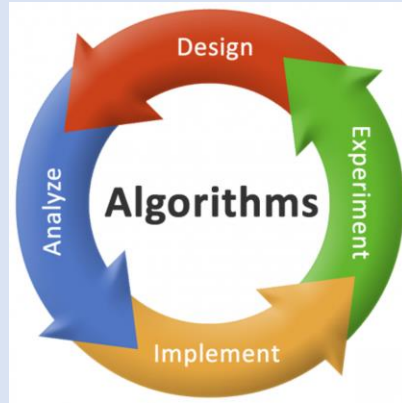


Dynamic Programming Sequence Alignment

COP 3503

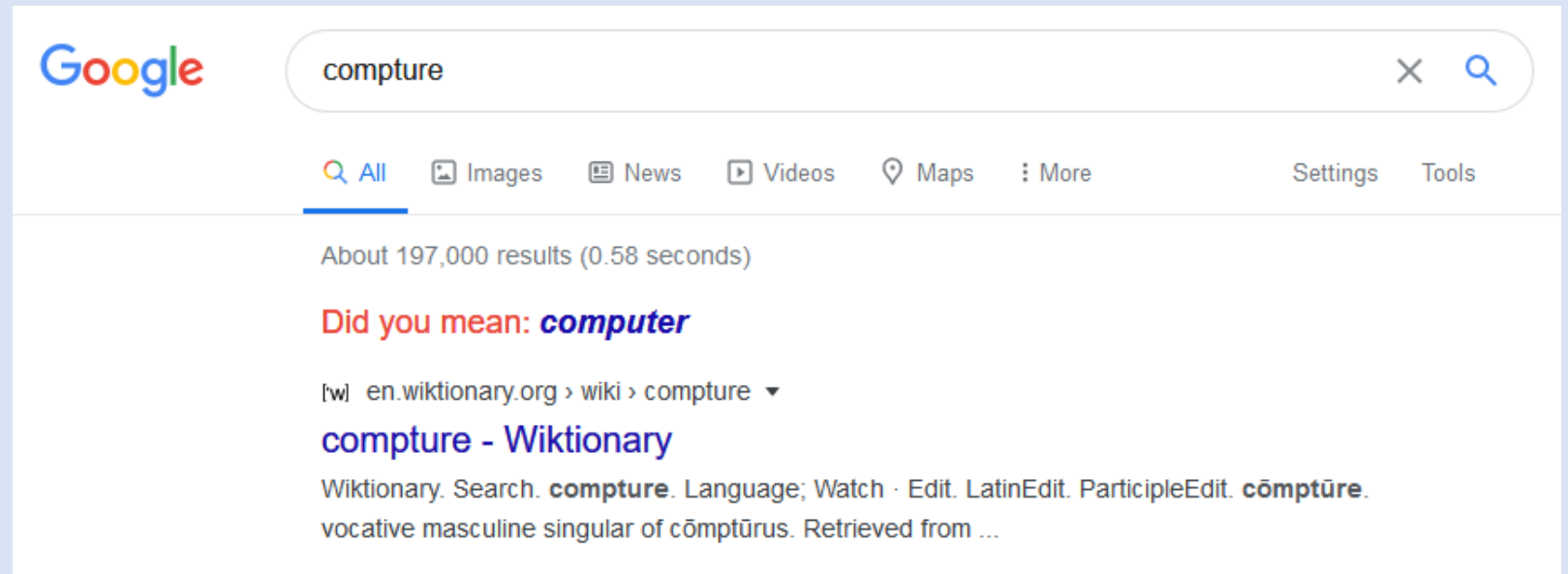
Fall 2021

Department of Computer Science
University of Central Florida
Dr. Steinberg



The Problem

- Dictionaries and websites are getting more and more useful.
- What happens when you type “ocurrance”?
 - “Perhaps you mean occurrence?”
 - The dictionary will search its entries for the word most “similar” to the one typed.



How can we model similarity between two strings?

o-currance
occurrence

one gap
one mismatch

o-curr-ance
occurre-nce

three gaps
No mismatch

Another Application – Molecular Biology

- The organism's genome
 - Divided into giant linear DNA molecules known as chromosomes
 - String over the alphabet {A, C, G, T} to determine the properties of the organism
 - Adenine, Cytosine, Guanine, Thymine
- Why is this important?
- Let's say we have two bacteria strains to observe. Sequence alignment allows scientists to observe any possible substrings that have similarities.

