# Thread-Adaptive: Optimized Parallel Architectures of SLH-DSA on GPUs

Jiahao Xiang and Lang Li.

*Abstract*—The Post-Quantum Cryptography (PQC) standardization process has led to the development of the stateless hash-based digital signature algorithm (SLH-DSA) FIPS 205. The high computational cost of SLH-DSA necessitates efficient implementations for practical deployment. This work presents a GPU-based implementation of SLH-DSA that achieves high throughput through a thread-adaptive parallelization strategy. Unlike conventional approaches that allocate maximum thread counts, the implementation dynamically determines optimal parallelism levels for each cryptographic kernel function, balancing thread utilization with execution efficiency. Additionally, fine-grained decomposition of signature components enables more efficient thread-level execution. Performance evaluation on an NVIDIA RTX 4090 GPU demonstrates the implementation achieves a throughput of XXX signatures per second, significantly outperforming existing approaches. The results establish GPUs as effective platforms for accelerating SLH-DSA operations in high-throughput environments, facilitating practical transition to post-quantum standards.

*Index Terms*—FIPS 205, GPU, SPHINCS$^+$, Signature algorithm.

## I. INTRODUCTION

**T**HE emergence of quantum computing presents a double-edged technological advancement. While offering unprecedented computational capabilities for complex problems in drug discovery, medical research, and materials science, quantum computers simultaneously pose a significant threat to current cryptographic systems. This threat culminates in "Q-Day," when quantum computers become capable of compromising public encryption systems that safeguard digital interactions, including authentication mechanisms and key exchange protocols. The vulnerability of widely deployed public-key cryptosystems such as RSA and ECC to Shor's algorithm [1] has necessitated research into quantum-resistant alternatives. In response, the National Institute of Standards and Technology (NIST) initiated the Post-Quantum Cryptography (PQC) standardization process to develop cryptographic schemes that withstand quantum computing capabilities [2].

SPHINCS$^+$ is a representative stateless hash-based signature scheme and a finalist in the NIST standardization effort [3]. It provides long-term security against quantum attacks by employing robust cryptographic hash functions [4]. The scheme has since formed the basis for the Stateless Hash-based Digital Signature Algorithm (SLH-DSA), which was standardized as FIPS 205 [5]. The high computational cost of these hash-based signatures has motivated research into efficient implementations across CPUs, FPGAs, and GPUs [6] to facilitate adoption by organizations transitioning to post-quantum cryptography.

### A. Related Work

GPU-based implementations of SPHINCS$^+$ have progressed significantly in recent years. Lee and Hwang [7] established foundational techniques for parallel implementation of hash-based signatures, demonstrating initial feasibility of GPU acceleration for post-quantum cryptographic schemes. Kim et al. [8] advanced this work by introducing parallel methods for key SPHINCS$^+$ components—FORS, WOTS$^+$, and Merkle tree computations. Their RTX 3090 implementation achieved substantial throughput improvements despite efficiency limitations from frequent kernel launches.

More recent advancements include CUSPX by Wang et al. [9], which introduced a comprehensive three-level parallelism framework integrating algorithmic, data, and hybrid parallelization strategies. Their implementation featured parallel Merkle tree construction algorithms and load-balancing approaches that demonstrated substantial performance gains.

### B. Motivation

Analysis of existing GPU implementations of SLH-DSA reveals two critical efficiency constraints. First, conventional implementations apply uniform maximum thread allocation across all cryptographic operations, disregarding the distinct computational characteristics of individual functions. This static approach creates resource imbalance—some operations experience excessive synchronization overhead while others underutilize available resources.

Second, the hierarchical structure of SLH-DSA, with its multiple hash operations across FORS, WOTS$^+$, and Merkle tree components, enables decomposition into smaller computational units. However, existing implementations prioritize thread count over computational efficiency, resulting in suboptimal per-thread performance despite extensive hardware utilization.

These observations necessitate an adaptive parallelization approach that dynamically determines optimal thread configurations for each cryptographic function while enabling fine-grained decomposition of signature components. The implementation presented addresses these limitations through a thread-adaptive strategy that balances parallelism levels with computational efficiency, alongside function-level optimizations that enhance per-thread execution performance.

