# 3000788 Intro to Comp Molec Biol
## Lecture 2: Sequence alignment

**Fall 2025**

### Sira Sriswasdi, PhD
- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)
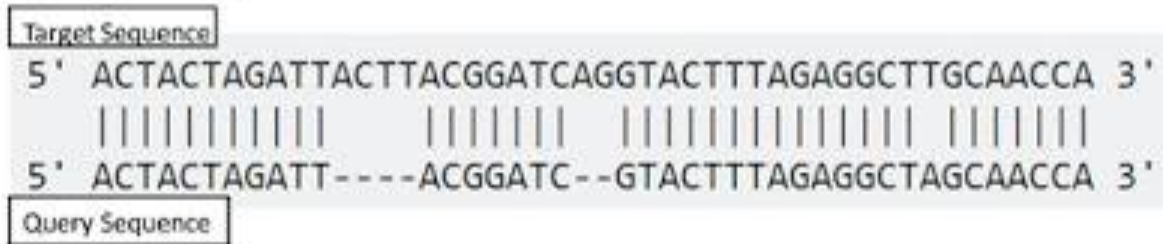
# Today's agenda

- What is sequence alignment?

- Why do we perform sequence alignment?

- Strategies/algorithms for sequence alignment

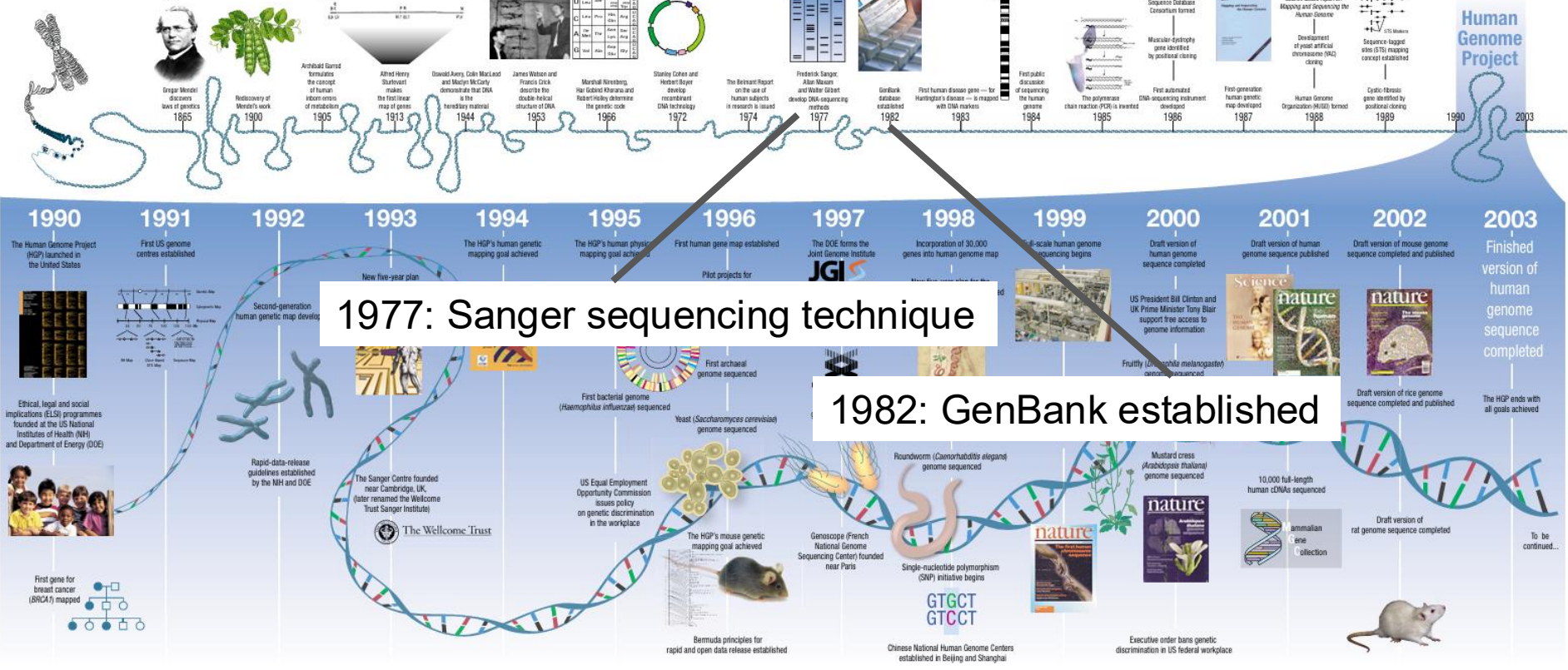# Global and local sequence alignments

## Local Alignment

Target Sequence

```
5'  ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA  3'
                 ||||  ||||||  ||||||||||||||||
Query Sequence 5'   TACTCACGGATGAGGTACTTTAGAGGC  3'
```

## Global Alignment

Target Sequence

```
5'  ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA  3'
    ||||||||||           |||||||   |||||||||||||||  |||||||
5'  ACTACTAGATT----ACGGATC--GTACTTTAGAGGCTAGCAACCA  3'
```

Query Sequence

Images generated from BABA http://baba.sourceforge.net/

Landmarks in genetics and genomics

1977: Sanger sequencing technique

1982: GenBank established

# Basic Local Alignment Search Tool

Stephen F. Altschul[1], Warren Gish[1], Webb Miller[2]
Eugene W. Myers[3] and David J. Lipman[1]

[1]*National Center for Biotechnology Information*
*National Library of Medicine, National Institutes of Health*
*Bethesda, MD 20894, U.S.A.*

[2]*Department of Computer Science*
*The Pennsylvania State University, University Park, PA 16802, U.S.A.*

[3]*Department of Computer Science*
*University of Arizona, Tucson, AZ 85721, U.S.A.*

# An early day of sequence alignment

- Found new cell extracts with interesting functions

- Cloned the genes, but which are responsible for the functions?

