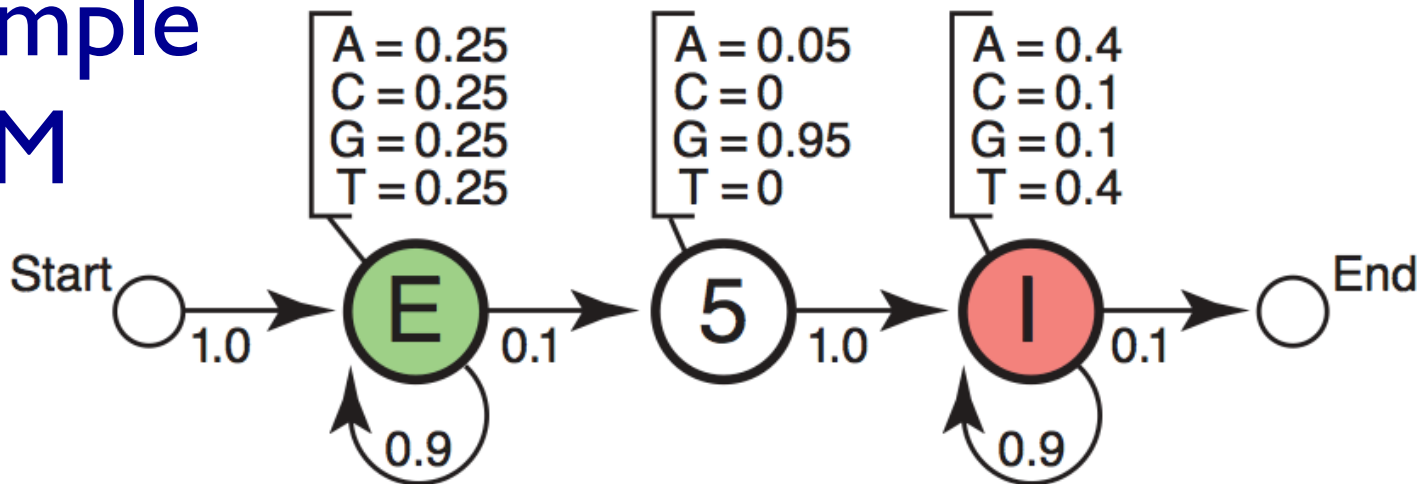


# Genome 559

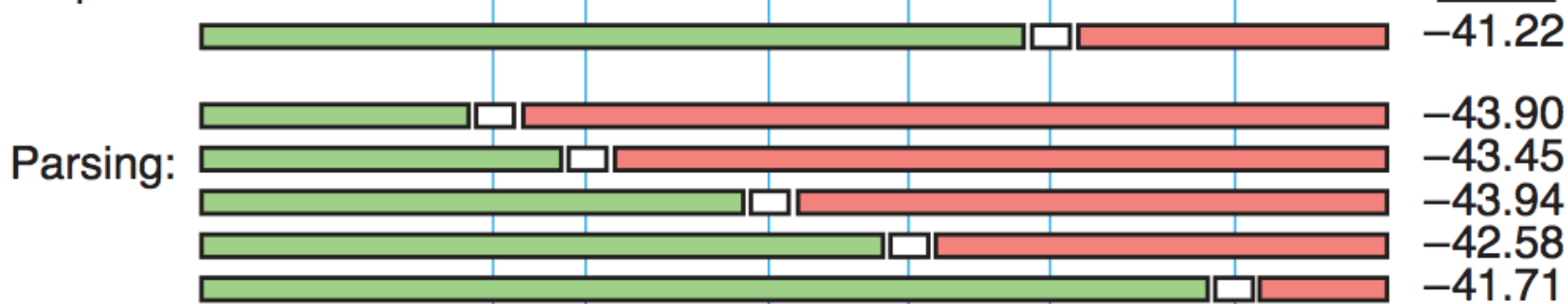
## Hidden Markov Models

# A simple HMM

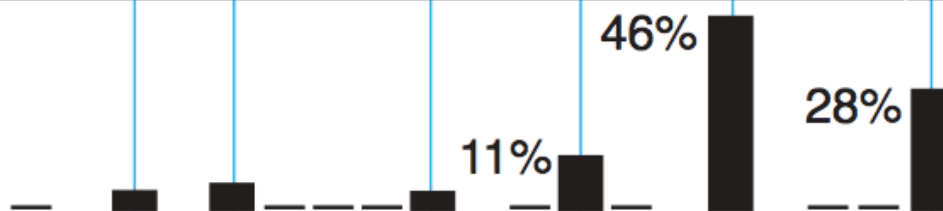


Sequence: **CTTCATGTGAAAGCAGACGTAAGTCA**

State path: **EEEEEEEEEEEEEEEEEE5IIIIIIII**  $\log P$



Posterior decoding:



# Notes

Probability of a given a state path and output sequence is just product of emission/transition probabilities

If state path is *hidden*, you need to consider *all* possible paths (usually exponentially many). E.g., find:

Total probability of a given seq (sum over all paths)

Probability of the most probable single path

“Dynamic programming” algorithms similar to seq alignment can solve these problems relatively quickly

# Viterbi: Most probable path



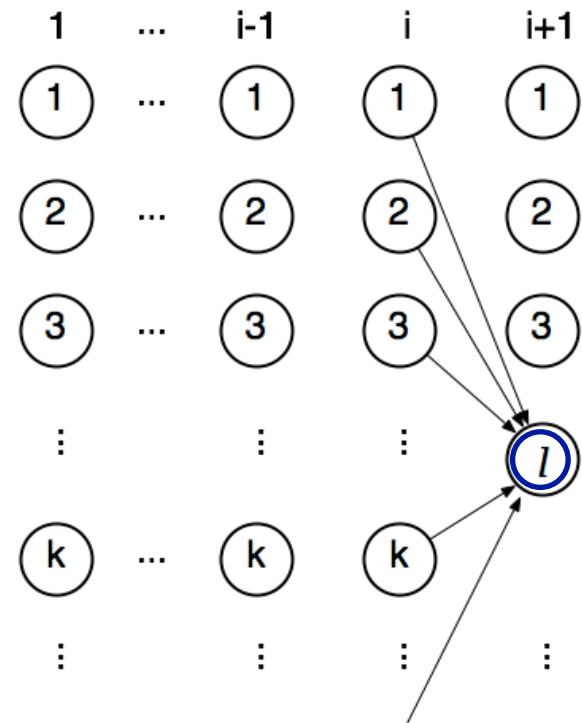
	C	A	G	A	T
E	.25	$.25^2 \cdot .9$	$.25^3 \cdot .9^2$	$.25^4 \cdot .9^3$	$.25^5 \cdot .9^4$
S	0	$.25 \cdot .1 \cdot .05$	$.25^2 \cdot .9 \cdot .1 \cdot .95$	$.25^3 \cdot .9^2 \cdot .1 \cdot .05$	0
I	0	0	$.25 \cdot .1 \cdot .05 \cdot .1$	$.25^2 \cdot .9 \cdot .1 \cdot .95 \cdot .9 \cdot .4$	...
				$.25 \cdot .1 \cdot .05 \cdot .1 \cdot .9 \cdot .4$	

# The Viterbi Algorithm

$v_l(i)$  = probability of the most probable path emitting  $x_1, x_2, \dots, x_i$  and ending in state  $l$

Initialize:

$$v_l(0) = \begin{cases} 1 & \text{if } l = \textit{Begin state} \\ 0 & \text{otherwise} \end{cases}$$



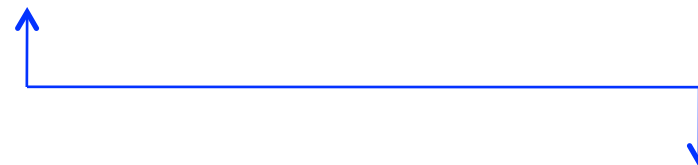
General case:

$$v_l(i + 1) = \underbrace{e_l(x_{i+1})}_{\text{emission}} \cdot \max_k (v_k(i) \underbrace{a_{k,l}}_{\text{transition}})$$

