

Global Multiple Sequence Alignment

```

HUMAN MKWVTFISLL FLFSSAYSRG V..FRRDA.H KSEVAHRFKD LGEENFKALV
RABIT MKWVTFISLL FLFSSAYSRG V..FRREA.H KSEIAHRFND VGEEHFGLV
PIG   ~-WVTFISLL FLFSSAYSRG V..FRRDT.Y KSEIAHRFKD LGEQYFKGLV
CHICK MKWVTLSIFI FLFSSATSRN LQRFARDAEH KSEIAHRYND LKEETFKAVA
  
```

Align k sequences, so that residues in each column share a property of interest:

- a common ancestor
- a structural or functional role

Global Multiple Sequence Alignment

Given sequences $s_1 \dots s_k$ of lengths $n_1 \dots n_k$

seek $s'_1 \dots s'_k$ of length $l \geq \max\{n_j\}$ such that

- Obtain s_i from s'_i by removing gaps
- No column contains all gaps
- The score of the alignment is optimal

Scoring function: Sum-of-Pairs

$$\text{Score} = \sum_{a=1}^k \sum_{b=1}^k \sum_{b > a} p(s'_a[i], s'_b[i])$$

(1) **A TT**
 (2) **A T _**
 (3) **ACAT**




$p[_,_] = 0$

$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= 0 + g + g = 2g \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

Scoring function: Sum-of-Pairs

$$\text{Score} = \sum_{a=1}^k \sum_{a=1}^k \sum_{b>a} p(s'_a[l], s'_b[l])$$




- (1) **A**TT  $p[_,_]=0$
 (2) **A**T  $p[_,_]=0$
 (3) **ACAT** 

$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= M + m + m = 2m + M \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

Scoring function: Sum-of-Pairs

$$\text{Score} = \sum_{a=1}^k \sum_{a=1}^k \sum_{b>a} p(s'_a[l], s'_b[l])$$

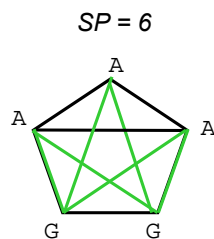
- (1) **A**TT  $p[_,_]=0$
 (2) **A**T  $p[_,_]=0$
 (3) **ACAT** 

$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= g + M + g = 2g + M \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

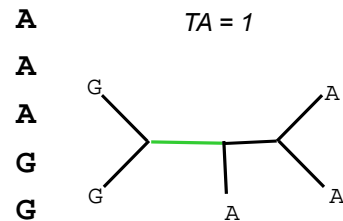
Scoring Multiple Alignments

Sum of Pairs

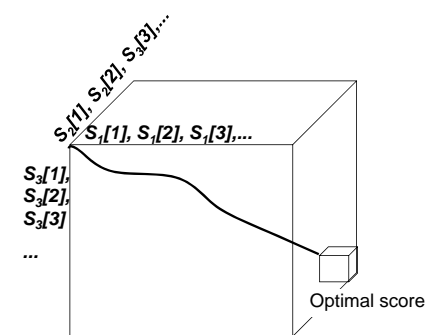


$$\text{Score} = \sum_{x=1}^k \sum_{y>x} d(s'_x[j], s'_y[j])$$

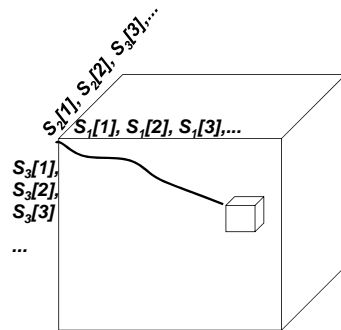
Tree alignment



Dynamic Programming for Multiple Alignment



Dynamic Programming for Multiple Alignment



Each cell has $O(2^k)$ neighboring cells
 Calculating the sum-of-pairs score for each neighbor is $O(k^2)$
 Number of cells in matrix: $O(n^k)$

Total computational complexity:
 $O(n^k 2^k k^2)$

Limits:

- ~ $k = 8 - 10$ sequences
- ~ $n = 500$ residues

MSA is NP-complete for Sum-of-Pairs scoring

Observations

1. A multiple alignment induces pairwise alignments
2. A column in the induced pairwise alignment may contain all gaps, even though no column in the MSA contains all gaps.

(1)

AG	CT
----	----

 (2)

AG	CT
----	----

 (3) **ACT**_T

3. The pairwise alignments induced by the *optimal multiple alignment* are *not* the same as the *optimal pairwise alignments*.

