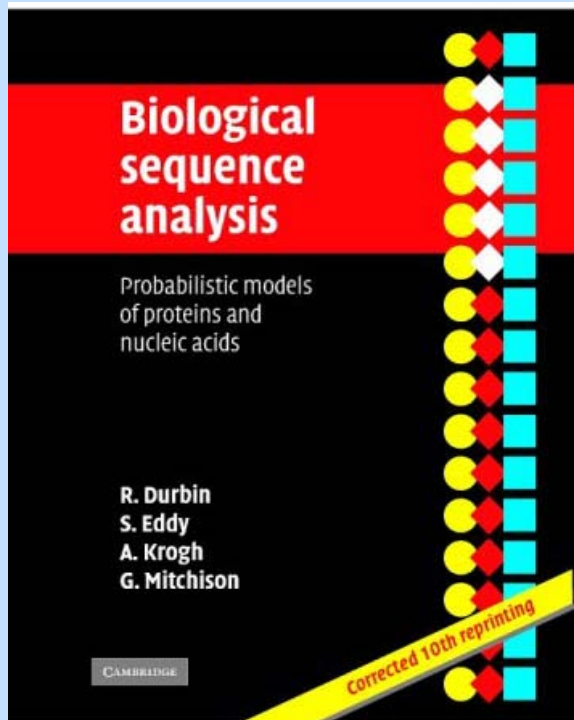


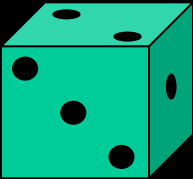
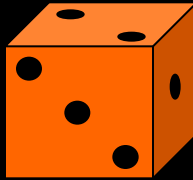
Hidden Markov Models



Main source: Durbin et al.,
“Biological Sequence Alignment”
(Cambridge, ‘98)




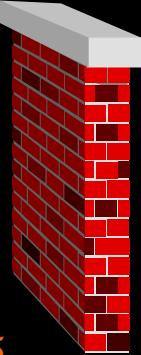
The occasionally dishonest casino

A   B

$P_A(1) =$
 $P_A(2) =$
 $\dots = 1/6$

$P_{A \rightarrow B} =$
 $P_{B \rightarrow A} =$
 $1/10$

$P_B(1) = 0.1$
 \dots
 $P_B(5) = 0.1$
 $P_B(6) = 0.5$



13652656643662612564
13652656643662612564

Can we tell when the loaded die is used?



Example - CpG islands

- CpG islands:
 - DNA stretches (100~1000bp) with frequent CG pairs (contiguous on same strand).
 - Rare, appear in significant genomic parts.
- Problem (1): Given a short genome sequence, decide if it comes from a CpG island.



Preliminaries: Markov Chains

$$(S, A, p)$$

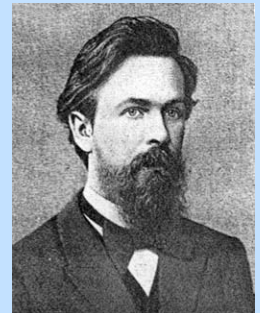
- S : State set
- p : Initial state prob. vector $\{p(x_1=s)\}$
[alternatively, use a begin state]
- A : Transition prob. matrix $a_{st} = P(x_i=t \mid x_{i-1}=s)$

Assumption: $X=x_1 \dots x_n$ is a random process with *memory length 1*, i.e.: $\forall s_i \in S$

$$P(x_i=s_i \mid x_1=s_1, \dots, x_{i-1}=s_{i-1}) = P(x_i=s_i \mid x_{i-1}=s_{i-1}) = a_{s_{i-1}, s_i}$$

- Sequence probability:

$$P(X) = p(x_1) \cdot \prod_{i=2 \dots L} a_{x_{i-1}, x_i}$$



Markov Models

- - Transition probs for non-CpG islands
- + Transition probs for CpG islands

-	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

+	A	C	G	T
A	0.180	0.274	0.425	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182



CpG islands: Fixed Window

- **Problem (1):** Given a short genome sequence X , decide if it comes from a CpG island.
- **Solution:** Model by a Markov chain. Let
 - a_{st}^+ : transition prob. in CpG islands,
 - a_{st}^- : transition prob. outside CpG islands.Decide by log-likelihood ratio score:

$$\text{score}(X) = \log \frac{P(X \mid \text{CpG-island})}{P(X \mid \text{non-CpG-island})} = \sum_{i=1}^n \log \frac{a_{x_{i-1}, x_i}^+}{a_{x_{i-1}, x_i}^-}$$

$$\text{bits_score}(X) = \frac{1}{n} \sum_{i=1}^n \log_2 \frac{a_{x_{i-1}, x_i}^+}{a_{x_{i-1}, x_i}^-}$$



Discrimination of sequences via Markov Chains

48 CpG islands, tot length ~60K nt. Similar non-CpG.

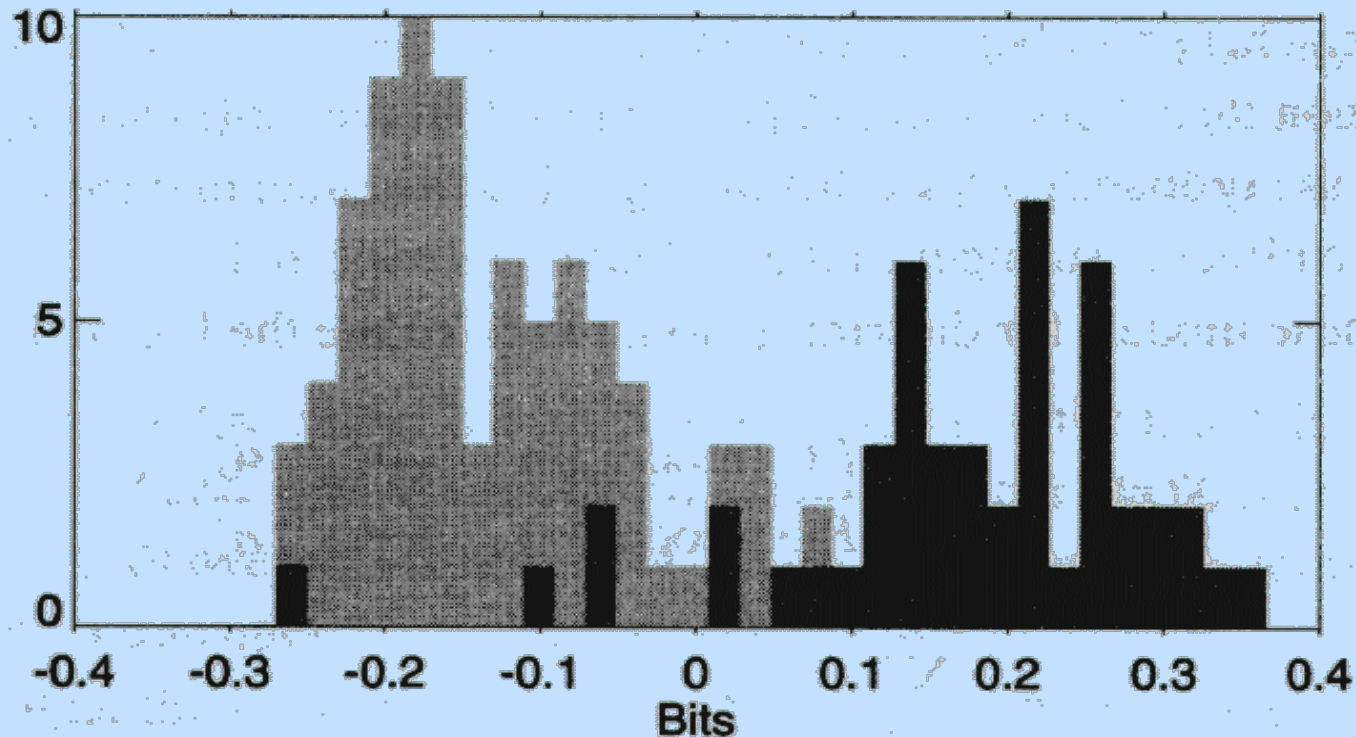


Figure 3.2 *The histogram of the length-normalised scores for all the sequences. CpG islands are shown with dark grey and non-CpG with light grey.*

Durbin et. al, Fig. 3.2



CpG islands - the general case

- **Problem(2)**: Detect CpG islands in a long DNA sequence.
- **Naive Solution - Sliding windows**: $\forall 1 \leq k \leq L-l$,
 - window: $X^k = (x_{k+1}, \dots, x_{k+l})$
 - score: $\text{score}(X^k)$
 - positive score \Rightarrow potential CpG island

Disadvantage: what is the length of the islands? How do we identify transitions?

