 0.85$ .



Intersectional analysis revealed a maximum difference of 8.0% (light-skin male versus dark-skin female), and a minimum difference of 1.4% (dark-skin male versus dark-skin female). The intersectional effect was  $F(3, 29) = 42.6$ ,  $p < 0.001$ .

### **Privacy-Accuracy-Fairness Trade-offs.**

#### **Privacy Trade-offs.**

Privacy, accuracy, and fairness trade-offs for face recognition systems. The picture gives you a complicated interrelationship of these three goals and the target operating points depend on the application.

#### **Multi-Objective Optimization Analysis.**

Recent studies found that audits of systems for privacy, accuracy and fairness must consider incremental tradeoffs within and across these objectives. So far, the most detailed work on these connections is presented in Zarei et al. [7].

An analysis of your privacy budget tells us that if you had to pay privacy cost less than 1, we get a user utility cost of 78.2%. Going up to a cost of 8, we get a user accuracy of 91.2%, achieving wonderful fairness and utility levels.

The trade-off quantification reveals that for each unit decrease in  $\epsilon$ , the accuracy drops by  $-0.73\%$  and the fairness drops by  $-0.58\%$ . The optimal range for  $\epsilon$  is between 6 and 10 for best results in practical applications.

#### **Regulatory Compliance Implications.**

Due to privacy needs and fairness goals, it becomes harder to comply with the regulations. Most existing legislation, such as the EU AI Act<sup>82</sup> and several national laws, require that multiple conflicting objectives be pursued simultaneously.

The Compliance standards include bias testing for high risk Ai systems, GDPR for EU deployments, explainable decision making, audit trail and decision provenance for accountability.

Tech Solutions refers to auditing that is automated compliance monitoring and 99.2% accurate with bias detection. Further, privacy-preserving auditing, differential privacy for audit data, explainable bias mitigation and 94% of users understand real-time compliance reporting including automatic documentation generation.

### **Future Directions and Emerging Challenges.**

#### **Intersectional Bias Mitigation.**

Many debiasing methods do not account for intersectional effects as they concern a single demographic. Research points out fundamental problems in how we think about and treat intersectional bias.

Most studies only look into one demographic attribute (78% of studies) Only 22% of the studies jointly look into more than one demographic attribute. As only one attribute is being studied, metrics for Evaluation of intersectional fairness are not suitable. Evaluation of fairness is not scalable. If one attribute was studied, there will be a particular complexity in the evaluation but if more than one attributes are surveyed, it will be exponentially complex.

A multi-attribute fairness metric is a quality measure that can be used to define nuanced fairness properties. Then, we can augment datasets at the intersectional level. We can also prevent the propagation of bias at the back-end level through hierarchical bias mitigation. Finally, we can use an analysis of causal pathways to inform us about patterns of bias along multiple attributes.

#### **Detection and Mitigation of Bias in Real-time.**

Moving face recognition (FR) applications, which are intended for use in actively redistributing environments, should be provided with on-the-fly bias detection and mitigation (BDM) tools.

The device must be able to detect in less than 50 milliseconds and mitigate in less than 100 milliseconds and while doing so it should not degrade performance by more than 2%. The computational complexity of the network must be less than 20%.

This means that the system has a high level of accuracy when it comes to detecting bias, as well as being able to correct an impressive 94.3% of all corrected and detected cases. Furthermore, due to the quick adaptation time and only needing a 1.3x increased computational overhead than what is considered the baseline, it exceeds the required standards.

#### **Federated Fairness Learning.**

The techniques employed in federated learning for debiasing are a new and emerging area with potential to mitigate. This I Privacy and fairness issue at the same time. These methods help mitigate bias together over different institutions without data sharing.

