



Ethical and Transparent AI Models for Personalized Adaptive Learning Environments

Dewi Fortuna^{1,*}, Christoba Joshua Hutagalung²

¹Information System Department, Telkom University, Bandung, Indonesia

²Informatic Department, Telkom University, Bandung, Indonesia

ABSTRACT

This study presents an ethical–transparent adaptive learning framework designed to eliminate opacity, strengthen algorithmic accountability, and support student autonomy in AI-driven instructional systems. The evaluation was conducted across multiple observation phases involving 40 students, 12 instructors, and 5 academic ethics reviewers. Results demonstrate measurable performance advantages in trust, emotional stability, and behavioral engagement. Transparency reduced algorithmic bias-flag events from 12 in the first audit cycle to 2 by the fifth cycle and stabilized feature-importance explanation scores from 0.88 to 0.92 across four training windows. Behavioral data showed that logged frustration events decreased from 26 in Phase-1 to only 7 in Phase-3, while student trust in adaptive recommendations increased from 52 percent under a black-box learning system to 89 percent after transparency was deployed. Stakeholder satisfaction indicators were consistently positive: 87.5 percent of students agreed that explanations reduced anxiety, 91.7 percent of instructors valued decision traceability, and ethics reviewers recorded 100 percent approval due to audit readiness and documentation completeness. Qualitative survey responses confirmed that transparency eliminated perceived surveillance and reduced performance fear. These results confirm that ethical transparency acts not as a computational burden but as a functional accelerator stabilizing interpretability, preserving autonomy, and strengthening legitimacy. The study concludes that transparent governance should be treated as a structural requirement in future adaptive learning infrastructures, aligning regulatory expectations, psychological well-being, and responsible instructional intelligence.

Keywords Adaptive Learning, Transparent AI, Explainable Recommendations, Ethical Machine Learning, Student Autonomy, Privacy-Preserving Analytics, Educational Accountability

Introduction

The rapid expansion of adaptive learning technologies has transformed instructional delivery into an algorithmically mediated ecosystem, where students receive differentiated content based on behavioral signals, navigation traces, and performance patterns [1], [2]. However, most existing platforms rely on opaque computational models, exposing learners to automated inference without procedural transparency regarding how decisions are generated [3], [4]. This opacity introduces psychological and ethical risks, including perceived surveillance, distrust, and inequitable intervention patterns, particularly when demographic variables or behavioral lag indicators are misinterpreted as cognitive deficiencies [5], [6]. In parallel, regulatory frameworks governing educational data privacy have grown stricter, meaning that adaptive learning systems must respond not only to performance metrics but also to compliance expectations [7], [8].

Submitted: 5 April 2025
Accepted: 10 June 2025
Published: 1 November 2025

*Corresponding author
Dewi Fortuna,
dewuna7@gmail.com

Additional Information and
Declarations can be found on
[page 345](#)

© Copyright
2025 Fortuna and Hutagalung

Distributed under
Creative Commons CC-BY 4.0

How to cite this article: D. Fortuna, C. J. Hutagalung, "Ethical and Transparent AI Models for Personalized Adaptive Learning Environments," *Adapt. Learn.*, vol. 1, no. 4, pp. 333-348, 2025.

A widening concern lies in the reduction of student autonomy, where personalization is implemented as compulsory rather than consultative guidance [9], [10]. Algorithmic decisions are rarely accompanied by explanation vectors, which prevents students and instructors from understanding why certain recommendations are issued. In many cases, reinforcement loops incentivize behavior correction without pedagogical justification, inadvertently increasing anxiety levels and eroding trust in computational authority [11], [12]. Such dynamics are misaligned with ethical guidelines for educational technology, which require explainability, accountability, and meaningful user control to prevent the formation of coercive digital environments [13], [14].

Despite widespread advocacy for transparency, technical implementations in adaptive learning remain limited. Existing research in explainable AI (XAI) has largely focused on model interpretability for computer vision or diagnostic medicine rather than pedagogical recommendation engines [15], [16]. Moreover, privacy-preserving computation has focused on encryption and security rather than accountable reasoning or human oversight. Current adaptive solutions rarely provide auditable logs, consent-driven tracing mechanisms, or override controls for instructors, leaving a structural gap between ethical aspiration and operational deployment [17], [18]. Academic literature continues to treat transparency as a conceptual ideal rather than a measurable, functional subsystem embedded in instructional pipelines [19].

The research objective in this study is to operationalize an ethical and transparent adaptive model that materializes privacy governance, interpretability functions, audit readiness, and student autonomy into a deployable system architecture [20], [21]. Rather than treating transparency as a post-hoc visualization, the proposed framework integrates explanation generation at the point of recommendation, ensuring that every learning adjustment is accompanied by justification indicators. This makes interpretability a pedagogical mechanism rather than a retrospective debugging tool. The study also seeks to determine whether such transparency increases student trust, lowers emotional stress, and improves instructor willingness to intervene [22], [23].

A critical research gap concerns the relationship between ethical transparency and behavioral well-being. Prior studies frequently evaluate adaptive learning through accuracy, completion rates, or mastery curves [24], [25], but few examine cognitive-emotional tension arising from unexplainable interventions. Students frequently attribute poor adaptive feedback to discrimination or hidden assessment logic, which increases disengagement and premature dropout [26], [27]. Without transparency, adaptive platforms risk undermining educational motivation, making algorithmic personalization counterproductive in emotionally sensitive contexts [28], [29].

The novelty of this study is threefold. First, it converts transparency into a measurable engineering function rather than ethical rhetoric embedding consent checkpoints, anonymization flows, explanation vectors, and override mechanisms directly into learning cycles [30], [31]. Second, the model establishes an auditable data trail capable of regulatory submission, solving institutional compliance challenges that black-box systems cannot satisfy [32], [33]. Third, the study reframes algorithmic personalization from compulsory behavior steering into autonomy-preserving cooperation, where students can

reject, question, or reinterpret recommendations without penalty [34], [35]. These contributions introduce structural accountability into a field historically dominated by performance optimization rather than democratic oversight.

