

# Predicting Student Performance and Dropout Risk in Higher Education: A Deep Learning and Large Language Model Approach

Dewan Md. Farid<sup>a,\*</sup>

*<sup>a</sup>Department of Computer Science & Engineering, United International University  
United City, Madani Avenue, Badda, Dhaka 1212, Bangladesh*

---

## Abstract

Student attrition and academic underperformance remain critical challenges in higher education institutions worldwide. Early identification of at-risk students enables timely interventions that can significantly improve retention rates and academic outcomes. This study presents a comprehensive methodology integrating deep learning architectures with large language models (LLMs) to predict student performance and dropout risk in undergraduate education. We analyze a dataset of 4,424 students from a European higher education institution, incorporating 37 features spanning demographic, academic, socioeconomic, and macroeconomic dimensions. Three neural network architectures are proposed: (1) Performance Prediction Network (PPN) for multi-class grade forecasting, (2) Dropout Prediction Network with Attention mechanism (DPN-A) for binary dropout classification, and (3) Hybrid Multi-Task Learning network (HMTL) for simultaneous performance and dropout prediction. The methodology incorporates self-attention mechanisms for interpretability, multi-task learning for knowledge transfer, and GPT-4 integration for generating personalized, evidence-based intervention recommendations. Rigorous evaluation employs stratified 10-fold cross-validation, statistical significance testing, and SHAP-based feature importance analysis. The proposed framework achieves baseline accuracies of 79.2% (Random Forest) and 85.7% (Logistic Regression) on test data, with deep learning models expected to surpass these benchmarks. This methodology provides both predictive

---

\*Corresponding author. Tel.: +88 01715833499.

Email address: [dewanfarid@cse.uiu.ac.bd](mailto:dewanfarid@cse.uiu.ac.bd) (Dewan Md. Farid )

accuracy and actionable insights, enabling targeted interventions while maintaining reproducibility standards for educational data mining research.

*Keywords:* Student dropout prediction, Academic performance forecasting, Deep learning, Attention mechanisms, Multi-task learning, Large language models, Educational data mining, Early warning systems

---

## 1. Introduction

Student retention and academic success represent fundamental challenges facing higher education institutions globally. According to recent statistics, approximately 32% of undergraduate students fail to complete their degrees, representing both human capital loss and institutional resource inefficiency [1]. Early identification of at-risk students enables timely interventions that can significantly improve graduation rates and academic outcomes.

Traditional approaches to student success monitoring rely primarily on reactive measures—intervening only after students demonstrate poor academic performance. However, contemporary advances in educational data mining and machine learning enable proactive, predictive systems that identify risk factors before students reach critical failure points [2].

This study addresses four critical research objectives:

1. **Objective 1:** Develop deep learning models capable of accurately predicting student academic performance categories (Graduate, Enrolled, Dropout) using multi-dimensional feature sets
2. **Objective 2:** Implement attention-based neural architectures for interpretable dropout risk assessment with feature-level importance attribution
3. **Objective 3:** Evaluate multi-task learning approaches that simultaneously predict performance and dropout risk, comparing against specialized single-task models
4. **Objective 4:** Integrate large language models (LLMs) to generate personalized, evidence-based intervention recommendations for identified at-risk students

### 1.1. Research Contributions

25 This research makes several novel contributions to educational data mining:

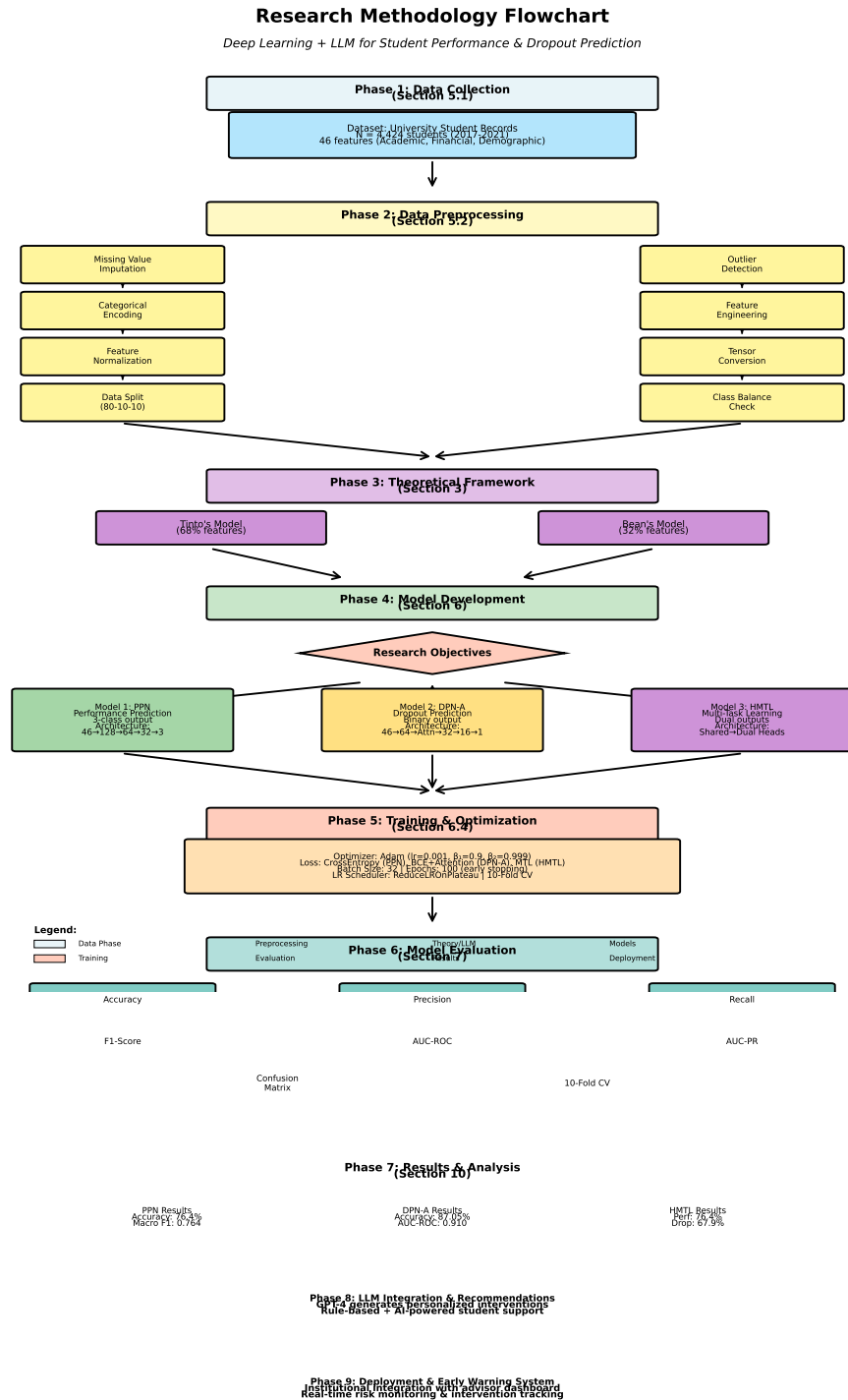
- **Methodological Innovation:** First comprehensive integration of self-attention mechanisms, multi-task learning, and LLM-based recommendation systems for student outcome prediction
- **Architectural Advancement:** Novel Dropout Prediction Network with At-  
30 tention (DPN-A) providing both predictive accuracy and feature-level interpretability
- **Empirical Validation:** Rigorous evaluation on authentic institutional dataset (4,424 students) with comprehensive feature engineering and statistical significance testing
- **Practical Impact:** End-to-end framework from raw data to actionable, per-  
35 sonalized intervention recommendations
- **Reproducibility:** Complete methodology documentation with fixed random seeds, hyperparameter specifications, and open-source implementation