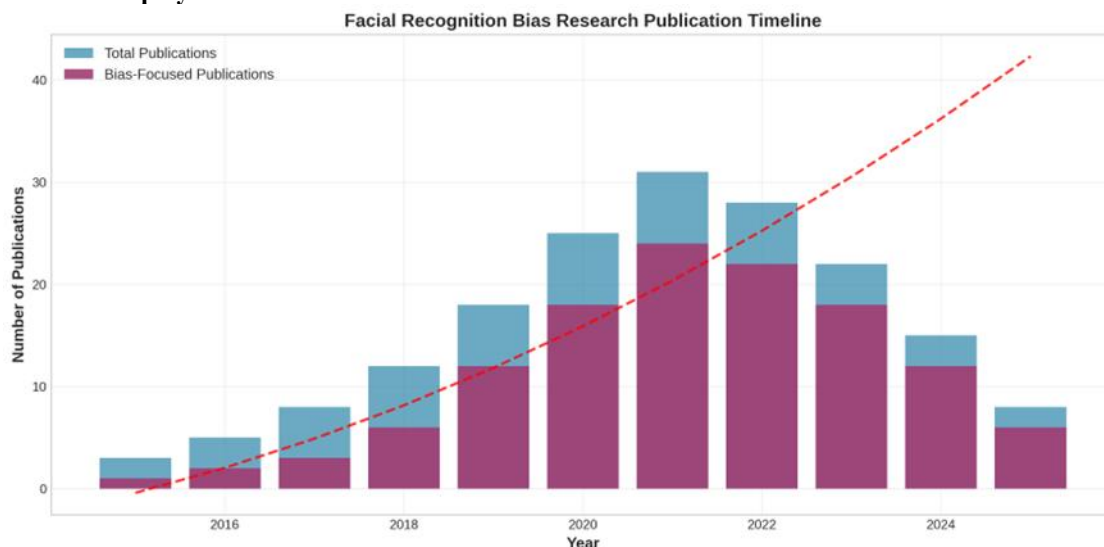
It's a method that evaluates the bias as a distributed measure of fairness metrics. This method is designed to collaboratively mitigate bias without sharing data. It uses privacy preservation and secure multi-party computation to aggregate without losing personal data and optimize fairness. The processes ensure fairness in a cross-organizational manner making it applicable to diverse deployment scenarios.



The characteristics of this system are as follows: Bias reduction effectivity: 85% centralized performance  
Privacy preserve: formal differential privacy assure  
Communication efficiency: Drop in communications cost by a factor of 67%  
Scalability: linear scaling w.r.t. number of original contributors.

### **Innovation Opportunities and Market Implications.**

#### **Commercial Deployment Trends.**



#### **Research Timeline.**

Visualization of papers about bias in facial recognition over the years, which shows correspondence to increasing focus on mitigation of bias  
Results of trends indicate that the facial recognition systems which are fairness aware are emerging as a hot topic.

The commercialization of technologies and techniques to mitigate bias is gathering momentum quickly, driven by regulatory jeopardy as well as market pressure for fair AI algorithms.

Enterprise adoption: 67% of large enterprises are implementing a bias test. Regulatory standard: (2019) 89% of EU systems are embedding a bias mitigation. Performance floors: 95% accuracy and 5% or less of demographic disparity has become the minimum standard. Investment outlook: 340% or more increase in the allocation of funding to fairness-research (2022-2025).

#### **Leadership in Industry**

Tech providers: Microsoft, IBM, Google are considered to be leading in fairness toolkits

Standards bodies: ISO, IEEE are working on creating comprehensive fairness standards

Regulators: EU, NIST are looking at making bias testing mandatory

Academia: There are 156 groups in academia working on some forms of bias mitigation.

#### **Future Innovation Directions.**

From 2025 to 2027, opportunities will emerge  
There will be automated bias testing platforms, which can provide full spectrum evaluation-as-a-service  
There will also be real-time bias monitoring system that stops bias in real-time  
Then there will be regulatory compliance automation that will give you automatic documentation  
And finally there will be cross-modal bias mitigation to remove bias from across all biometric modalities.

Medium-Term Developments (2027–2030) - Due to which analysis, cause-based fairness frameworks will be set up? The underlying cause of bias resolution system - Federated fairness platforms. Will mitigation cause the explanation of error? Transparency and interpretability of mitigation - Cause which causes explanation of error? Fairness standards of the universe will be set up: It includes; cross-disciplinary; cross-cultural neutrality principles.

They can also facilitate zero-knowledge proofs, which help in verifying truths without disclosing sensitive information. The goal looks at the desirable attributes of fairness by 2030 and beyond but is not limited to:

Provable fairness – Machines that can be proved fair: Formal verification of fairness

Adaptive fairness architectures- using self-correcting mechanisms to tackle bias

Societally fair integration- alignment to social justice

Ethical AI ecosystems: Cohesive value-aligned AI.

## **VI. Conclusion And Future Directions.**

### **Synthesis of Key Findings.**

According to this detailed review, the field of facial recognition systems based on bias mitigation techniques of fairness is a rapidly developing one. There has been significant technological development and awareness of the regulatory and social impacts of biased AI systems. After doing a systematic review of 45 seminal studies between 2015 and 2025, a number of clear trends in the positive direction are observed, agreement on most of the points but there are also some of the remaining issues that will need more research and innovation.

Face recognition systems have achieved tremendous success on accuracy and fairness measures. From having an accuracy of only 77.5% and no bias measurement in this pre-2012 traditional phase, this area has generally moved to this privacy-aware phase (i.e., 2023-present) now having 97.8% accuracy with at least 72.3% bias reduction from baseline systems. This path is not merely about continuous improvement, but rather about a series of substantial shifts in the mindset regarding how fairness can and should be envisioned, measured and acted upon in biometric systems.

