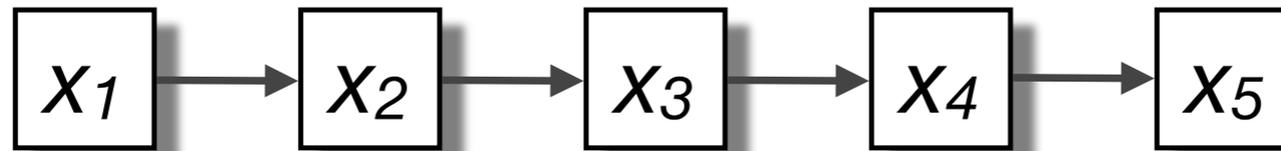


These shortcomings of PSSMs set the stage for a new kind of profile, based on **Markov chains**, called **Hidden Markov models (HMMs)**

- ▶ modeling positional dependencies
- ▶ recognizing pattern instances with indels
- ▶ modeling variable length patterns
- ▶ detecting boundaries

Markov chains

Markov chains are stochastic processes that undergo **transitions** between a finite series of **states** in a chainlike manner.



The system transverses states with probability

$$p(x_1, x_2, x_3, \dots) = p(x_1) p(x_2 | x_1) p(x_3 | x_2) p(x_4 | x_3) \dots$$

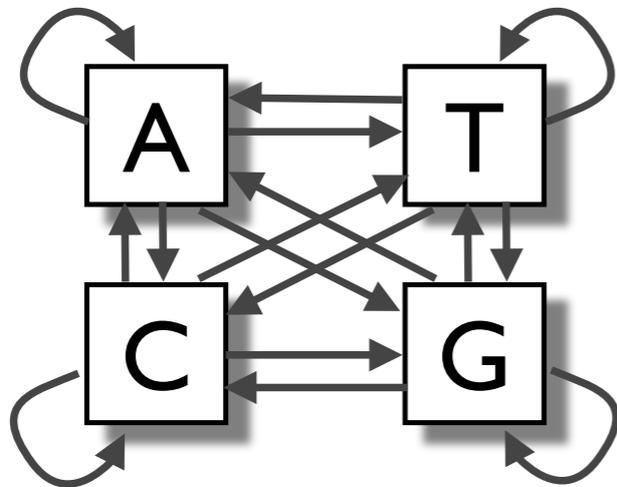
i.e. **Markov chains are memoryless**: the probability that the chain is in state x_i at time t , depends only on the state at the previous time step and not on the past history of the states visited before time $t-1$.

This specific kind of "*memorylessness*" is called the **Markov property**.

The **Markov property** states that the conditional probability distribution for the system at the next step (and in fact at all future steps) depends only on the current state of the system, and not additionally on the state of the system at previous steps.

Markov chains...

Markov chains, and their extension hidden Markov models (HMMs), are commonly represented by **state diagrams**, which consist of *states* and connecting *transitions*



E.g., A general Markov chain modeling DNA. Note that any sequence can be traced through the model by passing from one state to the next via the transitions.

A **transition probability** parameter (a_{ij}) is associated with each transition (arrow) and determines the probability of a certain state (S_j) following another state (S_i).

A Markov chain is defined by:

- a finite set of **states**, $S_1, S_2 \dots S_N$
- a set of **transition probabilities**: $a_{ij} = P(q_{t+1}=S_j|q_t=S_i)$
- and an **initial state probability distribution**, $\pi_i = P(q_0=S_i)$

Simple Markov chain example for $x=\{a,b\}$

Observed sequence: $x = \mathbf{abaaababbaa}$

Model:

transition probabilities

Prev i	Next j	Prob a_{ij}
a	a	0.7
a	b	0.3
b	a	0.5
b	b	0.5

initial state probability distribution

Start probs	π_i	a 0.5
		b 0.5

$$P(\mathbf{x}) = \mathbf{0.5 \times 0.3 \times 0.5 \times 0.7 \times 0.7 \times 0.3 \times 0.5 \times 0.3 \times 0.5 \times 0.5 \times 0.7}$$

Q. Can you sketch the state diagram with labeled transitions for this model?

Typical questions we can ask with Markov chains include:

- What is the probability of being in a particular state at a particular time?
(By time here we can read position in our query sequence)
- What is the probability of seeing a particular sequence of states?
(I.e., the score for a particular query sequence given the model)

Q. What do Markov chains add over the traditional PSSM approach?

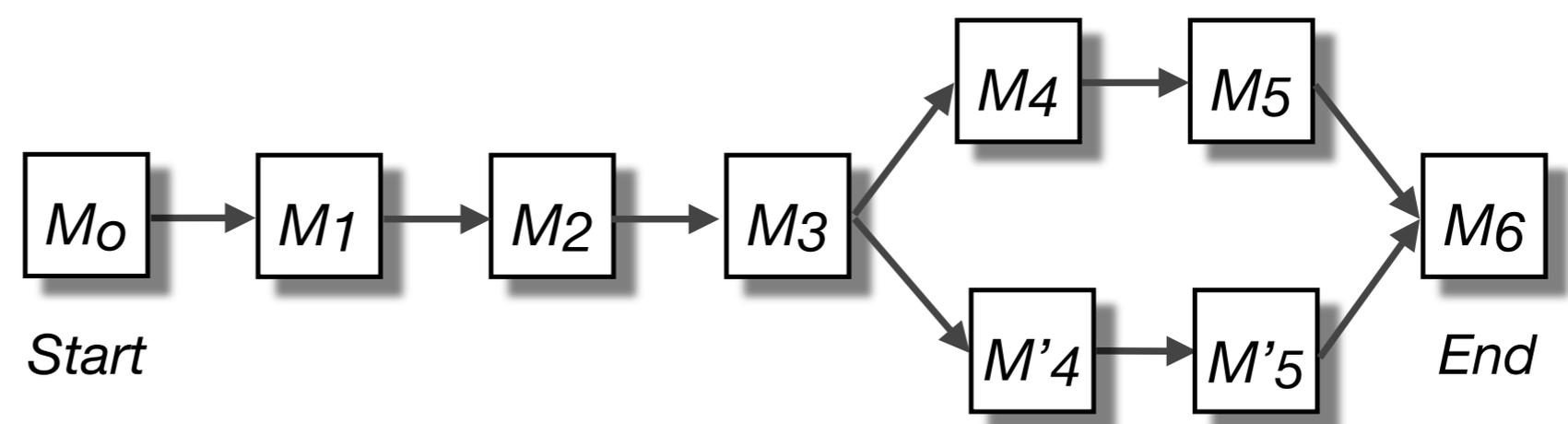
In particular how do Markov chains deal with the following PSSM weaknesses?

1. Positional dependencies
2. Pattern instances containing insertions or deletions
3. Variable length patterns, and
4. The detection boundaries (i.e. segmentation of sequences)

Markov chains: 1. Positional dependencies ✓

The connectivity or **topology** of a Markov chain can easily be designed to capture dependencies and variable length motifs.

WEIRD
WEIRD
WEIQH
WEIRD
WEIQH

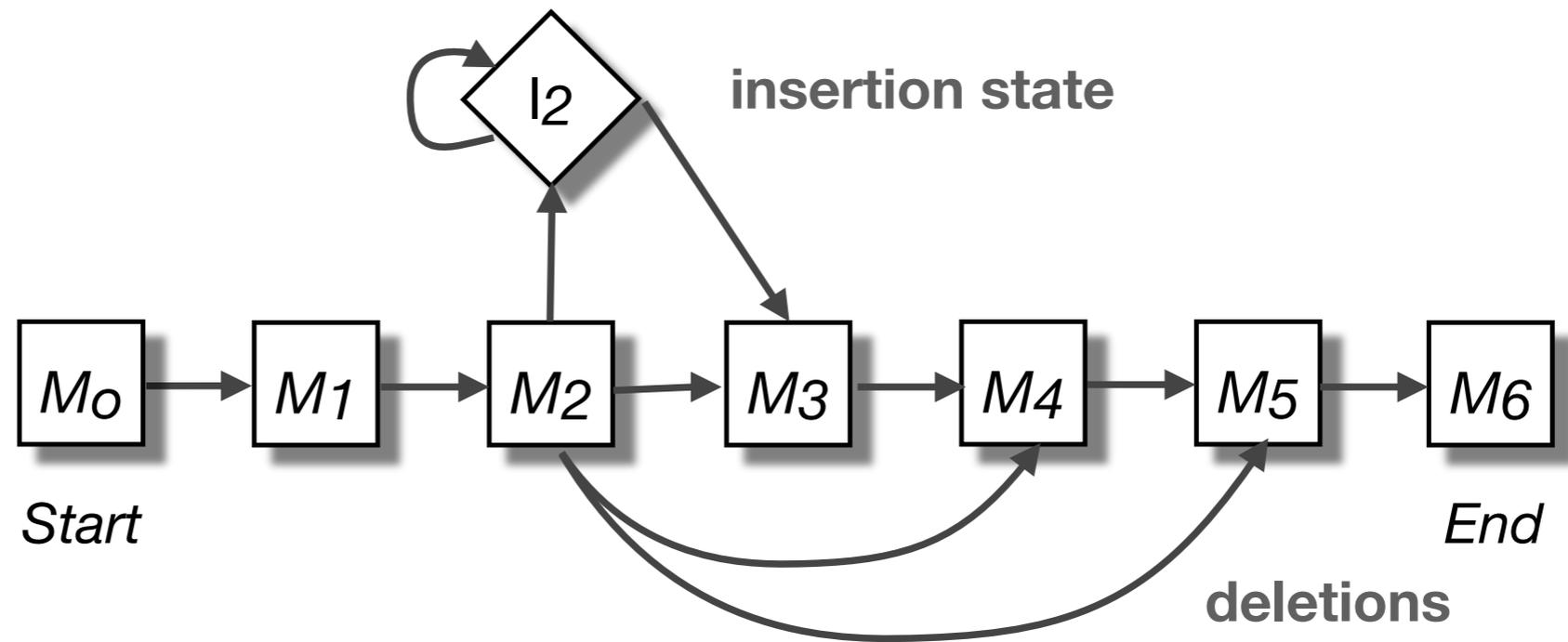


Recall that a PSSM for this motif would give the sequences **WEIRD** and **WEIRH** equally good scores even though the **RH** and **QR** combinations were not observed

Markov chains: 2. Insertions and deletions ✓

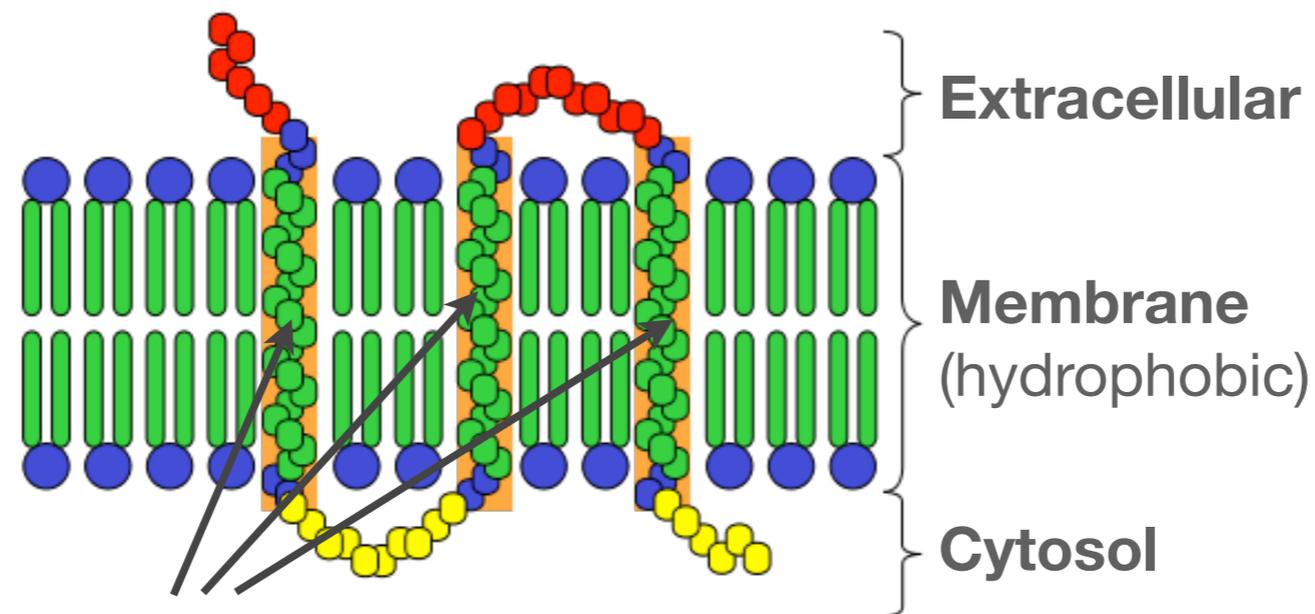
To address pattern instances with gaps and variable length motifs, we can construct a Markov chain to recognize a query sequences with insertions (via an extra insertion state) and deletions (via extra transitions (edges))

WETIRD
WE-IRD
WETIQH
WE-IRD
WE--QH



Markov chains: 3. Boundary detection ?

Giving a sequence we wish to label each symbol in the sequence according to its class (e.g. transmembrane regions or extracellular/cytosolic)



tend to be hydrophobic in composition

Given a training set of labeled sequences we can begin by modeling each amino acid as hydrophobic (**H**) or hydrophilic (**L**)

i.e. reduce the dimensionality of the 20 amino acids into two classes

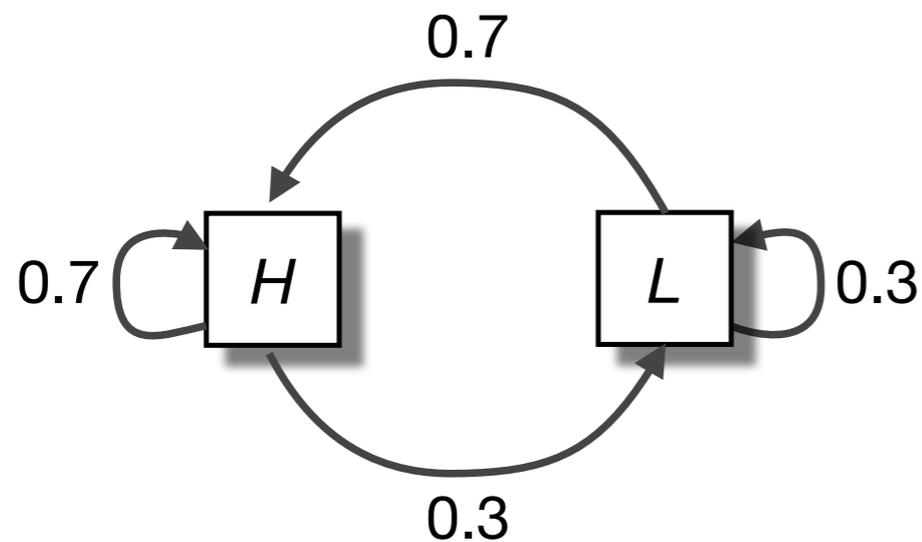
E.g., A peptide sequence can be represented as a sequence of Hs and Ls.

e.g. HHHLLHLHHLHL...

Markov chains: boundary detection...

A simpler question: **is a given sequence a transmembrane sequence?**

A Markov chain for recognizing transmembrane sequences



- States: S_H, S_L
- $\Sigma = \{H, L\}$
- $\pi(H) = 0.6, \pi(L) = 0.4$

Question: Is sequence **HHLHH** a transmembrane protein?

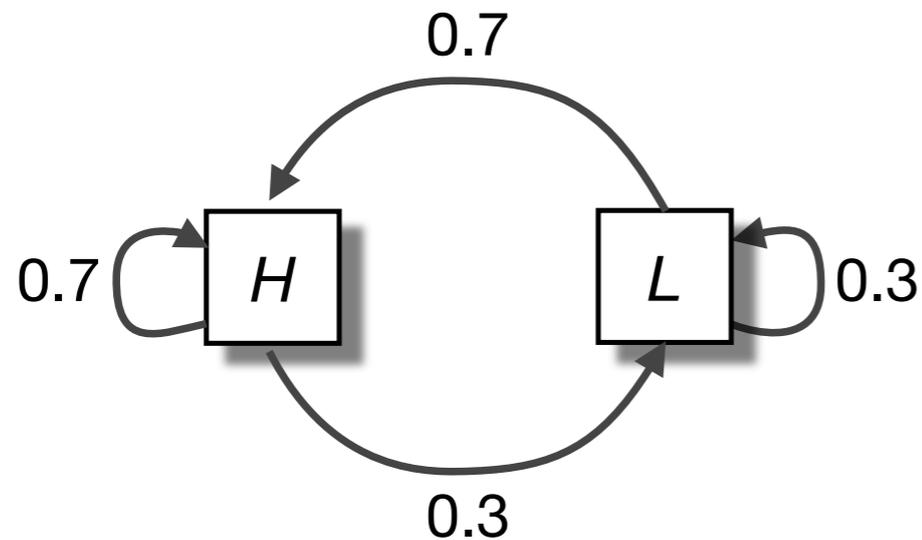
$$P(\text{HHLHH}) = 0.6 \times 0.7 \times 0.7 \times 0.3 \times 0.7 \times 0.7 = 0.043$$

Problem: need a threshold,
threshold must be length dependent

Markov chains: boundary detection

We can classify an observed sequence ($O = O_1, O_2, \dots$) by its log odds ratio

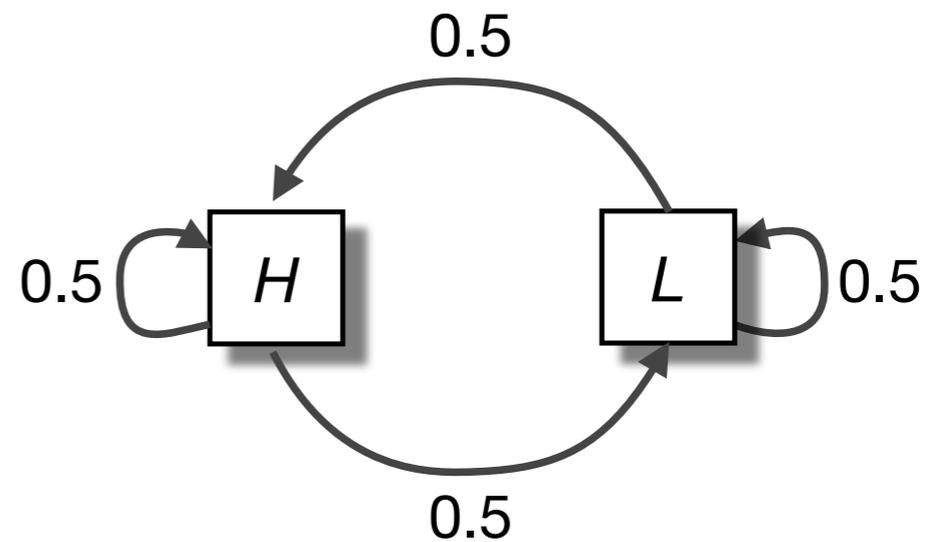
transmembrane model



Transmembrane (TM)

- $\pi(H) = 0.6, \pi(L) = 0.4$

null model



Extracellular/cytosolic (E/C)

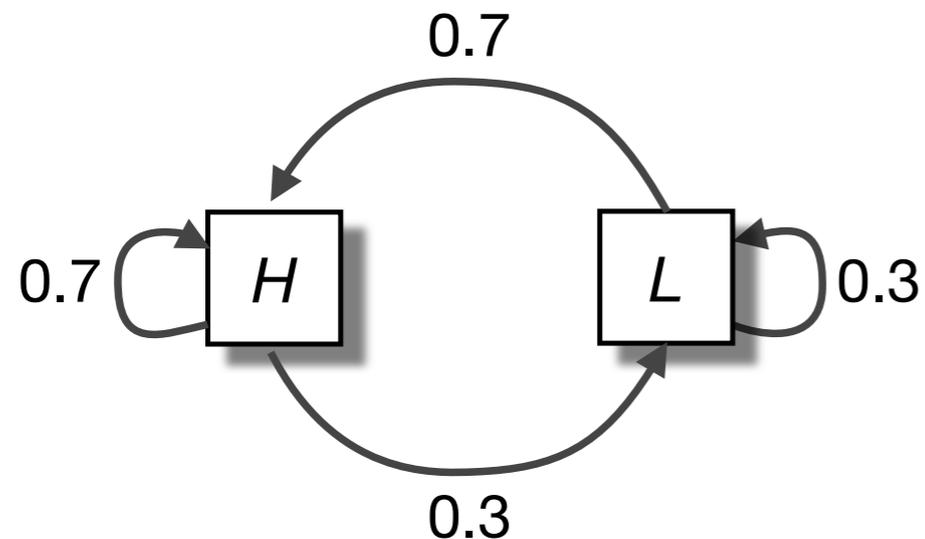
- $\pi(H) = 0.5, \pi(L) = 0.5$

$$\frac{P(\mathbf{HHLHH} \mid \text{TM})}{P(\mathbf{HHLHH} \mid \text{EC})} = \frac{0.6 \times 0.7 \times 0.7 \times 0.3 \times 0.7 \times 0.7}{0.5 \times 0.5 \times 0.5 \times 0.5 \times 0.5 \times 0.5} = \frac{0.043}{0.016} = 2.69$$

In other words, it is more than twice as likely that **HHLHH** is a transmembrane sequence. The log-odds score is: $\log_2(2.69) = 1.43$

Side note: Parameter estimation

Both initial probabilities ($\pi(i)$) and transition probabilities (a_{ij}) are determined from known examples of transmembrane and non-transmembrane sequences.



- initial probabilities $\pi(H)$, $\pi(L)$
- transition probabilities: a_{HH} , a_{HL} , a_{LH} and a_{LL} .

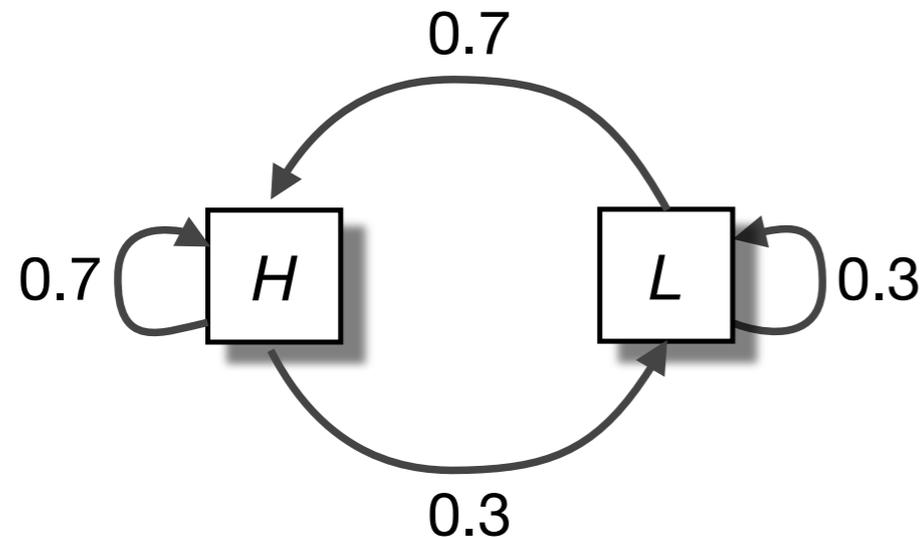
Given labeled sequences (TM and E/C), we determine the initial probabilities $\pi(i)$ by counting the number of sequences that begin with residue i .

