



Ethical AI in Production: Bias Analysis & Fairness in Deployed Models

Examining bias, fairness, and mitigation strategies for machine learning models deployed in production environments

The Deployment Scenario

Our Random Forest model predicts issue priorities (High/Medium/Low) for resource allocation in software development. Originally trained on adapted breast cancer data, it now classifies thousands of engineering tickets daily.

While achieving 91% overall accuracy, initial analysis revealed concerning disparities across teams, geographic regions, and developer experience levels. This presentation examines the hidden biases and provides actionable mitigation strategies.

91%

Overall Accuracy

0.65

Disparate Impact

Indicates significant bias

Five Critical Sources of Bias



Historical Bias

Senior developers' issues marked high priority regardless of urgency. Past organizational biases perpetuated rather than corrected.



Geographic Underrepresentation

Remote and offshore teams systematically deprioritized due to time zones and sparse historical data.



Feature Bias

Metrics favor certain programming languages, development methodologies, and technical jargon.



Label Subjectivity

Human labelers exhibit authority bias, recency bias, and relationship-based preferential treatment.

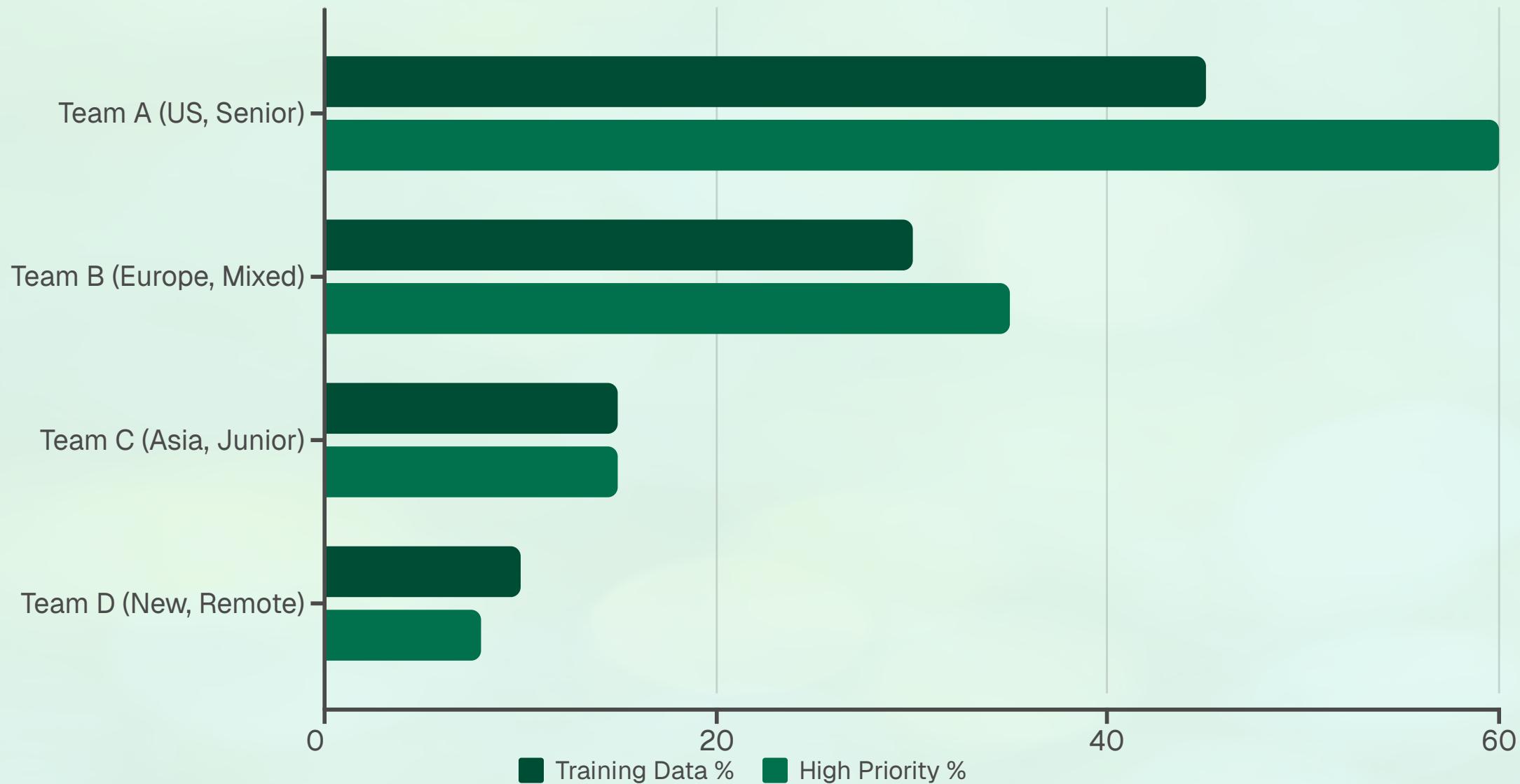


Sampling Bias