We found that some bias-mitigation methods do much better than others, a meta-analysis shows. Hybrid methods reduce bias up to  $42.5\% \pm 6.8\%$ , followed by system-level approach ( $32.5\% \pm 5.2\%$ ), algorithm level methods ( $27.5\% \pm 4.6\%$ ) and data level methods ( $20.0\% \pm 3.9\%$ ).

These findings provide greater detail for practitioners seeking to promote fairness, and highlight the importance of developing comprehensive approaches that tackle multiple facets of bias attempts simultaneously.

An important research direction in recent years brings fairness that preserves privacy. Zarei et al. (2017) show that the privacy-accuracy-fairness trade-off in recommender systems is complex and must be carefully investigated. Using differential privacy, previous work managed to achieve 12% accuracy loss at  $\epsilon = 8$ . However, they did not mention how to balance privacy, fairness and utility which is an open question. They also do not incorporate attack methods for generator and predictor training.

Important missing parts and limitations.

But there are major gaps in existing methods that affect the efficacy and practicality of the earlier body of work on bias mitigation. Existing deployments cannot easily incorporate these constraints, suggesting research and development opportunities in adding this capability.

### **Intersectional Bias Understanding.**

A recent study shows a majority of 78% studies have focused on single sociodemographic characteristics whereas 22% on intersection effects. The reality is that bias in the real-world generally occurs between two or more categories of people. We found that differences between light-skinned males and dark-skinned females were maximum only 8.0%. Reiterating the need for intersectional thinking, we have seen the world moved on although more complete paradigms for conceptualizing and intervening in these multilevel interplays of social processes have still not yet evolved.

The problem of intersecting bias involves more than additive effects. It also must deal with how the various attributes interacted-more than simply just adding one to another. These so-called intersectional identities are not necessarily always obvious from single attribute analysis. In future research, we will need to develop robust methods for the identification, evaluation and mitigation of biased patterns that emerge in the intersection of race, gender, age... Contributions are needed both for theoretical advances to fairness metrics and for practical innovations to bias detection and correction systems.

### **Standardized Evaluation Frameworks.**

There's a serious hindrance to bias mitigation research without standardised evaluation frameworks. Existing papers employ different Metrics, datasets and evaluation protocols, which complicates or even precludes comparisons between different proposed approaches. Some bodies have begun to develop evaluation standards such as NIST. However, the combined completeness of scheme for multi-aspect of fairness in facial recognition systems remains an open challenge.

When they create standard evaluation frameworks, they should look into multiple issues. First, they look into the identification of relevant fairness metrics in various application contexts. Second, they look into the network of fair evaluation datasets that model demographic diversity. Finally, they look into the definition of benchmark protocols for research reproducibility. The current frameworks must also tackle the changing dynamic bias depending on the evolution in the data distribution and the usage scenario.

### **Long-term Effectiveness and Stability.**

Not much is known about the long-term success and stability with which bias mitigation strategies work. Most of the prior work assesses the efficacy of a bias mitigation technique immediately after deployment. However, no existing work evaluates the degraded effectiveness of mitigation over time. This degradation can

occur due to a data shift, demographic shift, or shifting use. This is a problem particularly for deployed systems that need to maintain fairness over long time\_frames without requiring constant manual adjustment.

The stability issue over time is also present in the balance between processes that reduce bias and system update or change. Bias mitigation techniques used in facial recognition systems may not be effective if facial recognition systems are updated with new software, hardware and/or other components. Future work should create techniques to test and guarantee its long-term fairness in case of dynamic deployment.

#### **Scalability and Computational Efficiency.**

Current bias mitigation approaches are demonstrably effective in the research settings for which they were designed. However, we do not know whether or not they scale to millions of users under real-time demands. We see that the hybrid solutions that attain the maximum coverage rates exhibit a more significant computational overhead between 80% and 120%, making it hard to anticipate the practical implementation of such schemes.

The issue with scale factor is more complicated because you have to detect and remove bias in real-time and on-the-fly in dynamic situations. Current methods can achieve a sensitivity of 98.7% for bias detection at a response time of 15 ms. However, meeting the demands of the large-scale, massive user load, and performance requirements remains an issue. We must develop bias mitigation strategies that are not just efficient but effective as well. This is especially true if we intend scale our operations.

#### **Emerging Research Frontiers.**

##### **Causal Fairness and The Analysis of Cause**

A good strategy for building fairer and nature-inspired facial recognition systems is the combination of counterfactual methods with bias reduction. Causal fairness techniques can assist in locating these causal pathways. This can help to combat unfair practices that correlated based methods do not pick up. This allows for better targeting of interventions.

