

4005-735 Parallel Computing I
Project Report:
Longest Common Subsequence

Adam Risi
Rochester Institute of Technology
ajr7708@rit.edu

Richard Sarkis
Rochester Institute of Technology
res7008@rit.edu

February 14th, 2011

1 Abstract

The Longest Common Subsequence (LCS) problem is common to bioinformatics and a classic in computer science, generally considered NP-hard. The problem is commonly encountered in bioinformatics in the comparison of DNA which results in strings that are a combination of DNA base pairs representing an organism's DNA sequence. These sequences are typically very long in length resulting in a computationally intensive problem. By determining the longest common subsequence of DNA, biologists can determine the relative closeness of two organisms. The goal of this project was to research methods by which one could parallelize the LCS problem to speed up string comparison. Currently, the most common methods used by used in bioinformatics is BLAST and FASTA. Both algorithms emphasize speed of comparison at the expense of sensitivity of determining optimal alignments of sequences. Algorithms like Smith-Waterman provide exactness at the expense of speed, and in this problem area this project aims to research the current state of exact, parallelizable algorithms usable in DNA analysis.

2 Computational Problem

The Longest Common Subsequence (LCS) problem is common to bioinformatics and a classic in computer science. The problem is one in which given two sequences $X = \langle x_1, x_2, \dots, x_m \rangle$ and $Y = \langle y_1, y_2, \dots, y_n \rangle$ we attempt to find the maximum-length common subsequences of X and Y . A subsequence is not the same as a substring (subsequences do not need to be consecutive in the string). Given the string $X = \langle x_1, x_2, \dots, x_m \rangle$, string $Z = \langle z_1, z_2, \dots, z_k \rangle$ is then a subsequence of X if there is a strictly increasing sequence $\langle i_1, i_2, \dots, i_k \rangle$ of indices of X such that for all $j = 1, 2, \dots, k$ we have $x_{i_j} = z_j$ [3].

The problem is commonly encountered in bioinformatics in the comparison of DNA which is represented by the four DNA bases adenine, guanine, cytosine, and thymine represented relatively by the letters A, C, G, T . This results in strings that are a combination of these letters representing an organism's DNA sequence. These sequences are typically very long in length resulting in a computationally intensive problem.

The problem is generally considered NP-hard problem with an arbitrary number of input sequences but with a constant number of inputs the solution time is polynomial through the use of dynamic programming [2].

Parallelizing the LCS problem results in great increases in speed in the decoding of the genetic-derived character strings being compared. By determining the longest common subsequence of DNA, biologists can determine the relative closeness of two organisms.

The goal of this project is to determine the speed-up of this problem when using parallelization techniques and compare these performance metrics with a sequential single-CPU experiment.

3 Paper Analysis: Liu et al

3.1 Problem Description

With it being one of the primary tasks in bioinformatics, finding the LCS of a biosequence efficiently is imperative. Currently, there are various paralleled LCS implementations that vary in space and time complexity. They all seek to produce exact, correct results that sequential algorithms like Smith-Waterman and Needleman-Wunsch have been producing since the early 1980s. This paper presents an algorithm called *FAST_LCS* that promises better performance over other parallel LCS algorithms.

3.2 Contributions

Their contribution is their *FAST_LCS* algorithm which works in $max 8 * (n + 1) + 8 * (m + 1), L$ space, and $O(|LCS(X, Y)|)$ time, where m and n are the lengths of strings X and Y respectively, L is the number of identical character pairs and $|LCS(X, Y)|$ is the length of the LCS of X, Y [7]. A primary data structure for the algorithm builds a successor table TX and TY for the strings X and Y respectively, where each entry in these tables lists the successor of character $CH(i)$ in the string.

$$TX(i, j) = \begin{cases} \min\{k | k \in SX(i, j)\} & SX(i, j) \neq \phi \\ - & otherwise \end{cases}$$

[7]

Using the successor tables, a second table categorizing identical character pairs of X and Y is called $S(X)$ and $S(Y)$ respectively. A parallel computation, and the basic operation of the algorithm is producing the direct successors of each identical pair at each time step.

3.3 Investigative Use

The research into LCS algorithms, especially parallelizable ones is an ongoing area of active research; unfortunately it was clear that a unique solution, even a naive one was not approachable withing the context of this project. Instead, it was decided to choose the algorithm as described in [7] and implement it in Parallel Java.