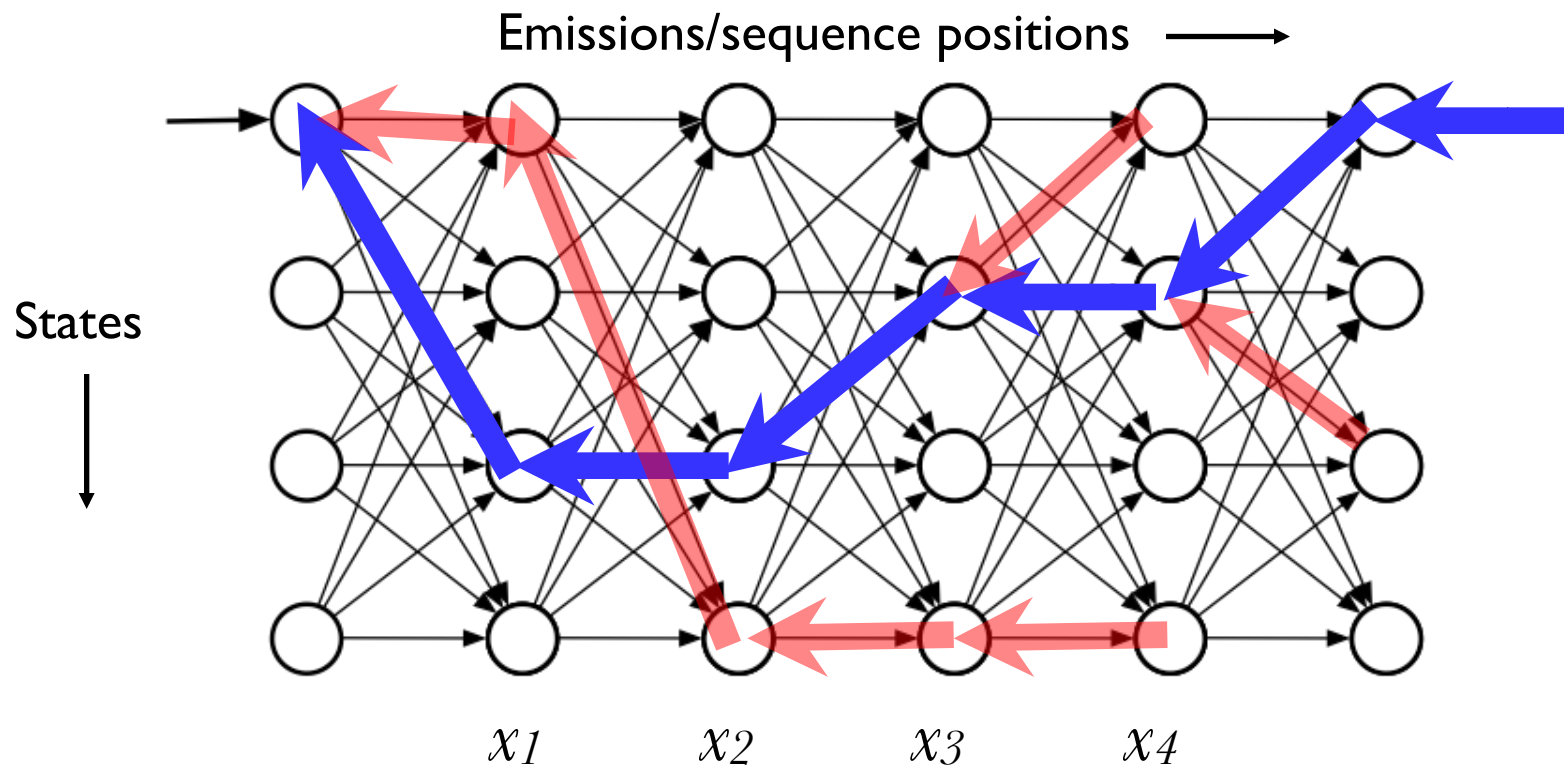
# Viterbi Traceback

Above finds *probability* of best path

To find the path itself, trace *backward* to the state  $k$  attaining the max at each stage


$$v_l(i + 1) = e_l(x_{i+1}) \cdot \max_k (v_k(i) a_{k,l})$$

# Viterbi Traceback



Viterbi score: 
$$v_l(i + 1) = e_l(x_{i+1}) \cdot \max_k (v_k(i) a_{k,l})$$

Viterbi path<sup>R</sup>: 
$$back_l(i + 1) = \arg \max_k (v_k(i) a_{k,l})$$

# An Application: Protein Alignments

```

Helix      AAAAAAAAAAAAAAAAAA      BBBBBBBBBBBBBBBBBBCCCCCCCCCCC
HBA_HUMAN  -----VLSPADKTNVKAAWGKVGGA--HAGEYGAEALERMFLSFPTTKTYFFHF
HBB_HUMAN  -----VHLTPEEKSAVTALWGKV----NVDEVGGEALGRLLVVYPWTQRFFESF
MYG_PHYCA  -----VLSSEGEWQLVLHVWAKVEA--DVAGHGQDILIRLFKSHPETLEKFDRF
GLB3_CHITP -----LSADQISTVQASFDKVKG-----DPVGILYAVFKADPSIMAKFTQF
GLB5_PETMA PIVDTGSVAPLSAAEKTIRSAWAPVYS--TYETS GVDILVKFFTSTPAAQEFFFKF
LGB2_LUPLU -----GALTESQAALVKSSWEEFNA--NIPKHTHRFFILVLEIAPAAKDLFS-F
GLB1_GLYDI -----GLSAAQRQVIAATWKDIAGADNGAGVGKDCLIKFLSAHPQMAAVFC-F
Consensus  LS.... v a w kv . . g . L . . F . F F
    
```

```

Helix      DDDDDDEEEEEEEEEEEEEEEEEEEEEEE      FFFFFFFF
HBA_HUMAN  -DLS-----HCSAQVKGHGKKVADALTNVAHV---D--DMPNALSALSDDLHAHKL-
HBB_HUMAN  GDLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHL---D--NLKGTFFATLSELHCDKL-
MYG_PHYCA  KHLKTEAEMKASEDLKKHGVTVLTAALGAILKK---K-GHHEAELKPLAQSHATKH-
GLB3_CHITP AG-KDLESIKCTAPFETHANRIVGFFSKIIGEL--P---NIEADVNTFVASHKPRG-
GLB5_PETMA KGLTTADQLKKSADVRWHAERIINAVNDAVASM--DDTEKMSMKLRDLGSKHAKSF-
LGB2_LUPLU LK-GTSEVPQNNPELQAHAGKVFKLVEAAIQLOQVTGVVVTDATLKNLGSVHVSKG-
GLB1_GLYDI SG----AS---DPGVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKA VGRHKGYGN
Consensus  . t . . . v . . Hg kv . a a . . l d . a l . l H .
    
```

