

## **Dynamic Programming**

- Simplex for LP: Greedy algorithm, makes a locally optimal choice.
- For many problems, we need a different approach called **Dynamic Programming**
- Finds solutions for problems with lots of **overlapping sub-problems.** Essentially, we try to solve each sub-problem **only once**.
- **Optimal substructure:** optimal solutions of **subproblems** can be used to find the optimal solutions of the **overall problem**.

**Example:** Finding the shortest path in a graph.



## **Dynamic Programming**

Typically, a dynamic programming solution is constructed using a series of steps:

- 1. Characterise the **structure** of an optimal solution.
- 2. Recursively define the value of an optimal solution.
- Compute the value of an optimal solution in a bottom-up (→ iteration) or top-down (→ recursion) fashion. That is, build it from the results of smaller solutions either iteratively from the bottom or recursively from the top.

## A Simple Example: Fibonacci numbers

**Fibonacci sequence:** The *n*-th number is the sum of the previous two. This can be implemented using a simple recursive algorithm:

#### function FIBONACCI(n)

if n = 0 then return 0 if n = 1 then return 1 return FIBONACCI(n - 1) + FIBONACCI(n - 2)

Problem: Overlapping sub-problems: Computing FIBONACCI(n-1) overlaps FIBONACCI(n-2)  $\rightsquigarrow$  exponential time complexity!



## A Simple Example (2)

Define **map object** m, maps each instance of FIBONACCI that has already been calculated to its result. **Modified recursion** requires only O(n) time:

var m; m[0] = 0; m[1] = 1function FIBONACCI(n) if m does not contain key nm[n] = FIBONACCI(n - 1) + FIBONACCI(n - 2)return m[n]

Or define array f and use **iteration:** f[0] = 0, f[1] = 1. FIBONACCI(n)

```
for i = 2 upto n step 1 do

f[i] = f[i-1] + f[i-2]

return f[n]
```

#### Another Example: Optimal Binary Search Trees

- **BST:** Tree where the key values are stored in the nodes, and the keys are ordered lexicographically.
- For each internal node all keys in the left subtree are less than the keys in the node, and all the keys in the right subtree are greater.
- Knowing the probabilities of searching each one of the keys makes it easy to compute the expected cost of accessing the tree.

An **OBST** is a BST with *minimal expected search costs*.



## OBST

- Keys  $k_1, \ldots, k_n$  in lexicographical order,
- **Probabilities** of accessing keys  $p_1, \ldots, p_n$ .
- Depth  $D_T(k_m)$  of node  $k_m$  in tree T.  $D_T(\text{root}) = 0$
- $T^{ij}$ : tree constructed from keys  $k_i, \ldots, k_j$
- **Costs:** number of comparisons done in a search.
- **Expected costs:** expected number of comparisons done during search in tree, given the acess probabilities  $p_i$

#### **OBST: Expected costs**

Definiton of expected costs of tree constructed from keys  $k_i, \ldots, k_j$ :

Probabilities 1/8 1/32 1/32 1/16 1/4 1/2  $C_{1,6} = 1 \cdot 1/16 + 2 \cdot (1/32 + 1/4) + 3 \cdot (1/8 + 1/32 + 1/2)$ = 85/32



## OBST

- Key observation: each subtree of an optimal tree is itself optimal (replacing a subtree with a better one lowers the costs of entire tree)
- Consider tree  $T^{ij}$  with root node  $r(T) = k_r$ .



## **Expected costs of tree** $T = T^{ij}$

$$\begin{split} C_{i,j} &= \sum_{m=i}^{j} p_m (D_T(k_m) + 1) \\ &= \sum_{m=i}^{r-1} p_m (D_T(k_m) + 1) + p_r + \sum_{m=r+1}^{j} p_m (D_T(k_m) + 1) \\ &= \sum_{m=i}^{r-1} p_m ((D_{T_L^r}(k_m) + 1) + 1) + \underbrace{p_r}_{\text{root}} + \underbrace{\sum_{m=r+1}^{j} p_m ((D_{T_R^r}(k_m) + 1) + 1)}_{\text{C(right subtree)|root=r}} \\ &= C(T_L^r) + \sum_{m=i}^{r-1} p_m + p_r + C(T_R^r) + \sum_{m=r+1}^{j} p_m \\ &= C_{i,r-1} + C_{r+1,j} + \sum_{m=i}^{j} p_m, \quad i \le r \le j. \end{split}$$

