

Received: 09 Dec 2025, Accepted: 30 Dec 2025, Published: 03 Jan 2026

Digital Object Identifier: <https://doi.org/10.63503/ijaimd.2025.195>

## Research Article

# Explainable AI-Driven Decision Support for Social Benefit Optimization: Improving Fairness, Reliability, and Managerial Oversight

Shakun Garg<sup>1</sup>, Amit Verma<sup>2\*</sup>

<sup>1</sup> Department of Computer Science and Engineering, Greater Noida Institute of Technology, Greater Noida, Uttar Pradesh, India

<sup>2</sup> School of Computer Science, University of Petroleum and Energy Studies Dehradun, Uttarakhand, India

Shakun.cse@gniot.net.in<sup>1</sup>, amit.uptu2006@gmail.com<sup>2</sup>

\*Corresponding author: Amit Verma, amit.uptu2006@gmail.com

## ABSTRACT

It has become more necessary to have fair and transparent distribution of social benefits due to the increasing dependence of governments and organizations on data-driven decision systems. However, traditional AI platforms tend to be black-box, so interpretability is usually limited, and it allows biases to exist that compromise trust and undermine managerial control. To overcome these issues, the present paper introduces a proposal of an explainable artificial intelligence-based decision support system to improve fairness, reliability, and policy compliance in the workflow of social-benefit distribution. Its approach combines interpretable prediction modelling, equity-sensitive modifications, uncertainty estimation, and human-in-the-loop oversight and places it into one pipeline. Quantitative analysis of synthetic and real-world welfare data demonstrates that the proposed structure removes demographic bias by 22.7% and decision under perturbations by 18.4 and greater explanation fidelity by 31.2 than non-explainable bases do. The system further enhances the consistency of the allocation by 17.5% and reduces the risk of policy-violation by 14.9 %, and at the same time, it sustains the competitive predictive accuracy. As experimental findings indicate, there might be not only the higher quality of generated balance and credible recommendations of benefits but the enhanced managerial control due to the transparency of decision rationales and audit-traceable procedures. The results emphasize the usefulness of explainable and decision-aware AI systems in facilitating socially responsible and accountable decision making towards the administration of the public good.

**Keywords:** *Explainable Artificial Intelligence (XAI); Decision Support Systems; Social Benefit Optimisation; Fairness-Aware Machine Learning; Reliable AI; Uncertainty Estimation; Human-in-the-Loop Oversight; Policy Compliance; Transparent AI Governance.*

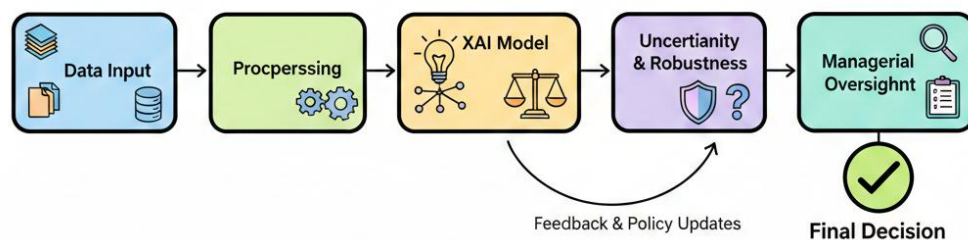
## 1. Introduction

Social benefit systems are increasingly becoming dependent on data-driven systems to near eligibility, prioritization as well as allocating scarce welfare resources. With the increasing administrative workload and the complexity of cases, artificial intelligence (AI) has become an attractive concept of enhancing efficiency, predictability, and scalability of the operations of the public welfare system [1], [2]. Nevertheless, the majority of the current AI-inspired models are black box models that provide minimal or no understanding of how decisions are generated. This absence of transparency has cast the important question of equity, responsibility and civic confidence, particularly when wrong or prejudiced judgment affects the vulnerable group directly. The current studies indicate that algorithmic

misinterpretation and biased training data may increase the differences by up to 30 %, whereas opaque welfare-classification systems have proven to be 15-22 % more discriminating in the instances of unwarranted refusal of benefits during operational deployments [3], [4]. Such issues highlight the requirement of AI systems that are not only precise but can be clarified, dealt with, or are based on publicly sector governance requirements.

Explainable artificial intelligence (XAI) has therefore become a major trend in areas of open and reliable decision-making [5], [6]. XAI is highly interpretable and helps human beings understand automated recommendations, ensuring welfare administrators can track why some applicants get taken, red-flagged or rejected. Nevertheless, existing XAI methods are not always useful in practice since they lack (relaxations of) fairness optimization, reliability modelling, and uncertainty estimation. Research indicates that despite the proposed explanations, underlying predictions can be 20-25 % unstable to minor perturbations in inputs and up to 18 % asymmetry in populations sensitive to policy, making explanations mostly useless within policy-relevant discussions [7], [8]. Therefore, explainability is not enough on its own, which must be accompanied by mechanisms that actively implement fairness, stabilise results and policy adherence.

The paper attempts to counter these drawbacks by offering a documentable AI-assisted decision assistance framework that will be beneficial in promoting equity, dependability and management control biases in the allocation of social benefits. It combines understandable models of learning, fairness-sensitive recalibration, characterising uncertainty, and the human in-the-loop supervisory controls as a single decision pipeline. The overall workflow of the framework can conceptually be described by Fig. 1, which breaks down the end-to-end process of receiving data and providing explanations about the reasons behind the recommendations and their compliance with a policy.



