

Fast and Cache-Oblivious Dynamic Programming with Local Dependencies

Philip Bille¹, Morten Stöckel²

¹*Technical University of Denmark*

²*IT University of Copenhagen*

March 9 2012

The Problem

What we have done

Dynamic Programming with Local Dependencies

Longest Common Subsequence

Standard Algorithms

Re-computations and I/O's

The FASTLSA Algorithm

The CO Algorithm

The FASTCO Algorithm

Experimental Results

Round Off

The Problem

- ▶ Contribution: A new approximate string matching algorithm, which is cache oblivious and provides a time/space tradeoff useful for large input strings.
- ▶ Theoretical bounds of the new method, named `FASTCO`, are at least as good as current state of the art methods.
- ▶ Practical experiments show our new method to be highly competitive for large input strings.

Alignments

- ▶ An alignment of strings $X = \{tcgacggc\}$ and $Y = \{tatcagta\}$:

<i>t</i>	<i>c</i>	<i>g</i>	<i>a</i>	-	<i>c</i>	<i>g</i>	<i>g</i>	<i>c</i>	-
<i>t</i>	-	-	<i>a</i>	<i>t</i>	<i>c</i>	<i>a</i>	<i>g</i>	<i>t</i>	<i>a</i>

- ▶ Score of an alignment: Summation of scores of aligned characters.
- ▶ Applications in bioinformatics, dictionaries, subversion etc. where n can be large.

The Problem

- ▶ Problem input: *Strings* X and Y , consisting of *characters* from Σ , where $|X| = |Y| = n$.
- ▶ Access to score function $S : (\sigma + 1) \times (\sigma + 1) \mapsto \mathbb{R}$.
- ▶ Problem output: An *alignment* of X and Y which is optimal wrt. score function S .

Alignments

- ▶ An alignment of strings $X = \{tcgacggc\}$ and $Y = \{tatcagta\}$:

<i>t</i>	<i>c</i>	<i>g</i>	<i>a</i>	-	<i>c</i>	<i>g</i>	<i>g</i>	<i>c</i>	-
<i>t</i>	-	-	<i>a</i>	<i>t</i>	<i>c</i>	<i>a</i>	<i>g</i>	<i>t</i>	<i>a</i>

- ▶ Score of an alignment: Summation of scores of aligned characters.
- ▶ Applications in bioinformatics, dictionaries, subversion etc. where n can be large.

Solving by Dynamic Programming

- ▶ Problem input: strings X and Y , where $|X| = |Y| = n$.
- ▶ Input strings form *DPM* c of $n + 1 \times n + 1$ values.
- ▶ Invariant of c : $c[i, j]$ for $0 \leq i, j \leq n$ holds the score of the optimal solution for strings $X[1 \dots i]$ and $Y[1 \dots j]$.

Filling the DPM

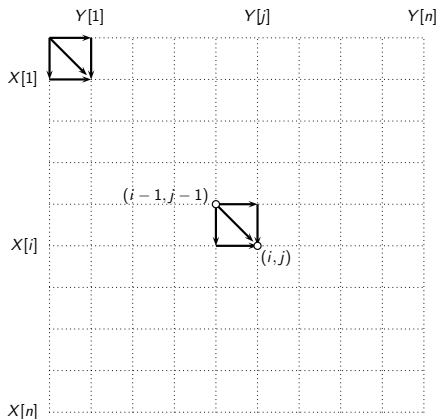
- ▶ Each value of c for $0 \leq i, j \leq n$ follows recurrence 1.

$$c[i, j] = \begin{cases} h((i, j)) & \text{if } j = 0 \vee i = 0, \\ f((i, j), X[i], Y[j]) & \text{if } i, j > 0 \end{cases} \quad (1)$$

- ▶ Functions f and h compute the values in constant time and make use of S .
- ▶ Function f derives the current value from one of three adjacent values.

