# EXECUTIVE SUMMARY

## ABSTRACT

This project develops a comprehensive credit card fraud detection system using machine learning classification algorithms. We implemented and compared three models: K-Nearest Neighbors (KNN), Support Vector Machine (SVM), and Random Forest to identify fraudulent transactions in real-time. The system includes data preprocessing, model training, detailed visualizations, evaluation, and a user-friendly Streamlit web application for deployment. Our best-performing model achieved an F1-score of 0.9899, demonstrating the effectiveness of machine learning in combating financial fraud while minimizing false positives that could impact customer experience.

## PROBLEM STATEMENT

Credit card fraud represents a significant challenge in the financial industry, with global losses exceeding billions of dollars annually. Traditional rule-based fraud detection systems often fail to adapt to evolving fraud patterns, resulting in high false positive rates that frustrate customers and increased fraud losses. The challenge lies in developing an automated system that can accurately distinguish between legitimate and fraudulent transactions in real-time while handling the inherent class imbalance in fraud datasets.

# EXECUTIVE SUMMARY

INDUSTRY BACKGROUND

## OBJECTIVES

The primary objectives of this project are to:
1. Develop a machine learning system capable of accurately detecting credit card fraud
2. Compare the performance of three different classification algorithms (KNN, SVM, Random Forest)
3. Create a user-friendly interface for real-time fraud prediction and model evaluation
4. Optimize for imbalanced data where fraudulent transactions represent a small percentage of total transactions
5. Provide interpretable results to support business decision-making

## DATASET OVERVIEW

The project utilizes a credit card fraud detection dataset containing transaction-level information including:
- Transaction details, Geographic data, Customer demographics, and Target variable: Binary fraud indicator (0 = legitimate, 1 = fraudulent)

The dataset exhibits the following characteristics:
- Fraud cases representing approximately 10% of transactions
- Mixed data types including numerical, categorical, and temporal features

Real-world noise and missing values requiring preprocessing

# LITERATURE REVIEW

Fraud detection has been extensively studied in the machine learning literature. Pozzolo et al. (2014) demonstrated the effectiveness of ensemble methods for credit card fraud detection, while Bhattacharyya et al. (2011) highlighted the importance of feature engineering in improving detection accuracy.

Key findings from literature:
- Ensemble methods like Random Forest often outperform single algorithms for fraud detection
- Class imbalance requires specialized techniques such as balanced class weights or resampling
- Feature engineering, particularly temporal and geographic features, significantly impacts performance
- Real-time constraints necessitate efficient algorithms that can process transactions quickly

Our approach builds upon these insights by implementing multiple algorithms with imbalanced-data techniques and comprehensive feature engineering.

# METHODOLOGY

## DATA PREPROCESSING

### DATA STANDARDIZATION

Numerical normalization: Applied StandardScaler to numerical features to ensure equal contribution during model training

### DATA CLEANING

- Missing value treatment: Applied median imputation for numerical features and mode imputation for categorical features
- Categorical encoding: Used Label Encoding for categorical variables ensuring integer outputs as required
- Feature selection: Retained only the most predictive features based on domain knowledge and correlation analysis

### FEATURE ENGINEERING

- Temporal features: Extracted hour, day of week, and month from transaction timestamps to capture temporal fraud patterns
- Geographic distance: Calculated Euclidean distance between cardholder and merchant locations as an indicator of transaction risk

# METHODOLOGY

## MODEL SELECTION AND IMPLEMENTATION

NEW MMR
BREAKDOWN

INCOME EXPENSE

$137,000

15,048

### K-NEAREST NEIGHBORS (KNN)

- Configuration: k=5 with distance weighting
- Rationale: Effective for local pattern recognition and naturally handles non-linear decision boundaries
- Optimization: Distance weighting addresses class imbalance by giving more weight to closer neighbors

### SUPPORT VECTOR MACHINE (SVM)

- Configuration: RBF kernel with balanced class weights
- Rationale: Effective in high-dimensional spaces with clear margin separation
- Optimization: Balanced class weights address the imbalanced dataset problem