Idea: Use Markov chains as before, with additional (hidden) states



Hidden Markov Model (HMM)

Finite set of **states**, capable of emitting symbols.

Example:

$Q = \{A_+, C_+, G_+, T_+, A_-, C_-, G_-, T_-\}$

Alphabet of **symbols**

Example: {A, C, G, T}

$M = (\Sigma, Q, \Theta)$

◆ $\Theta = (A, E)$

◆ **A: Transition**
prob. $a_{kl} \forall k, l \in Q$

◆ **E: Emission**
prob. $e_k(b) \forall k \in Q, b \in \Sigma$

path $\Pi = \pi_1, \dots, \pi_n$ (sequence of states - simple Markov chain; convention: π_0 - begin, π_{L+1} - end)

Given sequence $X = (x_1, \dots, x_L)$:

- $a_{kl} = P(\pi_i = l \mid \pi_{i-1} = k)$,
- $e_k(b) = P(x_i = b \mid \pi_i = k)$

$$P(X, \Pi) = a_{\pi_0, \pi_1} \cdot \prod_{i=1}^{L-1} e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}}$$

Ex.: express $P(X, \Pi)$ in terms of counts



Viterbi's Decoding Algorithm

(finding most probable state path)

Want: path Π maximizing $P(\mathbf{X}, \Pi)$

$v_k(i)$ = prob. of most probable path ending in state k at step i .

Init: $v_0(0) = 1$; $v_k(0) = 0 \quad \forall k > 0$

Step: $v_k(i+1) = e_k(x_{i+1}) \cdot \max_l \{v_l(i) \cdot a_{lk}\}$

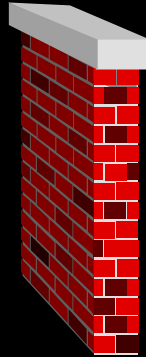
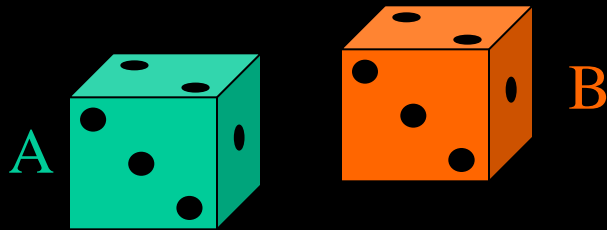
End: $P(\mathbf{X}, \Pi^*) = \max_l \{v_l(L) \cdot a_{l0}\}$

Time complexity: $O(Ln^2)$ for n states, L steps

Can find Π^* using back pointers.

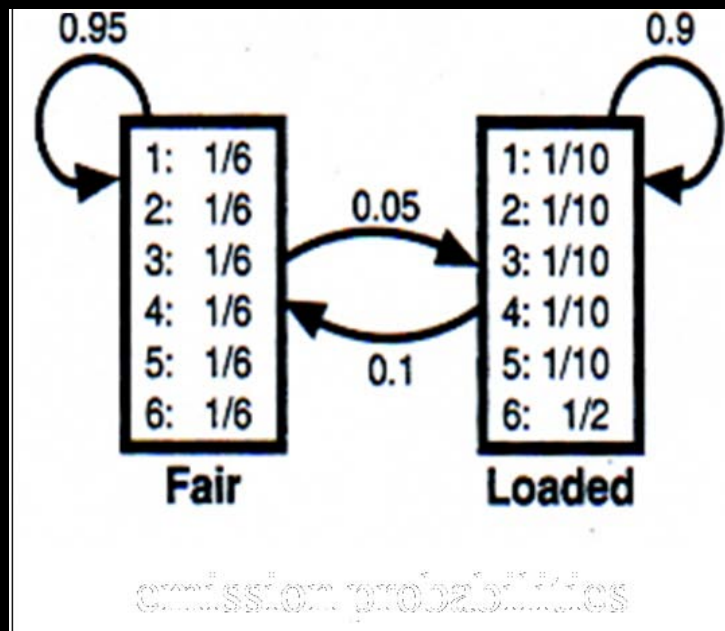


The occasionally dishonest casino (2)



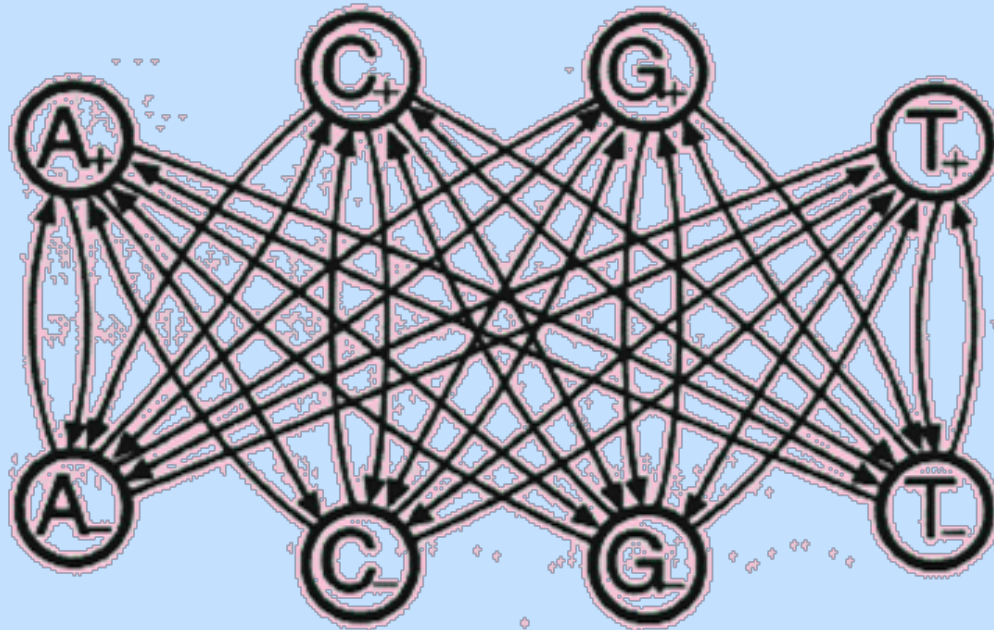
13652656643662612564

13652656643662612564



HMM for CpG Islands

- States: A_+ C_+ G_+ T_+ A_- C_- G_- T_-
- Symbols: A C G T A C G T
- Path $\Pi = \pi_1, \dots, \pi_n$: sequence of states



+	A	C	G	T
A	0.180	0.274	0.425	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

-	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

transition prob.

Posterior State Probabilities

Goal: calculate $P(\pi_i=k | X)$

- Our strategy:

- $P(X, \pi_i=k) =$

$$= P(x_1, \dots, x_i, \pi_i=k) \cdot P(x_{i+1}, \dots, x_L | x_1, \dots, x_i, \pi_i=k)$$

$$= P(x_1, \dots, x_i, \pi_i=k) \cdot P(x_{i+1}, \dots, x_L | \pi_i=k)$$

- $P(\pi_i=k | X) = P(\pi_i=k, X) / P(X)$

→ Need to compute these two terms - and $P(X)$



Forward Algorithm

Goal: calculate $P(X) = \sum_{\Pi} P(X, \Pi)$

Approximation: take max path Π^* from Viterbi alg.

