

# **Governance of Fairness in Data Annotation: A Study of Dual Paths in Law and Ethics**

**Shen Rongrong, Zhao Yaqin**

School of Law, Shandong Jianzhu University

## **Abstract**

As the foundational "meta-labor" for training artificial intelligence models, data annotation processes exhibiting fairness deficits and bias implantation have become significant sources of algorithmic discrimination, directly impeding the implementation of trustworthy artificial intelligence. This study systematically analyses the generation mechanisms of algorithmic bias within data annotation, revealing governance dilemmas arising from dual pathways: cognitive embedding by agents and structural exclusion by data objects. By comparing legislative approaches to data annotation fairness governance across the EU, US, and China, it identifies theoretical blind spots in current regulations concerning the formalization of value embedding. Building upon this, a dual governance framework of "rigid legal constraints coupled with flexible ethical guidance" is proposed. This framework outlines pathways for bias mitigation through dual dimensions of rule embedding and ethical review, offering a systematic solution to address the "inherent flaws" of algorithmic discrimination and achieve fairness in data annotation. It further drives the paradigm shift in AI governance from "algorithm explanation" towards "data traceability".

## **1. Introduction of the problem**

The rapid advancement of artificial intelligence has fostered global consensus on "Trustworthy AI". The report of the 20th National Congress of the Communist Party of China advocates "promoting the modernization of the national security system and capabilities, resolutely safeguarding national security and social stability", identifying artificial intelligence as a key domain for national security. It emphasizes balancing development with security, mitigating potential risks of AI, and ensuring its healthy development within the national security framework (Shan Zhigang 2024). <sup>[1]</sup> As the foundational labor in AI model training (Yao Jianhua 2020), <sup>[2]</sup> data annotation directly impacts algorithmic fairness, accuracy, and security, constituting a critical component of trustworthy AI governance. On 26 December 2024, the National Development and Reform Commission and other departments issued the Implementation Opinions on Promoting High-Quality Development of the

Data Annotation Industry, which stated that efforts should be guided by Xi Jinping Thought on Socialism with Chinese Characteristics for a New Era to advance the high-quality development of the data annotation industry.

Ethical trustworthiness, a core dimension of trustworthy artificial intelligence, explicitly requires algorithmic decisions to adhere to principles such as fairness and non-discrimination. As the starting point for algorithmic production, data annotation bears the crucial mission of translating human cognition into machine rules. Its quality and fairness directly determine the trustworthiness of algorithmic systems. However, recent years have witnessed frequent incidents of algorithmic bias stemming from data annotation: the Stable Diffusion model reproduced entrenched stereotypes in generated images due to unfiltered gender-discriminatory content in annotations (Hu Yong and Zhang Wenjie 2024); <sup>[3]</sup> sample imbalances in medical diagnostic datasets led to increased misdiagnosis rates for darker-skinned patients; and recruitment algorithms discriminated against female job seekers due to historical biases embedded in annotated data. These phenomena indicate that fairness deficits in the data annotation process have become a significant risk point in AI development.

## **2. Examination of Current Fairness Regulations for Data Annotation Domestically and Internationally**

### **2.1 Comparative Analysis of Legislative Models in Major Countries and Regions**

The European Union's Regulation (EU) 2024/1689 on Artificial Intelligence (AI Act) establishes the world's first risk-tiered governance framework for AI. Article 1 (EU) 2024/1689, establishing the world's first AI governance framework centered on risk classification. Article 10 explicitly requires training, validation, and testing datasets for high-risk AI systems to possess "representativeness, error-

free nature, and completeness,” while mandating the establishment of “bias detection and mitigation mechanisms” (Colmenarejo et al. 2022).<sup>[4]</sup> Furthermore, Article 10(5) permits the collection of sensitive data (such as race and gender) under specific conditions to correct historical biases. However, the bill does not directly define “fairness”, instead conveying it indirectly through concepts such as “representativeness” and “non-discrimination”, resulting in ambiguity in legal interpretation (Lund 2021).<sup>[5]</sup>

The United States has also implemented multiple measures concerning artificial intelligence governance. At the federal level, the AI Bill of Rights (2022) mandates that automated decision-making systems provide “algorithmic discrimination protections,” though it does not address oversight of annotation processes (Baumann et al. 2023).<sup>[6]</sup> As voluntary industry guidance, the NIST AI Risk Management Framework (AI RMF 1.0, 2023) lists “fairness” as a characteristic of trustworthy AI, requiring the identification of bias within training data. SP 1270 report (2022) systematically defined “label bias” for the first time, proposing mitigation through annotator training and multi-round reviews (McMillan-Major et al. 2024).<sup>[7]</sup> However, NIST standards lack enforceability and rely on voluntary corporate adoption. Additionally, several US states are actively exploring AI legislation. New York prohibits the use of AI recruitment tools trained on biased data, mandating independent audits of dataset representativeness (Dotan et al. 2023).<sup>[8]</sup> Colorado’s AI Act requires developers and deployers of high-risk AI systems to exercise reasonable care to prevent algorithmic discrimination, while publicly disclosing data sources and risk disclosures.