Optimal Pairwise Alignments

(1) ACT
 (2) AGT

1 substitution

(1) AC_T
 (2) A_GT
 (3) ACGT

2 indels

Although this costs more, it may be a biologically more realistic alignment

Since exact methods for MSA have exponential time complexity, heuristic approaches are used. Progressive alignment is the most commonly used.

Basic progressive alignment strategy:

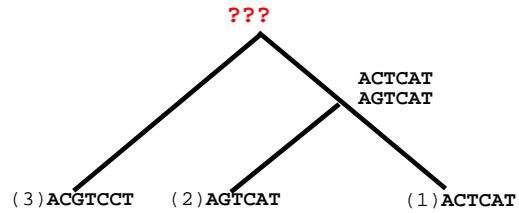
- Compute D , a matrix of distances between all pairs of sequences
- From D , construct a “guide tree” T
- Construct MSA by pairwise alignment of partial alignments (“profiles”) guided by T
- Improve alignment by postprocessing steps.

Optimal Pairwise Alignments

(1) ACTCAT	(1) ACTCAT	3
(2) AGTCAT	(2) AGTCAT	
(3) ACGTCCT	(2) A_GTCAT	5
	(3) ACGTCCT	
	(1) AC_TCAT	5
	(3) ACGTCCT	

$d(x, y) = 3$
 $d(x, _) = 2$

Progressive Alignment



- Use *profile alignment* to merge sequences according to a guide tree.
- Typically, most closely related sequences are merged first.

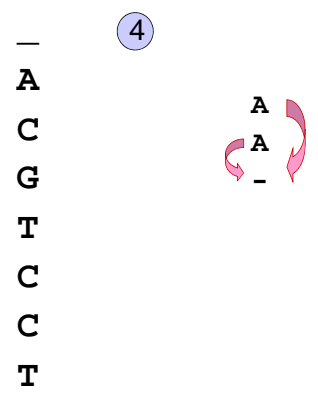
Merging strategy:

Align the profile (1,2) with sequence (3)

