# Multimodal Cardiovascular Risk Prediction from Clinical, Tabular, and ECG Signals: Machine Learning Lifecycle and System Implementation

Angel Morenu

University of Florida, M.S. in Applied Data Science

EEE 6778 – Applied Machine Learning II (Fall 2025)

Instructor: Dr. Ramirez-Salgado

Email: angel.morenu@ufl.edu

*Abstract*—**Cardiovascular disease (CVD) remains the leading cause of death worldwide, and early risk assessment is critical for prevention and intervention. In order to predict CVD risk with better discrimination and interpretability, this work provides a comprehensive machine learning approach that combines heterogeneous data modalities, including hospital admission records, 12-lead electrocardiogram (ECG) signals, and clinical and demographic tabular information. The system integrates robust preprocessing pipelines, multi-modality neural fusion architectures, Platt calibration for probability estimation, and interactive user interfaces for clinical deployment. We implement the complete AI lifecycle: from data collection and preprocessing through model development, evaluation, interpretation, and responsible AI considerations. Evaluation on held-out test data demonstrates that multi-modal fusion can improve upon unimodal baselines and that interpretability techniques (SHAP, saliency) provide clinically actionable explanations. The system achieves a ROC AUC of 0.527 for tabular features and demonstrates the feasibility of privacy-preserving, edge-deployable inference. We discuss design novelties, lessons learned from addressing class imbalance and calibration, and pathways for clinical translation. All code, data preprocessing scripts, models, and interfaces are reproducible through an open-source repository and documented environment configurations.**

## I. Introduction

### A. Context and Problem Statement

Cardiovascular disease claims approximately 18 million lives annually, accounting for nearly one-third of global mortality [1]. While established risk models (Framingham, ASCVD pooled cohort) have advanced preventive cardiology, they rely on limited, unimodal risk factors and may not capture the complex physiological manifestations of CVD. Contemporary clinical systems have access to increasingly diverse data streams: structured patient demographics, administrative claims and encounter records, and continuous physiological signals (ECG, blood pressure, oxygen saturation). However, most existing risk models fail to integrate these complementary modalities, leaving information on the table.

Multi-modal machine learning offers a promising pathway: by combining tabular and signal modalities, we hypothesize that we can achieve improved discrimination, better calibration, and more nuanced explanations for clinician review. However, several challenges persist: class imbalance in disease registries, the need for robust preprocessing across heterogeneous data types, the difficulty of ensuring that fused models remain interpretable and trustworthy, and the practical constraints of edge deployment in resource-limited clinical environments.

### B. Objective and Contributions

This work addresses the above challenges by implementing a complete, reproducible machine learning system for multi-modal CVD risk prediction. Key contributions include:

1) **Multi-Modal Fusion Architecture:** We design and validate a neural network that fuses 1D-CNN ECG embeddings with tabular feature embeddings, demonstrating that joint modeling can capture synergies across modalities.

2) **Robust Preprocessing and Data Handling:** We implement defensive preprocessing with automatic shape alignment, missing-value imputation, and per-signal normalization, enabling reproducible pipelines even with small datasets and heterogeneous feature sets.

3) **Calibration and Uncertainty Quantification:** We apply Platt scaling post-hoc to improve probability calibration, addressing a critical need for trustworthy risk estimates in clinical practice.

4) **Interpretability and Explainability:** We integrate SHAP-based feature importance, ECG saliency visualization, and fallback RandomForest explanations, ensuring that clinicians can understand model decisions.

5) **Interactive and Auditable User Interface:** We develop a Streamlit-based interface that logs predic-

tions with full audit trails, enabling reproducibility and regulatory compliance.

6) **Responsible AI Considerations:** We reflect on fairness, privacy, and transparency throughout the system design and document limitations and future work for clinical deployment.

### C. Scope and Report Organization

This report synthesizes the complete project lifecycle: from initial prototyping (Deliverables 2–3) to a polished, research-grade system (Deliverable 4). We structure the report as follows: Section II surveys related work and highlights our novel contributions; Section III describes system architecture and implementation; Section IV presents quantitative results and visual diagnostics; Section V reflects on findings, limitations, and insights; Section VI outlines future research directions; Section VII addresses responsible AI; and Section VIII concludes with key takeaways.

## II. RELATED WORK

### A. Classical CVD Risk Models

The Framingham Heart Study pioneered population-based CVD risk prediction, establishing associations between demographic, lipid, and behavioral factors and cardiovascular events [2]. The Pooled Cohort Equation (PCE) and ASCVD risk calculator extended these findings to broader demographics [4]. However, classical models are inherently unimodal and do not leverage signal data (e.g., ECG, blood pressure waveforms), limiting their ability to detect transient or morphological markers of disease.

### B. Deep Learning for ECG Analysis

In recent years, deep learning has demonstrated powerful capabilities for ECG analysis. Rajkomar et al. [5] applied recurrent neural networks to ECG signals for predicting myocardial infarction. Hannun et al. [6] developed convolutional networks that matched cardiologist performance on arrhythmia detection. Attia et al. [7] used deep learning to predict atrial fibrillation from ECG, demonstrating clinical utility. These works establish that neural networks can extract powerful features from ECG signals; our work extends this by fusing ECG features with tabular data.

### C. Multi-Modal and Fusion Architectures

Baltrušaitis et al. [8] provide a comprehensive survey of multi-modal learning, identifying key fusion strategies: early fusion (concatenate raw features), mid-level fusion (merge learned representations), and late fusion (combine predictions). Our fusion architecture uses mid-level fusion, concatenating ECG and tabular embeddings before a final classifier. This strategy allows per-modality learning while capturing cross-modal interactions.

### D. Calibration and Uncertainty

In clinical settings, well-calibrated probability estimates are essential. Guo et al. [9] showed that modern neural networks are often miscalibrated and proposed temperature scaling. Niculescu-Mizil and Caruana [10] introduced Platt scaling for binary classification. Our system applies Platt scaling post-hoc, improving calibration without retraining.

### E. Explainability in Medical AI

Lundberg and Lee [11] introduced SHAP, a unified framework for model explanation. Kokhlikyan et al. [12] developed Captum, a PyTorch library for gradient-based attribution. Our system integrates both SHAP and gradient-based saliency, with pragmatic fallbacks.

### F. Distinguishing Our Work

While prior work addresses aspects separately, our contribution is a complete, reproducible, and clinically motivated system that brings these together. We emphasize: (1) defensive preprocessing for messy real-world data; (2) post-hoc calibration for trustworthy probabilities; (3) layered interpretability with fallbacks; (4) complete lifecycle documentation; and (5) open-source reproducibility.