Beyond the EU and US, China is also advancing research and practice in AI legal governance. Regarding data annotation, China adopts a principle of “balancing development and security,” promoting the industry through policies such as the New Generation Artificial Intelligence Development Plan and the Implementation Opinions on Promoting High-Quality Development of the Data Annotation Industry. The Interim Measures for the Administration of Generative Artificial Intelligence Services further stipulate that data annotation must comply with principles of transparency and explainability to ensure algorithmic fairness. However, these regulations predominantly consist of general principles, lacking specific implementation rules and legal liability provisions. This results in issues such as low legal standing and insufficient enforceability.

Comparative analysis of three typical governance models—those of the EU, the US, and China—reveals that a single governance logic struggles to address the complex challenges of data annotation fairness. Consequently, this study proposes a dual “legal-ethical” governance framework. This framework seeks to draw upon the value-led approach of the

EU, the technological flexibility of the US, and the coordinated efficiency of China, thereby attempting to construct a more inclusive and systematic governance solution.

## 2.2 Current Research on Data Annotation Fairness Governance

Within the framework of trustworthy artificial intelligence, data annotation—as the core component of machine learning model training—has increasingly drawn academic attention to its fairness governance issues. Scholars globally have proposed diverse governance recommendations from dimensions including technological ethics, labor rights, and data quality assurance.

### 2.2.1 Domestic Scholars’ Governance Recommendations

Domestic scholars advocate strengthening the standardization and professionalization of annotation labor processes. Addressing the labor-intensive and low-barrier nature of data annotation, Yuan Xingyu (2025) proposes establishing occupational standards and training systems to transition annotation work from “low-skill repetitive labor” to “specialized knowledge services”. For instance, medical data annotation requires industry expert involvement and certification mechanisms to ensure annotators’ professional qualifications, thereby enhancing data quality and model credibility (Lu Gaofeng and Yao Zhiyu 2024).<sup>[9]</sup> Additionally, Fan Libo and Liu Changjiang propose improving efficiency and quality through refining annotation tasks, boosting annotation efficiency, and establishing industry standards (Fan Libo and Yu Xinyue 2022).<sup>[10]</sup>

Furthermore, scholars adopting a Marxist technological ethics perspective emphasize embedding “human-centered” ethical principles throughout data annotation processes to prevent technological alienation from exploiting workers’ physical and mental wellbeing (Han Lili and Ma Wanli 2020).<sup>[11]</sup> Concurrently, a tiered classification system for sensitive data must be established to restrict the circulation of raw, non-anonymized data during the annotation process, thereby mitigating privacy leakage risks (Lu Ruiheng et al. 2023).<sup>[12]</sup>

### 2.2.2 Governance Recommendations from International Scholars

The US-based Information Technology and Innovation Foundation advocates that data annotation governance must transcend national boundaries, promoting mutual recognition and coordination of governance standards through international organizations such as the United Nations and the OECD. Some scholars propose adopting open standards akin to the Content Origin and Authenticity Partnership (C2PA), embedding traceable encrypted metadata into the annotation process to ensure transparency and auditability across transnational data supply chains (John and Joseph 2024).<sup>[13]</sup>

Ioannis Pastatzidis (2022) and colleagues emphasize that annotation data must encompass diverse groups to prevent model bias arising from singular perspectives. They advocate for annotators' participation in stakeholder consultations to ensure equitable task allocation and remuneration mechanisms. Furthermore, they propose employing "fairness-aware data augmentation" during annotation to proactively balance sample distributions, thereby mitigating subsequent algorithmic bias (Pastatzidis et al. 2022).<sup>[14]</sup>

Beyond this, some scholars emphasize the concept of a "fair market", arguing that competitive mechanisms can naturally correct labelling inequities. However, they also note that ethical guidance remains essential in high-risk domains (such as law enforcement and finance) to prevent labelling biases from causing systemic discrimination. This research calls for the introduction of third-party audits and fairness assessment tools on labelling platforms to enhance transparency (Aaronson 2021).<sup>[15]</sup>