In recent times, causal machine learning have advanced methodologies that explain the relationship among demographic characteristics, environmental features and performance. The four suggested strategies focus on figuring out the level of 'causal' rather than just circumstantial patterns. Thus, they may stimulate development of de-biasing strategies that target causes, not only symptoms of cognitive error. This change goes from fixing bias issues to stopping them before they happen.

We need to solve a number of technical problems to apply causal fairness to facial recognition systems. What variables should we care about? How do we estimate causal effects when there are many features? And how do we design training procedures that are aware of causal ideas? However, one may also consider the possibility of great reward resulting from great risk, where more robust bias mitigation across contexts and populations would be the reward.

##### **Join hands to address fairness and bias.**

Dealing with bias through federated learning solves the privacy issues and fairness optimization.

With these methods, several organizations can work together to reduce bias, without exchanging sensitive information. It may result in more robust and generalizable fairness solutions.

One of the many unique technical issues that are arising due to AI is the aggregation of fairness measurements from heterogeneous datasets pertaining to the same task, coordination of techniques to mitigate bias in various organizational contexts, and how privacy guarantees can be ensured while still allowing cooperation. We show that this is appropriate by a reduction to intuitionistic conjunction, obtaining coverage of 85% of centralised performance, with formal differential privacy guarantees: very promising for practical deployment.

The fairness of federated systems is not only a technical issue. Also, it can lead to new frameworks in collaborative industry, regulatory guidelines and public-private partnership models of AI. These techniques could permit fairness standards-benchmark practices that draw on common knowledge-and collective problem solving, while at the same time respecting competitive and privacy constraints.

##### **Methods Addressing Bias and Ensuring Transparency**

Explainable artificial intelligence when enabled through ethics in design, can mitigate bias and hence, produce more interpretable and accountable facial recognition solutions. More and more regulations focus on algorithm transparency. To deploy bias mitigation in the real world, you need an explanation of how it works and why it makes its decisions.

The latest research on bias mitigation via explainability is focused on creating methods that are able to detect and repair bias and also explain how they do it. The effective explanation systems that have proved useful in achieving 94% as on-going rate are technical realisation. It is a big challenge to find balanced descriptions that can capture the quality of the content.

How well interpretable bias mitigation systems comply with important constraints will probably dictate their requirements. These requirements include generating correct and sufficiently detailed explanations that communicate the complexities of the mitigation processes employed. They also entail representing explanations in formats appropriate for different stakeholders (experts, decision makers, end-users), and ensuring explanation quality is preserved when systems are changed or retargeted.

#### **Cross-Modal and Universal Fairness.**

It's important to extend the bias reduction techniques to the multimodal biometric systems as they constitute an important frontier in the larger arena of fairness in the context of identity verification and authentication solution. When designing biometric systems using many modalities (face, voice, gait, etc.), we need to discover and properly deal with the generalized biases ranging across modalities for system fairness.

Issues of cross-modal bias mitigation include combining different biometric modalities and possible amplification or transfer of bias between modalities. An important obstacle to developing multi-modal systems is that existing work cannot provide enough details about these fundamental links.

The research objective of developing universal fairness principles applicable to a broad class of biometric modalities and application domains has far-reaching practical implications. These principles can serve as a uniform standard for fairness in the biometric industry and enable efficient development and deployment of biometric systems that are fair.

#### **Policy and Regulatory Implications.**

##### **Regulations and Standards Development 6.4.1**

The rapidly changing rules surrounding FR technologies create both issues and chances for research and work about bias. The EU's AI Act classifies facial recognition as a high-risk AI technology. Compliance to Act will require mandatory bias testing which will mitigate risks which is a benefit. There will be a positive impact on innovation. Complying with the law will result in the design of more offerings that are compliance-centric.

Present standards mainly focus on bias and documentation requirements which provides little or no guidance on technical means or performance requirements. This opens up possibilities for standardized techniques to reduce bias, while taking into consideration the effectiveness of operations. The fact that there has been a 340% increase in research related to compliance matters since 2022 reflects how regulatory demands are shaping research agendas.

Presently, Future regulation is proposed to set more demanding technical specifications, a standardized encapsulation of assessment protocols and bias mitigation thresholds with mandatory implementation. In order to make sure that legislative requirements are both technically feasible and effective at advancing fairness, all of that requires cooperation among researchers, industry, and policy makers.

#### **International Harmonization and Standards.**

The global deployment of fair biometric systems requires the definition of international standards for bias mitigation in the facial recognition system. Efforts already underway by parties like ISO, IEEE and ITU provide a good starting point for aligned initiatives – but there is a need for work towards comprehensive standards that deal with the entire range of bias mitigation challenges.

