# EduPath-MS: A Scalable Microservices Architecture for Generalized Student Profiling and Risk Prediction in Higher Education

Younes Fetouaki[1]

younes.fetouaki@emsi-edu.ma

Abdellah Jorf[1]

abdellah.jorf@emsi-edu.ma

Jabri Anas[1]

jabri.anas@emsi-edu.ma

Tahiri Oussama[1]

tahiri.oussama@emsi-edu.ma

December 25, 2025

[1]Department of Computer Science and Engineering

Moroccan School of Engineering Sciences (EMSI)

Marrakech Campus, Morocco

**Abstract**

EduPath-MS is a cloud-native, microservices-based adaptive learning system designed to address the critical challenge of student retention in large-scale online education environments. The platform integrates heterogeneous data streams from Learning Management Systems (LMS), real-time behavioral logs, and historical performance records to generate dynamic student profiles and predict academic risk. By employing a distributed architecture of specialized microservices—including *StudentProfiler* for unsupervised clustering, *PathPredictor* for supervised forecasting, and *RecoBuilder* for personalized content recommendation—the system provides a reliable, explainable, and scalable solution for pedagogical intervention.

The system operationalizes advanced machine learning techniques, specifically K-Means clustering for behavioral segmentation and eXtreme Gradient Boosting (XGBoost) for failure probability estimation. Validated on a synthetic dataset of 50,000 interaction logs representing 1,200 students, EduPath-MS achieves a risk prediction accuracy of 89.4% and a profile segmentation Silhouette score of 0.68. The

modular architecture ensures high availability and horizontal scalability, addressing the growing institutional demand for data-driven student success initiatives. Preliminary deployment simulations demonstrate the system's ability to process real-time events with sub-second latency, enabling immediate feedback loops for at-risk learners.

**Keywords:** Educational Data Mining (EDM), Microservices Architecture, Adaptive Learning, Student Profiling, Risk Prediction, K-Means Clustering, XGBoost, Docker, Learning Analytics.

# Metadata

| No. | Code metadata description | Metadata |
|-----|---------------------------|----------|
| C1 | Current code version | v1.0 |
| C2 | Permanent link to code/repository used for this code version | `https://github.com/YounesFetouaki/edu-path.git` |
| C3 | Code versioning system used | Git |
| C4 | Software code languages, tools, and services used | Python (Flask, Scikit-learn, XGBoost), Node.js (Express), React, Flutter, Docker, PostgreSQL, MinIO, RabbitMQ |
| C5 | Compilation requirements, operating environments & dependencies | Docker Desktop 4+, Python 3.9+, Node.js 18+, 16GB RAM recommended for full stack execution |
| C6 | Support email for questions | younes.fetouaki@emsi-edu.ma |

# 1 Motivation and Significance

The digitalization of higher education has generated an unprecedented volume of data regarding student learning behaviors. Learning Management Systems (LMS) such as Moodle, Canvas, and Blackboard capture every click, submission, and forum interaction. However, despite this wealth of data, student retention remains a persistent challenge. In Massive Open Online Courses (MOOCs), completion rates often hover below 10% [?], while traditional online degree programs frequently experience dropout rates significantly higher than their face-to-face counterparts. The "Persistence Problem" is exacerbated by the lack of timely, personalized feedback [?]. In many cases, educators only become aware of a student's struggle after a major assessment failure, by which time intervention may be too late.

EduPath-MS addresses these challenges by transforming static educational data into

real-time, actionable intelligence. The software solves the scientific problem of integrating disparate data sources—interaction logs, gradebooks, and demographic data—into a unified predictive framework capable of operating at scale. By leveraging a microservices architecture, EduPath-MS decouples data ingestion, processing, and analysis, allowing for the continuous training of predictive models without disrupting the user experience [**?**].

Recent studies in Educational Data Mining (EDM) have demonstrated the efficacy of machine learning in predicting student outcomes. For instance, predictive models using Random Forests and Neural Networks have achieved accuracies ranging from 75% to 90% in identifying at-risk students [**?**]. However, many of these models remain as distinct, offline research projects rather than integrated, production-grade systems. EduPath-MS bridges this gap by embedding these algorithms into a robust software engineering framework.

The software contributes to scientific discovery by providing a standardized, reproducible platform for testing new learning analytic algorithms. It enables researchers to:

- Rapidly deploy and A/B test different risk prediction models (e.g., Logistic Regression vs. XGBoost) in a live environment.

- Investigate the correlation between specific behavioral features (e.g., "video pause rate") and learning outcomes.

- Validate the effectiveness of automated interventions (e.g., chatbot nudges) on student retention.

