

A j -lanes tree hashing mode and j -lanes SHA-256

Shay Gueron^{1,2}

¹ Department of Mathematics, University of Haifa, Israel

² Intel Corporation, Israel Development Center, Haifa, Israel

August 21, 2012

Abstract. j -lanes hashing is a tree mode that splits an input message to j slices, computes j independent digests of each slice, and outputs the hash value of their concatenation. We demonstrate the performance advantage of j -lanes hashing on SIMD architectures, by coding a 4-lanes-SHA-256 implementation and measuring its performance on the latest 3rd Generation Intel[®] Core[™]. For message ranging 2KB to 132KB in length, the 4-lanes SHA-256 is between 1.5 to 1.97 times faster than the fastest publicly available implementation (that we are aware of), and between ~ 2 to ~ 2.5 times faster than OpenSSL 1.0.1c. For long messages, there is no significant performance difference between different choices of j . We show that the 4-lanes SHA-256 is faster than the two SHA3 finalists (BLAKE and Keccak) that have a published tree mode implementation. We explain why j -lanes hashing will be even faster on the future AVX2 architecture with 256 bits registers. This suggests that standardizing a tree mode for hash functions (SHA-256 in particular) would deliver significant performance benefits for a multitude of algorithms and usages.

Keywords: Tree mode hashing, SHA-256, SHA3 competition, SIMD architecture, Advanced Vector Extensions architectures, AVX, AVX2.

1 Introduction

The performance of hash functions plays an important role in various situations (e.g., for SSL/TLS connections that use HMAC for authenticated encryption). In particular, the performance of SHA-256 on high end processors is a performance baseline for the SHA3 competition [1].

Recently, [2] published a “Simultaneous Hashing” (S-HASH) method, for using SIMD architectures to speed up the computations of SHA-256 (and other hashes) over multiple messages. In this paper, we apply this technique to accelerate SHA-256 for a single message, using a tree mode that we call j -lanes hashing. We show that the resulting “ j -lanes SHA-256” is significantly faster than the standard (“linear” hereafter) SHA-256. This demonstrates the performance benefits of standardizing tree modes for hash functions (in particular, SHA-256 and SHA-512). It is interesting to compare our results to the two SHA3 finalists that already have a j -lanes tree mode implementation (see [3]): BLAKE ($j=2$ and $j=4$) and Keccak ($j=2$). We offer this comparison in Section 4.

2 *j*-lanes hashing and the special case of 4-lanes SHA-256

Tree hashing is a well known concept for speeding up hash functions computations, and is an efficient way for updating the hash value when only a portion of the message is changed. Some relevant references are [4], [5], [6], [7], [8]. We focus here on a specific tree construction, which is defined in the following section.

2.1 *j*-lanes hashing

Definition 1 (message *j*-Slicing): given a message m , its associated *j*-Sliced message is the permutation (not necessarily concatenation) of disjoint slices of m , namely $m = \text{permutation}(m_1 \parallel m_2 \parallel m_3 \parallel \dots \parallel m_j)$ under some agreed convention on how each slice is defined (for simplicity assume that m has at least j bits, to avoid empty string slices).

Definition 2 (*j*-lanes-hash): Let $h = h(\text{MESSAGE})$ be a hash function. Its associated *j*-lanes-hash, is a hash scheme that operates as follows:

1. *j*-Slicing the message to $m = \text{permutation}(m_1 \parallel m_2 \parallel m_3 \parallel \dots \parallel m_j)$.
2. Computing $t_1 = h(m_1)$, $t_2 = h(m_2)$, ..., $t_j = h(m_j)$.
3. Computing $t^* = h(t_1 \parallel t_2 \parallel \dots \parallel t_j)$.
4. Returning the digest t^* .

(hereafter we call Step 3 the “Wrapping” step).

j-lanes-hash is a special form of a tree mode (not a binary tree), where the number of nodes is $j+1$ and the height of the tree is 2. As a special case of a tree mode, the security properties of this construction follow from the more general theory on tree hashing (e.g., [5] and [6] discuss the security properties of a tree hash in the context of indifferenciability from an ideal hash function).

