Sequence Alignment



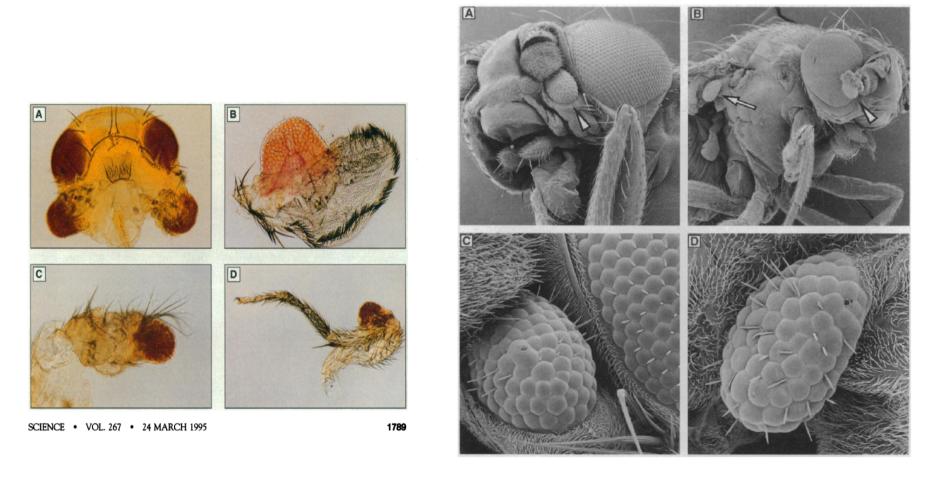
The Story of eyeless





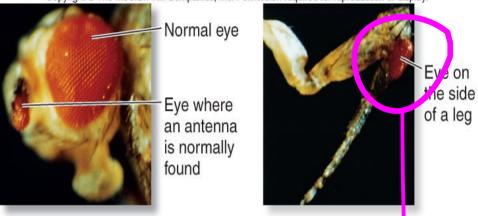
Wild type *eyeless*

Muahahaha....



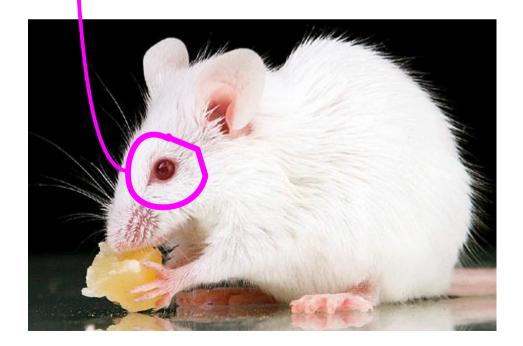
Ectopic expression of Drosophila eyeless gene

Lhahahahaha Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.

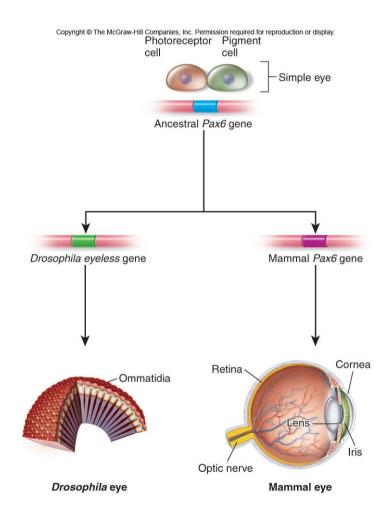


(a) Abnormal expression of Drosophila eyeless gene (b) Abnormal expression of mouse Pax6 gene in a fruit fly leg

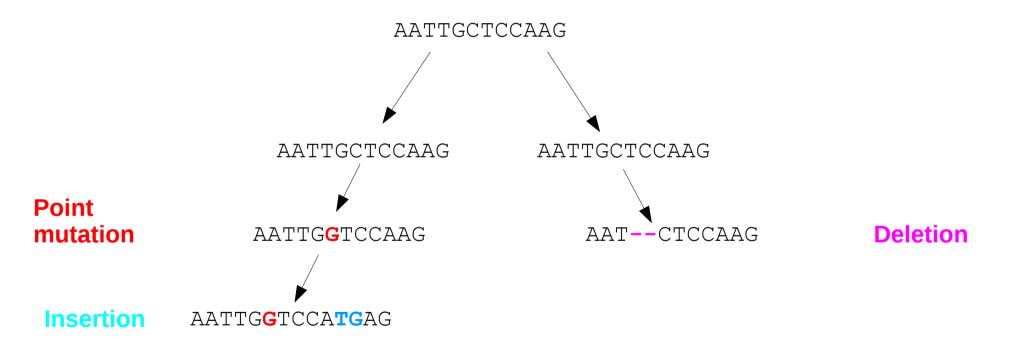
a,b: © Prof. Walter J. Gehring, University of Basel



