# A Comparative Study of Multi-source Data Fusion Approaches for Credit Default Early Warning

*Jiahui Han[1], Guanghe Cao[1,2]*

[1] *Master of Finance, MIT Sloan School of Management, MA, USA*
[1,2] *Computer Science, University of Southern California, CA, USA*

**A b s t r a c t**

This study presents a comparative analysis of multi-source data fusion approaches for early warning of credit defaults in financial institutions. The research integrates heterogeneous data sources, including credit bureau records, transaction behavior patterns, textual financial reports, and macroeconomic indicators. Three fusion strategies—early fusion, late fusion, and hybrid fusion—are systematically evaluated using ensemble machine learning algorithms, including XGBoost, LightGBM, and Random Forest. Experimental results on a real-world dataset comprising 125,847 credit records demonstrate that the hybrid fusion approach achieves the highest predictive performance with an AUC-ROC of 0.8934, outperforming the best single-source credit-bureau model (AUC-ROC 0.8234) by 7.0 percentage points (8.5% relative improvement). Feature importance analysis using SHAP values indicates that transaction behavior features account for 34.2% of the prediction, whereas NLP-extracted sentiment scores from financial texts account for 18.6%. Statistical tests (e.g., DeLong's test and bootstrap confidence intervals) indicate that the hybrid fusion configuration significantly outperforms the early-fusion baseline ($p < 0.001$ for AUC).

## 1. Introduction

### 1.1 Research Background and Motivation

#### 1.1.1 Growing Importance of Credit Risk Management in Financial Institutions

The global financial landscape has witnessed unprecedented growth in consumer lending activities over the past decade. Outstanding consumer credit in the United States reached $5.02 trillion by the end of 2024, representing a 4.7% annual increase from previous fiscal periods. This expansion has intensified the need for sophisticated risk assessment mechanisms capable of identifying potential defaults before they materialize into financial losses. Recent studies highlight the increasing exposure and complexity of consumer and enterprise credit markets, motivating the development of more accurate and robust early-warning models for defaults. In particular, multimodal learning that integrates structured financial variables with textual signals has demonstrated strong potential to improve credit risk prediction performance **Error! Reference source not found.**.

#### 1.1.2 Regulatory Requirements from the Federal Reserve and the Financial Stability Oversight Council

The Dodd-Frank Act established the Financial Stability Oversight Council with explicit mandates to monitor systemic risk. The Federal Reserve's Comprehensive Capital Analysis and Review program requires large banks to demonstrate robust stress testing capabilities. Regulatory guidance from the Office of the Comptroller of the Currency addresses the use of alternative data sources in credit underwriting, underscoring the need for financial institutions to validate predictive models using diverse data [1].

## 1.2 Problem Statement and Research Gaps

### 1.2.1 Limitations of Traditional Credit Scoring with Single-source Data

Conventional credit scoring relies predominantly on historical repayment records from credit bureaus. In practical lending and supervisory settings, credit risk models are expected to be stable under distribution shifts and class imbalance. Recent work, therefore, emphasizes handling imbalanced datasets and improving model robustness through resampling and ensemble strategies [3].

### 1.2.2 Challenges in Integrating Heterogeneous Financial Data Sources

Alternative data sources present both opportunities and technical challenges. Transaction-level data, social media activity, and textual information from financial disclosures offer complementary perspectives on creditworthiness. Integrating these heterogeneous streams requires addressing incompatibilities in data formats and feature representations[3].

### 1.2.3 Need for Improved Early Warning Mechanisms

Existing credit monitoring systems predominantly operate reactively. The economic cost of late default detection extends beyond direct losses to include collection expenses and reputational damage. Proactive early-warning mechanisms that identify deteriorating profiles 60-90 days before defaults could substantially reduce these costs [4].

## 1.3 Research Objectives and Contributions

### 1.3.1 Comparative Analysis Objectives and Scope

This research establishes a framework for evaluating multi-source data fusion approaches in credit default prediction. The primary objectives include a comparative assessment of early, late, and hybrid fusion strategies, performance benchmarking across ensemble learning algorithms, and interpretability analysis via SHAP-based feature-importance quantification.

### 1.3.2 Key Contributions of This Study

The contributions include a comprehensive taxonomy of data fusion strategies for credit risk applications, empirical evidence on the effectiveness of fusion approaches, a reproducible experimental framework, and practical recommendations for financial institutions.

