

# Reminders!

- PLEASE Register for the Virtual Bioinformatics Conference It's Free.
- September 24-26, 2002, Access Grid, Room ECS 212.
- Homework #1
- Run the BLAST Tutorial

# Useful Terms

- **E value:** Expected # of chance alignments with scores  $\geq S$ . The lower the E value, the more significant the score.
- **P value:** The probability of an alignment occurring with score  $\geq S$  for a random sequence. Calculated by relating the observed score,  $S$ , to the expected distribution of HSP scores from comparisons of random sequences of the same length and composition as the query to the database. The most highly significant P values will be those close to 0.
- **HSP:** High-scoring segment pair. Local alignments with no gaps that achieve high alignment scores

# Rules of Thumb

- Most sequences with significant similarity over their entire lengths are homologous.
- Matches that are > 50% identical in a 20-40 aa region occur frequently by chance.
- Distantly related homologs may lack significant similarity. Homologous sequences may have few absolutely conserved residues.
- A homologous to B & B to C  $\Rightarrow$  A homologous to C.
- Low complexity regions, transmembrane regions and coiled-coil regions frequently display significant similarity without homology.
- Greater evolutionary distance implies that length of a local alignment required to achieve a statistically significant score also increases.

# Rules of Thumb

- Results of searches using different scoring systems may be compared directly using normalized scores.
- If  $S$  is the (raw) score for a local alignment, the **normalized** score  $S'$  (in bits) is given by

$$S' = \frac{\lambda - \ln(K)}{\ln(2)}$$

The parameters depend on the scoring system.

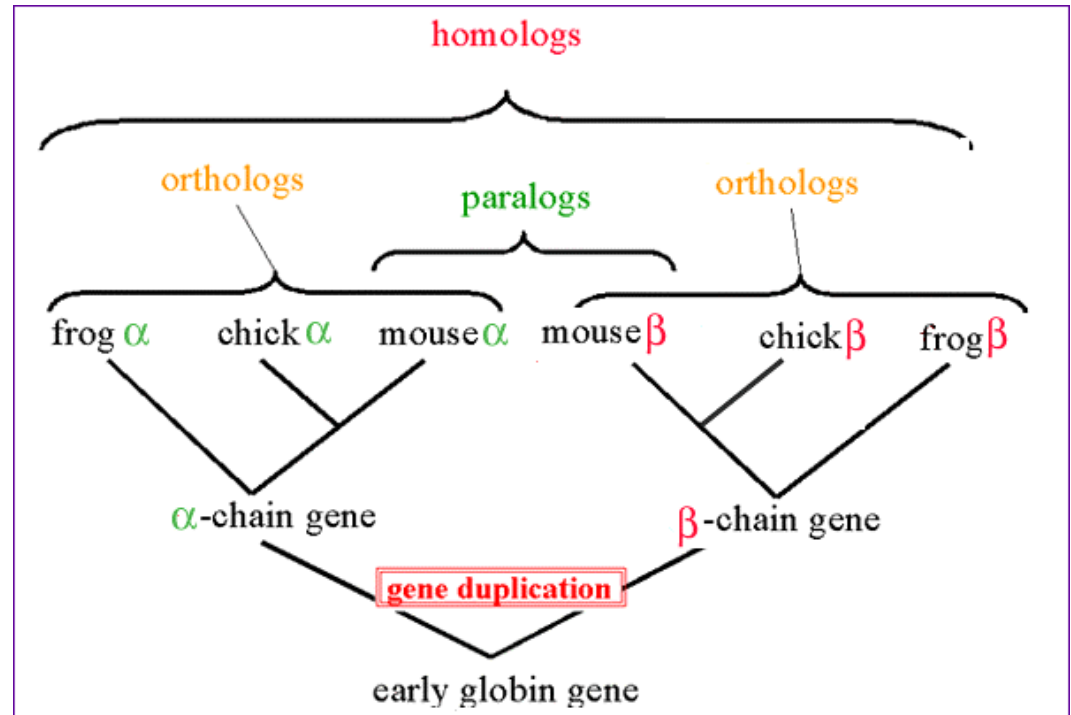
- **Statistically significant normalized score,**