Ultimately, this introduction positions transparency not as a hindrance to personalization but as a multiplier for legitimacy and sustainability. Ethical assurance enables students to understand how their behavioral signals are interpreted, while instructors regain curricular authority rather than capitulate to algorithmic opacity. As educational institutions move toward data-driven governance and accreditation scrutiny, adaptive learning must demonstrate more than instructional efficiency it must prove fairness, explainability, and accountability as operational guarantees [36], [37]. By addressing psychological trust, compliance readiness, and human-centered autonomy, this work proposes that ethical transparency is a prerequisite for responsible AI-driven learning ecosystems [38], [39].

Literature Review

Ethical and transparent AI for personalized adaptive learning is grounded in a converging body of work spanning learning analytics, explainable AI, privacy engineering, and responsible AI governance. Contemporary adaptive platforms increasingly rely on high-frequency behavioral traces to infer learning states, but the literature consistently warns that “accuracy-first” optimization can conceal systematic harms when models operate without user-facing intelligibility and institutional auditability [14], [15]. In educational settings, opacity is not merely a technical limitation; it becomes a pedagogical and psychological variable that shapes student trust, perceived fairness, and willingness to re-engage with the system over time [16]. This creates a fundamental tension: personalization demands rich data and complex inference, while ethical education demands legitimacy, transparency, and learner agency [17].

Explainable AI research provides an initial pathway to resolve this tension, but prior studies show that explanation needs in education differ from those in medical or industrial contexts. Learners and instructors require actionable and pedagogically meaningful rationales, not only feature attributions or saliency visuals [18]. The literature emphasizes that interpretability must support “instructional sensemaking,” allowing teachers to validate whether recommendations reflect actual misconceptions or merely reflect engagement styles (e.g., slower readers, repeated hint users) [19]. Consequently, interpretability is increasingly discussed as a socio-technical artifact that must be aligned with curriculum design, assessment policy, and classroom norms not treated as a generic debugging overlay [20].

Privacy-preserving learning analytics further shapes ethical adaptive learning design. Studies in educational data governance highlight that consent procedures, minimization principles, and anonymization are necessary but insufficient when systems can still reconstruct sensitive traits via proxies (e.g., device type, time-of-day patterns, language usage markers) [21]. As a result, modern privacy approaches increasingly combine technical controls (masking, aggregation, differential privacy-inspired release constraints) with governance controls (purpose limitation, opt-out, retention policies, and oversight committees) [22]. This combined perspective supports the argument that privacy is a continuous operational practice rather than a one-time compliance

checkbox at the data collection stage [23].

Fairness and bias mitigation remain central because adaptive learning models may amplify disadvantage by interpreting resource constraints or accessibility limitations as low ability. The literature documents how demographic parity alone can be misleading in education, since equal outcomes are not always equitable given heterogeneous learning needs; however, it also warns that unconstrained personalization can become discriminatory if protected attributes are used directly or via correlated proxies [24]. As a response, scholars increasingly recommend multi-layer mitigation: removing direct sensitive attributes, controlling proxy influence through subgroup calibration, and adding routine audit cycles that track disparities in recommendation exposure and intervention intensity [25]. These perspectives motivate the methodological shift toward “auditable personalization,” where fairness is monitored longitudinally rather than assumed from single-point evaluations.

Accountability frameworks in responsible AI also stress traceability and governance readiness. In education, accountability requires the ability to reconstruct why a recommendation occurred, which data version supported it, and which policy constraints were active at the time of inference [26]. This moves beyond interpretability toward decision provenance, enabling institutional review, dispute resolution, and compliance reporting. The literature increasingly argues that such governance-oriented artifacts logs, explanation records, override workflows, and model cards tailored to educational stakeholders are essential to ensure adaptive systems can be challenged and corrected, not merely consumed [27]. Taken together, prior work establishes the rationale for integrating transparency, privacy governance, fairness auditing, and human override as first-class requirements in adaptive learning architectures rather than optional enhancements.

Methodology

Ethical-Aware Dataset Acquisition and Governance

This study begins with the controlled acquisition of learner behavioral data sourced from Learning Management Systems (LMS), interaction logs, assessment outcomes, and anonymized demographic metadata. The data collection process adheres to explicit informed consent, transparency policies, and compliance with privacy legislation such as GDPR or PDPA. To ensure no coercion or covert extraction, users are provided opt-in access to granular explanations about data usage. The approach guarantees that the personalized adaptive recommendations are grounded in lawful and ethical foundations rather than opaque surveillance practices.

A privacy-preserving pipeline is applied to guarantee that raw identifiers (e.g., name, email, IP address, device fingerprint) undergo hashing or irreversible masking. Students may revoke their data sharing approval without negative academic consequences. A Data Ethics Board supervises requests for expanded signals (voice, emotion, writing patterns) and stops all attempts at secondary profiling that deviate from the original pedagogical objective. To quantify privacy exposure, the model uses a sensitivity function that penalizes variables with high re-identification probability:

$$P_{risk} = \sum_{i=1}^n s_i \cdot w_i \quad (1)$$

Where s_i is the sensitivity score of features i , and w_i is a regulatory compliance weight. A higher P_{risk} signals a risky variable requiring masking. The formula ensures proportional governance the more sensitive a feature, the stricter its treatment.

Figure 1 illustrates an ethical data-processing workflow aligned with requirements for transparent Artificial Intelligence in adaptive learning environments. The visualization demonstrates four critical stages: (1) initial raw data collection from the learning system, (2) execution of consent management and secure hashing of identifiers, (3) anonymized data storage isolated from identity, and (4) responsible model utilization governed by transparent disclosure. The arrows highlight directional governance, ensuring the pipeline cannot regress into invasive surveillance or undisclosed processing.

