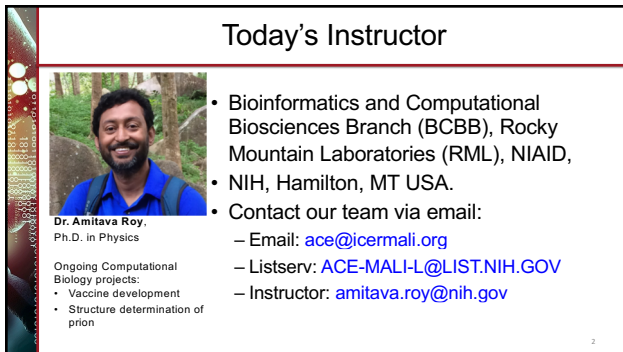AFRICAN CENTERS OF EXCELLENCE
IN BIOINFORMATICS

KAMPALA, UGANDA

**WEB COMPUTATIONAL BIOLOGY TRAINING**

1

## Today's Instructor
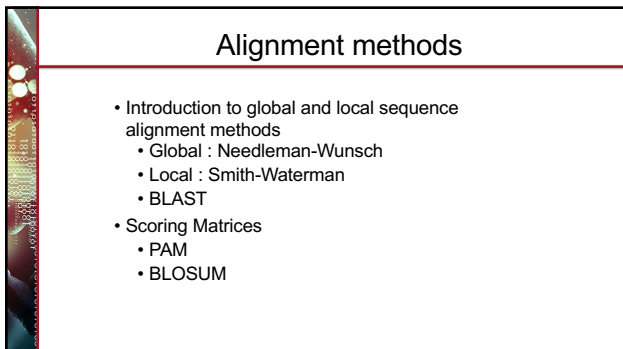
**Dr. Amitava Roy**,
Ph.D. in Physics

Ongoing Computational
Biology projects:
• Vaccine development
• Structure determination of prion

- Bioinformatics and Computational Biosciences Branch (BCBB), Rocky Mountain Laboratories (RML), NIAID,
- NIH, Hamilton, MT USA.
- Contact our team via email:
  – Email: ace@icermali.org
  – Listserv: ACE-MALI-L@LIST.NIH.GOV
  – Instructor: amitava.roy@nih.gov

2

## Alignment methods

- Introduction to global and local sequence alignment methods
  - Global : Needleman-Wunsch
  - Local : Smith-Waterman
  - BLAST
- Scoring Matrices
  - PAM
  - BLOSUM

## Function Prediction

- Multiple Sequence Alignment
  Dynamics Programming
  ClustlW, t-Coffee, Muscle

- Motif Search and Function Prediction
  Expectation Maximization, MEME

- Software

- Markov Model and Hidden Markov Model

## Alignment - Why search sequence databases?

- I have just sequenced something. What is known about the thing I sequenced?
- I have a unique sequence. Is there similarity to another gene that has a known function?
- I found a new protein in a lower organism. Is it similar to a protein from another species?

## Alignment - Perfect Searches

- First "hit" should be an exact match.
- Next "hits" should contain all of the genes that are related to your gene (homologs)
- Next "hits" should be similar but are not homologs

## How does one achieve the "perfect search"?

- Comparison Matrices (PAM vs. BLOSUM)
- Database Search Algorithms
- Databases
- Search Parameters
  - Expect Value-change threshold for score reporting
  - Translation-of DNA sequence into protein
  - Filtering-remove repeat sequences

## Alignment Algorithms

- Global : Needleman-Wunch
- Local : Smith-Watermann
- ➔ These two dynamic programming alignment algorithm are guaranteed to give OPTIMAL alignments
- ➔ But O(m*n) quadratic

## Scoring

- **Quality** = [10(match)] + [-1(mismatch)] - [(Gap Creation Penalty)(#of Gaps) +(Gap Ext. Pen.)(Total length of Gaps)]

Scoring scheme incorporates an evolutionary model--
- Matches are conserved
- Mismatches are divergences
- Gaps are more likely to disrupt function, hence greater penalty than mismatch.
Introduction of a gap (indel) penalized more than extension of a gap.

## Scoring - Estimating $p(\cdot,\cdot)$ for proteins

Generate a large diverse collection of accepted mutations. An *accepted mutation* is a mutation due to an alignment of closely related protein sequences. For example, Hemoglobin alpha chain in humans and other organisms (*homologous* proteins).

Let $p_a = n_a/n$ where $n_a$ is the number of occurrences of letter $a$ and $n$ is the total number of letters in the collection, so $n = \Sigma_a n_a$.

**Mutation counts**

$f_{ab} = f_{ba}$ be the number of mutations $a \Leftrightarrow b$,

$f_a = \sum_{b(b \neq a)} f_{ab}$ be the total number of mutations that involve $a$,

$f = \sum_a f_a$ be the total number of amino acids involved in a mutation.

Note that $f$ is twice the number of mutations.

