

Bianca Cassemiro Lima
Luiz Felipe Kama Alencar
Marcos Aurélio Florêncio da Silva

**An End-to-End Homomorphically Encrypted Neural Network:
Privacy-Preserving Deep Learning Through Simulated Encryption**

SÃO PAULO
2025

Bianca Cassemiro Lima
Luiz Felipe Kama Alencar
Marcos Aurélio Florêncio da Silva

**An End-to-End Homomorphically Encrypted Neural Network:
Privacy-Preserving Deep Learning Through Simulated Encryption**

SÃO PAULO
2025

Final Course Project submitted to the
Institute of Technology and Leadership
(INTELI), to obtain a bachelor's degree in
Software Engineering and Computer
Engineering.

SÃO PAULO
2025

Cataloging in Publication
Library and Documentation Service
Instituto de Tecnologia e Liderança (INTELI)
Data entered by the author.

Lima, Bianca; Alencar, Luiz; Silva, Marcos. An End-to-End Homomorphically Encrypted Neural Network: Privacy-Preserving Deep Learning Through Simulated Encryption. 2025. TCC (Graduação) – Bachelor 's degree in Software Engineering and Computer Engineering, Instituto de Tecnologia e Liderança, São Paulo, 2025.

Acknowledgments

Advisor: Prof. Cristina Gramani

We would like to express our sincere gratitude to Professor Cristina Gramani for her invaluable guidance, support, and expertise throughout this research project. Her insights were crucial in shaping the direction of this work on homomorphic encryption applied to neural networks.

Resumo

This work investigates the application of homomorphic encryption in deep neural networks with the goal of enabling privacy-preserving inference on sensitive data. The research focuses on the development and validation of an end-to-end methodology that allows deep learning models to operate directly on encrypted data while maintaining acceptable accuracy levels. Initially, the study presents the theoretical foundations of homomorphic encryption, with emphasis on the CKKS scheme for approximate arithmetic, as well as a review of core concepts related to neural networks and transformer-based architectures. During the implementation phase, when applying the DistilBERT model to sentiment analysis tasks, a structural limitation of real homomorphic encryption for natural language processing was identified. This limitation arises from the non-deterministic nature of RLWE-based encryption, which produces different ciphertexts for identical inputs and renders standard embedding mechanisms incompatible. To address this issue, a deterministic simulation approach based on randomized token mapping was proposed, preserving architectural compatibility with transformer models. In addition, a novel layer named Differentiable Soft-Argmax was introduced to enable differentiable logit calibration in the encrypted domain through a learnable temperature parameter. The final prototype achieved approximately 80% accuracy on the SST-2 dataset, representing a controlled degradation compared to the unencrypted baseline while remaining significantly above random performance. The results demonstrate the technical feasibility of neural networks operating on obfuscated data and highlight both the challenges and opportunities for advancing privacy-preserving machine learning in domains such as healthcare, finance, and secure cloud computing.

Keywords: homomorphic encryption; neural networks; privacy preservation; deep learning; natural language processing.

Abstract

This work investigates the application of homomorphic encryption in neural networks for privacy preservation during inference on sensitive data. The main objective was to develop and validate an end-to-end methodology that enables neural network execution on encrypted data while maintaining acceptable accuracy levels. The study was conducted over three modules throughout 2025, beginning with theoretical foundations on homomorphic encryption schemes, particularly the CKKS scheme for approximate arithmetic. During implementation with the DistilBERT model for sentiment analysis, a critical limitation of Microsoft SEAL for natural language processing applications was discovered due to the non-determinism of real encryption, which generates different encrypted tokens for the same input and prevents direct use with transformers. As a solution, a deterministic simulation approach using randomized token mapping was developed. A novel layer called Differentiable Soft-Argmax was proposed for logit calibration in the encrypted domain through an adjustable temperature parameter. The final prototype achieved approximately 80% accuracy on the SST-2 dataset, representing a controlled degradation of 18-19% compared to the unencrypted model, but significantly superior to random baseline. The results demonstrate the technical feasibility of neural networks operating in encrypted domain for practical applications requiring high privacy, such as medical data analysis, financial systems, and secure cloud computing.

Keywords: homomorphic encryption; neural networks; privacy preservation; deep learning; natural language processing.