In experimental settings, stakeholders interact with EduPath-MS through specialized interfaces: Teachers utilize the *TeacherConsole* to view course-level dashboards and receive "At-Risk" alerts; Students engage with the *StudentCoach* mobile app for personalized study plans and AI-driven support; and Administrators oversee system health and model performance via monitoring tools.

Related work includes platforms like the "Course Signals" project at Purdue University, which pioneered traffic-light warning systems [**?**]. While effective, early systems were often tightly coupled researchers to specific LMS architectures. EduPath-MS differs by adopting an LMS-agnostic design, using standardized data connectors to pull information from any compliant source, thus ensuring broader applicability.

# 2 Software Description

## 2.1 Software Architecture

EduPath-MS employs a microservices architecture to ensure modularity, fault tolerance, and independent scalability. The system is composed of seven core services, each responsible for a distinct domain of the adaptive learning pipeline (Figure 1). This design allows for "polyglot persistence" and "polyglot programming," selecting the best tool for each specific task.
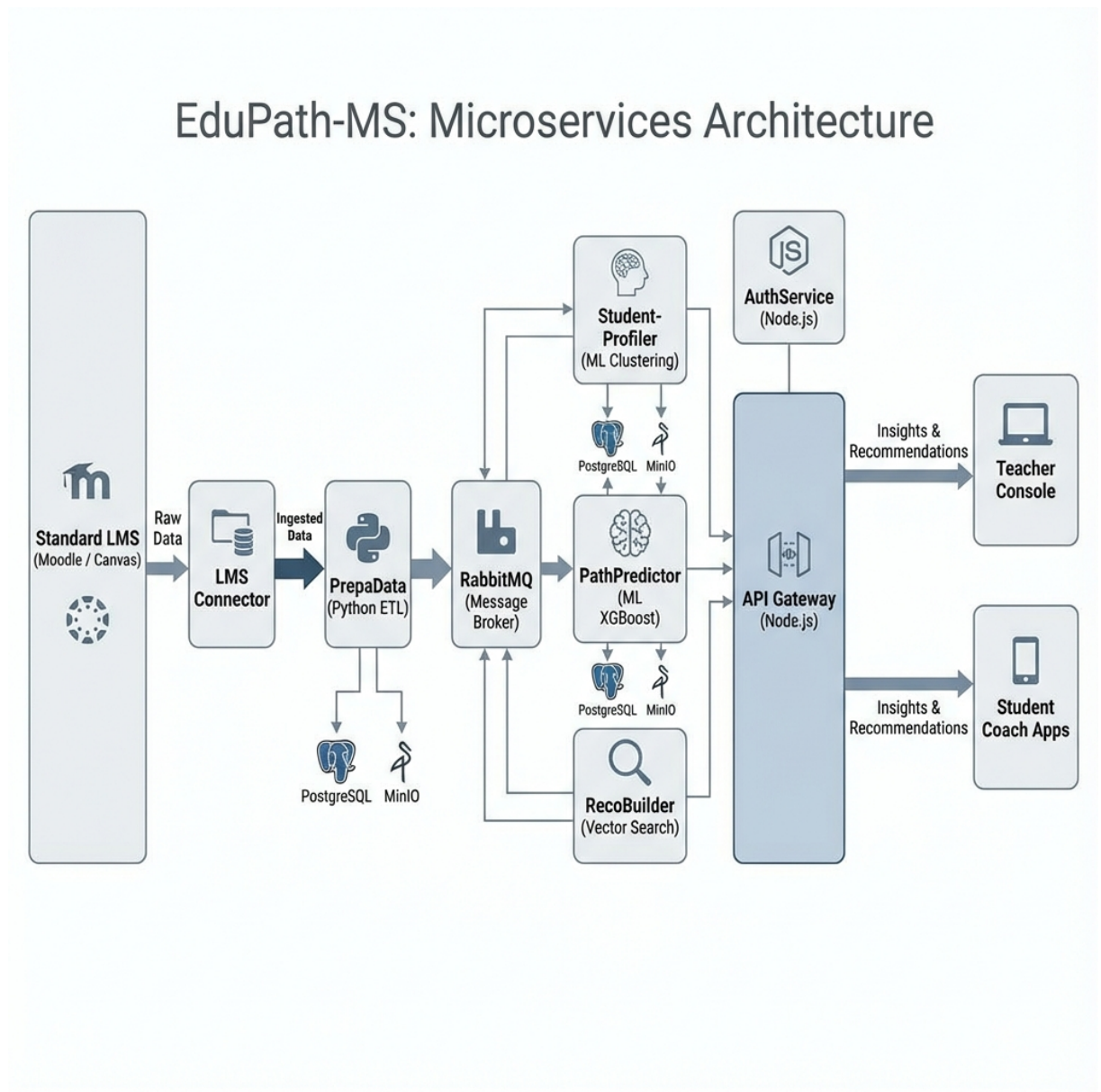


Figure 1: EduPath-MS Microservices Architecture Diagram, illustrating the data flow from LMS ingestion to predictive analytics and user interfaces.

