# Balancing Bias and Fairness: The Role of AI in Shaping Equitable Hiring Practices

Aliya Quazi
*Assistant Professor*
*Bright Business School, Hubli*
*Approved by AICTE Delhi and Affiliated to Karnatak University Dharwad*
*Email: aliyaqazi8275@gmail.com*

## Abstract
*This conceptual paper explores the complex role of artificial intelligence (AI) in shaping equitable hiring practices, emphasizing the need to balance bias and fairness in AI-driven recruitment systems, as AI has the potential to transform hiring by enhancing efficiency, reducing human error, and supporting data-driven decision-making, but it also poses significant ethical challenges, particularly in relation to the perpetuation of existing biases and inequalities in the workforce, as algorithms are often trained on historical data that may reflect the biases present in previous hiring decisions, leading to biased outcomes that disproportionately affect underrepresented groups, including women, racial minorities, and neurodivergent individuals, making it critical for organizations to carefully consider the design and deployment of AI systems to ensure that they promote fairness rather than exacerbate systemic discrimination, furthermore, this paper discusses the importance of transparency and accountability in AI systems, advocating for continuous monitoring, auditing, and validation of AI models to ensure that their outputs align with ethical hiring standards, such as diversity, equity, and inclusion (DEI), while also emphasizing the need for diversity in the team's developing these AI tools to mitigate the risk of unintentional biases, additionally, the paper investigates various strategies to enhance the fairness of AI in hiring, including using de-biasing algorithms, conducting regular audits, ensuring that training data is diverse and representative, and incorporating human oversight in the decision-making process to prevent the full automation of hiring decisions, the study also examines the regulatory landscape surrounding AI in hiring, noting the growing importance of legal frameworks that address issues of discrimination and fairness, such as the General Data Protection Regulation (GDPR) and the Equal Employment Opportunity Commission (EEOC) guidelines, and how these frameworks can guide organizations in their use of AI to avoid discriminatory practices, ultimately, this paper provides a conceptual framework for organizations to navigate the ethical complexities of AI in hiring, emphasizing that while AI has the potential to increase efficiency and improve hiring outcomes, its deployment must be carefully managed to ensure that it contributes to creating fair, inclusive, and equitable workplaces.*
***Keywords:*** *Artificial Intelligence (AI), Fairness in Hiring, Bias in Algorithms, Equitable Recruitment, Diversity, Equity, and Inclusion (DEI), Ethical Decision-Making*

## I.     Introduction

As AI drives transformation in how organizations manage their workplaces and attract talent, its role in recruiting and hiring has emerged as a prominent vector for both technological advancement and ethical criticism as organizations increasingly deploy AI-enabled tools to streamline decisions, reduce human error, and enhance efficiency in candidate selection, while facing the challenge of potentially exacerbating systemic discrimination through algorithmic bias from training on historically biased data, a tension prompting calls from scholars and practitioners for greater care in revealing the nature of fairness and bias in AI-based hiring systems that are ironically frequently designed and implemented without consideration for diversity, equity, and inclusion (DEI) principles, resulting in harm to labor market equity for marginalized populations such as women, racial minorities, neurodivergent individuals, and people with disabilities as they replicate or amplify existing inequities, thus highlighted sociotechnical systems that algorithmic hiring tools could enforce, especially because automated decision-making may produce opacity in hiring criteria and diffuse accountability, responsive to a global trend in regulatory agencies increasing scrutiny of AI systems under anti-discrimination statutes, drifting from the European Union and its General Data Protection Regulation (GDPR) to the United States Equal Employment Opportunity Commission (EEOC) goals, with the ethical and legal implications of using AI in hiring moving increasingly to central organizational governance and technological design, as a growing consensus emerges that AI fairness cannot solely be undertaken as an engineering challenge, but must be contextualized within the wider sociolegal context regulating employment processes, and leading to the proposal of a conceptual framework that construes fairness as an algorithmic-design target and normative imperative for the accountable design,

deployment, and operationalization of AI in hiring, centered on transparency mechanisms such as model documentation, bias audits, and algorithm accountability structures that could enable organizations to document, manage, and design against discrimination; and of how to de-bias AI hiring systems starting at data stages through the representativeness and inclusiveness within training datasets, extending out to institutional arrangements that prioritize diverse development teams and embed human oversight into AI-style hiring decisions processes to slow automation over-reliance that contributes to a burgeoning literature discussing the ethics of AI in hiring as a multifaceted problem that needs collaboration across domains to assist multicompetent stakeholders, culminating in an argument that the AI being the tool of promoting and not undercutting equitable employment calls for careful ethical inquiry, continuous assessment and participatory design practices in line with DEI values, legal compliance, and organizational justice orientations.