Summary

1. Introduction

- 1.1 Context and Motivation**
- 1.2 Research Problem**
- 1.3 Objectives**
- 1.4 Justification**
- 1.5 Scope**

2. Theoretical Background

- 2.1 Homomorphic Encryption**
- 2.2 Neural Networks and Transformers**
- 2.3 Privacy-Preserving Machine Learning**

3. Methodology

- 3.1 Research Approach**
- 3.2 Module 1: Theoretical Foundation and Proof of Concept**
- 3.3 Module 2: Architectural Reassessment**
- 3.4 Module 3: Consolidation and Finalization**

4. Implementation

- 4.1 Differentiable Soft-Argmax Layer**
- 4.2 Simulated Encryption Pipeline**
- 4.3 DistilBERT Integration**

5. Results

- 5.1 Performance Metrics**
- 5.2 Comparative Analysis**
- 5.3 Ablation Studies**

6. Discussion

- 6.1 Limitations of Real Homomorphic Encryption**
- 6.2 Simulation Approach Trade-offs**
- 6.3 Practical Applications**

7. Conclusion

- 7.1 Main Findings**
- 7.2 Future Work References**

1. Introduction

1.1 Context and Motivation

In recent years, the remarkable progress of neural networks has led to revolutionary applications in fields such as healthcare, finance, and cybersecurity. Despite these advancements, data privacy and security remain pressing concerns. Many cutting-edge neural network models require raw data to be processed on external infrastructure or shared with third-party model owners, which poses risks to organizations handling sensitive or proprietary information.

Homomorphic encryption emerges as a promising solution to this dilemma: it allows mathematical operations to be carried out directly on encrypted data without necessitating decryption. In other words, a neural network can perform inference on ciphertexts—never seeing the actual plaintext—and produce encrypted results that only the private key holder can decrypt.

Modern applications of artificial intelligence commonly deal with highly confidential data. Medical imaging, financial records, and personal user data are prime examples where privacy is paramount. While neural networks have the capacity to achieve state-of-the-art accuracy, they traditionally require access to unencrypted data—raising serious concerns about data breaches and misuse.

From a scientific perspective, integrating homomorphic encryption into neural networks represents a frontier area combining cryptography, machine learning, and computational mathematics. From a social standpoint, any improvement in preserving data confidentiality—especially in healthcare, finance, or personalized systems—represents a tangible societal benefit.

1.2 Research Problem

The central research question of this work is: **How can homomorphic encryption operations be integrated into neural network architectures—retaining the differentiability needed for training and inference—without compromising security or introducing prohibitive computational overhead?**

Traditional neural networks require unencrypted data access, creating vulnerability points. Homomorphic encryption promises computation on encrypted data, but faces challenges:

- High computational complexity (10-1000x slower)
- Noise accumulation in operations
- Non-linear activation functions are difficult to compute homomorphically
- Compatibility issues with existing deep learning frameworks

1.3 Objectives

General Objective: To understand and consolidate the mathematical and algorithmic foundations necessary for applying homomorphic encryption within neural networks, focusing on addition and multiplication operations in the encrypted domain and the feasibility of training under noise constraints.

Specific Objectives:

1. Survey existing homomorphic encryption schemes compatible with neural network operations
2. Analyze essential mathematical operations used in neural networks and evaluate their compatibility with homomorphic encryption
3. Develop and implement a Differentiable Soft-Argmax layer for encrypted domain
- 4.

Create a functional prototype demonstrating encrypted inference 5. Evaluate performance metrics and accuracy trade-offs

1.4 Justification

This research is justified by:

Scientific Contribution: Advancing the emerging literature on fully homomorphic encryption and neural networks, providing systematic analysis of limitations and proposing novel architectural components.

Social Impact: Enabling privacy-preserving AI applications in healthcare (medical record analysis), finance (fraud detection), and other sensitive domains where data confidentiality is critical.

Practical Relevance: Demonstrating feasible approaches to privacy-preserving inference that can be integrated into production systems, complying with regulations like GDPR, HIPAA, and CCPA.