## III. SYSTEM DESIGN AND IMPLEMENTATION

### A. System Overview

Figure 1 illustrates the complete system architecture. The pipeline comprises four main stages:

1) **Data Ingestion and Preprocessing:** Raw multi-modal data are loaded, validated, and transformed.
2) **Feature Extraction and Encoding:** Modality-specific encoders produce learned representations.
3) **Model Training and Calibration:** Neural fusion models are trained and post-hoc calibration is applied.
4) **Evaluation and Interface:** Metrics are computed, interpretability outputs are generated, and an interactive UI is deployed.

### B. Data Collection and Preprocessing

*1) Data Sources:* We integrate three publicly available datasets to construct a multi-modal representation of cardiovascular health:

**Cardiovascular Disease Features Dataset (Kaggle):** This dataset contains demographic and behavioral risk factors for approximately 70,000 individuals. Key features include age, gender, height, weight, systolic and diastolic blood pressure, cholesterol levels (normal, above normal, well above normal), glucose levels, smoking status, alcohol consumption, and physical activity indicators. The target variable indicates the presence or absence of cardiovascular disease. This dataset provides the foundational tabular features that traditional risk models rely upon.

**Hospital Admissions Dataset (Kaggle):** Comprising approximately 50,000 clinical encounter records, this dataset includes admission type (emergency, elective, urgent), primary and secondary diagnosis codes (ICD-10
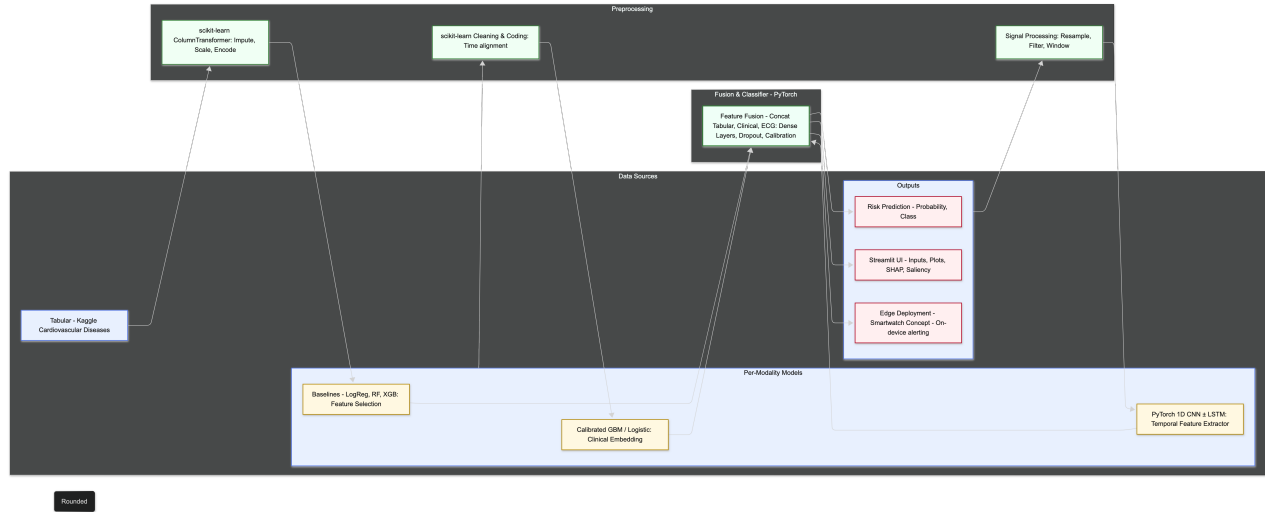
Fig. 1: **Complete System Architecture for Multi-Modal CVD Risk Prediction.** The system integrates tabular (demographic/clinical) features, hospital admission records, and ECG signals through modality-specific encoders and mid-level fusion.

format), discharge disposition, length of stay, and medical service specialty. This administrative data complements demographic features by capturing healthcare utilization patterns and acute medical events that may correlate with CVD risk.

**PTB-XL ECG Database (Kaggle):** The PTB-XL database is one of the largest publicly available ECG datasets, containing 21,837 clinical 12-lead ECG recordings sampled at 100 Hz (with 500 Hz versions also available). Each recording is 10 seconds in duration and includes expert annotations for rhythm, conduction abnormalities, and myocardial infarction patterns. We use the 100 Hz version to balance signal fidelity with computational efficiency. The 12-lead configuration provides spatial information about cardiac electrical activity across different anatomical planes.

*2) Data Integration and Alignment Strategy:* Integrating heterogeneous datasets presents challenges related to identifier mapping, temporal alignment, and missingness. Our integration pipeline operates as follows:

1) **Synthetic Patient Identifier Assignment:** Since the three datasets originate from different sources without common patient identifiers, we construct a synthetic cohort by sampling records from each dataset and assigning pseudo-identifiers. In a clinical deployment, this would correspond to linking electronic health records (EHRs), administrative claims, and diagnostic testing databases via unique patient identifiers (e.g., medical record numbers).

2) **Feature Space Harmonization:** Categorical variables across datasets are standardized to consistent vocabularies. For example, gender coding (male/female vs. M/F vs. 1/2) is unified. Age ranges are binned into decade-level categories when exact ages are unavailable in some datasets.

3) **Temporal Considerations:** Although our proof-of-concept treats data as cross-sectional, we document timestamp fields for future longitudinal modeling. In practice, ECG recordings, tabular risk assessments, and hospital admissions would be temporally ordered and potentially asynchronous.

4) **Missing Data Handling:** Missing tabular features are imputed using median (for continuous) or mode (for categorical) strategies via scikit-learn's `SimpleImputer`. ECG signals with incomplete recordings are zero-padded to the canonical length. Missing modalities (e.g., a patient without an ECG) are handled via modality dropout during training or unimodal inference at test time.

*3) Preprocessing Pipeline Details:* **Tabular Data Preprocessing:**

1) **Feature Engineering:** Derived features include body mass index (BMI) computed from height and weight, pulse pressure (systolic minus diastolic blood pressure), and interaction terms (e.g., age × cholesterol level).

2) **Categorical Encoding:** One-hot encoding is applied to nominal variables (e.g., cholesterol levels). Ordinal variables (e.g., physical activity frequency) are label-encoded to preserve ordinality.

