# Introduction to Bioinformatics
# 06 : Multiple Sequence Alignment

Instructor: Hossein Zeinali

Amirkabir University of Technology

# Overview

1. What is a multiple sequence alignment (MSA)?

2. Where/why do we need MSA?

3. What is a *good* MSA?

4. Algorithms to compute a MSA

# Multiple Sequence Alignment

- Generalize pairwise alignment of sequences to include **> 2** *homologous (related)* sequences

- Analyzing more than 2 sequences gives us <span style="color:red">much more information</span>:
  - Which amino acids are required? Correlated?
  - Evolutionary/phylogenetic relationships

- Similar to PSI-BLAST idea (not yet covered):
  - Use a set of homologous sequences to provide more "sensitivity"

# Multiple Sequence Alignments

# What is a MSA?

```
ATT-GC
ATTTGC
TTTTG
```

*Not* a MSA

```
AT-TGC
ATTTGC
ATTTG-
```

**MSA**

```
AT-T-GC
ATTT-GC
ATTT-G-
```

*Not* a MSA

Why?

# Definition of MSA

- Given a set of sequences, a multiple sequence alignment is an assignment of gap characters, such that
  - Resulting sequences have same length
  - No column contains only gaps

```
ATT-GC          AT-TGC          AT-T-GC
ATTTGC          ATTTGC          ATTT-GC
ATTTG           ATTTG-          ATTT-G-
```

*No*              Yes              *No*

# Applications of MSA

- Building phylogenetic trees and doing phylogenetic analysis of sequence families.
- Finding conserved patterns, e.g.:
  - Regulatory motifs
  - Protein domains
- Identifying and characterizing protein families
  - Find out which protein domains have same function
  - Prediction of protein secondary and tertiary structures
- DNA fragment assembly (in genomic sequencing)

**Goal:** Align homologous positions.

**But**: Without knowledge of phylogenetic tree is this very hard (sometimes impossible) to achieve!

# Scoring an Alignment

- In practice, simple scoring functions are used: usually, columns are scored independently, i.e.

$$S(m) = \sum_i S(m_i) + G$$

gap penalty

$i$th column of alignment $m$

# Scoring Function

- Sum of Pairs (SP) Score = sum of scores of all possible pairs of sequences in an MSA based on a particular scoring matrix

- Compute for each column $c$

$$S(m_i) = \Sigma_{k<l}\, s(m_i^k, m_i^l)$$

PAM or BLOSUM score

residue $l$

$m_i$

# Example

# How Score Gaps in MSAs?

- Want to align gaps with each other over all sequences.
- A gap in a pairwise alignment that "matches" a gap in another pairwise alignment should cost less than introducing a totally new gap.
  - Possible that a new gap could be made to "match" an older one by adjusting older pairwise alignment
  - Change gap penalty near conserved domains of various kinds (e.g. secondary structure elements, hydrophobic regions)

# Example of SP Score with Gap

| | F | Y | *G* | D |
|---|---|---|---|---|
| F | 5 | -2 | -2 | -1 |
| Y | | 7 | 1 | -5 |
| *G* | | | 4 | -3 |
| D | | | | 5 |

BLOSUM 60

$$m = \begin{bmatrix} F \\ F \\ F \end{bmatrix} \begin{bmatrix} - \\ - \\ Y \end{bmatrix} \begin{bmatrix} G \\ G \\ D \end{bmatrix}$$

Gap penalty = -8
$s(-,-) = 0$

$$\begin{bmatrix} G \\ G \\ D \end{bmatrix}$$

$$S(m) = S(m_1) + S(m_2) + S(m_3)$$
$$= 3s(F,F) + 2s(-,Y) + s(-,-) + s(G,G) + 2s(G,D)$$
$$= 15 -16 + 0 + 4 -6 = -3$$

# Algorithms for MSA

## *Exhaustive Methods*

- Multidimensional dynamic programming (DP)
  - <u>Divide-and-Conquer Alignment (DCA)</u> - "semi-exhaustive"
  - <u>Full DP Optimal Global Alignment?</u>
    - *Prohibitive in both time & space requirements for more than 10 sequences!!*

## *Heuristic Methods*

- Progressive alignments
  - We will cover Clustal, Star Alignment, T-Coffee, POA
  - Others: DbClustal and PRALINE -see text-book

# Algorithms for MSA (Cont.)

- Iterative methods
  - Idea: optimal solution can be found by repeatedly modifying existing suboptimal solutions (eg: PRRN)

- Block-based Alignment
  - Multiple re-building attempts to find best alignment (eg: DIALIGN2 & Match-Box)

- Local alignments
  - Profiles, Blocks, Patterns - more on these soon!

