



# Parallelization of the Smith-Waterman Algorithm with HPX

STELLAR

stellar.cct.lsu.edu



CENTER FOR COMPUTATION & TECHNOLOGY

Stephanie Crillo,<sup>1</sup> Bryce Adelstein-Lelbach,<sup>2,3</sup> Hartmut Kaiser<sup>2,3</sup>  
<sup>1</sup>Belmont Abbey College, <sup>2</sup>Louisiana State University, <sup>3</sup>LSU Center for Computation and Technology

## Abstract

Concepts introduced by the Smith-Waterman algorithm are utilized by various bioinformatics systems today, such as BLAST. The algorithm, although accurate, slows greatly when working with larger data sets, due to its  $O(n^2)$  complexity<sup>(1)</sup>. We wrote a simple Smith-Waterman algorithm, in both serial and parallel, in an attempt to improve the scalability of Smith-Waterman, as part of our larger effort to improve the scalability of BLAST through parallelization.

## Introduction

The Smith-Waterman algorithm, (Figure 1), is designed to find the optimal local alignment of two sequences.

$$H_{(i,j)} = \max \begin{cases} H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \\ 0 \end{cases}, 1 \leq i \leq n, 1 \leq j \leq m$$

**Figure 1: The Smith-Waterman Algorithm**  
 The first term describes a match/mismatch, the second term describes a deletion, and the third term describes an insertion.

The algorithm builds a matrix based on a specific scoring scheme (Figure 2). In our runs, matches were granted a score of +2, while deletions and insertions (gaps) were granted a score of -1.

	-	T	C	G	T	A	T	G	T
-	0	0	0	0	0	0	0	0	0
T	0	2	1	0	2	1	2	1	2
G	0	1	1	3	2	1	1	4	3
C	0	0	3	2	2	1	0	3	3
A	0	0	2	2	1	4	3	2	2
T	0	2	1	1	4	3	6	5	4
A	0	1	1	0	3	6	5	5	4
C	0	0	3	2	2	5	5	4	4
T	0	2	2	2	4	4	7	6	6

**Figure 2: Smith-Waterman Scoring Matrix**  
 This matrix was built in the comparison of the two sequences TGCATACT (our "database" sequence) and TCGTATGT (our "queried" sequence).

Upon completion of the matrix, the algorithm backtracks through the matrix to obtain an optimal alignment (Figure 3).