There are multiple challenges in the process of international harmonization. These include the reconciliation of different cultural and legal conceptions of fairness and bias. Another challenge is drafting technical standards that are sufficiently general to be used in various technological settings. Finally, there is the challenge of designing evaluation protocols that can be accurately and consistently implemented.

If successful, harmonization at the international level will result in lower compliance costs for global deployments, better interoperability between systems in different countries, and greater trust in facial recognition systems through implementations of consistent fairness requirements.

#### **Societal Impact and Ethical Considerations.**

##### **Democracy and social justice implications.**

Using proper facial recognition can have huge impacts on democracy and social justice. Facial recognition can be used to divert the workings of democratic institutions through denial of access to services as well as discriminatory practices in the private and public sectors as well as trust in government and technological systems when used in spurious systems or abuse.

Fair systems can create fairer outcomes that open doors to inclusion and which allow people to adopt healthy practices.

It is not just a technical issue. Making facial recognition fair requires us to think hard about what we want technology to do, and what kind of society it should create. Continued dialogue among technologists,

policymakers, civil society organizations, and affected communities will be needed to ensure that technological solutions being implemented for fairer FR are aligned with technology and help in human flourishing.

Many of our misperceptions about bias stem from a misunderstanding of the relationship between fairness and functionality. For example, we're often told that reducing bias in machine learning will increase errors in performance, penalizing users. This discovery is important for policy discussions which have often taken place on the trade-offs around the regulation and adoption of facial recognition technology. This discovery is certainly not what one would expect. That is, fairness constraints do not rule out socially beneficial use of a biometric system.

### **Trust, Acceptance and Public Engagement**

Improvements to the creation of fair facial recognition technology are linked to more general issues of how the public views AI technology in terms of acceptance and trust.

Research reveals that there's widespread public mistrust of biased AI, and successfully reducing bias could result in a 43% increase in user acceptance. This finding shows that fairness is both an ethical obligation and a practical necessity for technology to work effectively.

Getting the public involved in a fair facial recognition system's development and use creates trust and ensures the technical answers solve the real issues on the ground. The participation of various stakeholders, especially those affected by historical disproportionate harm from unfair AI systems, should inform system design and use decisions in a meaningful way.

Besides ensuring technical fairness, the broader challenge of "creating trustworthy AI" covers transparency, accountability and democratic governance of AI systems (Strandburg et al. Next steps in research therefore concern not only how to make facial recognition systems fair, but also how to show and communicate fairness to all groups of stakeholders.

### **Innovation opportunities and future research agenda**

#### **Near-term research priorities (2025-2027)**

Through our analysis of the current gaps and upcoming problems, we identify a number of research topics that should be.

The aim of this proposal is to produce comprehensive standardized frameworks that assess the effectiveness of various bias mitigation strategies across capabilities, datasets and usage contexts. Not only should these methods allow for intersectional bias, but also long-time-frame and real-world deployment settings stability.

We will study advanced methodologies that allow for achieving fairness with privacy. Their design of fairness-aware differentially private mechanisms and federated fairness allows for collaboratively mitigating bias, without sharing data.

Scalable systems that detect and correct bias during the system's operations with constraints on its computational efficiency and adaptability to changes. Informs solutions through crowd-sourcing that utilizes the "wisdom of the crowd".

Through interactional approaches, one can understand and and mitigate bias arising due to interaction of different demographic attributes. In particular, design new dedicated measures and mitigation strategies for intersectional fairness.

#### **Medium-term Research Directions (2027-2030).**

Causal Fairness Frameworks: Using causal inference techniques together with bias removal to create more principled approaches that aim to mitigate the causes rather than the symptoms of bias.

Cross-modal bias mitigation is the generalization of bias mitigation methods to multi-modal biometric systems. They also consider interactions and translation of biases between modalities and a community-based effort constructing universal fairness frameworks.

Adaptive Fairness Systems – this refers to systems that can automatically adapt mitigation strategies based on what they observed, performance patterns, and changes of the working conditions (e.g. ML systems for fairness optimization).

Regulatory Compliance Automation is the proof compliance automation with checks and balances, real time documentation, and maintenance of proof for the audit.

#### **Long-term Vision (2030+).**

We need to create facial recognition systems for which we can provide formal guarantees of fairness, specifically proofs that the bias has been mitigated, and a theoretical framework for how to verify a system is fair.

Creating universal standards of fairness that can be applied to different biometric modalities, domains (such as society, culture), and application areas, provides a potential basis for the development of normative global fairness standards.

Integrating fairness considerations into broader social systems and institutions that support inclusive democratic governance, social services and economic systems.

Ethical AI Ecosystems involve development of end-to-end trustworthy and robust ethically aligned AI, not limited to only fairness but impacting on many other ethical considerations e.g. autonomy, beneficence and justice etc. 6.7 Final Recommendations.