Not justified when \exists several near maximal paths

Exact alg : "Forward Algorithm"

$$f_k(i) = P(x_1, \dots, x_i, \pi_i = k)$$

- Init: $f_0(0) = 1$; $f_k(0) = 0 \quad \forall k > 0$
- Step: $f_k(i+1) = e_k(x_{i+1}) \cdot \sum_l f_l(i) \cdot a_{lk}$
- End: $P(X) = \sum_l f_l(L) \cdot a_{l0}$



Backward Algorithm

- $b_k(i) = P(x_{i+1}, \dots, x_L \mid \pi_i = k)$
- init: $\forall k, b_k(L) = a_{k0}$
- step: $b_k(i) = \sum_l a_{kl} \cdot e_l(x_{i+1}) \cdot b_l(i+1)$
- End: $P(X) = \sum_k a_{0k} \cdot e_k(x_1) \cdot b_k(1)$



Posterior State Probabilities (2)

Goal: calculate $P(\pi_i=k | X)$

- Recall:
 - $f_k(i) = P(x_1, \dots, x_i, \pi_i=k)$
 - $b_k(i) = P(x_{i+1}, \dots, x_L | \pi_i=k)$
 - Each can be used to compute $P(X)$
- $P(X, \pi_i=k) =$
 - $= P(x_1, \dots, x_i, \pi_i=k) \cdot P(x_{i+1}, \dots, x_L | x_1, \dots, x_i, \pi_i=k)$
 - $= P(x_1, \dots, x_i, \pi_i=k) \cdot P(x_{i+1}, \dots, x_L | \pi_i=k)$
 - $= f_k(i) \cdot b_k(i)$
- $P(\pi_i=k | X) = P(\pi_i=k, X) / P(X)$



Dishonest Casino (3)

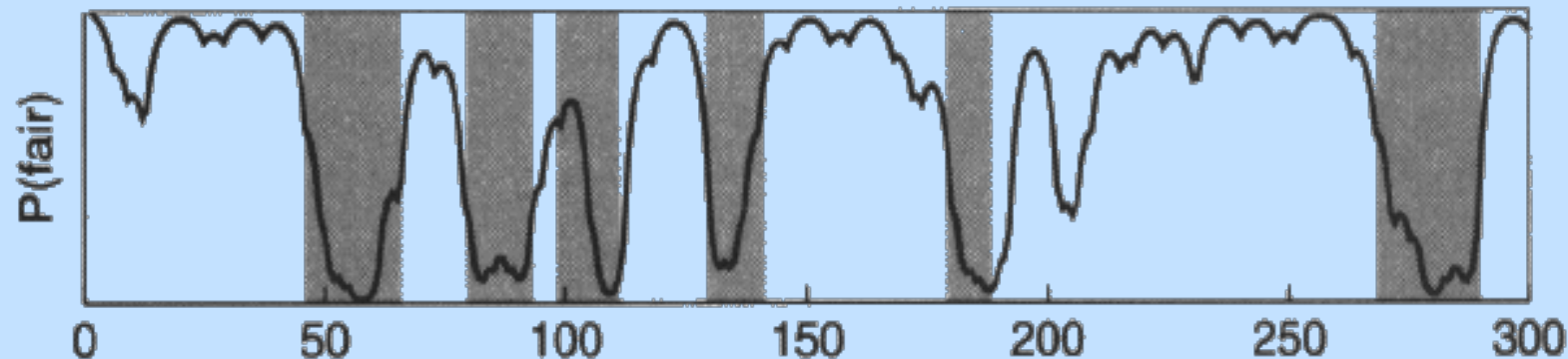


Figure 3.6 *The posterior probability of being in the state corresponding to the fair die in the casino example. The x axis shows the number of the roll. The shaded areas show when the roll was generated by the loaded die.*



Posterior Decoding

- Now we have $P(\pi_i=k | X)$. How do we decode?
1. $\pi_i^* = \operatorname{argmax}_k P(\pi_i=k | X)$
 - Good when interested in state at particular point
 - path of states π_1^*, \dots, π_L^* may not be legal
 2. Define a function of interest $g(i)$ on the states. Compute $G(i|X) = \sum_k P(\pi_i=k | X) \cdot g(k)$
 - E.g.: $g(i) = 1$ for states in S , 0 on the rest: $G(i|X)$ is posterior prob of symbol i coming from S

e.g., CpG island
 $S = \{A_+, C_+, G_+, T_+\}$



Parameter Estimation for HMMs

Log likelihood of model

Input: X^1, \dots, X^n independent **training sequences**

Goal: estimation of $\Theta = (A, E)$ (model parameters)

Note: $P(X^1, \dots, X^n | \Theta) = \prod_{i=1 \dots n} P(X^i | \Theta)$ (indep.)

$$l(x^1, \dots, x^n | \Theta) = \log P(X^1, \dots, X^n | \Theta) = \sum_{i=1 \dots n} \log P(X^i | \Theta)$$

Case 1 - Estimation When State Sequence is Known:

A_{kl} = #(occurred $k \rightarrow l$ transitions)

$E_k(b)$ = #(emissions of symbol b that occurred in state k)

Max. Likelihood Estimators:

- $a_{kl} = A_{kl} / \sum_{l'} A_{kl'}$
- $e_k(b) = E_k(b) / \sum_{b'} E_k(b')$

small sample, or
prior knowledge correction:

$$\begin{aligned} A'_{kl} &= A_{kl} + r_{kl} \\ E'_k(b) &= E_k(b) + r_k(b) \end{aligned}$$

- “Dirichlet priors”



Parameter Estimation in HMM

Case 2: -Estimation When States are Unknown

Input: X^1, \dots, X^n indep training sequences

Baum-Welch alg. (1972):

★ Expectation:

- compute expected no. of $k \rightarrow l$ state transitions: (ex.)
$$P(\pi_i = k, \pi_{i+1} = l \mid X, \Theta) = [1/P(x)] \cdot f_k(i) \cdot a_{kl} \cdot e_l(x_{i+1}) \cdot b_l(i+1)$$
- $A_{kl} = \sum_j [1/P(X^j)] \cdot \sum_i f_k^j(i) \cdot a_{kl} \cdot e_l(x_{i+1}^j) \cdot b_l^j(i+1)$
- compute expected no. of symbol b appearances in state k
$$E_k(b) = \sum_j [1/P(X^j)] \cdot \sum_{\{i \mid x_{i+1}^j = b\}} f_k^j(i) \cdot b_k^j(i) \text{ (ex.)}$$

★ Maximization:

- re-compute new parameters from A, E using max. likelihood.

- guarantees convergence
- monotone
- many local optima
- special case of EM

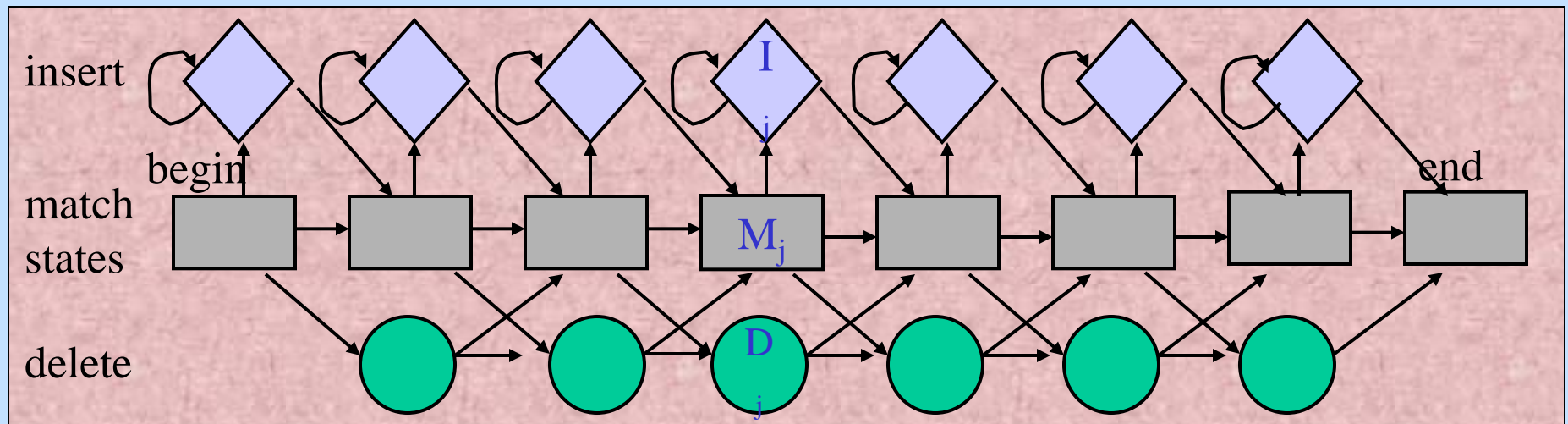
repeat (1)+(2) until improvement $\leq \epsilon$



Profile HMMs (Hausssler et al, 1993)

- Ungapped alignment of X against a profile M :
 - $e_i(a) = \text{prob. of observing } a \text{ at position } i.$
 - $P(X | M) = \prod_{i=1}^L e_i(x_i),$ or
 - $\text{Score}(X | M) = -\sum_{i=1}^L \log [e_i(x_i) / q_{x_i}]$
- indels: $AG---C$

background prob.