- Check the DNA sequences against known genes with validated functional experiments
  - Infer function
  - Infer taxonomy

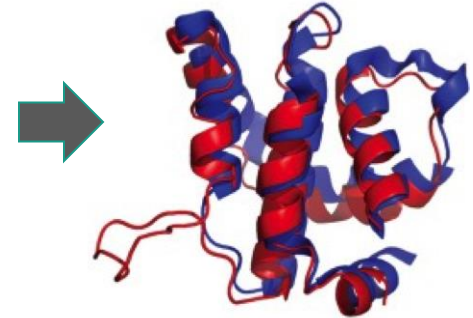# Implications of sequence homology

# Evolution occurs at the sequence level

## Histone H1 (residues 120-180)

```
HUMAN KKASKPKKAASKAPTKKPKATPVKKAKKKLAATPKKAKKPKTVKAKPVKASKPKKAKPVK
MOUSE KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKAKKPKVVKVKPVKASKPKKAKTVK
RAT   KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKAKKPKIVKVKPVKASKPKKAKPVK
COW   KKAAKPKKAASKAPSKKPKATPVKKAKKKPAATPKKTKKPKTVKAKPVKASKPKKTKPVK
CHIMP KKASKPKKAASKAPTKKPKATPVKKAKKKLAATPKKAKKPKTVKAKPVKASKPKKAKPVK
      ***.*********.***************  ******.****  **.***********.*  **
```

https://en.wikipedia.org/wiki/Homology_(biology)

- Evolutionarily-related genes/proteins have similar sequences
- Evolutionarily-related genes/proteins tend to have similar functions
- **Hence, genes/proteins with similar sequences may have similar functions, and may come from closely related organisms**

# Evolution tolerates conservation of function/structure



|  | | α1 | | α2 | | α3 |
|---|---|---|---|---|---|---|
| N1 | 1 | MRTLLIRYILWRNDNDQTQQNDDFKKLMLLDELVDDGDVCTLIKNMRMTL | | | | |
| N2 | 53 | IIAILNRFLTMNKDELNNTQCHIIKEFMTYEQMAIDHYGEYVNAILYQIR | | | | |

| | | α4 | | α5 | |
|---|---|---|---|---|---|
| N1 | 51 | SDGPLLDRLN-------------QPVNNIEDAKRMIAISAKVARDIGERSE | | | |
| N2 | 103 | KRPNQHHTIDLFKKIKRTPYDTFKVDPVEFVKKVIGFVSILNKYKPVYSY | | | |

| | | α6 | | α7 | |
|---|---|---|---|---|---|
| N1 | 90 | IRWEESFTILFRMIETYFDDLMIDLYG | | | |
| N2 | 153 | VLYENVLYDEFKCKINYVETKYF---- | | | |

Ferguson et al. J General Virology, 94: 2070-2081 (2013)

- Sequences can change if the function and structure remain intact
- **Hence, proteins with similar sequences may adopt similar 3D structure, down to the specific domains**

# Molecular probe design

- Sequence alignment confirm the **specificity** of your probes
  - Against targets and potential off-targets

# Sequence alignment predicts many things

- Infer evolutionary relationship across species
  - Many-to-many alignment between gene lists

- Identify the species of origin for a sequence
  - One-to-many alignment against a reference databa
  - Host vs pathogen

- Predict function and structure
  - Partial similarity is good enough
  - Locate conserved functional domain / motif

- Check the specificity of designed probes



https://sites.google.com/site/jkim339n/part2a

# Sequence alignment strategies

# Starting from exact match (seed / word)

..TTACGATAAGC**ATTTTCATAATA**CGACGTCA..

..AACCGACGT**ATTTTCATAATA**GCATAGCAT..

- **ATTTTCATAATA** pattern appears in both sequences
- Do you think the best alignment **must** have this at the center? Why?

....**ATTTTCATAATA**....
....**ATTTTCATAATA**....

# How long of a match can be expected?

- Input sequence length = 300
- Expected 95% similarity (genome re-sequencing)

- How many matches and mismatches to expect? **285 and 15**
    - We can expect to see at least (285/16) **18 matches in a row**

MM...MEM...MEM.........MEM...MM

    - NCBI's MEGABLAST searches for a run of 28 matches

# Dynamic programming

| No. of rocks | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Winner | | | | | | | | | | |
| No. of rocks | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| Winner | | | | | | | | | | ? |

- There is a pile of 20 rocks.
- Two players take turns by removing 1 or 2 rocks from the pile. Whoever removes the last rock(s) win.
- Who is the winner?