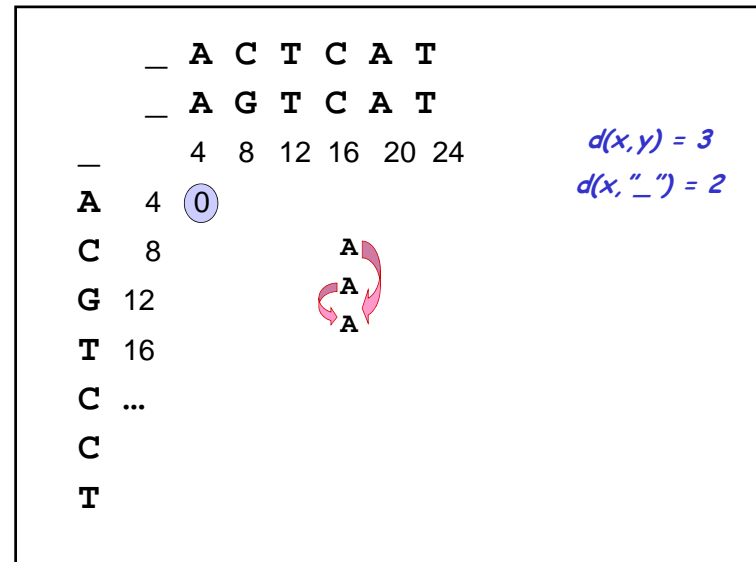
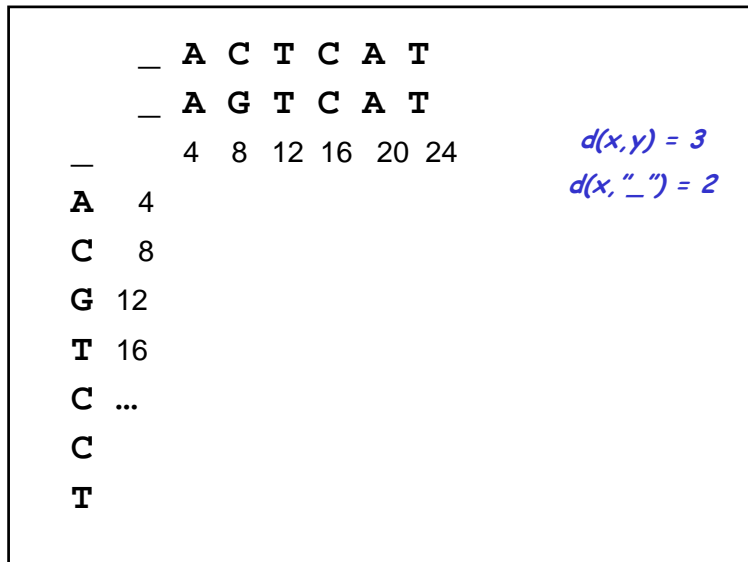
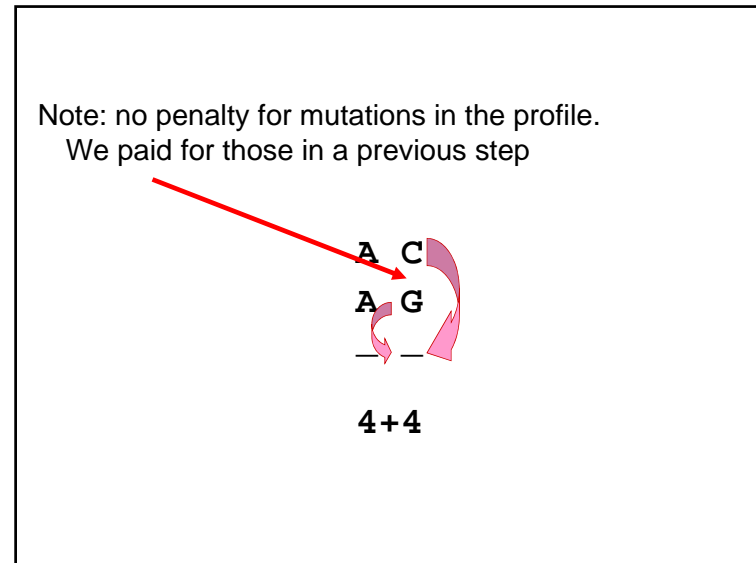
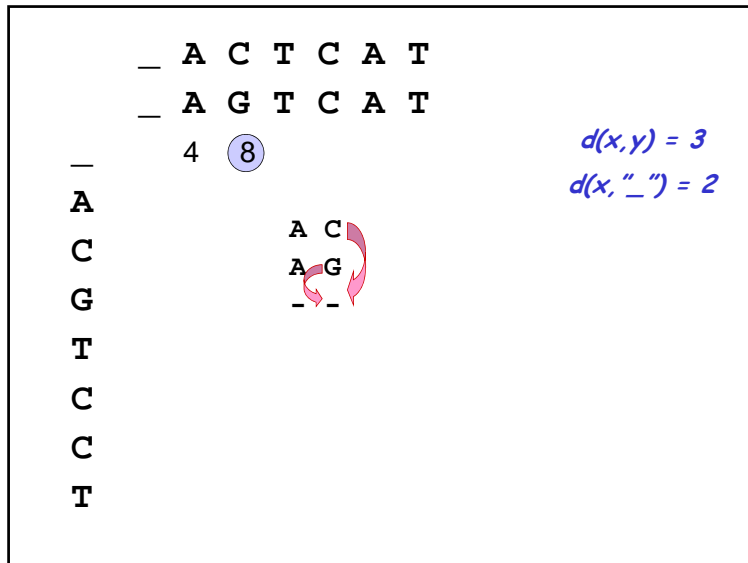
(1) ACTCAT	(1) ACTCAT	3
(2) AGTCAT	(2) AGTCAT	
(3) ACGTCCT	(2) A_GTCAT	5
	(3) ACGTCCT	
	(1) AC_TCAT	5
	(3) ACGTCCT	