# Dynamic Programming for MSA

- As with pairwise alignments, multiple sequence alignments can be computed by dynamic programming



2D

3D

# Generalized Needleman-Wunsch Algorithm

- Given 3 sequences x, y, and z:
- Main iteration loop:

$$
\begin{aligned}
F(i,j,k) = \max \; ( \; &F(i{-}1, j{-}1, k{-}1) + S(x_i, y_j, z_k), \\
&F(i{-}1, j{-}1, k\;\;\;) + S(x_i, y_j, -\;), \\
&F(i{-}1, j\;\;\;, k{-}1) + S(x_i, \;-, z_k), \\
&F(i{-}1, j\;\;\;, k\;\;\;) + S(x_i, \;-, -\;), \\
&F(i\;\;\;, j{-}1, k{-}1) + S(\;-, y_j, z_k), \\
&F(i\;\;\;, j{-}1, k\;\;\;) + S(\;-, y_j, -), \\
&F(i\;\;\;, j\;\;\;, k{-}1) + S(\;-, \;-, z_k) \; )
\end{aligned}
$$

3D

# What Happens to Computational Complexity?

- Given k sequences of length n:
  - Space for matrix: $O(n^k)$
  - Neighbors/cell: $2^k-1$
  - Time to compute SP score: $O(k^2)$
  - Overall runtime: $O(k^2 2^k n^k)$

3D

- So, full dynamic programming is limited to small datasets of less than ten short sequences.

# An Example of DP's Running Time

- Overall runtime: $O(k^2 2^k n^k)$

| # sequences | running time |
|---|---|
| 2 | 1 second |
| 3 | 2 minutes |
| 4 | 5 hours |
| 5 | 3 weeks |
| 6 | 9 years |

*Don't worry, there are fast heuristics*

Sequences: globins ($\approx$ 150 aa)

- Implementation example:
  - Divide-and-Conquer Alignment (DCA): semi-exhaustive
    - Breaking each of the sequences into two smaller sections.
  - http://bibiserv.techfak.uni-bielefeld.de/dca

# Progressive Alignment

Heuristic procedure:

1. Align *most* similar sequences first
2. Add sequences progressively

- Often **guide trees** is used to determine order of alignments.

**Examples:**    Star alignment

ClustalW

**Multiple Alignment by adding sequences**

1 ————————
2 ————————
3 ————————
4 ————————

# What is a Consensus Sequence?

- A single sequence that represents *most common* residue of each column in a MSA

- **Example**:

```
FGGHL-GF
F-GHLPGF
FGGHP-FG
```
FGGHL-GF

- **Steiner consensus sequence:** given sequences $s_1, \ldots, s_k$, find a sequence $s*$ that maximizes $\Sigma_i\ S(s*, s_i)$

# Guide Tree

Binary tree

- Leaves correspond to sequences

- Internal nodes represent alignments

- Root corresponds to final MSA

- Is created using neighbor-joining method

- Is a simple phylogenetic tree

# Example

S1: HVLIP
S2: HMIP
S3: HVLP
S4: LVLIP

# Progressive Alignment Steps



all individual pairwise alignment and construction of distance matrix

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A | – |   |   |   |   |
| B | 11 | – |   |   |   |
| C | 20 | 30 | – |   |   |
| D | 27 | 36 | 9 | – |   |
| E | 30 | 33 | 20 | 27 | – |

calculating a guide tree; C & D the closest pair; A & B the next closest pair

aligning C/D and A/B separately using dynamic programming

C/D and A/B alignments reduced to consensus sequences which are aligned to each other

creating a new consensus for C/D/A/B which aligns with E

completing alignment

# Clustal Program

- The most well-known progressive alignment program
  - [www.ebi.ac.uk/clustalw](http://www.ebi.ac.uk/clustalw)
- ClustalW and ClustalX: stand-alone programs which run on UNIX and Macintosh respectively.
- Does not rely on a single substitution matrix
  - Applies different scoring matrices depending on degrees of similarity.
- Uses of adjustable gap penalties
  - allow more insertions and deletions in regions that are outside the conserved domains, but fewer in conserved regions.
- Applies a weighting scheme to increase the reliability of aligning divergent sequences.