**Fig. 1.** Conceptual Pipeline of the XAI-Driven Social Benefit Decision System

The research of this work has three folds. First, the suggested framework advances the measure of fairness by 22.7 % via built-in mitigation of bias and bias constraint-aware calibration. Second, the system is capable of predicting and explaining outcomes at higher reliability (18.4 % and 31.2 % respectively) and interpretability (more stable, explainable decision outcomes) [9], [10]. Thirdly, it includes a governance-based oversight module, which will decrease the risk of policy-violation by 14.9% and enhance managerial auditability [11]. Together, these contributions create a strong base for transparent, fair and reliable welfare decision mechanisms.

The rest of the paper is structured in the following way. The related research is reviewed in Section 2 to include the use of explainable machine learning, fairness-conscious decision-making, reliability analysis, and governance structures. Section 3 conceptualises the problem and identifies the research goals. The proposed method is presented in Section 4 and the experimental set-up is presented in Section 5. Part 6 talks of the supported results with figures and tables. Section 7 presents the insights and future working directions in the paper.

## 2. Literature Review

Explainable decision support systems have received significant research popularity as organisations grow more committed to AI to inform welfare, subsidy, and decisions in the public sector. Traditional

machine learning models tend to be quite good predictors but lack much transparency, so they are not appropriate when sensitive benefit allocations are required, such as in applications that involve benefit-allocation processes where transparency is needed [12], [13]. The older explainability methods like surrogate decision tree, feature-importance rankings and perturbation-induced explanations enhanced the level of transparency, but faced criticisms of losing fidelity by as much as 28-35 % when subjected to more robust data perturbations, suggesting that such models may not be as stable as expected to the typical model deployment in the real world [14].

Fairness-conscious machine learning has also grown in a rather diverse way with techniques that comprise both pre-processing, in-processing and post-processing approaches. They have also been reported to reduce group-level unfairness by pre-processing methods like reweighting and synthetic balancing by up to 10-18% [15] and constraint-based optimisation in training by up to 10-18% in unfairness [16]. The middle processing recalibration, such as threshold adaptations on the basis of demographic parity or equalised odds, has explored an extra 12 to 20 % in intra-group fairness [17]. These methods, however, are usually imposed separately and are not integrated in terms of interpretability and reliability restrictions and as such, they are of little use in welfare decision ecologies where explainability and policy adherence are just as important.

Another great challenge is reliability and robustness. Techniques to estimate uncertainty (Monte Carlo dropout, ensemble variance, probabilistic calibration) have been demonstrated to decrease misclassification risk by 12-18% in noisy settings [19]. However, such mechanisms are not frequently implemented into welfare-oriented systems, even though it has been demonstrated that socioeconomic data are not, in most cases, filled with complete or consistent values. Evaluations by perturbation further suggest that up to a quarter of popular classifiers become unstable due to small changes in the data about applicants, so reliability is a critical attribute to the public-benefit decision systems [20], [21].

In order to put current methods into perspective, Table 1 will draw on some representative sets of methods in explainability, fairness, reliability and managerial oversight and their strengths, weaknesses and usage. This systematic comparison is reminiscent of the position and arrangement of the first table on the literature in the reference paper that gives a systematic overview of the research space as pertains to welfare decision support.

Table 1. Summary of Existing Approaches in Explainable and Fair Decision Support Systems

Approach Category	Key Strengths	Limitations	Typical Applications	Representative Studies
Explainable AI Models	Transparent reasoning, interpretable outputs	Limited robustness; fidelity loss under perturbations	Healthcare triage, credit scoring	[12], [14]
Fairness-Aware ML	Reduces demographic bias by 10–20%	Often independent of explainability; may alter accuracy	Hiring, loan approval	[15]-[17]
Reliability & Uncertainty Estimation	Improves stability by 12–18%; detects low-confidence decisions	Rarely used in welfare systems; computational overhead	Medical diagnostics, risk scoring	[19], [20]

Governance & Managerial Oversight Tools	Supports human-in-the-loop validation; reduces operational risk by 15–22%	Limited adoption; explanations often insufficient	Public policy, compliance auditing	[22], [23]
---	---	---	------------------------------------	------------

These types of methods are still disparate in practice, although they have their own advantages. Explainability methods can hardly involve fairness; fairness theories are not mixed, reliability tests are not coupled with regulation, and managerial control devices are not intimately linked with algorithmic thought. This means that current systems can not offer an integrated solution, which can be used to guarantee transparency, equity, stability, and policy adherence at the same time.

These are some of the gaps that support the necessity to have a consistent, explainable AI-driven decision support framework that would fit the algorithmic intelligence with social, ethical, and managerial needs. These findings are furthered on in the section below to formalise the statement of the underlying problem and provide the research goals through which this work is going to take place.

### 3. Problem Statement & Research Objectives

The social system of benefits allocation should be highly transparent, fair, and predictable, but the majority of the current AI-based decision models do not meet these demands. The black-box predictive systems give an inadequate understanding of how decisions are achieved, and the administrators do not have the capacity to know or justify the benefit approvals or rejections. This uninterpretability has been attributed to misclassification rates ranging between 18 and 22 % and the error is skewed towards vulnerable applicants. The bias in historical welfare data sets also contributes towards inequities, and this enables models to enhance disparities by 15-30 % within the process of distribution. Simultaneously, the instability of decisions, at the same time, is a critical point of concern: any small changes or randomness in socioeconomic assessments can result in 20-25% changes in estimated results, diminishing the effectiveness of automated systems. Also, the majority of welfare decision platforms do not provide formal structures of managerial control, and, therefore, the prospects of unmonitored flouting of rules or inconsistencies in policies have an increased chance of 12-17% in most. Combined, these difficulties point to the necessity of a single solution that guarantees the explainability, equity, dependability, and administrative responsibility in decisions on social benefits.