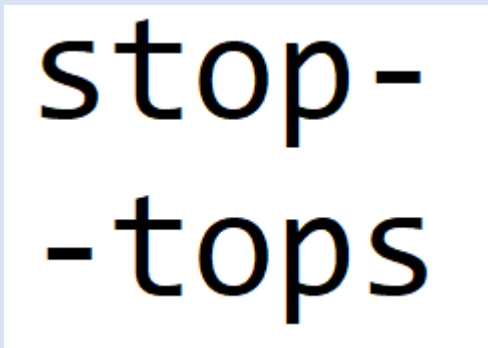
The Sequence Alignment Problem

Definition

- Let $X = x_1x_2\dots x_m$ and $Y = y_1y_2\dots y_n$
- $\{1, 2, \dots, m\}$ and $\{1, 2, \dots, n\}$ represents different positions in strings X and Y
- Matching is a set of ordered pairs with the property that each items occurs in at most one pair.
- We say that a matching M of these two sets is an alignment if there are no “crossing” pairs: if $(i,j), (i',j') \in M$ and $i < i'$ then $j < j'$

Example

- An alignment provides a way to line up two strings



stop-
-tops

- Corresponds $\{(2,1)(3,2)(4,3)\}$

Sequence Alignment Problem

- Suppose M is a given alignment between X and Y
 - *gap penalty*: each gap incurs at cost $\delta > 0$
 - *mismatch cost*: for each pair of letters p, q in the alphabet, there is a mismatch cost α_{pq} for lining up p and q
 - Assumption $\alpha_{pp} = 0$
 - The *cost* of M is the sum of gap and mismatch costs
- Objective: Find an *optimal alignment*, that means an alignment of minimum cost
- Values δ and α_{pq} are provided
- The lower the cost, the more similar X and Y are
- Looking back at *ocurrence* and *occurrence*, the first alignment is better if and only if $\delta + \alpha_{ae} < 3\delta$

Designing the Dynamic Programming Algorithm

- In the optimal alignment M :
 - either $(m,n) \in M$ or $(m,n) \notin M$
- In any alignment
 - if $(m,n) \notin M$, then either the m^{th} position of X or the n^{th} position of Y are not matched in M
 - Can be proven using proof by contradiction
- Property: 1 of the following is true for an alignment
 - $(m,n) \in M$
 - The m^{th} position of X is not matched
 - The n^{th} position of Y is not matched

Designing the Dynamic Programming Algorithm

- Let's define $\text{OPT}(i,j)$ as the minimum cost of an alignment between $x_1x_2\dots x_m$ and $y_1y_2\dots y_n$
- Recursively define $\text{OPT}(m,n)$

Case 1

- Pay $\alpha_{x_m y_n}$ then optimally align $x_1x_2\dots x_{m-1}$ and $y_1y_2\dots y_{n-1}$
 $\text{OPT}(m,n) = \alpha_{x_m y_n} + \text{OPT}(m-1,n-1)$

Case 2

- Pay gap cost δ , then optimally align $x_1x_2\dots x_{m-1}$ and $y_1y_2\dots y_n$
 $\text{OPT}(m,n) = \delta + \text{OPT}(m-1,n)$

Case 3

- Pay gap cost δ , then optimally align $x_1x_2\dots x_m$ and $y_1y_2\dots y_{n-1}$
 $\text{OPT}(m,n) = \delta + \text{OPT}(m,n-1)$

Designing the Dynamic Programming Algorithm

- Property: The minimum alignment costs satisfy the following recurrence, for $i \geq 1, j \geq 1$
- $\text{OPT}(i,j) = \min(\alpha_{x_i y_j} + \text{OPT}(i-1,j-1), \delta + \text{OPT}(i-1,j), \delta + \text{OPT}(i,j-1))$
- The minimal optimal alignment is achieved through these 3 values.

The Algorithm

Alignment(X,Y)

array $A[o..m, o..n]$

initialize array $A[i,0] = i\delta$ for each i

initialize array $A[0,j] = j\delta$ for each j

for $j = 1$ to n

 for $i = 1$ to m

$A[i,j] = \min\{\alpha_{x_i y_j} + A[i-1,j-1], \delta + A[i-1,j], \delta + A[i,j-1]\}$

return $A[m,n]$

Running Time: $O(mn)$