# Clustal

1. Perform pair-wise alignments between all pairs of sequences (n * (n-1)/2 possibilities)
2. Generate distance matrix
   - Distance between a pair = number of mismatched positions in alignment divided by total number of matched positions
3. Generate a Neighbor-Joining '**guide tree**' from distance table
4. Use guide tree to progressively align sequences in pairs from tips to root of tree

**1    2    3    4**

**Distance Matrix**

**1**
**2**
**3**
**4**

**Guide Tree**

1
2
3
4

**Progressive Alignment**

2
3
4
1

1 + 2
1 + 3
1 + 4
2 + 3
2 + 4
3 + 4

**Pairwise Alignments**

1. Compute pairwise alignments (DP)
2. Convert similarities into distances

   Distance between a pair = # of mismatched positions in alignment (divided by total # of matches)

3. Build guide tree from distances by *Neighbor Joining*
4. Align with respect to guide tree

# ClustalW

```
Hbb_Human   1   -
Hbb_Horse   2   .17    -
Hba_Human   3   .59    .60    -
Hba_Horse   4   .59    .59    .13    -
Myg_Whale   5   .77    .77    .75    .75    -
```

**CLUSTAL W**

Quick pair wise alignment
calculate distance matrix



Neighbour-joining tree
(guide tree)

Progressive alignment
following guide tree

alpha-helices

```
1   PEEKSAVTALWGKVN--VDEVGG
2   GEEKAAVLALWDKVN--EEEVGG
3   PADKTNVKAAWGKVGAHAGEYGA
4   AADKTNVKAAWSKVGGHAGEYGA
5   EHEWQLVLHVWAKVEADVAGHGQ
```

# Clustal Versions (http://www.clustal.org/)



**Clustal: Multiple Sequence Alignment**

Multiple alignment of nucleic acid and protein sequences

**Clustal Omega**

- Latest version of Clustal - fast and scalable (can align hundreds of thousands of sequences in hours), greater accuracy due to new HMM alignment engine
- Command line/web server only (GUI public beta available soon)

**ClustalW/ClustalX**

- "Classic Clustal"
- GUI (ClustalX), command line (ClustalW), web server versions available

# Displaying MSAs using ClustalW

```
CLUSTAL W (1.82) multiple sequence alignment

FOS_RAT     MMFSGFNADYEASSSRCSSASPAGDSLSYYHSPADSFSSMGSPVNTQDFCADLSVSSANF 60
FOS_MOUSE   MMFSGFNADYEASSSRCSSASPAGDSLSYYHSPADSFSSMGSPVNTQDFCADLSVSSANF 60
FOS_CHICK   MMYQGFAGEYEAPSSRCSSASPAGDSLTYYPSPADSFSSMGSPVNSQDFCTDLAVSSANF 60
FOSB_MOUSE  -MFQAFPGDYDS-GSRCSS-SPSAESQ--YLSSVDSFGSPPTAAASQE-CAGLGEMPGSF 54
FOSB_HUMAN  -MFQAFPGDYDS-GSRCSS-SPSAESQ--YLSSVDSFGSPPTAAASQE-CAGLGEMPGSF 54
            *...*  .:*::  .***** **:.:*    *  *..***.*    :.. :*: *:.*.   ...*
```

RED:       AVFPMILW (small)
BLUE:      DE    (acidic, negative chg)
MAGENTA:   RHK   (basic, positive chg)
GREEN:     STYHCNGQ (hydroxyl + amine + basic)

* entirely conserved column
: all residues have ~ same size *AND* hydropathy
. all residues have ~ same size *OR* hydropathy

# Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

**Important note:** This tool can align up to 4000 sequences or a maximum file size of 4 MB.

## STEP 1 - Enter your input sequences

Enter or paste a set of

| PROTEIN | ▼ |
|---|---|

sequences in any supported format:

```
>s1
NLFVALYDFVASG
>s2
KGVALIYALWDY
>s3
GQYRALYDYK
```