$$S' > \log\left(\frac{N}{E}\right)$$

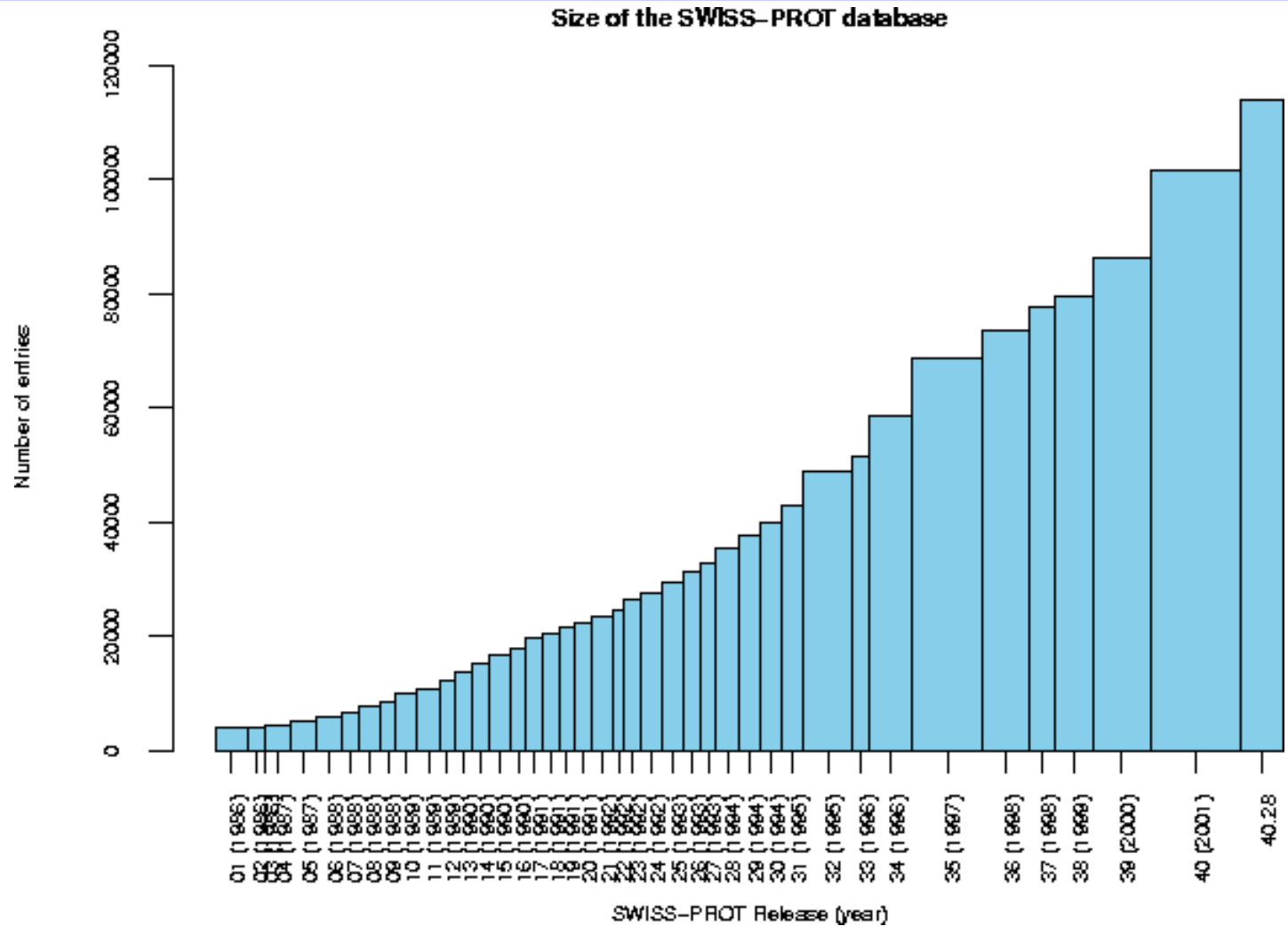
where E-value =  $E$ , and  $N$  = size of search space.

