# Algorithms for
# Approximate String Matching

**Yoan Pinzón**
Universidad Nacional de Colombia
Facultad de Ingeniería
Departamento de Ingeniería de Sistemas e Industrial
ypinzon@unal.edu.co
www.pinzon.co.uk

August 2006

## Levenshtein Distance

Levenshtein distance is named after the Russian scientist Vladimir Levenshtein, who devised the algorithm in 1965. If you cannot spell or pronounce Levenshtein, the metric is also called *edit distance*.

The edit distance $\delta(p, t)$ between two strings $p$ (pattern) and $t$ (text) ($m = |p|, n = |t|$) is the minimum number of insertions, deletions and replacements to make $p$ equal to $t$.

- **[Insertion]** insert a new letter $a$ into $x$. An insertion operation on the string $x = vw$ consists in adding a letter $a$, converting $x$ into $x' = vaw$.

- **[Deletion]** delete a letter $a$ from $x$. A deletion operation on the string $x = vaw$ consists in removing a letter, converting $x$ into $x' = vw$.

- **[Replacement]** replace a letter $a$ in $x$. A replacement operation on the string $x = vaw$ consists in replacing a letter for another, converting $x$ into $x' = vbw$.

► **Example 1:**

$p =$"approximate_matching"
$t =$"appropiate_meaning"

```
            1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
 String t:  a p p r o p r i a t  e     m  ϵ  e  a  n  i  n  g
            | | | | |     | | |  |     |        |        |  |  |
 String p:  a p p r o x i m a t  e     m  a  t  c  h  i  n  g
```

$\delta(t, p) = 7$

► **Example 2:**

$p =$"surgery"
$t =$"survey"

```
                   1 2 3 4 5 6 7
       String t:   s u r v e ϵ y
                   | | |   |   |
       String p:   s u r g e r y
```

$\delta(t, p) = 2$

► **Solution:** Using Dynamic Programing (DP): We need to compute a matrix $D[0..m, 0..n]$, where $D_{i,j}$ represents the minimum number of operations needed to match $p_{1..i}$ to $t_{1..j}$.

This is computed as follows:

$D[i, 0] = i$
$D[0, j] = j$
$D[i, j] = \min\{D[i-1, j] + 1, D[i, j-1] + 1, D[i, j] + \delta(p_i, t_j)\}$

$\delta(p, t) = D[m, n]$

► **Pseudo-code:**

```
1    procedure ED(p, t)    {m = |p|, n = |t|}
2    begin
3        for i ← 0 to m do D[i, 0] ← i
4        for j ← 0 to n do D[0, j] ← j
5        for i ← 1 to m do
6            for j ← 1 to n do
7                if p_i = t_j then D[i, j] ← D[i-1, j-1]
8                else
9                    D[i, j] ← min(D[i, j-1], D[i-1, j], D[i-1, j-1]) + 1
10           od
11       od
12       return D[m, n]
13   end
```

► **Running time:** $O(nm)$

## Example: $p=$"survey", $t=$"surgery"

| | | 0 ε | 1 s | 2 u | 3 r | 4 g | 5 e | 6 r | 7 y |
|---|---|---|---|---|---|---|---|---|---|
| 0 | ε | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 1 | s | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| 2 | u | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| 3 | r | 3 | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| 4 | v | 4 | 3 | 2 | 1 | 1 | 2 | 3 | 4 |
| 5 | e | 5 | 4 | 3 | 2 | 2 | 1 | 2 | 3 |
| 6 | y | 6 | 5 | 4 | 3 | 3 | 2 | 2 | **2** |

```
               1 2 3 4 5 6 7
   String t:   s u r g e r y
               | | |   |   |
   String p:   s u r v e ε y
```

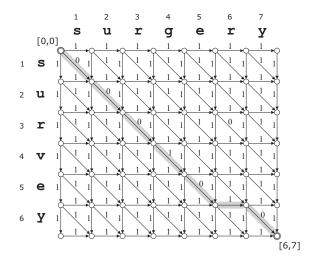## A Graph Reformulation
### for the Edit Distance problem

The Dynamic Programming matrix $D$ can be seen as a graph where the nodes are the cells and the edges represent the operations. The cost (weight) of the edges corresponds to the cost of the operations.