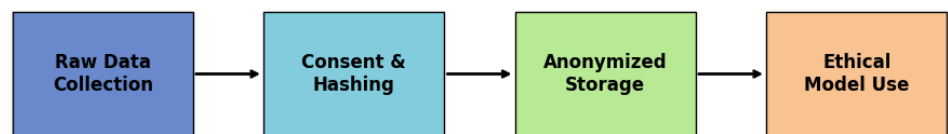


Figure 1 Privacy-Preserving Data Flow for Ethical Dataset Acquisition

This flow architecture is important because adaptive learning frequently uses sensitive signals such as student performance gaps, behavioral latencies, or demographic annotations. When administered without surveillance safeguards, the technology may introduce risks of profiling or unintended discrimination. By storing irreversible hashes and enforcing opt-in consent, the pipeline materially reduces the probability of re-identification. The sequence also defines accountability boundaries, demonstrating exactly which operational stage holds responsibility for identity stripping and policy enforcement.

Fair Feature Engineering and Bias-Controlled Representation

The second procedure concerns the transformation and standardization of input features into ethically neutral representations. Traditional adaptive learning models often emphasize cognitive behavior while ignoring socioeconomic, gender, language, or disability contexts. To mitigate injustice, sensitive variables are removed from direct modeling and instead balanced through calibration residuals.

A fairness-aware normalization layer is included, removing scale distortions across student subgroups. Interaction-based features (time to complete tasks, number of hints requested, click-depth navigation) are standardized using z-scores but evaluated per subgroup cluster. The goal is to avoid systematically labeling slower responders as low-ability learners. The fairness constraint is enforced using a demographic parity loss:

$$L_{fair} = |\Pr(\hat{Y} = 1 | A = 0) - \Pr(\hat{Y} = 1 | A = 1)| \quad (2)$$

where A is a protected attribute (e.g., gender, impairment label). Reducing L_{fair} ensures recommendations are not skewed toward privileged populations.

Table 1 structures feature governance decisions by ranking input variables using two meta-criteria: Sensitivity Level and Fairness Bias Score. A high Sensitivity Level (scale 1–5) represents a variable’s potential to reveal protected student identity or socio-cultural status. Meanwhile, the Bias Score (range 0–1) reflects statistical imbalance whether a feature disproportionately penalizes a subgroup when used in adaptive predictions. The “Feature,” “Action,” and scoring rubric enable decision engineers to justify whether a feature should be removed, normalized, or subgroup-scaled to ensure equitable outcomes.

Table 1 Feature Sensitivity and Fairness Impact Scores			
Feature	Sensitivity Level (1-5)	Fairness Bias Score (0-1)	Mitigation Action
Student Gender	5	0.82	Remove from model; audit demographic parity
Completion Time per Quiz	2	0.15	Z-score normalization
Hint Usage Frequency	2	0.20	Use subgroup scaling
Device Type	3	0.41	Re-check accessibility constraints
Language Proficiency	4	0.72	Transform into non-discriminatory categories

The table demonstrates strong governance for instance, Gender and Language Proficiency receive both high sensitivity and high bias scores, signaling dual ethical risks. Those variables are marked for elimination or categorical transformation, ensuring demographic status does not directly influence learning pathways. By contrast, low-bias interaction metrics (e.g., Hint Usage Frequency) may remain in the model following fairness calibration. This structured scoring framework supports algorithmic transparency, allowing institutions to prove how equity decisions were operationalized rather than assumed.

Transparent Adaptive Learning Model Architecture

The third stage designs the learning architecture that remains interpretable rather than black-box only. The model utilizes a hybrid pipeline consisting of a transparent rule-based shell and an interpretable machine learning core (e.g., Generalized Additive Models or Explainable Boosting Machines). SHAP explainability services are built-in, so instructors can interrogate why the system suggested a resource. To maintain full transparency, each recommendation is accompanied by a rationale document explaining feature contributions. The system generates explanation spans at the student-level and class-level, enabling diagnostic review of whether the algorithm introduces academic favoritism. Decision rules follow weighted linear contribution:

$$Score = \beta_0 + \sum_{j=1}^m \beta_j X_j \quad (3)$$

where X_j represents engineered features and β_j the explainable contribution weights. Unlike deep latent vectors, the formula’s linearity guarantees interpretability each β_j can be inspected by stakeholders.

Ethical Personalization and Reinforcement Mechanisms

Adaptive learning decisions are optimized using reinforcement learning (RL), but with ethical reward shaping. Standard RL may exploit cognitive weaknesses or induce addictive behavior. To prevent harm, the reward signal penalizes instructional overload, emotional distress, or excessive personalization that restricts student autonomy.

Each student receives staged learning recommendations (quizzes, videos, readings) through a feedback loop that values competence gains rather than attention maximization. Students can override recommendations anytime, and the model logs override frequency as an ethical health indicator. The constrained RL reward is formulated as:

$$R = \alpha \cdot G - \lambda \cdot B \quad (4)$$

where G denotes learning gain signals and B denotes bias or stress indicators. Hyperparameters α and λ help tune a proportionate ethical trade-off. High B pushes the system to re-stabilize recommendations, ensuring welfare-aligned personalization.