## 2. Related Work

## 2.1 Machine Learning Approaches for Credit Default Prediction

### 2.1.1 Traditional Machine Learning Algorithms (SVM, Random Forest, Logistic Regression)

Machine learning applications to credit risk have evolved substantially since logistic regression-based scorecards in the 1980s. Logistic regression remains prevalent in production credit scoring environments due to its interpretability, regulatory acceptance, and computational efficiency. The log-odds transformation provides intuitive probability estimates that align with risk-based pricing frameworks. Support Vector Machines gained attention in the 2000s due to improved classification accuracy on benchmark datasets. Random Forest algorithms achieve robust performance through ensemble aggregation of decision tree predictions while providing native feature-importance measures [5].

### 2.1.2 Ensemble Methods and Gradient Boosting (XGBoost, LightGBM, CatBoost)

Gradient boosting frameworks have emerged as dominant approaches. Because credit risk datasets often exhibit non-random missingness, effective reconstruction of missing data can materially affect downstream model performance. GAN-based multiple imputation has been explored to improve the reliability of credit risk assessment under missing data[7]. CatBoost addresses categorical feature handling through ordered target encoding, with AUC improvements of 2-8 percentage points over traditional algorithms[7].

### 2.1.3 Deep Learning Approaches (Neural Networks, LSTM, CNN)

Deep learning architectures offer representational flexibility for modeling complex feature interactions. LSTM networks are effective for sequential credit data, capturing temporal dependencies in transaction histories [9]. CNN-LSTM architectures combine spatial feature extraction with temporal modeling, achieving state-of-the-art results on credit prediction benchmarks**Error! Reference source not found.**.

## 2.2 Multi-source Data Fusion Techniques in Finance

### 2.2.1 Early Fusion, Late Fusion, and Hybrid Fusion Strategies

Data fusion paradigms follow three architectures. Early fusion concatenates features into unified representations before training. Late fusion trains separate models on individual sources and combine predictions at the decision level. Hybrid fusion balances these through intermediate integration points[9].

### 2.2.2 Multimodal Learning for Credit Risk Assessment

Beyond feature concatenation, graph-based and hybrid GNN approaches can model relational dependencies among entities and implicitly learn the relative importance connected signals, which is beneficial for credit risk analysis[10]

## 2.3 NLP and Text Analysis for Financial Applications

### 2.3.1 Financial Sentiment Analysis and Text Mining

NLP techniques enable the extraction of structured features from unstructured financial texts. Textual fields (e.g., loan descriptions or borrower-provided narratives) can be encoded into compact representations and fused with structured variables to enhance multi-source credit risk assessment[13].

### 2.3.2 Integration of Unstructured Text with Structured Data

Effective integration requires bridging representational differences between sparse text embeddings and dense feature vectors. An additional practical challenge is the presence of structured missingness in credit-scoring data. Recent techniques explicitly analyze missingness patterns and reconstruct incomplete variables to improve predictive stability**Error! Reference source not found.**.

## 3. Methodology

## 3.1 Data Sources and Preprocessing

### 3.1.1 Traditional Credit Data (Credit Bureau, Loan History)

The experimental dataset integrates information from multiple institutional sources spanning January 2019 to December 2024. Traditional credit bureau data comprises 47 features, including payment history indicators, outstanding debt levels, credit utilization ratios, account age metrics, and inquiry counts.

The dataset contains 125,847 unique credit records, with a default rate of 7.83%, defined as accounts that reach 90+ DPD within 24 months of the observation date. All features are constructed using information available up to that date. Table 1 provides descriptive statistics for key credit bureau features stratified by default status.

**Table 1:** Descriptive Statistics of Credit Bureau Features by Default Status

| Feature | Non-Default Mean | Non-Default SD | Default Mean | Default SD | p-value |
|---|---|---|---|---|---|
| Credit Score | 712.4 | 68.3 | 623.7 | 82.1 | <0.001 |
| Total Debt ($) | 45,672 | 38,291 | 67,834 | 52,147 | <0.001 |

| | | | | |
|---|---|---|---|---|
| Credit Utilization (%) | 31.2 | 24.7 | 58.4 | 31.3 | <0.001 |
| Accounts in Good Standing | 8.7 | 4.2 | 5.3 | 3.8 | <0.001 |
| Recent Inquiries (6mo) | 1.4 | 1.8 | 2.9 | 2.7 | <0.001 |
| Derogatory Marks | 0.3 | 0.7 | 1.8 | 2.1 | <0.001 |
| Average Account Age (mo) | 127.3 | 72.6 | 84.2 | 61.4 | <0.001 |

Following standard preprocessing practices in credit risk early-warning pipelines, extreme numerical values are capped to reduce model instability and prevent rare outliers from dominating training dynamics[12].

