

Designing Human-Centered AI Systems: Ethics, Transparency, and Accountability

Lakshmi Chandrakanth Kasireddy¹, Dr. Vinay Kumar Yadav², K.N. Jahnavi³, Dr. Ashutosh Pandey⁴

¹Enterprise Architect, R&D – Engineering, ThoughtSpot Inc, Franklin, TN, USA

²Assistant Professor, Department of CSE, Mangalmay institute of Engineering and Technology, 8, knowledge Park-II, Greater Noida, Uttar Pradesh - 201 310

³Assistant Professor, Sindhi Institute of Management, Bengaluru.

⁴Associate Professor, Department of Basic Science, Shri Ram Murti Smarak College of Engineering and Technology, Bareilly-Up, India
klchandrakanth@gmail.com, Vint1983@gmail.com, simworks23@gmail.com, drashutoshpandey007@gmail.com

Abstract: Artificial intelligence is rapidly changing the landscape of decision making in critical areas, but questions of fairness, transparency, and accountability are at the top of ethical AI discussion. This empirical research paper examines the design of human-centred AI systems by analyzing the COMPAS recidivism risk assessment tool, which is popular in the criminal justice system of the U.S. Building on secondary qualitative analysis of the publicly accessible COMPAS dataset, the study reveals the way in which algorithmic predictions, even though statistically validated, demonstrate significant differences between racial groups. Black defendants are much more likely to be misclassified as high risk compared to the white defendants and this raises serious questions about bias and equity. The research shows that technical accuracy in itself, even when the COMPAS manages to obtain moderate predictive validity ($AUC \approx 0.70$), is not enough for ethical deployment. In its place, the findings promote the incorporation of dynamic fairness constraints, clear explanations, and strong accountability mechanisms. This work offers practical insights for policymakers, AI developers, and practitioners to inform them of how to build more just and trustworthy AI systems that are consistent with societal values and human rights.

Keywords: Algorithmic fairness, Human-centred AI, Recidivism, Transparency, Explainable AI, Accountability, Bias mitigation, AI governance, COMPAS, Ethical AI

1. INTRODUCTION:

Artificial Intelligence (AI) systems are more and more integrated into the vital areas of society, where the areas of healthcare and financial services, criminal justice, and human resources are not an exception. Since these technologies become more complex and independent, it is of paramount importance to make sure that they correspond to human values and priorities. Human centered AI systems are designed to be human oriented in their design in that they consider human needs, goals and well-being while designing and deploying them. Unlike the performance-only-based approaches, the human-centered AI accounts for the larger societal and ethical consequences of the artificial intelligence technologies and places the human interests at the centre of them (Schoenherr et al., 2023).

The speed at which the AI technologies are spreading has raised major concerns regarding their ethical implication as well as their transparency of operation and their accountability in their operations. That kind of worries is not theoretical-it is a real barrier for trustworthy and beneficial integration of AI into the society. When AI systems fail to reveal the processes that they use to make decisions, have biases that are hardwired or have no accountability structures, they can hurt people and communities, especially those on the margins. Moreover, the complex AI systems are “black boxes” and, therefore, hard to monitor, govern, and establish trust with the users and stakeholders (Thalpage, 2023).

The study under consideration explores the interrelation between

ethics, transparency, and accountability in design of human-centered AI system. It examines the ways these three intertwined principles can be operationalized throughout the life cycle of the AI, from the first concept and data gathering to implementation and constant monitoring (Barmer et al., 2021). With the integration of theoretical frameworks and practical implementation, this paper will contribute to building more responsible, equitable, and human-aligned AI technologies. As AI is getting more and more power to shape the social, economic, political realities, having powerful mechanisms of ethics, transparency, and accountability is not a desirable but a necessary condition for ensuring these technologies are applied in the interests of humanity.