# Dynamic programming

|   |   |   |   | B |
|---|---|---|---|---|
|   |   |   |   |   |
|   |   |   |   |   |
| A |   |   |   |   |

- How many ways to travel from A to B using only **right** and **up**?
  - (7 choose 3) = 35 ways

| 1 | 4 | 10 | 20 | **35** |
|---|---|----|----|--------|
| 1 | 3 | 6  | 10 | 15     |
| 1 | 2 | 3  | 4  | 5      |
| 1 | 1 | 1  | 1  | 1      |

# Dynamic programming



- How many ways to travel from A to B using only **right** and **up**?
  - Without running into aliens

|   | 3 | 6 | 10 | **18** |
|---|---|---|----|----|
| 1 | 3 | 3 | 4  | 8  |
| 1 | 2 |   | 1  | 4  |
| 1 | 1 | 1 | 1  | 1  |

# Global sequence alignment



| | GAP | A | T | G | C | T |
|---|---|---|---|---|---|---|
| GAP | 0 | | | | | |
| A | | | | | | |
| G | | | | | | |
| C | | | | | | |
| T | | | | | | |

Match : 1
Mismatch : -1
GAP : -2

Seq1: ATGCT
      | |||
Seq2: A-GCT

# Local sequence alignment = reset when score <0

| | | A | T | G | C | T |
|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 1 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 0 | 1 | 0 | 0 |
| C | 0 | 0 | 0 | 0 | 2 | 0 |
| T | 0 | 0 | 0 | 0 | 0 | 3 |

Match       : 1
Mismatch : -1
GAP          : -2

Seq1: ~~AT~~GCT
           | | |
Seq2:    ~~A~~GCT

- Starting over is better than negative score

# Dynamic programming for sequence alignment

- The best alignment for TTCATA vs TGCTCGTA
  - T/T with the best alignment for TCATA vs GCTCGTA
  - T/– with the best alignment for TCATA vs TGCTCGTA
  - –/T with the best alignment for TTCATA vs GCTCGTA

- Different score for each possibility
  - Match
  - Mismatch
  - Gap

# Scoring of sequence alignment

# Nucleotide alignment scores

```
Ref:  ACCGTATCG
      ||    ||||
Query: AC---ATCG
```

**Scoring Parameters**

| Match/Mismatch Scores | 1,-2 |
| Gap Costs | Linear |

Score = +1+1-1-1-1+1+1+1+1 = +3

**Scoring Parameters**

| Match/Mismatch Scores | 2,-3 |
| Gap Costs | Existence: 5 Extension: 2 |

Score = +2+2-5-2-2-2+2+2+2+2 = +1

- Gap cost models
    - Constant = Same penalty regardless of length
    - Linear = Penalty x Length
    - Affine = **Existence** + (**Extension** x Length)

# Interpretation of match/mismatch scores

- ## Match / Mismatch = +1 / −2
    - A mismatch followed by two matches = no score gain
    - Get hits with **high similarity**
    - Re-sequencing, closely related species

- ## Match / Mismatch = +2 / −3
    - A mismatch followed by two matches +1 score
    - Get hits with **intermediate similarity**
    - More distance species

# Interpretation of gap scores

- Constant = An insertion/deletion can be of any length

- Linear = Long indel is less likely than short indel

- Affine = Existence + (Extension x Length)
  - Combination of constant and linear

# Basic Local Alignment Search Tool: BLAST

# NCBI's nucleotide BLAST interface

# Nucleotide BLAST algorithm



**Program Selection**

Optimize for
- ⦿ Highly similar sequences (megablast)
- ◯ More dissimilar sequences (discontiguous megablast)
- ◯ Somewhat similar sequences (blastn)

Choose a BLAST algorithm ❓

Megablast is intended for comparing a query to closely related sequences and works best if the target percent identity is 95% or more but is very fast.

Discontiguous megablast uses an initial seed that ignores some bases (allowing mismatches) and is intended for cross-species comparisons.

BlastN is slow, but allows a word-size down to seven bases.

- MEGABLAST: word size = 28, match/mismatch = +1/−2, linear gap
- BLASTN: word size = 11, match/mismatch score = +2/−3, affine gap

# MEGABLAST vs BLASTN

MEGABLAST = few, high-identity hits

BLASTN = lots of intermediate-identity hits

# BLAST result



**Query coverage** = % of input sequence used

**Identity** = % of identity in the aligned region

**E value** = expected number of hits with the same or higher score by chance

Typical cutoff is 1e-5

# Understanding E value

- Given an input sequence of length $N$ and a reference sequence of length $M$
- E value for a hit with score $S$ is proportional to $N$ x $M$ x $e^{-\lambda S}$

$N$

Matches with score ≥$S$

$M$

Matches with score ≥$S$

Match with score $2S$

Number of expected hits scales linearly

Matches with score ≥$S$

Number of expected hits scales linearly

Matches with score ≥$S$

$P(\text{score } 2S) = P(\text{score } S) \times P(\text{score } S)$

Two consecutive matches with score $S$

# E value as Poisson distribution

Sequence

Hits with score >*S*

- Event of interest = hits with score >*S* on the sequence of length *N*
- Expected number of events = E value
- Probability of observing *k* hits with score >*S* = $\dfrac{E^k e^{-E}}{k!}$

# Low complexity region

CG island

CCCGCGCGCCCCGGCGCCCGATGCAACTAGC

**Filters and Masking**

Filter ☑ Low complexity regions ❓

Mask regions of low compositional complexity that may cause spurious or misleading results.

