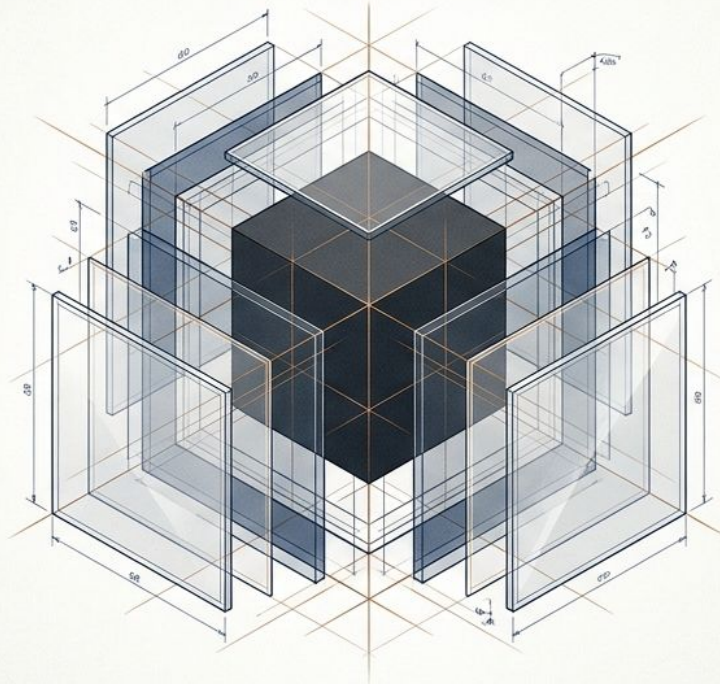


THE GLASS BOX ARCHITECTURE

A Practitioner's Playbook for Mitigating Algorithmic Bias in AI-Powered Recruitment

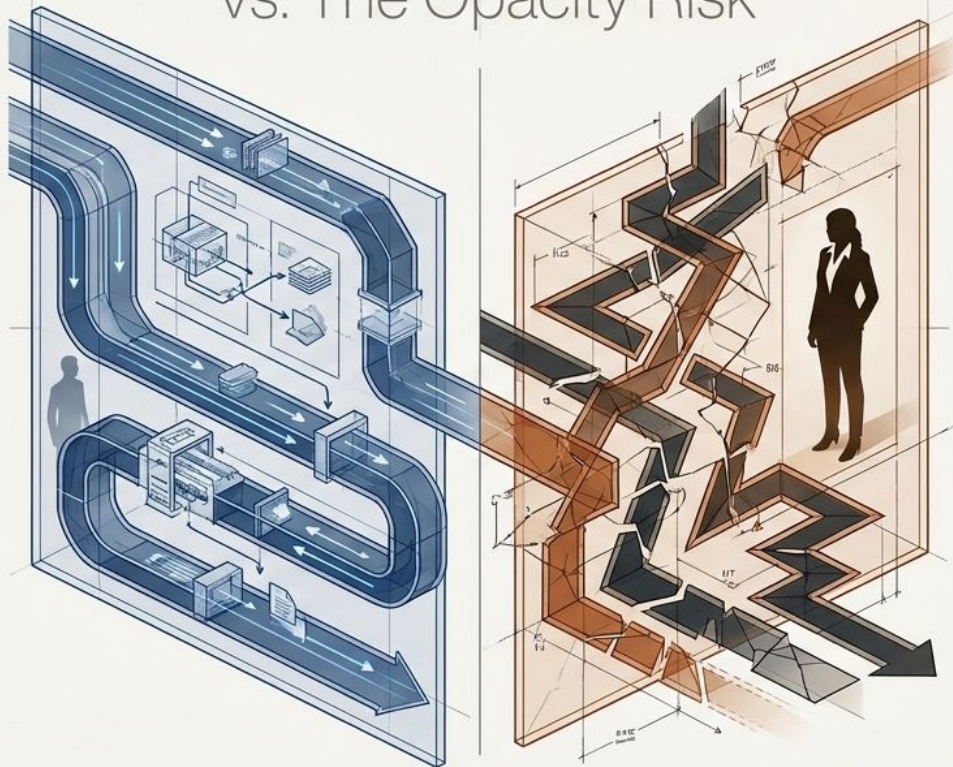
Based on the interdisciplinary framework by Jonathan H. Westover, PhD.