Note that the definition covers several setups. One example is “interleaving” segments of a given message (which we use here, for directly taking advantage of SIMD architectures). Another case is when the data is consumed from j locations (e.g., j pointers) of a message. This can occur in an application that hashes a file system (or a directory) where j is the number of files (and each file is a node in the tree). We assume hereafter that the processed messages are sufficiently long to gain performance advantage of the *j*-lanes tree mode (and ignore trivially short message).

2.2 Applying *j*-lanes hashing to derive a 4-lanes SHA-256

We use SHA-256 as the underlying hash algorithm, and generate a “*j*-lanes SHA-256”. Our motivation is the potential performance advantage that stems from the parallelization offered by SIMD architectures (or multithreaded implementations).

By splitting the message into j independent slices, the hash computations are reduced to the problem of hashing multiple independent messages, supplemented by the fixed-cost Wrapping step. Techniques for using SIMD architectures for hashing multiple independent messages (of different lengths), and the resulting performance speed-ups, are described in detail in [2]. We use these techniques here.

SHA-256 operates on 32-bit words. Therefore, on processors that support the AVX (or SSE) architecture that has 128-bit registers and the necessary integer instructions, a natural choice for j -lanes (SHA-256) hashing is $j=4$, with the obvious convention for slicing the message: consecutive 128-bit chunks of the message are treated as 4 consecutive 32-bit words, each one of a different slice. These 4 words fit in as 4 “elements” of a single AVX register (*xmm*), and the SHA-256 computations can therefore be parallelized using the SIMD architecture (see [2] for details).

If the byte-length of the message is divisible by 256, the slices have equal lengths. Otherwise, (at least) one of slice has a different length, and this situation requires different handling in the last Update (with negligible performance cost).

j -lanes hashing involves some overhead, and therefore, the performance gains are expected to be (fully) manifested only for sufficiently long messages.

To illustrate, we note that the performance of SHA-256 is closely proportional to the number of invocations of its compression function (“Update” hereafter). Consider a message whose byte-length l is divisible by 256, and write $l = 256x$ for some integer x . Hashing (with SHA-256) such a message requires $4x+1$ Updates, where the last one due to the padding block. On the other hand, 4-lanes SHA-256 for this message requires $4(x+1) + 3$ Updates, accounting for 1 padding block for each slice, and 3 Updates for the Wrapping step which requires hashing of a 128 bytes message. Comparing the linear (i.e., serial) SHA-256 to the 4-lanes SHA-256, we see that the latter involves 6 additional Updates. However, from the total of $4x+7$ Updates, $4x$ can be parallelized, in particular by using the AVX architecture. This is the reason why the overall performance is expected to improve.

3 Performance studies

This section discusses some performance studies for of j -lanes SHA-256. We first describe the measurement methodology.

- Each measured function was isolated, run 25,000 times (warm-up), followed by 100,000 iterations that were timed (using the RDTSC instruction) and averaged.
- To minimize the effect of background tasks running on the system, each experiment was repeated five times, and the minimum result was recorded.
- All the runs were carried out on a system where the Intel® Hyper-Threading Technology, the Intel® Turbo Boost Technology, and the Enhanced Intel Speedstep® Technology, were *disabled*.
- The runs were executed on the 3rd Generation Intel® Core™ i7-3770 processor (previously known as “Architecture Code name Ivy Bridge”).
- In all cases, the reported performance numbers account for the full computations (i.e., including the padding and, when relevant, the final hashing of the j digests).

In our studies, we used two SHA-256 and three j -lanes-SHA-256 ($j=4, 8, 16$) implementations as follows:

- OpenSSL (1.0.1c) linear: standard hashing using OpenSSL function.
- 4-SMS linear: standard hashing using the n -SMS ($n=4$) method (see [9], [10]; we used here an improved version of this implementation).