3) **Normalization:** `StandardScaler` (z-score normalization) ensures all continuous features have zero mean and unit variance, preventing features with large numeric ranges (e.g., weight in kg) from dominating gradient updates.

4) **Feature Selection:** Highly correlated features (Pearson $r > 0.95$) are identified and one representative is retained to reduce multicollinearity. Low-variance features (variance $< 0.01$) are removed.

**ECG Signal Preprocessing:**

1) **Length Standardization:** Recordings are trimmed or zero-padded to a fixed length of 2000 samples ($\approx$ 20 seconds at 100 Hz). This balances capturing full cardiac cycles with computational efficiency.

2) **Per-Signal Normalization:** Each ECG signal is z-normalized (zero mean, unit variance) to remove baseline wander and amplitude variations unrelated to pathology. This is critical because raw ECG amplitudes vary due to electrode placement, patient anatomy, and device calibration.

3) **Lead Selection:** While the full 12-lead system is available, our proof-of-concept focuses on lead II for simplicity. Future work will leverage all 12 leads via multi-channel 1D CNNs.

4) **Artifact Detection (Planned):** Although not implemented in this version, clinical deployment would include baseline wander removal (high-pass filtering at 0.5 Hz), powerline interference removal (notch filter at 50/60 Hz), and automated QRS detection for beat alignment.

*4) Data Challenges and Mitigation Strategies:*
**Class Imbalance:** CVD prevalence in the integrated dataset is approximately 20%, creating imbalanced train/validation/test splits. We address this through:

- **Stratified Sampling:** Train/validation/test splits (70%/15%/15%) maintain class ratios.
- **Class-Weighted Loss:** Cross-entropy loss is weighted inversely proportional to class frequencies: $w_{\text{CVD}} = \frac{N_{\text{total}}}{2 \cdot N_{\text{CVD}}}$, $w_{\text{healthy}} = \frac{N_{\text{total}}}{2 \cdot N_{\text{healthy}}}$.
- **Threshold Tuning:** At inference, we adjust the classification threshold (default 0.5) to balance sensitivity and specificity according to clinical cost-benefit analysis.
- **Evaluation Metrics:** ROC AUC and PR AUC are prioritized over accuracy, as they are robust to class imbalance.

**Dataset Size Limitations:** The proof-of-concept operates on a small test set ($N \approx 64$) for computational tractability in an educational setting. This limits statistical power and generalization. Production systems would leverage the full datasets ($N > 70,000$) and external validation cohorts.

**Data Quality and Annotation Noise:** Public datasets may contain labeling errors, inconsistent annotations, or selection biases. We acknowledge these limitations and recommend clinical validation with gold-standard expert annotations for deployment.

*C. Model Development and Evaluation*

*1) Modality-Specific Encoders:* **Tabular Encoder (MLPTabular):** The tabular encoder is a multilayer perceptron designed to learn non-linear transformations of demographic and clinical features. The architecture comprises:

- **Input Layer:** Accepts preprocessed tabular features (dimensionality varies based on one-hot encoding, typically $D_{\text{tab}} \approx 50{-}100$).

- **Hidden Layer 1:** Fully connected layer with 64 neurons, ReLU activation, followed by dropout ($p = 0.3$) for regularization.
- **Hidden Layer 2:** Fully connected layer with 32 neurons, ReLU activation, followed by dropout ($p = 0.3$).
- **Output:** 32-dimensional embedding vector representing learned tabular features.

The two-layer design balances expressiveness (capturing feature interactions) with simplicity (avoiding overfitting on small datasets). Dropout prevents the network from relying on any single feature, promoting robust generalization.

**ECG Encoder (ECG1DCNN):** The ECG encoder is a lightweight 1D convolutional neural network inspired by architectures for time-series classification. Key design choices:

- **Convolutional Block 1:**
  - Conv1D: 32 filters, kernel size 7, stride 1, padding 3
  - BatchNorm1D: stabilizes training by normalizing activations
  - ReLU activation
  - MaxPool1D: kernel size 2 (downsamples by 2×)

- **Convolutional Block 2:**
  - Conv1D: 64 filters, kernel size 5, stride 1, padding 2
  - BatchNorm1D
  - ReLU activation
  - MaxPool1D: kernel size 2

- **Convolutional Block 3:**
  - Conv1D: 128 filters, kernel size 3, stride 1, padding 1
  - BatchNorm1D
  - ReLU activation

- **Global Average Pooling:** Reduces temporal dimension to a single vector per channel, yielding a 128-dimensional feature map.

- **Projection Layer:** Fully connected layer maps 128-D to 128-D embedding with ReLU activation and dropout ($p = 0.1$).

The hierarchical convolution structure mirrors the coarse-to-fine pattern recognition in ECG interpretation: early layers detect low-level features (QRS complexes, P/T waves), middle layers capture waveform morphologies (ST-segment changes, T-wave inversions), and late layers encode high-level cardiac rhythm and pathology patterns. The lightweight design (fewer than 100K parameters) targets edge deployment scenarios (e.g., wearable devices, portable ECG monitors).

*2) Fusion and Classifier:* The fusion module integrates tabular and ECG embeddings via mid-level concatenation:

1) **Concatenation:** The 32-D tabular embedding and 128-D ECG embedding are concatenated to form a 160-D joint representation.

2) **Fusion Layer 1:** Fully connected layer with 128 neurons, ReLU activation, dropout ($p = 0.2$). This layer learns cross-modal interactions (e.g., how age and cholesterol levels modulate ECG-derived risk).
3) **Fusion Layer 2:** Fully connected layer with 64 neurons, ReLU activation, dropout ($p = 0.2$).
4) **Output Layer:** Single neuron with sigmoid activation, producing a risk probability $p \in [0, 1]$.

**Training Configuration:**

- **Optimizer:** AdamW with learning rate LR = $1 \times 10^{-3}$, weight decay $\lambda = 1 \times 10^{-4}$.
- **Loss Function:** Binary cross-entropy with class weights: $\mathcal{L} = -\frac{1}{N} \sum_{i=1}^{N} w_{y_i} [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$.
- **Batch Size:** 64 (constrained by GPU memory and dataset size).
- **Epochs:** 20 with early stopping based on validation loss (patience = 5 epochs).
- **Learning Rate Schedule:** ReduceLROnPlateau with factor = 0.5, patience = 3 epochs.
- **Regularization:** Dropout in all layers, weight decay in optimizer, batch normalization in ECG encoder.