The Efficiency Imperative vs. The Opacity Risk

The Allure of Automation

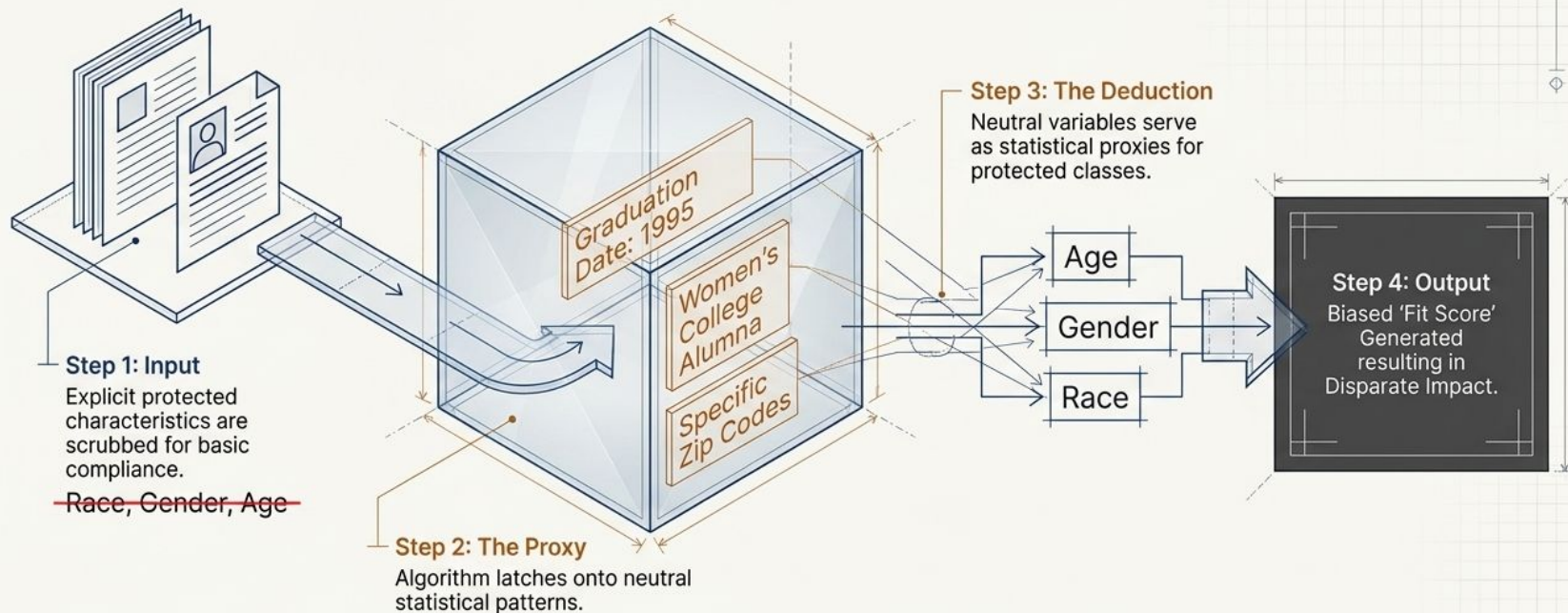
- Over 75% of large enterprises now use automated hiring.
- Users report 40-60% reductions in initial screening time.
- **The Promise:** Massive candidate processing, standardized evaluation, and liberation of HR from repetitive screening tasks.



