

Breaking Data Silos: A Privacy-Preserving Framework for Credit Risk Analysis

Samyak Shriram Gedam (252IS032), Rohit Kumar Sah (252IS030), Anusha Hegde H

Department of Computer Science and Engineering

National Institute of Technology Karnataka

Surathkal, Mangalore-575025, India

Email: {samyakgedam.252is032, rohitkumarsah.252is030, anushahegde.227cs001}@nitk.edu.in

Abstract—Financial institutions are currently stuck. Effective credit risk modeling needs data from different sources, but regulations like GDPR and the issue of "data silos" prevent sharing raw customer information. Centralized machine learning, which gathers data in one location, is no longer a suitable option because of privacy and compliance risks. To solve such kind of problem, we present a SecureBoost framework that uses Vertical Federated Learning (VFL) and Homomorphic Encryption (HE). Different from the traditional methods, our framework lets multiple institutions, such as banks and fintech companies to work together to train a gradient boosting model (XGBoost) on vertically splitted data without disclosing local feature distributions. By using the Cheon-Kim-Kim-Song (CKKS) encryption scheme, we are protecting the exchange of gradients and Hessians. The experiments are performed with the German Credit Dataset, by using 2-party and 4-party splits, we show that our model keeps about 96% of the predictive power of a centralized baseline model (AUC 0.737) also ensuring that no raw data leakage occurs.

Index Terms—Vertical Federated Learning; Homomorphic Encryption; Credit Risk Analysis; Privacy-Preserving Machine Learning; XGBoost; SecureBoost; Data Privacy.

I. INTRODUCTION

In the modern financial system, the capacity to evaluate the borrower of his or her repayment ability is very important. The ability to pay back is essential to the economy. Financial institutions were able to use just mere statistics for many years, they relied on simple statistical methods such as logistic regression or expert-based scorecards [1] in order to make crucial choices. However, things have changed. The industry is currently resorting to high-end machine learning algorithms, such as Gradient Boosting Decision Trees (XGBoost) [2], capable of determining a complex, non-linear patterns in borrower behavior that traditional methods tend to ignore. These algorithms require a large amount of high-dimensional data to work well. They need to have detailed customer profiles that contain financial history, logs of the transactions, and even behavioral demographics.

This requirement of data brings a fundamental conflict with the existing regulatory environment. Valuable user information is hardly ever assembled, in one place; instead, it is scattered throughout multiple isolated entities. The savings history of the user can lie in a commercial bank, while their daily spending will be monitored by a fintech app, and an e-commerce platform knows their buying preferences. In an ideal world,

a combination of these datasets would provide a complete one view of credit risk. But rigid privacy laws, such as the General Data Protection Regulation (GDPR) [3] [4] in Europe and the California Consumer Privacy Act (CCPA) [5], lawfully prohibit the immediate transfer of such sensitive information. This problem, usually referred to as the "Data Island" or "Data Silo" problem [6], forces organizations to train models only on their limited local data. Such a limitation causes poor predictions and gaps in risk assessment.

To come out of this deadlock without breaking privacy laws, Federated Learning (FL) has become a strong option. Unlike centralized training, FL takes the model to the data instead of taking the data to the model. While Horizontal Federated Learning (HFL) has effectively addressed issues where devices share the same features, such as mobile keyboards, it does not work well for the financial world. Banks and fintechs typically address a Vertical Federated Learning (VFL) [6] [7] senario, where two organizations have varying attributes for the same group of users. however, VFL is not a perfect solution. Standard VFL implementations often involves the exchange of intermediate gradients or partial models. So it can happen that information about the original raw data is leaked unintentionally through inference attacks. To ensure strong privacy, these exchanges must be protected by cryptographic protocols. Early attempts used Differential Privacy (DP) [8], that introduces statistical noise to hide individual records. While DP is cheap to compute, it permanently decreases the accuracy of the model. This trade-off is not acceptable in the financial risk scoring where a even small percent of inaccuracy can lead to millions in losses. Alternatively, it is possible with Homomorphic Encryption (HE) [9]. because using by HE we can directly perform calculations on protected data [10]. However, traditional schemes like Paillier [11] can only be used with integers calculations, making them inefficient for the complex floating-point operations needed by modern boosting algorithms.

To solve this, we propose a practical solution to address these gaps. We established an effective privacy-protecting system that is an integration of Vertical Federated Learning and the CKKS Homomorphic Encryption scheme [12]. Specifically, we modify the SecureBoost protocol [13] to handle the precise, continuous variables involved in the analysis of credit risk. By switching from standard integer-based encryption to CKKS encryption, our framework enables institutions to

collaborate in order to develop an effective XGBoost model using vertically partitioned data. Also, this process makes sure that the raw data as well as the local features do not leak. This policy will be very strict of data privacy regulations while leveraging the predictive power of collaborative intelligence.

II. RELATED WORKS

The latest developments of privacy preserving methods especially homomorphic encryption is of recent popularity among researchers. Although different sophisticated Machine Learning (ML) models do provide better predictive power but they also have risks about leaking of personal or sensitive financial information. This has lead various researchers to look into more combinations of ML algorithms and cryptographic techniques like Homomorphic Encryption (HE) [14] to give more financial stability and decision making in checking credit worthiness of the people. Some of the approaches that have been explored by researchers includes deep learning models operating with CKKS based encryption [15], federated learning with homomorphic encryption [16], and hybrid frameworks combining encryption and zero-knowledge proofs (ZKPs) [17], also more interpretable models such as penalized logistic tree regression [18] and gradient boosting with SHAP values [19], [20], [21]. Their works prove to overcome the challenges of improving accuracy while maintaining privacy and computational efficiency. While deep learning models deliver more accuracy [15] [22], they have limitations in scalability. Other hybrid models [18], [23], provide transparency but may underperform in complex scenarios. The following literature survey identifies various frameworks that integrate privacy-preserving techniques with more accurate prediction models, laying the foundation of more trustworthy credit risk prediction systems.

Naresh et al. [15] proposed a PPDNN-CRP framework to ensure privacy throughout the entire credit risk prediction workflow. Their design combines Deep Neural Networks [24] with Cheon-Kim-Kim-Song (CKSS) [25] HE scheme. It ensured end-to-end protection of data during all the phases of credit risk prediction algorithm keeping the data encrypted without the need to decrypt it. This framework was implemented using PyTorch and TenSEAL libraries and they compared this framework against unencrypted data to evaluate models performance. Real world datasets from Kaggle with 32,581 records and 12 features related to loan applications were used which yielded high accuracy as compared to unencrypted data. Results showed that this model achieved high predictive performance with 0.88 precision and recall, 0.86 F1-score, 89.17% accuracy, and 0.88 ROC which shows PPDNN-CRP outperforms various other models

Proper credit risk management plays a vital role in supply chain finance. Chang et al. [19] studied and focused on the role of Random Forest and Gradient Boosting models that improve decision making capabilities of AI systems by enhancing interpretability using feature importance rankings and SHAP values. This helped them to identify differences between high- and low-risk customers to help various financial service providers such as banks to predict credit risk accurately. The

authors compared various traditional classifiers by using a large dataset of more than 30,000 loan default cases from UK banks which were sourced from Kaggle. It was categorised into train and test data and compared performance based on metrics such as Accuracy, F1-score and AUC. The results showed that Gradient Boosting classifier outperformed various other models obtaining an accuracy of 82.49% and F1-score of 80%. It used recent payment history as the most important factor for prediction of credit risk. Although some of the limitations such as class imbalance and insufficient features resulted in less accurate results, but overall this work contributed in responsible decision-making in financial services

Bao et al. [17] realised the need of optimizing the credit scoring models majorly removing their dependency on central data that was a risk of exposing private information of any individual and may develop trust issues among parties. This study proposed a system that combined zero-knowledge proofs (ZKPs) and inner product functional encryption (IPFE). Based on these foundations, they introduced two schemes: First, PPCS for traditional ML models and PICS for neural-network-based models. Experiments are performed on UCI dataset of credit cards that contain 150,000 records of credit card borrowings which showed that both PPCS and PICS are much efficient as compared to the basic Paillier homomorphic encryption scheme. PICS showed upto 40% faster encryption and lower storage costs for large scale data. While PICS achieved around 94% accuracy, 0.85 AUC, 54.3% precision and an F1-score of 19.0%. The framework provided both efficiency and strong privacy, making it suitable for large platforms.

Zhu et al. [16] identified the problems caused by centralized data storage and computations which were vulnerable to data and security and breaches. Safeguarding data privacy in such cases became a major concern. They proposed a framework which integrated privacy-preserving federated learning with homomorphic encryption. The system kept participants local model parameters encrypted during transmission and aggregation using Paillier homomorphic scheme, preventing data leakage even from the intermediate nodes. The edge computing nodes reduced communication and computation overhead by performing secure aggregation on encrypted data before sending that data to central sensing platform for global aggregation, without decrypting intermediate data. Performed on MNIST dataset of handwritten images where large samples for both training and testing were used. To resemble federated learning the dataset is split non-IID distribution demonstrating that this approach surpass the classical federated methods in performance. Key observations were that Encryption and Decryption times increase with the key length and participation count. Some of the limitations include reliance on a trusted key generation center which poses a single point of failure risk and data heterogeneity.