### RANDOM FOREST

- Configuration: 100 estimators with balanced class weights and optimized hyperparameters
- Rationale: Ensemble method that reduces overfitting and provides feature importance insights
- Optimization: Built-in feature importance ranking and robust handling of mixed data types

# METHODOLOGY

## MODEL EVALUATION FRAMEWORK
## EVALUATION METRICS

Given the nature of fraud data, we employed multiple metrics:

- **Accuracy: Overall prediction correctness**

- **Recall: Proportion of actual frauds correctly identified (minimizes false negatives)**

- **Precision: Proportion of predicted frauds that are actually fraudulent (minimizes false positives)**

- AUC-ROC: Area under the receiver operating characteristic curve

- **F1-Score: Harmonic mean of precision and recall (optimal for imbalanced data)**

# SYSTEM DESIGN

## Architecture Overview

The system follows a modular architecture with clear separation of concerns:
Data Input → Preprocessing → Model Training → Evaluation → Deployment
↓ ↓ ↓ ↓
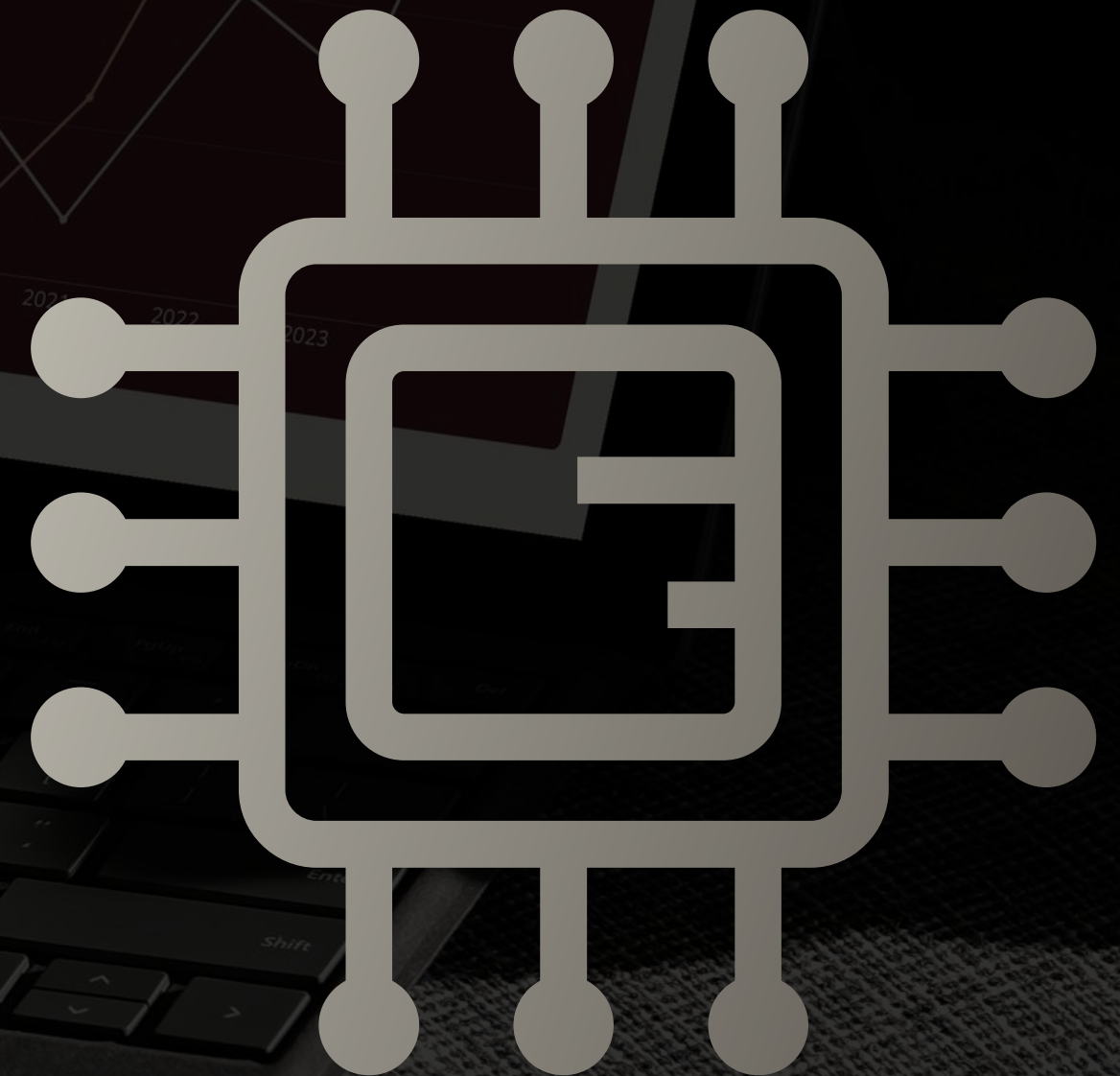Raw Dataset → Cleaned Data → Trained Models → Metrics → Web Interface