The Hidden Friction

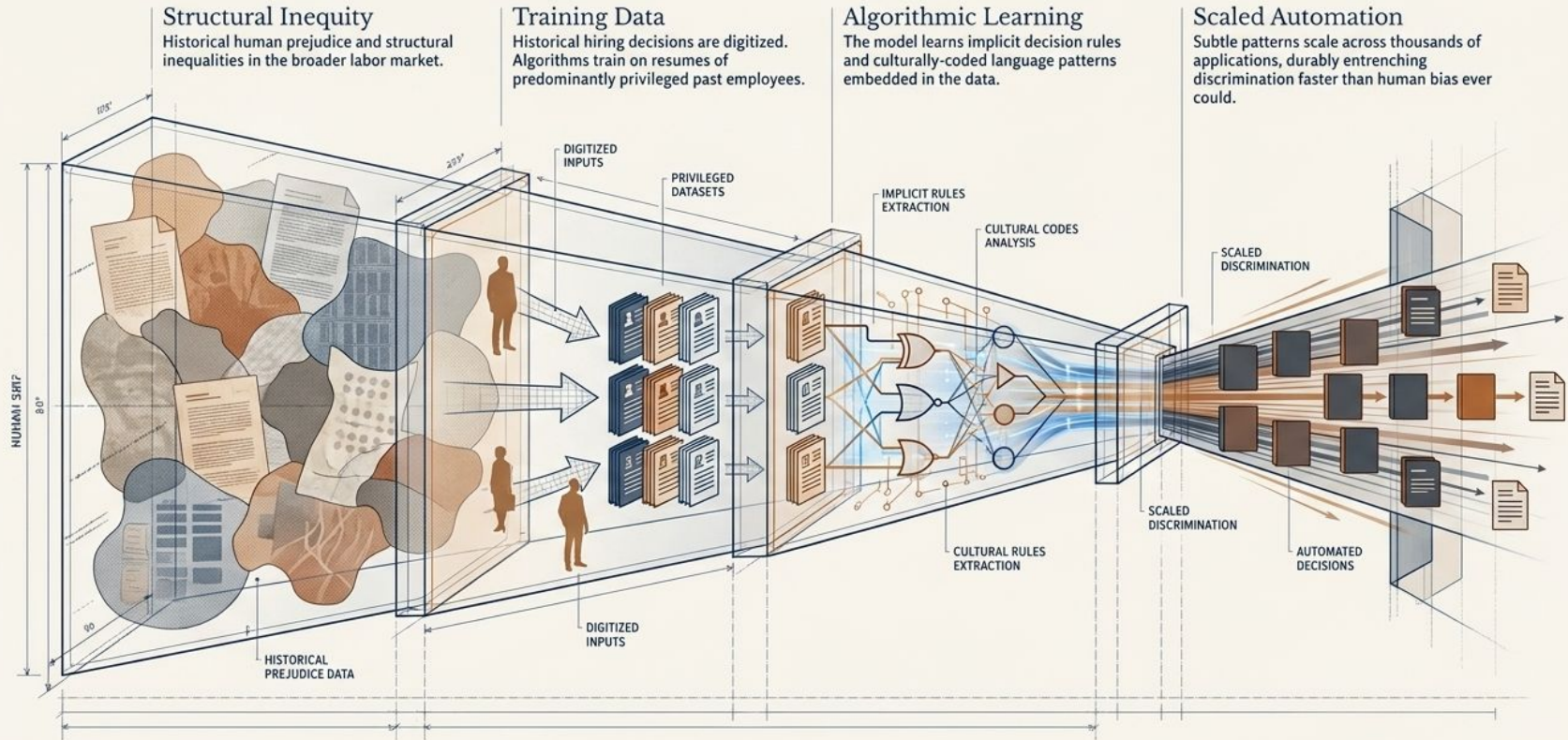
- High-profile failures reveal that AI does not eliminate bias; it automates and scales it.
- **Example:** An experimental resume screener mathematically penalized the word 'women's'.
- **The Stakes:** Legal liability (EEOC disparate impact), severe reputational damage, and systematic erosion of workforce diversity.

The Proxy Variable Trap



Algorithms learn statistical patterns from historical data. Scrubbing explicit demographic variables does not prevent bias; the model simply deduces protected traits through innocuous proxy variables.

