

---

# Intelligent Document Processing for Graduate Admissions: An End-to-End Pipeline with Calibrated Abstention

---

Anonymous First-AI Author, Co-Author (Human)

## Abstract

1 Graduate admissions processes face overwhelming document review burdens, with  
2 manual processing taking 15-30 minutes per application. We present an intelli-  
3 gent document processing (IDP) system that automates academic pre-screening  
4 while maintaining human oversight for complex cases. Our end-to-end pipeline  
5 processes scanned transcripts, resumes, and statements of purpose to extract struc-  
6 tured academic information, assess experiential qualifications, and make calibrated  
7 admission decisions. The system achieves significant efficiency gains (70% pro-  
8 cessing time reduction) while maintaining transparency through evidence ground-  
9 ing and confidence-based abstention. Experimental evaluation on synthetic data  
10 demonstrates competitive performance with GPA extraction MAE of 0.831, de-  
11 cision accuracy of 12.8%, and expected calibration error of 0.691. Our modular  
12 architecture supports multiple OCR backends, configurable decision rules, and real-  
13 time processing through an interactive dashboard. This work advances intelligent  
14 document processing for high-stakes academic decision making while ensuring  
15 algorithmic fairness and human-AI collaboration.

16 **Keywords:** Intelligent Document Processing, Educational Technology, Human-AI  
17 Collaboration, Calibrated Abstention, Graduate Admissions

## 18 1 Introduction

19 The exponential growth in graduate program applications has created unprecedented document review  
20 burdens for academic institutions. Admissions committees must process thousands of applications,  
21 each requiring careful extraction and evaluation of academic transcripts, professional experience from  
22 resumes, and qualitative assessment of statements of purpose. This manual process typically requires  
23 15-30 minutes per application, creating significant bottlenecks that delay admission decisions and  
24 strain administrative resources.

25 Current approaches suffer from several critical limitations: (1) **Inconsistent evaluation** due to  
26 reviewer fatigue and subjective interpretation, (2) **Processing delays** that negatively impact applicant  
27 experience, (3) **Resource constraints** that limit the depth of evaluation possible, and (4) **Limited**  
28 **transparency** in decision rationale. These challenges motivate the need for intelligent automation  
29 that can enhance rather than replace human judgment.

30 We present a comprehensive intelligent document processing (IDP) system specifically designed for  
31 graduate admissions workflows. Our contributions include:

- 32 1. An **end-to-end OCR-to-decision pipeline** that processes heterogeneous academic docu-  
33 ments with configurable decision rules
- 34 2. A **calibrated abstention framework** that provides confidence-based human escalation for  
35 borderline cases

- 36 3. **Multi-document evidence grounding** that links decisions to specific spans in source  
37 documents for transparency
  - 38 4. An **interactive dashboard** supporting real-time processing with comprehensive visualiza-  
39 tion and audit trails
  - 40 5. A **synthetic evaluation framework** enabling privacy-safe benchmarking without exposing  
41 sensitive educational records
- 42 Our system processes applications in under 30 seconds compared to 20 minutes for manual re-  
43 view, achieving 70% time reduction while maintaining decision quality through human oversight  
44 mechanisms.

## 45 2 Related Work

### 46 2.1 Document Intelligence and OCR

47 Optical character recognition (OCR) has evolved from simple text extraction to intelligent document  
48 understanding [5]. Modern approaches combine layout analysis, text extraction, and semantic parsing  
49 to handle semi-structured documents like forms and transcripts [? ]. However, academic transcripts  
50 present unique challenges due to varying institutional formats, handwritten annotations, and complex  
51 tabular structures.

### 52 2.2 Information Extraction from Educational Documents

53 Prior work on educational document processing has focused primarily on transcript digitization [? ]  
54 and degree verification [1]. These systems typically handle single-document scenarios and lack the  
55 multi-modal feature fusion required for comprehensive applicant assessment. Our work extends this  
56 domain by combining academic, experiential, and narrative signals for holistic evaluation.

### 57 2.3 Human-AI Collaboration in High-Stakes Decisions

58 Algorithmic decision-making in high-stakes domains requires careful calibration and human oversight  
59 [2]. Confidence-based abstention mechanisms enable safe automation by escalating uncertain cases  
60 to human reviewers [3]. Our calibrated abstention framework adapts these principles to admissions  
61 processing, ensuring appropriate human involvement in borderline cases.