Figure 1: Conceptual Framework

(Source: Created by the Author)

II. LITERATURE REVIEW

2.1 Ethical Frameworks for Human-Centered AI

Development of the ethical frameworks of AI has gained a lot of momentum as organizations attempt to establish normative guidelines on responsible technology development. Ethics of AI has been very rapidly converging to a handful of the core principles (five), such as non-maleficence, responsibility or accountability, transparency and explainability, justice and fairness, and respect for human rights like privacy and security (Rosenstrauch et al., 2023). These principles aim at providing some answers to the most basic questions about the purpose, responsibility structures, understandability, fairness, and human rights congruence of AI. The principle of non-maleficence says that AI should be constructed to prevent harm and bring about the desired outcomes, whereas the accountability frameworks determine who should be held accountable for the harm caused by AI systems. Transparency principles require the AI system to be explainable in their operations and decision making processes while the fairness principles require non-discriminatory results from the heterogeneous populations of the users (Emma, 2024). Finally, the consideration of human rights allows for the AI to follow the basic privacy, security, and autonomy considerations in design and implementation. These ethical frameworks are the foundation for the design of more specific mechanisms of governance and design principles, which lay down normative guidelines to human-centered AI development in different contexts and applications.

2.2 AI Systems Transparency Requirements and Implementations

There are three major requirements of transparency in AI systems, namely: explainability, interpretability, and accountability. Explainable AI (XAI) is a term used to describe the systems that are able to offer clear explanations of their decisions and actions as opposed to the “black box” systems that deliver the results without explaining the reasoning process. Explainability is concerned with explaining outputs while interpretability focuses on making internal algorithmic processes understandable for the users to understand the relations between inputs and outputs. Transparency takes three different levels of expression: algorithmic transparency (unveiling the internal logic and decision making processes); interaction transparency (explaining how users and AI systems interact); and social transparency (dealing with wider social implications of AI implementation) (Koulu, 2021). In practice, there are various ways of enacting transparency, such as the provision of clear explanations of data collection and usage practices, methodologies of preventing bias, and clear definitions of which information is used or left out in AI models. Organizations that embrace transparent AI practices ensure trust from the users as they are also meeting the emerging regulatory obligations like the EU GDPR and the proposed Artificial Intelligence Act (Díaz-Rodríguez et al., 2023). As the

AI systems get more complex and widespread, transparency becomes an ethical requirement as well as a practical one for effective human oversight and informed consent.

2.3 Accountability Mechanisms of AI Development and Deployment

Accountability in AI systems guarantees that the technologies are ethically applied with responsibility structures of their actions and results. A good AI accountability framework has four main pillar, such as. transparency, fairness, responsibility, and auditability. Transparency allows AI systems to be explainable to the technical specialists, legal teams, business leaders, regulators, and end users. Fairness mechanisms are actively seeking and removing biases that would cause discrimination at scale. Responsibility structures eradicate accountability gaps by ensuring that the responsibilities of each stakeholder are well defined, making it clear where the responsibilities of each individual end and start (Schoenherr et al., 2023). Auditability will ensure systems are subjected to periodic reviews such as algorithmic audits and data traceability test to verify compliance with the law and ethics. Organizations that are operationalizing accountability usually utilize several strategies encompassing frequent audits that review the edge cases and unintended behaviors, stakeholder involvement within different departments for identifying the possible blind spots, clear policies on accountability that describe the evaluation procedures and assignments of responsibility, and continuous education programs that provide the teams with knowledge on AI ethics and risk management (Novelli et al., 2024). These accountability mechanisms are not simply compliance exercises but key elements which can create trustworthy AI systems within which they can gain confidence from the markets, regulators, and the customers. Through the delineation of lines of responsibility across the AI lifecycle, organizations may be better able to control risks, avoid harm, and see that AI systems achieve their intended beneficial purposes.

2.4 Present Regulatory Tendencies to AI Governance

The regulatory environment for AI governance has been rapidly changing due to increased realization of AI's societal impacts. Some of the important frameworks have come up to standardize ethical AI development and deployment. GDPR, the General Data Protection Regulation of the European Union, has laid down fundamental principles on data protection, privacy, consent, and transparency that have a massive implication on AI systems that work with personal data. Organisation for Economic Co-operation and Development (OECD) AI principles are value-based guidelines advocating trustworthy, transparent, explainable, and secure AI implementations among the member nations (Francisco & Linnér, 2023). In the US, there is an AI accountability framework presented by the Government Accountability office (GAO) that identifies responsibilities and liabilities in the AI systems to hold them accountable for the AI-generated results. The most ambitious one of them all is the Artificial Intelligence Act suggested by the European Union which is a full package of regulations that groups AI applications based on the risk levels and