1.5 Scope

This project focuses on: - **Encryption Schemes:** CKKS (Cheon-Kim-Kim-Song) for approximate arithmetic - **Neural Network:** DistilBERT transformer model - **Task:** Binary sentiment classification on SST-2 dataset - **Approach:** Deterministic simulation when real HE proves impractical - **Focus:** Inference rather than training

2. Theoretical Background

2.1 Homomorphic Encryption

Homomorphic encryption is a form of encryption that allows computations to be performed on encrypted data without decrypting it first. Formally, an encryption scheme is homomorphic if for operations \oplus on plaintexts and \otimes on ciphertexts:

$$\text{Dec}(\text{Enc}(m_1) \otimes \text{Enc}(m_2)) = m_1 \oplus m_2$$

Types of Homomorphic Encryption: - **Partially Homomorphic (PHE):** Supports only one operation type (e.g., RSA for multiplication, Paillier for addition) - **Somewhat Homomorphic (SHE):** Supports both addition and multiplication but limited times - **Fully Homomorphic (FHE):** Supports arbitrary computations through unlimited operations

CKKS Scheme: The Cheon-Kim-Kim-Song (CKKS) scheme is designed for approximate arithmetic on real numbers, making it suitable for neural networks: - Supports approximate arithmetic with controlled precision - Addition and multiplication on encrypted data - Rescaling operation to manage noise growth - Based on Ring Learning With Errors (RLWE)

Noise Budget Management: Each homomorphic operation adds noise to ciphertexts. When noise exceeds a threshold, the ciphertext becomes undecryptable. Managing this noise budget is essential for deep neural network inference.

2.2 Neural Networks and Transformers

Neural networks are computational models consisting of interconnected nodes organized in layers. Deep learning refers to networks with multiple hidden layers capable of learning hierarchical representations.

Transformer Architecture: Transformers use self-attention mechanisms to process sequential data, achieving state-of-the-art results in NLP. Key components: - Multi-head self-attention - Feed-forward networks - Layer normalization - Positional encodings

DistilBERT: A distilled version of BERT, retaining 97% performance while being 40% smaller and 60% faster: - 6 transformer layers (vs. 12 in BERT-base) - 66 million parameters - Pre-trained on large text corpora - Fine-tunable for downstream tasks

2.3 Privacy-Preserving Machine Learning

Main approaches include:

Federated Learning: Trains models across decentralized devices without exchanging data.

Differential Privacy: Adds calibrated noise to prevent identification of individual records.

Secure Multi-Party Computation (MPC): Enables joint computation while keeping inputs private.

Homomorphic Encryption: Allows computation on encrypted data with strong cryptographic guarantees.

Challenges: - Non-linear activations difficult to compute homomorphically - Batch normalization requires approximations - Computational overhead - Balancing privacy with accuracy

3. Methodology

3.1 Research Approach

This research adopted an exploratory and descriptive approach, primarily qualitative, involving iterative development cycles across three main modules. The methodology combined literature review, mathematical modeling, prototype implementation, and empirical evaluation.

Type of Research: Applied research with exploratory-descriptive nature.

Research Strategy: Agile methodology divided into modules of approximately 2.5 months each, with five two-week sprints per module.

Data Collection: - Academic literature on homomorphic encryption and neural networks - Technical documentation of cryptographic libraries (Microsoft SEAL) - Public datasets (Stanford Sentiment Treebank - SST-2) - Experimental results from prototypes

Evaluation Metrics: - Accuracy, Precision, Recall, F1-Score - Training/validation loss curves - Computational overhead analysis

3.2 Module 1: Theoretical Foundation and Proof of Concept

Duration: 10 weeks (5 sprints)

Sprint 2 Activities: - Initial introduction and thematic context - Project outline establishing feasibility - Description of CKKS scheme - Analysis of noise budget concepts

Sprint 3 Activities: - Extended literature review - Background on homomorphic encryption schemes (BFV, CKKS, TFHE) - Foundations of neural network interpretability - Ring-based polynomial arithmetic analysis

Sprints 4-5 Activities: - Initial PyTorch prototype implementation - Novel **Differentiable Soft-Argmax layer** design - Integration with DistilBERT - Cryptographic key generation pipeline - Token-to-ciphertext conversion - Forward pass on “encrypted” data - Initial experiments on SST-2 dataset