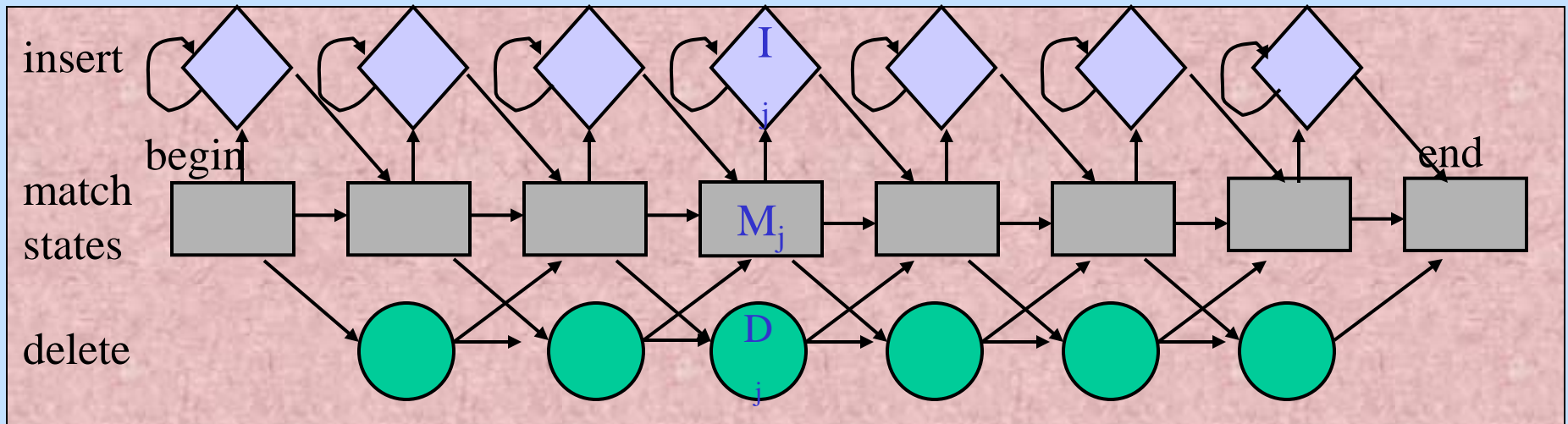
Profile HMMs

No log odds contribution from the emission

- Gapped alignment of X against a profile M :
assume $e_{I_j}(a) = q_a$ (background prob.)

\Rightarrow gap of length k contributes to log-odds:

$$\left\{ \begin{array}{l} \log(a_{M_j, I_j}) + \log(a_{I_j, M_{j+1}}) \\ \text{gap open} \end{array} \right\} + (k-1) \cdot \log(a_{I_j, I_j}) \quad (\text{gap extension})$$



Profile HMM

- Transition Probabilities

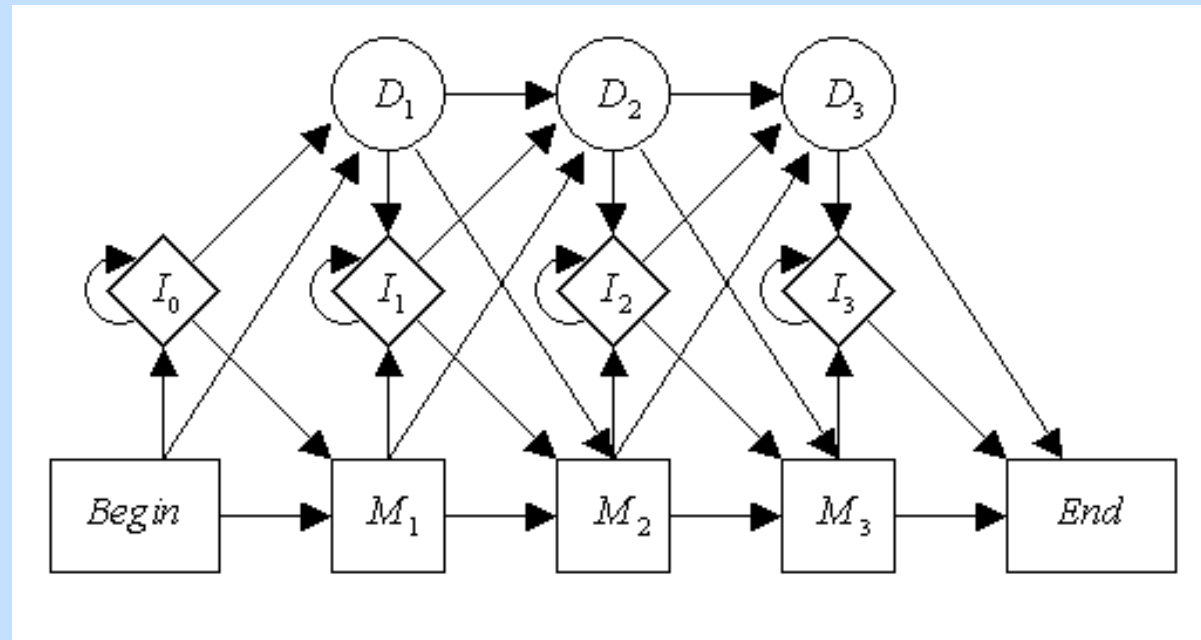
- $M_i \rightarrow M_{i+1}$
- $M_i \rightarrow D_{i+1}$
- $M_i \rightarrow I_i$

- $I_i \rightarrow M_{i+1}$
- $I_i \rightarrow I_i$
- $I_i \rightarrow D_{i+1}$

- $D_i \rightarrow D_{i+1}$
- $D_i \rightarrow M_{i+1}$
- $D_i \rightarrow I_i$

- Emission probabilities

- $M_i \rightarrow a$
- $I_i \rightarrow a$



Example

- Suppose we are given the aligned sequences

****---***

AG---C

A-AT-C

AG-AA-

--AAAC

AG---C

- Suppose also that the "match" positions are marked...