To determine transition probabilities, a_{ij} , we first determine A_{ij} (the number of transitions from state i to j in the training data, i.e. count the number of ij pairs in the training data). Then normalize by the number of i^* pairs.

$$a_{ij} = \frac{A_{ij}}{\sum_{j'} A_{ij'}}$$

Side note: Parameter estimation...

Both initial probabilities ($\pi(i)$) and transition probabilities (a_{ij}) are determined from known examples of transmembrane and non-transmembrane sequences.



$\pi(H)$ = # of sequences that begin with H, normalized by the total # of training sequences

- $\pi(H) = 0.6$, $\pi(L) = 0.4$

HHLLHHLLLLHLHLLHLLLLHLHHHL

HHHLHHLHLLLLLLLHHHHLLLLHHHHHL

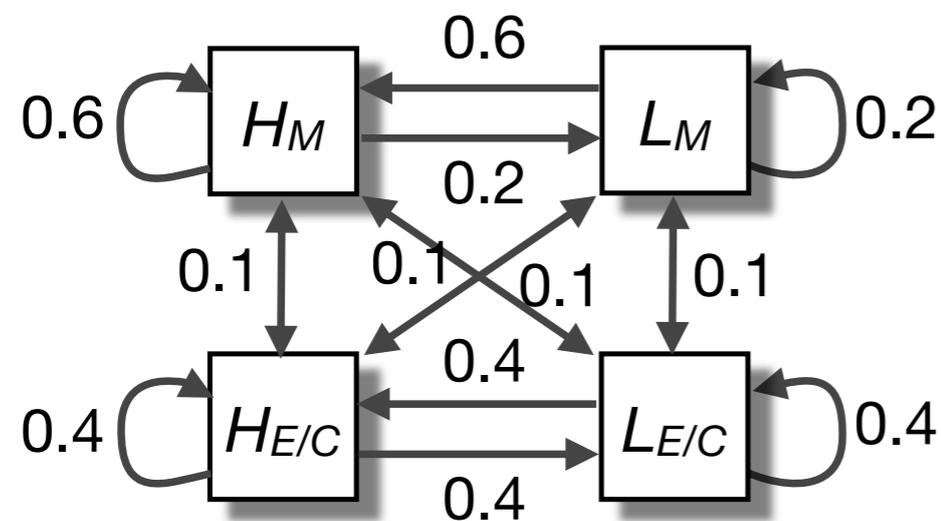
HH . . . ($A_{HL} = 12$, $A_{H^*} = 40$)

$$a_{HL} = \frac{A_{HL}}{\sum_i A_{Hi}} = \frac{\#HL \text{ pairs}}{\# H^* \text{ pairs}} = \frac{12}{40}$$

Boundary detection challenge

Given sequence of Hs and Ls, find all transmembrane regions:

To approach this question we can construct a new four state model by adding transitions connecting the TM and E/C models



Transitions between the *M* states and the *E/C* states indicate boundaries between membrane regions and cytosolic or extracellular regions.

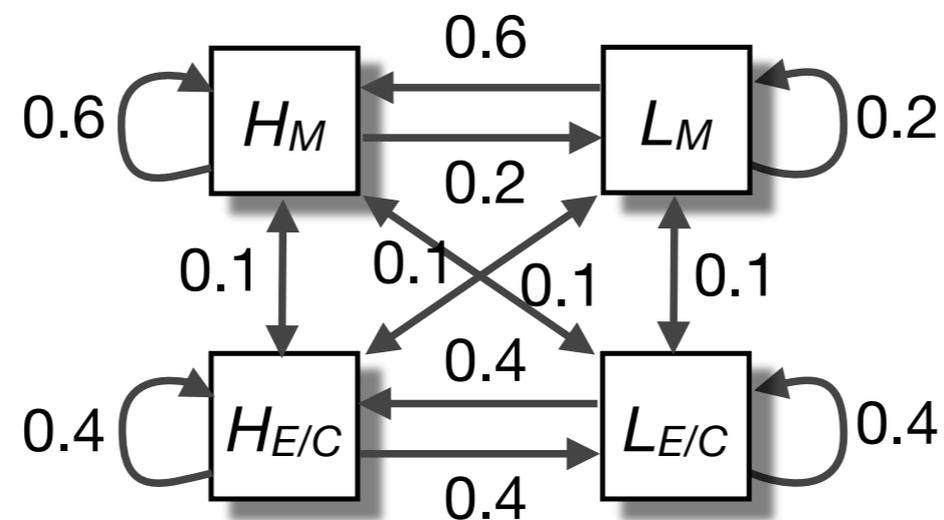
However this is no longer a standard Markov chain!

Boundary detection challenge...

In a Markov chain, there is a one-to-one correspondence between symbols and states, which is not true of our new merged four state, two symbol model.

For example, both H_M and $H_{E/C}$ are associated with hydrophilic residues.

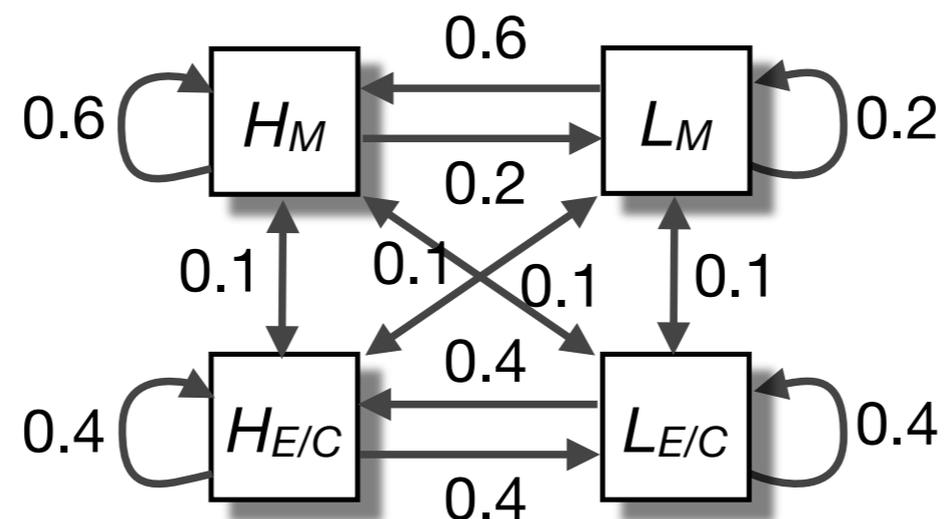
- This four-state transmembrane model is a **hidden Markov model**.



So whats hidden?

We will distinguish between the *observed* parts of the problem and the *hidden* parts

- In the Markov models we have considered previously it is clear which states account for each part of the observed sequence
Due to the one-to-one correspondence between symbols and states
- In our new model, there are multiple states that could account for each part of the observed sequence
i.e. we don't know which state emitted a given symbol from knowledge of the sequence and the structure of the model
 - ▶ This is the *hidden* part of the problem

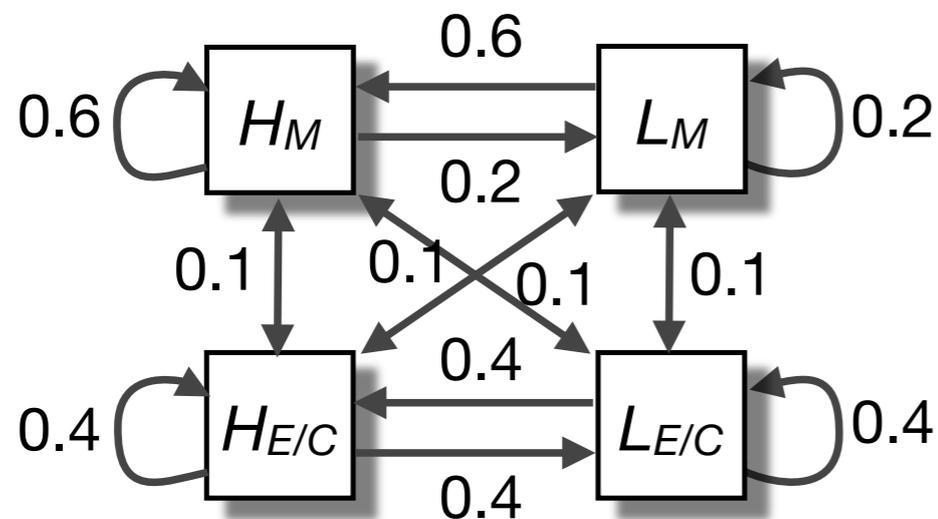


For our Markov models

- Given HLLH..., we know the exact state sequence ($q_0=S_H, q_1=S_L, q_2=S_L, \dots$)

For our HMM

- Given HLLH..., we must infer the most probable state sequence
- This HMM state sequence will yield the boundaries between likely TM and E/C regions



HM, LM, LM, HM
 HM, LM, LM, HE/C
 HM, LM, LH/C, HM
 HM, LM, LH/C, HE/C
 HM, LE/C, LM, HM
 HM, LE/C, LM, HE/C
 HM, LE/C, LH/C, HM,
 HM, LE/C, LH/C, HE/C,
 HE/C, LM, LM, HM
 HE/C, LM, LM, HE/C
 HE/C, LM, LH/C, HM
 HE/C, LM, LH/C, HE/C
 HE/C, LE/C, LM, HM
 HE/C, LE/C, LM, HE/C
 HE/C, LE/C, LH/CM, HM
 HE/C, LE/C, LH/CM, HE/C

Side note: HMM states as sequence emitters

It's useful to imagine HMM states **emitting symbols** each time they are visited

In this way, transversing the model will “generate” a sequence with a certain probability (i.e. “score”).

This probability is a product of the state path taken through the model

That is, it depends on *initial probabilities*, *transition probabilities* and ***emission probabilities*** (the probability that a visited state emits a particular symbol) along the path

There may be many possible paths that can generate the same sequence

An HMM is a **full probabilistic model** – the model parameters θ and the overall sequence “scores” $P(x, S | HMM, \theta)$ are all probabilities. As a result, we can use standard **Bayesian probability theory** to manipulate these numbers in powerful ways, including optimizing parameters, calculating confidence in predictions, and interpreting the statistical significance of scores.

Hidden Markov models (HMMs)

Markov Chains

- States: $S_1, S_2 \dots S_N$
- Initial probabilities: π_i
- Transition probabilities: a_{ij}

One-to-one correspondence
between states and symbols

Hidden Markov Models

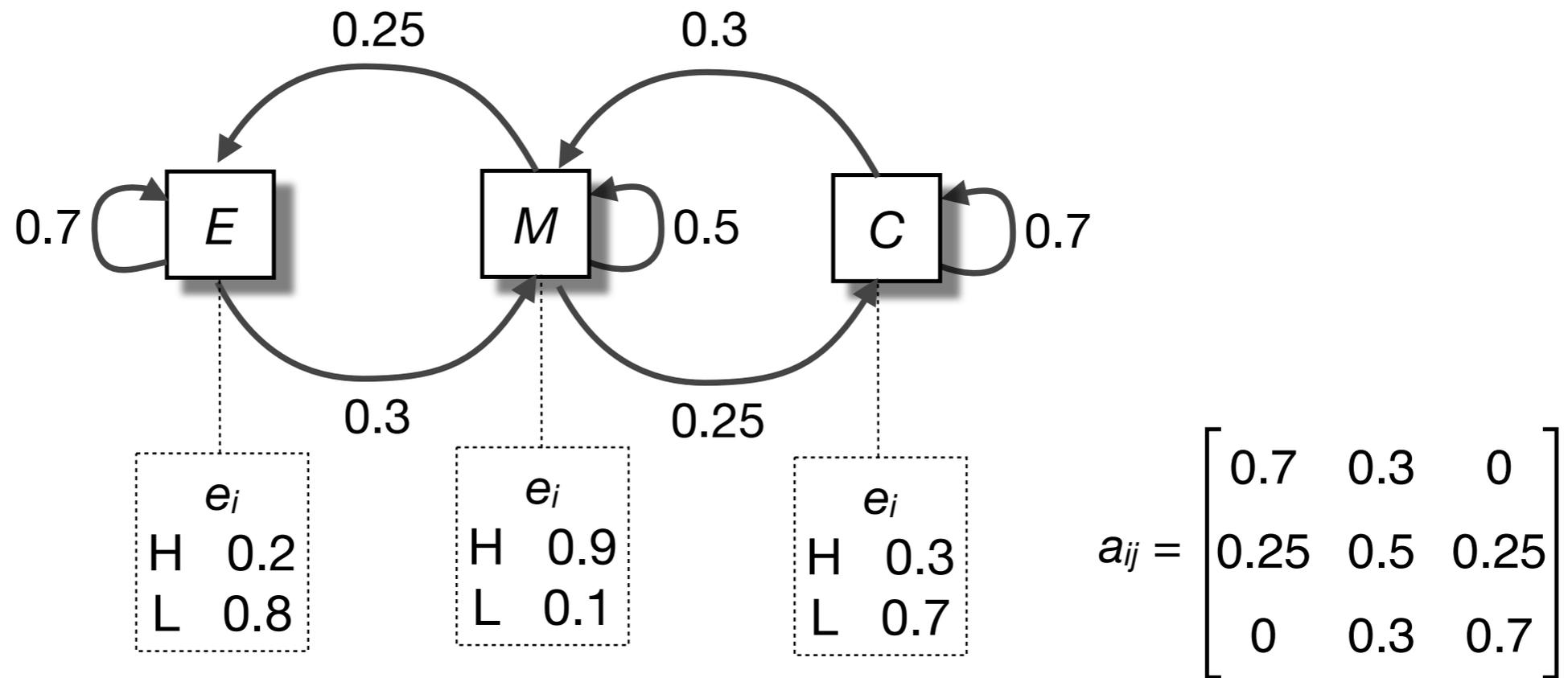
- States: $S_1, S_2 \dots S_N$
- Initial probabilities: π_i
- Transition probabilities: a_{ij}
- **Alphabet** of emitted symbols, Σ
- **Emission probabilities:** $e_i(a)$
probability state i emits symbol a

Symbol may be emitted by more
than one state

Similarly, a state can emit more
than one symbol

Example three state HMM

In this example we will use only one state for the transmembrane segment (M) and use emission probabilities to distinguish between H and L residues. We will also add separate E & C states with distinct emission probabilities.



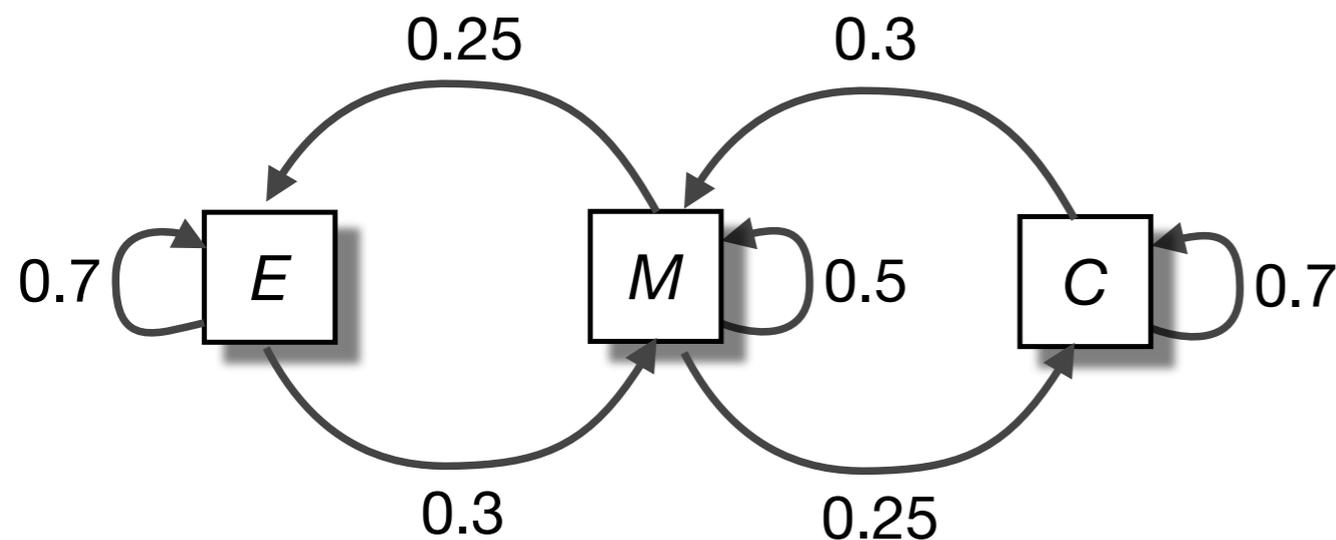
Side note: Parameter estimation

As in the case of Markov chains, the HMM parameters can be learned from labeled training data

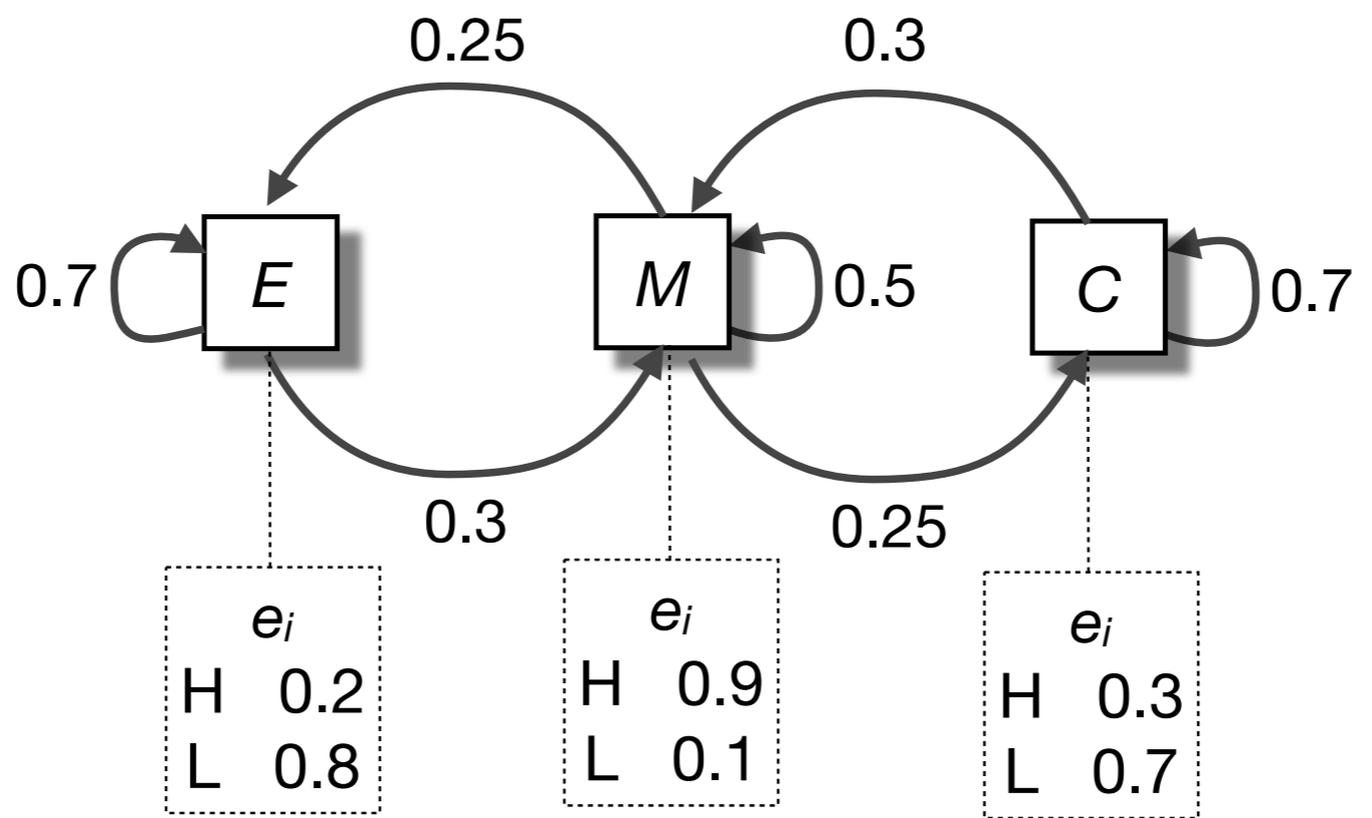
Note that we now have to learn the initial probabilities, transition probabilities and *emission probabilities*