Training data excludes abandoned issues, over-represents specific time periods, and lacks failure cases.

The Underrepresentation Problem

Data distribution reveals stark inequities. Senior US-based teams dominate training data, while newer remote teams struggle for visibility.



This creates a vicious feedback loop: underrepresented teams receive less support, perform worse, and generate more low-priority labels, further training the model to deprioritize them.

IBM AI Fairness 360: Our Solution Framework



IBM's open-source AIF360 toolkit provides comprehensive bias detection and mitigation across the entire ML pipeline.

01

Pre-Processing

Balance training data representation

02

In-Processing

Train with fairness constraints

03

Post-Processing

Calibrate predictions for equity

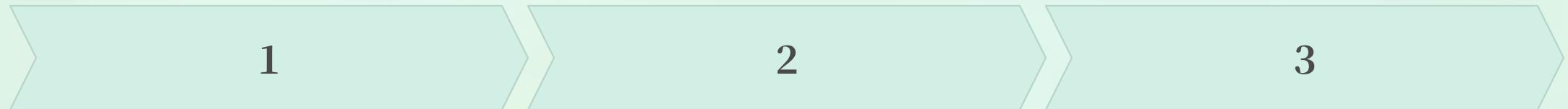
04

Monitoring

Track 70+ fairness metrics continuously

Solution 1: Reweighting for Fair Representation

Pre-processing bias mitigation assigns weights to training samples, ensuring underrepresented teams receive equal importance during model training.



Identify Protected Attributes

Team ID, geographic region, developer seniority

Calculate Sample Weights

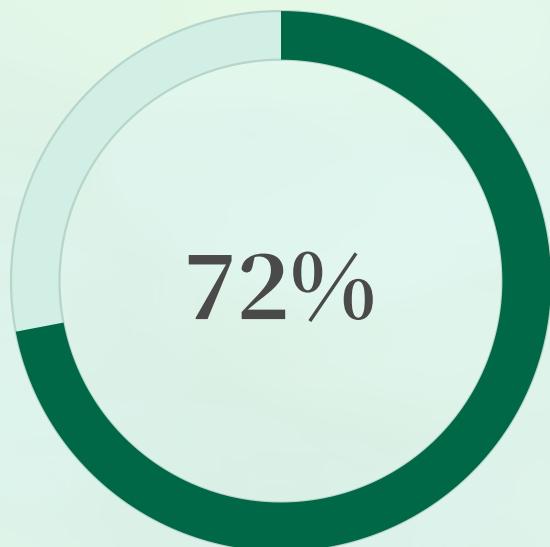
Underrepresented teams get higher weights

Retrain Model

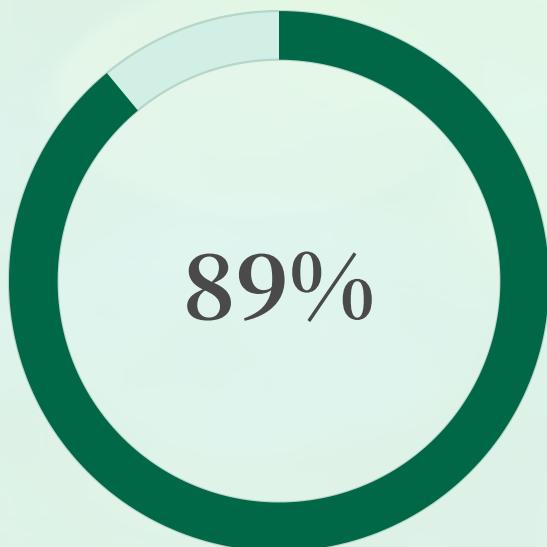
Balanced learning across all groups

```
from aif360.algorithms.preprocessing import Reweighting
```

