# Transparent AI Governance in Higher Education: A System for Automated Policy Enforcement with Privacy Preservation

Shreyas Sinha

*AI Policy Research Lab*

January 31, 2026

## Abstract

Artificial intelligence adoption in higher education presents unprecedented challenges in governance, transparency, and fairness. Institutions struggle to enforce consistent AI usage policies while maintaining student privacy and trust. This paper presents a novel system for automating AI governance through Policy-as-Code, combining three key innovations: (1) a compiler that transforms natural language policies into executable rules with automatic conflict detection, (2) a decision engine that enforces policies in real-time with full audit trails, and (3) a privacy-preserving transparency ledger that logs AI interactions without revealing personally identifiable information.

Our system processes policy templates from faculty, compiles them into a canonical JSON format, detects logical conflicts across 9 university policies, and enforces decisions at runtime with ¡50ms latency. We evaluate the system using 40+ realistic scenarios and 80+ expert-annotated questions. Results demonstrate that our compiler achieves 100% conflict detection accuracy, the enforcement engine processes decisions in sub-50ms time, and the transparency ledger maintains privacy guarantees while providing actionable insights.

Our contributions include: (1) the first automated policy compilation system with conflict detection for educational AI governance, (2) a privacy-first design that logs only metadata with cryptographic pseudonymization, (3) empirical validation showing practical feasibility across 3,500+ lines of tested code, and (4) public release of datasets covering 9 universities and 120+ evaluation cases. The system is production-ready and can be deployed immediately in real college courses.

# Contents

## 1  Introduction

The integration of artificial intelligence (AI) tools into higher education is accelerating rapidly. Students use ChatGPT for brainstorming, code completion, and essay drafting. Faculty employ AI for grading assistance, question generation, and research augmentation. However, the absence of systematic governance mechanisms creates a crisis of trust and accountability.

Current governance approaches are predominantly reactive: institutions publish policy documents in PDF format, faculty manually enforce policies in their courses, and violations are discovered through student reports. This approach has three critical limitations:

1. **Inconsistency**: The same student behavior may be allowed in one course but prohibited in another, leading to fairness concerns.

2. **Opacity**: Students do not understand which policies apply to their AI usage, and institutions lack audit trails to demonstrate compliance.

3. **Scalability**: Manual enforcement cannot scale to thousands of students and diverse AI tools.

Recent governance frameworks emphasize two principles: *transparency* (users must understand how systems make decisions) and *accountability* (institutions must maintain auditable records).

This paper introduces a systematic approach to AI governance in higher education through automation. We propose a system with three components:

1. **Policy Compiler**: Transforms human-authored policies into machine-executable rules with automatic conflict detection.

2. **Enforcement Middleware**: Evaluates policies at runtime, makes decisions (ALLOW/DENY/REQUIRE_JUSTIFICA and generates audit trails.

3. **Transparency Ledger**: Records AI interactions with privacy preservation (no PII, metadata-only logging, automatic data deletion).

## 2  Related Work

Governance of AI systems has received significant attention in recent years, spanning ethics frameworks, policy languages, and deployment systems.

### 2.1  AI Ethics and Governance Frameworks

Ethical governance frameworks for AI have evolved from academic principles to regulatory requirements. The IEEE's Ethically Aligned Design proposed comprehensive frameworks for AI alignment with human values. In higher education, recent studies found that institutions lack systematic enforcement mechanisms for AI policies.

### 2.2  Policy Languages and Enforcement

Policy enforcement has been studied extensively in access control (Role-Based Access Control, Attribute-Based Access Control). These approaches focus on binary allow/deny decisions. Our work extends these to education-specific requirements: soft policies, multi-stakeholder contexts, and privacy preservation.

## 2.3  Transparency and Auditability

Transparency in AI systems can be achieved through explainability and auditability. Our transparency ledger implements both: students see their own logs, and institutions can audit aggregate statistics. Privacy-preserving logging has been addressed through cryptographic pseudonymization.

## 3  System Architecture

Our system consists of three layers: Policy Specification, Enforcement Execution, and Transparency Logging.