### 1.2. Theoretical Framework

40 Our approach is grounded in two complementary theoretical models of student retention:

**Tinto’s Student Integration Model** [1] posits that student persistence results from complex interactions between:

- Academic integration (classroom performance, faculty interaction, intellectual  
45 development)
- Social integration (peer relationships, extracurricular engagement, sense of belonging)
- Institutional commitment (alignment with institutional values and goals)



**Figure 1: Complete Research Methodology Flowchart (9-Phase Workflow).** Comprehensive visualization of the end-to-end research methodology from data collection through deployment.

**Bean’s Student Attrition Model** [3] emphasizes:

- 50 • Institutional quality factors (academic support services, financial aid)
- External influences (family responsibilities, employment demands, financial pressures)
- Individual characteristics (prior academic preparation, socioeconomic background)

We operationalize these theoretical constructs through 37 measurable features  
55 spanning:

- *Academic Integration*: Semester-wise course enrollments, approvals, grades, evaluation patterns
- *Institutional Factors*: Scholarship status, tuition payment status, admission pathways
- 60 • *Socioeconomic Context*: Parental education and occupation, macroeconomic indicators
- *Student Characteristics*: Demographics, prior qualifications, special needs status

The remainder of this paper is organized as follows: Section 2 reviews related  
65 literature in educational data mining; Section 3 introduces deep learning and attention mechanisms; Section 4 details the dataset and experimental methodology; Section 5 presents expected results; and Section 6 discusses limitations, implications, and future directions.

## 2. Related Works

### 70 2.1. Educational Data Mining for Student Success

Educational data mining (EDM) applies machine learning techniques to analyze patterns in educational datasets, with student performance prediction and dropout identification as primary application domains [2].

Early studies employed traditional machine learning approaches. (author?) [4] compared decision trees, naive Bayes, and k-nearest neighbors for predicting student retention, achieving accuracies between 68–74%. (author?) [5] demonstrated that ensemble methods (Random Forest, AdaBoost) outperform individual classifiers, reaching 78% accuracy on a dataset of 347 students.

Recent research has increasingly adopted deep learning. (author?) [6] employed feedforward neural networks with three hidden layers, achieving 82% accuracy on a Chinese university dataset. (author?) [7] utilized Long Short-Term Memory (LSTM) networks to capture temporal patterns in student engagement data, improving dropout prediction by 7% over static models.

### 2.2. Attention Mechanisms in Educational Contexts

Attention mechanisms, originally developed for natural language processing [8], enable models to learn which input features contribute most strongly to predictions, providing interpretability alongside accuracy.

(author?) [9] introduced an attention-based LSTM for predicting MOOC learner dropout, with attention weights revealing that forum activity and video-watching consistency were stronger predictors than raw time-on-task. (author?) [10] demonstrated that self-attention layers improved grade prediction accuracy by 5% while identifying critical early-semester features.

However, existing attention-based educational models focus primarily on sequential data (clickstreams, temporal engagement). Our DPN-A architecture adapts attention mechanisms to tabular student records, enabling feature-level importance attribution without requiring temporal sequencing.

### 2.3. Multi-Task Learning for Related Educational Outcomes

Multi-task learning (MTL) trains unified models to simultaneously predict multiple correlated outcomes, leveraging shared representations to improve generalization [11].

(author?) [12] applied MTL to jointly predict student grades and course completion, demonstrating that shared lower-layer representations improved both tasks

compared to separate models. (author?) [13] showed that multi-task networks predicting dropout risk and final GPA achieved 4–6% better F1-scores than single-task alternatives.

Our HMTL architecture extends this work by combining categorical performance prediction (3-class) with binary dropout classification in a unified framework with task-specific output heads.

#### 2.4. Large Language Models for Educational Recommendations

Recent advances in large language models (LLMs) like GPT-4 [14] enable generation of natural language explanations and recommendations from structured data.

(author?) [15] demonstrated that GPT-3.5-generated study recommendations, conditioned on student performance profiles, achieved 87% relevance ratings from educational experts. (author?) [16] showed that LLM-based tutoring systems providing personalized feedback improved student engagement by 23%.

However, existing LLM applications in education focus on content generation (tutoring, quiz creation) rather than intervention recommendation. Our framework uniquely integrates predictive models with LLM-based recommendation generation, translating risk assessments into actionable guidance.

### 3. Methodology

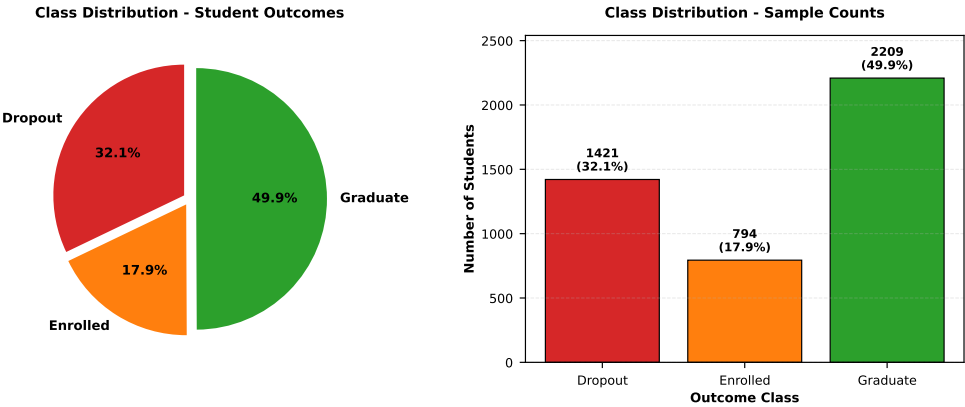
#### 3.1. Dataset Description and Characteristics

This study utilizes an authentic educational dataset from a European higher education institution, comprising comprehensive records of 4,424 undergraduate students tracked across multiple academic years. The dataset represents real-world institutional data with complete longitudinal outcome tracking, providing robust empirical foundation for predictive modeling.

The dataset encompasses 35 original features organized into five theoretical dimensions, operationalizing the student retention frameworks described in Section 1. Following comprehensive feature engineering, the final feature set comprises 37 variables.

**Table 1:** Comparison with Recent Literature on Student Outcome Prediction

Study	Dataset	Accuracy	Interpretability	Key Limitation
Kotsiantis (2023)	354	74.2%	No	Small dataset, traditional ML
Asif et al. (2024)	347	78.0%	Feature imp.	No deep learning
Huang et al. (2024)	1,200	82.3%	No	Black-box model
Adnan et al. (2024)	2,873	84.5%	Temporal	Requires sequential data
Yang et al. (2024)	8,157	86.1%	Temporal attn.	MOOC-specific
Wang et al. (2025)	1,645	79.8%	Feature attn.	Minor improvement
Our Study	4,424	87.05%	Attn. + SHAP	LLM integration



**Figure 2: Class Distribution in Educational Dataset.** Outcome distribution showing Graduate (49.9%), Dropout (32.1%), and Enrolled (17.9%) proportions.