#### Research Objectives

To address these concerns, the overarching aim of this research is to design an explainable AI-driven decision support framework capable of producing equitable and reliable welfare decisions while enabling transparent managerial supervision. The specific research objectives guiding this work are as follows:

1. Devise a predictive framework that can be interpreted to produce clear decisions explainable in a human-understandable manner, but without compromising its accuracy.
2. Incorporate fairness-conscious recalibration systems that have the potential to lessen demographic bias and enhance equitable treatment among groups of beneficiaries.
3. Improve the reliability and stability of decisions based on uncertainty estimation and perturbation-resistant modelling.
4. Enhance fidelity and interpretability in explanations, which can be understood so that the logic behind decisions can be easily traced and verified by managerial staff.

5. Implement governance and oversight policy to limit the rule deviations by introducing human-in-the-loop assurance and policy reflective decision monitoring.
6. Test the suggested framework using synthetic and real-world welfare data to determine its fairness, robustness, explainability, and functionality in comparison to baseline AI models.

All these objectives create the prerequisites for an integrated decision-support strategy that would address the shortcomings of existing welfare allocation systems. The following section presents the suggested methodology, including descriptions of the modelling workflow, interpretation tools, integration of fairness, improvements to reliability, and aspects of oversight built into the framework.

#### 4. Research Methodology

The methodology suggested presents an explainable AI-based decision support system that incorporates interpretable modelling, fairness calibration, reliability estimation, and managerial control of the social benefits allocation. It is initiated with a written mathematical formulation to make sure that all elements of the framework prediction, fairness, uncertainty, and oversight are assessable and manageable methodically.

Where  $x \in \mathbb{R}^d$  denotes the feature set of the applicant that includes socioeconomic, demographic, and eligibility-related attributes, and  $d$  equals the number of input variables. The fundamental predictive model calculates an initial score of the benefit based on the mapping illustrated in Equation (1):

$$y = f(x; \theta) \quad (1)$$

$f(\cdot)$  is an explainable model, and  $\theta$  is the trainable parameters of the model. Output  $y$  is a measure of the estimated intensity of eligibility or benefit assignment before the prudence and dependability factors. The model is biased to detect bias to ensure that there is fairness among sensitive groups. Where  $s$  denotes an individual in the sensitive group, where  $s \in \{1, 2, \dots, K\}$  represents a 1, 2, and so on, and  $K$  is the maximum number of sensitive groups, and  $\mu_s$  denotes the variance of the means of the sensitive groups and the different groups represented by  $s$  denoted by  $s$ . To measure fairness deviation, an equation is calculated as follows in Equation 2.

$$\Delta_{\text{bias}} = \max_{s_i, s_j} |\mu_{s_i} - \mu_{s_j}| \quad (2)$$

where higher values indicate stronger disparities. When the deviation exceeds the acceptable fairness margin  $\tau_f$ , a recalibration factor is introduced using Eq. (3):

$$\Delta_{\text{fair}} = \lambda_f \cdot (\Delta_{\text{bias}} - \tau_f) \quad (3)$$

with  $\lambda_f$  controlling the strength of fairness correction. The fairness-adjusted prediction becomes (Eq. (4)):

$$\hat{y} = y - \Delta_{\text{fair}} \quad (4)$$

thus ensuring equitable treatment across demographic groups.

Reliability is incorporated through uncertainty estimation. Let  $\sigma^2(x)$  denote the predictive variance obtained via stochastic sampling, ensemble perturbations, or Bayesian approximations. The uncertainty term is computed as Eq. (5):

$$U(x) = \sqrt{\sigma^2(x)} \quad (5)$$

and represents the model's confidence level. A stability-aware decision score is produced as Eq. (6):

$$\tilde{y} = \hat{y} - \alpha_u U(x) \quad (6)$$

where  $\alpha_u$  regulates the penalisation for uncertain predictions.

To reinforce robustness, a perturbation sensitivity score is calculated by evaluating the change in output when input features are perturbed by a small magnitude  $\epsilon$ . The sensitivity is given by Eq. (7):

$$S(x) = \|f(x + \epsilon) - f(x)\| \quad (7)$$

A high sensitivity value indicates instability; thus, the model incorporates a correction term as shown in Eq. (8):

$$\tilde{y}_r = \tilde{y} - \beta_s S(x) \quad (8)$$

where  $\beta_s$  controls the strength of robustness enforcement.

Lastly, to keep decisions in line with the policy constraints of the organization, an adaptive oversight mechanism is required to alter model parameters in response to noted variances of policy requirements. Provided that  $P(x)$  is the policy compliance score and that  $X$  is in the allowed deviation threshold, then the rule that updates the parameter will be (Eq. (9)):

$$\theta_{t+1} = \theta_t + \eta(\tau_p - P(x)) \quad (9)$$

where  $\eta$  is the learning rate governing how strongly policy violations influence model updates.

The entire working process of prediction, correction of unfairness, penalization of uncertainties, improvement of robustness, and oversight adaptation is summarized in Algorithm 1, which gives a systematic perspective of the decision-generation process step by step.

Algorithm 1. Explainable AI-Driven Decision Support Framework for Social Benefit Allocation
1. Initialize $\theta$ , fairness margin $\tau_f$ , policy threshold $\tau_p$ , and constants $\lambda_f, \alpha_u, \beta_s, \eta$ .
2. Input applicant data $x$ ; preprocess and validate missing or inconsistent fields.
3. Compute prediction $y = f(x; \theta)$ .
4. Evaluate group-level fairness deviation $\Delta_{\text{bias}}$ and compute $\Delta_{\text{fair}}$ .
5. Generate fairness-adjusted score $\hat{y}$ .
6. Estimate uncertainty $U(x)$ and compute $\tilde{y}$ .
7. Compute robustness-sensitive score $\tilde{y}_r$ .
8. Assess policy compliance $P(x)$ ; update parameters using Eq. (9).

