

Alignments & Co

(approx. pattern matching)

# Basics: Dynamic Programming (DP)

"programming" refers to "tabular" method, not to writing program code.

DP is a general, powerful alg. design technique for solving optimization problems

DP  $\hat{=}$  type of very smart exhaustive search that can be applied if the problem can be "subdivided" into overlapping subproblems.

Exmp! Fibonacci numbers:  $f(1) = f(2) = 1$   
 $f(n) = f(n-1) + f(n-2)$   
 $\Rightarrow 1, 1, 2, 3, 5, 8, 13, \dots$

Naive recursive way:

```
F(n)
{
  IF (n ≤ 2) f = 1
  ELSE f = F(n-1) + F(n-2)
  return f
}
```

Exponential time!

Why: Recurrence Nr:

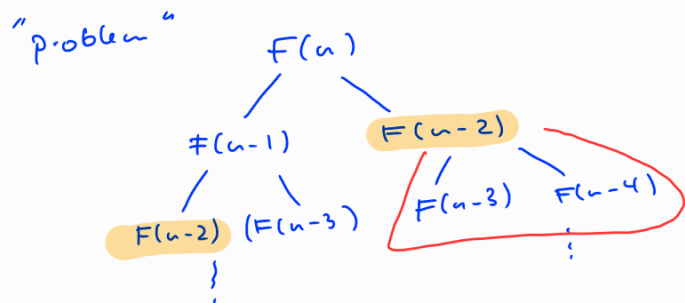
$T(n)$  = time to compute  $n$ -th Fibnr.

$$\Rightarrow T(n) = T(n-1) + T(n-2) + O(1)$$

$$\geq 2T(n-2)$$

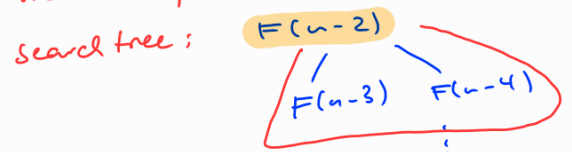
$$\begin{matrix} T(n-1) \\ \geq T(n-2) \end{matrix} \geq 2 \cdot 2 T(n-4)$$

$$\geq 2^{\frac{n}{2}} = O(2^{\frac{n}{2}})$$



We compute  $F(n)$  several times, although it is enough to compute it ones

We can prune the entire



$\Rightarrow$  we must "memorize" the intermediate results.

memo =  $\emptyset$  // dictionary

mF(n)

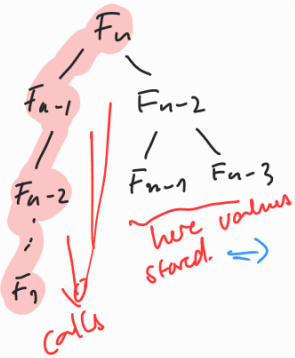
IF (n  $\in$  memo) return memo[n] // check if F(n) already computed.  
IF (n  $\leq$  2) f = 1  
ELSE f = mF(n-1) + mF(n-2) } still the same lines.  
memo[n] = f  
return f

Runtime

Sketch:  $\Rightarrow$

mF(k) only recurses the first time it's called  $\neq$  k  
& memoized calls cost  $O(1)$

the nr of non-memoized calls is n:  
mF(1), mF(2), ..., mF(k), ..., mF(n)



runtime  $O(n)$  ! (goes even faster)

iterativ\_ f(n)

list f(n) = 1xn empty list  
f(1) = f(2) = 1  
FOR (i = 3 ... n)  
[ f(i) = f(i-1) + f(i-2)  
return f

# The Longest Common Subsequence (LCS) problem

This is a classical problem in Bioinformatics, where one wants to understand how "close"

DNA/protein sequences are.

= simply a string  
some alphabet.

There are several ways to address this problem.  
Here we consider one of them: LCS.

Def: Let  $X = x_1 \dots x_n$  &  $Z = z_1 \dots z_k$  be two strings (= sequences)  
 $Z$  is a subsequence of  $X$  if  
 $\exists$  indices of  $X$   $i_1 \dots i_k$  st  $i_1 < i_2 < \dots < i_k$   
&  $z_j = x_{i_j}$

Exmpl:  $X = \text{A B C B D A B}$ ,  $Z = \text{B C D B}$

$\Rightarrow Z$  can be obtained from  $X$   
by removing elements from  $X$ .  
& keep order of remaining elements.

If  $Z$  is subsequence of  $X$  &  $Y \Rightarrow Z$  is common subsequ.  
of  $X$  &  $Y$ .

Exmpl:  $X = \text{A B C B D A B}$   
 $Y = \text{B D C A B A}$

$Z = \text{BCA}$ ,  $Z' = \text{BCBA}$ ,  $Z'' = \text{BDAB}$  are common subsequ.  
of  $X$  &  $Y$ .

Aim: Find longest common subsequ. (LCS) of  $X$  &  $Y$ .

Brute FORCE is not a good idea,  
since  $X = x_1 \dots x_n$  has  $2^n$  subsequences!

For  $X = x_1 \dots x_n$ , define the  $i$ -th prefix  $X_i$  of  $X$   
as  $X_i = x_1 \dots x_i$ ,  $1 \leq i \leq n$  & put  $X_0 = \epsilon$  "empty sequ."

Step 1: Characterization of LCS.

Theorem 4.2  
[opt. substructure of LCS]

Let  $X = x_1 \dots x_m$ ,  $Y = y_1 \dots y_n$  be two sequences.  
&  $Z = z_1 \dots z_k$  be some LCS of  $X$  &  $Y$ .

- (1)  $x_m = y_n \Rightarrow z_k = x_m = y_n$  &  $z_{k-1}$  is LCS of  $X_{m-1}$  &  $Y_{n-1}$
- (2)  $x_m \neq y_n$  &  $z_k \neq x_m \Rightarrow Z$  is LCS of  $X_{m-1}$  &  $Y$
- (3)  $x_m \neq y_n$  &  $z_k \neq y_n \Rightarrow Z$  is LCS of  $X$  &  $Y_{n-1}$ .

Illustration: (1)  
 $X = A A B A C D$   
 $Y = A B B C D$   
 LCS  $Z = A B C D$   
 (Note: A bracket underlines 'A B C' in Z, and a box encloses 'D' in both X and Y.)