$$a_{ij} = \frac{A_{ij}}{\sum_{j'} A_{ij'}}$$

$$e_i(x) = \frac{E_i(x)}{\sum_{x'} E_i(x')}$$



	<i>E</i>	<i>M</i>	<i>C</i>
π_i	0	0	1
$e_i(H)$	0.2	0.9	0.3
$e_i(L)$	0.8	0.1	0.7



$\pi(E) = 0$
 $\pi(M) = 0$
 $\pi(C) = 1$

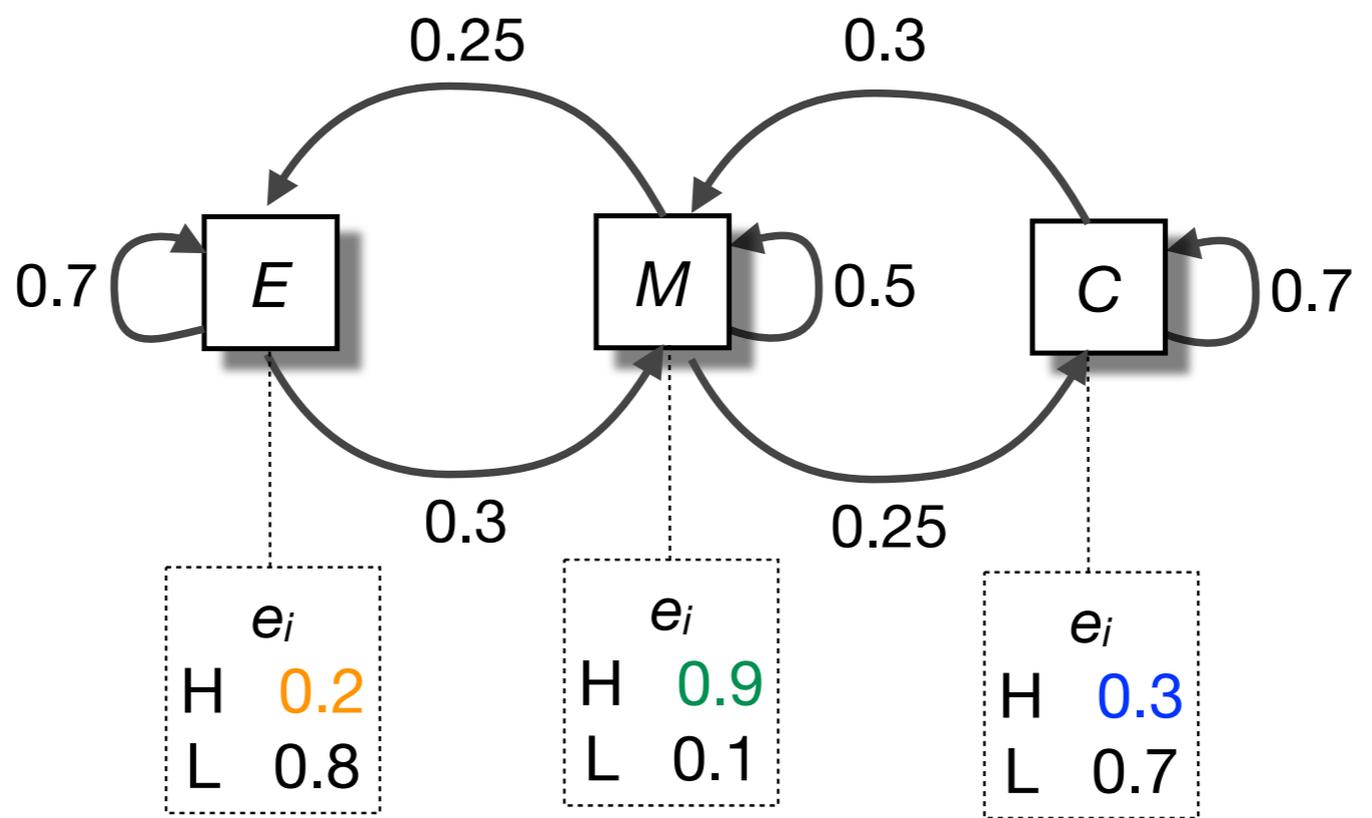
e_i
 H 0.2
 L 0.8

e_i
 H 0.9
 L 0.1

e_i
 H 0.3
 L 0.7

Query Sequence

States	H	H	L	L	H
<i>E</i>					
<i>M</i>					
<i>C</i>					
START					



$\pi(E) = 0$
 $\pi(M) = 0$
 $\pi(C) = 1$

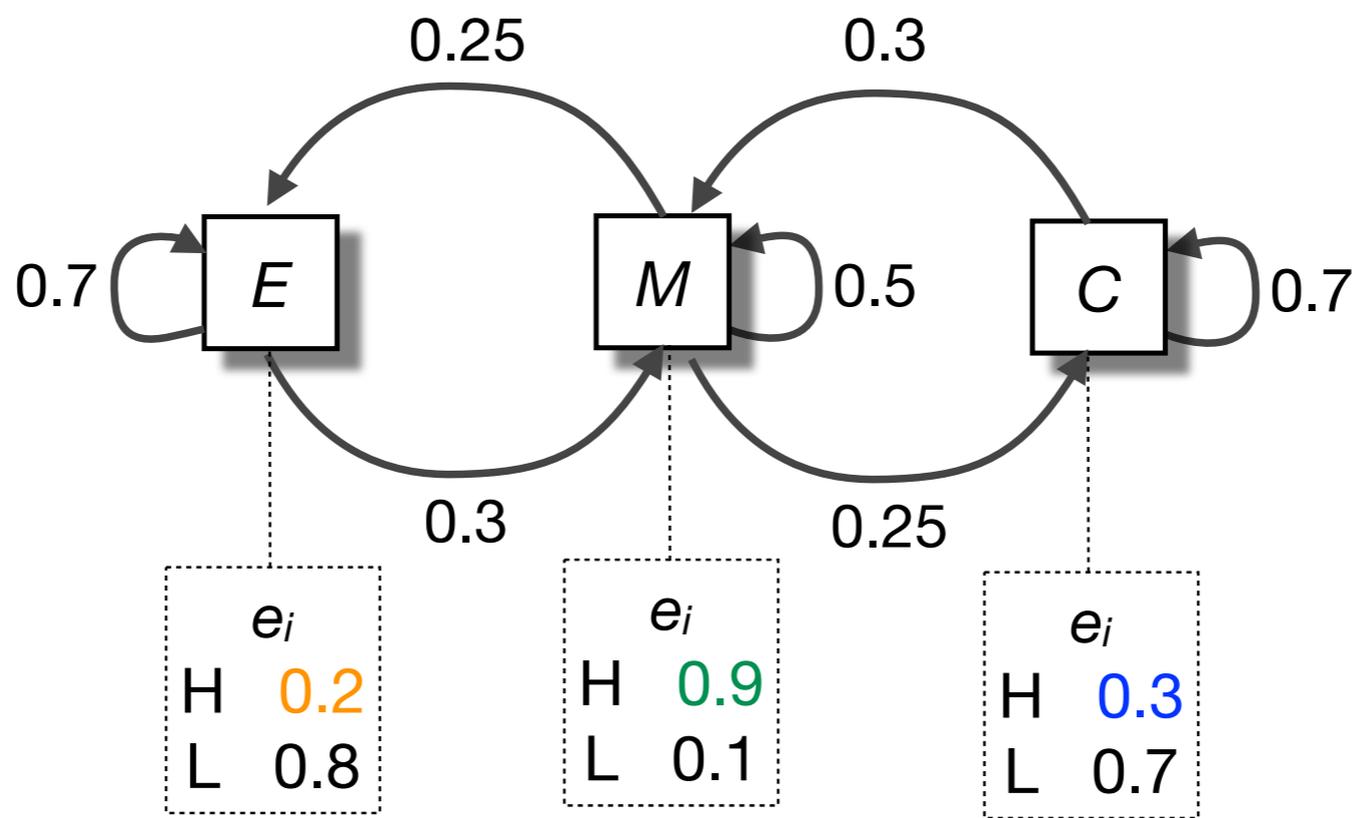
e_i
 H 0.2
 L 0.8

e_i
 H 0.9
 L 0.1

e_i
 H 0.3
 L 0.7

Query Sequence

States	H	H	L	L	H
E	0x0.2 =0				
M	0x0.9 =0				
C	1x0.3 =0.3				
START					



$\pi(E) = 0$
 $\pi(M) = 0$
 $\pi(C) = 1$

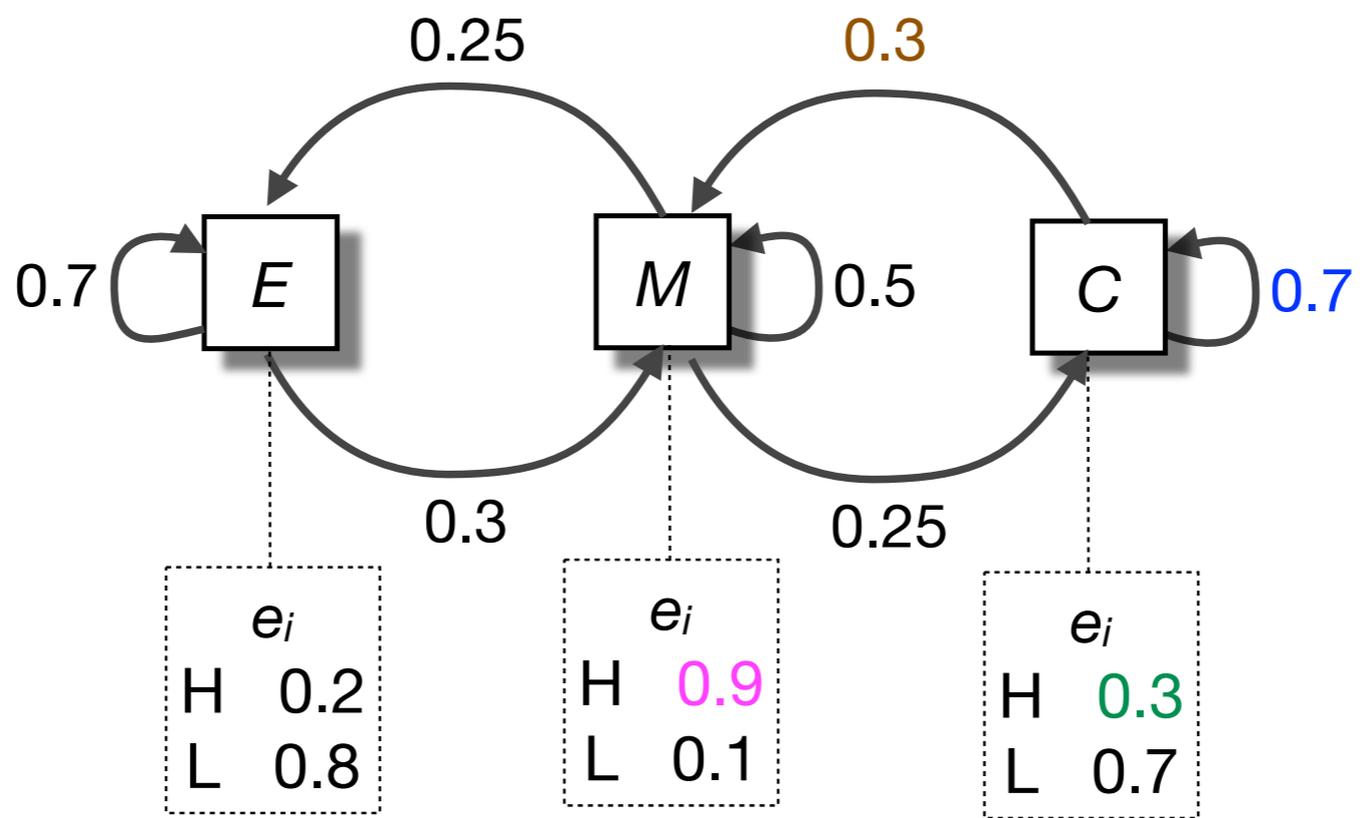
e_i
 H 0.2
 L 0.8

e_i
 H 0.9
 L 0.1

e_i
 H 0.3
 L 0.7

Query Sequence

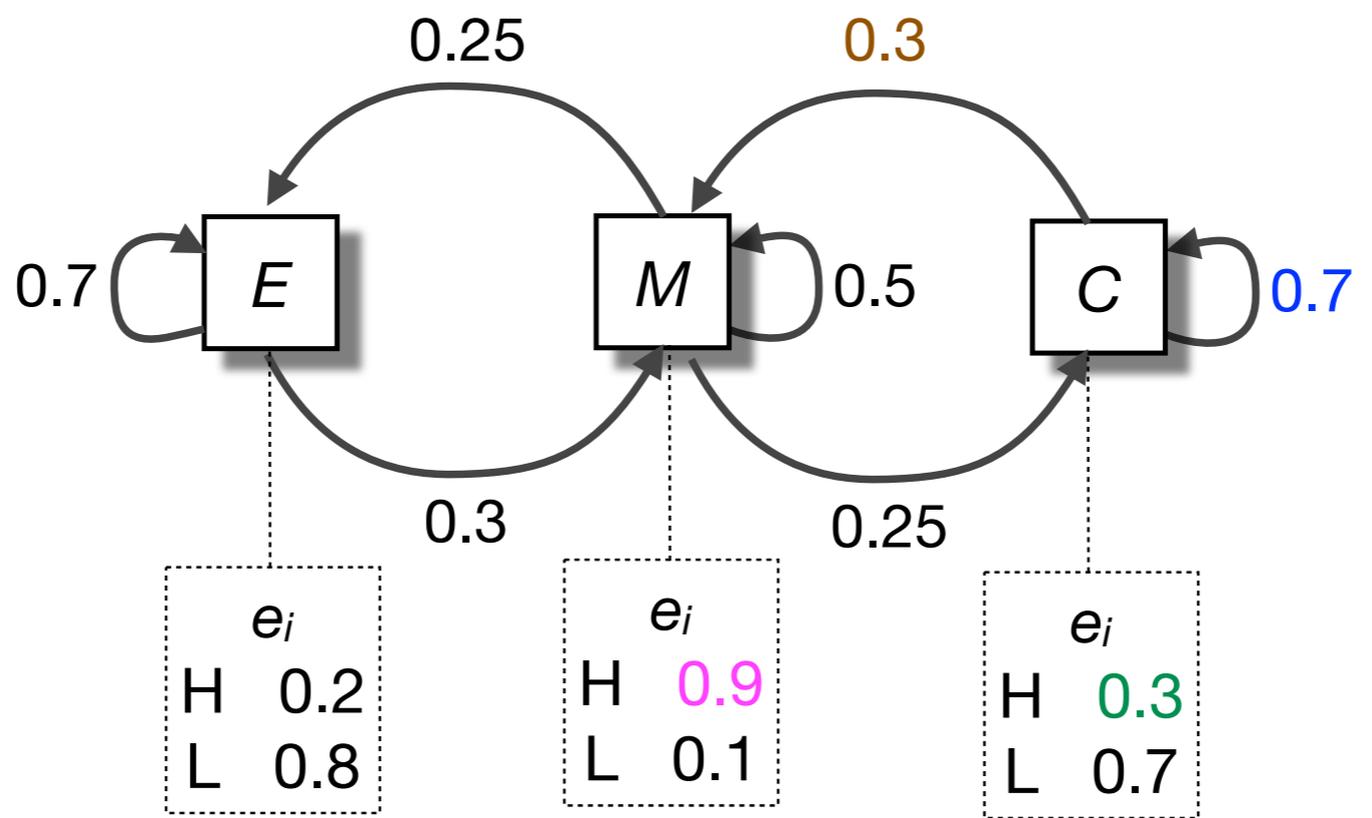
States	H	H	L	L	H
E	0x0.2 =0				
M	0x0.9 =0				
C	1x0.3 =0.3				
START					



Query Sequence

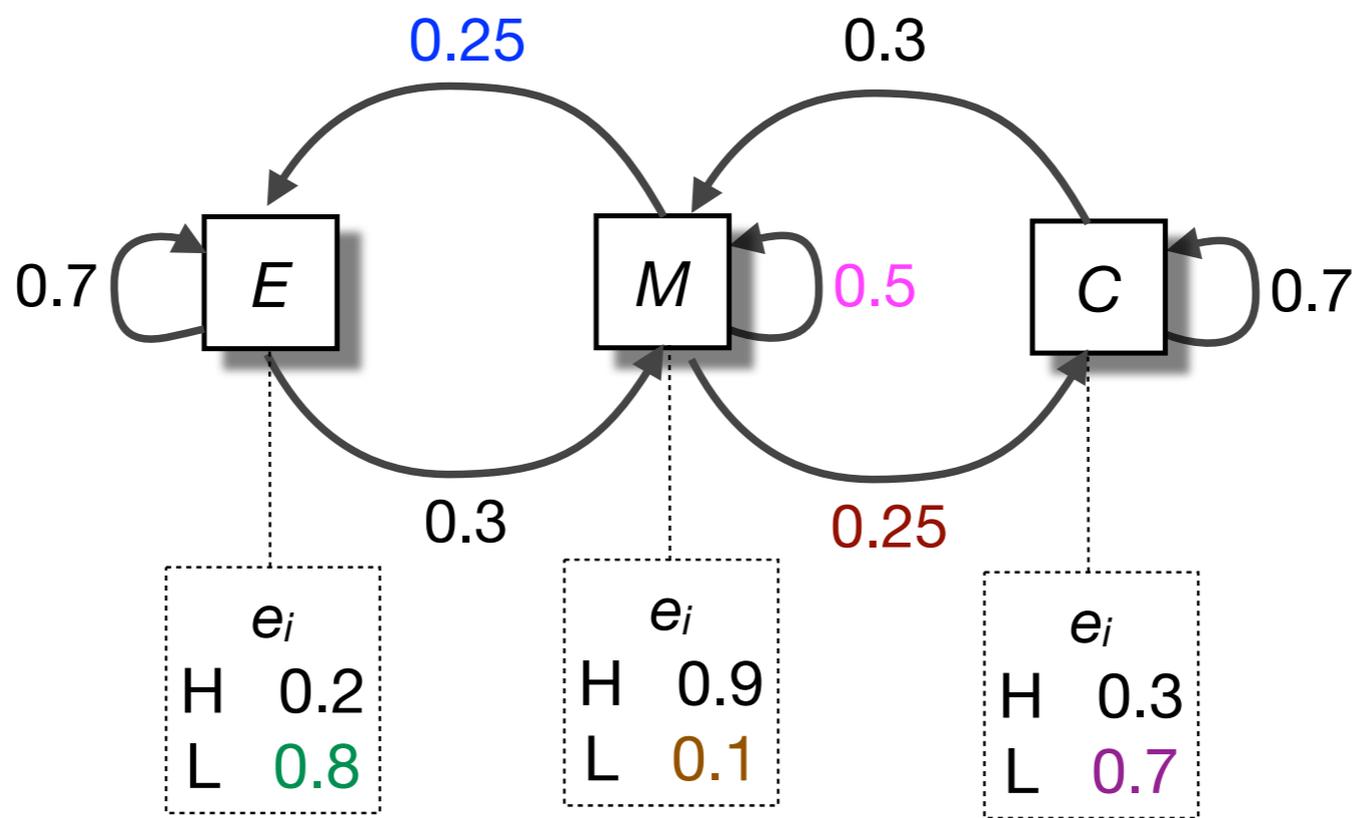
States	H	H	L	L	H
<i>E</i>	0×0.2 =0	-			
<i>M</i>	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081			
<i>C</i>	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063			

START



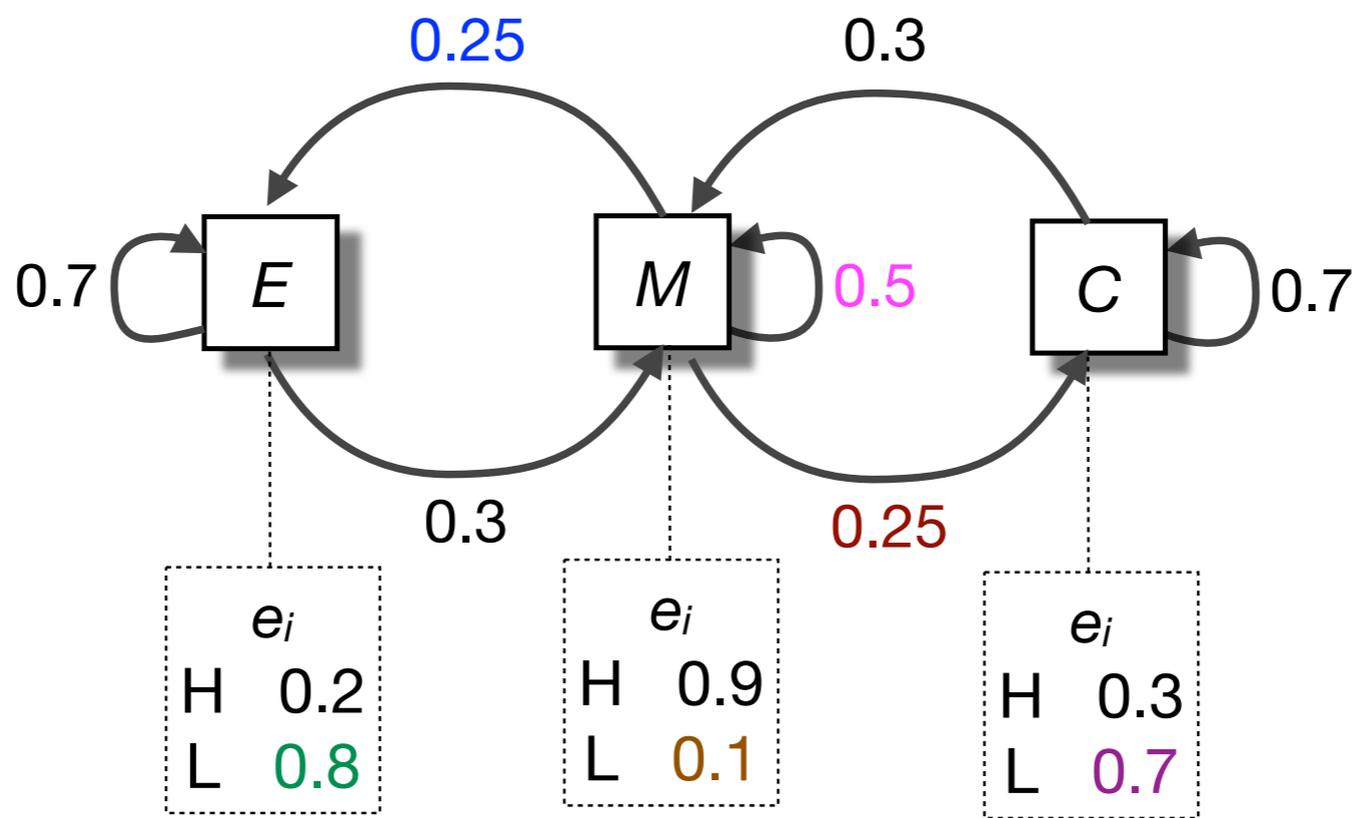
Query Sequence

States	H	H	L	L	H
<i>E</i>	0×0.2 =0	-			
<i>M</i>	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081			
<i>C</i>	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063			
START					



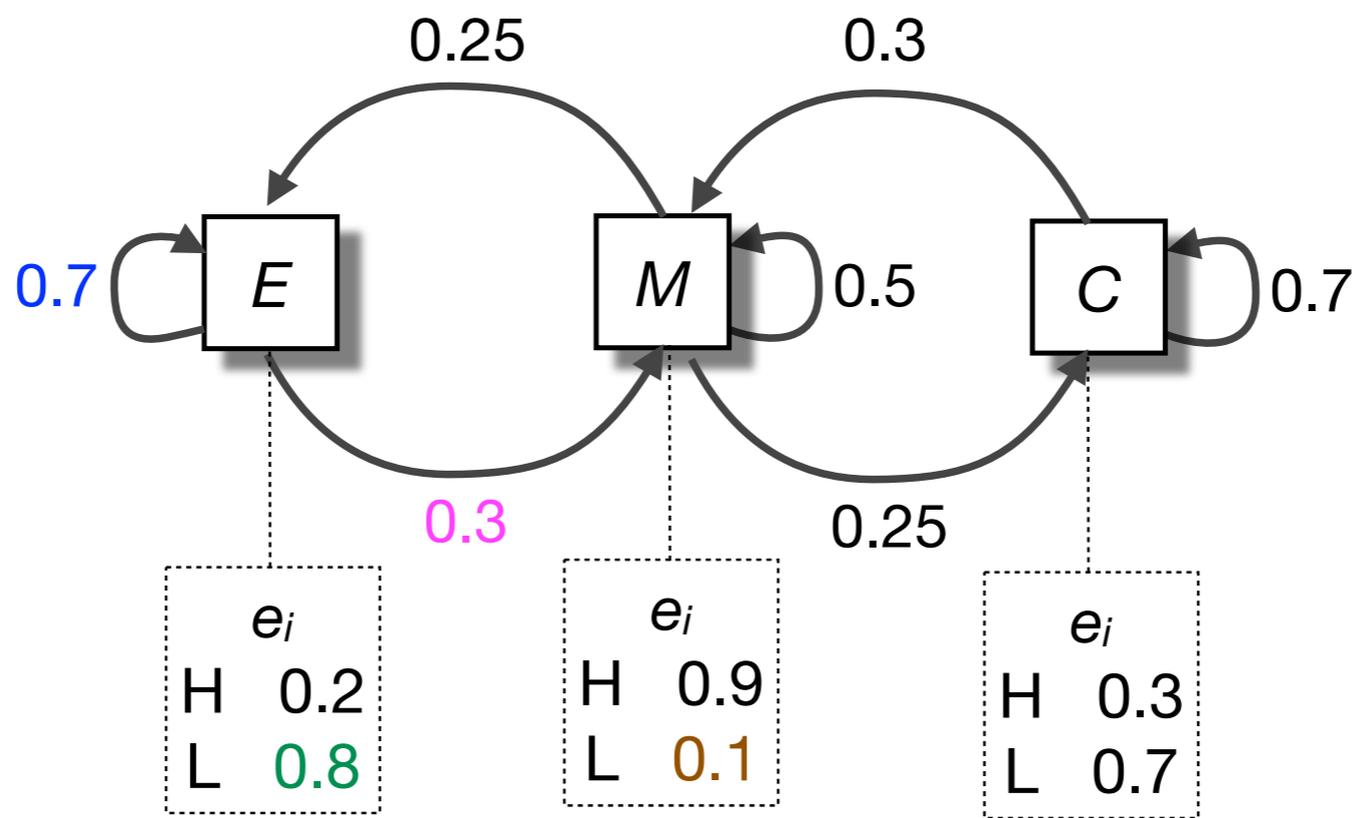
Query Sequence

States	H	H	L	L	H
<i>E</i>	0×0.2 =0	-	$0.25 \times 0.8 \times 0.081$ =0.016		
<i>M</i>	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081	$0.5 \times 0.1 \times 0.081$ =0.04		
<i>C</i>	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063	$0.25 \times 0.7 \times 0.081$ =0.014		
START					



