



# ALGORITHMS OF BIOINFORMATICS

## 3

## Comparing Sequences

*13 November 2025*

Prof. Dr. Sebastian Wild



# Outline

## 3 Comparing Sequences

- 3.1 Sequence Alignment
- 3.2 Dynamic Programming
- 3.3 Global – Local – Semilocal
- 3.4 General Scores & Affine Gap Costs
- 3.5 Bounded-Distance Alignments
- 3.6 Exhaustive Tabulation
- 3.7 Linear-Space Alignments
- 3.8 Multiple Sequence Alignment



## 3.1 Sequence Alignment



# Sequence Similarity

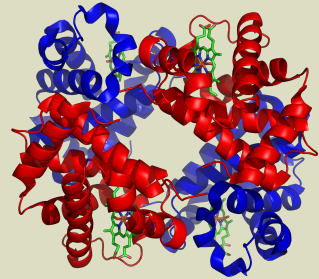
**Example:** two proteins from *human hemoglobin*

Human Hemoglobin  $\alpha$  globin subunit <https://www.uniprot.org/uniprotkb/P69905>

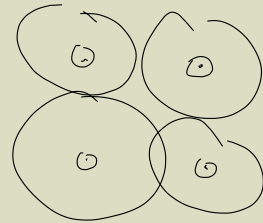
Human Hemoglobin  $\beta$  globin subunit <https://www.uniprot.org/uniprotkb/P68871>

~> *essentially symmetric copies with same function*

3D Structure of hemoglobin



[https://commons.wikimedia.org/wiki/File:1GZX\\_Haemoglobin.png](https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png)





# Sequence Similarity

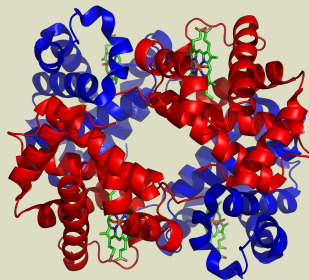
**Example:** two proteins from *human hemoglobin*

Human Hemoglobin  $\alpha$  globin subunit <https://www.uniprot.org/uniprotkb/P69905>

Human Hemoglobin  $\beta$  globin subunit <https://www.uniprot.org/uniprotkb/P68871>

~> *essentially symmetric copies with same function*

3D Structure of hemoglobin



[https://commons.wikimedia.org/wiki/File:1GZX\\_Haemoglobin.png](https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png)

Sequences of the subunits (142 resp. 147 amino acids):

MVLSPADKTNVKAANGKVGAHAGEYGAEALERMF<sup>1</sup>LSFPTTKTYFPHFDLSHGSAQVKGHGKKVADALTNAVAHVDDMPNALS<sup>2</sup>LSDLHAHKL<sup>3</sup>RVDPVNF<sup>4</sup>KLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR

MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVL<sup>1</sup>GAFSDGLAHLDNLKGT<sup>2</sup>FATLSELHCDKLHVDPENFRLLGNVLCVLAHFGKEFTPPVQAA<sup>3</sup>YQKVAVG<sup>4</sup>ANALAHKYH



# Sequence Similarity

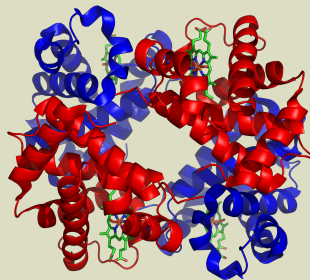
**Example:** two proteins from *human hemoglobin*

Human Hemoglobin  $\alpha$  globin subunit <https://www.uniprot.org/uniprotkb/P69905>

Human Hemoglobin  $\beta$  globin subunit <https://www.uniprot.org/uniprotkb/P68871>

~> *essentially symmetric copies with same function*

3D Structure of hemoglobin



[https://commons.wikimedia.org/wiki/File:1GZX\\_Haemoglobin.png](https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png)

Sequences of the subunits (142 resp. 147 amino acids):

MVLSPADKTNVKAANGKVGAHAGEYGAEALERMFLSPTTKTYFPHFDLSHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR

MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVLGAFSDGLAHLNLIKGTATLSELHCDKLHVDPENFRLLGNVLCVLAHFGKEFTPPVQAAQYQKVAVGANALAHKYH

*These are supposed to be “similar”!?*



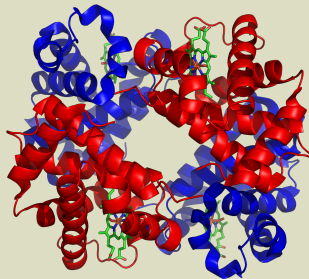
# Sequence Similarity

**Example:** two proteins from *human hemoglobin*

Human Hemoglobin  $\alpha$  globin subunit <https://www.uniprot.org/uniprotkb/P69905>Human Hemoglobin  $\beta$  globin subunit <https://www.uniprot.org/uniprotkb/P68871>

↪ essentially symmetric copies with same function

### 3D Structure of hemoglobin



[https://commons.wikimedia.org/wiki/File:1GZX\\_Haemoglobin.png](https://commons.wikimedia.org/wiki/File:1GZX_Haemoglobin.png)

Sequences of the subunits (142 resp. 147 amino acids):

MVLSPADKTNVKAAGWKVGAHAGEYGAEALERMFISFPTTKTYFPHFDLSHGSAOVKGGKVKDALTNAVAHVDDMPNALSALSDLHAHKL RVD PNVFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR

MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVVPWTORFFESFGDLSTPDVAMGNPKVKAHGKKVLGAFSDGLAHLNLLKGTFTLSELHCDKLHVDPENFRLLGNVLCVLAHFGKEFTPPVOAAYOKVVAGVANALAHKYH

*These are supposed to be “similar”!?*

*Alignment* by EMBOSS Needle <https://www.ebi.ac.uk/jdispatcher/psa>

MV - LSPADKTNVKAAMGKVGAHAGEYGAELERMFLSPTTKTYFPHF-DLS-----HGSAQVKGHGKKVADALTNAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR  
| | : | : | : | |||| : | ||| : : : | : : | ||| : : || |||| | :: : ||:: : ||:| | | || |:| : | | | |||| | : | : | : | |  
MVHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVLGAFSDGLAHLNDLKGTFATSELCHKDCLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYOKVAVAGVANALAHKYH

| = same amino acid (65x); : = similar amino acids (25x)

→ 60% same



# String Distances

*Mutations mean much in bioinformatics needs fuzzy comparisons . . .*

*How can we formally define these?*

- ▶ This unit studies wide class of options
- ▶ Algorithmically, all are similar to deal with
- ▶ Unfortunately, general case again hard . . .



# String Distances

*Mutations mean much in bioinformatics needs fuzzy comparisons . . .*

*How can we formally define these?*

- ▶ This unit studies wide class of options
- ▶ Algorithmically, all are similar to deal with
- ▶ Unfortunately, general case again hard . . .
- ▶ Simplest string distance function: **Hamming distance**  $d_H = \text{\#mismatches}$ 
  - ⚡ only defined for strings of same length



# String Distances

*Mutations mean much in bioinformatics needs fuzzy comparisons . . .*

*How can we formally define these?*

- ▶ This unit studies wide class of options
- ▶ Algorithmically, all are similar to deal with
- ▶ Unfortunately, general case again hard . . .
- ▶ Simplest string distance function: **Hamming distance**  $d_H = \text{\#mismatches}$ 
  - ⚡ only defined for strings of same length
  - ▶ How about strings like this:

$A = \text{alongsharedstring}$   
 $B = \text{longsharedstrings}$   $\rightsquigarrow d_H(A, B) = |A| = 17$       *These are maximally different!?*



# String Distances

*Mutations mean much in bioinformatics needs fuzzy comparisons . . .*

*How can we formally define these?*

- ▶ This unit studies wide class of options
- ▶ Algorithmically, all are similar to deal with
- ▶ Unfortunately, general case again hard . . .
- ▶ Simplest string distance function: **Hamming distance**  $d_H = \text{\#mismatches}$ 
  - ⚡ only defined for strings of same length
  - ▶ How about strings like this:

$A = \text{alongsharedstring}$   
 $B = \text{longsharedstrings}$   $\rightsquigarrow d_H(A, B) = |A| = 17$       *These are maximally different!?*

$\rightsquigarrow$  *Need a more flexible notion . . .*



# Edit Distance

Natural idea for distances: describe **how** to get from  $A$  to  $B$

$\rightsquigarrow$  *relative compression!*

$A[0..17) = \cancel{a}longsharedstrings$

$B[0..17) = longsharedstrings$

“Edit script”:

0. Start with  $\overset{A}{\cancel{M}}$ .
1. Delete  $\overset{A}{\cancel{M}}[0]$
2. Insert  $s$  at end of  $\overset{A}{\cancel{M}}$ .

$\rightsquigarrow$  2 character operations needed  $\rightsquigarrow d_{\text{edit}}(A, B) = 2$



# Edit Distance

Natural idea for distances: describe **how** to get from  $A$  to  $B$   $\rightsquigarrow$  *relative compression!*

$A[0..17) = \text{alongsharedstring}$

$B[0..17) = \text{longsharedstrings}$

“Edit script”:

0. Start with  $S_1$ .

1. Delete  $S_1[0]$

2. Insert  $s$  at end of  $S_1$ .

$\rightsquigarrow$  2 character operations needed  $\rightsquigarrow d_{\text{edit}}(A, B) = 2$

## Edit Distance Problem

► **Given:** String  $\underbrace{A[0..m)}$  and  $\underbrace{B[0..n)}$  over alphabet  $\Sigma = [0..\sigma)$ .