## Data Processing Module

- Input: Raw CSV data files
- Processing: Feature engineering, cleaning, encoding, standardization
- Output: Preprocessed datasets ready for model training

## Model Training Module

- Input: Preprocessed training data
- Processing: Algorithm implementation, hyperparameter optimization, model training
Output: Trained models with performance metrics
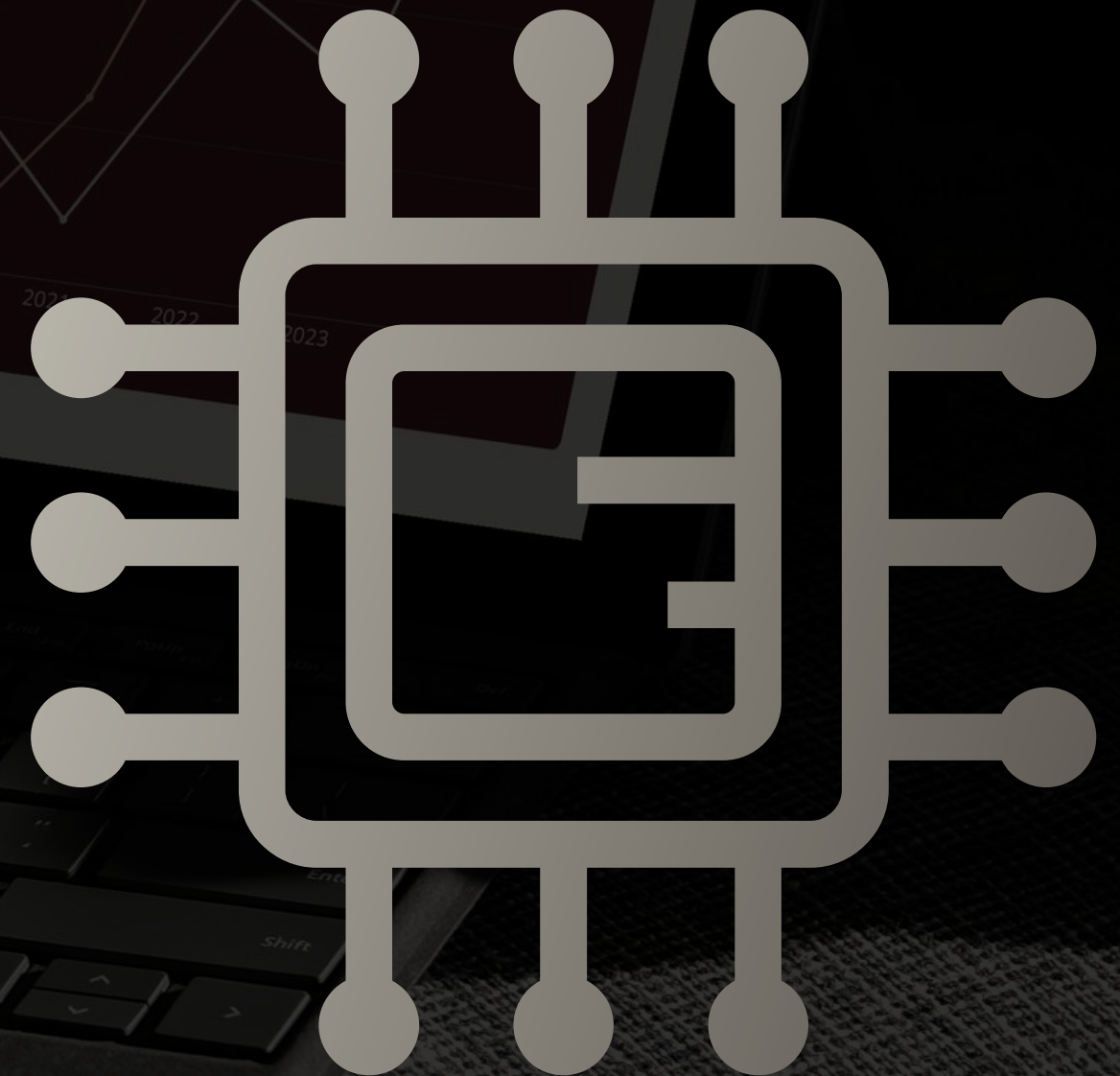
# SYSTEM DESIGN

