# Bioinformatics II
## Lecture 1: Multiple Sequence Alignment MSA using Dynamic Programming

**Dr Manaf A Guma**

**University Of Anbar- college of Applied sciences-Hit**

**Department of applied chemistry**

1

# What is the MSA?

- It is an alignment of more than 2 sequences.

- Why do we run MSA? Or what is the purposes of MSA?

1. *To highlight conservation and variation.*

2. *How? By identifying the regions of similarity among different species.*

3. *To find the relation among different species.*

4. *To find the profile of sequence from the database.*
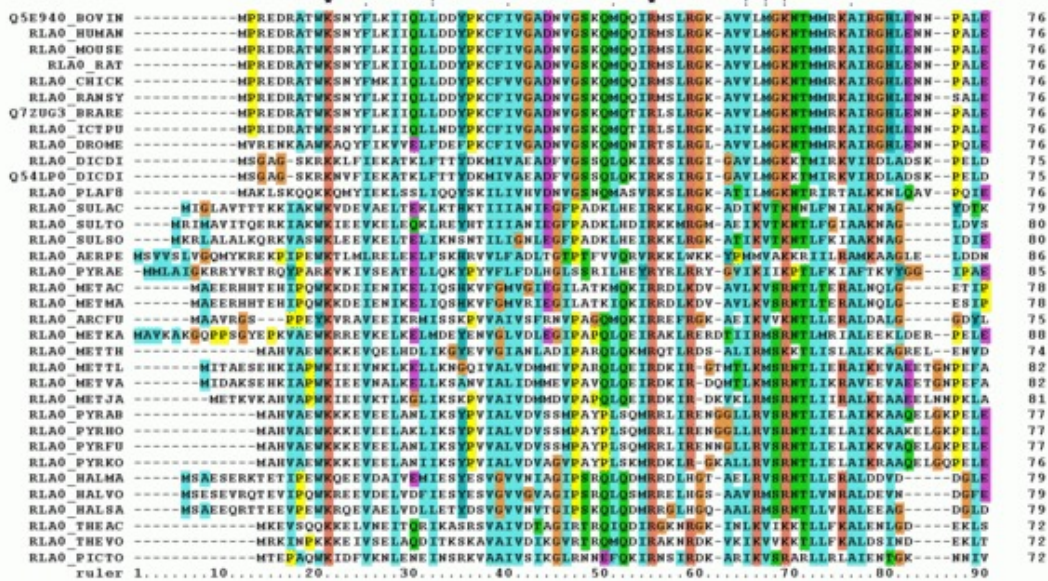
5. *To know how to draw phylogenetic trees.*

2

# The purpose of MSA? *To highlight conservation and variation.*



3

# When you have a huge data!



4

# Why do we use dynamic programming in MSA?

• Because there is a huge database which makes the comparison very difficult if we run MSA by hand.

• Which software and websites are commonly used to do MSA?

1. BLAST.

2. FASTA.

3. ClustalW.

FASTA format) do you remember it !

>AT1G09780|1|training
GTGGAGTAGAAGAATTGAGAGCCTTATCAG
TTTTTGAAGAGAGGGCTGAAACTCTCTAGT
TATCTTTTGTTGCTTTTCTAATAATAAGAG
TTTACACACAG

Part 1

Part 2

Part 3

5

# How do you use BLAST to run MSA? (Tutorial)

1. We need to have a specific sequence for (protein or DNA for a specific species) that we need to find the similarity with it.

2. If we do not have it, then we go to https://www.uniprot.org and then find the Protein seq.

3. Copy the seq (in a FASTA format) do you remember it !

4. Open https://blast.ncbi.nlm.nih.gov/Blast.cgi and find blast protein-protein.

5. Paste the seq in the box labeled with Enter Query Sequence:

6. Click on BLAST to find the similarities.

7. The result will show the comparison (the identity and the scoring of the similarity) of the protein to various proteins in the database.

8. It also shows you the matrices used to generate the comparison.

6

# Can we get MSA form BLAST? What can we get?

- We can get only pairwise alignment using BLAST. (what is pairwise? do you remember?)

- But we can not get all of the sequences aligned together in the same screen using BLAST.

- We can get the profile of each sequence (the type of the species, the gene name and gene number etc.)

7

# An example to see how BLAST works



8

# How do you use FASTA to run MSA?

1. Get the protein/DNA seq from https://www.uniprot.org.

2. copy the seq in FSATA format.

3. Open FASTA web page https://www.ebi.ac.uk/Tools/sss/fasta/.

4. Paste the seq.

5. The results will show different choses to get various bioinformatic analysis in a table.

6. You can show the MSA by clicking on visual output.

7. You can also download the seq by clicking on Download

9

# The tables of FASTA results: an example



You can download all the seq form here

10

# An example to see how FASTA works

# We will carry on with MSA in the next lecture!

• How to get such like this !

# Bioinformatics I
## Lecture 2: Multiple Sequence Alignment
## MSA using Dynamic Programming

**Dr Manaf A Guma**

**University Of Anbar- college of Applied sciences-Hit**

**Department of applied chemistry**

1

# Obtaining MSA using ClustalW databases?

- **ClustalW is the "classic" MSA tool using C++ programming made by JD Thompson, DG Higgins, and TJ Gibson.**

- The original publication describing ClustalW is one of the 100 most cited publications in 'web of science'.

- Why do we run ClustalW? Or what is the purposes of ClustalW?