$\Rightarrow z_3$  is LCS of  $X_5 = X_{n-1}$   
 $Y_4 = Y_{n-1}$   
 $m=6, n=5$   
 $\Rightarrow$  can reduce problem to find LCS of  $X_{m-1}$  &  $Y_{n-1}$

(2)  
 $X = A A B A C A$   
 $Y = A B B C D$   
 LCS  $Z = A B C$

! this is independent of letter  $y_n$   
 that is both cases may occur:  
 $z_k = y_n$  or  $z_k \neq y_n$

But since  $z_k$  never "matches"  $x_n$ ,  
 we can reduce the problem to find  
 LCS of  $X_{m-1}$  &  $Y$

(3) analog to (2).

Proof: (1)  $x_u = y_u \Rightarrow z_u = x_u$  [otherwise if  $z_u \neq x_u$  we can append  $x_u$  to  $z_1 \dots z_u$  & get larger subsequence since  $z$  is LCS.]  
[Thm 4.2] already CS

$\Rightarrow x_u = x_u = z_u$

Now prefix  $z_{u-1}$  is a common subseq. of  $x_{u-1}$  &  $y_{u-1}$

to show:  $z_{u-1}$  is LCS of  $x_{u-1}$  &  $y_{u-1}$

But this is clear, since if there is a longer one,

say  $W$  then we could append  $z_u + W$  to obtain a seqn. which is longer than  $z$   $\therefore z$  is LCS

(2) Clearly if  $z_u \neq x_u \Rightarrow z$  is common subseq. of  $x_{u-1}$  &  $y$ .

to show  $z$  is LCS of  $x_{u-1}$  &  $y$ .

Again, if there is a longer common subseq.  $W$  of  $x_{u-1}$  &  $y$

then  $W$  is also a common subseq. of  $x$  &  $y$ ,

but longer than  $z$   $\therefore z$  is LCS

(3) analog to (2)

□

Step 2: recursive solution of LCS

By Thm 4.2:  $x_u = y_u \Rightarrow$  Find LCS of  $x_{u-1}$  &  $y_{u-1}$  & append  $x_u = y_u$  to this LCS.

$x_u \neq y_u \Rightarrow$  Find LCS of  $x_{u-1}$  &  $y$  or  $x$  &  $y_{u-1}$  whichever of these is longer is an LCS of  $x$  &  $y$

Let  $c_{ij}$  = length (# of letters) of LCS of prefixes  $X_i$  &  $Y_j$

NOTE  $x_0 = \epsilon, y_0 = \epsilon \Rightarrow c_{ij} = 0$ , if  $i=0$  or  $j=0$

Recursive Formula:

$$c_{ij} = \begin{cases} 0 & , \text{if } i=0 \text{ or } j=0 \\ c_{i-1, j-1} + 1 & , i, j > 0 \text{ \& } x_i = y_j \\ \max \{c_{i, j-1}, c_{i-1, j}\} & , i, j > 0 \text{ \& } x_i \neq y_j \end{cases}$$

### Step 3: Computing length of LCS.

Based on recursive formula for  $c_{ij}$   
we get for free an exponential-time recursive alg.  
But DP can be used to get it in polynomial time.

Let  $X = x_1 \dots x_m$  &  $Y = y_1 \dots y_n$

LCS( $X, Y$ )

Let  $b[1 \dots m, 1 \dots n]$   
 $c[0 \dots m, 0 \dots n]$  be new arrays

//  $b$  used to construct LCS via backtracking  
//  $c$  stores length.

FOR ( $i = 0 \dots m$ ) DO  $c[i, 0] = 0$   
FOR ( $j = 0 \dots n$ ) DO  $c[0, j] = 0$

FOR ( $i = 1 \dots m$ ) DO

FOR ( $j = 1 \dots n$ ) DO

IF ( $x_i = y_j$ )

$c[i, j] = c[i-1, j-1] + 1$   
 $b[i, j] = "\swarrow"$

ELSE IF ( $c[i-1, j] \geq c[i, j-1]$ )

$c[i, j] = c[i-1, j]$   
 $b[i, j] = "\uparrow"$

ELSE

$c[i, j] = c[i, j-1]$   
 $b[i, j] = "\leftarrow"$

RETURN  $c$  &  $b$ .



### Step 4: construct LCS:

PRINT-LCS( $b, X, i, j$ ) // initial call: PRINT-LCS( $b, X, m, n$ )

IF ( $i = 0$  or  $j = 0$ ) RETURN.

IF ( $b[i, j] = "\swarrow"$ )

PRINT-LCS( $b, X, i-1, j-1$ )  
print  $x_i$

ELSE IF ( $b[i, j] = "\uparrow"$ )

PRINT-LCS( $b, X, i-1, j$ )

ELSE  
PRINT-LCS( $b, X, i, j-1$ )

RUNTIME:  $\Theta(m+n)$   
since it decrements at most  
one of  $i$  &  $j$  in each call

Exmpl:

$X = ABCB, Y = BDC$

	j	0	1	2	3
			B	D	C
i	0	0	0	0	0
1	A	0			
2	B	0			
3	C	0			
4	B	0			

	j	0	1	2	3
			B	D	C
i	0	0	0	0	0
1	A	0	0	0	0
2	B	0	1	1	1
3	C	0	1	1	2
4	B	0	1	1	2

now  $i = 1 \dots m / j = 1 \dots m$

$i = 1, j = 1, 2, 3$

$i = 2, j = 1, 2, 3$

$x_2 = y_1$  "↖"

$i = 3, j = 1, 2, 3$   $x_3 = x_2$  "↖"

$i = 4, j = 1, 2, 3$   $x_4 = y_1$  "↖"

	j	0	1	2	3
			B	D	C
i	0	0	0	0	0
1	A	0	0	0	0
2	B	0	1	1	1
3	C	0	1	1	2
4	B	0	1	1	2

Backtracking:

start at  $c[m, n]$

order "↖", "↑", "←"

but there might be further opt. solutions!

$b[3, 3] = \text{↖} \rightarrow \text{print 'C'}$

$b[2, 1] = \text{↖} \rightarrow \text{print 'B'}$

$\Rightarrow z = \underline{BC}$

↑  
this value is length of LCS

Theorem 4.3:  $LCS(\cdot) + \text{PRINT-LCS}(\cdot)$  correctly returns a length & LCS of given  $X = x_1 \dots x_m$  &  $Y = y_1 \dots y_n$  in  $\Theta(mn)$  time.

Proof: correctness by Thm 4.2  
runtime dominated by 2 FOR-loops ( $i = 1 \dots m, j = 1 \dots n$ )  
 $\Rightarrow \Theta(mn)$

□



## "Non-exact" matching (Alignments & Co)

so far: mostly "exact" matching (except LCS)

Differences due to sequ. error, natural variations

↓  
insertion, del,  
additions -

→ need "approx. matching"

want distance between 2 strings.