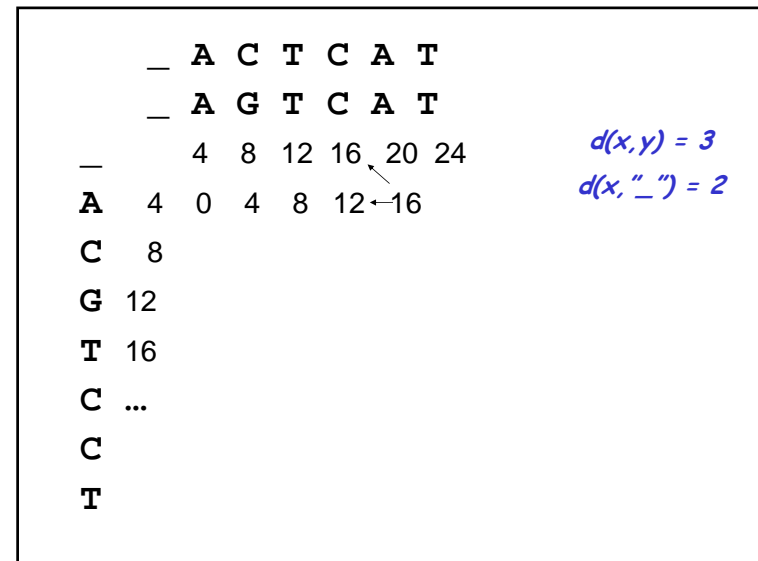
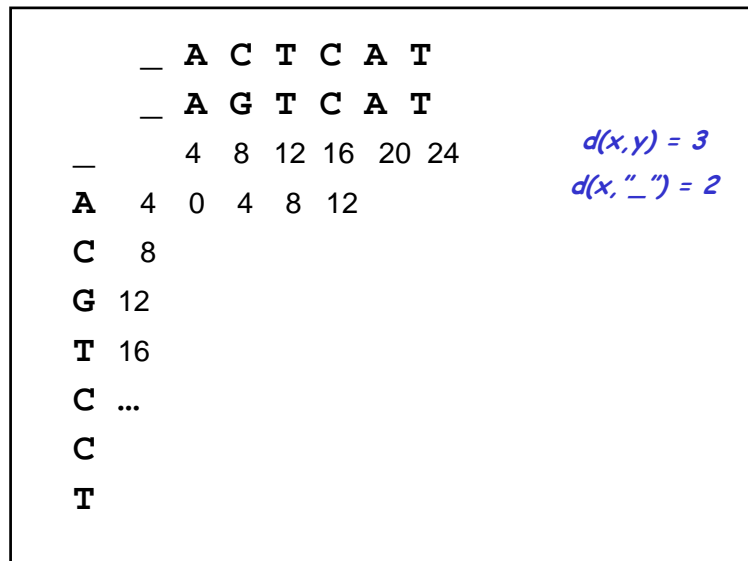
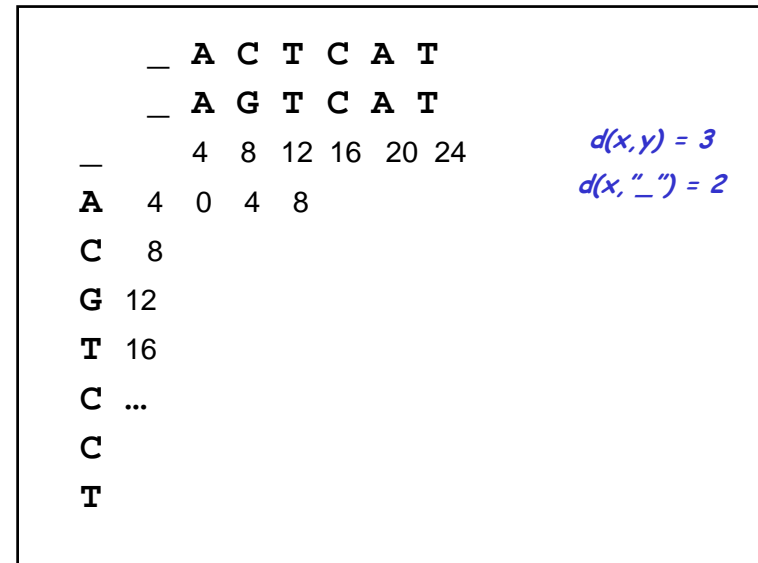
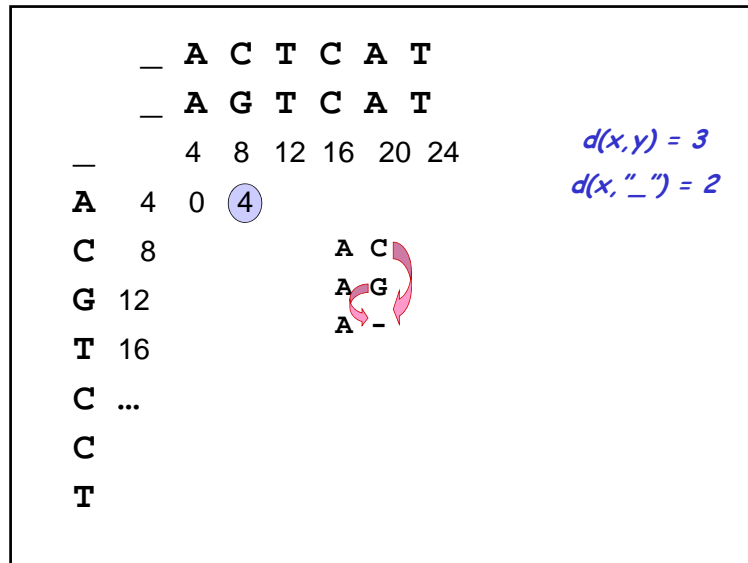
$d(x, y) = 3$
 $d(x, _) = 2$