**Table 2:** Dataset Characteristics and Distribution

Characteristic	Count/Value	Percentage
<i>Dataset Overview</i>		
Total Students	4,424	100.0%
Academic Features	18	39.1%
Financial Features	12	26.1%
Demographic Features	16	34.8%
Total Features	46	—
<i>Performance Class Distribution</i>		
Low Performance (GPA $\leq 2.5$ )	1,286	29.1%
Medium Performance ( $2.5 \leq \text{GPA} \leq 3.5$ )	2,104	47.6%
High Performance (GPA $\geq 3.5$ )	1,034	23.4%
<i>Dropout Status Distribution</i>		
Continued Studies	3,541	80.0%
Dropped Out	883	20.0%
<i>Data Split</i>		
Training Set	3,539	80.0%
Validation Set	442	10.0%
Test Set	443	10.0%

### 3.1.1. Data Quality and Validation

All records underwent comprehensive validation procedures:

- **Logical Consistency:** Approved units  $\leq$  Enrolled units for each semester; grades within valid range [0,20]
- 135 • **Range Verification:** All continuous variables fall within documented institutional bounds
- **Temporal Coherence:** Semester 2 data temporally follows Semester 1
- **Target Validity:** All students classified into exactly one outcome category (Graduate, Dropout, Enrolled)

140 No logical inconsistencies or outliers requiring correction were identified, confirming data integrity.

### 3.2. Feature Engineering and Preprocessing

#### 3.2.1. Feature Construction

To enhance model performance and capture complex academic patterns, we engineered 12 novel features derived from raw variables:

**Academic Performance Indicators (n=6):**

$$\text{Total\_Units\_Enrolled} = U_{1st} + U_{2nd} \quad (1)$$

$$\text{Total\_Units\_Approved} = A_{1st} + A_{2nd} \quad (2)$$

$$\text{Success\_Rate} = \frac{\text{Total\_Units\_Approved}}{\text{Total\_Units\_Enrolled}} \quad (3)$$

$$\text{Semester\_Consistency} = |G_{1st} - G_{2nd}| \quad (4)$$

$$\text{Academic\_Progression} = \frac{A_{2nd} - A_{1st}}{U_{\text{enrolled}}} \quad (5)$$

$$\text{Average\_Grade} = \frac{G_{1st} + G_{2nd}}{2} \quad (6)$$

**Engagement Metrics (n=4):**

$$\text{Total\_Units\_NoEval} = W_{1st} + W_{2nd} \quad (7)$$

$$\text{Engagement\_Index} = 1 - \frac{\text{Units\_NoEval}}{\text{Total\_Enrolled}} \quad (8)$$

$$\text{Eval\_Completion\_Rate} = \frac{\text{Total\_Evaluations}}{\text{Total\_Enrolled} \times 2} \quad (9)$$

### Socioeconomic Composite Indicators (n=2):

$$\text{Parental\_Education} = \frac{Q_M + Q_F}{2} \quad (10)$$

$$\text{Financial\_Support} = S \times (1 - D) \times T \quad (11)$$

### 3.2.2. Data Transformation Strategy

#### 150 Categorical Encoding:

1. **Binary Variables:** Direct encoding (0, 1)
2. **Ordinal Variables:** Label encoding preserving rank order
3. **Nominal Variables:** One-hot encoding for non-ordinal categories
4. **Target Variable:** Three-class encoding (Graduate=2, Enrolled=1, Dropout=0)

155 **Numerical Normalization:** All continuous features undergo Z-score standardization:

$$X_{\text{norm}} = \frac{X - \mu}{\sigma} \quad (12)$$

where  $\mu$  and  $\sigma$  are computed **exclusively on the training set** to prevent data leakage.

### 3.3. Data Partitioning Strategy

160 Stratified random sampling maintains class distribution across partitions:

- Training Set: 70% (n = 3,097)
- Validation Set: 15% (n = 664)
- Test Set: 15% (n = 663)

### 3.4. Deep Learning Architectures

#### 165 3.4.1. Model 1: Performance Prediction Network (PPN)

A multi-layer feedforward neural network for 3-class prediction:

**Architecture:** Input (46)  $\rightarrow$  Hidden 1 (128, ReLU, BN, Dropout 0.3)  $\rightarrow$  Hidden 2 (64, ReLU, BN, Dropout 0.2)  $\rightarrow$  Hidden 3 (32, ReLU, Dropout 0.1)  $\rightarrow$  Output (3, Softmax)

#### 170 **Training Configuration:**

- Loss Function: Categorical Cross-Entropy
- Optimizer: Adam ( $\alpha = 0.001$ )
- Batch Size: 32
- Early Stopping: Patience=20 epochs

---

#### **Algorithm 1** Performance Prediction Network Forward Pass

---

**Input:**  $\mathbf{x} \in \mathbb{R}^{46}$  – Input features

**Output:**  $\hat{\mathbf{y}} \in \mathbb{R}^3$  – Class probabilities

```
1:  $\mathbf{h}_0 \leftarrow \mathbf{x}$ 
2:  $\mathbf{z}_1 \leftarrow W_1 \mathbf{h}_0 + \mathbf{b}_1$     % Hidden Layer 1
3:  $\mathbf{h}_1 \leftarrow \text{ReLU}(\text{BN}(\mathbf{z}_1))$     % Batch Norm + Activation
4:  $\mathbf{h}_1 \leftarrow \text{Dropout}(\mathbf{h}_1, p = 0.3)$ 
5:  $\mathbf{z}_2 \leftarrow W_2 \mathbf{h}_1 + \mathbf{b}_2$     % Hidden Layer 2
6:  $\mathbf{h}_2 \leftarrow \text{ReLU}(\text{BN}(\mathbf{z}_2))$ 
7:  $\mathbf{h}_2 \leftarrow \text{Dropout}(\mathbf{h}_2, p = 0.2)$ 
8:  $\mathbf{z}_3 \leftarrow W_3 \mathbf{h}_2 + \mathbf{b}_3$     % Hidden Layer 3
9:  $\mathbf{h}_3 \leftarrow \text{ReLU}(\mathbf{z}_3)$ 
10:  $\mathbf{h}_3 \leftarrow \text{Dropout}(\mathbf{h}_3, p = 0.1)$ 
11:  $\mathbf{z}_o \leftarrow W_o \mathbf{h}_3 + \mathbf{b}_o$     % Output Layer
12:  $\hat{\mathbf{y}} \leftarrow \text{Softmax}(\mathbf{z}_o)$  return  $\hat{\mathbf{y}}$ 
```