To get an idea on "how similar" 2 strings are.

## Hamming distance

Simpler idea = Hamming distance, def for 2 strings of same length.

$$\text{for } \begin{array}{l} S = s_1 \dots s_n \\ S' = s'_1 \dots s'_n \end{array} \text{ def } d_{\text{Hammy}}(s, s') = \sum_{i=1}^n \delta(s_i, s'_i)$$

$$\text{where } \delta(s_i, s'_i) = \begin{cases} 0, & s_i = s'_i \\ 1 & \text{else.} \end{cases}$$

$$[d_{\text{Hammy}}(s, s') = \# \text{ mismatches}]$$

Exmpl:  $S = \text{ATCATGA}$   
 $S' = \text{ATGATCA} \Rightarrow d_{\text{Hammy}}(S, S') = 2$

problem, in general, this is not what is meant with similarity from a biol. POV.

$$\begin{array}{l} S = \text{ATATATATA} \\ S' = \text{TATATATAT} \end{array} \left. \begin{array}{l} \} \text{very similar in their} \\ \} \text{'characteristic'} \end{array} \right\} \text{but } d_{\text{Hammy}} = 9$$

need something else.

## LCS

$$s = GC$$

$$s_1 = GCGC \rightarrow LCS = 2$$

$$s = GC$$

$$s_1 = GCGCGCGCGC \rightarrow LCS = 2$$

Both  $s$  &  $s_1$ ,  $s$  &  $s_2$  have same LCS score,  
although  $s$  &  $s_1$  "are more similar" than  $s$  &  $s_2$

LCS is good for finding inexact "subpattern"  
but not for comparing similarity between sequ.

(LCS does not take into account insertion/deletions)

# Edit-Distance (Levenshtein-Distance) & Alignments

when genes get copied etc:

A T G C G T  
A G C G T } T on pos. 2 removed  
[deletion]

A T G C G T  
A T G G C G T } G on pos 4 added  
[insertion]

A T G C G T  
T T G C G T } A on pos 1 replaced by T  
[substitution (mismatch)]

DEF edit-distance = min Nr of deletion, insertion, substitutions to transform  $s$  into  $s'$   
 $d_{\text{Edit}}(s, s')$

## Observation:

- $s, s'$  don't need to be of same length
- $d_{\text{Edit}}(s, s') = d_{\text{Edit}}(s', s)$  (symmetric)
- $||s| - |s'|| \leq d_{\text{Edit}}(s, s') \leq d_{\text{Ham}}(s, s')$

- edits don't need to be unique.  $s = AT, s' = G$   
    remov A  
    repleu T by G  
    or remov T  
    repleu A by G

$s = \text{TGCATAT}$      $s' = \text{ATCCGAT}$

$\text{TGCATAT}$  } delete T at pos 7  
 $\text{TGCATA}$  } delete A at pos 6  
 $\text{TGCAT}$  } insert A before pos 4  
 $\text{ATGCAT}$  } replace G by C  
 $\text{ATCCAT}$  } replace G by C  
 $\text{ATCCGAT}$  } insert G before A

$\text{TGCATAT}$  } add A before pos 1  
 $\text{ATGCATAT}$  } delete T pos 6  
 $\text{ATGCATAT}$  } replace A pos 5 by G  
 $\text{ATCGAT}$  } replace G pos 3 by C

$\Rightarrow d_{edit}(s, s') \leq 5$

$\Rightarrow d_{edit}(s, s') \leq 4$

But what is min edit ???

↳ derive along alignments

DEF: Alignment  $\alpha(s, s')$  of  $s, s'$ ,  $s, s' \in \Sigma^*$  ( $\Sigma$  alphabet)

is  $2 \times n$  matrix,  $n \geq \max\{|s|, |s'|\}$

such that

- ▶ each row consist of characters in  $\Sigma$  & '-' gap symbol
- ▶ deletion of gaps in each row yields  $s$  resp  $s'$ .
- ▶ in no column are 2 gaps

Exmpl

$S = T G C A T A T$   
 $S' = A T C C G A T$

$$A(S, S') = \begin{pmatrix} - & T & G & C & A & T & A & T \\ A & T & C & C & G & - & A & T \end{pmatrix}$$

↑ insertion  
↑ match  
↑ replacement (unit match)  
↑ deletion

DEF: For given alphabet  $\Sigma$ , let

$$\delta: (\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \rightarrow \mathbb{R}$$

be a cost-function

[constant case:  $\delta(x, y) = \begin{cases} 1 & \text{if } x \neq y \\ 0 & \text{else} \end{cases}$  (unit score)]

$$A(S, S') = \begin{pmatrix} a_1 & \dots & a_\ell \\ b_1 & \dots & b_\ell \end{pmatrix} \Rightarrow \delta(A) = \sum_{j=1}^{\ell} \delta(a_j, b_j)$$

Exmpl.

$$\begin{pmatrix} - & T & G & C & A & T & A & T \\ A & T & C & C & G & - & A & T \end{pmatrix}$$

unit score:

$$1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 = 4$$

Aim:

Find alignment  $A$  that minimizes  $\delta(A)$   
(= optimal alignment)

Lemma 1  $\min d_{\text{Edit}}(s, s') = \min_{\text{ut}} \mathcal{J}(\text{ut}(s, s'))$  for  $\mathcal{J}$   
= unit costs.

proof: ut yields edit-operat. & vice versa

$$\min \mathcal{J}_{\text{edit}}(s, s') \stackrel{\geq}{=} \min_{\text{ut}} \mathcal{J}(\text{ut}(s, s')) \Rightarrow "=" \quad \square$$

In general, costs differ from unit costs.

e.g.

$$\begin{aligned} \mathcal{J}(x, x) &= -2 \quad \forall x \in \{A, C, T, G\} \\ \mathcal{J}(x, -) &= \mathcal{J}(-, x) = 1 \\ \mathcal{J}(A, C) &= 5 \\ &\vdots \end{aligned}$$

DEF: let  $s, T$  be strings,  $\mathcal{J}$  be cost-fct.