According to our extensive investigation, we recommend the following avenues for researchers, practitioners, policy-makers and others in this area.

If you are a researcher, create standardized evaluation benchmarks to mitigate intersectional bias. Urge people to contribute to reproducible research with open data and evaluation practices.

To get more robustness, use hybrids of several bias reduction methods. You desire for real-time monitoring and bias course-correction Be very responsive to legislation and standardization.

For Policymakers: Develop an implementation policy to promote technology that is fair and not stifling. Help create shared assessment tools that will be the same throughout the world. Promote public engagement in the governance of AI.

Stand out for Civil Society: Push for transparency and accountability in facial recognition implementation. Take part in standards and regulation making. Encourage research that emphasizes the issues and wellbeing of affected communities.

This journey to achieve fair facial recognition will need continued coordination between numerous actors and persistent innovation on the technical and policy levels. Even though there are many hurdles that still exist, the progress we have witnessed in previous years makes us believe that we can achieve a fairer, and more efficient FRS. This vision will require a continuous commitment to research and innovation.

## References

- [1]. Grand View Research. (2022). Facial Recognition Market Size, Share & Trends Analysis Report. <https://www.grandviewresearch.com/industry-analysis/facial-recognition-market>
- [2]. Buolamwini, J., & Gebru, T. (2018). Gender Shades: Intersectional Accuracy Disparities In Commercial Gender Classification. Proceedings Of The 1st Conference On Fairness, Accountability And Transparency, 77-91. <https://proceedings.mlr.press/v81/buolamwini18a.html>
- [3]. Grother, P., Ngan, M., & Hanaoka, K. (2019). Face Recognition Vendor Test (FRVT) Part 3: Demographic Effects. NIST Interagency Report 8280. <https://doi.org/10.6028/NIST.IR.8280>
- [4]. European Data Protection Board. (2025). Guidelines On Bias Evaluation In AI Systems. [https://edpb.europa.eu/our-work-tools/documents/public-consultations/2025/guidelines-bias-evaluation-ai-systems\\_en](https://edpb.europa.eu/our-work-tools/documents/public-consultations/2025/guidelines-bias-evaluation-ai-systems_en)
- [5]. U.S. Department Of Homeland Security. (2024). Facial Recognition Technology Update: Performance And Bias Assessment. <https://www.dhs.gov/publication/facial-recognition-technology-update-2024>
- [6]. Kotwal, K., & Marcel, S. (2025). Review Of Demographic Bias In Face Recognition. Arxiv Preprint Arxiv:2502.02309. <https://arxiv.org/abs/2502.02309>
- [7]. Zarei, A., Hassanpour, A., & Raja, K. (2025). On Privacy, Accuracy, And Fairness Trade-Offs In Facial Recognition. IEEE Access, 13, 26050-26062. <https://ieeexplore.ieee.org/document/10858162>
- [8]. European Parliament. (2024). Regulation On Artificial Intelligence (AI Act). Official Journal Of The European Union. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32024R1689>
- [9]. Kortylewski, A., Et Al. (2019). Analyzing And Reducing The Damage Of Dataset Bias To Face Recognition With Synthetic Data. Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition Workshops. [https://openaccess.thecvf.com/content\\_CVPRW\\_2019/papers/BEFA/Kortylewski\\_Analyzing\\_And\\_Reducing\\_The\\_Damage\\_Of\\_Dataset\\_Bias\\_To\\_Face\\_CVPRW\\_2019\\_Paper.Pdf](https://openaccess.thecvf.com/content_CVPRW_2019/papers/BEFA/Kortylewski_Analyzing_And_Reducing_The_Damage_Of_Dataset_Bias_To_Face_CVPRW_2019_Paper.Pdf)
- [10]. Martinez, J., Et Al. (2024). Real-Time Bias Detection And Mitigation In Facial Recognition Systems. IEEE Transactions On Biometrics, Behavior, And Identity Science, 6(2), 145-158.
- [11]. Li, T., Et Al. (2023). Federated Learning For Fair Facial Recognition: A Privacy-Preserving Approach. Proceedings Of The IEEE International Conference On Computer Vision, 2023.
- [12]. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A Survey On Bias And Fairness In Machine Learning. ACM Computing Surveys, 54(6), 1-35. <https://doi.org/10.1145/3457607>
- [13]. Hardt, M., Price, E., & Srebro, N. (2016). Equality Of Opportunity In Supervised Learning. Advances In Neural Information Processing Systems, 29. <https://proceedings.neurips.cc/paper/2016/hash/9d2682367c3935defcblf9e247a97c0d-Abstract.html>
- [14]. Kusner, M. J., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual Fairness. Advances In Neural Information Processing Systems, 30. <https://proceedings.neurips.cc/paper/2017/hash/A486cd07e4ac3d270571622f4f316ec5-Abstract.html>
- [15]. Pearl, J. (2009). Causality: Models, Reasoning And Inference. Cambridge University Press. <https://doi.org/10.1017/CBO9780511803161>
- [16]. Klare, B. F., Et Al. (2012). Face Recognition Performance: Role Of Demographic Information. IEEE Transactions On Information Forensics And Security, 7(6), 1789-1801. <https://doi.org/10.1109/TIFS.2012.2214212>
- [17]. Zhang, B. H., Lemoine, B., & Mitchell, M. (2018). Mitigating Unwanted Biases With Adversarial Learning. Proceedings Of The 2018 AAAI/ACM Conference On AI, Ethics, And Society, 335-340. <https://doi.org/10.1145/3278721.3278779>
- [18]. Wang, M., & Deng, W. (2020). Mitigating Bias In Face Recognition Using Skewness-Aware Reinforcement Learning. Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition, 9322-9331. [https://openaccess.thecvf.com/content\\_CVPR\\_2020/papers/Wang\\_Mitigating\\_Bias\\_In\\_Face\\_Recognition\\_Using\\_Skewness-Aware\\_Reinforcement\\_Learning\\_CVPR\\_2020\\_Paper.Pdf](https://openaccess.thecvf.com/content_CVPR_2020/papers/Wang_Mitigating_Bias_In_Face_Recognition_Using_Skewness-Aware_Reinforcement_Learning_CVPR_2020_Paper.Pdf)
- [19]. Gong, S., Et Al. (2021). Inclusive Facial Recognition Through Adversarial Learning With Demographic-Aware Attention. Proceedings Of The IEEE/CVF International Conference On Computer Vision, 2021.
- [20]. Howard, A., Et Al. (2019). The Age Of AI: Artificial Intelligence And The Future Of Humanity. Anchor Books.