► **Goal:**  $d_{\text{edit}}(A, B) =$  minimal #symbol operations to transform  $A$  into  $B$   
operations can be insertion/deletion/substitution of single character  
+ optimal edit script (with this number of operations)



# Edit Distance Example

**Example:** edit distance  $d_{\text{edit}}(A, B)$  with  $A = \text{algorithm}$ ,  $B = \text{logarithm}$ ?



## Edit Distance Example

**Example:** edit distance  $d_{\text{edit}}(A, B)$  with  $A = \text{algorithm}$ ,  $B = \text{logarithm}$ ?

012345678

algorithm

logarithm



# Edit Distance Example

**Example:** edit distance  $d_{\text{edit}}(A, B)$  with  $A = \text{algorithm}$ ,  $B = \text{logarithm}$ ?

0	1	2	3	4	5	6	7	8	
	a	l	g	o	r	i	t	h	m
	l	o	g	a	r	i	t	h	m

Edit script:

1. Delete  $A[0]$
2. Insert o after  $A[1] = \text{l}$
3. Replace  $A[3] = \text{o}$  by a



# Edit Distance Example

**Example:** edit distance  $d_{\text{edit}}(A, B)$  with  $A = \text{algorithm}$ ,  $B = \text{logarithm}$ ?

0	1	2	3	4	5	6	7	8
a	l	g	o	r	i	t	h	m
l	o	g	a	r	i	t	h	m

Edit script:

1. Delete  $A[0]$
2. Insert o after  $A[1] = \text{l}$
3. Replace  $A[3] = \text{o}$  by a

Compact representation of edit script: *String alignment*

0	1	2	3	4	5	6	7	8	9
a	l	-	g	o	r	i	t	h	m
-		+		x					
-	l	o	g	a	r	i	t	h	m



## Edit Distance Example

**Example:** edit distance  $d_{\text{edit}}(A, B)$  with  $A = \text{algorithm}$ ,  $B = \text{logarithm}$ ?

012345678  
algorithm  
logarithm

Edit script:

1. Delete  $A[0]$
2. Insert  $o$  after  $A[1] = u$
3. Replace  $A[3] = o$  by  $a$

### Compact representation of edit script: *String alignment*

```
0123456789
a l - g o r i t h m
- | + | x | | | |
- l o g a r i t h m
```

Formally: string over pairs of letters or *gap symbols*

$$\left\{ \begin{bmatrix} c \\ c \end{bmatrix} : c \in \Sigma \right\} \cup \left\{ \begin{bmatrix} c \\ - \end{bmatrix}, \begin{bmatrix} - \\ c \end{bmatrix} : c \in \Sigma \right\} \cup \left\{ \begin{bmatrix} c \\ c' \end{bmatrix} : c, c' \in \Sigma, c \neq c' \right\}$$

↪ Edit distance = #  $\begin{bmatrix} c \\ - \end{bmatrix}, \begin{bmatrix} - \\ c \end{bmatrix}, \begin{bmatrix} c \\ c' \end{bmatrix}$  with  $c \neq c'$



# Edit Distance and Longest Common Subsequence

- Note: close relation to *longest common subsequence*

Optimal edit script  $\approx$  maximal number of matches = longest common subsequence



# Edit Distance and Longest Common Subsequence

- Note: close relation to *longest common subsequence*  
Optimal edit script  $\approx$  maximal number of matches = longest common subsequence

- But: Optimal alignment may not contain any longest common subsequence

```

axxa  axxa  axxa
|   |   |   |   |
a  ayya  ayya  ayy

```

```

~ axxaaxxaaxxa
| 4 4 | 4 5 | 5 5      8
aayyaayyaayy ~

```

- LCS and edit distance are equivalent if we only allow insert and delete operations



## 3.2 Dynamic Programming



## Recap: The 6 Steps of Dynamic Programming

1. Define **subproblems** (and relate to original problem)
2. **Guess** (part of solution)  $\rightsquigarrow$  local brute force
3. Set up **DP recurrence** (for quality of solution)
4. Recursive implementation with **Memoization**
5. Bottom-up **table filling** (topological sort of subproblem dependency graph)
6. **Backtracing** to reconstruct optimal solution




# Recap: The 6 Steps of Dynamic Programming

↪ see *Efficient Algorithms*

1. Define **subproblems** (and relate to original problem)
2. **Guess** (part of solution) ↪ local brute force
3. Set up **DP recurrence** (for quality of solution)
4. Recursive implementation with **Memoization**
5. Bottom-up **table filling** (topological sort of subproblem dependency graph)
6. **Backtracing** to reconstruct optimal solution

► Steps 1–3 require insight / creativity / intuition;  
Steps 4–6 are mostly automatic / same each time

↪ Correctness proof usually at level of DP recurrence

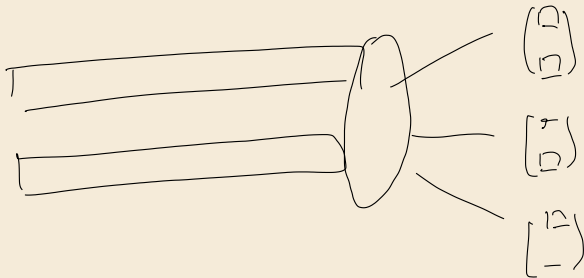
 running time too! worst case time = #subproblems · time to find single best guess



# Edit Distance by DP

original  $i = m, j = n$

1. **Subproblems:**  $(i, j)$  for  $0 \leq i \leq m, 0 \leq j \leq m$  compute  $d_{\text{edit}}(A[0..i], B[0..j])$
2. **Guess:** What to do with last positions? (insert/delete/(mis)match)






# Edit Distance by DP

- Subproblems:**  $(i, j)$  for  $0 \leq i \leq m, 0 \leq j \leq m$  compute  $d_{\text{edit}}(A[0..i], B[0..j])$
- Guess:** What to do with last positions? (insert/delete/(mis)match)
- Recurrence:**  $D(i, j) = d_{\text{edit}}(A[0..i], B[0..j])$

$$D(i, j) = \begin{cases} i & \text{// } i \text{ deletions} & \text{if } j = 0 \\ j & \text{// } j \text{ insertions} & \text{if } i = 0 \\ \min \begin{cases} D(i-1, j) + 1, & \text{delete} \\ D(i, j-1) + 1, & \text{insertion} \\ D(i-1, j-1) + [A[i-1] \neq B[j-1]] & \text{otherwise} \end{cases} & \end{cases}$$



// version bracket  $\begin{cases} 1 & \text{could be} \\ 0 & \text{else} \end{cases}$



# Edit Distance by DP

1. **Subproblems:**  $(i, j)$  for  $0 \leq i \leq m, 0 \leq j \leq n$  compute  $d_{\text{edit}}(A[0..i], B[0..j])$
2. **Guess:** What to do with last positions? (insert/delete/(mis)match)
3. **Recurrence:**  $D(i, j) = d_{\text{edit}}(A[0..i], B[0..j])$

$$D(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \begin{cases} D(i-1, j) + 1, \\ D(i, j-1) + 1, \\ D(i-1, j-1) + [A[i-1] \neq B[j-1]] \end{cases} & \text{otherwise} \end{cases}$$

$\rightsquigarrow O(nm)$  subproblems

►  $O(1)$  time to check all guesses (per subproblem)

$\rightsquigarrow O(nm)$  overall time and space

► An optimal *edit script* can be constructed by a *backtrace* (see below)



## Edit Distance – Step 4: Memoization

- ▶ Write **recursive** function to compute recurrence
- ▶ But *memoize* all results! (symbol table: subproblem  $\mapsto$  optimal cost )

$\rightsquigarrow$  First action of function: check if subproblem known

- ▶ If so, return cached optimal cost
- ▶ Otherwise, compute optimal cost and remember it!

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace



## Edit Distance – Step 4: Memoization

- ▶ Write **recursive** function to compute recurrence
- ▶ But *memoize* all results! (symbol table: subproblem  $\mapsto$  optimal cost)

~> First action of function: check if subproblem known

- ▶ If so, return cached optimal cost
- ▶ Otherwise, compute optimal cost and remember it!