Calculating A, E

count transitions and emissions:

transitions		**---*	emissions						
	0	1	2	3		0	1	2	3
M-M					A				
M-D					A-AT-C				
M-I					AG-AA-				
I-M					--AAAC				
I-D					AG---C				
I-I									
D-M									
D-D									
D-I									



Calculating A, E

count transitions and emissions:

	transitions			
	0	1	2	3
M-M	4	3	2	4
M-D	1	1	0	0
M-I	0	0	1	0
I-M	0	0	2	0
I-D	0	0	1	0
I-I	0	0	4	0
D-M	-	0	0	1
D-D	-	1	0	0
D-I	-	0	2	0

**** --- ***
AG --- C
A-AT-C
AG-AA-
--AAAC
AG --- C

	emissions			
	0	1	2	3
A	-	4	0	0
C	-	0	0	4
G	-	0	3	0
T	-	0	0	0
A	0	0	6	0
C	0	0	0	0
T	0	0	1	0
G	0	0	0	0

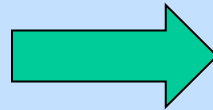


Estimating Maximum Likelihood probabilities using fractions

emissions

	0	1	2	3
A	-	4	0	0
C	-	0	0	4
G	-	0	3	0
T	-	0	0	0

A	0	0	6	0
C	0	0	0	0
T	0	0	1	0
G	0	0	0	0



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

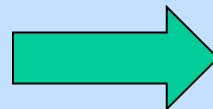
A	.25	.25	.86	.25
C	.25	.25	0	.25
T	.25	.25	.14	.25
G	.25	.25	0	.25



Estimating ML probabilities (contd)

transitions

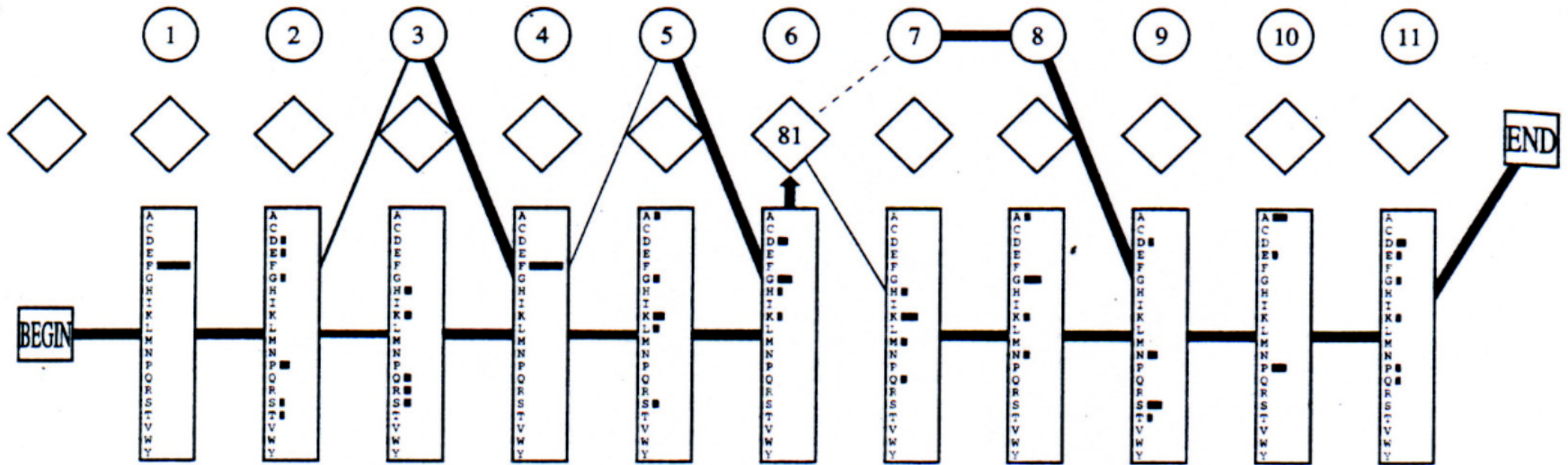
	0	1	2	3
M-M	4	3	2	4
M-D	1	1	0	0
M-I	0	0	1	0
I-M	0	0	2	0
I-D	0	0	1	0
I-I	0	0	4	0
D-M	-	0	0	1
D-D	-	1	0	0
D-I	-	0	2	0



	0	1	2	3
M-M	.8	.75	.66	1.0
M-D	.2	.25	0	0
M-I	0	0	.33	0
I-M	.33	.33	.28	.33
I-D	.33	.33	.14	.33
I-I	.33	.33	.57	.33
D-M	-	0	0	1
D-D	-	1	0	0
D-I	-	0	1	0



HMM from multiple alignment...



```

FPHF-DLS-----HGSAQ
FESFGDLSTPDAVMGNPK
FDRFKHLKTEAEMKASED
FTQFAG-KDLESIKGTAP
FPKFKGLTTADQLKKSAD
FS-FLK-GTSEVPQNNPE
FG-FSG-----AS---DPG
    
```

Shade:
insert



.. and multiple alignment from a given HMM

- Align each sequence to the profile separately

```
FPHF-Dls.....HGSAQ      FS-FLKngvdptaai--NPK
FESFGDlstpdavMGNPK      FESFGDlstpdav..MGNPK
FDRFKHlkteaemKASED      FDRFKHlkteaem..KASED
FTQFAGkdlesi.KGTAP      FTQFAGkdlesi...KGTAP
FPKFKGlttadqlKKSAD      FPKFKGlttadql..KKSAD
FS-FLKgtsevp.QNNPE      FS-FLKgtsevp...QNNPE
FG-FSGas.....--DPG      FG-FSGas.....--DPG
```

- Right: a new sequence realigned with the model
- Inserts are unaligned.



Searching with Profile HMMs

Compute: $\log \frac{P(X|Model)}{P(X|random)}$

$$\prod_i q_{x_i}$$

$V_j^M(i)$: log odds of best path matching x_1, \dots, x_i to submodel up to level j , ending with x_i emitted by state M_j
 $V_j^I(i)$: same, ending with x_i emitted by state I_j
 $V_j^D(i)$: score for best path ending in D_j after x_i has been emitted (and x_{i+1} has not been emitted yet)

- $$V_j^M(i) = \log [e_{M_j}(x_i) / q_{x_i}] + \max \left\{ \begin{array}{l} V_{j-1}^M(i-1) + \log(a_{M_{j-1}, M_j}), \\ V_{j-1}^I(i-1) + \log(a_{I_{j-1}, M_j}), \\ V_{j-1}^D(i-1) + \log(a_{D_{j-1}, M_j}) \end{array} \right\}$$
- $$V_j^I(i) = \log [e_{I_j}(x_i) / q_{x_i}] + \max \left\{ \begin{array}{l} V_j^M(i-1) + \log(a_{M_{j-1}, I_j}), \\ V_j^I(i-1) + \log(a_{I_j, I_j}), \\ V_j^D(i-1) + \log(a_{D_j, I_j}) \end{array} \right\}$$
- $$V_j^D(i) = \max \left\{ \begin{array}{l} V_{j-1}^M(i) + \log(a_{M_{j-1}, D_j}), \\ V_{j-1}^I(i) + \log(a_{I_{j-1}, D_j}), \\ V_{j-1}^D(i) + \log(a_{D_{j-1}, D_j}) \end{array} \right\}$$



Training a profile HMM from unaligned sequences

- choose length of the profile HMM, initialize parameters
- train the model using Baum-Welch
- obtain MA with the resulting profile as before.

(Formulas incl. forward/backward - in handout)



An Illustrative Study

"HMMs in Computational Biology", Krogh, Brown,
Mian, Sjoalnder, Haussler '93

- Globin experiment:
 - Heme-containing proteins, involved in the storage and transport of oxygen
 - 625 globins from Swissprot, ave. length 145AA
 - Training set: 400 sequences
 - Built model and ran it against test globins and all Swissprot proteins (25K).



