

## Amazon's AI Recruiting Tool Bias

The primary source of bias in Amazon's AI recruiting tool was biased training data. The model was trained on resumes submitted to Amazon over a 10-year period—most of which came from male applicants, reflecting historical gender imbalances in tech hiring. As a result, the AI learned to penalize resumes that included words like "women's," or that came from all-women's colleges, equating them with lower hiring suitability.

### Three Fixes to Make the Tool Fairer

- **Debias the Training Data**

Use a more diverse and balanced dataset with equal representation across gender, ethnicity, and educational background.

Remove or neutralize features that serve as proxies for gender (e.g., names, gendered terms, college names).

- **Use Fairness-Conscious Algorithms**

Apply fairness-aware machine learning methods such as adversarial debiasing or reweighing to reduce dependency on sensitive attributes.

Incorporate constraints during training to ensure equal opportunity or demographic parity.

- **Human-in-the-Loop Oversight**

Combine AI screening with human review, especially for underrepresented applicants.

Train HR personnel to audit and interpret AI recommendations, not follow them blindly.

### Fairness Evaluation Metrics (Post-Correction)

- **Demographic Parity:**  
Measures whether the selection rate is the same across different groups (e.g., male vs. female applicants).
- **Equal Opportunity:**  
Ensures true positive rates (qualified candidates correctly selected) are equal across groups.
- **Disparate Impact Ratio:**  
Compares the selection rates of protected and unprotected groups. A ratio  $< 0.8$  typically signals adverse impact under U.S. EEOC guidelines.
- **Calibration by Group:**  
Assesses whether predictions (e.g., scores for interview suitability) are equally reliable across demographic groups.

# Facial Recognition Misidentifies Minorities

## Ethical Risks

- **Wrongful Arrests & Legal Consequences**  
Facial recognition systems have been shown to have higher error rates for people of color—especially Black and Asian individuals—due to biased training data. This can lead to misidentification, resulting in unjust police actions, false arrests, and emotional or reputational damage.
- **Discrimination and Inequity**  
Disproportionate error rates perpetuate systemic biases, reinforcing racial profiling and deepening mistrust between law enforcement and marginalized communities.
- **Privacy Violations**  
Mass surveillance through facial recognition, especially without consent, raises serious concerns about loss of anonymity, constant tracking, and the chilling effect on free expression and movement.
- **Lack of Transparency and Accountability**  
Many facial recognition algorithms are proprietary “black boxes,” leaving affected individuals without recourse or understanding of how they were flagged, undermining trust in the system.

## Policy Recommendations for Responsible Deployment

- **Mandatory Bias Auditing**  
Require regular independent audits of accuracy across different demographic groups (race, age, gender), and prohibit deployment if performance is unequal.
- **Transparency Requirements**  
Vendors must disclose model performance, datasets used, and limitations. End users (e.g., police, agencies) should be transparent with the public when deploying such tools.
- **Informed Consent & Public Notice**  
Use facial recognition only in clearly marked, opt-in settings—never in public surveillance without strict legal oversight.
- **Human Oversight Mandate**  
Prohibit any use of facial recognition for automated decision-making. Always require human review before action is taken based on AI-generated matches.
- **Strict Use Cases & Moratoria**  
Limit use to only critical scenarios (e.g., finding missing persons). Consider temporary bans until fairness and privacy safeguards are fully enforced.