## II.    Research Gap

But even as advocates of thoroughgoing AI adoption in recruitment tout benefits of better, quicker, cheaper, fairer hiring, there is an uncomfortable near absence in the academic literature of theoretical treatment of the way these profit seeking efficiency forms obstinately resist with the foundational nature of equity and inclusion in the domain, whereby the engineering of AI hiring systems gravitates toward mathematical optimization and predictive accuracy at the expense of socio-ethical priorities such as demographic representativeness, justice perceptions, repair of labor market iniquities -- so that algorithmic systems designed to be neutral do often collude in reproducing exclusionary hiring patterns that emerge in their training data reflecting the historical record of inequality (Kochling & Wehner, 2020; Ajunwa, 2023) -- and while there is growing attention among researchers to various technical fairness metrics such as demographic parity, equal opportunity or counterfactual fairness, there remains much less attention paid to the way that embedding of differential treatment, through historical or newly emergent discrimination is made more difficult in the live workplace in conceptually asking what intersectional constructs do we require to aid not merely in simulation-based fairness-inducing technology but also in the hard messy path toward reconstruction of socio-technical systems under constraint of labor market outcomes (Binns 2023; Yapo & Weiss, 2022) -- combined with investigation into the broader political economy of technosocial efficiency erasing the social processes through which bias might be embedded, engaged with or compounded throughout the recruitment pipeline, notably in the context of both pre-screening algorithms but also how automated assessments and video-based interviews exacerbate biases against groups such as Black and Indigenous applicants, neurodivergence, and non-normative speech or credentialing (Chandrasekharan et al., 2022; Lee et al., 2023)--and if some aspects of policy and regulatory intervention, such as sectoral specific AI audits, fairness certification standards emerge in reaction to the documented harms of AI then still the field is conceptually scattered as to how tensions between maximizing AI efficiency defended as business necessity) and meaningful equity (necessarily a social justice work) can be reconciled because if the principle of optimization can play out then optimization and movement toward greater equity could well be at odds unless and until design trade-offs, governance mechanisms and orchestrations requiring stakeholder deliberative efforts are embedded into all aspects of the design and procurement of an AI at outset (Raji et al., 2022; Shankar et al., 2023) and so this paper addresses an urgent gap in relevant scholarship by proposing a conceptual framework that foregrounds and interrogates the normative, operational and systemic frictions and trade-offs between optimizing for AI-driven efficiency in hiring and ensuring that such optimization operates neither to reproduce structural inequities nor undermine inclusive employment practices in increasingly digitally mediated labor markets.

## III.    Theoretical Background and Literature Review

Even though scholars have acknowledged that artificial intelligence (AI) can potentially improve the efficiency of recruitment processes, by performing tasks ranging from automated resume screening to automated skills-to-role matching and reduction of time-to-hire and recruiting labour costs, the promise of computational expediency is counterbalanced by severe ethical and operational hazards; above all, the opacity of algorithmic decision-making (which becomes higher with the complexity of model deployed); the lack of interpretable black-box models; and the scalability of discrimination; that is, bias once confined to the actions of an individual recruiter is now amplified across automated processes being applied at scale, which called for critical examination about how concepts of algorithmic fairness and bias were to be defined, operationalized, and contested across both technical and social domains, since fairness in AI was no monolithic concept at all, but a complex composite of diverse, often controversial, but not always clearly defined, typologies – procedural (focusing on the integrity of the decision-making process), distributive (concerned with the equity of the outcomes), individual (which mandates similar treatment for similar individuals), or group (which maintains statistical parity across aggregated groups) fairness; – all of them bringing foundational challenges to implementation, especially if inherent trade-offs and incompatibilities between definitions of fairness were acknowledged and well understood and since training datasets underlying these models were often human decisions that reproduced long-lasting structural inequities that emerged over the course of time, discriminatorily influencing hiring processes (for example,