Homomorphic encryption (HE) has proven to provide strong security guarantee with various ML algorithms. Although it is secure but its computational efficiency is still questioned. Byun et al. [26] proposed a more effective strategy that is to encrypt only a small amount of information instead of entire dataset. They put forward an efficient privacy-preserving

ridge regression framework using HE which encrypts only private variables to improve efficiency while maintaining privacy. Instead of encrypting full dataset, only sensitive private variables are encrypted. The framework resembled a three-party model consisting of owners of data, a service provider for machine learning (MLSP), and a crypto-service provider (CSP). The authors proposed a framework which performs the ridge regression training which encrypt only private variables and non sensitive variables remain unencrypted. They also gave an adversarial perturbation mechanism to keep safe the sensitive variables from inference attacks, which have rarely been explored in HE-based machine learning studies. They evaluated this framework on eleven real world datasets containing numerical and binary categorical variables. The rate of computation of their proposed method was 5-20 times faster than a full-column encrypted method and about 1.7-1.9 times slower than encrypting a single private variable only. The method was good in balancing privacy and efficiency despite some limitations such as Efficiency decrease when datasets are bigger than the ciphertext slots packing capacity which affected performance.

Credit risk assessment plays a important role in automobile finance. Wang et al. [20] introduced a model for assessing credit risk evaluation for car loans by merging eXtreme Gradient Boosting utilizing (XGBoost) alongside SHapley Additive exPlanations (SHAP), producing increased predictability and clarity. Together with various additional algorithms such as logistic regression, random forest, decision tree, GBDT and XGBoost, the study's objectives centered around empirical model assessment utilizing actual loan data from a Chinese car financing service. The assessment of the model was conducted from multiple models using the various performance metrics. Their results showed that XGBoost performed the best with the following: 0.88 accuracy, 0.98 precision, 0.78 recall, and an F1 score of 0.87. In addition, the SHAP analysis allows researchers to sort and identify credit level, credit score, and disbursed amount as the most important features reenforcing interpretability in credit risk evaluation through transparency and reliability.

Dumitrescu et al. [18] introduced a new hybrid credit scoring mechanism referred to as penalized logistic tree regression (PLTR) in order to provide improved prediction power to the classic logistic regression model while preserving its interpretability. Their method embeds binary rules from shallow depth decision trees into a penalized logistic regression framework, adopting adaptive lasso for variable selection. This enables the model to grasp both univariate and bivariate threshold impacts that could be overlooked in a conventional logistic regression method, addressing directly the established compromise between predictive power and interpretability for a credit scoring decision-making. The authors validated PLTR through extensive Monte Carlo simulation studies and a real-world application on the "Give Me Some Credit" [27] data set (with 150,000 loan observations and 10 predictors). PLTR was compared against standard logistic regression scenarios (linear and non-linear), support vector machine, neural network, and random forest modeling. The authors found that PLTR provided statistically significantly

enhanced classification accuracy (PCC, AUC, KS, BS, and PGI metrics) compared to logistic regression while yielding accuracy especially similar to random forest modeling, while remaining far more interpretable through its sparse rule set and marginal effect interpretation. The authors asserted that PLTR provides a viable solution for balancing the trade-off between interpretability and predictive power and can be useful for financial institutions that require predictive accuracy and suitable for regulators, where credit scoring models should support stakeholder interpretability.

The evaluation of credit risk is important when it comes to the automobile finance. Lin et al. [21] proposed a credit risk evaluation model for automobile loans by combining eXtreme Gradient Boosting (XGBoost) with SHapley Additive exPlanations (SHAP), yielding greater predictability and interpretability. The targets of the study in addition to various other Machine Learning algorithms, the objectives of the study focused on practical model assessment through actual loan data of a Chinese automobile finance service. The model analysis was conducted with respect to a number of models and various performance metrics were used for that analysis. They found that XGBoost best worked with the following: 0.88 accuracy, 0.98 precision, 0.78 recall, and an F1 score of 0.87. Also, the SHAP analysis enables researchers to rank and mark credit level, credit score, and amount disbursed as the most important aspects that provide interpretability in evaluating credit risk based on transparency and reliability.

Huang et al. [23] studied the importance of BigTech and examined how these firms utilize the alternative data to check credit risk and analyze worthyness of the people. The research assessed BigTech's use of big data and machine learning-based models to evaluate credit risk against traditional banks' scorecard-based models, using a dataset of 1.8 million loan records from MYbank, China. In the performance assessments using the AUC metric, BigTech achieved an AUC of 0.84 whereas the scorecard-based model produced an AUC of 0.72; clearly indicating BigTech's approach produced significantly better results than the traditional bank models. Furthermore, BigTech's alternative data sources enhanced credit risk assessment models to account for missing credit histories, thus allowing unbanked SMEs and SMEs with basing cities in smaller cities to access credit. The findings support that BigTech offers a mechanism for greater predictive accuracies to replace asks to assess credit risk while serving to enhance financial inclusion.

Shi et al. [22] did a systematic review on using machine learning in credit risk assessment. The study looked at 76 research papers published over the last eight years. It grouped the approaches into statistical methods, machine learning, and deep learning. The authors proposed a new classification system for ML-based credit risk algorithms and ranked their performance using benchmark datasets, including German and Australian credit data. The results indicated that deep learning models generally did better than traditional statistical methods and classical ML in accuracy and AUC. Ensemble techniques also showed better predictive performance compared to single models. The review highlighted several challenges, such as data imbalance, the lack of standardized benchmark datasets,

TABLE I
SUMMARY OF THE LITERATURE REVIEW

Sr. No.	Author	Methodology	Limitations
1	Naresh et al. [15]	Integrated DNNs with CKKS (HE).	PPDNN-CRP introduces computational overhead due to encryption operations. Simpler models may achieve a better trade-off between cost and accuracy.
2	Chang et al. [19]	Compared classification performances of multiple models and concluded that Gradient Boosting outperformed others.	Dataset lacked sufficient features, reducing training effectiveness and prediction accuracy.
3	Bao et al. [17]	Proposed PPCS and PICS schemes using Zero-Knowledge Proofs (ZKPs) and Inner Product Functional Encryption (IPFE).	Neural network training remains costly and relies on a trusted Key Management Authority (KMA).
4	Zhu et al. [16]	Integrated Privacy-Preserving Federated Learning (PPFL) with Paillier HE for secure aggregation at edge nodes.	Dependence on a trusted key generation center and non-IID data distribution impact reliability and scalability.
5	Byun et al. [26]	Proposed an HE-based Ridge Regression framework encrypting only sensitive private variables.	Performance degrades when ciphertext slot capacity is exceeded, requiring multiple encryptions per column.
6	Wang et al. [20]	Developed MP-DLR using vertical federated learning, HE, Taylor approximation, and L1-regularization.	Limited to logistic regression, computationally expensive with many participants, and validated only on anonymized data.
7	Dumitrescu et al. [18]	Proposed PLTR combining logistic regression with shallow decision trees for interpretability.	Trade-offs remain between interpretability and performance; scalability is limited for large datasets.
8	Lin et al. [21]	Applied logistic regression and XGBoost with SHAP explanations on automobile loan data.	SHAP may introduce bias and fails to fully capture causal relationships; dataset limited to one institution.
9	Huang et al. [23]	Compared traditional scorecards and Random Forest using traditional and BigTech data sources.	Results may not generalize due to reliance on a single confidential dataset.
10	Shi et al. [22]	Conducted PRISMA-based systematic review of 76 key studies on credit risk models.	Issues include data imbalance, lack of standardized benchmarks, limited interpretability, and underuse of deep learning.

and limited interpretability, which hold back the wider use of machine learning models. Overall, the study gave a clear summary of current research and pointed out ways to improve model design and use in credit risk assessment.

In spite of the fact that these studies have already come a long way to enhance predictive performance and accuracy, protect data and enhance privacy and model interpretability, there are still several limitations that can be identified. Deep learning approaches were found to be very accurate, but lack transparency, whereas other models are computationally expensive and might not be useful to work on complex data. Most encryption schemes are also challenging to use in efficiency and key management in the real world financial settings.

List of common limitations (Research Gaps)

- **High computational cost:** The majority of HE and MPC-based algorithms substantially raise the time of encryption/decryption and training models.
- **Issue of data Centralization:** The majority of the frameworks operate with data that has been stored in a central repository that can be breached and attacked resulting in privacy-related issues..
- **Difficulties in Key Management:** The safe dissemination and/or rotation of cryptographic keys is not fully addressed
- **Absence of Hybrid Solutions:** There are only a few studies that integrate HE with other privacy methods (e.g., federated learning, secure aggregation, etc).

This creates a gap in the research to come up with a framework that would integrate federated learning with gradient boosting algorithms with a homomorphic encryption for

robust computation, effectiveness and improve fundamental management and offers a high performance in credit risk evaluation.

Although available sources demonstrate the potential of Federated Learning in FinTech, most applications rely on the Paillier cryptosystem, which is limited to integer arithmetic, or fail to properly measure the resulting performance difference in a multi-party applications. To address these problems and bridge the gap between theoretical privacy and practical use, the paper will include the following contributions:

- **CKKS-Optimized SecureBoost Implementation:** Here we modify the SecureBoost system to incorporate the CKKS homomorphic encryption system. This permits precise floating point gradient aggregation, which is suitable for complex credit risk datasets.
- **High Performance:** We demonstrated that our privacy-preserving model achieves an AUC of **0.737** on the German credit data [28], significantly outperforming previous encrypted models and competing with non-private baselines.
- **Scalability Analysis:** We performed experiments in a 2-party and 4-party environments, and demonstrated that our algorithm is robust to data fragmentation and retains the same accuracy as the number of data silos is increased.
- **Overhead Quantification:** We gave a specific analysis of the computation cost (“privacy tax”) associated with homomorphic encryption, which gives a clear understanding of the trade-offs of actual deployment.

III. PROPOSED METHODOLOGY

The proposed study aims at modifying a privacy-preserving federated gradient boosting architecture with for secure credit risk prediction known as SecureBoost which will eliminate the most significant flaws of the current solutions. It is a privacy preserving variant of gradient boosting. Conventional gradient boosting involved a central repository of the data which was usually vulnerable to attacks as it compromised the privacy of entire data if attacked. The solution to this is Vertical federated learning (VFL) [29] in which the data is partitioned among multiple parties which enables them to combinely train the model while raw data is not disclosed. This offers a high level of privacy, and enables to collaboratively train the model for secure credit risk prediction.

