

FHEMaLe: Framework for Homomorphic Encrypted Machine Learning

B. Pradeep Kumar Reddy¹,Sameeksha Goyal²,Ruchika Meel³, Ayantika Chatterjee⁴

pradeepkumarreddy.bukka@gmail.com, sameeksha.goyal.810@gmail.com, ruchikasingh758@gmail.com, cayantika@gmail.com

Indian Institute of Technology, Kharagpur, India-721302

Aim: The FHEMaLe framework enables privacypreserving machine learning (ML) on both edge and cloud environments using Fully Homomorphic Encryption (FHE). It addresses the privacy challenges of edge ML, where resource-constrained devices process sensitive data, and cloud ML, where data is outsourced to honest-but-curious servers. By performing computations on encrypted data, FHEMaLe ensures

Motivation

confidentiality while supporting scalable ML on diverse platforms.

- Edge and Cloud Synergy: Edge ML reduces latency and bandwidth needs, while cloud ML supports large-scale data processing. Both require robust privacy measures.
- Privacy Challenges: Edge devices face physical capture risks, and cloud servers are vulnerable to data breaches in honest-but-curious models.
- FHE Trade-offs: CKKS excels in approximate arithmetic for neural networks (NNs) but struggles with exact operations (e.g., comparisons in KNN, SVM). TFHE supports bit-level logic but is computationally intensive.