- j -lanes using OpenSSL: using OpenSSL's (1.0.1c) SHA-256 function to implement j -lanes-SHA-256.
- j -lanes using the n -SMS: using the n -SMS SHA-256 implementation ([9], [10]) to implement j -lanes SHA-256.
- AVX j -lanes hashing (j -lanes hashing for short): an optimized implementation of j -lanes SHA 256, using the S-HASH implementation of [2], and the AVX architecture.

The results are illustrated in Figures 1-3.

Figure 1 compares the different implementations for an 8KB message and $j=4$. Without parallelizing the hashing of the slices (as in j -lanes using OpenSSL and j -lanes using the n -SMS), the 4-lanes SHA-256 is slower than the linear implementation. This is due to the overheads of the j -lanes method. For example, OpenSSL (1.0.1c) uses 129 Updates and performs at 12.87 Cycles/Byte, while the 4-lanes SHA-256 implementation that simply calls the OpenSSL functions, uses 135 Updates, and performs at 13.57 Cycles/Byte. On the other hand, the optimized (using AVX) 4-lanes SHA-256 implementation is 2.45 times faster than OpenSSL.

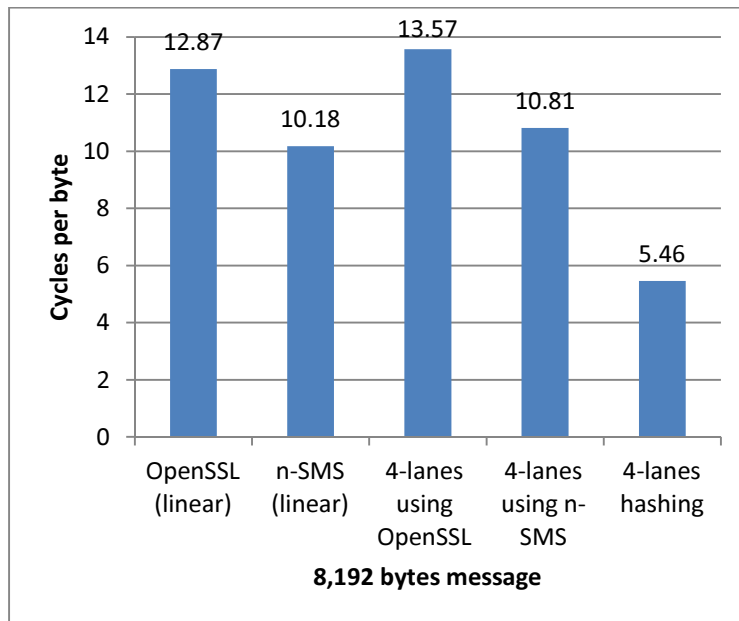


Fig. 1. Performance of different implementations of 4-lanes SHA-256, compared to linear SHA-256, for a 8192- bytes message. Measurements taken on the 3rd Generation Intel® Core™ Processor.

Figure 2 illustrates the effect of the choice of j ($= 4, 8, 16$). Obviously, increasing j involves additional overhead to the j -lanes hashing. For example, 16-lanes SHA-256 for an 8KB message involves 153 Updates, and is therefore slower than the 4-lanes SHA-256 that uses only 135 Updates (see top panel). However, both 8-lanes and 16-lanes SHA-256 are still significantly faster than the best performing linear

implementation. For long messages (see bottom panel), the relative impact of the overheads decreases, and we obtain roughly the same performance for $j = 4, 8, 16$.