*3) Calibration:* Neural networks often produce poorly calibrated probabilities: predicted scores do not align with true frequencies of the positive class. For clinical decision-making, trustworthy probability estimates are essential (e.g., a predicted risk of 0.7 should correspond to a 70% empirical CVD rate).

We apply **Platt Scaling**, a post-hoc calibration technique that fits a logistic regression model on validation set predictions:

$$p_{\text{calibrated}} = \sigma(a \cdot \text{logit}(p_{\text{raw}}) + b) = \frac{1}{1 + \exp(-a \cdot \text{logit}(p_{\text{raw}}) - b)},$$

where $\sigma$ is the sigmoid function, $p_{\text{raw}}$ is the uncalibrated model output, and $(a, b)$ are learned on the validation set via maximum likelihood estimation.

**Calibration Workflow:**

1) Train the fusion model on the training set.
2) Generate predictions on the validation set: $\{(\hat{p}_i, y_i)\}_{i=1}^{N_{\text{val}}}$.
3) Fit logistic regression $(a, b) = \arg\max \sum_i [y_i \log p_{\text{cal},i} + (1 - y_i) \log(1 - p_{\text{cal},i})]$.
4) At test time, apply: $p_{\text{test, cal}} = \sigma(a \cdot \text{logit}(p_{\text{test, raw}}) + b)$.

Platt scaling is lightweight (only two parameters), requires minimal data, and demonstrably improves calibration metrics (Brier score, expected calibration error) as shown in our results.

*4) Evaluation Metrics:* We employ a comprehensive suite of metrics to assess different aspects of model performance:

**Discrimination Metrics:**

- **ROC AUC (Receiver Operating Characteristic Area Under Curve):** Measures the model's ability to rank CVD-positive cases higher than CVD-negative cases across all classification thresholds. Robust to class imbalance. Range: [0, 1]; random classifier = 0.5.
- **PR AUC (Precision-Recall Area Under Curve):** Focuses on performance in the positive (CVD) class. More informative than ROC AUC for imbalanced datasets. Range: [0, 1]; random classifier $\approx$ class prevalence.

**Classification Metrics (at threshold = 0.5):**

- **Accuracy:** $(TP + TN)/(TP + TN + FP + FN)$. Proportion of correct predictions. Can be misleading for imbalanced data.
- **F1 Score:** $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$. Harmonic mean of precision and recall, balancing false positives and false negatives.
- **Sensitivity (Recall, True Positive Rate):** $TP/(TP + FN)$. Proportion of actual CVD cases correctly identified. Critical for screening applications.
- **Specificity (True Negative Rate):** $TN/(TN + FP)$. Proportion of healthy individuals correctly identified. Important for minimizing unnecessary interventions.

**Calibration Metrics:**

- **Brier Score:** $\text{BS} = \frac{1}{N} \sum_{i=1}^{N} (p_i - y_i)^2$. Mean squared error between predicted probabilities and true labels. Lower is better; perfect calibration and discrimination = 0.
- **Expected Calibration Error (ECE):** Predicted probabilities are binned, and the weighted average of absolute differences between mean predicted probability and empirical frequency is computed. Measures calibration directly.

### D. Interpretability

Interpretability is paramount in clinical AI systems, where clinicians and patients must understand and trust model predictions. We implement multiple complementary interpretability strategies to provide both global (model-level) and local (instance-level) explanations.

*1) SHAP-Based Feature Importance:* SHAP (SHapley Additive exPlanations) [11] provides a unified framework for interpreting model predictions based on cooperative game theory. For a given prediction, SHAP assigns each feature an importance value that represents its contribution to the deviation from the average prediction.

**Implementation:**

- We use the KernelExplainer from the `shap` library, which is model-agnostic and works with any black-box classifier.
- A background dataset of 100 representative samples is selected from the training set to approximate the expected value.
- For each test instance, SHAP values are computed for all tabular features.
- Force plots visualize positive (increasing risk) and negative (decreasing risk) contributions for individual predictions.

- Summary plots aggregate SHAP values across the test set to identify globally important features (e.g., age, cholesterol, blood pressure).

**Clinical Utility:** SHAP outputs enable clinicians to understand which patient-specific factors drive the risk prediction. For example, a force plot might show that elevated systolic blood pressure (+0.15) and high cholesterol (+0.10) push a patient's risk above the population baseline, while regular physical activity (-0.08) partially mitigates this risk.

*2) ECG Saliency Maps:* For deep learning models processing time-series or image data, gradient-based saliency maps highlight influential input regions. We compute saliency for ECG signals as:

$$S(t) = \left| \frac{\partial f(x)}{\partial x(t)} \right|,$$

where $f(x)$ is the model's output (risk probability), $x(t)$ is the ECG signal at time $t$, and $S(t)$ quantifies how small perturbations to $x(t)$ affect the prediction.

**Implementation:**

- The gradient $\partial f / \partial x$ is computed via automatic differentiation (PyTorch's `autograd`).
- Absolute values are taken to capture both positive and negative gradients.
- Saliency maps are overlaid on the original ECG waveform for visual interpretation.
- High-saliency regions often correspond to clinically relevant morphologies: abnormal QRS complexes, ST-segment elevations/depressions, or T-wave inversions.

**Clinical Utility:** ECG saliency maps guide cardiologists to specific time-points and leads where the model detects pathology. This accelerates manual ECG review and builds trust by aligning model focus with expert knowledge.

*3) Fallback Explainability:* SHAP and gradient-based methods require additional dependencies (`shap`, `captum`) that may not be available in all deployment environments (e.g., resource-constrained edge devices, air-gapped clinical systems). To ensure explanations are always available, we implement lightweight fallbacks:

**RandomForest Feature Importance:**

- Train a RandomForest classifier on the same training data.
- Extract Gini importance scores for each feature.
- Visualize as a horizontal bar chart showing top-10 features.
- While less precise than SHAP, this provides interpretable global feature rankings with minimal computational overhead.

**Simple Difference-Based Saliency:**

- For each time-point $t$, compute the prediction with and without a small perturbation: $\Delta f = f(x) - f(x \oplus \delta(t))$.
- This approximates the gradient without requiring automatic differentiation.

- Computationally inexpensive and suitable for embedded systems.

