

Spring 2024
BIOS 477/877
Bioinformatics and Molecular Evolution
Lecture 8

BIOS477/877 L8 - 1

1

TODAY'S TOPICS

- Gap Penalty
- Scoring (Substitution) Matrix
- Assignment 4

BIOS477/877 L8 - 2

2

Alignment Strategy

- **Protein alignment is easier** than DNA alignment
 - DNA has only 4 nucleotide types (they can match just by chance more easily)
 - Protein sequences evolve more slowly than DNA sequences (genetic code is redundant; nonsynonymous substitutions are less frequent than synonymous substitutions)
- **If DNA sequences are from coding regions:**
 - Translate them, and **align at the protein level** first
 - This ensures gaps inserted between codons (prevents insertion of frame-shifting gaps)
- **Do not blindly rely on the default parameter set**
 - Try various **scoring matrices, gap penalties, etc.**

BIOS477/877 L8 - 3

3

DNA alignment at protein level

1. Translate DNA sequences to amino acid sequences
2. **Align them at the protein level**
3. Reverse translate the protein alignment to DNA alignment

- **TranslatorX** <http://translatorx.co.uk/>
- **RevTrans 2.0** <https://services.healthtech.dtu.dk/services/RevTrans-2.0/>
- **PAL2NAL** <http://www.bork.embl.de/pal2nal/>
- **tranalign** <https://www.bioinformatics.nl/cgi-bin/emboss/tranalign>
(included in EMBOSS servers; see the course/Link page.)

BIOS477/877 L8 - 4

4

Global vs. local alignments

Global alignment
(Semi-global alignment)

Local alignment

$s(a_i, b_j) = 2$ where $a_i = b_j$
 $s(a_i, b_j) = -1$ where $a_i \neq b_j$

} **Match/Mismatch scores**

$w = -2$
 $(w=0 \text{ for free end gap})$

} **Gap penalty**

BIOS477/877 L8 - 5

5

Indel Evolution and Gap Penalty

➤ **A gap of length $k \neq k$ gaps of length 1**

H
A T T C C G
deletion ↓
F
A T C C G
deletion ↓
C
A C C G
deletion ↓
A C G

Multiple insertion/deletion events

or

A T T C C G
deletion ↓
A C G
T T C

Single insertion/deletion event

Which is more likely?
Which is biologically easier?

BIOS477/877 L8 - 6

6

Indel Evolution and Gap Penalty

- Indel mutations are often strongly deleterious
- Indel events are rare (less common than point mutations)
- Multi-residue indels are not uncommon (e.g., hotspot, repetitive DNA)

Replication slippage

From Human Molecular Genetics 2 (available in NCBI Bookshelf)

Unequal crossover

BIOS477/877 L8 - 7

7

Indel Evolution and Gap Penalty

- Indel mutations are often strongly deleterious
- Indel events are rare (less common than point mutations)
- Multi-residue indels are not uncommon
- **Fewest number of unlikely events**
→ most likely evolutionary hypothesis

Maximum parsimony

AATCTATA
AATCTATA
2 indels
1 indel
AA-G-ATA
AA--GATA
(more likely than 2 indel events)

BIOS477/877 L8 - 8

8

Indel Evolution and Gap Penalty

➤ Fewer, but longer, indel event is more likely than too many small indels

ATTCCG
↓ deletion
ACG

Single insertion/deletion event

During alignments

Extension of existing gaps
should cost less than
creation (or opening)
of a new gap

Two types of gap penalties are required

BIOS477/877 L8 - 9

9

Gap Penalty Functions

- **Linear (length-proportional) gap penalty:**
 $w(x) = gx$
 g : gap penalty
 x : length of a gap
- **Affine gap penalty:**
 $w(x) = \begin{cases} g_o + g_e(x - 1) & \text{when } x > 0 \\ 0 & \text{when } x = 0 \end{cases}$
 g_o : gap opening penalty
 g_e : gap extension penalty [usually $g_o > g_e$]
 x : length of a gap

BIOS477/877 L8 - 10

10

Gap Penalty Functions

Affine gap penalty:

Extension of existing gaps costs less than creation/opening of new gaps

$(g_o$: gap opening penalty, g_e : gap extension penalty)

BIOS477/877 L8 - 11

11

Simple Alignments

➤ Varied length & gaps considered

AAATCTATA
AA-G-ATA
↑↑↑

3

AAATCTATA
AA--GATA
↑↑↑

3

AAATCTATA
AAG-AT-A
↑↑↑

1

- Alignment Score =
(match score) x (the number of matched pairs) +
(mismatch score) x (the number of mismatched pairs) +
(gap penalty) x (the number of gaps)

If match score = 1, mismatch score = 0, gap penalty = -1
& if using linear gap penalty,

BIOS477/877 L8 - 12

12

Simple Alignments

➤ Varied length & gaps considered

2 gap events

1 gap event
Length 2

2 gap events

• Alignment Score =
 (match score) x (the number of matched pairs) +
 (mismatch score) x (the number of mismatched pairs) +
 Σ{(for each gap event) → {(gap opening penalty) +
 (gap extension penalty) x (gap length - 1)}}

If match score = 1, mismatch score = 0,
 & if using affine gap penalty
 gap opening penalty = -2, gap extension penalty = -1

BIOS477/877 L8 - 13

13

Simple Alignments

➤ Varied length & gaps considered

1

2

-1

• Alignment Score =
 (match score) x (the number of matched pairs) +
 (mismatch score) x (the number of mismatched pairs) +
 Σ{(for each gap event) → {(gap opening penalty) +
 (gap extension penalty) x (gap length - 1)}}

If match score = 1, mismatch score = 0,
 & if using affine gap penalty
 gap opening penalty = -2, gap extension penalty = -1

BIOS477/877 L8 - 14

14

Empirical Indel Distribution: DNA

Based on the comparisons of >1700 processed pseudogenes against their functional homologues in the human genome

number of insertions/deletions (n_k)

gap length (k)

Deletions are more frequent than insertions

Zhang and Gerstein (2003)

BIOS477/877 L8 - 15

15

Empirical Indel Distribution: DNA

Frequency

Indel length (bp)

Based on the comparisons of 23 noncoding region sequences between *Drosophila simulans* and *D. sechellia*

Keightley and Johnson (2004) *Genome Res* 14:442-450

This distribution was later used in MCALIGN2:
 Wang et al. (2006) *BMC Bioinformatics* 7: 292.

Figure 2 The empirical distribution of indel lengths in noncoding DNA between *D. simulans* and *D. sechellia*, and the indel length frequency distribution model assumed for the MC analysis.

BIOS477/877 L8 - 16

16

Empirical Indel Distribution: Protein

% of Pairs

Gap Length

Based on the comparisons of 4,952 protein pairs from human, mouse, and rat.

Sequences were aligned by a dynamic programming method...

Chang and Benner (2004) *J Mol Biol* 341:617-631

BIOS477/877 L8 - 17

17

Gap Penalty Function (more realistic)

➤ Empirical indel size distributions (both for DNA and proteins) can be described by a power law:
 $f_k = Ck^{-b}$ [k : indel size, b : the power parameter]

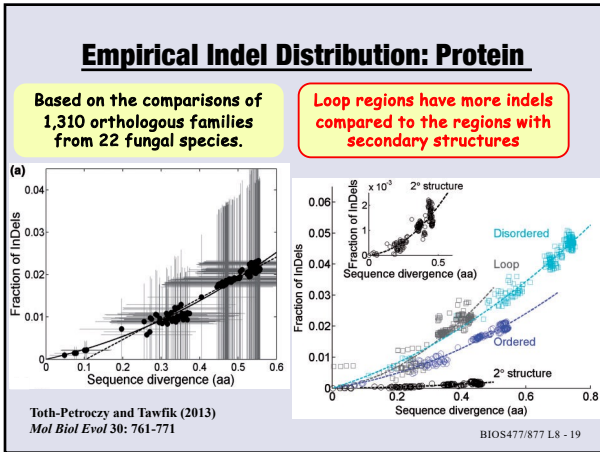
→ Corresponding gap penalty function
 $w = a + b \ln(k)$
 a : gap opening penalty
 b : gap extension penalty