# Homologs: Orthologs & Paralogs

- **Homology:** Similarity due to common ancestry.
- **Orthologs:** Homologous sequences in different species that arose from a common ancestral gene during speciation; may or may not be responsible for a similar function.
- **Paralogs:** Homologous sequences within a single species that arose by gene duplication.

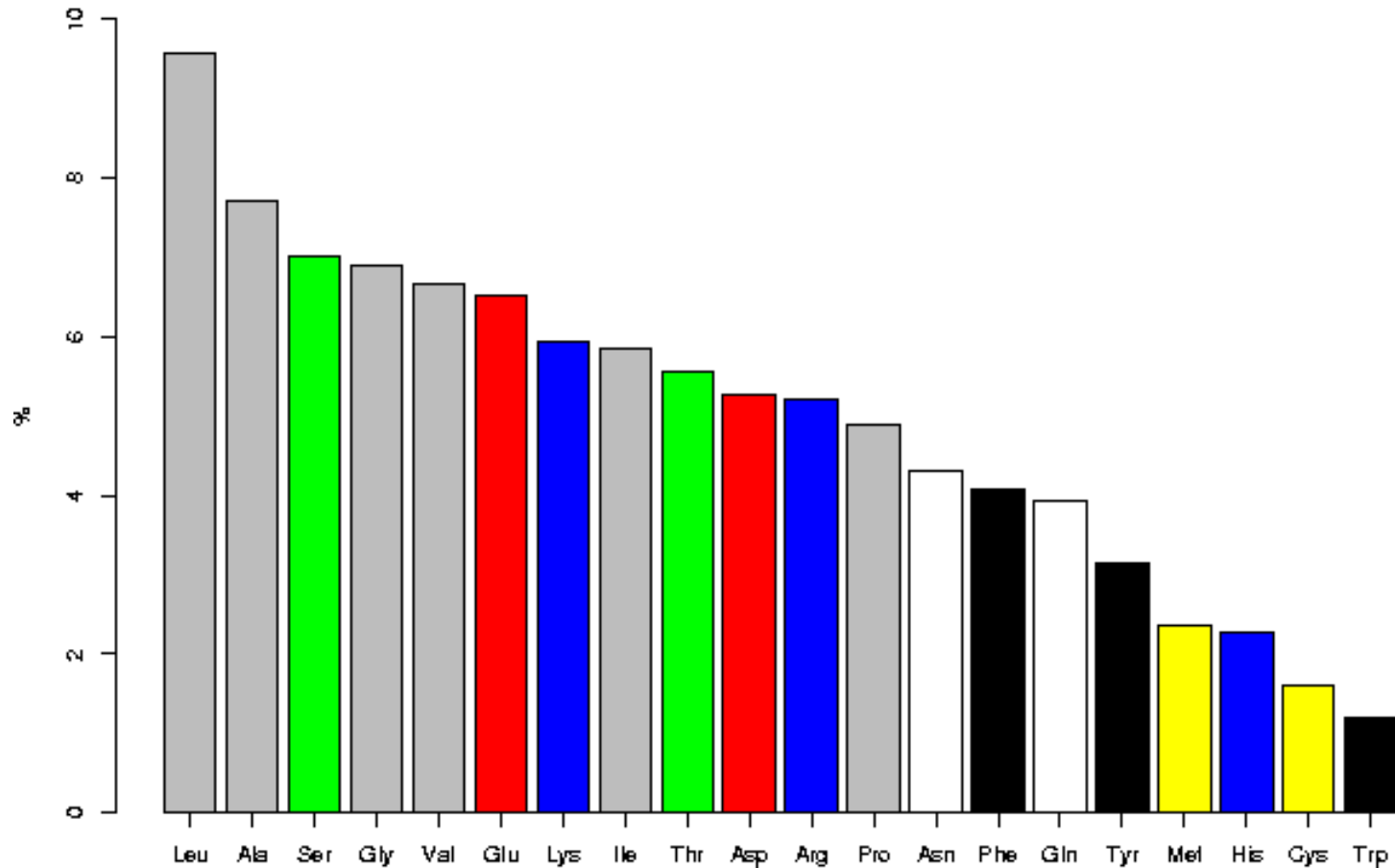


# Growth of SWISS-PROT

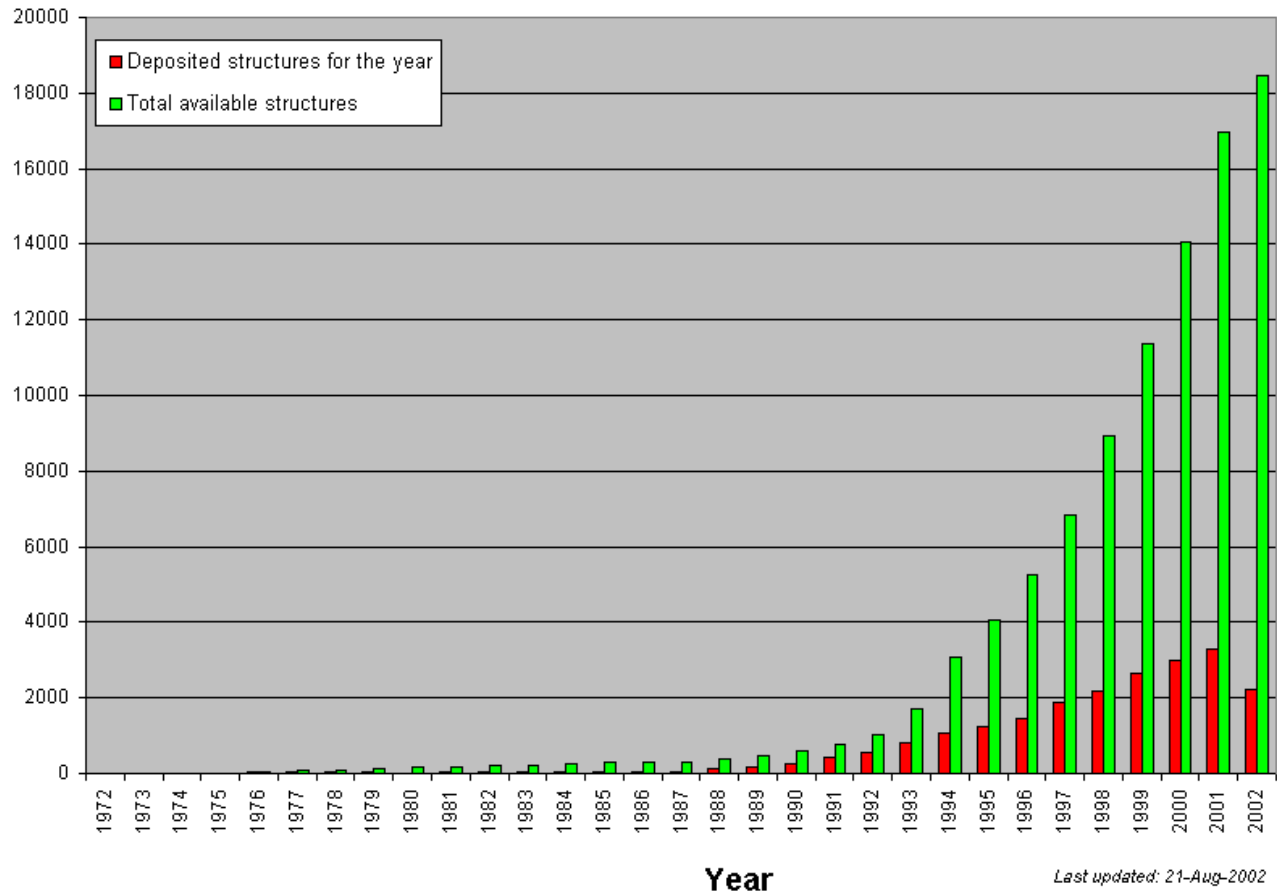


# Amino-acid composition from SWISS-PROT

Amino acid composition

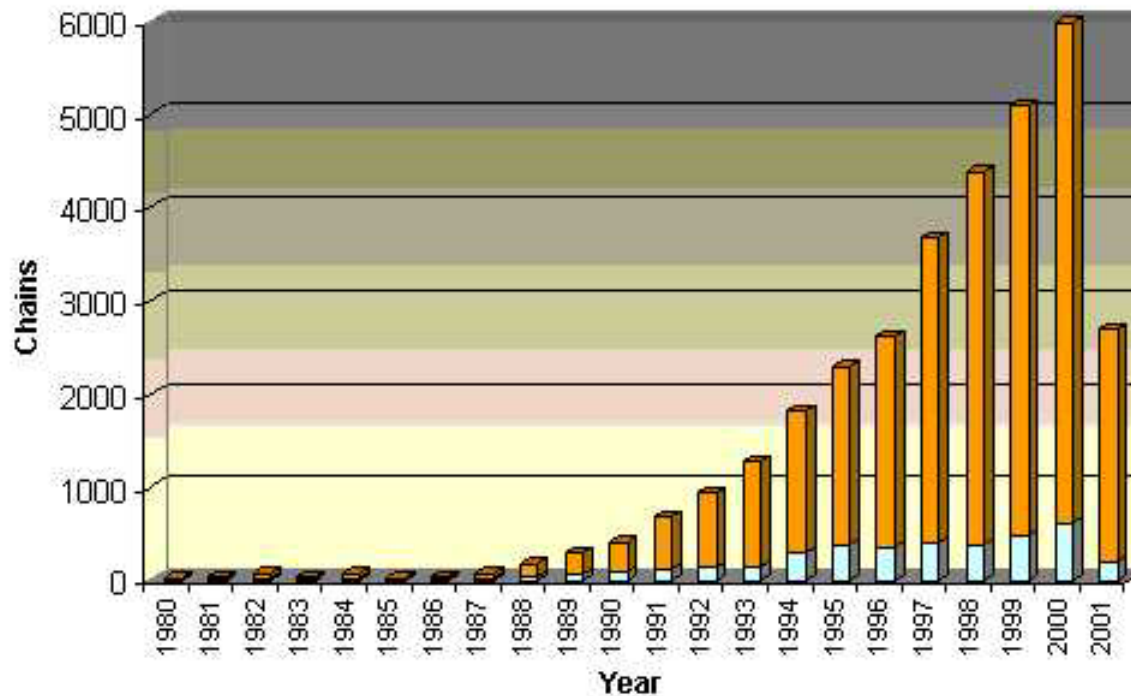