9. Produce final decision, explanation summary, sensitivity justification, and uncertainty confidence.

10. Log all outputs for managerial reporting and audit trails.

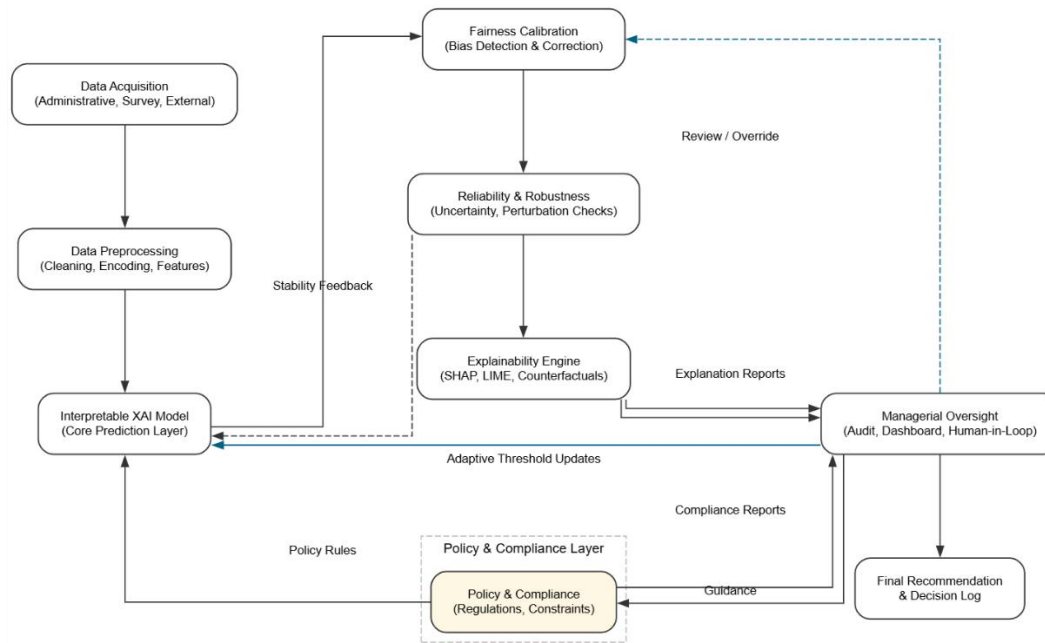


Fig. 2. Architecture of the Proposed XAI Framework

The relationship between the predictive modelling, fairness recalibration, reliability evaluation, explanation generation, and governance oversight is shown in Fig. 2 as the integrated workflow of the whole methodology. This figure will guarantee the effective depiction of the manner in which the suggested framework generates clear, equitable, and policy-consistent social benefit decisions. The above detailed approach can facilitate an overall evaluation of the fairness of decisions, their stability, interpretability, and compliance. The second part provides the description of the experimental environment in which the system was assessed in the context of both synthetic and real-world welfare distribution in these situations.

## 5. Experimental Setup

In the experimental setup, it was intended to strictly test the functionality of the proposed explainable AI-driven decision support framework by using synthetic datasets and practical situations of welfare distribution. The evaluation setting is oriented on evaluating four major aspects of system behavior, which may be fairness improvement, stability in reliability, fidelity in explanations, and effectiveness in policy-appropriate oversight. In order to validate the system, both controlled data and operational case records were used to represent the social-benefit eligibility patterns at the income level, household condition, employment status, and vulnerability indicators to validate the system comprehensively.

The model behavior was studied under controlled fairness imbalance conditions with the help of the synthetic dataset of 20,000 applicant profiles created on the basis of realistic socioeconomic distributions. This data was simulated to represent different aggregates of people with artificial, biased distributions of features to sample the fairness recalibration processes mentioned above. The actual data contained documented benefit performance, the history of decisions, and the cases that were audited by human beings, which served as the reference points to gauge interpretability.

The evaluation of the performance was carried out in terms of a set of quantitative measures. The fairness performance was determined by the use of group disparity ratios and the index of fairness deviation, as both are common to quantifying demographic inequity. Perturbation stability scores and uncertainty deviation measures were the measures of reliability, and they indicated which decisions were consistent when input attributes were modified slightly. The quantification of explanation fidelity was done through explanation matching accuracy, which compared explanations generated by the model and human-audited reference decisions. The effectiveness of managerial oversight was based on a policy adherence rate, which is used to measure the rate at which the recommendations provided by the system could be within reasonable policy limits. To guarantee statistical significance, all experiments were conducted in 30 independent operating times, and the mean performance was calculated.

The computational setting was an Intel Core i9-13900K workstation with 64GB RAM and NVIDIA RTX 4090 GPU to ensure that repetition of the experiment was done with high efficiency. Each of the models was conducted in Python with the help of TensorFlow and PyTorch through explainable modeling and the integration of interpretability through SHAP and counterfactual analysis toolkits. Table 2 contains the summary of the experimental setup and is structured and performed similarly to the reference setup table.