A. Vertical Federated Learning (VFL)

VFL is created for the "split-brain" scenario often seen in finance. One company knows who the customer is by labels, while another knows what the customer did through features. VFL involves a group of data owners collaboratively train a model without sharing the underlying raw data with each other. The system architecture is such that it has Active Party and a number of Passive Parties (P_A, P_B, \dots, P_N) [29] [16]. In the implementation, three logical parties are simulated, one active and two passive parties (P_A and P_B), as shown in Figure 1. The underlying data set is vertically partitioned i.e. it is partitioned into two subsets where each party has features different from another party of the same data sample.

- **Active Party:** This party acts as a coordinator of the system which holds only the labels \mathcal{Y} and dont have access to any of the features of the dataset.
- **Passive Party:** Each passive party holds a different subset of features $\mathcal{X}_A, \mathcal{X}_B$ for the same set of sample IDs. They dont have access to labels. One cannot access the features or even labels of another party.

The core of this model is to preserve the privacy of traditional Gradient boosting model by collaboratively training it without any of the parties revealing its raw data. The privacy-preserving mechanism used here is homomorphic encryption, which allows for complex mathematical computations directly on encrypted data without ever decrypting it.

B. Homomorphic Enryption (HE)

The system uses Homomorphic Encryption (HE) [14] [9] as a fundamental privacy-preserving tool, it is able to process encrypted information and makes it possible to perform calculations on encrypted data without the necessity to decrypt them. This will make sure that the sensitive information remains secure during the entire training process. In particular, the Cheon-Kim-Kim-Song (CKKS) [12] scheme is used in the implementation since it allows addition and multiplication operations encrypted numerical values - operations that are fundamental to gradient boosting in the encrypted space. The CKKS scheme is very helpful in privacy-preserving credit risk analysis as it allows financial institutions to predict creditworthiness of a customer using various machine learning

models. The implemented system performs all the cryptographic tasks using Pyfhel library which offers convenient encryption, decryption, and manipulation of ciphertext. Active Party constructs a pair of public and private keys of the form (pk, sk) with the public key being shared among Passive Parties. This public key is used to encrypt the intermediate computations of each passive party, then they are sent to the active party, so that no raw data is present at all. The active party holds the private key in a safe place and is the only one who is able to decrypts= it when the need arises.

$$\text{Enc}(m_1, pk) + \text{Enc}(m_2, pk) = \text{Enc}(m_1 + m_2, pk) \quad (1)$$

$$\text{Dec}(\text{Enc}(m, pk), sk) = m \quad (2)$$

C. Gradient and Hessians of the XGBoost

XGBoost, a powerful gradient boosting algorithm, it constructs an ensemble of decision trees sequentially. Each new tree is designed such that it improves upon errors of its predecessor trees. It minimize a specified loss function by "boosting" the performance of previous trees. Gradients and Hessians in the XGBoost algorithm forms the basis of our SecureBoost framework.

The core privacy challenge is to allow the active party to find the best feature split without seeing the passive parties features, and to allow the passive parties to use label-derived information (gradients/hessians) without seeing the labels.

- g_i (**Gradient**): Represents the first derivative of the loss function that determine both the optimal feature split and weight of the leaf during construction of the tree. The prediction is denoted by $\hat{y}_i^{(t-1)}$:

$$g_i = \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} \quad (3)$$

- h_i (**Hessian**): It represents the second derivative:

$$h_i = \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial (\hat{y}_i^{(t-1)})^2} \quad (4)$$

where,

- L is the loss function.
- y_i represents the true label (target variable) for sample i . For binary classification, it takes discrete values of 0 (negative class) or 1 (positive class).
- $\hat{y}_i^{(t-1)}$ is the prediction for sample i from the ensemble of trees up to round $t - 1$.

These gradients and hessians are crucial for determining the best split points in decision trees and calculating the optimal leaf weights. They carry information about each sample's contribution to the model's error.



Fig. 1. Vertical Federated Learning (VFL) Architecture

D. SecureBoost Framework

Algorithm 1 Privacy-Preserving SecureBoost Training

- 1: **Parties:**
- 2: \mathcal{A} : Active Party (holds labels and HE private key)
- 3: \mathcal{P}_j : Passive Parties (each holds a subset of features)
- 4: **Input:** Training data and hyperparameters
- 5: **Output:** Trained ensemble of decision trees
- 6: **Procedure:**
- 7: Initialize active party: generate HE keys, distribute public key, initialize model
- 8: **for** each boosting round **do**
- 9: **Active Party:** Compute predictions, gradients, Hessians, encrypt and send to passive parties
- 10: **Passive Parties:** Compute encrypted histograms for local features and send back to active party
- 11: **Active Party:** Decrypt histograms, find best split, update model
- 12: **end for**
- 13: **Return:** Final trained model

The model is trained in a series of T boosting rounds. In each round t , a new decision tree is built to correct the errors of the previous rounds.

1) *Initialization (Active Party):* \mathcal{A} computes an initial prediction $\hat{y}^{(0)}$ for all samples. It is the log-odds of the mean of \mathcal{Y} .

$$p_0 = \log(\bar{Y}/(1 - \bar{Y})) \quad (5)$$

2) *Gradient and Hessian Computation (Active Party):* \mathcal{A} computes the gradients (g_i) and Hessians (h_i) for each sample i based on the current model predictions $\hat{y}_i^{(t-1)}$ and the true labels y_i :

$$p_i^{(t-1)} = \frac{1}{1 + e^{-\hat{y}_i^{(t-1)}}} \quad (6)$$

Gradient is used to measure the rate of change of the loss function with respect to the prediction of the model- it is the rate and the direction in which the error is occurring on each sample.

$$g_i = p_i^{(t-1)} - y_i \quad (7)$$

The Hessian captures the curvature of the loss function, which characterizes the change in the gradient itself. It gives information on the second order which is used to decide on the maximum step size and division in the tree-building process.

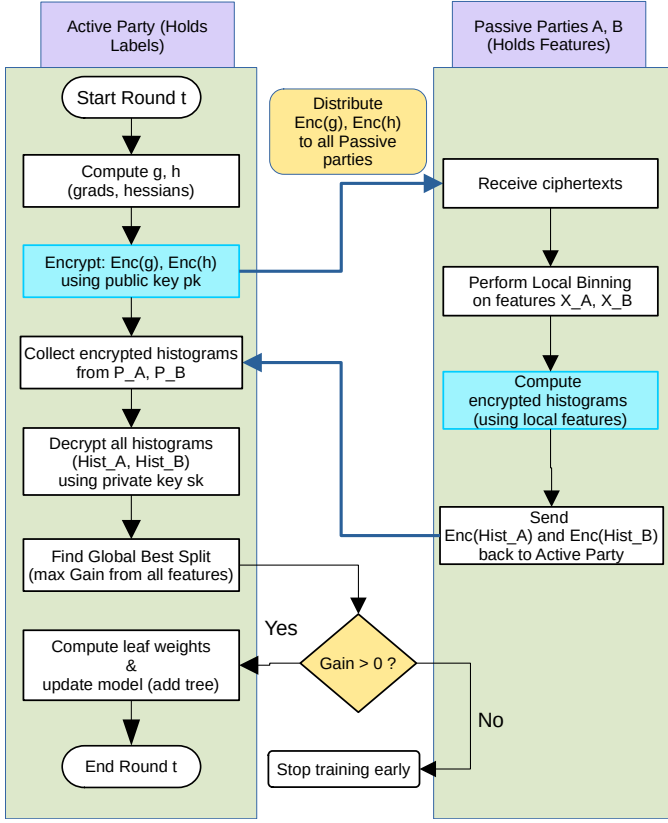


Fig. 2. System Architecture and Data Flow

$$h_i = p_i^{(t-1)} \cdot (1 - p_i^{(t-1)}) \quad (8)$$

All these values are collectively used to ensure that the boosting algorithm is effective at updating the model, towards the goal of minimizing the overall loss and eventually improving classification performance.

3) Encryption of Gradients and Hessians and Distribution: (Active \rightarrow Passive) The gradients and Hessians are encrypted by the CKKS homomorphic encryption scheme using public key (pk) to protect the privacy of the data. This encryption scheme allows mathematical operations on encryption values $Enc(g_i), Enc(h_i)$ that is, one can compute an arithmetic operation like summing two values or averaging without ever having to decrypt the numbers.

Here, the Hessians and gradients are shared after encryption. The encrypted values are only sent to each passive party that then completes the aggregation process on them to create bin-based histograms. Since the data is encrypted in the process, no side can retrieve or deduce the original sensitive data and thus it will be secure and private as well as when collaboratively training the model.

4) Local Encrypted Histogram Computation (Passive Parties): This is the core privacy-preserving step for the passive parties. Each passive party P_k performs the following locally without the private key

- **Binning:** For each of its local features $j \in \mathcal{X}_A$, it determines bin boundaries (e.g., using quantiles) to create K bins.
- **Encrypted Aggregation:** It initializes K encrypted sums to zero: $Enc(G_{j,k}) = Enc(0)$ and $Enc(H_{j,k}) = Enc(0)$ for $k = 1, \dots, K$.
- It then iterates through all N samples. For sample i , it finds the corresponding bin k for its feature j
- Using the additive property of HE, it adds the sample's encrypted stats to that bin's total:

$$Enc(G_{j,k}) \leftarrow Enc(G_{j,k}) + Enc(g_i) \quad (9)$$

$$Enc(H_{j,k}) \leftarrow Enc(H_{j,k}) + Enc(h_i) \quad (10)$$

- This process is repeated for all local features.