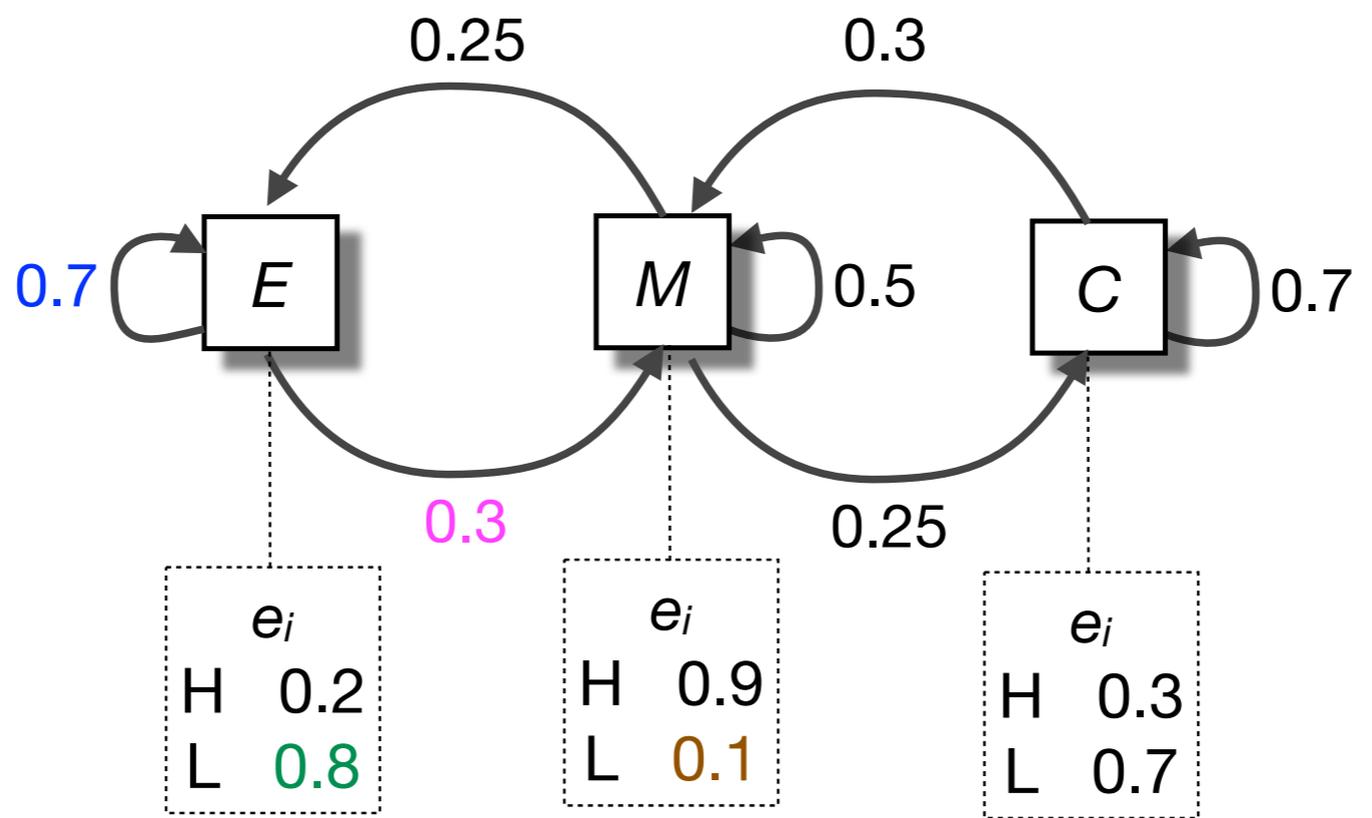
Query Sequence

States	H	H	L	L	H
<i>E</i>	0x0.2 =0	-	0.25x0.8x0.081 =0.016		
<i>M</i>	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04		
<i>C</i>	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014		
START					



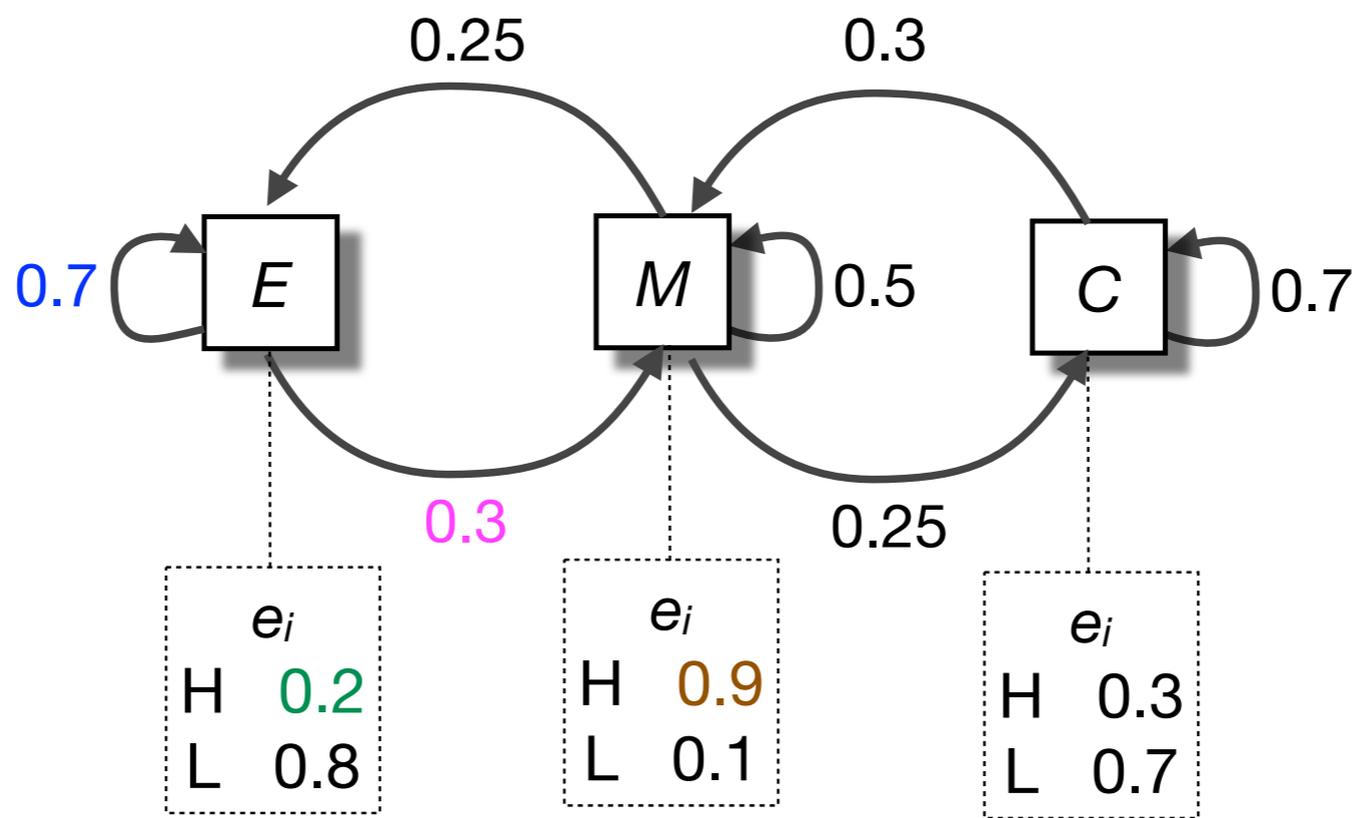
Query Sequence

States	H	H	L	L	H
<i>E</i>	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	
<i>M</i>	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	
<i>C</i>	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	
START					



Query Sequence

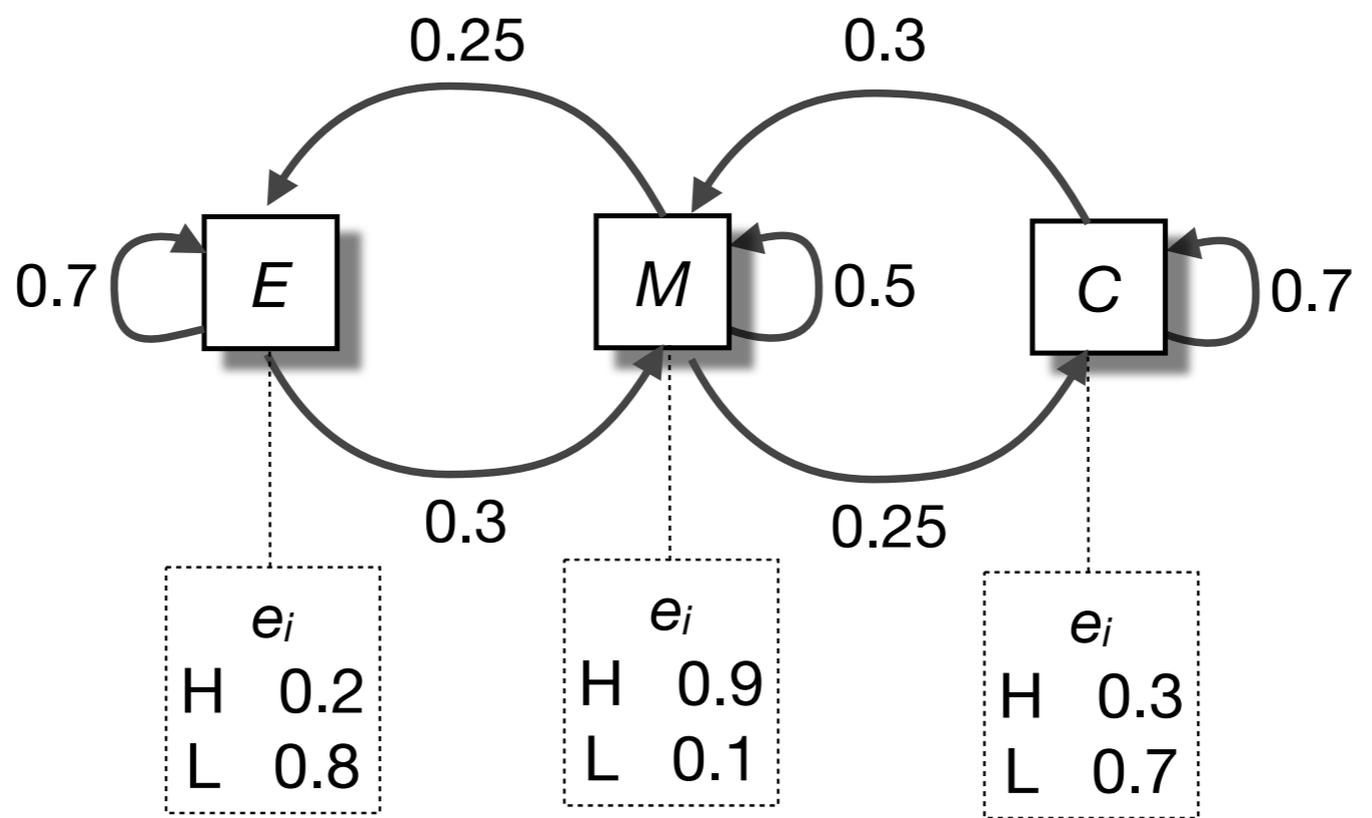
States	H	H	L	L	H
<i>E</i>	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	
<i>M</i>	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	
<i>C</i>	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	
START					



Query Sequence

States	H	H	L	L	H
<i>E</i>	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	0.7x0.2x0.009 =0.001
<i>M</i>	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	0.3x0.9x0.009 =0.002
<i>C</i>	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	-

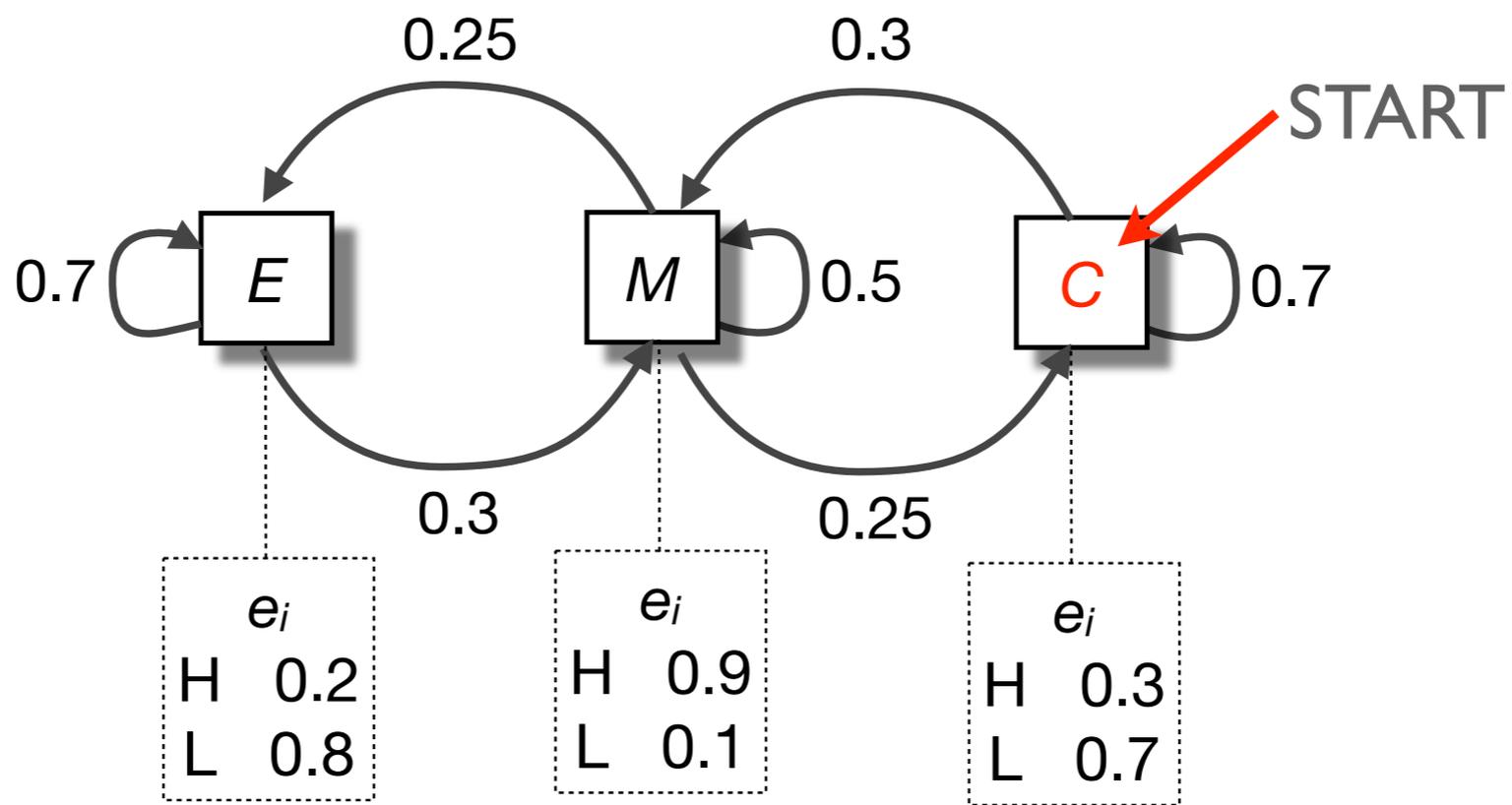
START



Query Sequence

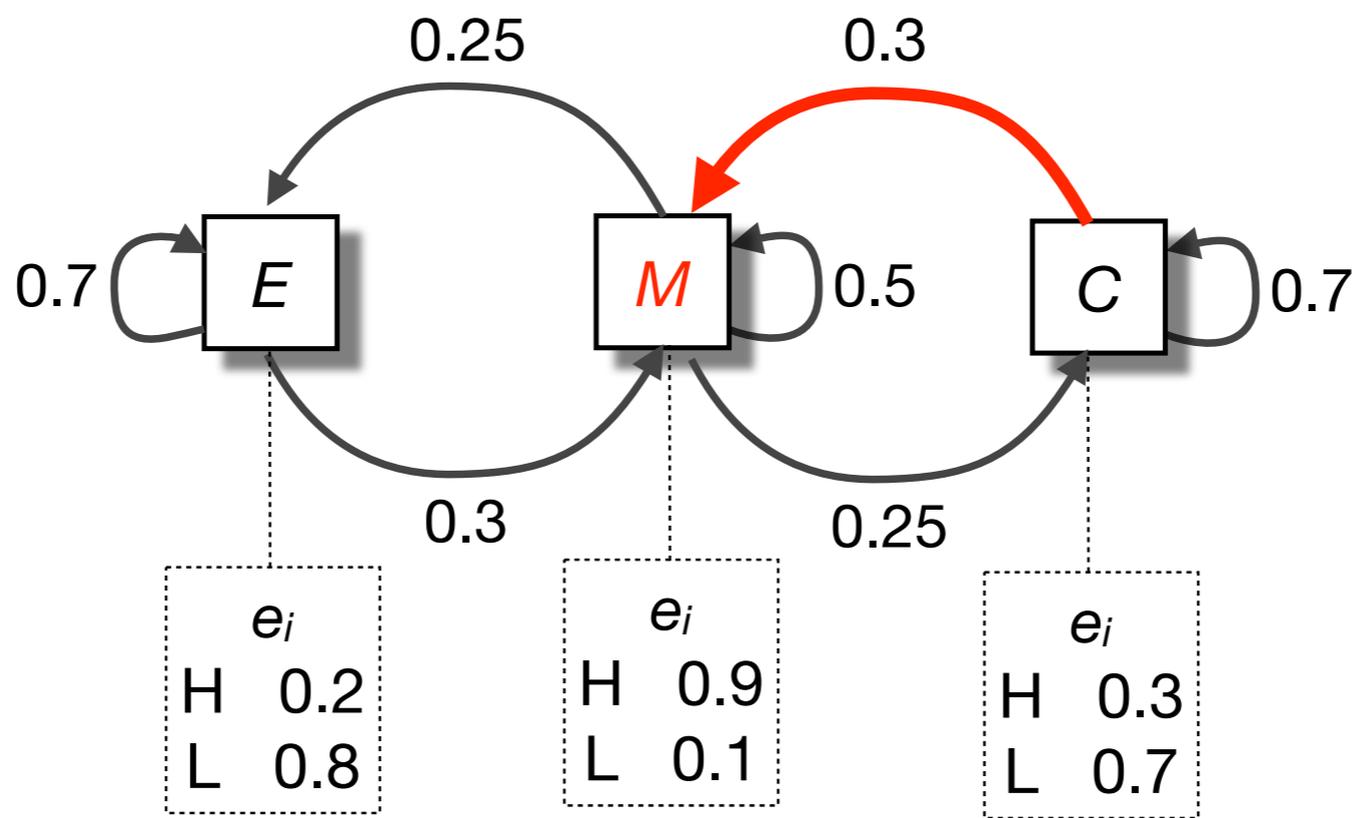
States	H	H	L	L	H
<i>E</i>	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	0.7x0.2x0.009 =0.001
<i>M</i>	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	0.3x0.9x0.009 =0.002
<i>C</i>	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	-

START



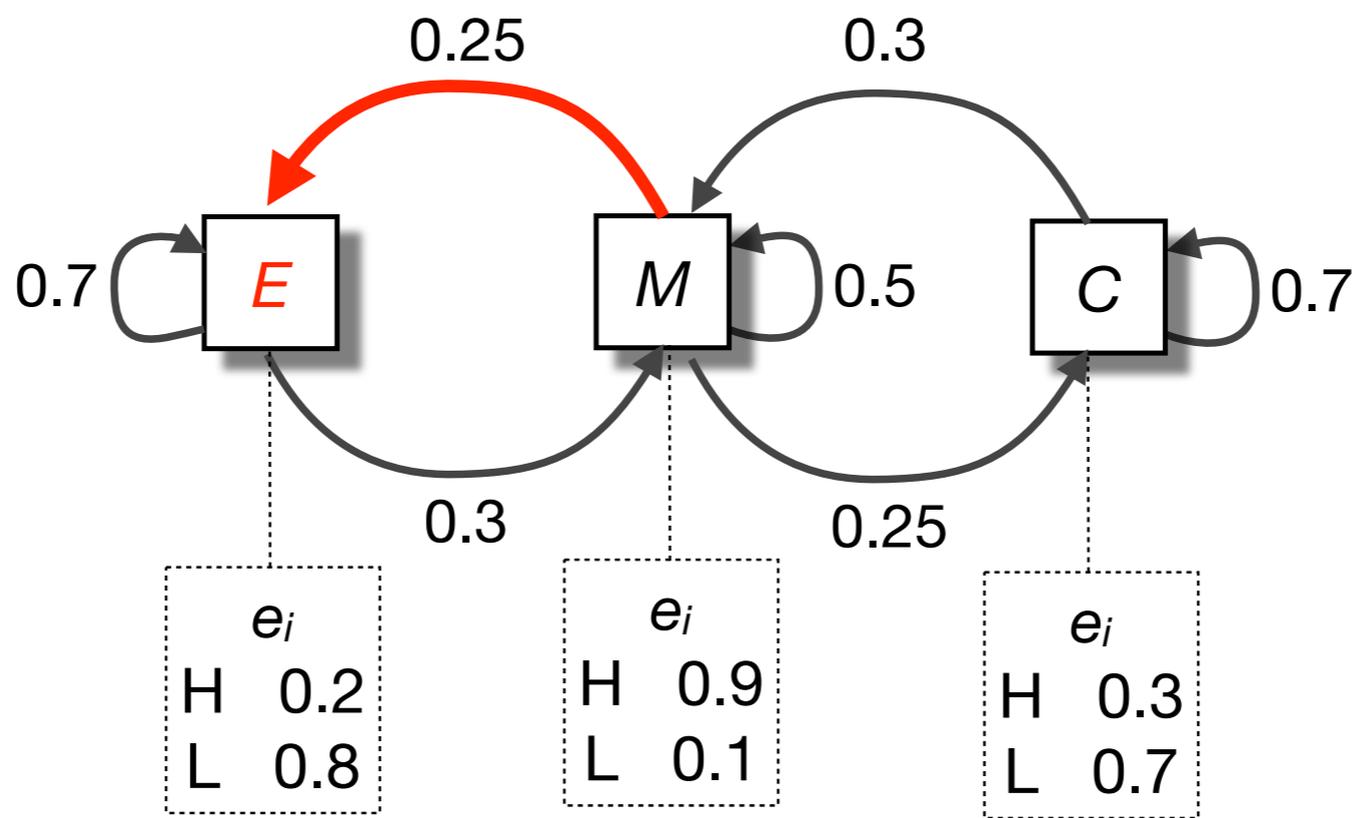
Query Sequence

States	H	H	L	L	H
E	0×0.2 =0	-	$0.25 \times 0.8 \times 0.081$ =0.016	$0.7 \times 0.8 \times 0.016$ =0.009	$0.7 \times 0.2 \times 0.009$ =0.001
M	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081	$0.5 \times 0.1 \times 0.081$ =0.04	$0.3 \times 0.1 \times 0.016$ =0.0005	$0.3 \times 0.9 \times 0.009$ =0.002
C	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063	$0.25 \times 0.7 \times 0.081$ =0.014	-	-
START	C				