### 3.1.2 Transaction Behavior Data Collection and Cleaning

Transaction-level banking data provides granular behavioral signals unavailable from aggregated credit bureau reports. The dataset incorporates 18 months of checking and savings account transactions, totaling approximately 47.3 million individual transactions. Feature engineering transforms raw records into 156 behavioral indicators organized across temporal, categorical, and statistical dimensions.

Temporal features capture spending velocity, payment timing patterns, and cash-flow volatility. The coefficient of variation of monthly income deposits serves as an indicator of stability. Categorical transaction analysis aggregates spending by merchant category codes, enabling identification of discretionary versus essential expenditure patterns. Table 2 summarizes the engineered transaction behavior features and their univariate predictive power, measured using Information Value.

**Table 2:** Transaction Behavior Feature Summary with Predictive Power Metrics

| Feature Category | Features | Mean IV | Max IV | Features with IV > 0.1 |
|---|---|---|---|---|
| Income Patterns | 23 | 0.087 | 0.234 | 8 |
| Spending Behavior | 42 | 0.065 | 0.189 | 12 |
| Cash Flow Metrics | 31 | 0.112 | 0.312 | 18 |
| Temporal Patterns | 28 | 0.054 | 0.156 | 7 |
| Account Balance | 19 | 0.143 | 0.287 | 14 |
| Transaction Frequency | 13 | 0.078 | 0.167 | 5 |

Missing transaction data arises from account dormancy periods. The Generative Adversarial Imputation Network approach addresses missing value imputation through adversarial training. From the 156 engineered transaction indicators, 89 were retained after IV/RFE-based feature selection.

### 3.1.3 Macroeconomic Indicators Integration

Systematic risk factors are incorporated through 24 macroeconomic time series. Federal Reserve Economic Data provides monthly observations for unemployment rates, consumer price indices, housing price indices, and yield curve spreads. Feature engineering emphasizes change rates and momentum indicators rather than level values.

### 3.2 Feature Extraction and Engineering

### 3.2.1 Numerical Feature Extraction from Structured Data

Structured data processing begins with comprehensive exploratory analysis. Numerical features undergo standardization using robust scaling methods. Log transformations are appropriate for right-skewed monetary variables. Polynomial feature expansion generates second-order interaction terms among the 20 most essential base features.

Table 3 presents the feature-extraction pipeline stages and the resulting feature counts.
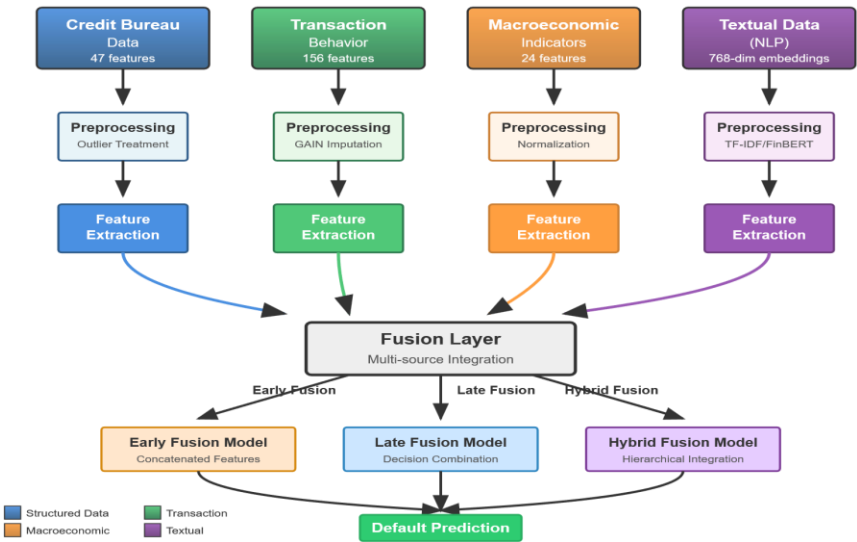
**Table 3:** Feature Extraction Pipeline Summary

| Processing Stage | Input Features | Output Features | Transformation Type |
|---|---|---|---|
| Raw Data Integration | - | 227 | Data collection |
| Missing Value Imputation | 227 | 227 | GAIN imputation |
| Outlier Treatment | 227 | 227 | Winsorization |
| Polynomial Expansion | 20 | 210 | Interaction terms |
| Ratio Features | 227 | 45 | Relationship encoding |
| Final Feature Set | - | 794 | Combined representation |

### 3.2.2 Text Feature Extraction Using NLP Techniques

Textual data sources include borrower-provided employment descriptions, loan-purpose statements, and industry-level financial news articles, which are used to construct monthly sentiment indices aligned by industry and month.TF-IDF vectorization converts preprocessed text into sparse numerical representations with approximately 8,500 unique terms.