### **2.2.3 Summary of Research Trends and Positioning of This Study**

In summary, current research broadly agrees that data annotation governance must balance technical reliability with social justice, exhibiting the following trends:

First, a shift from "algorithm-centric" to "data-centric" approaches: increasing recognition that algorithmic fairness stems from data quality necessitates advancing governance controls to the annotation stage. Second, a transition from "technical remediation" to "institutional prevention": relying solely on bias-mitigation technologies proves insufficient, requiring legal frameworks to embed fairness principles at source. Third, a shift from "platform responsibility" to "supply chain accountability": governance targets extend beyond AI service providers to encompass the entire chain of data collection, annotation, and circulation. Fourth, a shift from "domestic governance" to "global collaboration": the highly internationalized nature of the data annotation industry necessitates transnational legal coordination. However, the following limitations persist: Firstly, most proposals remain theoretical, lacking empirical research to validate their effectiveness. Secondly, insufficient attention is paid to the agency of annotators, particularly as the differentiated needs of data annotation personnel are inadequately incorporated into governance frameworks; Thirdly, interdisciplinary research requires further deepening, with the organic integration of legal, ethical, and computer science methodologies remaining a key future direction for breakthroughs. This study proposes a dual "legal-ethical" governance framework, distinct from existing "soft law-hard law" or "technology-norm" dichotomies, emphasizing the mechanism for translating ethical principles into legal rules.

### **2.3 Theoretical Blind Spots in Existing Research**

Current legal scholarship on trustworthy AI governance exhibits an "algorithm-centric" bias, revealing a theoretical blind spot in the formalization of value embedding.

The legal community's overreliance on the principle of "technological neutrality" has excluded subjective biases in the labelling process from legal evaluation systems, resulting in an "inherent flaw" of algorithmic discrimination. While some scholars advocate achieving fairness and impartiality through algorithmic technological neutrality (e.g., avoiding the parameterization of specific human characteristics) (Kanellopoulou-Botti et al. 2019),<sup>[16]</sup> technological neutrality does not equate to neutrality in algorithmic values or decision outcomes (Zhao Chao 2024).<sup>[17]</sup> Ethical norms, as crystallizations of societal value consensus, can establish boundaries for algorithmic annotation rules, preventing the encoding of unreasonable factors such as historical biases or group stereotypes into algorithms. While existing research widely acknowledges the importance of ethical principles, it largely remains at the level of advocating "soft law" or making declarations of principle. There is a lack of institutionalized, operational legal design for effectively embedding substantive value requirements such as fairness and non-discrimination into the specific rules of the pre-algorithmic data annotation stage.

In light of this, this study aims to transcend the "algorithm-centric" paradigm by constructing a dual "legal-ethical" governance framework. It focuses on resolving the core issue of how values can be embedded at the source, thereby providing a systematic pathway for governing fairness in data annotation.

## **3. Mechanisms Generating Algorithmic Discrimination in Data Annotation and Legal Dilemmas**

Within the framework of building trustworthy artificial intelligence, algorithmic discrimination arising from data annotation has become a fairness crisis that modern rule of law must confront. The algorithmic discrimination discussed herein primarily refers to "statistical discrimination" (Disparate Impact), wherein the annotation process, irrespective of subjective intent, leads to algorithmic decision outcomes disproportionately negatively affecting protected groups (such as specific genders or races). As the cognitive foundation of AI systems, data annotation fundamentally mediates the translation from human value judgements to machine decision-making rules. When this translation lacks necessary legal regulation and ethical constraints, latent subjective biases and sample imbalances within the annotation process can precipitate systemic discrimination in algorithmic decisions.

### 3.1 The Dual Pathways of Algorithmic Discrimination Generation

Algorithmic discrimination in data annotation follows a dual pathway. Firstly, the embedding pathway of cognitive biases. As social individuals, annotators' value preferences within their cognitive frameworks become unconsciously technologized during annotation. Research indicates annotators exhibit significant gender-occupation association biases, such as defaulting "nurse" to female and "engineer" to male. These micro-level, pervasive cognitive biases are captured and amplified through large-scale training datasets, ultimately solidifying into systemic biases within algorithmic decisions (Ye Qing and Liu Zongsheng 2023).<sup>[18]</sup> Secondly, the exclusionary pathway of objective data structures. Structural deficiencies in data samples create "cognitive blind spots" for algorithms regarding specific groups. Taking medical diagnostic systems as an example, a Stanford University research team discovered that when dermatology datasets predominantly collected samples from white individuals, it led to increased misdiagnosis rates for patients with darker skin tones (Wen et al. 2022).<sup>[19]</sup> This is not merely a technical issue; it further validates Bruno Latour's Actor-Network Theory—technology is not value-neutral but rather an 'enabler' of social relations and power structures. Actors within the data annotation network—annotators, platforms, standards, and data—collectively reproduce existing social inequalities as "objective" discrimination in the algorithmic era.