- Gap extension penalty is proportional to the logarithm of gap length k (logarithmic gap penalty system)

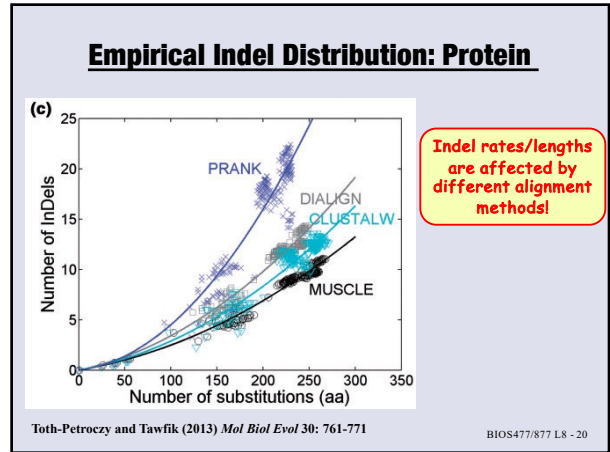
→ increases more slowly with gap length than in the affine gap penalty system (easier long gaps)
 (e.g., Cartwright 2006; Cartwright 2009; Loewenthal et al. 2021)

BIOS477/877 L8 - 18

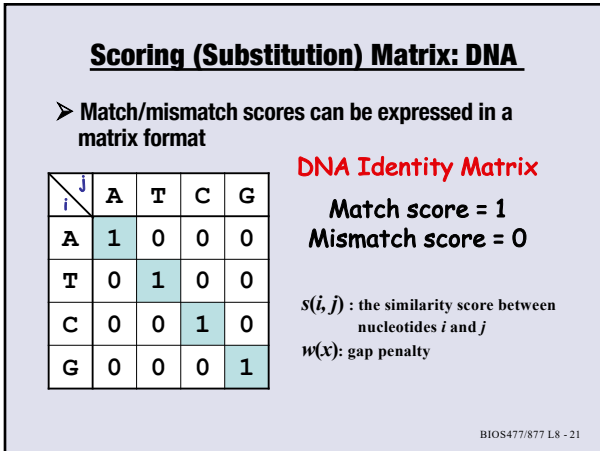
18



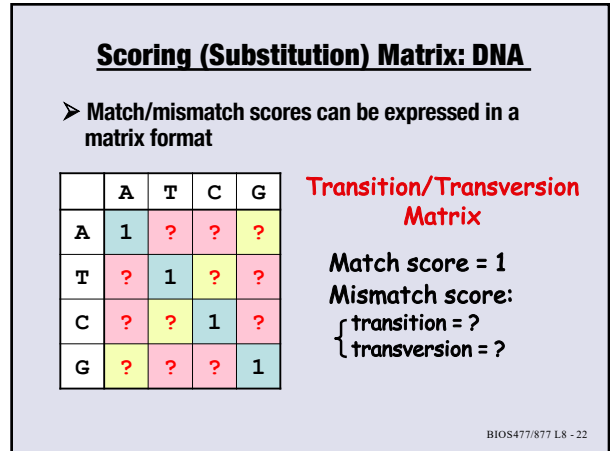
19



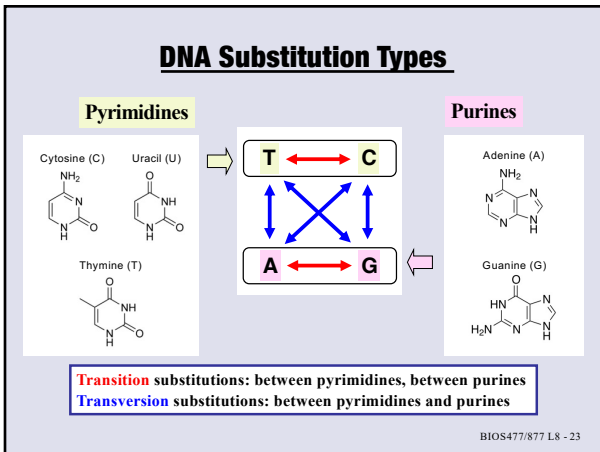
20



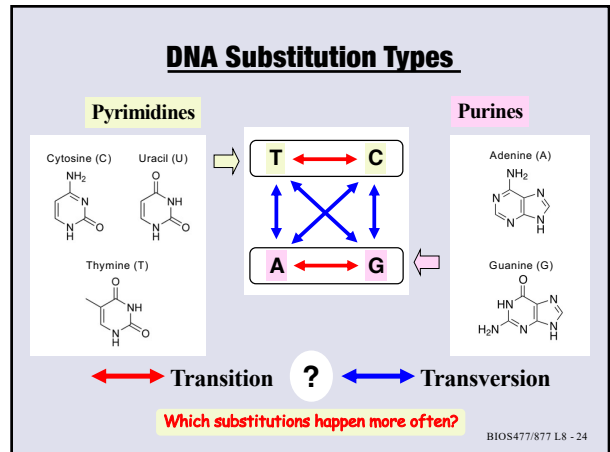
21



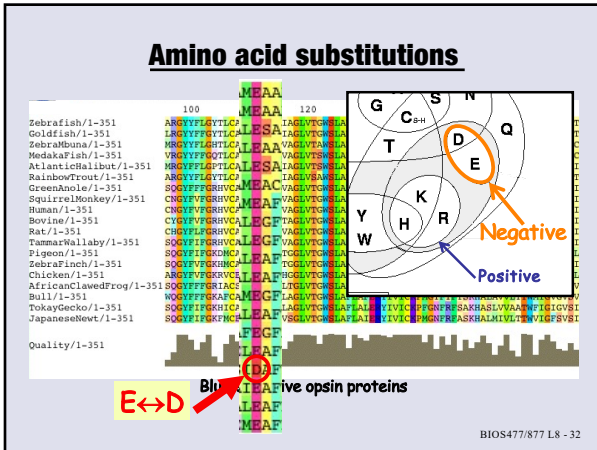
22