Or, upload a file: Browse... No file selected.     Use a example sequence | Clear sequence | See more example inputs

## STEP 2 - Set your parameters

OUTPUT FORMAT

| ClustalW with character counts | ▼ |
|---|---|

*The default settings will fulfill the needs of most users.*

More options...  *(Click here, if you want to view or change the default settings.)*

## STEP 3 - Submit your job

☐ Be notified by email *(Tick this box if you want to be notified by email when the results are available)*

Submit

# Clustal Omega

```
                                              s2 0.316667
                                              s1 0.3
                                              s3 0.3
```

```
CLUSTAL O(1.2.4) multiple sequence alignment


s2        KGVALIYALWDY----        12
s1        ---NLFVALYDFVASG        13
s3        ---GQYRALYDYK---        10
             **:*:
```

# Clustal Drawbacks

- Is not suitable for comparing sequences of different lengths because it is a global alignment–based method.
- The final alignment result is influenced by the order of sequence addition.
- The "greedy" nature of the algorithm: it depends on initial pairwise alignment.
- Any errors made in first steps cannot be corrected.
- To alleviate some of the limitations, a new generation of algorithms have been developed.

# Star Alignment

- Fast heuristic to compute MSA

- Good approximation of *optimal* MSA, if scoring scheme satisfies triangle inequality

**Algorithm:**

1. ***Compute pairwise similarities***

2. ***Select center*** $s_c$ that maximizes $\Sigma_{i \neq c} S(s_c, s_i)$

3. ***Add sequences in decreasing order of similarity*** to *center* $s_c$
   - Rule: "once a gap, always a gap"

Does that function look familiar?

- **Recall: Consensus sequence:** single sequence (*more accurately;* "model") that represents *most common* residue of each column in MSA.

- **Recall: Steiner consensus sequence:** Given sequences $s_1, \ldots, s_k$, find a sequence s* that maximizes $\Sigma_i\ S(s^*, s_i)$

```
FGGHL-GF
F-GHLPGF
FGGHP-FG
```
FGGHL-GF

$S_1=$ ATTCGGATT
$S_2=$ ATCCGGATT
$S_3=$ ATGGAATTTT
$S_4=$ ATGTTGTT
$S_5=$ AGTCAGG

|       | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ |
|-------|-------|-------|-------|-------|-------|
| $S_1$ |       | 14    | −4    | 0     | −6    |
| $S_2$ | 14    |       | −4    | 0     | −8    |
| $S_3$ | −4    | −4    |       | 0     | −14   |
| $S_4$ | 0     | 0     | 0     |       | −6    |
| $S_5$ | −6    | −8    | −14   | −6    |       |

34

# Step 3 - Add sequences in decreasing order

MPE
| |
MKE

MSKE
| | |
M-KE

$s_1$: MPE
$s_2$: MKE
$s_3$: MSKE
$s_4$: SKE

S₁

S₃

S₂

MKE
| |
SKE

S₄

MSKE
M-KE

M-PE
MSKE
M-KE

S-KE
M-PE
MSKE
M-KE

$S_2+S_3$          $+S_1$          $+S_4$

| S1 | A | T | T | G | C | C | A | T | T |
|----|---|---|---|---|---|---|---|---|---|
| S2 | A | T | G | G | C | C | A | T | T |
| S3 | A | T | C | C | A | A | T | T | T | T |
| S4 | A | T | C | T | T | C | T | T |
| S5 | A | C | T | G | A | C | C |

|       | S1 | S2 | S3 | S4 | S5 |     |
|-------|----|----|----|----|----|-----|
| S1    | -  | 7  | -2 | 0  | -3 | **2** |
| S2    | 7  | -  | -2 | 0  | -4 | **1** |
| S3    | -2 | -2 | -  | 0  | -7 | **-11** |
| S4    | 0  | 0  | 0  | -  | -3 | **-3** |
| S5    | -3 | -4 | -7 | -3 | -  | **-17** |
| **2** | **1** | **-11** | **-3** | **-17** |  |  |

