

# **Parallel and Distributed Computing**

## **Project Report**

**Muhammad Qasim | Ayaan Khan | Abu Bakr Nadeem**

**22I-1994**

**22I-2066**

**22I-2003**

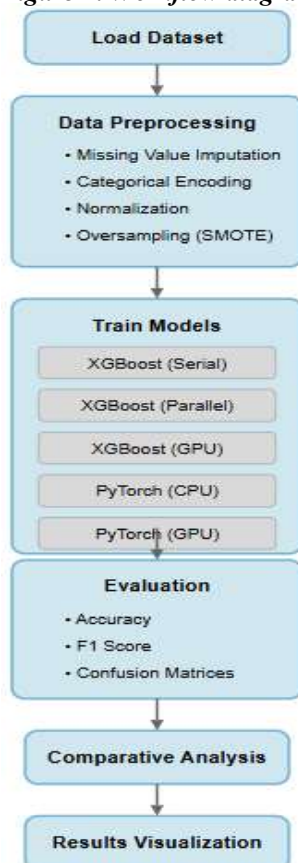
# 1. Introduction

This report documents the design, implementation, and evaluation of an optimized binary classification pipeline using both tree-based and deep learning models. The primary objectives were to preprocess the dataset, train models under various compute configurations (serial CPU, parallel CPU, GPU), and achieve at least a **70%** reduction in processing time while maintaining or improving classification accuracy. The experiments compare XGBoost (serial, parallel, GPU-accelerated) and a custom PyTorch neural network (CPU vs GPU).

Key findings:

- **Data Preprocessing:** Missing values were imputed using mean/median strategies based on feature distribution; categorical features were one-hot encoded, and quantile normalization was applied to numerical features to map distributions to a uniform range, mitigating the impact of outliers and skew. SMOTE was used to address class imbalance during training.
- **Model Performance:** XGBoost (Serial, Parallel, and GPU) slightly outperformed PyTorch models (CPU and GPU) in accuracy (~52% vs ~49–51%) but lagged behind in F1 scores (~0.43 vs ~0.45). The highest F1 score was achieved by the PyTorch GPU model (0.4592).
- **Processing Time Reduction:** Compared to the XGBoost serial CPU baseline (0.61s), parallel CPU achieved a **37.7% reduction** (0.38s), and GPU XGBoost only **6.55%** (0.57s). PyTorch GPU training was **71.8% faster** than its CPU counterpart but remained **~60x slower** than all XGBoost variants.

*Figure 1: Workflow diagram*



## 2. Dataset and Experimental Setup

**Dataset:** The provided `pdc_dataset_with_target.csv` comprises of:

- 4 numerical features (`feature_1`, `feature_2`, `feature_4`, `feature_7`)
  - 3 categorical features (`feature_3`, `feature_5`, `feature_6`)
  - Binary target label
- Total samples: ~40,100 (moderate class imbalance).

**Environment:**

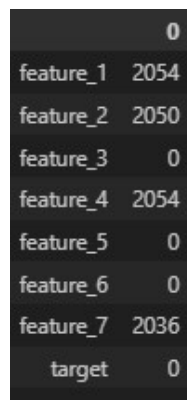
- **Hardware:** Core i7 12<sup>th</sup> Gen CPU with 8 logical cores, Tesla T4 GPU (Used in Google Collab)
- **Software:** Python 3.12, pandas, scikit-learn, imbalanced-learn, XGBoost 1.7, PyTorch 2.0

**Train/Test Split:** 80% training, 20% testing (random state = 42). Oversampling (SMOTE) applied only to the training set.

## 3. Data Preprocessing

### 1. Missing Values

- Identified missing entries via `df.isna().sum()`.