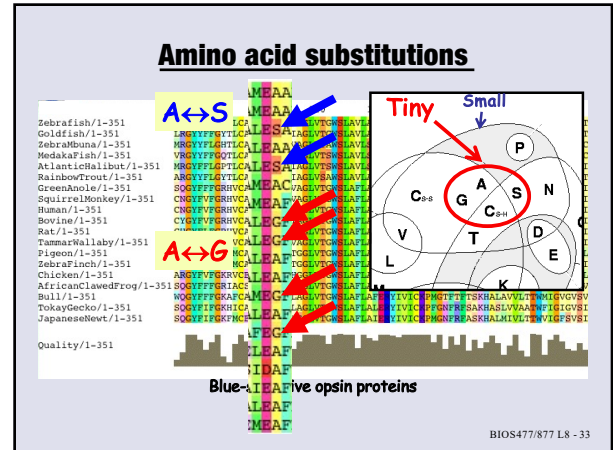
23



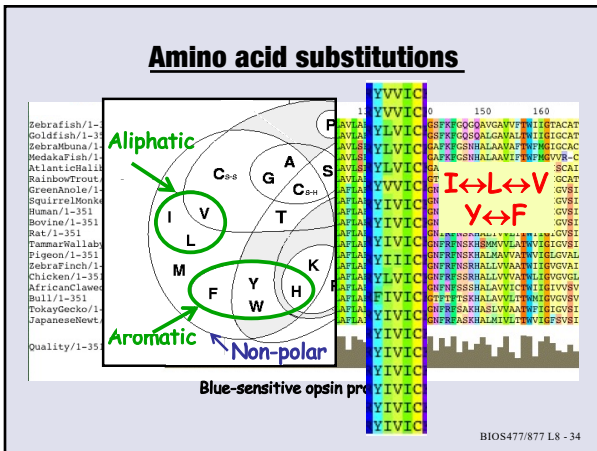
24



32



33



34

Amino acid substitution matrices

- Matrices based on various amino acid properties (hydrophobicity, charge, electronegativity, size, etc.)
 - Biologically meaningful matrix can be obtained by combining all of these matrices (including genetic code matrix). Not easy!
- Matrices based on empirical data
 - Alignments show the results of experiments done by the Nature
 - Capture the relative substitutability of amino acid pairs in the context of evolution
 - The model of protein evolution