### 3.2 Manifestations of Imbalanced Fairness

The current governance framework exhibits significant institutional vacuums and technical limitations in the data annotation phase. Firstly, governance focus is lagging. Existing regulations predominantly concentrate on transparency and interpretability at the algorithm application end (e.g., the Regulations on the Management of Algorithm Recommendations for Internet Information Services), which constitutes back-end technical rectification. However, bias is deeply ingrained at the front-end source of the data supply chain. Just as back-end model fine-tuning cannot fundamentally rectify deficiencies in front-end data representation, algorithmic optimization alone cannot eradicate biases embedded during data annotation. A prime example is Amazon's recruitment algorithm, which discriminated against female applicants. The historical gender bias inherent in its training data was amplified and entrenched through the annotation process. Secondly, technical fixes possess inherent limitations. Google Research experiments confirmed that after eliminating gender disparities through adversarial training, the model's regional discrimination index paradoxically increased (Hassani 2021).<sup>[20]</sup> This demonstrates that ex post compensatory measures at the algorithmic level struggle to resolve value conflicts at the data production

stage (Zanna and Sano 2024).<sup>[21]</sup> Also reveals the inherent limitations of technical corrective measures—when value conflicts are rooted in the original data production process, relying solely on algorithmic optimization cannot achieve substantive justice.

Furthermore, this mechanism of data annotation reproduction exposes the deep-seated paradox of algorithmic discrimination: algorithms, cloaked in the guise of "objectivity" and "science", legitimize and reinforce existing societal biases. Data annotation has transcended mere technical operations, evolving into a value-driven process through which power entities reshape societal cognitive orders. Resolving this predicament thus necessitates moving beyond traditional "algorithm-centric" governance approaches. Instead, we must establish a collaborative governance mechanism embedding rules and conducting ethical scrutiny at the data source, thereby laying a robust foundation of justice for trustworthy artificial intelligence.

## 4. Dual Legal-Ethical Governance Pathways for Ensuring Fairness in Data Annotation

### 4.1 Collaborative Mechanism for Rule Embedding and Ethical Review in Algorithmic Fairness

Data annotation is not merely a technical process but a value-embedded social practice. Fan Hongxia's empirical research demonstrates that annotators' cognitive biases can influence AI system decisions through a "bias leakage" mechanism (Fan Hongxia and Yu Luhong 2024).<sup>[22]</sup> Consequently, to address algorithmic discrimination at its source, a collaborative governance system must be established that combines the rigid constraints of legal rules with the flexible guidance of ethical review. The core of this system lies in transforming abstract ethical principles into concrete, executable, reviewable, and accountable rules throughout the entire data annotation process through standardized and proceduralized methods.

During the planning phase of data annotation, ethical review should be proactively integrated into the formulation of annotation rules. Legislatively, the National Standardization Administration of China could spearhead the development of a mandatory national standard titled "Guidelines for Fairness in Data Annotation". The drafting process itself should undergo interdisciplinary ethical hearings, with content extending beyond technical parameters to explicitly define minimum diversity thresholds and deviation tolerances for sensitive attributes such as gender, ethnicity, and geography. The UK's Digital Ethics Assessment Framework may be referenced to establish review metrics such as cultural sensitivity, value neutrality, and social inclusivity (Zhang Anqi 2025).<sup>[23]</sup>

During the annotation implementation phase, a certification mechanism linking ethical compliance to legal obligations should be established. The root causes of generative AI's hallucination and bias risks lie in data annotation quality defects (Zhang Xin 2023).<sup>[24]</sup> Therefore, a "Trusted Annotation Certification" system is recommended, requiring annotation platforms to obtain ISO 38507 certification (ISO/IEC 38507:2022 Information Technology — Governance of IT — Governance Implications of the Use of Artificial Intelligence by Organizations). Concurrently, the ethical competence of annotation personnel should be institutionalized, mandating completion of prescribed hours of ethics and legal training with successful assessment. Their training certification records shall serve as critical evidence in future judicial proceedings to determine whether platforms fulfilled their "duty of reasonable care".