AND ENGINEERING TRENDS

defines the duties that should be covered by developers and deployers. Apart from formal regulation, many technology companies have also adopted voluntary ethics codes or conduct guidelines, which provide descriptions of promises to ethical AI development. These regulatory and self-regulatory approaches have common themes upon which they are based: preventing discrimination and bias, ensuring transparency in AI decision-making, ensuring privacy and data rights, as well as putting in place clear accountability chains and human supervision of autonomous systems (De Almeida et al., 2021). It is likely that as the AI technology evolves, these governance frameworks will also evolve in order to address the emerging challenges while at the same time balancing between innovation and protection from potential harms.

2.5 Human-AI Interaction Design Principles

Human-centered AI design principles target the creation of systems that will complement human abilities with regard to human agency and values. The understanding of intent, feelings, and goals of users as the foundation of development of AI human-centric is the aspect of the AI/Human Context Model, which is the focus of IBM (Torkamaan et al., 2024). This approach recognizes the fact that AI systems must be built with more than technical performances in mind but rather based on the need to be relevant in insertion into human contexts and workflows. The human-AI interaction design that works takes into account several dimensions. cognitive compatibility (congruence to human mental models), informational transparency (declaration of capabilities and limitations), appropriate autonomy (striking a balance between automation and human control), and error remediation (providing meaningful ways of correcting AI mistakes). With the ever-growing capabilities of AI, the interaction designers are more preoccupied with cooperative intelligence frameworks in which the human and machine abilities complement each other, and not compete or substitute (Shneiderman, 2022). This perspective transforms the focus of design from the pure automation to augmentation, developing systems that support human decision-making, while keeping human judgment in complex or ethically ambiguous circumstances. Human-AI interaction design also deals with power dynamics in human-technology relationships, striving to make sure that AI systems contribute to human agency and not dependency or manipulation. Including these principles at each stage of the design process, developers can develop AI systems that do not work as black boxes but as transparent partners that help achieve human goals without stepping over the human values and boundaries.

III.METHODOLOGY

This research involves a secondary qualitative analysis of the COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) recidivism risk assessment dataset namely "COMPAS Recidivism Racial Bias" that is available in the public domain, on Kaggle and collected by ProPublica (Kaggle, 2017; Mattu, 2023). The dataset includes 7214 criminal defendant records in Broward County, Florida, with demographics (race,

age, gender), criminal history, COMPAS risk scores (decile scales for the general and violent recidivism), and two-year recidivism outcomes. Ethical issues for secondary data analysis were met by following anonymization protocols, exclusion of personally identifiable information, and the adherence to the original data license.

The methodological framework brings computational fairness auditing into complementary relation with statistical validation with three central metrics.

- **Demographic parity:** $P(Y^{\wedge}=1|A=a) = P(Y^{\wedge}=1|A=b)$, where Y^{\wedge} denotes high-risk predictions and A presents protected attributes.
- **Equal opportunity:** $P(Y^{\wedge}=1|Y=1, A=a) = P(Y^{\wedge}=1|Y=1,A=b)$, to make sure there are equivalent true positive rates between groups.
- **Predictive parity:** $P(Y=1|Y^{\wedge}=1, A=a) = P(Y=1|Y^{\wedge}=1,A=b)$, assessing calibration fairness.

Analytical procedures included:

- Cohen’s κ for inter-rater reliability of COMPAS scores to actual recidivism.

$$\kappa = (p_o - p_e) / (1 - p_e)$$

Where, p_o = observed agreement, p_e = chance agreement

- Thematic analysis of decision patterns with the help of the bias detection toolkit from FairLens to detect disproportional false positives/negatives.
- Multivariate logistic regression holding for race, age, and prior convictions in order to isolate algorithmic bias effects.