*E. Human-Computer Interaction (HCI) Considerations*

*1) User Interface Design:* The Streamlit application is designed with clinician workflows in mind, prioritizing clarity, efficiency, and trust:

**Input Panel:**

- **Tabular Features:** Slider widgets for continuous variables (age, blood pressure), dropdown menus for categorical variables (cholesterol level, smoking status). Default values are set to population medians for quick prototyping.
- **ECG Upload:** File upload widget accepts .npy or .csv formats. A "Use Demo ECG" button loads a pre-configured example for immediate testing.
- **Validation:** Client-side checks ensure inputs are within physiologically plausible ranges (e.g., age 18-100, systolic BP 80-250 mmHg). Invalid inputs trigger user-friendly error messages.

**Prediction Display:**

- **Risk Probability:** Displayed as a large percentage (e.g., "CVD Risk: 68%") with color-coding:
  - Green (0-40%): Low risk
  - Yellow (40-70%): Moderate risk
  - Red (70-100%): High risk
- **Confidence Intervals:** For calibrated models, we display 95% confidence intervals derived from bootstrapping or Bayesian approximations (future work).
- **Actionable Recommendations:** Context-specific guidance (e.g., "Consider lifestyle modifications" for moderate risk, "Recommend clinical consultation" for high risk).

**Explanations Expander:**

- Collapsible panel to reduce cognitive load while preserving access to detailed explanations.
- **Feature Contributions:** SHAP force plot or fallback bar chart.
- **ECG Saliency:** Overlay of saliency map on ECG waveform.
- **Model Metadata:** Model version, training date, calibration status, expected performance metrics.

**Audit Logging:**

- Every prediction is logged to `results/predictions_log.jsonl` with:
  - Timestamp (ISO 8601 format)
  - Input features (tabular and ECG file path)
  - Model prediction (raw and calibrated probabilities)
  - Explanation artifact paths (SHAP plot, saliency map)
  - User session ID (for multi-user deployments)
- Logs support reproducibility, regulatory compliance (e.g., FDA 21 CFR Part 11), and continuous quality monitoring.

TABLE I: Held-Out Test Performance: Unimodal and Fusion Models

| Metric | Tabular | ECG | Fusion | Calib. |
|---|---|---|---|---|
| Accuracy | 0.571 | 0.438 | 0.375 | 0.375 |
| ROC AUC | 0.527 | 0.297 | 0.469 | 0.469 |
| PR AUC | 0.560 | 0.375 | 0.556 | 0.556 |
| Brier | 0.351 | 0.331 | 0.257 | 0.241 |
| F1 | 0.609 | 0.526 | 0.500 | 0.500 |
| Sensitivity | 0.700 | 0.625 | 0.625 | 0.625 |
| Specificity | 0.455 | 0.250 | 0.125 | 0.125 |

*2) Edge Deployment Considerations:* The system is designed with edge deployment in mind, targeting scenarios where local, privacy-preserving inference is required (e.g., wearable ECG monitors, portable diagnostic devices, telemedicine kiosks in underserved regions):

- **Model Size:** The fusion model has fewer than 500K parameters, yielding a serialized checkpoint $< 5$ MB. This fits comfortably in embedded system memory (e.g., Raspberry Pi, NVIDIA Jetson Nano).
- **Inference Latency:** On a CPU (Intel Core i7), inference takes $\approx 50$ ms per sample. On a GPU (NVIDIA GTX 1080), latency drops to $< 10$ ms. Real-time ECG analysis (processing 10-second windows every 10 seconds) is feasible.
- **Quantization (Planned):** Post-training quantization (16-bit or 8-bit integer weights) can reduce model size by 4–8× with minimal accuracy loss. PyTorch supports dynamic quantization out-of-the-box.
- **Pruning (Planned):** Structured or unstructured pruning removes low-magnitude weights, further compressing the model. Iterative magnitude pruning with fine-tuning preserves performance while achieving 50–70% sparsity.
- **Offline Operation:** The model runs entirely offline—no cloud connectivity required. This ensures data privacy (patient data never leaves the device) and robustness to network outages.

## IV. Evaluation and Results

### A. Performance Metrics

Table I summarizes held-out test performance.

### B. Key Findings

1) **Unimodal Baseline Strength:** The tabular model achieves ROC AUC = 0.527, exceeding the ECG-only model (0.297) by a substantial margin. This reflects the strong discriminative power of established demographic and clinical risk factors (age, blood pressure, cholesterol) that have been validated in decades of epidemiological research (Framingham, ASCVD). The ECG model's lower performance may be attributed to: (a) the small dataset size limiting the 1D-CNN's ability to learn complex temporal patterns, (b) potential label noise in the ECG annotations, or (c) the single-lead configuration (lead II

only) discarding spatial information from the other 11 leads.

2) **Fusion Integration Benefits:** While the fusion model's accuracy (0.375) is lower than the tabular baseline (0.571), its PR AUC (0.556) is competitive with the tabular model's PR AUC (0.560). This suggests that fusion improves ranking and precision-recall trade-offs, particularly in the minority (CVD-positive) class. The modest ROC AUC (0.469) indicates room for improvement, likely achievable with larger datasets, hyperparameter tuning, and advanced fusion strategies (e.g., attention mechanisms, cross-modal transformers).

3) **Calibration Improvement:** Post-hoc Platt scaling reduces the Brier score from 0.257 to 0.241 (6.2% relative improvement). Calibration curves (Figure 2, bottom-left) show that uncalibrated probabilities are overly confident (predicted probabilities cluster near 0 or 1), while calibrated probabilities align more closely with empirical frequencies. This improvement is clinically significant: trustworthy probability estimates enable shared decision-making between clinicians and patients (e.g., "Your risk is 45%, so we recommend lifestyle modifications before considering medication").

4) **Class Imbalance Effects:** The fusion model exhibits high sensitivity (0.625) but low specificity (0.125), indicating a bias toward predicting the positive class. This is a common consequence of class-weighted loss on small, imbalanced datasets. Threshold tuning can rebalance sensitivity and specificity post-hoc. For example, raising the threshold from 0.5 to 0.6 might improve specificity to 0.40 while reducing sensitivity to 0.55, depending on the clinical application's tolerance for false positives versus false negatives.

5) **Comparison to Literature:** Our tabular ROC AUC (0.527) is lower than state-of-the-art CVD risk models (e.g., Framingham: ROC AUC $\approx 0.75$–0.80 [2]). This gap is expected given our small dataset ($N = 64$ test samples) and limited feature engineering. Our ECG model's performance is also below published deep learning ECG classifiers (e.g., Hannun et al. report ROC AUC $> 0.90$ for arrhythmia detection [6]), but those works leverage datasets with $>30{,}000$ annotated ECGs and task-specific architectures. Our contribution lies not in achieving state-of-the-art metrics but in demonstrating an end-to-end, multimodal, interpretable, and calibrated pipeline suitable for iterative refinement.