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace

---

```
1 procedure editDist(i, j):
2   if i == 0
3     return j
4   else if j == 0
5     return i
6   end if
7   best := +∞
8   Di := cachedED(i, j - 1) + 1
9   Dd := cachedED(i - 1, j) + 1
10  Dm := cachedED(i - 1, j - 1) + [A[i] ≠ B[j]]
11  best := min{Dd, Di, Dm}
12  return best
```

---

$$D(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \begin{cases} D(i, j - 1) + 1, \\ D(i - 1, j) + 1, \\ D(i - 1, j - 1) + [A[i - 1] \neq B[j - 1]] \end{cases} & \text{otherwise} \end{cases}$$

---

```
13 procedure cachedED(r[i..j], c[i..j]):
14   // D[0..m][0..n] initialized to NULL at start
15   if D[i][j] == NULL
16     D[i][j] := editDist(i, j)
17   return D[i][j]
```

---



## Edit Distance – Step 5: Table Filling

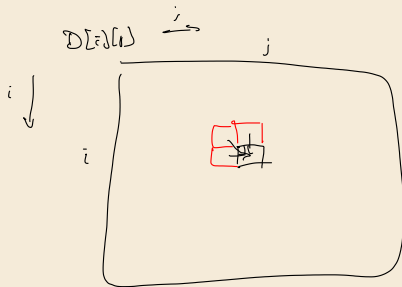
- ▶ Recurrence induces a DAG on subproblems (who calls whom)
  - ▶ Memoized recurrence traverses this DAG (DFS!)
  - ▶ We can slightly improve performance by systematically computing subproblems following a fixed topological order

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace



## Edit Distance – Step 5: Table Filling

- ▶ Recurrence induces a DAG on subproblems (who calls whom)
  - ▶ Memoized recurrence traverses this DAG (DFS!)
  - ▶ We can slightly improve performance by systematically computing subproblems following a fixed topological order
- ▶ **Topological order** here: lexicographic by  $(i, j)$



1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace



## Edit Distance – Step 5: Table Filling

- ▶ Recurrence induces a DAG on subproblems (who calls whom)
  - ▶ Memoized recurrence traverses this DAG (DFS!)
  - ▶ We can slightly improve performance by systematically computing subproblems following a fixed topological order
- ▶ **Topological order** here: lexicographic by  $(i, j)$

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace

---

```
1 procedure editDist(A[0..m], B[0..n]):  
2   D[0..m][0..n] := new array  
3   for i = 0, 1, ..., m // iterate over subproblems ...  
4     for j = 0, 1, ..., n // ... in topological order  
5       if i == 0  
6         D[i][j] := j  
7       else if j == 0  
8         D[i][j] := i  
9       else  
10        D[i][j] := min { D[i][j - 1] + 1,  
                          D[i - 1][j] + 1,  
                          D[i - 1][j - 1] + [A[i - 1] ≠ B[j - 1]]  
11      return D[m][n]
```

---



## Edit Distance – Step 5: Table Filling

- ▶ Recurrence induces a DAG on subproblems (who calls whom)
  - ▶ Memoized recurrence traverses this DAG (DFS!)
  - ▶ We can slightly improve performance by systematically computing subproblems following a fixed topological order
- ▶ **Topological order** here: lexicographic by  $(i, j)$

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace

```
1 procedure editDist(A[0..m], B[0..n]):  
2   D[0..m][0..n] := new array  
3   for i = 0, 1, ..., m // iterate over subproblems ...  
4     for j = 0, 1, ..., n // ... in topological order  
5       if i == 0  
6         D[i][j] := j  
7       else if j == 0  
8         D[i][j] := i  
9       else  
10        D[i][j] := min { D[i][j-1] + 1,  
                          D[i-1][j] + 1,  
                          D[i-1][j-1] + [A[i-1] ≠ B[j-1]]  
11      return D[m][n]
```

- ▶ Same  $\Theta$ -class as memoized recursive function
- ▶ In practice usually substantially faster
  - ▶ lower overhead
  - ▶ predictable memory accesses



## Edit Distance – Step 6: Backtracing

- ▶ So far, only determine the **cost** of an optimal solution
  - ▶ But we also want the solution itself
- ▶ By *retracing* our steps, we can construct optimal edit script

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace



## Edit Distance – Step 6: Backtracing

- ▶ So far, only determine the **cost** of an optimal solution
  - ▶ But we also want the solution itself
- ▶ By *retracing* our steps, we can construct optimal edit script

---

```
1 procedure editScript(A[0..m], B[0..n]):
2   D[0..m][0..n] := editDist(A[0..m], B[0..n])
3   return traceback(m, n)
4
5 procedure traceback(i, j):
6   if i == 0
7     return Insert(B[0]), ..., Insert(B[j - 1])
8   else if j == 0
9     return Delete(A[0]), ..., Delete(A[i - 1])
10  else if D[i][j] == D[i][j - 1] + 1
11    return traceback(i, j - 1), Insert(B[j - 1])
12  else if D[i][j] == D[i - 1][j] + 1
13    return traceback(i - 1, j), Delete(B[i - 1])
14  else if A[i - 1] == B[j - 1]
15    return traceback(i - 1, j - 1)
16  else return traceback(i - 1, j - 1), Replace(A[i - 1] → B[j - 1])
```

---

1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace

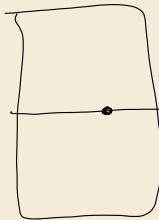
- ▶ follow recurrence a second time



## Edit Distance – Step 6: Backtracing

- ▶ So far, only determine the **cost** of an optimal solution
  - ▶ But we also want the solution itself
- ▶ By *retracing* our steps, we can construct optimal edit script

```
1 procedure editScript(A[0..m], B[0..n]):  
2   D[0..m][0..n] := editDist(A[0..m], B[0..n])  
3   return traceback(m, n)  
4  
5 procedure traceback(i, j):  
6   if i == 0  
7     return Insert(B[0]), ..., Insert(B[j - 1])  
8   else if j == 0  
9     return Delete(A[0]), ..., Delete(A[i - 1])  
10  else if D[i][j] == D[i][j - 1] + 1  
11    return traceback(i, j - 1), Insert(B[j - 1])  
12  else if D[i][j] == D[i - 1][j] + 1  
13    return traceback(i - 1, j), Delete(B[i - 1])  
14  else if A[i - 1] == B[j - 1]  
15    return traceback(i - 1, j - 1)  
16  else return traceback(i - 1, j - 1), Replace(A[i - 1] → B[j - 1])
```



1. Subproblems
2. Guess!
3. DP Recurrence
4. Memoization
5. Table Filling
6. Backtrace

- ▶ follow recurrence a second time
  - ▶ always have for running time:  
backtracing =  $O(\text{computing } M)$
- ~> computing optimal cost and  
computing optimal solution have  
same complexity



### **3.3 Global – Local – Semilocal**



# Local Alignment

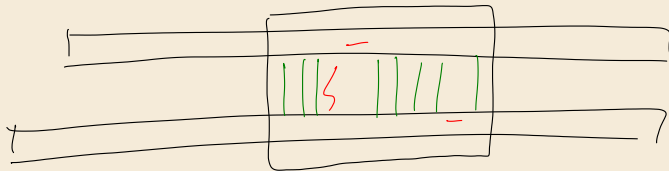
*So far, we assumed that we know similar regions.*

*How to detect significantly similar regions hidden in larger strings?*

↪ Allow new edit script operations (all cost 0):

- ▶  $\text{IgnorePrefix}(A[0..i])$       free deletes at beginning
- ▶  $\text{IgnorePrefix}(B[0..j])$       free inserts at beginning
- ▶  $\text{IgnoreSuffix}(A[i..m])$       free deletes at end
- ▶  $\text{IgnoreSuffix}(B[j..n])$       free inserts at end

↪ *Local Alignment*





# Local Alignment

*So far, we assumed that we know similar regions.*

*How to detect significantly similar regions hidden in larger strings?*

↪ Allow new edit script operations (all cost 0):

- ▶ IgnorePrefix( $A[0..i)$ )      free deletes at beginning
- ▶ IgnorePrefix( $B[0..j)$ )      free inserts at beginning
- ▶ IgnoreSuffix( $A[i..m)$ )      free deletes at end
- ▶ IgnoreSuffix( $B[j..n)$ )      free inserts at end

↪ *Local Alignment*

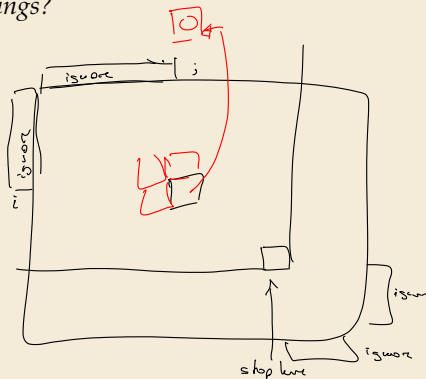
▶ Easy to incorporate in DP recurrence:

0. switch to **maximizing score** (instead min difference), otherwise empty substring is best

↪ Matches contribute +1 reward, rest penalty (negative score)

1. Always allow 4th option: **start** a **new** local alignment from here (at score 0)

2. Allow to finish at any  $D[i][j]$  ↪ free suffix





# Local Alignment Recurrence

$$D(i, j) = \begin{cases} 0 & \text{if } j = 0 \\ 0 & \text{if } i = 0 \\ \min \begin{cases} 0, \\ D(i-1, j) - 1, \\ D(i, j-1) - 1, \\ D(i-1, j-1) + [A[i-1] = B[j-1]] \\ \quad - [A[i-1] \neq B[j-1]] \end{cases} & \text{otherwise} \end{cases}$$

Optimal local alignment score:  $\max_{i \in [0..m], j \in [0..n]} D[i][j]$



## Semilocal Alignment a.k.a. Fitting Alignment

*Slight twist: We know conserved region, but need to find best match in larger sequence.*

*What substring of  $B[0..n]$  is the best match for  $A[0..m]$ ? (typically then  $m \ll n$ )*



# Semilocal Alignment a.k.a. Fitting Alignment

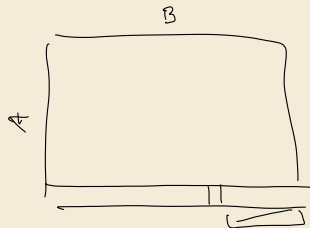
*Slight twist: We know conserved region, but need to find best match in larger sequence.*

*What substring of  $B[0..n]$  is the best match for  $A[0..m]$ ? (typically then  $m \ll n$ )*

$\rightsquigarrow$  only allow IgnorePrefix( $B[0..j)$ ) and IgnoreSuffix( $B[j..n)$ )

$$\rightsquigarrow D(i, j) = \begin{cases} -i & \text{if } j = 0 \\ \mathbf{0} & \text{if } i = 0 \\ \min \begin{cases} D(i-1, j) - \mathbf{1}, \\ D(i, j-1) - \mathbf{1}, \\ D(i-1, j-1) + [A[i-1] = B[j-1]] \end{cases} & \text{otherwise} \end{cases}$$

Optimal local alignment score:  $\max_{j \in [0..n]} D[m][j]$





## 3.4 General Scores & Affine Gap Costs



## General Scores

DP algorithm remains unchanged if we let contribution of (mis)match  $A[i - 1]$  vs  $B[j - 1]$  depend on used letters.