1. *To find the MSA.*

2. *Again, why do we find MSA? highlight conservation and variation.*

3. *How? By identifying the regions of similarity among different species.*

4. *To find the relation among different species.*

5. *To find the profile of sequence from the database.*

6. *To know how to draw phylogenetic trees.*

2

# How CLUSTAL W deals with MSA??

**CLUSTAL W: deals with multiple sequence alignment through:**

**1. Sequence weighting.**

**2. position-specific gap penalties**

**3. weight matrix choice.**

- **What is the lasts version of** ClustalW?

- ClustalW It is an old version, the version is Clustal Omega which is  much faster and better tools are available. http://www.ebi.ac.uk/Tools/msa/

3

# How do you use ClustalW to run MSA? (very common)

1. Get the protein/DNA seq from https://www.uniprot.org.

2.  copy the seq in FSATA to download multiple seq.

3. Open FASTA web page

https://www.ebi.ac.uk/Tools/msa/clustalo/

1. Paste the multiple seq in the box.

2. Run the FASTA omega. You can color it.

3. You see also the phylogenetic tree as well. p

4

# An example of ClustalW Omega



5

# The old version presentation of the ClustalW



6

# What other programs used for MSA?

Because Often multiple sequence alignments require manual editing:

1.  ***Jalview*** is a powerful MSA-editor for MSA. see

    http://www.jalview.org/index.html

2.  ***Muscle: https://www.ebi.ac.uk/Tools/msa/muscle/***
3.  PRANK: https://www.ebi.ac.uk/research/goldman/software/prank.
4.  MAFFT: https://mafft.cbrc.jp/alignment/software/

7

# What are the benefits of MSA?

1.  Find out which parts "do the same thing"

    Similar genes are conserved across widely divergent species, often performing similar functions.

2.  Structure prediction

    Use knowledge of structure of one or more members of a protein MSA to predict structure of other members

3.  Create "profiles" for protein families

    Allow us to search for other members of the family

4.  Genome assembly: how many gene in this genome.

5.  MSA is to build a phylogenetic analysis.

8

# How to find the most conservative amino acid in a seq among multiple species?

- **Sequence Logos and conservativity can be found using**

- http://weblogo.berkeley.edu/

- Sequence logos are based on **M**ultiple **S**equence **A**lignments
- Very useful to visualise Sequence profiles and motifs.



The DNA-binding helix-turn-helix motif of the CAP family

9

# Tutorial

- Find a TPM1 (tropomyosin) gene for human by typing it in www.uniport.com. Type the gene name

- Go inside the page do click alignment.

- The job will take time.

- Download seq, paste it in
  https://www.uniprot.org/blast/uniprot/B202003208BC4D7ADE02784B0C2481C7F3DE0963A0E5076S

- Download the whole seq, paste it in
  http://weblogo.berkeley.edu/logo.cgi



10

# Bioinformatics II
# Lecture 3: Phylogenetic trees

**Dr Manaf A Guma**

**University Of Anbar- college of Applied sciences-Hit**

**Department of applied chemistry**

1

---

# What is a phylogenetics tree!

- What is a phylogenetics tree?

- It is a diagram draws evolutionary relationships between a set of organism phylogenetic from Darwin notebook.

- It reflects the process of descent (from parents to children) with modification process ancestors to descendant.

- MSA multiped sequence alignment led us to reveal the evolutionary relationships between the species.

- Note:

- A little of sequences and species can be solved by hand.

- But we need a computational statistics to figure out the big groups.

Orangutan   Gorilla   Chimpanzee   Human

*From the Tree of the Life Website, University of Arizona*

2

**Some <mark>identifications</mark> and <mark>notes</mark> related to the topic.**

- **Before we study this topic, we need to know the following information:**

- **Phylogeny** is the study of relationships among different groups of organisms and their evolutionary development. **It** attempts to trace the evolutionary history of all life on the planet.

- <mark>**Taxonomy**</mark> is the science of naming, describing and classifying organisms and includes all plants, animals and microorganisms of the world.

- <mark>Phylogenetics</mark>: the study of evolutionary relatedness among organisms.

- The term living is called (extant) and dead is called (extinct).

- What are the following terms: <mark>Homologous, Paralogous, Analog ad Orthologous genes?</mark>

3

---

**What is a homologous gene (or homolog)?**

- It is a <mark>gene inherited in two species</mark> by a common ancestor.

- Note: the homologous genes can be similar in sequence, similar sequences are not necessarily homologous.

- <mark>What is an Analog?</mark>

- One of two organs or parts in different species of animals or plants which differ in structure or development but are similar in function.



4

# What are orthologous genes (or orthologs)?

- **Orthologous genes** (or **orthologs**) are a particular class of homologous **genes**.

- They are found **in** different species and have diverged following the speciation of the species hosting them.

- Therefore, **orthologous genes in** different species derive from a common ancestral **gene** found **in** the ancestor of those species.

5

# What are Paralogous genes (or paralogs)?

- Paralogous genes (or **paralogs**) are a particular class of homologous genes.

- They are the result of gene duplication.

- The gene copies resulting from the duplication are called paralogous of each other. ...

- Paralogous genes can be remined in the genome after their duplication, but some copies can also be lost.



6

3

**Flow chart to describe Orthologs - Paralogs**

Koonin EV. 2005. Annu Rev Genet. 39:309-38.

Tekaia F. Inferring Orthologs: Open Question and Perspectives (Review). *Genomics Insights* 2016:9 17-28.

7



**Homologs - Paralogs - Orthologs**

Homologs: $A_1$, $B_1$, $A_2$, $B_2$

Paralogs : $A_1$ vs $B_1$ and $A_2$ vs $B_2$

Orthologs: $A_1$ vs $A_2$ and $B_1$ vs $B_2$

Sequence analysis

8

# Phylogeny analyses

1. Starting point to draw the tree needs a set of ==homologous==, aligned DNA or protein sequences.

2. Result of the process: a tree describing evolutionary relationships between the considered sequences i.e.

9

# How many ways to describe the process of reconstructing a phylogeny?

- There have been two different ways and two different philosophies to the process of reconstructing a phylogeny:

- One approach is the **phenetic approach**. In this approach, a tree is constructed by considering the phenotypic similarities of the species without trying to understand the evolutionary pathways of the species.