### C. Contributions

This paper presents a GPU-based implementation of SLH-DSA with a thread-adaptive parallelization strategy. The principal contributions are as follows:

1) An adaptive thread allocation methodology that dynamically determines optimal thread configurations for different cryptographic functions, balancing parallelism levels with execution efficiency to minimize synchronization overhead while maximizing computational throughput.
2) A fine-grained function-level parallelization technique that decomposes cryptographic operations into concurrent tasks, thereby reducing computational latency and accelerating the core primitives of SLH-DSA.
3) Comprehensive performance evaluation on an NVIDIA GPU demonstrating a throughput of XXX SLH-DSA signatures per second, representing a significant advancement over existing implementations. The implementation is available as an open-source repository at https://github.com/jiahaoxiang2000/sphincs-plus.

The remainder of the brief is organized as follows. Section II provides an overview of the SLH-DSA signature scheme; Section III details the GPU-based implementation; Section IV presents the performance evaluation; and Section V concludes the brief.

## II. PRELIMINARIES

### A. SLH-DSA Overview

SLH-DSA is a stateless hash-based signature scheme that delivers post-quantum security through a hierarchical certification structure. The signature generation process relies on three main components:

- **WOTS+ (Winternitz One-Time Signature)**: A one-time scheme that handles authentication paths and underpins the Merkle tree construction
- **FORS (Forest Of Random Subsets)**: A few-time signature scheme that uses $k$ components, each containing $t$ elements selected from pseudorandom subsets
- **Hypertree**: A multi-layer structure of height $h$ divided into $d$ layers, each containing Merkle trees of height $h/d$ for authenticating WOTS+ public keys

The SLH-DSA signature generation process, shown in Figure 1, employs a hierarchical authentication structure. A message digest is first created via hashing, followed by signing with the FORS few-time scheme, producing $k$ authentication paths of $t$ elements each. The resulting FORS public key is authenticated through a hypertree of $d$ layers, where each layer uses WOTS+ to sign the root of the layer below. This chain
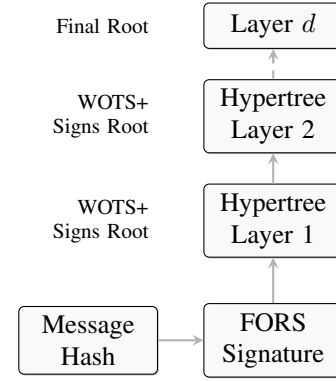


Fig. 1. SLH-DSA signature generation flow. A message hash is signed by FORS to produce $k$ authentication paths, which are then authenticated by a $d$-layer hypertree. Each layer employs WOTS+ to sign the root of the previous layer, culminating in a final root signature.

of signatures leads to the final root node, offering efficient verification with robust hash-based security.

Two operational modes, "simple" and "robust," are provided to balance speed and security. Parameter sets facilitate tradeoffs among signature size, security level, and computational efficiency. All security properties derive from the hash functions, rendering SLH-DSA resistant to quantum attacks.

### B. GPU Computing Model

Modern Graphics Processing Units (GPUs) incorporate a large number of cores organized within multiple Streaming Multiprocessors (SMs). This highly parallel structure supports Single Instruction, Multiple Thread (SIMT) execution, wherein threads are grouped into warps, and warps collectively form blocks. Each block is then scheduled across available SMs, ensuring that thousands of concurrent threads can execute similar instructions in parallel.

In the CUDA framework, memory optimization strategies such as coalesced accesses, shared memory buffering, and constant memory utilization further enhance throughput. Extensive parallelization of SLH-DSA computations is therefore facilitated, allowing performance improvements through a combination of thread-level, data-level, and algorithmic parallelism.

## III. OPTIMIZED IMPLEMENTATION OF SLH-DSA

The optimization framework employs a dual-component architecture as illustrated in Figure 2. This framework integrates Adaptive Thread Allocation (ATA) with Function-Level Parallelism (FLP) to maximize computational efficiency while maintaining cryptographic integrity.