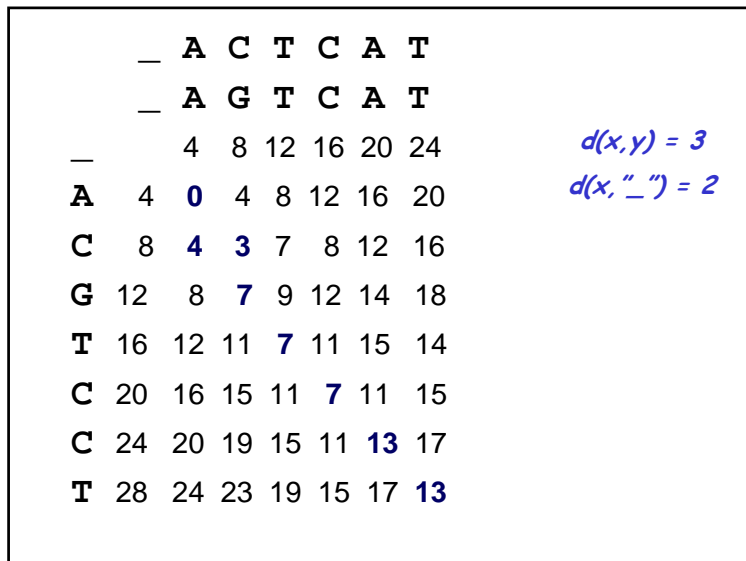
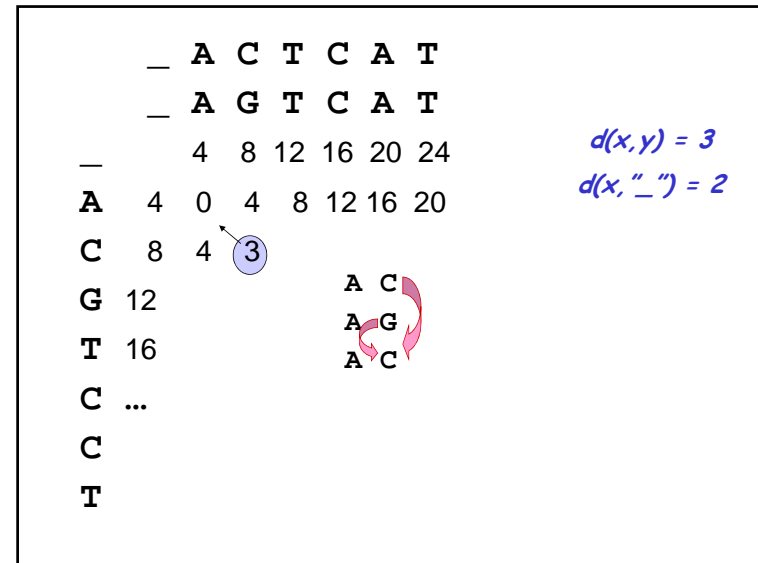
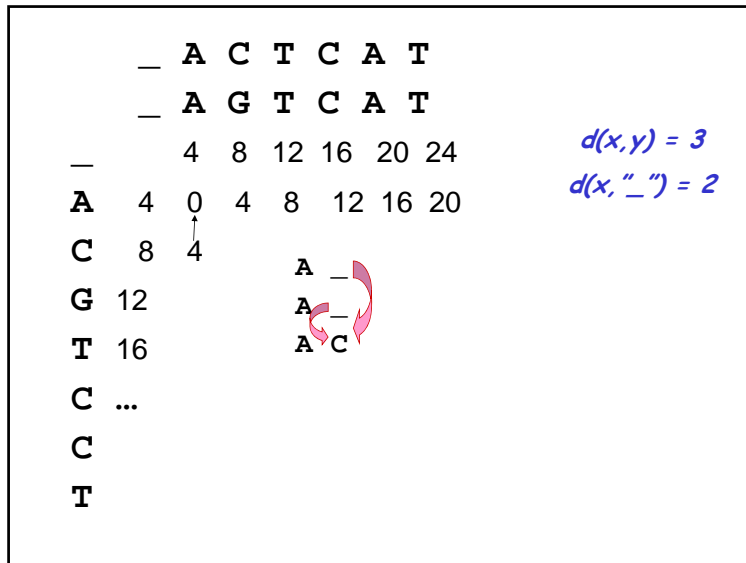
_ A C T C A T
 _ A G T C A T