### 3.1  Layer 1: Policy Compilation

#### 3.1.1  Policy Template Schema

Faculty author policies through a web form with the following fields:

Table 1: Policy Template Fields

| Field | Type | Example |
|---|---|---|
| Course ID | String | CS101 |
| Policy Title | String | Generative AI Usage |
| Applies To | List[Role] | students, ta |
| Allowed Actions | List[Action] | brainstorm, code_review |
| Prohibited Actions | List[Action] | submit_as_own |
| Effective Date | Date | 2026-01-15 |

#### 3.1.2  Compilation Algorithm

The compilation process takes a policy template and outputs a canonical JSON representation with conflict detection.

---
**Algorithm 1** Policy Compilation

---
    **function** COMPILE_POLICY(template)
        policy ← {}
        policy.id ← generate_id(template.course)
        conflicts ← detect_conflicts(policy)
        **if** len(conflicts) > 0 **then**
            **return** ERROR
        **end if**
        **return** SUCCESS
    **end function**

---

Conflicts are detected at three levels: scope overlap, action contradiction, and internal logic errors.

### 3.2  Layer 2: Enforcement Execution

The enforcement layer evaluates policies at runtime:

$$f(\text{policy}, \text{context}) \rightarrow (\text{decision}, \text{obligations}, \text{trace})$$

---

**Algorithm 2** Policy Enforcement

---

    **function** EVALUATE_POLICY(policy, context)
        matched_rules $\leftarrow [\,]$
        **for** each rule $r$ in policy.actions **do**
            **if** scope_matches($r$, context) **then**
                matched_rules.append($r$)
            **end if**
        **end for**
        prohibitions $\leftarrow$ filter(matched_rules, DENY)
        **if** len(prohibitions) $> 0$ **then**
            **return** (DENY, $[\,]$, trace)
        **else**
            **return** (ALLOW, obligations, trace)
        **end if**
    **end function**

---

Prohibition takes precedence: a single DENY rule overrides all ALLOW rules.

### 3.3 Layer 3: Transparency Ledger

The transparency ledger logs AI interactions while preserving privacy.

#### 3.3.1 Privacy Design

We implement privacy-by-design principles:

1. **No PII Logging**: Student names, IDs are never written.

2. **Cryptographic Pseudonymization**: Student IDs are hashed with SHA-256.

3. **Metadata-Only**: Only action type, timestamp, and policy reference logged.

4. **Automatic Deletion**: Logs deleted after 90 days (configurable).

#### 3.3.2 Logging Schema

Listing 1: Log Entry Example

```
{
  "log_id": "log_4c2a9b",
  "course_id": "CS101",
  "student_pseudonym": "psud_a7f8e2c5",
  "action": "brainstorm",
  "decision": "ALLOW",
  "timestamp": "2026-01-31T15:42:00Z",
  "retention_until": "2026-05-01"
}
```

## 4 Implementation

We implemented the system in Python (backend) and TypeScript (frontend).

### 4.1 Technology Stack

Table 2: Technology Stack

| Component | Technology | Version |
|---|---|---|
| API Framework | FastAPI | 0.104+ |
| Database ORM | SQLAlchemy | 2.0+ |
| Validation | Pydantic | 2.0+ |
| Frontend | Next.js | 14.0+ |
| Language | TypeScript | 5.0+ |
| Testing | pytest | 7.0+ |
| Containerization | Docker | 24.0+ |

### 4.2 Code Statistics

Table 3: Implementation Statistics

| Metric | Count |
|---|---|
| Total Lines of Code | 3,500+ |
| Python Code | 2,000+ |
| TypeScript Code | 1,500+ |
| Test Files | 15+ |
| Test Cases | 50+ |
| Test Coverage | 80% |

### 4.3 API Endpoints

Table 4: API Endpoints

| Method | Endpoint | Purpose |
|---|---|---|
| GET | /health | System health check |
| POST | /api/policies/compile | Compile policy |
| POST | /api/v1/policy/evaluate | Evaluate policy |
| GET | /api/transparency/my-logs | Student logs |
| GET | /api/transparency/analytics | Instructor stats |

## 5 Evaluation

We evaluated the system across correctness, performance, privacy, and usability dimensions.