FinBERT embeddings provide dense 768-dimensional representations capturing contextual semantics. Financial sentiment scores derived from FinBERT classification heads quantify positive, negative, and neutral sentiment intensities.

Figure 1: Multi-source Data Integration Architecture



This figure illustrates the complete data integration pipeline architecture. The visualization depicts four parallel data streams (credit bureau, transaction behavior, macroeconomic indicators, and textual data) flowing through source-specific preprocessing modules. Each stream passes through feature-extraction stages, represented as processing blocks. The streams converge at a central fusion layer that supports three branching paths for early, late, and hybrid fusion strategies. Color coding differentiates data modalities: blue for structured numerical data, green for transaction sequences, orange for macroeconomic time series, and purple for textual features. For fusion experiments, high-dimensional TF-IDF and FinBERT embeddings were summarized into compact indicators (e.g., sentiment and uncertainty scores), yielding 27 textual features listed in Table 4.

### 3.2.3 Feature Selection and Dimensionality Reduction Methods

The high-dimensional combined feature space necessitates systematic feature selection. A multi-stage selection pipeline applies filter, wrapper, and embedded methods in sequence. Information Value thresholds eliminate features with an IV below 0.02. The IV calculation follows the standard formula:

$$IV = \sum_{k=1}^{n} (Dist\_Good_k - Dist\_Bad_k) \times WoE_k$$

where WoE = ln (Distribution_Good / Distribution_Bad)

Recursive Feature Elimination with cross-validated performance evaluation serves as the wrapper method. The final reduced feature set contains 187 variables selected through consensus across multiple selection methods.

## 3.3 Data Fusion Strategies Comparison

### 3.3.1 Early Fusion Approach Implementation

Early fusion concatenates all preprocessed features into a single representation before model training. The combined feature vector comprises 187 selected features spanning credit bureau attributes, transaction behavior indicators, macroeconomic context variables, and NLP-derived text features. The implementation addresses modality imbalance by inversely weighting features based on the source dimensionality.

### 3.3.2 Late Fusion Approach Implementation

Late fusion trains independent models on each data source and combines predictions at the decision level. Four source-specific models operate on credit bureau data, transaction behavior features, macroeconomic indicators, and textual features, respectively. Prediction combination strategies include simple averaging, weighted averaging, and stacking through meta-learner training.

Table 4 compares the source-specific model performances used to inform late-fusion weight allocation.

**Table 4:** Source-Specific Model Performance for Late Fusion

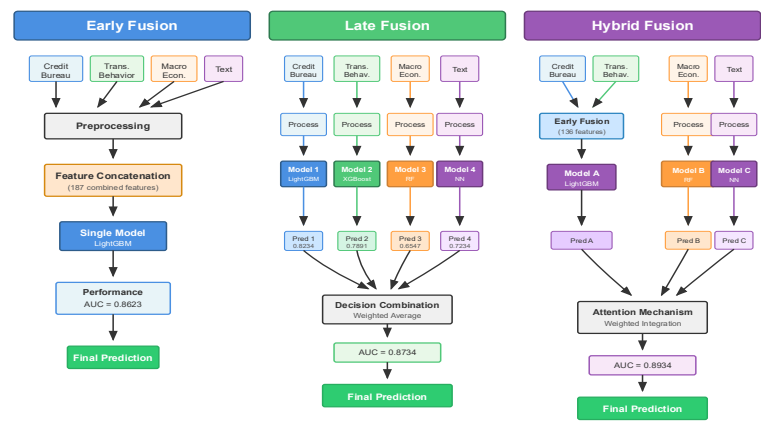| Data Source | Model Type | Features | AUC-ROC | Precision | Recall | Fusion Weight |
|---|---|---|---|---|---|---|
| Credit Bureau | LightGBM | 47 | 0.8234 | 0.724 | 0.687 | 0.312 |
| Transaction Behavior | XGBoost | 89 | 0.7891 | 0.692 | 0.654 | 0.271 |
| Macroeconomic | Random Forest | 24 | 0.6547 | 0.583 | 0.521 | 0.168 |
| Textual Features | Neural Network | 27 | 0.7234 | 0.647 | 0.612 | 0.249 |

### 3.3.3 Hybrid Fusion Approach Implementation

Hybrid fusion combines elements of early and late strategies through hierarchical integration. The implemented architecture groups credit bureau and transaction behavior features for early fusion, while macroeconomic and textual features maintain separate processing paths.