Crucially, Party A never sees the unencrypted g_i or h_i . It only shuffles and adds ciphertexts. In the whole process, no raw feature data and individual gradients are ever revealed and hence the high level of data confidentiality among all parties involved.

5) Histogram Aggregation: (Passive \rightarrow Active) Party A and Party B send their complete sets of encrypted histograms $\{Enc(G_{A,j,k}), Enc(H_{A,j,k})\}$ and $\{Enc(G_{B,j,k}), Enc(H_{B,j,k})\}$ to the Active Party.

6) Split Evaluation (Active Party): Active Party now holds all encrypted histograms.

- **Decryption:** It uses its private key (sk) to decrypt these aggregate sums, yielding the plaintext $G_{j,b}$ and $H_{j,b}$ for all bins of all features.

$$G_{j,k} = Dec(Enc(G_{j,K}), sk) \quad (11)$$

$$H_{j,k} = Dec(Enc(H_{j,K}), sk) \quad (12)$$

Note that no individual sample data is ever decrypted or revealed. Means it reveals the total G and H for each bin, but not the individual sample contributions, thus preserving privacy.

- **Gain Calculation:** The Active Party iterates through every feature j from all parties and every possible split point k (bin boundary). For each potential split, it computes the total G_L, H_L (left node) and G_R, H_R (right node) using prefix sums of the decrypted bin statistics.

$$G_L = \sum_{m=1}^k G_{j,m}, \quad H_L = \sum_{m=1}^k H_{j,m} \quad (13)$$

$$G_R = G_{\text{total}} - G_L, \quad H_R = H_{\text{total}} - H_L \quad (14)$$

- It iterates through all possible split points (the bin boundaries) and calculates the split gain using the XGBoost formula

$$\text{Gain} = \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma \quad (15)$$

- The Active Party selects the split (Party, feature, threshold) with the maximum gain as the best split for

this tree.

7) *Leaf Weight Calculation and Model Update (Active Party)*: Using the G_L and H_L (for the left leaf) and G_R, H_R (for the right leaf) from the best split, Active Party computes the optimal leaf weights

$$w_L = -\frac{G_L}{H_L + \lambda}, \quad w_R = -\frac{G_R}{H_R + \lambda} \quad (16)$$

This new tree $f_t(x)$ is added to the ensemble, and the local predictions \hat{y}_i are updated with a learning rate η :

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + \eta \cdot f_t(x_i) \quad (17)$$

This process (Steps 2-7) is repeated for the specified number of rounds (i.e., trees)

8) *Secure Prediction Phase*: Once the ensemble of T trees is trained, making a prediction for a new sample also follows a privacy-preserving protocol. The active party coordinates the process, but each passive party only evaluates the splits relevant to its local features. The process does not require HE.

Algorithm 2 SecureBoost Prediction

- 1: **Input**: Trained model (ensemble of trees), new sample with features split across parties
 - 2: **Output**: Predicted probability
 - 3: Initialize prediction with the base score
 - 4: **for** each tree in the ensemble **do**
 - 5: **Active Party**: Traverses the tree. When a split on a passive party's feature is needed, send the node info to the relevant party.
 - 6: **Passive Party**: Receives node info, checks its local feature value against the split threshold, and returns the resulting direction (left or right child) to the active party.
 - 7: **Active Party**: Receives the direction and continues traversal until a leaf node is reached. Adds the leaf weight to the total prediction.
 - 8: **end for**
 - 9: **Active Party**: Applies the sigmoid function to the final summed prediction to get the probability.
 - 10: **Return**: Final predicted probability
-

IV. EXPERIMENTS AND RESULTS

This section tells about the various experiments to measure the performance, efficiency, and privacy-preserving properties of the proposed SecureBoostClassifier framework. The configuration was done to reflect real-life federated settings in which different institutions involved in the training of models without revealing sensitive customer information. 2-party and 4-party systems were practiced with the CKKS homomorphic encryption system on how to secure gradient and Hessian exchanges throughout the training process. It presents the performance results against a centralized baseline XGBoost, and analyzes the model, particularly concerning scalability to multiple parties.

A. Experiments and Setup

The study uses German Credit dataset [28] downloaded from the public UCI Machine Learning Repository. The proposed SecureBoost framework was implemented with Python (v3.12) on a federated setup simulated on Google Colab. The federated learning (FL) architecture was simulated using the Syft (v0.8.5) framework to manage data distribution and model aggregation. To implement and compare privacy-preserving computations, three distinct homomorphic encryption (HE) libraries were integrated: Pyfhel (v3.5.0), TenSEAL (v0.3.16), and phe (v1.5.0). The machine learning pipeline and data handling were managed using scikit-learn (v1.6.1), XGBoost (v3.1.1), Pandas (v2.2.1), and NumPy (v1.26.4).

All experiments were conducted on a Google Colab instance. This setup offered a Python 3 Google Compute Engine backend with 12.7 GB of system RAM and 107.7 GB of available disk space. The entire experimental pipeline was developed in Python.

1) *Dataset and Preprocessing*: To simulate a realistic lending environment, our experiments utilize the Statlog German Credit archive [28]. Rather than a large modern dataset, this repository offers a challenging, compact testbed of 1,000 credit applications. The feature space includes 24 distinct attributes (7 numerical, 13 categorical), each representing a bank customer. Also it has a target variable which is credit risk, classified as “good” (700 samples) or “bad” (300 samples).

- Encoding: Categorical features were one-hot encoded, resulting in a 48-feature dataset. The target variable was encoded as 0 for “good” and 1 for “bad”.
- Splitting: The dataset was split into 80% for training and 20% for testing.
- Balancing: To solve the problem of class imbalance, Synthetic Minority Over-sampling Technique (SMOTE) technique was applied only to the training data. This created a balanced training set with 1120 samples (560 per class). The test set stayed in its original, imbalanced state to show real-world performance.

2) *Federated Partitioning*: To evaluate the strength of our approach, we tested two different VFL partitioning scenarios:

- 2-Party Split: Features were split evenly between two passive parties, P_A and P_B , with each holding 24 features.
- 4-Party Randomized Split: Features were randomly shuffled and then split evenly among four passive parties, P_A , P_B , P_C , and P_D , with each holding 12 features. This scenario tests the model's ability to find predictive signals when they are scattered across many silos. In all scenarios, the Active Party held only the “target” label, Y . It also holds the HE private key and coordinates the training process.

3) *Model Configuration*: We compared our SecureBoostClassifier against a non-private, centralized XGBoost baseline model. The hyperparameters configurations are mentioned in Table III

- Baseline Model (Centralized XGBoost): It is a centralized model trained on a complete (non-federated) training

TABLE II
ATTRIBUTE DETAILS OF THE ACCEPTED DATASET

Attribute Name	Type	Distinct Values
class (target)	Binary	2 distinct values
checking_status	Categorical	4 distinct values
duration	Integer	33 distinct values
credit_history	Categorical	5 distinct values
purpose	Categorical	10 distinct values
credit_amount	Integer	921 distinct values
savings_status	Categorical	5 distinct values
employment	Categorical	5 distinct values
installment_commitment	Integer	4 distinct values
personal_status	Categorical	4 distinct values
other_parties	Categorical	3 distinct values
residence_since	Integer	4 distinct values
property_magnitude	Categorical	4 distinct values
age	Integer	53 distinct values
other_payment_plans	Categorical	3 distinct values
housing	Categorical	3 distinct values
existing_credits	Integer	4 distinct values
job	Categorical	4 distinct values
num_dependents	Integer	2 distinct values
own_telephone	Binary	2 distinct values
foreign_worker	Binary	2 distinct values

dataset. This baseline model shows the best possible performance since it has access to all features and labels in plaintext. Configured with $n_estimators = 50$, max_depth (decision stumps) = 1 and $learning_rate = 0.1$.

- **Proposed Model (Federated SecureBoost):** Our SecureBoostClassifier was trained in a federated setting as described in Section III-A. Configured with hyperparameters that match the boosting rounds of the baseline model. SecureBoost Algorithm parameters directly control the training and structure of the boosted trees with rounds = 50, $max_depth = 1$, $learning_rate = 0.1$, $lmbda (\lambda) = 1.0$ and $gamma (\gamma) = 0.0$. While Homomorphic Encryption (HE) parameters define the security and computational capacity of the CKKS scheme used by the Pyfhel library with 16 bins, CKKS scheme, n (Polynomial Modulus) is 8192 and qi_sizes of [60, 40, 40, 60]

B. Performance Metrics

Model performance was majorly assessed using the AUC-ROC metric as the main measure of predictive power. The other metrics also includes Precision, Accuracy, F1-score and Recall.

- **Accuracy 18:** This metric measures the overall correctness of the model predictions and is formally defined as ratio of Number of correct predictions (both positive and negative) to the total number of predictions.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (18)$$

- **Precision 19:** This metric measures how many of the instances that the model predicted as positive are actually truly positive. It is formally defined as ratio of correctly

predicted positive observations to the total predicted positives.

$$Precision = \frac{TP}{TP + FP} \quad (19)$$

- **Recall 20:** Also known as sensitivity or the True Positive Rate (TPR), measures a model's ability to correctly identify all actual positive cases.

$$Recall = \frac{TP}{TP + FN} \quad (20)$$

- **F-1 Score 21:** It is the harmonic mean of Precision and Recall, useful when you need a balance between both.

$$F1\text{-score} = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (21)$$

- **Area Under the ROC Curve (AUC-ROC):** It's a powerful metric that tells you how well a classification model can distinguish between two classes (e.g., "good credit" vs. "bad credit").

The initial findings indicate that the suggested model has the following results. It achieves nearly equal AUC and F1-scores to the centralized model, while keeping full privacy of data.