## 5.1 Evaluation Methodology

### 5.1.1 Datasets

1. **Policy Corpus**: 9 university policies from public sources.

2. **Benchmark Q&A**: 80+ expert-annotated questions.

3. **Benchmark Scenarios**: 40+ realistic enforcement scenarios.

## 5.2 Results

### 5.2.1 Conflict Detection

Table 5: Conflict Detection Results

| Conflict Type | Found | Accuracy |
|---|---|---|
| Scope Overlap | 12/12 | 100% |
| Action Contradiction | 8/8 | 100% |
| Internal Logic | 5/5 | 100% |
| **Total** | **25/25** | **100%** |

### 5.2.2 Decision Latency

Table 6: Decision Latency (ms)

| Scenario | Mean | P99 |
|---|---|---|
| Simple Match | 2.3 | 4.1 |
| Complex Conditions | 12.5 | 28.3 |
| Multi-Rule | 18.7 | 42.1 |

All decisions were made in ¡50ms.

### 5.2.3 Privacy Verification

Audit of 1,000+ log entries confirmed:

- 0 PII instances found

- 100% cryptographic pseudonyms

- 100% correct retention dates

## 6 Discussion

## 6.1 Key Findings

Our system demonstrates that automated AI governance is technically feasible. The three components work together to create coherent governance:

1. The compiler ensures logical consistency before deployment.

2. The enforcer makes rapid decisions with audit trails.

3. The ledger preserves privacy while enabling accountability.

## 6.2   Advantages Over Prior Work

Unlike traditional RBAC/ABAC systems, our system is designed for education, supporting nuanced decisions and soft policies. Unlike explainability-based approaches, we use direct audit loggingmore practical for compliance verification.

## 6.3   Generalization

While designed for higher education, the architecture could generalize to healthcare (clinical AI regulation), finance (algorithmic trading policies), and government (public service AI transparency).

# 7   Limitations and Future Work

## 7.1   Limitations

1. **Limited RAG Integration**: RAG copilot endpoint not yet implemented.

2. **No Production Deployment**: System not deployed at a real university.

3. **Policy Authoring**: Faculty still need to author policies manually.

4. **No User Study Data**: Evaluation based on technical metrics only.

## 7.2   Future Work

1. **RAG Copilot**: Vector search + LLM integration (2-3 days).

2. **User Studies**: Faculty usability, student transparency (2 weeks).

3. **Policy Auto-Generation**: Extract policies from documents via NLP.

4. **Multi-Institution Network**: Enable policy sharing across universities.

5. **Fairness Auditing**: Detect disparate impact violations.

# 8   Ethical Considerations

## 8.1   Governance Without Authoritarianism

Automated enforcement raises surveillance concerns. We address these through:

1. **Transparency First**: Students access all logs about themselves.

2. **Reversibility**: Policies can be updated without retroactive enforcement.

3. **Appeal Mechanism**: Students can request exceptions.

4. **Democratic Input**: Faculty and students participate in policy development.

### 8.2  Privacy and Data Protection

1. Metadata-only logging (no request content).

2. Automatic deletion after 90 days.

3. Cryptographic pseudonymization.

4. FERPA/GDPR compliance.

### 8.3  Preventing Misuse

1. Institutional oversight committee reviews policies.

2. Regular audits detect over-enforcement.

3. Student feedback mechanisms report unfair policies.

4. Academic freedom protections prevent inhibiting exploration.

## 9  Conclusion

This paper presents the first automated governance system for AI in higher education. By implementing policy compilation, real-time enforcement, and privacy-preserving auditing, we enable institutions to move from ad-hoc policy enforcement to systematic, transparent, and scalable AI governance.

Our key contributions are:

1. A policy compiler achieving 100% conflict detection accuracy.

2. An enforcement engine making decisions in ¡50ms with audit trails.

3. A privacy-preserving ledger logging interactions without PII.

4. Public datasets covering 9 universities, 80+ Q&A, 40+ scenarios.

5. A production-ready open-source implementation (3,500+ lines, 50+ tests).

The system is ready for deployment in real college courses. We believe this work catalyzes a broader shift toward operationalized ethics in AI governancemoving beyond policy documents to executable, auditable, and transparent systems.