Representative globins

```

Helix      AAAAAAAAAAAAAAAAAA  BBBBBBBBBBBBBBBBCCCCCCCCC  DDDDDDEE
HBA_HUMAN  -----VLSPADKTNVKAAWGKVGA--HAGEYGAEALERMFLSFPTTKTYFPHF-DLS-----HGSA
HBB_HUMAN  -----VHLTPEEKSAVTALWGKV----NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
MYG_PHYCA  -----VLSEGEWQLVLHVWAKVEA--DVAGHGQDILIRLFKSHPETLEKFRDFKHLKTEAEMKASE
GLB3_CHITP -----LSADQISTVQASFDKVKG-----DPVGILYAVFKADPSIMAKFTQFAG-KDLESIKGTA
GLB5_PETMA PIVDTGSVAPLSAAEKTIRSAPVYS--TYETSGVDILVKFFTSTPAAQEFPKFKGLTTADQLKSA
LGB2_LUPLU -----GALTESQAALVKSSWEEFNA--NIPKHTHRFFILVLEIAPAAKDLFS-FLK-GTSEVPQNNP
GLB1_GLYDI -----GLSAAQRQVIAATWKDIAGADNGAGVVKDCLIKFLSAHPQMAAVFG-FSG----AS---DP
    
```

```

Helix      EEEEEEEEEEEEEEEEE  FFFFFFFF  FGGGGGGGGGGGGGGGGGG
HBA_HUMAN  QVKGHGKQVADALTNVAHV---D--DMPNALSALSDLHAHKL--RVDPVNFKLLSHCLLVTLAAHLP
HBB_HUMAN  KVKAHGKQVLAGFSDGLAHL---D--NLKGTATLSELHCDKL--HVDPENFRLLGNVLCVLAHFGKE
MYG_PHYCA  DLKKHGVTVLTALGAILKK---K-GHHEAELKPLAQSHATKH--KIPIKYLEFISEAIIHVLHSRHPG
GLB3_CHITP PFETHANRIVGFFSKIIGEL--P---NIEADVNTFVASHKPRG---VTHDQLNFRAGFVSYMAHT--D
GLB5_PETMA DVRWHAERIINAVNDAVASM--DDTEKMSMKLRDLGKHAHSF--QVDPQYFKVLAAVIADTVAAG----
LGB2_LUPLU ELQAHAGKVFKLVEAAIQLQVTGVVTDATLKNLGSVHVSKG---VADAHFPVVKAILKTIKEVVGAK
GLB1_GLYDI GVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKAVGVRHKGYGKHIKAQYFEPLGASLLSAMEHRIGK
    
```

```

Helix      HHHHHHHHHHHHHHHHHHHHHHHHHHH
HBA_HUMAN  FTPAVHASLDKFLASVSTVLTISKYR-----
HBB_HUMAN  FTTPVQAAYQKVAVAGVANALAHKYH-----
MYG_PHYCA  FGADAQGAMNKALELFRKDIAAKYKELGYQG
GLB3_CHITP FA-GAEAAWGATLDTFFGMIFSKM-----
GLB5_PETMA -----DAGFEKLMSMICILLRSAY-----
LGB2_LUPLU WSEELNSAWTIAAYDELAIVIKKEMNDAA---
GLB1_GLYDI MNAAAKDAAWAAAYADISGALISGLQS-----
    
```

Alignment from
Bashford et al 87
A-H: alpha helices



Realignment by trained HMM

```

Helix          AAAAAAAAAAAAAAAAAA  BBBB BBBB BBBB BBBB BBBBBB CCCCCCCCCC  DDDDDDDDEE
                *****+  ++++++*****+  +
HBA_HUMAN  V.....LSPADKTNVKAAWGKVG..HAGEYGAEALERMFLSFPTTKTYFPHF--DLSHGSAQ----
HBB_HUMAN  Vh.....LTPEEKSAVTALWGKV--.NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
MYG_PHYCA  V.....LSEGEWQLVLHVWAKVEA..DVAGHGQDILIRLFKSHPETLEKFDKFRFKHLKTEAEMKASE
GLB3_CHITP -.....LSADQISTVQASFDKV--.KGDVPG--ILYAVFKADPSIMAKFTQF-AGKDLESIKGTA
GLB5_PETMA PivdtgsvapLSAAEKTKIRSAWAPVYS..TYETSGVDILVKFFTSTPAAQEFFPKFKGLTTADQLKKS
LGB2_LUPLU Ga.....LTESQAALVKSSWEEFNA..NIPKHTRFFILVLEIAPAAKDLF-SFLKGTSEVPQ--NNP
GLB1_GLYDI G.....LSAAQRQVIAATWKDIAGadNGAGVVGKDCLIKFLSAHPQMAAVF-GF----SGASD---P
    
```

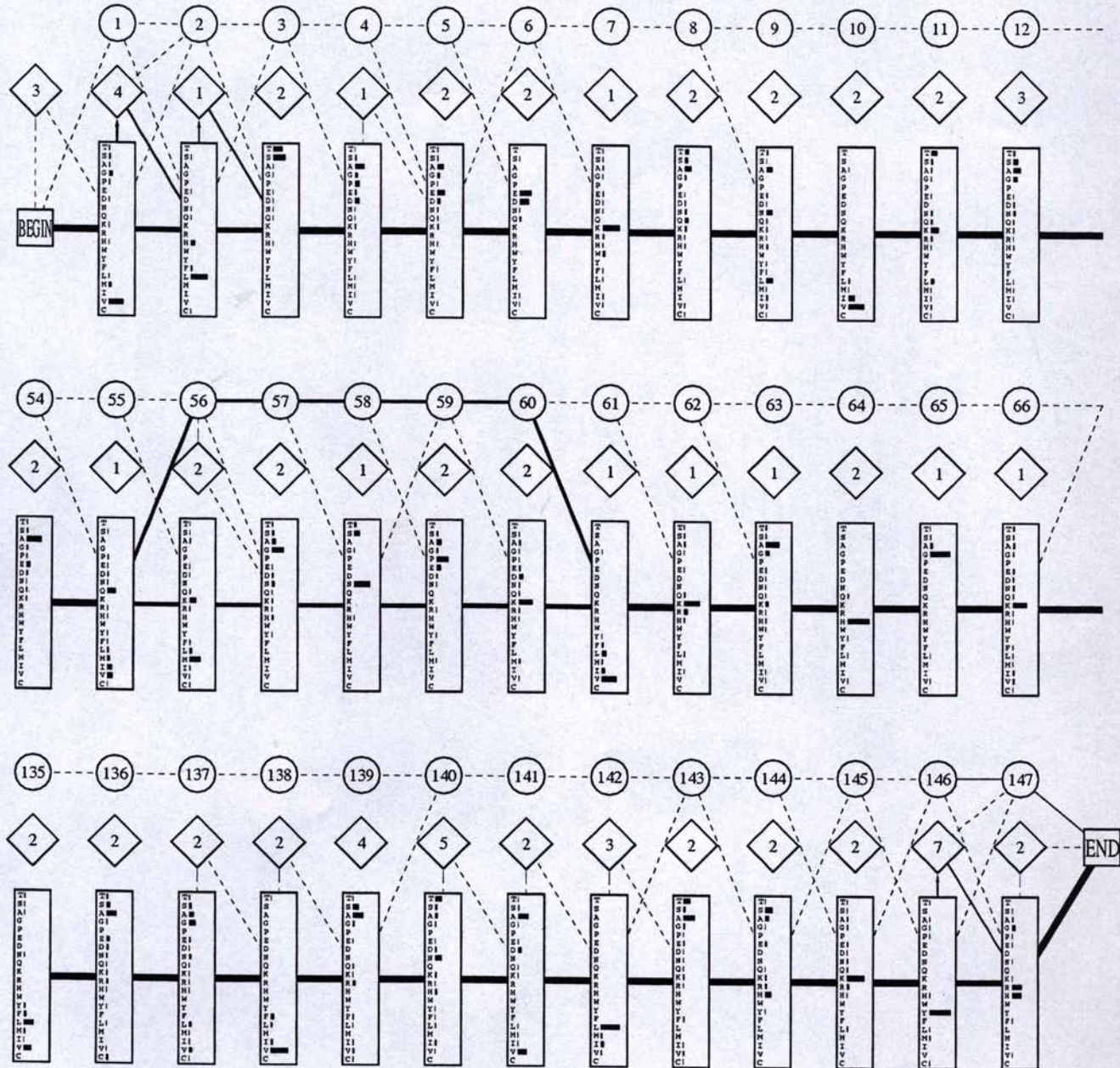
```

Helix          EEEEEEEEEEEEEEEEEEE  FFFFFFFF  FFFFFGGG  GGGGGGGGGGGGGGGG
                +*****+  *****  *****+
HBA_HUMAN  -VKGHGKKVADALTNVAHVDD....MPNALSALSDLHA...HKLRVDPV.NFKLLSHCLLVTLAAHLP
HBB_HUMAN  KVKAHGKKVLGAFSDGLAHLDN....LKGTFATLSELHC...DKLHVDPE.NFRLLGNVLCVLAHHFG
MYG_PHYCA  DLKKHGVTVLTAALGAILKKKGH....HEAELKPLAQSHA...TK-HKIPIkYLEFISEAIHVLHSRHP
GLB3_CHITP PFETHANRIVGFFSKIIIGELPN....IEADVNTFVASHK...PR-GVTHD.QLNFRAGFVSYMKAH--
GLB5_PETMA DVRWHAERIINAVNDAVASMDDtek..MSMKLRDLSGKHA...KSFQVDPQ.YFKVLAAVIADTVAA---
LGB2_LUPLU ELQAHAGKVFKLVYEAAIQLQVtgvvvTDTLKNLGSVHV...SK-GVADA.HFPVVKEAILKTIKEVVG
GLB1_GLYDI GVAALGAKVLAQIGVAVSHLGDegk..MVAQMKAVGVRHKgygNK-HIKAQ.YFEPLGASLLSAMEHRIG
    
```

