## CSE 421 Algorithms

Richard Anderson Lecture 19 Longest Common Subsequence

### Longest Common Subsequence

- C=c<sub>1</sub>...c<sub>g</sub> is a subsequence of A=a<sub>1</sub>...a<sub>m</sub> if C can be obtained by removing elements from A (but retaining order)
- LCS(A, B): A maximum length sequence that is a subsequence of both A and B

ocurranec occurrence

attacggct

tacgacca

Determine the LCS of the following strings

#### BARTHOLEMEWSIMPSON

KRUSTYTHECLOWN

# String Alignment Problem

Align sequences with gaps
 CAT TGA AT

#### CAGAT AGGA

- Charge  $\delta_{\textbf{x}}$  if character x is unmatched
- Charge  $\gamma_{xy}$  if character x is matched to character y

Note: the problem is often expressed as a minimization problem, with  $\gamma_{xx}=0$  and  $\delta_x>0$ 

### LCS Optimization

- $A = a_1 a_2 ... a_m$
- $B = b_1 b_2 \dots b_n$
- Opt[ j, k] is the length of LCS(a<sub>1</sub>a<sub>2</sub>...a<sub>j</sub>, b<sub>1</sub>b<sub>2</sub>...b<sub>k</sub>)

### Optimization recurrence

If  $a_j = b_k$ , Opt[ j,k ] = 1 + Opt[ j-1, k-1 ]

If  $a_j != b_k$ , Opt[j,k] = max(Opt[j-1,k], Opt[j,k-1])

### Give the Optimization Recurrence for the String Alignment Problem

- Charge  $\delta_x$  if character x is unmatched
- Charge  $\gamma_{xy}$  if character x is matched to character y

Opt[ j, k] =

Let  $a_j = x$  and  $b_k = y$ Express as minimization



Code to compute Opt[j,k]



else if Opt[i-1, j] >= Opt[i, j-1]

if A[i] = B[j] { Opt[i,j] := 1 + Opt[i-1,j-1]; Best[i,j] := Diag; }

 $\{ \ Opt[i, j] := Opt[i-1, j], Best[i, j] := Left; \}$ else  $\{ \ Opt[i, j] := Opt[i, j-1], Best[i, j] := Down; \}$ 

# How good is this algorithm?

• Is it feasible to compute the LCS of two strings of length 100,000 on a standard desktop PC? Why or why not.



- The computation can be done in O(m+n) space if we only need one column of the Opt values or Best Values
- The algorithm can be run from either end of the strings

# Computing LCS in O(nm) time and O(n+m) space

- Divide and conquer algorithm
- Recomputing values used to save space



### **Constrained LCS**

- LCS<sub>i,j</sub>(A,B): The LCS such that

   a<sub>1</sub>,...,a<sub>i</sub> paired with elements of b<sub>1</sub>,...,b<sub>j</sub>
   a<sub>i+1</sub>,...,a<sub>m</sub> paired with elements of b<sub>i+1</sub>,...,b<sub>n</sub>
- LCS<sub>4,3</sub>(abbacbb, cbbaa)

### A = RRSSRTTRTS B=RTSRRSTST

Compute LCS<sub>5,0</sub>(A,B), LCS<sub>5,1</sub>(A,B),...,LCS<sub>5,9</sub>(A,B)

### A = **RRSSRTTRTS** B=RTSRRSTST

Compute LCS<sub>5,0</sub>(A,B), LCS<sub>5,1</sub>(A,B),...,LCS<sub>5,9</sub>(A,B)

j	left	right
0	0	4
1	1	4
2	1	3
3	2	3
4	3	3
5	3	2
6	3	2
7	3	1
8	4	1
9	4	0



### Divide and Conquer

- $A = a_1, \dots, a_m$   $B = b_1, \dots, b_n$
- Find j such that
  - LCS(a1...am/2, b1...bj) and
  - $-LCS(a_{m/2+1}...a_m,b_{j+1}...b_n)$  yield optimal solution
- Recurse



Prove by induction that T(m,n) <= 2cmn

# Memory Efficient LCS Summary

- We can afford O(nm) time, but we can't afford O(nm) space
- If we only want to compute the length of the LCS, we can easily reduce space to O(n+m)
- Avoid storing the value by recomputing values
  - Divide and conquer used to reduce problem sizes