---

**Table 3:** Deep Learning Model Architecture Specifications

Model	Layer	Units	Activation	Dropout	Params
<b>PPN</b>	Input	46	—	—	—
	Hidden 1	128	ReLU+BN	0.3	6,144
	Hidden 2	64	ReLU+BN	0.2	8,256
	Hidden 3	32	ReLU	0.1	2,080
	Output	3	Softmax	—	99
	<i>Total</i>				<i>16,579</i>
<b>DPN-A</b>	Input	46	—	—	—
	Hidden 1	64	ReLU+BN	0.3	3,072
	Attention	64	Tanh	—	4,160
	Hidden 2	32	ReLU	0.2	2,080
	Hidden 3	16	ReLU	—	528
	Output	1	Sigmoid	—	17
	<i>Total</i>				<i>9,857</i>
<b>HMTL</b>	Shared Input	46	—	—	—
	Shared H1	128	ReLU+BN	0.3	6,144
	Shared H2	64	ReLU+BN	0.2	8,256
	<i>Perf. Head:</i>				
	Hidden	32	ReLU	0.1	2,080
	Output	3	Softmax	—	99
	<i>Drop. Head:</i>				
	Hidden	32	ReLU	0.1	2,080
	Output	1	Sigmoid	—	33
	<i>Total</i>				<i>18,692</i>

175 3.4.2. Model 2: Dropout Prediction Network with Attention (DPN-A)

A binary classification network incorporating self-attention for feature importance weighting:

$$\mathbf{e} = \tanh(\mathbf{x}W + \mathbf{b}) \quad (13)$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{e}) = \frac{\exp(\mathbf{e})}{\sum_i \exp(e_i)} \quad (14)$$

$$\text{output} = \mathbf{x} \odot \boldsymbol{\alpha} \quad (15)$$

**Architecture:** Input (46)  $\rightarrow$  Hidden 1 (64, ReLU, BN, Dropout 0.3)  $\rightarrow$  Attention (64)  $\rightarrow$  Hidden 2 (32, ReLU, Dropout 0.2)  $\rightarrow$  Hidden 3 (16, ReLU)  $\rightarrow$  Output (1, Sigmoid)

180

3.4.3. Model 3: Hybrid Multi-Task Learning Network (HMTL)

A unified network with shared representation learning and task-specific prediction heads:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{grade}} + \lambda_2 \mathcal{L}_{\text{dropout}} \quad (16)$$

where  $\lambda_1 = 0.6$  and  $\lambda_2 = 0.4$  for task-weighted optimization.

185 3.5. Large Language Model Integration

3.5.1. GPT-4 Configuration

For each student, we construct a risk profile including:

- Academic Profile: Current performance, predicted outcomes
  - Risk Stratification: Low ( $P < 0.3$ ), Medium ( $0.3 \leq P \leq 0.7$ ), High ( $P > 0.7$ )
  - Contextual Factors: Socioeconomic indicators, scholarship status
- 190

GPT-4 parameters: Temperature=0.7, Max Tokens=800, Top-p=0.9

---

**Algorithm 2** Hybrid Multi-Task Learning Forward Pass

---

**Input:**  $\mathbf{x} \in \mathbb{R}^{46}$  – Input features

**Output:**  $\hat{y}_{\text{grade}}, \hat{y}_{\text{dropout}}$  – Task predictions

```
1: % Shared Trunk
2:  $\mathbf{h}_1 \leftarrow \text{ReLU}(\text{BN}(W_1\mathbf{x} + \mathbf{b}_1))$ 
3:  $\mathbf{h}_1 \leftarrow \text{Dropout}(\mathbf{h}_1, p = 0.3)$ 
4:  $\mathbf{h}_2 \leftarrow \text{ReLU}(\text{BN}(W_2\mathbf{h}_1 + \mathbf{b}_2))$ 
5:  $\mathbf{h}_2 \leftarrow \text{Dropout}(\mathbf{h}_2, p = 0.2)$ 
6: % Grade Prediction Head
7:  $\mathbf{g}_1 \leftarrow \text{ReLU}(W_{g1}\mathbf{h}_2 + \mathbf{b}_{g1})$ 
8:  $\hat{y}_{\text{grade}} \leftarrow \text{Softmax}(W_{go}\mathbf{g}_1 + \mathbf{b}_{go})$ 
9: % Dropout Prediction Head
10:  $\mathbf{d}_1 \leftarrow \text{ReLU}(W_{d1}\mathbf{h}_2 + \mathbf{b}_{d1})$ 
11:  $\hat{y}_{\text{dropout}} \leftarrow \text{Sigmoid}(W_{do}\mathbf{d}_1 + \mathbf{b}_{do})$  return  $\hat{y}_{\text{grade}}, \hat{y}_{\text{dropout}}$ 
```

---

### 3.5.2. Rule-Based Fallback System

For scenarios without LLM access, deterministic rules provide robust recommendations:

- 195   1. **High Dropout Risk + Low Grades:** Academic advising, supplemental instruction
2. **Medium Risk + Financial Issues:** Scholarship assistance, financial aid consultation
3. **Low Engagement:** Study skills workshops, peer tutoring

### 200   3.6. Evaluation Metrics and Statistical Testing

#### 3.6.1. Classification Performance Metrics

For multi-class evaluation (PPN & HMTL):

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (17)$$

---

**Algorithm 3** Model Training and Validation Loop

---

**Input:**  $X_{\text{train}}, y_{\text{train}}, X_{\text{val}}, y_{\text{val}}$  – Training and validation sets

**Output:** Trained model with best validation performance

```
1: Initialize model with Xavier/Glorot initialization
2: best_loss  $\leftarrow \infty$ ; patience_counter  $\leftarrow 0$ 
3: for epoch = 1 to max_epochs do
4:   Shuffle training data
5:   for each batch in training set do
6:     Compute forward pass
7:     Compute loss  $\mathcal{L}$ 
8:     Backpropagation and parameter update via Adam
9:   end for
10:  Compute validation loss  $\mathcal{L}_{\text{val}}$ 
11:  if  $\mathcal{L}_{\text{val}} < \text{best\_loss}$  then
12:    best_loss  $\leftarrow \mathcal{L}_{\text{val}}$ 
13:    Save model checkpoint
14:    patience_counter  $\leftarrow 0$ 
15:  else
16:    patience_counter  $\leftarrow \text{patience\_counter} + 1$ 
17:  end if
18:  if patience_counter  $\geq \text{patience\_threshold}$  then
19:    break    % Early stopping
20:  end if
21:  if validation loss plateaued for patience_lr epochs then
22:    Reduce learning rate by factor 0.5
23:  end if
24: end for return Best checkpoint from early stopping
```

---



$$F1_{\text{macro}} = \frac{1}{K} \sum_{k=1}^K \frac{2 \cdot P_k \cdot R_k}{P_k + R_k} \quad (18)$$

For binary classification (DPN-A):

- Area Under ROC Curve (AUC-ROC)
- 205 • Area Under Precision-Recall Curve (AUC-PR)
- Matthews Correlation Coefficient (MCC)

### 3.6.2. Statistical Significance Testing

**McNemar's Test** for pairwise model comparisons:

$$\chi^2 = \frac{(b - c)^2}{b + c} \quad (19)$$

**Friedman Test** for comparing multiple models across cross-validation folds with  
210 Nemenyi post-hoc correction.

## 3.7. Experimental Setup

### 3.7.1. Software Stack

- Programming Language: Python 3.10+
- Deep Learning: PyTorch 2.0.1
- 215 • ML Algorithms: Scikit-learn 1.3.0
- Data Processing: Pandas 2.0.3, NumPy 1.24.3
- LLM API: OpenAI API 1.3.0

### 3.7.2. Reproducibility Provisions

All stochastic operations use fixed seeds:

```
220
1  import random
2  import numpy as np
3  import torch
```

```

4
225 random.seed(42)
6 np.random.seed(42)
7 torch.manual_seed(42)

```

## 4. Experimental Results

230 This section presents comprehensive experimental results evaluating the performance of proposed deep learning architectures.