## **OBST:** algorithm

#### **Recursive algorithm**:

- consider every node as being the root
- split rest of the keys into left and right subtrees and recursively calculate their costs.

$$C_{i,i} = p_i$$

$$C_{i,j} = 0 \forall j < i \quad (\text{tree with no nodes})$$

$$C_{i,j} = \sum_{m=i}^{j} p_m + \min_{i \le r \le j} [C_{i,r-1} + C_{r+1,j}]$$

Use **memoization** to avoid solving the same problem over and over. Or use **iterative** algorithm.

## DP for an OBST

- Precompute  $P_{ij} = \sum_{m=i}^{j} p_m$ .
- Fill C-matrix by diagonals (start with main diagonal, move up-right)
- $\bullet\,$  Store "winning" root index in matrix R



Find tree by backtracking: start in upper right corner R<sub>1,n</sub>
→ root of full tree, say root = k.
Right subtree: proceed with R(k + 1, n)
→ root of right subtree T<sub>k+1,n</sub>, say R(k + 1, n) = r.
Draw edge k → r.
Left subtree: R(1, k - 1) = l → root of left subtree, edge k → l.
Recurse.

#### Computations

In our case:



 $E[cost] = \mathbf{1} \cdot 1/4 + \mathbf{2} \cdot (1/2 + 1/8) + \mathbf{3} \cdot (1/32) + \mathbf{4} \cdot (1/16 + 1/32) \\= 1/32[8 + 2(16 + 4) + \mathbf{3} + 4(2 + 1)] = 63/32.$ 

#### **Chained Matrix Multiplication**

- Problem: Given a series of n arrays (of appropriate sizes) to multiply:  $A_1 \times A_2 \times \cdots \times A_n$ .
- Determine where to place parentheses to minimize the number of multiplications.
- Matrix multiplication is associative:  $A_1(A_2A_3) = (A_1A_2)A_3$ , so all placements give same result.
- Formal problem:

Given a sequence of matrices  $A_1, A_2, \ldots, A_n$ , insert parentheses so that the product of the matrices needs the **minimal number of multiplications.** 

## Number of Multiplications / Parenthesizations

- Multiplying an i × j and a j × k matrix requires ijk multiplications: each element of the product requires j multiplications, and there are ik elements
- Given the matrices  $A_1, A_2, A_3, A_4$ , assume the dimensions of  $A_1 = d_0 \times d_1$ ,  $A_2 = d_1 \times d_2$  etc.
- Below are the five possible parenthesizations of these arrays, along with the number of multiplications:

 $(A_1A_2)(A_3A_4) : d_0d_1d_2 + d_2d_3d_4 + d_0d_2d_4$  $((A_1A_2)A_3)A_4 : d_0d_1d_2 + d_0d_2d_3 + d_0d_3d_4$  $(A_1(A_2A_3))A_4 : d_1d_2d_3 + d_0d_1d_3 + d_0d_3d_4$  $A_1((A_2A_3)A_4) : d_1d_2d_3 + d_1d_3d_4 + d_0d_1d_4$  $A_1(A_2(A_3A_4)) : d_2d_3d_4 + d_1d_2d_4 + d_0d_1d_4$ 

## Number of Multiplications / Parenthesizations

- The number of parenthesizations is at least  $T(n) \ge T(n-1) + T(n-1)$ :
  - Since the number with the first element removed is T(n-1), which is also the number with the last removed
  - Thus the number of parenthesizations is  $\Omega(2^n)$

- The number is actually 
$$T(n) = \sum_{k=1}^{n-1} T(k)T(n-k)$$

- This is because the original product can be split into 2 subproducts in k places. Each split is to be parenthesized optimally.