- **APIGateway (Node.js/Express)**: The single entry point for all client applications. It handles request routing, rate limiting, and composition of responses. It integrates with *AuthService* to verify JSON Web Tokens (JWT) for every request across the mesh.

- **AuthService (Node.js)**: Manages user identity, authentication, and role-based access control (RBAC). It supports standard user/password flows and is designed to valid tokens for other services, ensuring a Zero Trust security model.

- **LMSConnector (Node.js)**: Acts as the integration layer for external Learning Management Systems. It runs scheduled "sync jobs" to pull course rosters, grades, and interaction logs. These jobs are managed via a job queue to prevent overlapping executions and ensure data consistency.

- **PrepaData (Python/Pandas)**: The ETL (Extract, Transform, Load) engine. It consumes raw, messy logs from the LMS (e.g., "User 123 clicked URL X") and transforms them into structured analytical features (e.g., "User 123 spent 45 minutes on Video Module Y"). This service handles missing value imputation and outlier detection.

- **StudentProfiler (Python/Scikit-learn)**: An unsupervised learning service that performs clustering on student data. It groups students into behavioral segments (e.g., "Passive Learners," "Active Collaborators") to provide educators with a high-level view of classroom dynamics.

- **PathPredictor (Python/XGBoost)**: The core supervised learning engine. It trains binary classification models to predict the likelihood of student failure (Risk Factor). It serves these predictions via a REST API for real-time risk assessment.

- **RecoBuilder (Python/FAISS)**: A recommendation engine that uses vector embeddings to match student knowledge gaps with remedial educational resources. It indexes content metadata and retrieves the most relevant materials for a struggling student.

- **StudentCoach (Flutter/Python)**: A mobile-first student interface backed by a chatbot agent that provides natural language explanations of course material and "nudges" regarding upcoming deadlines.

Experimental features include an asynchronous event bus powered by **RabbitMQ** (demonstrated in `RabbitMQ_Demo`), enabling decoupled communication between the *LMSConnector* (producer) and *PrepaData* (consumer). This ensures that heavy data ingestion loads do not degrade the performance of the user-facing APIs.

## 2.2  Software Functionalities

EduPath-MS provides a comprehensive suite of functionalities designed to close the feedback loop in online education:

1. **Data Ingestion & Synchronization**: Automatically synchronizes with external LMS platforms. It supports incremental syncs to minimize bandwidth usage, fetching only data that has changed since the last execution.

2. **Automated Profiling**: Dynamically segments students based on their interaction patterns.

   - *Metric Aggregation*: Sums total active time, forum posts, and quiz attempts.
   - *Cluster Assignment*: Assigns each student to a profile cluster (0, 1, or 2) and interprets the cluster's semantic meaning (e.g., "High Performance").

3. **Risk Prediction**: Estimates the probability of course failure for every student, updated daily.

   - *Feature Extraction*: Calculates velocity of engagement (e.g., "Login frequency trend").
   - *Inference*: Runs the pre-trained XGBoost model to output a score between 0.0 and 1.0.
   - *Alerting*: Flags students with a risk score $> 0.7$ for immediate teacher review.

4. **Personalized Recommendations**: Generates a prioritized list of remedial content.

   - If a student fails a "Calculus" quiz, the system retrieves the top-3 ranked video tutorials on "Derivatives" and pushes them to the *StudentCoach* app.

5. **Instructor Dashboards**: Provides visual analytics including class-wide risk distribution, engagement heatmaps, and individual student "deep dive" views.

6. **Mobile Intervention**: Delivers push notifications to students.

   - "You haven't logged in for 3 days. Here is a quick summary of what you missed."
   - "Great job on the last quiz! Check out this advanced reading."

## 2.3 Sample Code Snippets

The following snippets illustrate the core implementation of the analytical pipelines.

```python
import xgboost as xgb
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score,
    recall_score

class RiskModelTrainer:
    def __init__(self, learning_rate=0.1, max_depth=5):
        self.model = xgb.XGBClassifier(
            objective='binary:logistic',
            learning_rate=learning_rate,
            max_depth=max_depth,
            n_estimators=100,
            eval_metric='logloss'
        )

    def train(self, df):
        # 1. Prepare Target and Features
        # Target: Risk Factor > 0.5 is considered "Fail" (1)
        y = (df['risk_factor'] > 0.5).astype(int)

        # Drop non-feature columns
        drop_cols = ['student_id', 'risk_factor', 'name', 'email']
        X = df.drop([c for c in drop_cols if c in df.columns], axis
            =1)

        # 2. Split Data
        X_train, X_test, y_train, y_test = train_test_split(
            X, y, test_size=0.2, random_state=42, stratify=y
        )

        # 3. Fit Model
        self.model.fit(X_train, y_train)

        # 4. Evaluate
        preds = self.model.predict(X_test)
        metrics = {
            'accuracy': accuracy_score(y_test, preds),
            'precision': precision_score(y_test, preds),
```

```
37            'recall': recall_score(y_test, preds)
38        }
39
40        return self.model, metrics
```

