

Multiple Alignments

Laboratory of Bioinformatics I
Module 2

Emidio Capriotti

<http://biofold.org/>



Biomolecules
Folding and
Disease

Department of Pharmacy and
Biotechnology (FaBiT)
University of Bologna



Multiple Structure Alignment

Align the structure of 5 structures of Cytochromes

```
3zcf:A --GDVEKGKKIFIMKCsQCHTVEkgg-----khKTG--PNLHG--LfgRKTgqapgysyt----aank
3o20:A --GDVEKGKKIFVQKCaQCHTVEkgg-----khKTG--PNLHG--LfgRKTgqapgftyt----dank
2ce0:A --LDIQRGATLFNRACaACHDTGg-----nIIQpgATLFTkdL--ERN-----
1cxc:A qeGDPEAGAKAFNQCQ-TCHVIVddsgettiagrnaKTG--PNLYG--VvgRTAgtqadfkgygegmkag
1i8o:A --EDAKAGEAVFKQCM-TCHRADk-----nMVG--PALAG--VvgRKAgtaagftysp-lnhnsq
```

```
3zcf:A nkgiIW-GEDTLMEYLENPKkyi-----pgTKMiFvGiK-----KKEERAD
3o20:A nkgiTW-KEETLMEYLENPKkyi-----pgTKMiFaGiK-----KKTERED
2ce0:A ----GVdTEEEIYRVTYFGK-----GRM-PgF-GekctprgqctfgprlQDEEIKL
1cxc:A akglAW-DEEHFVQYVQDPTkflkeyt-----gdakakGKMtF-KlK-----KEADAHN
1i8o:A eaglVW-TADNIVPYLADPNafllkfltekqkadqavgvTKMtF-KlA-----NEQQRKD
```

```
3zcf:A LIAYLKKATne----
3o20:A LIAYLKKATne----
2ce0:A LAEFVKFQAdqgwpt
1cxc:A IWAYLQQVAvrp---
1i8o:A VVAYLATLK-----
```

Functional sites

Conserved sites

Similar substitutions



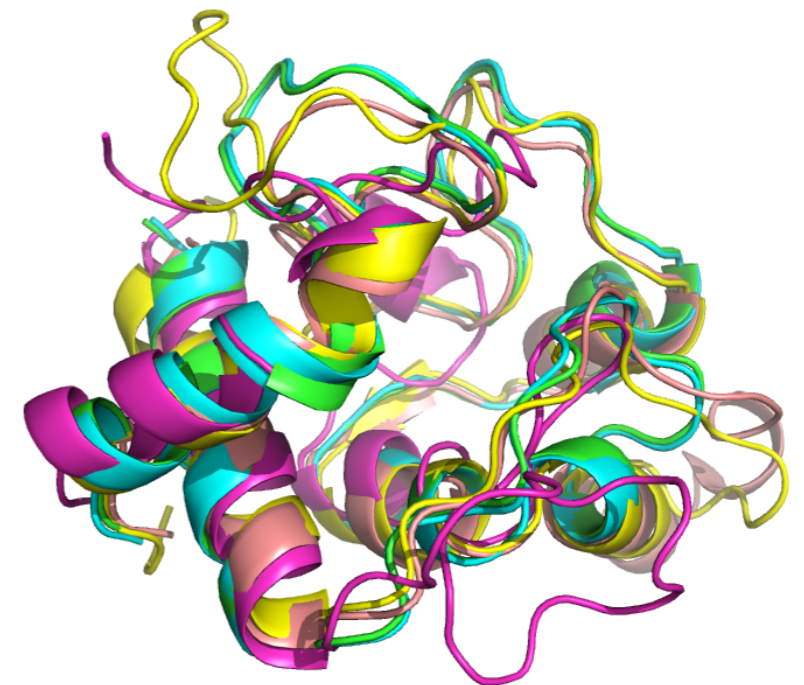
Important Information

The comparison among multiple structures **highlight the most conserved and the most variable sites.**

Conserved sites could be functionally or structurally important.

For each site the residue distribution is estimated

The information is not general, but family specific



Multiple Alignment

A representation of a **set of sequences, where equivalent residues** (e.g. functional, structural) are aligned in columns.

Part of an alignment of SH2 domains from 14 sequences

```

          *      *      :      * **:*      *      :      :      :
lnk_rat   -----Y P W F H G P I S R V R A A Q L V Q L Q G P D A H G V F L V R Q S E S R R - G E Y V L T F N L Q -----G R A K H L R L V L T E R G Q C R V Q H --L H F P S V V D M L
crk1_mouse -----S A W Y M G P V T R Q E A Q T R L Q G Q R ---H G M F L V R D S S T C P - G D Y V L S V S E N -----S R V S H Y I I N S L P N R R F K I G D --Q E F D H L P A L L
nck_human -----W Y Y G K V T R H Q A E M A L N E R G H --E G D F L I R D S E S S P - N D F S V S L K A Q -----G K N K H F K V Q L K - E T V Y C I G Q --R K F S T M E E L V
ht16_hydat -----W Y H G K I T R E V A V Q V L L R K G G R - D G F F L I R D C G N A P - E D Y V L S M M F R -----S Q I L H F Q I N C L G D N K F S I D N G - P I F Q G L D M L I
pip5_human -----K P W Y Y D S L S R G E A E D M L M R I P R --D G A F L I R K R E G S --D S Y A I T F R A R -----G K V K H C R I N R D G - R H F V L G T S - A Y F E S L V E L V
fer_human  -----W Y H G A I P R I E A Q E L L K K -----Q G D F L V R E S H G K P - G E Y V L S V Y S D -----G Q R R H F I I Q Y V - D N M Y R F E G --T G F S N I P Q L I
1ab2      -----E E W F H G V L P R E E V V R L L N N -----D G D F L V R E T I R N E E S Q I V L S V C W -----N G H K H F I V Q T T G E G N F R F E G --P P F A S I Q E L I
1mil      -----H S W Y H G P V S R N A A E Y L L S S G I ---N G S F L V R E S E S S P - G Q R S I S L R Y E -----G R V Y H Y R I N T A S D G K L Y V S S E - S R F N T L A E L V
1blj      -----E P W F H G K L S R R E A E A L L Q L -----N G D F L V R E S T T T P - G Q Y V L T G L Q S -----G Q P K H L L L V D P - E G V V R T K D --H R F E S V S H L I
1shd      G S V A P V E T L E V E K W F F R T I S R K D A E R Q L L A P M N K - A G S F L I R E S E S N K - G A F S L S V K D I T T Q - G E V V K H Y K I R S L D N G G Y Y I S P R - I T F P T L Q A L V
1lkkA     -----S I Q A E E W Y F G K I T R R E S E R L L L N A E N P - R G T F L V R E S E A -----Y C L S V S D F D N A K G L N V K H Y K I R K L D S G G F Y I T S R - T Q F N S L Q Q L V
1csy      -----L E P E P W F F K N L S R K D A E R Q L L A P G N T - H G S F L I R E S E S T A - G S F S L S V R D F D Q N Q G E V V K H Y K I R N L D N G G F Y I S P R - I T F P G L H E L V
1bfi      -----S H E K M P W F H G K I S R E E S E Q I V L I G S K T - N G K F L I R A R D N N --G S Y A L C L L H E -----G K V L H Y R I D K D K T G K L S I P E G - K K F D T L W Q L V
1gri      -----H H D E K T W N V G S S N R N K A E N L L R G K R ---D G T F L V R E S S K Q --G C Y A C S V V V D -----G E V K H C V I N K T A T G - Y G F A E P Y N L Y S S L K E L V
          -----E M K P H P W F F G K I P R A K A E E M L S K Q R H --D G A F L I R E S E S A P - G D F S L S V K F G -----N D V Q H F K V L R D G A G K Y F L W V --V K F N S L N E L V

```

* conserved identical residues
: conserved similar residues

Sequence Logo

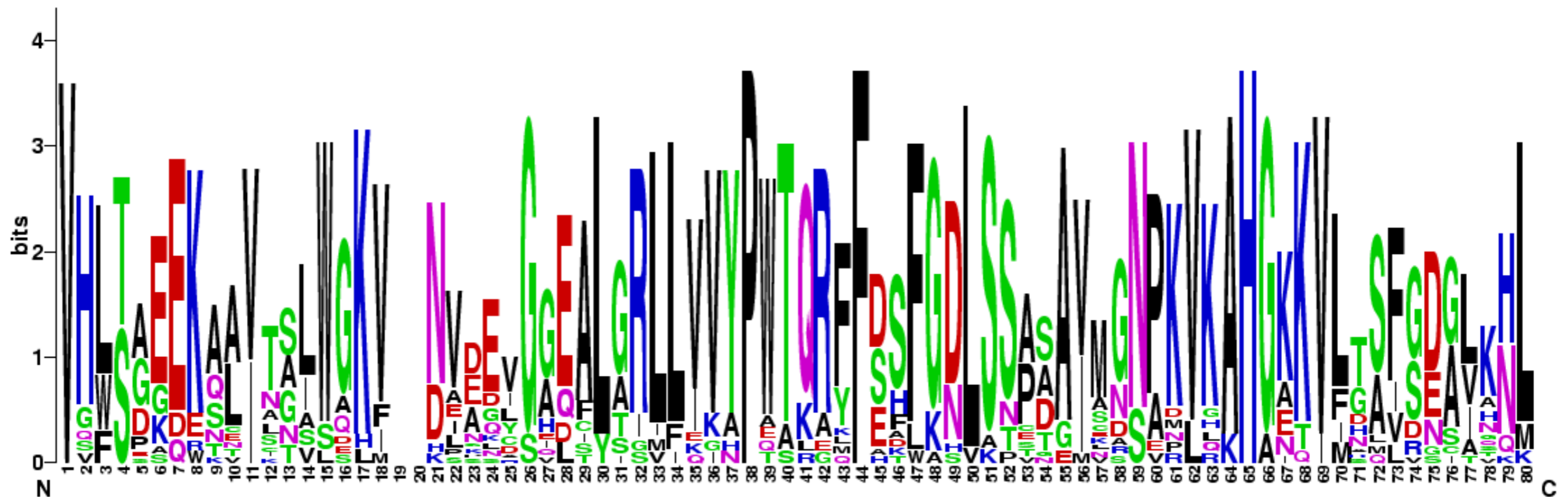
Plot drawn using a score derived from the information entropy (S) calculating the information content (I) in each position

$$S(p) = \sum_{i=1}^{20} -p_i \log_2 p_i$$

$$I = \log_2 20 - S(p) - e_n$$

$$h_i = p_i \times I \quad h_i = \text{height } aa_i$$

$$e_n = \frac{k-1}{2n \times \ln 2} \quad \begin{array}{l} k = \# \text{ aa types} \\ n = \# \text{ sequences} \end{array}$$



weblgo.berkeley.edu

Alignment Scoring

How to score an alignment of many sequences?