n = 1

#### **Characterizing the Optimal Parenthesization**

- An optimal parenthesization of  $A_1 \dots A_n$  must break the product into two expressions, each of which is parenthesized or is a single array
- Assume the break occurs at position  $\boldsymbol{k}$
- In the optimal solution, the solution to the product  $A_1 \dots A_k$  must be optimal:
  - Otherwise, we could improve  $A_1 \ldots A_n$  by improving  $A_1 \ldots A_k$
  - But the solution to  $A_1 \dots A_n$  is known to be optimal
  - Contradiction, thus the solution to  $A_1 \dots A_n$  is optimal
- Use Dynamic Programming: Consider a recursive solution, then improve it's performance with memoization or by rewriting bottom up.
- Dimensions of matrix  $A_i$  is  $d_{i-1} \times d_i$  $\rightsquigarrow A_i \times \cdots \times A_j$  is of size  $d_{i-1} \times d_j$

#### **Recursive Solution**

- $M_{ij}$  = number of multiplies in best way to parenthesize arrays  $A_i, \ldots A_j$
- M[i, i] = 0 since no product is required
- $\bullet$  The optimal solution of  $A_i \times A_j$  must break at some point, k, with  $i \leq k < j$
- $M[i,j] = M[i,k] + M[k+1,j] + d_{i-1}d_kd_j$

• 
$$M[i,j] = \begin{cases} 0 & \text{if } i = j \\ \min_{i \le k < j} \{ M[i,k] + M[k+1,j] + d_{i-1}d_kd_j \} & \text{if } i < j \end{cases}$$

• This is easily expressed as a recursive function (with exponential complexity)

#### **Efficient Computation**

- $\bullet$  We must find a way to calculate this bottom up. Which values does M[i,j] depend on?
- Consider a  $n \times n$  matrix of values M[i, j]: Diagonal is 0, build upper triangular table by diagonals
- Example: Array dimensions:

A\_1: 2 x 3 A\_2: 3 x 5 A\_3: 5 x 2 A\_4: 2 x 4 A\_5: 4 x 3

Array sizes: 
$$A_i = d_{i-1}$$
 by  $d_i$ :  $\frac{j \ 0 \ 1 \ 2 \ 3 \ 4 \ 5}{d_j \ 2 \ 3 \ 5 \ 2 \ 4 \ 3}$ 

#### **Example Showing Tables and Calculations**

$M_{ij}$	= r	umt	oer o	f mu	ltipl	ies in best way to parenthesize arrays $A_i, \ldots A_j$ :
i,j	1	2	3	4	5	
1	0	30	42	58	78	
2		0	30	54	72	
3			0	40	54	
4				0	<b>24</b>	
5					0	

Calculating  $M_{25}$  = number of multiplies in the optimal way to parenthesize  $A_2A_3A_4A_5$ :

$$M_{25} = \min \begin{cases} M_{22} + M_{35} + d_1 d_2 d_5, \\ M_{23} + M_{45} + d_1 d_3 d_5, \\ M_{24} + M_{55} + d_1 d_4 d_5 \end{cases}$$

#### **Example Showing Tables and Calculations**

$$= \min \begin{cases} 0 + 54 + 45 = 99\\ 30 + 24 + 18 = 72\\ 54 + 0 + 36 = 90 \end{cases}$$
$$= \min \begin{cases} (A_2)(A_3A_4A_5)\\ (A_2A_3)(A_4A_5)\\ (A_2A_3A_4)(A_5) \end{cases}$$

Optimal locations for parentheses:

20

# DP for Aligning Biological Sequences

Histone H1 (residues 120-180)



Thomas Shafee - Own work, CC BY 4.0, https://commons.wikimedia.org/w/index.php?curid=37188728

By

## **Mutations**

- **Mutation**: Heritable change in the DNA sequence. Occur due to exposure to **ultra violet radiation** or other **environmental conditions**.
- **Two levels** at which a mutation can take place:
  - **Point mutation:** within a single gene.
    - **substitution** (change of one nucleotide),
    - insertion (addition of nucleotides),
    - deletion.
  - Chromosomal mutation: whole segments interchange, either on the same chromosome, or on different ones.

#### **Point Mutations**

- May arise from spontaneous mutations during DNA replication.
- The rate of mutation increased by mutagens
   (physical or chemical agent that changes the genetic material).
- Mutagens: Physical (UV-, X-rays or heat), or chemical (molecules misplace base pairs / disrupt helical shape of DNA).