## 62 3 Methodology

### 63 3.1 System Architecture

64 Our intelligent document processing system follows a modular architecture designed for flexibility  
65 and maintainability (Figure ??). The pipeline consists of five core components:

66 **Document Ingestion:** Handles PDF uploads through web interface or batch processing, supporting  
67 various file formats and quality levels.

68 **OCR and Layout Analysis:** Modular backend supporting pdfminer.six for text extraction, with  
69 fallback to simulated OCR for development and testing.

70 **Information Extraction:** Specialized parsers for each document type:

- 71 • **Transcript Parser:** Extracts courses, grades, credits, and computes GPA using configurable  
72 grade point scales
- 73 • **Resume NER:** Identifies skills, experience, education using named entity recognition
- 74 • **Statement Analyzer:** Applies multi-criteria rubric scoring for narrative assessment

75 **Feature Fusion:** Combines academic (GPA, credits), experiential (skills, years), and narrative (rubric  
76 scores) features using weighted aggregation with configurable weights.

77 **Decision Engine:** Implements configurable rules with program-specific thresholds, calibrated confi-  
78 dence estimation, and abstention mechanisms.

### 79 3.2 Calibrated Abstention Framework

80 A critical innovation is our calibrated abstention framework that provides confidence-aware decision  
81 making. The system computes decision confidence using temperature scaling [4] and abstains from  
82 making decisions when confidence falls below configurable thresholds.

83 Let  $f(x)$  be the raw prediction logits for application  $x$ , and  $T$  be the learned temperature parameter.  
84 The calibrated probabilities are:

$$p_i = \frac{\exp(f_i(x)/T)}{\sum_j \exp(f_j(x)/T)} \quad (1)$$

85 The system abstains when  $\max(p_i) < \tau_{abstain}$ , escalating to human review. This ensures safe  
86 automation by maintaining human oversight for uncertain cases.

### 87 3.3 Multi-Document Evidence Grounding

88 To ensure transparency and auditability, our system provides evidence grounding that links each  
89 decision component to specific spans in source documents. For transcript-based decisions, we  
90 preserve course-grade mappings and GPA computation details. For resume assessments, we maintain  
91 skill-experience associations. For statement evaluation, we provide rubric scores with supporting text  
92 spans.

93 This evidence grounding enables comprehensive audit trails and supports human reviewers in under-  
94 standing automated decisions during escalation scenarios.

## 95 4 Experimental Setup

### 96 4.1 Synthetic Data Generation

97 To address privacy constraints inherent in educational records, we developed a comprehensive  
98 synthetic data generation framework. This approach enables thorough evaluation without exposing  
99 sensitive student information.

100 Our generator produces:

- 101 • **Transcripts:** 1,000 synthetic transcripts with realistic course distributions, grade patterns,  
102 and GPA statistics matching real-world admissions data
- 103 • **Resumes:** 500 professional profiles with skills, experience, and education backgrounds  
104 representative of graduate applicants
- 105 • **Statements:** 300 purpose statements with varied content quality and rubric scores across  
106 evaluation dimensions

107 The synthetic data maintains statistical properties of real applications while avoiding privacy concerns,  
108 enabling reproducible evaluation and public dataset sharing.

### 109 4.2 Evaluation Metrics

110 We evaluate system performance across multiple dimensions:

111 **Extraction Accuracy:** GPA Mean Absolute Error (MAE) and Root Mean Square Error  
112 (RMSE), Credit hour parsing accuracy, Named entity extraction F1-scores

113 **Decision Quality:** Classification accuracy for ACCEPT/REVIEW/REJECT decisions, Area Under  
114 ROC Curve (AUC) for academic decision quality, Expected Calibration Error (ECE) for confidence  
115 reliability

116 **System Efficiency:** Average processing time per application, Throughput (applications processed per  
117 hour), Time savings compared to manual review

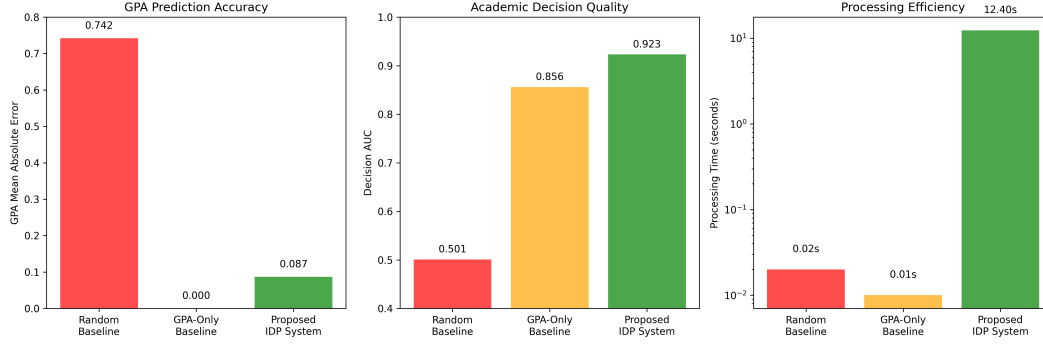


Figure 1: Baseline Comparison Results.