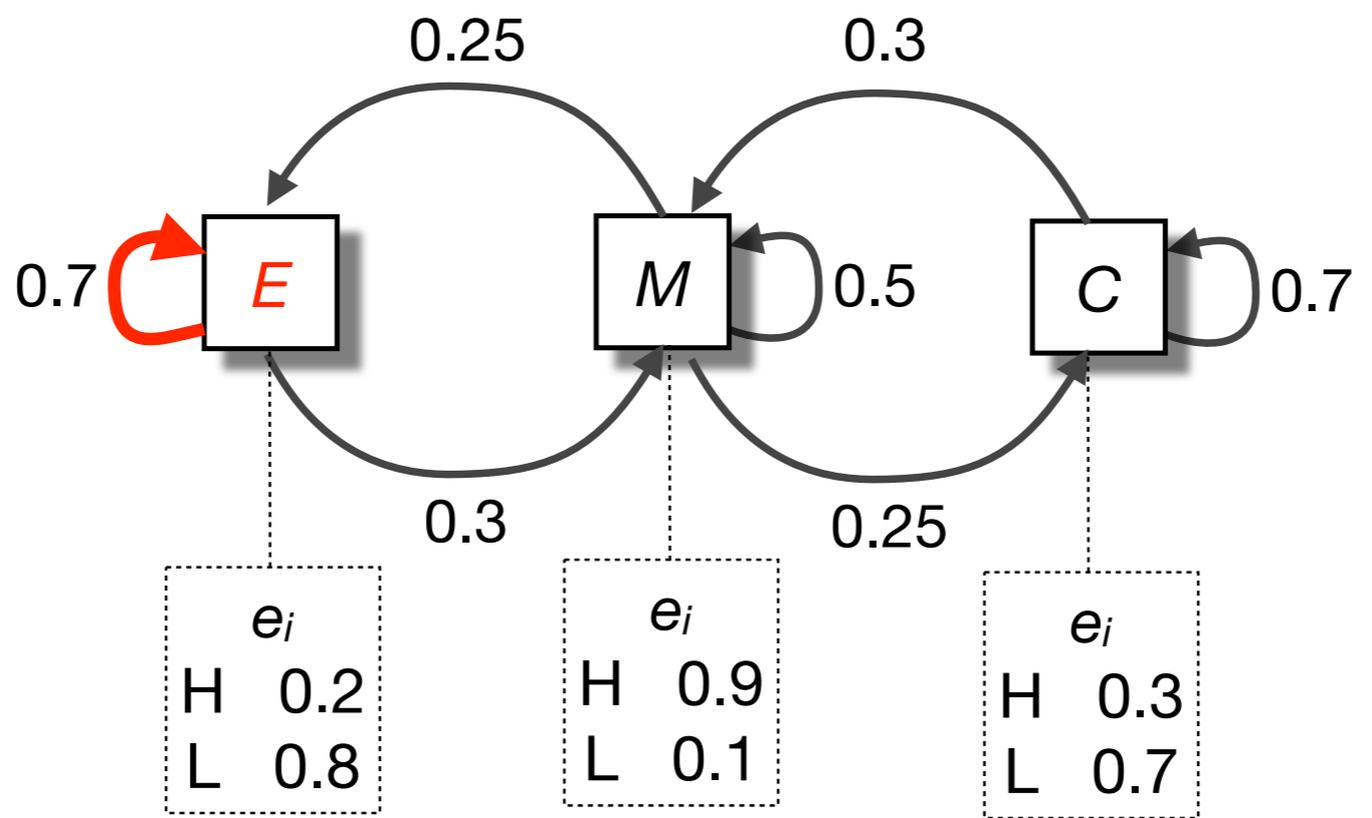
Query Sequence

States	H	H	L	L	H
E	0×0.2 =0	-	$0.25 \times 0.8 \times 0.081$ =0.016	$0.7 \times 0.8 \times 0.016$ =0.009	$0.7 \times 0.2 \times 0.009$ =0.001
M	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081	$0.5 \times 0.1 \times 0.081$ =0.04	$0.3 \times 0.1 \times 0.016$ =0.0005	$0.3 \times 0.9 \times 0.009$ =0.002
C	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063	$0.25 \times 0.7 \times 0.081$ =0.014	-	-
START	C	M			



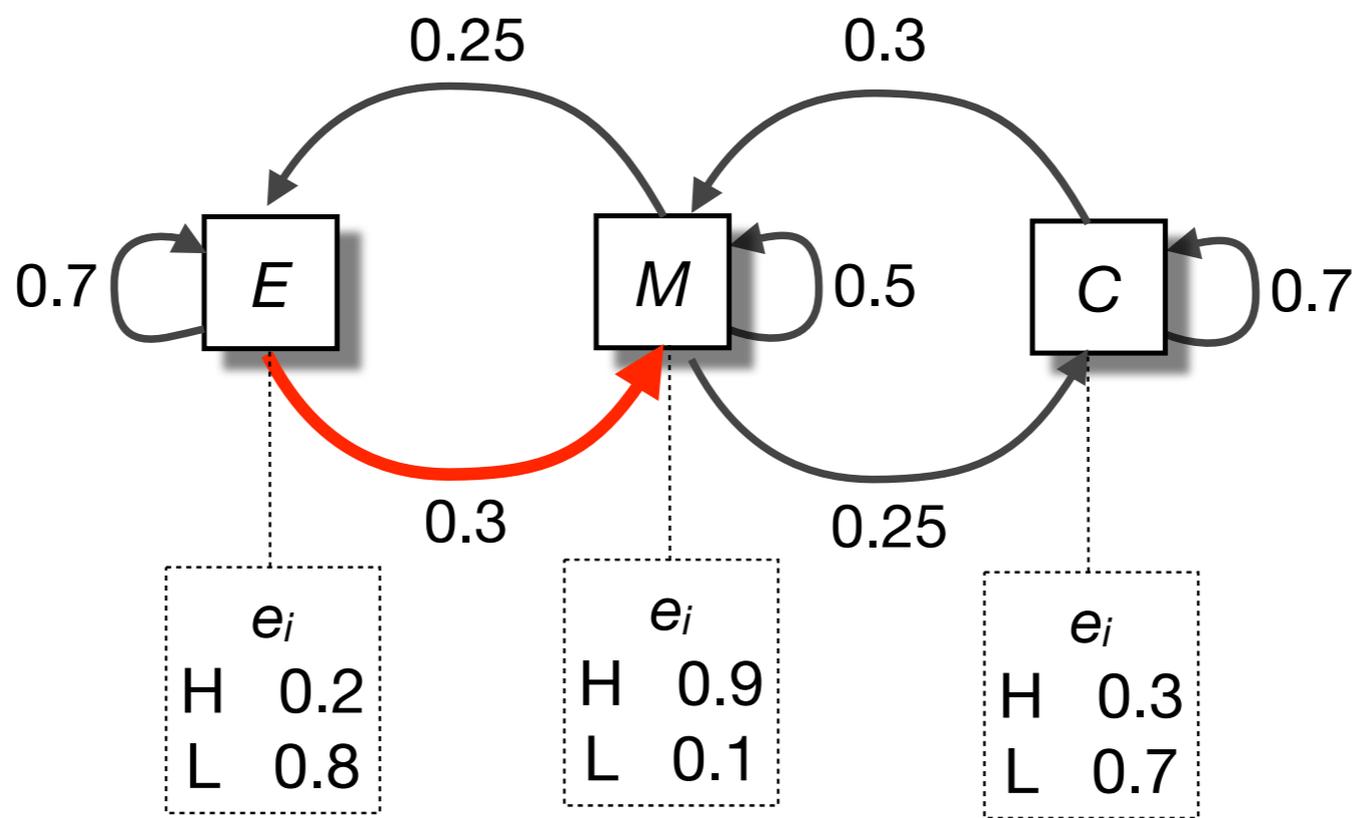
Query Sequence

States	H	H	L	L	H
E	0×0.2 =0	-	$0.25 \times 0.8 \times 0.081$ =0.016	$0.7 \times 0.8 \times 0.016$ =0.009	$0.7 \times 0.2 \times 0.009$ =0.001
M	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081	$0.5 \times 0.1 \times 0.081$ =0.04	$0.3 \times 0.1 \times 0.016$ =0.0005	$0.3 \times 0.9 \times 0.009$ =0.002
C	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063	$0.25 \times 0.7 \times 0.081$ =0.014	-	-
START	C	M	E		



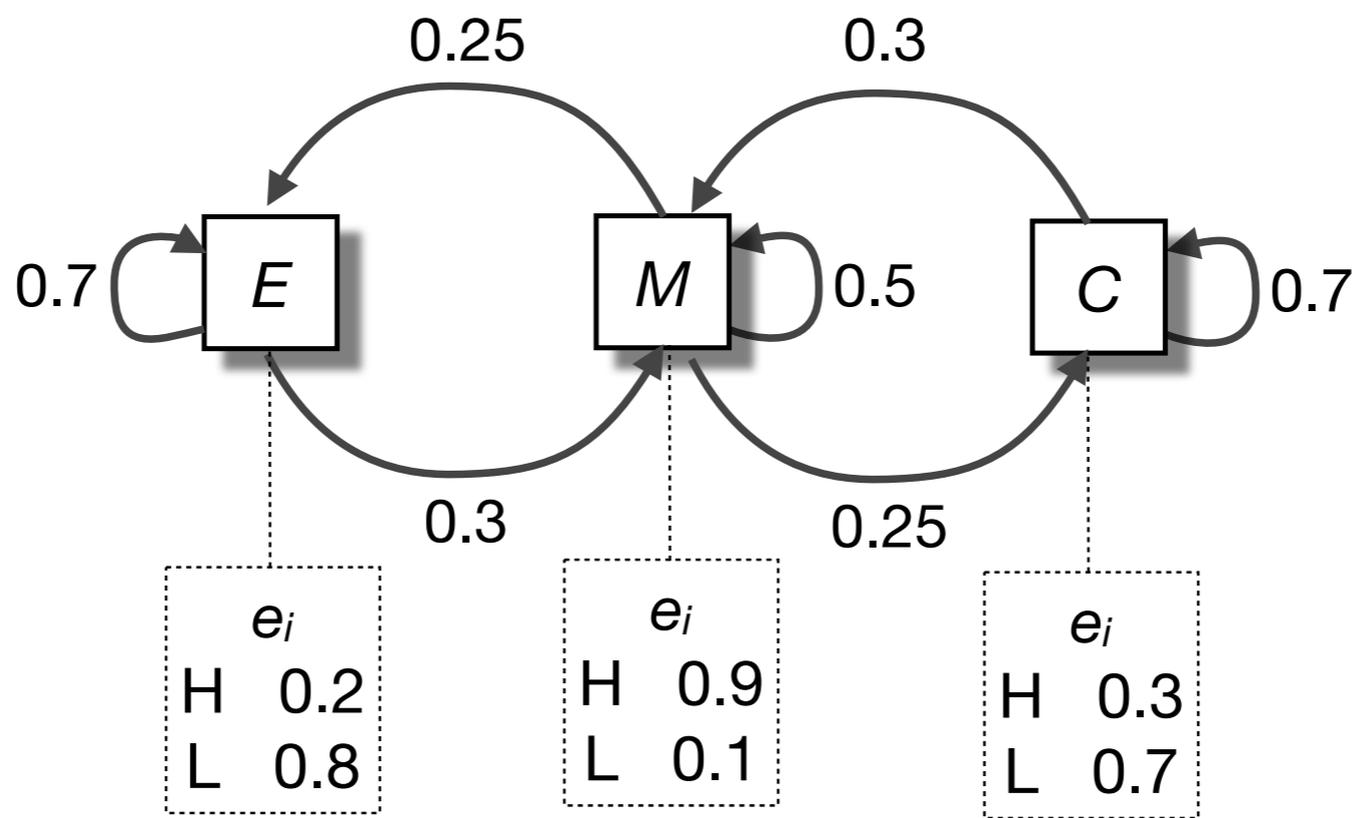
Query Sequence

States	H	H	L	L	H
E	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	0.7x0.2x0.009 =0.001
M	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	0.3x0.9x0.009 =0.002
C	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	-
START	C	M	E	E	



Query Sequence

States	H	H	L	L	H
E	0x0.2 =0	-	0.25x0.8x0.081 =0.016	0.7x0.8x0.016 =0.009	0.7x0.2x0.009 =0.001
M	0x0.9 =0	0.3x0.9x0.3 =0.081	0.5x0.1x0.081 =0.04	0.3x0.1x0.016 =0.0005	0.3x0.9x0.009 =0.002
C	1x0.3 =0.3	0.7x0.3x0.3 =0.063	0.25x0.7x0.081 =0.014	-	-
START	C	M	E	E	M



Query Sequence

States	H	H	L	L	H
E	0×0.2 =0	-	$0.25 \times 0.8 \times 0.081$ =0.016	$0.7 \times 0.8 \times 0.016$ =0.009	$0.7 \times 0.2 \times 0.009$ =0.001
M	0×0.9 =0	$0.3 \times 0.9 \times 0.3$ =0.081	$0.5 \times 0.1 \times 0.081$ =0.04	$0.3 \times 0.1 \times 0.016$ =0.0005	$0.3 \times 0.9 \times 0.009$ =0.002
C	1×0.3 =0.3	$0.7 \times 0.3 \times 0.3$ =0.063	$0.25 \times 0.7 \times 0.081$ =0.014	-	-
START	C	M	E	E	M

Most Probable State Sequence

We have just used the Viterbi algorithm

The **Viterbi algorithm** finds the most probable “state path” (S^*) (i.e. sequence of hidden states) for generating a given sequence ($x = x_1, x_2, \dots, x_N$)

$$S^* = \operatorname{argmax} P(x, S)$$

This process is often called **decoding** because we “decode” the sequence of symbols to determine the hidden sequence of states

HMMs were originally developed in the field of speech recognition, where speech is “decoded” into words or phonemes to determine the meaning of the utterance

Note that we could have used brute force by calculating $P(x|S)$ for all paths but this quickly becomes intractable for longer sequences or HMMs with a large number of states

The Viterbi algorithm is guaranteed to find the most probable state path given a sequence and an HMM

See Durbin *et al.* *Biological Sequence Analysis*

Three key HMM algorithms

- **Viterbi algorithm**

Given observed sequence x and an HMM M , composed of states S , calculate the most likely state sequence, S^*

- ▶ $S^* = \operatorname{argmax} P(x, S)$

- **Forward algorithm**

Given observed sequence x and an HMM composed of states S , calculate the probability of the sequence for the HMM, $P(x|M)$

- ▶
$$P(x) = \sum_S P(x, S)$$

- **Baum-Welch algorithm**

Given many observed sequences, estimate the parameters of the HMM

- ▶ heuristic expectation maximization method to optimize of a_{ij} and $e_i(a)$

The forward algorithm

Another important question is how well does a given sequence fit the HMM?

To answer this question we must sum over all possible state paths that are consistent with the sequence in question

(Because we don't know which path emitted the sequence)

The number of paths can quickly become intractable. The **forward algorithm** is a simple dynamic programming solution that makes use of the Markov property so that we don't have to explicitly enumerate every path.

The **forward algorithm** basically replaces the maximization step of the Viterbi algorithm with sums to calculate the probability of the sequence given a HMM.

$$P(x) = \sum_s P(x, S)$$

See Durbin *et al.* *Biological Sequence Analysis*

The Baum-Welch algorithm

The **Baum-Welch algorithm** is an **heuristic optimization** algorithm for learning probabilistic models in problems that involve hidden states

If we *know* the state path for each training sequence (i.e. no hidden states with respect to the training sequences), then learning the model parameters is simple (just like it was for Markov chain models)

- count how often each transition and emission occurs
- normalize to get probabilities

If we *don't know* the path for each training sequence, we can use the **Baum-Welch algorithm**, an expectation maximization method, which estimates counts by considering every path weighted by its probability

- start from a given initial guess for the parameters
- perform a calculation which is guaranteed to improve the previous guess
- run until there is little change in parameters between iterations

For sequence profile-HMMs we train from a MSA and hence we can *estimate* our probabilities from the observed sequences

Segmentation/boundary detection

Given: A test sequence and a HMM with different sequence classes

Task: Segment the sequence into subsequences, predicting the class of each subsequence

Question: What is the most probable “path” (sequence of hidden states) for generating a given sequence from the HMM?

Solution: Use the **Viterbi algorithm**

Classification/sequence scoring

Given: A test sequence and a set of HMMs representing different sequence classes

Task: Determine which HMM/class best explains the sequence

Question: How likely is a given sequence given a HMM?

Solution: Use the **Forward algorithm**

Learning/parameterization

Given: A model, a set of training sequences

Task: Find model parameters that explain the training sequences

Question: Can we find a high probability model for sequence characterization

Solution: Use the **Forward backward algorithm**

Segmentation/boundary detection

Question: What is the most probable “path” (sequence of hidden states) for generating a given sequence from the HMM?

HMMER: **hmmalign** - align sequences to our HMM

Classification/sequence scoring

Question: How likely is a given sequence given a HMM?

HMMER: **hmmsearch** - find sequences that match our HMM

Learning/parameterisation

Question: Can we find a high probability model for sequence characterization

HMMER: **hmmbuild** - setup our HMM parameters

Half time break...

Questions:

For what kinds of motifs are PSSMs not well suited?

What is the Markov property?

In what important ways do HMMs differ Markov chains?

What is the Viterbi algorithm used for?

How does the Forward algorithm differ from the Viterbi algorithm?

For what kinds of motifs are PSSMs not well suited?

PSSMs are not well suited to pattern instances containing insertions or deletions, variable length patterns and those with positional dependencies.

What is the Markov property?

The Markov property states that the conditional probability distribution for the system at the next step (and in fact at all future steps) depends only on the current state of the system, and not additionally on the state of the system at previous steps.

In what important ways do HMMs differ Markov chains?

HMMs differ from Markov chains in a number of ways:

- In HMMs, the sequence of states visited is hidden. Unlike Markov Chains, there is no longer a one-to-one correspondence between states and output symbols.
- In a HMM the same symbol may be emitted by more than one state.
- In a HMM a state can emit more than one symbol.

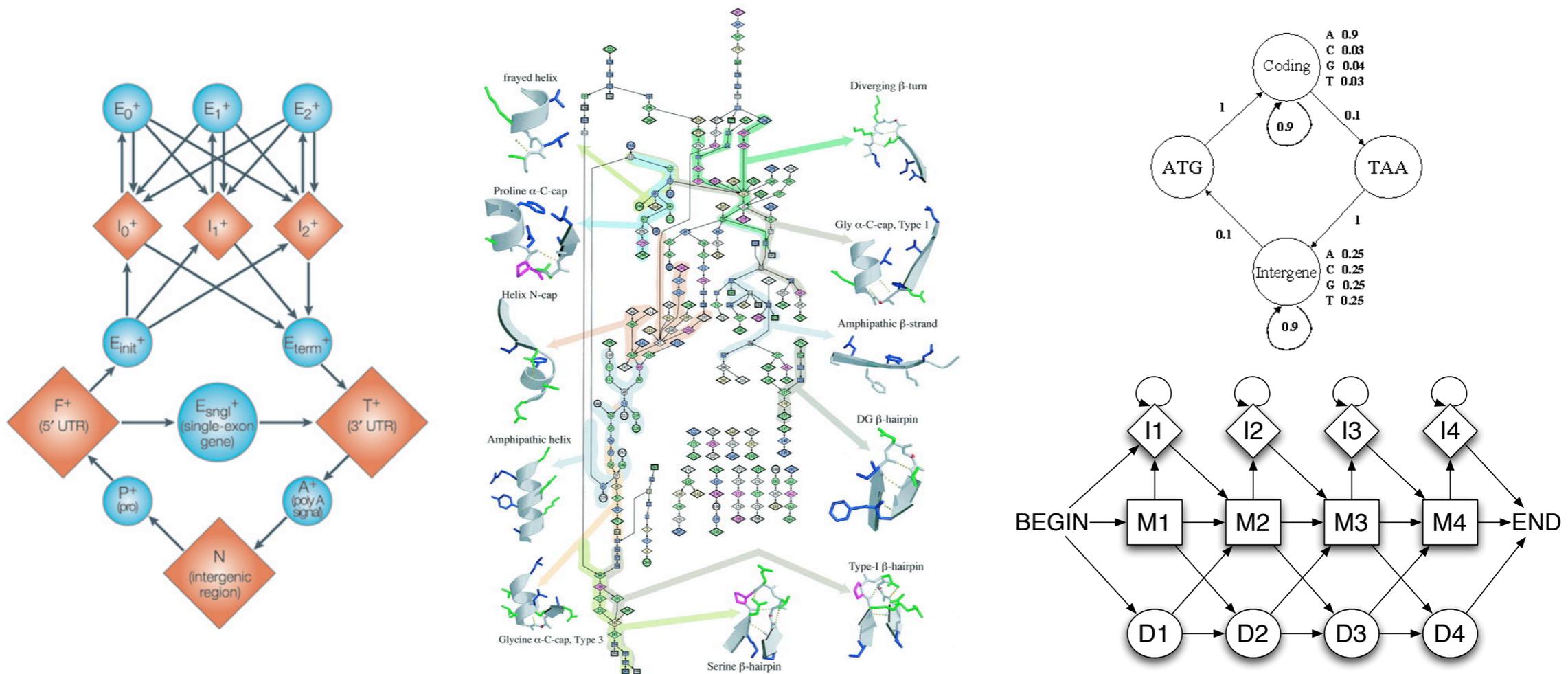
What is the Viterbi algorithm used for?