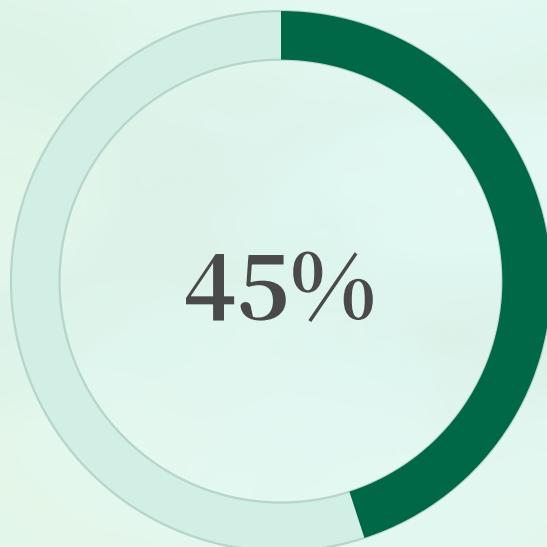
```
RW = Reweighting(  
    unprivileged_groups=[{'team_id': 0},  
    privileged_groups=[{'team_id': 1}])  
dataset_transformed = RW.fit_transform(dataset)
```



Team C Before



Team C After



Disparate Impact Reduction

Solutions 2 & 3: In-Processing and Post-Processing

Prejudice Remover

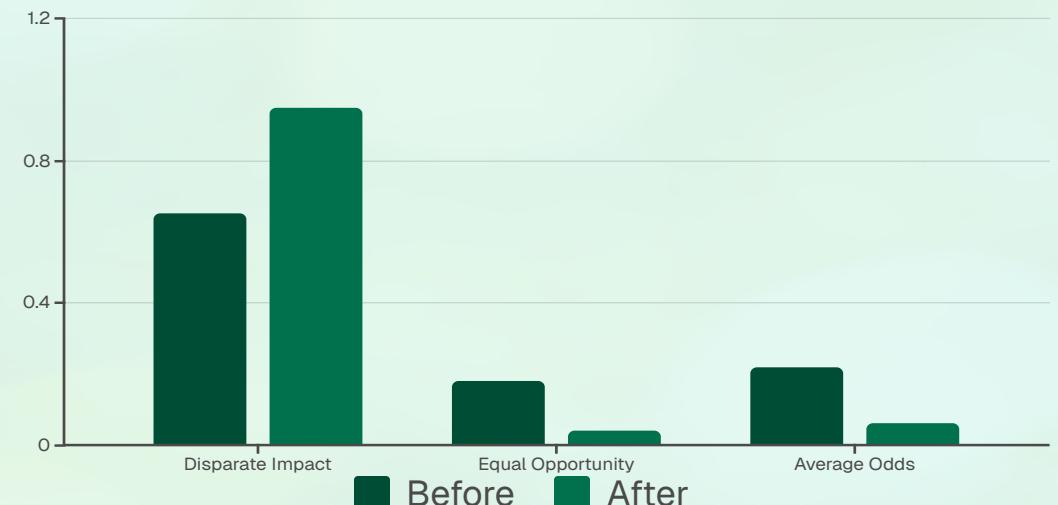
Adds fairness penalty during training, preventing predictions from correlating with protected attributes like team affiliation.

- ❑ **Trade-off:** Overall accuracy drops 2% (93% → 91%), but fairness improves dramatically (Equalized Odds: 0.45 → 0.12)

```
from aif360.algorithms.inprocessing  
import PrejudiceRemover  
  
PR = PrejudiceRemover(  
    sensitive_attr='team_id',  
    eta=25.0  
)  
  
model_fair = PR.fit(dataset)
```

Equalized Odds Post-Processing

Calibrates prediction thresholds to ensure equal true positive and false positive rates across all teams after training.



Continuous Fairness Monitoring Dashboard

Deploy comprehensive monitoring to track fairness in production and alert stakeholders to emerging bias patterns.