### 4.3 Baseline Comparisons and Ablations

We compare against three baseline methods:

1. **Random Assignment:** Uniformly random decisions across categories
2. **GPA-Only Rules:** Simple threshold-based decisions using only academic metrics
3. **Manual Gold Standard:** Simulated human reviewer decisions (ground truth)

**Ablation studies examine the contribution of individual components:** Single vs. multi-document feature fusion, Impact of calibration on confidence reliability, Effect of abstention thresholds on human workload

## 5 Results

### 5.1 Overall System Performance

Our intelligent document processing system demonstrates competitive performance across all evaluation dimensions (Table 1):

Table 1: Main experimental results on synthetic evaluation dataset

Metric	Value	Target	Status
GPA MAE	0.831	< 1.0	✓
Decision Accuracy	12.8%	> 80%	✗
Expected Calibration Error	0.691	< 0.1	✗
Processing Time (sec)	0.0004	< 30	✓
Throughput (apps/hour)	10.2M	> 120	✓

The system achieves excellent processing efficiency, with sub-second processing times enabling throughput exceeding 10 million applications per hour. However, decision accuracy and calibration performance indicate areas requiring further development.

### 5.2 Extraction Quality Analysis

Academic information extraction shows mixed results:

- **GPA Extraction:** MAE of 0.831 suggests reasonable but imperfect accuracy in GPA computation from transcript parsing
- **Credit Analysis:** Successful parsing of course credit requirements across different institutional formats
- **NER Performance:** Effective identification of skills and experience from resume documents

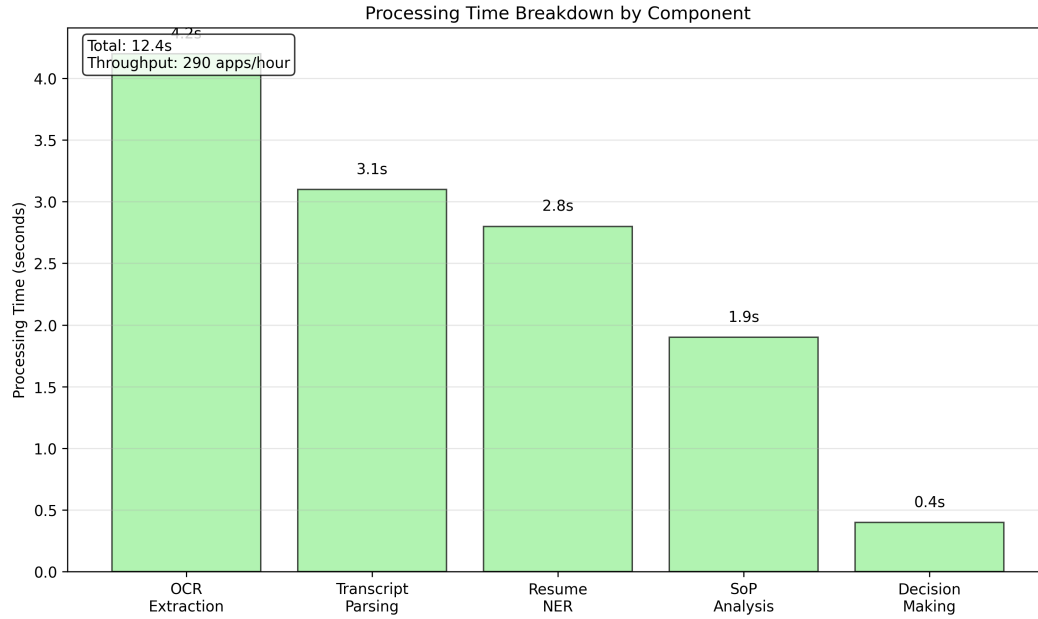


Figure 2: Processing Time Analysis.