Given M sequences A_i , we can define a score for the multiple sequence alignment as the **sum of the scores of all the induced pair alignments**

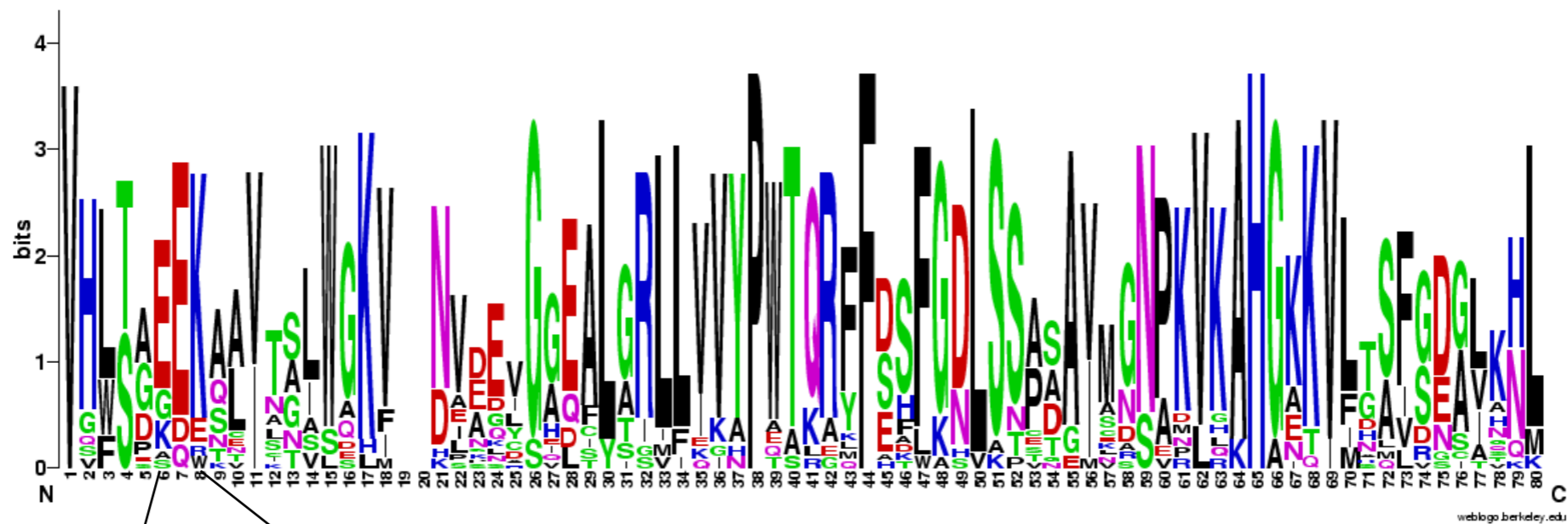
$$S = \sum_{i < j} S(A_i, A_j)$$

$$S \begin{bmatrix} 1 > \text{ASPTLPLSLA} \\ 2 > \text{SS-TLPA--A} \\ 3 > \text{SSPTLPA--A} \end{bmatrix} = S \begin{bmatrix} 1 > \text{ASPTLPLSLA} \\ 2 > \text{SS-TLPA--A} \end{bmatrix} + S \begin{bmatrix} 1 > \text{ASPTLPLSLA} \\ 3 > \text{SSPTLPA--A} \end{bmatrix} + S \begin{bmatrix} 2 > \text{SS-TLPA--A} \\ 3 > \text{SSPTLPA--A} \end{bmatrix}$$

Entropy Score

The multiple sequence alignment can be obtained minimizing the entropy

$$S = \sum_{j=1}^{N_{\text{columns}}} \sum_{i=1}^{20} -p_{ji} \log_2 p_{ji}$$



The substitution score may depend on the position.

Profile-Based Alignment

Given the position i along a sequence profile, it is represented by a 20-element vector $P_i = P_i(A) P_i(C) \dots P_i(Y)$

A	0
C	85
D	0
E	0
F	5
G	0
H	0
I	0
K	0
L	2
M	0
N	8
P	0
Q	0
R	0
S	0
T	0
V	0
W	0
Y	0

Given the residue in position j along the sequence to align: S_j
The score for aligning S_j to the vector P_i is:

$$Score(i, j) = \sum_{k=1}^{20} P_i(r_k) \cdot M(r_k, s_j)$$

where M is a matrix score (BLOSUM or PAM)

The score can be used in dynamic programming procedures (Needleman-Wunsch, Smith-Waterman)

Sequence to Profile Score

Alignment score between P_i and S_i is

$$= 0.85 \cdot M(C,C) + 0.05 \cdot M(C,F) + 0.02 \cdot M(C,L) + 0.08 \cdot M(C,N) =$$

$$= 0.85 \cdot (9) + 0.05 \cdot (-2) + 0.02 \cdot (-1) + 0.08 \cdot (-3) = 7.29$$

P_i	
A	0
C	85
D	0
E	0
F	5
G	0
H	0
I	0
K	0
L	2
M	0
N	8
P	0
Q	0
R	0
S	0
T	0
V	0
W	0
Y	0

$S_i = \text{"C"}$

$$Score(i, j) = \sum_{k=1}^{20} P_i(r_k) \cdot M(r_k, s_j)$$

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

Alignment Strategies

Multiple sequence alignment (MSA)

The algorithmic problem is to find the alignment with the maximum score

Exact algorithms

Algorithms based of **multi-dimensional dynamic programming** have been implemented. However they are **too slow** when many sequences have to be compared.

Progressive alignments

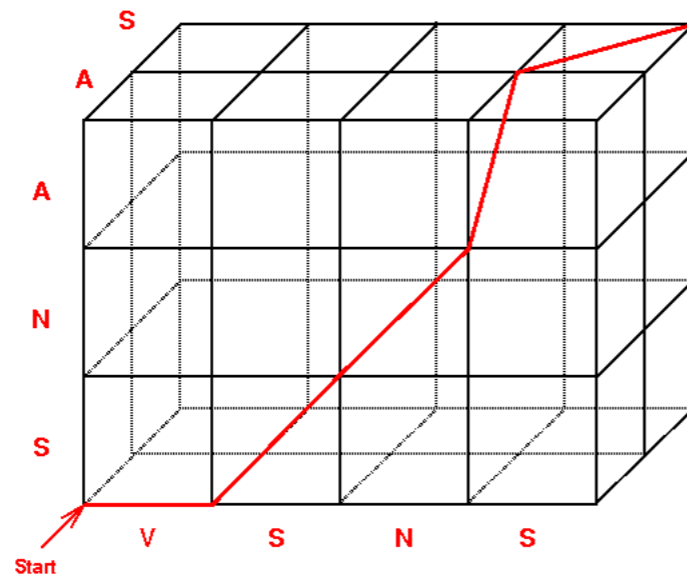
Iterative algorithms

Consistency-based algorithms

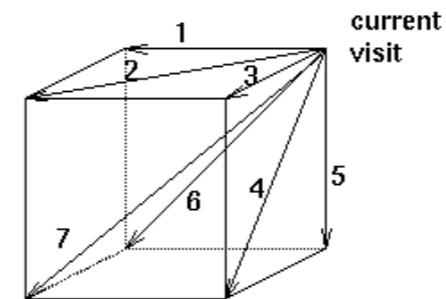
Optimal Alignment

Optimal Multiple Alignment: MSA (Lipman et al. 1989, Gupta et al. 1995)

Extension of dynamic programming for 2 sequences => N dimensions



V S N - S
- S N A -
- - - A S



Problem: calculation time and memory requirements

Time proportional to N^k for k sequences of length N => limited to less than 10 sequences

Progressive MSA

Idea: Progressively **align pairs** of sequences (or groups of sequences)

Problem:

Start with which sequences? How to decide order of alignment?

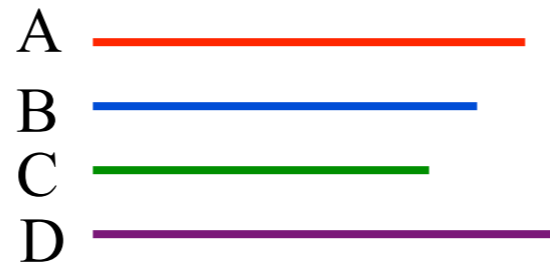
First align the most closely related sequences

How to measure the similarity of the sequences?

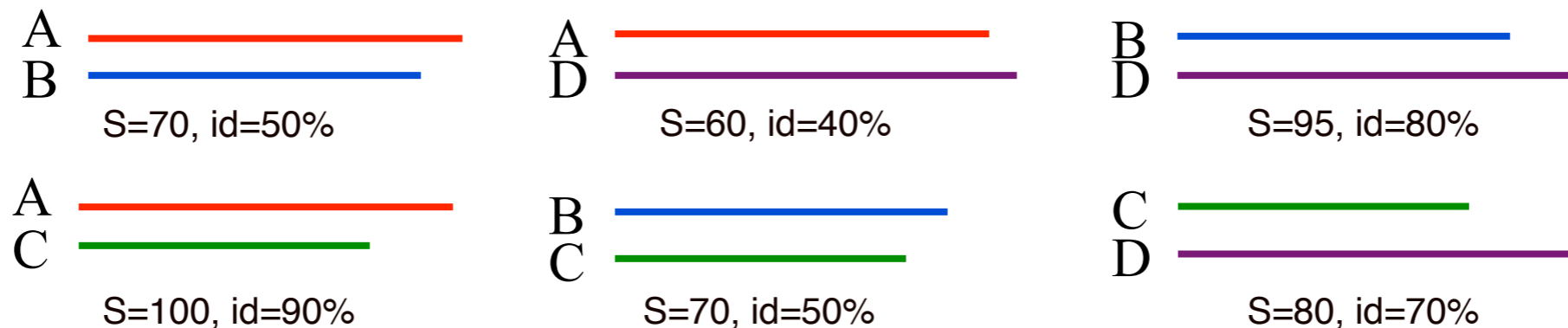
align all the sequences pairwise

calculate the similarity between each pair from the alignment

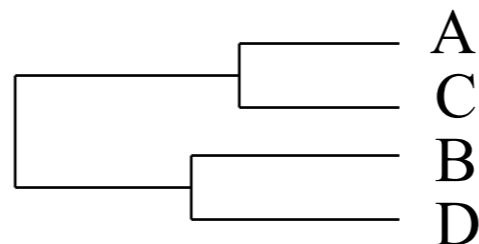
Progressive MSA - Start



Step 1: Pairwise sequence alignment: exact, all-against-all



Step 2: Build a similarity tree



Progressive MSA - End

Step 3: Exact alignment of the most similar sequences, following the tree



Step 4: Build the profile from the sub alignments

Step 5: Perform profile-to-profile alignment following the similarity tree, until comprising all the sequences