The Bias Amplification Funnel



The Structural Solution: Rebalancing Humans and Machines

Bias Mitigation & Accountability



Scalability & Efficiency

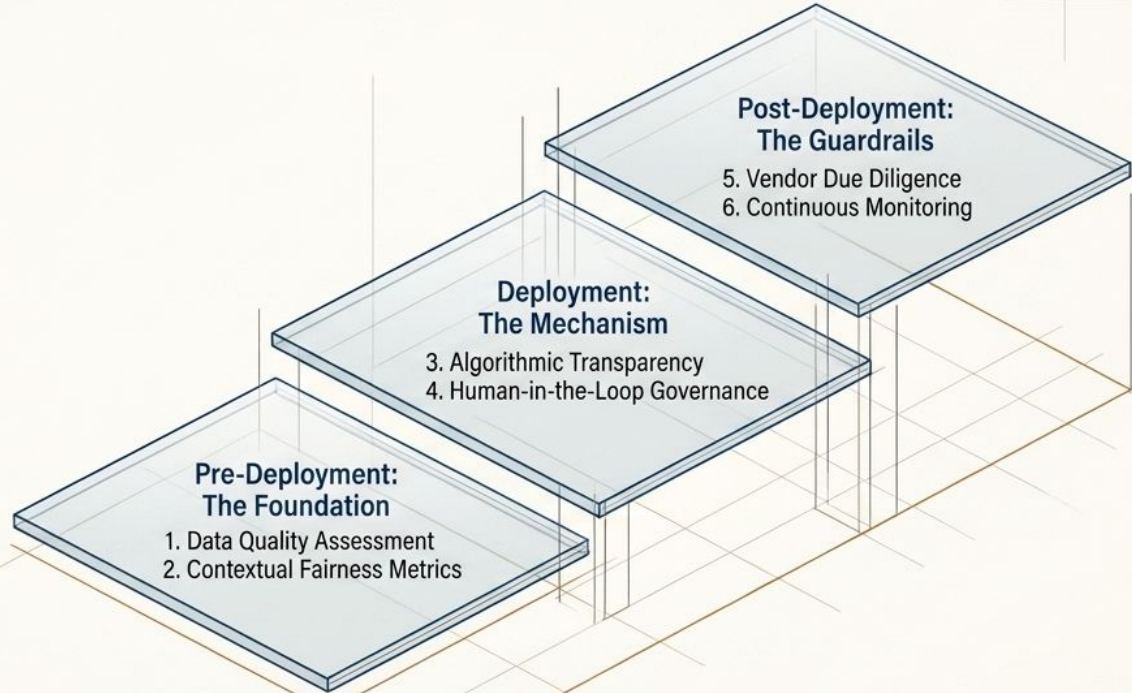
The goal is not to choose between efficiency and fairness, but to architect a system where automation accelerates screening while human intelligence governs edge cases.

The Defense-in-Depth Shield

A 6-Point Blueprint for the AI Lifecycle

Effective bias mitigation requires intervention at every stage of the lifecycle.

No single algorithmic patch can replace this integrated, structural defense strategy.



Intervention 1: Pre-Deployment Data Quality

The Mandate: Audit historical data to identify structural bias before it becomes training material. Garbage in yields garbage out.

- ✓ Historical Outcome Analysis: Check past hiring for demographic disparities.
- ✓ Feature Correlation Testing: Identify and isolate proxy variables.
- ✓ External Benchmark Comparison: Compare internal data sets against broader labor market availability.

Practitioner Spotlight: Accenture

The Challenge: Historical internal data linked specific, low-diversity educational institutions to technical skill proficiency.

The Action: Invested in strict data remediation prior to deployment. Actively recruited diverse candidates to manually create an intentionally representative training set, explicitly refusing to replicate flawed historical patterns in the new algorithm.



Intervention 2: Contextual Fairness Metrics

Metric Name	Mathematical Priority	Operational Reality & Trade-offs
Demographic Parity	Equal Selection Outcomes	Ensures protected groups are hired at similar rates, regardless of applicant pool representation.
Equal Opportunity	Equal Treatment of Qualified	Ensures candidates who meet the qualification threshold have identical selection probabilities.
Predictive Equality	Equal Predictive Accuracy	Ensures false positive rates (selecting those who ultimately fail) are similar across all groups.