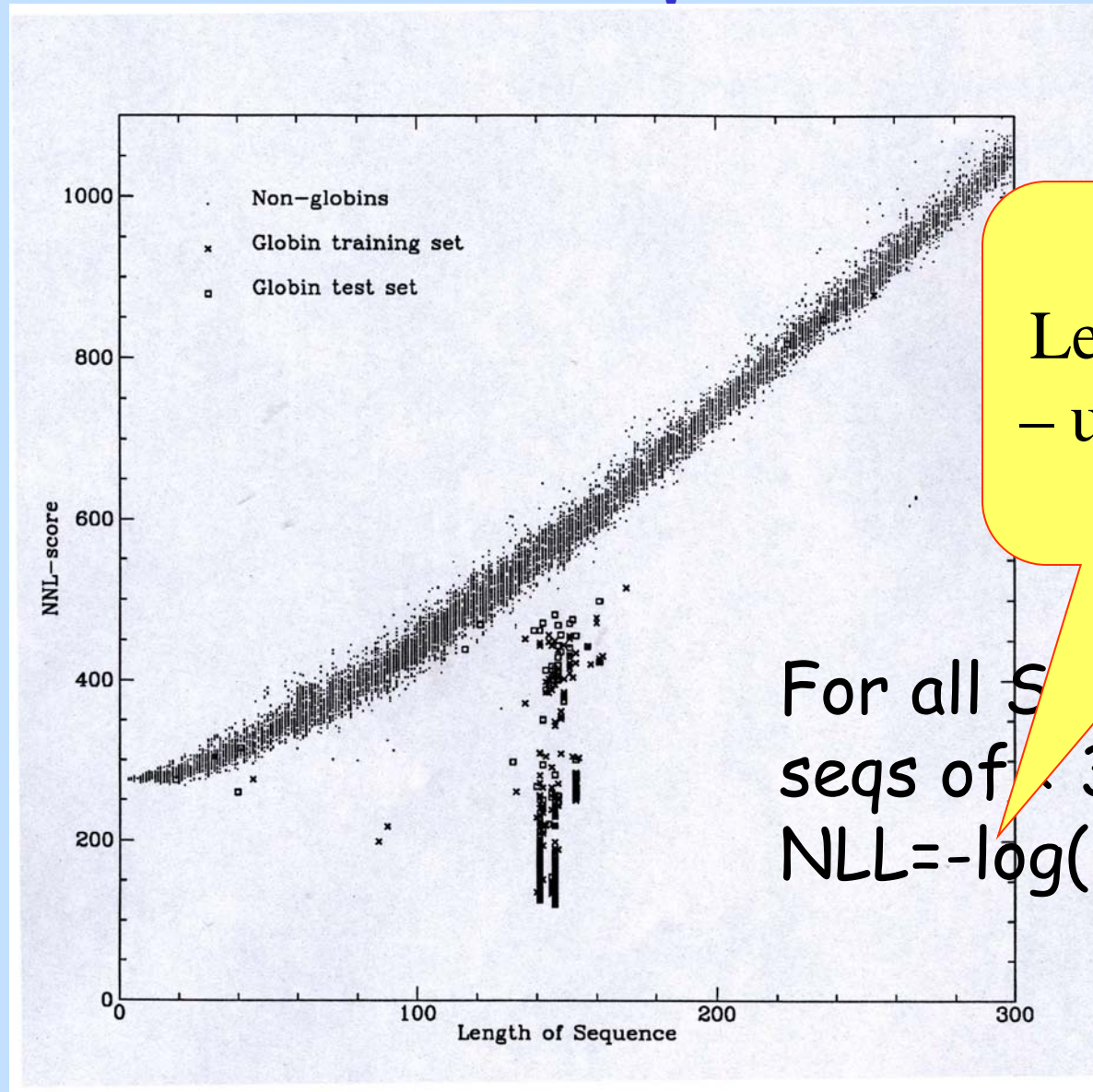
```

Helix          HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
                +*****+
HBA_HUMAN  AEFTP AVHASL DKFLASVSTVLT SKY.....R
HBB_HUMAN  KEFTPPVQAAYQKV VAGVANALAHKY.....H
MYG_PHYCA  GDFGADAQGAMNKALELFRKDIAAKYkelgyqG
GLB3_CHITP TDF-AGAEAAWGATLD TFFGMIFSKM.....-
GLB5_PETMA GD-----AGFEKLMSMICILLRSAY.....-
LGB2_LUPLU AKWSEELNSAWTIAYDELAIVIKKEMnda...A
GLB1_GLYDI GKMNAAAKDAWAAAYADISGALISGLq.....S
    
```