The first-stage fusion concatenates credit bureau and transaction behavior features into a behavioral representation of 136 features. Second-stage fusion combines first-stage model predictions via a stacking mechanism weighted by a meta-learner. The hybrid fusion loss function incorporates both prediction accuracy and weight smoothness regularization:

$$L_{\text{hybrid}} = L_{\text{BCE}}(y, \hat{y}) + \lambda \sum_{i=1}^{n} (|w_i - w_{\text{mean}}|^2)$$

Where L  BCE denotes binary cross-entropy, w  i represents the meta-learner–estimated fusion weights for source i, and lambda controls regularization strength.

Figure 2: Fusion Strategy Comparison Workflow



Comparison of fusion strategies showing data flow, model architecture, and performance metrics

This figure presents a comparative visualization of the three fusion strategies, arranged in a three-column layout: early fusion on the left, late fusion in the center, and hybrid fusion on the right. Each column depicts the data processing flow from raw sources at the top to final predictions at the bottom. Early fusion shows all sources converging immediately after preprocessing. Late fusion displays parallel vertical paths converging only at the decision combination. Hybrid fusion entails partial early convergence of related sources, followed by late combination. Performance metrics annotate each model block.

## 4. Experimental Design and Results

### 4.1 Experimental Setup

#### 4.1.1 Dataset Description and Characteristics

The experimental dataset encompasses 125,847 credit records from regional banking institutions in the northeastern United States. Temporal coverage spans January 2019 through December 2024, capturing pre-pandemic, pandemic, and post-pandemic conditions. Class distribution exhibits moderate imbalance with 7.83% positive instances and 92.17% negative instances.

Data partitioning employs temporal stratification. Records from January 2019 through December 2022 (78,234 observations) constitute the training set, January 2023 through December 2023 (31,456 observations) form the validation set, and January 2024 through December 2024 (16,157 observations) comprise the test set.

#### 4.1.2 Handling Imbalanced Data and Missing Values

Class-imbalance mitigation employs SMOTE-NC, which generates synthetic minority-class observations via interpolation. The oversampling ratio is designed to achieve approximately 1:1 class balance in the augmented training set. Table 5 summarizes the patterns of missing data and the imputation strategies. SMOTE-NC was used because the selected feature set contains categorical variables (e.g., encoded merchant-category groups); otherwise, standard SMOTE was applied.

**Table 5:** Missing Data Patterns and Imputation Strategies

| Feature Category | Missing Rate Range | Missing Pattern | Imputation Strategy |
|---|---|---|---|
| Credit Bureau Core | 0.2% - 3.4% | MCAR | MICE |

| Credit Bureau Supplementary | 5.7% - 18.2% | MAR | MICE |
| Transaction Statistics | 0% - 8.3% | MCAR | MICE |
| Transaction Temporal | 12.4% - 47.8% | MNAR | GAIN |
| Textual Features | 23.6% - 31.2% | MAR | Zero imputation |

### 4.1.3 Evaluation Metrics (AUC-ROC, Precision, Recall, F1-Score, KS Statistic)

Model evaluation employs a comprehensive suite of metrics. AUC-ROC is the primary metric for discrimination. Precision quantifies the proportion of predicted defaults corresponding to actual defaults. The Kolmogorov-Smirnov statistic measures maximum separation between cumulative distribution functions:

$$KS = \max_{s}|F_{\text{default}}(s) - F_{\text{non-default}}(s)|$$

where F denotes the cumulative distribution function, and s represents the model score.
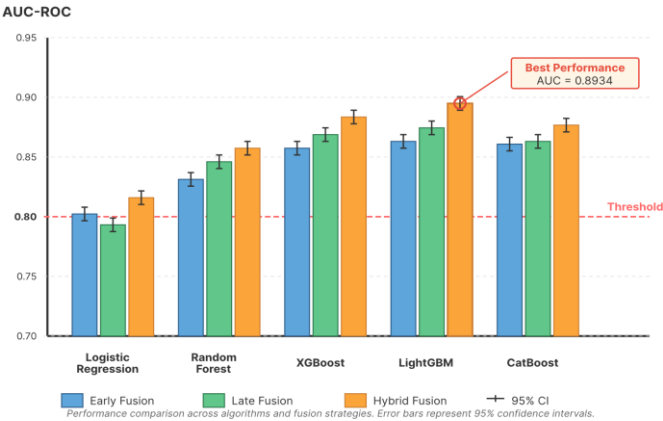
## 4.2 Comparative Analysis of Algorithms

### 4.2.1 Baseline Algorithm Performance

Baseline models establish benchmarks using traditional algorithms. Logistic regression with L2 regularization provides an interpretable linear baseline. A Random Forest with 500 trees serves as the ensemble baseline. Baseline performance indicates that the Random Forest achieves the highest AUC (0.8156) among traditional methods, followed by logistic regression (0.7923).