Prior to the delivery and application of annotation outputs, a firewall linking technical verification with legal accountability must be established. Developers of AI systems in high-risk domains are mandated to conduct fairness assessments on annotated datasets using standardized tools such as Fairlearn or IBM's AI Fairness 360 before model training. These assessment reports must be submitted as mandatory documentation for algorithmic registration. Should discriminatory incidents subsequently occur, these reports will serve as core evidence in determining developer liability. Concurrently, the Dutch "Data Donation" model permits users to set usage timeframes and scenario restrictions when annotating personal data. This approach could be referenced in legislation to explicitly stipulate that usage limitations set by data subjects when providing data constitute inviolable legal boundaries for annotation tasks. Violations of these boundaries should be directly recognized as infringements.

## 4.2 Theoretical Innovation of Dual Legal-Ethical Governance

Existing AI governance research predominantly focuses on the dichotomy of "soft law versus hard law" or the co-governance of "technology and norms." The "legal-ethical" dual governance framework proposed in this study introduces three innovations:

Firstly, governance focus shifts downstream: from algorithmic application (e.g., autonomous driving, facial recognition) to data production (annotation phase), achieving "source governance". Second, integration of normative systems: overcoming the dualistic opposition between law and ethics by translating ethical principles such as fairness and non-discrimination into operational legal rules, forming a normative hierarchy of "ethics → soft law → hard law"; Third, reconfiguration of principal responsibilities: transcending the traditional "platform-user" dual responsibility

framework by incorporating annotators, annotation platforms, algorithm developers, and data users into a unified system of rights and obligations, thereby establishing a "full-chain accountability mechanism".

## Conclusion

As the bedrock of AI trust systems, fairness governance in data annotation concerns not only technological credibility but also social justice and legal order. By analyzing the mechanisms generating algorithmic bias and the legal dilemmas in data annotation, this paper proposes a governance framework centered on "legal-ethical" synergy, shifting the focus from back-end algorithmic explanation to front-end data traceability. This framework institutionalizes fairness values within data annotation through rule embedding and ethical review pathways, providing foundational governance for trustworthy AI. Looking ahead, advancements in generative AI and multimodal large models will pose dual challenges of technological iteration and international coordination. Continuous refinement of institutional design is essential to ensure AI evolves fairly, reliably, and benevolently within the rule of law.

## REFERENCES

- [1] Economic Daily. 2024. Shan Zhigang: Consolidating the Security Foundation for Artificial Intelligence Development. [http://paper.ce.cn/pc/content/202401/03/content\\_287202.html](http://paper.ce.cn/pc/content/202401/03/content_287202.html). Accessed: 2024-01-03.
- [2] Yao Jianhua. 2020. Research on the Operational Mechanisms and Labor Control of Online Crowdsourcing Platforms: Taking Amazon Mechanical Turk as an Example. *Journal of Journalism and Communication Studies* (07), 17-32+121-122. DOI: 10.20050/j.cnki.xwdx.2020.07.005.
- [3] Hu Yong and Zhang Wenjie. 2024. Data Annotation Governance: Backstage Risks and Governance Shifts in Trustworthy Artificial Intelligence. *Yunnan Social Sciences* (06), 29-36.
- [4] Alejandra Bringas Colmenarejo, Luca Nannini, Alisa Rieger, Kristen M. Scott, Xuan Zhao, Gourab K Patro, Gjergji Kasneci, and Katharina Kinder-Kurlanda. 2022. Fairness in Agreement With European Values: An Interdisciplinary Perspective on AI Regulation. *Computing Machinery* 107-118. doi.org/10.1145/3514094.3534158.
- [5] European Tech Alliance. 2021. Lund, K.S., Piech, M., & Piech, M. European Tech Alliance Position on the European Commission's Artificial Intelligence Act proposal. <https://eutechalliance.eu/position-on-the-european-commission-artificial-intelligence-act-proposal/> Accessed: 2021-11-10.
- [6] Joachim Baumann, Alessandro Castelnovo, Riccardo Crupi, Nicole Inverardi, and Daniele Regoli. 2023. Bias on