- Probability of getting a hit with score >*S* depends on the nucleotide composition (easier with only C and G)
- BLAST withholds these regions from score calculation

# Protein sequence alignment

# Amino acid side chains



https://www.technologynetworks.com/applied-sciences/articles/essential-amino-acids-chart-abbreviations-and-structure-324357

# Similar side chains are tolerated by evolution


wikipedia.com

- D, E have –COOH groups
- K, R have charged –$NH_2$ groups
- A, V, I, L are small hydrocarbon
- F, Y, W have benzene rings

- Alignment score must reflect these

# Block Substitution Matrix (BLOSUM)

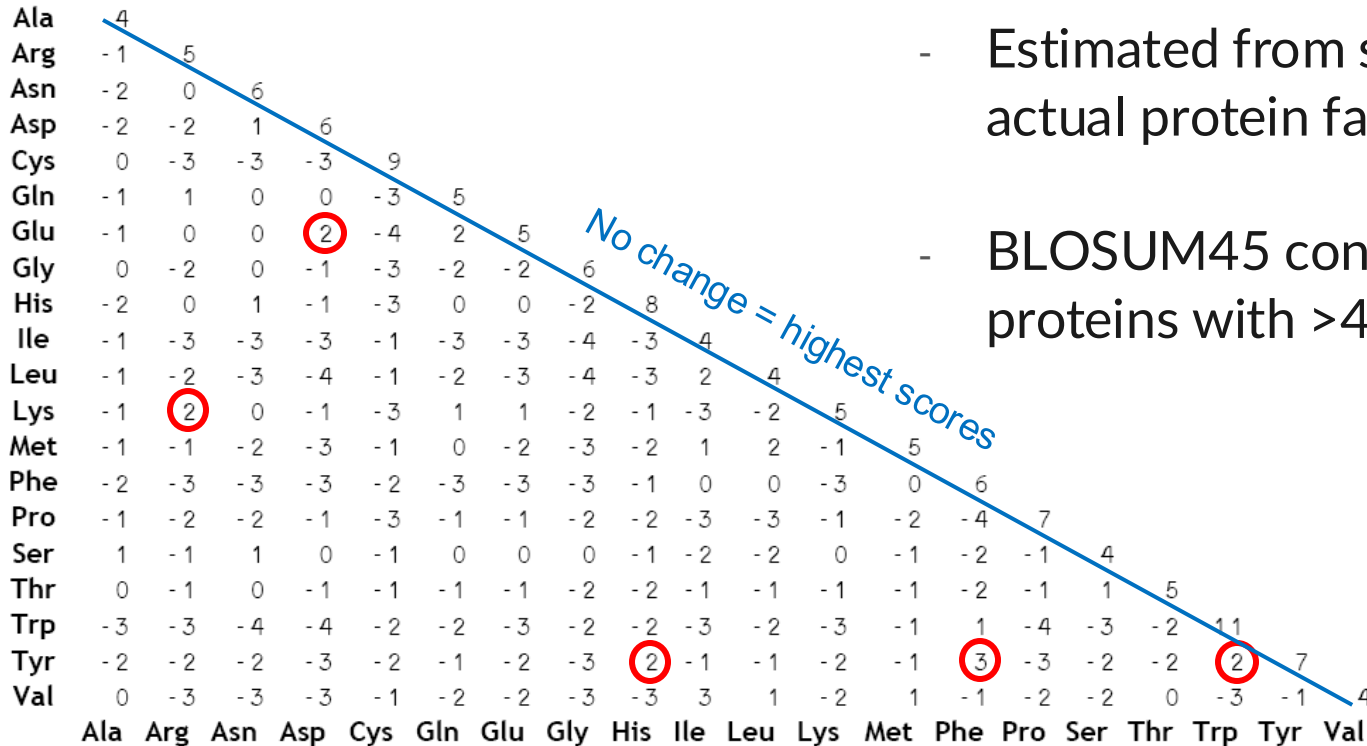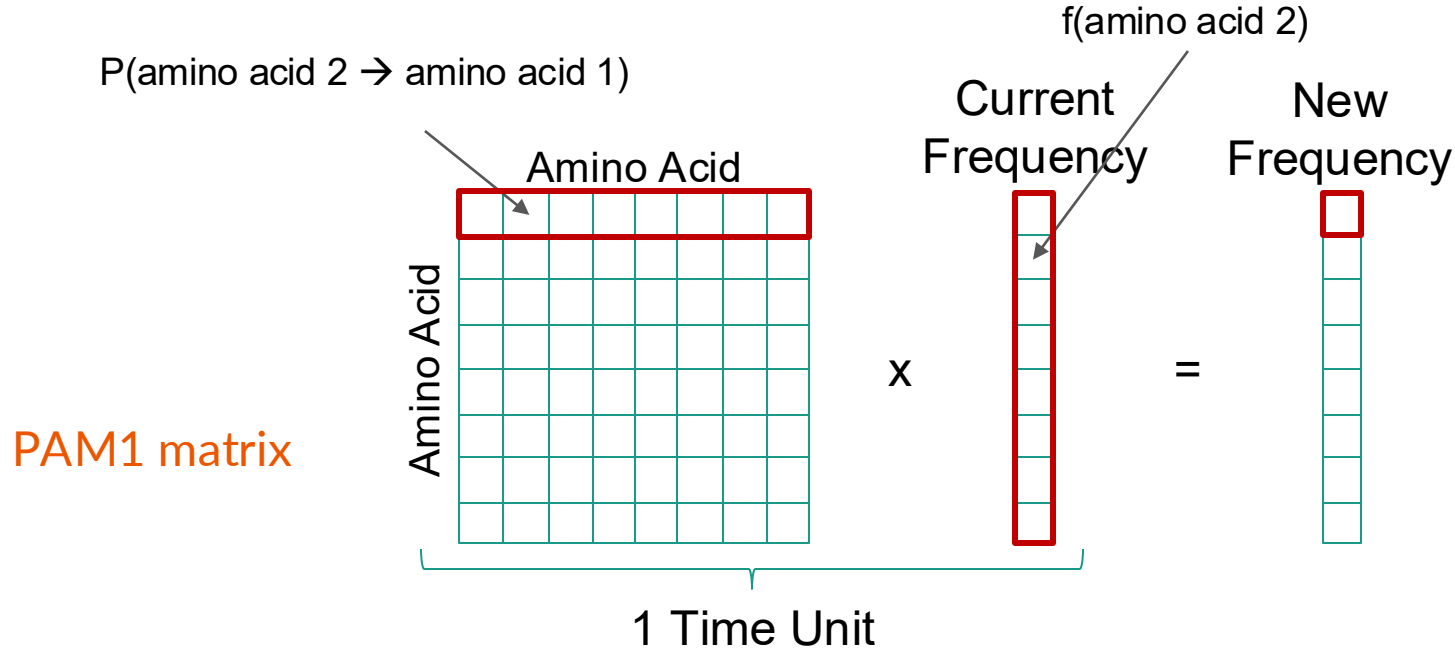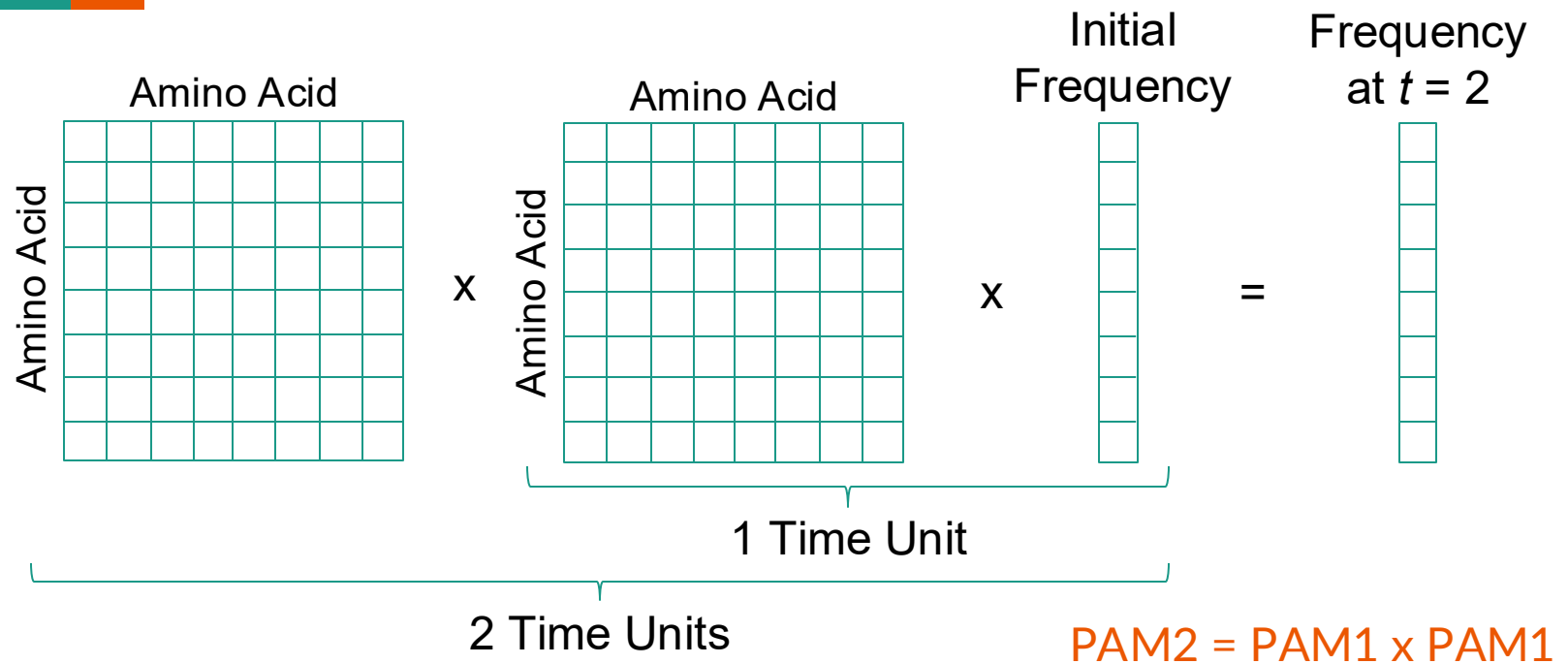|     | Ala | Arg | Asn | Asp | Cys | Gln | Glu | Gly | His | Ile | Leu | Lys | Met | Phe | Pro | Ser | Thr | Trp | Tyr | Val |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Ala | 4   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Arg | -1  | 5   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Asn | -2  | 0   | 6   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Asp | -2  | -2  | 1   | 6   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Cys | 0   | -3  | -3  | -3  | 9   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Gln | -1  | 1   | 0   | 0   | -3  | 5   |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Glu | -1  | 0   | 0   | 2   | -4  | 2   | 5   |     |     |     |     |     |     |     |     |     |     |     |     |     |
| Gly | 0   | -2  | 0   | -1  | -3  | -2  | -2  | 6   |     |     |     |     |     |     |     |     |     |     |     |     |