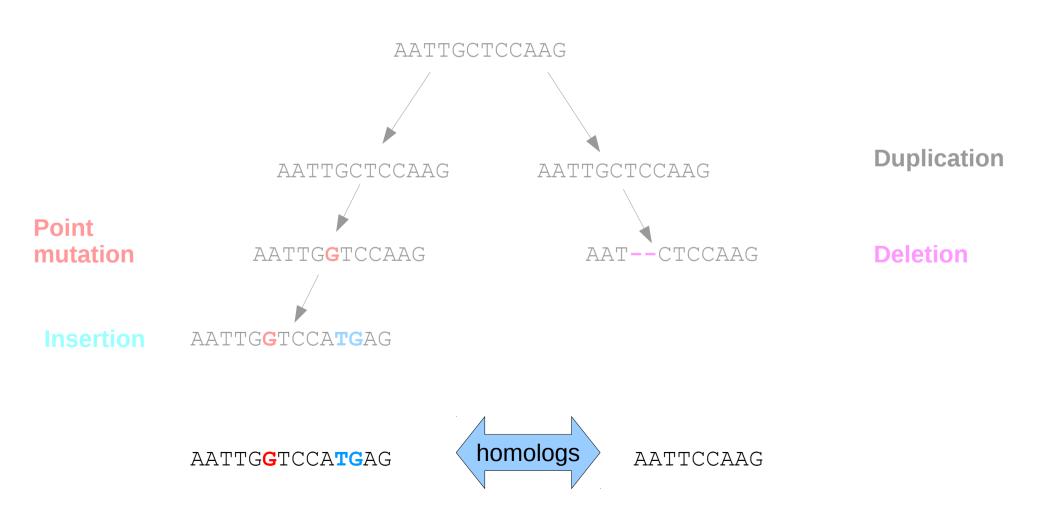
Homology: Common ancestry



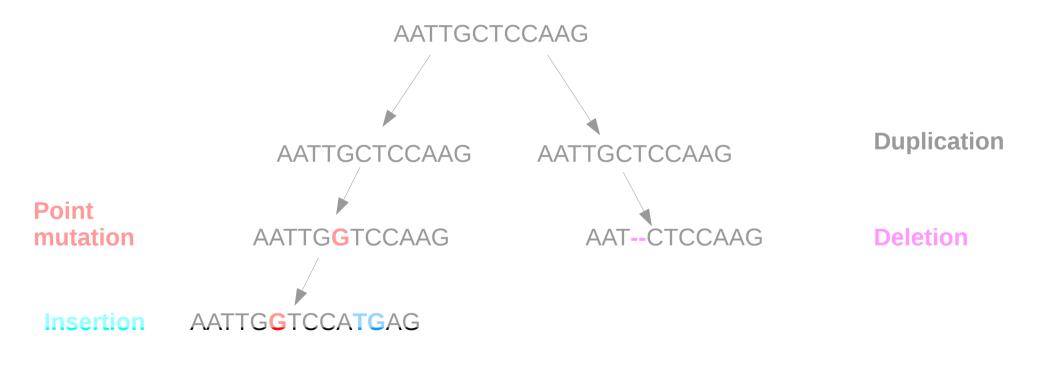
Sequence Evolution Events



Sequence Evolution Events



Sequence Evolution Events



AATTGGTCCATGAG AAT - - CTCCA - - AG

Sequence alignment: reconstructing evolutionary events since LUCA

Sequence Alignment for

- Predicting function
- Searching sequence databases
- Gene finding (not by length alone...)
- Genome construction / sequence assembly

Local & Global Alignments

Global: indicative of whole gene/protein homology

Local: Common substructure in two different genes / proteins

```
----QTGKGS-SRIWDN
| | | | | | |
----SAGKGAIMRLGDA
```

Cost of Evolutionary Events

$$c(-, a) = c(a,-) = -1$$

 $c(a,b) = -1$

$$c(a,b) = 1$$

insertion/deletion (indel)

when: a≠b mut

when: a=b

mutation

Cost of Evolutionary Events

$$c(-, a) = c(a,-) = -1$$
 $c(a,b) = -1$ $c(a,b) = 1$ when: $a \neq b$ mutation $c(a,b) = 1$ when: $a = b$