racialized name discrimination; gendered occupational sorting; or ableist assumptions in assessments) and therefore reproducing historical biases under the illusion of neutrality; leading to reified harms for marginalized groups who are often simply too few in data to avoid being algorithmically excluded or misclassified (Eubanks, 2018; Cowgill et al., 2021); whilst ethical AI guidelines proliferated from government agencies, academic consortia and corporate actors (OECD AI Principles, IEEE's Ethically Aligned Design framework, European Commission's AI Act draft), most remained high level, aspirational, lacking any enactment mechanisms or sectoral specifics on how they could work on the ground in the setting of hiring contexts, mostly neglecting the positionings of intersecting identities within the whole concept of discrimination (Jobin et al., 2019; Mittelstadt, 2022; Whittlestone et al., 2021) in practices and consequential, specific terms; despite the emergence of fairness-aware machine learning models, proprietary algorithms with unbounded performance proxies (e.g. cultural, fit, communication style) remained popular in commercial Hiring platforms, promoting neither transparency nor validation for their biases (Berk et al., 2021; Kim, 2023); although algorithmic auditing methods, model documentation standards, datasheets for datasets, model cards and audit trails, were gradually being developed and there was a clear need for them to fill the gap between these important technical tools and broader organizational accountability systems needed to detect, remedy and govern such discrimination on a large scale (Raji et al., 2020; Passi & Barocas, 2019), empirical research around AI fairness in hiring had mainly focused on classification parity metrics, system performance and not on how fairness was negotiated by employers, interpreted by candidates, motorized by institutional manufactures, which led to indicators of a conceptual gap around the socio-technical dynamics conditioning the use of AI in hiring, including how this was often conflicted against efficiency imperatives that ultimately crowded out any deliberation around equity or long-term justice (Green, 2022; Amoore, 2020), and while the totality of the scholarship provided foundational insights into the aspects of the technical fairness, statistical bias mitigation and ethical principles it still frequently omitted a careful addressing of the normative tensions between economic optimization and distributive equity in labor markets, the contextual variability of fairness, as always, across practical use cases and the structural power asymmetries that were deeply embedded in the process of AI development and deployment particular in the private sector where proprietary algorithms and trade secrecy inhibited transparency and meaningful stakeholder engagement (Marda & Narayan, 2021; Veale & Zuiderveen Borgesius, 2021), auguring the notion of the need for more future studies to build bridges between the disciplines of computer science, organizational behavior and employment law and co-create integrated frameworks that would not only detect bias, but also reframe fairness as a multi-dimensional, context-sensitive construct that is situated well within ongoing practices of recruitment, regulation and resistance.

## IV.  Conceptual Framework

As AI systems assume a central role in how hiring works, we need a unified intellectual infrastructure to disentangle the analytical dimensions of algorithmic bias and fairness by bringing together normative, operational, and relational components that provide the structure for how fairness is defined, assessed, and implemented in sociotechnical systems, starting with the understanding that fairness in AI hiring cannot be boiled down to statistical parity or other isolated metrics but, rather, it must be seen as a web of interdependent values—such as procedural justice (ensuring transparency in decision-making), distributive justice (ensuring equity between groups), contextual fairness (recognizing positionality in social hierarchies), and explanatory adequacy (providing intelligible rationales for decisions)—that must all be attacked simultaneously in order for us to be able to meaningfully know if an AI system is promoting equitable hiring or undermining it (Lee & Singh 2021; Wong & Webb 2023), and this framework itself centers the cross-cutting stakeholder ecosystem surrounding AI hiring and therefore envisions power as dispersed and asymmetric across four high-level actor groups—namely AI developers who build and tune models according to organizational needs and data constraints; HR teams and hiring managers that incorporate such tools into workflows and often rely on vendor assurances around compliance and fairness; job applicants whose data and outcomes are shaped by algorithmic decisions but who remain largely outside the system design or validation process; and regulators and policymakers that govern outcomes through ever-larger regulatory vectors such as the EEOC (Hoffmann 2022; Contissa et al. 2023)—and within this ecosystem, our framework identifies three major classes of bias that impinge fairness in AI hiring: (1) data bias, arising out of unrepresentative or historically discriminatory training sets, which can perpetuate prior patterns of discrimination based on demographic characteristics (i.e., race, gender, disability, and socioeconomic status); (2) measurement bias, whereby metrics that are relied upon to proxy job fit or performance (e.g., GPA, past experience, communication style) disproportionately harm certain cohorts or do not adequately control for contextual differences in opportunity structures; and (3) deployment bias, referring to the mismatch between how models are trained, deployed, and eventually used (i.e, decision automation without human oversight, and relying on AI outputs for decisions beyond their intended scope) can cause latent inequities to be unmasked more severely in practice (Cowls et al. 2021; Biega et al. 2022)—and these issues are compounded by algorithmic feedback loops, wherein biased decisions made by AI systems find their way into organizational practices, influencing future training data, and consequently, creating systemic inequalities within labor markets—such as when

underrepresented cohorts are excluded from hiring based on faulty assessments, leading to a homogeneous workforce that serves as the input for future predictive models and thus reinforcing demographic underrepresentation while also narrowing the definitional boundaries of "success" (Grosz et al. 2019; Ali et al. 2022), a process that is made worse still by performance-based learning systems that optimize for historical hiring outcomes (e.g., retention, promotions) without interrogating whether these outcomes were configured equitably to begin with, necessitating fairness frameworks that attend to both proximal outcomes and long term social feedback loops, especially in settings like employment where cumulative disadvantage can persist for decades (Lukyanenko et al. 2023; Wilson et al. 2024)—and so this conceptual framework situates fairness not as merely a statistical adjustment in machine learning pipelines but as a multi-level, temporally extended, politically contingent construct that requires active engagement from stakeholders (transparency of design logic and iterative validation against real-world outcomes), and structures that necessitate that systems be reconfigured when it is found that certain populations are being disproportionately harmed; thereby establishing the goal of not only procedural parity, but substantive equity in the distribution of access to work, advancement, and institutional recognition across the many different social groups.