### 4.2.2 Ensemble Method Performance Comparison

Gradient-boosting algorithms undergo extensive hyperparameter optimization via Bayesian search. XGBoost hyperparameters include learning rate (0.01-0.3), maximum depth (3-10), and subsample ratio (0.6-1.0). LightGBM optimization additionally considers the number of leaves (20-100). Optimization uses 50 Bayesian iterations, with 5-fold cross-validation AUC as the objective metric.

Figure 3: Algorithm Performance Comparison Across Fusion Strategies



This figure displays a grouped bar chart comparing algorithm performance across fusion strategies. The x-axis arranges algorithms (Logistic Regression, Random Forest, XGBoost, LightGBM, CatBoost) in groups, with each group containing three bars representing early fusion (blue), late fusion (green), and hybrid fusion (orange). The y-axis shows AUC-ROC values ranging from 0.70 to 0.95. Error bars indicate 95% confidence intervals computed through bootstrap resampling. A horizontal dashed line at 0.80 marks the strong performance threshold. Annotations highlight the highest-performing combination (LightGBM with hybrid fusion; AUC = 0.8934).

## 4.3 Analysis of Experimental Results

### 4.3.1 Data Fusion Strategy Effectiveness Comparison

Comprehensive evaluation reveals consistent patterns in relative performance. Hybrid fusion achieves superior results for all tested algorithms, with average AUC improvements of 4.2% over early fusion and 2.8% over late fusion. The advantage of hybrid fusion is most pronounced for gradient-boosting algorithms, suggesting synergistic interactions between the fusion architecture and ensemble-learning mechanisms.

Early fusion demonstrates competitive performance for linear models, including logistic regression, where the unified feature space aligns naturally with the algorithm's assumption of additive feature contributions. The performance gap between fusion strategies widens for nonlinear algorithms that model complex feature interactions. Late fusion exhibits the highest variance across experimental replications, attributable to the propagation of source-specific model uncertainty through the prediction combination stage.

Table 6 presents comprehensive performance metrics across all fusion strategies and algorithm combinations.

**Table 6:** Comprehensive Performance Metrics by Fusion Strategy and Algorithm

| Fusion Strategy | Algorithm | AUC-ROC | Precision | Recall | F1-Score | KS Statistic |
|---|---|---|---|---|---|---|
| Early | Logistic Regression | 0.8023 | 0.698 | 0.652 | 0.674 | 0.412 |
| Early | Random Forest | 0.8312 | 0.741 | 0.698 | 0.719 | 0.467 |
| Early | XGBoost | 0.8567 | 0.773 | 0.724 | 0.748 | 0.512 |
| Early | LightGBM | 0.8623 | 0.782 | 0.736 | 0.758 | 0.523 |
| Early | CatBoost | 0.8534 | 0.768 | 0.718 | 0.742 | 0.498 |
| Late | Logistic Regression | 0.7934 | 0.687 | 0.643 | 0.664 | 0.398 |
| Late | Random Forest | 0.8456 | 0.756 | 0.712 | 0.733 | 0.487 |
| Late | XGBoost | 0.8678 | 0.789 | 0.741 | 0.764 | 0.534 |
| Late | LightGBM | 0.8734 | 0.798 | 0.752 | 0.774 | 0.547 |
| Late | CatBoost | 0.8623 | 0.779 | 0.732 | 0.755 | 0.518 |
| Hybrid | Logistic Regression | 0.8156 | 0.712 | 0.667 | 0.689 | 0.434 |
| Hybrid | Random Forest | 0.8567 | 0.769 | 0.723 | 0.745 | 0.512 |
| Hybrid | XGBoost | 0.8823 | 0.812 | 0.763 | 0.787 | 0.567 |
| Hybrid | LightGBM | 0.8934 | 0.824 | 0.778 | 0.800 | 0.589 |
| Hybrid | CatBoost | 0.8756 | 0.798 | 0.751 | 0.774 | 0.543 |

The optimal configuration—LightGBM with hybrid fusion—achieves an AUC-ROC of 0.8934, which is 7.0 percentage points (8.5% relative) higher than the best single-source credit-bureau baseline (AUC-ROC 0.8234), and a KS statistic of 0.589 that substantially exceeds standard industry thresholds. The KS statistic of 0.589 substantially exceeds industry thresholds.