| His | -2  | 0   | 1   | -1  | -3  | 0   | 0   | -2  | 8   |     |     |     |     |     |     |     |     |     |     |     |
| Ile | -1  | -3  | -3  | -3  | -1  | -3  | -3  | -4  | -3  | 4   |     |     |     |     |     |     |     |     |     |     |
| Leu | -1  | -2  | -3  | -4  | -1  | -2  | -3  | -4  | -3  | 2   | 4   |     |     |     |     |     |     |     |     |     |
| Lys | -1  | 2   | 0   | -1  | -3  | 1   | 1   | -2  | -1  | -3  | -2  | 5   |     |     |     |     |     |     |     |     |
| Met | -1  | -1  | -2  | -3  | -1  | 0   | -2  | -3  | -2  | 1   | 2   | -1  | 5   |     |     |     |     |     |     |     |
| Phe | -2  | -3  | -3  | -3  | -2  | -3  | -3  | -3  | -1  | 0   | 0   | -3  | 0   | 6   |     |     |     |     |     |     |
| Pro | -1  | -2  | -2  | -1  | -3  | -1  | -1  | -2  | -2  | -3  | -3  | -1  | -2  | -4  | 7   |     |     |     |     |     |
| Ser | 1   | -1  | 1   | 0   | -1  | 0   | 0   | 0   | -1  | -2  | -2  | 0   | -1  | -2  | -1  | 4   |     |     |     |     |
| Thr | 0   | -1  | 0   | -1  | -1  | -1  | -1  | -2  | -2  | -1  | -1  | -1  | -1  | -2  | -1  | 1   | 5   |     |     |     |
| Trp | -3  | -3  | -4  | -4  | -2  | -2  | -3  | -2  | -2  | -3  | -2  | -3  | -1  | 1   | -4  | -3  | -2  | 11  |     |     |
| Tyr | -2  | -2  | -2  | -3  | -2  | -1  | -2  | -3  | 2   | -1  | -1  | -2  | -1  | 3   | -3  | -2  | -2  | 2   | 7   |     |
| Val | 0   | -3  | -3  | -3  | -1  | -2  | -2  | -3  | -3  | 3   | 1   | -2  | 1   | -1  | -2  | -2  | 0   | -3  | -1  | 4   |