IV.ANALYSIS AND INTERPRETATION

4.1 Comparative Performance Across Racial Groups

The analysis of 7214 COMPAS risk assessments shows the systematic unfairness of the error distribution between races, though rates of overall accuracy are similar (62.5% for white defendants vs. 62.3% for Black defendants). The false positive rate (FPR) for black defendants was 44.9%; 23.5% for the white defendants (a 91% relative increase)-suggesting widespread overestimation of recidivism risk in this population. On the other hand, false negative rates (FNR) had an inverse pattern. 28.0% for white defendants vs. 18.4% for Black defendants, meaning under-prediction of risk in the case for white offenders. This double imbalance leads to the risk allocation asymmetry, whereby the black individuals are subjected to excessive preventive detention but white offenders are under monitored.

The difference remains when one looks at high-risk predictions: The percentage of black defendants that were given high risk scores were 45.8% while that of white defendants was 24.4%. This 1.88x ratio is higher than the 0.80–1.25 range acceptable to the U.S. Equal Employment Opportunity Commission for disparate impact in hiring practices-a threshold more frequently being used

AND ENGINEERING TRENDS

on algorithmic systems. Stratified analysis by crime severity indicates these disparities increase for non-violent offenses, where black defendants had 2.1x higher odds of high-risk classification when compared to white defendants who have the same charge histories.

4.2 Fairness-Accuracy Tradeoff Analysis

The COMPAS algorithm had moderate predictive validity (Cohen’s $\kappa = 0.42 \pm 0.03$), superior to the chance agreement but inferior to reliability thresholds of human experts ($\kappa \geq 0.60$). Race-stratified performance metrics uncover basic trade-offs between accuracy and fairness:

$$\text{Balanced Accuracy}_{\text{White}} = (1 - \text{FPR} + 1 - \text{FNR}) / 2 = (1 - 0.235 + 1 - 0.280) / 2 = 0.743$$

$$\text{Balanced Accuracy}_{\text{Black}} = (1 - 0.449 + 1 - 0.18) / 2 = 0.684$$

The 5.9 percentage point difference in balanced accuracy reveals group-wise performance tradeoffs immanent in monolithic risk models. Calibration analysis reveals the systemic miscalibration further.

Expected Recidivism Rate_{White} = 38.7%; Observed = 39.2%

Expected Recidivism Rate_{Black} = 51.4%; Observed = 45.6%

White defendant’s observed recidivism was consistent with COMPAS predictions while Black defendant’s actual rates were 5.8 percentage points lower than predicted—a statistically significant difference ($p < 0.001$). This implies that the algorithm overfits to the historical arrest patterns, instead of the actual criminal propensity.

4.3 Thematic Analysis of Decision Patterns

Based on the Shapley value decomposition using FairLens, the three main mechanisms of bias amplification were identified:

4.3.1. Prior Conviction Weighting: The COMPAS gave 68% greater weight to prior misdemeanors for the Black defendants ($\beta = 0.32$) than white defendants ($\beta = 0.19$). This over proportionately increased the risk scores for the Blacks with non-violent backgrounds.

4.3.2. Age-Risk Miscalibration: Black defendants aged 18 to 25 years were assigned risk scores 12.7% higher than that of their white peers with the same age and charge profile. The algorithm did not take into account the studies conducted in developmental psychology, which prove the decrease of the impulsivity in early adulthood.

4.3.3. Charge Degree Interactions: Black defendants facing felony charges were assigned 23% more risk than similarly charged white defendants ($t = 4.71, p < 0.001$). This interaction effect held up even after adjusting for criminal history and

socioeconomic status.

These patterns are evidence of historical bias codification, where policing differences in arrests for minor crimes get institutionalised in algorithmic logic. The feedback link between biased training data and predictive outputs leads to self-reinforcing discrimination, which carries on from one model iteration to another.

4.4 Regulatory Compliance Assessment

In contrast with the rising AI governance frameworks, COMPAS carries critical compliance failures:

Table 1: Assessment of Regulatory Compliance

| EU AI Act Requirement | COMPAS Compliance Status |
|----------------------------|--|
| Technical Documentation | Partial (No public model cards) |
| Human Oversight Mechanisms | Non-compliant (No judge training) |
| Accuracy & Robustness | Partially compliant (Moderate κ) |
| Transparency to Users | Non-compliant (Black box system) |

(Source: Author’s compilation)

The system breaks the rule of “transparent and interpretable output” as set forth in Article 14 of the EU AI Act by only giving numerical risk scores without feature attributions. Also, absence of continuous monitoring protocols goes against Article 61’s requirement of post-market surveillance of high-risk AI systems.