- ▶ For example, replacing amino acid with chemically similar one might not affect function  
     $\rightsquigarrow$  contributes small positive score
- ▶ replacing amino acid with dissimilar one  $\rightsquigarrow$  negative score



# General Scores

DP algorithm remains unchanged if we let contribution of (mis)match  $A[i - 1]$  vs  $B[j - 1]$  depend on used letters.

- ▶ For example, replacing amino acid with chemically similar one might not affect function  
     $\rightsquigarrow$  contributes small positive score
- ▶ replacing amino acid with dissimilar one  $\rightsquigarrow$  negative score

Formally, any function giving additive scores for columns  $S : (\Sigma \cup \{-\})^2 \setminus \{[-]\} \rightarrow \mathbb{R}$  works.

## General Alignment Score $S$ :

- ▶ symmetric matches/substitutions matrix  $p : \Sigma \times \Sigma \rightarrow \mathbb{R}$        $(p(a, b) = p(b, a))$
- ▶ gap penalty  $g \in \mathbb{R}$   
     $\rightsquigarrow S\left(\begin{bmatrix} c \\ c' \end{bmatrix}\right) = p(a, b), S\left(\begin{bmatrix} c \\ - \end{bmatrix}\right) = S\left(\begin{bmatrix} - \\ c \end{bmatrix}\right) = g$   
     $\rightsquigarrow$  score of alignment sum of scores of columns



# BLOSUM Matrices

	C	S	T	A	G	P	D	E	Q	N	H	R	K	M	I	L	V	W	Y	F	
C	9																				C
S	-1	4																			S
T	-1	1	5																		T
A	0	1	0	4																	A
G	-3	0	-2	0	6																G
P	-3	-1	-1	-1	-2	7															P
D	-3	0	-1	-2	-1	-1	6														D
E	-4	0	-1	-1	-2	-1	2	5													E
Q	-3	0	-1	-1	-2	-1	0	2	5												Q
N	-3	1	0	-2	0	-2	1	0	0	6											N
H	-3	-1	-2	-2	-2	-2	-1	0	0	1	8										H
R	-3	-1	-1	-1	-2	-2	-2	0	1	0	0	5									R
K	-3	0	-1	-1	-2	-1	-1	1	1	0	-1	2	5								K
M	-1	-1	-1	-1	-3	-2	-3	-2	0	-2	-2	-1	-1	5							M
I	-1	-2	-1	-1	-4	-3	-3	-3	-3	-3	-3	-3	-3	1	4						I
L	-1	-2	-1	-1	-4	-3	-4	-3	-2	-3	-3	-2	-2	2	2	4					L
V	-1	-2	0	0	-3	-2	-3	-2	-2	-3	-3	-3	-2	1	3	1	4				V
W	-2	-3	-2	-3	-2	-4	-4	-3	-2	-4	-2	-3	-3	-1	-3	-2	-3	11			W
Y	-2	-2	-2	-2	-3	-3	-3	-2	-1	-2	2	-2	-2	-1	-1	-1	-1	2	7		Y
F	-2	-2	-2	-2	-3	-4	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	1	3	6	F
	C	S	T	A	G	P	D	E	Q	N	H	R	K	M	I	L	V	W	Y	F	



## Affine Gap costs

*In sequence evolution, insertions of single stretch of  $k$  characters much more likely than  $k$  isolated (single-character) insertions*

*So far, we score these the same.*



# Affine Gap costs

*In sequence evolution, insertions of single stretch of  $k$  characters much more likely than  $k$  isolated (single-character) insertions*

*So far, we score these the same.*

↪ *affine gap costs:*

score  $k$  contiguous insertions (or  $k$  contiguous deletions) instead as  $g_0 + k \cdot g$   
(usually then  $g_0 \gg g$ )

► If we represent contiguous insertions as  $\begin{bmatrix} + \\ c_1 \end{bmatrix} \begin{bmatrix} - \\ c_2 \end{bmatrix} \cdots \begin{bmatrix} - \\ c_k \end{bmatrix}$   
can assign  $S\left(\begin{bmatrix} + \\ c \end{bmatrix}\right) = g_0 + g$  and  $S\left[\begin{bmatrix} - \\ c \end{bmatrix}\right] = g$ .

► DP algorithm can be extended to handle these refined scores

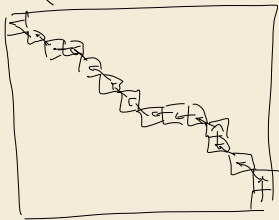
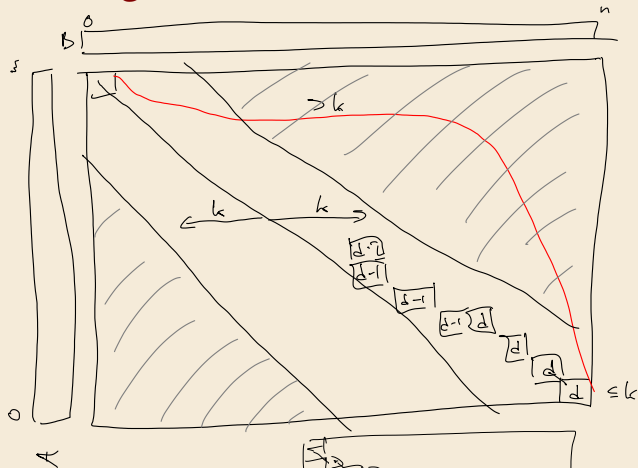
↪ exercises



## 3.5 Bounded-Distance Alignments



# Good Alignment or Abort



traceback path

Given:  $A[0..m]$

$B[0..n]$

ED score

global alignment

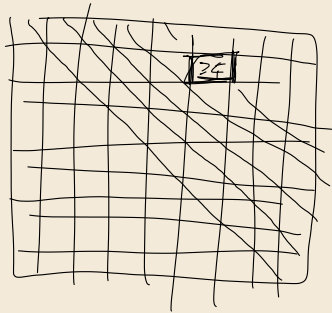
+  $k \in \mathbb{N}$

Goal: If  $d_{\text{edit}}(A, B) \leq k$

output alignment

otherwise "large dist"





• off diagonal by  $d$  steps  $\Rightarrow D \geq d$

$\Rightarrow$  only need to consider band matrix  
of size  $n \cdot 2k$  instead of full  $m \cdot n$  matrix

$\Rightarrow$  in DP recurrence, skip any option  
where  $|j-i| > k$

if  $D[m][n] > k$  don't get correct distance!

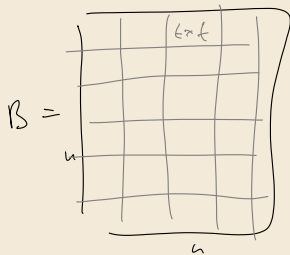
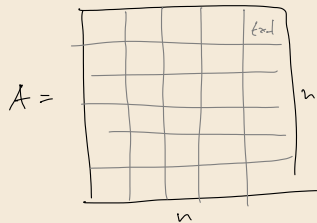
if  $\leq k$  yes!



## 3.6 Exhaustive Tabulation



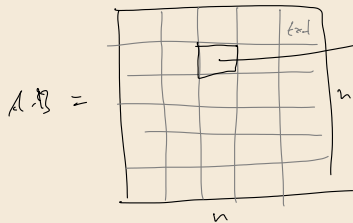
# Boolean matrix multiplication



$$\frac{n^3}{\log^2 n}$$

$$A \cdot B \quad O(n^3)$$

(ignore DnC)



$\Sigma$  of products of  $t \times t$  matrices from  $A$  &  $B$

$2^{t^2}$  different  $\square$

$\rightarrow$  compute ALL products  
up front!

+ use indirect addressing



## Four Russians?

The *exhaustive-tabulation technique* to follow is often called “Four Russians trick” . . .

- ▶ The algorithmic technique was published 1970 by V. L. Arlazarov, E. A. Dinitz, M. A. Kronrod, and I. A. Faradžev
- ▶ all worked in Moscow at that time . . . but not even clear if all are Russians  
(Arlazarov and Kronrod are Russian)



## Four Russians?

The *exhaustive-tabulation technique* to follow is often called “Four Russians trick” . . .

- ▶ The algorithmic technique was published 1970 by V. L. Arlazarov, E. A. Dinitz, M. A. Kronrod, and I. A. Faradžev
- ▶ all worked in Moscow at that time . . . but not even clear if all are Russians  
(Arlazarov and Kronrod are Russian)
- ▶ American authors coined the othering term “Method of Four Russians”  
. . . name in widespread use



# A Trick for Matrix Multiplication

Suppose we want to multiply two  $n \times n$  Boolean matrices  $C = A \cdot B$ .