# Star-Alignment Example (Cont.)

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| S₁ | A | T | T | G | C | C | A | T | T | |
| S₂ | A | T | G | G | C | C | A | T | T | |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| S₁ | A | T | T | G | C | C | A | T | T | - | - |
| S₃ | A | T | C | - | C | A | A | T | T | T | T |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| S₁ | A | T | T | G | C | C | A | T | T | |
| S₄ | A | T | C | T | T | C | - | T | T | |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| S₁ | A | T | T | G | C | C | A | T | T | |
| S₅ | A | C | T | G | A | C | C | - | - | |

# Star-Alignment Example (Cont.)

Let's use the alignment of $S_1$ and $S_2$.

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

$S_1$ and $S_2$ are aligned

Now, let's add $S_3$, using its alignment to $S_1$.

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

$S_1$, $S_2$, and $S_3$ are aligned

Then, let's add $S_4$, using its alignment to $S_1$.

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

$S_1$, $S_2$, $S_3$, and $S_4$ are aligned

38

Finally, let's add $S_5$, using its alignment to $S_1$.

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $S_1$ | A | T | T | G | C | C | A | T | T | - | - |
| $S_2$ | A | T | G | G | C | C | A | T | T | - | - |
| $S_3$ | A | T | C | - | C | A | A | T | T | T | T |
| $S_4$ | A | T | C | T | T | C | - | T | T | - | - |
| $S_5$ | A | C | T | G | A | C | C | - | - | - | - |

$S_1$, $S_2$, $S_3$, $S_4$ and $S_5$ are aligned

For consistency, once a gap is added, it is never removed.

(Tehran Polytechnic)

# Complexity of Star Alignment?

Given *k* sequences of length *n*, and an upper bound *l* for alignment length. We need:

- $O(k^2n^2)$ to compute the alignments
- $O(k^2)$ to compute the center
- $O(k^2l)$ to build multiple alignment

Overall: **$O(k^2n^2)$**

Is this really much better than **$O(k^22^kn^k)$?**

YES! *Remember:    k = # of sequences*

*n = length of sequences*

# T-Coffee Program

- T-Coffee (Tree-based Consistency Objective Function for alignment Evaluation) performs progressive sequence alignments as in Clustal.
  - [www.ch.embnet.org/software/TCoffee.html](www.ch.embnet.org/software/TCoffee.html)
- The main difference: T-Coffee performs both global and local pairwise alignment for all possible pairs involved.
- The global pairwise alignment is performed using the Clustal program.
- The local pairwise alignment is generated by the *Lalign* program
  - The top ten scored alignments are selected.
- T-Coffee outperforms Clustal when aligning moderately divergent sequences. However, it is also slower than Clustal.

# https://www.ebi.ac.uk/Tools/msa/tcoffee/

# POA (Partial Order Alignments)

- POA is a progressive alignment program that does not rely on guide trees.
  - The multiple alignment is assembled by adding sequences in the order they are given.
  - A partial order graph is used to represent a growing multiple alignment.
  - [www.bioinformatics.ucla.edu/poa/](www.bioinformatics.ucla.edu/poa/)
- Each time a new sequence is added, it is aligned with every sequence within the partial order graph individually using the Smith–Waterman algorithm.
- POA is local alignment-based and has been shown to produce **more accurate** alignments than Clustal. It is also **faster** than Clustal.
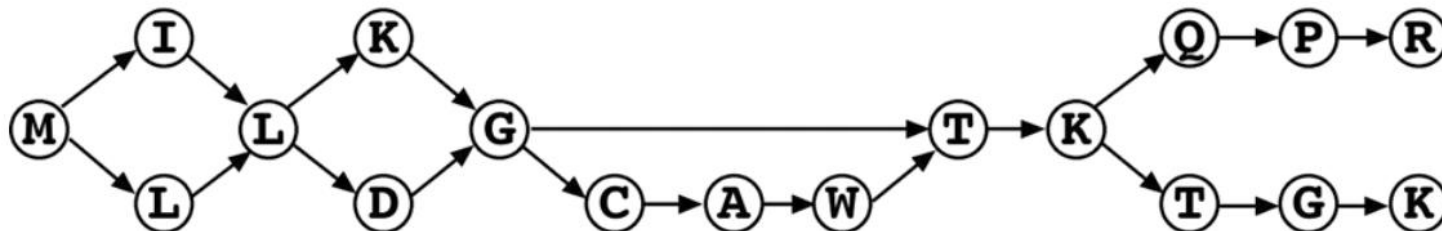
M I L K G T K Q P R
M V L D G C A W T K T G K