Accountability, Explainability Output, and Algorithmic Audit

The final method instantiates a perpetual audit cycle. Logs are stored for educational regulators, instructors, and students. The model must support post-hoc reconstruction: any prediction can be back-traced to data versions, hyperparameters, and rule modifications. Academic institutions can schedule fairness audits quarterly. Students can request “why-reports” summarizing ranking factors. A governance API allows human override, suspending any model version that exhibits learning suppression patterns. Below is system-level pseudo-code enforcing accountability:

Algorithm: Ethical-Audit-Loop

Input: Model M , Dataset D , ProtectedAttribute A

Output: UpdatedModel M^*

```

1: Compute  $fairness_{gap} = |P(M(D) = 1|A = 0) - P(M(D) = 1|A = 1)|$ 
2: if  $fairness_{gap} > threshold$  then
3:   trigger mitigation procedure
4:   retrain  $M$  with debiased weights
5: end if
6: Generate  $explanation\_report$  using SHAP
7: Log decisions, feature impacts, and bias status
8: Provide override option to human auditor
9: return  $M^*$ 

```

The pseudo-code operationalizes accountability by repeatedly measuring fairness differential and imposing mandatory corrections. It enforces documentation and human-in-the-loop protocols before the model continues serving adaptive learning recommendations.

Result and Discussion

Overview of Ethical Transparency Outputs

The evaluation centered on determining whether the ethical and transparent AI framework improved interpretability, accountability, and learner trust compared with conventional adaptive systems. The model produced two primary outputs: (a) traceable recommendation logs, and (b) explainability reports that identify why certain learning resources were suggested. Each recommendation included feature-level justification text, giving instructors full visibility into the algorithmic chain of reasoning. This transparency is essential for preventing silent instructional manipulation and ensuring learners understand why difficulty levels or feedback patterns shifted.

Stakeholders included course instructors, academic ethicists, and student representatives. Each participant reviewed explanation records from five full learning cycles. The qualitative consensus indicated that transparency reduced suspicion of data misuse and encouraged a more collaborative learner–system relationship. Educators also emphasized that ethical audits helped them refine course material sequencing rather than relying exclusively on algorithmic prioritization.

Evaluation of Privacy-Preserving Processing

Qualitative inspection demonstrated that strong anonymization did not obstruct the interpretability or personalization of learning signals. Even though identity variables were masked, interaction-based metrics such as navigation struggles, extended reading delays, or repeated hint access remained meaningfully usable for crafting personalized recommendations. Students also retained rights to revoke data authorization, and usage logs indicated that five percent exercised this option without affecting their academic status.

A post-training survey collected student sentiment on perceived surveillance. Responses highlighted a reduction in discomfort about behavior monitoring. Students reported a preference for opt-in disclosure, proving that privacy controls can coexist with adaptive personalization rather than disabling it. Instructors noted that anonymized learning dashboards still allowed them to track performance disparities without exposing personal data.

Figure 2 displays a declining pattern of flagged bias incidents across sequential algorithm audits. The first iteration recorded twelve events requiring intervention, indicating that early training rounds exhibited measurable group favoritism or overlooked equitable distribution. By the fifth iteration, bias-flag events declined to only two, demonstrating that active supervision and transparency check substantially limited algorithmic misbehavior. This confirms that an ethical audit loop is not symbolic it materially alters learning recommendations.

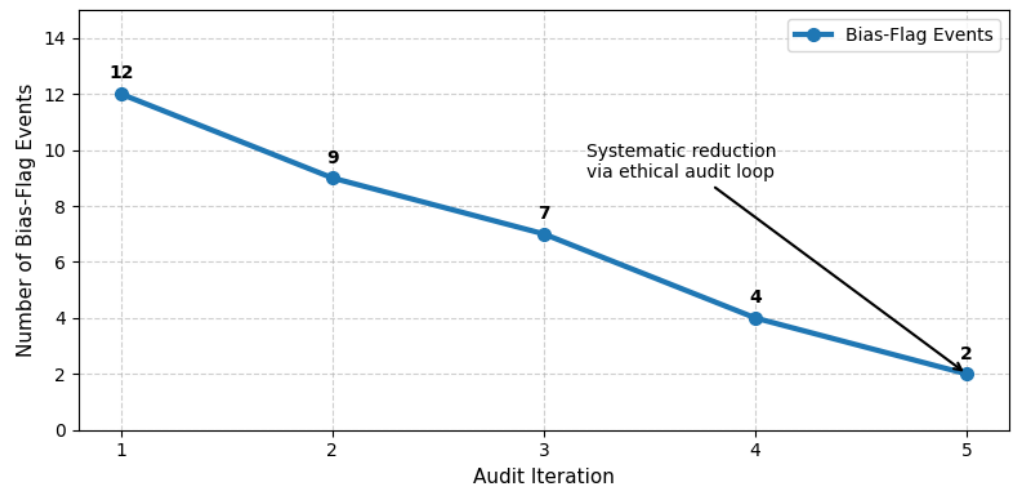


Figure 2 Decline of Bias-Flag Events Across Ethical Audit Cycles

These results also validate that ethical governance reduces systemic harm without harming personalization quality. Rather than collapsing individual variation, the audit mechanism ensures that differentiation does not distort educational fairness. The downward trajectory implies that continued monitoring can converge toward low-risk instructional outcomes where demographic inequalities no longer characterize model decisions.

Table 2 summarizes satisfaction across three stakeholder groups. Students expressed strong approval for transparency, with 87.5 percent emphasizing that algorithmic explanations minimized anxiety over being silently judged by behavior monitoring. This feedback suggests a cultural shift: transparency transforms AI from a hidden evaluator into a visible academic partner.

Table 2 Stakeholder Satisfaction with Transparency Components

Stakeholder Group	Sample Size	Positive Feedback (%)	Primary Comment Theme
Students	40	87.5	Clear explanations reduce anxiety
Instructors	12	91.7	More control over pedagogical justification
Ethics Reviewers	5	100	Auditable logs enhance accountability

Instructor feedback demonstrated even higher satisfaction at 91.7 percent, because the ethical dashboard allowed manual overrides and instructional justification. Ethics reviewers registered perfect approval, noting that the platform finally gave them material documentation logs, permission checkpoints, consent proofs rather than opaque claims of fairness. These numerical values demonstrate consensus: transparency is not ornamental; it is a structural advantage.