$d(x, y) = 3$
 $d(x, _) = 2$







Optimal Pairwise Alignments		Progressive alignment
(1) ACTCAT	(2) AGTCAT	(1,2) + (3)
(2) A_GTCAT	(3) ACGTCCT	(1) AC_TCAT $4m+2g$
(3) ACGTCCT	(2) AG_TCAT	(2) AG_TCAT
(1) AC_TCAT	(2) A_GTCAT	An alternate alignment
(3) ACGTCCT	(3) ACGTCCT	(1) AC_TCAT
		(2) A_GTCAT $2m+4g$
		(3) ACGTCCT

Optimal Pairwise Alignments

- (1) ACTCAT
 (2) AGTCAT
 (2) A_GTCAT
 (3) ACGTCCT
 (1) AC_TCAT
 (3) ACGTCCT

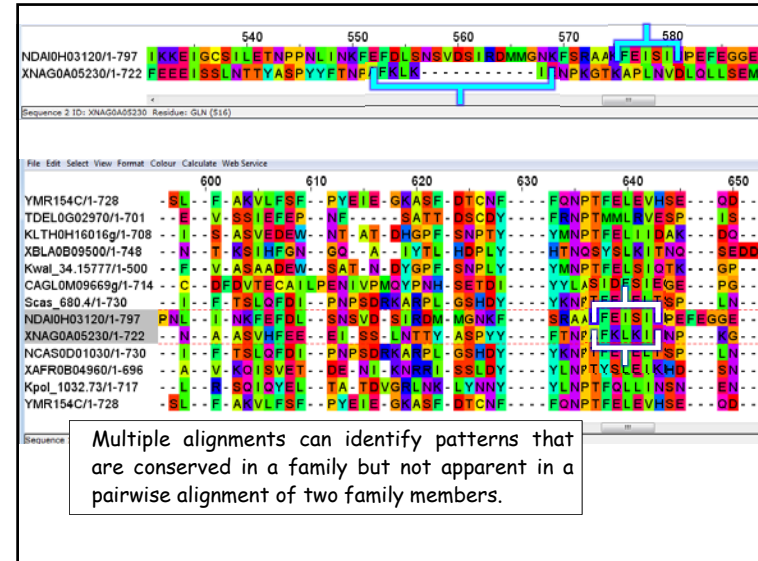
Progressive alignment

(1,2) + (3)

- (3) ACGTCCT
 (1) AC_TCAT 16
 (2) AG_TCAT

An alternate alignment

- (1) AC_TCAT
 (2) A_GTCAT 14
 (3) ACGTCCT



Progressive alignment

- “Once a gap, always a gap”
 - You can’t go back and correct a bad decision at an earlier step.
- Progressive alignment is not guaranteed to give the optimal alignment.
- But it does have better complexity...

Complexity of progressive alignment

- Distance matrix
 - Each pairwise alignment $O(n^2)$
 - Number of pairwise alignments $O(k^2)$
 - Iterative construction of MSA
 - Number of merge steps $O(k)$
 - Each pairwise alignment $O(k^2n^2)$
- Entire method $O(k^2n^2)$

Summary: Progressive alignment heuristics

- Not guaranteed to give the optimal MSA
- Bad choice of gaps propagates
- Complexity
 - Progressive: $O(k^2n^2)$
 - versus DP: $O(n^k 2^k k^2)$
- Typically, merge the most closely related sequences first.

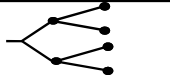


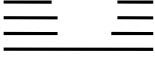
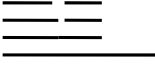
Mathematical correctness is not a guarantee of biological accuracy. The performance of MSA programs is typically evaluated using benchmarks based on biological data:

- Curated structural alignment
- Automated structural alignment
- Real or simulated sequence

Various benchmarks are designed to mimic properties of different types of data sets encountered in practice, especially those that are challenging to align:

- Highly divergent sequences, e.g., <50% or <30% identity
- A family of related sequences plus several outliers, or “orphan” sequences
- Related sequences that differ due to large N or C terminal extensions or large internal insertions or deletions

Benchmark challenges

PROBLEM	Description
	Even Phylogenetic Spread.
	One Outlier Sequence
	Two Distantly related Groups
	Long Internal Indel
	Long Terminal Indel

Source: BaliBase, Thompson et al, NAR, 1999.

