



# FraudFusion

DS 5500 – Capstone Project  
Anish Rao, Raghu Ram Sattanapalle

March 18<sup>th</sup>, 2025



# Objectives

## Goal:

- Our main goal is to **improve the synthetic fraud** data generation process to more accurately represent the real data

## Main Problem:

- **Severe class imbalance** in credit card fraud detection.
- **Need for high-quality synthetic fraud data** to improve model learning.

## Significance:

- Our 7th model iteration increased fraud detection recall by **5.75 percentage points** (82.75% to 88.50%).
- Potentially reduce financial losses by capturing more fraudulent transactions.



# Hypothesis

## Key hypotheses:

- Augmenting training data with these synthetic samples will **improve the recall** of fraud detection classifiers.
- **Custom loss functions** & feature-weighted learning will result in **more realistic** synthetic fraud samples

## Expected outcomes:

- **Improved alignment** between synthetic and real fraud transaction distributions
- **Measurable improvements in fraud detection** sensitivity with acceptable precision trade-offs

## Importance:

- In real-world fraud scenarios, **missing fraudulent transactions (false negatives) is costly**, leading to financial loss and reputational damage.

# Personal Learning Objectives

## Anish:

- I was not really familiar with generative models before working on this project, so increasing the depth of my knowledge on this topic is my primary objective
- More experience building advanced machine learning models from scratch
- Develop expertise in designing custom loss functions for specialized domains

## Raghu:

- Enhance proficiency in statistical analysis and data visualization methods to assess the quality of generated data
- Develop and refine expertise in the practical application of generative AI techniques
- Mastery in implementing and evaluating diffusion models for complex data



# Team Learning Objectives

## **Collective Technical Goals:**

- Master implementation of specialized diffusion models for tabular data with mixed feature types
- Develop a systematic approach for evaluating synthetic data quality beyond basic distribution matching
- Create reproducible methodologies for fraud detection that balance precision and recall trade-offs

## **Collaboration Enhancements:**

- Strengthen capability to iterate rapidly through model versions with structured evaluation criteria
- Develop a shared vocabulary and framework for discussing generative model performance

## **Expected Takeaways:**

- Comprehensive workflow for synthetic data generation in highly imbalanced domains
- Transferable approaches for custom loss function design in financial applications



# Challenges & Solutions

## Data Challenges

**Complex Feature Distributions:** Bimodal transaction amounts, cyclical patterns in temporal features

- *Solution:* Implemented specialized amount distribution loss and cyclic encoding for temporal features

## Technical Challenges

**Distribution Matching:** Difficulty capturing bimodal fraud patterns

- *Solution:* Developed custom loss functions:
  - Amount distribution loss with weighted quantile matching
  - Feature-weighted learning (amt: 1.8, time: 1.3)
  - Post-processing distribution matching