Wikipedia. By Jonsta247 - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=12481467

#### **Importance of Mutations**

 Mutations are responsible for inherited disorders & diseases.
 Sickle-cell anemia caused by missense point mutation in hemoglobin (in blood cells, responsible for oxygen transport.)
 Hydrophilic glutamic acid replaced with hydrophobic valine.
 ~> deformed red blood cells.

Sequence for Normal Hemoglobin: 6th codon: adenine (A)

AUG	GUG	CAC	CUG	ACU	CCU	G <mark>A</mark> G	GAG	AAG	UCU	GCC	GUU	ACU
START	Val	His	Leu	Thr	Pro	Glu	Glu	Lys	Ser	Ala	Val	Thr

Sickle Cell Hemoglobin: ~> thymine (DNA), uracil (RNA)

AUG	GUG	CAC	CUG	ACU	CCU	G <mark>U</mark> G	GAG	AAG	UCU	GCC	GUU	ACU
START	Val	His	Leu	Thr	Pro	Val	Glu	Lys	Ser	Ala	Val	Thr

• Mutations are the source of **phenotypic variation** 

 $\Rightarrow$  **new species** and **adaption** to environmental conditions.

#### **Sequence Comparison: Motivation**

Basic idea: similar sequences  $\rightsquigarrow$  similar proteins. Protein folding: 30 % sequence identity  $\Rightarrow$  structures similar.



Rout et al., Scientific Reports, vol 8, no 7002 (2018)

## **Comparing sequences**

**Theory:** during evolution **mutations** occurred, creating differences between families of contemporary species.

Missense mutation



U.S. National Library of Medicine

 $https://commons.wikimedia.org/w/index.php?curid{=}25399199$ 

#### **Comparing sequences**

Comparing two sequences: looking for **evidence** that they have **diverged from a common ancestor** by a **mutation process**.

Histone H1 (residues 120-180)



Thomas Shafee - Own work, CC BY 4.0, https://commons.wikimedia.org/w/index.php?curid=37188728

By

## **Sequence Alignment**

#### Informal definition:

Alignment of sequences  $x = x_1 \dots x_n$  and  $y = y_1 \dots y_m$ :

#### (i) insert spaces,

(ii) place resulting sequences **one above the other** so that every character or space has a counterpart.

**Example:** ACBCDDDB and CADBDAD. Possible alignments:

## **Optimal Alignment**

**Given:** two sequences x and y over alphabet  $\mathcal{A}$ .

 $\mathcal{A} = \{ \texttt{A},\texttt{G},\texttt{C},\texttt{T} \} (\mathsf{DNA}) \\ \mathcal{A} = \{ \texttt{A},\texttt{R},\texttt{N},\texttt{D},\texttt{C},\texttt{Q},\texttt{E},\texttt{G},\texttt{H},\texttt{I},\texttt{L},\texttt{K},\texttt{M},\texttt{F},\texttt{P},\texttt{S},\texttt{T},\texttt{W},\texttt{Y},\texttt{V} \} (\mathsf{proteins})$ 

Formalizing **optimality of an alignment**: define

- the costs for substituting a letter by another letter
   ⇒ substitution matrix;
- the costs for **insertion**  $\Rightarrow$  **gap penalties**.

## The Scoring Model

- Idea: assign a score to each alignment, choose best one.
- Additive scoring scheme: Total score = sum of all scores for pairs of letters + costs for gaps.
   Implicit assumption:

Mutations at different sites have occurred **independently**. (In most cases) reasonable for DNA and protein sequences.

- All common algorithms use additive scoring schemes.
- Modeling dependencies is possible, but at the price of significant computational complexities.

#### **Substitution Matrices**

#### • Expectation:

Identities in real alignments are more likely than by chance.

- Derive score for aligned pairs from a **probabilistic model**.
- Score: relative likelihood that two sequences are evolutionary related as opposed to being unrelated

 $\rightarrow$  score = ratio of probabilities.