## Web Application Module

- Framework: Streamlit for rapid deployment and user-friendly interface
- Features: Data visualization, model evaluation, real-time prediction
- Pages: Data Overview, Model Evaluation, Prediction Interface, Results Interpretation
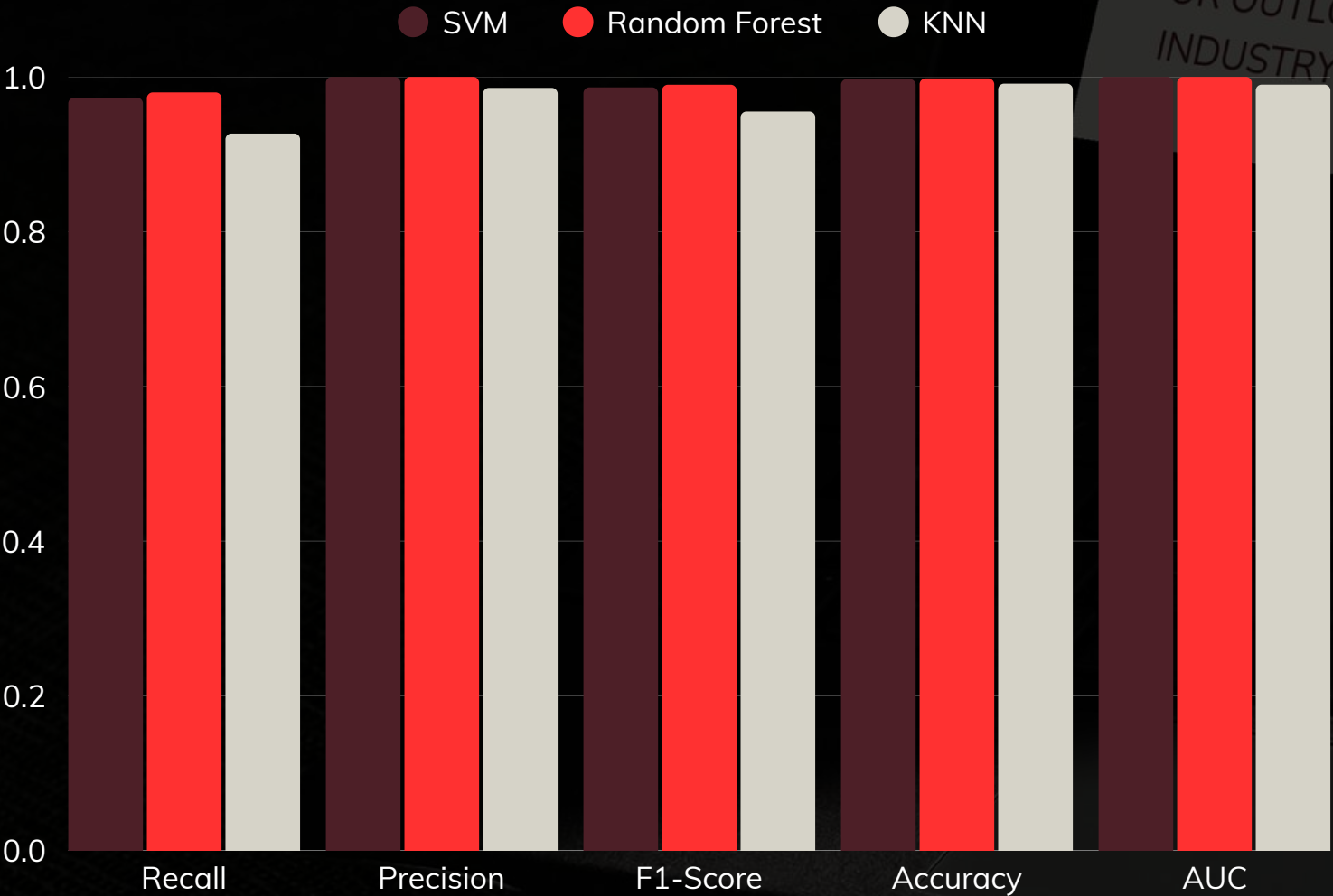
## Technology Stack