$\delta(t,p) =$ shortest-path from node [0,0] to the node [n,m].

## Running time: $O(nm \log(nm))$

## Example: $p=$"survey", $t=$"surgery"

## Hamming Distance

The Hamming distance H is defined only for strings of the *same* length. For two strings $p$ and $t$, $H(p, t)$ is the number of places in which the two strings differ, i.e., have different characters.

▶ **Examples:**

$H(\texttt{"pinzon"}, \texttt{"pinion"}) = 1$
$H(\texttt{"josh"}, \texttt{"jose"}) = 1$
$H(\texttt{"here"}, \texttt{"hear"}) = 2$
$H(\texttt{"kelly"}, \texttt{"belly"}) = 1$
$H(\texttt{"AAT"}, \texttt{"TAA"}) = 2$
$H(\texttt{"AGCAA"}, \texttt{"ACATA"}) = 3$
$H(\texttt{"AGCACACA"}, \texttt{"ACACACTA"}) = 6$

▶ **Pseudo-code:** too easy!!

▶ **Running time:** $O(n)$

## Approximate String Matching with $k$ Differences

▶ **Problem:** The *k-differences approximate string matching problem* is to find all occurrences of the *pattern* string $p$ in the *text* string $t$ with at most $k$ differences (substitution, insertions, deletions).

▶ **Solution:** Using DP

$D[i, 0] = i$
$D[0, j] = 0$
$D[i, j] = \min\{D[i-1, j] + 1, D[i, j-1] + 1, D[i, j] + \delta(p_i, t_j)\}$

if $D[m, j] \leq k$ then we say that $p$ occurs at position $j$ of $t$.

## ▶ Pseudo-code:

```
1   procedure KDifferences(p, t, k)    {m = |p|, n = |t|}
2   begin
3       for i ← 0 to m do D[i, 0] ← i
4       for j ← 0 to n do D[0, j] ← 0
5       for i ← 1 to m do
6           for j ← 1 to n do
7               if pᵢ = tⱼ then D[i, j] ← D[i − 1, j − 1]
8               else
9                   D[i, j] ← min(D[i, j − 1], D[i − 1, j], D[i − 1, j − 1]) + 1
10          od
11      od
12      for j ← 0 to n do
13          if D[m][j] ≤ k then Output(j)
14      od
15  end
```

## ▶ Running time: $O(nm)$

## ▶ Example:

$p =$ "CDDA",
$t =$ "CADDACDACDBACBA"
$k = 1$

|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
|   |   | $\epsilon$ | C | A | D | D | A | C | D | A | C | D | B | A | C | B | A |
| 0 | $\epsilon$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | C | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 |
| 2 | D | 2 | 1 | 1 | 1 | 1 | 2 | 1 | 0 | 1 | 1 | 0 | 1 | 2 | 1 | 1 | 2 |
| 3 | D | 3 | 2 | 2 | 1 | 1 | 2 | 2 | 1 | 1 | 2 | 1 | 1 | 2 | 2 | 2 | 2 |
| 4 | A | 4 | 3 | 2 | 2 | 2 | **1** | 2 | 2 | **1** | 2 | 2 | 2 | **1** | 2 | 3 | 2 |

⇑ (position 5) ⇑ (position 8) ⇑ (position 12)

$p$ occurs in $t$ ending at positions 5, 8 and 12.

## Longest Common Subsequence
### Preliminaries

For two sequences $x = x_1 \cdots x_m$ and $y_1 \cdots y_n$ $(n \geq m)$

we say that $x$ is a *subsequence* of $y$ and equivalently, $y$ is a *supersequence* of $x$, if for some $i_1 < \cdots < i_p$, $x_j = y_{i_j}$.