**Training Stability:** Diffusion models sensitive to mixed data types

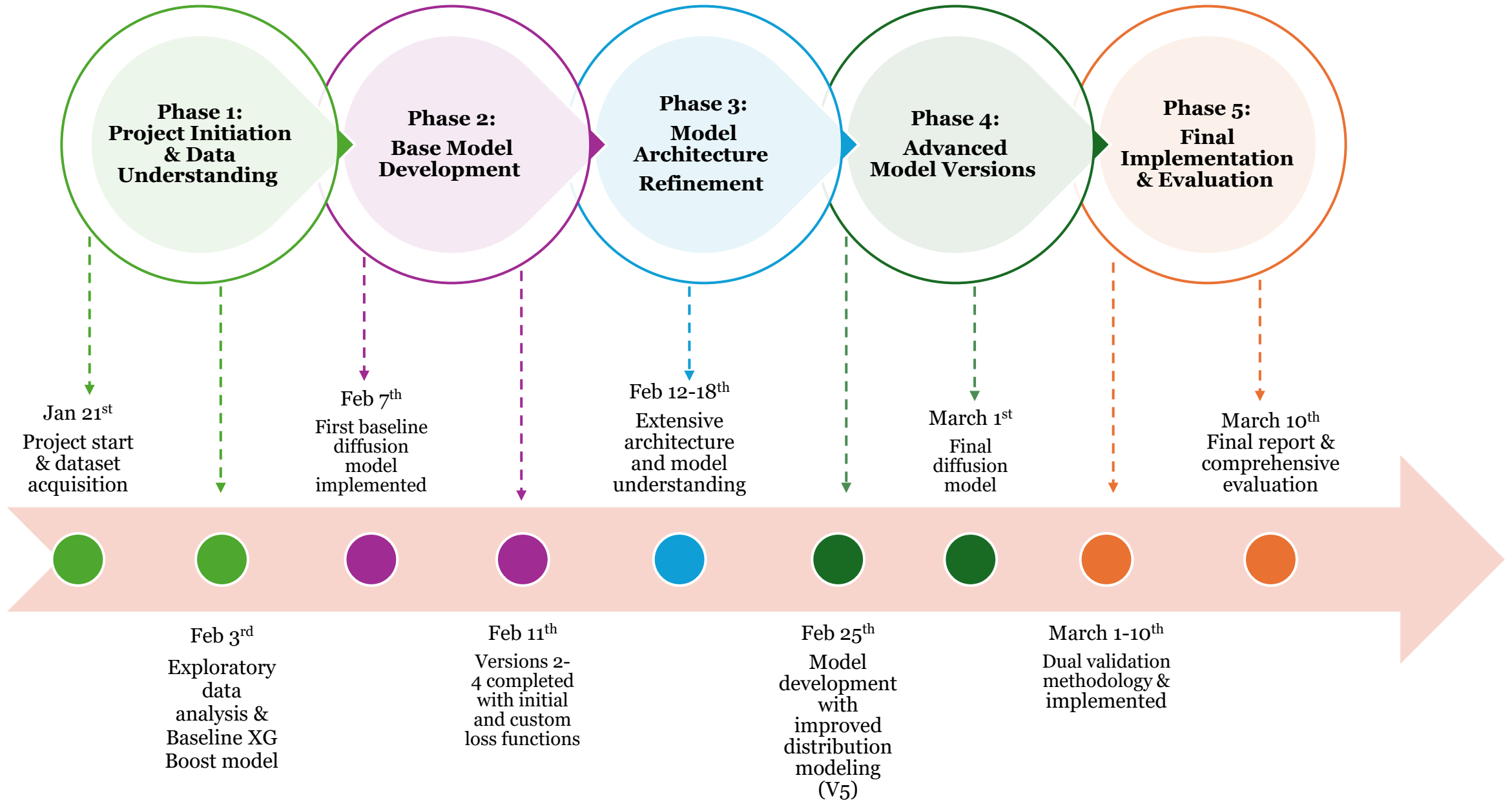
- *Solution:*
  - Gradient clipping (max\_norm=0.5)
  - Adaptive learning rate scheduling
  - Comprehensive NaN detection and recovery
  - Reduced batch size (32) for stability

## Validation Challenges

**Performance Assessment:** Risk of overfitting to synthetic patterns

- *Solution:* Implemented dual validation streams:
  - Pure validation set (real data only)
  - Synthetic validation set (combined real/synthetic)
  - Controlled synthetic sample allocation

# Timeline & Key Milestones





# Future Direction

## Model Enhancements:

- **Temporal Modeling:** Incorporate sequence modeling for transaction streams to detect complex fraud patterns over time
- **Transfer Learning:** Pre-train diffusion models on larger financial datasets, then fine – tune for specific types

## Expanded Applications:

- **Cross-Domain:** Apply to insurance fraud, anti-money laundering, and other financial crimes
- **Multimodal Fusion:** Combine transaction data with text (e.g., transaction descriptions)

## Further Research:

- **Adaptive Generation:** Design systems that dynamically adjust synthetic data characteristics based on classifier feedback
- **Advanced Distribution Matching:** Explore adversarial distribution matching for even more realistic synthetic data