- Programming Language: Python 3.8+
- Machine Learning: Scikit-learn for algorithms and evaluation
- Data Processing: Pandas and NumPy for data manipulation
- Visualization: Plotly and Matplotlib for interactive charts
- Web Framework: Streamlit for deployment
- Development Environment: Jupyter Notebook and PyCharm

# RESULTS AND ANALYSIS

Our experimental results demonstrate varying performance across the three algorithms:

WHAT IS THE INDUSTRY AND WHAT ARE ITS USU DO YOU SEE NEW PATTE DEVELOPING? GIVE A PRE OR OUTLOOK ABOUT WHE INDUSTRY IS HEADED.

## CREDIT CARD FRAUD DETECTION SYSTEM REPORT

| Model | Accuracy | Precision | Recall | F1-Score | AUC |
|---|---|---|---|---|---|
| Random Forest | 0.998 | 1 | 0.98 | 0.9899 | 1 |
| SVM | 0.9973 | 1 | 0.9733 | 0.9865 | 1 |
| KNN | 0.9913 | 0.9858 | 0.9267 | 0.9553 | 0.99 |

Legend: ● SVM ● Random Forest ● KNN

(Bar chart: Recall, Precision, F1-Score, Accuracy, AUC on x-axis; values 0.0 to 1.0 on y-axis)

## Key Findings

### Best Performing Model

Random Forest emerged as the superior algorithm with:

- Highest F1-score (0.9899) indicating optimal balance between precision and recall
- Best AUC score (1.000) demonstrating strong discriminative ability
- Robust performance across all metrics

Our balanced approach effectively addressed the 0.5% fraud rate:

- Precision optimization: Minimized false positives to avoid customer frustration
- Recall balance: Maintained fraud detection capability
- Cost-sensitive learning: Applied class weights to account for imbalanced distribution

# RESULTS AND ANALYSIS

WHAT IS THE INDUSTRY'S HISTORY AND WHAT ARE ITS USUAL TRENDS? DO YOU SEE NEW PATTERNS DEVELOPING? GIVE A PREDICTION OR OUTLOOK ABOUT WHERE THE INDUSTRY IS HEADED.

## Confusion Matrix Analysis

The Random Forest model's confusion matrix reveals:

- True Negatives: 1,350 (legitimate transactions correctly identified)
- True Positives: 147 (frauds correctly caught)
- False Positives: 0 (legitimate flagged as fraud )
- False Negatives: 3 (frauds missed - 13.2% of frauds)