The Viterbi algorithm is used to find the most probable state path given a sequence and an HMM

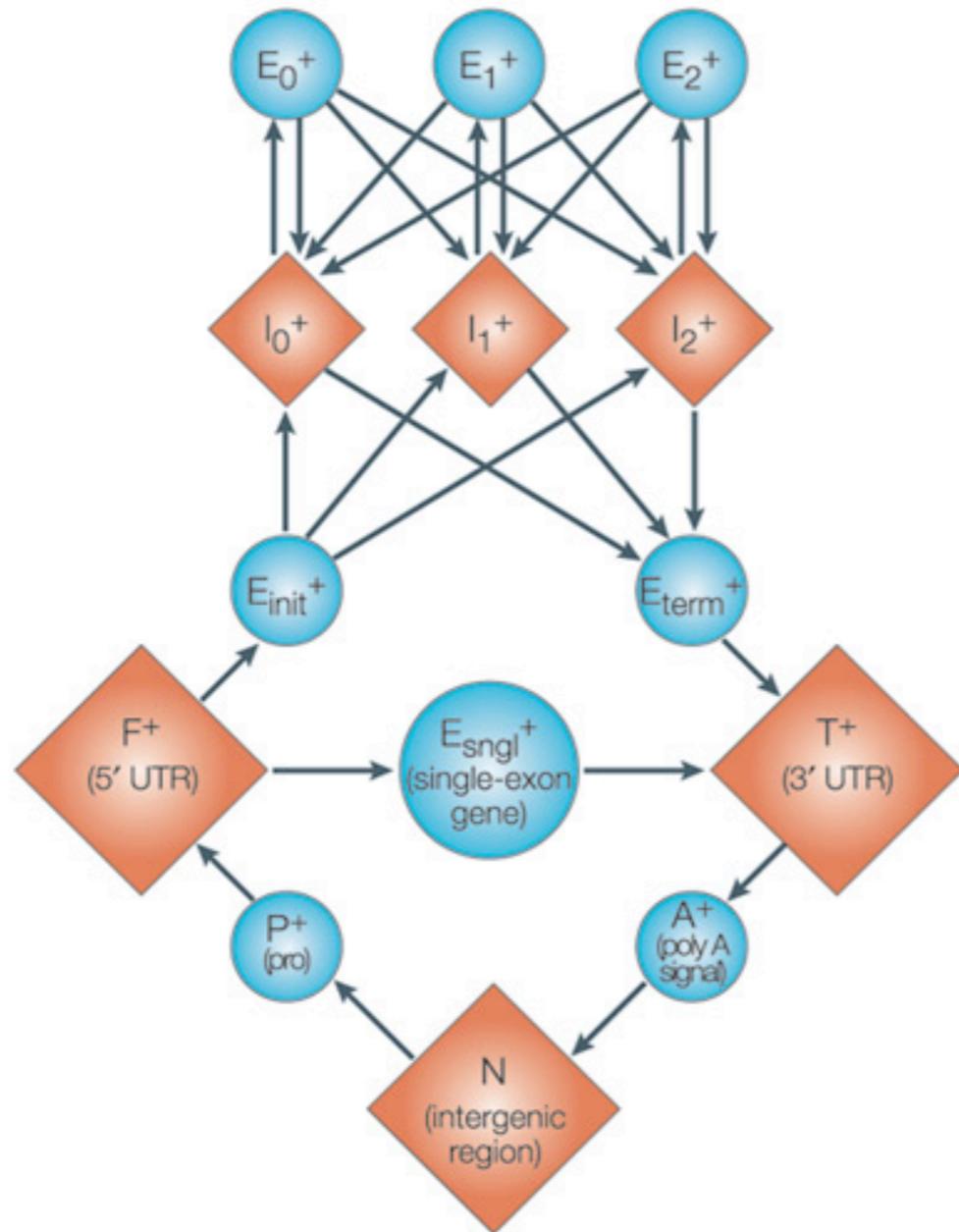
HMM network structure is hand tailored to the problem

No algorithm for the prediction of optimal HMM network structure and probabilities has yet been able to beat simple hand-built topologies

These topologies are tailored to the problem at hand - exon/intron detection, transmembrane regions, secondary structure elements, protein families...



GenScan - gene-prediction HMM

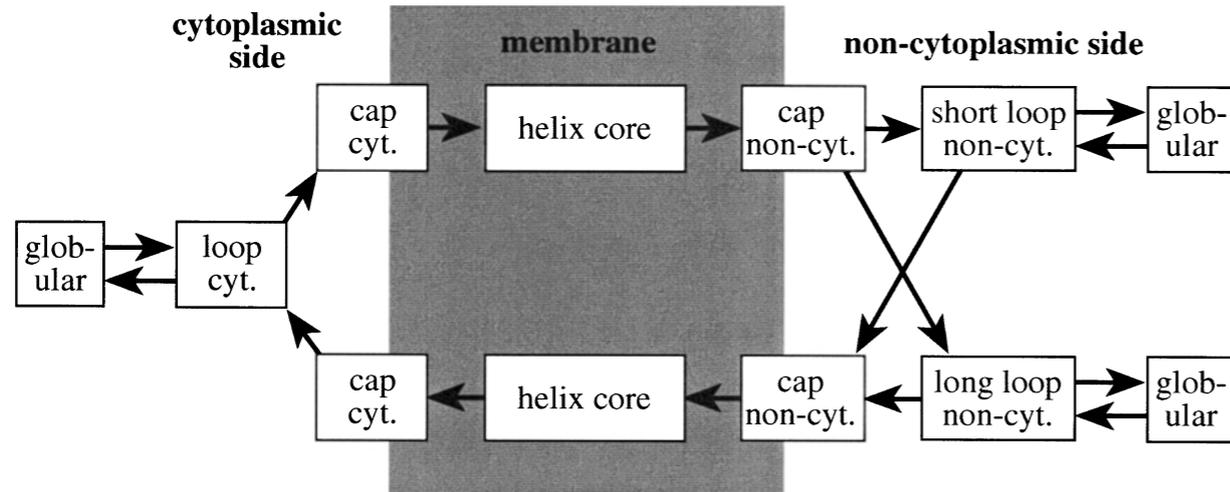


Here, each circle or square represents a functional unit (a state) of a gene on its forward strand (for example, E_{init} is the 5' coding sequence (CDS) and E_{term} is the 3' CDS, and the arrows represent the transition probability from one state to another. The GenScan HMM is trained by pre-computing the transition probabilities from a set of known gene structures.

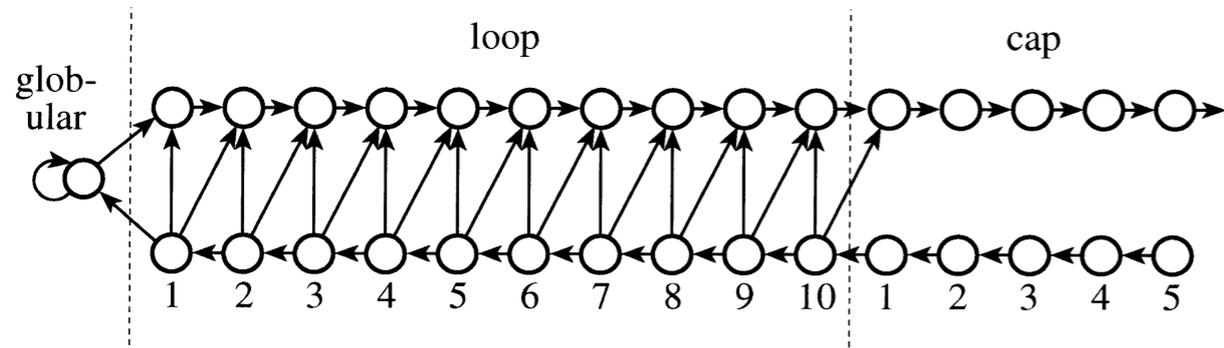
See: Zhang et al. (2002) Nature Reviews Genetics 3, 698-709

Reverse strand: mirror reflection of above

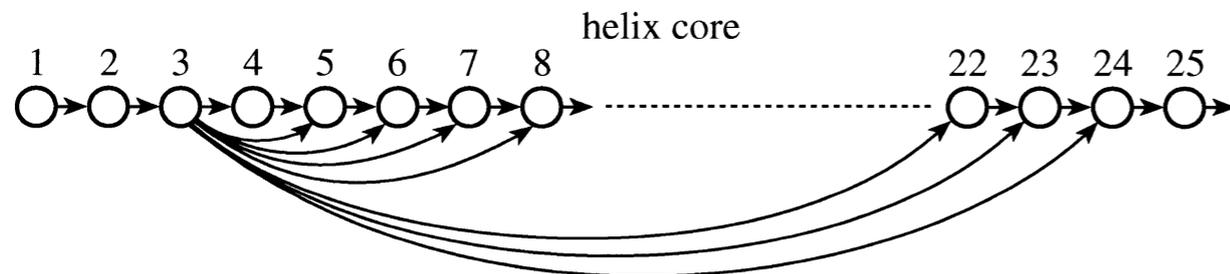
TMHMM - transmembrane protein topology prediction



(b)



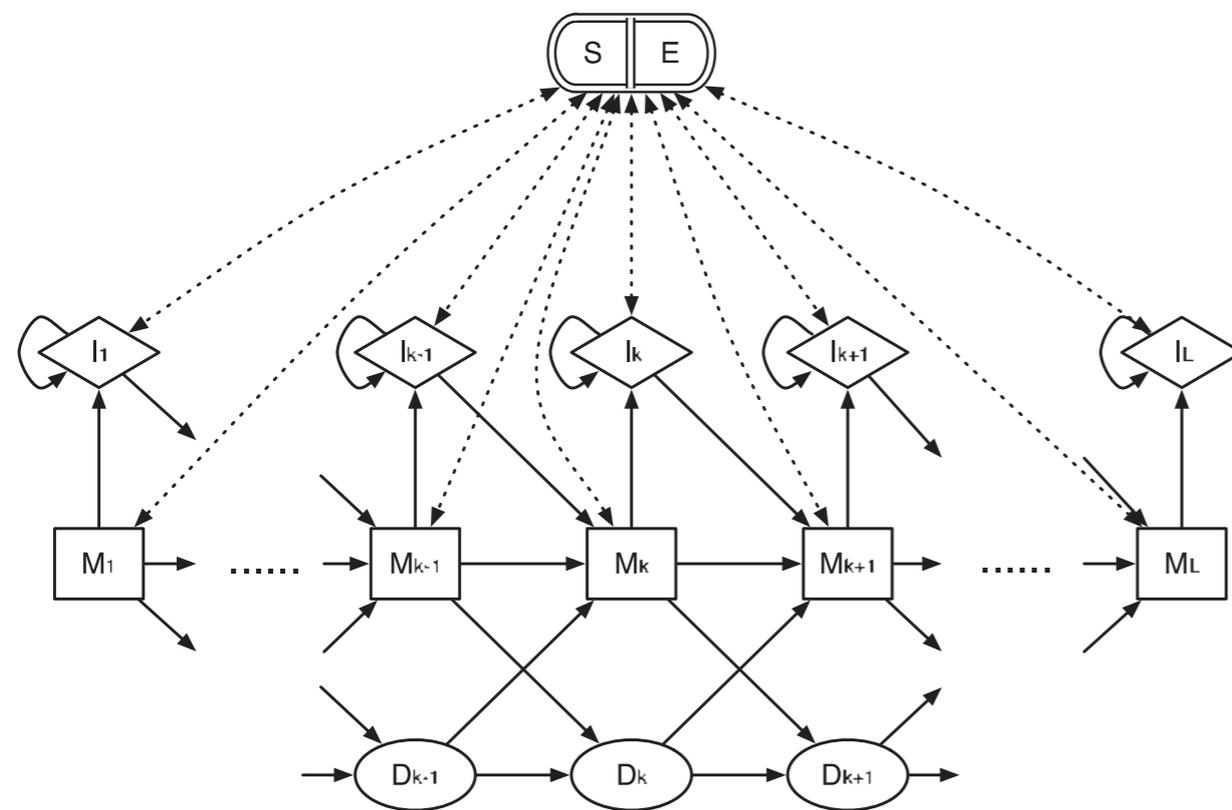
(c)



Each box corresponds to one or more states in the HMM. Cyt. represents the cytoplasmic side of the membrane and non-cyt. the other side. (b) The detailed structure of the inside and outside loop models and helix cap models. (c) The structure of the model for the helix core modeling lengths between 5 and 25, which translates to helices between 15 and 35 when the caps are included.

See: Krogh et al. (2001) JMB 305, 567-580

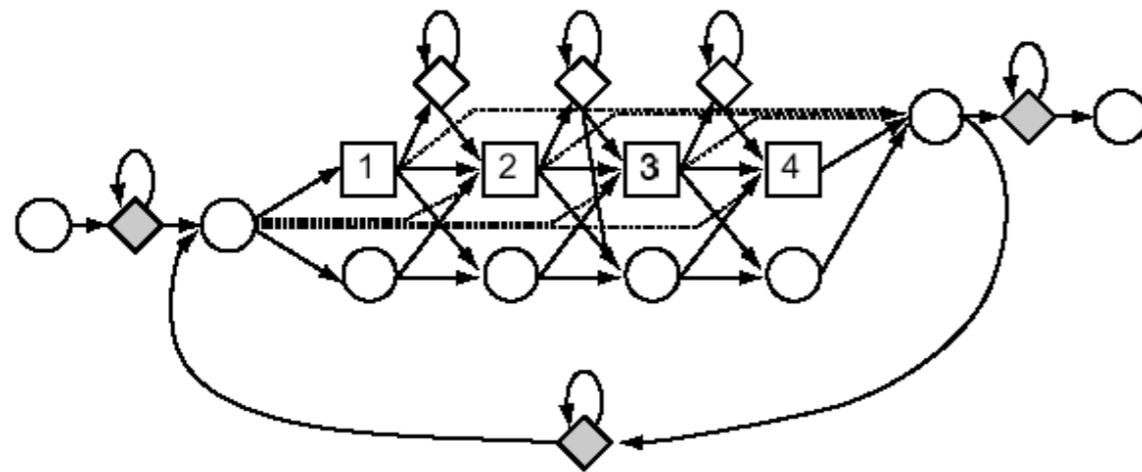
SAMTOOLS - SNP calling in NextGen sequencing data



Application of HMMs in the area of SNP discovery from NextGen sequencing data, to greatly reduce false SNP calls caused by misalignments around insertions and deletions (indels). The central concept is per-Base Alignment Quality, which accurately measures the probability of a read base being wrongly aligned.

See: Li et al. (2011) *Bioinformatics* 27, 1157–1158

HMMER - protein homology detection and alignment



Profile HMM architecture used in HMMER2, SAM and PFTOOLS protein homology detection and alignment packages. *Match states* carry position-specific emission probabilities for scoring residues at each consensus position. *Insert states* emit residues with emission probabilities identical to a background distribution. We will describe this in more detail shortly...

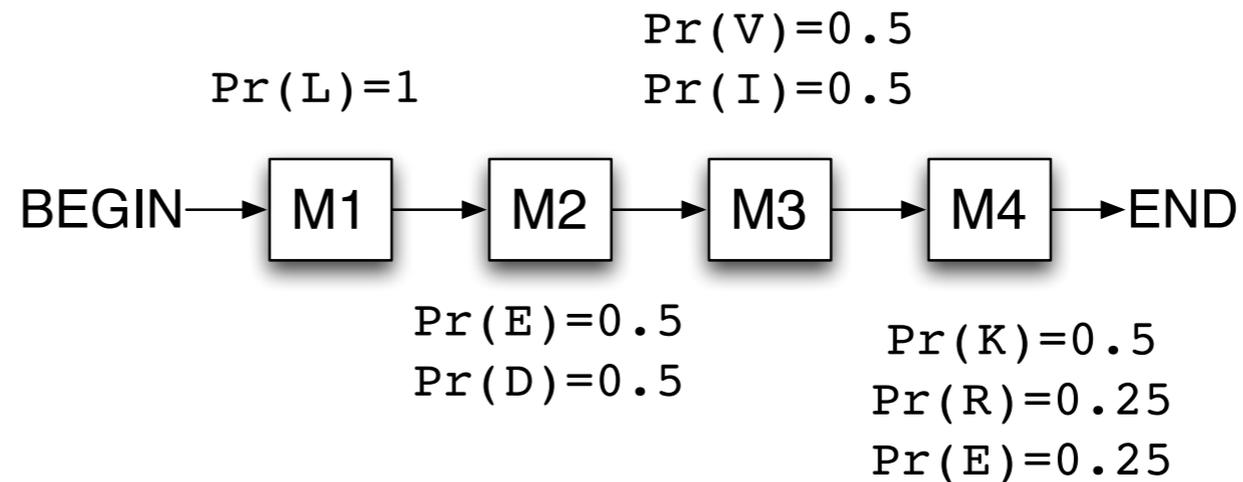
See: Eddy (1998) *Bioinformatics* 14, 755–763

Building sequence profile-HMMs: Match states

How do the above HMMs relate to profiles? Let's see how we can use the HMM framework to build **profile HMMs** that describe families of related sequences.

In the last lecture, we built a profile for the alignment:

s1 -
 LEVK
s2 LDIR
s3 LEIK
s4 LDVE



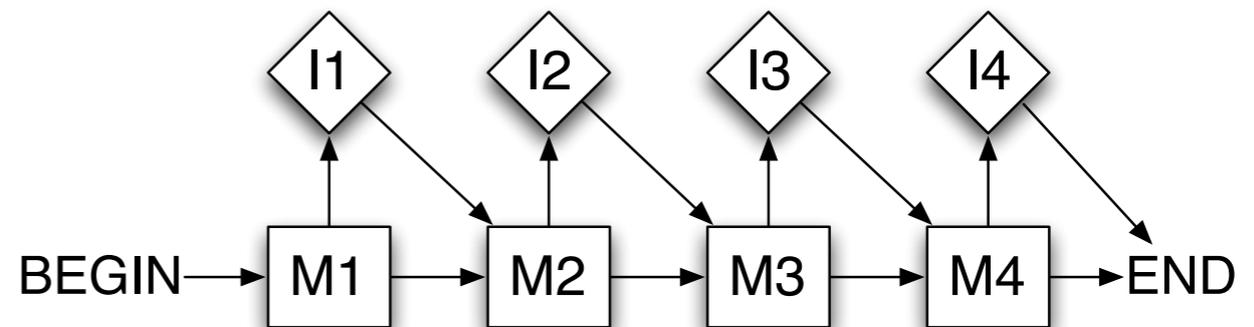
Ignoring the “background” frequencies for now, a profile for this alignment can be viewed as a simple HMM with one “match” state for each column, where consecutive match states are separated by transitions of probability 1.

Q. Why is this not a Markov chain?

Building profile-HMMs: Insert states

Introduce **insert states** (I_j), which will model inserts after the j th column in our alignment.

s1	LE-VK
s2	LD-IR
s3	LE- <u>IK</u>
s4	LD-VE
query1	LD <u>AV</u> K

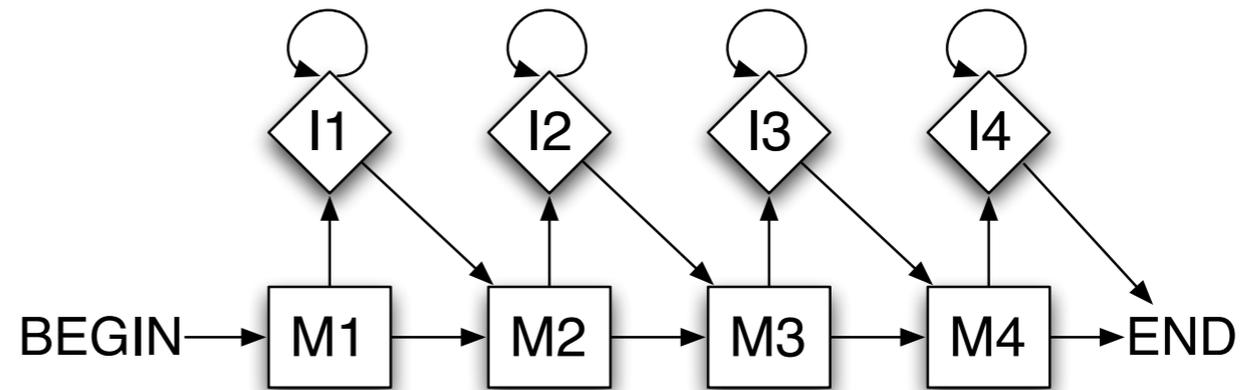


Typically, the output probabilities for insert states are set equal to the background probabilities. Note that we can have different probabilities for entering different insert states, and this models the fact that insertions may be less well-tolerated in certain portions of the alignment.

Building profile-HMMs: Insert states + *affine* gaps

For any particular insert state, we may have different transition probabilities for entering it for the first time vs. staying in the insert state; this models *affine* gap penalties.