# PDB Growth





# Growth in New Folds - PDB



# Multiple Alignments

- **Family alignment for the ITAM domain**

- CD3D\_MOUSE/1-2    **E**Q**L**Y**Q**P**L**R**D**R    **E**D**T**Q-**Y**S**R**L**G**    GN
- Q90768/1-21      **D**Q**L**Y**Q**P**L**G**E**R    **N**D**G**Q-**Y**S**Q**L**A**    TA
- CD3G\_SHEEP/1-2    **D**Q**L**Y**Q**P**L**K**E**R    **E**D**D**Q-**Y**S**H**L**R**    KK
- P79951/1-21      **N**D**L**Y**Q**P**L**G**Q**R    **S**E**D**T-**Y**S**H**L**N**    SR
- FCEG\_CAVPO/1-2    **D**G**I**Y**T**G**L**S**T**R    **N**Q**E**T-**Y**E**T**L**K**    HE
- CD3Z\_HUMAN/3-0    **D**G**L**Y**Q**G**L**S**T**A    **T**K**D**T-**Y**D**A**L**H**    MQ
- C79A\_BOVIN/1-2    **E**N**L**Y**E**G**L**N**L**D    **D**C**S**M-**Y**E**D**I**S**    RG
- C79B\_MOUSE/1-2    **D**H**T**Y**E**G**L**N**I**D    **Q**T**A**T-**Y**E**D**I**V**    TL
- CD3H\_MOUSE/1-2    **N**Q**L**Y**N**E**L**N**L**G    **R**R**E**E-**Y**D**V**L**E**    KK
- CD3Z\_SHEEP/1-2    **N**P**V**Y**N**E**L**N**V**G    **R**R**E**E-**Y**A**V**L**D**    RR
- CD3E\_HUMAN/1-2    **N**P**D**Y**E**P**I**R**K**G    **Q**R**D**L-**Y**S**G**L**N**    QR
- CD3H\_MOUSE/2-0    **E**G**V**Y**N**A**L**Q**K**D    **K**M**A**E**A**Y**S**E**I**G    TK
- Consensus/60%    -.lYpsLspc    pcsp.YspLs    pp

