

# BIOS 477/877 Bioinformatics and Molecular Evolution

Instructor: Etsuko Moriyama  
(School of Biological Sciences)

Spring 2026

Lecture 8

BIOS477/877 L8 - 1

1

## Today's topics

- Alignment strategy: a good practice
- Gap Penalty
- Scoring (Substitution) Matrix
- Graduate Only Assignment
- 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

[Some useful tools]

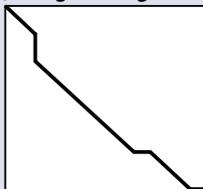
- [TranslatorX](#)
- [RevTrans 2.0](#)
- [PAL2NAL](#)
- [tranalign](#)  
(included in EMBOSS servers; see the course/Link page.)

BIOS477/877 L8 - 4

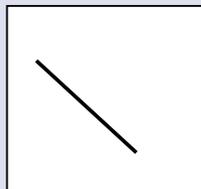
4

## 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  $a_i \neq b_j$   
 $w = -2$   
 $(w=0$  for free end gap)

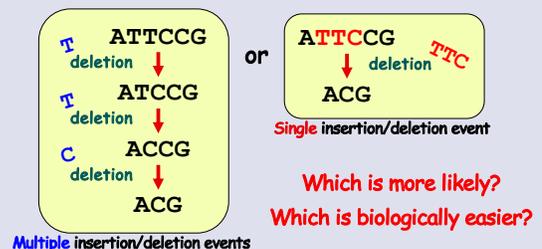
**Match/Mismatch scores**  
**Gap penalty**

BIOS477/877 L8 - 5

5

## Insertions and deletions: mechanisms

- A gap of length  $k \neq k$  gaps of length 1



BIOS477/877 L8 - 6

6

### Insertions and deletions: mechanisms (continued)

- 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)

Redeling et al. (2024)

BIOS477/877 L8 - 7

7

### Insertions and deletions: mechanisms (continued)

- 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**

BIOS477/877 L8 - 8

8

### Indel evolution and gap penalty

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

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

Score = 3      3      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)

Use match score = 1, mismatch score = 0, gap penalty = -1 using linear gap penalty

BIOS477/877 L8 - 12

12

### Simple alignments (with affine gap penalty)

➤ Varied length & gaps considered

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)}]

Match score = 1, mismatch score = 0, & using affine gap penalty:  
 gap opening penalty = -2, gap extension penalty = -1

BIOS477/877 L8 - 13

13

### Simple alignments (with affine gap penalty)

➤ Varied length & gaps considered

Score = 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)}]

Match score = 1, mismatch score = 0, & 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

Deletions are more frequent than insertions

Zhang and Gerstein (2003)

BIOS477/877 L8 - 15

15

### Empirical indel distribution: DNA

Based on the comparisons of 23 noncoding region sequences between *Drasophila simulans* and *D. sechella*

Keyhtley and Johnson (2004)  
 This distribution was later used in MCALIGN2 (Wang et al. 2006)

BIOS477/877 L8 - 16

16

### Empirical indel distribution: protein

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)

BIOS477/877 L8 - 17

17

### Empirical indel distribution: protein

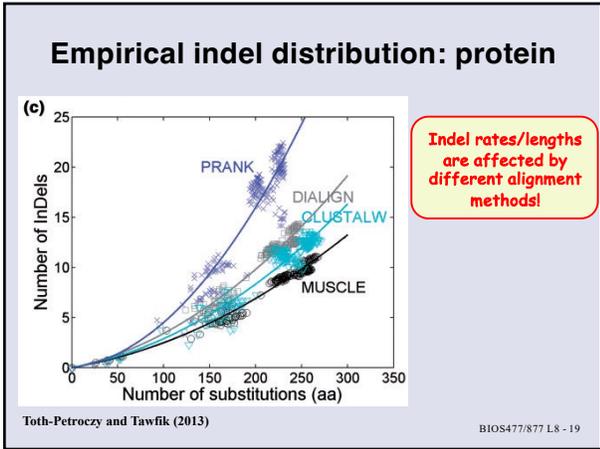
Based on the comparisons of 1,310 orthologous families from 22 fungal species.