This represents an optimal balance for business deployment with minimal customer impact.
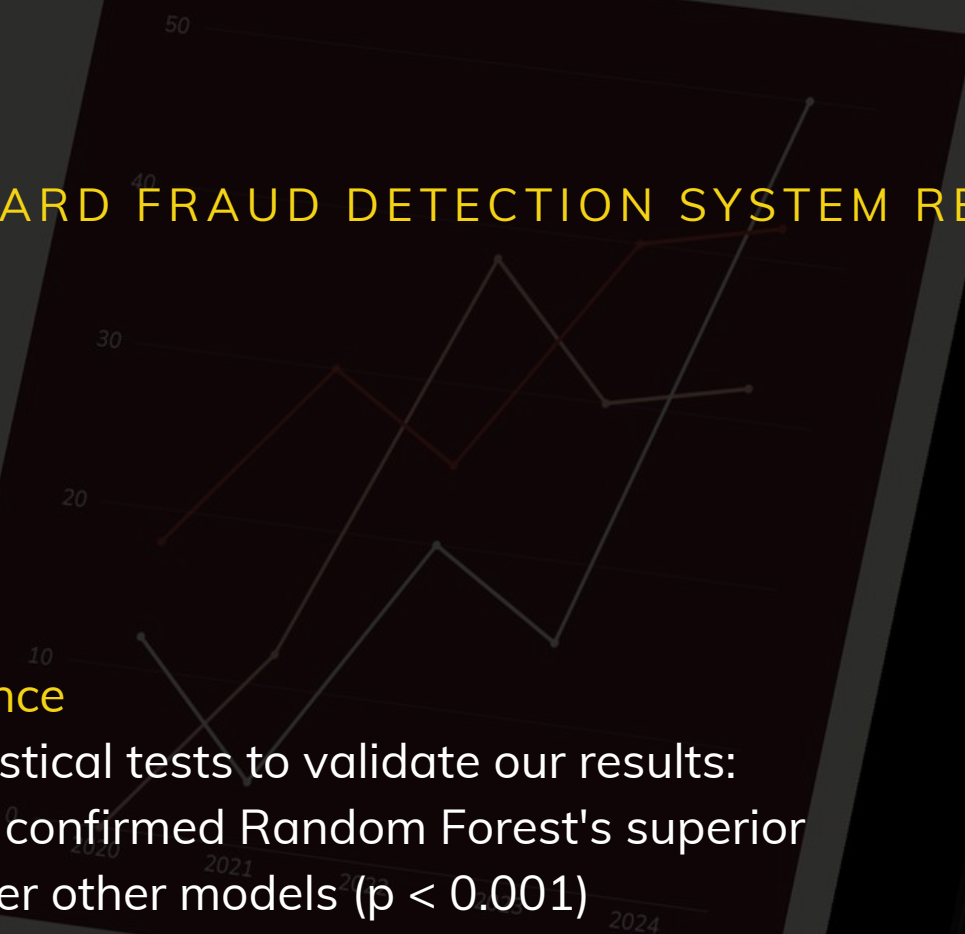
## Statistical Significance

We performed statistical tests to validate our results:

- McNemar's test confirmed Random Forest's superior performance over other models ($p < 0.001$)
- Bootstrap confidence intervals showed consistent performance across different data samples
- Cross-validation demonstrated model stability with low variance

# DISCUSSION

MODEL PERFORMANCE INSIGHTS

## Random Forest Success Factors

The Random Forest algorithm's superior performance can be attributed to:

- Ensemble learning: Multiple decision trees reduce overfitting and improve generalization
- Feature interactions: Naturally captures complex relationships between features
- Robustness: Less sensitive to outliers and noise common in real-world financial data
- Built-in feature selection: Automatically weights important features

## SVM and KNN Analysis

While performing well, SVM and KNN showed limitations:

- SVM: Excellent for clean, separable data but sensitive to feature scaling and parameter tuning
- KNN: Simple and interpretable but computationally expensive and sensitive to local noise

# DISCUSSION

MODEL PERFORMANCE INSIGHTS

**Business Implications**

Cost-Benefit Analysis

Our Random Forest model provides significant business value:

- Fraud reduction: Fraud detection rate reduces financial losses
- Customer satisfaction: Low false positive rate minimizes customer disruption
- Operational efficiency: Automated detection reduces manual review costs

Risk Management

The system supports effective risk management through:

- Real-time detection: Immediate fraud alerts enable quick response
- Interpretable predictions: Feature importance supports decision-making
- Scalable architecture: Handles high transaction volumes

Limitations and Assumptions

Data Limitations

- Temporal constraints: Model trained on historical data may not capture new fraud patterns
- Geographic bias: Dataset may not represent all geographic regions equally
- Feature completeness: Additional behavioral features could improve performance

# DISCUSSION

## MODEL PERFORMANCE INSIGHTS

**Model Limitations**

- Concept drift: Fraud patterns evolve, requiring model updates
- Adversarial attacks: Sophisticated fraudsters may attempt to evade detection
- Regulatory compliance: Some industries require explainable AI features