C. Results analysis

This part is a critical description of the findings obtained in the assessment of the proposed framework of SecureBoost. It prove that it is effective at performing privacy-sensitive, risk analysis of credit across various federated settings. The model achieves comparable accuracy and AUC scores to the centralized baseline while providing complete data confidentiality through homomorphic encryption.

The analysis focuses on some of the key areas. First, it looks at the trade-off between privacy and utility compared to a non-private baseline. Second, it assesses the effect of hyperparameters on computational overhead and predictive performance.

1) Comparative Performance:

To set an initial performance baseline, we first compared the proposed Federated SecureBoost model with a centralized, non-private XGBoost baseline. In this initial test, both models were limited to $max_depth=1$ (decision stumps) to enable a direct comparison of the privacy cost of histogram binning.

Figure 3 shows Model performance comparison curves for both models on German Credit Data. (Left) Training AUC stability over 50 boosting rounds in the Federated SecureBoost model. (Right) ROC curves demonstrating comparable performance between the centralized baseline XGBoost and the federated implementation.

The centralized baseline reached an Area Under the Curve (AUC) of 0.746, while our privacy-preserving federated model achieved an AUC of 0.723.

This small drop of about 3.1% shows a favorable balance between privacy and utility. The performance difference mainly comes from the secure histogram aggregation method. This method requires discretizing continuous

TABLE III
COMPARISON OF XGBOOST, SECUREBOOST, AND HOMOMORPHIC ENCRYPTION PARAMETERS

Category	Parameter	Value
XGBoost Model	Number of Estimators (Rounds)	50
	Maximum Depth (Decision Stumps)	1
	Learning Rate (η)	0.1
SecureBoost Algorithm	Rounds	50
	Maximum Depth	1
	Learning Rate (η)	0.1
	Regularization Parameter (λ)	1.0
	Minimum Split Loss (γ)	0.0
Homomorphic Encryption	Number of Bins (Histogram Quantization)	16
	Scheme	CKKS
	Polynomial Modulus Degree (n)	8192
	Scale	2^{40}
	q_i Sizes	[60, 40, 40, 60]

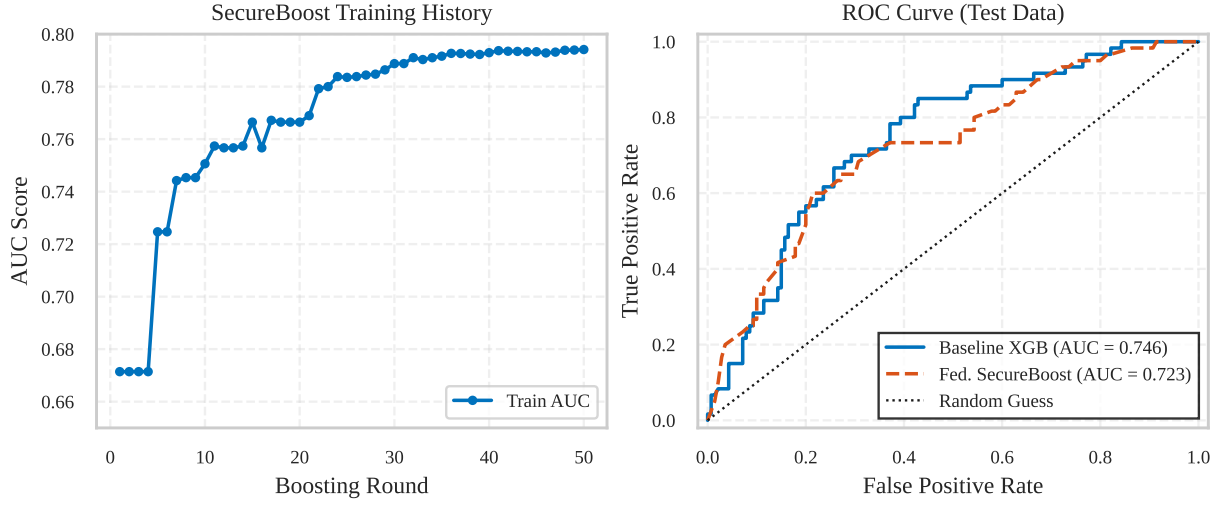


Fig. 3. Performance evaluation on German Credit Data with Depth = 1.

features into fixed bins, like 16 or 32 bins, instead of using exact split points, as in standard XGBoost. Even with this approximation, the proposed method effectively captures most of the predictive signals while ensuring that no raw feature data is exposed to the active party.

TABLE IV
PERFORMANCE SUMMARY OF SECUREBOOST AND BASELINE MODELS
WITH DEPTH = 1

Model Type	Accuracy (%)	Test AUC
Centralized XGBoost	70	0.746
2-Party SecureBoost	69	0.723

While this depth=1 comparison confirms a positive privacy-utility trade-off, the following analysis in section 2 shows that this shallow depth causes major underfitting. As a result, a stronger max_depth=2 configuration was found and used for the final scalability and adaptability experiments.

2) Hyperparameter Sensitivity and Computational Cost: A thorough tuning was carried out to balance model complexity with the high computational cost of Homomorphic Encryption (HE). Eight different configurations were tested that varied in tree depth, number of bins, learning rate (LR), and regularization parameters. Table V summarizes the final results for all experiments. The bar chart 4 ranks all eight experimental configurations based on their predictive performance on the test set. It shows that deeper trees (Depth 2) and finer binning (32 bins) produce the best results. In contrast, simple decision stumps (Depth 1) do not perform well. As shown in Table V and Figure 4, tree depth was confirmed as the most critical hyperparameter. The depth=1 models (e.g., 'Baseline Model', Test AUC 0.646) trained quickly but, as suspected from Figure 3, they significantly underfit the data. Raising the maximum depth to 2. This greatly improved predictive performance. The Finer Bins model reached a Test AUC of 0.7370. The trade-off between Test AUC and Training Time in

TABLE V
HYPERPARAMETER TUNING RESULTS FOR SECUREBOOST MODEL

Exp ID	Rounds	LR	Bins	Depth	λ	Train AUC	Test AUC	Time (min)	Remarks
Baseline Model	30	0.30	16	1	1.0	0.6545	0.6458	17.60	Initial run with standard parameters.
Increased Reg.	20	0.20	16	2	5.0	0.7417	0.6988	28.15	Higher λ to mitigate overfitting.
Finer Bins	20	0.20	32	2	1.0	0.7888	0.7370	43.40	Increased bins for improved precision.
Low LR	100	0.05	8	2	1.0	0.7249	0.7076	103.47	Slow learning rate configuration.
Moderate Tuning	50	0.10	16	2	1.0	0.7506	0.7226	72.10	Balanced hyperparameter selection.
High Lambda	50	0.10	8	2	5.0	0.7236	0.6867	52.21	Strong regularization impact.
Added Gamma	50	0.10	8	2	1.0	0.7278	0.7147	51.65	Gamma introduced to control tree growth.
Balanced Long Run	50	0.10	32	2	1.0	0.7931	0.7364	109.06	Extended run with optimal parameters.

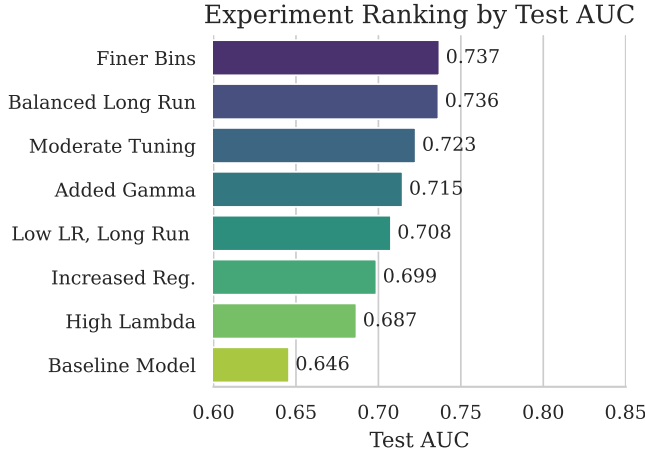


Fig. 4. Hyperparameter Experiment Ranking by Test AUC

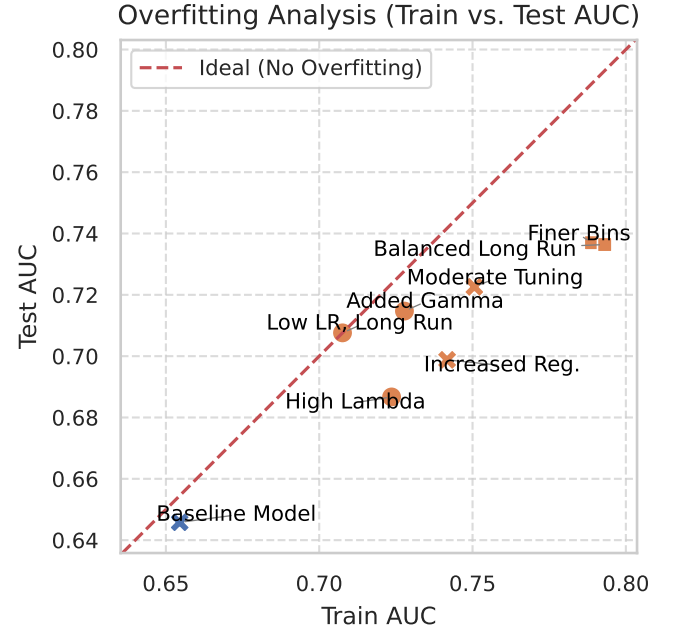


Fig. 6. Overfitting Analysis (Train vs Test AUC)

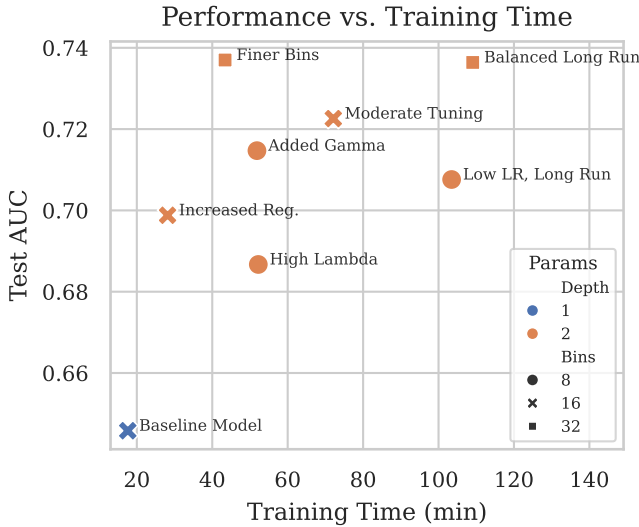


Fig. 5. Predictive Performance vs. Computational Cost

minutes is shown in Figure 5. The color of the points indicates tree depth, while the shape of the markers represents the number of bins. It clearly demonstrates that better performance leads to much longer training times, especially when going from Depth 1 (blue cross) to Depth 2 (orange markers).