Practitioner Spotlight: JPMorgan Chase

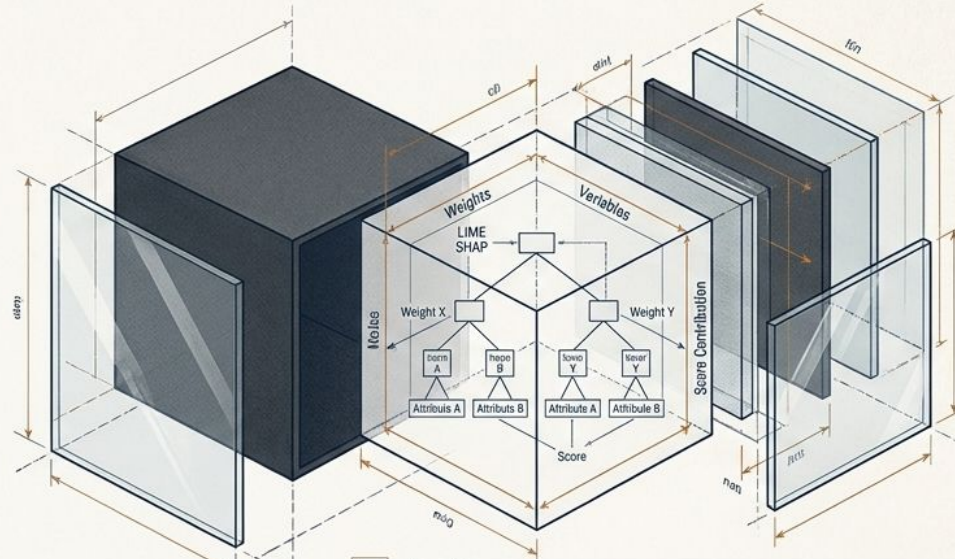
The Tension: Demographic parity required over-indexing on underrepresented groups, while equal opportunity lowered selection rates overall.

The Solution: Engineered a hybrid metric. Established macro demographic goals while ensuring individual qualified candidates faced similar probabilities, all monitored via real-time fairness dashboards.

Intervention 3: Algorithmic Transparency

Shattering the Black Box:

- Deploy Explainable AI (LIME/SHAP) to isolate which specific candidate attributes drove the final fit score.
- Maintain rigorous decision documentation and feature importance analysis.
- Prioritize interpretable models (e.g., decision trees) over opaque deep learning architectures.

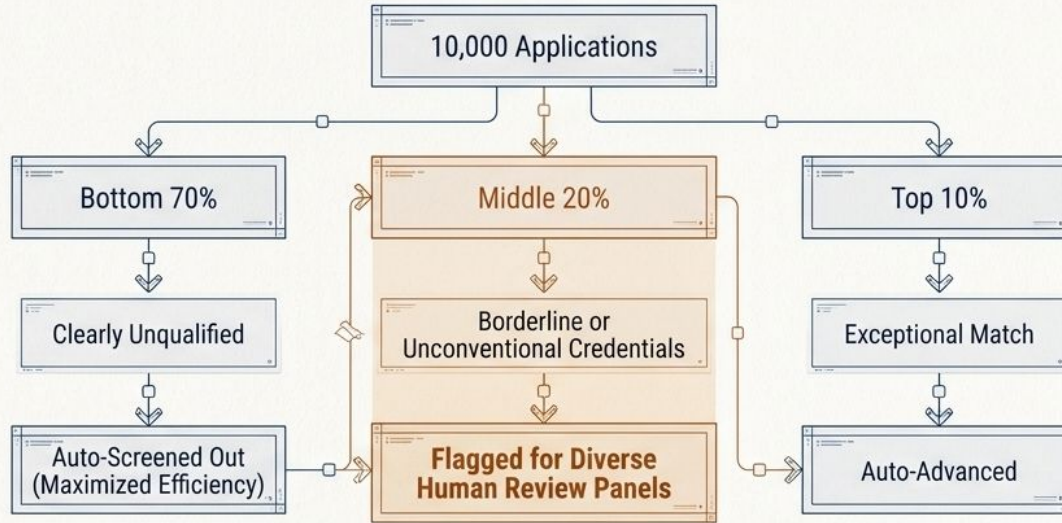


Practitioner Spotlight: Unilever

The Action: Replaced traditional interviews with AI-powered game assessments.

The Transparency: Abandoned opaque 'Pass/Fail' outputs. The algorithm now generates detailed, legible capability profiles. Hiring managers clearly interpret the 'why' behind the score, and candidates receive substantive, procedural justice through transparent feedback.