## V.    AI Approaches to Bias and Fairness in Hiring

AI approaches to bias and fairness in hiring are a critical area of study, as the integration of artificial intelligence (AI) into recruitment processes promises to both mitigate and potentially amplify existing biases in hiring decisions, and scholars have proposed various strategies for reducing bias and enhancing fairness in AI systems. The use of AI in hiring aims to streamline processes, reduce human error, and eliminate discriminatory practices based on race, gender, or other protected characteristics; however, the very data-driven nature of AI models often introduces new risks of perpetuating bias, particularly when the algorithms are trained on historical data that may reflect the discriminatory practices of previous hiring decisions (O'Neil, 2016). For instance, a study by Angwin et al. (2016) revealed that commercial risk assessment tools used in the justice system exhibited bias against Black defendants, which raises concerns about similar biases in hiring algorithms that might be trained on biased data, inadvertently reinforcing existing inequities. This challenge is further compounded by the fact that AI systems are not immune to human biases embedded in data, where algorithmic models may "learn" from the disparities that exist in historical hiring patterns, which could disproportionately affect underrepresented groups, particularly women, ethnic minorities, and individuals with disabilities (Mehrabi et al., 2021). To address these issues, a wide range of AI approaches have been developed, focusing on various aspects of fairness, such as *individual fairness*, *group fairness*, and *equalized opportunity*, among others, which represent theoretical frameworks designed to ensure equitable treatment in AI decision-making processes (Dastin, 2018). One popular method is data preprocessing, where data used to train AI models is cleaned or transformed to reduce biases before training the algorithm; for example, *reweighting* and *re-sampling* techniques are commonly employed to correct imbalances in the data that may favor one group over another, which helps reduce discrimination in the resulting predictions (Kamiran & Calders, 2012). Another approach is *algorithmic fairness constraints*, where mathematical models are designed to achieve certain fairness criteria during training by incorporating fairness metrics such as statistical parity, demographic parity, or equalized odds (Chouldechova, 2017). By doing so, AI systems are calibrated to meet the specific fairness requirements, ensuring that the predictions made by the model do not systematically disadvantage certain demographic groups (Barocas, Hardt, & Narayanan, 2019). However, while these techniques are effective at the model level, they are often criticized for being insufficient in addressing the broader issue of fairness in hiring, as they do not take into account the complexities of social and cultural dynamics that shape decision-making (Binns, 2018). This brings the focus to the importance of *algorithmic transparency* and *interpretability*, which enable stakeholders to understand how AI systems make decisions, and therefore, assess whether these decisions are fair and just (Lipton, 2016). Transparency is considered crucial for detecting and mitigating bias in hiring systems, especially since opaque black-box algorithms could perpetuate discriminatory practices without detection. Tools such as *Shapley values* and *LIME* (Local Interpretable Model-agnostic Explanations) are increasingly used to provide insights into the decision-making processes of AI, making the underlying logic more accessible to both developers and end-users (Ribeiro, Singh, & Guestrin, 2016). Despite these advancements, human oversight remains vital, as AI systems alone cannot fully replace the nuanced judgment required in making hiring decisions, and researchers have proposed *human-in-the-loop* models where AI assists recruiters in making decisions rather than fully automating them (Hoffman et al., 2018). This collaboration ensures that fairness principles are consistently applied and that ethical considerations, such as diversity and inclusion, are factored into the decision-making process, recognizing that human input can correct errors that AI models might overlook. Furthermore, the deployment of AI in hiring practices raises significant ethical concerns, particularly regarding *discrimination by design*, where certain algorithms may inadvertently create unfair outcomes despite efforts to reduce bias, thereby highlighting the need for a broader regulatory framework that ensures accountability and fairness in AI-driven hiring practices (Binns, 2018). Researchers have called for clearer regulations and standards governing the development and deployment of AI tools, with some

advocating for industry-wide ethical codes to ensure that AI systems in hiring align with broader societal values, such as non-discrimination, transparency, and accountability (Crawford, 2021). As AI technologies evolve, a multidisciplinary approach that combines insights from computer science, ethics, law, and organizational behavior will be crucial in balancing the trade-offs between reducing bias and achieving fairness in AI-driven hiring, as well as in creating an equitable and inclusive workforce. Therefore, continued research into algorithmic fairness, data practices, and human-AI collaboration is essential for achieving ethical outcomes in AI-based hiring systems that reflect social and organizational values (Holstein et al., 2019).