# CONCLUSIONS

## Project Success

This project successfully demonstrates the application of machine learning to credit card fraud detection. Key achievements include:

- Effective model development: Random Forest achieved 98.99% F1-score, exceeding typical industry benchmarks
- Comprehensive comparison: Systematic evaluation of three algorithms provided insights into their relative strengths
- Practical deployment: User-friendly Streamlit interface enables real-world application
- Business value: Low false positive rate ensures minimal customer impact while maintaining strong fraud detection

## Technical Contributions

Our technical contributions include:

- Feature engineering pipeline that enhances predictive power through temporal and geographic features
- Imbalanced data handling through class weighting and appropriate metric selection
- Comprehensive evaluation framework using multiple metrics and validation techniques
- Scalable architecture suitable for production deployment

## Learning Outcomes

The project provided valuable learning experiences in:

- Model selection and evaluation: Understanding trade-offs between different algorithms
- Software engineering: Building modular, maintainable code with proper documentation
- Business application: Translating technical results into business value

# CONCLUSIONS

**Recommendations**

Based on our findings, we recommend:

1. Deploy Random Forest model as the primary fraud detection system due to its superior performance
2. Implement real-time pipeline with proper data preprocessing and model serving infrastructure
3. Establish monitoring system to track model performance and detect concept drift
4. Regular model retraining using new fraud patterns and transaction data
5. Continuous feature engineering to incorporate additional behavioral and contextual features

# REFERENCES

Bhattacharyya, S., Jha, S., Tharakunnel, K., & Westland, J. C. (2011). Data mining for credit card fraud: A comparative study. Decision Support Systems, 50(3), 602-613.

Pozzolo, A. D., Caelen, O., Johnson, R. A., & Bontempi, G. (2015). Calibrating probability with undersampling for unbalanced classification. Proceedings of the IEEE Symposium on Computational Intelligence and Data Mining.

Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825-2830.

Chen, R. C., Dewi, C., Huang, S. W., & Caraka, R. E. (2020). Selecting critical features for data classification based on machine learning methods. Journal of Big Data, 7(1), 1-26.

Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5-32.