4 Paper Analysis: Liu et al

4.1 Problem Description

4.2 Contributions

4.3 Investigative Use

5 Problem Solution

This project involves two sets of deliverables, a sequential version utilizing Hirschberg's LCS algorithm and a parallel LCS algorithm provided developed by [7].

5.1 Sequential

We investigated the single-processor, single-threaded version of an LCS algorithm using Hirschberg's algorithm[4]. Refer to section 6 for benchmarks recorded for this sequential algorithm, compared against our parallelized algorithm. This code was written in Java using the Parallel Java libraries [5].

Hirschberg's algorithm, developed by Dan Hirschberg in 1975, utilizes a dynamic programming model for finding the least cost sequence alignment between two string inputs. It is a divide and conquer version of the Needleman-Wunsch algorithm, and is used in bioinformatics for finding exact, maximal global alignments of DNA base pair sequences.[1]. The algorithm is considered superior to heuristic methods like BLAST or FASTA, but results in a space complexity of $O(\min(n, m))$ space and time complexity of $O(mn)$ time, where X and Y are two input strings and $|x| = n$ and $|y| = m$.

Algorithm Description

Hirschberg's algorithm is divided into three distinct parts, Algorithm A, B and C. Algorithm B is an optimization of Algorithm A, thus our sequential version of the algorithm relies only on Algorithms B and C directly.

We introduce string $C = c_1c_2 \cdots c_p$ that is a subsequence of string A , where $A = a_1a_2 \cdots a_m$ if and only if there is a mapping $F1, 2, \cdots, p \rightarrow 1, 2, \cdots, m$ such that $f(i) = k$ when $c_i = a_k$ and F is a monotonically increasing function. Therefore, C is a common subsequence of input strings $A = a_1a_2 \cdots a_m$ and $B = b_1b_2 \cdots b_n$ if C is a subsequence of A and C is a subsequence of B [4].

Algorithm B produces an output vector LL , which is a derivation of Algorithm A in [4] such that the matrix L as defined in that algorithm is reduced to requiring only rows i and $i - 1$, hence the use of a 2-row array K of j -elements wide.

The space complexity of Algorithm B as implemented in our Java version is $O(3n)$ deriving from the total allocated space of $3(n + 1)$ integers for the three local storage arrays. The input arrays are pre-allocated and passed by reference. The time complexity of Algorithm B is simply $O(mn)$ reflecting the nested for-loops executed mn times.

ALGORITHM-B(m, n, A, B, LL)

```

1  ▷ Initialization
2   $K[1][j] \leftarrow 0[j = 0 \dots n]$ 
3  for  $i \leftarrow 1$  to  $m$ 
4      do
5           $K[0][j] \leftarrow K[1][j](j = 0 \dots n)$ 
6          for  $j \leftarrow 1$  to  $n$ 
7              do
8                  if  $A[i] = B[j]$ 
9                      then
10                          $K[1][j] \leftarrow K[0][j - 1] + 1$ 
11                     else
12                          $K[1][j] \leftarrow \max(K[1][j - 1], K[0][j])$ 
13   $LL[j] \leftarrow K[1][j](j = 0 \dots n)$ 

```

Algorithm C recursively produces a solution to the LCS sub-problem, where at its finish returns the full longest common subsequence string. In our Java implementation of Algorithm C , the space complexity is proportional to $O(m + n)$ space. The time complexity of this algorithm is $O(mn)$, factoring in linear invocations of Algorithm- B and - C inside this function.

ALGORITHM-C(m, n, A, B, C)

```

1  if  $n = 0$ 
2    then
3       $\triangleright$  If the problem is trivial, solve it
4       $C \leftarrow \epsilon$   $\triangleright$  Empty string
5  elseif  $m = 1$ 
6    then
7      for  $j \leftarrow 0$  to  $n$ 
8        do
9           $C \leftarrow A[1]$ 
10
11 else
12    $\triangleright$  Otherwise, divide and conquer
13    $i \leftarrow \lfloor m/2 \rfloor$ 
14    $\triangleright$  Evaluate
15   ALGORITHM-B( $i, n, A_{1,i}, B_{1,n}, L1$ )
16   ALGORITHM-B( $m - i, n, A_{n,i+1}, B_{n,1}, L2$ )
17    $\triangleright$  Find  $j$ 
18   for  $j \leftarrow 0$  to  $n$ 
19     do
20        $M \leftarrow \text{MAX}(L1[j] + L2[n - j])$ 
21        $k \leftarrow \text{MIN}(j)$  such that  $L1(j) + L2(n - j) = M$ 
22
23    $\triangleright$  Solve simpler sub-problems
24   ALGORITHM-C( $i, k, A_{1,i}, B_{1,k}, C1$ )
25   ALGORITHM-C( $m - i, n - k, A_{i+1,m}, B_{k+1,n}, C2$ )
26    $\triangleright$  Resulting output
27    $C \leftarrow \text{STRCAT}(C1, C2)$ 
28
```