Generating a Solution

- ▶ The solution can be seen as a *path* through c , consisting of all c -value derivations from upper left to bottom right.
- ▶ Graph theoretical view: Optimal path in weighted grid graph with $(n + 1) \times (n + 1)$ nodes.

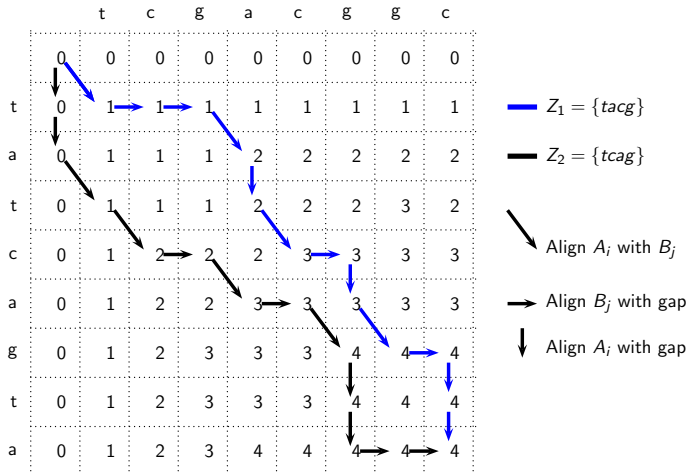


Chosen problem: LCS

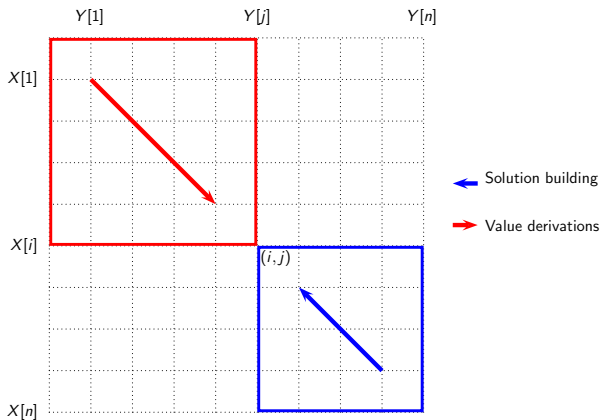
- ▶ LCS: Score function yields 1 for equal characters, 0 otherwise and goal is to maximize the alignment score.
- ▶ DPM values follow recurrence 2.

$$c[i, j] = \begin{cases} 0 & \text{if } j = 0 \vee i = 0, \\ c[i-1, j-1] + 1 & \text{if } i, j > 0 \wedge X[i] = Y[j], \\ \max \begin{cases} c[i, j-1] \\ c[i-1, j] \end{cases} & \text{if } i, j > 0 \wedge X[i] \neq Y[j] \end{cases} \quad (2)$$

An LCS example



Why is this difficult?



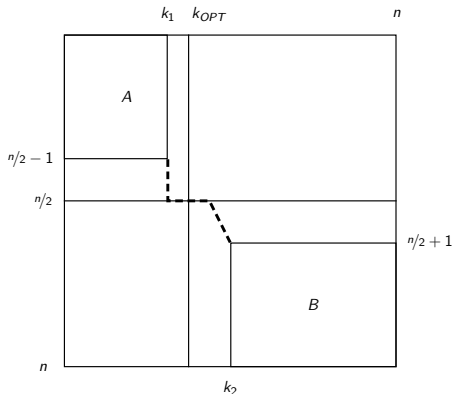
Dynamic Programming

1. Forward pass: Compute and store all $n + 1 \times n + 1$ values in a row wise manner using recurrence 2 for each value.
2. Backward pass: Starting from $c[n, n]$, follow derivations to achieve solution path through c .
3. Time: $O(n^2)$
4. Space: $O(n^2)$
5. I/O's: $O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$

Hirschberg's Algorithm

- ▶ Divide and conquer: Calculate a split such that the current problem can be computed by two smaller disjoint problems.
- ▶ Time:

$$\sum_{i=1}^{\log n} cn^2/2^{i-1} \leq 2cn^2 = O(n^2)$$
- ▶ Space: $O(n)$
- ▶ I/O's: $O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$