- Demand: A Modelling Framework That Generates Synthetic Data With Bias. *Computing Machinery* 1002–1013. doi.org/10.1145/3593013.3594058.
- [7] Angelina McMillan-Major, Emily M. Bender, & Batya Friedman. 2024. Data Statements: From Technical Concept to Community Practice. *ACM J. Responsib. Comput* 1(1), 17. doi.org/10.1145/3594737.
- [8] Dotan, R., Kriebitz, A., Lütge, C., & Max, R. 2023. US Regulation of Artificial Intelligence. *Handbook on Applied AI Ethics*.
- [9] Lu Gaofeng and Yao Zhiyu. 2024. Algorithm-Relationship-Intermediary: Constructing a Hybrid Control Framework for Platform Labour Processes – An Embedded Study of AI Data Annotators. *Modern Communication (Journal of Communication University of China)* 46(08), 38-47.doi.org/10.19997/j.cnki.xdcb.2024.08.005.
- [10] Fan Libo and Yu Xinyue. 2022. Strategic Pathways for Data Factories to Move Beyond Contract Manufacturing: The Case of the Data Annotation Industry, *Science and Technology Management Research* 42(24),125-136.
- [11] Han Lili and Ma Wanli. An Exploration of the Humanistic Effects of Technology Ethics in the Context of Technological Alienation. *People's Forum: Academic Frontiers* (06),92-95. doi.org/10.16619/j.cnki.rmltxsqy.2020.06.012.
- [12] Lu Ruiheng, Xu Xiaogeng, Bai Xuejun et al. 2023. Research on Classification and Grading of Sensitive Personal Information. *Information Security and Communications Confidentiality* (04),46-56.
- [13] John C, Joseph M. 2024. Interoperable Provenance Authentication of Broadcast Media using Open Standards-based Metadata. *Watermarking and Cryptography*. arXiv. 2405.12336.
- [14] Ioannis Pastaltzidis, Nikolaos Dimitriou, Katherine Quezada-Tavarez, Stergios Aidinlis, Thomas Marquenie, Agata Gurzawska, and Dimitrios Tzovaras. 2022. Data augmentation for fairness-aware machine learning: Preventing algorithmic bias in law enforcement systems. *Computing Machinery* 2302–2314. doi.org/10.1145/3531146.3534644.
- [15] Aaronson, S.A. 2021. Transatlantic Priorities: Data Governance. *Intereconomics* 56, 59–60. doi.org/10.1007/s10272-021-0952-2
- [16] Kanellopoulou-Botti M, Panagopoulou F, Nikita M et al. 2019. The Right to Human Intervention, Law, Ethics and Artificial Intelligence. 2019 Computer Ethics Philosophical Enquiry. Norfolk. doi.org/10.2139/ssrn.3430075
- [17] Zhao Chao. 2024. Approaches to Governing Algorithmic Discrimination from an Equal Protection Perspective. *Journal of Shanghai Institute of Economic Management Cadres* 22(01),53-63. doi.org/10.19702/j.cnki.jsemc.2024.01.005.
- [18] Ye Qing and Liu Zongsheng. 2023. Causes and Governance Strategies for Algorithmic Gender Bias in Artificial Intelligence Scenarios. *Journal of Guizhou Normal University (Social Sciences Edition)* (05),54-63. doi.org/10.16614/j.gznuj(skb).2023.05.006.
- [19] Wen D, Khan S, Ibrahim H et al. 2022. Characteristics of publicly available skin cancer image datasets: a systematic review. *The Lancet Digital Health* 4(1): 64 – 74.
- [20] Hassani BK. 2021. Societal bias reinforcement through machine learning: a credit scoring perspective. *AI Ethics*, 1: 239-247.
- [21] Zanna K and Sano A. 2024. Enhancing Fairness and Performance in Machine Learning Models: A Multi-Task Learning Approach with Monte-Carlo Dropout and Pareto Optimality. arXiv:2404.08230.
- [22] Fan Hongxia and Yu Luhong. 2024. Bias Infiltration in Digital Annotation Behaviour by Crowdsourcing Platform Users. *Journal of Journalism* (11),101-112. doi.org/10.16057/j.cnki.31-1171/g.2024.11.006.
- [23] Personal library. 2025. Zhang Anqi: A Review of Artificial Intelligence Governance in the United Kingdom. [http://www.360doc.com/content/25/0411/16/12810717\\_1151033973.shtml](http://www.360doc.com/content/25/0411/16/12810717_1151033973.shtml). Accessed: 2025-04-11.
- [24] Zhang Xin. 2023. Data Risks and Governance Pathways for Generative Artificial Intelligence. *Legal Science (Journal of Northwest University of Political Science and Law)* 41(05),42-54. doi.org/10.16290/j.cnki.1674-5205.2023.05.006.