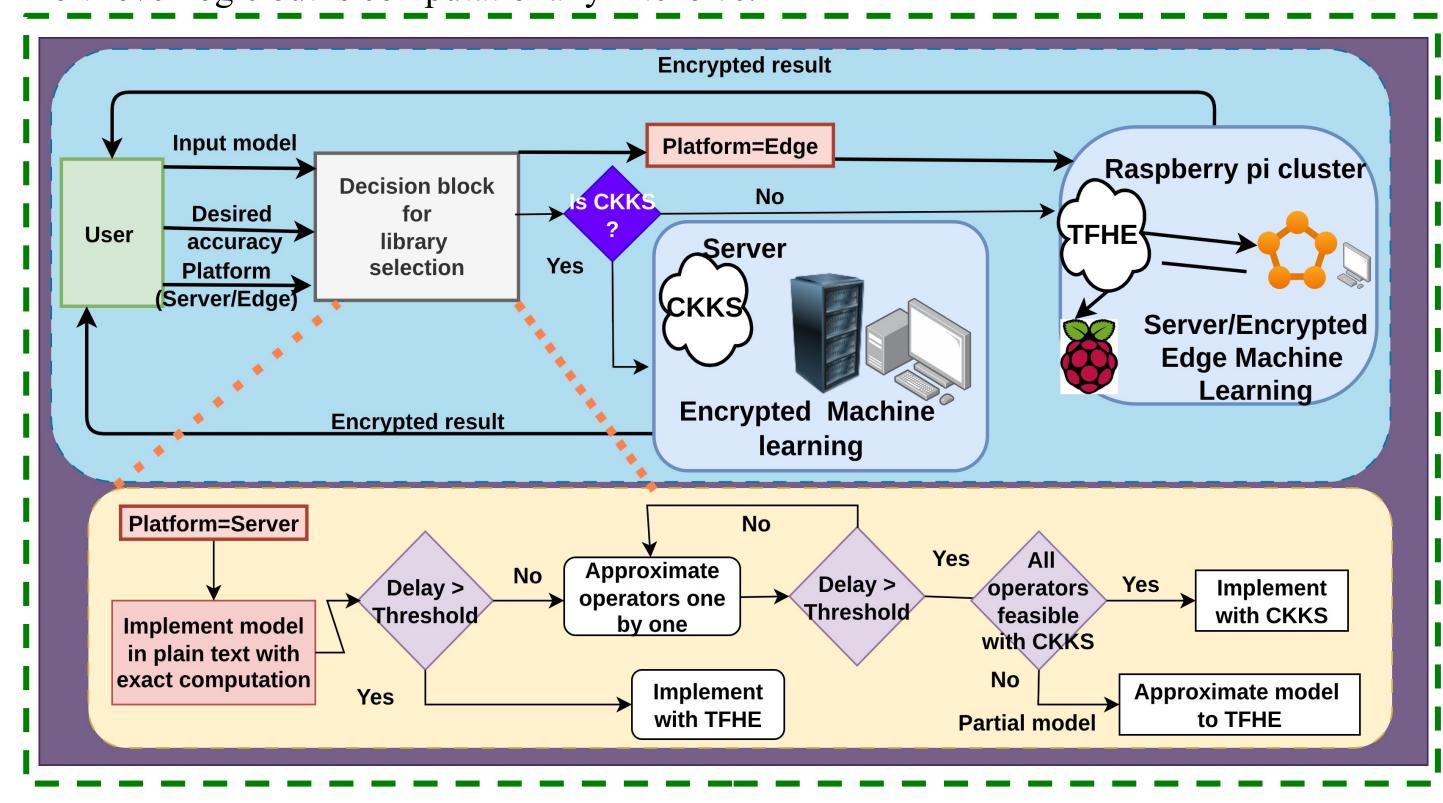


Fig. 1: Encrypted Machine Learnig as a Service

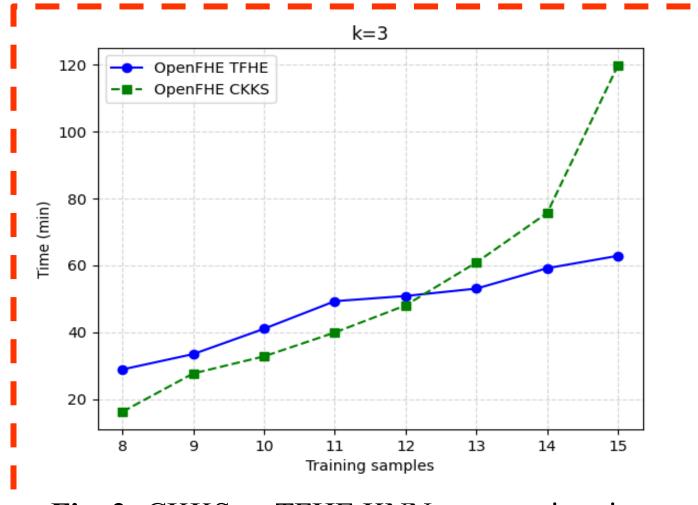


Fig. 2: CKKS vs TFHE KNN computation time on

Why KNN is Not Suitable in **Encrypted Domain?**

- KNN requires exact comparisons and distance-based sorting for classification, which are challenging
- CKKS enables homomorphic distance computations for KNN efficiently, secure and efficient neighbor selection (sorting of encrypted data) scales poorly remains a critical challenge.

cloud

- in FHE.

Enc_data Node_1: in n groups Enc_data Node_1: Node 2: ompute Partia ompute Parti mpute Partial DF (PDF)DF (PDF)DF (PDF)PDFNode2:

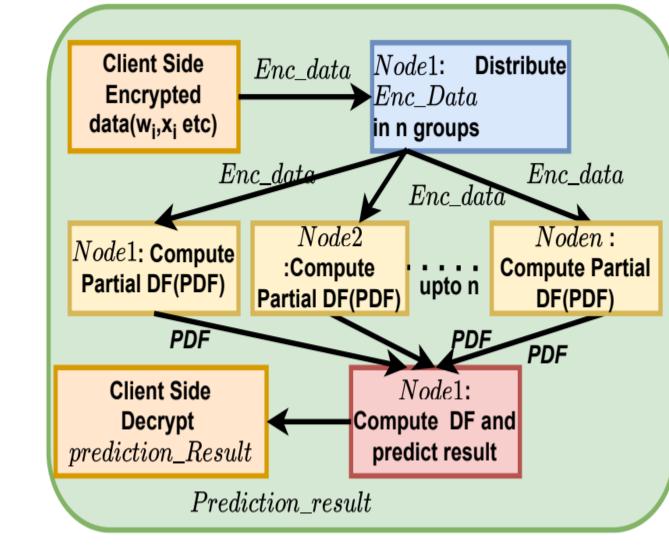


Fig 4: Encrypted Distributed LR Framework.

Fig 5: Encrypted Distributed SVM Framework.

Why CKKS is Suitable for Neural Networks?

CKKS is ideal for NNs due to its support for approximate arithmetic and SIMD operations:

- Vectorized Operations: CKKS efficiently handles matrix multiplications and dot products, core to NN layers (e.g., linear, attention heads).
- SIMD Packing: Enables parallel evaluation of multiple slots, improving throughput for highdimensional embeddings.
- Tolerable Approximations: NNs tolerate minor inaccuracies in operations like ReLU or sigmoid, unlike KNNs need for exact comparisons.
- Cloud Efficiency: CKKS achieves faster NN processing on high-resource servers.