Loop regions have more indels compared to the regions with secondary structures

Toth-Petroczy and Tawfik (2013)

BIOS477/877 L8 - 18

18



19

### 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

$w = a + b \ln(k)$  ← Affine gap penalty:  $w = a + bk$

$a$ : 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 and 2009); Loewenthal *et al.* 2021; Wygoda *et al.* (2024)  
 Review: Redelings *et al.* (2024)

BIOS477/877 L8 - 20

20

### Scoring (substitution) matrix: DNA identity matrix

$i \backslash j$	A	T	C	G
A	1	0	0	0
T	0	1	0	0
C	0	0	1	0
G	0	0	0	1

Match/mismatch scores can be expressed in a matrix format

Match score = 1  
 Mismatch score = 0

$s(i, j)$ : the similarity score between nucleotides  $i$  and  $j$

BIOS477/877 L8 - 21

21

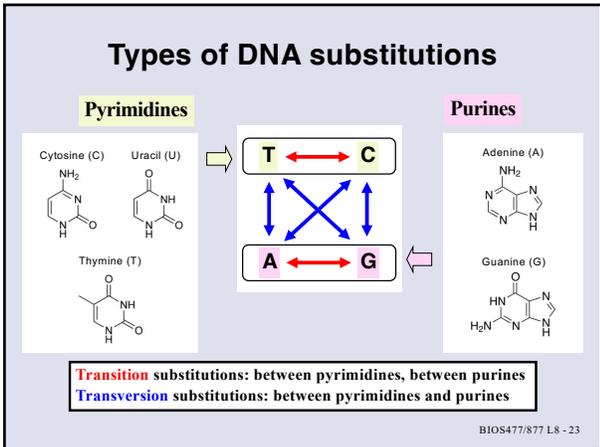
### Scoring (substitution) matrix: Transition/transversion matrix

	A	T	C	G
A	1	?	?	?
T	?	1	?	?
C	?	?	1	?
G	?	?	?	1

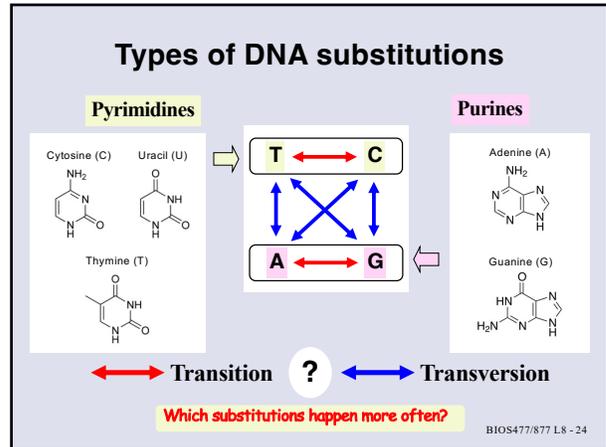
Match score = 1  
 Mismatch score:  
 { Transition = ?  
 { Transversion = ?