Also known as **scoring** or **cost function**: mutations and indels are more **Costly** than evolutionary "inertia"

Homologous Camels GAMAL JAMAL



BIOLOGY LIMNOLOGY

--BIOLOGY LIMNOLOGY

-1+-1+-1+-1+1+1+1+1=1

BI--OLOGY LIMNOLOGY

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

BI--OLOGY LIMNOLOGY

--BIOLOGY LIMNOLOGY

-1+-1+-1+-1+1+1+1+1+1=1

BP-0L0GY LIMNOLOGY

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

--BIOLOGY

-1+-1+-1+-1+1+1+1+1+1=1

Improving the cost function

$$c(-, a) = c(a,-) = -1$$

insertion/deletion (indel)

$$c(-) = -0.2$$

extension

$$c(a,b) = -1$$

when: a≠b n

mutation

$$c(a,b) = 1$$

when: a=b

Also known as cost function:

Point mutations are costly

Starting indels is costly

Extending indels cost less than starting them. (Why?)

The Dotplot

```
Y
G
X
X
O
X
L
X
O
X
N
N
N
M
L
X
L
X
B
I
O
L
O
G
Y
```

Back to DNA

 Scoring matrix or substitution matrix describes the cost of mutations

	Α	Т	G	С	-
Α	1	-1	-1	-1	-1
Т	-1	1	-1	-1	-1
G	-1	-1	1	-1	-1
С	-1	-1	-1	1	-1
-	-1	-1	-1	-1	N/D



The simplest matrix: IDENTITY MATRIX

```
Human Chr. 21 37218151-37247016, +
Chimp Chr. 22 37250843-37279887.+
 37250843 GGAGCAATGCACAAGCTCTCTAGTTCGCAAGGTAAAGGTGCGACCTTTCT 37250892
 37218177 -
 37218357 TGCCCCACACACACAGATTCAGGTCCCAGCCGTCCTAACCCCCCAGCTGCC 37218406
```

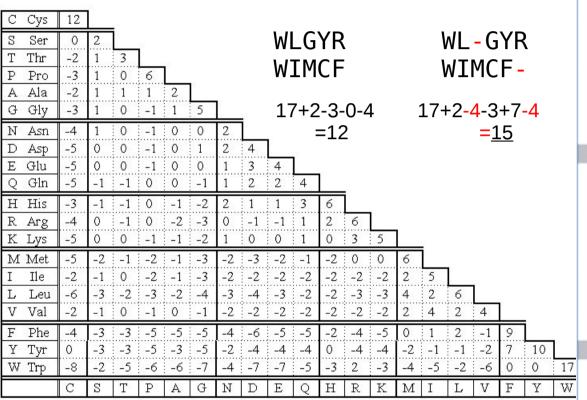
Proteins: serious business

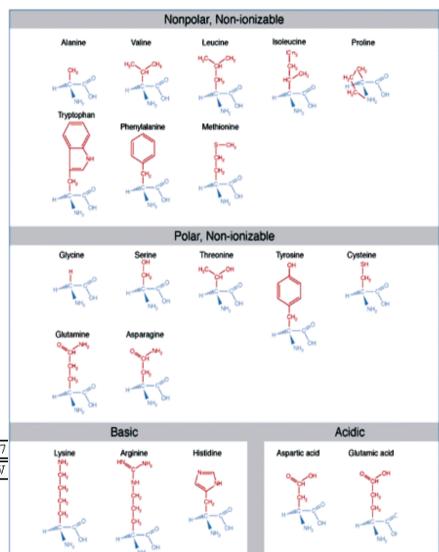
- Protein sequences: 20 (21) amino acids
- Some amino acids tend to be more conserved
- Some amino acids substitute others more readily. (Why?)
- Hence: protein substitution matrices are <u>quantitative</u> rather than <u>qualitative</u>

Scoring matrices for Proteins

- Simulate the evolutionary cost of mutations and indels in proteins
- Generated by a combination of:
 - Observations
 - Statistical model
- Often "repeats" the known biochemistry