- The trees constructed via this method are called **phenograms**.

- The second approach is called the **cladistic approach**. via these methods, a tree is reconstructed by considering the various possible pathways of evolution and choosing from amongst these the best possible tree.

- Trees reconstructed via these methods are called **cladograms**.

10

## How does the phylogenetic tree describe the evolutionary relationship?

- By using large DNA sets to build tree underlying principle and logic sequences.

- The sequences are separated by shorter evolutionary distance which are expected to be similar to one another than sequences separated over longer distances.

- How about trees with long evolutionary distances ?

11

---

**Bayesian phylogenetic tree of hominin mitochondrial relationships based on the Sima de los Huesos mtDNA sequence determined using the inclusive filtering criteria.**

▲ Sima de los Huesos
■ Denisovans
■ Neanderthals
■ Africans
■ Asians and Europeans

0.0090

nature

12

## Examples of phylogenetic trees

13

---



**Pace (2001) described a tree of life based on small subunit rRNA sequences.**
Pace, N. R. (1997) *Science* **276**, 734-740

**This tree shows the main three branches described by Woese and colleagues.**

**This tree is referred to as the tree of life or the universal tree.**

14

## Chlamydiae



**Fig. 1.** Phylogeny of chlamydiae. 16*S* rRNA-based neighbor-joining tree showing the affiliation <mark>of environmental and pathogenic chlamydiae with major bacterial phyla</mark>. Arrow, to outgroup. Scale bar, 10% estimated evolutionary distance.

Science 304:728-30.2004.

15

## Eukaryotes
(Baldauf et al., 2000)



16

**Mammalian phylogeney**

**Science 2011; 334:521-24.**

17



18

Chen et al. NAR 34: D363-D368 (2006)

19

**How the analyses of predicted proteomes revealed significant evolutionary processes?**

**Expansion**, **Exchange** and **Reduction**.



**Evolutionary processes include**

**Expansion***
duplication

Ancestor

**Phylogeny***

genesis

**HGT**
Horizontal gene transfer

species genome

**Exchange***Rearrangements***loss** **Reduction***

20

# Bioinformatics II:
## Lecture 4: Building Phylogenies using different methods

**Dr Manaf A Guma**

**University of Anbar- college of applied sciences-Heet.**

**Department of chemistry**

1

---

# What does tree consist of:



1. Root: where tree begins.

2. branch represent different species.

3. Node to separate the branches.

4. Tips of the tree (called leaves as in a plant tree).

5. There are 2D of the tree:

6. The length of the branch to the tips tells the time over.

7. The other dimension does not tell any things about the relationships.



2

## Key features of phylogenetic trees

• **An unrooted tree**

external nodes
branches
external nodes
internal nodes
Hypothetical ancestor

A, B, C, D

• **Rooted trees**

1  2  3  4  5

3

## Rooted and Unrooted trees

• An important feature in phylogenetics:

➡ **trees that make an effect about a common ancestor and the direction of evolution, and those that do not.**

A, B, C, D

A, B, C, D

• In rooted trees a single node is set as a **common ancestor, and a unique path leads from it through evolutionary time to any other node.**

• Unrooted trees **only specify the relationships between nodes and say nothing about the direction in which evolution occurred.**

• Roots can usually be assigned to unrooted trees through the use of an outgroup.

4

2

# Reconstructing trees: A simple example: step 1

- **Choose the taxa.** You decide to study the major clades of vertebrates shown in the leftmost column of the table below.

- (Note that many vertebrate lineages are excluded from this example for the sake of simplicity.)

**Taxa: a classification group.**

Or Taxon.

Moisture ندي او رطب
Amphibians البرمائيات
Retaining الاحتفاظ
Amniotic الذي يحيط بالجنين
Fenestrae النوافذ
Rodents القوارض
Primates الرئيسيات such as human, monkeys ets
ray finned fishes أسماك الراي ذات الزعانف

| | Vertebrae? | Bony skeleton? | Four limbs? | Amniotic egg?* | Hair? | Two post-orbital fenestrae?** |
|---|---|---|---|---|---|---|
| Sharks and relatives | YES | no | no | no | no | no |
| Ray-finned fishes | YES | YES | no | no | no | no |
| Amphibians | YES | YES | YES | no | no | no |
| Primates | YES | YES | YES | YES | YES | no |
| Rodents and rabbits | YES | YES | YES | YES | YES | no |
| Crocodiles and relatives | YES | YES | YES | YES | no | YES |
| Dinosaurs and birds | YES | YES | YES | YES | no | YES |

*amniotic egg: an egg in which the embryo is surrounded by the moisture-retaining amnion membrane

**post-orbital fenestrae: holes in the skull behind the eye

# Step 2 & step 3

- **Determine the characters.** After studying the vertebrates, you select a set of traits, which seem to be homologies, and build the following data table to record your observations.

- (Note that many relevant vertebrate characters are excluded from this example for the sake of simplicity.

- **Determine the polarity of characters.** From studying fossils and outgroups closely related to the vertebrate clade, you hypothesize that the ancestor of vertebrates had none of these features.