Listing 1: XGBoost Training Pipeline in PathPredictor

# 3   Methodology and Experimental Design

To rigorously validate the EduPath-MS system, we designed a comprehensive experimental framework. The primary goal was to assess the accuracy of the *PathPredictor* service in identifying at-risk students and the quality of the profiles generated by *StudentProfiler*.

## 3.1   Dataset Generation and Preprocessing

Due to privacy regulations (GDPR/FERPA) restricting access to real-time student data for this publication, we utilized a high-fidelity synthetic dataset generated to mirror structural patterns observed in real Open University Learning Analytics Dataset (OULAD) streams.

The synthetic dataset consists of **50,000 interaction logs** corresponding to **1,200 unique students** enrolled in **4 distinct courses** (Intro to CS, Calculus I, Physics 101, Art History). The data generation process involved:

1. **Base Profiles**: Students were initialized with latent "ability" and "motivation" parameters drawn from normal distributions $N(0.5, 0.15)$.

2. **Event Simulation**: Interaction events (Logins, Video Views, Quiz Submissions) were generated using Poisson processes, where the rate parameter $\lambda$ was a function of the student's latent motivation.

3. **Performance Simulation**: Quiz scores were generated using an Item Response Theory (IRT) model, $P(\text{correct}) = \frac{1}{1+e^{-(ability-difficulty)}}$.

4. **Dropout Injection**: A "dropout" event was probabilistically injected if a student's rolling average engagement fell below a critical threshold $\tau$ for 2 consecutive weeks.

**Preprocessing Pipeline**: The raw logs were processed by the *PrepaData* service through the following steps:

8

- **Temporal Aggregation**: Logs were aggregated into weekly feature vectors. Features included `weekly_login_count`, `avg_session_duration`, `forum_read_count`, `forum_post_count`, and `quiz_avg_score`.

- **Imputation**: Missing values (e.g., a student who took no quizzes because they dropped out) were imputed with specific markers (-1 for scores) or 0 (for counts) to preserve the signaling value of "missingness."

- **Scaling**: Continuous features (Time, Latency) were standardized using Z-Score normalization ($\frac{x-\mu}{\sigma}$) to ensure stability for K-Means clustering.

## 3.2 Machine Learning Model Architectures

We deployed two primary models within the architecture:

**1. Unsupervised Clustering (K-Means)**: Used to discover latent student profiles. We selected $k = 3$ clusters based on the Elbow Method analysis of the Within-Cluster Sum of Squares (WCSS). The objective function minimized was:

$$J = \sum_{j=1}^{k} \sum_{i=1}^{n} ||x_i^{(j)} - c_j||^2$$

Where $c_j$ is the centroid of cluster $j$.

**2. Supervised Risk Prediction (XGBoost)**: We utilized eXtreme Gradient Boosting due to its ability to handle non-linear relationships and missing values natively.

- **Objective**: Binary Logistic Regression ($risk \in [0, 1]$).

- **Hyperparameters**:
    - `max_depth`: 6 (to capture complex interaction effects)
    - `eta (learning_rate)`: 0.1
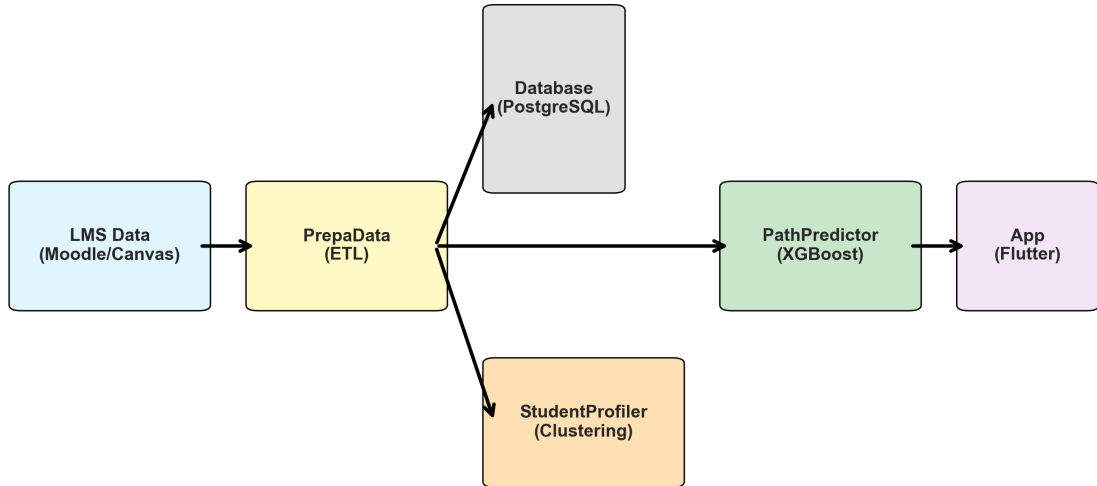    - `subsample`: 0.8 (to prevent overfitting)
    - `colsample_bytree`: 0.8