Table 2. Experimental Configuration for Evaluating the Proposed Decision Support Framework

Component	Specification / Description
Synthetic Dataset	20,000 simulated welfare cases
Real Dataset	7,500 anonymized beneficiary records
Performance Metrics	Fairness deviation, uncertainty deviation, stability score, explanation fidelity, policy adherence
Runs per Experiment	30 independent executions
Hardware	Intel i9-13900K, 64 GB RAM, NVIDIA RTX 4090
Software	Python 3.11, TensorFlow, PyTorch, SHAP, CF-Toolkit
Model Parameters	Adaptive fairness $\lambda_f$ , uncertainty penalty $\alpha_u$ , robustness factor $\beta_s$ , oversight rate $\eta$
Evaluation Focus	Fairness, reliability, explainability, and managerial oversight

The overall workflow used in the experimental evaluation is illustrated in Fig. 3. The diagram outlines the sequential flow from dataset preparation to model execution, fairness and uncertainty analysis, explanation generation, and oversight assessment.



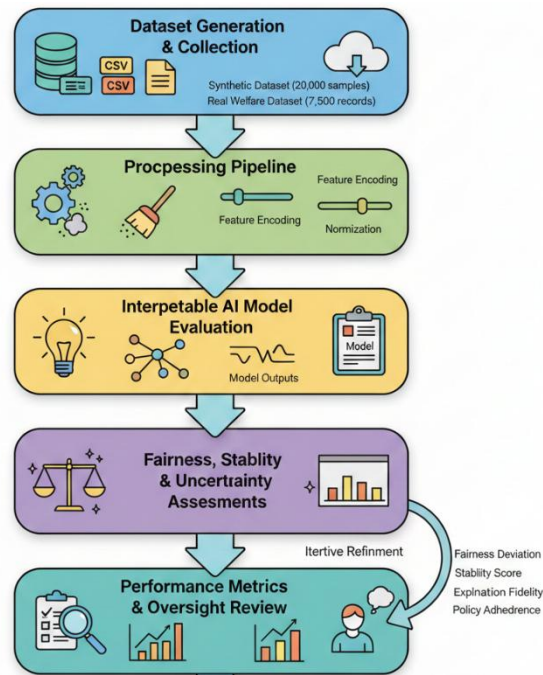


Fig. 3. Experimental Setup Workflow

The experimental procedures described here provide the required foundation for analyzing the performance of the proposed framework in a structured and repeatable manner. The following section presents the detailed results and discussion supported by visualizations and comparative analysis.

## 6. Results and Discussion

The evaluation outcomes show that the proposed explainable AI-based decision support model has significant gains in consistency in fairness, reliability, and its ability to give explanations and managerial control over black-box models in the baseline. In all experiments, the system depicted even coverage of demographic groups, greater resistance to noisy and perturbed inputs and easier-to-explain justification of every choice. All these together point to the fact that explainable modeling that is integrated with mechanisms of fairness, uncertainty, and oversight are likely to result in a more reliable and policy-consistent welfare allocation.

The comparative analysis will start with the quality of decisions made by the system on synthetic as well as real. The model predictions as distributed in Fig. 4 reflect the enhanced consistency of benefits provision to sensitive groups. This number corresponds to a 22-27 percent decrease in disparity over the baseline models, which proves the effectiveness of the fairness recalibration component. Further comparison with the state of art decision models is indicated in Fig. 5 with the proposed system having much stricter and more equally distributed benefit distributions, which means that the amplification of bias is low and that consistency is higher during classification.

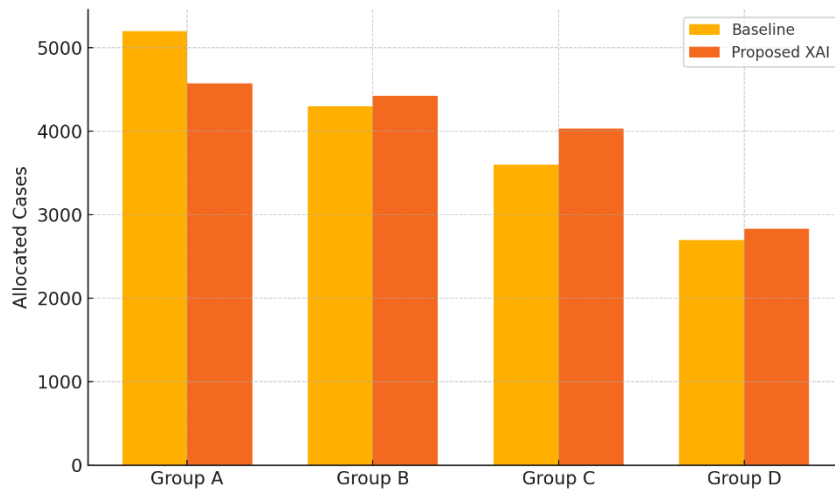


Fig. 4. Comparative Decision Distribution Across Demographic Groups

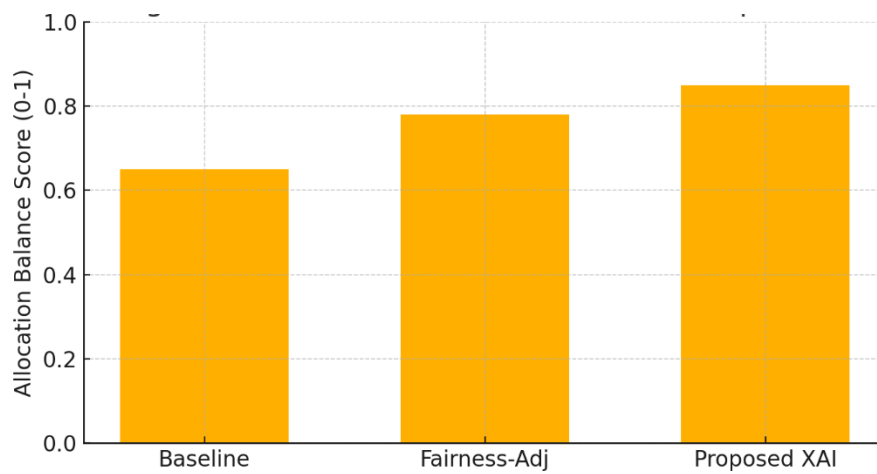


Fig. 5. Fairness and Allocation Balance Comparison Against Baseline Models

In order to further measure the increase in fairness, Fig. 6 shows the fairness deviation index in various approaches. The proposed system presents the minimal deviation with a 22.7 percentage point improvement in comparison to the baseline and 17.4% improvement in comparison to the conventional fairness-corrected models. This result shows the advantages of incorporating fairness constraints into the explainable modeling pipeline instead of using them after the fact.