Profile-Profile Alignment

The position i along the first sequence profile, it is represented by a 20-element vector
 $P^1_i = P^1_i(A) P^1_i(C) \dots\dots P^1_i(Y)$

The position j along the second sequence profile, it is represented by a 20-element vector
 $P^2_j = P^2_j(A) P^2_j(C) \dots\dots P^2_j(Y)$

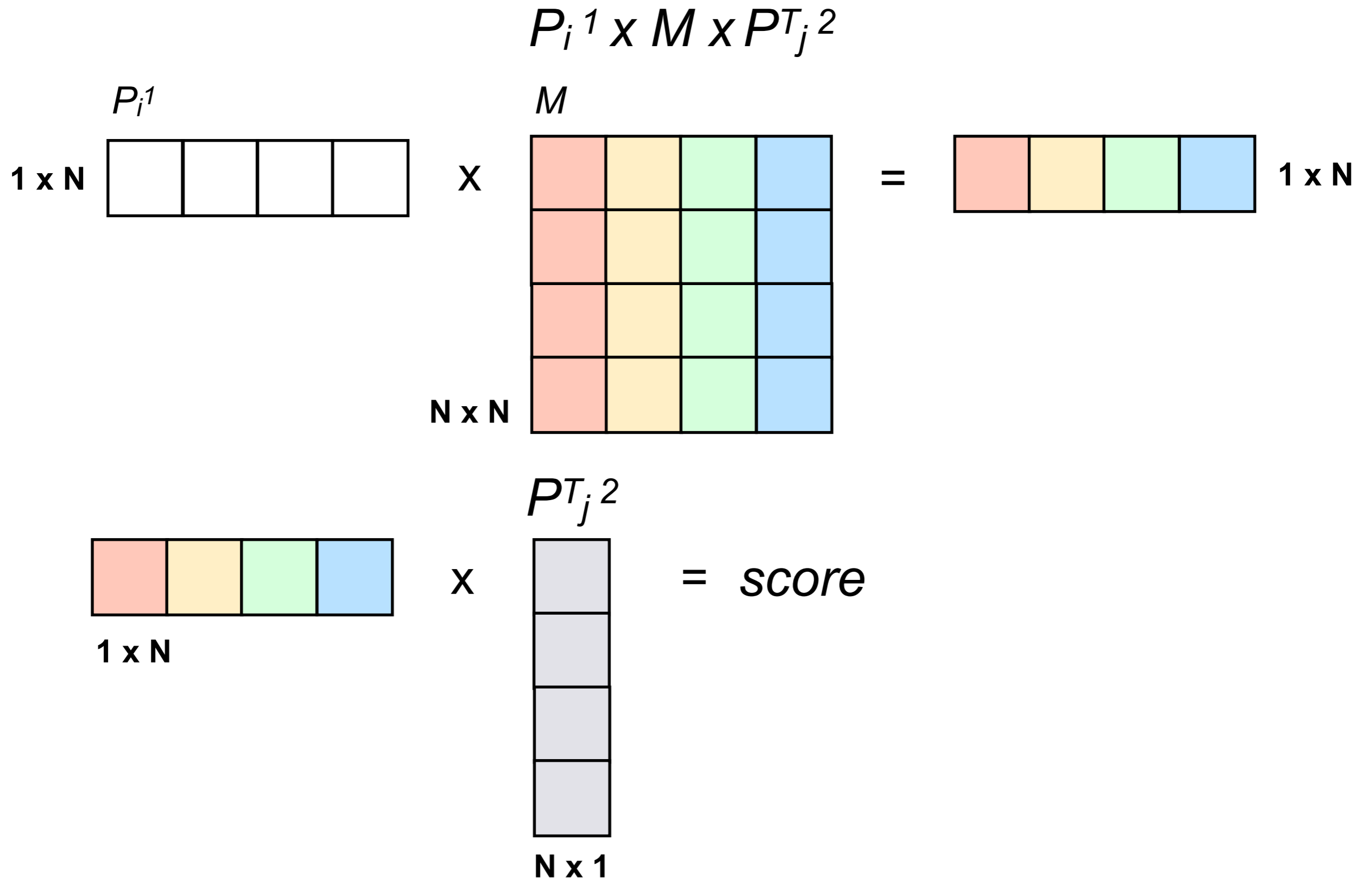
The score for aligning the two positions is:

$$Score(i, j) = \sum_{m=1}^{20} \sum_{k=1}^{20} P^1_i(r_m) P^2_j(r_k) \cdot M(r_m, r_k)$$

where M is a matrix score (BLOSUM or PAM)

The score can be used in dynamic programming procedures
(Needleman-Wunsch, Smith-Waterman)

Matrix representation



Adding Gaps

- Where gaps are added is a critical question
- Gaps are often added to the **first two (closest) sequences**
- To **change the initial gap** choices later on corresponds to give more **weight to distantly related sequences**
- To maintain the **initial gap choices** means that the initial gaps are the most believable

Limitations

- Dependence of the final MSA on the initial pairwise sequence alignment with the highest score
- Errors in initial alignments are propagated
- Gaps can proliferate, if not carefully evaluated
- Gaps can be amino-acid specific, so that you penalize introduction of gaps into segments that are less likely to have gaps (e.g. hydrophobic core)

Alignment Evaluation

How many conserved sites?

CLUSTAL 2.1 multiple sequence alignment

```
sp|P99999|CYC_HUMAN      -----MGDVEKGKKIFIMKCS-----QCHT  20
sp|P00004|CYC_HORSE     -----MGDVEKGKKIFVQKCA-----QCHT  20
sp|P0C0X8|CYC2_RHOSH    -----QEGDPEAGAKAFNQCQTCHVIVDDSGT  27
sp|P00091|CYC22_RHOPA   -----MVKKLLTILSIAATAGSLSIGTASAQDAKAGEAVFKQCMT  40
sp|Q93VA3|CYC6_ARATH_   MRLVLSGASSFTSNLFCSSQQVNGRGKELKNPISLNHNKDLDFLLKKLAP  50
                               * . . . .

sp|P99999|CYC_HUMAN      VEKGGKHKTGPNLHG--LFGRKTGQAPGYS-YTAANKN---KGI IWGEDT  64
sp|P00004|CYC_HORSE     VEKGGKHKTGPNLHG--LFGRKTGQAPGFT-YTDANKN---KGITWKEET  64
sp|P0C0X8|CYC2_RHOSH    TIAGRNAKTGPNLYG--VVGRTAGTQADFKGYGEGMKEAGAKGLAWDEEH  75
sp|P00091|CYC22_RHOPA   CHRADKNMVG PALGG--VVG RKAGTAAGFT-YSP LNHN SGEAGLVWTADN  87
sp|Q93VA3|CYC6_ARATH_   PLTAVLLAVSPICFPPESLGQTLDIQRGATLFNRACIGCHDTGGNIIQPG  100
                               . . . * . * : . . . :
                               . . . * . * : . . . :

sp|P99999|CYC_HUMAN      LMEYLENP-----KKYIPG-----TKMIFVGI  86
sp|P00004|CYC_HORSE     LMEYLENP-----KKYIPG-----TKMIFAGI  86
sp|P0C0X8|CYC2_RHOSH    FVQYVQDPTK-----FLKEYTGD-----AKAKGKMTFK-L  104
sp|P00091|CYC22_RHOPA   IINYLN DPNA-----FLKKFLTDK GKADQAVGVTKMTFK-L  122
sp|Q93VA3|CYC6_ARATH_   ATLFTKDLERNGV DTEEEIYRVTYFGKGRMPGFG--EKCTPRGQCTFGPR  148
                               : : : * . : *

sp|P99999|CYC_HUMAN      KKKEERADLIAYLKKATNE-----  105
sp|P00004|CYC_HORSE     KKKTEREDLIAYLKKATNE-----  105
sp|P0C0X8|CYC2_RHOSH    KKEADAHNIWAYLQQVAVRP-----  124
sp|P00091|CYC22_RHOPA   ANEQQRKDVVAYLATLK-----  139
sp|Q93VA3|CYC6_ARATH_   LQDEEIKLLAEFVKFQADQGWPTVSTD  175
                               : . : : : :
```

Alternative Alignment

How many conserved sites?

CLUSTAL 2.1 multiple sequence alignment

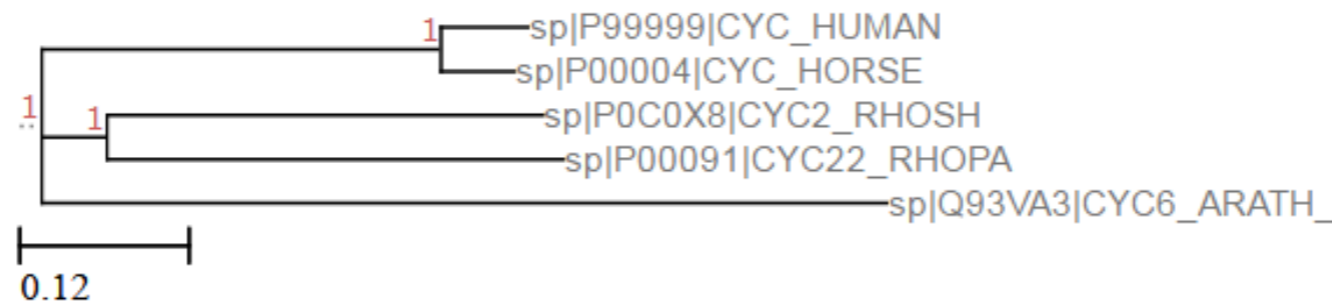
```
sp|P99999|CYC_HUMAN_Cytochrome      -MGDVEKGKKIFIMKCSQCHTVE-----KGGKHKTGPNLHGLFGRKTG 42
sp|P00004|CYC_HORSE_Cytochrome      -MGDVEKGKKIFVQKCAQCHTVE-----KGGKHKTGPNLHGLFGRKTG 42
sp|P0C0X8|CYC2_RHOSH_Cytochrom      QEGDPEAGAKAFN-QCQTCHVIVDDSGTTIAGRNAKTGPNLYGVVGRTAG 49
sp|P00091|CYC22_RHOPA_Cytochro     --QDAKAGEAVFK-QCMTCHR-----ADKN-MVGPALGGVVGRKAG 37
sp|Q93VA3|CYC6_ARATH_Cytochrom     QTLDIQRGATLFRACIGCHDTG-----GNIIQPGATLFTKDLERNG 42
          * : * * * ** . * . * . *

sp|P99999|CYC_HUMAN_Cytochrome      QAPGYS-YTAANKN---KGIIWGEDTLMEYLENPKKYIP----- 77
sp|P00004|CYC_HORSE_Cytochrome      QAPGFT-YTDANKN---KGITWKEETLMEYLENPKKYIP----- 77
sp|P0C0X8|CYC2_RHOSH_Cytochrom      TQADFKGYGEGMKEAGAKGLAWDEEHFVQYVQDPTKFLKEYTGD----AK 95
sp|P00091|CYC22_RHOPA_Cytochro     TAAGFT-YSPLNHNHSGEAGLVWTADNIINYLNPNFLKKFLTDK GKADQ 86
sp|Q93VA3|CYC6_ARATH_Cytochrom     VDTEEEIYRVTYFGKG-RMPGFGEKCTPRGQCTFGPRLQ----- 80
          . * : . . :

sp|P99999|CYC_HUMAN_Cytochrome      --GTKMIFVGIKKKEERADLIAYLKKATNE- 105
sp|P00004|CYC_HORSE_Cytochrome      --GTKMIFAGIKKKTEREDLIAYLKKATNE- 105
sp|P0C0X8|CYC2_RHOSH_Cytochrom      AKG--KMTFKLKKKEADAHNIWAYLQQVAVRP 124
sp|P00091|CYC22_RHOPA_Cytochro     AVGVTKMTFKLANEQQRKDVVAYLATLK--- 114
sp|Q93VA3|CYC6_ARATH_Cytochrom     -----DEEIKLLAEFVKFQADQGWPTVSTD- 105
```