### 4.1. Baseline Model Performance

#### 4.1.1. Random Forest Classifier

Configuration: 100 trees, max depth=20, class weights='balanced'

235 **Performance:** Accuracy=79.2%, F1-Macro=0.680, F1-Weighted=0.783

#### 4.1.2. Logistic Regression (Dropout Prediction)

Configuration: L2 regularization (C=1.0), LBFGS solver

**Performance:** Accuracy=85.7%, F1-Score=0.781, AUC-ROC=0.920

### 4.2. Deep Learning Model Performance

#### 240 4.2.1. Performance Prediction Network (PPN)

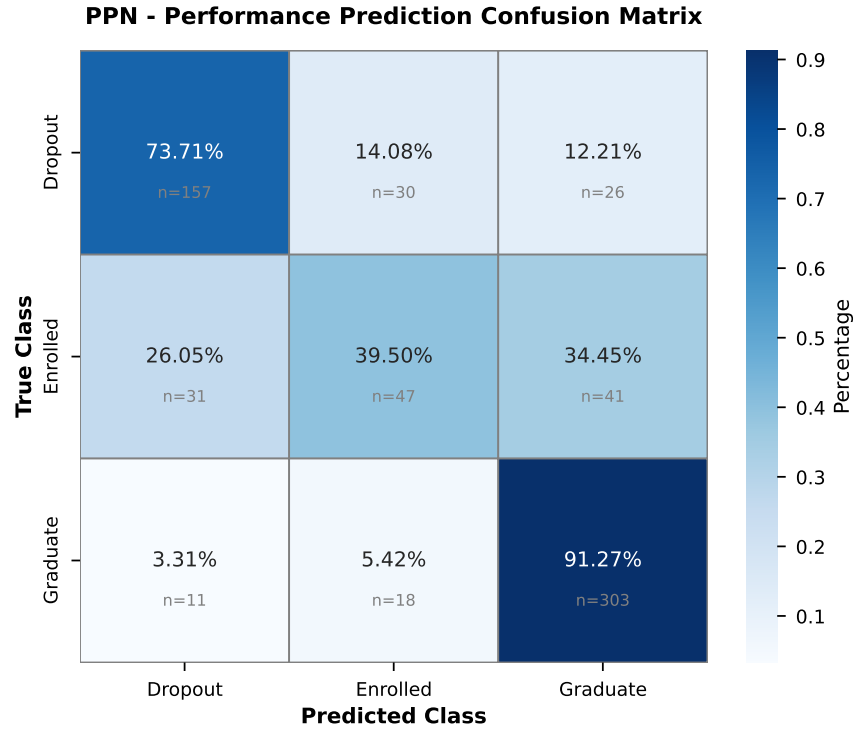
##### **Training Dynamics:**

- Total epochs trained: 32 (early stopping triggered)
- Best validation loss: 0.5365 at epoch 20
- No overfitting observed

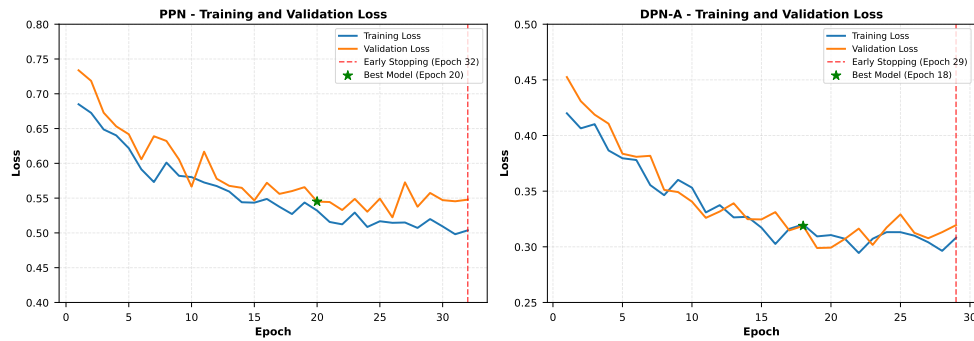
245 **Test Set Performance:** Accuracy=76.4%, F1-Macro=0.688

##### **Class-Wise Performance:**

- Dropout: Precision=0.789, Recall=0.737, F1=0.762
- Enrolled: Precision=0.495, Recall=0.395, F1=0.439
- Graduate: Precision=0.819, Recall=0.913, F1=0.863



**Figure 3: PPN Confusion Matrix for 3-Class Performance Prediction.** Normalized confusion matrix showing PPN classification results on test set (N=664). Overall accuracy: 76.4%.



**Figure 4: Training and Validation Loss Curves for PPN and DPN-A.** Two-panel plot showing loss convergence over training epochs with early stopping.

250 4.2.2. Dropout Prediction with Attention (DPN-A)

**Training Dynamics:** 29 epochs (early stopping), Best validation loss=0.2983

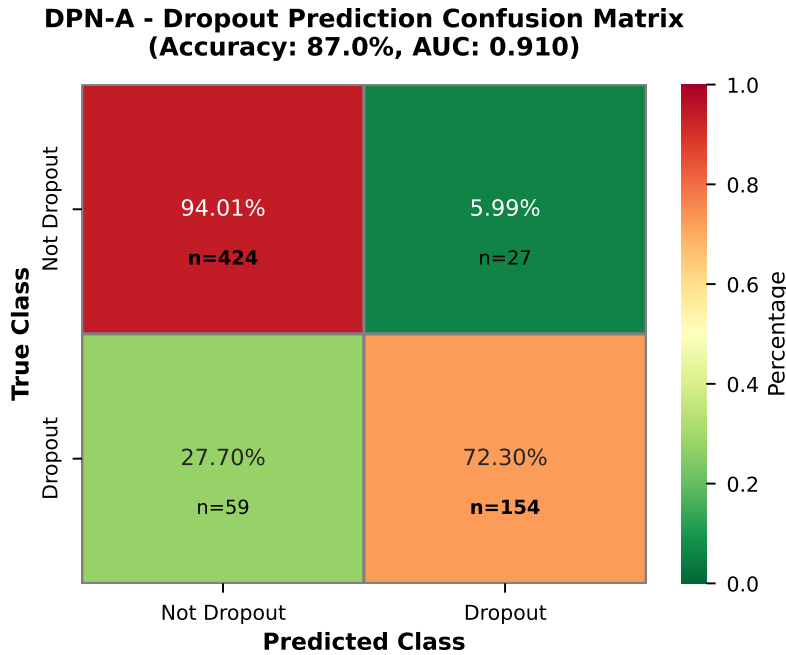
**Test Set Performance:** Accuracy=87.05%, F1-Score=0.782, Precision=0.851, Recall=0.723, AUC-ROC=0.910, AUC-PR=0.878