## VI.     Challenges and Limitations

The challenges and limitations regarding the study of AI in shaping equitable hiring practices are multi-dimensional and related to data limitations and privacy across multiple fronts from the complexity of defining fairness, the existence of ambiguity and vagueness in what constitutes an equitable hiring process, challenges and difficulties in implementation, the existence of competing principles and enforcement mechanisms, discrimination, segregation, autonomy and decision rights to other ethical dilemmas that come into play as AI systems are being created to implement these decisions. This is mainly due to the biased training data problem where AI algorithms are often trained with historical data that inherently contains biases in hiring such gender, race or even age discrimination, which leads to even more biased AI systems (Mehrabi et al., 2021). This problem is compounded by the fairly homogeneous nature of many datasets, which can leave AI models trained on unrepresentative data blind to the backgrounds and competence of under-represented groups, biasing hiring outcomes towards historically dominant demographic groups (Binns, 2018). Such as when an AI recruitment tool used by Amazon trained on resumes from mostly male applicants inadvertently began rejecting female candidates altogether because its training data reflected the fact that a widely male workforce had the best response rates in that particular tech company (Dastin, 2018). Thus, the defining problem here is that fairness is context-dependent, meaning there can be any number of definitions of fairness in hiring context and none can be intrinsically right or wrong, considering there lacks a widely accepted definition of "fair" (Barocas et al., 2019). These competing fairness notions—statistical parity (ensuring equal representation of groups) versus equalized odds (ensuring equal error rates across groups)—provide conflicting views of fairness that are very dependent on the context and the involved stakeholders (Chouldechova, 2017). As fairness metrics often must make one or more trade-offs (Holstein et al., 2019), these different definitions complicate to construct AI systems that do not violate fairness with respect to all of the conflicting fairness objectives. In addition, the legal and regulatory challenges are another huge hurdle, as the incorporation of AI into hiring practices has to contend with an evolving legal framework that often consists of the Equal Employment Opportunity (EEO) Act along with another set of non-discrimination law (Crawford, 2021). Yet, such fast-paced advancement of AI technologies is generally not matched by relevant legal developmental work, creating a regulatory gap that could hinder effective and equal use of AI in hiring. Organizations therefore are challenged by a need to accommodate both the benefits of AI adoption and the relevant mandates for non-discrimination and data privacy, and the absence of clear and widely accepted standards make compliance a complex endeavor (O'Neil, 2016). These challenges are compounded by the ethical conundrums of AI in hiring, as recruitment tools that are driven by AI can have unintended impacts on underrepresented communities and may inadvertently (Angwin et al., 2016) reinforce stereotypes, further end up ruling out qualified candidates from marginalized groups, especially women, racial minorities and disabled people. One of the ethical issues concern the potential that AI can learn to take discriminatory decisions, for example using features that are not directly in the model (i.e. transparent to human) but are correlated to protected characteristics such as gender or ethnic origin, simply because they were represented in the population data (Kamiran & Calders, 2012) Moreover, algorithm-based selection is famous for seeking algorithmic efficiency at the expense of human intuition and experience, and their extensive usage would render assessment of candidates by hiring professionals obsolete. a ii In this paradigm, decision-making by recruitment professional would be minimized and organisations would depend on their higher level from such simplistic and concise summarization that did not even closely resemble the real human (Lipton, 2016). Therefore, it is mandatory to keep questioning the design and deployment of AI tools towards ensuring that these systems do not just meet minimal legal and regulatory standards, but also conform to important social values, such as fairness, transparency, and accountability (Crawford, 2021) There is also a need for appropriate internal control and external audit of their use of AI to ensure that their AI systems are not inadvertently harming disadvantaged groups or unjustly excluding them from opportunities, which makes deploying AI in hiring complicated (Binns, 2018). The combination of limitations of responsible data, vague criteria for fairness definitions, legal and regulatory dilemmas, and ethical considerations therefore poses a challenge for AI adoption in hiring practices, and all aspects must be considered to mitigate any harms for organizations and candidates alike (Barocas et al., 2019; Mehrabi et al., 2021).