Improve the Alignment

The alignment is based on a guide tree computed on the basis of the pairwise distances (guide tree).



The sequence distances computed starting from the MSA can be different (“phylogenetic” tree)



Phylogenetic Tree

If the **trees are very different**, the final MSA is somehow incoherent with respect of the procedure used to derive it.

It is then possible to iterate the progressive alignment procedure, **using the “phylogenetic” tree as guide**.

Iterative Alignment Method

Iterative Methods: MUSCLE

MUSCLE (MULTiple Sequence Comparison by Log Expectation), 3 steps:

draft progressive:

consists of a progressive sequence alignment

- I (accuracy) it uses log-expectation score instead of PPS score in profile-profile alignment;
- I (efficiency) uses k-mer distance instead of alignment score for sequence similarity (a k-mer is a substring of length k)
- I instead of neighbour joining, it uses UPGMA (Unweighted Pair Group Method with Arithmetic Mean)

improved progressive:

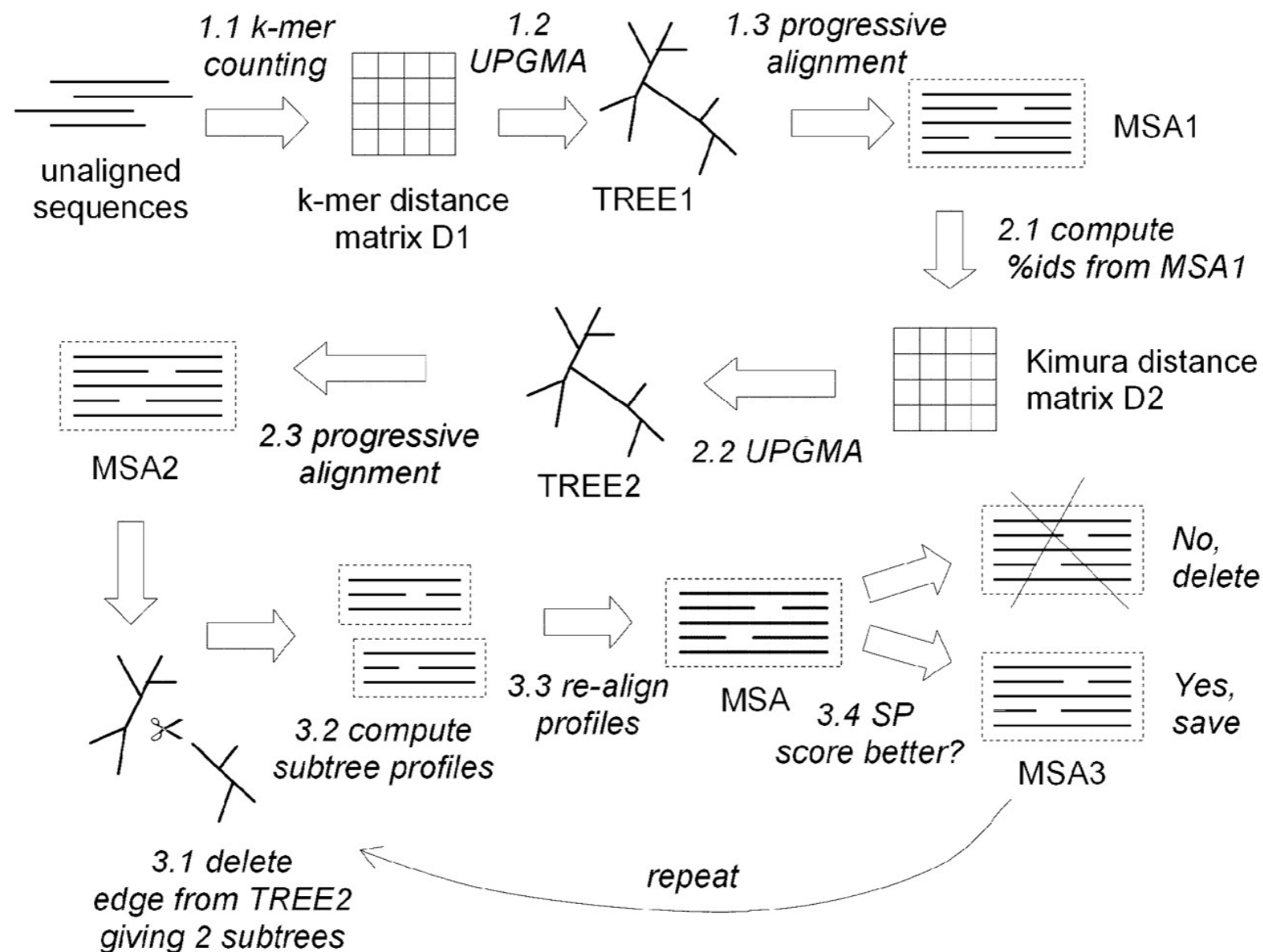
- use alignment to compute more accurate pairwise distance between sequences, Kimura distance: $-\ln(1 - D - D^2/5)$, where D is the fraction of identical bases between the pair of sequences.
- from new distance matrix, build the guide tree and a new alignment.

refinement: tries to improve alignment

refines multiple alignment using the tree-dependent restricted partition technique - a process of deleting edges of guide tree, and re-combine the alignment of the disjoint trees, if better.

Muscle

The first guide tree is not based on pairwise alignment, but on the comparison between the vectors containing the k-mer compositions of each sequence (faster)



Alignment with Muscle

Calculate the alignment of five sequences of the Cytochromes used before with Muscle

```

sp|Q93VA3|CYC6_ARATH      MRLVLSGASSEFTSNLFCSSQQVNGRKGKELKNPISLNHNKDLDFLLKKLAPPLTAVLLAVS
sp|P99999|CYC_HUMAN      -----MG-----
sp|P00004|CYC_HORSE      -----MG-----
sp|P0C0X8|CYC2_RHOSH     -----QEG-----
sp|P00091|CYC22_RHOPA    -----MVKKLLTILSIAATAGSLSIGTASAQ-----
                                     *
```

```

sp|Q93VA3|CYC6_ARATH      PICFPPEESLGQTLDIQRGATLFNRACIGCH-----DTGGNI-----
sp|P99999|CYC_HUMAN      -----DVEKGKKIFIMKCSQCH-----TVEKGGKHKTGPNLHGLFGRKTG
sp|P00004|CYC_HORSE      -----DVEKGKKIFVQKCAQCH-----TVEKGGKHKTGPNLHGLFGRKTG
sp|P0C0X8|CYC2_RHOSH     -----DPEAGAKAFNQ-CQTCHVIVDDSGTTIAGRNAKTGPNLYGVVGRTAG
sp|P00091|CYC22_RHOPA    -----DAKAGEAVFKQ-CMTCH-----RADKNMVGPPALGGVVGRKAG
                                     * : * * * ** . * :
```

```

sp|Q93VA3|CYC6_ARATH      IQPGATLFTKDLER---NGVDTEEEIYRVTYFGKGRM-----PGFGEKCTPRGQCTF
sp|P99999|CYC_HUMAN      QAPGYS-YTAANKN---KGIIWGEDTL-MEYLENPKKYI-----PG-----TKMIF
sp|P00004|CYC_HORSE      QAPGFT-YTDANKN---KGITWKEETL-MEYLENPKKYI-----PG-----TKMIF
sp|P0C0X8|CYC2_RHOSH     TQADFKGYGEGMKEAGAKGLAWDEEHF-VQYVQDPTKFL-----KEYTGDAKAKGKMTF
sp|P00091|CYC22_RHOPA    TAAGFT-YSPLNHNSGEAGLVWTADNI-INYLNDPNAFLKKFLTDKGKADQAVGVTKMTF
                                     . . . : * : : * . . : *
```

```

sp|Q93VA3|CYC6_ARATH      -GPRLQDEEIKLLAEFVKFQADQGWPVSTD
sp|P99999|CYC_HUMAN      VGIKKKEERADLIAYLKKATNE-----
sp|P00004|CYC_HORSE      AGIKKKTEREDLIAYLKKATNE-----
sp|P0C0X8|CYC2_RHOSH     -KLKKEADAHNIWAYLQQVAVRP-----
sp|P00091|CYC22_RHOPA    -KLANEQQRKDVVAYLATLK-----
                                     : : . : * :
```

Consistency

For any multiple alignment, the induced pairwise alignments are necessarily consistent;

given a multiple alignment containing **three sequences x , y , and z** , if position **x_i aligns with position z_k** in the projected x - z alignment and position **z_k aligns with y_j** in the projected z - y alignments, then **x_i must align with y_j** in the projected x - y alignment.

Consistency-based techniques apply this principle in reverse, using evidence from intermediate sequences to guide the pairwise alignment of x and y , such as needed **during the steps of a progressive alignment**.