BIOS477/877 L8 - 35

35

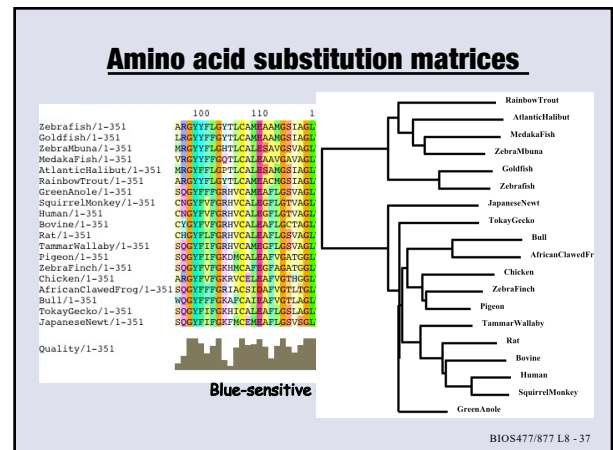
Amino acid substitution matrices

- Identity matrix
- Genetic code matrix
- Matrices based on AA properties
- **Matrices based on empirical data**
 - Dayhoff matrices (PAM120 etc.)
 - BLOSUM matrices (BLOSUM62 etc.)
 - Gonnet matrices (Gonnet 250 etc.)
 - JTT matrices
 - and more ...

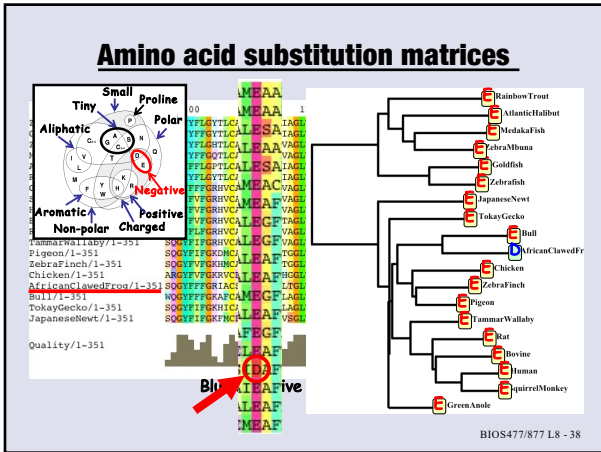
AAindex (566 indices + ~150 matrices): <https://www.genome.jp/dbget/aaindex.html>