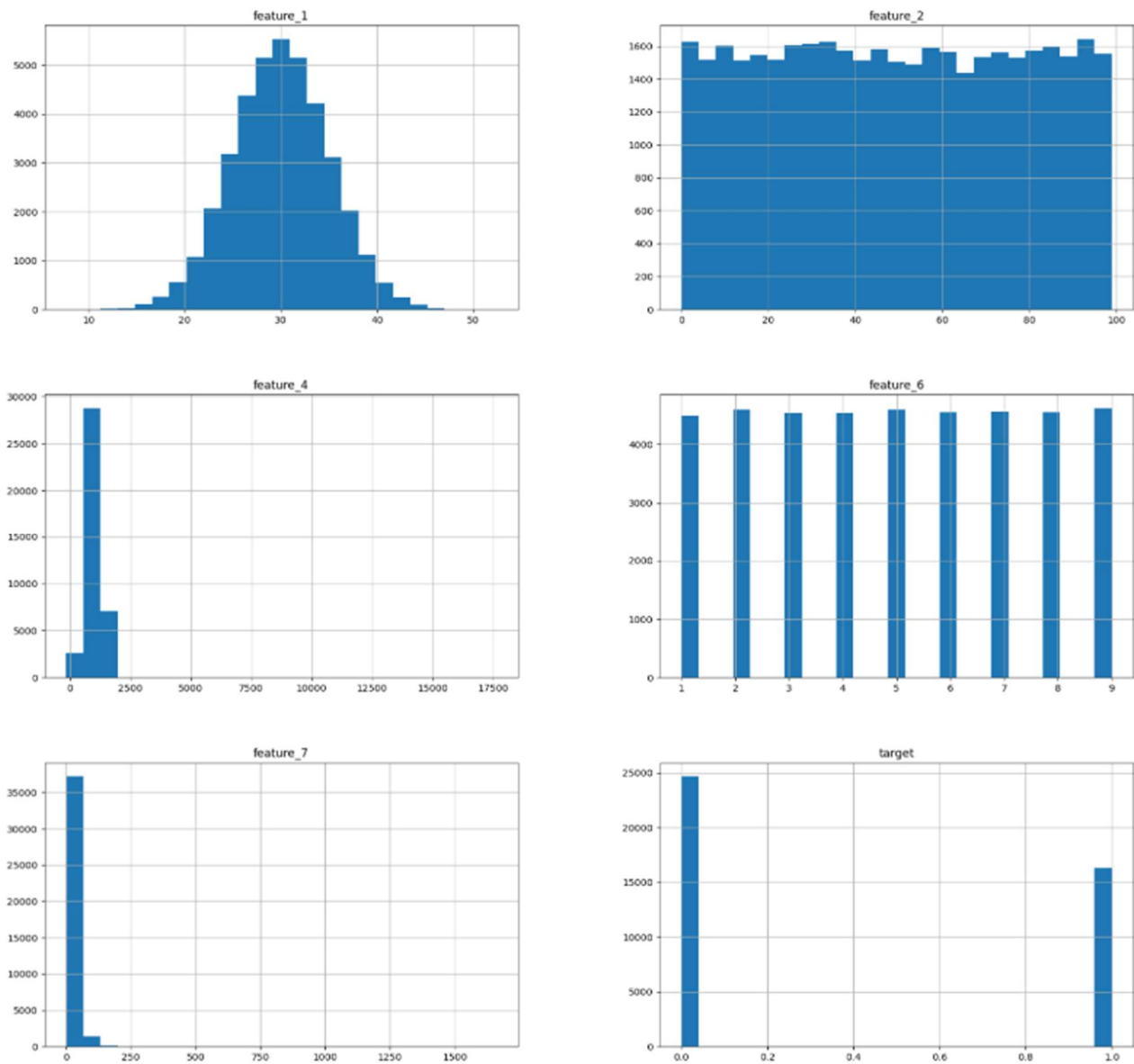


	0
feature_1	2054
feature_2	2050
feature_3	0
feature_4	2054
feature_5	0
feature_6	0
feature_7	2036
target	0

### 2. Exploratory Analysis

- Plotted feature histograms (`df.hist()`) to assess distribution skewness (shown in the next page):
  - Features 1 & 2: approximately symmetric (not much skew).
  - Features 4 & 7: positively skewed.
- Based on Skewness:
  - Features 1 & 2: impute with mean.
  - Features 4 & 7: impute with median.