Then  $D(i, j) = \text{min costs of alignments of } s[1..i] \text{ \& } T[1..j]$

[ that is  $D(|s|, |s'|) = \mathcal{J}(\text{ut}(s, s'))$  ]

Lemma 2 let  $S = s_1 \dots s_k$ ,  $T = t_1 \dots t_\ell$  strings

Then, for all  $i, j \geq 1$  it holds that

$$D(i, j) = \min \begin{cases} D(i-1, j) + \delta(s_i, -) \\ D(i, j-1) + \delta(-, t_j) \\ D(i-1, j-1) + \delta(s_i, t_j) \end{cases}$$

where

$$D(0, 0) = 0$$

$$D(0, j) = D(0, j-1) + \delta(-, t_j)$$

$$D(i, 0) = D(i-1, 0) + \delta(s_i, -)$$

$D(0, 0) = 0 \iff$  align empty string  $\epsilon$  with  $\epsilon$

$D(i, 0) \iff$  align  $S[1..i]$  with  $\epsilon$   
 $\iff$  delete all characters  $s_1 \dots s_i$  to get  $\epsilon$   
 $\iff \sum_{r=1}^i \delta(s_r, -)$

$D(0, j) \iff$  align  $\epsilon$  with  $T[1..j]$   
 $\iff$  add  $t_1 \dots t_j$  to  $\epsilon$  to get  $T[1..j]$   
 $\iff \sum_{r=1}^j \delta(-, t_r)$

in case of unit costs:  $D(0, 0) = 0$   
&  $D(i, 0) = i$  for  $i \geq 1$   
&  $D(0, j) = j$  for  $j \geq 1$



# proof of lemma 2

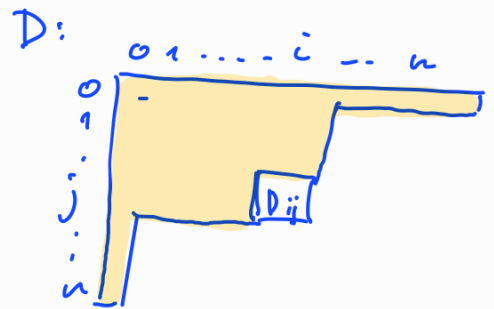
Induction.

$$D(0,0)$$

$$D(0,i) \quad i \geq 1$$

$$D(j,0) \quad j \geq 1 \quad \text{are correct.}$$

$\Rightarrow$  when computing  $D(i,j)$  we assume that all previous values have been correctly computed:



Consider  $D(i,j)$ .

let  $\sigma_{ij}$  be opt alignment for  $s[1..i]$  &  $T[1..j]$

$\Rightarrow$  3 cases how  $\sigma_{ij}$  "ends"

$$(a) \begin{pmatrix} \boxed{A} & s_i \\ & - \end{pmatrix} \quad (b) \begin{pmatrix} \boxed{B} & - \\ & t_j \end{pmatrix} \quad (c) \begin{pmatrix} \boxed{C} & s_i \\ & t_j \end{pmatrix}$$

Suppose case (a)

since  $\sigma_{ij}$  opt alignment for  $s[1..i]$  &  $T[1..j]$

$\Rightarrow A$  opt alignment for  $s[1..i-1]$  &  $T[1..j]$

(if not  $\Rightarrow \exists A''$  better alignment for  $s[1..i-1]$  &  $T[1..j]$

now add column  $\begin{pmatrix} s_i \\ - \end{pmatrix}$  to  $A''$

to get better alignment than  $\sigma_{ij}$   $\downarrow$ )

$$\Rightarrow \delta(A) \stackrel{\text{ind hyp}}{=} D(i-1,j)$$

$$\Rightarrow \delta(\sigma_{ij}) \stackrel{\text{Def}}{=} D(i,j) \stackrel{A \text{ is}}{=} D(i-1,j) + \delta(s_i, -)$$

is correctly computed

analog for case (b) & (c)

- (b) B opt alignment for  $S[1..i]$  &  $T[1..j-1]$
- (c) C opt alignment for  $S[1..i-1]$  &  $T[1..j]$ .

&  $\delta(s_i, t_j)$  correctly computed.

under the assumption that we have case a) b) c)

However at least one case must occur & no further case can occur.

$$\Rightarrow \delta(s_i, t_j) = \min \begin{cases} D(i-1, j) + \delta(s_i, -) \\ D(i, j-1) + \delta(-, t_j) \\ D(i-1, j-1) + \delta(s_i, t_j) \end{cases}$$

"def  
 $D(i, j)$

$\Rightarrow \delta(S, T) = D(|S|, |T|)$  correctly determined  $\checkmark$

globalAlign( $S, T, \delta$ ) //  $s = s_1 \dots s_k, t = t_1 \dots t_l$

init:  $D(0, 0) = 0$

FOR ( $i = 1 \dots k$ )

$D(i, 0) = D(i-1, 0) + \delta(s_i, -)$

FOR ( $j = 1 \dots l$ )

$D(0, j) = D(0, j-1) + \delta(-, t_j)$

FOR ( $i = 1 \dots k$ )

FOR ( $j = 1 \dots l$ )

$D(i, j) = \min \begin{cases} D(i-1, j) + \delta(s_i, -) \\ D(i, j-1) + \delta(-, t_j) \\ D(i-1, j-1) + \delta(s_i, t_j) \end{cases}$

Return D

By the latter arguments & lemma 2,

globalAlign correctly computes min-edit cost  
(opt. alignment score)

Runtime:  $O(|S| |T|)$

Exmpl,  $J = \text{unit costs}$ ,  $S = T b C A$   
 $t = A T C$

D	0	1	2	3	4
$\Sigma$	0	1	2	3	4
A	1	1	2	3	3
T	2	1	2	3	4
C	3	2	2	2	3

$\textcircled{3} = D(|S|, |T|)$   
 $= \text{min-edit costs.}$

So far only know value of min-edit costs,  
But not the alignment.