BIOS477/877 L8 - 36

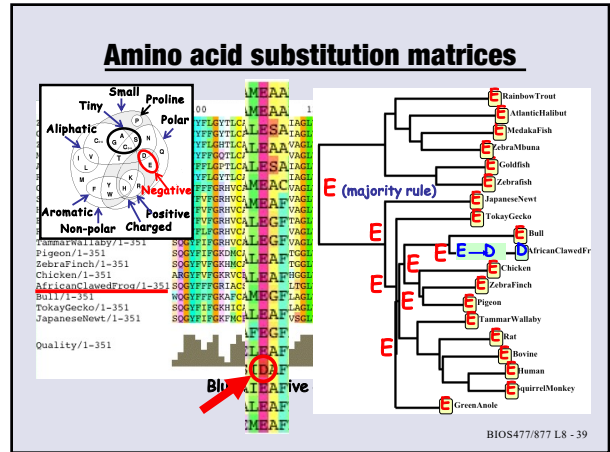
36



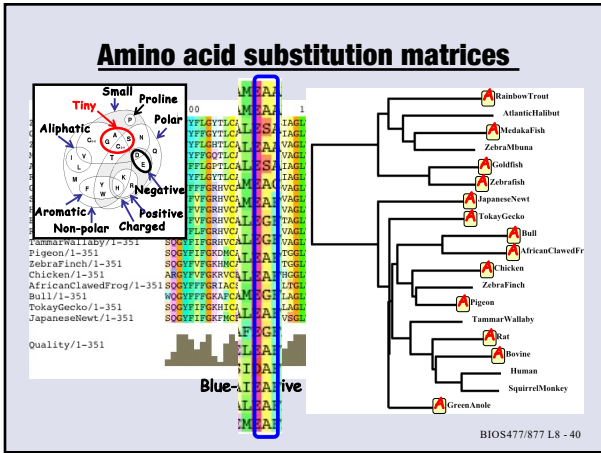
37



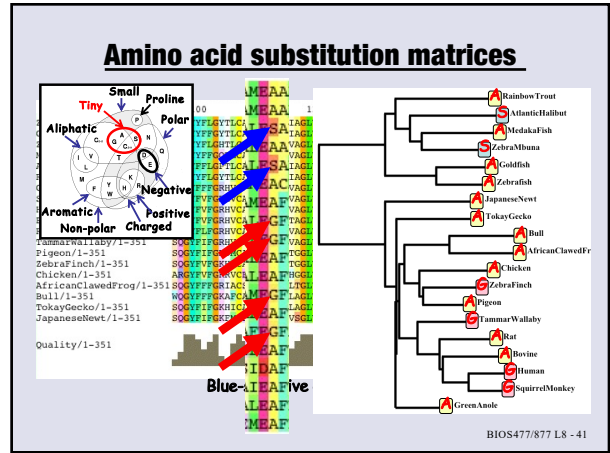
38



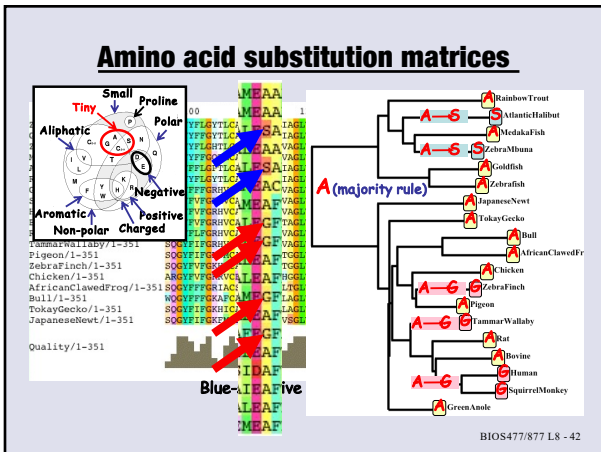
39



40



41



42

Substitution matrices based on empirical data

- **PAM matrices**
 - Dayhoff, Schwartz, and Orcutt (1978)
- **BLOSUM matrices**
 - Henikoff and Henikoff (1992)

Also see Eddy (2004) Nature Biotechnology 22: 1035-36

BIOS477/877 L8 - 43

43

Margaret O. Dayhoff
(1925-1983)

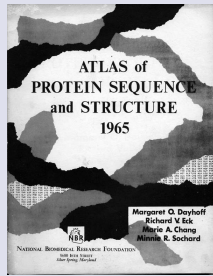


Founder of
the field of Bioinformatics
The first woman in the field

Dayhoff developed a single letter code for the amino acids

Read the [Smithsonian website](https://www.smithsonianmag.com/science-nature/margaret-dayhoff-180) also in Strasser (2010)

BIOS477/877 L8 - 44



Collection of all known protein sequences
1st *Atlas* contained 65 proteins
Developed into PIR (Protein Information
Resource), a brain-child of Dayhoff

PAM matrices (Dayhoff *et al.* 1978)

- **Accepted point mutations** (point accepted mutations, percent accepted mutations)
 - **accepted by selection**: no (or very weak) deleterious effect, maintaining the function
- **Based on 1,572 changes in 71 groups of closely related proteins (34 protein families)**
 - at least 85% identical
 - no ambiguity in alignments, no gap
 - most likely observed substitutions do not affect protein functions (accepted by selection, close to neutral)
 - successive (multiple) substitutions at one site are minimal (no hidden substitution)

BIOS477/877 L8 - 45