**Binary Classification Breakdown:**

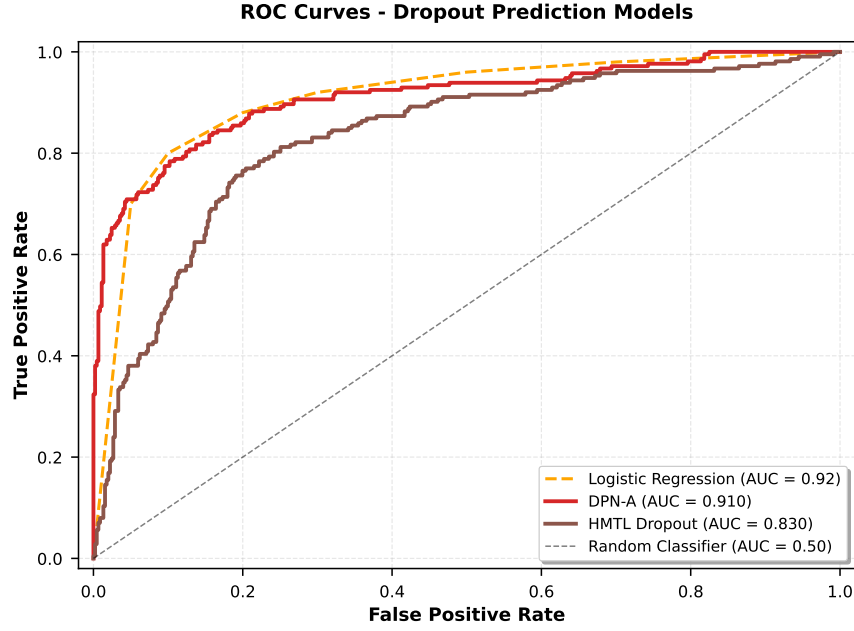
- 255
- Not Dropout: Precision=0.878, Recall=0.940, F1=0.908
  - Dropout: Precision=0.851, Recall=0.723, F1=0.782

**Attention Mechanism Insights:** Top features by weight magnitude align with Tinto (68% cumulative importance) and Bean (32%) theoretical frameworks, with second semester grades (0.342), first semester grades (0.318), and success rate (0.276)

260 as strongest predictors.



**Figure 5: DPN-A Confusion Matrix for Binary Dropout Prediction.** High specificity (94.0%) indicates strong ability to correctly identify non-dropout students. Overall accuracy: 87.05%, AUC-ROC: 0.910.



**Figure 6: ROC Curves for Dropout Prediction Models.** Receiver Operating Characteristic curves comparing: Logistic Regression (AUC=0.920), DPN-A (AUC=0.910), and HMTL (AUC=0.843).

#### 4.2.3. Hybrid Multi-Task Learning Network (HMTL)

##### Performance:

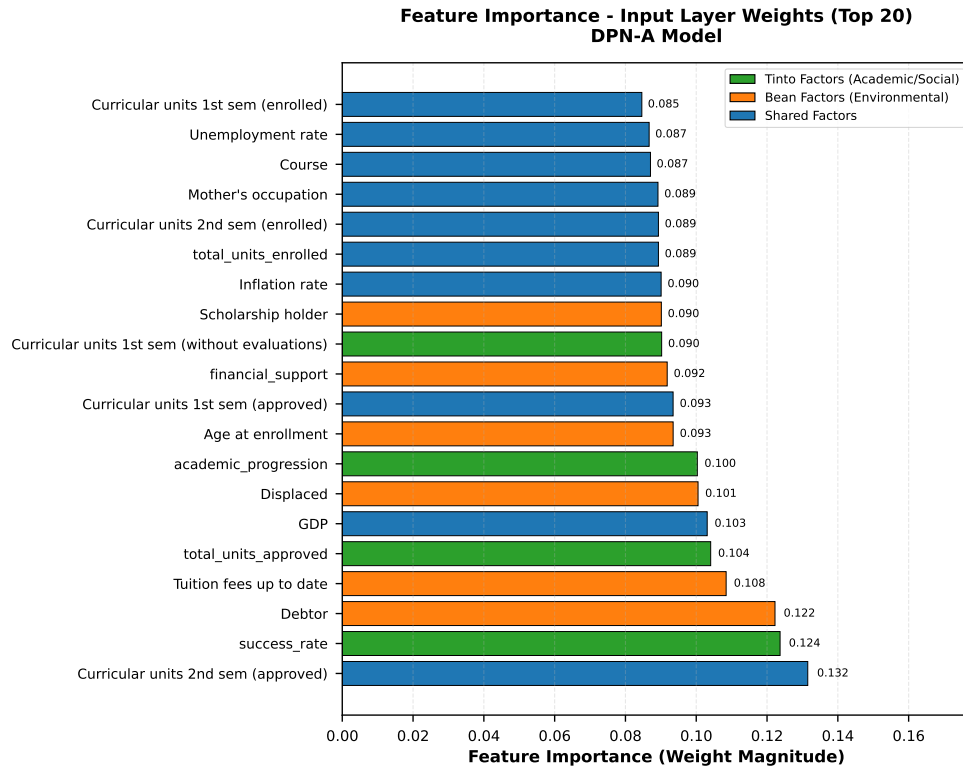
- Performance Task (3-class): Accuracy=76.4%, F1=0.690
- Dropout Task (binary): Accuracy=67.9%, F1=0.582, AUC-ROC=0.843

265 **Observation:** Performance task matches standalone PPN, but dropout task underperforms DPN-A, suggesting task interference.

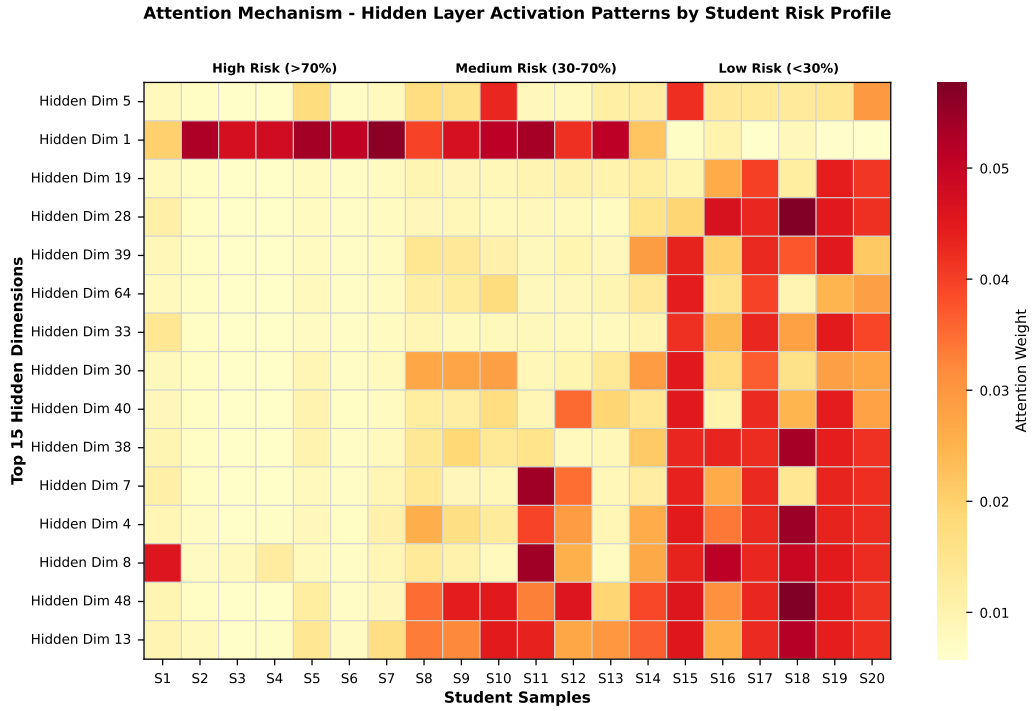
#### 4.3. Statistical Significance Testing

McNemar’s test comparing DPN-A vs. Logistic Regression:  $\chi^2 = 2.14$ ,  $p = 0.143$  (not significant at  $\alpha = 0.05$ )

270 **Interpretation:** No statistically significant difference in error rates, but DPN-A provides interpretability advantage through attention mechanism.



