# SafePayAI - Semi-Technical Pitch

*"Every 8 seconds, someone in India falls victim to a UPI fraud. As digital payments skyrocket to billions of transactions daily, fraudsters are evolving faster than traditional rule-based detection systems can keep up. SafePayAI changes that game."*

## The Problem

Traditional fraud detection relies on **static rules** — if amount > ₹50,000, flag it. But fraudsters adapt instantly.

**Core challenges:**

- **Class imbalance:** Legitimate transactions outnumber fraudulent ones 1000:1
- **Evolving patterns:** Fraud techniques change weekly
- **False positives:** Blocking legitimate users damages trust
- **Real-time demands:** Decisions must happen in milliseconds

## Our Solution

SafePayAI uses a **two-stage AI architecture**:

### 1️⃣ GAN-Powered Data Augmentation

We use **Generative Adversarial Networks** to solve the class imbalance problem:

- The **Generator** creates synthetic transaction data from random noise
- The **Discriminator** learns to distinguish real vs synthetic data
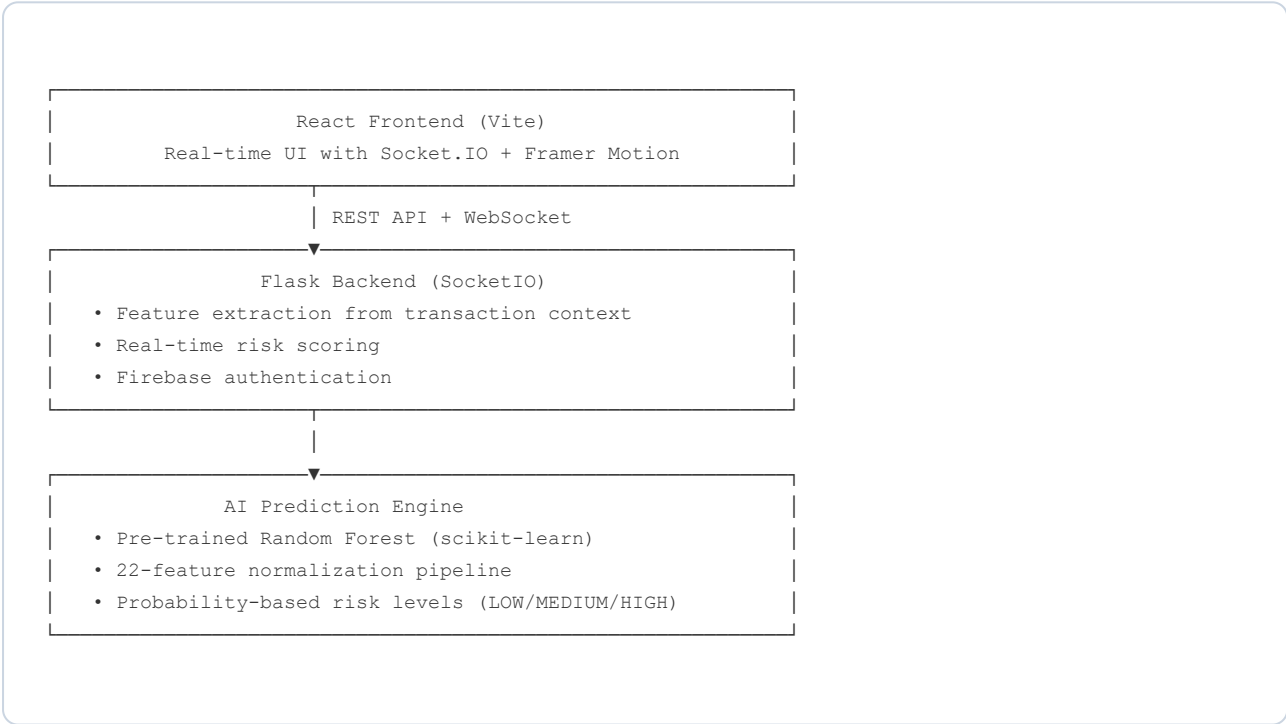- This adversarial training produces high-quality synthetic fraud data