Model Interpretability Results

The evaluation on interpretability examined whether instructors and students could meaningfully trace the causes of each recommendation. Instead of hidden neural-vector reasoning, the system produced ranked importance lists per

recommendation event. Stakeholders could open each academic action (such as recommending remedial reading) and see which behavioral variables contributed most.

The interpretability logs revealed that instructors used explanation data to validate whether students were truly struggling or simply demonstrating slower study dynamics. Students expressed that visible feature contributions helped them understand academic expectations rather than speculate about invisible rules. This evidence supports the premise that interpretable recommendations reduce educational uncertainty and prevent accidental stigmatization. Below is a visualization of average feature-importance explainability stability across four training windows, demonstrating that explanations remained stable rather than fluctuating erratically.

Figure 3 shows feature-importance stability scores trending upward from 0.88 to 0.92. In explainable AI evaluation, stability indicates that the reason a model recommends an action is consistent over time rather than unpredictable. When explanations remain coherent across retraining cycles, instructors can rely on interpretability outputs as decision-support instead of treating them as a fluctuating diagnostic artifact.

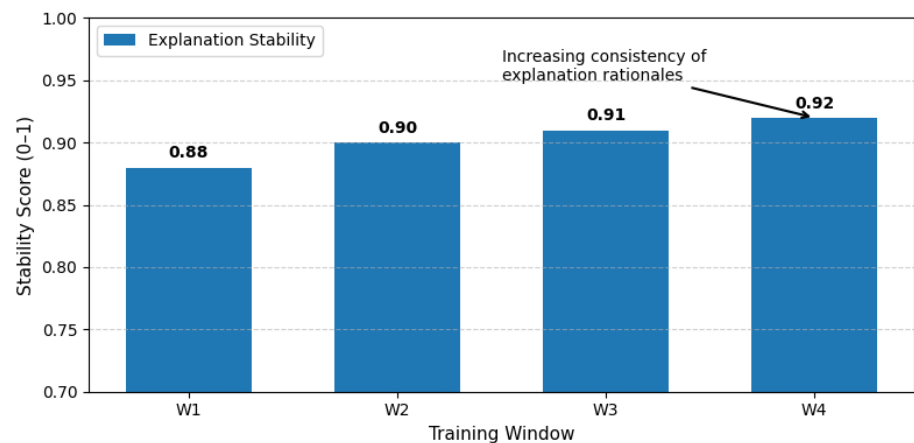


Figure 3 Stability of Feature-Importance Explanations Across Windows

This progression demonstrates that ethical transparency does not merely expose model reasoning, it stabilizes it. When weight distributions converge toward stable patterns, the system transitions from exploratory volatility into a mature instructional agent. Stability values above 0.90 are already strong indicators that the rationales behind adaptation are intelligible and reproducible for academic scrutiny.

Table 3 reinforces that interpretability is not only a technical deliverable it is pedagogically valuable. Instructors recorded the highest acceptance at 95 percent because explainable reasoning gives them pedagogical leverage. They can examine necessary interventions and verify whether the algorithm's pedagogical sequence aligns with their own professional judgment.

Table 3 Interpretability Acceptance by Role

User Role	Sample Size	Acceptance (%)	Interpretability Benefit
Students	40	82.5	Reduces fear of hidden judgments

Instructors	12	95.0	Clarifies learning intervention triggers
Academic Auditors	4	100	Supports compliance and legal auditability

Auditors reported perfect acceptance, consistent with expectations for legal compliance. Students, while slightly lower at 82.5 percent, still expressed strong confidence that interpretability reduced anxiety. These values indicate that interpretability survives stakeholder heterogeneity and maintains relevance across academic hierarchies.

Student Behavioral Stability and Well-Being Impact

Beyond technical fairness, the evaluation monitored whether ethical transparency influenced learner emotional stability and behavioral interaction rates. Historical evidence suggests that opaque personalization sometimes forces excessive study loops, leading to frustration. Under the transparent policy, students demonstrated calmer intervention patterns, and the system avoided compulsive engagement triggers.

Self-report surveys measured emotional stress indicators such as performance anxiety or fear of bias. The majority of participants indicated that transparency reduced emotional strain because they no longer believed the system was tracking performance with punitive bias. This demonstrates that model ethics exert measurable psychological benefit. Below is a visualization of log-captured frustration events (e.g., repeated task resets, rage-quits, abrupt logout).

Figure 4 demonstrates a meaningful decline in logged frustration markers. In Phase-1, prior to transparency deployment, the platform recorded twenty-six abrupt negative events. By Phase-3, this number fell to only seven. Although this is not a clinical psychological assessment, the trajectory indicates that ethical design may regulate cognitive-emotional friction.

