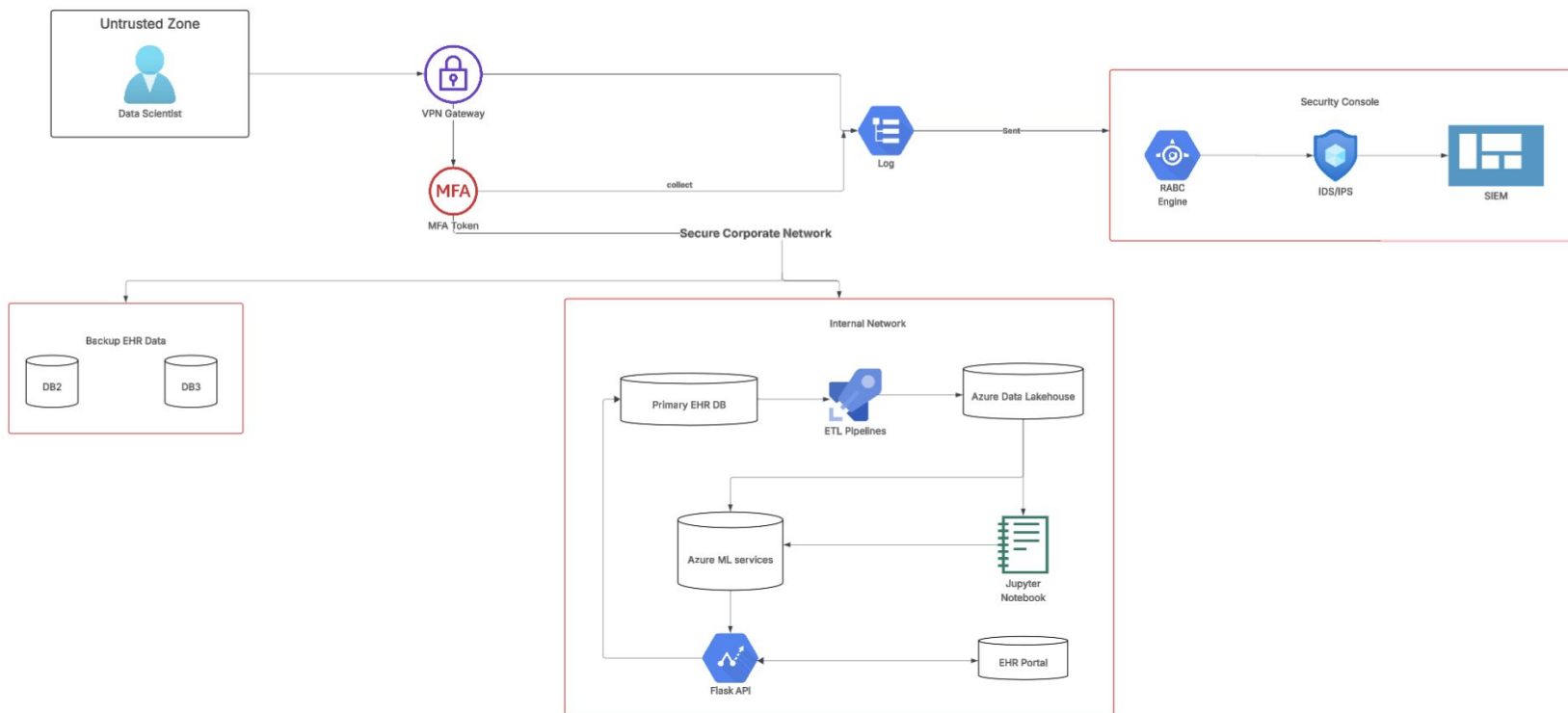# Diabetes Readmission

Group 1: Yicheng Pan, Prabhakar Pandey,
Miaoxuan Zhang, Danny Chen, Vanessa Chen
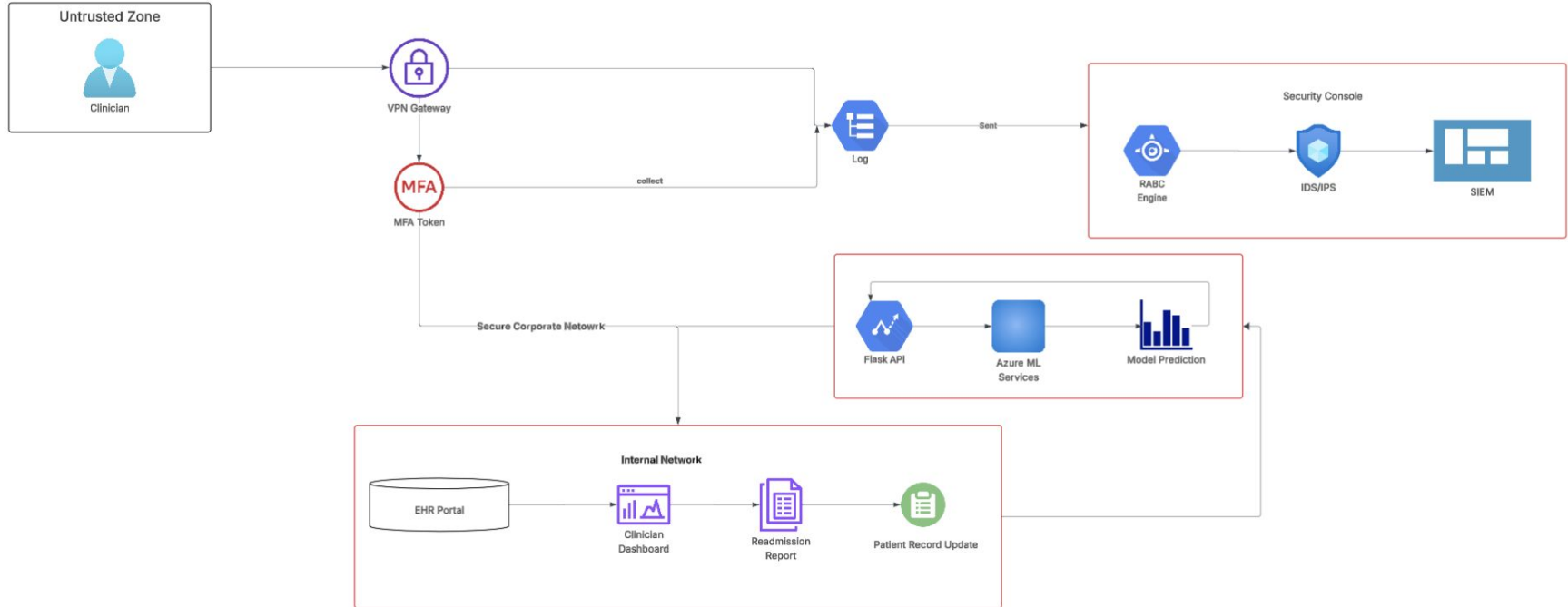
# Tools and Technologies

- Lucidchart

- Excel

- Jupyter Notebook

- Python

- Tableau

- Canva

- ChatGPT

# Network Diagram (Data Scientist)

# Network Diagram (Clinician)

# 🔒 Secure Entry: Identity Verification and Role Control

🛡️ **VPN + MFA**

Encrypted tunnel plus identity check — like a security gate with both badge and fingerprint
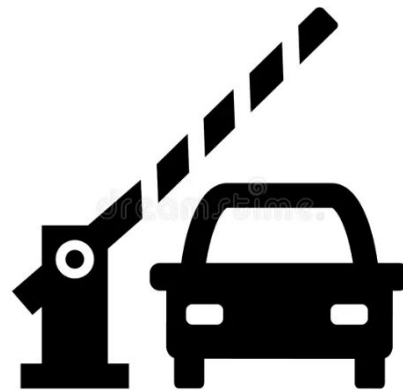
👥 **RBAC (Role-Based Access Control)**

Only lets you into rooms you're authorized for — like a hotel key card that opens only your room

📉 **Principle of Least Privilege**

Just enough access to do your job — no access to unnecessary areas, reducing blast radius