No change = highest scores

- Estimated from substitution rates of actual protein families

- BLOSUM45 constructed using proteins with >45% conservation

https://en.wikipedia.org/wiki/BLOSUM

# Point Accepted Mutation (PAM)



- Mimic one step of evolution

# Point Accepted Mutation (PAM)

Amino Acid

Amino Acid
<br>(vertical label)

x

Amino Acid

Amino Acid
<br>(vertical label)

x

Initial
Frequency

=

Frequency
at $t = 2$

1 Time Unit

2 Time Units

PAM2 = PAM1 x PAM1

- Extrapolate evolution via multiplication

# PAM vs BLOSUM

| PAM | BLOSUM |
|---|---|
| PAM100 | BLOSUM90 |
| PAM120 | BLOSUM80 |
| PAM160 | BLOSUM60 |
| PAM200 | BLOSUM52 |
| PAM250 | BLOSUM45 |

Data from https://en.wikipedia.org/wiki/BLOSUM



- BLOSUM for low identity, PAM for high identity

# Protein BLAST algorithms

**Program Selection**

Algorithm
- ○ Quick BLASTP (Accelerated protein-protein BLAST)
- ⦿ blastp (protein-protein BLAST)
- ○ PSI-BLAST (Position-Specific Iterated BLAST)
- ○ PHI-BLAST (Pattern Hit Initiated BLAST)
- ○ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

- BLASTP assumes that all amino acid residues are the same
- But each protein domain & motif evolve differently
  - Unknown pattern: PSI-BLAST
  - Known pattern: PHI-BLAST