# Contributions (Anish)

## **Feature Engineering & Data Preprocessing:**

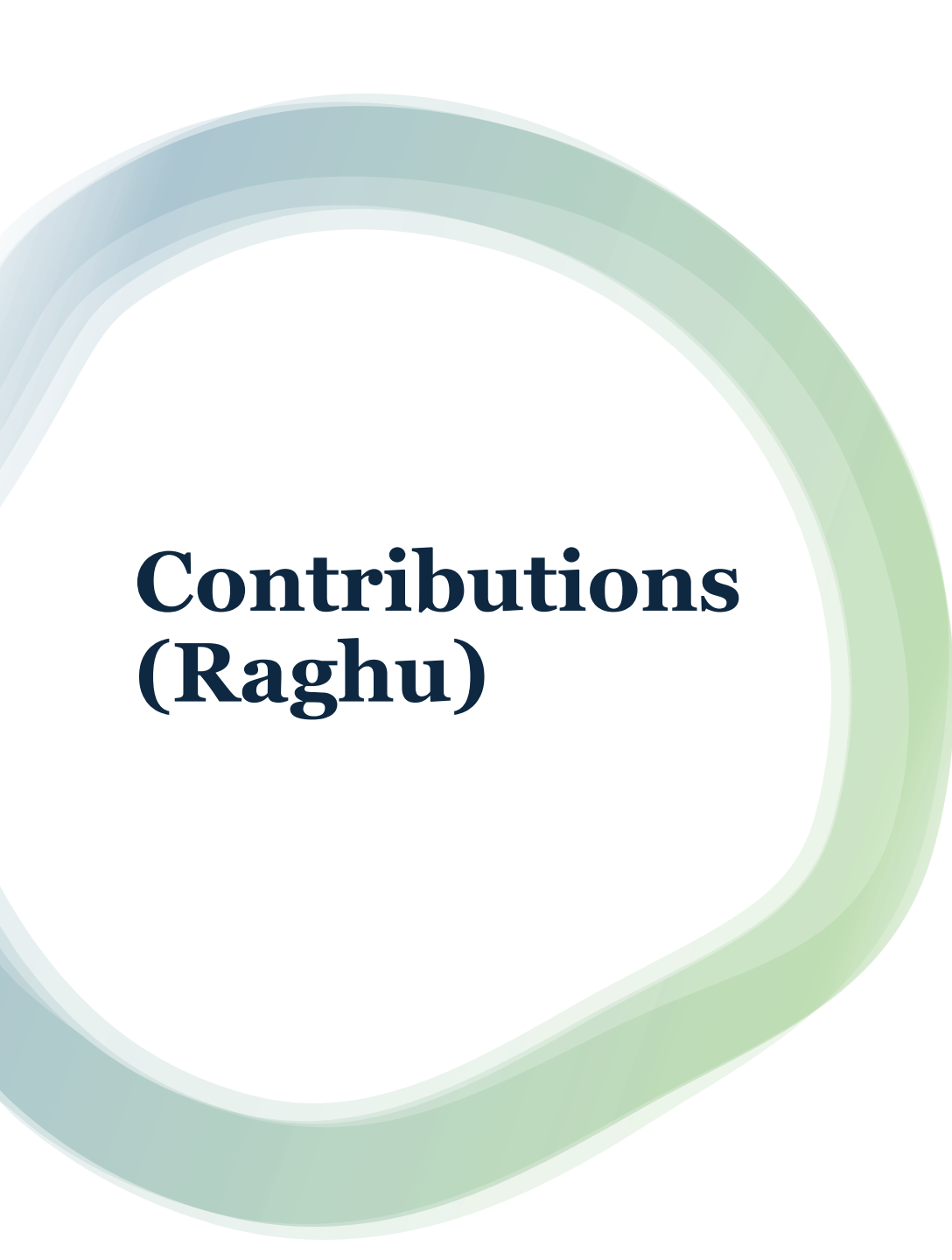
- Developed comprehensive pipeline for cleaning and preparing transaction data for diffusion model iterations
- Created specialized transformations for temporal and categorical features

## **XGBoost Implementation:**

- Designed and optimized the fraud detection classifier pipeline
- Created controlled synthetic data augmentation framework
- Conducted performance analysis across precision-recall spectrum

## **Diffusion Model development & testing:**

- Architecture & Hyper-parameter testing for diffusion model improvements



# Contributions (Raghu)

## **Feature Engineering:**

- Developed cyclic encoding for temporal fraud patterns

## **Diffusion Model development & testing:**

- Developed & enhanced FraudDiffuse architecture with custom components for mixed data types
- Designed & implemented custom loss functions (Amount distribution & Feature weighted MSE loss)
- Created framework for synthetic data quality assessment
- Implemented specialized bimodal distribution matching for transaction amounts
- Implemented post –processing distribution transformation techniques



**Thank you!**