**ACCESS GRANTED**

# Surveillance Watch and Respond

📊 **SIEM (Security Information and Event Management)**

Centralizes logs from all systems for unified threat visibility

🛡️ **IDS/IPS (Intrusion Detection/Prevention Systems)**

Blocks malicious traffic like a digital gatekeeper

📝 **Audit Logging**
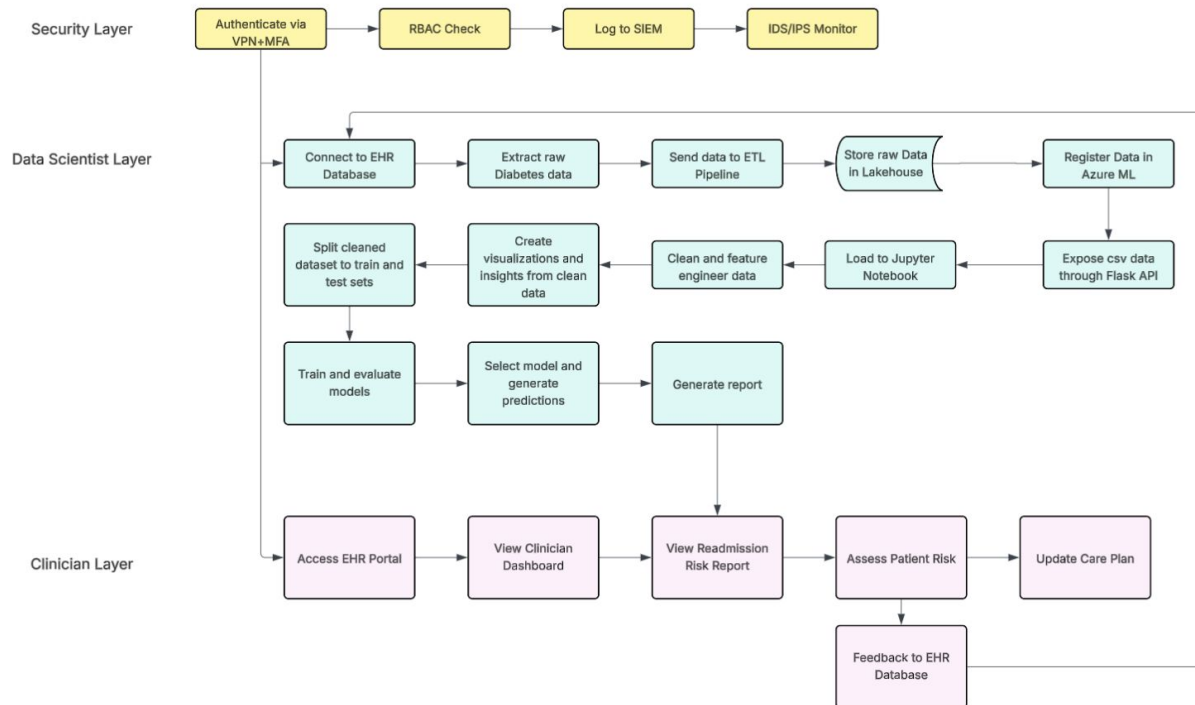
Tracks all access and activity in the environment

24 HR VIDEO SURVEILLANCE

# Following The Trails: A Three Layered Approach

- **Security Layer**
  - Protect sensitive health data

- **Data Scientist Layer**
  - Process and analyze data

- **Clinician Layer**
  - Translate analysis into care decisions

# Data Flow Diagram

# Data Retrieval

1.  Downloaded the CSV file, "<u>Diabetes 130-US Hospitals for Years 1999-2008</u>"

2.  Saved it to a folder in our local device

3.  Copied the file's path

4.  Used pd.read_csv() to access the data using the path

# Exploratory Data Analysis

1. Insights into Target Variable
2. Regarding Metrics
3. Checking for Missing Data
4. Feature Categorization
5. Bivariate Analysis
6. Correlation Between Numerical Variables

# Data Cleaning

**Cleaning**

1. Drop unnecessary columns + create new columns

2. Get rid of duplicates

3. Check and address missing data

4. Combine NO and >30 together in the readmitted column

**Regression**

1. Convert categorical features to numeric

2. Used a 70-15-15 data split percentage for training, validation, and test

3. Over-sampled the data as that led to the better model

# 9 Models Tested

| Logistic Regression | Decision Tree | Random Forest |
|---|---|---|
| KNN | Linear SVC | Gradient Boosting |
| Catboost | Stochastic Gradient Descent | XGB |

# Final Model – Random Forest

## Top 3 Models

|  | 1) Random Forest | 2) Logistic Regression | 3) Linear SVC |
|---|---|---|---|
| AUC | **0.687** | 0.686 | 0.684 |
| Accuracy | **0.634** | 0.684 | 0.882 |
| Recall | **0.641** | 0.579 | 0.109 |
| Precision | **0.177** | 0.191 | 0.373 |
| Specificity | **0.633** | 0.697 | 0.977 |
| Prevalence | **0.110** | 0.110 | 0.110 |

# Random Forest – Hyperparameter Tuning

| AUC | 0.692 |
|------|-------|
| Accuracy | 0.667 |
| Recall | 0.610 |
| Precision | 0.188 |
| Specificity | 0.674 |
| Prevalence | 0.110 |

# Top 20 Most Influential Predictive Features



Top 20 Most Influential Features (Post-Tuning)

# Visualizing Results

# Real World Implications

- Prioritize high risk patient populations

- Strengthen preventive care in outpatient settings

- Identify barriers to metformin treatment

- Promote diabetes medication optimization

- Improve comprehensive discharge planning

# Thank You!

Any Questions?