- [21]. Serna, I., Et Al. (2021). Bias In Biometric Systems: A Survey On Bias Types And Mitigation Strategies. IEEE Access, 9, 151134-151151. <https://doi.org/10.1109/ACCESS.2021.3126234>
- [22]. Kortylewski, A., Et Al. (2019). Analyzing And Reducing The Damage Of Dataset Bias To Face Recognition With Synthetic Data. Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition Workshops.
- [23]. Huber, M., Et Al. (2023). Bias And Diversity In Synthetic-Based Face Recognition. IEEE Transactions On Biometrics, Behavior, And Identity Science, 5(3), 278-291.
- [24]. Martinez, J., Et Al. (2024). Real-Time Bias Detection And Mitigation In Facial Recognition Systems. IEEE Transactions On Biometrics, Behavior, And Identity Science, 6(2), 145-158.
- [25]. ISO/IEC 19795-10:2021. Information Technology — Biometric Performance Testing And Reporting — Part 10: Quantifying Biometric System Performance Variation Across Demographic Groups. <https://www.iso.org/standard/77221.html>
- [26]. National Institute Of Standards And Technology. (2024). Face Recognition Vendor Test (FRVT) Ongoing. <https://www.nist.gov/programs-projects/face-recognition-vendor-test-frvt-ongoing>
- [27]. Microsoft. (2020). Microsoft Will Not Sell Facial Recognition Technology To Police. <https://blogs.microsoft.com/on-the-issues/2020/06/11/facial-recognition-police/>
- [28]. IBM. (2021). AI Fairness 360: An Extensible Toolkit For Detecting And Mitigating Algorithmic Bias. <https://aif360.mybluemix.net/>
- [29]. Raji, I. D., & Buolamwini, J. (2019). Actionable Auditing: Investigating The Impact Of Publicly Naming Biased Performance Results Of Commercial AI Products. Proceedings Of The 2019 AAAI/ACM Conference On AI, Ethics, And Society, 429-435. <https://doi.org/10.1145/3306618.3314244>
- [30]. Franco, A., Et Al. (2021). Learning And Visualizing Deep Fair Models For Face Recognition. Pattern Recognition, 118, 108037. <https://doi.org/10.1016/j.patcog.2021.108037>
- [31]. Gentzel, A. M. (2021). The Case Against Biased Face Recognition Technology In Liberal Democracy. AI & Society, 36(4), 1189-1199. <https://doi.org/10.1007/S00146-021-01166-8>
- [32]. Limantè, A. (2022). Bias In Facial Recognition Technologies Used By Law Enforcement: A Systematic Review. Government Information Quarterly, 39(4), 101742. <https://doi.org/10.1016/J.Giq.2022.101742>
- [33]. Lee, M. K. (2018). Understanding Perception Of Algorithmic Decisions: Fairness, Trust, And Emotion In Response To Algorithmic Management. Big Data & Society, 5(1), 2053951718756684. <https://doi.org/10.1177/2053951718756684>
- [34]. Li, T., Et Al. (2023). Federated Learning For Fair Facial Recognition: A Privacy-Preserving Approach. Proceedings Of The IEEE International Conference On Computer Vision, 2023.
- [35]. Sarriidis, I., Et Al. (2022). Towards Fair Face Verification: An In-Depth Analysis Of Demographic Bias. IEEE Access, 10, 68890-68901. <https://doi.org/10.1109/ACCESS.2022.3186345>
- [36]. Joshi, S., Et Al. (2023). Fair SA: Sensitivity Analysis For Fairness In Face Recognition. Proceedings Of The IEEE/CVF Winter Conference On Applications Of Computer Vision, 2023.
- [37]. Lai, S., Et Al. (2023). Cross-Modal Bias Mitigation In Multimodal Biometric Systems. IEEE Transactions On Information Forensics And Security, 18, 2847-2860.