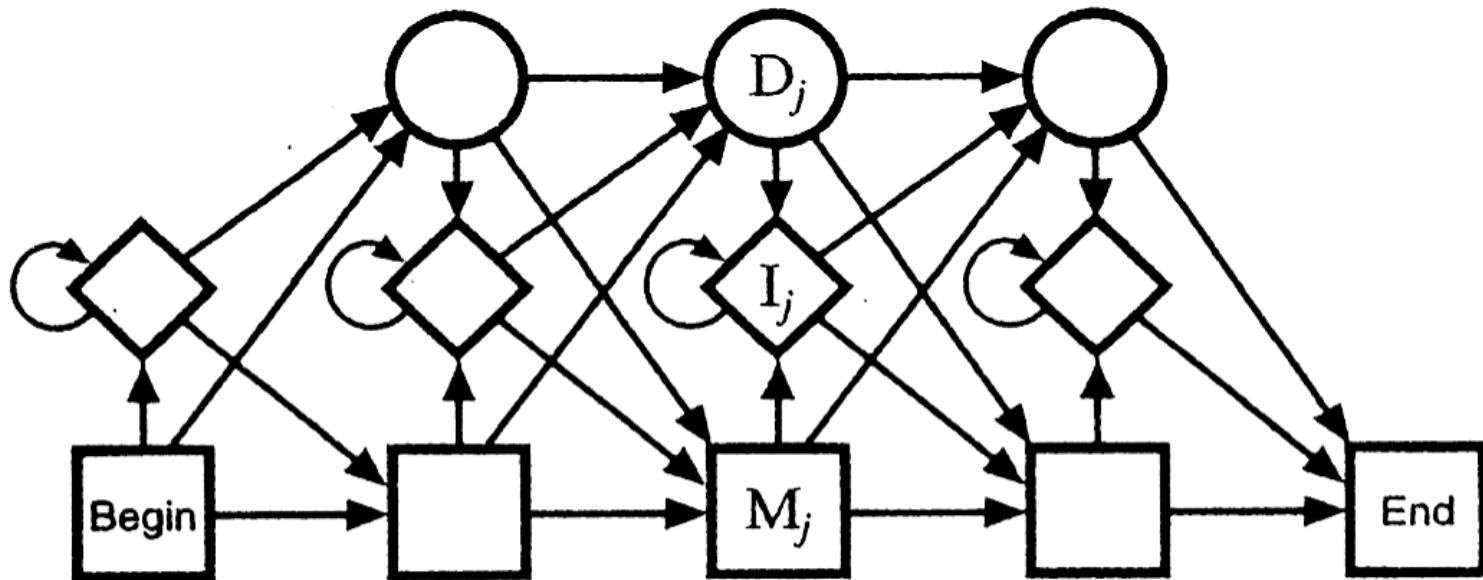
```

Helix      FGGGGGGGGGGGGGGGGGGGGGGGGGGGG      HHHHHHHHHHHHHHHHHHHHHHHHHHHHH
HBA_HUMAN  -RVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLTKFLASVSTVLTSKYR-----
HBB_HUMAN  -HVDPENFRLLGNVLCVLAHIFGKEFTPPVQAAAYQKV VAGVANALAHKYH-----
MYG_PHYCA  -KIPIKYLEFISEAIIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
GLB3_CHITP --VTHDQLNFRAGFVSYMKAHT--DFA-GAEAAWGATLDTFFGMIFSKM-----
GLB5_PETMA -QVDPQYFKVLA AVIADTVAAG-----DAGFEKLMSMICILLRSAY-----
LGB2_LUPLU --VADAHFPVVK EAILKTIKEVVGAKWSEELNSAWTIA YDELAIVIKKEMNDAA---
GLB1_GLYDI KH IKAQYFEPLGASLLSMEHRIGGKMNAAKDAWAAAYADISGALISGLQSQS-----
Consensus  v . f l . . . . . f . aa . k . . l sky
    
```

WMM might be a good model of the marked blocks, but what about the gaps??



# Profile HMM Structure



**Figure 5.2** *The transition structure of a profile HMM.*

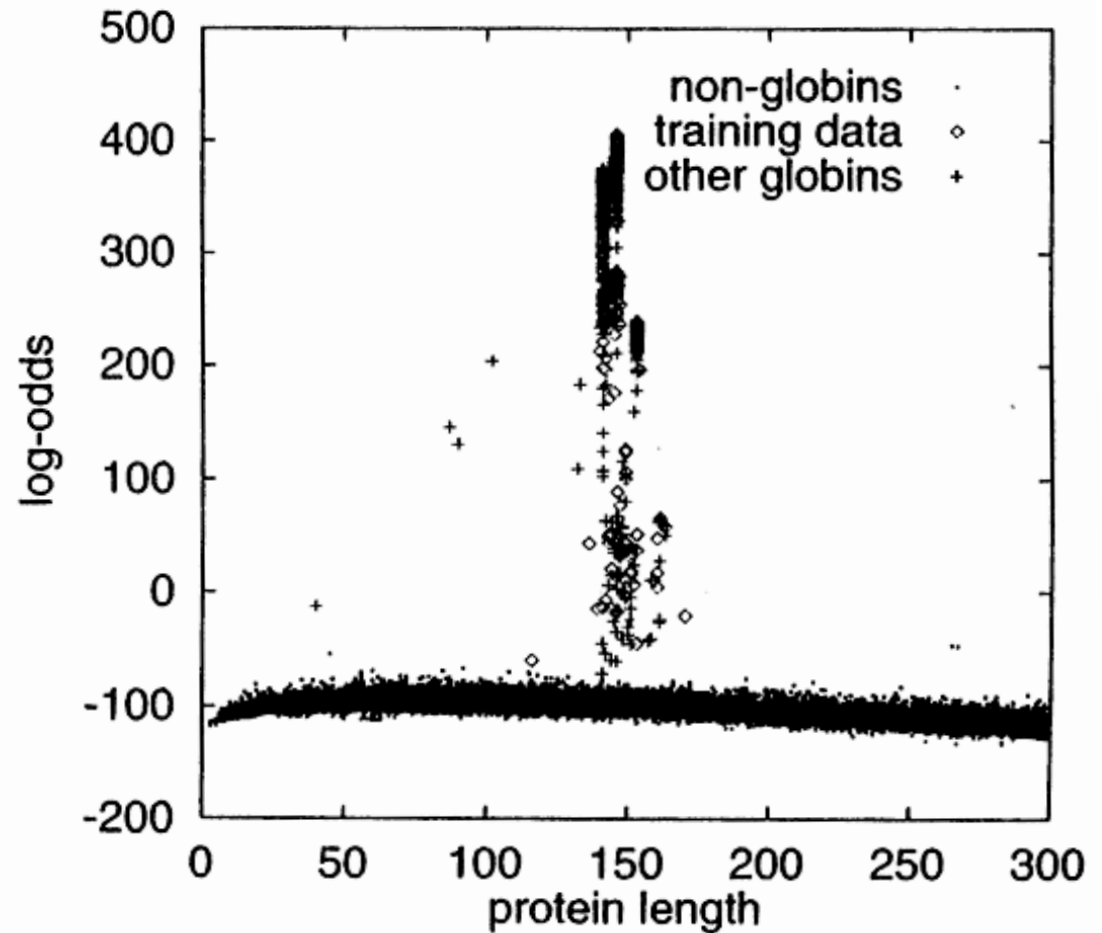
M<sub>j</sub>: Match states (20 emission probabilities)

I<sub>j</sub>: Insert states (Background emission probabilities)

D<sub>j</sub>: Delete states (silent - no emission)

# Odds Scores

Length-normalized log odds scores, globin model



# HMMs in Action: Pfam

<http://pfam.sanger.ac.uk/>

Hand-curated “seed” multiple alignments  
(domains, not full-length proteins)