$\Rightarrow$  TRACEBACK (can be compute  
when computing D)

	0	1	2	3	4
D	Σ	T	b	c	A
0	Σ	0	1	2	3
1	A	1	1	2	3
2	T				
3	c				

want to compute  $D(1,4)$

$$D(1,4) = \min ( D(1,3) + 1, D(0,4) + 1, D(0,3) + 0 )$$

$$= \min ( 4, 5, 3 )$$

in this case "unique" path path that  
 $\rightarrow$   $D(1,4)$  can only be computed by  
 using entry  $D(0,3)$  & this info can be stored  
 when computing  $D(1,4)$

	0	1	2	3	4
D	Σ	T	b	c	A
0	Σ	0	1	2	3
1	A	1	1	2	3
2	T	2	1	2	
3	c				

$$D(2,2) = \min ( D(2,1) + 1, D(1,2) + 1, D(1,1) + 1 )$$

$$= \min ( 2, 3, 2 )$$

$\Rightarrow$   $D(2,2)$  can be computed  
 using "coming from"  $D(2,1)$  or  $D(1,1)$

	0	1	2	3	4
D	ε	T	G	C	A
0	0	1	2	3	4
1	1	1	2	3	3
2	2	1	2	3	4
3	3	2	2	2	3

Now follow one "path" from  $D(1,1)$  back to  $D(0,0)$  to construct alignment as follow:

$$\begin{aligned} \text{IF } \begin{matrix} (i-1, j) \\ \uparrow \\ (i, j) \end{matrix} &\Rightarrow D(i-1, j) + \delta(s_i, -) \\ &\Rightarrow \left( \dots \begin{matrix} s_i \\ - \end{matrix} \dots \right) \end{aligned}$$

$$\begin{aligned} \text{IF } \begin{matrix} (i-1, j-1) \\ \swarrow \\ (i, j) \end{matrix} &\Rightarrow D(i-1, j-1) + \delta(s_i, t_j) \\ &\Rightarrow \left( \begin{matrix} s_i \\ t_j \end{matrix} \right) \end{aligned}$$

$$\begin{aligned} \text{IF } (i, j-1) \rightarrow (i, j) & \\ \Rightarrow D(i, j-1) + \delta(-, t_j) & \\ \Rightarrow \left( \begin{matrix} - \\ t_j \end{matrix} \right) & \end{aligned}$$

	0	1	2	3	4
D	ε	T	G	C	A
0	0	1	2	3	4
1	1	1	2	3	3
2	2	1	2	3	4
3	3	2	2	2	3

Here 2 possible paths

$\hat{=}$  2 possible opt. alignments.

$$\left( \begin{matrix} T G C A \\ A T C - \end{matrix} \right)$$

$$\left( \begin{matrix} - T G C A \\ A T - C - \end{matrix} \right)$$

Example shows, since no penalty for gaps, we may have a lot of them.

Alg globalAlign also better known as  
Needleman-Wunsch-algorithm  
(1970)

(also globalAlign slightly derivate, since  
original Needleman-Wunsch Alg had runtime  
 $O(|S|^2|T|^2)$   
 $\rightarrow O(|S||T|)$  is runtime of Wegner/Isidor 1974)

---


GAPS:    ACCGTCTGCT            ACCGTCTGCT  
          A-C--C-G-T            ACCGT-----

in evolution, "consecutive seqs" of gaps  
easier to realize than a lot of  
single gaps.

& also to avoid "unnecessary" gaps

as in

  $\begin{pmatrix} T G C A \\ A T C - \end{pmatrix}$

  $\begin{pmatrix} - T G C A \\ A T - C - \end{pmatrix}$

Def: gap penalty function  $g: \mathbb{N} \rightarrow \mathbb{R}$

$$\text{st } g(k+n) \leq g(k) + g(n) \quad (\text{"Sub-additive"})$$

$g(i)$  = penalty for having  $i$  consecutive gaps.

$g(k+n) \leq g(k) + g(n)$  ensures that it is better to have  $k+n$  "consecutive" gaps than one somewhere of length  $k$  and one of length  $n$ .

Exmpl  $g(i) = i + 1$ , unit-cost  $\delta$ .

A C C G T C T G C T  
A - C - - C - G - T  
 $0 + 2 + 0 + 3 + 0 + 2 + 0 + 2 + 0 = 9$

A C C G T C T G C T  
A C C G T - - - - -  
0 +  $g(5) = 0 + 6 = 6$

## Alignment with gap penalties

(as above but modified 'D')

$$D(0,0) = 0, \quad D(0,k) = D(k,0) = g(k), \quad k \geq 1$$

$$D(i,j) = \min \begin{cases} D(i-1, j-1) + \delta(u_i, v_j) \\ \min_{1 \leq k \leq i} D(i-k, j) + g(k) \\ \min_{1 \leq k \leq j} D(i, j-k) + g(k) \end{cases}$$

Problem:  $O(|S||T|(|S|+|T|))$  "Cubic"

Solution: affine gap costs

$$g(k) = \alpha + \beta \cdot k$$

$\alpha$  = cost for "opening" gaps

$\beta$  = costs for extending gaps.

БОТОН (1982),

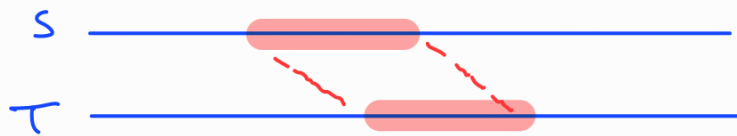
Altschul & Erickson (1986)]

↳  $O(|S||T|)$  - time Alg



So far, what we have done is global Align,  
 that is,  
 Compare entire S with entire T

local: Find most similar pair of  
 substrings of S & T



in practice, matches get pos. weight (high similarity)  
 deletion/insertion neg-weight

in this case we call  $\delta$  scoring fct.

Exmpl

$\delta$	-	A	C	G	T
-	-6	-6	-6	-6	-6
A	-6	2	-4	-4	-4
C	-6	-4	2	-4	-4
G	-6	-4	-4	2	-4
T	-6	-4	-4	-4	2

gaps

matches

mismatches

XX

L-D (local D):

$$L-D(0, i) = L-D(i, 0) = 0 \quad \forall i \geq 0$$

[ but substring of S & E  
(resp E & T) is  
the empty string  $\epsilon$   
= "0" ]

$$L-D(i, j) = \max \begin{cases} L-D(i-1, j) + \delta(s_i, -) \\ L-D(i, j-1) + \delta(-, t_j) \\ L-D(i-1, j-1) + \delta(s_i, t_j) \\ 0 \end{cases} \quad \text{with } \delta = \text{scoring function.}$$

Exmpl:

S = G C C G , S' = A C C A ,  $\delta$  as above

L-D	$\epsilon$	G	C	C	G
$\epsilon$	0	0	0	0	0
A	0	0	0	0	0
C	0	0	2	0	0
C	0	0	2	4	0
A	0	0	0	0	0

max + traceback to first '0'

$\Rightarrow$  get CC as best local align.