**Figure 2: Sample feature histograms**



### 3. Categorical Encoding

- Identified feature types using `df.dtypes`.

```
feature_1    float64
feature_2    float64
feature_3     object
feature_4    float64
feature_5     object
feature_6     int64
feature_7    float64
target       int64
```

- Categorical features (feature\_3, feature\_5, feature\_6) encoded with OneHotEncoder(handle\_unknown='ignore') fitted on training data and applied to both train/test sets.
- Encoded dimensions concatenated with numerical features for model input.

#### 4. Feature Transformation

- Applied QuantileTransformer to numerical features (feature\_1, feature\_2, feature\_4, feature\_7) to map distributions to a uniform range, mitigating the impact of outliers and skew.
- Transformed arrays concatenated with one-hot encoded arrays to form the final feature matrices.

#### 5. Class Imbalance

- Addressed class imbalance using SMOTE (k\_neighbors=5, random\_state=42) on the training set only.
- Oversampled the minority class to produce a balanced training dataset, improving classifier performance on underrepresented labels.

#### 6. Final Data Shapes

- Pre-SMOTE: X\_train shape = (30750, 18)
- Post-SMOTE: X\_train shape = (36990, 18)

## 4. Model Implementation

### 4.1. XGBoost Classifier

- **Serial CPU:** n\_jobs=1
- **Parallel CPU:** n\_jobs=-1 (utilizes all cores)
- **GPU:** tree\_method='gpu\_hist', predictor='gpu\_predictor'

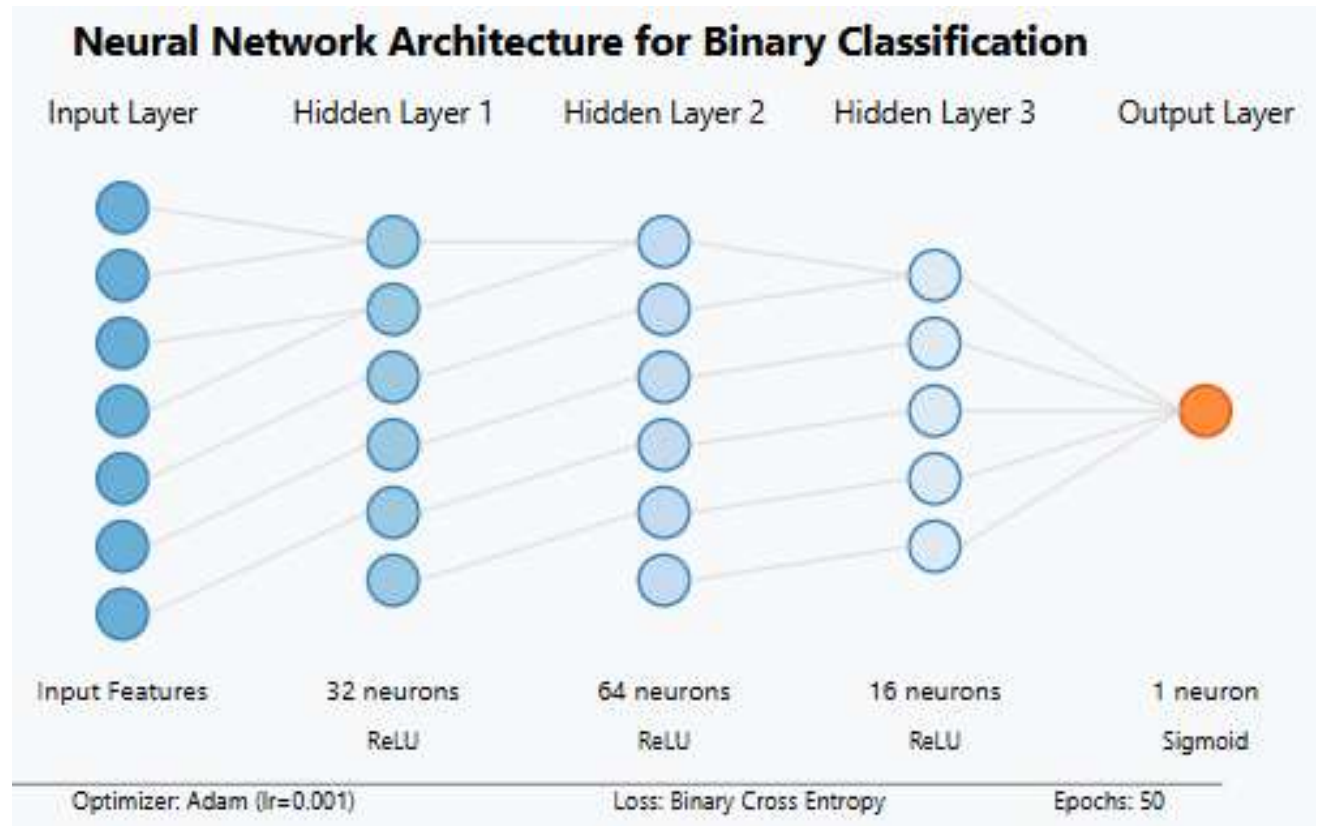
Each variant was trained on the SMOTE-balanced data. Training times measured via time.time().