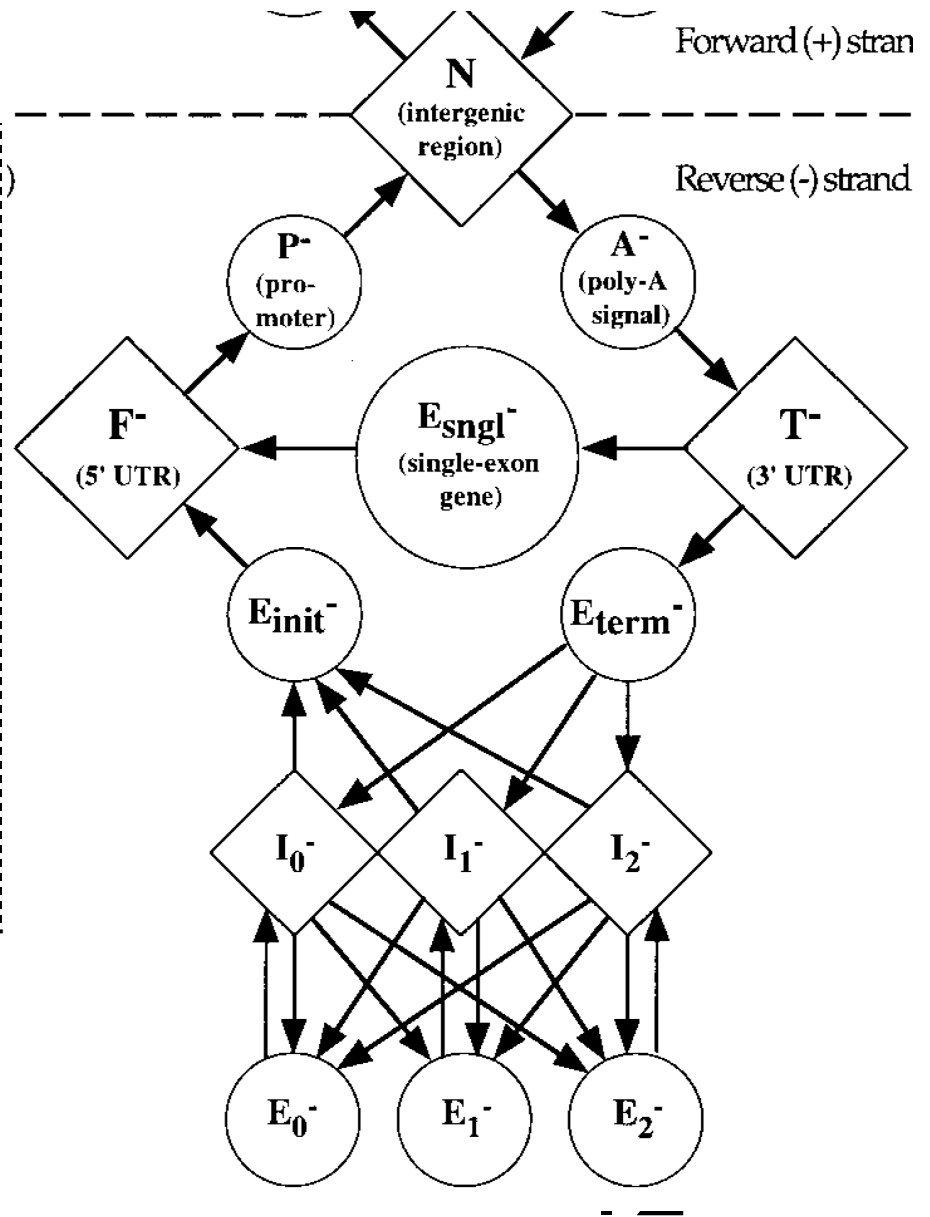
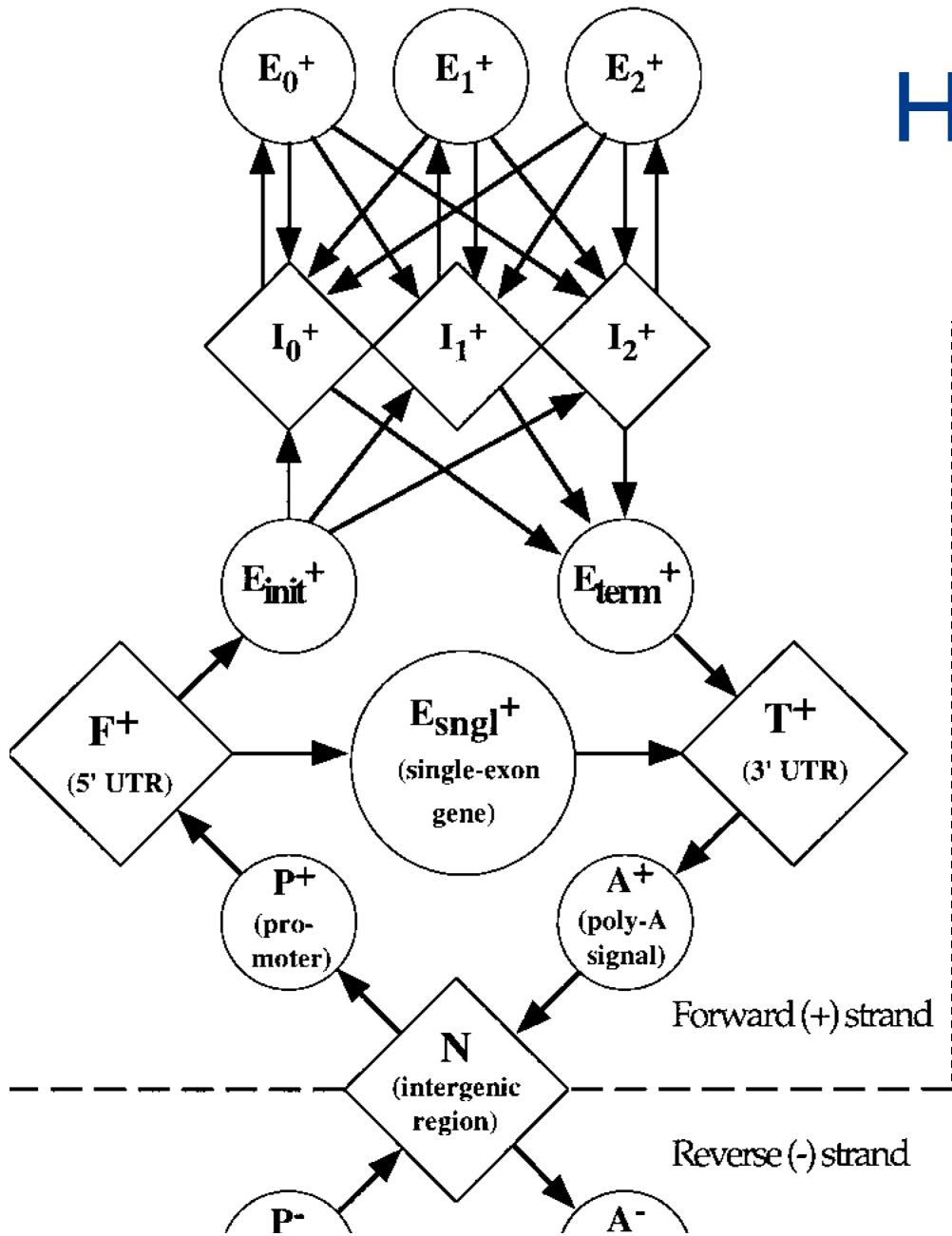
Train profile HMM from seed alignment

Hand-chosen score threshold(s)

Automatic classification/alignment of all other  
protein sequences

11912 families in Pfam 24.0, 10/2009  
(covers ~75% of proteins)

# HMM Gene Finder



# HMM Summary

## Search

Viterbi – best single path (max of products)

Forward – sum over all paths (sum of products)

Posterior decoding

## Model building

Typically fix architecture (e.g. profile HMM), then

Learn parameters – the Baum-Welch Algorithm

## Scoring

Odds ratio to background

Excellent tools available (SAM, HMMer, Pfam, ...)

*A very widely used tool for biosequence analysis*

Many variants on all of these points



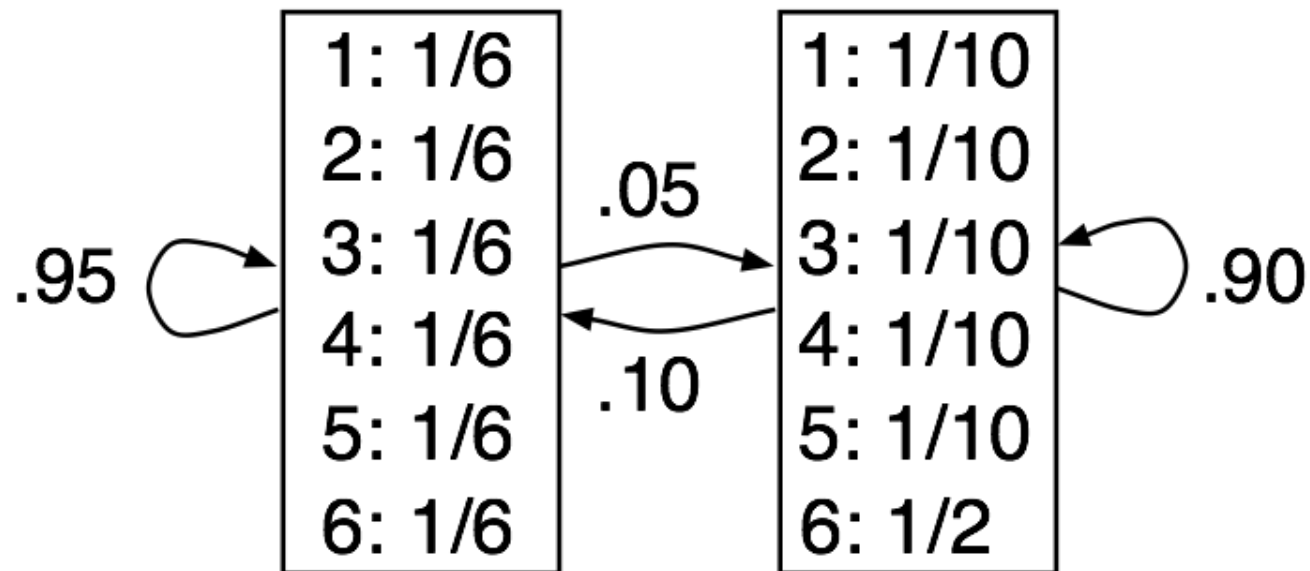
# Hidden Markov Models

(HMMs; Claude Shannon, 1948)

States:	$1, 2, 3, \dots$
Paths:	sequences of states $\pi = (\pi_1, \pi_2, \dots)$
Transitions:	$a_{k,l} = P(\pi_i = l \mid \pi_{i-1} = k)$
Emissions:	$e_k(b) = P(x_i = b \mid \pi_i = k)$
Observed data:	emission sequence
Hidden data:	state/transition sequence