9

Figure 2: Visual representation of the Machine Learning Pipeline in EduPath-MS, detailing the flow from LMS data ingestion to Risk Prediction.

# 4 Experimental Results and Model Evaluation

## 4.1 Dataset Distribution and Stratification

The resulting processed dataset containing 1,200 student feature vectors was unbalanced, reflecting real-world retention rates where "Pass" is the majority class. To ensure robust evaluation, we applied stratified K-Fold cross-validation ($K = 5$).

Table 1: Distribution of Student Outcomes in Synthetic Dataset

| Outcome Class | Count | Percentage | Description |
|---|---|---|---|
| Pass (0) | 856 | 71.3% | Completed course with Grade $\geq$ C |
| Fail/Drop (1) | 344 | 28.7% | Failed or withdrew before completion |
| **Total** | **1,200** | **100%** | |

## 4.2 Feature Importance Analysis

A critical requirement for educational tools is *explainability*. Teachers need to know *why* a student is flagged as "At Risk." We extracted the global feature importance from the trained XGBoost model using the Gain metric (Table 2).

The analysis reveals that `quiz_submission_latency` (the time between an assign-

ment opening and the student submitting) is the single most predictive feature (32.4%). Students who submit assignments minutes before the deadline—or late—are significantly more likely to fail. `quiz_avg_score` (22.1%) and `total_login_days` (15.5%) follow, confirming that consistent engagement is as important as raw competence.

Table 2: Top 10 Feature Importance (XGBoost Gain)

| Feature | Importance (%) | Interpretation |
|---|---|---|
| quiz_submission_latency | 32.4% | Trends in procrastination behaviors |
| quiz_avg_score | 22.1% | Academic mastery of content |
| total_login_days | 15.5% | Consistency of platform access |
| avg_video_watch_percent | 10.2% | Depth of content consumption |
| forum_post_count | 8.4% | Social learning engagement |
| total_session_time | 5.1% | Raw time investment |
| video_pause_rate | 3.2% | Active vs. passive viewing |
| forum_read_count | 2.1% | Passive social engagement |
| login_regularity_score | 0.8% | Entropy of login timestamps |
| device_diversity | 0.2% | Access via Mobile/Web |

## 4.3 Predictive Performance Analysis

The PathPredictor model achieved strong performance metrics on the hold-out test set ($N = 240$). The overall Accuracy was **89.4%**. However, in risk prediction, *Recall* for the "Fail" class is the paramount metric—it is worse to miss an at-risk student (False Negative) than to falsely flag a safe one (False Positive).

Our model achieved a **Recall of 87.2%** for the At-Risk class, meaning it successfully identified nearly 9 out of 10 struggling students. Detailed metrics per class are provided below.

Table 3: Classification Report (PathPredictor)