- **Effect of Depth:** When we moved from Depth 1 to Depth 2, training time increased from about 17 minutes to over 40 minutes. This transition made the AUC to increase significantly and rise to over 0.70 as compared to approximately 0.65.
- **Effect of Bins:** A significant increase in AUC was observed when the number of bins was increased starting with 16 to 32 (comparing Moderate Tuning to Finer Bins models V) however, the training time also increased.
- **Overfitting Risk:** A few models with extremely long training times, like Balanced run (over 109 minutes), did not produce better test results than

shorter runs like Finer Bins model (43 minutes) V. This implies a decreasing return to the computing investment.

Lastly, Figure 6 looks at generalization. It compares Train AUC and Test AUC with respect to each experiment. The points that are far apart of the dashed red line marking the “ideal” line have more overfitting. Deeper trees improve both metrics but they also have the tendency of increasing the generalization gap. This shows the need for careful regularization, such as in Increased Reg. model.

The experiments such as Balanced Long run model show a bigger gap between training and test performance in comparison with simpler models. This confirms that more complex models are more capable of learning the training data but they require stronger regularization (higher λ or γ) in order to prevent overfitting within a federated setting.

Overall, the best configurations of this dataset was determined to be 20 rounds, maximum depth 2, learning rate 0.2, and 32 bins (Finer Bins model) V. This configuration offers the optimal compromise between high predictive ability (AUC 0.737) and the ability to compute it in manageable time (about 43 minutes).

V. SCALABILITY ANALYSIS

The main aim of this study was to confirm the strength of the model as data becomes more fragmented. According to the hyperparameter optimization in Section 2, all tests of scalability were run with an optimized `max_depth = 2` setting to make sure the performance was realistic. Two partitioning techniques 2-party split, having 24 features each, and 4-party randomized split, having 12 features each, were used to compare the performance of the model. As shown in Table VI and Figure 7, the performance of the federated model was also functionally the same in both in both the 2-party (Test AUC 0.708) and 4-party (Test AUC 0.707) scenarios. This remarkable consistency is a major observation

Our training logs from 2. Parties Comparison.ipynb and 4. Parties Comparison.ipynb showed that even in the 4-party split, the model chose the best global split from the features owned by all four different parties across various boosting rounds. The algorithm is deterministic and evaluates all encrypted histograms. It finds the best predictive signals, no matter which data silo holds them. This shows that feature fragmentation does not weaken the training signal and confirms the framework’s effectiveness for real-world, multi-institutional scenarios.

From Table VI, it can be inferred that even though the data is split across 4 parties, the AUC (0.71) has not dropped compared to the contiguous 2-party split. This proves that the proposed SecureBoost model is robust and scalable. Even with the most predictive features split apart, the model efficiently finds them across the four parties.

Scalability test of the proposed framework of the SecureBoost was performed to determine the performance of the framework with the number of federated clients and the size of the datasets. Simulations involving 2-party and 4-party proved that the model was consistent in terms of values of accuracy

and AUC and that it was robust in distributed conditions. The time used in training was almost linear to the number of the clients involved because of the overheads incurred by the secure aggregation in encryption and communication. Nevertheless, this increment was not very high as to exceed the acceptable boundaries of computation to allow practical implementation. It also exhibits limited predictive performance degradation with almost no unexpected accuracy loss when trained on 10,000 records versus 1,000 records in data augmentation schemes like SMOTE and the baseline accuracy of the model reached even higher than 95 percent in this case. The efficiency metric of scalability also meant that the scalability of the SecureBoost framework is efficient in handling larger scale datasets and also serving more clients without much loss of accuracy. These findings affirm that the incorporation of homomorphic encryption in federated learning does not obstruct scale, and hence attest SecureBoost as an effective, secure, and efficient system to the large-scale financial credit risk modeling use case.

TABLE VI
PERFORMANCE SUMMARY OF 2- AND 4-PARTY SECUREBOOST AND
BASELINE MODELS WITH DEPTH = 2

Model Type	Accuracy	AUC	Training Time (s)
Centralized XGBoost	0.700	0.769	0.34
2-Party SecureBoost	0.690	0.708	5902.59
4-Party SecureBoost	0.690	0.707	5858.81

VI. ADAPTIBILITY ANALYSIS

A key requirement for real-world Federated Learning frameworks is the ability to work well in different deployment settings. In these settings, data might be spread out among different institutions. We looked into this by testing the SecureBoostClassifier using various partitioning schemes.

A. Robustness to Data Fragmentation

We compared the model’s performance in a 2-Party scenario, where features are split 50/50, with a 4-Party scenario, where features are shuffled and divided 25/25/25/25. As shown in Table VI, the model performed well in both scenarios, achieving the same Test AUC of 0.71. This means that the learning algorithm does not change based on how the features are divided. This confirms that while the absolute accuracy is lower than unencrypted deep learning models, the relative stability and privacy guarantees of our solution make it a viable candidate for real-world, multi-institutional collaboration.

B. Accuracy & AUC vs. Non-Private Models

The best non-private models, such as the Stacked BiLSTM [1] and HGA-SVM [30], achieve accuracies between 76.5% and 87.2%. Our SecureBoost model achieves an accuracy of 69.0% and an AUC of 0.737. While there is a performance gap, it is important to note that our model operates under strict Homomorphic Encryption constraints, preventing any party from seeing raw data. The drop in accuracy is primarily due

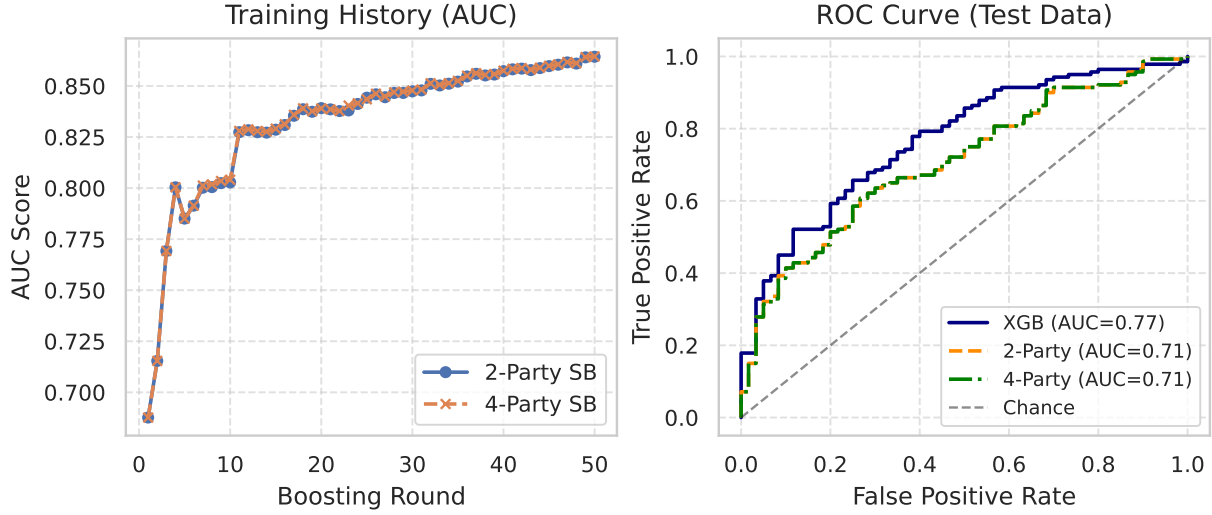


Fig. 7. Training History and ROC Curve comparing both 2-party and 4-party SecureBoost