Smith-Waterman
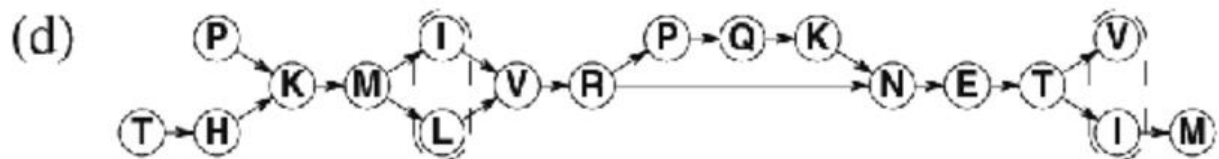alignment

M I L K G - - - T K Q P R
M V L D G C A W T K T G K
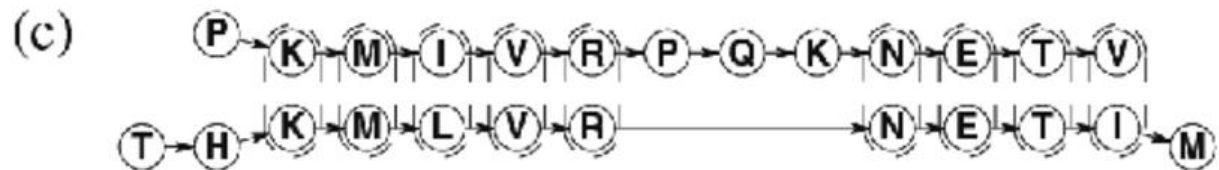
Condense to a
graph profile

- **Conversion of a sequence alignment into a graphical profile in the POA algorithm. Identical residues in the alignment are condensed as nodes in the partial order graph.**
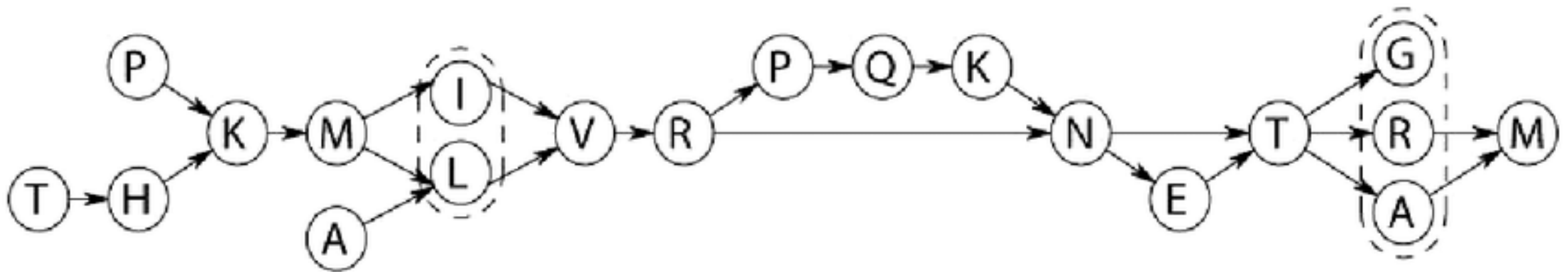
```
.  .  P  K  M  .  I  V  R  P  Q  K  N  E  T  G  .

.  .  .  .  .  A  L  V  R  P  Q  K  N  .  T  R  M

T  H  .  K  M  .  L  V  R  .  .  .  N  E  T  A  M
```

# Iterative Alignment

- Idea: an optimal solution can be found by repeatedly modifying existing suboptimal solutions.

- The procedure starts by producing a low-quality alignment and gradually improves it by iterative realignment.

- The method may reduce the "greedy" problem of the progressive strategy:
  - Because the order of the sequences used for alignment is different in each iteration.

- The method is also heuristic in nature and does not have guarantees for finding the optimal alignment.