# The Occasionally Dishonest Casino

1 fair die, 1 “loaded” die, occasionally swapped





Rolls	315116246446644245311321631164152133625144543631656626566666
Die	FFL
Viterbi	FFL
Rolls	651166453132651245636664631636663162326455236266666625151631
Die	LLLLLLFFL
Viterbi	LLLLLLFFL
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFLL
Viterbi	FFL
Rolls	366163666466232534413661661163252562462255265252266435353336
Die	LLLLLLLLLFF
Viterbi	LLLLLLLLLLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Rolls	233121625364414432335163243633665562466662632666612355245242
Die	FFL
Viterbi	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL

**Figure 3.5**

*Rolls: Visible data—300 rolls of a die as described above.*

*Die: Hidden data—which die was actually used for that roll (F = fair, L = loaded).*

*Viterbi: the prediction by the Viterbi algorithm is shown.*

# Inferring hidden stuff

Joint probability of a given path  $\pi$  & emission sequence  $x$ :

$$P(x, \pi) = a_{0, \pi_1} \prod_{i=1}^n e_{\pi_i}(x_i) \cdot a_{\pi_i, \pi_{i+1}}$$

*But  $\pi$  is hidden*; what to do? Some alternatives:

Most probable single path

$$\pi^* = \arg \max_{\pi} P(x, \pi)$$

Sequence of most probable states

$$\hat{\pi}_i = \arg \max_k P(\pi_i = k \mid x)$$

# The Viterbi Algorithm: The most probable path

Viterbi finds:  $\pi^* = \arg \max_{\pi} P(x, \pi)$

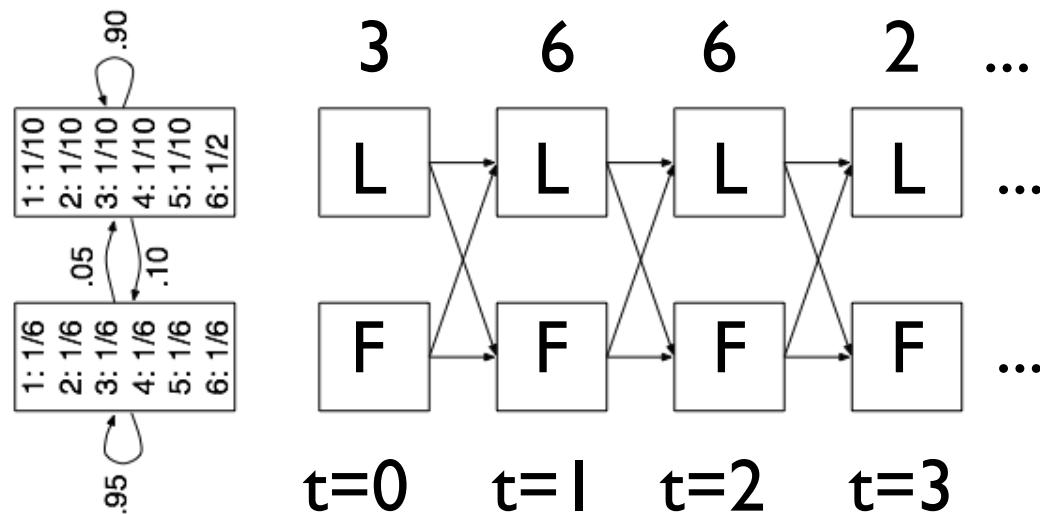
Possibly there are  $10^{99}$  paths of prob  $10^{-99}$

More commonly, one path (+ slight variants) dominate others.

(If not, other approaches may be preferable.)

Key problem: exponentially many paths  $\pi$

# Unrolling an HMM



Conceptually, sometimes convenient

Note exponentially many paths

Rolls	315116246446644245311321631164152133625144543631656626566666
Die	FFL
Viterbi	FFL
Rolls	651166453132651245636664631636663162326455236266666625151631
Die	LLLLLLFFL
Viterbi	LLLLLLFFL
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFLLL
Viterbi	FFL
Rolls	366163666466232534413661661163252562462255265252266435353336
Die	LLLLLLLLLFF
Viterbi	LLLLLLLLLLLLLLLLLFF
Rolls	233121625364414432335163243633665562466662632666612355245242
Die	FFL
Viterbi	FFL

### Figure 3.5

*Rolls: Visible data—300 rolls of a die as described above.*

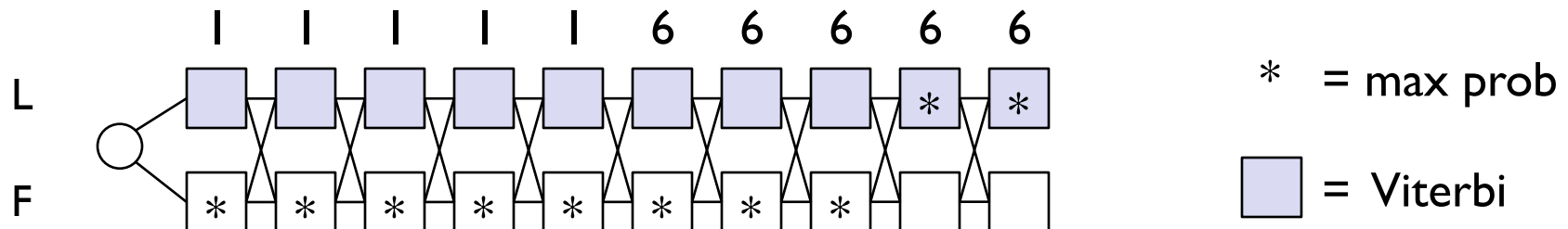
*Die: Hidden data—which die was actually used for that roll (F = fair, L = loaded).*

*Viterbi: the prediction by the Viterbi algorithm is shown.*

# Most probable path $\neq$ Sequence of most probable states

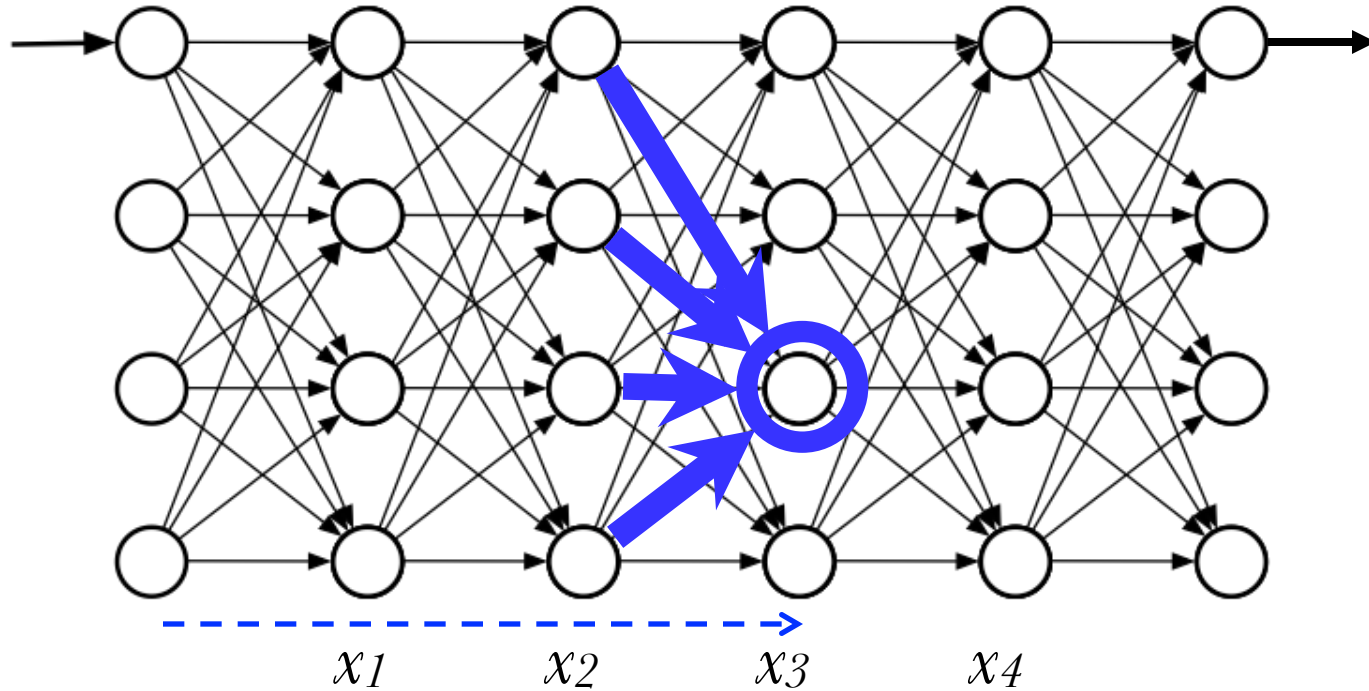
Another example, based on casino dice again

Suppose  $p(\text{fair} \leftrightarrow \text{loaded})$  transitions are  $10^{-99}$  and roll sequence is  $11111\dots 66666$ ; then fair state is more likely all through 1's & well into the run of 6's, but eventually loaded wins, and the improbable  $F \rightarrow L$  transitions make Viterbi = *all L*.