We divide  $A$ ,  $B$ , and  $C$  into  $\ell \times \ell$  *micro matrices*.

$\rightsquigarrow$   $C$  consists of  $\left(\frac{n}{\ell}\right)^2$  micro matrices, each of which is the sum of  $\frac{n}{\ell}$  micro-matrix products.



# A Trick for Matrix Multiplication

Suppose we want to multiply two  $n \times n$  Boolean matrices  $C = A \cdot B$ .

We divide  $A$ ,  $B$ , and  $C$  into  $\ell \times \ell$  *micro matrices*.

$\rightsquigarrow$   $C$  consists of  $\left(\frac{n}{\ell}\right)^2$  micro matrices, each of which is the sum of  $\frac{n}{\ell}$  micro-matrix products.

The number of *different* possible micro matrix products is  $L = 2^{\ell^2} \cdot 2^{\ell^2}$ .

If we pick  $\ell = \frac{1}{4}\sqrt{\lg n}$ , we have only  $L = 2^{2\ell^2} = \sqrt{n}$  different products.

$\rightsquigarrow$  **Exhaustive Tabulation:** Can precompute all  $\sqrt{n}$  possible micro-matrix sums/products!



# A Trick for Matrix Multiplication

Suppose we want to multiply two  $n \times n$  Boolean matrices  $C = A \cdot B$ .

We divide  $A$ ,  $B$ , and  $C$  into  $\ell \times \ell$  *micro matrices*.

↪  $C$  consists of  $\left(\frac{n}{\ell}\right)^2$  micro matrices, each of which is the sum of  $\frac{n}{\ell}$  micro-matrix products.

The number of *different* possible micro matrix products is  $L = 2^{\ell^2} \cdot 2^{\ell^2}$ .

If we pick  $\ell = \frac{1}{4}\sqrt{\lg n}$ , we have only  $L = 2^{2\ell^2} = \sqrt{n}$  different products.

↪ **Exhaustive Tabulation:** Can *precompute* all  $\sqrt{n}$  possible micro-matrix sums/products!

For two micro matrices  $a$  and  $b$ , we store  $a \cdot b$  at the offset  $a_{1,1} \dots a_{\ell,\ell} b_{1,1} \dots b_{\ell,\ell}$ , where we interpret this bitstring as a binary number.

On a word RAM, we can use this as indirect memory access in  $O(1)$  time.

↪ Any micro matrix sum/product takes  $O(1)$  time  
after a total of  $O(\sqrt{n} \cdot \log^{3/2} n)$  preprocessing.



# A Trick for Matrix Multiplication

Suppose we want to multiply two  $n \times n$  Boolean matrices  $C = A \cdot B$ .

We divide  $A$ ,  $B$ , and  $C$  into  $\ell \times \ell$  *micro matrices*.

↪  $C$  consists of  $(\frac{n}{\ell})^2$  micro matrices, each of which is the sum of  $\frac{n}{\ell}$  micro-matrix products.

The number of *different* possible micro matrix products is  $L = 2^{\ell^2} \cdot 2^{\ell^2}$ .

If we pick  $\ell = \frac{1}{4}\sqrt{\lg n}$ , we have only  $L = 2^{2\ell^2} = \sqrt{n}$  different products.

↪ **Exhaustive Tabulation:** Can precompute all  $\sqrt{n}$  possible micro-matrix sums/products!

For two micro matrices  $a$  and  $b$ , we store  $a \cdot b$  at the offset  $a_{1,1} \dots a_{\ell,\ell} b_{1,1} \dots b_{\ell,\ell}$ , where we interpret this bitstring as a binary number.

On a word RAM, we can use this as indirect memory access in  $O(1)$  time.

↪ Any micro matrix sum/product takes  $O(1)$  time  
after a total of  $O(\sqrt{n} \cdot \log^{3/2} n)$  preprocessing.

The total time to compute one micro matrix in  $C$  is thus  $O(\frac{n}{\ell})$ .

So the total time to compute  $C$  is  $O(n^3/\ell^3) = O(n^3/\log^{3/2} n)$ .



# A Trick for Matrix Multiplication

Suppose we want to multiply two  $n \times n$  Boolean matrices  $C = A \cdot B$ .

We divide  $A$ ,  $B$ , and  $C$  into  $\ell \times \ell$  *micro matrices*.

↪  $C$  consists of  $\left(\frac{n}{\ell}\right)^2$  micro matrices, each of which is the sum of  $\frac{n}{\ell}$  micro-matrix products.

The number of *different* possible micro matrix products is  $L = 2^{\ell^2} \cdot 2^{\ell^2}$ .

If we pick  $\ell = \frac{1}{4}\sqrt{\lg n}$ , we have only  $L = 2^{2\ell^2} = \sqrt{n}$  different products.

↪ **Exhaustive Tabulation:** Can precompute all  $\sqrt{n}$  possible micro-matrix sums/products!

For two micro matrices  $a$  and  $b$ , we store  $a \cdot b$  at the offset  $a_{1,1} \dots a_{\ell,\ell} b_{1,1} \dots b_{\ell,\ell}$ , where we interpret this bitstring as a binary number.

On a word RAM, we can use this as indirect memory access in  $O(1)$  time.

↪ Any micro matrix sum/product takes  $O(1)$  time  
after a total of  $O(\sqrt{n} \cdot \log^{3/2} n)$  preprocessing.

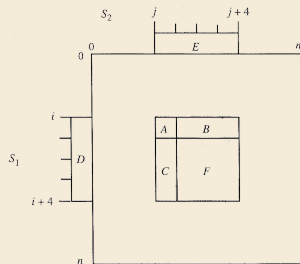
The total time to compute one micro matrix in  $C$  is thus  $O(\frac{n}{\ell})$ .

So the total time to compute  $C$  is  $O(n^3/\ell^3) = O(n^3/\log^{3/2} n)$ .

Note: By taking  $n \times \ell$  resp.  $\ell \times n$  “micro strips” instead of squares, we can choose  $\ell = \Theta(\log n)$  and obtain final time  $O(n^3/\log^2 n)$ .



# Exhaustive Tabulation for Edit Distance



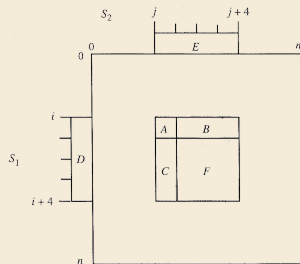
## Micro matrix

- ▶ Split  $D(i, j)$  matrix Again  $\ell \times \ell$  submatrices corresponding to  $\ell$ -char substrings of  $S_1$  and  $S_2$
- ▶ values in  $F$  only depend on  $A, B, C, D$ , and  $E$ !
- ~> can make progress micro matrix by micro matrix

Gusfield, *Algorithms on Strings, Trees, and Sequences*, Fig. 12.21



# Exhaustive Tabulation for Edit Distance



## Micro matrix

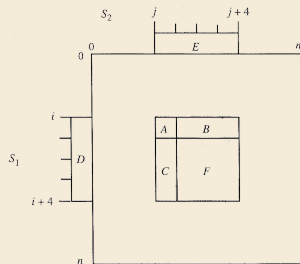
- ▶ Split  $D(i, j)$  matrix Again  $\ell \times \ell$  submatrices corresponding to  $\ell$ -char substrings of  $S_1$  and  $S_2$
- ▶ values in  $F$  only depend on  $A, B, C, D$ , and  $E$ !
- ~> can make progress micro matrix by micro matrix

Gusfield, *Algorithms on Strings, Trees, and Sequences*, Fig. 12.21

*But ... exhaustive tabulation doesn't seem to work! The values of  $D(i, j)$  keep increasing!  
How shall we bound the number of possible micro matrices?*



# Exhaustive Tabulation for Edit Distance



## Micro matrix

- ▶ Split  $D(i, j)$  matrix Again  $\ell \times \ell$  submatrices corresponding to  $\ell$ -char substrings of  $S_1$  and  $S_2$
- ▶ values in  $F$  only depend on  $A, B, C, D$ , and  $E$ !
- ~> can make progress micro matrix by micro matrix

Gusfield, *Algorithms on Strings, Trees, and Sequences*, Fig. 12.21

*But ... exhaustive tabulation doesn't seem to work! The values of  $D(i, j)$  keep increasing! How shall we bound the number of possible micro matrices?*

- ▶ **Observation:** The difference between neighboring cells  $D(i, j)$  and  $D(i, j + 1)$  respectively  $D(i, j)$  and  $D(i + 1, j)$  is in  $\{-1, 0, +1\}$ .
  - ▶  $D(i, j + 1) \leq D(i, j) + 1$  is trivial from recurrence
  - ▶  $D(i, j) \leq D(i, j + 1) + 1$  needs closer look / case distinction



~> Apply tabulation for offset, not actual values in  $D(i, j)$



## Putting the Micro Matrices together

	A	T	T	C	A
G	1	-1	0	1	1
C	-1				0
A	-1				0
T	1				-1
T	1	1	-1	0	1

... A T T C A ...

.				
G			-1 -1 0 1	
C			-1            0	
A			-1            0	
T			1             -1	
.			1 1 -1 0 1	
.				
.				