- Uses a double nested iterative strategy.
- Two sets of iterations:
- *Outer iteration:*
  – An initial random alignment is generated that is used to derive a guide tree.
  – Weights are subsequently applied to optimize the alignment.

http://prrn.ims.u-tokyo.ac.jp/

48

- *Inner iteration*:
  – The sequences are randomly divided into two groups.
  – Randomized alignment is used for each group in the initial cycle, after which the alignment positions in each group are fixed.
  – The two groups, each treated as a single sequence, are then aligned to each other using global dynamic programming.

http://prrn.ims.u-tokyo.ac.jp/

49

Guide tree

Application of weights

Outer iteration

X

Alignment using dynamic programming

Inner iteration

Alignment score converged

- *Inner iteration (Cont.):*
  - The process is repeated through many cycles until the total SP score no longer increases.
  - At this point, the resulting alignment is used to construct a new guide tree. New weights are applied to optimize alignment scores.
  - The newly optimized alignment is subject to further realignment in the inner iteration.
  - This process is repeated over many cycles until there is no further improvement in the overall alignment scores.

http://prrn.ims.u-tokyo.ac.jp/

# MAFFT (Multiple Alignment using FFT)

A progressive alignment is made then divided into sub-alignments by tree-dependent partitioning. Partitions are re-aligned, then subgroups are aligned. If an objective score improves, this new alignment replaces the initial one and the process may be repeated.



Sequences → Initial alignment → Tree-dependent partitioning → Divide into subalignments → Group-to-group alignment → Replace (if score is improved)

# Block-Based Alignment

- The progressive and iterative alignment strategies are largely global alignment based:
  - May fail to recognize conserved domains and motifs among highly divergent sequences of varying lengths.

- For such divergent sequences that share only regional similarities, a local alignment based approach has to be used.

- The strategy identifies a **block of ungapped alignment** shared by all the sequences, hence, the block-based local alignment strategy.

# DIALIGN

- It does not apply gap penalties and thus is not sensitive to long gaps.
- The method breaks each of the sequences down to smaller segments and performs all possible pairwise alignments between the segments.
- High-scoring segments, called *blocks*, among different sequences are then compiled in a progressive manner to assemble a full multiple alignment.
- The program has been shown to be especially suitable for aligning *divergent sequences* with only local similarity.

# Example

## Dialign-Pfam

### Dialign-Pfam

Dialign-Pfam identifies possible domains in protein sequences by scanning the input sequences using HMMER against PFAM database. It then uses this information to align protein sequences using Dialign.

### Input Sequences

Paste your sequences in multiple FASTA format:

Or, upload your sequences file in multiple FASTA format:

Choose File   No file chosen

Insert Sample    Reset

### Thresholds

HMMER assigns quality scores to matches between sequences and models of proteins and domains in a database. In order to control which hits are used by our algorithm, we use two threshold values for E-values of HMMER hits.

$E_m$: 0.005      $E_d$: 0.0001

Submit

This website is free and open to all users and there is no login requirement.

54

# Protein-Coding DNA Sequences

- Alignment at the protein level is more sensitive than at the DNA level.
- Sequence alignment directly at the DNA level can often result in **frameshift** errors
  – in DNA alignment gaps are introduced irrespective of codon boundaries.
- In the process of achieving maximum sequence similarity at the DNA level, mismatches of genetic codons occur that violate the *accepted evolutionary scenario*
  – insertions or deletions occur in units of codons
- There are occasions when sequence alignment at the DNA level is often necessary, for example, in constructing DNA-based molecular phylogenetic trees.

**Protein alignment**

```
Ser  Ala  Glu
Thr   -   Asp
```

```
AGT  GCA  GAA          AGT  GCA  GAA
ACA  ---  GAT          A--  -CA  GAT

  correct                 incorrect
```

**DNA alignment**

- DNA can be translated into an amino acid sequence before carrying out alignment to avoid the errors of inserting gaps within codon boundaries.

- After alignment of the protein sequences, the alignment can be converted back to DNA alignment.

# Example: RevTrans Program

## RevTrans 1.4 Server

http://www.cbs.dtu.dk/services/RevTrans/

[NOTICE: New improved version is now open for testing: RevTrans 2.0]

RevTrans takes a set of DNA sequences, virtually translates them, aligns the peptide sequences, and uses this as a scaffold for constructing the corresponding DNA multiple alignment.