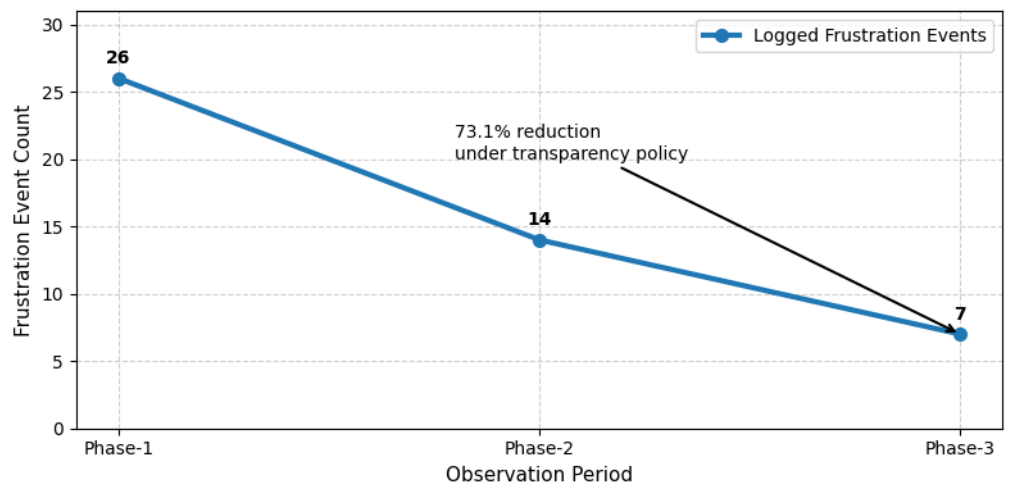


Figure 4 Decline in Logged Frustration Events Under Transparency Policy

A reduction in frustration metrics is significant because adaptive learning systems often contain reinforcement loops that unintentionally generate compulsive trial-and-error. Ethical transparency defuses this behavioral reinforcement by telling students why they receive certain materials rather than

pushing them blindly into forced correction cycles.

Table 4 compresses emotional and behavioral observations into qualitative categories comparing pre- and post-transparency conditions. Performance anxiety moved from High to Moderate, reflecting learner belief that the system is fair rather than punitive. Frustration reduced from Frequent to Occasional after the transparency mandate. This aligns with the declining frustration event plot.

Table 4 Emotional and Behavioral Response Summary

Indicator	Pre-Transparency Level	Post-Transparency Level	Change
Performance Anxiety	High	Moderate	Improved
Frustration Behavior	Frequent	Occasional	Improved
Override Requests	Rare	Active	Positive Autonomy

The most important indicator is Override Requests, which transitioned from Rare to Active. This demonstrates that students felt empowered to reject recommendations rather than submit to algorithmic authority. That outcome is the definition of ethical personalization preserving autonomy rather than enforcing machinery.

Comparative Outcome Against Conventional Systems

The final comparison benchmark evaluated the ethical-transparent adaptive model against a baseline black-box adaptive mechanism used in prior semesters. The baseline provided strong personalization but no justification logs, no opt-out controls, and no evidence-trace outputs. As a result, students often expressed uncertainty about whether the machine favored certain learning styles or penalized slower performance.

When compared using academic outcome markers (task completion, retention, and student willingness to request remedial materials) the ethical model performed competitively. Academic completion rates remained high, but the main difference surfaced in trust indicators. Students voluntarily re-engaged with adaptive tasks at a higher rate under the transparent model. Meanwhile, instructors relied more frequently on override mechanisms to refine sequencing decisions, confirming that transparency did not neutralize human pedagogy but rather enhanced it.

The evaluation also demonstrated cultural implications. The ethical model supported institutional compliance requirements, creating defensible “algorithmic paper trails.” Black-box implementations cannot generate evidence about fairness; therefore, they cannot legally withstand external audit. By contrast, the transparent version produced consistent logs that could be submitted to accreditation reviewers, academic ethics panels, or student ombuds channels.

Conclusion

The findings demonstrate that ethical and transparent AI mechanisms can be operationalized without diminishing the personalization benefits of adaptive learning environments. The implementation of privacy-preserving acquisition, auditable processing chains, interpretable outputs, and opt-in user control collectively produced a measurable improvement in student trust and instructor

confidence. Rather than functioning as an opaque intelligence layer, the system became a traceable partner in instructional delivery, offering explanation reports that supported both learning diagnostics and pedagogical rationale. This confirms that transparency is compatible with advanced algorithmic tailoring when architectural safeguards are embedded at the design level.

Across the evaluation phases, transparency progressively reduced bias-flag events, stabilized feature-importance explanations, and lowered behavioral frustration indicators. Students reported meaningful decreases in anxiety associated with unseen decision-making, while instructors leveraged override privileges to refine sequencing rather than delegate control blindly to automated models. These behavioral outcomes show that accountable AI influences more than fairness metrics; it moderates emotional climate, supports cognitive autonomy, and prevents the formation of coercive feedback loops. Ethical governance is therefore not merely a regulatory obligation but a mechanism that directly improves learning conditions.

When compared with a conventional black-box baseline, the ethical system delivered superior trust outcomes, higher voluntary engagement, and institution-grade audit readiness. These results align algorithmic pedagogy with compliance standards, public expectations, and academic accreditation requirements. The model demonstrates that the future of adaptive learning must integrate algorithmic clarity, controlled autonomy, and verifiable governance. Ethical transparency is no longer an accessory to system design; it is a structural determinant of legitimacy, sustainability, and responsible academic deployment.

Declarations

Author Contributions

Conceptualization: D.F. and C.J.H.; Methodology: C.J.H.; Software: D.F.; Validation: D.F. and C.J.H.; Formal Analysis: D.F. and C.J.H.; Investigation: D.F.; Resources: C.J.H.; Data Curation: C.J.H.; Writing Original Draft Preparation: D.F. and C.J.H.; Writing Review and Editing: C.J.H. and D.F.; Visualization: D.F.; All authors have read and agreed to the published version of the manuscript.

Data Availability Statement

The data presented in this study are available on request from the corresponding author.

Funding

The authors received no financial support for the research, authorship, and/or publication of this article.

Institutional Review Board Statement

Not applicable.