5.2 Parallelized

Research into the problem will help us determine the best LCS algorithm for parallelization. A modification of one of the algorithms chosen in the sequential portion of the project would be preferred for the most accurate analysis. The code will be written in C and parallelized using MPI (either OpenMPI or MPICH). Again, metrics described in section 6 will be used for our analysis of the chosen algorithm in its sequential and parallelized forms. It's not yet determined if we will investigate a pure MPI solution or consider a hybrid approach using MPI and PThreads.

6 Performance Metrics

The following table provides reported metrics for various test conditions. The variables are the length of input strings A and B , resulting in increasing memory and time requirements.

A (bytes)	B (bytes)	K	T (msec)	Speedup	Eff	EDSF
512	512	seq	6	xxx	xxx	xxx
512	512	1				xxx
512	512	2				
512	512	4				
512	512	8				
512	512	16				
1024	1024	seq		xxx	xxx	xxx
1024	1024	1				xxx
1024	1024	2				
1024	1024	4				
1024	1024	8				
1024	1024	16				
2048	2048	seq		xxx	xxx	xxx
2048	2048	1				xxx
2048	2048	2				
2048	2048	4				
2048	2048	8				
2048	2048	16				
4096	4096	seq		xxx	xxx	xxx
4096	4096	1				xxx
4096	4096	2				
4096	4096	4				
4096	4096	8				
4096	4096	16				
8192	8192	seq		xxx	xxx	xxx
8192	8192	1				xxx
8192	8192	2				
8192	8192	4				
8192	8192	8				
8192	8192	16				
16384	16384	seq		xxx	xxx	xxx
16384	16384	1				xxx
16384	16384	2				
16384	16384	4				
16384	16384	8				
16384	16384	16				
32768	32768	seq		xxx	xxx	xxx
32768	32768	1				xxx
32768	32768	2				
32768	32768	4				
32768	32768	8				
32768	32768	16				
65536	65536	seq		xxx	xxx	xxx
65536	65536	1				xxx
65536	65536	2				
65536	65536	4				
65536	65536	8				
65536	65536	16				

7 Conclusion

Bibliography

- [1] Hirschberg's algorithm. *Wikipedia*. <http://goo.gl/RyxLH>.
- [2] Longest common subsequence problem. *Wikipedia*. <http://goo.gl/5bMjT>.
- [3] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. The MIT Press, 2nd edition, 2001.
- [4] D. S. Hirschberg. A linear space algorithm for computing maximal common subsequences. *Commun. ACM*, 18:341–343, June 1975.
- [5] Alan Kaminsky. Parallel java library. *Rochester Institute of Technology*. <http://www.cs.rit.edu/~ark/pj.shtml>.
- [6] Jinxian Lin and Yiqing Lv;. Computing all longest common subsequences on mpi cluster. *Computational Intelligence and Natural Computing Proceedings (CINC), 2010 Second International Conference on*, 1:386 – 389, 2010.
- [7] Wei Liu, Ling Chen, and Lingjun Zou. A parallel lcs algorithm for biosequences alignment. *InfoScale '07: Proceedings of the 2nd international conference on Scalable information systems*, Jun 2007.
- [8] B Ben Mabrouk, H Hasni, and Z Mahjoub. Parallelization of the dynamic programming algorithm for solving the longest common subsequence problem. *Computer Systems and Applications (AICCSA), 2010 IEEE/ACS International Conference on*, pages 1 – 8, 2010.
- [9] Nuraini Rashid, Rosni Abdullah, and Abdullah Talib. Parallel homologous search with hirschberg algorithm: a hybrid mpi-pthreads solution. *ICCOMP'07: Proceedings of the 11th WSEAS International Conference on Computers*, Jul 2007.