

# Academic Paper Review Report

Review Date: 2025-07-13

## Paper Information

**Paper Title:** cuTraNTT: GPU-based Transposed Number Theoretic Transform with Low Latency Homomorphic Encryption for IoT Applications  
**Journal/Conference:** Cluster Computing

## Review Summary

This paper presents cuTraNTT, a GPU-based implementation of the Number Theoretic Transform (NTT) optimized for fully homomorphic encryption (FHE) in IoT applications. The work introduces several novel techniques including skipping two levels of radix-4 NTT, transposed polynomial arrangement, and optimized Toom-Cook 4-way multiplication. While the technical contributions are solid and performance improvements are demonstrated, the gains are modest and some claims require clarification.

## Detailed Review

### Innovation Assessment

The main innovations of this paper include:

1. **cuTraNTT technique:** Novel approach to skip the last two levels of radix-4 NTT and arrange polynomials in transposed manner, eliminating block synchronization and reducing from two kernels to one
2. **Single-kernel NTT implementation:** Significant contribution as previous GPU implementations (Jung et al., Yang et al.) required two separate kernels with mandatory synchronization
3. **32-bit arithmetic optimization:** Leverages GPU's native 32-bit architecture instead of 64-bit arithmetic used in existing implementations
4. **Optimized  $16 \times 16$  polynomial multiplication:** Fast parallel implementation using Toom-Cook 4-way algorithm to handle overhead from skipped NTT levels
5. **IoT-focused FHE:** First work to specifically target embedded GPU devices (Jetson Orin Nano) for FHE deployment

### Areas needing improvement:

- Limited novelty in individual components - mostly engineering optimizations of existing techniques
- Performance gains are modest (6% improvement over state-of-the-art)
- Trade-offs not thoroughly analyzed (NTT speedup vs. PWM overhead)

## Technical Quality Assessment

The technical quality of the paper is generally good:

### Strengths:

- **Solid theoretical foundation:** Well-grounded in RLWE-based FHE theory and radix-4 Cooley-Tukey algorithm
- **Comprehensive algorithm analysis:** Detailed breakdown of synchronization requirements in radix-4 NTT (2 sync points vs 5 in radix-2)
- **Sound mathematical approach:** Correct application of incomplete NTT theory and proper handling of polynomial arithmetic
- **Well-structured implementation:** Clear separation of concerns between NTT optimization and point-wise multiplication handling
- **Appropriate parameter selection:** Focuses on  $N=2^{16}$  which enables bootstrapping in CKKS

### Weaknesses:

- **Limited theoretical novelty:** Mainly combines existing techniques (incomplete NTT + transposed layout)
- **Incomplete complexity analysis:** Missing detailed computational complexity comparison with baseline methods
- **Insufficient memory analysis:** No thorough evaluation of memory usage patterns and bandwidth requirements
- **Unclear algorithmic details:** Some implementation specifics (e.g., exact transposition scheme) are not fully described
- **Missing error analysis:** No discussion of numerical stability or precision requirements for 32-bit arithmetic

## Experimental Assessment

The experimental design is comprehensive and covers major validation scenarios across multiple GPU platforms (RTX 4060, Jetson Orin Nano, V100, A100). The authors conducted sufficient comparative experiments with state-of-the-art implementations.

### Experimental Strengths:

- **Multi-platform validation:** Testing on both high-end (V100, A100) and embedded (Jetson Orin Nano) GPUs demonstrates broad applicability
- **Comprehensive benchmarking:** Comparison against Jung et al. and Yang et al. provides proper baseline evaluation
- **Relevant metrics:** Focus on latency measurements aligns well with IoT application requirements
- **Parameter consistency:** Uses  $N=2^{16}$  consistently for meaningful comparisons
- **Performance breakdown:** Separate analysis of HENC and HMULT operations provides insight into trade-offs

### Suggested Improvements:

- **Missing energy consumption analysis:** Critical for IoT/embedded applications but not evaluated
- **Limited parameter exploration:** Only  $N=2^{16}$  tested; scalability to other sizes unclear
- **Insufficient memory usage analysis:** No detailed memory footprint comparison
- **Missing accuracy verification:** No correctness testing or numerical precision analysis
- **Lack of real-world workloads:** All benchmarks are synthetic; no actual FHE application scenarios

## Writing Quality Assessment

The overall writing quality of the paper is good, with clear structure and fluent expression.