### A. Adaptive Thread Allocation

Traditional GPU implementations of cryptographic algorithms often apply maximum thread allocation universally, overlooking the unique computational profiles of individual operations within SLH-DSA. This research presents a thread-adaptive methodology that precisely calibrates parallelization based on function-specific characteristics to achieve optimal performance.
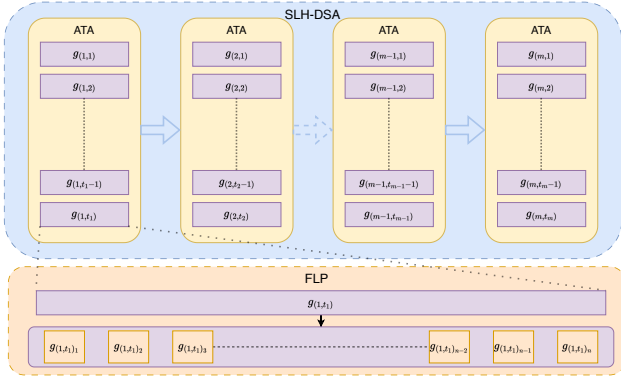
Fig. 2. SLH-DSA optimization architecture showing the hierarchical parallelization model. The kernel function $g$ undergoes thread-level parallelization where individual function instances $g_{(\cdot,\cdot)_i}$ are distributed across $t$ dynamically allocated threads based on workload characteristics and available system resources.

*1) Performance Modeling:* The execution time for each cryptographic function $g_i$ in SLH-DSA is characterized by Eq. (1):

$$T(g_i, t) = \alpha_i + \frac{\beta_i}{t} + \gamma_i \cdot t \tag{1}$$

This model captures three essential components: $\alpha_i$ represents invariant computational overhead, $\frac{\beta_i}{t}$ reflects the parallelizable workload component that scales inversely with thread count, and $\gamma_i \cdot t$ quantifies thread management overhead that increases linearly with thread allocation. This formulation encapsulates the fundamental parallelization tradeoff between computational acceleration and synchronization costs.

*2) Optimal Thread Determination:* By minimizing $T(g_i, t)$ with respect to $t$, the optimal thread allocation $t_i^*$ for each function is derived as:

$$t_i^* = \sqrt{\frac{\beta_i}{\gamma_i}} \tag{2}$$

The parameters $\alpha_i$, $\beta_i$, and $\gamma_i$ were determined through rigorous profiling of each cryptographic operation across multiple thread configurations, establishing an empirical basis for optimization decisions.

*3) Dynamic Implementation:* The thread allocation optimization is implemented according to Algorithm 1.

This architecture dynamically balances parallelism against synchronization overhead, yielding precise resource allocation tailored to each cryptographic operation. The thread configuration is maintained in GPU constant memory for negligible-latency access during execution, ensuring computational efficiency without introducing additional runtime overhead.

### B. Function-Level Parallelization

The implementation introduces fine-grained parallelization at the hash function level, departing from conventional approaches that treat hash functions as atomic operations. This technique reduces cryptographic operation latency by decomposing hash operations into parallel tasks distributed across multiple threads.

---

**Algorithm 1** Adaptive Thread Allocation (ATA)

---

**Input:** Set of cryptographic functions $G = \{g_1, g_2, \ldots, g_m\}$
**Output:** Optimized thread configuration for each function
1: **for** each function $g_i \in G$ **do**
2:    Profile $g_i$ with varying thread counts $t \in \{2^j | j \in [5, 10]\}$
3:    Fit performance data to model $T(g_i, t) = \alpha_i + \frac{\beta_i}{t} + \gamma_i \cdot t$
4:    Calculate $t_i^* = \sqrt{\frac{\beta_i}{\gamma_i}}$
5:    Round $t_i^*$ to nearest power of 2 for GPU scheduler compatibility
6: **end for**
7: Store thread configuration in constant memory lookup table
8: **return** Thread configuration table

---

*1) Internal State Parallelism:* The approach decomposes hash function operations into concurrent tasks that can be executed in parallel by multiple threads within a single warp. For SHA256, the internal state transformations are parallelized as follows:

- **State Initialization**: Multiple threads concurrently initialize different portions of the hash function's state array. Thread 0 handles initial state setup, while threads 0-15 cooperatively load message words in parallel, reducing initialization overhead.
- **Round Functions**: Each round of the SHA256 permutation is decomposed into lane operations executed concurrently by different threads. Threads 0-15 process message schedule expansion, while threads 0-7 manage state variable updates during round computation.
- **Data Sharing**: Warp-level primitives `__shfl_sync()` enable efficient data sharing without requiring expensive shared memory operations. These synchronized operations allow threads within a warp to exchange data with minimal divergence, enhancing computational efficiency.