> **Result:** Our model learns from 10x more fraud examples without collecting more real fraud cases.
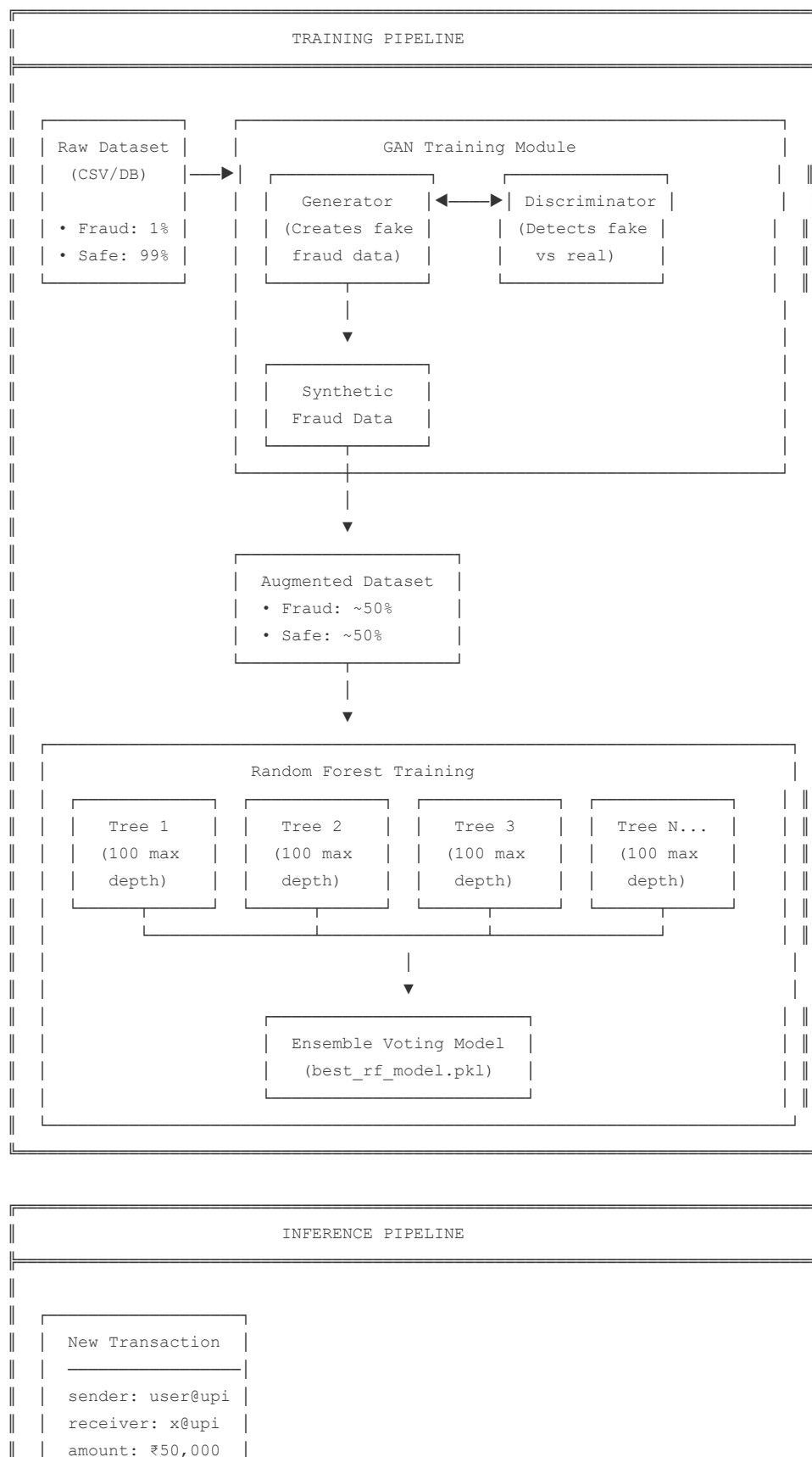
### 2️⃣ Random Forest Classifier

Our production model analyzes **22 real-time features** per transaction:

| Feature Category | Examples |
|---|---|
| **Behavioral** | Transaction frequency, amount deviation from history |
| **Recipient Risk** | Trust score, verification status, blacklist status |
| **Contextual** | Transaction hour, geo-location anomalies |
| **Device** | Device fingerprint, VPN/proxy detection |
| **Historical** | Past fraud complaints, account age |