# The Forward Algorithm

For each state/time, want *total* probability of all paths leading to it, with given emissions



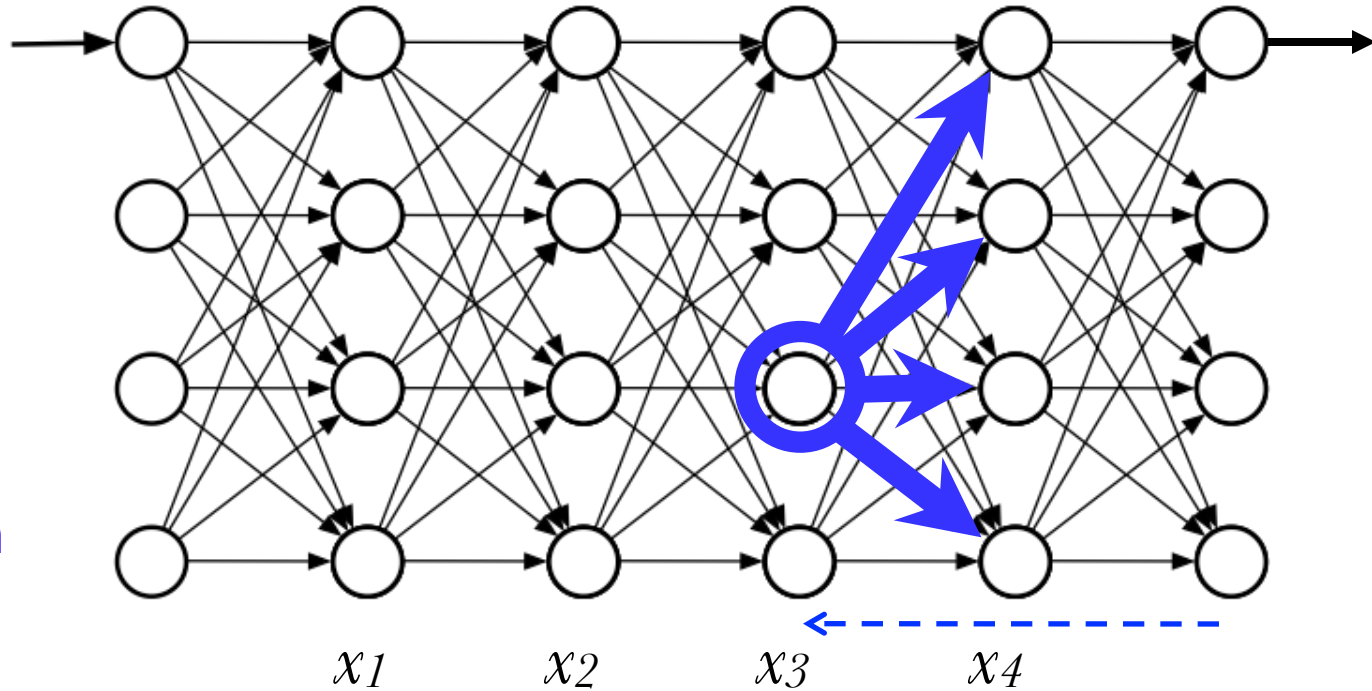
$$f_k(i) \triangleq P(x_1 \dots x_i, \pi_i = k)$$

$$f_l(i+1) = e_l(x_{i+1}) \sum_k f_k(i) a_{k,l}$$

$$P(x) = \sum_{\pi} P(x, \pi) = \sum_k f_k(n) a_{k,0}$$

# The Backward Algorithm

Similar:  
for each  
state/time,  
want total  
probability  
of all paths  
from it, with  
given  
emissions,  
conditional  
on that  
state.



$$b_k(i) \triangleq P(x_{i+1} \cdots x_n \mid \pi_i = k)$$

$$b_k(i) = \sum_l a_{k,l} e_l(x_{i+1}) b_l(i+1)$$

$$b_k(n) = a_{k,0}$$



# In state $k$ at step $i$ ?

$$P(x, \pi_i = k)$$

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

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

$$= f_k(i) \cdot b_k(i)$$

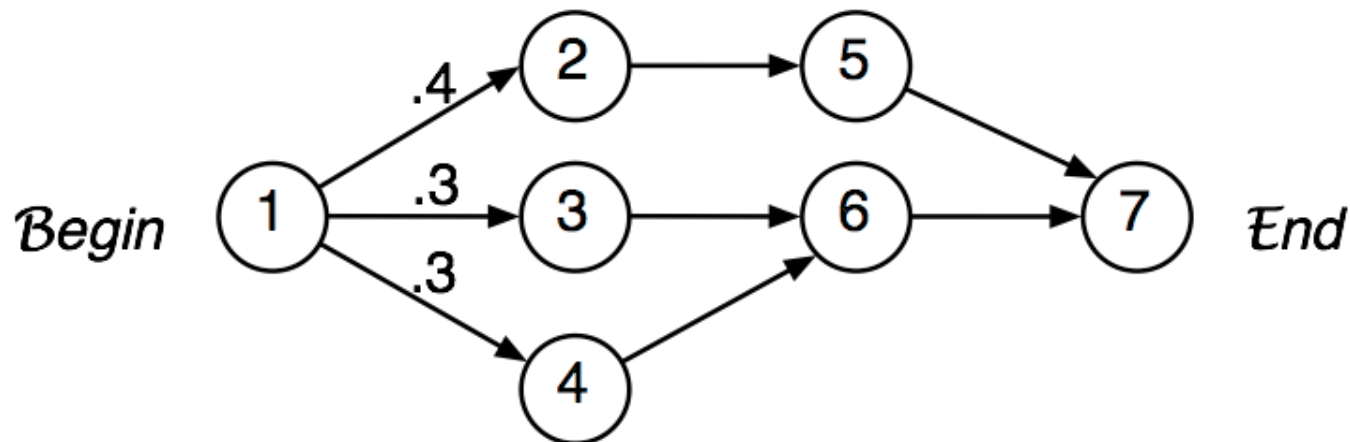
$$P(\pi_i = k \mid x) = \frac{P(x, \pi_i = k)}{P(x)} = \frac{f_k(i) \cdot b_k(i)}{P(x)}$$

# Posterior Decoding, I

Alternative 1: what's the most likely state at step  $i$ ?

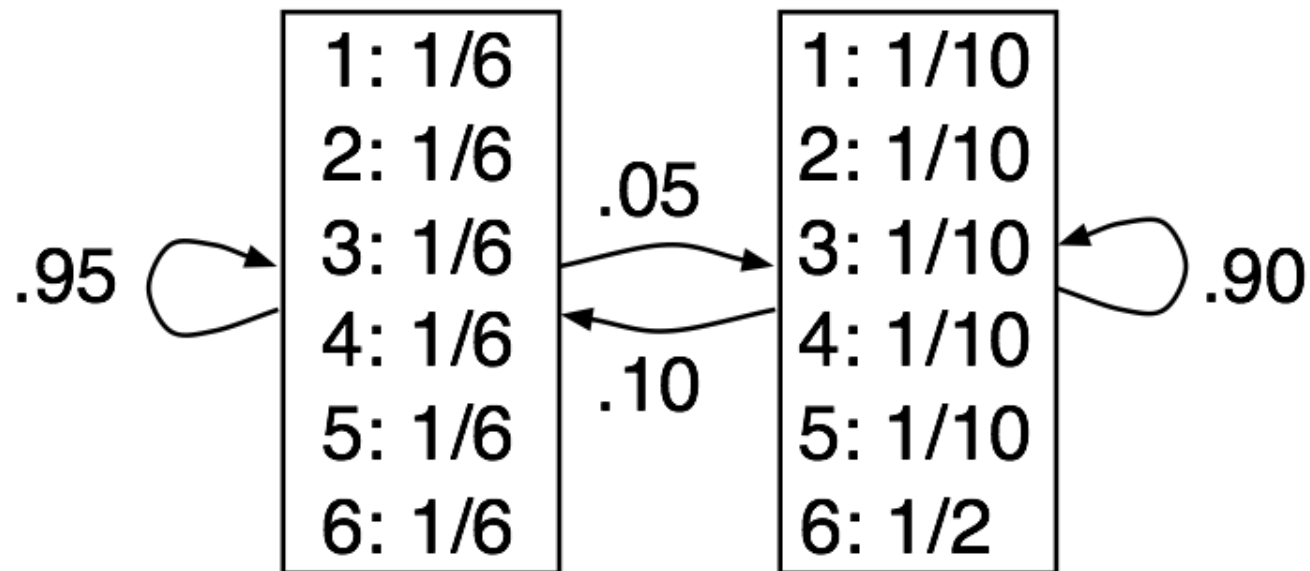
$$\hat{\pi}_i = \arg \max_k P(\pi_i = k \mid x)$$