Transitive Relation

In mathematics, a binary relation R over a set X is transitive if whenever an element a is related to an element b , and b is in turn related to an element c , then a is also related to c .

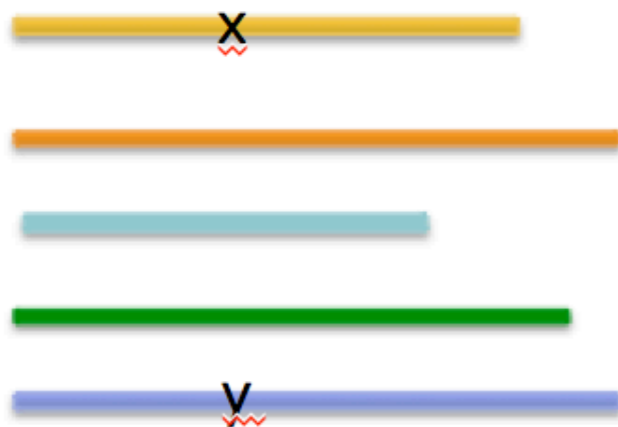
$$\forall a, b, c \in X : (aRb \wedge bRc) \Rightarrow aRc$$

Transitivity in Alignments

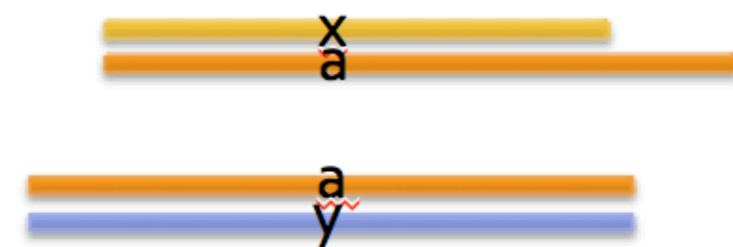
$$\forall a, b, c \in X : (aRb \wedge bRc) \Rightarrow aRc$$

$$\forall x, y, z \in \text{aligned} : (xA \ln z \wedge zA \ln y) \Rightarrow xA \ln y$$

Multiple Sequence Alignment

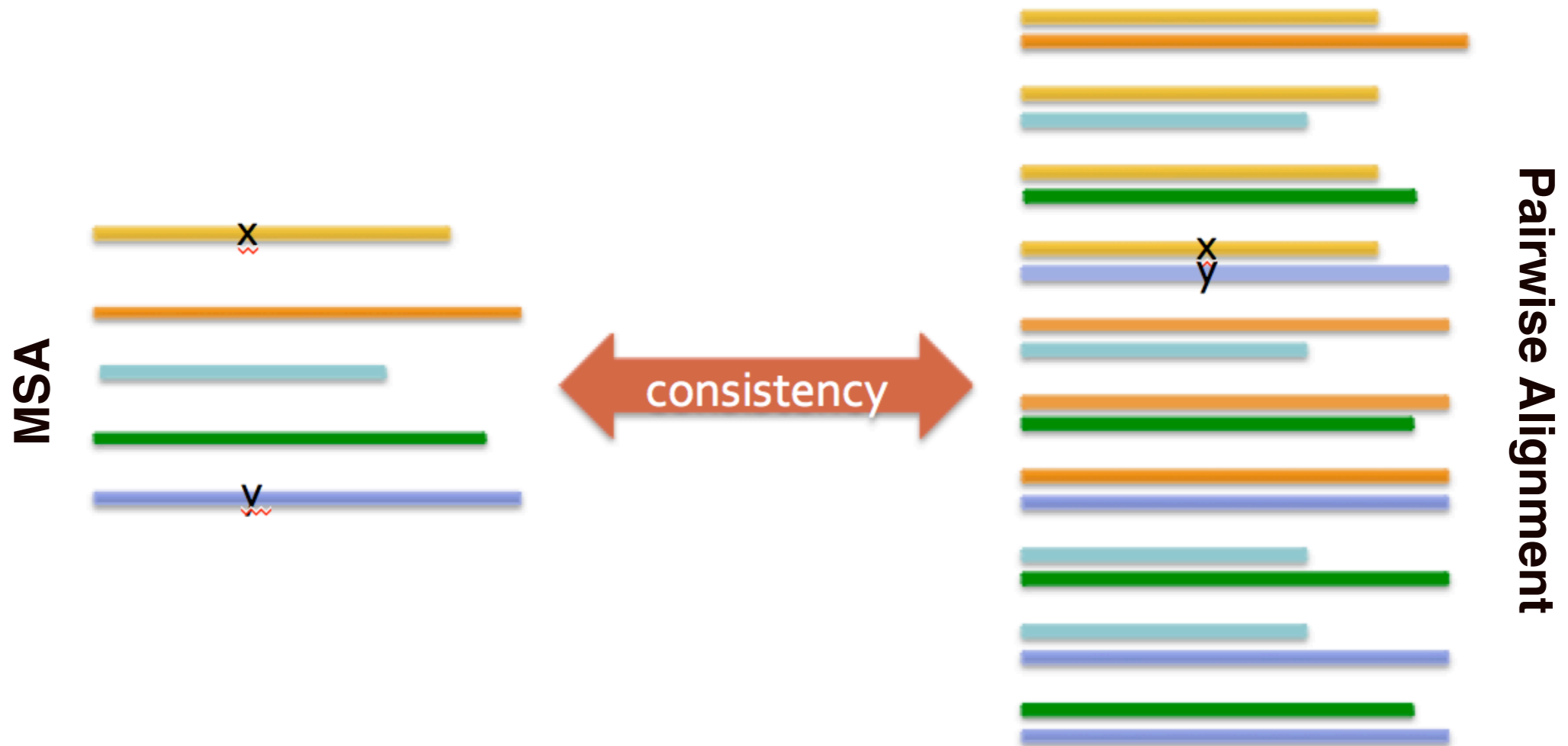


Pairwise Alignment



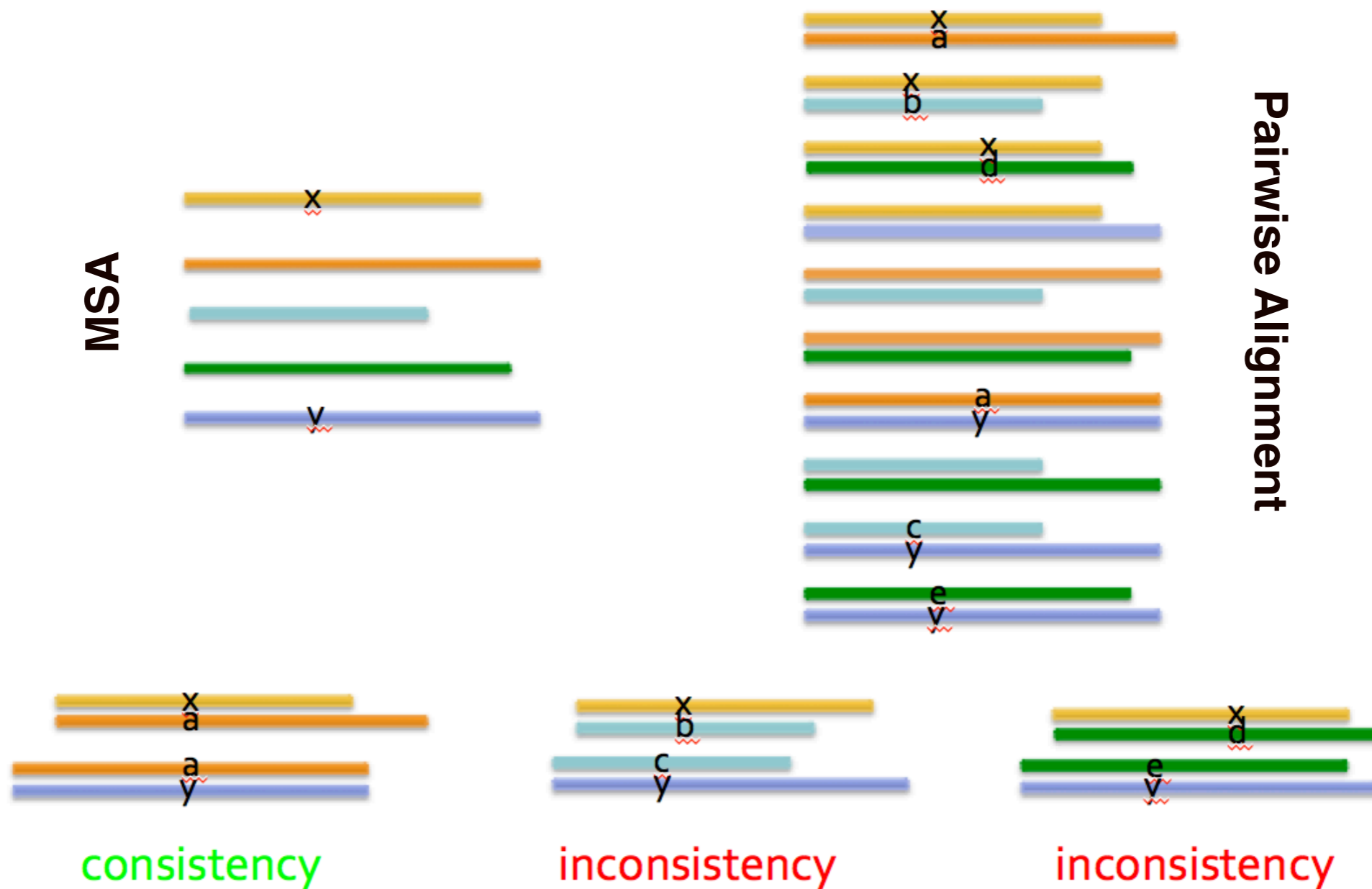
How it can be applied?

Consistency between MSA & pairwise alignment : 0/1
How can we increase the resolution of confidence?