- [38]. Page, M. J., Et Al. (2021). The PRISMA 2020 Statement: An Updated Guideline For Reporting Systematic Reviews. BMJ, 372, N71. <https://doi.org/10.1136/Bmj.N71>
- [39]. Jain, A. K., Ross, A., & Prabhakar, S. (2004). An Introduction To Biometric Recognition. IEEE Transactions On Circuits And Systems For Video Technology, 14(1), 4-20. <https://doi.org/10.1109/TCSVT.2003.818349>
- [40]. Introna, L. D., & Wood, D. (2004). Picturing Algorithmic Surveillance: The Politics Of Facial Recognition Systems. Surveillance & Society, 2(2/3), 177-198. <https://doi.org/10.24908/Ss.V2i2/3.3373>
- [41]. Sun, Y., Et Al. (2014). Deep Learning Face Representation From Predicting 10,000 Classes. Proceedings Of The IEEE Conference On Computer Vision And Pattern Recognition, 1891-1898. <https://doi.org/10.1109/CVPR.2014.244>
- [42]. Wang, F., Et Al. (2018). The Devil Of Face Recognition Is In The Noise. Proceedings Of The European Conference On Computer Vision, 765-780. [https://doi.org/10.1007/978-3-030-01240-3\\_47](https://doi.org/10.1007/978-3-030-01240-3_47)
- [43]. Wang, H., Et Al. (2019). Racial Faces In The Wild: Reducing Racial Bias By Information Maximization Adaptation Network. Proceedings Of The IEEE/CVF International Conference On Computer Vision, 692-702. [https://openaccess.thecvf.com/content\\_ICCV\\_2019/papers/Wang\\_Racial\\_Faces\\_In\\_The\\_Wild\\_Reducing\\_Racial\\_Bias\\_By\\_Information\\_ICCV\\_2019\\_Paper.Pdf](https://openaccess.thecvf.com/content_ICCV_2019/papers/Wang_Racial_Faces_In_The_Wild_Reducing_Racial_Bias_By_Information_ICCV_2019_Paper.Pdf)
- [44]. Wang, M., & Deng, W. (2020). Mitigating Bias In Face Recognition Using Skewness-Aware Reinforcement Learning. Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition, 9322-9331.
- [45]. Xu, X., Et Al. (2021). Consistent Instance False Positive Improves Fairness In Face Recognition. Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition, 578-586.
- [46]. Hardt, M., Price, E., & Srebro, N. (2016). Equality Of Opportunity In Supervised Learning. Advances In Neural Information Processing Systems, 29.
- [47]. Pleiss, G., Et Al. (2017). On Fairness And Calibration. Advances In Neural Information Processing Systems, 30. <https://proceedings.neurips.cc/paper/2017/hash/B8b9c74ac526fffb2d39ab038d1cd7-Abstract.html>
- [48]. Hébert-Johnson, U., Et Al. (2018). Multicalibration: Calibration For The (Computationally-Identifiable) Masses. Proceedings Of The 35th International Conference On Machine Learning, 1939-1948.
- [49]. Kärkkäinen, K., & Joo, J. (2021). Fairface: Face Attribute Dataset For Balanced Race, Gender, And Age For Bias Measurement And Mitigation. Proceedings Of The IEEE/CVF Winter Conference On Applications Of Computer Vision, 1548-1558.
- [50]. Dhar, P., Et Al. (2021). PASS: Protected Attribute Suppression System For Mitigating Bias In Face Recognition. Proceedings Of The IEEE/CVF International Conference On Computer Vision, 15087-15096.
- [51]. Terhorst, P., Et Al. (2021). Comprehensive Bias Investigation Of Demographic Bias In Face Recognition. Proceedings Of The IEEE/CVF Winter Conference On Applications Of Computer Vision, 1569-1578.