Given a finite set of sequences, $S$, a *longest common subsequence* (LCS) of $S$ is a longest possible sequence $s$ such that each sequence in $S$ is a supersequence of $s$.

Example: $y=$"longest", $x=$"large"

```
String y:  l o n g e s t
           |     | | |
String x:  l a r g e
```

LCS$(y,x)=$"lge"

---

▶ **Problem:** The *Longest Common Subsequence* (LCS) of two strings, $p$ and $t$, is a subsequence of both $p$ and of $t$ of maximum possible length.

▶ **Solution:** Using Dynamic Programing: We need to compute a matrix $L[0..m, 0..n]$, where $L_{i,j}$ represent the LCS for $p_{1..i}$ and $t_{1..j}$.

This is computed as follows:

$$L[i,j] = \begin{cases} 0, & \text{if either } i = 0 \text{ or } j = 0 \\ L[i-1, j-1] + 1, & \text{if } p_i = t_j \\ \max\{L[i-1,j], L[i,j-1]\}, & \text{if } p_i \neq t_j \end{cases}$$

▶ **Pseudo-code:**

```
1   procedure LCS(p, t)    {m = |p|, n = |t|}
2   begin
3       for i ← 0 to m do L[i, 0] ← 0
4       for j ← 0 to n do L[0, j] ← 0
5       for i ← 1 to m do
6           for j ← 1 to n do
7               if pi = tj then L[i, j] ← L[i − 1, j − 1] + 1
8               else
9                   if L[i, j − 1] > L[i − 1, j] then L[i, j] ← L[i, j − 1]
10                  else L[i, j] ← L[i − 1, j]
11          od
12      od
13      return L[m][n]
14  end
```

▶ **Running time:** $O(nm)$

---

▶ **Example 1:** $p =$"survey" and $t =$"surgery".

|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|
|   |   | $\epsilon$ | s | u | r | g | e | r | y |
| 0 | $\epsilon$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | s | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | u | 0 | 1 | 2 | 2 | 2 | 2 | 2 | 2 |
| 3 | r | 0 | 1 | 2 | 3 | 3 | 3 | 3 | 3 |
| 4 | v | 0 | 1 | 2 | 3 | 3 | 3 | 3 | 3 |
| 5 | e | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 |
| 6 | y | 0 | 1 | 2 | 3 | 3 | 4 | 4 | **5** |

```
              1   2   3   4   5   6   7
String t:     s   u   r   g   e   r   y
              |   |   |       |    /
String p:     s   u   r   v   e   y
```

LCS$(p, t) =$"surey"

LLCS$(p, t) = L[6, 7] = 5$

▶ **Example 2:**

$p =$"ttgatacatt"
$t =$"gaataagacc"

|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | $\epsilon$ | g | a | a | t | a | a | g | a | c | c |
| 0 | $\epsilon$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | t | 0 | 0 | 0 | 0 | **1** | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | t | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 3 | g | 0 | **1** | 1 | 1 | 1 | 1 | 1 | **2** | 2 | 2 | 2 |
| 4 | a | 0 | 1 | **2** | 2 | 2 | 2 | 2 | 2 | **3** | 3 | 3 |
| 5 | t | 0 | 1 | 2 | 2 | **3** | 3 | 3 | 3 | 3 | 3 | 3 |
| 6 | a | 0 | 1 | 2 | **3** | 3 | **4** | 4 | 4 | 4 | 4 | 4 |
| 7 | c | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 4 | **5** | 5 |
| 8 | a | 0 | 1 | 2 | 3 | 3 | 4 | **5** | 5 | 5 | 5 | 5 |
| 9 | t | 0 | 1 | 2 | 3 | **4** | 4 | 5 | 5 | 5 | 5 | **5** |

LCS$(p, t) =$?

LLCS$(p, t) = L[9, 10] = 5$

# Part II

A Fast and Practical Bit-Vector Algorithm for the Longest Common Subsequence Problem

## Some More Definitions

The ordered pair of *positions* $i$ and $j$ of $L$, denoted $[i, j]$, is a *match* iff $x_i = y_j$.