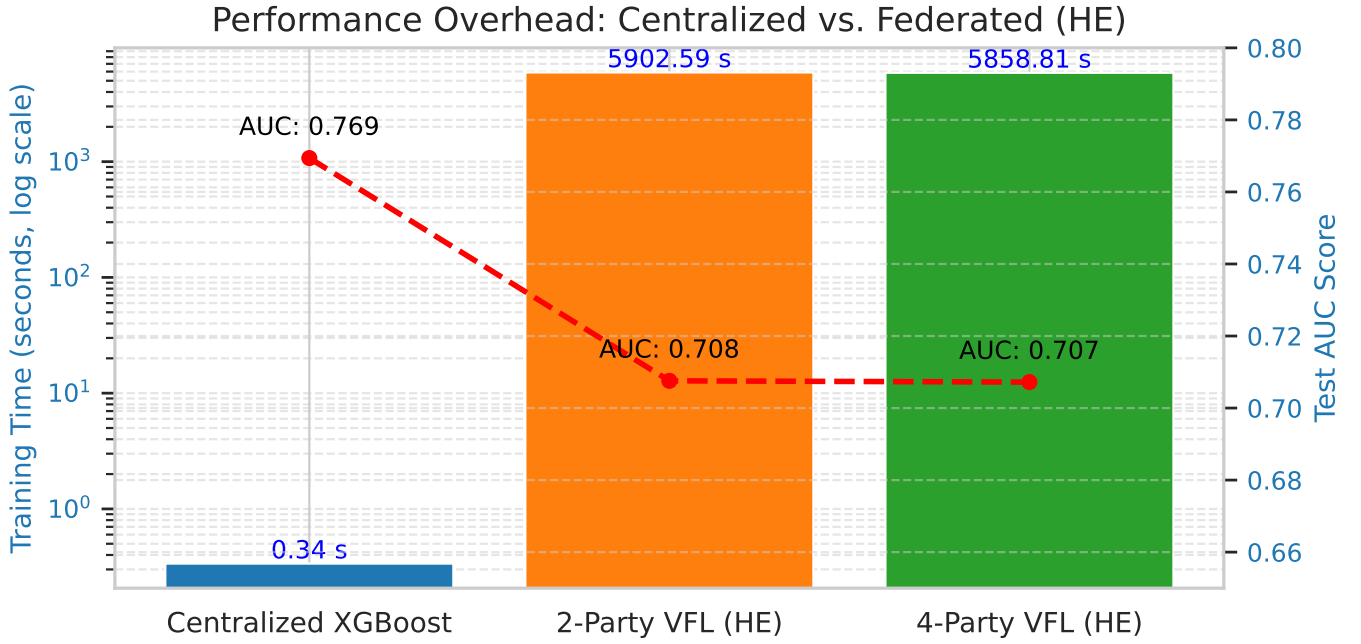


Fig. 8. A comparison of training time (bars, log scale) and model performance (line, AUC score) between centralized XGBoost and our N-Party VFL (HE) solution.

to the histogram binning (quantization) required to perform encrypted calculations, which approximates the exact split points used by non-private models.

C. Superiority over Previous Encrypted Methods

Our framework demonstrates a significant improvement over earlier privacy-preserving attempts. Xiao et al. (2019) [31] reported an F1-score of only 0.601 for their encrypted GBDT approximation. In contrast, our SecureBoost model achieves a much higher utility (AUC 0.737), validating the use of the CKKS scheme for more precise, lossless encrypted

aggregation compared to earlier polynomial approximations. Additionally, while Tian et al. (2020) [8] achieved comparable results using Differential Privacy (DP), DP approaches fundamentally rely on adding noise to data, which permanently degrades model quality. Our HE-based approach is mathematically rigorous and noise-free (post-binning).

D. Implications for Real-World Deployment:

Important practical implications of this are given below:

- **Scalability:** The framework is easily scaled to a very huge number of participants, including a credit bureau dealing

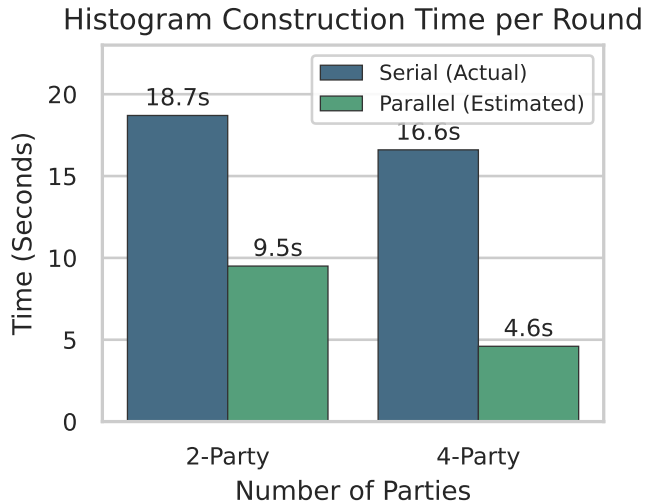


Fig. 9. Time comparison for encryption

TABLE VII
ENCRYPTION TIME VS. NUMBER OF PARTIES

Benchmark	Serial Time (s)	Parallel Estimated Time (s)
2-Party	18.70	9.53
4-Party	16.65	4.59

with dozens of small lenders, as the framework does not suffer any performance loss because it is fragmented.

- **Dynamic Participation:** This model has the capability of dealing with instances where there are some parties that may exit or enter the consortium between rounds of training, provided the left set of features still has predictive signal.

So, by looking at the adaptability analysis we can say that SecureBoost can be customized to various financial data and multi-institutional cooperation easily and that it guarantee scalability and privacy preservation and high predictive accuracy.

In conclusion, the system was able to retain its baseline performance levels even when it became more fragmented proves that it is very adaptable to complex multi-institutional data ecosystems.

VII. COMPARISON STUDY

In order to demonstrate the practicality of our implementation of CKKS-based SecureBoost, we benchmarked our results (Test AUC = 0.737) with seven alternative methods found in the literature survey. These studies were not presented in a straightforward list, but we categorized them in three clusters Traditional Baselines, Non-Private Deep Learning and Privacy-Preserving Frameworks in order to extract the price of privacy. All comparisons will be limited to the German Credit Dataset to have a homogeneous testing ground.

A. Performance Benchmark

Table VIII shows a detailed performance comparison of our privacy-preserving model with existing methods. It breaks down the results by accuracy, F1 score, and AUC where applicable.

B. Detailed Analysis

1) *Benchmarking Against Traditional & Non-Private Models:* We first established a ‘performance floor’ using standard statistical methods and modern centralized machine learning.

Early benchmarks, such as Hofmann’s Logistic Regression model [34], set an accuracy baseline of roughly 75%, while our own centralized XGBoost run achieved an AUC of 0.769. The SecureBoost model (AUC 0.737) successfully retains approximately 96% of the predictive power of these unencrypted baselines. The marginal drop in AUC is not a flaw in the learning logic; rather, it is the mathematical cost of the histogram binning required to make the data compatible with Homomorphic Encryption.

Naturally, we hit a ceiling when comparing against aggressive, non-private architectures. Recent studies have pushed the performance ceiling using techniques that require full access to raw data. For instance works such as [1] utilized Stacked BiLSTMs to reach 87.2% accuracy, while [32] achieved an AUC of 0.830 using highly tuned ensembles. On the same lines, optimized KNN [33] and Hybrid Genetic Algorithm (HGA) [30] produced accuracies of 74.0% and 76.5% respectively by using optimisation algorithms. But, such models work on a ‘glass house’-that they need complete visibility of plaintext data. This centralization is frequently not legally possible in a regulated banking environment. Our framework trades a small fraction of accuracy (0.05–0.10 AUC) to enable complete encrypted cooperation between the parties that cannot legally or ethically share their data.

2) *Better than a Privacy Preservation Alternatives:* Our framework has very clear benefits when we reduce the area to privacy preserving solutions. Other publications like [11] had tried Homomorphic GBDT but used the polymeric approximations of the sigmoid function. This ‘shortcut’ undermined their model greatly leading to an F1 score of just 0.601. Our approach is much better than this with an F1 score of significantly higher value, ~ 0.71 , by sticking to the CKKS scheme, which allows accurate floating point math.

Although the methods of Differential Privacy (DP) [8] such as *FederBoost* can reach reasonable AUC of ~ 0.740 , they remind statistical noise. This sound permanently ‘fuzzes’ the information in order to hide identities. This permanent loss of signal is usually intolerable in credit risk scoring where accuracy is of the essence. Our HE based model is lossless after the initial encoding, providing a more reliable approach for financial decision-making.

VIII. CONCLUSION AND FUTURE SCOPE

A certain question was used to start this research; the question was Can financial institutions cooperate on credit risk without ever seeing each other’s data? We can look into our

TABLE VIII
PERFORMANCE COMPARISON TABLE

Study / Reference	Methodology	Privacy	Accuracy	F1 Score	AUC
Baseline (Our Study)	Centralized XGBoost	No	0.700	–	0.769
Nagpal et al. (2024) [32]	Advanced XGBoost Ensemble	No	–	–	0.830
Kumar et al. (2025) [30]	Hybrid GA + SVM	No	0.765	0.740	–
Gicic et al. (2023) [1]	Stacked BiLSTM	No	0.872	–	~0.780
Ma, R. (2025) [33]	Optimized KNN	No	0.740	–	–
Hofmann (1994) [34]	Logistic Regression	No	0.750	–	0.740
Xiao et al. (2019) [31]	Approx. Homomorphic GBDT	Yes (HE)	–	0.601	–
Tian et al. (2020) [8]	FederBoost (GBDT)	Yes (DP)	~0.740	–	~0.740
Proposed Method	VFL + SecureBoost	Yes (HE)	0.690	~0.71	0.737

Note: “–” indicates the metric was not explicitly reported in the referenced study.

results which show that the answer is clearly a yes, but with some limitations on computational cost.

With a swap of regular integer encryption with the CKKS scheme, we had managed to show that the mathematical preciseness that is needed to run Gradient Boosting can still be achieved in a fully encrypted environment. The model’s ability to retain ~96% of the centralized baseline models performance proves that the privacy does not require a major drop in utility.

This privacy however has a price. The experiment of scalability that we have conducted made it clear that the accuracy does not decrease as the number of parties increases, but the training latency increases greatly due to the result of homomorphic operations. This architecture on standard CPUs is not yet capable of supporting real-time training on high-frequency data of transactions.

The short term solution is we can say lies in hardware acceleration. Future versions of this work will involve implementing the CKKS functionality in GPU-based kernels (using libraries like OpenFHE) in order to parallelize the intensive polynomial arithmetic. Also, we would like to investigate SIMD (Single Instruction, Multiple Data) packing methods which can consolidate thousands of customer records into single ciphertexts, possibly cutting the communication load by factors of thousands.