Note: the sequence of most likely states  $\neq$  the most likely sequence of states. May not even be legal!



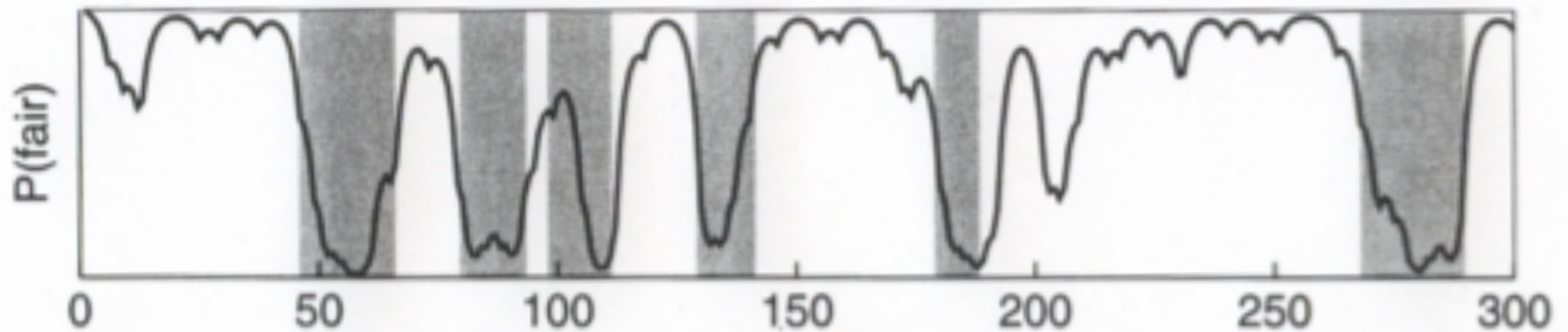
# The Occasionally Dishonest Casino

1 fair die, 1 “loaded” die, occasionally swapped





# Posterior Decoding



**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, II

Alternative 1: what's most likely state at step  $i$ ?

$$\hat{\pi}_i = \arg \max_k P(\pi_i = k \mid x)$$

Alternative 2: given some function  $g(k)$  on states, what's its expectation. E.g., what's probability of “+” model in CpG HMM ( $g(k)=1$  iff  $k$  is “+” state)?

$$G(i \mid x) = \sum_k P(\pi_i = k \mid x) \cdot g(k)$$

# CpG Islands again

Data: 41 human sequences, totaling 60kbp,  
including 48 CpG islands of about 1kbp each

Viterbi:

Found 46 of 48  
plus 121 “false positives”

Post-process:

46/48  
67 false pos

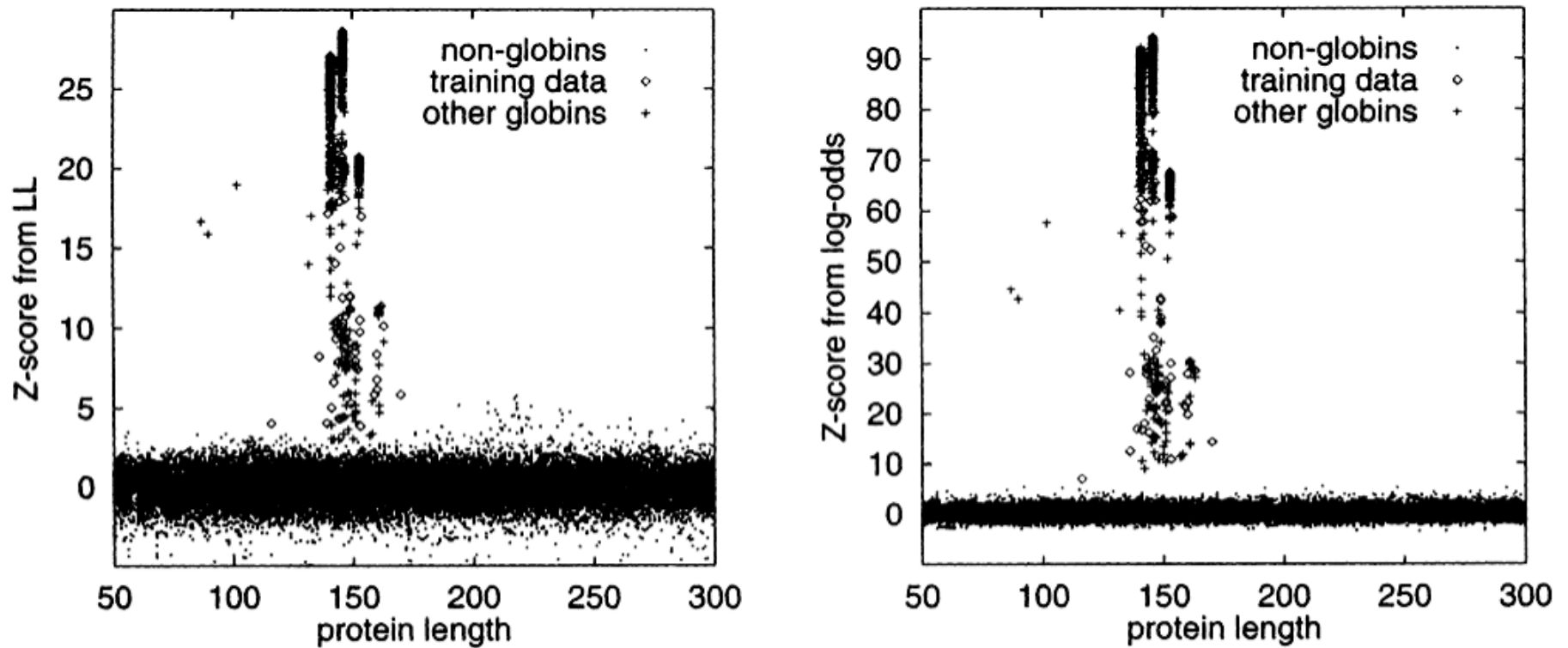
Posterior Decoding:

same 2 false negatives  
plus 236 false positives

46/48  
83 false pos

Post-process: merge within  
500; discard < 500

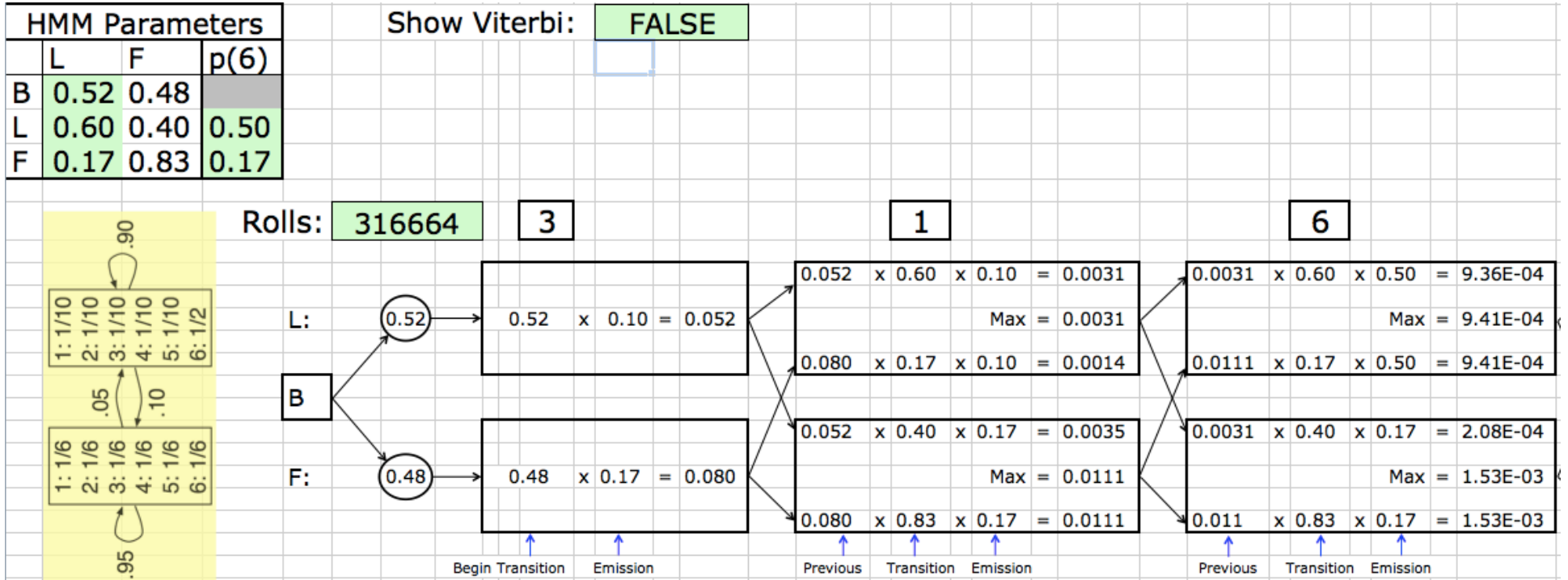
# Z-Scores



**Figure 5.6** *The Z-score calculated from the LL scores (left) and the log-odds (right).*

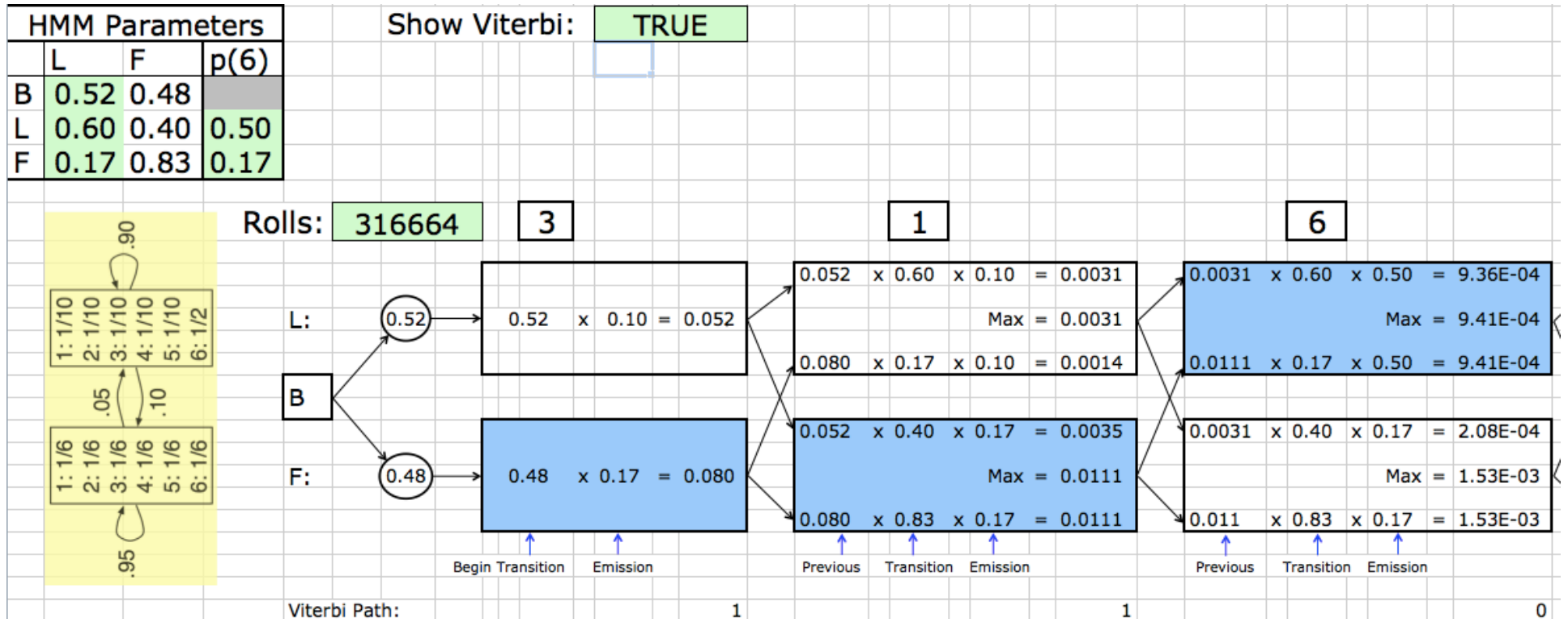


# HMM Casino Example



(Excel spreadsheet on web; download & play...)

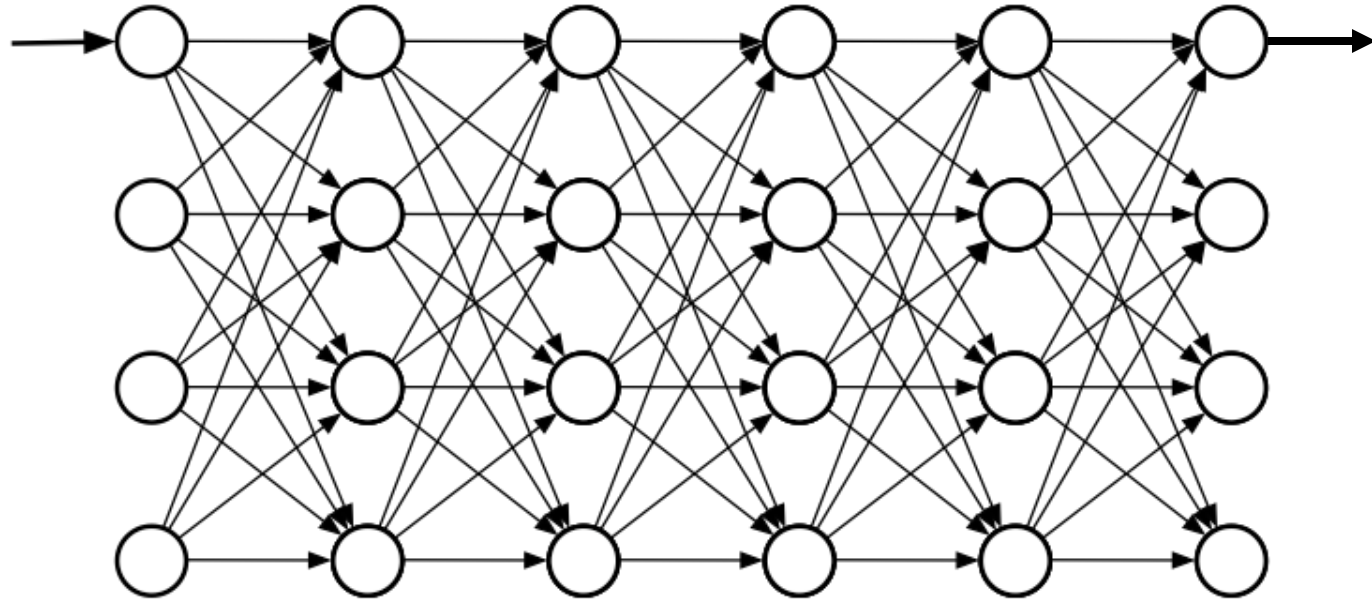
# HMM Casino Example



(Excel spreadsheet on web; download & play...)

# An HMM (unrolled)

States



$x_1$

$x_2$

$x_3$

$x_4$

Emissions/sequence positions  $\longrightarrow$

# HMMs in Action: Pfam

<http://pfam.sanger.ac.uk/>

Proteins fall into families, both across & within species

Ex: Globins, GPCRs, Zinc fingers, Leucine zippers,...

Identifying family very useful: suggests function, etc.

So, search & alignment are both important

One very successful approach: profile HMMs