# Pattern hit initiated (PHI-BLAST)
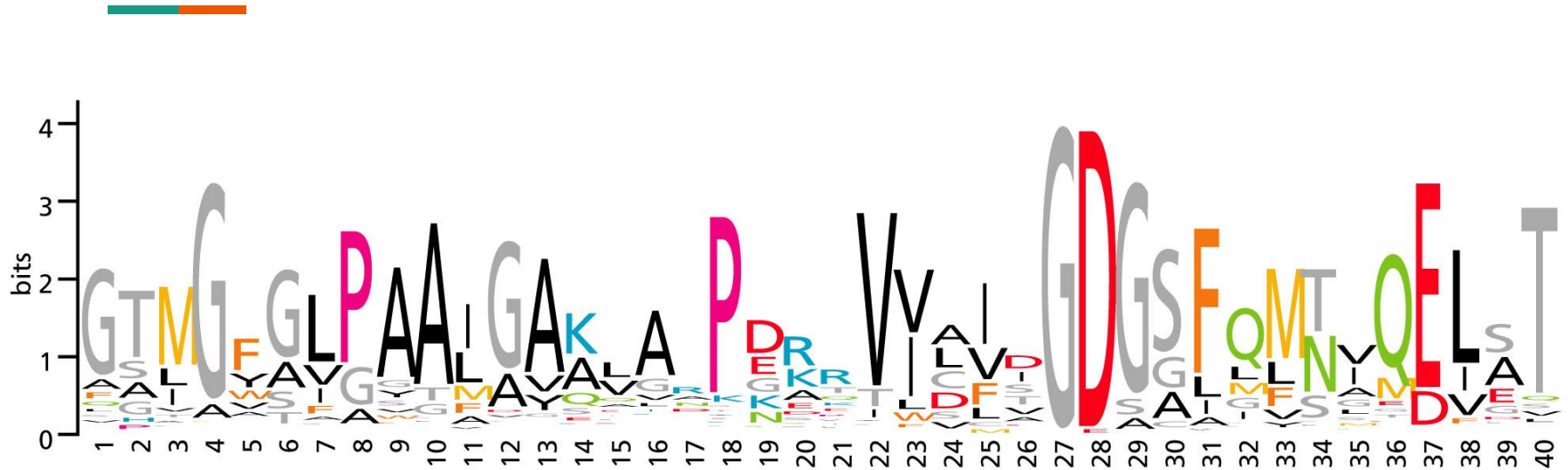
x = any amino acid

`[LIVMF]-G-E-x-[GAS]-[LIVM]-x(5,11)-R-[STAQ]`

L, I, V, M, or F

any sequences of 5-11 amino acids

- Combine regular BLASTP with user-specified pattern
- Hits must be similar to the input sequence AND match the pattern
- Search for known protein domain
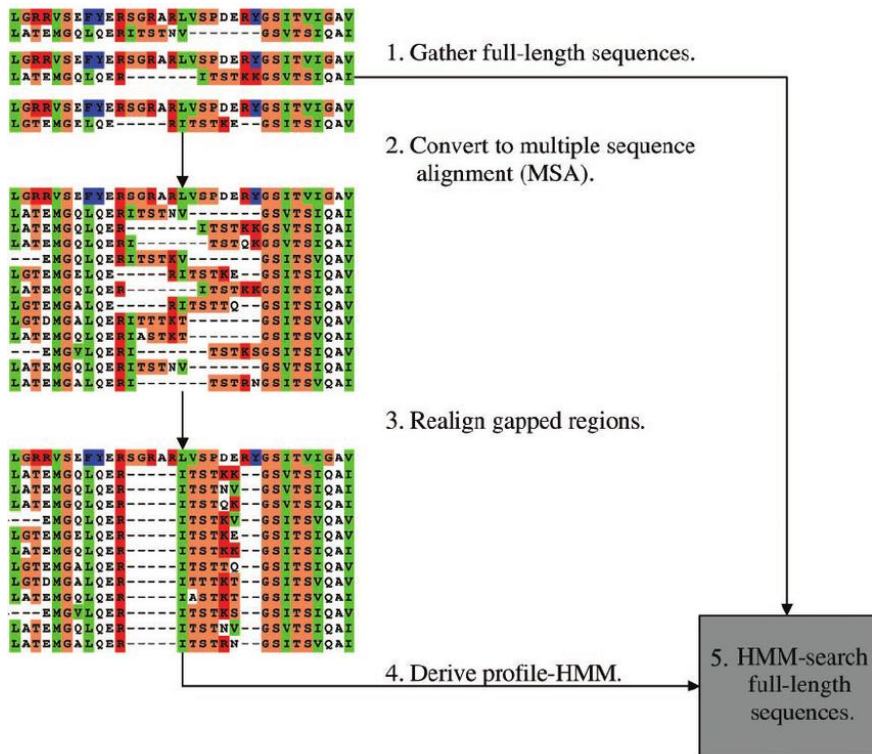
# Position-specific scoring matrix (PSSM)



www.nemates.org/uky/520/Lecture/Lect6/BIO520_2010_Lect6.pp

weblogo.berkeley.edu

- Different scoring for each position in the motif/domain
- How do we know?

# Position-specific iteratred (PSI-BLAST)



- Start from input sequences

- First BLASTP search

- Construct PSSM from hits

- Re-search using the PSSM

- Repeat

Frickey, T. and Lupas, A. NAR 32:5231-8 (2004)

# Using BLASTP to annotate protein function

| Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|---|---|---|---|---|---|---|---|---|
| hypothetical protein JCGZ_15894 [Jatropha curcas] | Jatropha curcas | 1161 | 1161 | 99% | 0.0 | 89.37% | 689 | KDP41487.1 |
| NADPH--cytochrome P450 reductase [Manihot esculenta] | Manihot esculenta | 1159 | 1159 | 100% | 0.0 | 86.98% | 691 | XP_021601058.2 |
| NADPH--cytochrome P450 reductase [Manihot esculenta] | Manihot esculenta | 1145 | 1145 | 100% | 0.0 | 86.25% | 690 | XP_021601060.1 |
| NADPH--cytochrome P450 reductase-like [Hevea brasiliensis] | Hevea brasiliensis | 1130 | 1130 | 99% | 0.0 | 85.59% | 689 | XP_021642755.1 |
| NADPH--cytochrome P450 reductase [Ricinus communis] | Ricinus communis | 1124 | 1124 | 99% | 0.0 | 84.64% | 692 | XP_002514049.1 |
| LOW QUALITY PROTEIN: NADPH--cytochrome P450 reductase-like [Hevea brasilien… | Hevea brasiliensis | 1120 | 1120 | 100% | 0.0 | 84.81% | 698 | XP_021660128.1 |
| hypothetical protein COLO4_35252 [Corchorus olitorius] | Corchorus olitorius | 1111 | 1111 | 100% | 0.0 | 82.08% | 1505 | OMO57587.1 |
| Flavodoxin [Corchorus capsularis] | Corchorus capsularis | 1093 | 1093 | 100% | 0.0 | 82.08% | 692 | OMO50775.1 |
| NADPH--cytochrome P450 reductase-like [Hibiscus syriacus] | Hibiscus syriacus | 1085 | 1085 | 100% | 0.0 | 81.24% | 693 | XP_039050423.1 |
| hypothetical protein CXB51_011412 [Gossypium anomalum] | Gossypium anomalum | 1083 | 1083 | 100% | 0.0 | 81.10% | 694 | KAG8494022.1 |
| NADPH:cytochrome P450 reductase [Gossypium hirsutum] | Gossypium hirsutum | 1083 | 1083 | 100% | 0.0 | 81.24% | 693 | ACN54323.1 |
| NADPH--cytochrome P450 reductase-like [Gossypium hirsutum] | Gossypium hirsutum | 1083 | 1083 | 100% | 0.0 | 81.10% | 693 | NP_001313876.2 |