Part of the final globin model



Score vs sequence length

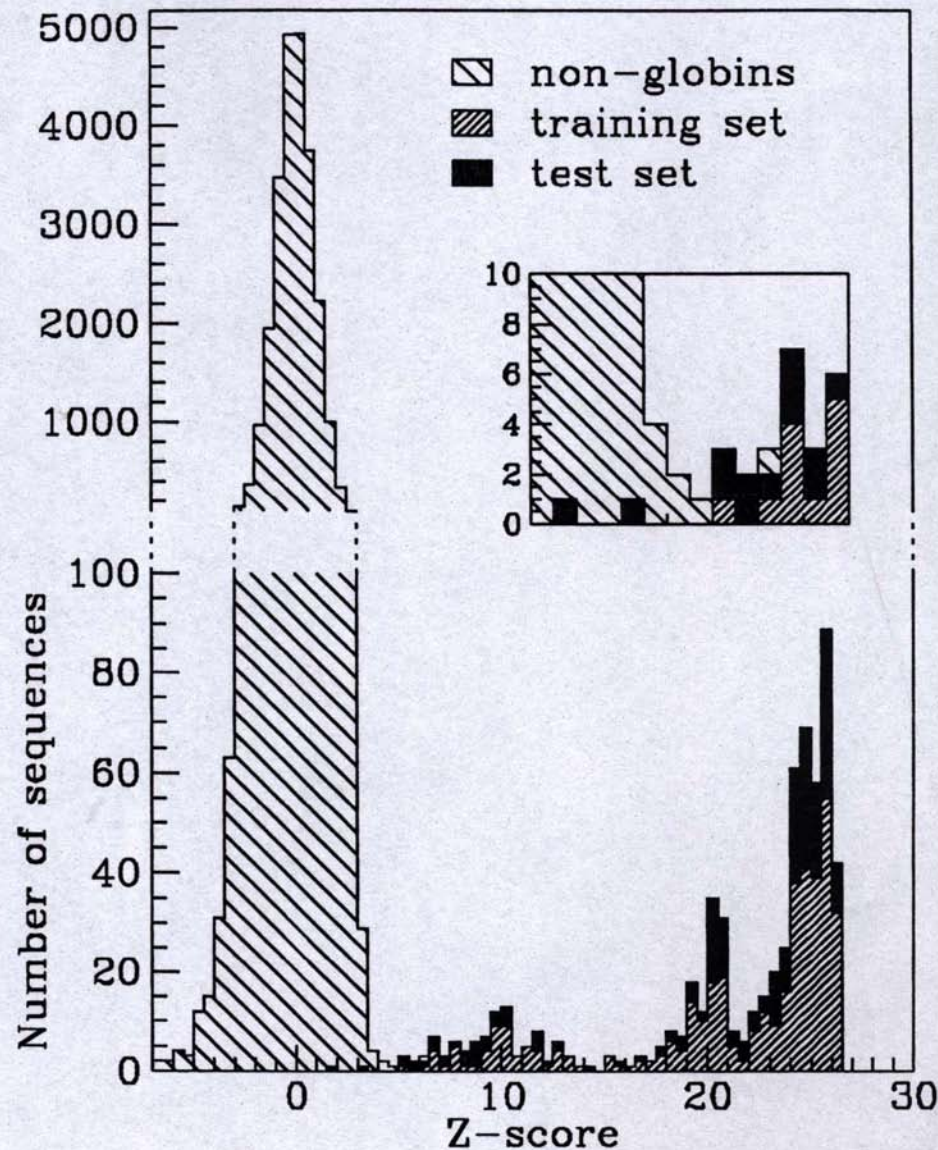


Negative LL.
Length dependent
– using log odds is preferable

For all Swissprot
seqs of < 300AA
 $LL = -\log(P(\text{seq}|\text{model}))$



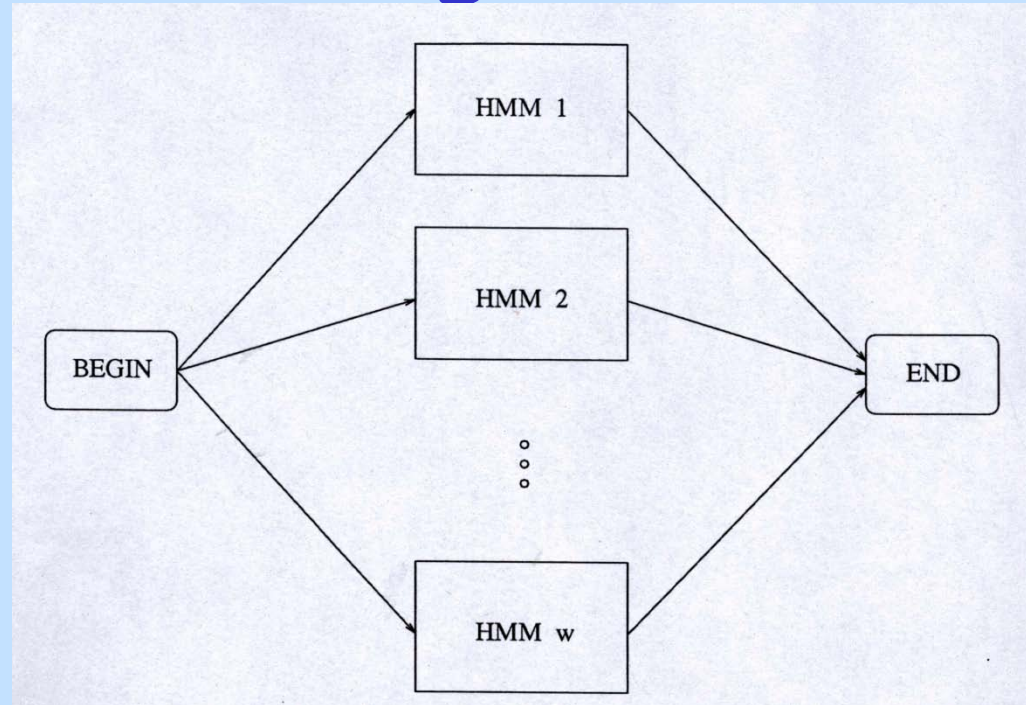
Z-score distribution



- Z-score:
 $(S - E(S)) / sd(S)$
- On ~25K proteins, cutoff 5 misses 2/628 globins with no fp



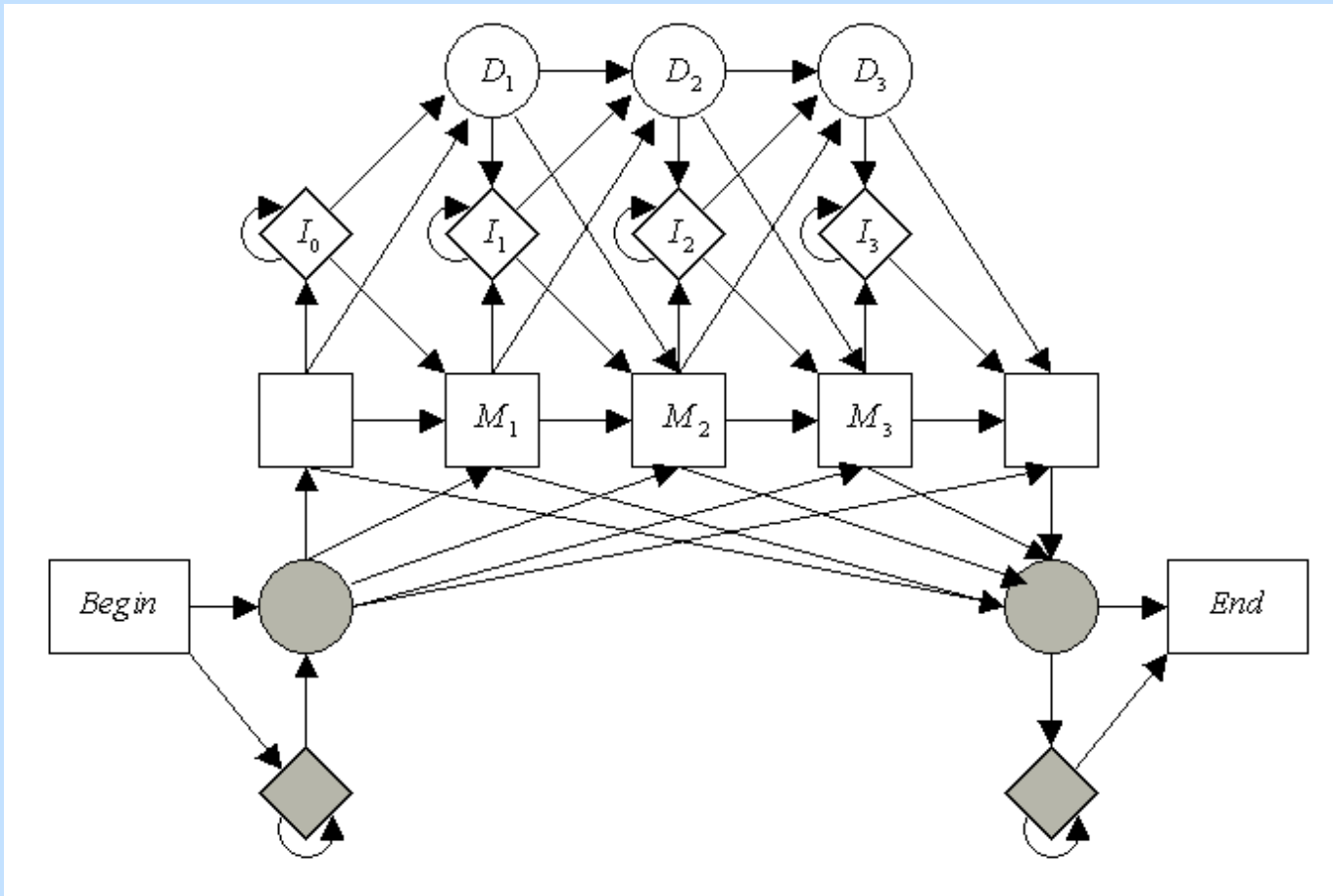
Discovering subfamilies



- 10 component HMM. Training set: 628 globins. Classified each sequence using the model.
- Generated 7 nonempty clusters, with ~560 falling into 4 "biologically meaningful" clusters.



Local alignment in HMM



Add **flanking states** for regions of unaligned seq.





Pfam 22.0 (July 2007, 9318 families)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

USING PFAM	YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...
SEQUENCE SEARCH	Analyze your protein sequence for Pfam matches
VIEW A PFAM FAMILY	View Pfam family annotation and alignments
VIEW A CLAN	See groups of related families
VIEW A SEQUENCE	Look at the domain organisation of a protein sequence
VIEW A STRUCTURE	Find the domains on a PDB structure
KEYWORD SEARCH	Query Pfam by keywords
	Or view the help pages for more information