If $[i, j]$ is a match, and an LCS $s_{i,j}$ of $x_1 x_2 ... x_i$ and $y_1 y_2 ... y_j$ has length $k$, then $k$ is the *rank* of $[i, j]$.

The match $[i, j]$ is *$k$-dominant* if it has rank $k$ and for any other pair $[i', j']$ of rank $k$, either $i' > i$ and $j' \leq j$ or $i' \leq i$ and $j' > j$.

Computing the $k$-dominant matches is all that is needed to solve the LCS problem, since the LCS of $x$ and $y$ has length $p$ iff the maximum rank attained by a dominant match is $p$.

A match $[i, j]$ *precedes* a match $[i', j']$ if $i < i'$ and $j < j'$.

Let $r$ be the total number of matches points, and $d$ be the total number of dominant points (all ranks). Then $0 \leq p \leq d \leq r \leq nm$.

Let $\mathcal{R}$ denote a partial order relation on the set of matches in $L$.

A set of matches such that in any pair one of the matches always precedes the other in $\mathcal{R}$ constitutes a *chain* relative to the partial order relation $\mathcal{R}$.

A set of matches such that in any pair neither element of the pair precedes the other in $\mathcal{R}$ constitutes an *antichain*.

Sankoff and Sellers (1973) observed that the LCS problem translates to finding a longest *chain* in the *poset* of matches induced by $\mathcal{R}$.

A decomposition of a poset into antichains partitions the poset into the minimum possible number of antichains.

$t =$ "aaagtgacctagcccg"
$p =$ "tccagatg"

$LCS(p, t) =$ "tccagg"

$t$ = "aaagtgacctagcccg"
$p$ = "tccagatg"

$t$

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| | a | a | a | g | t | g | a | c | c | t | a | g | c | c | c | g |

$p$: 1 t, 2 c, 3 c, 4 a, 5 g, 6 a, 7 t, 8 g

## A Simple Bit-Vector Algorithm

Here we will make use of word-level parallelism in order to compute the matrix $L$ more efficiently.

The algorithm is based on the $O(1)$-time computation of each column in $L$ by using a bit-parallel formula under the assumption that $m \leq w$, where $w$ is the number of bits in a machine word or $O(nm/w)$-time for the general case.

An interesting property of the LCS allows to represent each column in $L$ by using $O(1)$-space. That is, the values in the columns (rows) of $L$ increase by at most one. i.e. $\Delta L[i,j] = L[i,j] - L[i-1,j] \in \{0,1\}$ for any $(i,j) \in \{1..m\} \times \{1..n\}$.

In other words $\Delta L$ will use the relative encoding of the dynamic programming table $L$.

$\Delta L'$ is defined as NOT $\Delta L$.

Example: $x=$ "ttgatacatt" and $y=$ "gaataagacc".

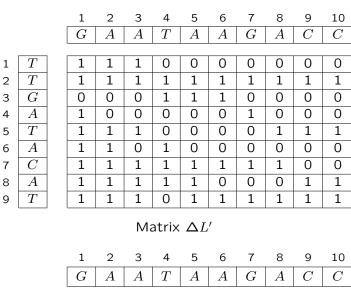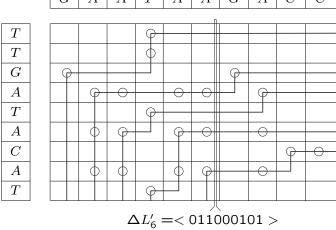| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\epsilon$ | $G$ | $A$ | $A$ | $T$ | $A$ | $A$ | $G$ | $A$ | $C$ | $C$ |
| 0 | $\epsilon$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | $T$ | 0 | 0 | 0 | 0 | **1** | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | $T$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 3 | $G$ | 0 | **1** | 1 | 1 | 1 | 1 | 1 | **2** | 2 | 2 | 2 |
| 4 | $A$ | 0 | 1 | **2** | 2 | 2 | 2 | 2 | 2 | **3** | 3 | 3 |
| 5 | $T$ | 0 | 1 | 2 | 2 | **3** | 3 | 3 | 3 | 3 | 3 | 3 |
| 6 | $A$ | 0 | 1 | 2 | **3** | 3 | **4** | 4 | 4 | 4 | 4 | 4 |
| 7 | $C$ | 0 | 1 | 2 | 3 | 3 | 4 | 4 | 4 | 4 | **5** | 5 |
| 8 | $A$ | 0 | 1 | 2 | 3 | 3 | 4 | **5** | 5 | 5 | 5 | 5 |
| 9 | $T$ | 0 | 1 | 2 | 3 | **4** | 4 | 5 | 5 | 5 | 5 | **5** |