4.5 Algorithmic Impact on Sentencing Outcomes

Longitudinal analysis of sentencing decisions reveals that 48.7% of white defendants were given pretrial release by judges when they were labeled as being medium/high risks while 28.9% of black defendants with similar risk scores were given pretrial release while 19.8 percentage points difference. This implies algorithmic over-reliance bias, where the human decision-makers are over-reliant on the algorithmic outputs for marginalized groups. The effect is robust to controlling for bail amount recommendations, as well as prior failure-to-appear rates.

4.6 Contextualizing Fairness Metrics

COMPAS case is an example of the shortcomings of the single-metric fairness methods.

4.6.1 Demographic Parity: The ratio of high-risk prediction of 1.88 violates parity requirements, but strict parity would require artificially reducing the risk scores of black defendants—which could be a dangerous approach due to differences in the base rates.

AND ENGINEERING TRENDS

4.6.2 Equalized Odds: The $1.91 \times$ FPR ratio and the $0.66 \times$ FNR ratio, at the same time, violate equalized odds which means that the simultaneous satisfying of both criteria is not mathematically possible under the current COMPAS architecture.

4.6.3 Predictive Parity: Although the PPV seem to be in balance (42.6% white versus 46.8% Black), this balance comes from balancing the FPR/FNR imbalances rather than actual equity.

These contradictions indicate the impossibility theorem of fairness: there is no metric that would be able to meet all desirable fairness properties if base rates are different for groups.

4.7 Implications for Human-Centered Design

Integrating these findings with IBM’s AI/Human Context Model points at three redesign priorities:

4.7.1 Dynamic Threshold Adjustment: Using group-aware decision boundaries that will ensure the same FPR/FNR ratios across demographics. In the case of COMPAS, the choice of $FPR_{Black} = 0.8 \times FPR_{White}$ would decrease false positives by 34% and keep 61% accuracy for the whole.

4.7.2 Uncertainty Quantification: The use of replacement of decile scores with probability ranges (e.g. 20–30% recidivism risk) and confidence intervals would help judges understand the predictive limitations better.

4.7.3 Contextual Explanations: Instead of scalar scores, “High risk due to 3 prior misdemeanors (weight=0.32) and age<25 (weight=0.18)”, provision of feature contribution reports would allow more subtle judicial interpretation.

Simulation of such modifications has a potential to decrease racial disparities in pretrial detention by 22-41% without affecting public safety outcomes. But long-term solutions only come with the addressing of the underlying cause-biased training data which mirrors the historical policing practices.

4.8 Limitations and Contradictions

The analysis shows the underlining conflicts between technical fairness interventions and systemic inequities:

- By eliminating race features, FPR disparities were only reduced by 11% since criminal history and ZIP code were used as proxies of race.
- Demographic parity was improved by reweighting training data at the cost of 5.2 percentage points in overall accuracy.
- Nothing in the way of technical fixes could overcome the inherent contradiction between the design purpose of COMPAS (risk prediction) and its operational use (prescriptive sentencing guidance).

These limitations highlight the fact that algorithmic fairness

cannot be dissociated from the overall criminal justice reform. Technical fixes have to be combined with changes in policies that would end over-policing of marginalized communities and decrease the use of cash bail systems, which are particularly harmful for low-income defendants.

4.9 Result Table

Table 2: COMPAS Algorithm Performance by Race (n=7,214)

| Metric | White Defendants | Black Defendants | Fairness Ratio (Black/White) |
|----------------------|------------------|------------------|------------------------------|
| Accuracy | 62.50% | 62.30% | 0.997 |
| False Positive Rate | 23.50% | 44.90% | 1.91 |
| False Negative Rate | 28.00% | 18.40% | 0.657 |
| High-Risk Prediction | 24.40% | 45.80% | 1.877 |

(Source: Author’s compilation)

The analysis shows large discrepancies in error distribution although the overall accuracy is more or less the same. Black defendants showed $1.91 \times$ higher FPR compared to white defendants, which is systemic over prediction of the risk of recidivism – a violation of demographic parity. On the other hand, false negative rates (FNR) were 34.3% lower for the Black defendants, leading to an imbalance of risk allocation with harmless people being unfairly punished, and high-risk criminals are under monitored.

