# ABSTRACT

This project explores how machine learning can be used to detect fraudulent financial transactions. Using a real-world dataset of over 1.2 million records, a predictive pipeline was developed combining feature engineering, anonymization, and an XGBoost classifier optimized for recall.
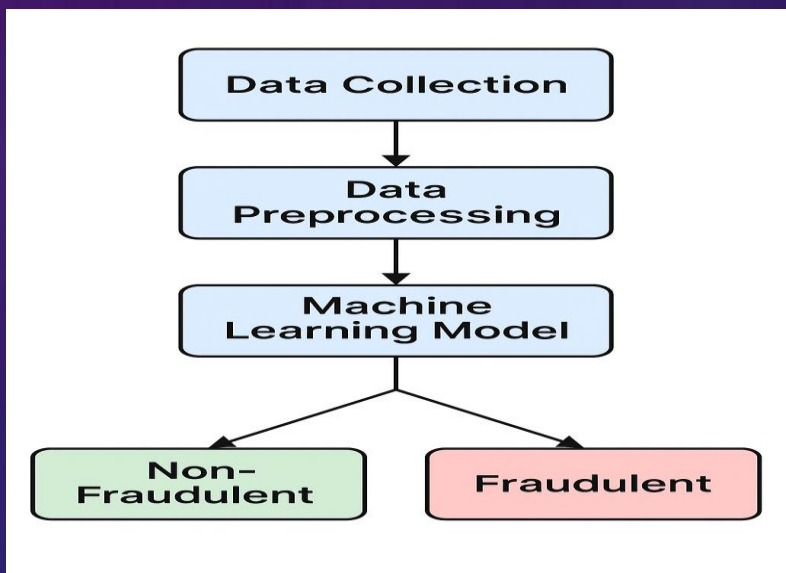
## Introduction & Background

Financial fraud detection is a critical application of machine learning. With fraud cases being rare but highly damaging, this project focuses on developing an effective classifier that prioritizes recall. The dataset was sourced from a real-world transactional system and includes time, location, demographic, and transactional metadata. The aim is to detect fraud while minimizing false negatives.

## Diagram / Design

The project uses a modular pipeline including:
- Data loading and cleaning
- Feature engineering (log transforms, interactions)
- SHA-256 anonymisation
- Preprocessing (scaling, encoding)
- XGBoost model training
- Threshold tuning for optimal Fβ-score



## Specification & Implementation

• Data: 1.3M transactions, sampled to 50K for efficient training
• Sensitive fields hashed for privacy (SHA-256)
• Feature engineered: amt_log, city_pop_log, amt_category
• XGBoost classifier tuned via GridSearchCV
• Metrics: F1, Fβ (β=2), Confusion Matrix
• Tools: Python, Pandas, Scikit-learn, XGBoost

## Testing & Evaluation

Model performance was evaluated over thresholds from 0.2 to 0.5.
At threshold 0.5:
- Accuracy: 99.45%
- Precision: 51.85%
- Recall: 72.41%
- F1 Score: 0.60
- Fβ Score (β=2): 0.67
Model generalised well across cross-validation folds and proved effective for fraud detection.

## Conclusions & Future Work

This project successfully developed a fraud detection model that balances high recall with usable precision. The model can serve as the foundation for production-ready systems.
Future improvements include integrating real-time streaming data, expanding feature sets, and deploying the model as a service.