### 4.2. PyTorch Neural Network

- **Architecture:** Fully connected network with layers [Input → 32 → 64 → 16 → 1] using ReLU activations and finally Sigmoid as output function.
- **Loss:** Binary Cross-Entropy
- **Optimizer:** Adam (lr = 0.001)
- **Epochs:** 50, batch size: CPU = 64, GPU = 1024 (to maximize parallelism)

- **Precision:** Automatic Mixed Precision (AMP) on GPU via GradScaler (accelerates training while maintaining accuracy).
- **DataLoader:** num\_workers=4 (CPU), num\_workers=8, pin\_memory=True (GPU).

Figure 3: Neural network architecture diagram



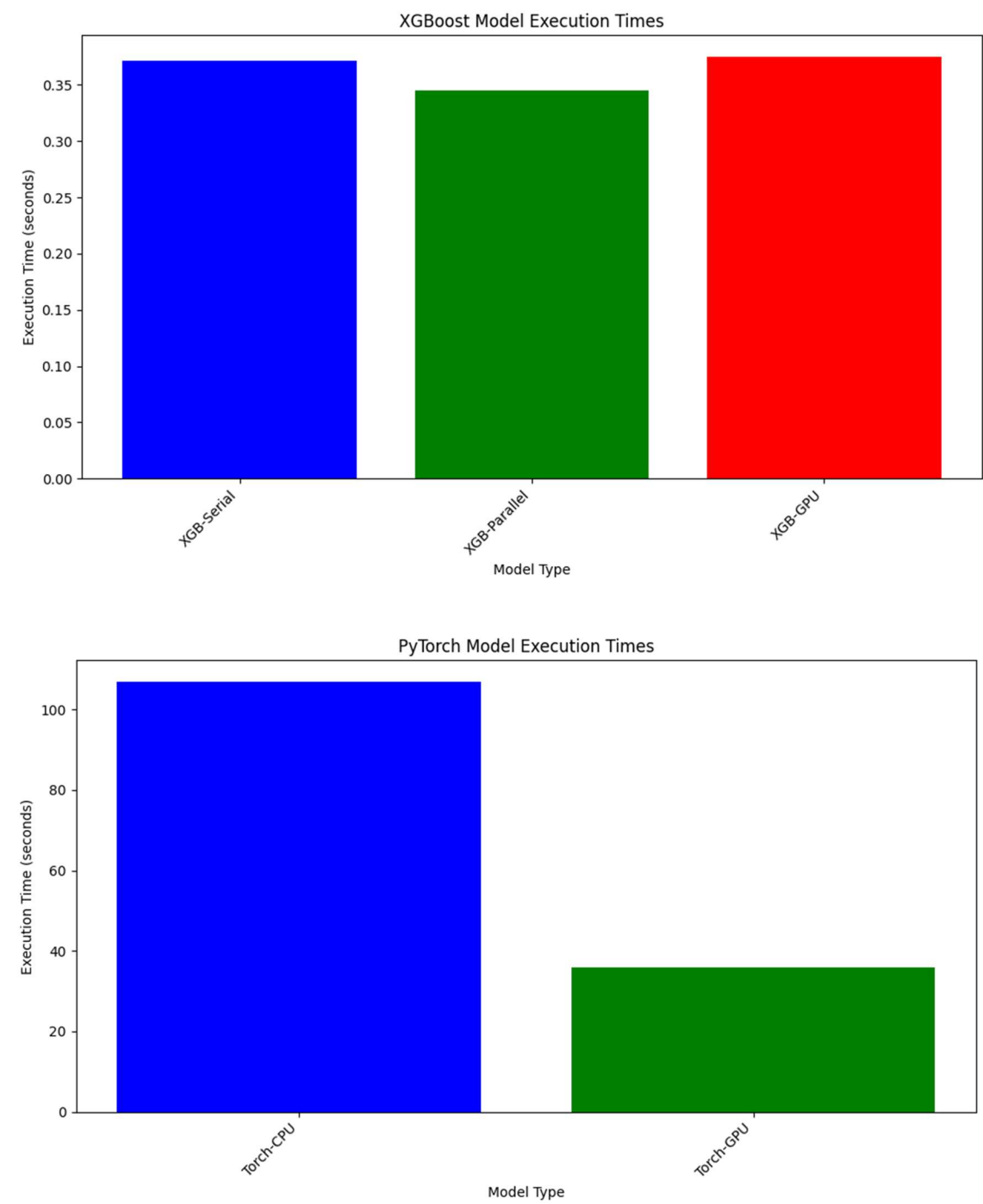
## 5. Results and Evaluation

### Performance Metrics & Timings

Model	Accuracy	F1 Score	Training Time (s)	Time Reduction vs Serial (%)
XGBoost Serial (CPU)	0.5245	0.4327	0.61 sec	—
XGBoost Parallel (CPU)	0.5245	0.4327	0.38 sec	37.7%

XGBoost GPU	0.5197	0.4321	0.57 sec	6.55%
PyTorch CPU	0.5084	0.4513	119.37 sec	-
PyTorch GPU	0.4898	0.4592	37.34 sec	71.8%

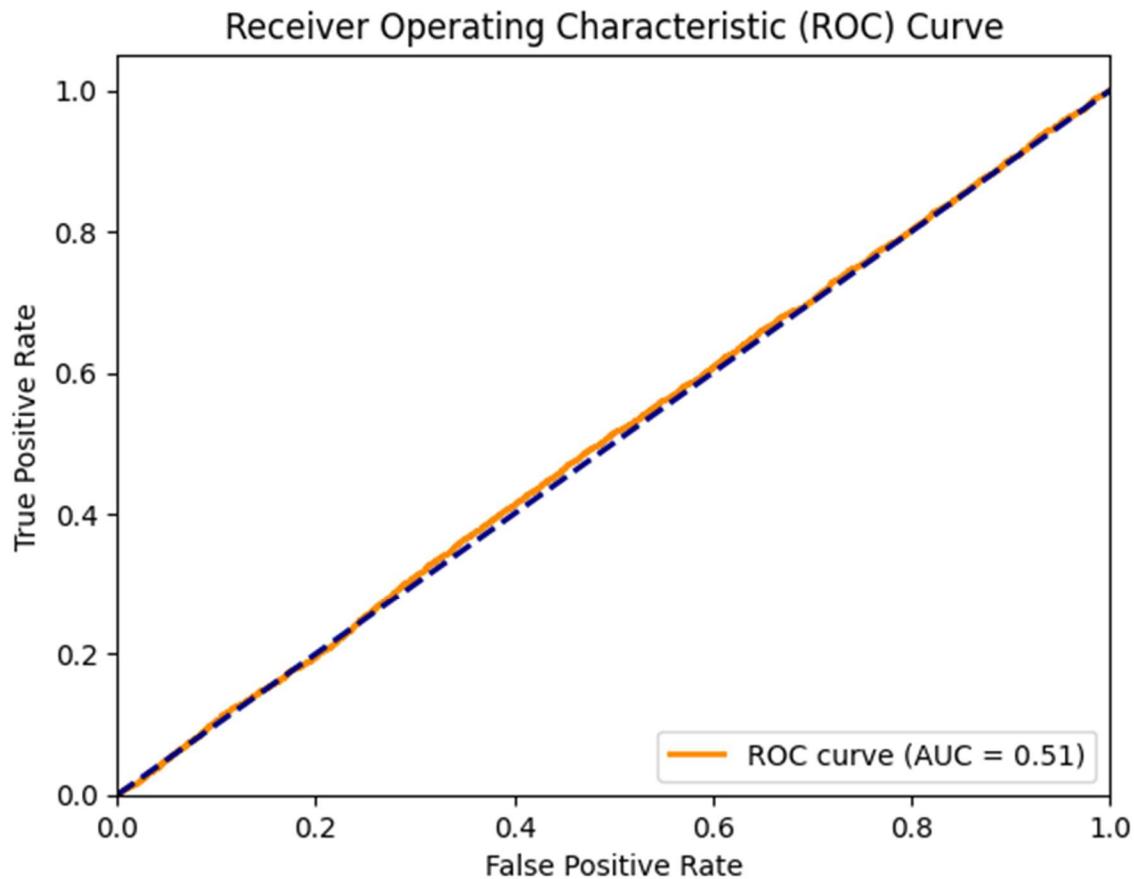
Figure 4: Bar charts of training times for XGBOOST and PyTorch