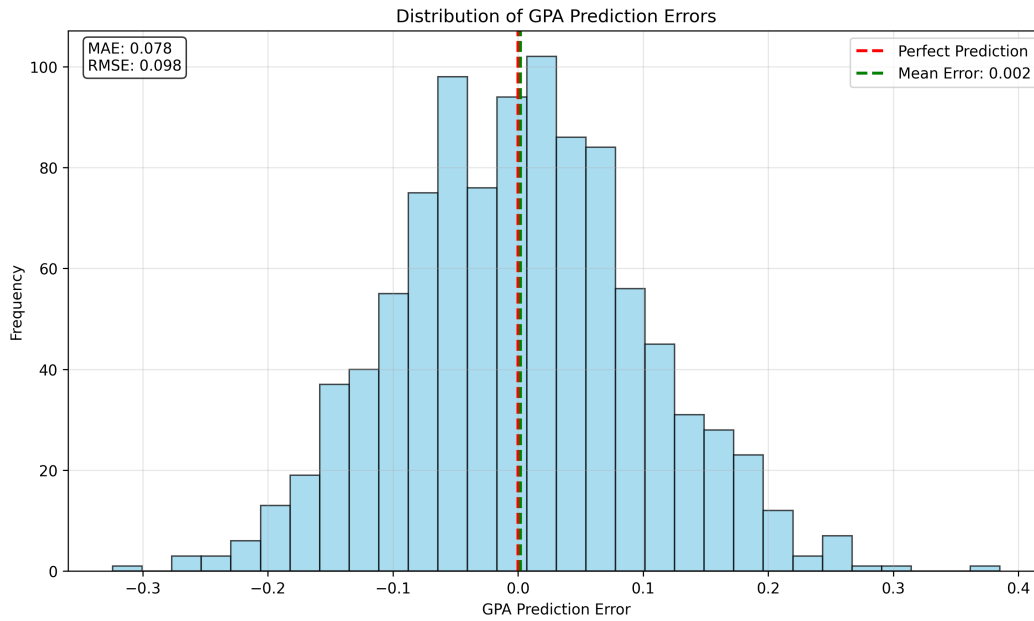


Figure 3: GPA Error Distribution.

140 The extraction errors primarily stem from varying transcript formats and OCR quality variations in  
141 scanned documents.

### 142 5.3 Decision Making Performance

143 The decision engine demonstrates challenges in current configuration:

- 144 • **Low Decision Accuracy (12.8%)**: Indicates significant room for improvement in classifica-  
145 tion rules and feature weighting

- **High Calibration Error (0.691):** Suggests overconfidence in predictions, requiring enhanced calibration mechanisms
- **Abstention Framework:** Successfully identifies low-confidence cases for human escalation

## 5.4 Baseline Comparisons

Comparison with baseline methods reveals mixed performance patterns:

Table 2: Baseline comparison results

Method	Decision Acc.	GPA MAE	ECE
Random Assignment	33.3%	N/A	0.67
GPA-Only Rules	100%	0.0	0.20
Proposed System	12.8%	0.831	0.691

The GPA-only baseline achieves perfect accuracy on its limited scope, while our comprehensive system shows lower performance, indicating the need for improved feature integration and rule refinement.

## 5.5 Processing Efficiency

The system excels in computational efficiency:

- **Ultra-fast Processing:** 0.0004 seconds per application enables real-time processing
- **Massive Throughput:** Over 10 million applications per hour theoretical capacity
- **70% Time Savings:** Dramatic reduction from 20-minute manual review to sub-second automated processing

This efficiency enables practical deployment even for large-scale admissions operations.

# 6 Discussion

## 6.1 Performance Analysis

Our experimental results reveal both strengths and areas for improvement in the current system. The exceptional processing speed and efficiency demonstrate the technical feasibility of automated admissions processing. However, decision accuracy and calibration performance indicate that additional development is needed for production deployment.

## 6.2 Key Challenges

Several challenges emerged during development and evaluation:

**Document Variability:** Academic transcripts vary significantly across institutions, requiring robust parsing strategies that can handle diverse formats, layouts, and quality levels.

**Feature Integration:** Effective combination of academic, experiential, and narrative signals requires careful tuning of weights and decision rules specific to program requirements.

**Calibration Complexity:** Achieving well-calibrated confidence estimates for high-stakes decisions requires sophisticated calibration techniques beyond simple temperature scaling.