- Clade: a group of organisms believed to have evolved from a common ancestor, according to the principles of cladistics.

| | Vertebrae? | Bony skeleton? | Four limbs? | Amniotic egg? | Hair? | Two post-orbital fenestrae? |
|---|---|---|---|---|---|---|
| Ancestor | no | no | no | no | no | no |

# Step 4

- **Group taxa by traits .** We might start out by examining the egg character.
- We focus in on the group of lineages that share the **traits** form of this character, an amniotic egg (A,) and hypothesize that they form a clade (B):

| A | Vertebrae? | Bony skeleton? | Four limbs? | Amniotic egg? | Hair? | Two post-orbital fenestrae? |
|---|---|---|---|---|---|---|
| Sharks and relatives | YES | no | no | no | no | no |
| Ray-finned fishes | YES | YES | no | no | no | no |
| Amphibians | YES | YES | YES | no | no | no |
| Primates | YES | YES | YES | YES | YES | no |
| Rodents and rabbits | YES | YES | YES | YES | YES | no |
| Crocodiles and relatives | YES | YES | YES | YES | no | YES |
| Dinosaurs and birds | YES | YES | YES | YES | no | YES |

B

Sharks  Ray-finned fish  Amphibians  Primates  Rodents & rabbits  Crocodiles  Dinosaurs & birds

Amniotic egg evolved

7

# Then,

- We go through the whole table like this, grouping clades according to traits

| C | Vertebrae? | Bony skeleton? | Four limbs? | Amniotic egg? | Hair? | Two post-orbital fenestrae? |
|---|---|---|---|---|---|---|
| Sharks and relatives | YES | no | no | no | no | no |
| Ray-finned fishes | YES | YES | no | no | no | no |
| Amphibians | YES | YES | YES | no | no | no |
| Primates | YES | YES | YES | YES | YES | no |
| Rodents and rabbits | YES | YES | YES | YES | YES | no |
| Crocodiles and relatives | YES | YES | YES | YES | no | YES |
| Dinosaurs and birds | YES | YES | YES | YES | no | YES |
| | 1 | 2 | 3 | 4 | 5 | 6 |

8

## Step 5

- **Work out conflicts that arise.** There are no conflicts here. Every group is a subset of another group (see C).

- **Build your tree.** Based on the groupings above, you produce this tree:

## Step 7

- **You have made a phylogeny.**

- Of course, this was just an example of the tree-building process.

- Phylogenetic trees are generally based on many more characters and often involve more lineages.

**What are the common Methods to build a phylogenetic tree?**

1. Distance-based measures: (UPGMA and NJ methods)

2. Character Based Methods: (Parsimony: straightforward method and Maximum likelihood). They are statistical measures.

3. Additional Method (Quartets Based and Disc Covering) (we are NOT going to study).

11

# Bioinformatics II:
## Lecture 5: Building a phylogenetic tree using Parsimony methods

**Dr Manaf A Guma**

**University of Anbar- college of applied science-Heet.**

**Department of chemistry**

1

# What is Parsimony methods

- **It is a method taxa (for classification) in ways that minimize the number of evolutionary changes.**

- **What is the Goal of the method?**

- Find the tree that allows evolution of the sequences with the fewest changes.

- This is called a ***most parsimonious*** (MP) tree

- **NOTE: Taxonomists : scientists who do classification through evaluations.**



2

# What is the idea of the Parsimony method?

- It is based on, all other species being equal.

- For e.g: 4 nucleotide changes, is more likely to be true than a complex one. (sub-grouping).

- The method assumes that change in characteristic occurs in lineage over time.

- But, the rate of mutation is NOT constant over time.

- So, leaves (descendent) can have different distances from the root.



3

# Example?  Construct a tree that represent minimum evolutionary changes

- For site 1, the tree that represents minimum changes groups all the (parrots) -a type of birds-showing ''T'' in one group and all those showing ''C'' in another group.



| Species | Sites | | | Species | Sites | | |
|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | | 1 | 2 | 3 |
| | T | A | G | | T | A | G |
| | C | A | A | | C | A | A |
| | T | C | G | | T | C | G |
| | T | C | G | | T | C | G |
| | C | A | A | | C | A | A |

4

**The question: How thee taxa divergence started from an ancestor?**

• The two trees are topologically identical, but they assume different.

ancestral states



5

# 2

• For site 2, there are three possible parsimonious trees.



6

# 2

- All these trees agree on:

- How to cluster parrots in different groups.

- The minimum number of substitution.

- But they disagree on:

- How this evolutionary divergence has occurred in the course of evolution.

- Again, trees are topologically identical, but they assume different ancestral state

7

# 3

- For site 3, there are 2 possible parsimonious tree.



8

## Calculate the number of substitutions steps on each branch?

• So, this depends on the ancestral states.



9



10

## The maximum parsimony tree is the one with minimum total number of substitution.

# Fitch's algorithm: a computer work

- The maximum parsimony tree depends on the sequences of ancestral nodes.

- Fitch's algorithm is used used to build ancestral nodes.

- To count the number of steps of a tree requires at given site:

- Construct a set of nucleotides that are observed there.

- Go down the tree.

- For each ancestral node (say A).

- Consider its two immediate descendent sets (say D1 and D2).

- Nucleotide set in ancestor A

- A=D1 U D2, if D1 and D2 are dis-join (U means Union)

- Else, A=D1 ∩ D2 means **Intersection**



https://www.youtube.com/watch?v=gLM5N__kPiU

13

# Tutorial

- Make a Phylogenetic tree using NCBI?