(a) Matrix $L$

| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\epsilon$ | $G$ | $A$ | $A$ | $T$ | $A$ | $A$ | $G$ | $A$ | $C$ | $C$ |
| 0 | $\epsilon$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | $T$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | $G$ | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 4 | $A$ | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| 5 | $T$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| 6 | $A$ | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 1 |
| 7 | $C$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| 8 | $A$ | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| 9 | $T$ | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

(b) Matrix $\Delta L$

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $G$ | $A$ | $A$ | $T$ | $A$ | $A$ | $G$ | $A$ | $C$ | $C$ |
| 1 | $T$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | $T$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 3 | $G$ | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 4 | $A$ | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 5 | $T$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 6 | $A$ | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | $C$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| 8 | $A$ | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 |
| 9 | $T$ | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 1 |

Matrix $\Delta L'$



$\Delta L'_6 = <011000101>$

20

21

First we compute the array $M$ of the vectors that result for each possible text character. If both the strings $x$ and $y$ range over the alphabet $\Sigma$ then $M[\Sigma]$ is defined as $M[\alpha]_i = 1$ if $y_i = \alpha$ else 0.

Example: $x=$ "ttgatacatt" and $y=$ "gaataagacc".

| | | $A$ | $C$ | $G$ | $T$ |
|---|---|---|---|---|---|
| 1 | $T$ | 0 | 0 | 0 | 1 |
| 2 | $T$ | 0 | 0 | 0 | 1 |
| 3 | $G$ | 0 | 0 | 1 | 0 |
| 4 | $A$ | 1 | 0 | 0 | 0 |
| 5 | $T$ | 0 | 0 | 0 | 1 |
| 6 | $A$ | 1 | 0 | 0 | 0 |
| 7 | $C$ | 0 | 1 | 0 | 0 |
| 8 | $A$ | 1 | 0 | 0 | 0 |
| 9 | $T$ | 0 | 0 | 0 | 1 |

(a) Matrix M

| | | $A$ | $C$ | $G$ | $T$ |
|---|---|---|---|---|---|
| 1 | $T$ | 1 | 1 | 1 | 0 |
| 2 | $T$ | 1 | 1 | 1 | 0 |
| 3 | $G$ | 1 | 1 | 0 | 1 |
| 4 | $A$ | 0 | 1 | 1 | 1 |
| 5 | $T$ | 1 | 1 | 1 | 0 |
| 6 | $A$ | 0 | 1 | 1 | 1 |
| 7 | $C$ | 1 | 0 | 1 | 1 |
| 8 | $A$ | 0 | 1 | 1 | 1 |
| 9 | $T$ | 1 | 1 | 1 | 0 |

(a) Matrix M'

---

$\boxed{\textbf{Basic steps of the algorithm}}$

1. Computation of $M$ and $M'$

2. Computation of matrix $\Delta L'$ $(L)$ as follows:

$$L = \begin{cases} 2^m - 1, & \text{for } j = 0 \\ (L_{j-1} + (L_{j-1} \text{ AND } M(y_j))) \text{ OR } (L_{j-1} \text{ AND } M'(y_j)), & \text{for } j \in \{1..n\} \end{cases}$$