BIOS477/877 L8 - 22

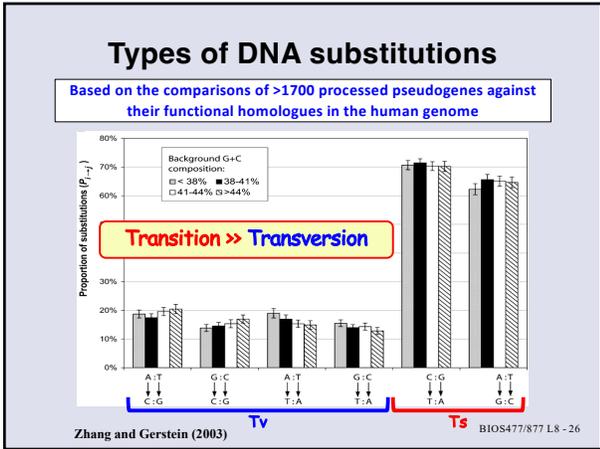
22



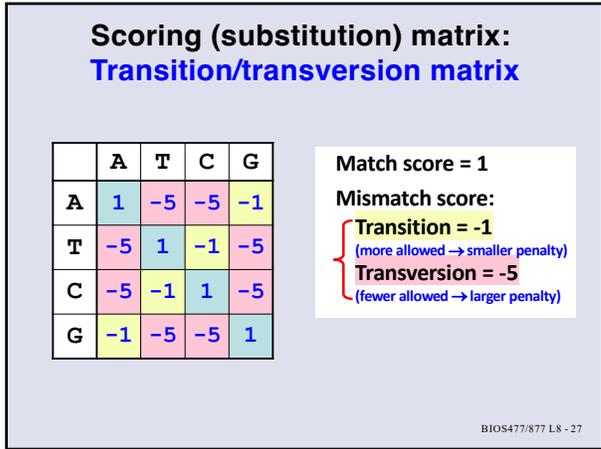
23



24



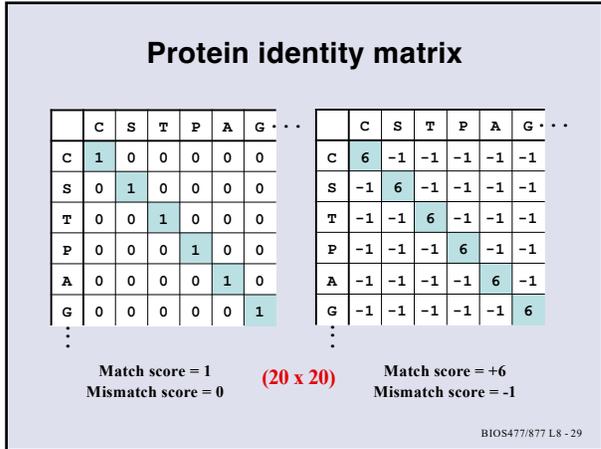
26



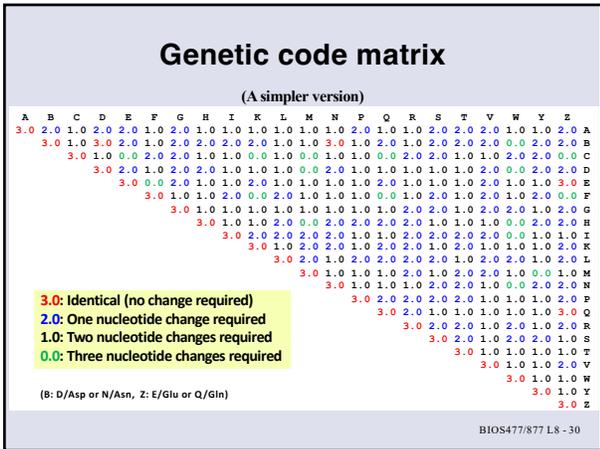
27

- ### 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)
- BIOS477/877 L8 - 28