### 4.3.2 Feature Importance Analysis Using SHAP Values

TreeSHAP algorithms enable efficient computation for gradient boosting models, generating instance-level feature attributions. Analysis of the optimal LightGBM hybrid fusion model reveals that transaction behavior features collectively account for 34.2% of total SHAP importance, followed by credit bureau features (31.8%), textual sentiment features (18.6%), and macroeconomic indicators (15.4%).

Table 7 presents the top 20 features by SHAP importance.

**Table 7:** Top 20 Features by SHAP Importance in Optimal Model

| Rank | Feature Name | Source Category | SHAP Importance | Cumulative % |
|------|--------------|-----------------|-----------------|--------------|
| 1 | Credit Utilization Ratio | Credit Bureau | 8.7% | 8.7% |
| 2 | Payment History Percentage | Credit Bureau | 7.2% | 15.9% |
| 3 | Cash Flow Volatility (3mo) | Transaction | 6.3% | 22.2% |
| 4 | Derogatory Mark Count | Credit Bureau | 5.1% | 27.3% |
| 5 | Monthly Spending Variance | Transaction | 4.8% | 32.1% |
| 6 | Total Debt to Income Ratio | Credit Bureau | 4.2% | 36.3% |
| 7 | Account Balance Trend | Transaction | 3.9% | 40.2% |
| 8 | Recent Credit Inquiries | Credit Bureau | 3.6% | 43.8% |
| 9 | Negative Sentiment Score | Textual | 3.2% | 47.0% |
| 10 | Discretionary Spending Ratio | Transaction | 2.9% | 49.9% |
| 11 | Average Account Age | Credit Bureau | 2.7% | 52.6% |
| 12 | Income Deposit Regularity | Transaction | 2.5% | 55.1% |
| 13 | Unemployment Rate (Regional) | Macroeconomic | 2.3% | 57.4% |
| 14 | Late Payment Frequency | Transaction | 2.1% | 59.5% |
| 15 | Housing Price Index Change | Macroeconomic | 1.9% | 61.4% |
| 16 | Uncertainty Language Score | Textual | 1.8% | 63.2% |
| 17 | Weekend Transaction Ratio | Transaction | 1.6% | 64.8% |
| 18 | Open Account Count | Credit Bureau | 1.5% | 66.3% |
| 19 | Consumer Confidence Index | Macroeconomic | 1.4% | 67.7% |
| 20 | ATM Withdrawal Frequency | Transaction | 1.3% | 69.0% |

### 4.3.3 Statistical Significance Testing

For McNemar's test, probabilities were converted to class labels using the threshold that maximizes F1 on the validation set. McNemar's test assesses classification agreement between model pairs. DeLong's test provides a direct comparison of AUC values. The comparison between optimal hybrid fusion (AUC=0.8934) and early fusion (AUC=0.8623) yields a z-statistic of 4.23 ($p<0.001$). Bootstrap resampling with 1000 iterations yields a 95% confidence interval for the optimal model AUC of [0.8856, 0.9012]. The AUC difference ($\Delta$AUC) between hybrid and early fusion is 0.0311.

## 5. Conclusion

### 5.1 Summary of Findings

### 5.1.1 Optimal Data Fusion Strategy Identification

This research provides empirical evidence establishing hybrid fusion as the optimal integration approach for multi-source credit default prediction. The systematic comparison across early, late, and hybrid strategies demonstrates consistent

superiority of hybrid fusion across all tested algorithms, with average AUC improvements of 4.2% over early fusion and 2.8% over late fusion. The hybrid architecture's ability to balance cross-source interaction modeling with source-specific optimization proves advantageous for heterogeneous financial data integration.

### 5.1.2 Best-Performing Algorithm Combinations

LightGBM emerges as the optimal algorithm across fusion strategies, with the LightGBM-hybrid combination achieving the highest overall performance (AUC = 0.8934, KS = 0.589). The algorithm's histogram-based optimization and leaf-wise growth strategy are well suited to high-dimensional, mixed-type feature spaces characteristic of multi-source credit data.