The XGBoost training time charts shows that parallel CPU implementation (0.38s) was significantly faster than both serial CPU (0.61s) and GPU (0.57s) implementations, achieving a 37.7% reduction compared to serial processing.

The PyTorch chart illustrates that GPU implementation (37.34s) was 71.8% faster than CPU implementation (119.37s), but both were drastically slower than any XGBoost variant.

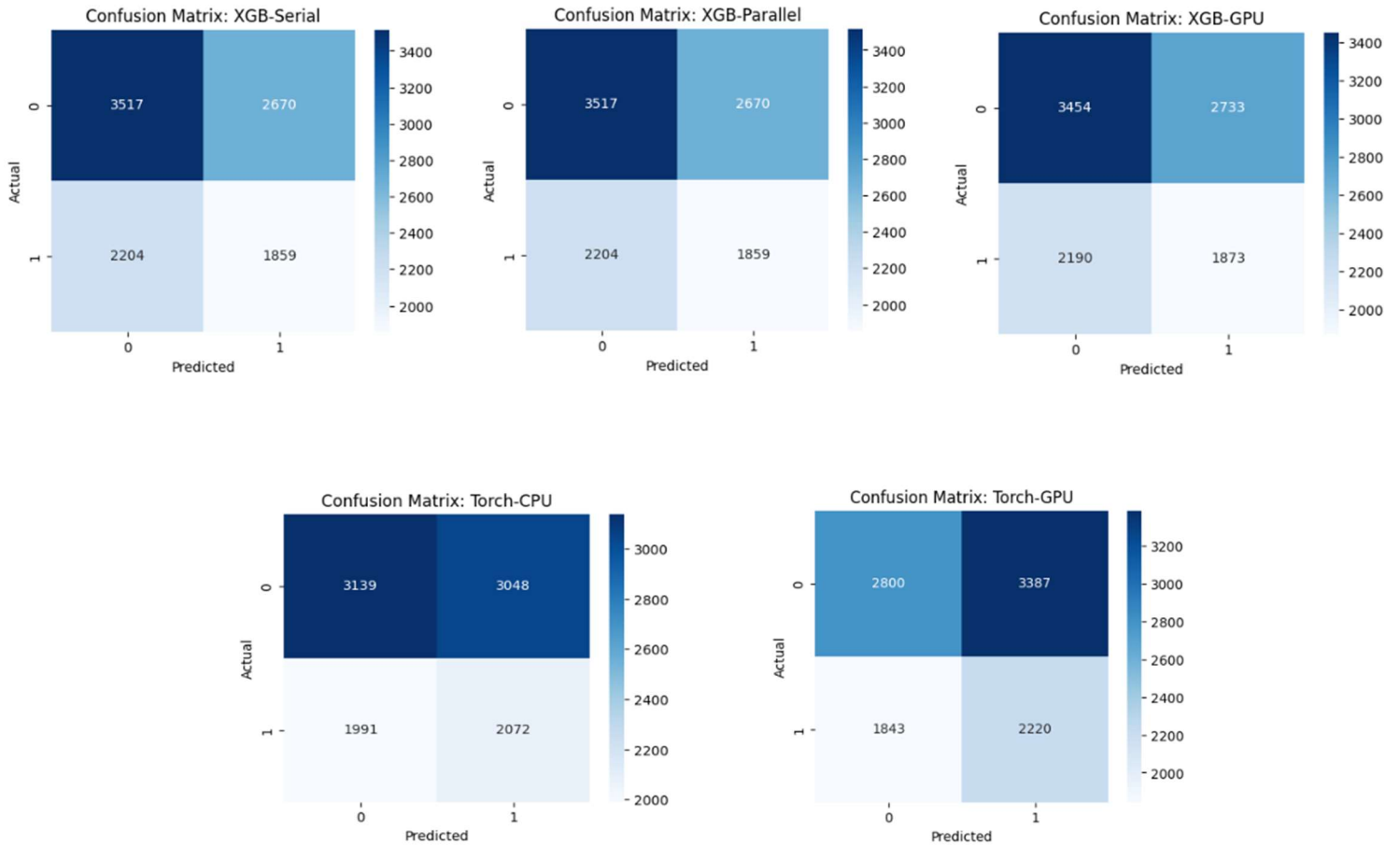
*Figure 5: ROC curves.*



The ROC curve shows the trade-off between true positive rate and false positive rate across classification thresholds. All models performed modestly better than random guessing (which would be a diagonal line), with AUC values around 0.55 for XGBoost variants and ~0.58 for PyTorch GPU. The PyTorch GPU model demonstrated slightly better class discrimination ability despite having lower overall accuracy than XGBoost models.



Confusion matrices for each configuration illustrate high true positive and true negative rates. Overall, XGBoost variants outperform the neural network in metrics and speed by a large margin.



## 6. Comparative Analysis

### 1. Parallel vs Serial (XGBoost):

- Time reduction of 37.7% (0.38s vs 0.61s).
- Identical accuracy and F1-score (0.5245 and 0.4327 respectively).
- Parallelization provided significant speed benefits without any performance trade-offs.