( This is the basic idea that together with gap-penalty function is known as Smith-Waterman-Alg )

Another large Example:

		Y															
		ε	T	A	T	A	T	G	C	G	G	C	G	T	T	T	
X	ε	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	G	0	0	0	0	0	0	2	0	2	2	0	2	0	0	0	
	G	0	0	0	0	0	0	2	0	2	4	0	2	0	0	0	
	T	0	2	0	2	0	2	0	0	0	0	0	0	4	2	2	
	A	0	0	4	0	4	0	0	0	0	0	0	0	0	0	0	
	T	0	2	0	6	0	6	0	0	0	0	0	0	0	2	2	2
	G	0	0	0	0	2	0	8	2	2	2	0	2	0	0	0	
	C	0	0	0	0	0	0	2	10	4	0	4	0	0	0	0	
	T	0	2	0	2	0	2	0	4	6	0	0	0	0	2	2	2
	G	0	0	0	0	0	0	4	0	6	8	2	2	0	0	0	
	G	0	0	0	0	0	0	2	0	2	8	4	4	0	0	0	
	C	0	0	0	0	0	0	0	4	0	2	10	4	0	0	0	
	G	0	0	0	0	0	0	2	0	6	2	4	12	6	0	0	
	C	0	0	0	0	0	0	0	4	0	2	4	6	8	2	0	
	T	0	2	0	2	0	2	0	0	0	0	0	0	8	10	4	
	A	0	0	4	0	4	0	0	0	0	0	0	0	0	2	4	6

max.

		Y																
		ε	T	A	T	A	T	G	C	G	G	C	G	T	T	T		
X	ε	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	G	0	0	0	0	0	0	2	0	2	2	0	2	0	0	0		
	G	0	0	0	0	0	0	2	0	2	4	0	2	0	0	0		
	T	0	2	0	2	0	2	0	0	0	0	0	0	0	4	2	2	
	A	0	0	4	0	4	0	0	0	0	0	0	0	0	0	0	0	
	T	0	2	0	6	0	6	0	0	0	0	0	0	0	0	2	2	2
	G	0	0	0	0	2	0	8	2	2	2	0	2	0	0	0	0	
	C	0	0	0	0	0	0	2	10	4	0	4	0	0	0	0	0	
	T	0	2	0	2	0	2	0	4	6	0	0	0	0	2	2	2	
	G	0	0	0	0	0	0	4	0	6	8	2	2	0	0	0	0	
	G	0	0	0	0	0	0	2	0	2	8	4	4	0	0	0	0	
	C	0	0	0	0	0	0	0	4	0	2	10	4	0	0	0	0	
	G	0	0	0	0	0	0	2	0	6	2	4	12	6	0	0	0	
	C	0	0	0	0	0	0	0	4	0	2	4	6	8	2	0	0	
	T	0	2	0	2	0	2	0	0	0	0	0	0	0	8	10	4	
	A	0	0	4	0	4	0	0	0	0	0	0	0	0	2	4	6	

traceback

Y T A T G C - G G C G  
 | | | | | | | |  
 X T A T G C T G G C G

# Protein Comparison & Scoring Functions

(max instead of min in Alg above)

$$\text{LCS} \hat{=} \sigma(x, y) = \begin{cases} 1 & x=y \text{ (match)} \\ -\infty & \text{mismatch (not allowed)} \\ 0 & \text{insert, deletion,} \\ & \text{ie } x=- \text{ or } y=- \end{cases}$$

simple idea:

eg. In Humans:

A ↔ G, C ↔ T substitution

~2 times more frequent than other

A, G ∈ Purines  
C, T ∈ Pyrimidines

& deletion/insertion less frequent than substitution

	-	A	C	G	T
-		-6	-6	-6	-6
A		2	-4	-2	-4
C			2	-4	-2
G				2	-4
T					2

Scores:

(substitution ~1 in 1000)

Indels: ~1 in 3000

(take two human genomes & compare)

Often protein-seq. instead of DNA-seq.  
are compared.

Why?

- AA<sup>G</sup>  
AA<sup>A</sup>  $\leadsto \neq 0$  Align. score

But both AA<sup>G</sup> & AA<sup>A</sup> encode Lysine  
same amino acid.

"Silent mutations"

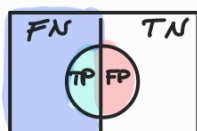
AA<sup>G</sup>  $\leadsto \neq$  Align score but possibly same score  
as AA<sup>G</sup>, AAA (similar)

But AA<sup>G</sup>  $\leadsto$  Lysine  
TAA  $\leadsto$  STOP  $\leadsto$  yields possibly fatal  
protein

comparing protein sequences "more robust"

- due to genetic code, very similar proteins  
might be encoded by pair of DNA  
sequences that share only limited  
similarity

- 20 Aa instead 4 Bases  $\Rightarrow$  higher sensitivity &  
shorter sequences (runtime)



$$\text{sensitivity} = \frac{TP}{FN+TP} = \text{true pos. rate}$$

# How to obtain "good" scoring?

As an example "simple idea",

Observation: Asn, Asp, Glu, Ser most "mutable" amino acid  
Cys & Trp are the "least" mutable

eg. Ser  $\rightarrow$  Phe is  $\sim 3$  times higher than Trp  $\rightarrow$  Phe

Based on this knowledge one can start to design scoring matrices.

Scoring matrix  $\delta$

$\delta(i,j)$  usually reflects how often amino acid  $i$  substitutes amino acid  $j$

if we have large sets of alignments

$\delta(i,j)$  simply counts how many times  $i$  is replaced by  $j$ .

Problem

to get alignments, we need already some scoring function.

highly ( $> 90\%$ )

But for similar sequences, unit-score can be used first & then  $\delta$  can be established & used to construct less "obvious" alignments.

# Multiple sequ. alignments (MSA)

So far 2 sequences only.

## Formal

set of  $k$  sequences  $S = \{S_1 \dots S_k\}$ ,  $S_i \in \Sigma^*$

multiple sequ. alignment  $A(S_1 \dots S_k)$  is  $k \times n$  matrix  
( $n \geq \max\{|S_1| \dots |S_k|\}$ )

It:

- each row consist of characters in  $\Sigma$  & '-' gap symbol
- deletion of gaps in each row  $i$  yields  $S_i$
- in no column all rows are filled with gaps.