Consistent Alignments

The information are in the pairwise alignments



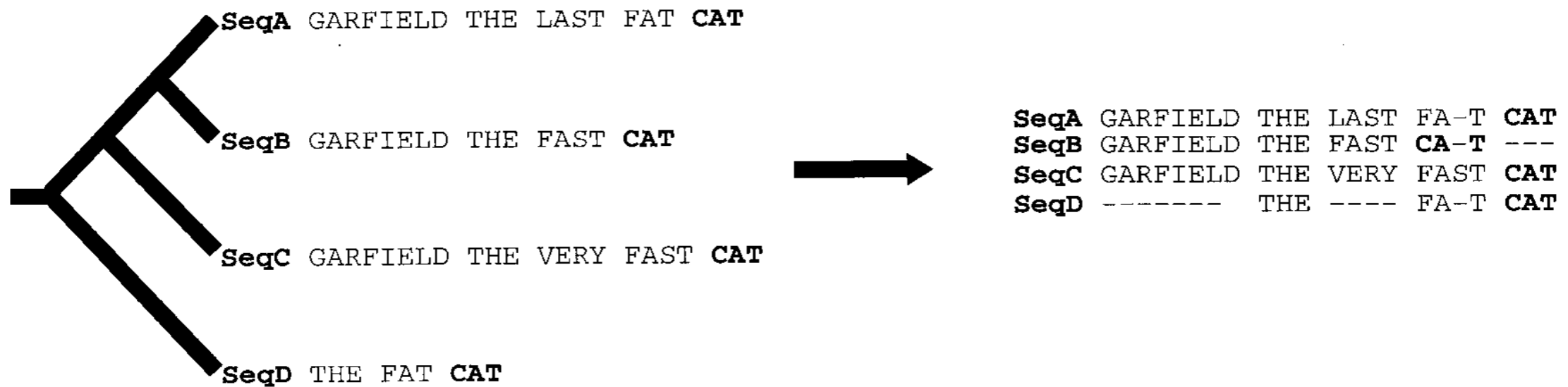
How to Improve?

- MSA from **progressive alignments can be largely inconsistent** with respect to the set of pairwise alignments used to build the guide tree
- Consistency-based methods try **to build the tree in a more consistent way**

T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation

SeqA GARFIELD THE LAST FAT CAT	Prim. Weight = 88	SeqB GARFIELD THE ---- FAST CAT	Prim Weight = 100
SeqB GARFIELD THE FAST CAT ---		SeqC GARFIELD THE VERY FAST CAT	
SeqA GARFIELD THE LAST FA-T CAT	Prim. Weight = 77	SeqB GARFIELD THE FAST CAT	Prim. Weight = 100
SeqC GARFIELD THE VERY FAST CAT		SeqD ----- THE FA-T CAT	
SeqA GARFIELD THE LAST FAT CAT	Prim. Weight =100	SeqC GARFIELD THE VERY FAST CAT	Prim. Weight = 100
SeqD ----- THE ---- FAT CAT		SeqD ----- THE ---- FA-T CAT	



This would be the ClustalW alignment of the four sequences.

CAT is evidently misaligned

T-Coffee - Start

The T-Coffee strategy starts from pairwise alignments as well.

Each pair of aligned residues is associated with a weight equal to the average identity among matched residues (gapped positions are neglected).

Identity values are used instead of alignment scores.

SeqA GARFIELD THE **LAST** **FAT** CAT **Prim. Weight = 88**
SeqB GARFIELD THE **FAST** **CAT** ---

SeqA GARFIELD THE **LAST** FA-T CAT **Prim. Weight = 77**
SeqC GARFIELD THE **VERY** FAST CAT

SeqA GARFIELD THE LAST FAT CAT **Prim. Weight = 100**
SeqD ----- THE ---- FAT CAT

SeqB GARFIELD THE ---- FAST CAT **Prim Weight = 100**
SeqC GARFIELD THE VERY FAST CAT

SeqB GARFIELD THE FAST CAT **Prim. Weight = 100**
SeqD ----- THE FA-T CAT

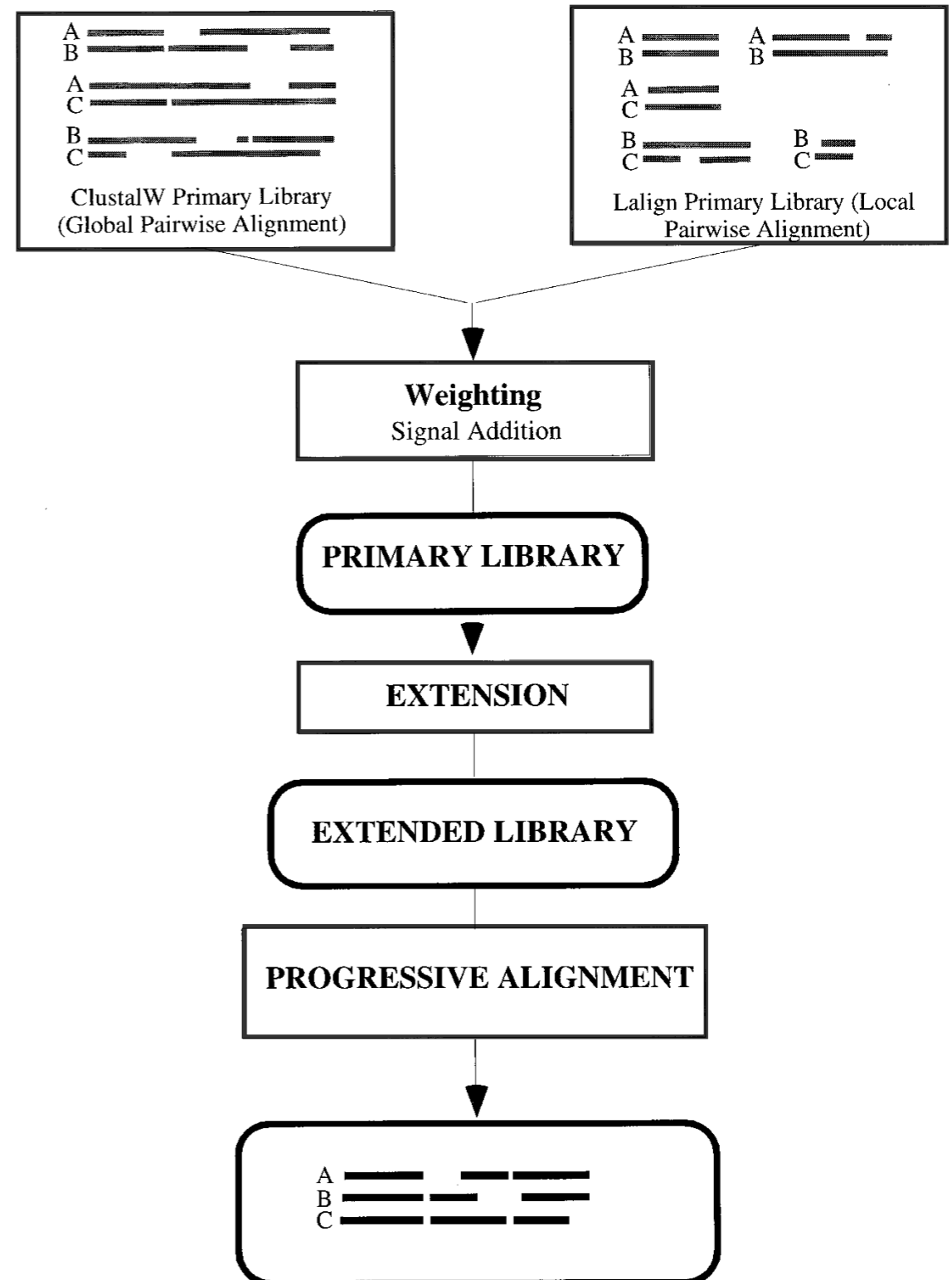
SeqC GARFIELD THE VERY FAST CAT **Prim. Weight = 100**
SeqD ----- THE ---- FA-T CAT

T-Coffee Flowchart

The **extended pairwise alignments** are used to build a guide tree and **progressive alignment procedure** is then applied.

T-Coffee **considers both global and local pairwise alignments**.

It can also **add supplementary information** (domain, motifs....)



Alignment with T-Coffee

Calculate the alignment of five sequences of the Cytochromes used before with T-Coffee

CLUSTAL W (1.83) multiple sequence alignment

```
sp|P00004|CYC_HORSE      -----
sp|P00091|CYC22_RHOPA  MVK-----KLLTILSIAATAGS
sp|P0C0X8|CYC2_RHOSH   Q-----
sp|P99999|CYC_HUMAN    -----
sp|Q93VA3|CYC6_ARATH  MRLVLSGASSFTSNLFCSSQQVNGRGKELKNPISLNHNKDLDFLLKKLAPPLTAVLLAVS
```

```
sp|P00004|CYC_HORSE      -----MGDVEKGKKIFVQKCAQCHTVE-----KGGKHKTGPNLHGLFGRK
sp|P00091|CYC22_RHOPA  -----LSIGTASAQDAKAGEAVFK-QCMTCHRA-----DKNMVG PALGGVVGRK
sp|P0C0X8|CYC2_RHOSH   -----EGDPEAGAKAFN-QCQTCHVIVDDSGTTIAGRNAKTGPNLYGVVGRK
sp|P99999|CYC_HUMAN    -----MGDVEKGKKIFIMKCSQCHTVE-----KGGKHKTGPNLHGLFGRK
sp|Q93VA3|CYC6_ARATH  PICFPPESLG--QTLDIQRGATLFRACIGCHDTG-----GNIIQPG-----
                        * : * * * ** *

```

```
sp|P00004|CYC_HORSE      TGQAPGFT-YTDANKN---KGITWKEETL-MEYLENPKKYI-----PGTKM
sp|P00091|CYC22_RHOPA  AGTAAGFT-YSPLNHNSGEAGLVWTADNI-INYLNDPNAFLKKFLTDKGKADQAVGVTKM
sp|P0C0X8|CYC2_RHOSH   AGTQADFKGYGEGMKEAGAKGLAWDEEHF-VQYVQDPTKFLKEYTG-----DAKAKGKM
sp|P99999|CYC_HUMAN    TGQAPGYS-YTAANKN---KGI IWGEDTL-MEYLENPKKYI-----PGTKM
sp|Q93VA3|CYC6_ARATH  -----ATLFT---KDLERNVDTEEEIYRVTYFGK--GRMPGFGE-----KCTPRGQC
                        . : : * : : * . . :

```

```
sp|P00004|CYC_HORSE      IFAGIKKKTEREDLIAYLKK-----ATNE
sp|P00091|CYC22_RHOPA  TFK-LANEQQRKDVVAYL-----ATLK
sp|P0C0X8|CYC2_RHOSH   TFK-LKKEADAHNIWAYLQQ-----VAVRP
sp|P99999|CYC_HUMAN    IFVGIKKKEERADLIAYLKK-----ATNE
sp|Q93VA3|CYC6_ARATH  TFG-PRLQDEEIKLLAEFVKFQADQGWPTVSTD
                        * : : . : * :

```

Alignment Benchmark

BAiBASE was the first large scale benchmark specifically designed for MSA, providing high quality manually refined reference **alignments based on 3D structure** superpositions.