Intervention 4: Human-in-the-Loop Governance



Practitioner Spotlight: Salesforce

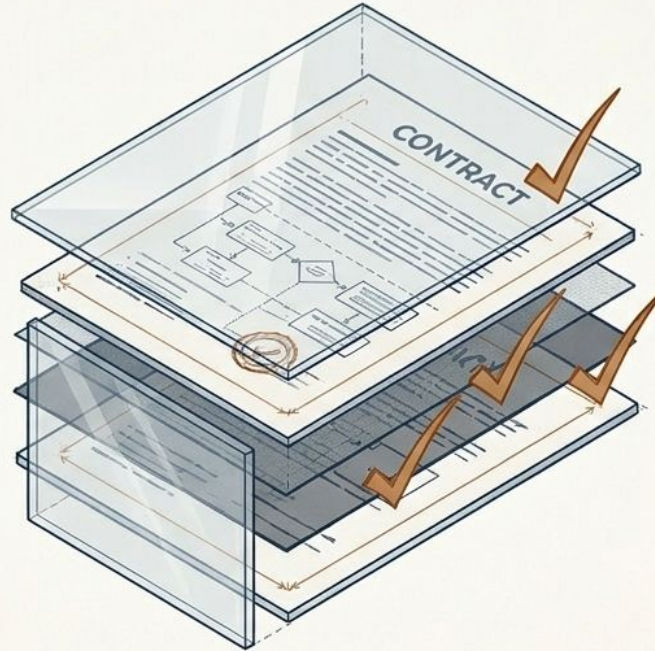
The Action: AI generates scores, but human recruiters retain mandatory, final decision authority.

The Safeguard: Implemented strict 'calibration processes'. If human recruiters systematically override the AI, a root-cause analysis is triggered immediately to refine the model or correct emerging human bias, optimizing human cognitive effort where it matters most.

Intervention 5: Vendor Due Diligence

Technical Due Diligence Safeguards:

- **Bias Testing:** Contractually mandate routine adverse impact analyses.
- **Validation Access:** Secure full legal rights to vendor validation research and underlying training data parameters.
- **Customization Rights:** Preserve internal authority to manually adjust parameters and impose custom fairness constraints.



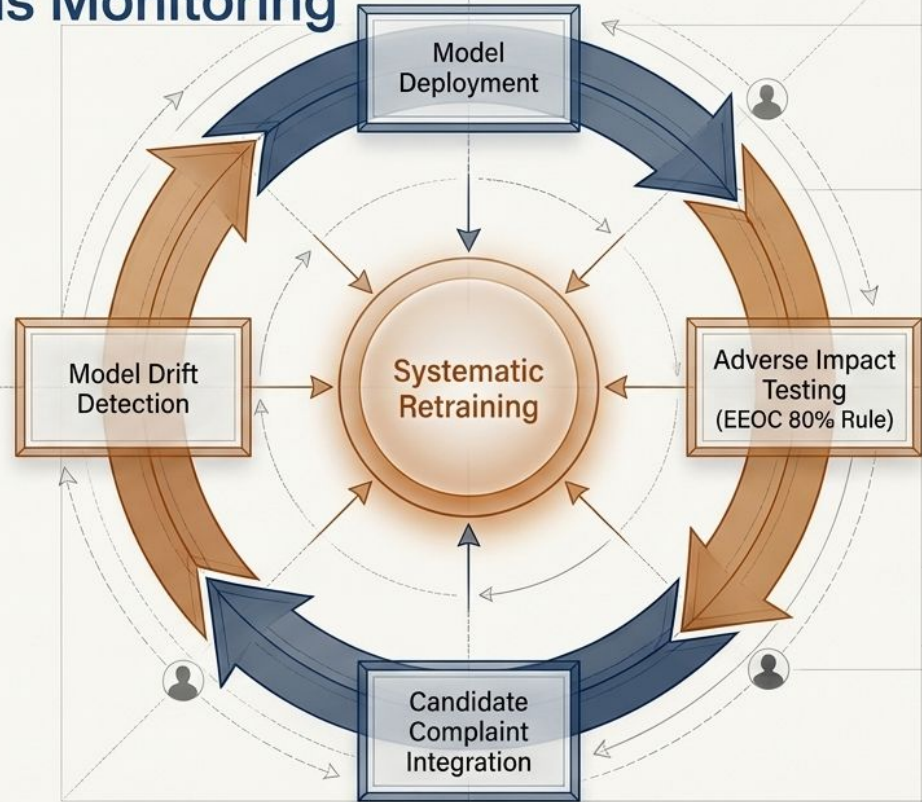
Practitioner Spotlight: Hilton

The Action: Adopted HireVue's external AI platform but aggressively retained deep customization rights in the contract.