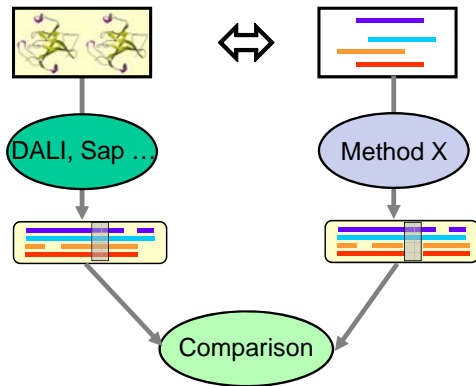
Mathematical correctness is not a guarantee of biological accuracy. The performance of MSA programs is typically evaluated using benchmarks based on biological data:

- Curated structural alignment
- Automated structural alignment
- Real or simulated sequence

Various benchmarks are designed to mimic properties of different types of data sets encountered in practice, especially those that are challenging to align:

- Highly divergent sequences, e.g., <50% or <30% identity
- A family of related sequences plus several outliers, or “orphan” sequences
- Related sequences that differ due to large N or C terminal extensions or large internal insertions or deletions

BaliBase: Reference MSAs based on structural alignment.



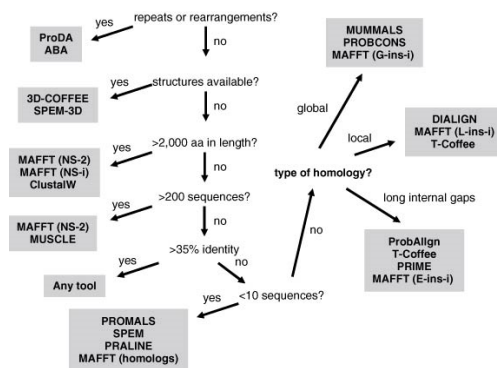
Note that implementation choices result in substantial differences in running time:

Aligner	Performance*	Time
DIALIGN	57.2	12 h, 25 min
CLUSTALW	58.9	2 h, 57 min
T-Coffee	63.6	144 h, 51 min
MUSCLE	64.8	3 h, 11 min
MAFFT	64.8	2h,36min
ProbCons	66.9	19 h, 41 min
ProbCons-ext	68.0	37 h, 46 min

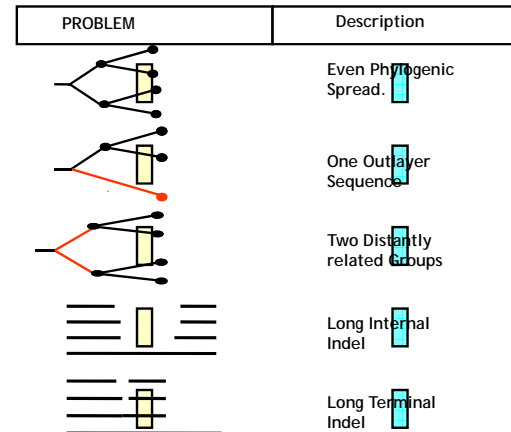
* Fraction of correctly aligned residue pairs

Do et al, Genome Research, 2005

Which program to choose?



Do and Katoh, 2008



Source: BaliBase, Thompson et al, NAR, 1999,

Approaches for improving MSA (Speed or accuracy)

- Iterative refinement of the MSA
- Faster estimation of the guide tree
- Better scoring
 - Combining information from various sources
 - Consistency in alignments of 3 sequences
 - Weighting sequences pairs
- Position specific gap penalties

Iterative refinement

Progressive “alignment suffers from its greediness”
Notredame et al, JMB 2000

1. Randomly select one sequence, remove it and realign it with the rest of the alignment
2. Remove each sequence in turn and realign with the remaining alignment. Select the best of these as the new alignment.
3. Randomly split into two sub alignments and realign them.

Apply strategy repeatedly until convergence or out of computer time

Approaches for improving MSA (Speed or accuracy)

- Iterative refinement of the MSA
- Faster estimation of the guide tree
- Better scoring
 - Combining information from various sources
 - Consistency in alignments of 3 sequences
 - Weighting sequences pairs
- Position specific gap penalties

Combining information from multiple sources

T. Coffee, Notredame, Higgins, Heringa, JMB 2000

```
DRHNSNIKV
DLKPENLLI
```

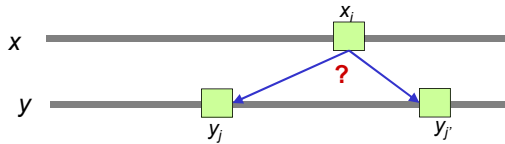
22% identity

```
DRHNSNIKVDDG_QLFHIDFGHFLD
YLHSLDIYRDLKPENLIDQQGYIQV
```

