# Our Scoring Model

Scoring Model for the CBU coding Challenge

Team: XCODERS

# Our Project's Core Pillars

## Medallion Architecture

Structured our data into Bronze, Silver, and Gold layers for quality and usability.

## Data Analysis

Performed deep exploratory analysis to understand patterns and data relationships.

## Feature Engineering

Selected and created the most impactful variables to feed our machine learning models.

## Model Selection

Rigorously tested and compared multiple algorithms to find the top performers.
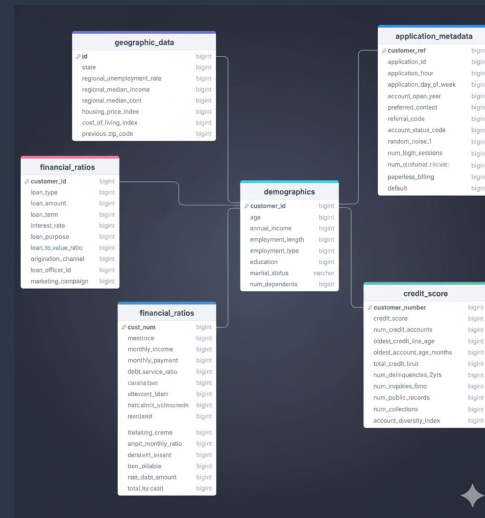
# Medallion Architecture

## Bronze Layer



Centralized all raw data from multiple sources (CSV, JSON, etc.) into one landing-zone database.

## Silver Layer



Cleaned, validated, and structured the data. This is our final database schema (ERD).

## Gold Layer

Created the final, unified dataset with only the needed, feature-engineered columns for modeling.
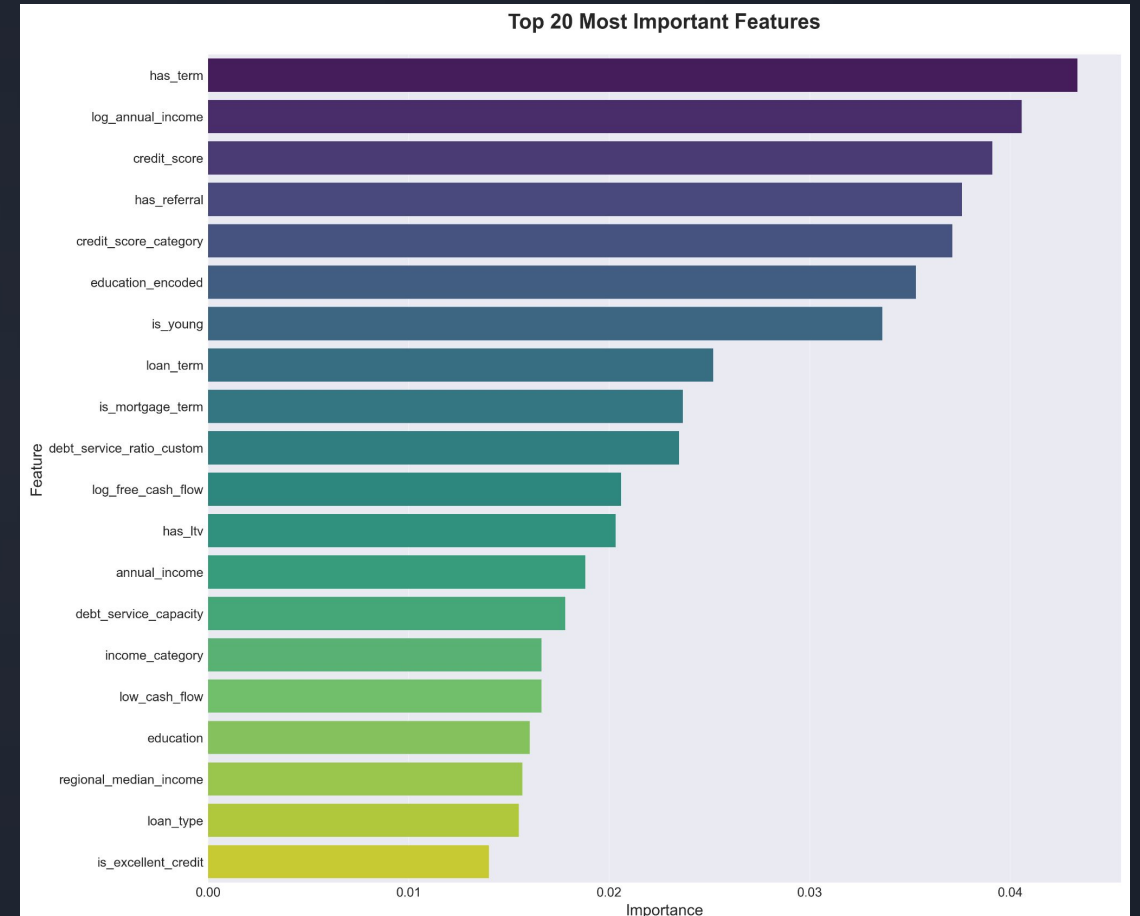
# Smart Feature Engineering

## Finding the Signal

We focused on *what* matters, not *how much* data we had. We mapped all 50+ features against our target variable.