- First assumption: Ungapped alignment, n = m.
- *R*: Random model:

Letter a occurs **independently** with some frequency  $q_a$ 

 $\Rightarrow$  Pr(two sequences) = product of probabilities for each letter:

$$P(x, y|R) = \prod_{i} q_{x_i} \prod_{i} q_{y_i}.$$

#### **Substitution Matrices**

• M (match): aligned pairs occur with joint probability

$$P(x, y|M) = \prod_{i} p_{x_i y_i}$$

• Ratio ~> "odds ratio":

$$\frac{P(x, y|M)}{P(x, y|R)} = \prod_{i} \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}}$$

• To arrive at an **additive** scoring system  $\rightarrow$  **log-odds ratio**:

$$S = \sum_{i} \log \left( \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} \right) = \sum_{i} s(x_i, y_i)$$

 s(a, b): log-likelihood ratio of pair (a, b) occurring as an aligned pair as opposed to an unaligned pair → substitution matrix.

#### **BLOSUM62** substitution matrix

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	lle	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val
Val	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
Tyr	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
Trp	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Thr	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Ser	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Pro	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
Phe	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Met	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Lys	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Leu	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
lle	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
His	-2	0	1	-1	-3	0	0	-2	8											
Gly	0	-2	0	-1	-3	-2	-2	6												
Glu	-1	0	0	2	-4	2	5													
Gln	-1	1	0	0	-3	5														
Cys	0	-3	-3	-3	9															
Asp	-2	-2	1	6																
Asn	-2	0	6																	
Arg	-1	5																		
Ala	4																			

Wikipedia

#### **Gap penalties**

**Gap penalty types** for a gap of length g:

- Linear:  $\gamma(g) = -gd$ , with d being the gap weight.
- Affine:  $\gamma(g) = -d (g 1)e$ , gap-open penalty d, gap-extension penalty e. Usually e < d.
- Convex: e.g.  $\gamma(g) = -d \log(g)$ . Each additional space contributes less than the previous space.



## **Global Alignment: Needleman-Wunsch algorithm**

#### The Global Alignment problem:

**INPUT**: two sequences  $x = x_1 \dots x_n$  and  $y = y_1 \dots y_m$ .

**TASK**: Find optimal alignment for linear gap penalties  $\gamma(g) = -gd$ .

Let F(i, j) be the optimal alignment score of the **prefix sequences**  $x_{1...i}$  and  $y_{1...j}$ . A zero index i = 0 or j = 0 refers to an **empty sequence**. F(i, j) has following properties:

Base conditions: 
$$F(i,0) = \sum_{k=1}^{i} -d = -id$$
  
 $F(0,j) = \sum_{k=1}^{j} -d = -jd, \quad F(0,0) = 0.$ 

Recurrence relation:

for 
$$1 \leq i \leq n, \ 1 \leq j \leq m$$
:

$$F(i,j) = \max \begin{cases} F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

## **Tabular Computation of Optimal Alignment**

Starting from F(0,0) = 0, fill the whole matrix  $(F)_{ij}$ :

for i = 0 or j = 0, calculate new value from left-hand (upper) value.

for  $i, j \ge 1$ , calculate the bottom righthand corner of each square of 4 cells from one of the 3 other cells:

keep a pointer in each cell back to the cell from which it was derived  $\Rightarrow$  **traceback pointer**.





#### **Global Alignment: Example**

x = HEAGAWGHEE, y = PAWHEAE. Linear gap costs d = 8. Scoring matrix: BLOSUM50



Durbin et al., Cambridge University Press

#### **Example: traceback procedure**



## **Time and Space Complexity**

**Theorem.** The time complexity of the Needleman-Wunsch algorithm is O(nm). Space complexity is O(m), if only F(x, y) is required, and O(nm) for the reconstruction of the alignment.

**Proof:** 

**Time:** when computing F(i, j), only cells (i - 1, j - 1), (i, j - 1), (i - 1, j) are examined  $\rightsquigarrow$  constant time. There are (n + 1)(m + 1) cells  $\rightsquigarrow O(nm)$  **time complexity.** 



**Space :** row-wise computation of the matrix: for computing row k, only row k - 1 must be stored  $\rightsquigarrow O(m)$  **space. Reconstructing** the alignment: all traceback pointers must be stored  $\rightsquigarrow O(nm)$  **space complexity.** 

## **Local Alignments**

#### The Local Alignment problem:

**INPUT**: two sequences  $x = x_1, \ldots, x_n$  and  $y = y_1, \ldots, y_m$ . **TASK**: find subsequences a of x and b of y, whose similarity (=optimal global alignment score) is maximal (over all such pairs of subsequences). Assume linear gap penalties  $\gamma(g) = -gd$ .