The Problem

What we have done

Dynamic Programming with Local Dependencies

Longest Common Subsequence

Standard Algorithms

Re-computations and I/O's

The FASTLSA Algorithm

The CO Algorithm

The FASTCO Algorithm

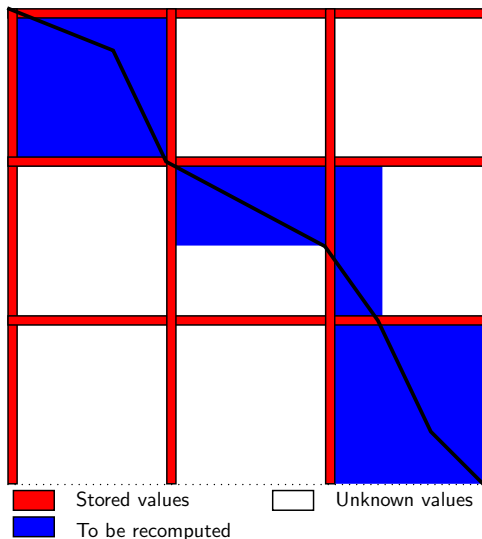
Experimental Results

Round Off

FASTLSA Intuition

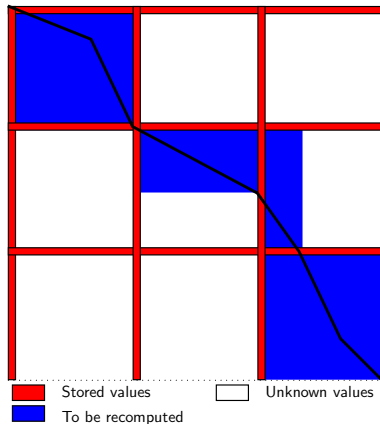
- ▶ Inefficiency of Hirschberg: Linear space means quadratic re-computations.
- ▶ FastLSA by Driga et al.: Split DPM in k parts in both directions. Store k "important" rows and columns to achieve a faster backward pass. If current problem matrix size is less than D compute directly, recurse otherwise.
- ▶ Provides a time/space tradeoff through parameter k .

FASTLSA time/space tradeoff



FASTLSA Complexities

- ▶ Forward: $O(n^2)$
- ▶ Backward: $O(n^2/k)$
- ▶ Space: $O(kn + D)$
- ▶ I/O's: $O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$



Wasting I/O's

- ▶ For each character in X , Y is accessed from left to right.
- ▶ Incurs $O(\frac{n}{B})$ I/O's per row.
- ▶ $O(\frac{n^2}{B})$ total I/O's.
- ▶ Issue: Y continuously accessed and thus brought into memory.

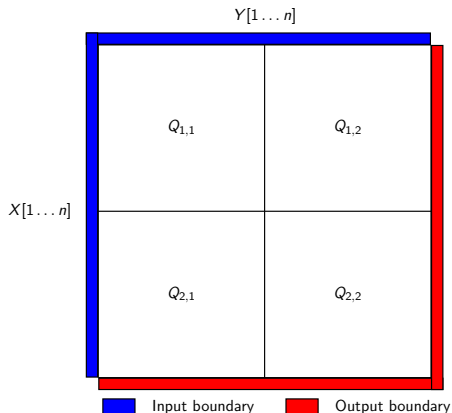
	t	c	g	a	c	g	g	c
	0	0	0	0	0	0	0	0
t	0	1	1	1	1	1	1	1
a	0	1	1	2	2	2	2	2
t	0	1	1	2	2	2	3	2
c	0	1	2	2	2	3	3	3
a	0	1	2	2	3	3	3	3
g	0	1	2	3	3	3	4	4
t	0	1	2	3	3	3	4	4
a	0	1	2	3	4	4	4	4

CO Intuition

- ▶ Deals with I/O inefficiency of left-right computation.
- ▶ Algorithm CO by Chowdhury et al.: If input matrix is constant size compute directly, otherwise split DPM in 4 quadrants and recurse. Compute boundaries of quadrants in a cache efficient manner to achieve better I/O complexity.