## Scoring - PAM-1 matrices

Define $M_{ab}$ to be the symmetric probability matrix for switching between $a$ and $b$. We set, $M_{aa} = 1 - m_a$, so that $m_a$ is the probability that $a$ is involved in a change.

$$M_{ab} = \Pr(a \to b) = \Pr(a \to b \mid a \text{ changed}) \cdot \Pr(a \text{ changed}) = \frac{f_{ab}}{f_a} m_a$$

We define $M_{ab}$, such that only 1% of amino acids change according to this matrix or 99% don't. Hence the name, **1-P**ercent **A**ccepted **M**utation (**PAM**). In other words,
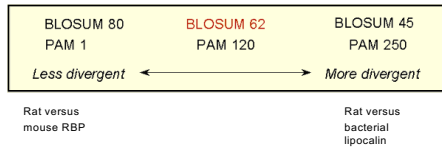
$$\sum_a p_a M_{aa} = \sum_a p_a (1 - m_a) = 1 - \sum_a p_a m_a = 0.99$$

## Scoring - BLOSUM Outline

- **Idea**: use aligned ungapped regions of ***protein families***. These are assumed to have a common ancestor. Similar ideas but better statistics and modeling. It uses 2000 conserved blocks from 500 families.
- **Procedure:**
  - Cluster together sequences in a family whenever more than L% identical residues are shared, for BLOSUM-L.
  - Count number of substitutions across different clusters (in the same family).
  - Estimate frequencies using the counts.
- **Practice**: BIOSUM-50 and BLOSOM62 are widely used.

  Considered the state of the art nowadays.

## Scoring - BLOSUM vs PAM

| BLOSUM 80 | BLOSUM 62 | BLOSUM 45 |
|---|---|---|
| PAM 1 | PAM 120 | PAM 250 |

*Less divergent* ← → *More divergent*

Rat versus
mouse RBP

Rat versus
bacterial
lipocalin

## MSA

- Multiple sequence alignment (MSA)
- Generalize DP to 3 sequence alignment
  - Impractical
- Heuristic approaches to MSA
  - Progressive alignment – ClustalW (using substitution matrix based scoring function)
  - Consistency-based approach – T-Coffee (consistency-based scoring function)
  - MUSCLE (MUSCLE-fast, MUSCLE-prog): reduces time and space complexity

## MSA -From pairwise to multiple alignment

- Alignment of 2 sequences is represented as a 2-row matrix
- In a similar way, we represent alignment of 3 sequences as a 3-row matrix

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

- Score: more conserved columns, better alignment

## What's multiple sequence alignment (MSA)

- A model
- Indicates relationship between residues of different sequences
- Reveals similarity/disimilarity
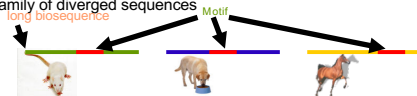
## Why we need MSA

- MSA is central to many bioinformatics applications
- Phylogenetic tree
- Motifs
- Patterns
- Structure prediction (RNA, protein)

## Motif

- Set of similar substrings,
  - within a single long sequence



  - or a family of diverged sequences

## Protein Motif: Activity Sites



..**Y**KFST**Y**AT**WW**IR**Q**AIT**R**..

## Example: Globin Motifs



**Hemoglobin alpha subunit**

## Motif discovery problem

- Given sequences

seq. 1
seq. 2
seq. 3

- Find motif
  - the number of motifs
  - the width of each motif
  - the locations of motif occurrences

**IGR**G**GFGE**V**Y** at position 515
**IG**E**GCFGQ**V**V** at position 430
**VGS**G**GFGQ**V**Y** at position 682

## Why find motifs?

In proteins—may be a critical component
- Find similarities to known proteins
- Find important areas of new protein family

In DNA—may be a *binding site*
- Discover how the gene expression is regulated

## Why is this hard?

Input sequences are long (thousands or millions of residues)

Motif may be *subtle*
- Instances are short.
- Instances may be only slightly similar.

?

?

## Computational problems for *in silico* motif detection

- Extract a motif model based on (experimentally) identified motifs

    **Supervised learning**

- Search for motif instances based on given motif model(s)

    **Prediction**

- Uncover novel motifs computationally from genomic sequences

    **Unsupervised learning**

## Profile based predictors of protein domains / motifs

**prosite**

Motif database in form of regular expressions. Not necessarily the whole domain.

K-x(12)-[DE]   = lysine, any 12, Aspartic acid or Glutamic acid.
Returns 1 or 0, i.e. very rigid and can be very inaccurate for small simple motifs

**PRINTS**

Motif search tools based on Prosite but with multiple alignment profiling

**SMART**

**Pfam**

Collection of HMM's usually covering the whole domain

---

## Exercises:

**Section A:**
- Sequence retrieval of a *P. falciparum* protein (cyclophilin) using SRS
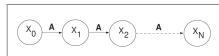- BLAST and Fasta searches by cutting & pasting the sequence.

**Section B:**
**Exercise 1 Part I (row 1):**
- Search PROSITE server by cutting & pasting the cyclophilin sequence
**Exercise 1 Part II (row2):**
- Pfam server
**Exercise 1 Part III (row3):**
- SMART server
**Exercise 1 Part IV (row4):**
- InterPro server
**Exercise 2:**
- Sequence retrieval of *P. falciparum* PFC0125w protein using SRS.
- TMHMMv2.0 server.
- SignalPv3.0 server.
**Section C:**
- Other web resources.

---

## Introducing Hidden Markov Models
## First – a Markov Model

**A Markov Model** is a chain-structured process where future states depend only on the present state, not on the sequence of events that preceded it.

$$X_0 \xrightarrow{A} X_1 \xrightarrow{A} X_2 \xrightarrow{A} \cdots \xrightarrow{A} X_N$$

The X at a given time is called the **state**.
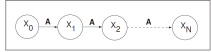The value of Xn depends only on Xn-1.

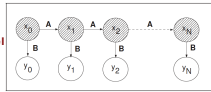**State** :  sunny  cloudy  rainy       sunny ?

## The Hidden Markov Model

**A Hidden Markov Model** is a Markov chain for which the state is only partially observable.

**A Markov Model**



**A Hidden Markov Model**



**Hidden states** : the (TRUE) states of a system that can be described by a Markov process (e.g., the weather).

**Observed states** : the states of the process that are `visible' (e.g., umbrella).