BAiBASE is divided into several reference datasets:

1. cases with small numbers of equidistant sequences, and was further subdivided by percent identity;
2. families with one or more “orphan” sequences;
3. a pair of divergent subfamilies, with less than 25% identity between the two groups;
4. sequences with large terminal extensions (N/C-terminal);
5. sequences with large internal insertions and deletions.

Benchmark Evaluation

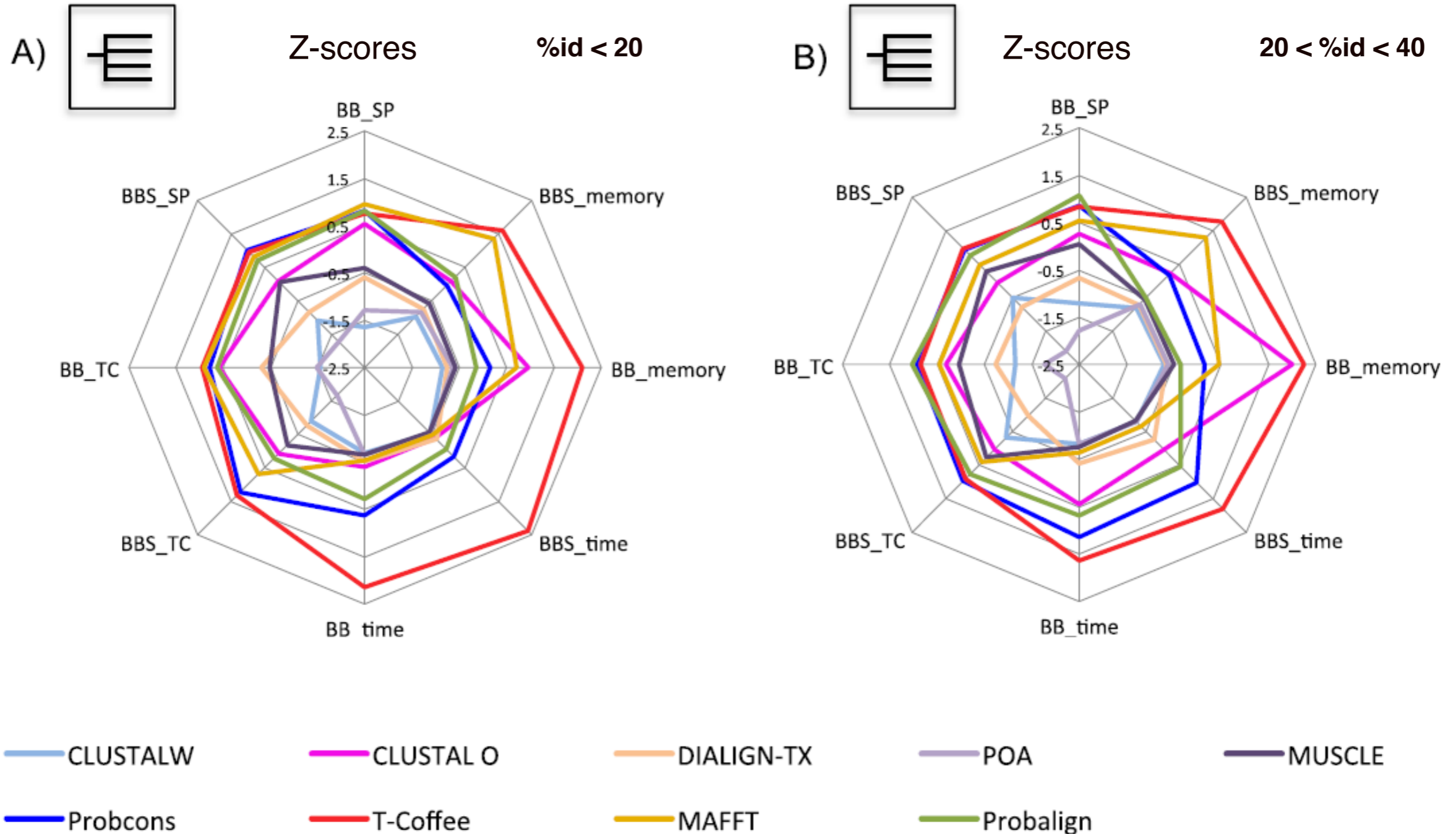
SP score determines the extent to which the programs succeed in aligning input sequences in an MSA. It is calculated as the ratio of the sum of scores p for all pairs of residues in every column of the alignment by the sum of scores in the reference alignment; $p = 1$ if the pair of compared residues is aligned identically in the reference alignment, otherwise $p = 0$.

The TC score is calculated considering the ratio of the sum of scores c by the number of columns in the alignment, being $c = 1$ if all residues in the column are aligned identically in the reference alignment, otherwise $c = 0$

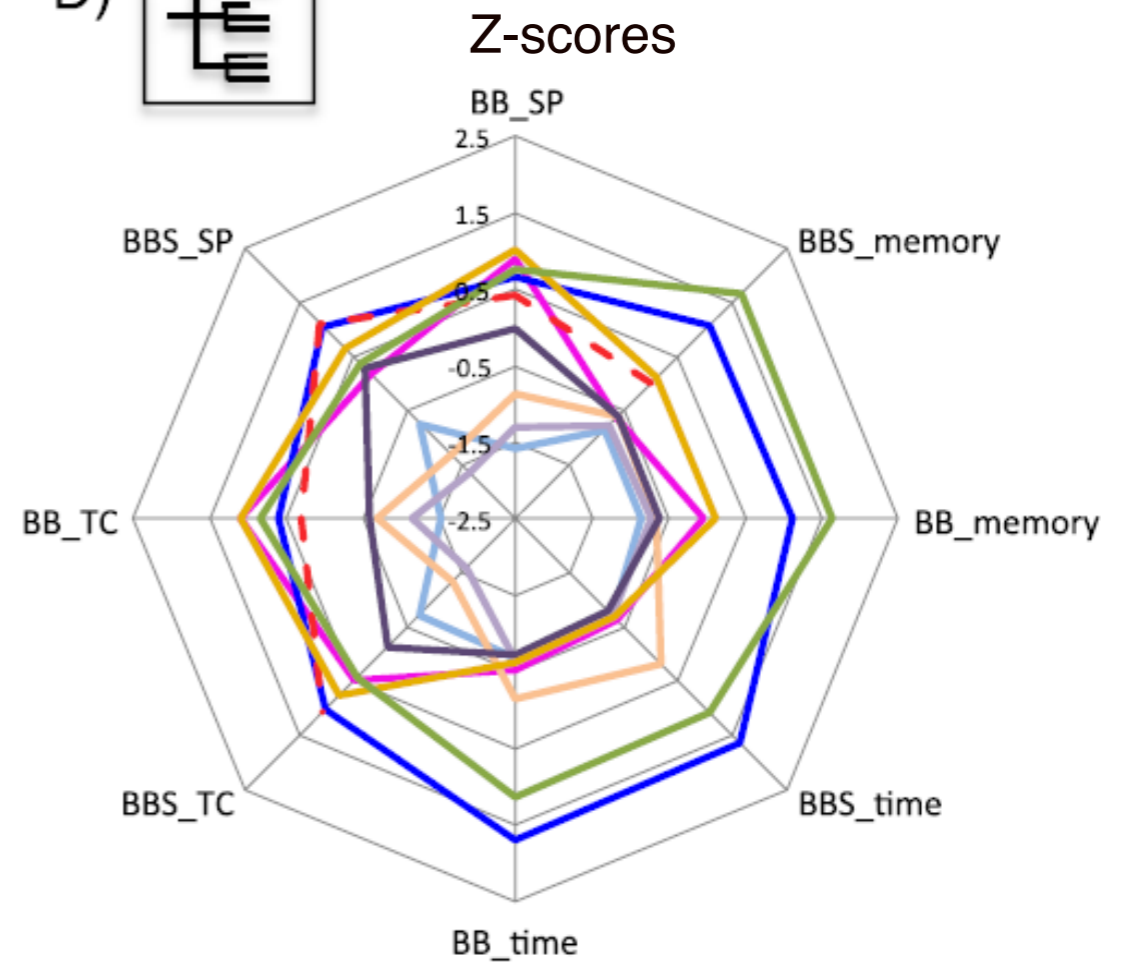
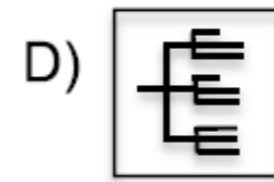
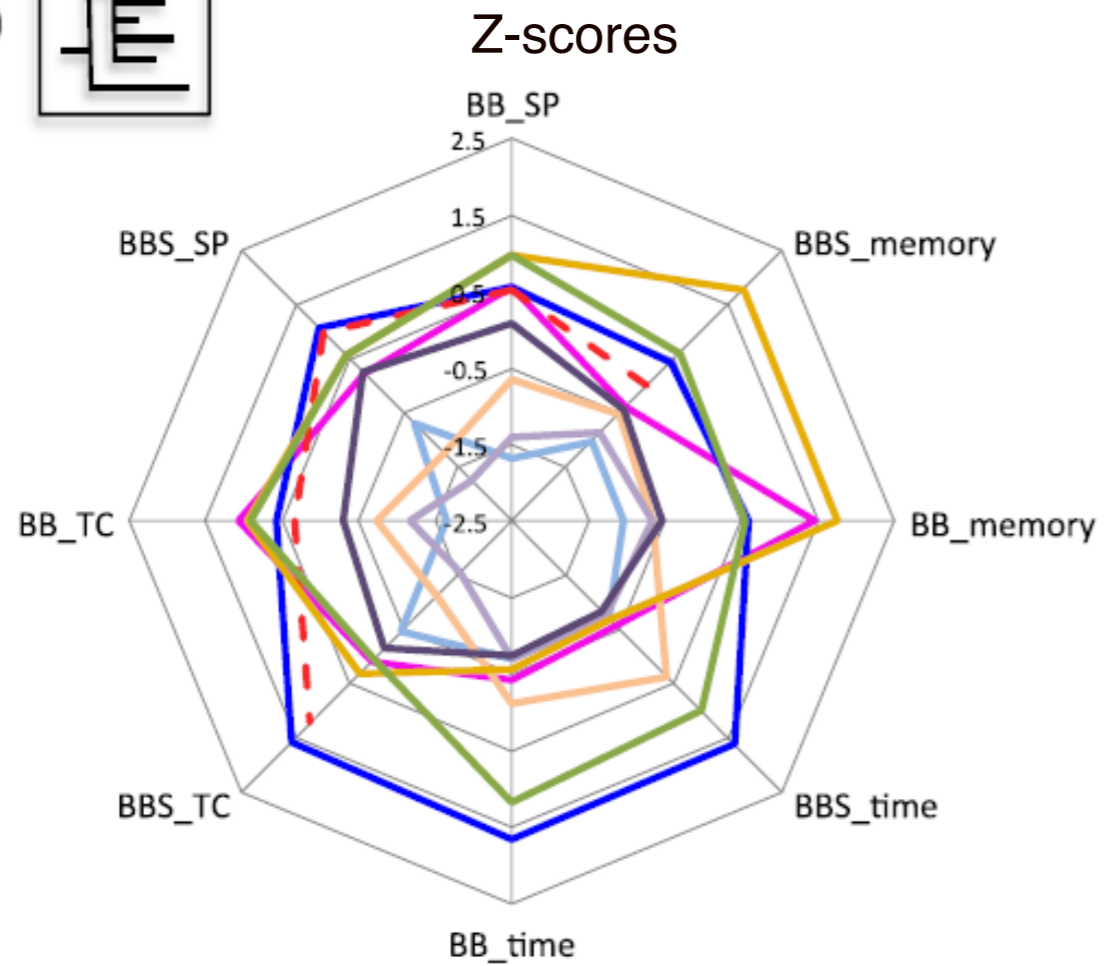
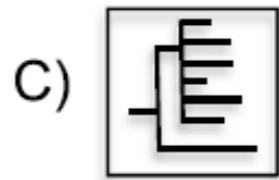
Time of execution