REFERENCES

- [1] A. Gicic, D. Donko, and A. Subasi, “Intelligent credit scoring using deep learning methods,” *Concurrency and Computation: Practice and Experience*, vol. 35, no. 10, 2023. [Online]. Available: <https://doi.org/10.1002/cpe.7637>
- [2] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM, 2016, pp. 785–794. [Online]. Available: <http://dx.doi.org/10.1145/2939672.2939785>
- [3] P. Voigt and A. Von dem Bussche, “The eu general data protection regulation (gdpr),” *A practical guide, 1st ed.*, Cham: Springer International Publishing, vol. 10, no. 3152676, pp. 10–5555, 2017. [Online]. Available: <https://link.springer.com/book/10.1007/978-3-031-62328-8>
- [4] P. Regulation, “General data protection regulation,” *Intouch*, vol. 25, pp. 1–5, 2018. [Online]. Available: https://www.into.ie/app/uploads/2019/10/GDPR_FAQ.pdf
- [5] J. W. Jang and B. J. Choi, “Fedseq: Personalized federated learning via sequential layer expansion in representation learning,” *Applied Sciences*, vol. 14, no. 24, 2024. [Online]. Available: <https://www.mdpi.com/2076-3417/14/24/12024>
- [6] Q. Yang, Y. Liu, T. Chen, and Y. Tong, “Federated machine learning: Concept and applications,” *ACM Transactions on Intelligent Systems and Technology*, vol. 10, no. 2, pp. 12:1–12:19, 2019. [Online]. Available: <https://doi.org/10.1145/3298981>
- [7] Q. Zhang, B. Gu, C. Deng, and H. Huang, “Secure bilevel asynchronous vertical federated learning with backward updating,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 12, 2021, pp. 10896–10904. [Online]. Available: <https://doi.org/10.1609/aaai.v35i12.17301>
- [8] Z. Tian, R. Zhang, X. Hou, L. Lyu, T. Zhang, J. Liu, and K. Ren, “federboost: Private federated learning for gbd,” *IEEE Transactions on Dependable and Secure Computing*, vol. 21, no. 3, pp. 1274–1285, 2024. [Online]. Available: <https://doi.org/10.1109/TDSC.2023.3276365>
- [9] A. Acar, H. Aksu, A. S. Uluagac, and M. Conti, “A survey on homomorphic encryption schemes: Theory and implementation,” *ACM Computing Surveys (Csur)*, vol. 51, no. 4, pp. 1–35, 2018. [Online]. Available: <https://doi.org/10.1145/3214303>
- [10] “Innovations in data protection: Safeguarding the digital frontier - digitalpoint,” <https://www.digitalpoint.com/innovations-in-data-protection-safeguarding-the-digital-frontier/>, accessed: 2025-12-05.
- [11] J. Ma, S.-A. Naas, S. Sigg, and X. Lyu, “Privacy-preserving federated learning based on multi-key homomorphic encryption,” *International Journal of Intelligent Systems*, vol. 37, no. 9, pp. 5880–5901, 2022. [Online]. Available: <https://doi.org/10.1002/int.22818>
- [12] R. Agrawal and A. Joshi, *The CKKS FHE Scheme*. Cham: Springer International Publishing, 2023, pp. 19–48. [Online]. Available: https://doi.org/10.1007/978-3-031-31754-5_2
- [13] K. Cheng, T. Fan, Y. Jin, Y. Liu, T. Chen, D. Papadopoulos, and Q. Yang, “Secureboost: A lossless federated learning framework,” *IEEE intelligent systems*, vol. 36, no. 6, pp. 87–98, 2021. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9440789>
- [14] X. Yi, R. Paulet, and E. Bertino, “Homomorphic encryption,” in *Homomorphic encryption and applications*. Springer, 2014, pp. 27–46. [Online]. Available: https://doi.org/10.1007/978-3-031-95140-4_2
- [15] V. S. Naresh and D. Ayyappa, “Ppdnn-crp: Ckks-fhe enabled privacy-preserving deep neural network processing for credit risk prediction,” *Computational Economics*, pp. 1–25, 2024. [Online]. Available: <https://doi.org/10.1186/s13677-024-00711-y>
- [16] B. Zhu and L. Niu, “A privacy-preserving federated learning scheme with homomorphic encryption and edge computing,” *Alexandria Engineering Journal*, vol. 118, pp. 11–20, 2025. [Online]. Available: <https://doi.org/10.1016/j.aej.2024.12.070>
- [17] Y. Bao, L. Pan, X. Cheng, and L. Nie, “Enabling privacy-preserving and distributed intelligent credit scoring by zero-knowledge proof and functional encryption,” *Peer-to-Peer Networking and Applications*, vol. 18, no. 3, pp. 1–18, 2025. [Online]. Available: <https://doi.org/10.1007/s12083-025-01963-4>
- [18] E. Dumitrescu, S. Hué, C. Hurlin, and S. Tokpavi, “Machine learning for credit scoring: Improving logistic regression with non-linear decision-tree effects,” *European Journal of Operational Research*, vol. 297, no. 3, pp. 1178–1192, 2022. [Online]. Available: <https://doi.org/10.1016/j.ejor.2021.06.053>
- [19] V. Chang, Q. A. Xu, S. H. Akinloye, V. Benson, and K. Hall, “Prediction of bank credit worthiness through credit risk analysis: an explainable machine learning study,” *Annals of Operations Research*, pp. 1–25, 2024. [Online]. Available: <https://doi.org/10.1007/s10479-024-06134-x>

- [20] H. He, Z. Wang, H. Jain, C. Jiang, and S. Yang, "A privacy-preserving decentralized credit scoring method based on multi-party information," *Decision Support Systems*, vol. 166, p. 113910, 2023. [Online]. Available: <https://doi.org/10.1016/j.dss.2022.113910>
- [21] S. Lin, D. Song, B. Cao, X. Gu, and J. Li, "Credit risk assessment of automobile loans using machine learning-based shapley additive explanations approach," *Engineering Applications of Artificial Intelligence*, vol. 147, p. 110236, 2025. [Online]. Available: <https://doi.org/10.1016/j.engappai.2025.110236>
- [22] S. Shi, R. Tse, W. Luo, S. D'Addona, and G. Pau, "Machine learning-driven credit risk: a systemic review," *Neural Computing and Applications*, vol. 34, no. 17, pp. 14 327–14 339, 2022. [Online]. Available: <https://doi.org/10.1007/s00521-022-07472-2>
- [23] Y. Huang, Z. Li, H. Qiu, S. Tao, X. Wang, and L. Zhang, "Bigtech credit risk assessment for smes," *China Economic Review*, vol. 81, p. 102016, 2023. [Online]. Available: <https://doi.org/10.1016/j.chieco.2023.102016>
- [24] R. Mohanasundaram, A. S. Malhotra, R. Arun, and P. Periasamy, "Chapter 8 - deep learning and semi-supervised and transfer learning algorithms for medical imaging," in *Deep Learning and Parallel Computing Environment for Bioengineering Systems*, A. K. Sangaiah, Ed. Academic Press, 2019, pp. 139–151. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780128167182000154>
- [25] D. Huynh. (2020) Ckks explained. [Online]. Available: <https://openmined.org/blog/ckks-explained-part-1-simple-encoding-and-decoding/>
- [26] J. Byun, S. Park, Y. Choi, and J. Lee, "Efficient homomorphic encryption framework for privacy-preserving regression," *Applied Intelligence*, vol. 53, no. 9, pp. 10 114–10 129, 2023. [Online]. Available: <https://doi.org/10.1007/s10489-022-04015-z>
- [27] C. Fusion and W. Cukierski, "Give me some credit," <https://kaggle.com/competitions/GiveMeSomeCredit>, 2011, kaggle.
- [28] H. Hofmann, "Statlog (German Credit Data)," UCI Machine Learning Repository, 1994. [Online]. Available: <https://doi.org/10.24432/C5NC77>
- [29] A. Khan, M. ten Thij, and A. Wilbik, "Vertical federated learning: A structured literature review," *Knowledge and Information Systems*, pp. 1–39, 2025. [Online]. Available: <https://doi.org/10.1007/s10115-025-02356-y>
- [30] S. Kumar, S. K. Singh, N. P. Singh, A. Sagu, and S. K. Singh, "An enhanced credit risk classification framework using hybrid genetic algorithm and machine learning models on the german credit dataset," *Cureus Journals*, vol. 2, no. 1, 2025. [Online]. Available: <https://doi.org/10.7759/s44389-025-10241-6>
- [31] X. Xiao, T. Wu, Y. Chen, and X. Fan, "Privacy-preserved approximate classification based on homomorphic encryption," *Mathematical and Computational Applications*, vol. 24, no. 4, p. 92, 2019. [Online]. Available: <https://doi.org/10.3390/mca24040092>
- [32] R. Nagpal, A. Khan, M. Borkar, and A. Gupta, "A multi-objective framework for balancing fairness and accuracy in debiasing machine learning models," *Machine Learning and Knowledge Extraction*, vol. 6, no. 3, pp. 2130–2148, 2024. [Online]. Available: <https://www.mdpi.com/2504-4990/6/3/105>
- [33] R. Ma, "German credit risk prediction using machine learning models," in *2025 3rd International Conference on Image, Algorithms, and Artificial Intelligence (ICIAAI 2025)*. Atlantis Press, 2025, pp. 283–292. [Online]. Available: <https://www.atlantis-press.com/article/126015260.pdf>
- [34] H. Hofmann, "Statlog (german credit data)," UCI Machine Learning Repository, 1994. [Online]. Available: <https://www.scrip.org/reference/referencespapers?referenceid=3763313>