- Suspected novel CYP reductase from an indigenous plant
  - >80% similarity to known and predicted CYP reductase class I
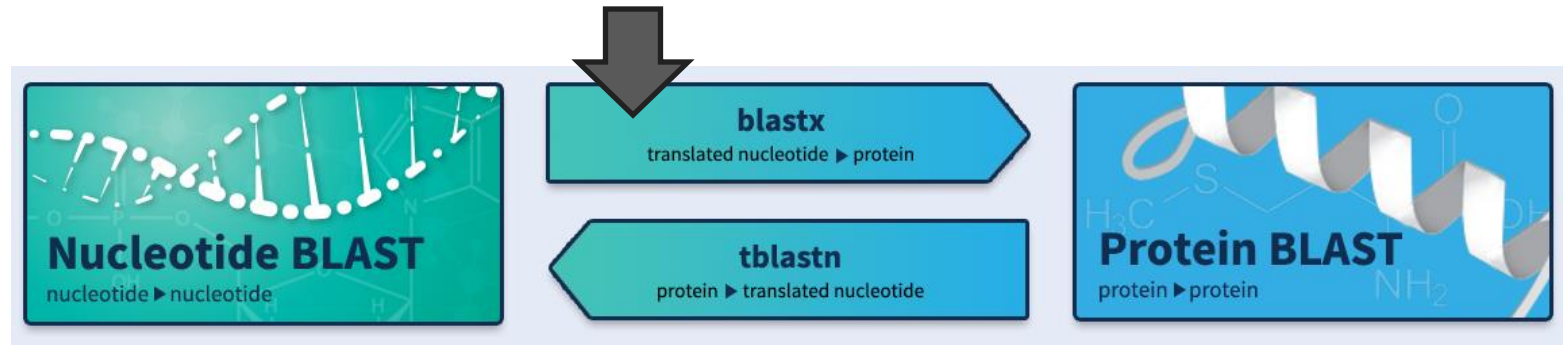
# InterPro: Protein domain search



- Use patterns from known protein functional domains

# Mixed protein-nucleotide alignment

# BLASTX and TBLASTN



- For coding DNA, alignment at protein level is more informative
  - Codon structure, not all amino acid change is equally likely
- **BLASTX** = align translated DNA to protein database
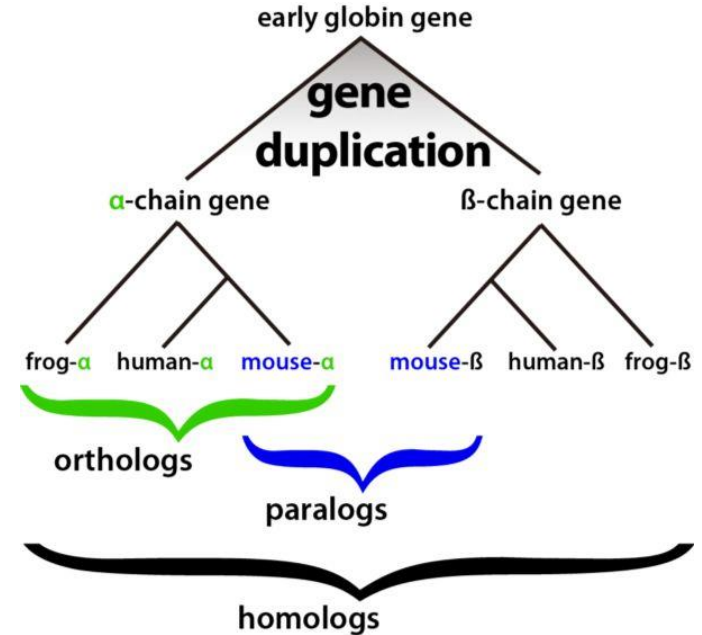- **tBLASTN** = align protein to translated DNA database

# Example use cases

- **BLASTX**
    - Get RNA sequence from RNA-seq
    - Unsure which **open reading frame** is translated
    - Does this RNA translate to a known protein?

- **tBLASTN**
    - Identified novel protein, with no similarity to protein database
    - Is there a transcriptomics study that have identified the RNA of a related protein?

# More advanced alignments

# All-vs-all search for finding orthologs across genomes

- Evolutionarily related genes
  - Group 1 = {Mouse-a, Human-a}
  - Group 2 = {Mouse-b, Human-b}

- BLAST mouse to human
- BLAST human to mouse

- Reciprocal best hit:
  - Human-a is the best hit for Mouse-a
  - Mouse-a is the best hit for Human-a



https://sites.google.com/site/jkim339n/part2a

# Multiple sequence alignment (MSA)



Edgar, BMC Bioinformatics, 5, 113 (2014)

- Dynamic programming is not feasible. Too many ways of grouping sequences (on top of aligning positions).

# Heuristic algorithms

Growing possibilities

- Greedy algorithm optimize the next step and/or the future steps

- Randomized algorithm makes a lot of random decisions and keeps the best one found

# Alignment output



ClustalW

Aligned FASTA

- Dashes "-" are added to indicate insertions/deletions
- All output sequences have the same length

# Any question?

- See you next time