- ▶ Choose micro matrices with one row/col overlapping
- ▶ initialize first row and col (as per recurrence)

Brubach Ghurye, *A Succinct Four Russians Speedup for Edit Distance Computation and One-against-many Banded Alignment*, Fig. 1



## Putting the Micro Matrices together

	A	T	T	C	A
G		1	-1	0	1
C	-1				0
A	-1				0
T	1				-1
T	1	1	-1	0	1

[illegible]

Brubach Ghurye, *A Succinct Four Russians Speedup for Edit Distance Computation and One-against-many Banded Alignment*, Fig. 1

- ▶ Choose micro matrices with one row/col overlapping
- ▶ initialize first row and col (as per recurrence)
- ▶ number of different micro matrices:  
 $\leq \sigma^{2\ell} \cdot 3^{2(\ell-1)}$

$$\rightsquigarrow \ell \leq \frac{1}{4} \log_{3\sigma}(n) \quad \text{for } O(\sqrt{n}) \text{ micro matrices}$$



Brubach Ghurye, *A Succinct Four Russians Speedup for Edit Distance Computation and One-against-many Banded Alignment*, Fig. 1

### My Banded Alignment, Fig. 1

- $$\rightsquigarrow \ell \leq \frac{1}{4} \log_{3\sigma}(n) \quad \text{for } O(\sqrt{n}) \text{ micro matrices}$$

- $\leadsto O(1)$  time per micro matrix

 $\rightsquigarrow O(n^2/\log^2 n)$  time overall



# Can we do better?

## Theorem 3.1 (Conditional Lower Bound for Edit Distance)

An algorithm for computing the edit distance of any two strings of length  $n$  in time  $O(n^{2-\delta})$  for constant  $\delta > 0$  would refute the Strong Exponential-Time Hypothesis. ◀



Backurs, Indyk: *Edit Distance Cannot Be Computed in Strongly Subquadratic Time (unless SETH is false)*, STOC 2015

## Definition 3.2 (Exponential-Time Hypothesis)

The *Exponential-Time Hypothesis (ETH)* asserts that there is a constant  $\delta > 0$  so that every algorithm for 3SAT requires  $\Omega(2^{\delta k})$  time, where  $k$  is the number of variables. ◀

## Definition 3.3 (Strong Exponential-Time Hypothesis)

The *Strong Exponential-Time Hypothesis (SETH)* asserts that for every  $\varepsilon > 0$  there is a  $k$  such that  $k$ SAT requires  $\Omega(2^{(1-\varepsilon)k})$  time, where  $k$  is the number of variables. ◀



# Can we do better?

## Theorem 3.1 (Conditional Lower Bound for Edit Distance)

An algorithm for computing the edit distance of any two strings of length  $n$  in time  $O(n^{2-\delta})$  for constant  $\delta > 0$  would refute the Strong Exponential-Time Hypothesis. ◀



Backurs, Indyk: *Edit Distance Cannot Be Computed in Strongly Subquadratic Time (unless SETH is false)*, STOC 2015

## Definition 3.2 (Exponential-Time Hypothesis)

The *Exponential-Time Hypothesis (ETH)* asserts that there is a constant  $\delta > 0$  so that every algorithm for 3SAT requires  $\Omega(2^{\delta k})$  time, where  $k$  is the number of variables. ◀

## Definition 3.3 (Strong Exponential-Time Hypothesis)

The *Strong Exponential-Time Hypothesis (SETH)* asserts that for every  $\varepsilon > 0$  there is a  $k$  such that  $k$ SAT requires  $\Omega(2^{(1-\varepsilon)k})$  time, where  $k$  is the number of variables. ◀

*Unlikely to see “truly subquadratic” algorithms (even for constant alphabets)*



## 3.7 Linear-Space Alignments



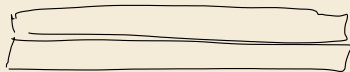
# Saving Space is Easy for Score

Assume here that  $n \leq m$ .

DP for  $D[i][j]$ ,

only need  $O(n)$  space:

- ▶  $D[i][j]$  depends on  $D[i-1][j]$ ,  $D[i][j-1]$ , and  $D[i-1][j-1]$ .
- ▶ clearly enough to keep previous and current row of  $D$





# Saving Space is Easy for Score

Assume here that  $n \leq m$ .

DP for  $D[i][j]$ ,  
only need  $O(n)$  space:

- ▶  $D[i][j]$  depends on  $D[i-1][j]$ ,  $D[i][j-1]$ , and  $D[i-1][j-1]$ .
- ▶ clearly enough to keep previous and current row of  $D$
- ▶ actually, can even *overwrite* as we go along  
     $\rightsquigarrow$  single row sufficient

```
1 procedure Score(A[0..m], B[0..n])
2   D := ScoresRow(A, B)
3   return D[n]
4
5 procedure ScoresRow(A[0..m], B[0..n])
6   D[0..n] := new array
7   for j := 0, ..., n
8     D[j] := j · g
9   for i := 1, ..., m
10    match := (i - 1) · g
11    for j = 1, ..., n
12      new := min { match + p(A[i - 1], B[j - 1])
13                  D[j] + g
14                  D[j - 1] + g
15      match := D[j]
16      D[j] := new
```



# The Middle-Point Problem

*To reconstruct alignment/edit script **using standard backtrace**, need full table  $D[0..n][0..m]$ .*

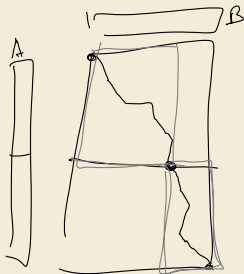


# The Middle-Point Problem

To reconstruct alignment/edit script **using standard backtrace**, need full table  $D[0..n][0..m]$ .

But can also reconstruct edit script using Divide & Conquer DP approach!

- **Idea:** Construct edit script for turning  $A[0..m/2)$  into  $B[0..j^*)$  and for turning  $A[m/2..m)$  into  $B[j^*..n)$
- But we don't know *middle point*  $j^*$  ... so need to **guess** it!  $\rightsquigarrow$  use DP!





# The Middle-Point Problem

To reconstruct alignment/edit script *using standard backtrace*, need full table  $D[0..n][0..m]$ .

But can also reconstruct edit script using Divide & Conquer DP approach!

- ▶ **Idea:** Construct edit script for turning  $A[0..m/2)$  into  $B[0..j^*)$   
and for turning  $A[m/2..m)$  into  $B[j^*..n)$
- ▶ But we don't know *middle point*  $j^* \dots$  so need to **guess** it!  $\rightsquigarrow$  use DP!

*Hold on, are we running in circles?*



# The Middle-Point Problem

To reconstruct alignment/edit script **using standard backtrace**, need full table  $D[0..n][0..m]$ .

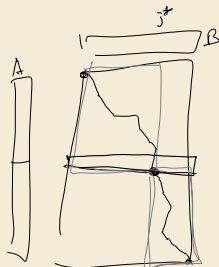
But can also reconstruct edit script using Divide & Conquer DP approach!

- **Idea:** Construct edit script for turning  $A[0..m/2)$  into  $B[0..j^*)$  and for turning  $A[m/2..m)$  into  $B[j^*..n)$
- But we don't know *middle point*  $j^*$  ... so need to **guess** it!  $\rightsquigarrow$  use DP!

*Hold on, are we running in circles?*

No!  $j^*$  optimizes **sum of scores** of  $A[0..m/2) \rightarrow B[0..j^*)$  and  $A[m/2..m) \rightarrow B[j^*..n)$

→ Can use linear-space ScoresRow!





# The Middle-Point Problem

To reconstruct alignment/edit script **using standard backtrace**, need full table  $D[0..n][0..m]$ .

But can also reconstruct edit script using Divide & Conquer DP approach!

- ▶ **Idea:** Construct edit script for turning  $A[0..m/2)$  into  $B[0..j^*)$   
and for turning  $A[m/2..m)$  into  $B[j^*..n)$
- ▶ But we don't know *middle point*  $j^* \dots$  so need to **guess** it!  $\rightsquigarrow$  use DP!

*Hold on, are we running in circles?*

No!  $j^*$  optimizes **sum of scores** of  $A[0..m/2) \rightarrow B[0..j^*)$  and  $A[m/2..m) \rightarrow B[j^*..n)$

$\rightsquigarrow$  Can use linear-space ScoresRow!

- ▶ Score for  $A[0..m/2) \rightarrow B[0..j^*)$  is  $D[m/2][j^*]$



# The Middle-Point Problem

To reconstruct alignment/edit script **using standard backtrace**, need full table  $D[0..n][0..m]$ .

But can also reconstruct edit script using Divide & Conquer DP approach!

- ▶ **Idea:** Construct edit script for turning  $A[0..m/2)$  into  $B[0..j^*)$   
and for turning  $A[m/2..m)$  into  $B[j^*..n)$
- ▶ But we don't know *middle point*  $j^* \dots$  so need to **guess** it!  $\rightsquigarrow$  use DP!

*Hold on, are we running in circles?*

No!  $j^*$  optimizes **sum of scores** of  $A[0..m/2) \rightarrow B[0..j^*)$  and  $A[m/2..m) \rightarrow B[j^*..n)$

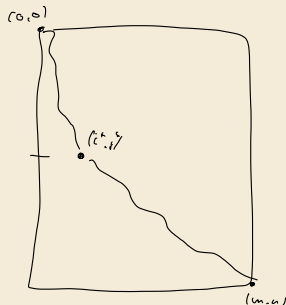
$\rightsquigarrow$  Can use linear-space ScoresRow!