#### **Subsequence** = **contiguous** segment of a sequence.

Consider first a simpler problem by **fixing the endpoint** of the subsequences at index pair (i, j):

**Local suffix alignment problem:** given x, y, i, j, find suffixes  $\alpha$  of  $x_{1,...,i}$ and  $\beta$  of  $y = y_{1,...,j}$  such that their global alignment score is maximal.

$$(x_1, \ldots, \underbrace{x_k, \ldots, x_i}_{\alpha}), \quad (y_1, \ldots, \underbrace{y_l, \ldots, y_j}_{\beta})$$

## Local suffix alignments

Consider global alignment path to cell (i, j). Where to start? Intuition: Indices (k, l) found by following the path back to (0, 0), but stopping at the first negative value.



**Remark**: If we consider all solutions (i.e. for all (i, j) pairs), we look at all possible subsequences (no restrictions on  $\alpha, \beta$ )

Maximal solution of local suffix alignment over all pairs (i, j)= solution of local alignment problem.

#### **Smith-Waterman Algorithm**

F(i, j): optimal local suffix alignment for indices i, j.

**Global alignment** with one **modification**: Prefixes whose scores are  $\leq 0$  are **discarded**  $\rightsquigarrow$  alignment can **start anywhere**.

Recurrence relation: 
$$F(i,j) = \max \begin{cases} 0\\F(i-1,j-1) + s(x_i,y_j)\\F(i-1,j) - d\\F(i,j-1) - d \end{cases}$$

Finally, find indices  $i^*$  and  $j^*$  after which the similarity only decreases. Stop the alignment there.

$$F(i^*, j^*) = \max_{i,j} F(i, j)$$

#### Traceback...

...starts at highest value until a cell with 0 is reached.



Adapted from Durbin et al., Cambridge University Press. https://doi.org/10.1017/CBO9780511790492.004

#### Local vs. Global Alignment: Biological Considerations

- Many proteins have **multiple domains**, or modules.
- Some domains are present (with high similarity) in many other proteins
- **Local** alignment can detect similar regions in otherwise dissimilar proteins.



Durbin et al., Cambridge University Press. https://doi.org/10.1017/CBO9780511790492.004

#### Other gap models

• So far: linear gap model. Not ideal for biological sequences, since it penalizes additional gap steps as much as the first. But in reality: When gaps do occur, they are often longer than one character.

HBA\_HUMAN GSAOVKGHGKKVADALTNAVAHV---D--DMPNALSALSDLHAHKL ++ ++++H+ KV + +A ++ +L+ L+++H+ K LGB2\_LUPLU NNPELOAHAGKVFKLVYEAAIOLOVTGVVVTDATLKNLGSVHVSKG

Durbin et al., Cambridge University Press. https://doi.org/10.1017/CBO9780511790492.004

• For a general gap cost function  $\gamma(g)$ , we can still use the standard dynamic programming recursion with slight modifications:

$$F(i,j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(k, j) + \gamma(i-k), & k = 0, \dots, i-1, \\ F(i, k) + \gamma(j-k), & k = 0, \dots, j-1. \end{cases}$$

• **Problem:** requires  $O(n^3)$  operations to align two sequences of length n, rather than  $O(n^2)$ . Why?

#### Alignment with affine gap costs

For affine gap costs,  $\gamma(g) = -d - (g - 1)e$ , there exists a **solution**: Modify recurrence by introducing another two "states". Denote by

- M(i,j) the best score given that  $x_i$  is aligned to  $y_j$ ,
- $I_x(i,j)$  the best score given that  $x_i$  is aligned to a gap,
- $I_y(i,j)$  the best score given that  $y_j$  is aligned to a gap.