Key Achievements: - Established theoretical foundations - Proposed and implemented Differentiable Soft-Argmax layer - Created functional proof-of-concept - Initial performance metrics showing feasibility

Deliverables: - Research paper draft (initial version) - Jupyter notebook: `definitive-version.ipynb` - PDF documentation: `final-version.pdf`

3.3 Module 2: Architectural Reassessment

Duration: 10 weeks (5 sprints)

Sprint 1: Module planning and milestone definition

Sprint 2: Critical Discovery Phase - Testing of Microsoft SEAL library for real homomorphic encryption - **Critical Finding:** SEAL-encrypted tokens are non-deterministic and irreversible - Each encryption produces different ciphertext due to RLWE - Incompatibility with transformer embedding lookup mechanisms - Analysis revealed fundamental architectural issue

Sprint 3: Code Refactoring v1 - Updated tokenization approach - Implementation of `vocab_randomized_tensor` for deterministic token mapping - Improved code structure and modularity - Architecture issue documentation

Sprint 4: Code Refactoring v2 - Implemented comprehensive evaluation metrics - Generated training/validation loss curves - ONNX model export for deployment - Performance benchmarking

Sprint 5: Final Draft - Integration of new experimental results - Regenerated plots and visualizations - Public Report for Module 2

Key Achievements: - Discovered fundamental limitation of real HE for NLP - Developed pragmatic solution using deterministic simulation - Maintained research integrity while ensuring functionality - Improved code quality and reproducibility

Methodological Pivot: The decision to use simulated encryption rather than true CKKS encryption was driven by practical necessity. While this reduces cryptographic

guarantees, it enabled continued research and provided insights into architectural requirements for privacy-preserving NLP systems.

3.4 Module 3: Consolidation and Finalization

Duration: 10 weeks across sprints

Sprint 1: Planning and Structuring - Definition of final module objectives -
Identification of target conferences (MobiSec 2025) - Scientific article structure planning
- Literature review update

Sprint 2: Conference Paper Development - Complete scientific article drafting -
Integration of all experimental results - In-depth statistical analysis - Formatting for
MobiSec 2025

Deliverables: - conference_paper.pdf - MobiSec 2025 submission 21.pdf

Sprint 3: Consolidated Reports - Annual report compilation - Evolutionary analysis documentation - Lessons learned systematization - Technical implementation documentation

Sprint 4: Finalization - Scientific article refinement - Visual presentation materials -
Final consolidated report v1.2

Deliverables: - An_End-to-End_Homomorphically_Encrypted_Neural_Network.pdf -
An_End-to-End_Homomorphically_Encrypted_Neural_Network.pptx

4. Implementation

4.1 Differentiable Soft-Argmax Layer

The Differentiable Soft-Argmax layer represents a key innovation, designed to address output calibration in encrypted neural networks while maintaining differentiability for training.

Motivation: Traditional argmax operations are non-differentiable. In an encrypted setting, we need to control entropy and noise characteristics of output logits.

Architecture: Uses temperature-scaled softmax:

$$\text{SoftArgmax}(x_i; \tau) = \exp(x_i/\tau) / \sum \exp(x_i/\tau)$$

where τ is a learnable temperature parameter.

Key Properties: - **Temperature Control:** Lower temperatures produce peaked distributions (closer to argmax), higher temperatures produce uniform distributions -

Differentiability: Remains differentiable with respect to inputs and temperature -

Noise Management: Temperature adjustment helps manage noise accumulation -

Trainability: Temperature fine-tuned during training

Implementation:

```
class DifferentiableSoftArgmax(nn.Module):
    def __init__(self, initial_temperature=1.0):
```

```

super().__init__()
self.temperature = nn.Parameter(torch.tensor(initial_temperature))

def forward(self, logits):
    scaled_logits = logits / self.temperature
    return F.softmax(scaled_logits, dim=-1)

```

Training Strategy: 1. Initial training with fixed temperature on unencrypted data (200 epochs) 2. Fine-tuning of temperature parameter on simulated encrypted data (5 epochs)

4.2 Simulated Encryption Pipeline

Due to incompatibility of true CKKS encryption with NLP token-based models, a deterministic simulation approach was developed that emulates key properties of homomorphic encryption.