# CLUSTALW

- \* identical
- : conserved substitutions
- . semi-conserved substitutions

```

gi|2213819   CDN-ELKSEAIIEHLCASEFALR-----MKIKEVKKENGDKK 223
gi|12656123  ----ELKSEAIIEHLCASEFALR-----MKIKEVKKENGD-- 31
gi|7512442   CKNKNDDDDNDIMETLCKNDFALK-----IKVKEITYINRDTK 211
gi|1344282   QDECKFDYVEVYETSSSGAFSLGRFCGAEPPLVSSHHELAVLFRTDH 400
                : .      : *   . .  *:*          .  **:
  
```

Red: AVFPMLW (Small & hydrophobic)

Blue: DE (Acidic)

Magenta: RHK (Basic)

Green: STYHCNGQ (Hydroxyl, Amine, Basic)

Gray: Others

# How to Score Multiple Alignments?

- Sum of Pairs Score (SP)
  - Optimal alignment:  $O(d^N)$  [Dynamic Prog]
  - Approximate Algorithm: **Approx Ratio 2**
    - Locate Center:  $O(d^2N^2)$
    - Locate Consensus:  $O(d^2N^2)$

**Consensus char**: char with min distance sum

**Consensus string**: string of consensus char

**Center**: input string with min distance sum

# Multiple Alignment Methods

- Phylogenetic Tree Alignment (NP-Complete)
  - Given tree, task is to label leaves with strings
- Iterative Method(s)
  - Build a MST using the distance function
- Clustering Methods
  - Hierarchical Clustering
  - K-Means Clustering

# Multiple Alignment Methods (Cont'd)

- Gibbs Sampling Method
  - Lawrence, Altschul, Boguski, Liu, Neuwald, Winton, *Science*, 1993
- Hidden Markov Model
  - Krogh, Brown, Mian, Sjolander, Haussler, *JMB*, 1994

# Profile Method

PROFILE METHOD, [M. Gribskov et al., '90]

Location in Seq.	Sequence							Protein Name
	1	2	3	4	5	6	7	
14	G	V	S	A	S	A	V	Ka RbtR
32	G	V	S	E	M	T	I	Ec DeoR
33	G	V	S	P	G	T	I	Ec RpoD
76	G	A	G	I	A	T	I	Ec TrpR
178	G	C	S	R	E	T	V	Ec CAP
205	C	L	S	P	S	R	L	Ec AraC
210	C	L	S	P	S	R	L	St AraC
13	G	V	N	K	E	T	I	Br MerR

FREQUENCY TABLE

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	0	2	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	4	0	0
3	0	0	0	0	0	1	0	0	0	0	1	0	0	0	6	0	0	0	0	0
4	1	0	0	1	0	0	0	1	1	0	0	0	3	0	1	0	0	0	0	0
5	1	0	0	2	0	1	0	0	0	1	0	0	0	0	3	0	0	0	0	0
6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	5	0	0	0
7	0	0	0	0	0	0	4	0	2	0	0	0	0	0	0	0	0	2	0	0