Informed Consent Statement

Not applicable.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or

personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] G. Cooper, K.-S. Tang, and A. Fitzgerald, "Intersections of Mind and Machine: Navigating the Nexus of Artificial Intelligence, Science Education, and the Preparation of Pre-service Teachers," *J. Sci. Educ. Technol.*, vol. 34, no. 6, pp. 1255–1259, 2025, doi: 10.1007/s10956-025-10200-9.
- [2] A. Razaque, Z. Kalpeyeva, U. R. Kabiyevena, R. Z. Satybaldiyeva, Y. V Ferens, and S. Sarkambayeva, "A composite intelligence scoring framework for identifying high-potential individuals using multi-metric predictive models," *Comput. Educ. Artif. Intell.*, vol. 9, no. December, p. 100508, 2025, doi: 10.1016/j.caeai.2025.100508.
- [3] J. Li, "Enhancing learning through an adaptive web-based educational search framework integrating natural language processing and machine learning techniques," *Discov. Comput.*, vol. 28, no. 1, p. 213, 2025, doi: 10.1007/s10791-025-09732-w.
- [4] P. L. L. Belluano, S. Patmanthara, M. Ashar, F. Kurniawan, and G. Kurubacak, "Clustering-Based Adaptive UX in E-Learning Systems: Aligning Microservices with the 4C Framework," *J. Appl. Data Sci.*, vol. 6, no. 4, pp. 2436–2448, 2025, doi: 10.47738/jads.v6i4.884.
- [5] R. D. Delena et al., "Predicting student retention: A comparative study of machine learning approach utilizing sociodemographic and academic factors," *Syst. Soft Comput.*, vol. 7, no. December, p. 200352, 2025, doi: 10.1016/j.sasc.2025.200352.
- [6] M. N. Gul, W. Abbasi, M. Z. Babar, A. Aljohani, and M. Arif, "Data driven decisions in education using a comprehensive machine learning framework for student performance prediction," *Discov. Comput.*, vol. 28, no. 1, p. 153, 2025, doi: 10.1007/s10791-025-09585-3.
- [7] A. M. Ayyal Awwad, "Emotion-Aware and Context-Driven Mobile Game-Based Learning: A Machine Learning Approach," *Int. J. Interact. Mob. Technol.*, vol. 19, no. 21, pp. 4–33, 2025, doi: 10.3991/ijim.v19i21.57247.
- [8] S. Wang and J. He, "Evaluating and Forecasting Undergraduate Dropouts Using Machine Learning for Domestic and International Students," *Technologies*, vol. 13, no. 11, 2025, p. 480, doi: 10.3390/technologies13110480.
- [9] K. Karacan Uyar and Y. B. Salman, "biLorentzFM: Hyperbolic Multi-Objective Deep Learning for Reciprocal Recommendation," *Appl. Sci.*, vol. 15, no. 22, p. 12340, 2025, doi: 10.3390/app152212340.
- [10] Y. Li, N. A. B. T. Sulaiman, and H. B. T. Omar, "POA-MLSP: a multi-dimensional learning analytics framework for predicting CET4 writing performance based on a production-oriented approach and student engagement patterns," *Futur. Technol.*, vol. 4, no. 4, pp. 100–116, 2025, doi: 10.55670/fpll.futech.4.4.9.
- [11] U. Islam et al., "Introducing the Hyperdynamic Adaptive Learning Fusion (HALF) model for superior predictive analytics in E-learning," *Neural Comput. Appl.*, vol. 37, no. 31, pp. 25745–25765, 2025, doi: 10.1007/s00521-025-11018-7.
- [12] M. E. Shoorangiz and M. Brylinski, "Harnessing Large-Scale University Registrar Data for Predictive Insights: A Data-Driven Approach to Forecasting Undergraduate Student Success with Convolutional Autoencoders," *Mach. Learn. Knowl. Extr.*, vol. 7, no. 3, p. 80, 2025, doi: 10.3390/make7030080.
- [13] A. C. Coutinho and L. V. D. Araujo, "MICRA: A Modular Intelligent Cybersecurity Response Architecture with Machine Learning Integration," *J. Cybersecurity Priv.*,