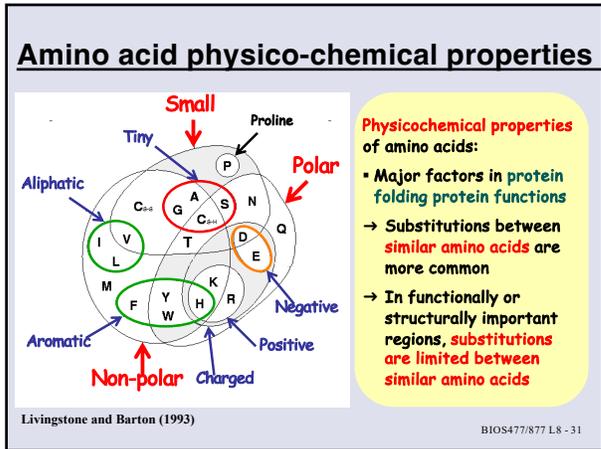
28



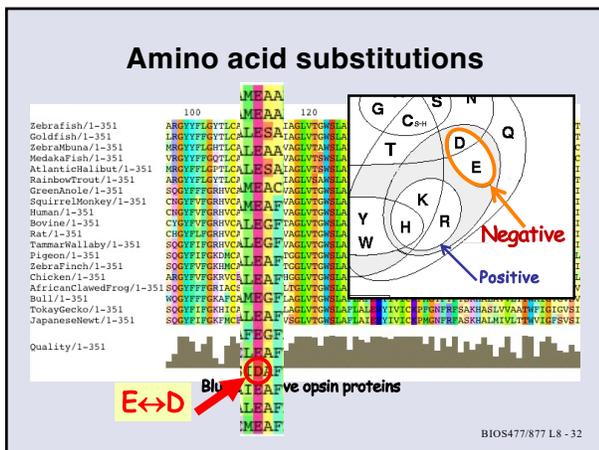
29



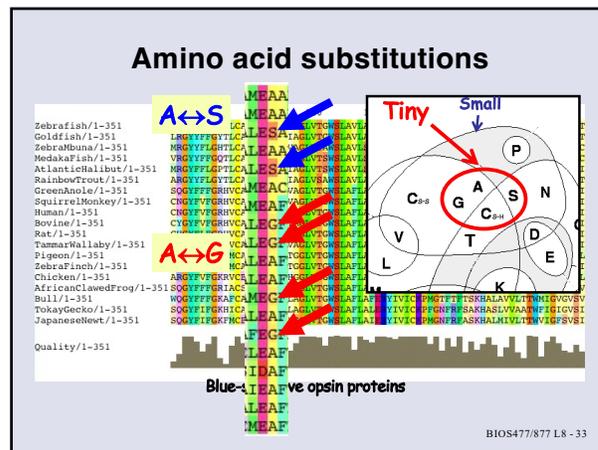
30



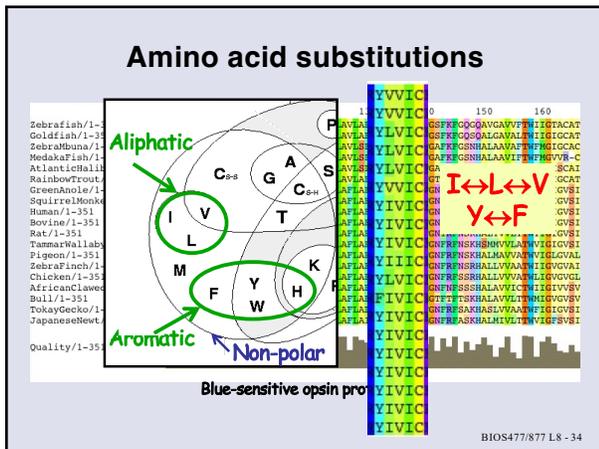
31



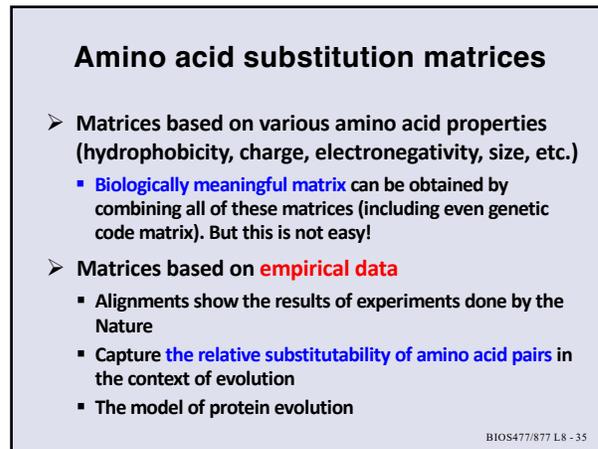
32



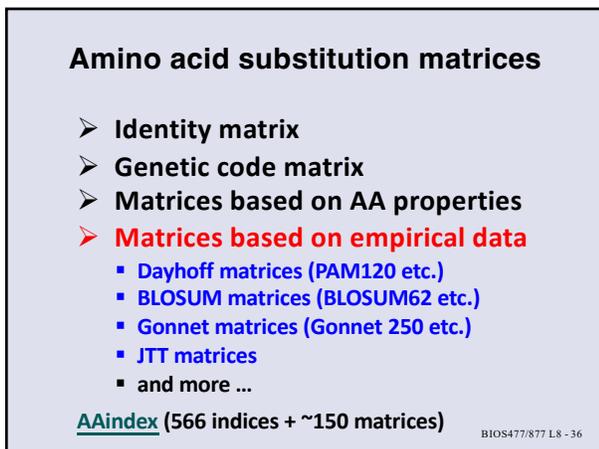
33



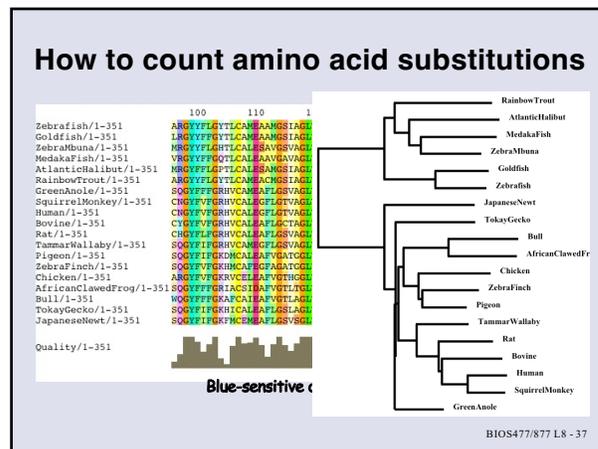
34



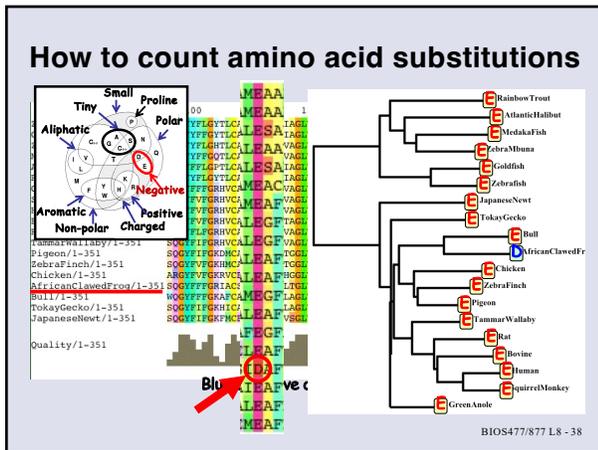
35



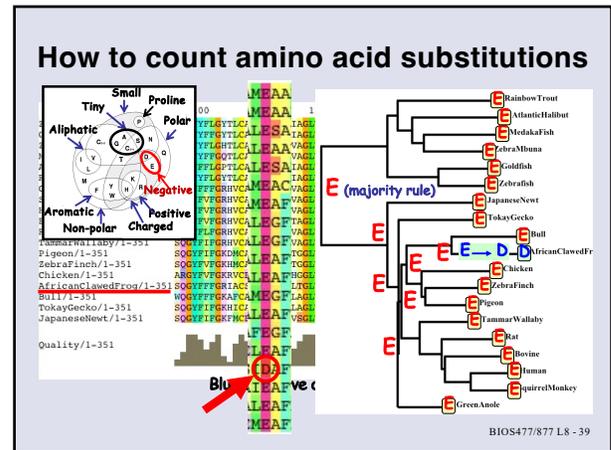
36



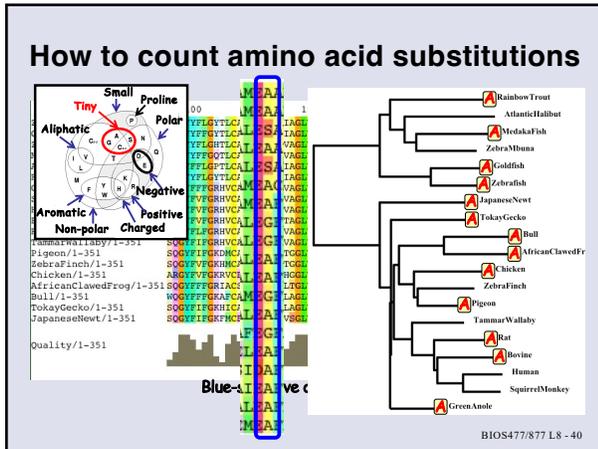
37



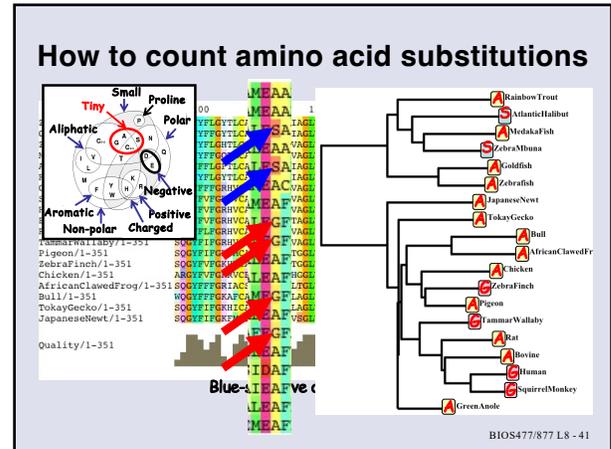
38



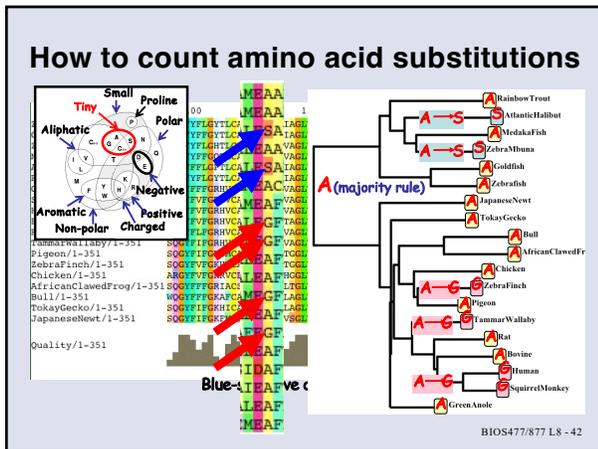
39



40



41



42

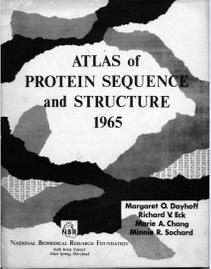
### Amino acid substitution matrices based on empirical data

- PAM matrices  
Dayhoff, Schwartz, and Orcutt (1978)
- BLOSUM matrices  
Henikoff and Henikoff (1992)
- Also see Eddy (2004)

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

Collection of all known protein sequences  
1st *Atlas* contained 65 proteins

Developed into PIR (Protein Information Resource), a brain-child of Dayhoff

Read the [Smithsonian website](#); also in Strasser (2010)

BIOS477/877 L8 - 44

44

### 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

45

### Graduate Only Assignment

- **Report 1: March 8**
  - Short proposal (2-3 pages)
- **Report 2: May 3**
  - Final report (5-7 pages)
- Comparative analysis of bioinformatics methods
  - e.g., Multiple alignment methods:  
Clustal W vs. MUSCLE vs. MAFFT
  - Compare three methods
  - Choose appropriate input datasets (3 or more)
  - Read the Guidelines!
- Undergraduate students can do these assignments for bonus points (**both reports need to be submitted**)

BIOS477/877 L8 - 46

46