## Technical Architecture

```
    ┌──────────────────────────────────────────────────┐
    │                React Frontend (Vite)             │
    │      Real-time UI with Socket.IO + Framer Motion │
    └──────────────────────────────────────────────────┘
                      │ REST API + WebSocket
    ┌──────────────────▼───────────────────────────────┐
    │              Flask Backend (SocketIO)            │
    │  • Feature extraction from transaction context   │
    │  • Real-time risk scoring                        │
    │  • Firebase authentication                       │
    └──────────────────────────────────────────────────┘
                      │
    ┌──────────────────▼───────────────────────────────┐
    │                AI Prediction Engine              │
    │  • Pre-trained Random Forest (scikit-learn)      │
    │  • 22-feature normalization pipeline             │
    │  • Probability-based risk levels (LOW/MEDIUM/HIGH)│
    └──────────────────────────────────────────────────┘
```

# ML Pipeline Architecture

```
╔═══════════════════════════════════════════════════════════════╗
║                      TRAINING PIPELINE                        ║
╠═══════════════════════════════════════════════════════════════╣
║                                                               ║
║  ┌───────────────┐   ┌─────────────────────────────────────┐ ║
║  │ Raw Dataset   │   │         GAN Training Module          │ ║
║  │  (CSV/DB)     │──▶│                                     │ ║
║  │               │   │  ┌─────────────┐   ┌──────────────┐ │ ║
║  │ • Fraud: 1%   │   │  │  Generator  │◀─▶│ Discriminator│ │ ║
║  │ • Safe: 99%   │   │  │ (Creates fake│   │ (Detects fake│ │ ║
║  │               │   │  │  fraud data) │   │   vs real)   │ │ ║
║  └───────────────┘   │  └─────────────┘   └──────────────┘ │ ║
║                      │         │                            │ ║
║                      │         ▼                            │ ║
║                      │  ┌─────────────┐                     │ ║
║                      │  │  Synthetic  │                     │ ║
║                      │  │ Fraud Data  │                     │ ║
║                      │  └─────────────┘                     │ ║
║                      └─────────────────────────────────────┘ ║
║                                │                              ║
║                                ▼                              ║
║                      ┌─────────────────────┐                 ║
║                      │  Augmented Dataset   │                 ║
║                      │  • Fraud: ~50%       │                 ║
║                      │  • Safe: ~50%        │                 ║
║                      └─────────────────────┘                 ║
║                                │                              ║
║                                ▼                              ║
║  ┌──────────────────────────────────────────────────────┐   ║
║  │              Random Forest Training                  │   ║
║  │  ┌─────────┐ ┌─────────┐ ┌─────────┐ ┌──────────┐   │   ║
║  │  │ Tree 1  │ │ Tree 2  │ │ Tree 3  │ │ Tree N...│   │   ║
║  │  │(100 max │ │(100 max │ │(100 max │ │(100 max  │   │   ║
║  │  │ depth)  │ │ depth)  │ │ depth)  │ │ depth)   │   │   ║
║  │  └─────────┘ └─────────┘ └─────────┘ └──────────┘   │   ║
║  │       └──────────┴────────┴──────────┘               │   ║
║  │                      │                               │   ║
║  │                      ▼                               │   ║
║  │            ┌─────────────────────┐                   │   ║
║  │            │ Ensemble Voting Model│                   │   ║
║  │            │  (best_rf_model.pkl) │                   │   ║
║  │            └─────────────────────┘                   │   ║
║  └──────────────────────────────────────────────────────┘   ║
╚═══════════════════════════════════════════════════════════════╝


╔═══════════════════════════════════════════════════════════════╗
║                     INFERENCE PIPELINE                        ║
╠═══════════════════════════════════════════════════════════════╣
║                                                               ║
║  ┌───────────────────┐                                        ║
║  │ New Transaction   │                                        ║
║  │ ───────────────── │                                        ║
║  │ sender: user@upi  │                                        ║
║  │ receiver: x@upi   │                                        ║
║  │ amount: ₹50,000   │                                        ║
```

```
                        Feature Extraction Engine

   ┌─────────────┐  ┌─────────────┐  ┌─────────────┐  ┌─────────────┐
   │ Behavioral  │  │  Recipient  │  │ Contextual  │  │ Historical  │
   │  Features   │  │   Profile   │  │   Signals   │  │    Data     │
   │ (5 values)  │  │ (6 values)  │  │ (5 values)  │  │ (6 values)  │
   └─────────────┘  └─────────────┘  └─────────────┘  └─────────────┘

                      ┌─────────────────────┐
                      │  22-Feature Vector   │
                      │  [0.2, 0.8, 0.1, ...] │
                      └─────────────────────┘


                        Random Forest Prediction

                  ┌──────────────────────────────┐
                  │      best_rf_model.pkl        │
                  │  model.predict_proba(features) │
                  └──────────────────────────────┘

                  ┌──────────────────────────────┐
                  │  Probability: [0.15, 0.85]    │
                  │  → 85% chance of fraud         │
                  └──────────────────────────────┘


                        Risk Level Classification

   Probability < 30%  ──────▶  ✅ LOW RISK     ──────▶  Approve
   Probability 30-70% ──────▶  ⚠️ MEDIUM RISK  ──────▶  Warn + Proceed
   Probability > 70%  ──────▶  🔴 HIGH RISK    ──────▶  Block + Explain
```