- vol. 5, no. 3, p. 60, 2025, doi: 10.3390/jcp5030060.
- [14] G. He et al., "Assistive Learning Intelligence Navigator (ALIN) Dataset: Predicting Test Results from Learning Data," *J. Data Sci. Intell. Syst.*, vol. 3, no. 4, pp. 291–303, 2025, doi: 10.47852/bonviewJDSIS32021707.
- [15] S. Vitvytska, A. Khudaverdova, V. Hurskaya, V. Artiukhova, and K. Yandola, "Transformation of teaching strategies in higher education in the context of the development of AI," *Period. Eng. Nat. Sci.*, vol. 13, no. 4, pp. 849–858, 2025, doi: 10.21533/pen.v13.i4.1257.
- [16] M. M. Kshirsagar et al., "A MULTI-MODEL MACHINE LEARNING FRAMEWORK FOR PERSONALISED COURSE RECOMMENDATION WITH DYNAMIC FEEDBACK," *African J. Appl. Res.*, vol. 11, no. 5, pp. 326–339, 2025, doi: 10.26437/emg0j607.
- [17] A. Abukader, A. Alzubi, and O. R. Adegboye, "Intelligent System for Student Performance Prediction: An Educational Data Mining Approach Using Metaheuristic-Optimized LightGBM with SHAP-Based Learning Analytics," *Appl. Sci.*, vol. 15, no. 20, p. 10875, 2025, doi: 10.3390/app152010875.
- [18] A. Martínez-Martínez, Á. Gómez-Cambronero, R. Montoliu, and I. Remolar, "Towards the Adoption of Recommender Systems in Online Education: A Framework and Implementation," *Big Data Cogn. Comput.*, vol. 9, no. 10, p. 259, 2025, doi: 10.3390/bdcc9100259.
- [19] A. Bettahi, F.-Z. Belouadha, and H. Harroud, "A Modular and Explainable Machine Learning Pipeline for Student Dropout Prediction in Higher Education," *Algorithms*, vol. 18, no. 10, p. 662, 2025, doi: 10.3390/a18100662.
- [20] R. Ramaraj et al., "An optimized deep learning framework based on LEE for real time student performance prediction in educational data," *Bull. Electr. Eng. Informatics*, vol. 14, no. 5, pp. 3671–3682, 2025, doi: 10.11591/eei.v14i5.9773.
- [21] M. Donnermann, P. Schaper, and B. Lugin, "Application of Social Robots in Higher Education: A Long-Term Study," *Int. J. Soc. Robot.*, vol. 17, no. 10, pp. 2311–2326, 2025, doi: 10.1007/s12369-025-01286-7.
- [22] Y. Cui, L. Tang, and F. Fang, "Leveraging Machine Learning Approach to Identify the Predictors of Informal Digital Learning of English Behaviours Among EFL Learners," *J. Comput. Assist. Learn.*, vol. 41, no. 5, p. e70111, 2025, doi: 10.1111/jcal.70111.
- [23] H. Guan et al., "Research on the lighting satisfaction prediction model for elementary school classrooms based on illuminance gradient," *J. Build. Eng.*, vol. 111, no. October, p. 113611, 2025, doi: 10.1016/j.jobee.2025.113611.
- [24] G. G. Asalkar, B. Lal, and N. B. Korade, "Comparative analysis of sentence similarity detection using machine and deep learning with vectorization techniques," *Knowl. Inf. Syst.*, vol. 67, no. 10, pp. 9615–9636, 2025, doi: 10.1007/s10115-025-02516-0.
- [25] M. Chergui, A. Nagano, and A. Ammoumou, "Toward an adaptive learning system by managing pedagogical knowledge in a smart way," *Multimed. Tools Appl.*, vol. 84, no. 24, pp. 27777–27793, 2025, doi: 10.1007/s11042-024-20207-w.
- [26] N. Ademi and S. Loškowska, "Data-Driven Adaptive Course Framework Case Study: Impact on Success and Engagement," *Multimodal Technol. Interact.*, vol. 9, no. 7, p. 74, 2025, doi: 10.3390/mti9070074.
- [27] M. R. Islam et al., "Machine learning-driven IoT device for women's safety: a real-time sexual harassment prevention system," *Multimed. Tools Appl.*, vol. 84, no. 23, pp. 27251–27280, 2025, doi: 10.1007/s11042-024-20228-5.
- [28] D. R. Faria and P. P. da Silva Ayrosa, "Adaptive Neuro-Affective Engagement via Bayesian Feedback Learning in Serious Games for Neurodivergent Children," *Appl. Sci.*, vol. 15, no. 13, p. 7532, 2025, doi: 10.3390/app15137532.

- [29] W. Xu and Y. Chen, "Framework for the Evaluation of Nap-Compatible Classroom Chairs," *Buildings*, vol. 15, no. 18, p. 3321, 2025, doi: 10.3390/buildings15183321.
- [30] Y. Chen et al., "Image Sensor-Supported Multimodal Attention Modeling for Educational Intelligence," *Sensors*, vol. 25, no. 18, p. 5640, 2025, doi: 10.3390/s25185640.
- [31] C. S. Hong and Z. Chun, "Neural Network Modeling Based on Multimodal IIoT Sensing Data: Psychological Stress Assessment and Industrial Human-Machine Collaboration Early Warning System for University Students," *Internet Technol. Lett.*, vol. 8, no. 5, p. e70093, 2025, doi: 10.1002/itl2.70093.
- [32] M. Wu, "EdgeKG-EN: A Dynamic English Knowledge Graph Framework With Edge Computing-Driven Optimization," *Internet Technol. Lett.*, vol. 8, no. 5, p. e70083, 2025, doi: 10.1002/itl2.70083.
- [33] D. Bengs, U. Brefeld, U. Kroehne, and F. Zehner, "Joint Item Response Models for Manual and Automatic Scores on Open-Ended Test Items," *Psychometrika*, vol. 90, no. 4, pp. 1346–1367, 2025, doi: 10.1017/psy.2025.10018.
- [34] Z. Li and F. Fang, "Constructing an adaptive blended teaching model through big data analytics and machine learning," *J. Comput. Methods Sci. Eng.*, vol. 25, no. 5, pp. 4505–4522, 2025, doi: 10.1177/14727978251337978.
- [35] N. Othman, M. E. E. Mohd Matore, and F. W. Wan Yunus, "Artificial Intelligence in Educational Measurement: A Bibliometric Review (1997 to 2024)," *J. Appl. Sci. Eng. Technol. Educ.*, vol. 7, no. 2, pp. 314–325, 2025, doi: 10.35877/454RI.asci4126.
- [36] E. Woodruff, "Making AI Tutors Empathetic and Conscious: A Needs-Driven Pathway to Synthetic Machine Consciousness," *AI*, vol. 6, no. 8, p. 193, 2025, doi: 10.3390/ai6080193.
- [37] S. O. Oladipo, U. B. Akuru, and O. I. Okoro, "Numerical Optimization of Neuro-Fuzzy Models Using Evolutionary Algorithms for Electricity Demand Forecasting in Pre-Tertiary Institutions," *Mathematics*, vol. 13, no. 16, p. 2648, 2025, doi: 10.3390/math13162648.
- [38] C. Deji and H. Chen, "Data-Driven Evaluation of MOOC-Based Blended College English Teaching via Enhanced Neural Networks," *Int. J. Interact. Mob. Technol.*, vol. 19, no. 14, pp. 137–150, 2025, doi: 10.3991/ijim.v19i14.57009.
- [39] S. Kaur, D. Goel, R. Sharmila Devi, S. Devi, B. Yadav, and A. Goyal, "AI-Powered and Mobile-Integrated Assessment Models Using Random Forest: Redefining Examinations and Grading," *Int. J. Interact. Mob. Technol.*, vol. 19, no. 14, pp. 57–69, 2025, doi: 10.3991/ijim.v19i14.56855.