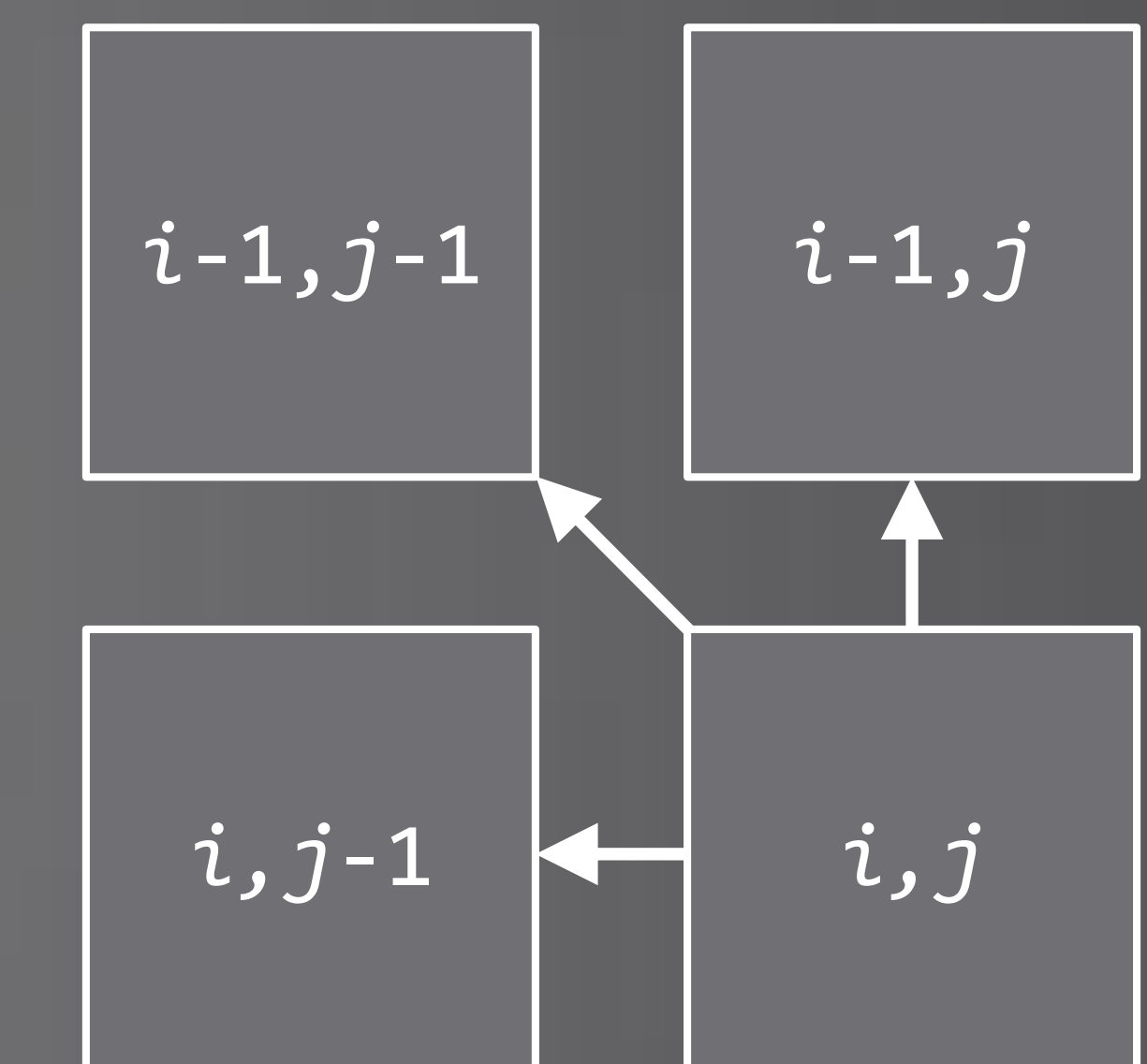
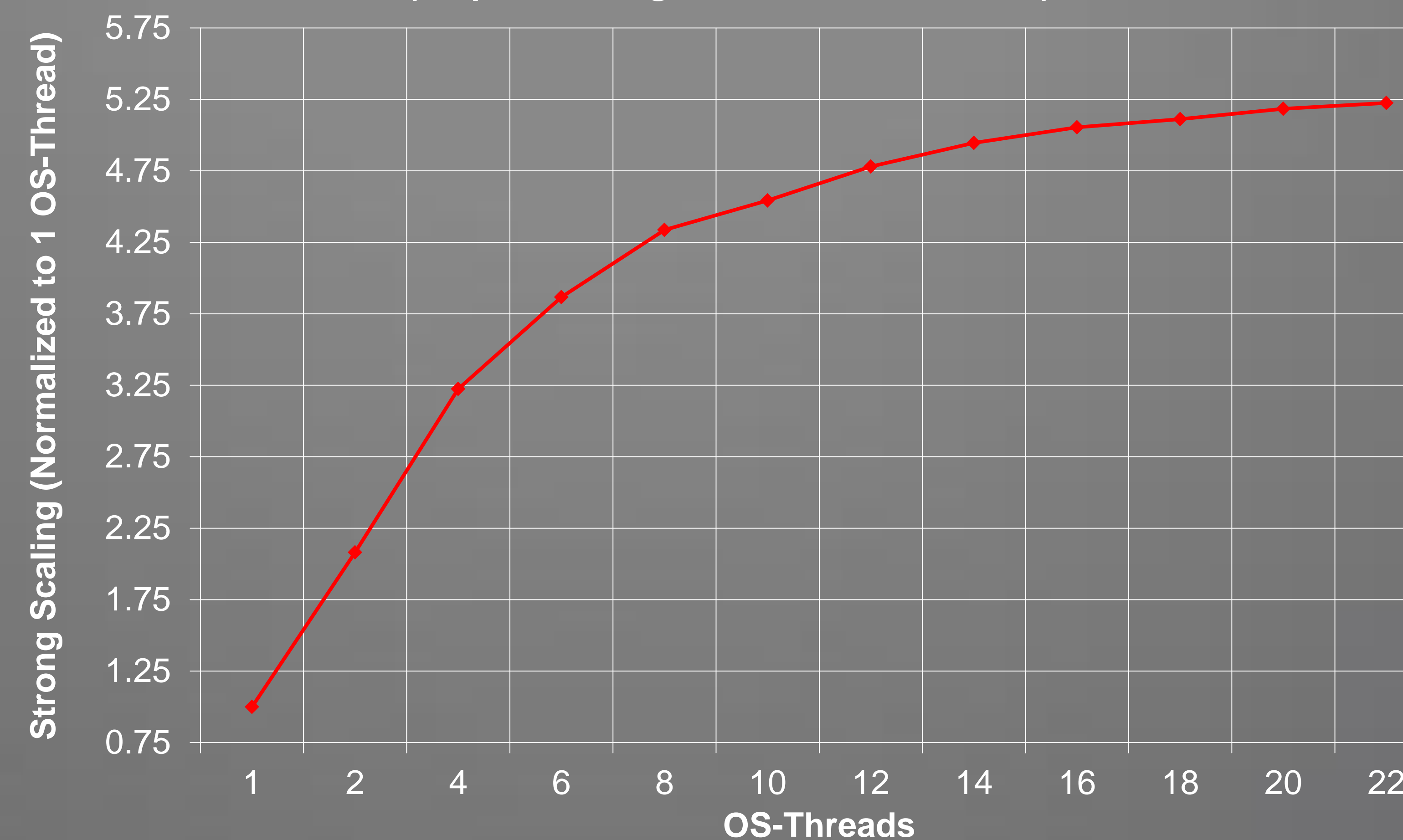
	-	T	C	G	T	A	T	G	T
-	0	0	0	0	0	0	0	0	0
T	0	2	1	0	2	1	2	1	2
G	0	1	1	3	2	1	1	4	3
C	0	0	3	2	2	1	0	3	3
A	0	0	2	2	1	4	3	2	2
T	0	2	1	1	4	3	6	5	4
A	0	1	1	0	3	6	5	5	4
C	0	0	3	2	2	5	5	4	4
T	0	2	2	2	4	4	7	6	6