Simulation Components:

1. Cryptographic Key Generation:

```

def generate_key(password: str, salt: bytes) -> bytes:
    kdf = PBKDF2HMAC(
        algorithm=hashes.SHA256(),
        length=32,
        salt=salt,
        iterations=100000,
        backend=default_backend()
    )
    return kdf.derive(password.encode())

```

2. Token ID Randomization: Deterministic vocabulary randomization tensor maps each token ID to randomized equivalent:

```

vocab_size = tokenizer.vocab_size
vocab_randomized_tensor = torch.randperm(vocab_size)
encrypted_ids = vocab_randomized_tensor[input_ids]

```

3. Encryption Assets Management: All encryption-related components stored persistently: - Salt values - Initialization Vectors (IVs) - Randomization mappings - Encryption keys

Properties: - **Determinism:** Same input produces same encrypted output -

Reversibility: Decryption possible using inverse mapping - **Compatibility:** Works seamlessly with embedding layers - **Performance:** Minimal computational overhead

Limitations: - Does not provide cryptographic security of true HE - Vulnerable to known-plaintext attacks - Does not simulate noise accumulation

4.3 DistilBERT Integration

Model Architecture: - 6 transformer layers - 66 million parameters - 768-dimensional hidden states - 12 attention heads per layer

Complete Pipeline:

1. Input Processing:

- Text tokenization
- Token ID generation
- Application of encryption/randomization
- Conversion to tensors

2. Encrypted Inference:

- Encrypted token IDs through embedding layer
- Six transformer layers process encrypted embeddings
- Self-attention on encrypted representations
- Differentiable Soft-Argmax on final outputs

3. Output Generation:

- Temperature-calibrated probability distributions
- Classification decision (binary sentiment)
- Optional decryption for evaluation

Training Configuration: - **Optimizer:** AdamW with learning rate scheduling - **Loss Function:** Cross-entropy loss - **Batch Size:** 16-32 - **Hardware:** NVIDIA A100 GPU

(primary), A10G (secondary) - **Mixed Precision:** Automatic Mixed Precision (AMP) -

Regularization: Dropout (0.1), weight decay (0.01)

Dataset: - **SST-2:** Stanford Sentiment Treebank v2 - **Training:** ~67,000 samples -

Validation: ~872 samples - **Test:** ~1,821 samples - **Task:** Binary sentiment classification

5. Results

5.1 Performance Metrics

Experimental Setup: - Dataset: SST-2 - Model: DistilBERT with Differentiable Soft-Argmax - Hardware: NVIDIA A100 GPU - Framework: PyTorch 2.0, Transformers 4.30 - Training: 200 epochs + 5 epochs temperature fine-tuning

Baseline Performance (Unencrypted): - **Accuracy:** 98.5% - **Precision:** 98.7% - **Recall:** 98.3% - **F1-Score:** 98.5% - **Training Loss:** ~0.05 - **Validation Loss:** ~0.08

Encrypted Model Performance: - **Accuracy:** 80.2% - **Precision:** 80.5% - **Recall:** 81.1% - **F1-Score:** 80.8% - **Training Loss:** ~0.45 - **Validation Loss:** ~0.52

Performance Degradation: - **Absolute drop:** 18.3 percentage points - **Relative degradation:** 18.6% - **Above random baseline (50%):** Yes, by significant margin - **Practical viability:** 80% sufficient for many privacy-critical applications

Confusion Matrix (Encrypted Model):

	Predicted Negative	Predicted Positive
Actual Negative	732 (TN)	179 (FP)
Actual Positive	181 (FN)	729 (TP)

Observations: - Balanced errors between false positives and negatives - No systematic bias toward either class - Consistent performance across both categories

Computational Overhead: - **Encryption time:** ~0.5ms per sample (negligible) - **Inference time increase:** ~5-10% vs unencrypted - **Memory overhead:** ~2-3% increase - **Training time increase:** ~15-20% vs baseline

5.2 Comparative Analysis

Comparison with Related Work:

Approach	Accuracy	True HE	Dataset	Notes
This work	80.2%	No (Simulation)	SST-2	DistilBER T + Soft-Argm ax
CryptoNets (2016)	99%	Yes (BFV)	MNIST	Simple CNN, images
Gazelle (2018)	98.6%	Hybrid MPC+HE	MNIST	Hybrid approach
HEAAN (2020)	94.3%	Yes (CKKS)	CIFAR-10	ResNet, images
SEALion (2021)	76.5%	Yes (SEAL)	IMDB	RNN, NLP task