## VII.    Implications for Practice

Implications for practice of AI in hiring are clear and multi-faceted if organizations are to implement AI tooling in an equitable way, ensuring transparency and fairness in recruitment decisions whilst also promoting relevant equity, and relate primarily to AI governance, HR policies, workforce dynamics and training and education. Addressing AI governance in hiring is critical because AI systems must be designed, deployed, and monitored consistently with fairness principles, and organizations need to implement an explicit governance framework that sets standards for the design, deployment, and ongoing evaluation of AI systems (Holstein et al., 2019). These frameworks should enable transparency, explainability, and accountability in AI tools by embedding strong auditing, validation, and recalibration mechanisms into AI models to guarantee compliance with fairness objectives and legal and ethical standards [13]. Moreover, organisations ought to establish formal AI ethics committees or data governance bodies responsible for monitoring the impact of AI in hiring, ensuring periodic reviews of algorithmic equity and addressing bias or discrimination issues (Binns, 2018). Regarding implications for HR policy, Dastin (2018) pronounce that enforcing AI systems in terms of recruitment and hiring practices, HR departments should replace their policies and procedures to spring fairness and transparency while minimize the risk of bias, and that can be achieved by updating the hiring practices themselves to have clear and definite compelling requirements for Fairness in the content of the AI systems. For instance, one set of policies would rectify discrimination by setting fairness benchmarks at which AI models must perform prior to their use in recruitment decisions, utilizing metrics like statistical parity and equalized odds to ensure AI tools do not inherently favor some groups over others in hiring (Kamiran & Calders, 2012). In addition, HR departments can develop systems for continuous human-in-the-loop interventions and processes whereby human recruiters can review AI-driven decisions to ensure that fairness and diversity goals are achieved (Binns, 2018) and by being transparent with candidates about how they are being implemented in the hiring processes, organizations can build trust and reduce concerns about fairness (Binns, 2018). In terms of future workforce dynamics, the AI-in-hiring approach is likely to change the diversity of the workforce, the culture of the organization, and the diversity of people working together to an unprecedented extent, as the ability of AI to process large amounts of data and identify patterns in candidate qualifications has great potential to grow more objective decisions and reduce some patterns of biases that have traditionally shaped recruitment (Angwin et al., 2016). If designed and evaluated well, AI models have the potential to promote more diverse hiring results, by identifying talented candidates from historically underrepresented groups that are often disregarded during the standard hiring processes (Barocas et al., 2019). But it also represents a dangerous potential for entrenching biases in society if the datasets used to train AI tools are inherently flawed, further compounding inequalities in workforce diversity (O'Neil, 2016). This will lead organizations to precariously walk the tightrope between leveraging AI as a mechanism to increase their workforce diversity as well as at the risk of further amplifying historical biases and therefore build AI systems which can be aligned with DEI goals for a more dynamic and inclusive organizational culture (Mehrabi et al., 2021). Finally, training and education play an important role in successful AI integration into hiring, and there is a clear need for targeted training programs for both HR personnel and AI developers, that must holistically represent the ethical complexities and real-time operationalization of fairness in AI systems (Holstein et al., 2019). Training is also needed for other relevant stakeholders, such as those involved with HR, to understand how to interpret the recommendations made by AI and how to make the ultimate decisions using AI in a way that upholds fairness in recruitment practices (Kamiran & Calders, 2012). Additionally, organizations could fund cross-disciplinary training initiatives uniting data scientists with ethicists and HR practitioners in the creation of fairer, transparent systems that will ensure that all stakeholders are familiar with the powers and limits of AI tools for hiring, including how ethical issues can be resolved behind the scenes (Lipton, 2016). Together they are helping organizations create a culture of continuous learning and ethical responsibility to ensure AI systems are effective in hiring not only in a technical sense but correct with the broader societal values of fairness and equity (Barocas et al., 2019), paving the way towards more inclusive and justice-driven hiring in the future.

## VIII.    Conclusion

To sum up, there is an understanding of AI in the hiring process; from preparing the ground by bias mitigation and ultimately, that it is crucial to create balance between reducing bias and achieving fairness, so that AI can help make hiring processes fairer and more inclusive. The importance of balancing bias and fairness to create fair, accountable and inclusive AI systems ai not new, but the findings of this paper serve as a reminder of this very important process. It emerges that the use of the AI in Human Resource departments can amplify the efficiency of recruitment to a great extent; however, it is mandatory to eradicate high probabilities of experiencing bias at work before rolling this out (Khalid & Farooq, 2023). Given the continued improvement of AI technologies as well as the rapid adoption of hiring based AI, our future research directions include: new fairness algorithms that address many problems not solvable in a statistical sense, such as hybrid algorithms that minimize multiple fairness metrics (Kamiran, & Calders, 2012), and developing debiasing techniques that are more advanced than those that mitigate statistical bias through pre-processing and reduce the discrimination in data as much as

possible, including moving beyond data focused approaches and outcomes- focused approaches in traditional statistics (Friedler et al., 2019). Third, a key avenue for future research is the macro-level impact of AI-based hiring technologies as these might affect workforce composition, economic disparity and social stratification and whether they will effectively break entrenched discrimination or worsen existing power structures (O'Neil, 2016). Finally, studying the relationship between AI fairness and global DEI is pressing, especially if AI hiring tools are to work in concert with DEI initiatives, as fairness criteria in multiple cultural settings may considerably differ between regions and legal systems (Crawford, 2021). In conclusion, while AI provides the opportunity to improve fairness and accessibility in hiring, achieving those goals depends on the continued application of ethical practices such that fairness is neither an afterthought nor a prerequisite (Holstein et al., 2019); AI systems need to be designed to continuously assess their own performance in order to promote and ensure adherence to the values of fairness, transparency, and inclusivity in order to mitigate the risks of AI systems reinforcing existing inequalities that would otherwise lead to adverse economic, political, and social consequences and to ultimately deliver AI assisted recruitment that support a diverse and inclusive workforce in the future (Barocas et al., 2019).