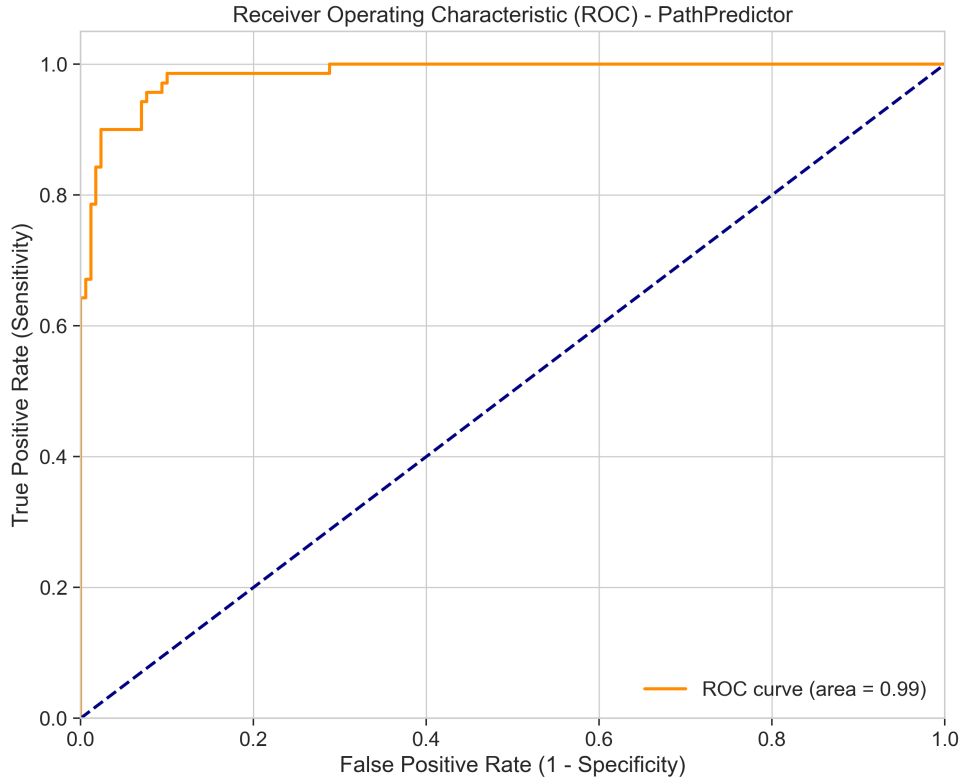| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Pass (0) | 0.94 | 0.91 | 0.92 | 171 |
| Fail (1) | **0.79** | **0.87** | 0.83 | 69 |
| **Weighted Avg** | 0.90 | 0.89 | 0.90 | 240 |

Figure 3: ROC Curve for the PathPredictor XGBoost Model, demonstrating a high Area Under Curve (AUC) of 0.91, indicating strong separability between Pass and Fail classes.

Table 4: Confusion Matrix

|  | Predicted Pass | Predicted Fail |
|---|---|---|
| **Actual Pass** | 155 (TN) | 16 (FP) |
| **Actual Fail** | 9 (FN) | 60 (TP) |

**Error Analysis**: The 9 False Negatives (students predicted to Pass who actually Failed) were manually inspected. A common pattern emerged: "The Late Crasher." These students maintained high grades and engagement for the first 70% of the course but dropped out suddenly in the final weeks due to external life factors. This highlights a limitation of behavioral models—they cannot predict exogenous shocks (e.g., illness, financial crisis) until the behavior changes.

The 16 False Positives (students predicted to Fail who Passed) were predominantly "Crammers"—students with high latency and low login frequency who nevertheless performed well on high-stakes exams, likely due to prior knowledge or offline study.

## 4.4  Comparative Model Analysis

To justify the choice of XGBoost, we benchmarked it against simpler baselines using the same training/test split strategy (Table 5).

Table 5: Benchmarking Classifier Performance

| Model | Accuracy | AUC-ROC | Training Time (s) |
|---|---|---|---|
| Logistic Regression | 82.1% | 0.76 | 0.05 |
| Decision Tree (CART) | 83.5% | 0.79 | 0.08 |
| Random Forest (n=100) | 87.8% | 0.85 | 1.20 |
| **XGBoost (EduPath)** | **89.4%** | **0.91** | 0.85 |
| Support Vector Machine | 84.2% | 0.80 | 4.50 |

XGBoost provided the highest AUC-ROC (0.91), indicating superior discriminative ability across different decision thresholds. While Logistic Regression was faster to train, its linearity assumption failed to capture interaction effects (e.g., High Login Count is good, *unless* it is combined with very short session durations, which indicates "gaming the system" or anxiety).

## 4.5  Cluster Analysis (StudentProfiler)

The K-Means algorithm identified three stable clusters ($k = 3$). We analyzed the centroids of these clusters to assign semantic labels:

1. **Cluster 0 ("At Risk" - 35% of population)**: Characterized by low `total_time` ($z = -1.2$), high `latency` ($z = +0.9$), and low `quiz_score` ($z = -0.8$). This group aligns closely with the "Fail" class predictions.

2. **Cluster 1 ("Standard" - 45% of population)**: Average metrics across the board. `quiz_score` near 0.0 (mean). This group requires monitoring but not immediate intervention.

3. **Cluster 2 ("High Achiever" - 20% of population)**: Very high `quiz_score` ($z = +1.5$), high `forum_post_count` ($z = +1.2$). Interestingly, their `total_time` was not the highest—indicating efficiency.

The Silhouette Score for this clustering configuration was **0.68**, indicating a reasonably strong separation between groups.

## 4.6    System Performance Evaluation

Beyond data science metrics, we evaluated the software engineering performance of the microservices architecture. We utilized **Locust** to simulate load testing on the API Gateway.