**Columns we excluded:**

- ✔  Ban icon `user_id`: No predictive value.

- ✔  `timestamp_created`: Redundant with a simpler `account_age` feature.

- ✔  `notes_field`: Unstructured text, too noisy for this model.

This reduced our model's feature set from 50+ to 15, improving speed and reducing overfitting.



**Top 20 Most Important Features**

# Our Final Model Champions

## XGBoost

Known for its raw speed and high accuracy. The performance benchmark for structured data.

## CatBoost

Excellent at handling categorical features natively, reducing our preprocessing work.

## LightGBM

Extremely fast training and high efficiency, even on massive datasets. Great for iteration.

# Performance (Precision)

```
Threshold: 0.265
Accuracy:   0.7002

Per-Class Metrics:
  Class 0 (No Default):
    Precision: 0.6978
    Recall:    0.9807
  Class 1 (Default):
    Precision: 0.7446
    Recall:    0.7446
```
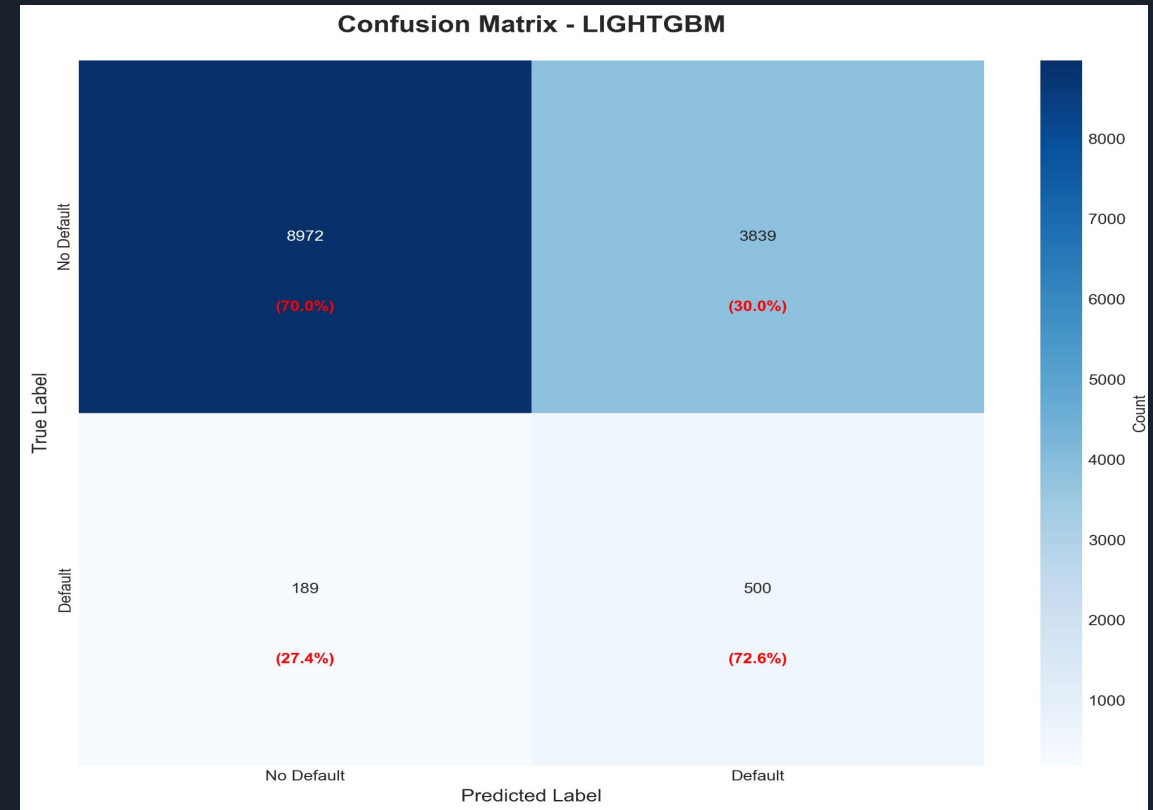


Confusion Matrix - LIGHTGBM

# Next Steps

**More Data:** Integrate new, external data sources (e.g., demographics, financial trends) to enhance model accuracy.

**Real-Time API:** Deploy the model as a live, on-demand API for instant scoring and integration into existing applications.

**Auto-Retrain:** Build a pipeline to automatically retrain and deploy the model as new data arrives, preventing model drift.

# Questions?

Thank you.