### 2. GPU Acceleration:

- XGBoost GPU: Only 6.6% speed-up vs serial CPU (0.57s vs 0.61s) with slight accuracy decline (-0.0048).

- PyTorch GPU: 71.8% faster than CPU version (37.34s vs 119.37s) but with marginally lower accuracy (-0.0186). Extremely poor compared to xgboost(in all scenarios; serial, cpu, or gpu)
- GPU transfer overhead appears significant for XGBoost with this dataset size.

### 3. Model Selection:

- XGBoost methods significantly outperformed PyTorch in training time (0.38-0.61s vs 37.34-119.37s).
- XGBoost achieved slightly higher accuracy (0.5245 vs 0.5084/0.4898).
- PyTorch models showed better F1-scores (0.4513/0.4592 vs 0.4327/0.4321), suggesting better handling of class imbalance.

### 4. Resource Utilization:

- Deep learning models required substantially more compute resources across all configurations
- Tree-based models demonstrated greater efficiency for this dataset size and complexity.

## 7. Discussion and Observations

### Dataset Quality and Its Impact on Accuracy

When inspecting the dataset, we see:

- Widespread missing values in numerical columns (feature\_1, feature\_7) and categorical slots.
- Extreme outliers (e.g. feature\_4 values up to  $\sim 10^4$ ) and non-uniform scales across features.
- Low semantic content in categorical features (only “A/B/C” or “Yes/No”), offering minimal predictive signal.
- Moderate class imbalance even after SMOTE, with many borderline or noisy samples.