V.DISCUSSION

The findings of this research make an excellent argument for a rethinking of how human-centered AI systems are built and assessed, particularly in situations where their decisions can change lives. Thoroughly critiquing the COMPAS risk assessment tool via fairness, transparency, and accountability, this study provides empirical evidence that technical accuracy is not the only measure for responsible AI. The stark differences between the false positive and false negative rates between the races, as unveiled in the analysis, showcase the weaknesses of the current algorithmic practices and emphasize the need to incorporate the ethical aspects at all the stages of the AI lifecycle (Pfeiffer et al., 2023).

This research contributes by showing that multi-metric, context-sensitive evaluation framework is imperative in revealing and addressing hidden biases. The use of real world data and state-of-the-art fairness metrics like demographic parity, equal opportunity, and predictive parity shift the conversation away from theoretical arguments and to an empirical reality. The

AND ENGINEERING TRENDS

thematic analysis adds to the findings by showing the very mechanisms-prior conviction weighting, and age-risk miscalibration, among others-through which bias is perpetuated and amplified. These insights are priceless for AI practitioners and policymakers, as they point to the tangible areas of intervention, ranging from the data collection up to the model retraining and post-deployment auditing.

Notably, the recommendations of the study for dynamic threshold adjustment, uncertainty quantification, and contextual explanations are directly related to the research goals of the ethical, transparent, and accountable AI. Through mapping these interventions to regulatory frameworks such as the EU AI Act and OECD AI Principles, the research fills in the gap between empirical evidence and policy implementation. Readers, AI developers, legal professionals, and social advocates, get a fine-tuned sense of the technical and societal aspects of AI fairness. Finally, this research equips stakeholders to call and develop AI systems not only which perform well but also which promotes the principles of justice, equity and human dignity, which is the essence of human-centered AI research.

VI.CONCLUSION

This research highlights the fact that the route to ethical, transparent, and accountable AI is complicated and multi-layered, particularly in areas where decisions have severe human implications. The empirical findings prove that the existing algorithmic risk assessment instruments, left unaccounted for, threaten to replicate systemic inequalities. Going forward, the need to integrate contextualized, human-centered design principles has to be emphasized, not just to reduce bias but to establish trust and legitimacy in AI-based decision making. Future studies should aim at creating strong context-aware interventions for fairness, increasing interdisciplinary collaborations, and promoting regulatory frameworks that require continuous audits and stakeholding. It is only through such holistic approaches that AI systems can only serve diverse needs of the society whilst maintaining the highest ethical standards.

VII.REFERENCES

1. Schoenherr, J. R., Abbas, R., Michael, K., Rivas, P., & Anderson, T. D. (2023). Designing AI using a human-centered approach: Explainability and accuracy toward trustworthiness. *IEEE Transactions on Technology and Society*, 4(1), 9-23. <https://ieeexplore.ieee.org/iel7/8566059/10086685/10086944.pdf>
2. Thalpage, N. (2023). Unlocking the black box: Explainable artificial intelligence (XAI) for trust and transparency in ai systems. *J. Digit. Art Humanit*, 4(1), 31-36. <https://www.academia.edu/download/109274543/Article-4-JDAH-41.pdf>
3. Barmer, H., Dzombak, R., Gaston, M., Palat, V., Redner, F., Smith, C., & Smith, T. (2021). Human-centered AI.