*New in RevTrans 1.4:* Improvements in the transcription model, restriction on 75 sequences removed, more alignments programs: Dialign 2, Dialign-T and ClustalW, - [Previous version: RevTrans 1.3]

| Instructions | Output format | Background | Software download | Article abstract |
|---|---|---|---|---|

**Paste in DNA sequences**

**Optional: Paste in peptide alignment**

**Upload file containing DNA sequences**

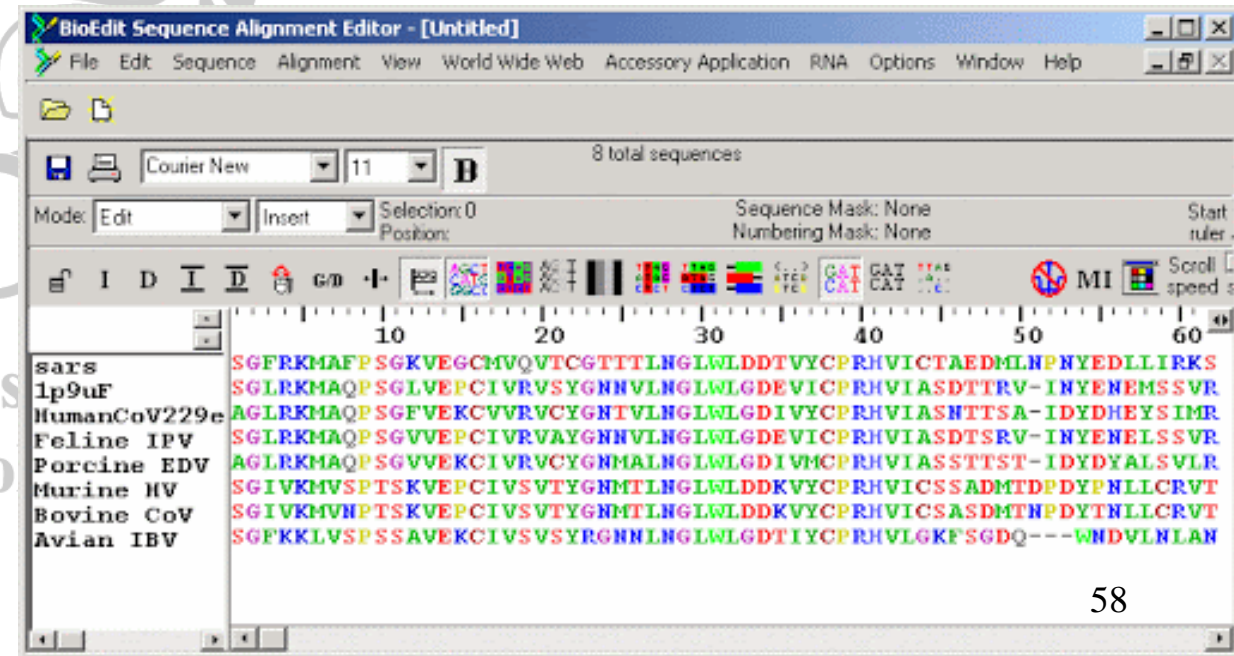Browse...   No file selected.

**Optional: Upload peptide alignment**

Browse...   No file selected.

*Valid formats: FASTA, MSF and ALN (Clustal) - any gaps will be removed from DNA sequences*

Submit query   Clear fields

Translate only   57

# Editing Alignments

- The automated alignment often contains misaligned regions.
  - the user should check the alignment carefully for biological relevance and edit the alignment if necessary.
  - This involves introducing or removing gaps to maximize biologically meaningful matches.

- *BioEdit* is a multifunctional sequence alignment editor for Windows.

- *Rascal* is a web-based program that automatically refines a multiple sequence alignment.

# References

- Mostly used:
  - Essential bioinformatics, Chapter 4 (Multiple Sequence Alignment)
- Second reference:
  - Bioinformatics and functional genomics, Chapter 6 (Multiple Sequence Alignment)

- IP notice: some slides were selected from Drena Dobbs' slides.

# Thanks for your attention