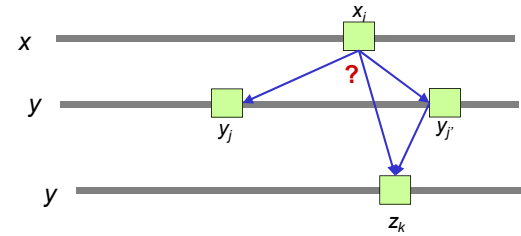
12% identity

Construct a *library* of pairwise alignments

Consistency



Consistency



Sequence 2 ID: XNAG0A05230 Residue: GLN (516)

File Edit Select View Format Colour Calculate Web Service

600 610 620 630 640 650

YMR154C/1-728 SL F AKVLFSE PYEIE GKASF DTCNF FONPTFELEVHSE OD

TDEL0G02970/1-701 E V SSIEFEP NF SATT DSCNY FRNPTMMLVESP IS

KLTH0H16016g/1-708 I S ASVEDEW NT AT DHGPF SNPTY YMNPTFELIDAK DQ

XLBA0B09500/1-748 N T KSIHFGN GQ A IYTL HDPLY HTNOSYSLK TNQ SEDD

Kwal_34.15777/1-500 F V ASAADRW SAT N YGPF SNPLY YMNPTFELS OTK GP

CAGL0M09669g/1-714 C DEVTCAI LPENI VPMQYPNH SETDI YYLASIDFSIEGE PG

Scas_680.4/1-730 I F TSLQFDI PNPSDRKAAPL GSHDY YKNPTFELTSP LN

NDAIH03120/1-797 PNL I NKFEFDL SNSVD SIRDM MGNKF SAALFEISIIPEFEGGE

XNAG0A05230/1-722 N A ASVHFEE EI SS LNTTY ASPYY FTNP FVKL I INP KG

NCAS0D01030/1-730 I F TSLQFDI PNPSDRKAAPL GSHDY YKNPTFELTSP LN

XAFR0B04960/1-696 A V KOISVET DE NI KNRR I SSLDY YLNP TYSLEIKHD SN

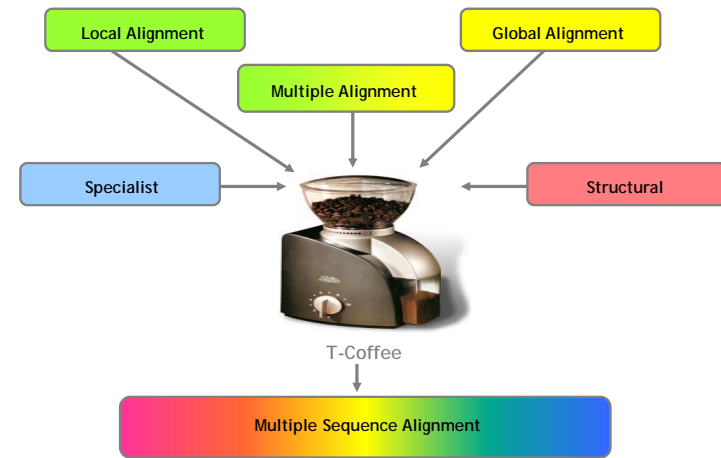
Kpol_1032.73/1-717 L R SOIQYEL TA TDVGR LNK LYNNY YLNP TOLLINSN EN

YMR154C/1-728 SL F AKVLFSE PYEIE GKASF DTCNF FONPTFELEVHSE OD

Sequence

Multiple alignments can identify patterns that are conserved in a family but not apparent in a pairwise alignment of two family members.

Combining information from multiple sources:

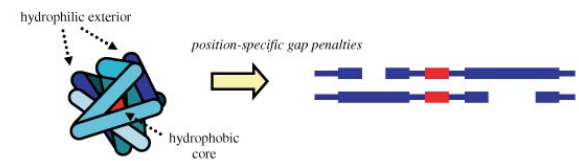


Approaches for improving MSA (Speed or accuracy)

- Iterative refinement of the MSA
- Faster estimation of the guide tree
- Better scoring
 - Combining information from various sources
 - Consistency in alignments of 3 sequences
 - Weighting sequences pairs
- Position specific gap penalties

Position specific gap penalties

Penalize gaps in hydrophobic and hydrophilic regions differently



Do and Katoh, 2008

Other improvements

- Sequence weighting



Assign weights so that **these sequences** do not dominate.

Do and Katoh, 2008