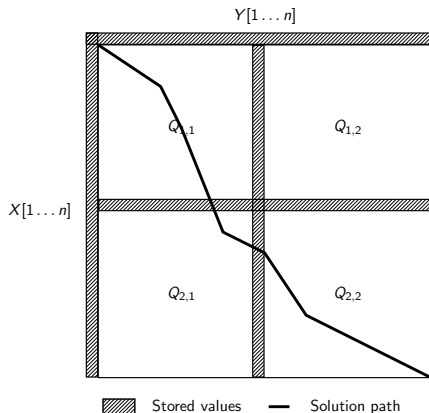
Computing Boundaries

- ▶ Problem: Given input boundary and strings, compute output boundary.
- ▶ Divide DPM into 4 quadrants and recursively compute quadrants in a left-right top-down order. Pass and combine boundaries.
- ▶ I/O's: $O\left(\frac{n^2}{BM}\right)$



CO Complexities

- ▶ Forward: $O(n^2)$
- ▶ Backward: $O(n^2)$
- ▶ Space: $O(n)$
- ▶ I/O's: $O\left(\frac{n^2}{BM} + \frac{n}{B} + 1\right)$



The Problem

What we have done

Dynamic Programming with Local Dependencies

Longest Common Subsequence

Standard Algorithms

Re-computations and I/O's

The FASTLSA Algorithm

The CO Algorithm

The FASTCO Algorithm

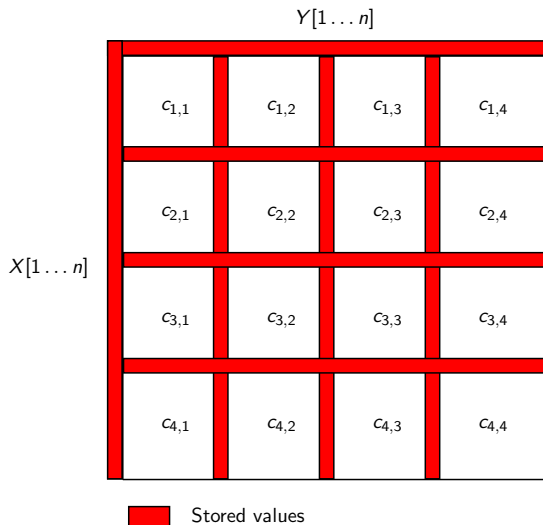
Experimental Results

Round Off

FASTCO Intuition

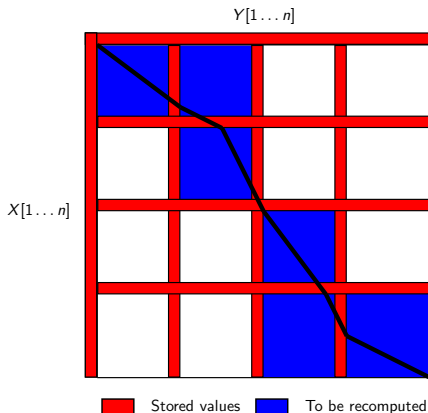
- ▶ FASTCO: Combine the two ideas from CO and FASTLSA to create a method better than both of them.
- ▶ Generalize CO to use the time/space tradeoff through parameter k .
- ▶ Result: A cache oblivious time/space tradeoff for dynamic programming with local dependencies.