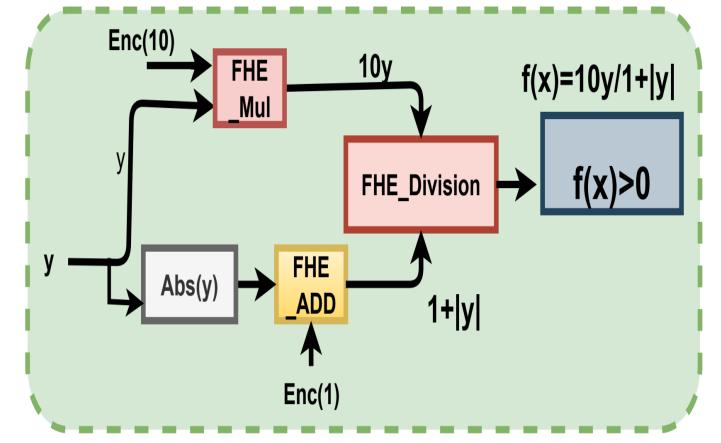
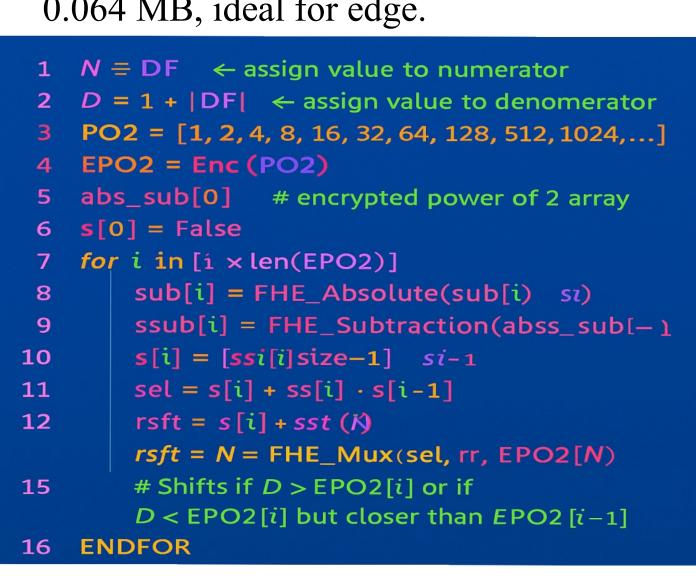


Fig 6: Fastsigmoid with Encrypted Gates.

Results and discussions

- Cloud Performance: On an HP Z240 (Intel Xeon, 64 GB RAM), CKKS achieves faster KNN for small datasets (248s vs. 780s for TFHE, 4 samples) but slows for larger datasets (7182s vs. 3768s, 15 samples).
- Edge Performance: Distributed processing on 11 Raspberry Pi nodes reduces KNN time (37 min), SVM (4.15 min), and LR (7.82 min) using OpenFHE.
- Memory Usage: CKKS requires 107 MB per element, limiting edge use; TFHE uses 0.064 MB, ideal for edge.





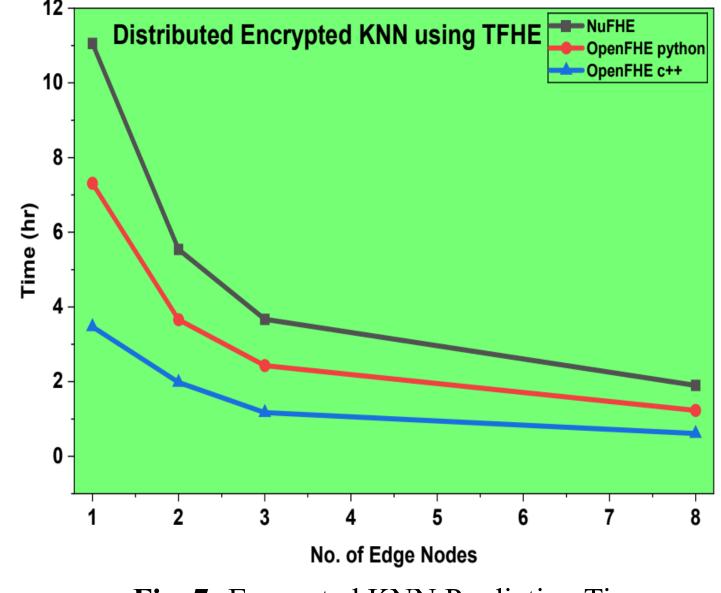
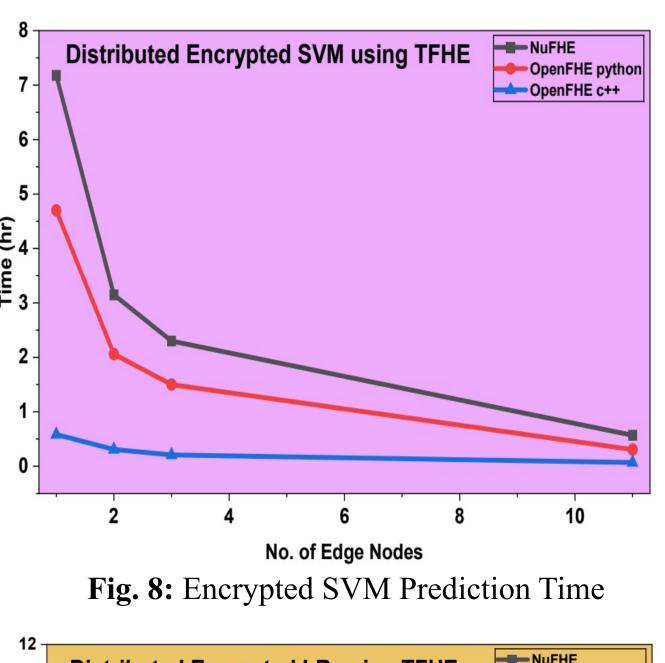