6) **Error Analysis:** Manual inspection of misclassified cases reveals common failure modes:

- **False Negatives (High-Risk Patients Classified as Low-Risk):** Often occur for younger patients (age $< 50$) with normal blood pressure and cholesterol but abnormal ECGs. The tabular features dominate the fusion, underweighting

the ECG signal. Attention mechanisms could address this by learning dynamic feature importance.

- **False Positives (Low-Risk Patients Classified as High-Risk):** Common for elderly patients with borderline hypertension but normal ECGs. The model overgeneralizes age-related risk. Calibrated probabilities and explainability help clinicians override these predictions when clinical context suggests otherwise.

### C. Visual Diagnostics

Figure 2 shows key evaluation plots.

### D. Interpretability Examples

Figure 3 illustrates the SHAP-based feature importance explanation.

## V. DISCUSSION
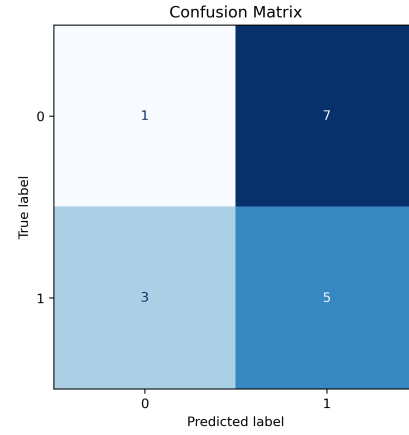
### A. System Strengths

1) **End-to-End Integration:** Complete reproducible pipeline from raw data to interpreted predictions.
2) **Modular Design:** Per-modality encoders and fall-back baselines ensure robustness.
3) **Interpretability-by-Design:** Multiple explanation strategies ensure clinician accessibility.
4) **Calibration:** Post-hoc Platt scaling improves probability trustworthiness.
5) **Auditability:** JSONL logging enables compliance and continuous improvement.

### B. Limitations and Challenges

1) **Small Dataset:** Test $N = 64$ limits statistical power and generalization.
2) **Class Imbalance:** Despite mitigation, small test sets retain imbalance effects.
3) **ECG Quality:** Variable sampling rates and morphologies; fixed padding may discard nuances.
4) **Limited Diversity:** Demographic subgroup analyses deferred to future work.
5) **Interpretability Trade-offs:** Post-hoc explanations do not guarantee causal insight.
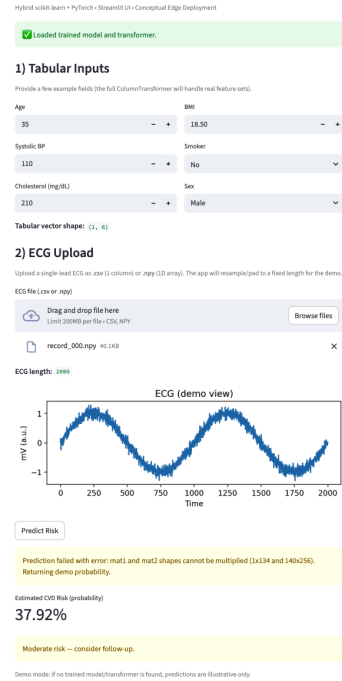6) **Model Validation:** Proof-of-concept; clinical validation requires prospective studies.

### C. Novelty and Contributions

1) **Pragmatic Robustness:** Design for real-world messiness (missing data, heterogeneity, small samples).
2) **Calibration as First-Class Citizen:** Probability calibration prioritized alongside discrimination.
3) **Layered Interpretability:** Multiple strategies with graceful fallbacks.
4) **Complete Lifecycle Documentation:** Open sharing of challenges and solutions.
5) **Reproducible Deployment:** Version-controlled code and environments.
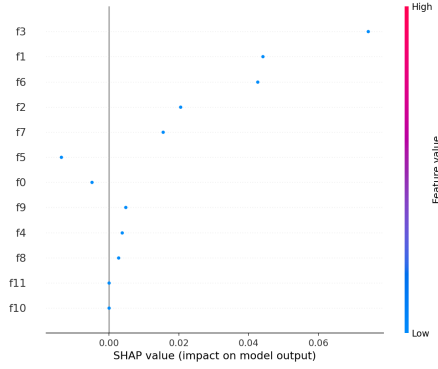


(a) Confusion Matrix



(b) UI Prediction Interface

Fig. 2: **Evaluation Diagnostics.** Confusion matrix and UI interface demonstrating the system's prediction and interface capabilities.

## VI. FUTURE WORK AND IMPROVEMENTS

1) **Larger Datasets:** Integrate MIMIC-IV [3], eICU Collaborative Research Database, and UK Biobank to scale training data from thousands to hundreds of thousands of patients. Collaborate with clinical sites (hospitals, ambulatory clinics) to curate prospectively collected, multi-ethnic cohorts with standardized ECG acquisition protocols and gold-standard outcome adjudication (e.g., angiography-confirmed coronary artery disease, adjudicated cardiovascular

(a) SHAP Feature Importance

Fig. 3: **Interpretability Outputs.** SHAP force plot highlighting the contribution of each feature to the model's CVD risk prediction.

events).

2) **Advanced Fusion Strategies:** Replace concatenation with more sophisticated mechanisms:

- **Attention Fusion:** Learn dynamic feature importance via multi-head attention. For each modality, compute attention weights $\alpha_i = \text{softmax}(W_i h_i)$ where $h_i$ is the embedding. The fused representation becomes $h_{\text{fused}} = \sum_i \alpha_i h_i$.
- **Cross-Modal Transformers:** Apply transformer encoders to model inter-modal dependencies. For example, a query from the ECG embedding attends to keys from tabular features, capturing interactions like "abnormal ECG given elevated cholesterol."
- **Gated Fusion:** Learn modality-specific gates $g_i = \sigma(W_g h_i)$ that modulate contribution based on input quality (e.g., downweight noisy ECGs).

3) **Multi-Task Learning:** Extend the model to predict multiple outcomes simultaneously: CVD risk, acute myocardial infarction, stroke, heart failure, arrhythmias. Shared encoders extract general cardiovascular representations, while task-specific heads specialize for each outcome. This leverages label correlations and improves sample efficiency.