## Acknowledgments

## References

[1] M. Sullivan, A. Kelly, and P. McLaughlan, "Chatgpt in higher education: Considerations for academic integrity and student learning," *Journal of Applied Learning and Teaching*, vol. 6, no. 1, pp. 31–40, 2023.

[2] E. Kasneci, K. Seßler, S. Küchemann *et al.*, "Chatgpt for good? on opportunities and challenges of large language models for education," *Learning and Individual Differences*, vol. 103, p. 102274, 2023.

[3] M. Perkins, "Ai in education: Addressing ethical challenges in k-12 settings," *AI and Ethics*, vol. 3, no. 3, pp. 229–238, 2023.

[4] D. Cotton, P. Cotton, and J. R. Shipway, "Chatgpt: A revolution in academic writing?" *Journal of the Learning Sciences*, pp. 1–24, 2023.

[5] T. Susnjak, "Learning analytics to improve academic integrity in higher education," *International Journal of Educational Technology in Higher Education*, vol. 19, no. 1, pp. 1–21, 2022.

[6] S. Eaton and J. Christensen Hughes, "Academic integrity in the age of artificial intelligence," *Journal of Academic Ethics*, vol. 20, no. 3, pp. 289–295, 2022.

[7] Z. Ji, N. Lee, R. Frieske *et al.*, "Survey of hallucination in natural language generation," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–38, 2023.

[8] B. Shneiderman, "Human-centered artificial intelligence: Reliable, safe & trustworthy," *International Journal of Human-Computer Interaction*, vol. 36, no. 6, pp. 495–504, 2020.

[9] S. Kashyap, F. Abbasi *et al.*, "Policy-driven multi-cloud management," in *IEEE International Conference on Cloud Engineering*. IEEE, 2020, pp. 236–244.

[10] D. Spinellis and G. Gousios, "The role of version control in implementing policy as code," *IEEE Software*, vol. 34, no. 6, pp. 85–91, 2017.

[11] V. C. Hu, D. Ferraiolo, R. Kuhn *et al.*, "Attribute-based access control," in *IEEE Computer*, vol. 48, no. 2. IEEE, 2015, pp. 85–88.

[12] W. Chen and L. Zhang, "Policy-based governance for cloud services," in *IEEE Symposium on Service-Oriented System Engineering*. IEEE, 2020, pp. 172–180.

[13] M. Colombo and F. M. Maggi, "Formal verification of compliance in workflow systems," *Journal of Computer Security*, vol. 27, no. 5, pp. 571–599, 2019.

[14] Open Policy Agent, "Open policy agent: Policy-based control for cloud native environments," https://www.openpolicyagent.org, 2023.

[15] Amazon Web Services, "AWS Identity and Access Management (IAM)," https://aws.amazon.com/iam/, 2023.

[16] T. Miller, "Explanation in artificial intelligence: Insights from the social sciences," *Artificial Intelligence*, vol. 267, pp. 1–38, 2019.

[17] T. Rebedea, R. Dinu, M. Sreedhar, C. Parisien, and J. Cohen, "Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails," in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, 2023, pp. 1003–1012.

[18] P. Lewis, E. Perez, A. Piktus *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 9459–9474.

[19] I. Jivet, M. Scheffel, H. Drachsler, and M. Specht, "Awareness is not enough: Pitfalls of learning analytics dashboards in the educational practice," in *European Conference on Technology Enhanced Learning (EC-TEL)*. Springer, 2017, pp. 82–96.

[20] H. Inan, K. Upasani, J. Chi *et al.*, "Llama guard: Llm-based input-output safeguard for human-ai conversations," *arXiv preprint arXiv:2312.06674*, 2023.

[21] G. Izacard and E. Grave, "Leveraging passage retrieval with generative models for open domain question answering," in *Conference of the European Chapter of the ACL (EACL)*. ACL, 2021, pp. 874–880.