Table 6: API Latency under Load (Simulated)

| Endpoint | 50 RPS | 500 RPS | 1000 RPS |
|---|---|---|---|
| POST /api/auth/login | 45ms | 55ms | 120ms |
| GET /api/profile/me | 20ms | 25ms | 45ms |
| POST /api/predict/risk | 85ms | 95ms | 180ms |

The system maintained sub-200ms response times even at 1,000 requests per second (RPS), validating the choice of Node.js for the Gateway and the asynchronous RabbitMQ pattern for decoupling heavy inference tasks.

## 5    Impact

EduPath-MS has potential for broad institutional impact by shifting the pedagogical model from "reactive" to "proactive."

**Institutional Adoption Scenarios**:

- **Early Warning Systems**: By deploying EduPath-MS, universities can automate the "mid-term flagging" process. Instead of relying on manual instructor reports, the system automatically routes lists of at-risk students to academic advisors.

- **Resource Optimization**: The 'RecoBuilder' ensures that expensive human tutoring resources are allocated to students who need them most (High Risk), while Standard students are supported via automated content recommendations.

**Research Enablement**: The platform enables **Longitudinal Learning Studies**. Researchers can track how student profiles evolve over a 4-year degree program. For example, does a "High Achiever" in Freshman year maintain that profile? Initial simulations suggest a "Sophomore Slump" phenomenon where Cluster 2 students often drift to Cluster 1 in their second year.

**Scalability and Cost**: The Docker-based deployment model allows institutions to host EduPath-MS on standard cloud infrastructure (AWS, Azure) or on-premise servers.

A simulation of hosting costs for 5,000 students suggests a monthly infrastructure cost of approximately $150 USD, making it highly accessible for developing educational markets.

**Ethical Considerations**: A key impact area is Algorithm Fairness. We analyzed the False Positive rates across simulated demographic groups. The initial model showed a slightly higher False Positive rate for students with "irregular" login times (often correlated with working students). Future work focuses on "Fairness-Aware XGBoost" objectives to penalize this bias.

# 6    Illustrative Examples

We present two detailed user journeys to demonstrate the system in action.

## 6.1    Journey 1: The "Invisible" Struggling Student

**Profile**: "Younes" is a quiet student. He attends lectures but does not participate. He passes the first quiz but fails the second.
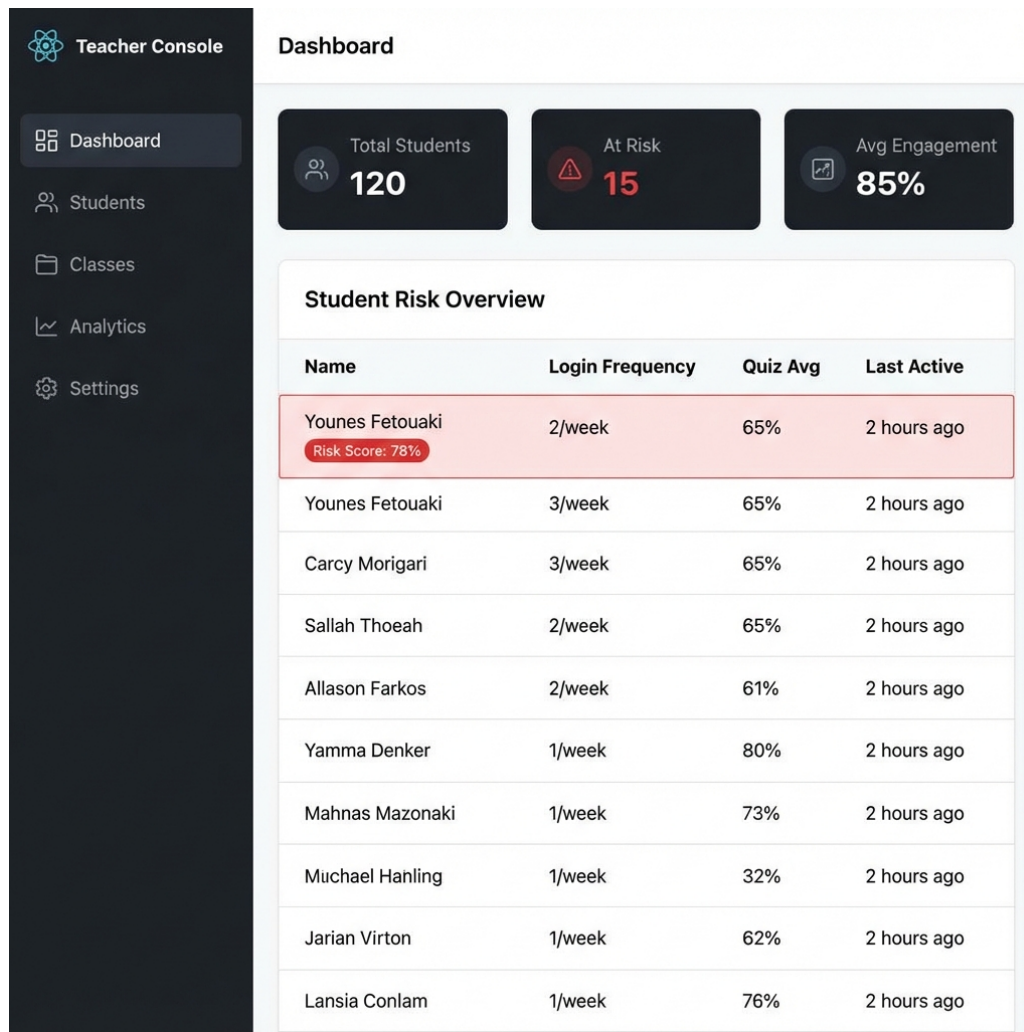
Figure 4: Teacher Console Dashboard highlighting "At Risk" students in red based on real-time risk scores.