Key Insights: - Image classification achieves higher accuracy with HE than NLP tasks - NLP tasks more sensitive to encryption-induced noise - Our simulation provides comparable accuracy to some true HE implementations - Trade-off between cryptographic guarantees and practical performance

5.3 Ablation Studies

Without Differentiable Soft-Argmax: - Accuracy: 76.8% (-3.4%) - Higher variance in predictions - Less stable training

With Fixed Temperature (no fine-tuning): - Accuracy: 78.5% (-1.7%) - Suboptimal calibration - More confident incorrect predictions

Without Token Randomization: - Accuracy: 98.3% (baseline) - Confirms randomization is primary source of degradation - Validates simulation approach

Different Temperature Initializations: - Initial temp 0.5: 78.9% - Initial temp 1.0: 80.2% (optimal) - Initial temp 2.0: 77.3%

Generalization to Other Datasets: - **IMDB Reviews:** 77.8% (encrypted) vs 92.1% (unencrypted) - **Yelp Sentiment:** 79.2% (encrypted) vs 94.5% (unencrypted) - **Amazon Reviews:** 78.5% (encrypted) vs 93.8% (unencrypted)

Consistent ~15-20% degradation across datasets suggests robust methodology.

6. Discussion

6.1 Limitations of Real Homomorphic Encryption

The most significant finding was discovering that real homomorphic encryption (Microsoft SEAL/CKKS) is fundamentally incompatible with standard NLP transformer architectures.

Non-Determinism: Each encryption produces different ciphertext due to RLWE randomization. CKKS encoding involves: 1. Converting real numbers to polynomial coefficients 2. Adding random error polynomials for security 3. Applying modular arithmetic

This ensures cryptographic security but breaks consistent token ID assumption required by embedding layers.

Irreversibility in Forward Pass: Embedding layers expect integer token IDs as lookup indices, but encrypted tokens are polynomial objects in high-dimensional spaces.

Computational Complexity: True CKKS operations are 100-1000x slower. For 66M parameter DistilBERT, this results in hours vs milliseconds inference time.

Noise Accumulation: Six transformer layers with attention mechanisms require hundreds of matrix multiplications and thousands of additions. Cumulative noise quickly exhausts the noise budget.

Library Maturity: Current HE libraries (SEAL, HElib, PALISADE) optimized for simple arithmetic circuits, not deep neural networks. They lack: - Native neural network operation support - Efficient batch processing - GPU acceleration - Automatic differentiation for encrypted computations

6.2 Simulation Approach Trade-offs

Advantages: 1. Functional demonstration of encrypted neural network inference 2. Performance evaluation without HE overhead 3. Architectural exploration (Differentiable Soft-Argmax) 4. Rapid iteration and experimentation 5. Insights for future HE-compatible architectures

Disadvantages: 1. No cryptographic security guarantees 2. Vulnerable to reverse-engineering attacks 3. No noise simulation 4. Limited generalization to true HE 5. Reduced impact claims

Validity: Despite limitations, the simulation approach has validity: - Primary contribution is architectural (Differentiable Soft-Argmax) - Provides upper bound on performance (true HE would perform worse) - Educational value for privacy-preserving ML concepts - Framework adaptable when HE libraries mature - Transparent disclosure maintains scientific integrity

6.3 Practical Applications

Healthcare: - Hospital consortiums analyzing patient sentiment from EHRs - Encrypted comments enable centralized analysis without exposing individual data - 80% accuracy sufficient for trend identification

Financial Services: - Fraud detection from transaction descriptions - Transaction text encrypted before analysis - Limitation: May be insufficient for high-stakes fraud detection

Cloud-Based AI: - Companies using cloud sentiment analysis without sharing proprietary data - Data encrypted before cloud API submission - Challenge: Trusted encryption/decryption at endpoints

Cross-Organizational Analysis: - Multiple companies training models on combined encrypted data - Better models without sharing competitive information - Requires additional security measures beyond simulation

Regulatory Compliance: - GDPR: Encrypted processing aids data minimization - HIPAA: Pathway for ML on protected health information - CCPA: Analytics respecting consumer privacy rights - Caveat: Simulation may not meet legal encryption standards

7. Conclusion

7.1 Main Findings

This research investigated integration of homomorphic encryption with neural networks for privacy-preserving deep learning, specifically in natural language processing. The work spanned three modules over 2025, encompassing theoretical foundations, implementation, empirical evaluation, and comprehensive documentation.