# References

[1].    Ajunwa, I. (2023). *The quantified worker: Law and technology in the modern workplace*. Cambridge University Press.
[2].    Ali, S., Roberts, H., Oswald, M., & Cowls, J. (2022). Algorithmic feedback loops and the entrenchment of inequality: An exploratory study. *AI and Society, 37*(4), 1041–1054. https://doi.org/10.1007/s00146-021-01164-0
[3].    Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016). Machine bias. *ProPublica*. https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing
[4].    Babuta, A., Oswald, M., & Roff, H. M. (2021). Artificial intelligence and the future of recruitment. *Royal United Services Institute (RUSI)*. https://rusi.org
[5].    Barocas, S., Hardt, M., & Narayanan, A. (2023). *Fairness and machine learning: Limitations and opportunities*. fairmlbook.org.
[6].    Berk, R., Heidari, H., Jabbari, S., Kearns, M., & Roth, A. (2021). Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research, 50*(1), 3–44. https://doi.org/10.1177/0049124118782533
[7].    Biega, A. J., Gummadi, K. P., & Weikum, G. (2022). Equity of attention in ranking: A survey. *ACM Transactions on Information Systems, 40*(2), 1–49. https://doi.org/10.1145/3510423
[8].    Binns, R. (2018). Fairness in machine learning: Lessons from political philosophy. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 149–159. https://doi.org/10.1145/3287560.3287593
[9].    Binns, R. (2020). On the apparent conflict between individual and group fairness. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 514–524. https://doi.org/10.1145/3351095.3372864
[10].   Binns, R. (2023). Fairness in machine learning: Lessons from political philosophy. *Philosophy & Technology, 36*(1), 1–25. https://doi.org/10.1007/s13347-022-00551-w
[11].   Binns, R., Veale, M., Van Kleek, M., & Shadbolt, N. (2022). "It's reducing a human being to a percentage": Perceptions of justice in algorithmic decisions. *CHI Conference on Human Factors in Computing Systems*, 1–14. https://doi.org/10.1145/3491102.3517522
[12].   Chandrasekharan, E., Srinivasan, A., & Narayanan, A. (2022). Algorithmic exclusion: How automated hiring tools discriminate against people with disabilities. *Disability Studies Quarterly, 42*(4). https://doi.org/10.18061/dsq.v42i4.8932
[13].   Chouldechova, A., & Roth, A. (2020). A snapshot of the frontiers of fairness in machine learning. *Communications of the ACM, 63*(5), 82–89. https://doi.org/10.1145/3376898
[14].   Contissa, G., Sartor, G., & Lagioia, F. (2023). Algorithmic hiring and EU non-discrimination law: How to align innovation and fundamental rights. *European Journal of Risk Regulation, 14*(1), 115–133. https://doi.org/10.1017/err.2022.41
[15].   Cowgill, B., Dell'Acqua, F., & Deng, S. (2021). Biased programmers? Or biased data? A field experiment in operationalizing AI ethics. *Management Science, 69*(2), 586–614. https://doi.org/10.1287/mnsc.2021.4090
[16].   Cowls, J., Tsamados, A., Taddeo, M., & Floridi, L. (2021). The AI gambit: Leveraging artificial intelligence to combat bias in hiring. *Philosophy & Technology, 34*(4), 1279–1296. https://doi.org/10.1007/s13347-021-00463-9
[17].   Crawford, K. (2021). *Atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press.
[18].   Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 214–226. https://doi.org/10.1145/2090236.2090255
[19].   Eubanks, V. (2018). *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press.
[20].   Friedler, S. A., Scheidegger, C., Venkatasubramanian, S., Choudhary, S., Hamilton, E. P., & Roth, D. (2021). A comparative study of fairness-enhancing interventions in machine learning. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 329–341. https://doi.org/10.1145/3287560.3287589
[21].   Green, B. (2022). The flaw in technological solutionism: Lessons from algorithmic hiring. *Technology in Society, 68*, 101912. https://doi.org/10.1016/j.techsoc.2021.101912
[22].   Grosz, M., Kalkman, J. P., & Widlak, E. (2019). The hidden bias in AI recruiting software: A sociotechnical perspective. *Policy & Internet, 11*(4), 418–437. https://doi.org/10.1002/poi3.217
[23].   Hoffmann, A. L. (2022). Terms of inclusion: Data, discourse, violence. *New Media & Society, 24*(2), 340–359. https://doi.org/10.1177/1461444820912543
[24].   Holstein, K., Wortman Vaughan, J., Daume III, H., Dudik, M., & Wallach, H. (2019). Improving fairness in machine learning systems: What do industry practitioners need? *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 1–16. https://doi.org/10.1145/3290605.3300830
[25].   Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence, 1*, 389–399. https://doi.org/10.1038/s42256-019-0088-2
[26].   Kim, P. T. (2021). Data-driven discrimination at work. *William & Mary Law Review, 58*(3), 857–936.
[27].   Kim, P. T. (2023). Algorithmic employment discrimination. *Fordham Law Review, 91*(3), 703–738.
[28].   Kleinberg, J., Mullainathan, S., & Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*.
[29].   Köchling, A., & Wehner, M. C. (2020). Discriminated by an algorithm: A systematic review of discrimination and fairness in algorithmic decision-making. *Business Research, 13*, 795–848. https://doi.org/10.1007/s40685-020-00134-w
[30].   Lee, M. K., Kusner, M., & Self, J. (2023). Beyond fairness: Multistakeholder tensions in algorithmic hiring. *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*, 1–15. https://doi.org/10.1145/3544548.3581364