FASTCO Forward Pass



FASTCO Backward Pass

- ▶ Compute $O(k)$ submatrices each of size $O\left(\frac{n^2}{k^2}\right)$.
- ▶ Backward: $O\left(\frac{n^2}{k}\right)$

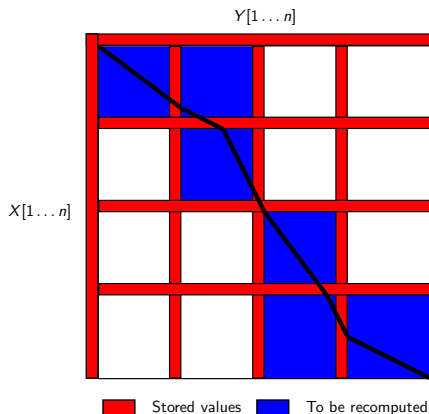


FASTCO I/O's

- ▶ I/O complexity now dependant on k .

$$I(n, k) = O\left(\frac{n^2 k}{MB} + \frac{n}{B} + 1\right)$$

- ▶ Typically $k = O(1)$



Algorithm Overview

Algorithm	Forward Pass	Backward Pass	Space	I/O	CO
FULLMATRIX	$O(n^2)$	$O(n)$	$O(n^2)$	$O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$	Yes
HIRSCHBERG	$O(n^2)$	$O(n^2)$	$O(n)$	$O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$	Yes
FASTLSA _k	$O(n^2)$	$O\left(\frac{n^2}{k} + n\right)$	$O(nk + D)$	$O\left(\frac{n^2}{B} + \frac{n}{B} + 1\right)$	No
CO	$O(n^2)$	$O(n^2)$	$O(n)$	$O\left(\frac{n^2}{BM} + \frac{n}{B} + 1\right)$	Yes
FASTCO _k	$O(n^2)$	$O\left(\frac{n^2}{k} + n\right)$	$O(nk)$	$O\left(\frac{n^2 k}{BM} + \frac{n}{B} + 1\right)$	Yes

The Problem

What we have done

Dynamic Programming with Local Dependencies

Longest Common Subsequence

Standard Algorithms

Re-computations and I/O's

The FASTLSA Algorithm

The CO Algorithm

The FASTCO Algorithm

Experimental Results

Round Off

Overview

- ▶ All algorithms were implemented in C++.
- ▶ Standardized DNA strings of sizes up to 2^{21} as input.
- ▶ Three architectures: Intel M, Intel I7, AMD X2.
- ▶ FASTLSA: $D = 1000 \times 1000$ for Intel I7 and AMD X2, 500×500 for Intel M.
- ▶ k values used: 2^i for $i \in \{3, 4, 5\}$.

Intel I7 Results

- ▶ 2.66GHz. 32KB L1, 256KB L2, 8MB L3 cache. 4GB memory

n	HB	CO	FLSA ₈	FLSA ₁₆	FLSA ₃₂	FCO ₈	FCO ₁₆	FCO ₃₂
2^{16}	0.016 <i>h</i>	0.012 <i>h</i>	0.009 <i>h</i>	0.009 <i>h</i>	0.008<i>h</i>	0.009 <i>h</i>	0.009 <i>h</i>	0.008<i>h</i>
2^{17}	0.063 <i>h</i>	0.049 <i>h</i>	0.0387 <i>h</i>	0.036 <i>h</i>	0.033<i>h</i>	0.036 <i>h</i>	0.035 <i>h</i>	0.034 <i>h</i>
2^{18}	0.251 <i>h</i>	0.194 <i>h</i>	0.150 <i>h</i>	0.144 <i>h</i>	0.132<i>h</i>	0.143 <i>h</i>	0.136 <i>h</i>	0.134 <i>h</i>
2^{19}	1.003 <i>h</i>	0.775 <i>h</i>	0.584 <i>h</i>	0.559 <i>h</i>	0.529<i>h</i>	0.565 <i>h</i>	0.539 <i>h</i>	0.530 <i>h</i>
2^{20}	4.059 <i>h</i>	3.129 <i>h</i>	2.320 <i>h</i>	2.290 <i>h</i>	2.127 <i>h</i>	2.238 <i>h</i>	2.258 <i>h</i>	2.100<i>h</i>
2^{21}	16.105 <i>h</i>	12.297 <i>h</i>	9.544 <i>h</i>	9.022 <i>h</i>	8.741 <i>h</i>	9.036 <i>h</i>	8.611 <i>h</i>	8.355<i>h</i>