3. Let LLCS be the number of times a carry took place.

## Pseudo-code

```
LLCS(x, y)    ▷ n = |y|, m = |x|, p = 0
 1 begin
 2      ▷ Preprocessing
 3      for i ← 1 until m do
 4          M[α](i) ← yᵢ = α
 5          M'[α](i) ← yᵢ ≠ α
 6      ▷ Initialization
 7      L₀ = 2ᵐ − 1
 8      ▷ TheMainStep
 9      for j ← 1 until n do
10          Lⱼ ← (Lⱼ₋₁ + (Lⱼ₋₁ AND M[yⱼ])) OR (Lⱼ₋₁ AND M'[yⱼ])
11          if Lⱼ(m + 1) = 1 then p++
12      return p
13 end
```

$$\text{LLCS}(x, y) \quad \triangleright\ n = |y|,\ m = |x|,\ p = 0$$
$$L_0 = 2^m - 1$$
$$M[\alpha](i) \leftarrow y_i = \alpha$$
$$M'[\alpha](i) \leftarrow y_i \neq \alpha$$
$$L_j \leftarrow (L_{j-1} + (L_{j-1}\ \text{AND}\ M[y_j]))\ \text{OR}\ (L_{j-1}\ \text{AND}\ M'[y_j])$$
$$\text{if}\ L_j(m + 1) = 1\ \text{then}\ p{+}{+}$$

## Illustration of $\triangle L'_4$ Computation

for $x=$ "gaataagacc" and $y=$ "ttgatacatt".



(a) Matrix $\triangle L'$        (b) Matrix M

## Page 26

$$L_4 \leftarrow \overbrace{(\; \underbrace{L_3 + \overbrace{\overline{(L_3 \& M_T)}}^{(1)}\;}_{(3)})\;|\; \overbrace{\overline{(L_3 \& M'_T)}}^{(2)}}^{(4)}$$

|       |   |   |   |   |   |   |   |   |   |   |
|-------|---|---|---|---|---|---|---|---|---|---|
| $L_3$ | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | & |
| $M_T$ | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |   |
| (1)   |   |   |   |   |   |   |   |   |   |   |

|        |   |   |   |   |   |   |   |   |   |   |
|--------|---|---|---|---|---|---|---|---|---|---|
| $L_3$  | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | & |
| $M'_T$ | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 |   |
| (2)    |   |   |   |   |   |   |   |   |   |   |

|        |   |   |   |   |   |   |   |   |   |   |
|--------|---|---|---|---|---|---|---|---|---|---|
| $L_3$  | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | + |
| (1)    | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |   |
| (3)    |   |   |   |   |   |   |   |   |   |   |

|      |   |   |   |   |   |   |   |   |   |   |   |
|------|---|---|---|---|---|---|---|---|---|---|---|
| (3)  | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | \| |
| (2)  |   | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |   |
| (4)  |   |   |   |   |   |   |   |   |   |   |   |

## Page 27

$$L_4 \leftarrow \overbrace{(\; \underbrace{L_3 + \overbrace{\overline{(L_3 \& M_T)}}^{(1)}\;}_{(3)})\;|\; \overbrace{\overline{(L_3 \& M'_T)}}^{(2)}}^{(4)}$$

|       |   |   |   |   |   |   |   |   |   |   |
|-------|---|---|---|---|---|---|---|---|---|---|
| $L_3$ | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | & |
| $M_T$ | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |   |
| (1)   | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |   |

|        |   |   |   |   |   |   |   |   |   |   |
|--------|---|---|---|---|---|---|---|---|---|---|
| $L_3$  | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | & |
| $M'_T$ | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 |   |
| (2)    | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |   |

|        |   |   |   |   |   |   |   |   |   |   |
|--------|---|---|---|---|---|---|---|---|---|---|
| $L_3$  | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | + |
| (1)    | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |   |
| (3)    | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |

|      |   |   |   |   |   |   |   |   |   |   |   |
|------|---|---|---|---|---|---|---|---|---|---|---|
| (3)  | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | \| |
| (2)  |   | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |   |
| (4)  | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |   |

## Automata for Addition



## Experimental Results