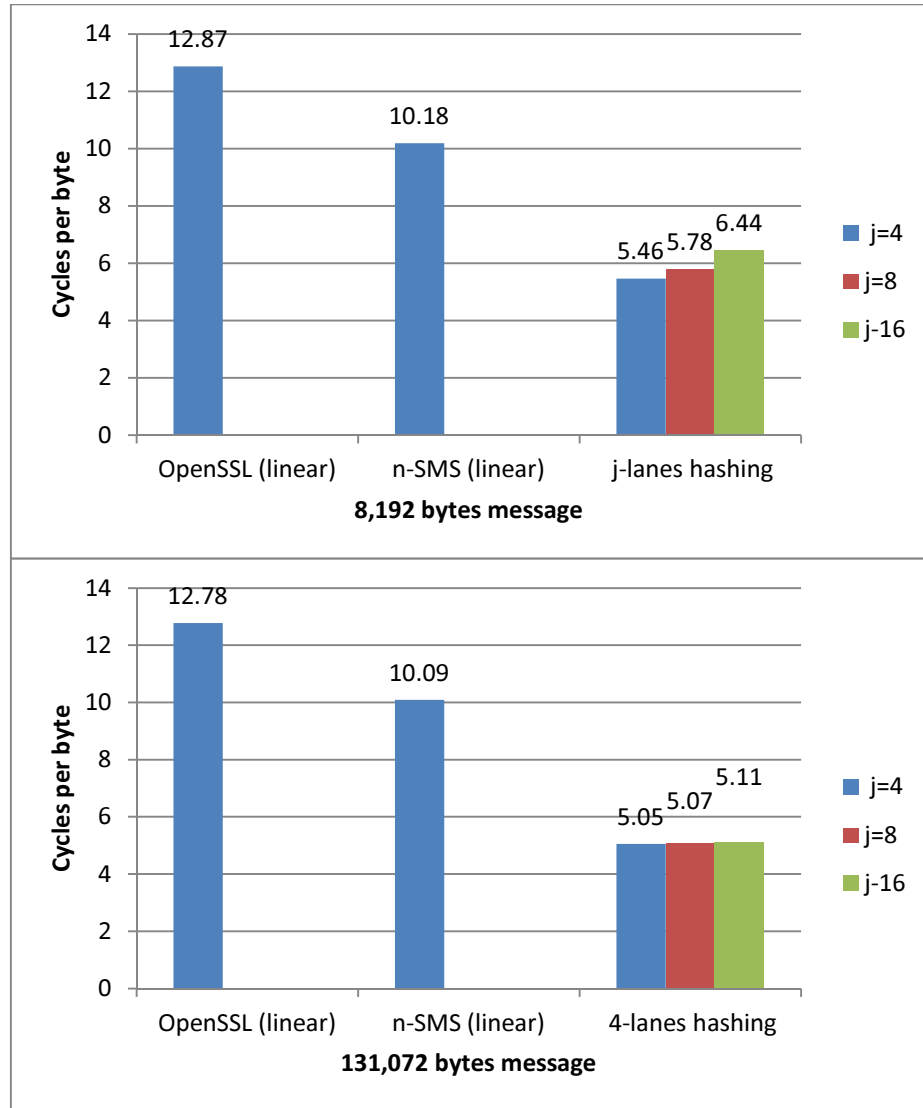


Fig. 2. Performance of j -lanes-SHA-256 for $j=4, 8, 16$, compared to linear SHA-256. The message length is 8,192 bytes (top panel) and 131,072 bytes (bottom panel). Measurements taken on the 3rd Generation Intel® Core™ Processor.

Figure 3 shows the performance advantage of the 4-lanes SHA-256 for messages of lengths varying from 2KB to 128KB: 4-lanes SHA-256 is between 1.55 to 2 times

faster than the best serial implementation (and $1.97 - 2.53$ times faster than OpenSSL 1.0.1c).