- https://www.youtube.com/watch?v=ilE-SkJCFOc

14

# Bioinformatics II:
## Lecture 6: RNA structure and function (structure prediction)

**Dr Manaf A Guma**

**University of Anbar- college of applied science-Heet.**

**Department of chemistry**

1

# What is the role of RNA in the life?

- The main cause of the divergence of species and the evolution is the RNA processes e.g splicing etc.
- So, it important to study such process on RNA and the structure and the function of the RNA
- **RNA structure**
  - Levels of organization
  - Energetics
- **RNA types & functions**
  1. Genomic information storage/transfer
  2. Structural
  3. Catalytic
  4. Regulatory

2

# RNA structure

- **We do know in biology that:**

- **One sequence leads to ------A unique structure (usually).**

- **This is true for proteins and RNA molecules? how?**

- Some of RNA molecules are relatively unstructured ( part of mRNA).

- But other RNA have structures such as tRNA ribosomes, enzymes ribo-switches.

- RNA molecules are often high complex and involved in all aspects of gene expression. See examples:

3

# Some examples



Natural Biosensors
(G riboswitch)

- Some examples natural biosensors:

- 1- Its function is to sense the Metabolite in the environment.

- 2- catalytic RNA can cut and stitch RNA and DNA molecules.

- 3-Ribosome; the protein synthesis factor in your body

- See other examples:



**Protein Synthesis Factories**
(ribosomes)

4

## RNA types & functions

| Types of RNAs | Primary Function(s) |
|---|---|
| mRNA - messenger | translation (protein synthesis) regulatory |
| rRNA - ribosomal | translation (protein synthesis)     <catalytic> |
| t-RNA - transfer | translation (protein synthesis) |
| hnRNA - heterogeneous nuclear | precursors & intermediates of mature mRNAs & other RNAs |
| scRNA - small cytoplasmic | signal recognition particle (SRP) tRNA processing                    <catalytic> |
| snRNA - small nuclear snoRNA - small nucleolar | mRNA processing, poly A addition <catalytic> rRNA processing/maturation/methylation |
| regulatory RNAs (siRNA, miRNA, etc.) | regulation of transcription and translation, other?? |

5

# Structured and unstructured RNA

• Some mRNA contains structured regions:

•  this is an example of a mRNA and a (middle part of mRNA molecule)
   contains cording regions: start and stop codons.

## The chemical structure of RNA

### The Chemical Structure of RNA

- It contains 2-OH group would you give Aaron a special power.

- 2-OH makes RNA backbone sticky.

- Also, sugars are not the flat because they are five membered rings.

5' end

The 2'-OH group gives RNA **Special Powers**

*acceptor*     *donor*

**It can make H-bonds**

Put simply....
The 2'-OH Makes the
RNA Backbone **Sticky**

3' end

7

---

## What are the types of bonds in the RNA structure?
## Covalent & non-covalent bonds in RNA

- Mainly there are three structures of RNA molecules:

- **1-Primary:**

-  Covalent bonds

- **2-Secondar and 3-Tertiary both have:**

-  Non-covalent bonds

  - H-bonds

  - (base-pairing)

  - Base stacking

5'AAUUGCGGG
AAAGGGGUCA
ACAGCCGUUC
AGUACCAAGU
CUCAGGGGAA
ACUUUGAGAU
GGCCUUGCAA
AGGGUAUGGU
AAUAAGCUGA
CGGACAUGGU
CCUAACCACG
CAGCCAAGUC
CUAAGUCAAC
AGAUCUUCUG
UUGAUAUGGA
UGCAGUUCA3'

Primary      Secondary      Tertiary

8

## What are the Common structural motifs in RNA?

**Helices**

**Loops**
- Hairpin
- Internal
- Bulge
- Multibranch



9

---

# Some kinds of RNA secondary structure

- In RNA the bases do not stick like spines of cactus

  But some are in a structure or dynamic they can be unpaired in the structure.

- Terminal leg of RNA hairpin have a specific structure



Hairpin loop

Bulge loops

10

# tRNA secondary structure

## 2D



## 3D

11

---

# RNA secondary structure

- Regions of all in a base pairing as it shown in the figure

Definition: Regions of base-pairing

Example and Parts List:



Base-Pairs: The usual suspects (A-U, G-C), PLUS some others

12

## Additional base baring are common in RNA structure

Important example:
**The G-U "wobble" pair**



Compare the G-U pair with an A-U pair:

## RNA Tertiary structure

• It is a specific arrangement of helices and of structure units in 3 dimension.

• RNA has multi-base interactions.

• It shows the minor and major grooves in the 3D.



MINOR GROOVE

MAJOR GROOVE

• 2.6 Å rise/bp
• 11 bp/turn

**RNA**

# What are the RNA structure prediction strategies?

## Secondary structure prediction

1.  **Energy minimization (thermodynamics)**

2.  **Comparative sequence analysis (co-variation).**

3.  **Combined experimental & computational**

15

15

# We will talk about: Energy minimization method

**What are the assumptions?**

**Native tertiary structure or "fold" of an RNA molecule is** (one of) **its "lowest" free energy configuration**(s)

Gibbs free energy $= \Delta G$ **in kcal/mol at 37°C** = equilibrium stability of structure.

lower values (negative) are more favorable

16

16

# What is Gibbs Free energy (G)?

**Gibbs Free energy (G)** is formally defined in terms of state functions **enthalpy** & **entropy,** & state variable, **temperature**

$$G = H - TS$$

$$\Delta G = \Delta H - T\Delta S \quad \text{(for constant temp)}$$

**Enthalpy (H)** = amount of heat absorbed by a system at constant pressure

**Entropy (S)** = measure of the amount of disorder or randomness in a system

Note = this is not the same as "entropy" in information theory, but is related

17

# Free energy minimization
## What are the rules?

```
A    U  Basepair    A=U
A    U  ――――→        A=U
     ΔG = -1.2 kcal/mole
```

**What gives here?**

**Why 1.2 vs 1.6?**

Because it is backward!

```
A    U              A=U
          Basepair
U    A  ――――→        U=A
     ΔG = -1.6 kcal/mole
```

The best configuration when the free energy value les than 0

18

- https://www.youtube.com/watch?v=WCrlm18KQ48&t=26s

- https://www.youtube.com/watch?v=YB13BkjN8RA

- https://www.youtube.com/watch?v=encRU80nOHg

19

# Bioinformatics II:
# Lecture 7: Protein Structure,
# Function and predication

**Dr Manaf A Guma**

**University Of Anbar- college of applied sciences-Heet.**

**Department of chemistry**

# What is a protein?

• Proteins is a biological molecules that our body is built of.