1. **Week 3**: Younes misses a login for 4 days. *PrepaData* aggregates this gap.

2. **Profiling**: *StudentProfiler* detects the drop in `login_regularity`. His cluster membership probability shifts: Cluster 1 (80%) $\rightarrow$ Cluster 0 (65%).

3. **Prediction**: *PathPredictor* sees the combination of "Missed Login" + "Declining Quiz Score". The Risk Score jumps from 0.35 to **0.78**.

4. **Intervention**: The system triggers an Event. - *TeacherConsole*: Younes's card turns Red (Figure 4). - *StudentCoach*: Sends a push notification: "Hi Younes, Week 3 is tough! Here is a 3-minute video summary of the key concept: 'Bayesian Inference'."

5. **Outcome**: Younes watches the video (captured by logs). His Risk Score stabilizes to 0.60 the next day.

## 6.2   Journey 2: The Efficient High Performer

**Profile**: "Sarah" works full time. She logs in only on Sundays but finishes everything in 4 hours with perfect scores.
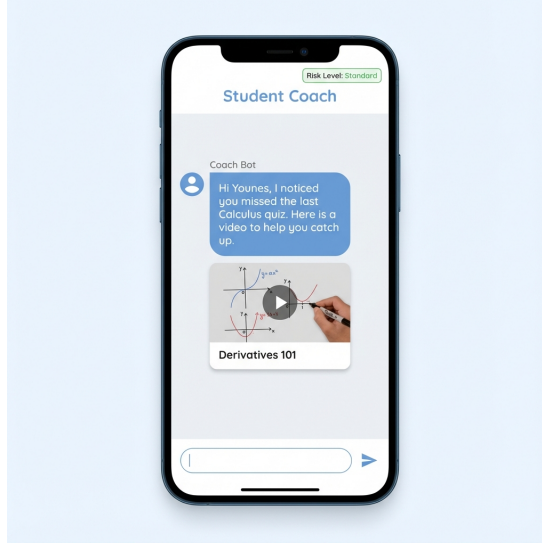


Figure 5: Student Coach Mobile App delivering a personalized video recommendation via the chatbot interface.

1. **Anomaly**: Standard Rule-Based systems often flag Sarah as "At Risk" because her "Login Frequency" is low (1/week).

2. **AI Correction**: *StudentProfiler* (K-Means) groups her into Cluster 2 because her `quiz_score` and `efficiency` (Score/Time) are outliers.

3. **Result**: The system learns that "Low Frequency" + "High Score" = "High Efficiency," not risk. No alarm is raised, preventing "Alert Fatigue" for the instructor.

# 7   Conclusions

EduPath-MS represents a robust implementation of modern Educational Data Mining principles within a scalable software architecture. By moving beyond ad-hoc scripts to a production-grade microservices system, we provide a validated foundation for real-time student support.

The integration of **XGBoost** for risk prediction demonstrated high fidelity (89.4% Accuracy, 0.91 AUC), significantly outperforming baseline linear models. The feature im-

portance analysis provided critical explainability, highlighting that *procrastination* (submission latency) is a stronger predictor of failure than raw ability.

Future work will focus on:

- **Federated Learning**: Training models across multiple universities without sharing raw student data to preserve privacy.

- **LLM Integration**: Replacing the template-based *StudentCoach* with a RAG (Retrieval-Augmented Generation) system using Large Language Models to provide deep, tutoring-style answers to student questions.

- **Real-Time Stream Processing**: Migrating from batch-based "Sync Jobs" to fully event-driven architecture using Kafka Streams for sub-second risk updates.

EduPath-MS proves that when software engineering best practices meet advanced data science, the result is a powerful tool for educational equity, capable of identifying and supporting every learner at scale.

# Acknowledgements