**Figure 3: Backtracking**  
 Backtracking begins at the highest number<sup>(2)</sup> (here, 7), and continues through the matrix until a point surrounded by zeros is reached<sup>(2)</sup>. Diagonal movement indicates a match or mismatch, top-down movement indicates a deletion, and left-right movement (not shown) indicates an insertion.

## Results

Parallelization of the Smith-Waterman code was done by dividing the matrix into smaller submatrices and computing the scores of these submatrices in parallel where possible (see Figure 4 and 5). Additionally, backtracking of each of these components was computed as the matrix was being constructed.

### Strong Scaling of the HPX Smith-Waterman Algorithm (Sequence Length: 4096, Grain Size: 64)

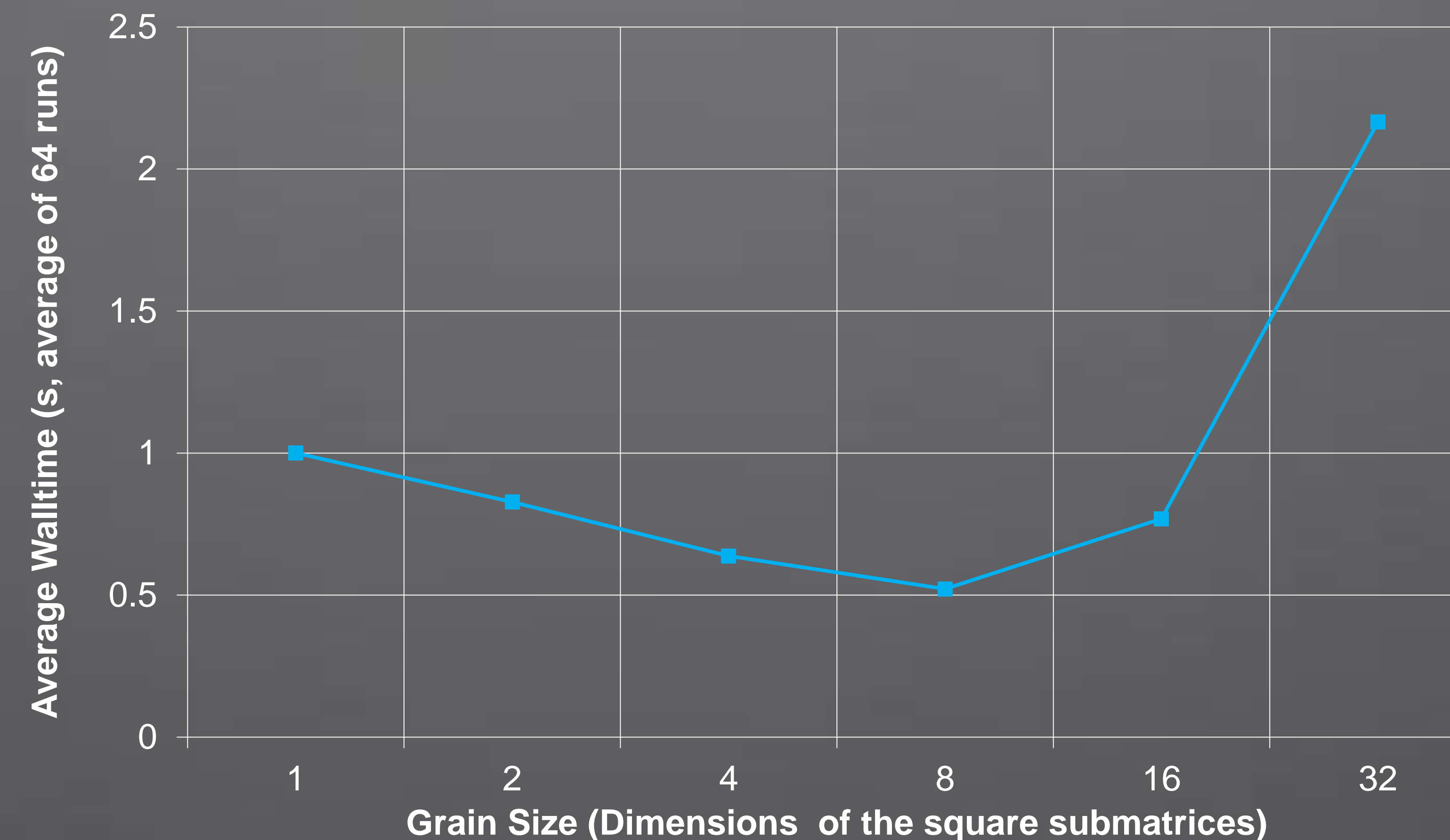


**Figure 4: Dependencies of the Algorithm**  
 To compute the value  $(i, j)$ , the values immediately above, diagonal and left of  $(i, j)$  are needed.



**Figure 5: Hot Zones**  
 Darker, or "hotter" cells represent regions in the matrix which have a higher potential for parallelization, because these regions are independent of each other.

### Grain Size Parameter Sweep (Sequence Length: 512, OS-Threads: 22)



## References

- (1) Torbjørn Rognes, 2011. Faster Smith-Waterman database searches with inter-sequence SIMD parallelisation. BMC Bioinformatics. [Internet]. [cited 2012 July 12]; 12:221. Available from: <http://www.biomedcentral.com/1471-2105/12/221>.
- (2) Smith T.F., Waterman M. S., 1981. Identification of Common Molecular Subsequences. Journal of Molecular Biology. [Internet]. [cited 2012 July 12]; 147:195-197. Available from: [http://dx.doi.org/10.1016/0022-2836\(81\)90087-5](http://dx.doi.org/10.1016/0022-2836(81)90087-5).

## Acknowledgements

- NSF Grants 1240655, 1117470, 1048019 and 1029161
- DOE ASCR X-Stack Program
- Center for Computation and Technology