• It consists of sequence of amino acids connected as a polymer '' multi mere'' by peptide bonds

• It is folded in a specific 3D shape to do a specific function in the living cells.

• Proteins have different structures and functions based on the places where they exist in.



Primary Structure of Protein

**Protein: amino acid →peptide → polypeptide →protein**

3



# CLASSIFICATION OF PROTEINS

## Classification Based on Structure
- *Fibrous Proteins*
- *Globular Proteins*
- *Intermediate Proteins*

## Classification Based on Composition
- *Simple Proteins*
- *Conjugated Proteins*

## Classification Based on Functions
- *Structural Proteins, Enzymes, Hormones*
- *Pigments, Transport Proteins, Contractile Proteins*
- *Storage Proteins, Toxins*

www.easybiologyclass.com

4

**How can we study both: Protein Structure & Function of the protein?**

• *Protein structure* - primarily determined by sequence.

• *Protein function* - primarily determined by structure.

• What can that (protein structure and function) be useful in bioinformatics?

• Most of amino acid sequences for specifics regions of a protein are similar to any other protein.

• So, by identifying protein sequences we can predict the structure and the function of a protein.

5

**What are the different structures of a protein?**

1. Primary structures  (simple form)

2. Secondary structures (alpha helix and beta sheets) (3D)

3. Tartary structures

4. Quaternary structures.

They differ by the types of bonds that connect each others.



LEVELS OF PROTEIN STRUCTURE

Primary Structure

Secondary Structure

β-Sheet

α-Helix

Tertiary Structure

Quaternary Structure

Lubrizol Life Science

6

# What are the basic Levels of Protein Structure?



7

# What is the Primary structure of a protein?

- The Primary structures of protein are:

- Linear sequence of amino acids

- This linear sequence is referred to as a polypeptide chain. The amino acids in the **primary structure** are held together by covalent bonds.



8

# What is the secondary structure of protein?

- The Secondary structures of protein are:

- Mostly consists of 2 types which are the α helix and the β pleated sheet.

- Other 2$^{nd}$ str:

- Both **structures** are held in shape by hydrogen bonds, which form between the carbonyl O of one amino acid and the amino H of another

β-pleated sheet

α helix

9

# What are basic differences between a-helix and b-sheet?

β-pleated sheet

α helix

| Alpha Helix | Beta Sheet |
|---|---|
| 1. It is rod like structure, coiled peptide chain arranged in spiral structure. | 1. It is sheet like structure, composed of 2 or more peptide chains. |
| 2. All bindings are intra-chain (on the inside). | 2. Inter-chain between separate polypeptide chains. The intra-chain in a single poly peptide chain folding back on itself. |
| 3. The spiral of a-helix prevents the chain from being fully extended. | 3. The chains are mostly extended and flat, Parallel and anti-parallel. |
| 4. All peptide bonds participate in H bonds between the C=O and N-H. | 4. The H bonds between N backbone and carbonyl group of **adjacent** chain. |
| 5. 3.6 per turn. | 5. number of a.a are Not fixed. |

10

5

# How many types can b-sheet form?



Parellel form

Anti-parellel

11

# What are the tertiary & quaternary structures of protein?

• Tertiary: It is 3D structures shape of protein.

• It has a single polypeptide chain "backbone" with one or more protein secondary structures that form the protein domain.

• Bonds?



12

# What is the Quaternary structure of proteins?

- Quaternary:

- It is an arrangement of multiple folded protein subunits in a multi-subunit complex.

- It involves at least 2 polypeptides (domains).

- It can be a dimer, tetramer, homo or hetero protein.



13

# what are the "Additional" Structural Levels of proteins?

1. **Coils:** type of 2$^{nd}$ structure that are not helices, sheets, or recognizable turns.

2. **Loops**: type of 2$^{nd}$ structure.

3. **Motifs**: combinations of 2$^{nd}$ structural elements

4. **Domains**: combinations of motifs.

- e.g: Globular proteins are built from recurring structural patterns

- Please look for their structure?

14

## What are the types of protein structures based on their folding?

- Folding protein is a process by which a polypeptide chain folds to become a biologically active protein in its native 3D structure.

- Types?

1.   Folded proteins: they are typically stable.

2.   Unfolded proteins: very little.

3.   Partially folded.

4.   *Intrinsically disordered proteins:  unstable protein because they are* dynamic.

- Predicting protein structure and function can be very hard *& fun!*

15

## What are Loops?



1. **Type of 2nd structure.**

2. **Connect helices and sheets.**

3. **They are usually located on surface of structure**

4. **Are more flexible and can adopt multiple conformations**

5. **Tend to have charged and polar amino acids**

6. **Usually involve active sites such as Calcium binding sites**

7. **Some fall into distinct structural families** (e.g., hairpin loops, reverse turns)



16

# What are the most common structural motifs?

1. Helix-turn-helix    e.g., DNA binding **H-T-H**

2. Helix-loop-helix **H-L-H**  e.g., Calcium binding sites in troponin.

3. Coiled-coil: two coils twist over each other.

**H-T-H**

**H-L-H**

17

# Why do we study the protein structure in bioinformatics?

- To be able to predict the structure of the protein by simulation software's using the sequences of known proteins.