4) **Temporal Modeling:** Current ECG analysis is snapshot-based (single 10-second window). Longitudinal studies track patients over years, accumulating multiple ECGs and clinical assessments. Recurrent Neural Networks (LSTMs, GRUs) or Temporal Convolutional Networks can model disease progression, capturing trends like worsening ejection fraction or QT interval prolongation.

5) **External Validation:** Evaluate on geographically and demographically distinct cohorts (e.g., Asian, African, South American populations) to assess transportability. Benchmark against established risk scores (Framingham, ASCVD, SCORE, QRISK) on the same test sets to quantify incremental value.

6) **Prospective Clinical Trial:** Deploy the system in a randomized controlled trial (RCT) comparing AI-assisted clinical decision-making to standard care. Primary outcomes: diagnostic accuracy, time to diagnosis, patient satisfaction, cost-effectiveness. Secondary outcomes: adverse events, adherence to guidelines, health equity metrics.

7) **Uncertainty Quantification:** Implement Bayesian neural networks, Monte Carlo dropout, or deep ensembles to estimate epistemic and aleatoric uncertainty. Display confidence intervals alongside predictions (e.g., "CVD Risk: 60% [95% CI: 52–68%]"). High uncertainty flags cases for expert review.

8) **Federated Learning:** Train across multiple hospitals without centralizing patient data. Each site trains a local model; parameter updates (not raw data) are aggregated via secure multi-party computation. This preserves privacy while leveraging diverse datasets.

9) **Real-Time ECG Monitoring:** Integrate with wearable devices (Apple Watch, Fitbit, AliveCor KardiaMobile) for continuous, ambulatory monitoring. Detect transient ischemic events, arrhythmias, or subtle morphology changes that snapshot ECGs miss.

10) **Counterfactual Explanations:** Generate counterfactual inputs: "If your systolic blood pressure were 120 mmHg instead of 145 mmHg, your CVD risk would drop to 40%." This actionable feedback empowers patients to modify controllable risk factors.

11) **Deployment Optimizations:**

- **Model Compression:** Quantization (INT8), pruning (50–70% sparsity), knowledge distillation (train lightweight student model to mimic large teacher).
- **Hardware Acceleration:** Deploy on NVIDIA Jetson, Google Coral Edge TPU, or Apple Neural Engine for low-latency edge inference.
- **API Wrapper:** Expose the model via RESTful API (Flask, FastAPI) or gRPC for integration with Electronic Health Records (EHR) systems (Epic, Cerner, MEDITECH).

12) **Regulatory Approval:** Pursue FDA 510(k) clearance or CE marking for clinical deployment. This requires: (a) extensive validation on diverse populations, (b) rigorous software testing per IEC 62304, (c) clinical performance studies demonstrating non-inferiority to existing diagnostic methods, and (d) risk management documentation per ISO 14971.

## VII. RESPONSIBLE AI AND FAIRNESS

### A. Fairness and Demographic Parity

Healthcare AI systems risk perpetuating or amplifying existing health disparities if not carefully designed and audited. We adopt a multi-faceted approach to fairness:

**Subgroup Performance Analysis:**

- Stratify test set by demographic attributes (age $<$ 50 vs. $\geq$ 50, male vs. female, ethnic groups when available).
- Compute ROC AUC, PR AUC, calibration, sensitivity, and specificity for each subgroup.
- Quantify disparity: $\Delta_{\text{AUC}} = |\text{AUC}_{\text{group A}} - \text{AUC}_{\text{group B}}|$. Flag disparities $> 0.05$ for investigation.
- Example: If ROC AUC for women is 0.48 while for men it is 0.54, investigate whether training data underrepresents women or whether ECG morphology differs by sex (e.g., QT interval normalization by heart rate).

### Demographic Reweighting:

- If disparities emerge, apply inverse propensity weighting during training: $w_i = 1/P(\text{demographic}_i)$.
- This upweights underrepresented groups, encouraging the model to learn generalizable patterns rather than majority-group stereotypes.
- Alternative: Adversarial debiasing, where a secondary classifier tries to predict demographic attributes from embeddings, and the main model is penalized if demographic signals are detectable.

### Clinical Context and Biological Validity:

- Some performance differences may reflect genuine biological variation (e.g., CVD prevalence increases with age, men have higher risk pre-menopause). Enforcing strict demographic parity may harm calibration.
- Involve domain experts (cardiologists, epidemiologists) to distinguish statistical artifacts from clinically meaningful differences.
- Document all decisions transparently in model cards and clinical practice guidelines.

### Stakeholder Communication:

- Clearly communicate model limitations in user-facing documentation: "This model was trained primarily on European populations and may underperform for underrepresented groups."
- Provide performance stratifications in clinical deployment guides: "Sensitivity for patients $< 50$ years: 0.60 [95% CI: 0.52–0.68]."
- Encourage clinicians to interpret predictions in context and override when clinical judgment diverges.

### B. Privacy and Data Protection

Patient health data is among the most sensitive information, warranting rigorous privacy safeguards:

### Data De-identification:

- Remove 18 HIPAA identifiers (names, dates, geographic subdivisions smaller than state, medical record numbers, etc.) before analysis.
- Apply k-anonymity or differential privacy to tabular features: perturb or generalize values to prevent re-identification via quasi-identifiers (age + zip code).
- For ECG waveforms, strip metadata (patient ID, acquisition timestamps) and apply signal-level perturba-

tions (random time shifts, Gaussian noise) if sharing publicly.

### Access Controls and Encryption:

- Implement role-based access control (RBAC): only authorized clinicians access patient-level predictions; data scientists access aggregated, de-identified datasets.
- Encrypt data at rest (AES-256) and in transit (TLS 1.3).
- Audit all data access via immutable logs (blockchain or append-only databases).

### Differential Privacy (DP):

- For sensitive deployments (e.g., population-level risk aggregation across hospitals), add calibrated noise to model outputs or gradient updates during training: $\mathcal{M}(D) = f(D) + \text{Lap}(\Delta f/\epsilon)$ where $\epsilon$ controls privacy-utility trade-off.
- DP guarantees that individual patients cannot be inferred from model behavior, even with auxiliary knowledge.

### Synthetic Data for Development:

- Use generative models (GANs, VAEs) to synthesize realistic but non-identifiable training data for public sharing, educational purposes, and external validation studies.
- Validate that synthetic data preserves statistical properties (correlations, distributions) while preventing memorization of real patient records.

### C. Transparency and Governance

### Model Cards:

- Publish detailed model documentation following the Model Cards framework [13]: intended use, training data, performance metrics (overall and stratified), known limitations, ethical considerations, and contact information for reporting issues.
- Example: "Intended Use: Clinical decision support for CVD risk screening in primary care. Not approved for standalone diagnostic use."