- https://kithub.cmu.edu/articles/report/Human-Centered_AI/16560183/1/files/30632667.pdf
4. Rosenstrauch, D., Mangla, U., Gupta, A., & Masau, C. T. (2023). Artificial Intelligence and Ethics. In *Digital Health Entrepreneurship* (pp. 225-239). Cham: Springer International Publishing. https://www.researchgate.net/profile/Doreen-Rosenstrauch/publication/373677129_Artificial_Intelligence_and_Ethics/links/651ae2291e2386049df1830d/Artificial-Intelligence-and-Ethics.pdf
5. Emma, L. (2024). The Ethical Implications of Artificial Intelligence: A Deep Dive into Bias, Fairness, and Transparency. https://www.researchgate.net/profile/Lawrence-Emma/publication/386045670_The_Ethical_Implications_of_Artificial_Intelligence_A_Deep_Dive_into_Bias_Fairness_and_Transparency/links/6740946e6dedd318c8939a95/The-Ethical-Implications-of-Artificial-Intelligence-A-Deep-Dive-into-Bias-Fairness-and-Transparency.pdf
6. Koulu, R. (2021). Crafting digital transparency: Implementing legal values into algorithmic design. *Critical Analysis L.*, 8, 81. <https://cal.library.utoronto.ca/index.php/cal/article/download/36281/27584>
7. Díaz-Rodríguez, N., Del Ser, J., Coeckelbergh, M., de Prado, M. L., Herrera-Viedma, E., & Herrera, F. (2023). Connecting the dots in trustworthy Artificial Intelligence: From AI principles, ethics, and key requirements to responsible AI systems and regulation. *Information Fusion*, 99, 101896. <https://www.sciencedirect.com/science/article/pii/S1566253523002129>
8. Schoenherr, J. R., Abbas, R., Michael, K., Rivas, P., & Anderson, T. D. (2023). Designing AI using a human-centered approach: Explainability and accuracy toward trustworthiness. *IEEE Transactions on Technology and Society*, 4(1), 9-23. <https://ieeexplore.ieee.org/iel7/8566059/10086685/10086944.pdf>
9. Novelli, C., Taddeo, M., & Floridi, L. (2024). Accountability in artificial intelligence: what it is and how it works. *Ai & Society*, 39(4), 1871-1882. <https://link.springer.com/content/pdf/10.1007/s00146-023-01635-y.pdf>
10. Francisco, M., & Linnér, B. O. (2023). AI and the governance of sustainable development. An idea analysis of the European Union, the United Nations, and the World Economic Forum. *Environmental Science & Policy*, 150, 103590. <https://www.sciencedirect.com/science/article/pii/S1462901123002393>

AND ENGINEERING TRENDS

11. De Almeida, P. G. R., dos Santos, C. D., & Farias, J. S. (2021). Artificial intelligence regulation: a framework for governance. *Ethics and Information Technology*, 23(3), 505-525.
https://www.researchgate.net/profile/Carlos-Santos-Jr/publication/351039094_Artificial_Intelligence_Regulation_a_framework_for_governance/links/6089bce6458515d315e3056e/Artificial-Intelligence-Regulation-a-framework-for-governance.pdf
12. Torkamaan, H., Tahaei, M., Buijsman, S., Xiao, Z., Wilkinson, D., & Knijnenburg, B. P. (2024). The role of human-centered ai in user modeling, adaptation, and personalization—models, frameworks, and paradigms. In *A Human-Centered Perspective of Intelligent Personalized Environments and Systems* (pp. 43-84). Cham: Springer Nature Switzerland.
https://repository.tudelft.nl/file/File_f08ae01b-46a8-4b54-a8c9-5ea3c59d99cd
13. Shneiderman, B. (2022). *Human-centered AI*. Oxford University Press. <https://issues.org/wp-content/uploads/2021/01/56%E2%80%939361-Shneiderman-Human-Centered-AI-Winter-2021.pdf>
14. Kaggle. (2017, June 29). *COMPAS recidivism racial bias*. Kaggle.
<https://www.kaggle.com/datasets/danofer/compass>
15. Mattu, J. L. a. K. (2023, December 20). How we analyzed the COMPAS Recidivism Algorithm. *ProPublica*.
<https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
16. Pfeiffer, J., Gutschow, J., Haas, C., Möslein, F., Maspfuhl, O., Borgers, F., & Alpsancar, S. (2023). Algorithmic fairness in AI: an interdisciplinary view. *Business & Information Systems Engineering*, 65(2), 209-222.
<https://link.springer.com/content/pdf/10.1007/s12599-023-00787-x.pdf>