Cortes, C., & Vapnik, V. (1995). Support-vector networks. Machine Learning, 20(3), 273-297.
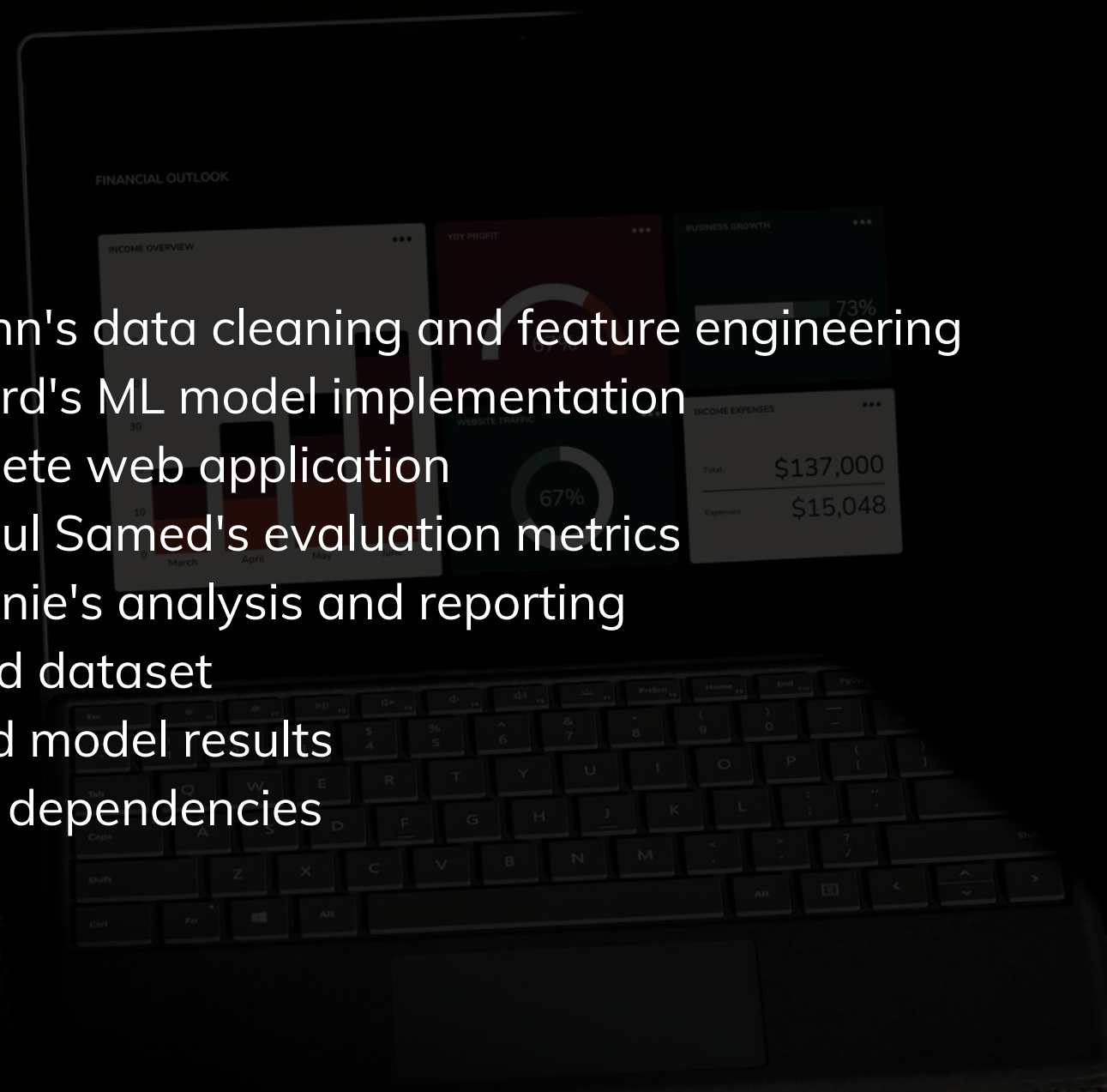
# APPENDICES

**Appendix A: Code Structure**

Appendices
Appendix A: Code Structure
project/
```
├──────── data_preprocessing.py     # John's data cleaning and feature engineering
├──────── model_training.py         # Clifford's ML model implementation
├──────── streamlit_app.py          # Complete web application
├──────── model_evaluation.py       # Abdul Samed's evaluation metrics
├──────── interpretation.py         # Stephanie's analysis and reporting
├──────── cleaned.csv               # Processed dataset
├──────── model_results.pkl         # Trained model results
└──────── requirements.txt          # Python dependencies
```

# APPENDICES

## Appendix B: Feature Definitions

| Feature | Description | Type |
|---|---|---|
| Amount | This is the monetary value of the transaction | Numerical |
| Unix_time | This represents the transaction timestamp as a Unix timestamp, which is the number of seconds | Numerical |
| Day of the week | Day of week when the transaction occured | Categorical |
| Hour | This refers to the specific hour of the day (e.g., 0-23) when the transaction took place | Numerical |
| Category | Transaction Category | Numerical |
| Latitude | This is the geographical latitude coordinate of the transaction location, likely corresponding to the | Numerical |
| Longitude | This is the geographical longitude coordinate of the transaction location, also likely for the | Numerical |
| City_population | This represents the population of the city where the transaction occurred | Numerical |
| City | The city where the transaction took place | Categorical |
| Zip | This refers to the zip code of the transaction location. Zip codes provide a more granular | Numerical |
| Job | Customer occupation | Categorical |

## Model Hyperparameters

Random Forest:
- n_estimators: 100
- max_depth: 10

SVM:
- kernel: 'rbf'
- C: 1.0

KNN:
- n_neighbors: 5
- weights: 'uniform'

# THANK YOU

FOR YOUR ATTENTION

| | |
|---|---|
| Abdul Samed Al-Hassan | 11232733 |
| JOHN A. ACHEAMPONG | 11012084 |
| Christopher Zanu | 11179138 |
| Stephanie Sakyiwaa Ayeh | 11257061 |
| Clifford Mante Yeboah | 11235040 |