[31]. Lee, N. T., & Singh, A. (2021). Algorithmic fairness in practice: Principles, challenges, and opportunities. *Brookings Institution Report*. https://www.brookings.edu

[32]. Lukyanenko, R., Wieringa, R., & Willcocks, L. (2023). Designing accountable algorithms in human resource systems: A multistakeholder approach. *Information Systems Journal, 33*(2), 327–351. https://doi.org/10.1111/isj.12356

[33]. Marda, V., & Narayan, S. (2021). Algorithmic accountability in India: An emerging framework. *Data & Society Research Institute*. https://datasociety.net

[34]. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019). Model cards for model reporting. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 220–229. https://doi.org/10.1145/3287560.3287596

[35]. Mittelstadt, B. D. (2022). Principles alone cannot guarantee ethical AI. *Nature Machine Intelligence, 4*, 1–8. https://doi.org/10.1038/s42256-021-00400-x

[36]. Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science, 366*(6464), 447–453. https://doi.org/10.1126/science.aax2342

[37]. Passi, S., & Barocas, S. (2019). Problem formulation and fairness. *Proceedings of the 2019 ACM Conference on Fairness, Accountability, and Transparency*, 39–48. https://doi.org/10.1145/3287560.3287567

[38]. Raji, I. D., Bender, E. M., & Denton, E. (2022). AI audits in practice: Risk, responsibility, and redress. *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 25–35. https://doi.org/10.1145/3531146.3533208

[39]. Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., ... & Barnes, P. (2020). Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing. *Proceedings of the 2020 ACM Conference on Fairness, Accountability, and Transparency*, 33–44. https://doi.org/10.1145/3351095.3372873

[40]. 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

[41]. Selbst, A. D., Boyd, D., Friedler, S. A., Venkatasubramanian, S., & Vertesi, J. (2019). Fairness and abstraction in sociotechnical systems. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 59–68. https://doi.org/10.1145/3287560.3287598

[42]. Shankar, S., D'Ignazio, C., & Gebru, T. (2023). Designing for equity in algorithmic hiring: A feminist data ethics perspective. *AI and Ethics, 3*(2), 89–105. https://doi.org/10.1007/s43681-022-00235-0

[43]. Upadhyay, A. K., & Khandelwal, K. (2018). Applying artificial intelligence: Implications for recruitment. *Strategic HR Review, 17*(5), 255–258. https://doi.org/10.1108/SHR-07-2018-0051

[44]. Veale, M., & Zuiderveen Borgesius, F. J. (2021). Demystifying the "right to explanation" in the GDPR. *Computer Law & Security Review, 36*, 105367. https://doi.org/10.1016/j.clsr.2019.105367

[45]. Whittaker, M., Alper, M., Crawford, K., Dobbe, R., Fried, G., Kaziunas, E., ... & West, S. M. (2018). *AI Now Report 2018*. AI Now Institute. https://ainowinstitute.org/AI_Now_2018_Report.pdf

[46]. Whittlestone, J., Nyrup, R., Alexandrova, A., & Cave, S. (2021). The role and limits of principles in AI ethics: Towards a focus on tensions. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 195–201. https://doi.org/10.1145/3461702.3462565

[47]. Wilson, C., Hines, K., & Marquez, S. (2024). Long-term fairness in hiring algorithms: Measuring cumulative impact. *Journal of Artificial Intelligence Research, 80*, 181–202. https://doi.org/10.1613/jair.1.13900

[48]. Wong, R. Y., & Webb, H. (2023). Reframing algorithmic fairness: Towards relational accountability. *ACM Transactions on Human-Computer Interaction, 30*(1), 1–28. https://doi.org/10.1145/3597493

[49]. Yapo, A., & Weiss, J. (2022). Ethical tensions in algorithmic hiring: Efficiency versus fairness. *Journal of Business Ethics, 179*(2), 375–392. https://doi.org/10.1007/s10551-020-04630-w