$S_1$	A	C	G	-	-	G	A	G	A
$S_2$	-	C	G	T	T	G	A	C	A
$S_3$	A	C	-	T	-	G	A	-	A
$S_4$	C	C	G	T	T	C	A	C	-

possible scores: Sum-of-Pair (SP) score:

$$SP(S_1 \dots S_n) = \sum_{1 \leq i < j \leq n} \delta(S_i, S_j),$$

where  $\delta(S_i, S_j)$  = "classical" align. score between  $S_i$  &  $S_j$

$S_1$     A C G - - G A G A  
 $S_2$     - C G T T G A C A  
 $S_3$     A C - T - G A - A  
 $S_4$     C C G T T C A C -

match = 2  
 mismatch = -2  
 gap open = -2  
 gap extension = -1

$\Rightarrow \delta(S_1, S_2)$ :
 

$S_1$	A	C	G	-	-	G	A	G	A
$S_2$	-	C	G	T	T	G	A	C	A
3									

  
-2 2 2 -2 -1 2 2 -2 2  
= 3

$\delta(S_1, S_3) = 4$   
 $\delta(S_1, S_4) = -5$   
 $\delta(S_2, S_3) = 2$   
 $\delta(S_2, S_4) = 6$   
 $\delta(S_3, S_4) = -6$

$\Rightarrow SP(S_1 \dots S_n) = 3 + 4 - 5 + 2 + 6 - 6 = 4$



Problem: finding one or multiple alignment with opt score is NP-hard.

Plenty of several methods to compute optimal MSA exists.

- exponential runtime algorithm  
(in part with runtime improvements via e.g. simulate annealing, searchspace-pruning (bread & bud))
- approx. alg. (Guaranteed: polytime MSA of solution (distance) is at most twice the opt one)
- heuristics → progressive approaches (eg. ClustalW)  
→ iterative — (eg. MUSCLE)  
(most commonly used)
- probabilistic methods  
(assumption on certain evol. models)  
(eg. markov models)  
profile HMM.

Classical Alg Based or concerned  
or with

Alignments:

---

- Blant
- Clustal
- MUSCLE.

# BLAST = Basic Local Alignment Search Tool

Umbrella term for a collection of the world's most widely used programs for analyzing biological sequence data

- ▶ BLAST is used to compare experimentally determined DNA or protein sequences with sequences already existing in a database.
- ▶ Basic idea: BLAST divides query sequences into short strings and initially only looks for (exact) matches of those strings in database strings. This is afterwards extended to get the entire alignment.
- ▶ very fast local alignment heuristic, but no optimality guarantee
- ▶ output: series of local alignments, i.e. comparisons of pieces of the searched sequence with similar pieces from the database. In addition, BLAST indicates how significant of the hits that have been found.

Databases e.g. for nucleotide sequences (Genbank of NCBI, EMBL, ...) or protein databases (SwissProt, RefSeq, Pfam, ...).

BLAST homepage: [blast.ncbi.nlm.nih.gov](http://blast.ncbi.nlm.nih.gov)

Tutorial: [digitalworldbiology.com/BLAST](http://digitalworldbiology.com/BLAST)

20/24

## BLAST "Types"

type	query	target
blastn	nucleotide	nucleotide
blastp	protein	protein
blastx	nucleotide (transl)	protein
tblastn	protein	nucleotide (transl)
tblastx	nucleotide (transl)	nucleotide (transl)

21/24

## Example

<https://blast.ncbi.nlm.nih.gov/Blast.cgi> → nucleotide blast → copy&paste → press button BLAST

>Sequence\_experimental

```
GACATTACGGCGACCCAGTCTCCCCGGTGTGTGTCAGTGGGACTGGGCC
AGACCGCAACCATCACTTGTACGGCCAGTCAAAGCATCTACAGTAACCT
TGCTTGGTACCAGCAGAGAGAAGGACAGAAGCCCTCTCTCCTGATCTAT
GCTGCGACAACGCGATACGAAGGAGTCTCCGAGCGATTTCAGCGGCAGTG
GATCAGGGACCAGTTTCACCCTGACAATCAGCAACGTTTCAGAATGAGGA
TGTCGCTGACTATTACTGTGATCGCATATTCGATCTACTCCGGTTCC
GTTGTTTTTCGGTGAAGGAACCAAGCTCAGACTGAGCCGT
```

mRNA for some (specified) protein of a nurse shark.

22/24

Clustal is a series of computer programs used in bioinformatics for multiple sequence alignment.

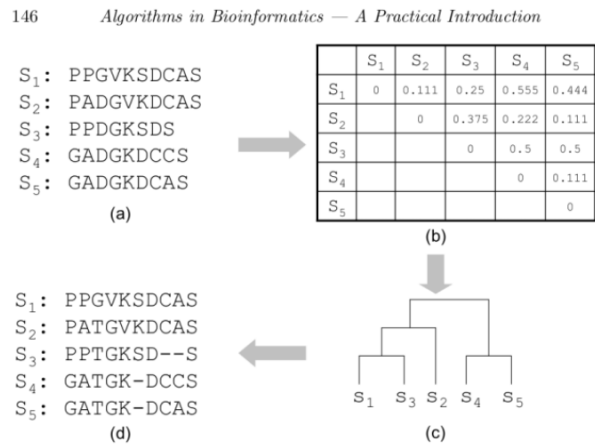
Brief History:

- ▶ Clustal (1981, first version)
- ▶ CLustalW (1994, great improvements)
- ▶ ClustalX (1997, first time with GUI)
- ▶ ClustalΩ (latest standard version, 2011)

## Clustal

Basic idea explained on ClustalW (3 steps for input  $\zeta = \text{set } \{S_1, \dots, S_k\}$  of sequences):

- W1** Compute for all pairs  $S_i, S_j \in \zeta$  a pairwise alignment  $\implies$  pairwise distances  $D(S_i, S_j)$
- W2** Use distance matrix  $D$  to compute phylogenetic tree  $T$  (via NeighborJoining-method)
- W3** Use  $T$  to carry out a multiple alignment



**FIGURE 6.6:** The three steps of ClustalW (a progressive alignment method). Five input sequences are given in (a). Step 1 computes the pairwise distance scores for these five sequences (see (b)). Then, Step 2 generates the guide tree such that similar sequences are grouped together first (see (c)). Step 3 aligns the sequences one by one according to the branching order of the guide tree, yielding the multiple alignment of all input sequences (see (d)).

Cluster W: 3 steps  $W1, W2, W3$

(W1)  $\forall$  pairs  $S_i, S_j$ : compute optimal global alignment.