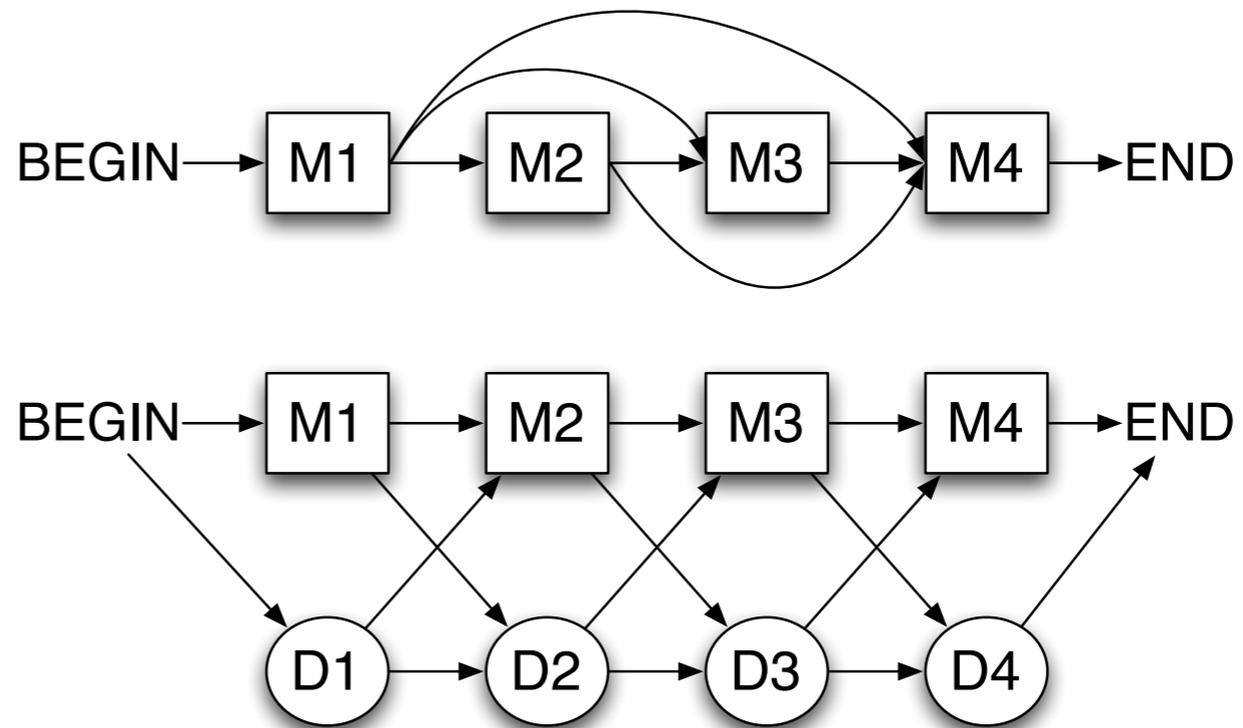
s1	LE	---	VK
s2	LD	---	IR
s3	LE	---	IK
s4	LD	---	VE
query1	LDA	--	VK
query2	LDAAV		VK



Building profile-HMMs: Delete states

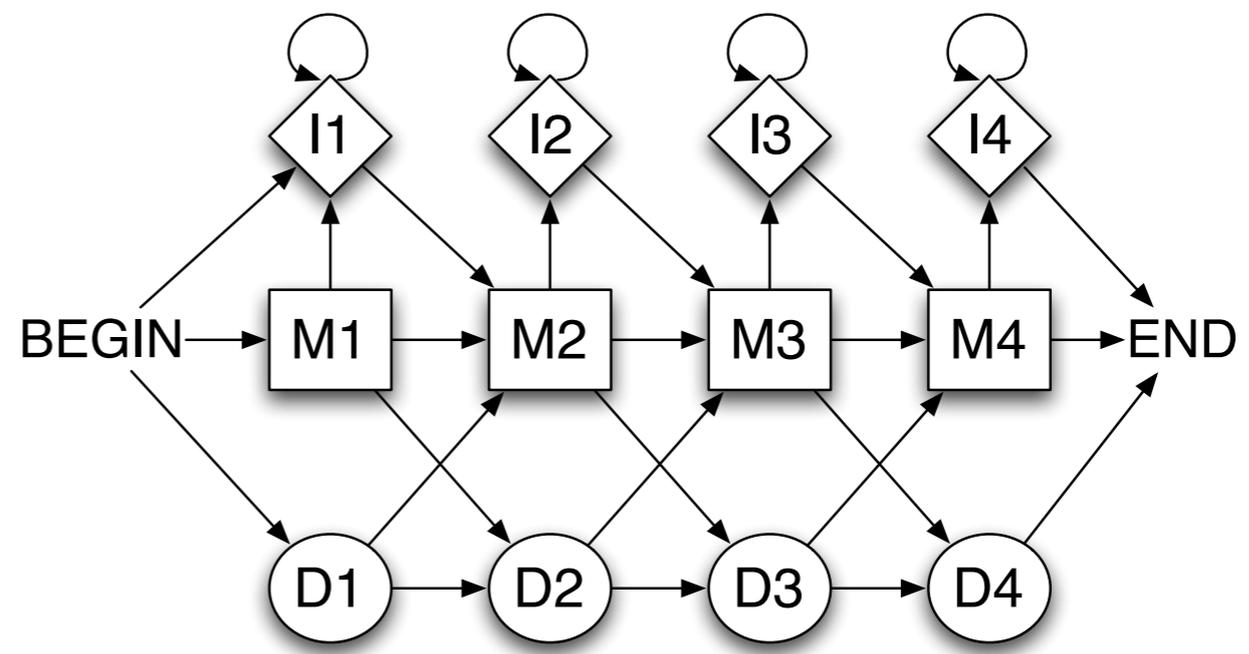
One could model deletions with additional transitions between match states. However, arbitrarily long gaps would introduce lots of transitions in the model. Instead, we will introduce **delete states** that do not emit any symbols

s1	LEVK
s2	LDIR
s3	LEIK
s4	LDVE
query3	L-VK



Building profile-HMMs

Putting it all together we get a complete profile HMM topology with match, insert and delete states.



However we still need to decide how many states our HMM has, what the transition probabilities are, etc.

Example profile-HMM building

- How do we pick the length of the HMM?
Common heuristic is to include only those columns that have > 50% occupancy
- How do we pick emission probabilities for match states?

$$b_{m1}(V) = 5/7$$

$$b_{m1}(F) = 1/7$$

$$b_{m1}(I) = 1/7$$

s1	VGA--NAGRPY
s2	VG---NVDKPV
s3	VGA--NVAHPH
s4	VAA-----PH
s5	VGS--TYEKPS
s6	FGA--NFEKPH
s7	IGAADNGARPY

How do we pick transition probabilities?

- We let the transition probability of going from state i to state j , a_{ij} be equal to:

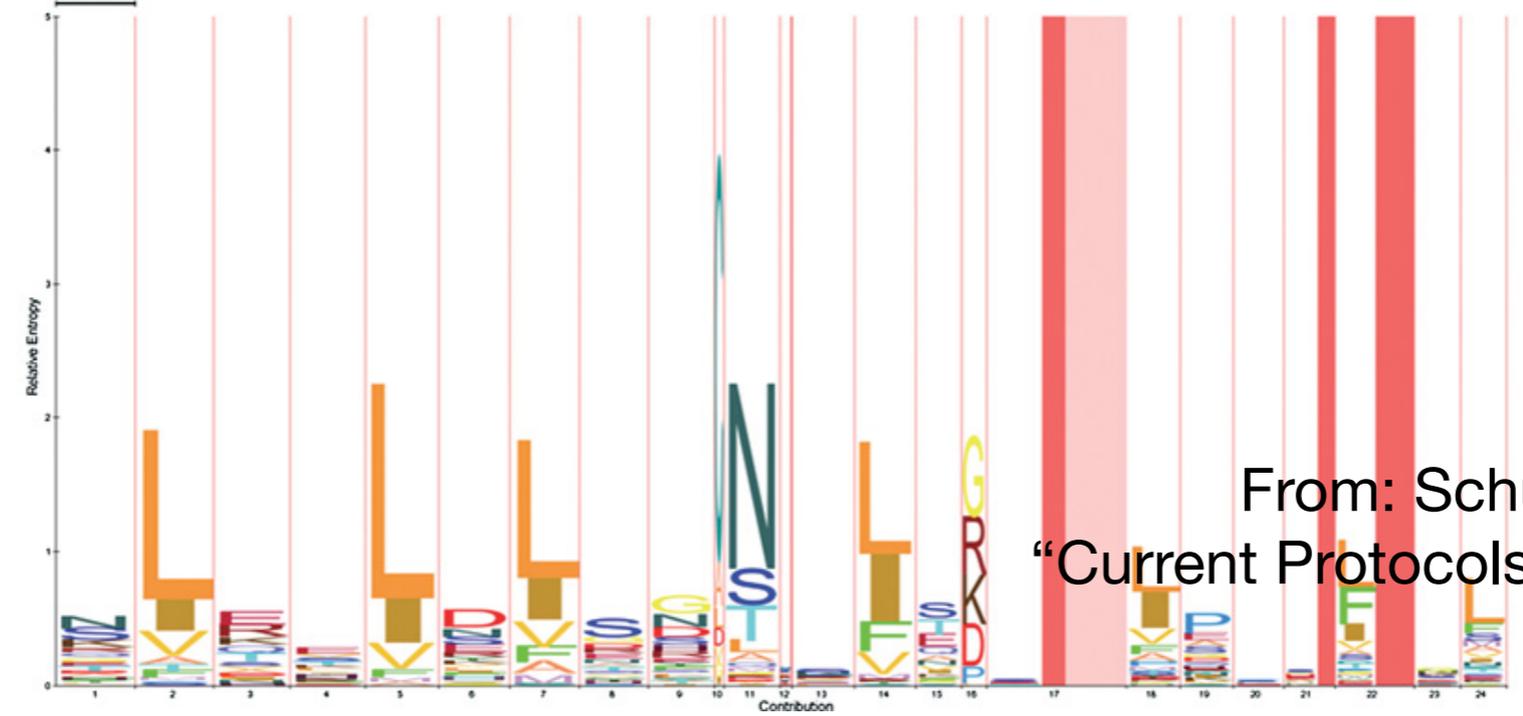
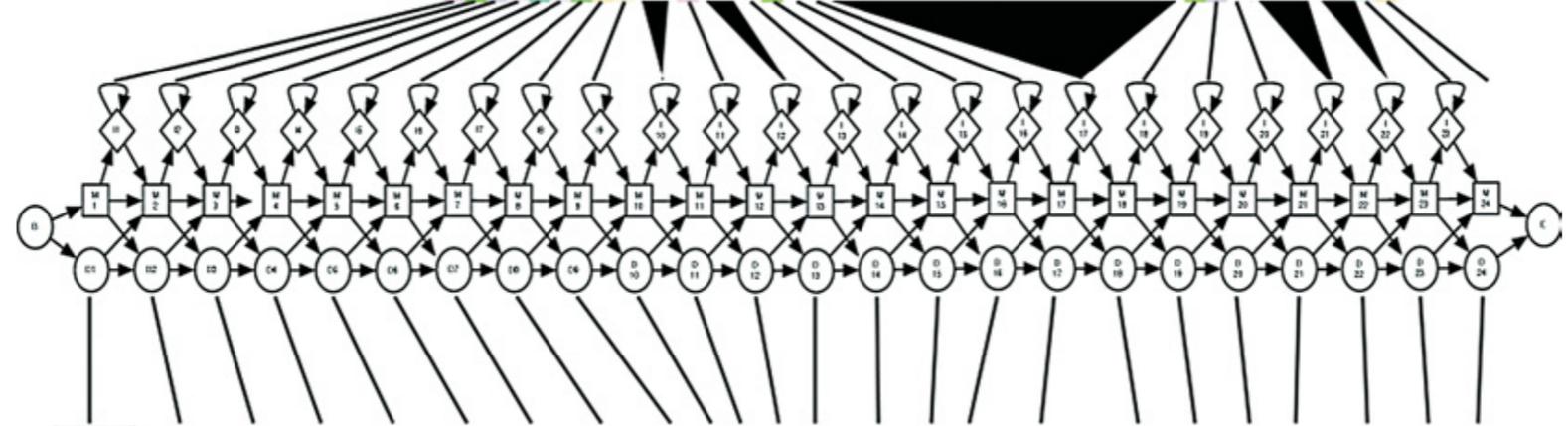
$$\frac{\text{No. of transitions from state } i \text{ to state } j}{\text{No. of transitions from state } i \text{ to any other state}}$$

$$a_{M2M3}(V) = 6/7 \quad \text{No. of matches (=6)}$$

$$a_{M2D3}(F) = 1/7 \quad \text{No. of gaps (=1)}$$

$$a_{M2I2}(I) = 0/7 \quad \text{No. of insertions (=0)}$$

Q9ARB2_LINUS/823-844	M L E Y L D I G R A . . P . R I V . H	L D G . . . L E N L
Q9M8N0_ARATH/320-341	R L T F L N L S F C . . S . K L T . G	L A F . . . F S I I
FLJ1_HUMAN/318-339	N L E E F M A A N . . N . N L E . L	V P E S . . L C R C
Q9VN74_DROME/90-112	A L H S L V I E N C T I V . H	I N D A A . F N Q E
Q8L8I7_PNTA/792-814	N L Q T I Q M Y R X . . E . S L Q . V	L P D S . . F G N L
Q9FHL8_ARATH/301-324	N L W S L N L S R . . N . . L F S D P	L P V V G . A R G F
SLIK6_MOUSE/65-87	R P F H L S L L N . . N . . G L T . M	L H T N D . F S G L
Q8NJJ8_EMENI/978-1000	T L T S L N I A S . . A . . K L V . Q	F R D T L . F D S L
Q9LUQ2_ARATH/92-113	A M K S L D V S F . . N . . S I S . E	L P E Q . . I G S A
Q9FH93_ARATH/169-188	R L T S L N L D F . . N . . R F N G T	L P S L N
Q898G0_CLOTE/268-288	Y L E R I N L D K . . N . . K I . K N	I E E . . . L E A N
Q8H6V2_MAIZE/678-699	N L R I L S I V D C . . V . S L Q . K	L P P . . . S D S F
Q9AR40_LINUS/692-713	D L K V L D I N Q . . T . . E I T . T	L K G E . . V E S L
Q9LE82_ARATH/350-377	H L T E I Y M S Y . . L . . N L E D E G T	E A L S E A L . L K S A
Q9H5N5_HUMAN/255-278	H L Q V L D L H Q C S L T . A D	D V M S L . . . T Q V I
Q8L4C7_ARATH/185-207	K L E Y L D I W G . . S . . N V T . N	Q G A V S . I L K F
Q9VSA4_DROME/1115-113E	Q L K A L R L Q C . . N . . A I . G S H	G L E A L . . L C G Q
TLR1_MOUSE/376-398	R L K T L S L Q K . . N . . Q L . K N	L E N I I . L T S A
Q9TXJ6_LEMA/445-465	G L R D I D L S H . . T . . K V H . N	I D A . . . L Q A S
FXL13_MOUSE/409-448	K L I Y L D L S G C . . T . Q V L . V E K C P R I S S V V L I G S P H I	S D S A . F K A L
Q9TXJ6_LEMA/927-948	A L T V V N A N S C . . V . N L T . S	I E A . . . L E S A
Q9M4X9_CHLRE/1417-1444	L L A V L H L H D . . N P . R L A . A D G	V A G L A A A . . L P G L
Q945S6_LYCPM/656-677	N L R H L D V S N . . T . . R L . K	M P L H . . L S R L



From: Schuster-Bockler *et al.*
 "Current Protocols in Bioinformatics"
 Supplement 18.

Side note: Weighting the training sequences

If there is a high degree of **redundancy** in our initial MSA (i.e. it contains a large group of very closely related sequences and a small number of more distantly related sequences) the resulting HMM will over represent the similar sequences and adversely effect our ability to detect distantly related sequences when searching databases

Sequences weighting attempts to compensate for this *sequence sampling bias* by differentially weighting sequences to reduce redundancy prior to model building

By default HMMER uses a sequence clustering tree as a guide to weight each sequence by its distance to other sequences. This approach will effectively down-weight the influence of redundant sequences.

A number of other approaches have been developed (Voronoi algorithm, maximum entropy, etc.)

See: Karchin *et al.* (1998) *Bioinformatics* 14, 772-778

Side note: Pseudocounts and Dirichlet distributions

Unfortunately, for alignments containing a small number of sequences the observed counts may not be representative of the family as a whole.

In such cases we must adjust the probabilities to account for our under-sampling (i.e. unobserved residues)

One common approach is to add **pseudocounts** to the observed counts so that no zero probabilities can occur.

Simplest approach is to just add one to all counts. More accurate adjustments consider prior knowledge about the behavior of sequence families adjusting counts according to pre-tabulated **Dirichlet distributions** - which are rather like protein comparison matrixes used in profile methods

Such information is often called **prior information**, indicating that it is known before any sequence data is seen

See: Durbin et al. "Biological Sequence Analysis"

Generating multiple sequence alignments

Large MSAs can be generated very quickly by using the Viterbi algorithm to find the most likely path through the HMM for a set of unaligned sequences

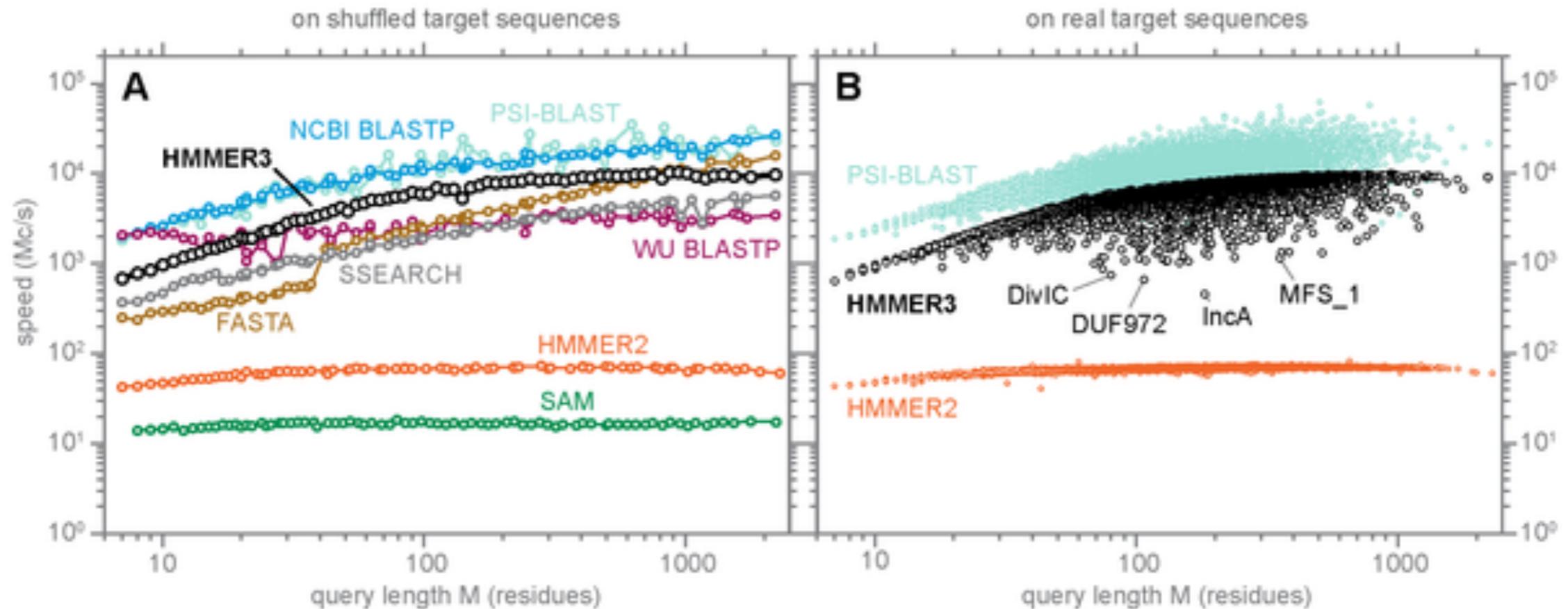
This is the basis of the PFAM database which uses the HMMER software package. Namely, HAMMER's *hmmalign* from the results of *hmmsearch*

MSA produced by HMMs are not true MSAs in the way that those produced by ClustalW are. ClustalW compares every sequence to every other sequence, whereas HMM aligning compares every sequence to the model independently so that the alignment between sequences is by proxy. Adding new sequences to the ClustalW alignment will add new information which may alter the alignment of existing sequences; adding new sequences to the HMM alignment never changes the alignment of any sequences relative to each other.

As an alternative to HMMER, you can use the *Sequence Alignment and Modeling Software System (SAM)*

<http://compbio.soe.ucsc.edu/sam.html>

HMM sequence searching performance

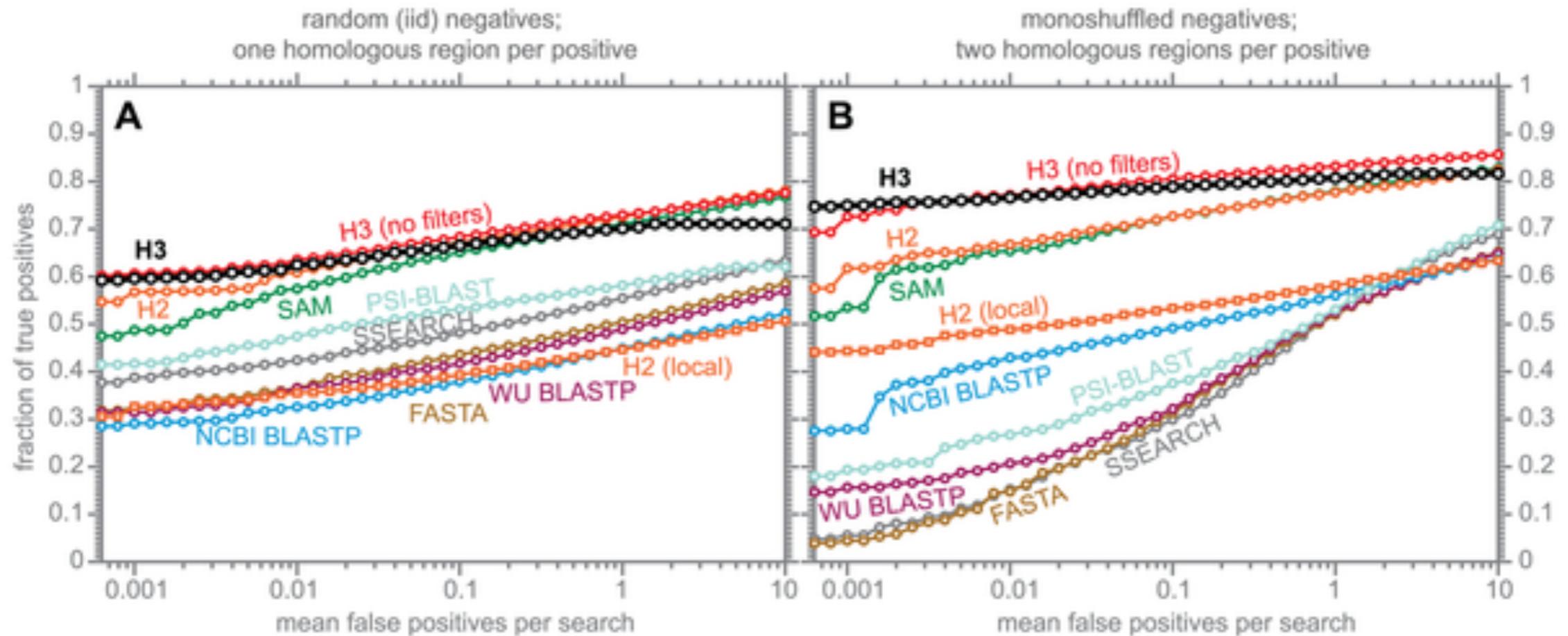


Recent speed benchmarks indicate that HMMER3 is approaching BLAST speed

Each point represents a speed measurement for one search with one query against target sequences. Both axes are logarithmic, for speed in millions of dynamic programming cells per second (Mc/s) on the y-axis and query length on the x-axis.

See: Eddy (2011) PLoS Comp Biol 7(10): e1002195

HMM sequence searching performance...



However HMMER3 has a much higher search sensitivity and specificity

In each benchmark, true positive subsequences have been selected to be no more than 25% identical to any sequence in the query alignment ... (see paper for details).