- ▶ Score for  $A[0..m/2) \rightarrow B[0..j^*)$  is  $D[m/2][j^*]$
- ▶ For  $A[m/2..m) \rightarrow B[j^*..n)$  we don't have an entry in  $D$ !
- ▶ But we can **reverse**  $A$  and  $B$



# Linear-Space Alignment

```
1 procedure editScript( $A[0..m]$ ,  $B[0..n]$ )
2   if  $m == 0$  then return Insert( $B[0]$ ), ..., Insert( $B[n - 1]$ )
3   else if  $n == 0$  then return Delete( $A[0]$ ), ..., Delete( $A[m - 1]$ )
4   else if  $m == 1$ 
5      $j := \arg \min_{0 \leq j < n} p(A[0], B[j])$ 
6     return Insert( $B[0..j]$ ), Replace( $A[0]$ ,  $B[j]$ ), Insert( $B[j + 1..n]$ )
7   else
8      $i^* := \lfloor \frac{m}{2} \rfloor$ 
9      $D_{top} := \text{ScoresRow}(A[0..i^*], B)$ 
10     $D_{bottom} := \text{ScoresRow}(A[i^*..m]^R, B^R)$  //  $s^R$  is  $s$  reversed
11     $j^* := \arg \min_{0 \leq j \leq n} D_{top}[j] + D_{bottom}[n - j]$ 
12    return editScript( $A[0..i^*]$ ,  $B[0..j^*]$ ), editScript( $A[i^*..m]$ ,  $B[j^*..n]$ )
13  endif
```



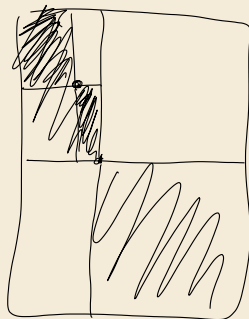
$(i^*, j^*)$  middle point



# Linear-Space Alignment

$$m \geq n$$

```
1 procedure editScript(A[0..m], B[0..n])
2   if  $m == 0$  then return Insert( $B[0]$ ), ..., Insert( $B[n - 1]$ )
3   else if  $n == 0$  then return Delete( $A[0]$ ), ..., Delete( $A[m - 1]$ )
4   else if  $m == 1$ 
5      $j := \arg \min_{0 \leq j < n} p(A[0], B[j])$ 
6     return Insert( $B[0..j]$ ), Replace( $A[0], B[j]$ ), Insert( $B[j + 1..n]$ )
7   else
8      $i^* := \lfloor \frac{m}{2} \rfloor$ 
9      $D_{top} := \text{ScoresRow}(A[0..i^*], B)$ 
10     $D_{bottom} := \text{ScoresRow}(A[i^*..m]^R, B^R) // s^R \text{ is } s \text{ reversed}$ 
11     $j^* := \arg \min_{0 \leq j \leq n} D_{top}[j] + D_{bottom}[n - j]$ 
12    return editScript( $A[0..i^*]$ ,  $B[0..j^*]$ ), editScript( $A[i^*..m]$ ,  $B[j^*..n]$ )
13  endif
```



- ▶ Non-recursive cost  $\Theta(n \cdot m)$  for ScoresRow
- ▶ “Area”  $n \cdot m$  in recursive calls is **halved** in each step.
- ↪ Total time  $\Theta(nm)$ , but using only  $\Theta(\min n, m)$  space

conservative estimate     $\log m$  rec. calls  
each cost  $\Theta(nm)$



## 3.8 Multiple Sequence Alignment



# Multiple-Sequence Alignment

*Biological sequences are often too noisy to recognize preserved regions from pairwise alignments.*



# Multiple-Sequence Alignment

*Biological sequences are often too noisy to recognize preserved regions from pairwise alignments.*

A shared region between two sequences could be random coincidence.

A shared region between **many** sequences hardly are.



# Multiple-Sequence Alignment

*Biological sequences are often too noisy to recognize preserved regions from pairwise alignments.*

A shared region between two sequences could be random coincidence.

A shared region between **many** sequences hardly are.

*“One or two homologous sequences whisper . . . a full multiple alignment shouts out loud”*

(Arthus Lesk)



## Multiple-Sequence Alignment

*Biological sequences are often too noisy to recognize preserved regions from pairwise alignments.*

A shared region between two sequences could be random coincidence.

A shared region between **many** sequences hardly are.

*"One or two homologous sequences whisper . . . a full multiple alignment shouts out loud"*

(Arthus Lesk)

**Example:**  *$\beta$ -globin* in different species:

[illegible]African Clawed Frog (*Xenopus laevis*): P02133

Zebrafish (*Danio rerio*): Q90486

Chicken (*Gallus gallus*): P02112Human (*Homo sapiens*): P68871

Mouse (*Mus musculus*): P02088

<https://www.ebi.ac.uk/jdispatcher/msa/clustalo>



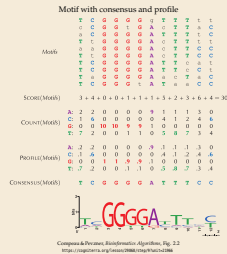
# Scoring Multiple Alignments

- ▶ Given sequences  $A_1[0..n_1), \dots, A_k[0..n_k)$  over common alphabet  $\Sigma$
- ▶ alignment is sequence of *columns* in  $(\Sigma_-)^k$  with  $\Sigma_- = \Sigma \cup \{-\}$



# Scoring Multiple Alignments

- ▶ Given sequences  $A_1[0..n_1), \dots, A_k[0..n_k)$  over common alphabet  $\Sigma$
- ▶ alignment is sequence of *columns* in  $(\Sigma_-)^k$  with  $\Sigma_- = \Sigma \cup \{-$
- ▶ going from 2 to  $k$  sequences requires score for  $k$ -columns
  - ▶ different options
  - ▶ One option: total Hamming distance (see Unit 2 for motifs)

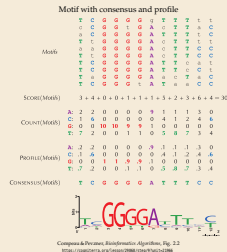




## Scoring Multiple Alignments

- ▶ Given sequences  $A_1[0..n_1), \dots, A_k[0..n_k)$  over common alphabet  $\Sigma$
- ▶ alignment is sequence of *columns* in  $(\Sigma_-)^k$  with  $\Sigma_- = \Sigma \cup \{-\}$
- ▶ going from 2 to  $k$  sequences requires score for  $k$ -columns
  - ▶ different options
  - ▶ One option: total Hamming distance (see Unit 2 for motifs)
  - ▶ Here: **SP-Score** (*sum-of-pairs score*) w.r.t.  $S$

$$d_{SP}\left(\begin{bmatrix} c_1 \\ \vdots \\ c_k \end{bmatrix}\right) = \sum_{1 \leq i < j \leq k} S\left(\begin{bmatrix} c_i \\ c_j \end{bmatrix}\right) \quad \text{for } S \text{ any pairwise-alignment score}$$



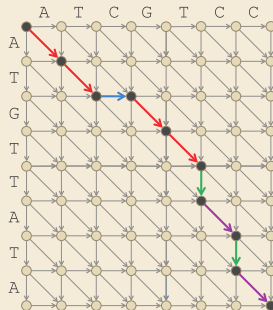


# Dynamic Programming Solution

Pairwise alignment = path in grid graph; optimal alignment = shortest path between corners



match/mismatch (/) , insertion () , or deletion () .



Compeau & Pevzner, *Bioinformatics Algorithms*, Fig. 5.5 & 5.6  
<https://cogniterra.org/lesson/29932/step/1?unit=22029>

DP solution with 2D matrix  $D[0..m][0..n]$

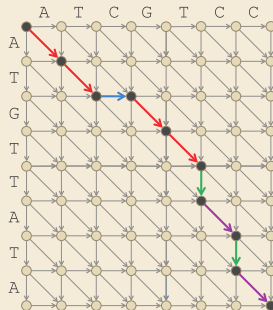


# Dynamic Programming Solution

Pairwise alignment = path in grid graph; optimal alignment = shortest path between corners



match/mismatch (↘/↘), insertion (→), or deletion (↓).



Compeau & Pevzner, *Bioinformatics Algorithms*, Fig. 5.5 & 5.6  
<https://cogniterra.org/lesson/29932/step/1?unit=22029>

↪ DP solution with 2D matrix  $D[0..m][0..n]$

For  $k$  strings, shortest path in  $k$ -dimensional grid graph

↪  $n_1 \cdot n_2 \cdot \dots \cdot n_k$  vertices to consider



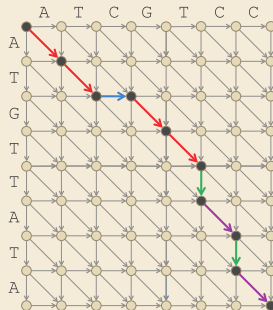


# Dynamic Programming Solution

Pairwise alignment = path in grid graph; optimal alignment = shortest path between corners

0	1	2	2	3	4	5	6	7	8
	A	T	-	G	T	T	A	T	A
	A	T	C	G	T	-	C	-	C
0	1	2	3	4	5	5	6	6	7
	↘	↘	→	↘	↘	↓	↘	↓	↘

match/mismatch (↘/↘), insertion (→), or deletion (↓).