### Writing Strengths:

- **Clear structure:** Well-organized with logical flow from background to implementation to results
- **Comprehensive related work:** Thorough coverage of GPU-based FHE implementations and lattice cryptography
- **Good technical exposition:** Complex algorithms (NTT, CKKS) explained with appropriate detail
- **Effective use of figures:** Algorithm descriptions and performance comparisons well-illustrated
- **Consistent notation:** Mathematical symbols and terminology used consistently throughout

### Areas for Improvement:

- **Incomplete IoT motivation:** Claims IoT focus but limited discussion of actual IoT constraints and requirements
- **Missing implementation details:** Some algorithmic specifics (transposition scheme) insufficiently described
- **Unclear performance claims:** Some comparisons lack proper context (different hardware, parameters)
- **Grammar and style issues:** Minor language improvements needed throughout
- **Insufficient future work discussion:** Limited exploration of research directions and limitations

## Specific Revision Suggestions

### Major Issues

1. **Insufficient theoretical analysis:** The paper lacks a rigorous complexity analysis comparing the proposed method against baselines. While performance improvements are claimed, the theoretical justification for why skipping two NTT levels should be beneficial is not thoroughly established.
2. **Limited evaluation scope:** The experimental evaluation is restricted to  $N=2^{16}$  parameters only. For a technique claiming broad applicability, testing with different polynomial sizes ( $N=2^{15}, 2^{17}$ ) is essential to demonstrate scalability.
3. **Missing energy consumption analysis:** Given the IoT focus, energy efficiency is crucial but completely absent from the evaluation. This is a significant oversight for embedded applications where power consumption is often more critical than raw performance.
4. **Incomplete algorithmic description:** The exact transposition scheme and memory layout optimizations are not sufficiently detailed for reproducibility. Section 3.1 mentions “transposed manner” but lacks concrete implementation details.
5. **Questionable IoT applicability claims:** While the paper claims IoT relevance, the evaluation focuses primarily on high-end GPUs (V100, A100). The single embedded platform (Jetson Orin Nano) is still relatively powerful and may not represent typical IoT constraints.

## Minor Issues

1. **Grammar and style inconsistencies:** Several grammatical errors throughout, e.g., “The field of FHE is highly evolving” should be “rapidly evolving”
2. **Notation inconsistencies:** Some mathematical symbols introduced but not consistently used (e.g.,  $\Phi_N$  notation)
3. **Missing related work:** Recent advances in FHE acceleration (2023-2024) not covered in related work section
4. **Figure quality:** Figure 2 could benefit from clearer labeling and explanation of the block-thread mapping
5. **Bibliography formatting:** Some references incomplete or inconsistently formatted

## Technical Suggestions

### Algorithm Improvement Suggestions

1. **Extend theoretical analysis:** Provide rigorous complexity analysis comparing computational and memory costs of cuTraNTT vs. conventional two-kernel approach
2. **Multi-parameter evaluation:** Test with different polynomial sizes ( $N=2^{15}, 2^{17}$ ) to demonstrate scalability and identify optimal parameter ranges
3. **Energy efficiency metrics:** Add power consumption measurements, especially for embedded platforms, as this is critical for IoT deployment
4. **Memory optimization details:** Provide detailed analysis of memory access patterns, cache efficiency, and bandwidth utilization
5. **Real-world benchmarks:** Include evaluation on actual FHE applications (e.g., private inference, secure aggregation) rather than synthetic workloads
6. **Adaptive implementation:** Consider dynamic switching between cuTraNTT and conventional approaches based on polynomial size or hardware characteristics

## Review Conclusion

### Final Recommendation:

#### Major Revision Required

This paper presents a technically sound approach to optimizing NTT for GPU-based FHE with some interesting engineering contributions. The cuTraNTT technique that eliminates kernel synchronization is noteworthy, and the focus on IoT applications addresses an important gap. However, several significant issues limit the impact:

1. **Limited theoretical depth:** The work lacks rigorous complexity analysis and theoretical justification for the proposed optimizations
2. **Narrow evaluation scope:** Testing only  $N=2^{16}$  parameters and missing critical IoT metrics like energy consumption significantly weakens the claims
3. **Modest performance gains:** The 6% improvement over state-of-the-art, while positive, is not substantial enough to justify publication without stronger theoretical contributions
4. **Incomplete IoT validation:** Despite claiming IoT focus, evaluation on actual IoT constraints and broader embedded platforms is insufficient

The paper would benefit from major revisions addressing theoretical analysis, comprehensive evaluation across multiple parameters, energy efficiency measurements, and clearer algorithmic descriptions before being suitable for publication.

Reviewer: isomo  
Review Date: 2025-07-13