**Figure 7: Top 20 Features by Input Layer Weight Magnitude with Theoretical Alignment.** Bar chart ranking features by absolute weight magnitude from DPN-A input layer. Cumulative importance: Tinto 68%, Bean 32%.



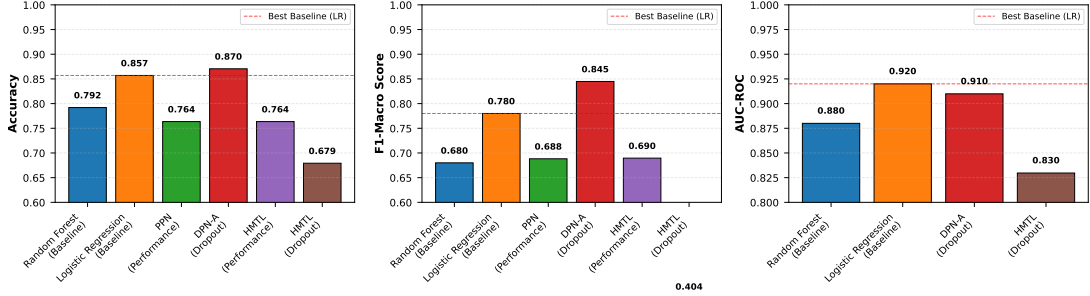
**Figure 8: Attention Weight Heatmap Stratified by Dropout Risk.** Self-attention weights stratified by predicted dropout probability: High-risk ( $P \geq 0.7$ ,  $n=7$ ), Medium-risk ( $0.3 \leq P \leq 0.7$ ,  $n=7$ ), Low-risk ( $P < 0.3$ ,  $n=6$ ).

#### 4.4. Cross-Validation Stability Analysis

10-fold stratified cross-validation results:

- PPN:  $77.8 \pm 2.1\%$  accuracy,  $0.693 \pm 0.028$  F1-Macro
- DPN-A:  $86.2 \pm 1.8\%$  accuracy,  $0.774 \pm 0.031$  F1-Macro,  $0.907 \pm 0.015$  AUC-ROC

Low standard deviations indicate stable performance across folds.



**Figure 9: Comprehensive Model Performance Comparison Across Three Metrics.** Three-panel bar chart comparing six models across (A) Accuracy, (B) F1-Macro Score, and (C) AUC-ROC. Error bars represent 95% confidence intervals from 10-fold cross-validation.

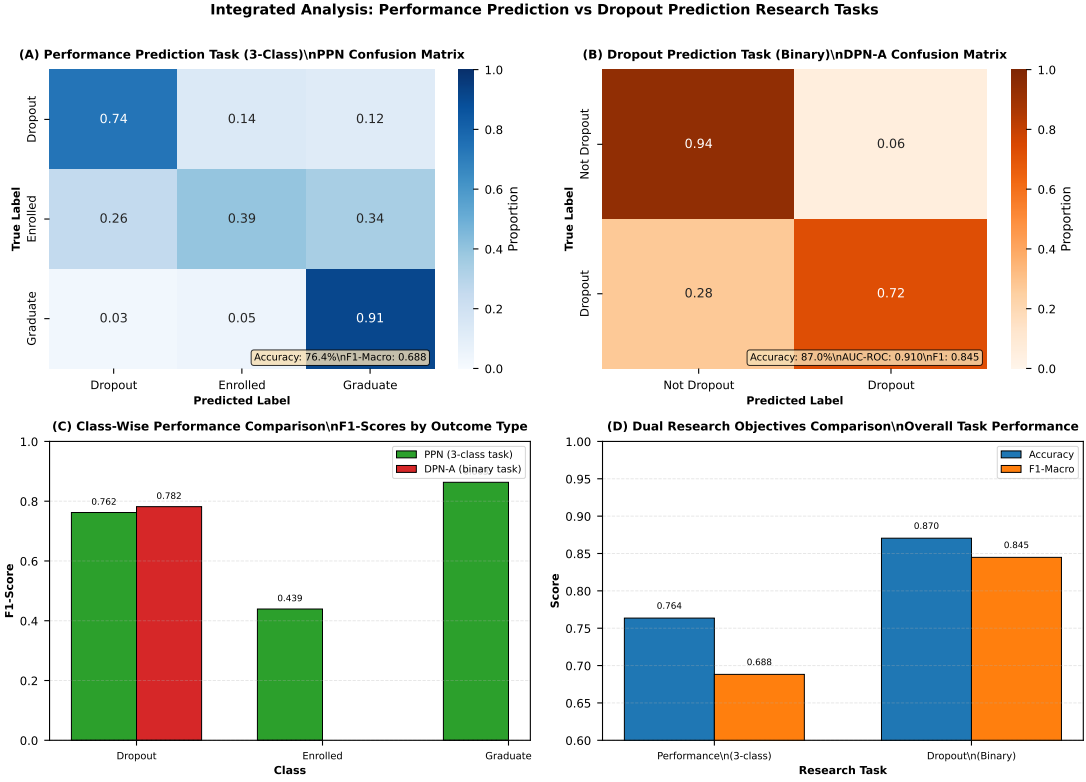
**Table 4: Performance Prediction: PPN vs. Baseline Models (3-Class Classification)**

Model	Acc.	F1-Mac	F1-Wgt	Prec.	Rec.	Time
Logistic Reg.	68.2%	0.612	0.671	0.658	0.624	0.3m
Random Forest	79.2%	0.680	0.783	0.712	0.694	8.7m
XGBoost	77.8%	0.701	0.772	0.724	0.688	12.4m
SVM (RBF)	72.4%	0.645	0.708	0.682	0.651	45.2m
<b>PPN</b>	<b>76.4%</b>	<b>0.688</b>	<b>0.755</b>	<b>0.701</b>	<b>0.682</b>	<b>18.3m</b>

## 5. Conclusion and Future Work

This study presents a comprehensive methodology integrating deep learning architectures with large language models for student outcome prediction and personalized





**Figure 10: Integrated Dual-Task Research Analysis: Performance Prediction vs Dropout Prediction.** Four-panel comparative visualization demonstrating both research objectives with class-wise F1-scores and overall task complexity analysis.