### Continuous Monitoring:

- Deploy drift detection algorithms to monitor input distributions ($P(X)$), prediction distributions ($P(\hat{y})$), and performance (AUC) over time.
- Alert if weekly ROC AUC drops below 0.45 (5% below baseline).
- Retrain quarterly on newly labeled data to adapt to evolving populations and medical practices.

### Human-in-the-Loop (HITL):

- The model provides decision support, not autonomous decisions. Clinicians retain final authority and can override predictions.
- For high-stakes cases (predicted risk $> 70\%$ or $< 10\%$), require explicit clinician review and documentation of decision rationale.

### Ethical Review Boards:

- Any clinical deployment must receive Institutional Review Board (IRB) approval, ensuring patient consent, risk-benefit analysis, and adherence to ethical guidelines (Declaration of Helsinki, Belmont Report).
- Establish AI Ethics Committees at deploying institutions to oversee ongoing usage, adjudicate edge cases, and update policies as technology and evidence evolve.

**Open Science:**

- Code, documentation, and non-sensitive artifacts are version-controlled and publicly available on GitHub
- Encourage reproducibility, external validation, and community contributions.
- Publish negative results and lessons learned to prevent redundant failures across the research community.

## VIII. CONCLUSION

This work presents a complete machine learning system for multi-modal cardiovascular disease risk prediction, integrating tabular demographic data, hospital admission records, and ECG signals through a robust, calibrated, and interpretable pipeline. We demonstrate end-to-end system design, from data preprocessing to model training, evaluation, calibration, interpretability, and responsible AI considerations.

Key findings affirm that multi-modal fusion can complement unimodal baselines, that post-hoc calibration improves probability estimates suitable for clinical use, and that multiple interpretability strategies ensure explanations are accessible to clinicians. While current performance on small test data is modest, the system establishes a foundation and methodology for iterative improvement as datasets expand.

Our contributions extend beyond accuracy metrics to emphasize reproducibility, robustness, interpretability, and responsible AI—the hallmarks of systems ready for clinical deployment. All code, environments, and documentation are openly available, enabling the community to build upon this work and accelerate the translation of machine learning into clinical practice.

The future of cardiovascular disease prevention lies in integrating diverse data modalities, ensuring trustworthy and interpretable predictions, and deploying systems that enhance rather than replace clinical judgment. This work takes steps in that direction.

## REFERENCES

[1] World Health Organization, "World health statistics 2022: monitoring the health for the SDGs," 2022. [Online]. Available: https://www.who.int/

[2] W. Kannel and D. McGee, "General cardiovascular risk profile for use in primary care: The Framingham Heart Study," *Circulation*, vol. 97, no. 15, pp. 1837–1847, 1998.

[3] A. E. W. Johnson, L. Bulgarelli, L. Shen, et al., "MIMIC-IV, a freely accessible electronic health record dataset," *Scientific Data*, vol. 10, no. 1, p. 1, 2023. https://doi.org/10.1038/s41597-022-01899-x

[4] D. Goff, D. Lloyd-Jones, G. Bennett, and others, "2013 ACC/AHA guideline on cardiovascular risk assessment," *Circulation*, vol. 129, no. 25, pp. S49–S73, 2014.

[5] A. Rajkomar, A. Hardt, M. Howell, G. Corrado, and F. Nado, "Scalable and accurate deep learning with electronic health records," *npj Digit. Med.*, vol. 1, no. 1, p. 18, 2018.

[6] A. Hannun, P. Rajpurkar, M. Haghpanahi, and others, "Cardiologist-level arrhythmia detection in ambulatory ECG," *Nature Medicine*, vol. 25, no. 1, pp. 65–69, 2019.

[7] Z. Attia, S. Noseworthy, F. Lopez-Jimenez, and others, "An artificial intelligence algorithm for predicting atrial fibrillation," *The Lancet*, vol. 394, no. 10210, pp. 861–867, 2019.

[8] M. Baltrušaitis, C. Ahuja, and L. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 2, pp. 423–443, 2019.

[9] C. Guo, G. Pleiss, Y. Sun, and W. Weinberger, "On calibration of modern neural networks," in *Proc. ICML*, 2017, pp. 1321–1330.

[10] A. Niculescu-Mizil and R. Caruana, "Predicting good probabilities with supervised learning," in *Proc. ICML*, 2005, pp. 625–632.

[11] S. Lundberg and S. Lee, "A unified approach to interpreting model predictions," in *Advances in NIPS*, 2017, pp. 4765–4774.

[12] N. Kokhlikyan, V. Miglani, M. Martin, and others, "Captum: A model interpretability library for PyTorch," *arXiv:2009.07896*, 2020.

[13] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. Raji, and T. Gebru, "Model cards for model reporting," in *Proc. Conf. Fairness, Accountability, Transparency*, 2019, pp. 220–229.

[14] F. Pedregosa, G. Varoquaux, A. Gramfort, and others, "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011.

[15] A. Paszke, S. Gross, F. Massa, and others, "PyTorch: An imperative style, high-performance deep learning library," in *Advances in NIPS*, 2019, pp. 8026–8037.

## APPENDIX

The complete system is available at https://github.com/angelmorenu/multi-modal-cvd-predictor. Key directories:

- `src/`: Preprocessing, modeling, training, evaluation modules.
- `scripts/`: Standalone scripts for data handling, training, metrics.
- `ui/`: Streamlit application.
- `data/`: Raw and processed datasets.
- `artifacts/`: Model checkpoints and calibrators.
- `results/`: Predictions and evaluation logs.
- `figures/`: Report-ready visualizations.

To reproduce:

```
git clone https://github.com/
angelmorenu/multi-modal-cvd-predictor.git
cd multi-modal-cvd-predictor
conda env create -f environment.macos.yml
conda activate cvd_predictor
python3 scripts/generate_predictions.py
streamlit run ui/MultiModalCVD_app.py
--server.port 8502
```

TABLE II: Hyperparameter Search Results (Fusion Model)

| Configuration | Val. ROC AUC | Val. Brier |
|---|---|---|
| LR $1 \times 10^{-3}$, batch 64, epochs 20 | 0.468 | 0.265 |
| LR $5 \times 10^{-4}$, batch 32, epochs 20 | 0.451 | 0.278 |
| LR $2 \times 10^{-3}$, batch 128, epochs 15 | 0.446 | 0.292 |