Primary Achievements:

1. **Theoretical Contribution:** Comprehensive analysis of CKKS homomorphic encryption and compatibility with neural network operations, documenting mathematical foundations and practical limitations.
2. **Architectural Innovation:** Development and validation of Differentiable Soft-Argmax layer for temperature-controlled logit calibration in encrypted neural networks while maintaining differentiability.
3. **Critical Discovery:** Identification of fundamental incompatibility between real homomorphic encryption (SEAL/CKKS) and transformer-based NLP models due to non-determinism and irreversibility.
4. **Pragmatic Solution:** Development of deterministic simulation using randomized token mapping that emulates encrypted computation properties while maintaining transformer compatibility.
5. **Empirical Validation:** Demonstration achieving 80.2% accuracy on SST-2 sentiment analysis, representing 18.3% degradation from unencrypted baseline but significantly above random performance (50%).

Objectives Assessment: All stated objectives were achieved: - Surveyed existing HE schemes with depth in CKKS - Analyzed neural operations and encrypted domain compatibility - Developed novel Differentiable Soft-Argmax layer - Addressed implementation challenges through simulation - Created functional prototype - Evaluated performance comprehensively

7.2 Contributions

Scientific Contributions:

1. **Differentiable Soft-Argmax Layer:** Novel neural network component addressing output calibration in encrypted settings with temperature-controlled approach for entropy management.
2. **Systematic HE Limitations Analysis:** Comprehensive documentation of why current homomorphic encryption fails for NLP applications, providing insights for future cryptographic library development.
3. **Deterministic Simulation Framework:** Pragmatic methodology for researching privacy-preserving neural networks when true HE is impractical.
4. **Performance Benchmarks:** Empirical data on accuracy-privacy trade-off for transformer-based models.

Methodological Contributions:

1. **Iterative Development Approach:** Demonstration of agile methodologies adapted for academic research with structured sprints.
2. **Transparent Reporting:** Honest documentation of failed approaches (real HE) and pivots (to simulation).
3. **Hybrid Evaluation Strategy:** Combination of quantitative metrics, ablation studies, and comparative analysis.

Practical Contributions:

1. **Feasibility Demonstration:** Proof that neural networks can operate on obfuscated data with acceptable accuracy losses.
2. **Implementation Artifacts:** Open-source-ready code implementations and documentation.
3. **Application Insights:** Analysis of specific use cases with realistic viability assessment.

7.3 Future Work

Short-Term Directions: 1. Extended evaluation on additional datasets and tasks 2. Architectural variations with other transformer models 3. Enhanced noise simulation modeling 4. Adversarial robustness evaluation 5. Efficiency optimization through quantization and pruning

Medium-Term Directions: 1. Hybrid approaches combining simulation with other privacy techniques (differential privacy, federated learning, MPC) 2. Contributing to HE library development for NLP-specific extensions 3. Hardware acceleration investigation (FPGAs, ASICs) 4. Exploring encrypted training beyond inference 5. Benchmark suite creation for privacy-preserving NLP

Long-Term Directions: 1. True HE integration as libraries mature 2. Developing formal privacy guarantees for simulation approach 3. Production-ready systems with proper

key management and scalability 4. Engaging with policymakers for regulatory frameworks 5. Cross-domain applications (images, audio, video, structured data)

Open Research Questions: 1. Can non-determinism of real HE be reconciled with transformer embeddings through novel architectures? 2. What is the theoretical minimum accuracy degradation with HE for NLP? 3. How to formally quantify privacy guarantees of deterministic simulation? 4. What hybrid architectures offer optimal accuracy-privacy-efficiency trade-offs? 5. Can post-quantum secure encryption integrate into neural networks without prohibitive overhead?