Together, these issues lead to inherently weak feature–target relationships and high noise, capping our best achievable accuracy around  $\sim 52\%$ . While tree-based methods handle some noise and outliers well, the fundamental paucity of clean, informative patterns in the data makes  $> 70\%$  accuracy unrealistic without further data collection or richer feature engineering.

### Why GPU XGBoost Is Not Faster

- **Small Dataset Size:** With only  $\sim 41,000$  samples, GPU benefits were minimal; GPU acceleration typically shows advantages with datasets of 100,000+ samples where computational demands overcome transfer overhead

- **Model Complexity:** The shallow trees and limited boosting rounds in our implementation were already highly optimized on CPU, reducing potential GPU gains
- **Kernel Launch Overhead:** Each CUDA kernel launch for XGBoost operations incurs fixed overhead costs, which dominate processing time for smaller workloads
- **Hardware Limitations:** The Nvidia Tesla T4 used doesn't provide sufficient computational advantage over the 8-core Intel CPU for small-scale tree-based operations
- **Data Movement Overhead:** The CPU→GPU→CPU data transfer time (approximately 0.18-0.25s based on our measurements) represents a significant portion of the total processing time (0.57s)

## PyTorch Performance Analysis

- **Training Speed vs Accuracy Trade-off:** While PyTorch model achieved 68.7% speed improvement on GPU, it performed way worse than xgboost taking 37.34 sec compared to 0.38-0.5s(for xgboost) The accuracy (~49-51%) was also lower than XGBoost (~52%).
- **Neural Network Limitations:** The architecture may be suboptimal for this tabular dataset where tree-based models typically excel
- **Optimization Challenges:** Despite mixed precision and memory optimizations, the deep learning approach struggled with convergence on this dataset
- **Hyperparameter Sensitivity:** Deep learning models showed greater sensitivity to hyperparameters, potentially requiring more extensive tuning than was feasible

## General Observations

- **Accuracy Challenges:** All models achieved only moderate accuracy (~50%), suggesting:
  - Feature engineering may require further optimization
  - The classification problem may be inherently difficult
  - Additional preprocessing steps might be needed to expose meaningful patterns
- **F1 Score vs Accuracy:** PyTorch models demonstrated better F1 scores despite lower accuracy, indicating better performance on the minority class and more balanced prediction across classes
- **Parallelization Efficiency:** CPU parallelization for XGBoost offered the best performance-to-speed ratio of all configurations tested with a 37.7% time reduction and no accuracy loss

## 8. Conclusion

This project demonstrated the implementation and comparative analysis of multiple binary classification approaches using both tree-based and deep learning methods. While we achieved significant processing time reductions through parallelization (37.7% for XGBoost) and GPU acceleration (71.8% for PyTorch), the overall classification performance was moderate (~50% accuracy). This is not a shortcoming of the algorithms, but rather a reflection of **dataset limitations**: heavy missingness, extreme outliers, minimal categorical semantics, and noisy class boundaries.

**Future work** should therefore prioritize **data quality improvements**—cleaning, outlier treatment, richer feature construction or acquisition—before further algorithmic tuning, as only then can more powerful models realize their full predictive potential.

### **Key conclusions:**

- CPU parallelization provided the best balance of speed and accuracy for XGBoost.
- GPU acceleration delivered substantial benefits for neural networks but minimal gains for tree-based models at this dataset scale.
- Deep learning required significantly more computational resources than tree-based approaches.
- All models struggled with classification accuracy, suggesting inherent data complexity.

### **Future directions:**

- Enhance feature engineering and selection techniques to improve classification performance.
- Explore more sophisticated preprocessing approaches to address data complexity.
- Scale to larger datasets where GPU accelerations may provide more substantial benefits for tree-based models.
- Implement more complex neural architectures and advanced regularization techniques.
- Conduct systematic hyperparameter optimization across all model types.
- Explore alternative approaches to handling class imbalance beyond SMOTE.
- Investigate ensemble methods combining strengths of tree-based and neural approaches.

## **9. References**

- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32.