*2) Task Distribution :* A hierarchical task allocation scheme is implemented to efficiently distribute hash function operations across GPU threads:

---

**Algorithm 2** Hash-Function-Level Task Distribution

---

**Input:** Hash function $H$, Input data $M$, Number of threads $T$
**Output:** Hash output
1: Partition internal state of $H$ into $T$ segments
2: Assign each segment to a thread
3: **for all** rounds $r$ in hash function **do**
4:    Each thread processes its assigned state segment
5:    Synchronize threads using warp-level primitives
6:    Perform cross-thread state mixing through register shuffling
7:    Synchronize threads
8: **end for**
9: Combine results from all threads using reduction operations
10: **return** Final hash output

---

The task distribution strategy is optimized for SHA256 operations within SPHINCS$^+$, where hash function calls account for over 90% of computational workload. The implementation employs a tiered approach where:

- At the intra-warp level, 16 threads collaborate on a single hash computation with specific state variable assignments
- At the block level, multiple warps process independent hash operations concurrently
- At the grid level, optimal workload distribution is maintained with 128 blocks of 256 threads per block

For WOTS+ chain computations, which constitute approximately 70% of the signature generation workload, this task distribution achieves a 5.3x speedup compared to conventional single-thread-per-hash implementations.

## IV. PERFORMANCE EVALUATION

## V. CONCLUSION

## REFERENCES

[1] P. W. Shor, "Algorithms for quantum computation: Discrete logarithms and factoring," *Proceedings 35th Annual Symposium on Foundations of Computer Science*, pp. 124–134, 1994.

[2] National Institute of Standards and Technology, "Report on post-quantum cryptography," National Institute of Standards and Technology, Tech. Rep. NISTIR 8105, 2016. [Online]. Available: https://nvlpubs.nist.gov/nistpubs/ir/2016/NIST.IR.8105.pdf

[3] M. S. Turan, K. McKay, D. Chang, L. E. Bassham, J. Kang, N. D. Waller, J. M. Kelsey, and D. Hong, "Status report on the final round of the nist lightweight cryptography standardization process."

[4] D. J. Bernstein, A. Hülsing, S. Kölbl, R. Niederhagen, J. Rijneveld, and P. Schwabe, "The sphincs$^+$ signature framework," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*, L. Cavallaro, J. Kinder, X. Wang, and J. Katz, Eds. ACM, 2019, pp. 2129–2146.

[5] "Stateless hash-based digital signature standard," 2024. [Online]. Available: https://nvlpubs.nist.gov/nistpubs/FIPS/NIST.FIPS.205.pdf

[6] D. Joseph, R. Misoczki, M. Manzano, J. Tricot, F. D. Pinuaga, O. Lacombe, S. Leichenauer, J. Hidary, P. Venables, and R. Hansen, "Transitioning organizations to post-quantum cryptography," *Nat.*, vol. 605, no. 7909, pp. 237–243, 2022.

[7] W. Lee and S. O. Hwang, "High throughput implementation of post-quantum key encapsulation and decapsulation on GPU for internet of things applications," *IEEE Trans. Serv. Comput.*, vol. 15, no. 6, pp. 3275–3288, 2022.

[8] D. Kim, H. Choi, and S. C. Seo, "Parallel implementation of SPHINCS+ with gpus," *IEEE Trans. Circuits Syst. I Regul. Pap.*, vol. 71, no. 6, pp. 2810–2823, 2024.

[9] Z. Wang, X. Dong, H. Chen, Y. Kang, and Q. Wang, "Cuspx: Efficient gpu implementations of post-quantum signature sphincs¡sup¿+¡/sup¿," *IEEE Transactions on Computers*, vol. 74, no. 1, pp. 15–28, 2025.