AMD X2 Results

- ▶ 2.5GHz. 64KB L1, 512KB L2 cache. 4GB

n	HB	CO	FLSA ₈	FLSA ₁₆	FLSA ₃₂	FCO ₈	FCO ₁₆	FCO ₃₂
2^{16}	0.017 <i>h</i>	0.009 <i>h</i>	0.010 <i>h</i>	0.010 <i>h</i>	0.010 <i>h</i>	0.007<i>h</i>	0.007<i>h</i>	0.007<i>h</i>
2^{17}	0.069 <i>h</i>	0.037 <i>h</i>	0.041 <i>h</i>	0.039 <i>h</i>	0.038 <i>h</i>	0.028 <i>h</i>	0.027 <i>h</i>	0.026<i>h</i>
2^{18}	0.278 <i>h</i>	0.149 <i>h</i>	0.169 <i>h</i>	0.159 <i>h</i>	0.156 <i>h</i>	0.114 <i>h</i>	0.108 <i>h</i>	0.104<i>h</i>
2^{19}	1.123 <i>h</i>	0.597 <i>h</i>	0.685 <i>h</i>	0.640 <i>h</i>	0.624 <i>h</i>	0.455 <i>h</i>	0.430 <i>h</i>	0.418<i>h</i>
2^{20}	4.474 <i>h</i>	2.389 <i>h</i>	2.752 <i>h</i>	2.574 <i>h</i>	2.498 <i>h</i>	1.846 <i>h</i>	1.721 <i>h</i>	1.671<i>h</i>
2^{21}	17.949 <i>h</i>	9.442 <i>h</i>	11.007 <i>h</i>	10.337 <i>h</i>	9.950 <i>h</i>	7.278 <i>h</i>	6.873 <i>h</i>	6.685<i>h</i>

Cachegrind, L1 and L2 misses

- ▶ 64KB L1, 512KB L2 cache

L1 cache misses $\times 10^6$

n	HB	CO	FCO ₈	FCO ₁₆	FCO ₃₂	FLSA ₈	FLSA ₁₆	FLSA ₃₂
2 ¹⁶	1,090	0.855	1.682	1.980	2.070	1,293	1,162	1,154
2 ¹⁷	4,639	1.952	3.823	4.566	7.49	5,156	4,834	4,626
2 ¹⁸	18,866	5.916	11.23	12.29	17.41	20,640	19,4654	18,789
2 ¹⁹	76,654	19.85	37.27	39.59	46.55	83,201	77,630	75,355

L2 cache misses $\times 10^6$

2 ¹⁶	604.6	0.345	1.038	1.373	1.711	1,151	1,146	1,150
2 ¹⁷	3,207	0.594	1.949	2.49	5.2	4,575	4,579	4,581
2 ¹⁸	16,517	1.312	4.385	5.067	9.373	18,741	18,303	18,290
2 ¹⁹	71,629	3.416	11.689	15.12	18.792	82,131	75,219	73,117

The Problem

What we have done

Dynamic Programming with Local Dependencies

Longest Common Subsequence

Standard Algorithms

Re-computations and I/O's

The FASTLSA Algorithm

The CO Algorithm

The FASTCO Algorithm

Experimental Results

Round Off

- ▶ Two ideas taken from each their algorithm were used to create an algorithm theoretically better than both of them.
- ▶ Experiments showed `FASTCO` performing equally good or better across all three sample architectures.
- ▶ Cache efficiency once again shown to be important when dealing with large data.

Fast and Cache-Oblivious Dynamic Programming with Local Dependencies

Philip Bille¹, Morten Stöckel²

¹*Technical University of Denmark*

²*IT University of Copenhagen*

March 9 2012