$\Rightarrow$   $x = \text{nongap pos. } (\bar{e}, e)$   
 $y = \text{\#matches}$

$$d(S_i, S_j) = 1 - \frac{y}{x}$$

PP6VKSDCAS      opt alignment  
PPD6KSD--S

$$x = 8, y = 6 \rightarrow d(S_i, S_j) = 1 - \frac{6}{8} = 0.25.$$

$\rightarrow$  this gives distance matrix.  $D$

(W2) use  $D$  to compute Phyl. tree (Later in course)

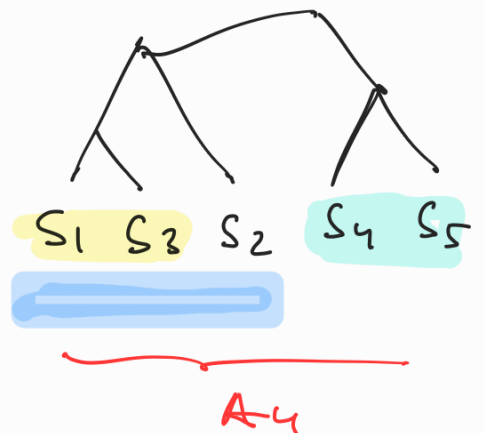
(W3)

align  $(S_1, S_3) \rightarrow A_1$

align  $(S_4, S_5) \rightarrow A_2$

align  $(A_1, S_2) \rightarrow A_3$

align  $(A_3, A_2) \rightarrow A_4.$



RUNTIME  $O(k^2 n^2 + k^3)$

$n = \max_{i=1..k} \{ |S_i| \}$

MUSCLE (multiple seq. comparison by  
log expectation)

## MUSCLE

**M**Ultiple **S**equence **C**omparison by **L**og-**E**xpectation (2004)

computer software used in bioinformatics for multiple sequence alignment.

Online available via <https://www.ebi.ac.uk/Tools/msa/muscle/>

Basic idea for input  $\zeta = \text{set } \{S_1, \dots, S_k\}$  of sequences (2nd and 3rd steps similar to ClustalW):

- 1** Compute  $k$ -mer distances  
=(dis)similarities  $D(S_i, S_j)$  between the sets of  $k$ -mers for all pairs  $S_i, S_j \in \zeta$   
Much(!) faster than [W1] in Clustal
- W2** Use distance matrix  $D$  to compute phylogenetic tree  $T$  (via UPGMA-method)
- W3** Use  $T$  to carry out a multiple alignment
- 4** Several re-iteration and refinement steps follow

for  $S = \{S_1 \dots S_n\}$  set of sequ.  
compute  $k$ -mers "distances"

(not pairw. alignment)

$k$ -mers = contiguous substring of length  $k$ .

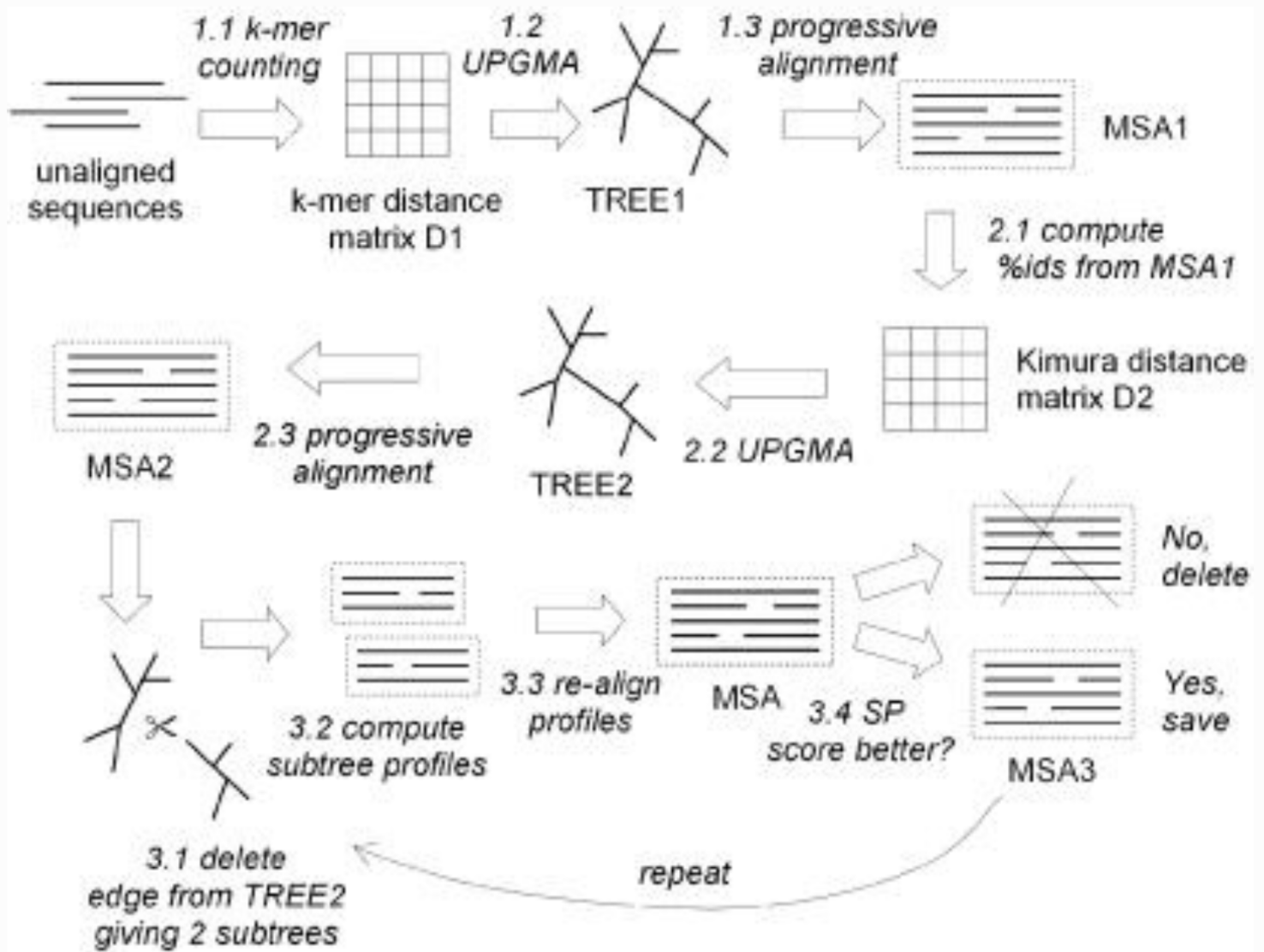
$\Rightarrow$  distance matrix.

$\downarrow$   
w2/w3 as in ClustalW

this step is considerably faster than WA is ClustalW

After this step several "recomputing"  
& "refinement"  
steps follow.

[ClustalW uses neighbor joining ( $O(n^3)$ )  
Muscle UPGMA ( $O(n^2)$ ) to compute trees]



In general more accurate & faster than ClustalW.