7

# Profile Method

FREQUENCY TABLE

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	0	2	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	4	0
3	0	0	0	0	0	1	0	0	0	0	1	0	0	0	6	0	0	0	0	0
4	1	0	0	1	0	0	0	1	1	0	0	0	3	0	1	0	0	0	0	0
5	1	0	0	2	0	1	0	0	0	0	1	0	0	0	0	3	0	0	0	0
6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	5	0	0	0
7	0	0	0	0	0	0	0	4	0	2	0	0	0	0	0	0	0	2	0	0

WEIGHT MATRIX

	A	C	E	G	I	K	L	M	N	P	R	S
1	0	108	0	101	0	0	0	0	0	0	0	0
2	21	78	0	0	0	0	44	0	0	0	0	0
3	0	0	0	23	0	0	0	0	46	0	0	102
4	21	0	32	0	38	32	0	0	0	86	39	0
5	21	0	62	23	0	0	0	74	0	0	0	72
6	21	0	0	0	0	0	0	0	0	0	69	0
7	0	0	0	0	98	0	44	0	0	0	0	0

$$Weight[i, AA] = \log \left( \frac{Freq[i, AA]}{p[AA] \cdot N} \right) \cdot 100$$

8



# Profile Method

WEIGHT MATRIX

	A	C	E	G	I	K	L	M	N	P	R	S
1	0	108	0	101	0	0	0	0	0	0	0	0
2	21	78	0	0	0	0	44	0	0	0	0	0
3	0	0	0	23	0	0	0	0	46	0	0	102
4	21	0	32	0	38	32	0	0	0	86	39	0
5	21	0	62	23	0	0	0	74	0	0	0	72
6	21	0	0	0	0	0	0	0	0	0	69	0
7	0	0	0	0	98	0	44	0	0	0	0	0

Given the following protein sequence:

```
M T E D L F G D L Q D D T I L A H L D N
P A E D T S R F P A L L A E L N D L L R
G E L S R L G V D P A H S L E I V V A I
C K H L G G G Q V Y I P R G Q A L D S L
I R D L R I W N D F N G R N V S E L T T
R Y G V T F N T V Y K A I R R M R R L K
```

# CpG Islands

- Regions in DNA sequences with increased occurrences of substring “CG”
- Rare: typically C gets methylated and then mutated into a T.
- Often around promoter or “start” regions of genes
- Few hundred to a few thousand bases long

## Problem 1:

- **Input:** Small sequence **S**
- **Output:** Is **S** from a CpG island?
  - Build Markov models:  $M^+$  and  $M^-$
  - Then compare

# Markov Models

<b>+</b>	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	0.180	0.274	0.426	0.120
<b>C</b>	0.171	0.368	0.274	0.188
<b>G</b>	0.161	0.339	0.375	0.125
<b>T</b>	0.079	0.355	0.384	0.182

<b>—</b>	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	0.300	0.205	0.285	0.210
<b>C</b>	0.322	0.298	0.078	0.302
<b>G</b>	0.248	0.246	0.298	0.208
<b>T</b>	0.177	0.239	0.292	0.292

# How to distinguish?

- Compute

$$S(x) = \log\left(\frac{P(x | M+)}{P(x | M-)}\right) = \sum_{i=1}^L \log\left(\frac{p_{x(i-1)x_i}}{m_{x(i-1)x_i}}\right) = \sum_{i=1}^L r_{x(i-1)x_i}$$

<b>r</b>	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	-0.740	0.419	0.580	-0.803
<b>C</b>	-0.913	0.302	1.812	-0.685
<b>G</b>	-0.624	0.461	0.331	-0.730
<b>T</b>	-1.169	0.573	0.393	-0.679

## Problem 1:

- **Input:** Small sequence **S**
- **Output:** Is **S** from a CpG island?
  - Build Markov Models:  $M^+$  &  $M^-$
  - Then compare

## Problem 2:

- **Input:** Long sequence **S**
- **Output:** Identify the CpG islands in **S**.
  - Markov models are inadequate.
  - Need Hidden Markov Models.