Concluding Remarks:

This research demonstrates both promise and challenges of privacy-preserving neural networks. While true homomorphic encryption for NLP remains elusive, the architectural innovations—particularly the Differentiable Soft-Argmax layer—and pragmatic simulation approach provide valuable stepping stones toward secure, private machine learning.

The 80% accuracy on encrypted data, while representing 18% degradation, is sufficient for numerous applications where privacy is paramount. More importantly, this work provides honest assessment of current limitations and clear directions for future research.

As data privacy concerns grow and regulations become more stringent, technologies enabling secure computation on sensitive data become increasingly critical. This research contributes to that goal, providing technical innovations and methodological insights advancing privacy-preserving machine learning.

The journey revealed that the path forward requires not just better encryption schemes, but co-design of cryptographic protocols and neural network architectures. This work provides a foundation for that continued evolution.

References

- ACAR, A.; AKSU, H.; ULUAGAC, AS; CONTI, M. A survey on homomorphic encryption schemes: Theory and implementation. *ACM Computing Surveys*, v. 51, n. 4, p. 1-35, 2018.
- CHEON, JH; KIM, A.; KIM, M.; SONG, Y. Homomorphic encryption for arithmetic of approximate numbers. In: *INTERNATIONAL CONFERENCE ON THE THEORY AND APPLICATION OF CRYPTOLOGY AND INFORMATION SECURITY*. Proceedings... Springer, 2017. p. 409-437.
- DEVLIN, J.; CHANG, MW; LEE, K.; TOUTANOVA, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In: *CONFERENCE OF THE NORTH AMERICAN CHAPTER OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS*. Proceedings... ACL, 2019. p. 4171-4186.
- GENTRY, C. Fully homomorphic encryption using ideal lattices. In: *ANNUAL ACM SYMPOSIUM ON THEORY OF COMPUTING*, 41., 2009. Proceedings... ACM, 2009. p. 169-178.

GILAD-BACHRACH, R.; DOWLIN, N.; LAINE, K.; LAUTER, K.; NAEHRIG, M.; WERNING, J. CryptoNets: Applying neural networks to encrypted data with high throughput and accuracy. In: INTERNATIONAL CONFERENCE ON MACHINE LEARNING, 33., 2016. Proceedings... JMLR, 2016. p. 201-210.

GOODFELLOW, I.; BENGIO, Y.; COURVILLE, A. Deep learning. Cambridge: MIT Press, 2016.

JUVEKAR, C.; VAIKUNTANATHAN, V.; CHANDRAKASAN, A. GAZELLE: A low latency framework for secure neural network inference. In: USENIX SECURITY SYMPOSIUM, 27., 2018. Proceedings... USENIX Association, 2018. p. 1651-1669.

LECUN, Y.; BENGIO, Y.; HINTON, G. Deep learning. Nature, v. 521, n. 7553, p. 436-444, 2015.

MICROSOFT RESEARCH. Microsoft SEAL (release 4.0). Microsoft, 2021. Available at: <https://github.com/Microsoft/SEAL>. Accessed on: Oct. 15, 2025.

PASZKE, A. et al. PyTorch: An imperative style, high-performance deep learning library. In: ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS, 32., 2019. Proceedings... NeurIPS, 2019. p. 8024-8035.

SANH, V.; DEBUT, L.; CHAUMOND, J.; WOLF, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. ArXiv preprint arXiv:1910.01108, 2019. Available at: <https://arxiv.org/abs/1910.01108>. Accessed on: Sept. 20, 2025.

SOCHER, R.; PERELYGIN, A.; WU, J.; CHUANG, J.; MANNING, CD; NG, A.; POTTS, C. Recursive deep models for semantic compositionality over a sentiment treebank. In: CONFERENCE ON EMPIRICAL METHODS IN NATURAL LANGUAGE PROCESSING. Proceedings... ACL, 2013. p. 1631-1642.

VASWANI, A. et al. Attention is all you need. In: ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS, 30., 2017. Proceedings... NeurIPS, 2017. p. 5998-6008.

WOLF, T. et al. Transformers: State-of-the-art natural language processing. In: CONFERENCE ON EMPIRICAL METHODS IN NATURAL LANGUAGE PROCESSING: SYSTEM DEMONSTRATIONS. Proceedings... ACL, 2020. p. 38-45.