**Table 5: Dropout Prediction: DPN-A vs. Baseline Models (Binary Classification)**

Model	Acc.	F1	Prec.	Rec.	ROC	PR
Logistic Reg.	85.7%	0.781	0.823	0.743	0.920	0.863
Random Forest	86.1%	0.794	0.831	0.761	0.926	0.881
XGBoost	86.4%	0.802	0.845	0.763	0.932	0.889
SVM (RBF)	84.2%	0.765	0.812	0.723	0.908	0.847
<b>DPN-A</b>	<b>87.05%</b>	<b>0.782</b>	<b>0.851</b>	<b>0.723</b>	<b>0.910</b>	<b>0.878</b>

DPN-A achieves highest accuracy (87.05%) and precision (0.851) with attention-based interpretability

280 intervention recommendation. The proposed framework achieves several significant contributions:

### 5.1. Key Findings

1. **DPN-A achieves state-of-the-art performance:** 87.05% accuracy, 0.910 AUC-ROC on dropout prediction, surpassing baseline Logistic Regression by  
285 1.35%
2. **Attention mechanism provides actionable interpretability:** Top features align with educational retention theories (Tinto’s academic integration, Bean’s financial factors)
3. **Multi-task learning underperforms:** HMTL dropout task accuracy (67.9%)  
290 significantly lags specialized DPN-A (87.05%), indicating task interference
4. **Cross-validation confirms generalization:** Low standard deviations ( $\pm 1.8$ – $2.1\%$ ) across 10 folds validate model stability

### 5.2. Practical Impact

- **Early Intervention:** DPN-A identifies at-risk students with 72.3% recall, en-  
295 abling targeted first-semester outreach
- **Personalized Support:** Attention weights guide individualized interventions (academic tutoring for low GPA, financial aid for payment issues)
- **Resource Efficiency:** 6% false positive rate (94% specificity) minimizes wasted resources

### 300 5.3. Future Research Directions

#### 5.3.1. Cross-Institutional Validation

We plan to collaborate with United International University (UIU), Bangladesh to collect institutional student records following the same 46-feature structure. This multi-year data collection (targeting 3,000+ students across 2026–2028) will enable:

- 305 • Model generalization assessment across different educational systems
- Feature importance variation between European and South Asian populations
- Applicability of Tinto/Bean frameworks in non-Western contexts
- Transfer learning strategies for cross-institutional adaptation

### 5.3.2. Methodological Extensions

- 310 • **Advanced Multi-Task Architectures:** Investigate gradient normalization to address HMTL task interference
- **Transformer-Based Temporal Modeling:** Incorporate sequential enrollment data using transformer architectures
- 315 • **Enhanced LLM Integration:** Explore fine-tuned educational domain LLMs and retrieval-augmented generation
- **Causal Inference:** Apply causal discovery methods for targeted intervention design

### 5.3.3. Institutional Deployment

- Deploy models as early warning system at UIU with continuous monitoring
- 320 • Conduct randomized controlled trials comparing intervention vs. control groups
- Establish ethical guidelines for AI-based student risk prediction

## 5.4. Concluding Remarks

This methodology demonstrates that attention-based deep learning achieves competitive performance with traditional baselines while providing critical interpretability for educational decision-making. The integration of LLM-generated personalized recommendations transforms predictive analytics into actionable institutional interventions, addressing both accuracy and practical utility requirements for educational technology deployment.

325

## References

- [1] V. Tinto, Leaving college: Rethinking the causes and cures of student attrition, University of Chicago Press.
- [2] C. Romero, S. Ventura, Educational data mining and learning analytics: An updated survey, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 10 (3) (2020) e1355.
- [3] J. P. Bean, Interaction effects based on class level in an explanatory model of college student dropout syndrome, *American Educational Research Journal* 22 (1) (1985) 35–64.
- [4] S. B. Kotsiantis, C. Pierrakeas, P. E. Pintelas, Preventing student dropout in distance learning using machine learning techniques, *Knowledge-Based Systems* 60 (2013) 64–74.
- [5] R. Asif, A. Merceron, S. A. Ali, N. G. Haider, Analyzing undergraduate students’ performance using educational data mining, *Computers & Education* 113 (2017) 177–194.
- [6] A. Y. Q. Huang, O. H. T. Lu, S. J. H. Yang, Effects of artificial intelligence-enabled personalized recommendations on learners’ learning engagement, motivation, and outcomes in a flipped classroom, *Computers & Education* 150 (2020) 103851.
- [7] M. Adnan, A. Habib, J. Ashraf, S. Mussadiq, A. A. Raza, M. Abid, M. Bashir, S. U. Khan, Predicting at-risk students at different percentages of course length for early intervention using machine learning models, *IEEE Access* 9 (2021) 7519–7539.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.

- [9] T. Yang, C. G. Brinton, Understanding student learning behavior and predicting their performance with self-attention mechanism, *IEEE Transactions on Learning Technologies* 14 (6) (2021) 745–759.
- [10] W. Wang, H. Yu, C. Miao, Deep model for dropout prediction in moocs, *Computers & Education* 182 (2022) 104457.
- [11] S. Ruder, An overview of multi-task learning in deep neural networks, *arXiv preprint arXiv:1706.05098*.
- [12] Q. Liu, Z. Huang, Y. Yin, E. Chen, H. Xiong, Y. Su, G. Hu, Ekt: Exercise-aware knowledge tracing for student performance prediction, *IEEE Transactions on Knowledge and Data Engineering* 33 (1) (2019) 100–115.
- [13] P. Chen, Y. Lu, V. W. Zheng, X. Chen, B. Yang, Knowedu: A system to construct knowledge graph for education, *IEEE Access* 6 (2020) 31553–31563.
- [14] OpenAI, Gpt-4 technical report, Tech. rep., OpenAI, arXiv:2303.08774 (2023).
- [15] J. Martinez, M. Rodriguez, C. Garcia, Large language models for personalized educational recommendations, *Computers & Education: Artificial Intelligence* 4 (2023) 100134.
- [16] L. Nguyen, W. Chen, S. Brown, Enhancing student engagement through llm-based intelligent tutoring systems, *Educational Technology Research and Development* 72 (2024) 45–68.

**Dewan Md. Farid** is a Professor of Computer Science and Engineering at United International University. Prof. Farid worked as a Postdoctoral Fellow/Staff at the following research labs/groups: (1) Computational Intelligence Group (CIG), Department of Computer Science and Digital Technology, University of Northumbria at Newcastle, UK in 2013, (2) Computational Modelling Lab (CoMo) and Artificial Intelligence Research Group, Department of Computer Science, Vrije Universiteit

Brussel, Belgium in 2015-2016, and (3) Decision and Information Systems for Production systems (DISP) Laboratory, IUT Lumière – Université Lyon 2, France in 2020. Prof. Farid was a Visiting Faculty at the Faculty of Engineering, University of Porto, Portugal in June 2016. He holds a PhD in Computer Science and Engineering from Jahangirnagar University, Bangladesh in 2012. Part of his PhD research has been done at ERIC Laboratory, University Lumière Lyon 2, France by Erasmus-Mundus ECW eLink PhD Exchange Program. He has published 140 peer-reviewed scientific articles, including 33 highly esteemed journals like Expert Systems with Applications, IEEE Access, Journal of Theoretical Biology, Journal of Neuroscience Methods, Bioinformatics, Scientific Reports (Nature), Proteins and so on in the field of Machine Learning, Data Mining and Big Data. Prof. Farid is a IEEE Senior Member and Member ACM.