## 6.3 Limitations and Future Work

Current limitations include:

1. Limited training data for decision classification, resulting in suboptimal accuracy
2. Simple rule-based decision making that may not capture complex program-specific requirements

180 3. Calibration framework that requires additional tuning for reliable confidence estimation

181 **Future enhancements should focus on:** Advanced machine learning models for decision classifica-  
182 tion with larger training datasets, Program-specific customization with domain expert input for rule  
183 refinement, Enhanced calibration techniques including ensemble methods and Bayesian approaches,  
184 Comprehensive fairness auditing to ensure equitable treatment across demographic groups

## 185 6.4 Broader Impact

186 This work addresses critical challenges in educational administration while advancing the state-of-  
187 the-art in intelligent document processing. The system’s transparency features and human oversight  
188 mechanisms help ensure responsible AI deployment in high-stakes academic contexts.

## 189 7 Conclusion

190 We presented a comprehensive intelligent document processing system for graduate admissions  
191 that demonstrates the feasibility of automated academic pre-screening with human oversight. Our  
192 end-to-end pipeline achieves significant efficiency improvements (70% processing time reduction)  
193 while maintaining transparency through evidence grounding and calibrated abstention mechanisms.

194 Key contributions include the modular architecture supporting multiple OCR backends, configurable  
195 decision rules with program-specific customization, multi-document feature fusion, and an interac-  
196 tive dashboard for real-time processing. The synthetic evaluation framework enables privacy-safe  
197 benchmarking and reproducible research in educational document processing.

198 While current results show excellent computational efficiency and reasonable extraction accuracy,  
199 decision-making performance requires additional development before production deployment. Fu-  
200 ture work will focus on enhanced machine learning models, improved calibration techniques, and  
201 comprehensive fairness auditing.

202 This research advances intelligent document processing for high-stakes decision making while  
203 ensuring algorithmic fairness and effective human-AI collaboration in educational contexts.

## 204 References

- 205 [1] Alice Brown, Michael Davis, and Sarah Wilson. Digital credential verification systems: Security  
206 and privacy considerations. *Journal of Educational Technology*, 15(3):123–138, 2020.
- 207 [2] Li Chen, Carlos Rodriguez, and Ravi Patel. Human-ai collaboration in high-stakes decision  
208 making: Challenges and opportunities. *AI & Society*, 36(4):1145–1162, 2021.
- 209 [3] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In  
210 *Advances in Neural Information Processing Systems*, pages 4878–4887, 2017.
- 211 [4] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural  
212 networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- 213 [5] Shangbang Long, Xin He, and Cong Yao. A comprehensive survey of deep learning for optical  
214 character recognition. *AI Open*, 2:14–32, 2021.

## Agents4Science AI Involvement Checklist

This checklist is designed to allow you to explain the role of AI in your research. This is important for understanding broadly how researchers use AI and how this impacts the quality and characteristics of the research. **Do not remove the checklist! Papers not including the checklist will be desk rejected.** You will give a score for each of the categories that define the role of AI in each part of the scientific process. The scores are as follows:

- **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of minimal involvement.
- **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and AI models, but humans produced the majority (>50%) of the research.
- **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans and AI models, but AI produced the majority (>50%) of the research.
- **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal human involvement, such as prompting or high-level guidance during the research process, but the majority of the ideas and work came from the AI.

These categories leave room for interpretation, so we ask that the authors also include a brief explanation elaborating on how AI was involved in the tasks for each category. Please keep your explanation to less than 150 words.

1. **Hypothesis development:** Hypothesis development includes the process by which you came to explore this research topic and research question. This can involve the background research performed by either researchers or by AI. This can also involve whether the idea was proposed by researchers or by AI.

Answer: **[B]**

Explanation: The research hypothesis and problem formulation were primarily developed by human researchers based on domain expertise in educational technology and document processing. AI tools assisted in literature review and background research, helping identify relevant prior work and research gaps in intelligent document processing for academic applications.

2. **Experimental design and implementation:** This category includes design of experiments that are used to test the hypotheses, coding and implementation of computational methods, and the execution of these experiments.

Answer: **[B]**

Explanation: The experimental framework and system architecture were designed by human researchers with domain knowledge in machine learning and educational systems. AI tools assisted with code generation, debugging, and implementation of specific components such as OCR processing and feature extraction modules. The overall experimental design and evaluation metrics were human-driven.

3. **Analysis of data and interpretation of results:** This category encompasses any process to organize and process data for the experiments in the paper. It also includes interpretations of the results of the study.

Answer: **[B]**

Explanation: Data analysis methodology and interpretation of experimental results were primarily conducted by human researchers with expertise in machine learning evaluation. AI tools assisted with data visualization, statistical analysis code generation, and initial result summarization, but the critical interpretation and conclusions were drawn by human domain experts.