Peak memory

Performance (I)



Performance (II)



CLUSTALW

CLUSTAL O

DIALIGN-TX

POA

MUSCLE

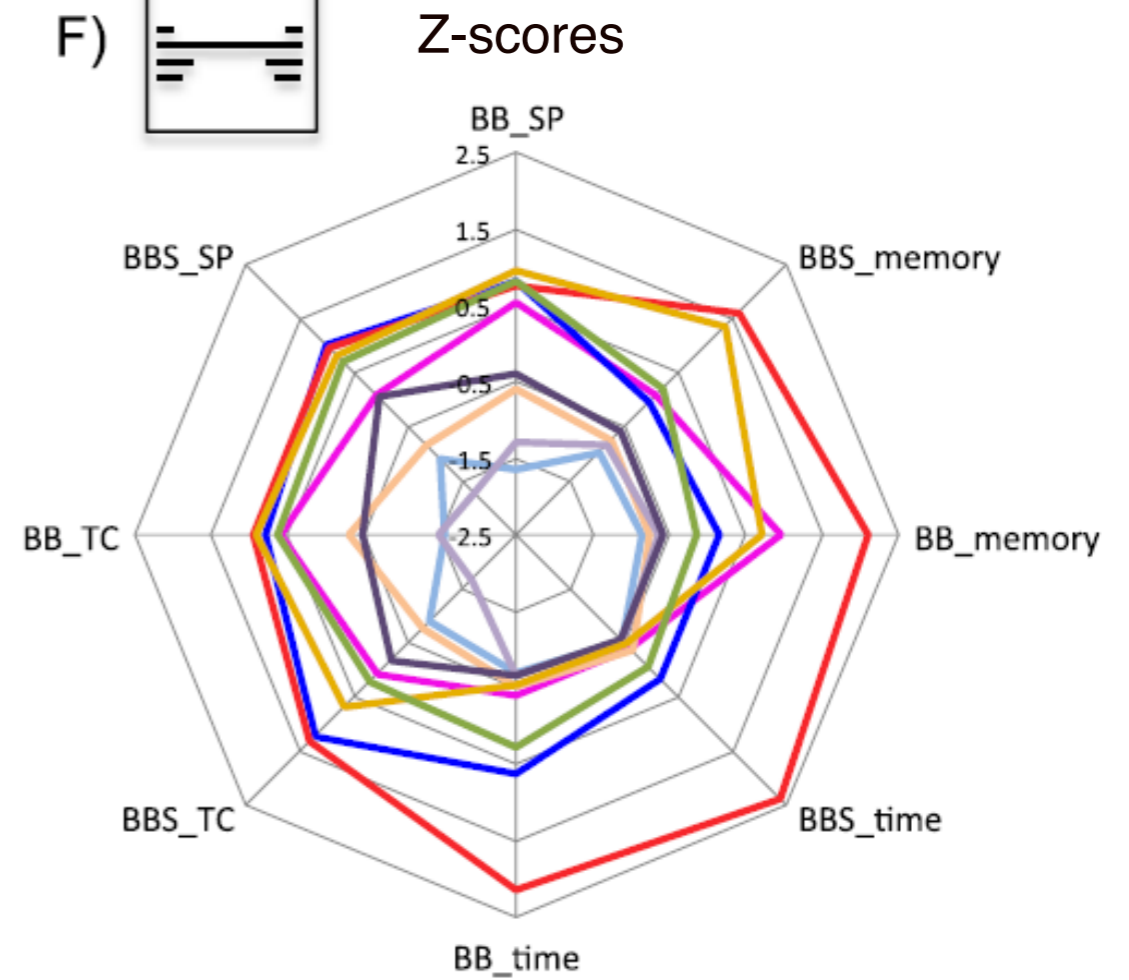
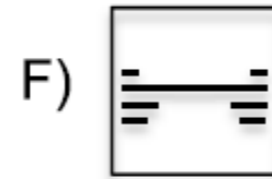
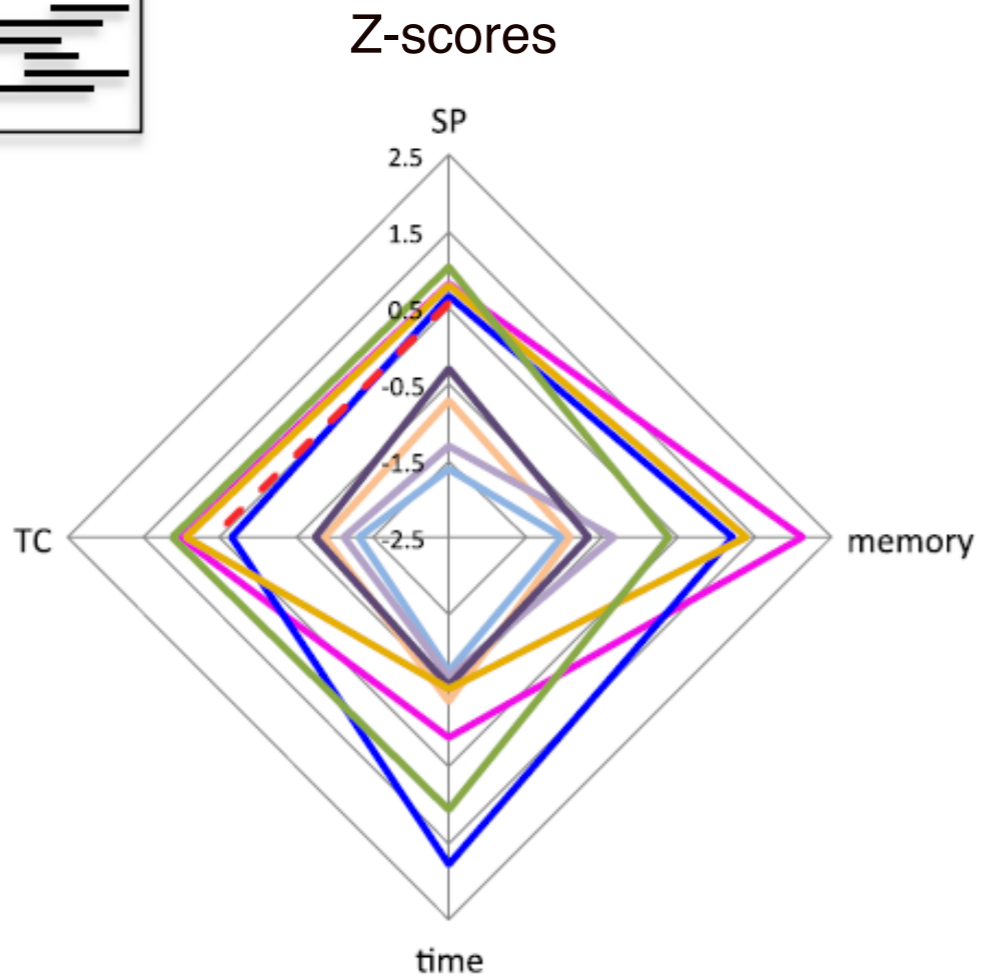
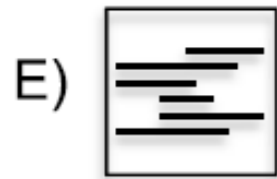
Probcons

T-Coffee

MAFFT

Probalign

Performance (III)



CLUSTALW

CLUSTAL O

DIALIGN-TX

POA

MUSCLE

Probcons

T-Coffee

MAFFT

Probalign

Results and Conclusions

Results: Our results indicate that mostly the consistency-based programs Probcons, T-Coffee, Probalign and MAFFT outperformed the other programs in accuracy. Whenever sequences with large N/C terminal extensions were present in the BALiBASE suite, Probalign, MAFFT and also CLUSTAL OMEGA outperformed Probcons and T-Coffee. The drawback of these programs is that they are more memory-greedy and slower than POA, CLUSTALW, DIALIGN-TX, and MUSCLE. CLUSTALW and MUSCLE were the fastest programs, being CLUSTALW the least RAM memory demanding program.

Conclusions: Based on the results presented herein, all four programs Probcons, T-Coffee, Probalign and MAFFT are well recommended for better accuracy of multiple sequence alignments. T-Coffee and recent versions of MAFFT can deliver faster and reliable alignments, which are specially suited for larger datasets than those encountered in the BALiBASE suite, if multi-core computers are available. In fact, parallelization of alignments for multi-core computers should probably be addressed by more programs in a near future, which will certainly improve performance significantly.

Exercise

Download from UniProtKB the sequences of the following proteins (in FASTA format)

P99999 (human)

P00004 (horse)

P0C0X8 (Rhodobacter)

P00091 (Rhodopseudomonas)

Q93VA3 (Arabidopsis)

Align with ClustalW @

<http://clustalw.ddbj.nig.ac.jp/>

<http://www.ch.embnet.org/software/ClustalW.html>

Write a script to calculate the information entropy of the MSA and for each column the most conserved residue and its frequency.

Exercise

Using the BLAST tool at Uniprot, retrieve all the SwissProt sequences that are similar with an E-value <0.001 to the Rhodospseudomonas cytochrome C (P00091).

Download the sequences in Fasta format and align with ClustalW, Muscle or T-Coffee

- Analyze the conserved positions in the alignments
- Repeat with the Arabidopsis (Q93VA3) and the human (P99999) sequences
- Compare the results and in particular the pattern of conserved residues