$$M(i,j) = \max \begin{cases} M(i-1,j-1) + s(x_i, y_j) \\ I_x(i-1,j-1) + s(x_i, y_j) \\ I_y(i-1,j-1) + s(x_i, y_j) \end{cases} \overset{s(x_i, y_j)}{\underset{(+1,+1)}{}} \overset{s(x_i, y_j)}{\underset{(+1,+1)}{}} \overset{-e}{\underset{(x_i, y_j)}{}} \\ I_x(i,j) = \max \begin{cases} M(i-1,j) - d \\ I_x(i-1,j) - e \\ I_y(i,j-1) - e \end{cases} \overset{s(x_i, y_j)}{\underset{(+1,+1)}{}} \overset{-e}{\underset{(x_i, y_j)}{}} \overset{f_x(x_i, y_j)}{\underset{(+1,+1)}{}} \overset{-e}{\underset{(x_i, y_j)}{}} \end{aligned}$$

#### **Example FSA alignment**



Durbin et al., Cambridge University Press. https://doi.org/10.1017/CBO9780511790492.004

#### FSA alignment corresponds to path through states.

Probabilistic version ~> Hidden Markov models

## **Global Alignment in Linear Space**

- Problem: genomic scale sequence analysis: comparing two large genomic sequences:  $m, n \approx 10^6 \Rightarrow$  space complexity  $10^{12}$  is clearly unacceptable!
- Solution: linear space algorithms with space complexity O(m+n).
- Basic idea: divide and conquer. Let  $u = \lfloor \frac{n}{2} \rfloor$  be the integer part of  $\frac{n}{2}$ .
  - Let v be a row index such that the cell (u, v) is on the optimal alignment.
  - Split dynamic programming problem into two parts:  $(0,0) \rightarrow (u,v)$  and  $(u,v) \rightarrow (n,m)$ .

Optimal alignment will be concatenation of individual sub-alignments.

– Repeat splitting until until u = 0: trivial



## **Global Alignment in Linear Space**

• For  $i \ge u$  define c(i, j) such that (u, c(i, j)) is on the optimal path from  $(1, 1) \rightarrow (i, j)$ .



• Let (i', j') be the preceding cell to (i, j) from which F(i, j) is derived. Update c(i, j) as:

$$c(i,j) = \begin{cases} j & , \text{ if } i = u, \\ c(i',j') & , \text{else} \end{cases}$$



- Local operation  $\rightsquigarrow$  need to store only the previous row of c().
- Finally, v = c(n, m).

#### **Global Alignment in Linear Space: Example**

Computing the c matrix for the first step (i = n = 6, j = m = 4, u = 3). The c values are written as subscripts. BLOSUM62, linear gap costs d = 8.

		0		1		2		3		4		5		6
		•		Н		Е		А		G		А		W
0	٠	0	$\leftarrow$	-8	$\leftarrow$	-16	$\leftarrow$	$-24_{0}$	$\leftarrow$	$-32_{0}$	$\leftarrow$	$-40_{0}$	$\leftarrow$	$-48_{0}$
		$\uparrow$	ĸ		$\overline{\mathbf{x}}$		ĸ				$\overline{\mathbf{x}}$			
1	Р	-8		-2		-9		-17 <mark>1</mark>	$\leftarrow$	$-25_{1}$		$-33_{0}$	$\leftarrow$	$-41_{0}$
		$\uparrow$	ĸ	$\uparrow$	$\overline{\mathbf{x}}$		ĸ				ĸ			
2	А	-16		-10		-3		$-4_{2}$	$\leftarrow$	$-12_{2}$		$-20_{1}$	$\leftarrow$	$-28_{1}$
		$\uparrow$		$\uparrow$			ĸ		×,		$\overline{\mathbf{x}}$		ĸ	
3	W	-24		-18		-11		$-6_{3}$		$-7_{2}$		-152		$-5_{1}$
		$\uparrow$	ĸ		$\overline{\mathbf{x}}$		ĸ		×,		$\overline{\mathbf{x}}$			$\uparrow$
4	Н	-32		-14		-18		$-13_{4}$		$-8_{3}$		$-9_{2}$		-13 <sub>1</sub>

Every c(i, j) defines a row index v such that (u, c(i, j)) is on the optimal path from (1, 1) to  $(i, j) \rightsquigarrow v = c(6, 4) = 1$ , so (3,1) is our desired element on the optimal path form (1,1) to (6,4).