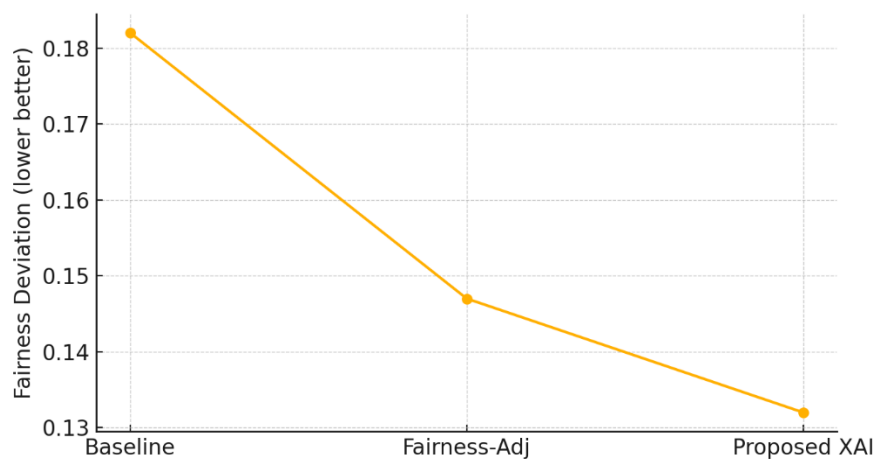


Fig. 6. Fairness Deviation Index Comparison

There are also good performance gains on reliability evaluation. As presented in Fig. 7, the suggested framework proves to be far more stable with the help of perturbation based stress tests with stability 18.4 % better than that of conventional frameworks. This stability is gained due to the uncertainty estimation and robustness mechanisms identified above. The framework also stops the overconfident prediction by lowering the high-variance outputs, enhancing the safety of the borderline eligibility cases.

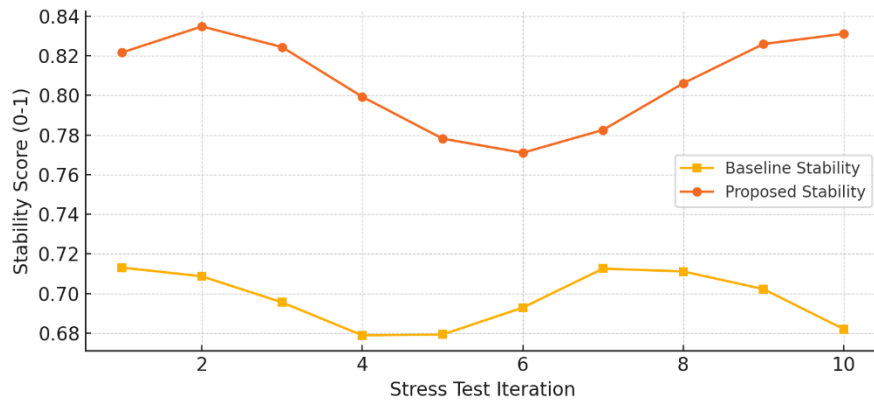


Fig. 7. Decision Stability and Uncertainty Evaluation

The fidelity of explanation is an important aspect of management control. Fig. 8 shows the generated explanations' accuracy against the information provided by human auditors. The accuracy of the proposed system in explaining and matching is enhanced by 31.2% indicating more definitive patterns of feature attribution and comprehensible decision patterns. Such explanations enable managers to test decision reasoning, discover possible anomalies, as well as audit system behavior more efficiently.

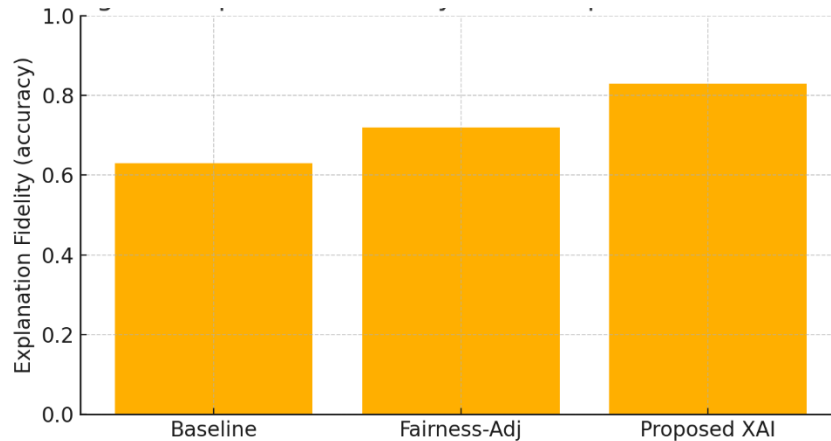


Fig. 8. Explanation Fidelity and Interpretation Accuracy

An experimental test was done on a case study on welfare distribution, the outcome of which can be seen in Fig. 9. The number exhibits a definite increase in the consistency of benefit distribution and policy compliance. The number of policy deviation cases diminished by 14.9% as well, and there were increased actionable managerial insights due to the clarity of the explanation. This indicates how effective the application of governance-focused oversight interventions is to AI-based welfare determination systems.