## 5.2 Practical Implications

### 5.2.1 Recommendations for Financial Institutions

Financial institutions implementing multi-source credit assessment should prioritize hybrid fusion architectures given demonstrated performance advantages. The modular structure of hybrid fusion facilitates phased deployment, enabling organizations to integrate new data sources incrementally. Data infrastructure investments should emphasize transaction-level data capture capabilities, given the substantial predictive contribution of transaction behavior features.

### 5.2.2 Applications in Credit Approval and Risk Pricing

The optimized fusion model supports multiple credit lifecycle applications beyond binary default prediction. Risk-based pricing calibration can leverage predicted default probabilities directly, with well-calibrated probability estimates enabling actuarially appropriate interest rate determination. Portfolio monitoring applications employ the model as an early warning trigger, flagging accounts exceeding dynamic risk thresholds for proactive intervention.

## 5.3 Limitations and Future Work

### 5.3.1 Study Limitations

Geographic concentration in the northeastern United States may limit generalizability to other regional markets with different economic characteristics and borrower populations. The consortium data source may not fully represent national credit market diversity. The 24-month default definition may not capture shorter-term liquidity crises or longer-term gradual deterioration patterns equally well.

### 5.3.2 Future Research Directions

Extensions should explore graph neural network architectures modeling network relationships between borrowers, building on emerging work demonstrating predictive value of social and financial network structures. Federated learning approaches could enable multi-institution model training while preserving data privacy. Causal inference methods applied to treatment effects of credit interventions would enhance model utility for proactive portfolio management.

## References

[1]. Abedin, M. Z., Guotai, C., Hajek, P., & Zhang, T. (2023). Combining weighted SMOTE with ensemble learning for the class-imbalanced prediction of small business credit risk. Complex & Intelligent Systems, 9(4), 3559-3579.

[2]. Wang, Y., Xu, Z., Ma, K., Chen, Y., & Liu, J. (2024, December). Credit Default Prediction with Machine Learning: A Comparative Study and Interpretability Insights. In 2024 4th International Conference on Communication Technology and Information Technology (ICCTIT) (pp. 496-500). IEEE.

[3]. Das, S., Huang, X., Adeshina, S., Yang, P., & Bachega, L. (2023). Credit risk modeling with graph machine learning. INFORMS Journal on Data Science, 2(2), 197-217.

[4]. Yao, G., Hu, X., Song, P., Zhou, T., Zhang, Y., Yasir, A., & Luo, S. (2024). AdaFNDFS: An AdaBoost ensemble model with fast nondominated feature selection for predicting enterprise credit risk in the supply chain. International Journal of Intelligent Systems, 2024(1), 5529847.

[5]. Shi, S., Tse, R., Luo, W., D'Addona, S., & Pau, G. (2022). Machine learning-driven credit risk: a systemic review. Neural Computing and Applications, 34(17), 14327-14339.

[6]. Zhao, F., Lu, Y., Li, X., Wang, L., Song, Y., Fan, D., ... & Chen, X. (2022). Multiple imputation method of missing credit risk assessment data based on generative adversarial networks. Applied Soft Computing, 126, 109273.

[7]. Zhu, M., Zhang, Y., Gong, Y., Xing, K., Yan, X., & Song, J. (2024, May). Ensemble methodology: Innovations in credit default prediction using lightgbm, xgboost, and localensemble. In 2024 IEEE 4th International Conference on Electronic Technology, Communication and Information (ICETCI) (pp. 421-426). IEEE.

[8]. Li, J., Xu, C., Feng, B., & Zhao, H. (2023). Credit risk prediction model for listed companies based on CNN-LSTM and attention mechanism. Electronics, 12(7), 1643.

[9]. Li, Y., & Chen, W. (2020). A comparative performance assessment of ensemble learning for credit scoring. Mathematics, 8(10), 1756.

[10]. Sun, M., Sun, W., Sun, Y., Liu, S., Jiang, M., & Xu, Z. (2024, October). Applying Hybrid Graph Neural Networks to Strengthen Credit Risk Analysis. In 2024 3rd International Conference on Cloud Computing, Big Data Application and Software Engineering (CBASE) (pp. 373-377). IEEE.

[11]. Wang, T., Liu, R., & Qi, G. (2022). Multi-classification assessment of bank personal credit risk based on multi-source information fusion. Expert systems with applications, 191, 116236.

[12]. Wang, D., Feng, J., Zou, W., & Chen, H. (2023, October). Credit risk assessment and early warning of supply chain finance based on xgboost-lstm-a model. In Proceedings of the 2023 4th International Conference on Computer Science and Management Technology (pp. 444-449).