- Also, we can do docking, which is used to design drugs, that shows the ability of the protein to interact to the drug.

- Same sequences of a protein may have same function in other proteins.

18

# How can we predict the 3D structure of a protein?

- Protein Structure Prediction  or "Protein Folding" Problem

1.  shows the amino acid **sequence** of a protein.

2.  predict its **3-dimensional structure** (**fold**).

The opposite of this is:

- "Inverse Folding" Problem

- **predict a protein fold.**

- identify every  amino acid sequence that can adopt its 3-D structure.

- Note: The are many websites to predict the 3D Structure **of the protein based on the database.**



19

# Protein structure databases, structural classification & visualization

**Please search for types of proteins in the following webs.**

1.   PDB = Protein Data Bank http://www.rcsb.org/pdb/

2.   (RISC) - several different structure viewers
     http://www.pdg.cnb.uam.es/cursos/Barcelona2002/pages/Farmac/CATH/index.html

3.   MMDB = Molecular Modeling Database

4.   http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Structure
                 **(NCBI Entrez) - Cn3D viewer**

5.   SCOP = Structural Classification of Proteins http://scop.mrc-lmb.cam.ac.uk
           Levels reflect both evolutionary and structural relationships

6.   CATH = Classification by Class, Architecture, Topology and Homology http://www.cathdb.info

**The next steps is to predict the protein structure, please follow…..**

20

- http://www.ebi.ac.uk/interpro/result/InterProScan/iprscan5-R20200704-091515-0179-41313245-p1m/

- http://smart.embl-heidelberg.de

- https://pfam.xfam.org

- For membrane

- http://www.cbs.dtu.dk/services/TMHMM/

- https://embnet.vital-it.ch/software/TMPRED_form.html

- Watch video  https://www.youtube.com/watch?v=8XS0cqxD5XU

- https://www.youtube.com/watch?v=i-u1kJPKUQs&fbclid=IwAR0nhkndR4Hs_SPu5QFHSK6j4MR1FRuH-3uMUNPksbMfpnb06iqlevSw7oA

21

11

# Bioinformatics II:
# Lecture 8: Protein docking

**Dr Manaf A Guma**

**University Of Anbar- college of applied sciences-Heet.**

**Department of chemistry**

---

## Basics Of Molecular Docking

- Docking is a structure-based technique which attempts to find the "best" match, between two molecules.

- WHAT IS MOLECULAR DOCKING?

- In the field of molecular modelling, docking is a method which predicts the preferred orientation of one molecule to a second when bound to each other to form a stable complex.



Target    Ligand    Complex

docking

docking

# Ligand-receptor Complex

- Knowledge of the preferred orientation in turn may be used to predict the strength of <mark>association or binding affinity between two molecules</mark> using for example scoring functions.



3

# Types of docking?

- 1. <mark>Protein-protein docking:</mark>

- Both molecules usually considered rigid

- 2<mark>. Protein-ligand docking:</mark>

- Flexible ligand.

- Rigid-receptor



4

# WHY is Docking Important?

- 1. Signal Transduction:

- The associations between biologically relevant molecules such as proteins, nucleic acids, carbohydrates, and lipids play a central role in signal transduction.

- 2. Drug-Designing and discovery:

- Docking is frequently used to predict the binding orientation of small molecule drug candidates to their protein targets in order to in turn predict the affinity and activity of the small molecule.

- 3. The behaviour of diseases:

- Understanding the mechanism of interaction between proteins molecules which could explain the cause of diseases.

5

# Parts of docking

- 1. Receptor or host or lock: The "receiving" molecule, most commonly a protein or other biopolymer.

- 2. Ligand or guest or key: The complementary partner molecule which binds to the receptor. Ligands are most often small molecules but could also be another biopolymer.

- 3. Docking: Computational simulation of a candidate ligand binding to a receptor.

6

# Parts of docking

- 4. <mark>Binding mode</mark>: The orientation of the ligand relative to the receptor as well as the conformation of the ligand and receptor when bound to each other.

- 5. <mark>Pose</mark> : A candidate binding mode.

- 6. <mark>Scoring</mark> : The evaluating of a particular pose by counting the number of favourable intermolecular interactions such as hydrogen bonds and hydrophobic contacts.

- 7. <mark>Ranking</mark> : the interaction is based on the predicted free-energy of binding.

7

# Introduction to protein-ligand docking

- Protein–ligand docking is a <u>molecular modelling</u> technique.

- The goal of protein–ligand docking is to predict the position and orientation of a <mark><u>ligand</u> (a small molecule, chemical molecule)</mark> when it is bound to a <u>protein</u> receptor or enzyme.

- For example ligand: Paracetamol or any.

- Example of enzyme or protein: lipase, globulin etc.

8

# Example



9

# Computer-aided drug design (CADD)

|  | Known ligand(s) | No known ligand |
|---|---|---|
| **Known protein structure** | **Structure-based drug design (SBDD)**<br><br>Protein-ligand docking | ***De novo* design** |
| **Unknown protein structure** | **Ligand-based drug design (LBDD)**<br>*1 or more ligands*<br>• Similarity searching<br>*Several ligands*<br>• Pharmacophore searching<br>*Many ligands (20+)*<br>• Quantitative Structure-Activity Relationships (QSAR) | **CADD of no use**<br>Need experimental data of some sort |

10

# Protein-ligand docking

- A Structure-Based Drug Design (SBDD) method
  - "structure" means "using protein structure"
- Computational method that mimics the binding of a ligand to a protein
- **Given...**



- **Predicts...**
  - The **pose** of the molecule in the binding site
  - The binding affinity or a **score** representing the strength of binding

11

# Pose vs. binding site

- **Binding site** (or "active site")

  - the part of the protein where the ligand binds

  - generally a cavity on the protein surface

  We must have the crystal structure of the protein bound with a known inhibitor.