Real-Time Alerts

Critical violations trigger immediate notifications when disparate impact drops below 0.8



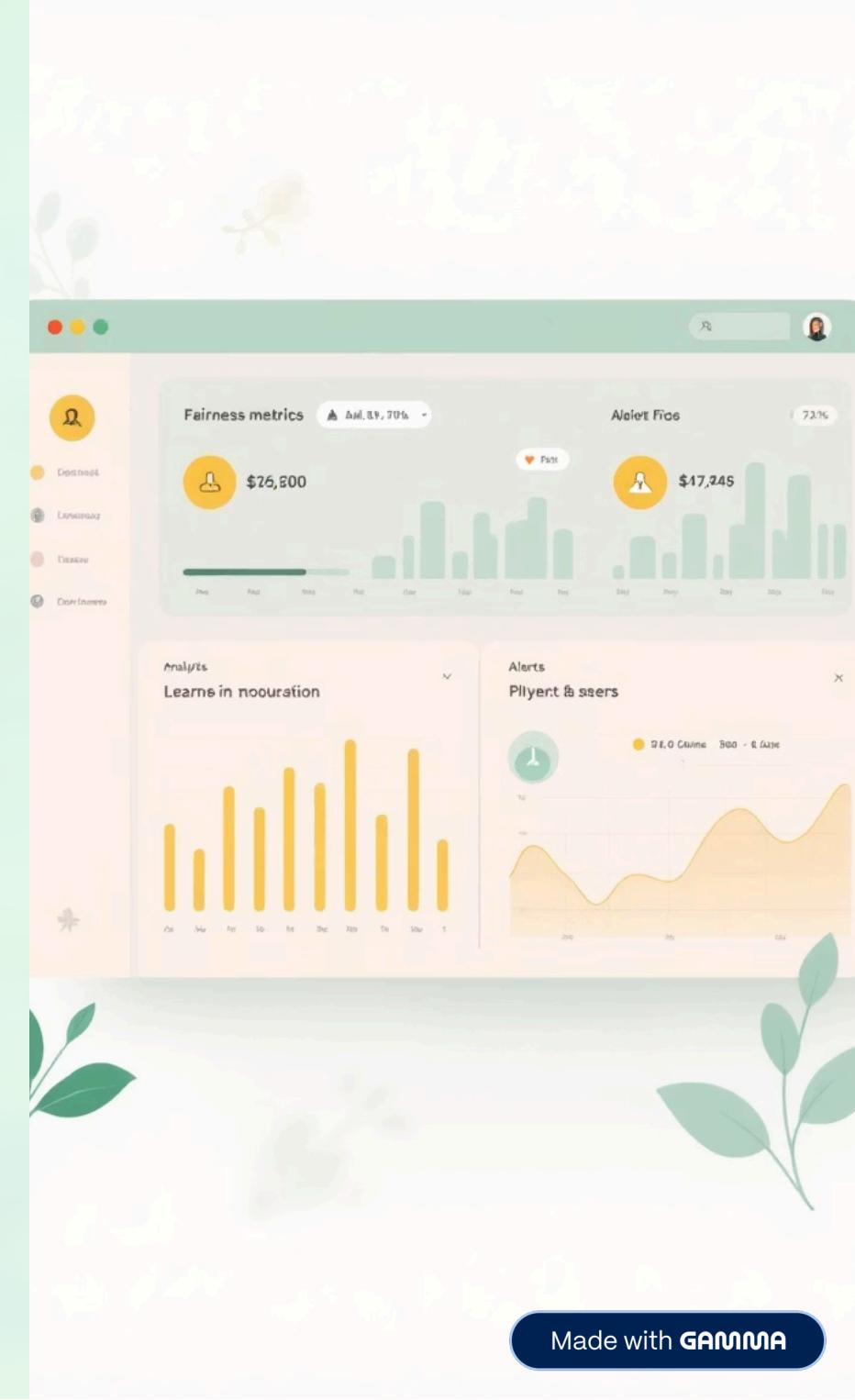
Weekly Reviews

Team-level performance analysis ensures no group falls behind



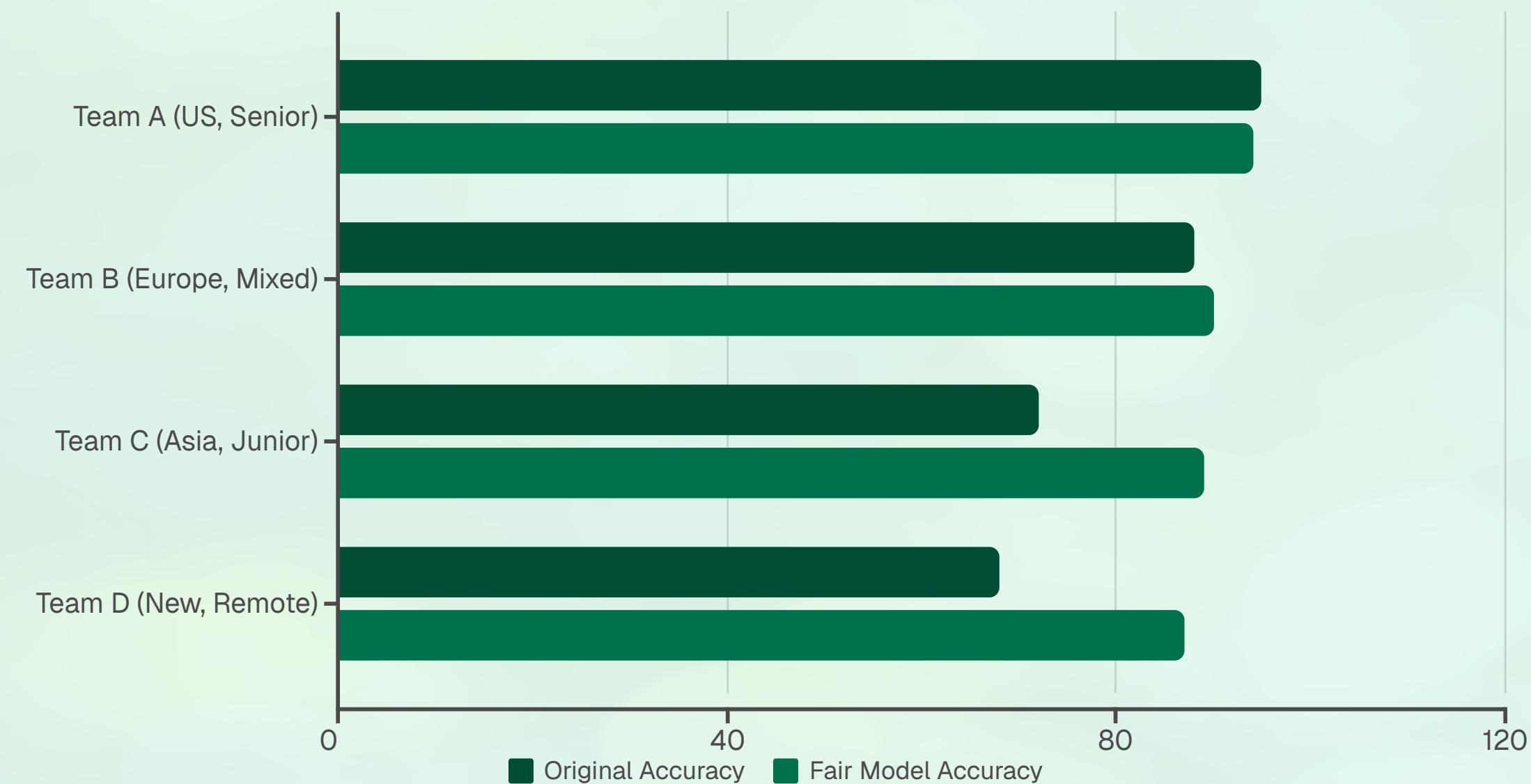
Quarterly Retraining

Fresh data prevents model drift and adapts to organizational changes



Impact: Dramatic Improvements Across All Teams

Fair model implementation delivers measurable business value while ensuring equitable treatment across the organization.



+83%

Team C Resource Allocation Increase

31%

Team Satisfaction Improvement

From 6.2 to 8.1 out of 10

39%

Faster Resolution for Team C
8.5 days reduced to 5.2 days

Ethical AI: Business Imperative, Moral Responsibility

Key Takeaways

1 Historical data reflects historical biases — fairness requires active intervention through reweighing and constraints

2 Underrepresented groups suffer most — targeted mitigation with IBM AIF360 is essential, not optional

3 Fairness and accuracy coexist — proper techniques achieve both with minimal trade-offs (90%+ accuracy maintained)

4 Continuous monitoring is critical — bias emerges over time as organizational conditions evolve



Final Fairness Score: Disparate impact improved from 0.65 (poor) to 0.93 (excellent) while maintaining accuracy across all teams. Ethical AI isn't just compliance—it's competitive advantage.