Fig. 7: Encrypted KNN Prediction Time

Table 1: Ckks vs TFHE Time (single prediction) and Memory per element (one element 16 bit) on cloud

ML Model	OpenFHE CKKS Time (sec)	OpenFHE TFHE Time (sec)	OpenFHE TFHE Memory (MB)	OpenFHE CKKS Memory (MB)
KNN	7182	3768	0.064	107
SVM	18	946	0.064	107
LR	17	650	0.064	107
NN	57	3982	0.064	107



Distributed Encrypted LR using TFHE OpenFHE python OpenFHE c++ No. of Edge Nodes

Fig. 9: Encrypted LR Prediction Time

Edge Nodes(ENs)

Methodology

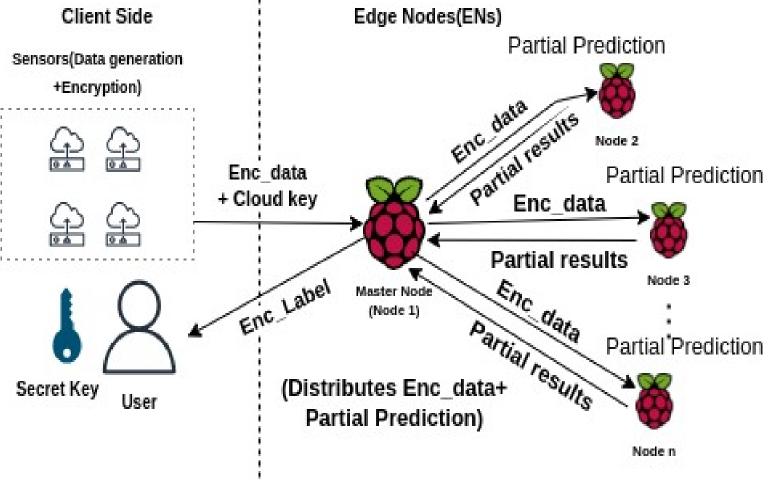


Fig. 3: Distributed Framework.

FHEMaLe evaluates ML models in plaintext to select CKKS or TFHE based on required operators and platform (edge or

cloud). Key components include:

- Library Selection: CKKS for SIMD-based operations (e.g., matrix multiplications in NNs) for exact in cloud environments; TFHE comparisons (e.g., KNN, SVM) on edge devices.
- Distributed processing for edge (Raspberry Pi clusters) and cloud (highresource servers) to optimize performance.
- Distributed Processing: Edge nodes compute partial decision functions (PDFs), aggregated by a master node (Node1). Cloud servers handle largescale encrypted computations.
- Encrypted Operations: FHE circuits (e.g., FHE Adder, FHE Mux) support secure dot products and label predictions.
- Novel encrypted approximate division algorithm for efficient sigmoid computation in logistic regression (LR).
- Evaluation and minimization of computational overhead introduced by encrypted data operations with distributed and concurrent computing on the edge devices network is also explored.

Future Directions

- > Optimize FHE circuits for additional ML algorithms (e.g., Decision Trees).
- Scale cloud and edge clusters for larger datasets and complex models.
- Develop hybrid CKKS-TFHE approaches for balanced performance.
- **Explore energy-efficient FHE for low-power** edge devices.

References

- A., & Sengupta, (2015).Chatterjee, algorithms Translating fully handle to homomorphic encrypted the cloud. IEEE **Transactions** Cloud Computing, 6(1), 287-300.
- 2. Reddy, B. P. K., Meel, R., & Chatterjee, A. (2024). Encrypted KNN Implementation on Distributed Edge Device Network. IACR Cryptol. ePrint Arch., 2024, 648.
- 3. Agrawal, R., & Joshi, A. (2023). On architecting fully homomorphic encryption-based computing systems. Springer International Publishing.
- 4. Pradeep Kumar Reddy, B., & Chatterjee, A. (2025). SMLaaS: Secure Machine Learning as a Service Ensuring Data and Model Parameter Privacy. Security and Privacy, 8(4), e70047.

Connect With Us