- **Pose** (or "binding mode")

  - The *geometry* of the ligand in the binding site

  - Geometry = **location, orientation and conformation**

- *Protein-ligand docking is **not** about identifying the binding site*



12

# Uses of docking

- The main uses of protein-ligand docking are for
  - Virtual screening,
  - Pose prediction

  - **Pose prediction**
  - If we know exactly where and how a known ligand binds...
    – We can see which parts are important for binding
    – We can suggest changes to improve affinity
    – Avoid changes that will 'clash' with the protein



13

# 2 Types of ligand-protein docking

- 2. Protein-ligand docking:

- Flexible ligand.

- Rigid-receptor

- We can classify the various search algorithms according to the degrees of freedom that they consider

- Rigid docking or flexible docking
  - With respect to the ligand structure



14

# Rigid docking

- The ligand is treated as a rigid structure during the docking.

- The DOCK algorithm developed by Kuntz and co-workers is generally considered one of the major advances in protein–ligand docking [Kuntz et al., *JMB*, **1982**, *161*, 269].

- Most docking software treats the protein as rigid

15

# Flexible docking

- **Flexible docking：**

- is the most common form of docking today

    - Conformations of each molecule are generated on-the-fly by the search algorithm during the docking process

    - The algorithm can avoid considering conformations that do not fit

16

# Böhm's empirical scoring function

- In general, scoring functions assume that the free energy of binding can be written as a linear sum of terms to reflect the various contributions to binding.

- Bohm's scoring function included contributions from hydrogen bonding, ionic interactions, lipophilic interactions and the loss of internal conformational freedom of the ligand.

$$\Delta G_{bind} = \Delta G_0 + \Delta G_{hb} \sum_{h-bonds} f(\Delta R, \Delta \alpha)$$
$$+ \Delta G_{ionic} \sum_{\substack{ionic \\ interactions}} f(\Delta R, \Delta \alpha)$$
$$+ \Delta G_{lipo} \left| A_{lipo} \right| + \Delta G_{rot} NROT$$

17

# $\Delta G_o$ calculation

- In analogy with any spontaneous process, protein–ligand **binding** occurs only when the change in Gibbs **free energy** ($\Delta G$) of the system is **negative** when the system reaches an equilibrium state at constant pressure and temperature.

- The **binding free energy** can be **calculated** using the rate constants $k_{on}$ and $k_{off}$ as

- $\Delta G = G$ **bound**-G unbound $= -kTln\ KeqC_0 = -kTln\ C_0 k_{on}/k_{off}$

-  where $K_{eq}$ is the **binding** equilibrium constant, $C_0$ is the reference concentration of 1 mol/L, k is Boltzmann's constant and T is the temperature in Kelvin
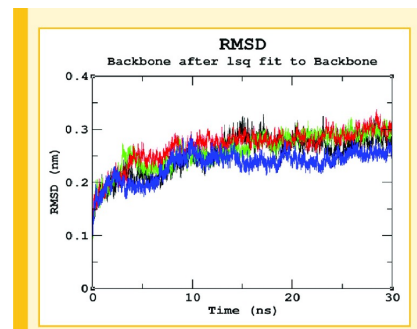
18

## $\Delta G_o$ calculation reflect the various contributions to binding.

- The $\Delta G$ values on the right of the equation are all constants.

- **$\Delta G_o$** is a contribution to the binding energy that does not directly depend on any specific interactions with the protein

- The **hydrogen bonding** and **ionic terms** are both dependent on the geometry of the interaction, with large deviations from ideal geometries (ideal distance R, ideal angle $\alpha$) being penalized.

- The **lipophilic term** is proportional to the contact surface area (Alipo) between protein and ligand involving non-polar atoms.

- The **conformational entropy term** is the penalty associated with freezing internal rotations of the ligand. It is largely entropic in nature. Here the value is directly proportional to the number of rotatable bonds in the ligand (NROT).

19

## Accuracy measured by RMSD

- Accuracy od docking is measured by RMSD (root mean squared deviation) compared to known crystal structures

  - RMSD = square root of the average of (the difference between a particular coordinate in the crystal and that coordinate in the pose)$^2$

  - Within 2.0Å RMSD considered cut-off for accuracy

  - More sophisticated measures have been proposed, but are not widely adopted



20

# Tutorial A: online ligand protein docking

- Step to find the ligand-protein docking: we can use CB-dock a web server for cavity detection guided protein ligand blind docking.

- You need to choose the pdb structure that aim to target the ligand

- from https://www.rcsb.org

- You need also to find the ligand that aim to target your protein from https://pubchem.ncbi.nlm.nih.gov

- Then you can use the server to dock http://clab.labshare.cn/cb-dock/php/

- Follow the link below:

- https://www.youtube.com/watch?v=xFVglxbkoSQ&t=42s

21

# docking programs available

- Large number of docking programs available
  - AutoDock, DOCK, e-Hits, FlexX, FRED, Glide, GOLD, LigandFit, QXP, Surflex-Dock…among others
  - Different scoring functions, different search algorithms, different approaches
  - See Section 12.5 in DC Young, Computational Drug Design (Wiley 2009) for good overview of different packages

- Note: protein-ligand docking is not to be confused with the field of protein-protein docking ("protein docking")

22

# Tutorial B: online protein- protein docking

- You can to the tutorial following the link below:

- https://www.youtube.com/watch?v=8-IPJqXYQ3Q

23