Compeau & Pevzner, *Bioinformatics Algorithms*, Fig. 5.5 & 5.6  
<https://cogniterra.org/lesson/29932/step/1?unit=22029>

↪ DP solution with 2D matrix  $D[0..m][0..n]$

For  $k$  strings, shortest path in  $k$ -dimensional grid graph

↪  $n_1 \cdot n_2 \cdot \dots \cdot n_k$  vertices to consider for  $k$  strings of  $n$  characters  $\Theta(n^k)$  time ⚡



## Bad News (Again)

MULTIPLE ALIGNMENT WITH SP-SCORE is NP-hard for any  $\sigma \geq 2$  and any metric  $S$



Elias: *Settling the Intractability of Multiple Alignment*, J. of Computational Biology 2006







# Bounding SP-scores

*Not all hope is lost.*

SP-score can be bounded by optimal pairwise alignments and heuristic for some alignment:

$$\sum_{1 \leq i < j \leq k} d_S(A_i, A_j) \leq d_{SP}(A_1, \dots, A_k) \leq d_{SP}(\text{some alignment})$$
$$\sum_{i \in I} d_S(M[:, i], [i, :])$$

- ▶ can be the basis for a Branch & Bound algorithm
- ▶ but: need efficient approximation algorithm for MULTIPLE ALIGNMENT WITH SP-SCORE



# Bounding SP-scores

*Not all hope is lost.*

SP-score can be bounded by optimal pairwise alignments and heuristic for some alignment:

$$\sum_{1 \leq i < j \leq k} d_S(A_i, A_j) \leq d_{SP}(A_1, \dots, A_k) \leq d_{SP}(\text{some alignment})$$

- ▶ can be the basis for a Branch & Bound algorithm
- ▶ but: need efficient approximation algorithm for MULTIPLE ALIGNMENT WITH SP-SCORE

↪ *Can we build a multiple alignment by successively adding in one new sequence at a time?*



## Extending Pairwise Alignments is tricky

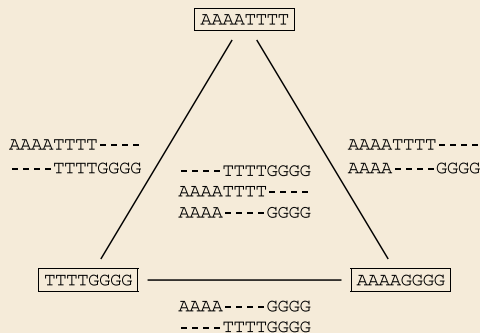
*Can we combine optimal pairwise alignment into a multiple alignment?*



# Extending Pairwise Alignments is tricky

*Can we combine optimal pairwise alignment into a multiple alignment?*

Sometimes Yes!



(a) Compatible pairwise alignments

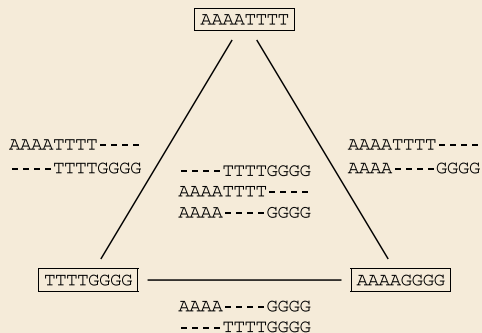
Jones & Pevzner, *Bioinformatics Algorithms*, Fig 6.22a



# Extending Pairwise Alignments is tricky

*Can we combine optimal pairwise alignment into a multiple alignment?*

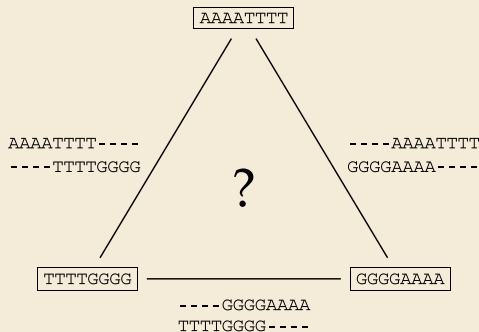
Sometimes Yes!



(a) Compatible pairwise alignments

Jones & Pevzner, *Bioinformatics Algorithms*, Fig 6.22a

But No in general ...



(b) Incompatible pairwise alignments

Jones & Pevzner, *Bioinformatics Algorithms*, Fig 6.22b



# Alignment Trees

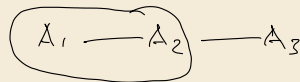
*Problem in example comes precisely from cycle!*



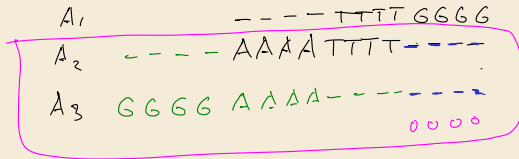
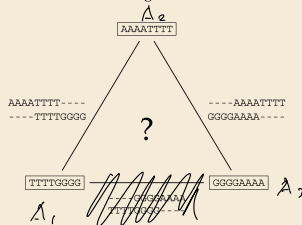
# Alignment Trees

*Problem in example comes precisely from cycle!*

- ▶ Given a *tree* over sequences  $A_1, \dots, A_k$
- ▶ Compute optimal *pairwise* alignments along all  $k - 1$  tree edges
- ▶ Build multiple alignment one edge at a time
  - ▶ Here, use  $[-]$  for every gap symbol in either endpoint of an edge  
We always assume  $S([-]) = 0$



But No in general ...





# Alignment Trees

*Problem in example comes precisely from cycle!*

- ▶ Given a *tree* over sequences  $A_1, \dots, A_k$
- ▶ Compute optimal *pairwise* alignments along all  $k - 1$  tree edges
- ▶ Build multiple alignment one edge at a time
  - ▶ Here, use  $[-]$  for every gap symbol in either endpoint of an edge  
We always assume  $S([-]) = 0$
- ▶ **Notation:**
  - ▶  $M \in (\Sigma_-^k)^N$  multiple alignment of length  $N \geq \max n_j$
  - ▶  $d_{SP}$  SP-Score w.r.t. pairwise score  $S$
  - ▶  $d_S(A, B)$  score of optimal pairwise alignment of  $A$  and  $B$
  - ▶  $M$  induces pairwise alignment  $M[:, [i, j]]$  for  $A_i$  and  $A_j$   
Note:  $S(M[:, [i, j]]) \geq d_S(A_i, A_j)$  and in general not optimal

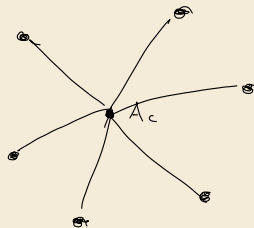


# Center-Star Approximation

Use simplest possible tree: A *star*!

## Center-Star Multiple Sequence Alignment

1. Compute all pairwise distances  $d_S(A_i, A_j)$
2. Find  $c \in [k]$  that minimizes  $\sum_j d_S(A_c, A_j)$
3. Construct  $M$  as alignment consistent with star alignment with center  $S_c$ .





# Center-Star Approximation – Analysis

## Theorem 3.4

Assume  $d_S$  is a metric for pairwise alignments. The center-star alignment for  $k$  strings is a  $(2 - \frac{2}{k})$ -approximation w.r.t. to the SP-score of the multiple sequence alignment. ◀

$M^*$  optimal MSA

$M_c$  center-star alignment

$$d_{SP} \triangleq \sum_{1 \leq i < j \leq k} S(M[:, \{i, j\}])$$

$$v \triangleq \sum_{1 \leq i, j \leq k} S(M[:, \{i, j\}])$$

$$v(M) = 2 \cdot d_{SP}(M)$$

$$\textcircled{1} \quad \underset{\substack{\text{OPT} \\ \text{OPT}}}{v(M^*)} = \sum_{i,j} S(M^*[:, \{i, j\}]) \geq \sum_{i,j} d_S(A_i, A_j) \geq \underset{\substack{\text{best center} \\ |}}{k \cdot \sum_i d_S(A_c, A_i)} = k \cdot m$$

$$\textcircled{2} \quad v(M_c) \leq 2(k-1)m$$

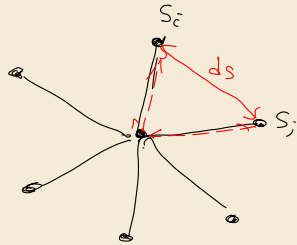
- ▶  $M \in (\Sigma^k)^N$  multiple alignment of length  $N \geq \max n_j$
  - ▶  $d_{SP}$  SP-Score w.r.t. pairwise score  $S$
  - ▶  $d_S(A, B)$  score of optimal pairwise alignment of  $A$  and  $B$
  - ▶  $M$  induces pairwise alignment  $M[:, \{i, j\}]$  for  $A_i$  and  $A_j$
- Note:  $S(M[:, \{i, j\}]) \geq d_S(A_i, A_j)$  and in general not optimal

Fact: for  $j$

$$S(M_c[:, \{c, j\}]) = d_S(A_c, A_j)_{34}$$



$$\begin{aligned}
 v(M_c) &= \sum_{i,j} S(M_c[i:j][i_{ij}]) \\
 &= 2(k-1) \cdot \sum_i S(M_c[i:j][c, i]) \\
 &\stackrel{=}{=}_{\text{fact}} 2(k-1) \underbrace{\sum_i d_S(A_c, A_i)}_{=: m}
 \end{aligned}$$



$$d_S(i, j) \leq d_S(S_i, S_c) + d_S(S_c, S_j)$$