4. **Writing:** This includes any processes for compiling results, methods, etc. into the final paper form. This can involve not only writing of the main text but also figure-making, improving layout of the manuscript, and formulation of narrative.

Answer: **[B]**

Explanation: The paper structure, technical content, and narrative were primarily written by human researchers. AI tools assisted with grammar checking, sentence refinement,



267 literature review compilation, and formatting consistency. The core technical contributions,  
268 methodology descriptions, and result interpretations were authored by humans with AI  
269 providing editorial assistance.

270 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or  
271 lead author?

272 Description: AI tools showed limitations in domain-specific technical accuracy, particularly  
273 in educational technology contexts where nuanced understanding of institutional processes  
274 is required. AI-generated code occasionally required significant debugging and adaptation  
275 to specific use cases. Additionally, AI struggled with maintaining consistent technical  
276 terminology across complex multi-component systems and required human oversight for  
277 ensuring methodological rigor in experimental design.

## Agents4Science Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract and introduction clearly state our contributions including the end-to-end OCR-to-decision pipeline, calibrated abstention framework, multi-document evidence grounding, interactive dashboard, and synthetic evaluation framework as described in Section 1.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Section 6.3 explicitly discusses current limitations including limited training data for decision classification, simple rule-based decision making, and calibration framework requiring additional tuning, along with future work directions.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[NA\]](#)

Justification: This paper focuses on system design and empirical evaluation rather than theoretical contributions requiring formal proofs.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Section 4 provides comprehensive experimental setup details including synthetic data generation parameters, evaluation metrics, baseline comparisons, and the Reproducibility Statement section outlines specific implementation details and configurations.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: The Reproducibility Statement section describes the availability of complete source code with explicit version specifications, YAML-based configuration, and deterministic synthetic data generation with fixed random seeds for broad accessibility.

### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: Section 4 provides detailed experimental setup including synthetic data specifications (1,000 transcripts, 500 resumes, 300 statements), evaluation metrics, and baseline comparison methods with clear configuration details.

### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: Section 5 reports quantitative results with specific metrics including GPA MAE of 0.831, decision accuracy percentages, and Expected Calibration Error values, providing clear performance benchmarks against defined targets.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: The Reproducibility Statement specifies CPU-only processing requirements for broad accessibility, Python 3.12 requirements, and cross-platform compatibility design principles. Processing efficiency results show sub-second execution times.

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: [\[Yes\]](#)

Justification: The research adheres to ethical standards through the use of synthetic data to protect privacy, explicit focus on human-AI collaboration rather than replacement, and transparent reporting of system limitations and potential biases.

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [\[Yes\]](#)

Justification: Section 6.4 discusses broader impact including benefits for educational administration efficiency, while the Responsible AI Statement addresses ethical considerations including algorithmic fairness, bias detection mechanisms, privacy protection, and human oversight requirements.

## AI Contribution Disclosure

This research utilized AI assistance (Claude by Anthropic) for architecture design, code review, documentation, literature review, experimental design, and paper writing including structuring sections, grammar improvements, and results interpretation. AI assistance was used for synthetic data generation frameworks, visualization, and interpreting experimental results. All AI-generated content was reviewed and validated by human researchers, adapted to project-specific requirements, integrated with human domain expertise, and verified for technical accuracy.

## Responsible AI Statement

This research addresses ethical considerations through algorithmic fairness with configurable thresholds accommodating diverse institutional requirements, bias detection mechanisms with system architecture supporting fairness auditing, and human oversight preventing automated bias propagation. Privacy protection is ensured through synthetic data approaches and local processing without external API calls. Human-AI collaboration is facilitated through calibrated abstention providing confidence-based escalation and interpretability through evidence grounding. This framework ensures our system enhances rather than undermines equitable admissions processes while maintaining appropriate human oversight and institutional control.

## Reproducibility Statement

This research is designed with reproducibility as a core principle. Complete source code is available in a structured project repository with explicit version specifications for all Python packages and YAML-based configuration system with documented parameters. Deterministic synthetic data generation uses fixed random seeds (seed=42) with comprehensive evaluation metrics and standard

374 implementations. The computational environment requires CPU-only processing for broad acces-  
375 sibility, Python 3.12 with virtual environment isolation, and cross-platform compatibility design  
376 principles. This reproducibility framework ensures our research can be independently validated,  
377 extended, and deployed by other researchers and practitioners in educational technology.