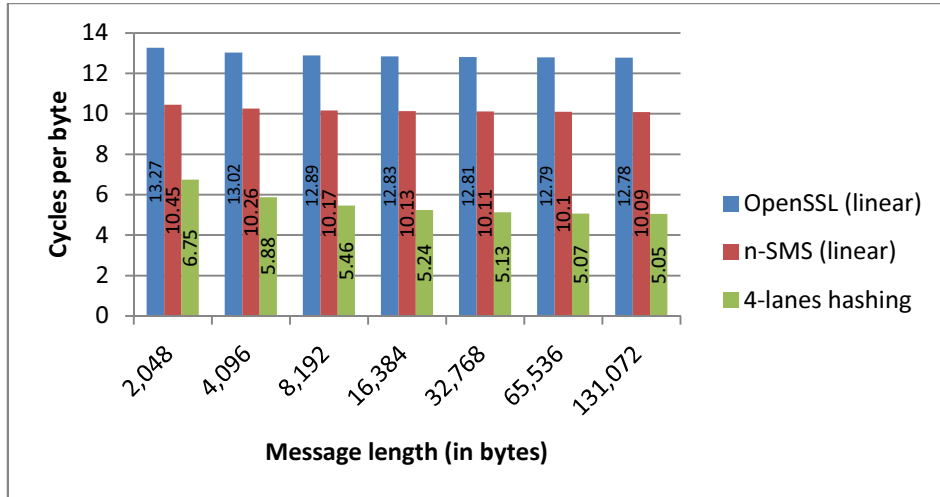


Fig. 3. Performance of 4-lanes SHA-256, compared to linear SHA-256, for different message lengths. Measurements taken on the 3rd Generation Intel[®] Core[™] Processor.

4 Conclusion

We demonstrated the performance gains of j -lanes hashing, using SHA-256 as the underlying hash algorithm. On the 3rd Generation Intel[®] Core[™] Processor (with AVX architecture) selecting $j=4$, gives speedup factors between $1.55x$ to $2x$, compared to the best available implementation (up to $2.53x$ when comparing to OpenSSL 1.0.1c). We focused on $j=4$, as the natural choice for the current AVX (and SSE) architectures. Interestingly, although $j=4$ yields the best results (the Wrapping overhead is the smallest among the tested cases), we note that the performance with all the studies choices $j=4, 8, 16$ is roughly the same for long messages.

We also comment that with the near future AVX2 architecture [11], a natural choice would be $j=8$ for SHA-1, SHA-256 and $j=4$ for SHA-512, and the j -lanes implementations will be significantly faster. Therefore, if a j -lanes hashing mode is adopted, and the ecosystem would prefer to support only a single value of j (to reduce the interoperability complexities), it seems that selecting $j=8$ would be a good choice.

In general, the j -lanes-hash can be useful in other scenarios, and with different values of j . One example mentioned about is hashing a file system, where j is the number of files (and each file is a node in the tree). Such computations can be accelerated not only by using SIMD architectures, but also by using the processing power of multi-cores systems.

We conclude that the j -lanes-hash could alleviate computational bottlenecks, and recommend that this mode (or a general tree mode) is standardized. To this end, we comment that standardization of a j -lanes (or any tree) mode should also properly define different initialization vectors (depending also on the value of j) in order to

distinguish the resulting digests from outputs of the linear SHA-256 (analogously to the how a digest truncation (e.g., SHA224) is defined).

4.1 A comment on the SHA3 finalists

We expect that the SHA3 finalists [1] could also gain from using the j -lanes-hash, at least to some extent, and the performance gains will further increase when the AVX2 architecture becomes available. However, at this point, it is hard to tell if these algorithms would outperform the j -lanes-SHA-256 and/or j -lanes-SHA-512, and by what margin.

Since the two finalists BLAKE and Keccak already have a tree mode implementation ($j=2$ and $j=4$ for BLAKE, and $j=2$ for Keccak; see [3]), we show the performance comparisons of SHA-256, BLAKE, and Keccak in linear and in j -lanes mode in Figure 4.