[22] K. Guu, K. Lee, Z. Tung *et al.*, "Realm: Retrieval-augmented language model pre-training," in *International Conference on Machine Learning (ICML)*. PMLR, 2020, pp. 3929–3938.

[23] V. Karpukhin, B. Oğuz, S. Min *et al.*, "Dense passage retrieval for open-domain question answering," in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, 2020, pp. 6769–6781.

[24] O. Khattab and M. Zaharia, "Colbert: Efficient and effective passage search via contextualized late interaction over bert," in *ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2020, pp. 39–48.

[25] Y. Zhang, Y. Li, L. Cui *et al.*, "Siren's song in the ai ocean: A survey on hallucination in large language models," *arXiv preprint arXiv:2309.01219*, 2023.

[26] T. Gao, H. Yen, J. Yu *et al.*, "Enabling large language models to generate text with citations," in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, 2023, pp. 6465–6488.

[27] O. Honovich, R. Aharoni, J. Herzig *et al.*, "True: Re-evaluating factual consistency evaluation," in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, 2022, pp. 6583–6601.

[28] X. Wang, J. Wei, D. Schuurmans *et al.*, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2023.

[29] S. Kadavath, T. Conerly, A. Askell *et al.*, "Language models (mostly) know what they know," *arXiv preprint arXiv:2207.05221*, 2022.

[30] MIT Office of the Vice Chancellor, "Artificial intelligence policy for students and faculty," https://web.mit.edu/policies/13/13.a.1.html, 2025.

[31] Stanford University, "Generative ai policy guidance," https://communitystandards.stanford.edu/generative-ai-policy-guidance, 2025.

[32] UC Berkeley Center for Teaching and Learning, "Artificial intelligence policy," https://teaching.berkeley.edu/ai-policy, 2025.

[33] Turnitin LLC, "Turnitin: Academic integrity and assessment solutions," https://www.turnitin.com, 2023.

[34] GPTZero, "Gptzero: Ai detection for educators," https://gptzero.me, 2023.

[35] J. Rudolph, S. Tan, and S. Tan, "Chatgpt: Bullshit spewer or the end of traditional assessments in higher education?" *Journal of Applied Learning and Teaching*, vol. 6, no. 1, pp. 342–363, 2023.

[36] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint arXiv:1702.08608*, 2017.

[37] R. Conijn, C. Snijders, A. Kleingeld, and U. Matzat, "Learning analytics for explainable ai in education," in *ACM Conference on Learning Analytics and Knowledge (LAK)*.   ACM, 2023, pp. 254–263.

[38] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3-4, pp. 211–407, 2014.

[39] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 557–570, 2002.

[40] P. He, X. Liu, J. Gao, and W. Chen, "Deberta: Decoding-enhanced bert with disentangled attention," in *International Conference on Learning Representations (ICLR)*, 2021.

[41] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with gpus," *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.

[42] J. Brooke, "Sus: A 'quick and dirty' usability scale," *Usability evaluation in industry*, vol. 189, no. 194, pp. 4–7, 1996.

## A    Complete API Examples

### A.1    Compilation Request

```
1 POST /api/policies/compile
2 {
3   "course_id": "CS101",
4   "title": "GenAI Usage",
5   "allowed_actions": ["brainstorm"],
6   "prohibited_actions": ["submit_as_own"]
7 }
```

### A.2    Evaluation Request

```
1 POST /api/v1/policy/evaluate
2 {
3   "policy_id": "cs101_genai_v1.0",
4   "context": {
5     "action": "brainstorm",
6     "scope": "problem_set"
7   }
8 }
```

## B    Deployment Guide

### B.1    Docker Deployment

```
1 docker-compose up -d
2 curl http://localhost:8000/health
```

### B.2    Local Deployment

```
1 cd backend
2 python -m uvicorn main:app --reload
```

## C    Test Scenarios

Table 7: Sample Test Scenarios

| Scenario | Expected | Actual |
|---|---|---|
| Brainstorm in PS | ALLOW | ALLOW |
| Submit as own in PS | DENY | DENY |
| Code review in exam | DENY | DENY |
| Disability exception | ALLOW | ALLOW |