# Key Differentiators

| What Others Do | What We Do |
| --- | --- |
| Rule-based thresholds | ML-driven pattern recognition |
| Block & frustrate users | Risk levels with explanations |
| Batch processing | Real-time prediction (<100ms) |
| Fixed fraud patterns | GAN-generated evolving patterns |
| Binary yes/no | Probability scores + risk factors |

# Demo Highlights

| Scenario | Recipient | Result |
| --- | --- | --- |
| ✅ Send to trusted merchant | trusted.merchant@upi | LOW risk, instant approval |
| ⚠️ Send to suspicious account | suspicious.account@upi | WARNING with risk factors displayed |
| 🔴 Send to known fraud actor | fraud.actor@upi | BLOCKED with explanation |

> *"The user sees WHY their transaction is flagged — building trust, not frustration."*

# Impact & Metrics

| Model Accuracy | False Positive Rate |
| --- | --- |
| **~94%** | **<5%** |
| on test data | minimal user friction |

| Inference Time | Features Analyzed |
| --- | --- |
| **<50ms** | **22** |
| per transaction | real-time signals |

# Vision

Integrate SafePayAI as a **middleware layer** for any payment platform — protecting millions of users with AI-powered, explainable fraud detection.

# Q&A Preparation

| Potential Question | Answer |
|---|---|
| *Why Random Forest over Deep Learning?* | Interpretability + speed. We can explain which features triggered a flag, and inference is instant without GPU. |
| *How does the GAN help?* | It synthetically generates realistic fraud scenarios, solving the 1:1000 class imbalance without overfitting. |
| *What about privacy?* | All predictions happen locally on transaction metadata — no card/account numbers are stored or transmitted. |
| *How scalable is this?* | Flask + pre-loaded model can handle thousands of requests/sec. For production, we'd add Redis caching and load balancing. |

> **Closing Statement:** *"SafePayAI doesn't just detect fraud — it **anticipates** it. By combining generative AI for training and ensemble models for prediction, we've built a system that evolves as fast as the fraudsters do. The future of payment security isn't just reactive. It's predictive."*

# Research References

1. Zong, K. et al. (2025). *Detection of AI Deepfake and Fraud in Online Payments Using GAN-Based Models*. arXiv:2501.07033.
   arxiv.org/abs/2501.07033

2. Gao, Y. et al. (2023). *Advancing Financial Fraud Detection: Self-Attention Generative Adversarial Networks*. ScienceDirect.
   sciencedirect.com/science/article/abs/pii/S1544612323012151

3. Taha, A. A. & Malebary, S. J. (2023). *A Survey on GAN Techniques for Data Augmentation in Credit Card Fraud Detection*. MDPI.
   mdpi.com/2504-4990/5/1/19

4. Zhang, X. et al. (2024). *Utilizing GANs for Fraud Detection: Model Training with Synthetic Transaction Data*. arXiv:2402.09830.
   arxiv.org/abs/2402.09830

5. Chen, Y. et al. (2025). *Semi-Supervised Bayesian GANs with Log-Signatures for Uncertainty-Aware Credit Card Fraud Detection*. MDPI Mathematics. mdpi.com/2227-7390/13/19/3229

6. Kumar, R. et al. (2025). *Optimizing Credit Card Fraud Detection with Random Forests and SMOTE*. Nature Scientific Reports. nature.com/articles/s41598-025-00873-y

7. Bhattacharyya, S. et al. (2018). *Ensemble Learning for Credit Card Fraud Detection.* ACM CODS-COMAD. dl.acm.org/doi/10.1145/3152494.3156815

8. Lee, J. et al. (2024). *An Integrated Multistage Ensemble Machine Learning Model for Fraudulent Transaction Detection*. Journal of Big Data. journalofbigdata.springeropen.com