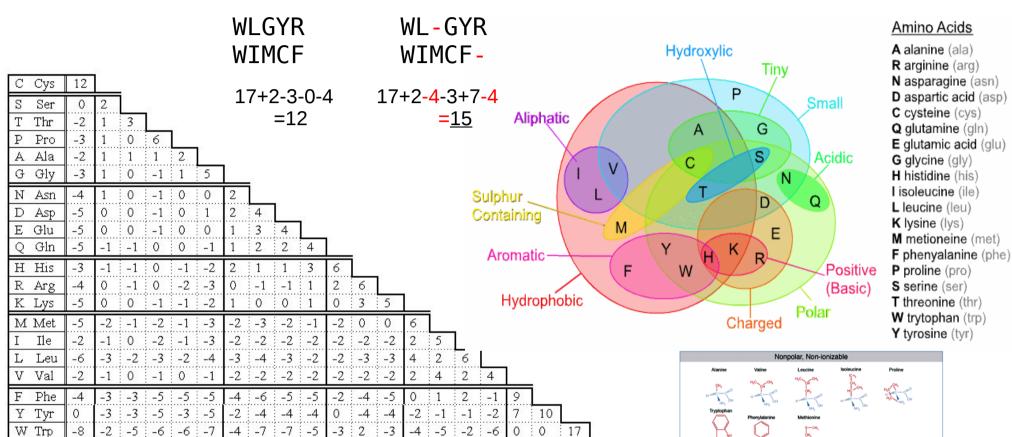
PAM Matrix





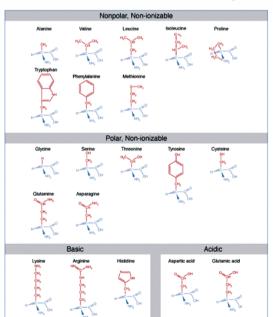


PAM Matrix





Margaret Dayhoff, 1978



PAM (Percent Accepted Mutations)

- Based on a database of 1572 changes in 71 groups of closely related proteins
- PAM1: Probability of 1 amino acid out of 100 will mutate

AA Change	PAM1
Phe ↔ Phe	0.9946
Phe ↔ Tyr	0.0021
Phe ↔ Leu	0.0013
Phe ↔ Ala	0.0002
Phe ↔ Arg	0.0001
•••	
SUM	1.000

PAM (Percent Accepted Mutations)

- Based on a database of 1572 changes in 71 groups of closely related proteins
- PAM1: Probability of 1 amino acid out of 100 will mutate
- PAM250 = PAM1 times itself 250 times

AA Change	PAM1	PAM250
Phe ↔ Phe	0.9946	0.32
Phe ↔ Tyr	0.0021	0.15
Phe ↔ Leu	0.0013	0.13
Phe ↔ Ala	0.0002	0.04
Phe ↔ Arg	0.0001	0.01
SUM	1.000	1,00

PAM

- Scoring matrix: 20x20
- Matrix entry M(i,j) = probability of i substituting j
- Diagonally symmetric M(i,j) = M(j,i)
- The higher the PAM number, the higher the evolutionary distance we can assess

PAM: Assumptions and Errors

Assumptions

- Replacement at any site depends only on the amino acid at that site and the probability given by the table (Markov model).
- Sequences that are being compared have a representative amino acid composition.

• Errors

- Many sequences depart from average composition.
- Rare replacements were observed too infrequently to resolve relative probabilities accurately (for 36 pairs no replacements were observed!).
- Errors in PAM1 are magnified in the extrapolation to PAM 250.
- Distantly related sequences usually have stretches of conserved residues. This implies that replacement is not equally probable over entire sequence.

BLOSUM Matrix

- Generated from BLOCKS database
- Initial alignments determined by %ID in BLOCK
- BLOSUM85: alignments of 85% ID

```
Block BP00180A
     LIGASE SYNTHETASE CARBAMOYL
ACCC BACSU P49787
                           GPSADAISKMG
                                         39
ACCC METJA | 058626
                                         52
                           GPNPDAIEAMG
                                         57
CARB BACSU P25994
                      276)
                           GIEGGCNVOLA
CARY BACSU P18185
                      954)
                           GTFASWMEOEG
                                         51
DUR1 YEAST
           P32528
                      735)
                           GPSGDIIRGLG
                                         17
PCCA HUMAN P05165
                           GSVGYDPNEKT
                                        100
PUR2 ECOLI P15640
                           GPTAGAAOLEG
                                         37
PURK PSEAE P72158
                           GOLGRMLALAG
                                         26
PURT PASHA P46927
                      261) GIFGVELFVCG
                                         46
PYR1 DROME | P05990
                           GVGGEVVFOTG
                                         18
SUCC METJA 057663
                       52)
                           GKAGGILFASN
                                         54
YFIO ECOLI P76594
                                         59
                      125) NSLGLLAPWOG
ACCC ECOLI
           P24182
                      104) GPKAETIRLMG
                                         27
ACCC HAEIN P43873
                                         23
                      104) GPTADVIRLMG
ACCC PSEAE P37798
                                         23
                      104)
                           GPTAEVIRLMG
CARB ECOLI P00968
                      397) ALRGLEVGATG
                                         28
                      429) GSGGLSIGOAG
CPSM HUMAN P31327
CPSM RAT | P07756
                      429) GSGGLSIGOAG
PYR1 DICDI P20054
                      372)
                           GSGGLSIGOAG
PYR1 HUMAN P27708
                      400)
                           GSGGLSIGQAG
PYR1 YEAST P07259
                           GSGGLSIGOAG
                      445)
COA1 HUMAN Q13085
                           SDLGISALODG
                       60)
```