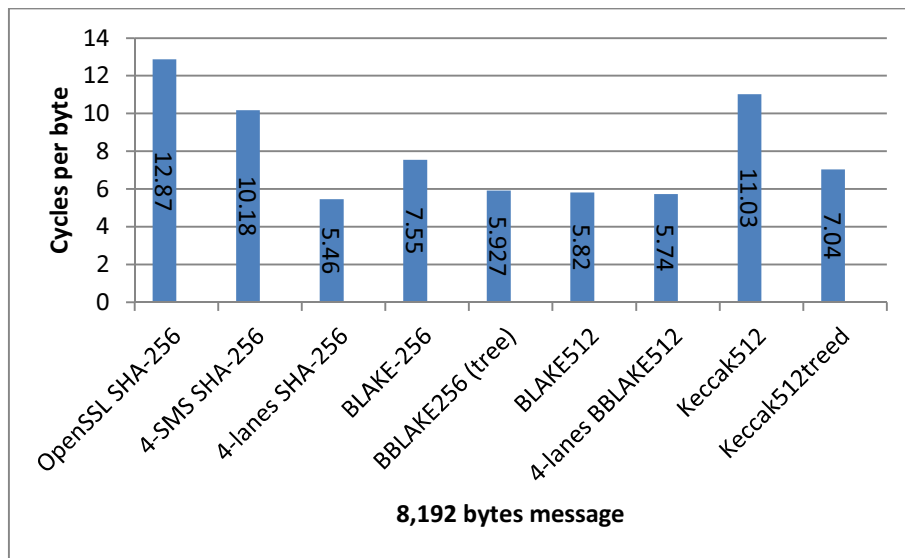


Fig. 4. Performance of SHA-256, BLAKE256, BLAKE512, and Keccak, in “linear” mode and in tree mode (for a 8,192 bytes message). Measurements taken on the 3rd Generation Intel[®] Core[™] Processor.

As expected, the j -lanes (tree mode) implementation improves the performance of all three algorithms. The results show that the j -lanes SHA-256 implementation is the fastest one of these three.

Recalling that SHA-256 (and SHA-512) is the performance baseline for SHA3, we conclude (from the currently available information) that considering the j -lanes mode still *does not* offer a performance advantage for SHA3 over SHA-256. This is consistent with the findings of [6]: migration to a new SHA3 standard could not be motivated by performance advantages on the high end platforms.

5 Acknowledgements

I thank Jean-Philippe Aumasson, Bart Preneel and Jesse Walker for helpful discussions.

6 References

- [1] NIST, cryptographic hash Algorithm Competition. <http://csrc.nist.gov/groups/ST/hash/sha-3/index.html>
- [2] Gueron, S., Krasnov, V.: Simultaneous hashing of multiple messages (2012), <http://eprint.iacr.org/2012/371.pdf>
- [3] SUPERCOP, <http://bench.cr.yp.to/supercop.html>
- [4] Bertoni, G., Daemen, J., Peeters, M., Van Assche, G.: Keccak sponge function family main document. Submission to NIST; updated (2009) <http://cuda-keccak.googlecode.com/svn/trunk/docs/Keccak-main-2.1.pdf>
- [5] Bertoni, G., Daemen, J., Peeters, M., Van Assche, G.: Sufficient conditions for sound tree and sequential hashing modes (2009) <http://eprint.iacr.org/2009/210>
- [6] Dodis, Y., Reyzin, L., Rivest, R.L., Shen, E.: Indifferentiability of Permutation-Based Compression Functions and Tree-Based Modes of Operation, with Applications to MD6. Proceedings of FSE 2009, Lecture Notes in Computer Science, 5665 104-121 (2009).
- [7] Merkle, R. C.: A certified digital signature. Advances in Cryptology. Proceedings of CRYPTO '89, Lecture Notes in Computer Science, 435: 218-238 (1990).
- [8] P. Sarkar, P. Schellenberg, P. J.: A parallelizable design principle for cryptographic hash functions. Cryptology ePrint Archive (2002), <http://eprint.iacr.org/2002/031>
- [9] Gueron, S., Krasnov, V.: Parallelizing message schedules to accelerate the computations of hash functions (2012), <http://eprint.iacr.org/2012/067.pdf>
- [10] Gueron, S., Krasnov, V.: [PATCH] Efficient implementations of SHA256 and SHA512, using the Simultaneous Message Scheduling method, <http://rt.openssl.org/Ticket/Display.html?id=2784&user=guest&pass=guest>
- [11] Intel (M. Buxton): Haswell New Instruction Descriptions Now Available! <http://software.intel.com/en-us/blogs/2011/06/13/haswell-new-instruction-descriptions-now-available/>