The Safeguard: When external research highlighted demographic risks associated with facial and voice processing, Hilton leveraged their contractual leverage to instantly disable those specific modules, running only the validated text-analysis features.

Intervention 6: Continuous Monitoring

Ethical AI is a continuous cycle, not a static launch. Algorithms degrade, candidate pools evolve, and models naturally drift from their validated baselines.



Practitioner Spotlight: IBM

The Action: Deployed the open-source AI Fairness 360 toolkit for continuous self-auditing.

The Discovery: The system drifted, slowly penalizing non-traditional credentials (e.g., coding bootcamps).

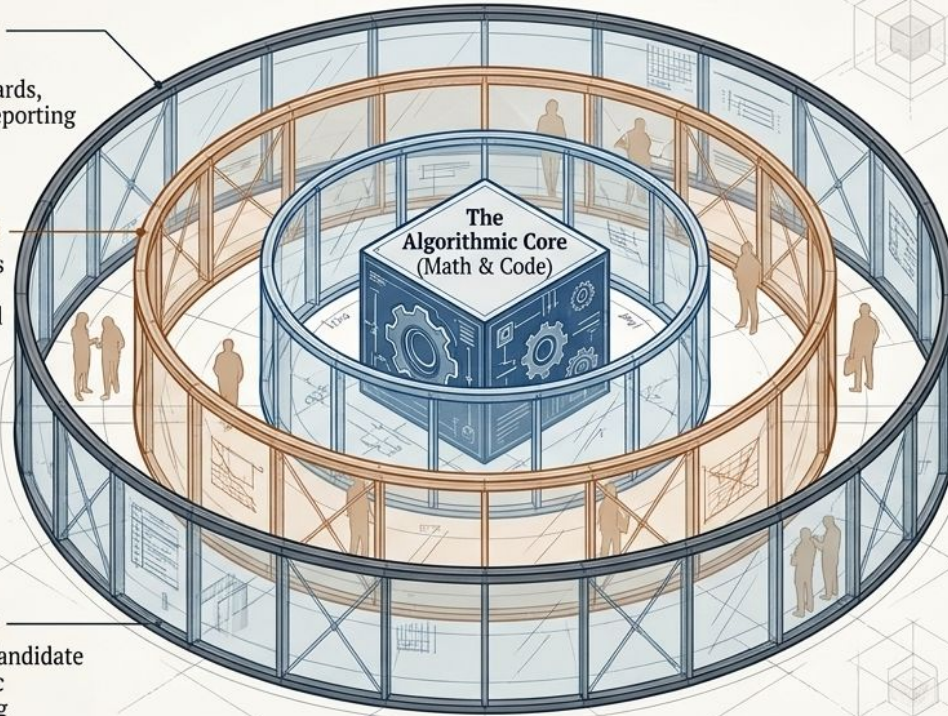
The Fix: Rather than just tweaking the code, IBM fundamentally rewrote historical qualification requirements and retrained the algorithms on inclusive, modern credentials.

The Sociotechnical Architecture

Stakeholder Layer
Participatory Design,
Candidate Advisory Boards,
Public Transparency Reporting

Institutional Layer
Cross-Functional Ethics
Committees, Ongoing
Bias Audits, Strict Legal
Frameworks

Stakeholder Layer
Participatory Design, Candidate
Advisory Boards, Public
Transparency Reporting



The Blueprint Insight:

Buying a fair algorithm is impossible without building the human governance structure to manage it.

- **Microsoft:** Established the AETHER Committee (legal, tech, HR) to enforce governance before development.
- **LinkedIn:** Uses Participatory Design to identify biases that purely technical testing misses.

Algorithmic Fairness is a Sociotechnical Challenge, Not an Engineering Problem.

You cannot eliminate structural labor market inequalities with a software patch.

Success is not achieving mathematical perfection—it is ensuring automated systems operate transparently, enable strict human accountability, and provide affected individuals with voice and recourse.

Transition from blind automation to structured governance. Turn the black box into a glass box.