Deriving BLOSUM

BABA AAAC AACC AABA AACC AABC

aa	Observed (qi) frequency
А	14/24
В	4/24
С	6/24

Mutation type	Observed mutations (qij)	Expected mutations (eij)
$A \leftrightarrow A$	26/60 = 0.433	(14/24)x(14/24) = 0.34
A ↔ B	8/60 = 0.1333	(14/24)x(4/24) = 0.0972
A ↔ C		(14/24)x(6/24) = 0.1459
B ↔ B	3/60 = 0.05	
B ↔ C		
C ↔ C		

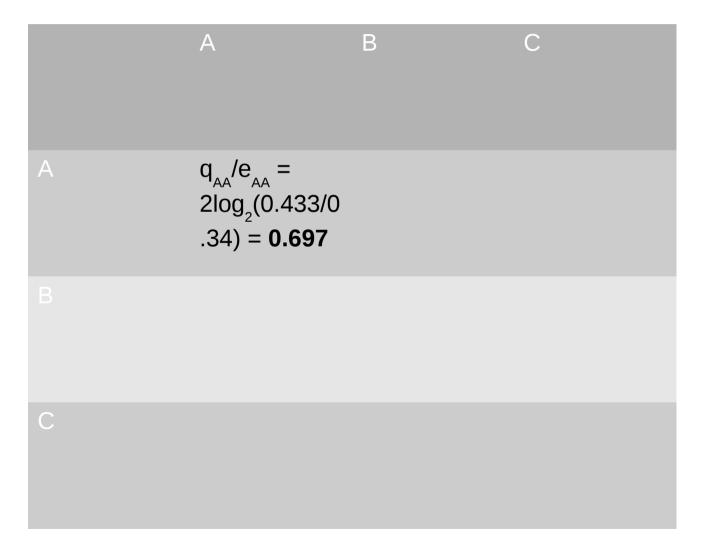
Deriving BLOSUM

$$q_{ij} = \frac{f_{ij}}{\sum_{i}^{20} (f_{ij})}$$

$$e_{ij} = q_i \times q_j$$

$$e_{ij} = q_i \times q_j$$

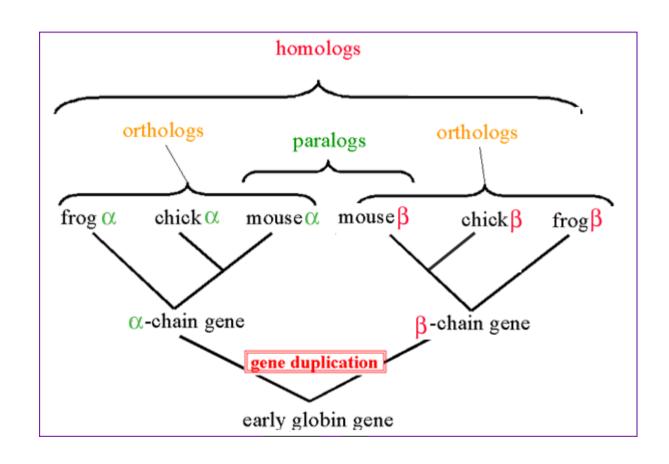
$$M_{ij} = 2 \log_2(\frac{q_{ij}}{e_{ij}})$$



BLOSUM vs. PAM

	PAM	BLOSUM
Derived from	Full-length proteins	Local blocks
Evolutionary model	Markovian & extrapolation	None implicit

Homology, Orthology, Paralogy



My Brain Hurts...