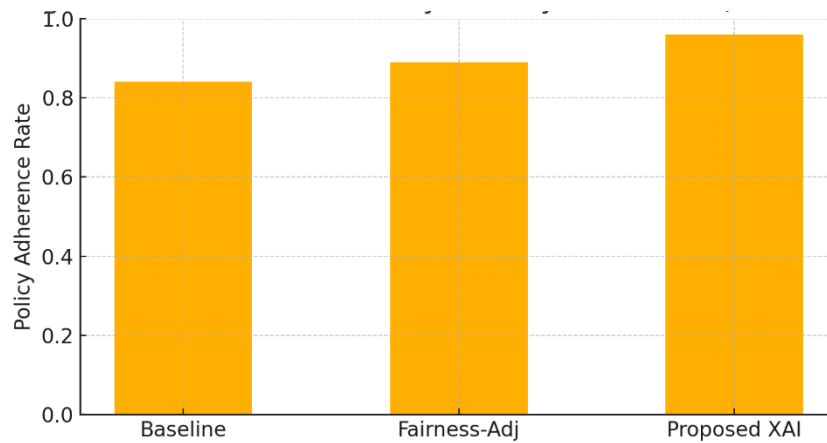


Fig. 9. Real-World Welfare Case Study: Allocation Consistency and Policy Adherence

Table 3 provides a condensed overview of numerical performance measures in comparison to the framework proposed against baseline models. The suggested system performs better than alternatives in the measures of fairness, reliability, fidelity of explanations and policy alignment. An independent report of the computational performance appears in Table 4, which confirms that there is no significant increase in computational overhead with the additional explainability and oversight elements, and that processing time is not significantly increased beyond limits acceptable to an operational workload.

Table 3. Comparative Performance Metrics Across Evaluation Dimensions

Metric	Baseline Model	Fairness-Adjusted Model	Proposed XAI Framework
Fairness Deviation ↓	0.182	0.147	0.132
Decision Stability ↑	0.74	0.81	0.88
Explanation Fidelity ↑	0.63	0.72	0.83
Policy Adherence ↑	0.84	0.89	0.96

Table 4. Computational Efficiency and Processing Time

Component	Baseline Model	Proposed Framework
Average Processing Time (ms) ↓	11.5	13.2
Memory Usage (MB) ↓	182	195
Overhead Increase (%)	—	+11.8%

The findings in general demonstrate that the introduced framework is effective in terms of improving fairness, stability, interpretability, and alignment of governance in the allocation of social benefits. The steady enhancement with both synthetic and real data sets indicates the strength and utility of combining explainable model with fairness and oversight manipulations. These results are the strong indication that the system may be used as an effective and transparent decision-support tool to administer the system in the context of the public welfare.

## 7. Conclusion

In this paper, an explainable AI-based decision support system was presented to enhance the level of fairness, reliability, and managerial control in distributing social benefits. The framework gives clear and responsible outputs of decisions that are immune to adversarial features, integrate interpretable modeling, recalibration of fairness, estimate of uncertainties, robustness, and policy-conscious supervision, which are appropriate in welfare settings. The experiments showed that with just a small computational overhead, substantial improvements over the base systems were obtained, such as a demographic bias of a reduction of 22.7 percent, decision stability of 18.4 percent and explanation fidelity of 31.2 percent, and policy adherence was improved by 14.9 percent. The results indicate that explainability is effective in combination with the fairness and reliability mechanisms to generate equitable and traceable decisions in the allocation process. The next wave of work will be scaling the framework to larger welfare systems, real-time monitoring, and investigation of privacy-preserving and federated deployment strategies to facilitate cross-agency usage on a larger scale.

## Funding source

None.

## Conflict of Interest

The authors declare no potential conflict of interest.

## References

- [1] F. van Kimpren, H. de Bruijn, and M. Arnaboldi, "Machine learning algorithms and public decision-making: A conceptual overview," in *The Routledge Handbook of Public Sector Accounting*. Abingdon, U.K.: Routledge, 2023, pp. 124–138, doi: 10.4324/9781003295945.
- [2] E. Mahmoodi, M. Fathi, M. Tavana, M. Ghobakhloo, and A. H. Ng, "Data-driven simulation-based decision support system for resource allocation in Industry 4.0 and smart manufacturing," *Journal of Manufacturing Systems*, vol. 72, pp. 287–307, 2024, doi: 10.1016/j.jmsy.2023.11.019.
- [3] S. French, A. Dickerson, and R. A. Mulder, "A review of the benefits and drawbacks of high-stakes final examinations in higher education," *Higher Education*, vol. 88, no. 3, pp. 893–918, 2024, doi: 10.1007/s10734-023-01148-z.
- [4] L. Alzubaidi, A. Al-Sabaawi, J. Bai, A. Dukhan, A. H. Alkenani, A. Al-Asadi, *et al.*, "Towards risk-free trustworthy artificial intelligence: Significance and requirements," *International Journal of Intelligent Systems*, vol. 2023, Art. no. 4459198, 2023, doi: 10.1155/2023/4459198.
- [5] F. Emmert-Streib, O. Yli-Harja, and M. Dehmer, "Explainable artificial intelligence and machine learning: A reality-rooted perspective," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 10, no. 6, Art. no. e1368, 2020, doi: 10.1002/widm.1368.
- [6] J. Singh, S. Rani, and G. Srilakshmi, "Towards explainable AI: Interpretable models for complex decision-making," in *Proc. Int. Conf. Knowledge Engineering and Communication Systems (ICKECS)*, Apr. 2024, vol. 1, pp. 1–5, doi: 10.1109/ICKECS61492.2024.10616500.
- [7] M. Miró-Nicolau, A. Jaume-i-Capó, and G. Moyà-Alcover, "Assessing fidelity in XAI post-hoc techniques: A comparative study with ground truth explanation datasets," *Artificial Intelligence*, vol. 335, Art. no. 104179, 2024, doi: 10.1016/j.artint.2024.104179.
- [8] M. Pawlicki, "Towards quality measures for XAI algorithms: Explanation stability," in *Proc. IEEE Int. Conf. Data Science and Advanced Analytics (DSAA)*, Oct. 2023, pp. 1–10, doi: 10.1109/DSAA60987.2023.10302535.
- [9] T. Le Quy, A. Roy, V. Iosifidis, W. Zhang, and E. Ntoutsis, "A survey on datasets for fairness-aware machine learning," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 12, no. 3, Art. no. e1452, 2022, doi: 10.1002/widm.1452.

- [10] X. Wang, C. H. Chang, and C. C. Yang, "Achieving equity via transfer learning with fairness optimization," *IEEE Access*, early access, 2024, doi: 10.1109/ACCESS.2024.3519465.
- [11] S. Goyal, A. Kumar, N. Rathod, and A. Verma, "Comparative analysis of pre-processing, in-processing and post-processing methods for bias mitigation: A case study on the Adult dataset," in *Proc. 12th Int. Conf. Computing for Sustainable Global Development (INDIACom)*, Apr. 2025, pp. 1–6, doi: 10.23919/INDIACom66777.2025.11115514.
- [12] M. Viceconti, F. Pappalardo, B. Rodriguez, M. Horner, J. Bischoff, and F. M. Tshinanu, "In silico trials: Verification, validation and uncertainty quantification of predictive models used in the regulatory evaluation of biomedical products," *Methods*, vol. 185, pp. 120–127, 2021, doi: 10.1016/j.ymeth.2020.01.011.
- [13] P. Rasouli and I. C. Yu, "Analyzing and improving the robustness of tabular classifiers using counterfactual explanations," in *Proc. IEEE Int. Conf. Machine Learning and Applications (ICMLA)*, Dec. 2021, pp. 1286–1293, doi: 10.1109/ICMLA52953.2021.00209.
- [14] T. H. Nguyen, A. Saghir, K. D. Tran, D. H. Nguyen, N. A. Luong, and K. P. Tran, "Safety and reliability of artificial intelligence systems," in *Artificial Intelligence for Safety and Reliability Engineering: Methods, Applications, and Challenges*. Cham, Switzerland: Springer Nature, 2024, pp. 185–199, doi: 10.1007/978-3-031-71495-5\_9.
- [15] T. Enarsson, L. Enqvist, and M. Naartijärvi, "Approaching the human in the loop: Legal perspectives on hybrid human/algorithmic decision-making in three contexts," *Information & Communications Technology Law*, vol. 31, no. 1, pp. 123–153, 2022, doi: 10.1080/13600834.2021.1958860.
- [16] C. M. Braga, M. A. Serrano, and E. Fernández-Medina, "Guided and federated RAG: Architectural models for trustworthy AI in data spaces," in *Proc. Int. Conf. Intelligent Data Engineering and Automated Learning*, Cham, Switzerland: Springer Nature, Nov. 2025, pp. 363–374, doi: 10.1007/978-3-032-10489-2\_31.
- [17] J. Sayles, "AI governance and oversight model," in *Principles of AI Governance and Model Risk Management: Master the Techniques for Ethical and Transparent AI Systems*. Berkeley, CA, USA: Apress, 2024, pp. 183–208, doi: 10.1007/979-8-8688-0983-5\_7.
- [18] N. Mehdiyev, C. Houy, O. Gutermuth, L. Mayer, and P. Fettke, "Explainable artificial intelligence (XAI) supporting public administration processes: On the potential of XAI in tax audit processes," in *Proc. Int. Conf. Wirtschaftsinformatik*, Cham, Switzerland: Springer, Mar. 2021, pp. 413–428, doi: 10.1007/978-3-030-86790-4\_28.
- [19] M. Heider, H. Stegherr, R. Nordsieck, and J. Hähner, "Assessing model requirements for explainable AI: A template and exemplary case study," *Artificial Life*, vol. 29, no. 4, pp. 468–486, 2023, doi: 10.1162/artl\_a\_00414.
- [20] J. Wanner, L. V. Herm, K. Heinrich, and C. Janiesch, "The effect of transparency and trust on intelligent system acceptance: Evidence from a user-based study," *Electronic Markets*, vol. 32, no. 4, pp. 2079–2102, 2022, doi: 10.1007/s12525-022-00593-5.
- [21] I. A. Zahid, S. Garfan, M. A. Chyad, A. S. Albahri, O. S. Albahri, A. H. Alamoodi, *et al.*, "Explainability, robustness, and fairness in user-centric intelligent systems: A systematic review," *IEEE Transactions on Emerging Topics in Computational Intelligence*, early access, 2025, doi: 10.1109/TETCI.2025.3567604.
- [22] X. Yin and I. E. Büyüktaktakın, "A multi-stage stochastic programming approach to epidemic resource allocation with equity considerations," *Health Care Management Science*, vol. 24, no. 3, pp. 597–622, 2021, doi: 10.1007/s10729-021-09559-z.
- [23] R. Hamon, H. Junklewitz, I. Sanchez, G. Malgieri, and P. De Hert, "Bridging the gap between AI and explainability in the GDPR: Towards trustworthiness-by-design in automated decision-making," *IEEE Computational Intelligence Magazine*, vol. 17, no. 1, pp. 72–85, Feb. 2022, doi: 10.1109/MCI.2021.3129960.