See: Eddy (2011) PLoS Comp Biol 7(10): e1002195

HMM limitations

HMMs are linear models and are thus **unable to capture higher order correlations** among positions (e.g. distant cysteins in a disulfide bridge, RNA secondary structure pairs, etc).

Another flaw of HMMs lies at the very heart of the mathematical theory behind these models. Namely, that the probability of a sequence can be found from the product of the probabilities of its individual residues.

This claim is only valid if the probability of a residue is independent of the probabilities of its neighbors. In biology, there are frequently **strong dependencies between these probabilities** (e.g. hydrophobic residues clustering at the core of protein domains).

These biological realities have motivated research into new kinds of statistical models. These include hybrids of HMMs and neural nets, dynamic Bayesian nets, factorial HMMs, Boltzmann trees and stochastic context-free grammars.

See: Durbin et al. "Biological Sequence Analysis"

PFAM: Protein Family Database of Profile HMMs

Comprehensive compilation of both multiple sequence alignments and profile HMMs of protein families.

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

PFAM consists of two databases:

- **Pfam-A** is a manually curated collection of protein families in the form of multiple sequence alignments and profile HMMs. HMMER software is used to perform searches.
- **Pfam-B** contains additional protein sequences that are automatically aligned. Pfam-B serves as a useful supplement that makes the database more comprehensive.
- Pfam-A also contains higher-level groupings of related families, known as **clans**

Pfam 25.0 (March 2011, 12273 families)

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

QUICK LINKS

[SEQUENCE SEARCH](#)

[VIEW A PFAM FAMILY](#)

[VIEW A CLAN](#)

[VIEW A SEQUENCE](#)

[VIEW A STRUCTURE](#)

[KEYWORD SEARCH](#)

[JUMP TO](#)

YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...

Analyze your protein sequence for Pfam matches

View Pfam family annotation and alignments

See groups of related families

Look at the domain organisation of a protein sequence

Find the domains on a PDB structure

Query Pfam by keywords

Enter any type of accession or ID to jump to the page for a Pfam family or clan, UniProt sequence, PDB structure, etc.

Or view the [help](#) pages for more information

Citing Pfam

If you find Pfam useful, please consider [citing](#) the reference that describes this work:

[The Pfam protein families database](#): R.D. Finn, J. Mistry, J. Tate, P. Coghill, A. Heger, J.E. Pollington, O.L. Gavin, P. Gunesekaran, G. Ceric, K. Forslund, L. Holm, E.L. Sonnhammer, S.R. Eddy, A. Bateman

Nucleic Acids Research (2010) Database Issue 38:D211-222

Mirrors

The following are official Pfam [mirror](#) sites:

 [WTSI, UK](#)

 [SBC, Sweden](#)

 [JFRC, USA](#)

Family: *Kinesin* (PF00225)

 126 architectures
  4150 sequences
  6 interactions
  248 species
  114 structures

Summary

[Domain organisation](#)

[Clans](#)

[Alignments](#)

[HMM logo](#)

[Trees](#)

[Curation & models](#)

[Species](#)

[Interactions](#)

[Structures](#)

Jump to... 

 enter ID/acc

Summary

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

Wikipedia: Kinesin

Pfam

Interpro

The Pfam group coordinates the annotation of Pfam families in [Wikipedia](#). This family is described by a Wikipedia entry entitled "[Kinesin](#)". [More...](#)

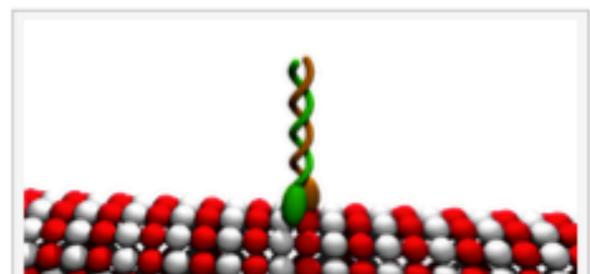
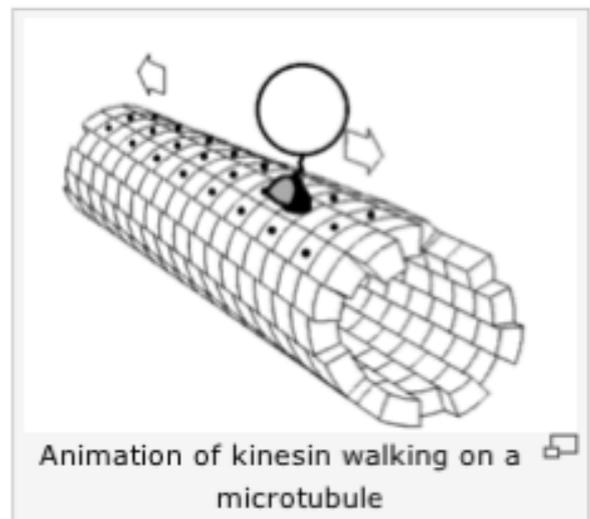
Kinesin

[Edit Wikipedia article](#)

A **kinesin** is a protein belonging to a class of [motor proteins](#) found in [eukaryotic](#) cells. Kinesins move along [microtubule](#) filaments, and are powered by the hydrolysis of [ATP](#) (thus kinesins are [ATPases](#)). The active movement of kinesins supports several cellular functions including [mitosis](#), [meiosis](#) and transport of cellular cargo, such as in [axonal transport](#). Most kinesins walk towards the plus end of a microtubule, which, in most cells, entails transporting cargo from the centre of the cell towards the periphery. This form of transport is known as [anterograde transport](#).

Contents [\[show\]](#)

- 1 Structure
 - 1.1 Overall structure
 - 1.2 Kinesin motor domain
- 2 Cargo transport
- 3 Direction of motion
- 4 Proposed mechanisms of movement
- 5 Theoretical Modeling of Kinesin
- 6 Kinesin and mitosis
- 7 Family members
- 8 See also
- 9 References
- 10 External links



Family: *Kinesin* (PF00225)

126 architectures 4150 sequences 6 interactions 248 species 114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

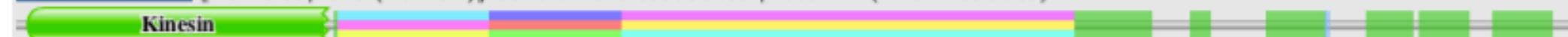
Jump to...

Domain organisation

Below is a listing of the unique domain organisations or architectures in which this domain is found. [More...](#)

There are 3185 sequences with the following architecture: **Kinesin**

[CENPE_HUMAN](#) [Homo sapiens (Human)] Centromere-associated protein E (2701 residues)



[Show](#) all sequences with this architecture.

There are 139 sequences with the following architecture: **Kinesin x 2**

[CIN8_YEAST](#) [Saccharomyces cerevisiae (Baker's yeast)] Kinesin-like protein CIN8 (1000 residues)



[Show](#) all sequences with this architecture.

There are 56 sequences with the following architecture: **Kinesin, FHA**

[KIF14_HUMAN](#) [Homo sapiens (Human)] Kinesin-like protein KIF14 (1648 residues)



[Show](#) all sequences with this architecture.

There are 54 sequences with the following architecture: **CH, Kinesin**

[Q9SS42_ARATH](#) [Arabidopsis thaliana (Mouse-ear cress)] Kinesin-like protein (897 residues)



[Show](#) all sequences with this architecture.

There are 54 sequences with the following architecture: **Kinesin, DUF3490**

[Q8LNZ2_ARATH](#) [Arabidopsis thaliana (Mouse-ear cress)] Kinesin-like protein (938 residues)



[Show](#) all sequences with this architecture.

There are 44 sequences with the following architecture: **Kinesin, FHA, KIF1B, DUF3694, PH**

[KIF1A_HUMAN](#) [Homo sapiens (Human)] Kinesin-like protein KIF1A (1690 residues)



[Show](#) all sequences with this architecture.

Family: *Kinesin* (PF00225)

 126 architectures
  4150 sequences
  6 interactions
  248 species
  114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Jump to...

Pfam Clans

This family is a member of clan [AAA \(CL0023\)](#), which contains the following 157 members:

6PF2K	AAA	AAA-ATPase like	AAA_10	AAA_2
AAA_3	AAA_5	AAA_6	AAA_7	AAA_8
AAA_9	AAA PrkA	ABC ATPase	ABC tran	ABC tran_2
Adeno IVa2	Adenylsucc synt	ADK	AFG1 ATPase	AIG1
APS kinase	Arch ATPase	Arf	ArgK	ArsA ATPase
ATP-synt ab	ATP bind 1	ATP bind 2	Bac DnaA	CbiA
CoaE	CobA CobO BtuR	CobU	cobW	CPT
CTP synth_N	Cytidylate kin	DAP3	DEAD	DEAD_2
DLIC	DNA pack_C	DNA pack_N	DNA_pol3_delta	DnaB_C
dNK	DUF1253	DUF1611	DUF2075	DUF2478
DUF258	DUF2813	DUF463	DUF699	DUF815
DUF853	DUF87	DUF927	Dynamin_N	Exonuc_V_gamma
FeoB_N	Fer4 NifH	Flavi DEAD	FTHFS	FtsK_SpoIIIE
G-alpha	Gal-3-0 sulfotr	GBP	GSPII_E	GTP EFTU
Gtr1 RagA	Guanylate kin	GvpD	HDA2-3	Helicase_C
Herpes Helicase	Herpes ori bp	Herpes TK	IIGP	IPPT
IPT	IstB IS21	KaiC	KAP_NTPase	Kinesin
Kinesin-relat_1	Kinesin-related	KTI12	LpxK	MCM
Mg chelatase	MipZ	Miro	MMR_HSR1	MobB
MukB	MutS_V	Myosin head	NACHT	NB-ARC
NOG1	NTPase_1	ParA	Parvo NS1	PAXNEB
PduV-EutP	PhoH	PIF1	Podovirus Gp16	Polyoma Ig T_C
Pox_A32	PPK2	PPV_E1_C	PRK	Rad17
Rad51	Ras	RecA	Rep fac_C	ResIII
RHD3	RHSP	RNA12	RNA_helicase	RuvB_N
SecA DEAD	Septin	Sigma54 activat	SKI	SMC_N
SNF2_N	Spore IV_A	SRP54	SRPRB	Sulfotransfer_1
Sulfotransfer_2	Sulphotransf	T4SS-DNA transf	Terminase_1	Terminase_3
Terminase_6	Terminase GpA	Thymidylate kin	TIP49	TK
TniB	Torsin	TraG-D_C	TrwB AAD bind	UPF0079
UvrD-helicase	Viral helicase1	VirC1	VirE	YhjQ
Zeta toxin	Zot			

Family: *Kinesin* (PF00225)

126 architectures 4150 sequences 6 interactions 248 species 114 structures

- Summary
- Domain organisation
- Clans

Alignments

- HMM logo
- Trees
- Curation & models
- Species
- Interactions
- Structures

Jump to...

enter ID/acc Go

Alignments

There are various ways to view or download the sequence alignments that we store. You can use a sequence viewer to look at either the seed or full alignment for the family, or you can look at a plain text version of the sequence in a variety of different formats. [More...](#)

View options

Alignment:	<input checked="" type="radio"/> Seed (87)	<input type="radio"/> Full (4150)
	<input type="radio"/> NCBI (6110)	<input type="radio"/> Metagenomics (525)
Viewer:	HTML	

View

Formatting options

Alignment:	<input checked="" type="radio"/> Seed (87)	<input type="radio"/> Full (4150)
Format:	Selex	
Order:	<input checked="" type="radio"/> Tree	<input type="radio"/> Alphabetical
Sequence:	<input checked="" type="radio"/> Inserts lower case	<input type="radio"/> All upper case
Gaps:	Gaps as "." or "-" (mixed)	
Download/view:	<input checked="" type="radio"/> Download	<input type="radio"/> View

Generate

Download options

Very large alignments can often cause problems for the formatting tool above. If you find that downloading or viewing a large alignment is problematic, you can also download a [gzip](#)-compressed, Stockholm-format file containing the [seed](#) or [full](#) alignment for this family.

You can also [download](#) a FASTA format file containing the **full-length sequences** for all sequences in the full alignment.

The main seed and full alignments are generated using sequences from the UniProt sequence database. However, we also generate

Table of protein sequence alignments for Kinesin (PF00225). Columns include protein name (e.g., KIF7_DICDI/34-349), seed sequence (e.g., DSKSISIRANGPQFTTDRIFGY..QET), and Pfam seed sequence (e.g., QSQIFEDV.AEPIVNDFL.DGYHGTIIAYG.QTAS). Sequences are color-coded by conservation.

Family: *Kinesin* (PF00225)

 126 architectures
  4150 sequences
  6 interactions
  248 species
  114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

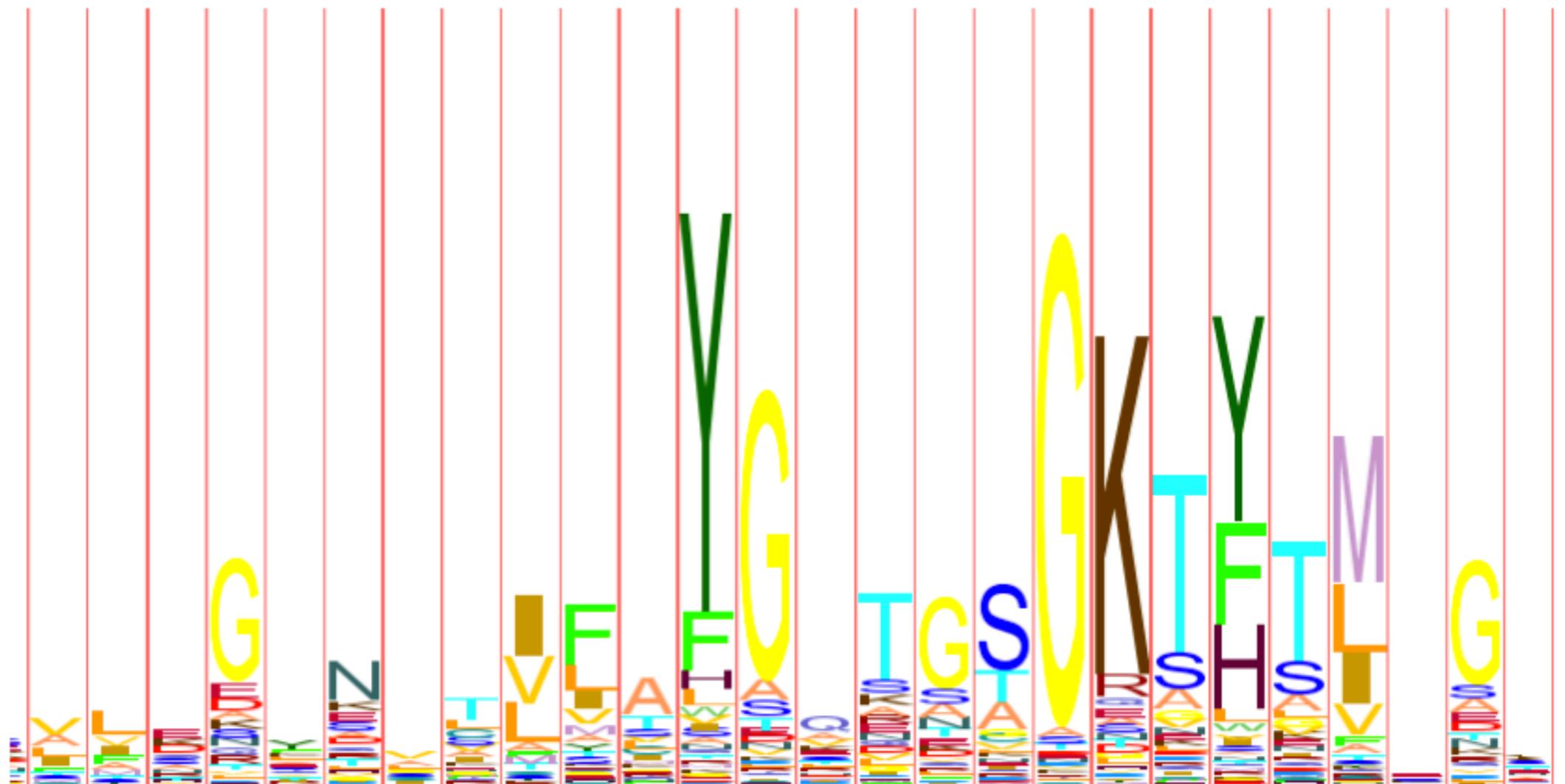
Structures

Jump to... ↓

enter ID/acc Go

HMM logo

HMM logos is one way of visualising profile HMMs. Logos provide a quick overview of the properties of an HMM in a graphical form. You can see a more detailed description of HMM logos and find out how you can interpret them [here](#). **More...**



Family: *Kinesin* (PF00225)

126 architectures 4150 sequences 6 interactions 248 species 114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Jump to... 

enter ID/acc

Go

Curation and family details

This section shows the detailed information about the Pfam family. You can see the definitions of many of the terms in this section in the [glossary](#) and a fuller explanation of the scoring system that we use in the [scores](#) section of the help pages.

Curation

Seed source:	Prosite
Previous IDs:	kinesin;
Type:	Domain
Author:	Bateman A, Finn RD
Number in seed:	87
Number in full:	4150
Average length of the domain:	298.60 aa
Average identity of full alignment:	31 %
Average coverage of the sequence by the domain:	34.30 %

HMM information

HMM build commands:	<i>build method:</i> hmmbuild -o /dev/null HMM SEED <i>search method:</i> hmmsearch -Z 11384036 -E 1000 --cpu 4 HMM pfamseq		
Model details:	Parameter	Sequence	Domain
	Gathering cut-off	22.5	22.5
	Trusted cut-off	22.5	22.5
	Noise cut-off	22.4	22.4
Model length:	333		
Family (HMM) version:	17		



Family: *Kinesin* (PF00225)

Loading page components (1 remaining)...

 126 architectures
  4150 sequences
  6 interactions
  248 species
  114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Jump to...

enter ID/acc

Go

Interactions

There are 6 interactions for this family. [More...](#)

[Tubulin](#)
[Tubulin_C](#)

[Tubulin_C](#)

[Kinesin](#)

[Tubulin](#)

[Kinesin](#)



Family: *Kinesin* (PF00225)

126 architectures
 4150 sequences
 6 interactions
 248 species
 114 structures

Summary

Domain organisation

Clans

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Jump to...

enter ID/acc

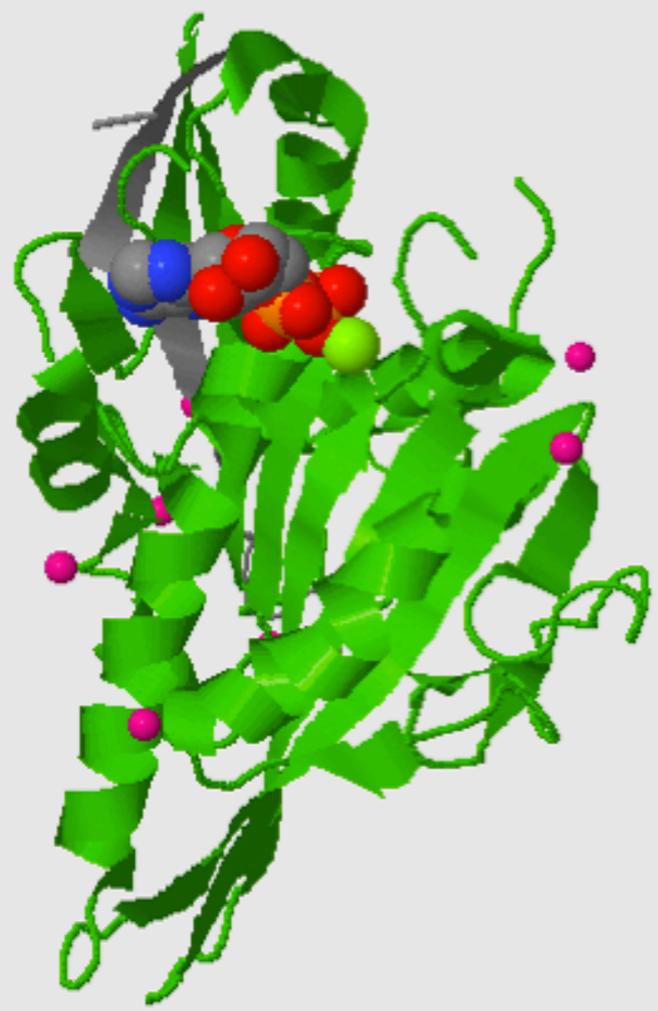
Structures

For those sequences which have a structure in the [Protein DataBank](#), we use the mapping between [UniProt](#), PDB and Pfam coordinate systems from the [PDBer](#) group, to allow us to map Pfam domains onto UniProt sequences and three-dimensional protein structures. The table below shows the structures on which the **Kinesin** domain has been found.

UniProt entry	UniProt residues	PDB ID	PDB chain ID	PDB residues	View			
A8BKD1_GIALA	11 - 335	2vvg	A	11 - 335	Jmol AstexViewer SPICE			
			B	11 - 335	Jmol AstexViewer SPICE			
CENPE_HUMAN	12 - 329	1t5c	A	12 - 329	Jmol AstexViewer SPICE			
			B	12 - 329	Jmol AstexViewer SPICE			
KAR3_YEAST	392 - 723		1f9t	A	392 - 723	Jmol AstexViewer SPICE		
			1f9u	A	392 - 723	Jmol AstexViewer SPICE		
			1f9v	A	392 - 723	Jmol AstexViewer SPICE		
			1f9w	A	392 - 723	Jmol AstexViewer SPICE		
				B	392 - 723	Jmol AstexViewer SPICE		
			3kar	A	392 - 723	Jmol AstexViewer SPICE		
KI13B_HUMAN	11 - 352	3qbj	A	11 - 352	Jmol AstexViewer SPICE			
			B	11 - 352	Jmol AstexViewer SPICE			
			C	11 - 352	Jmol AstexViewer SPICE			
			1ii6	A	24 - 359	Jmol AstexViewer SPICE		
				B	24 - 359	Jmol AstexViewer SPICE		
			1q0b	A	24 - 359	Jmol AstexViewer SPICE		
				B	24 - 359	Jmol AstexViewer SPICE		
			1x88	A	24 - 359	Jmol AstexViewer SPICE		
				B	24 - 359	Jmol AstexViewer SPICE		
						A	24 - 359	Jmol AstexViewer SPICE



PDB entry 3bfn



Jmol

PDB			UniProt			Pfam family	Colour
Chain	Start	End	ID	Start	End		
A	49	368	KIF22_HUMAN	49	368	Kinesin (PF00225)	

Close window



HMNER

biosequence analysis using profile hidden Markov models



Home Search Results Software Help About

HMNER3: a new generation of sequence homology search software

HMNER is used for searching sequence databases for homologs of protein sequences, and for making protein sequence alignments. It implements methods using probabilistic models called **profile hidden Markov models** (profile HMMs).

Compared to BLAST, FASTA, and other sequence alignment and database search tools based on older scoring methodology, HMNER aims to be significantly *more* accurate and *more* able to detect remote homologs because of the strength of its underlying mathematical models. In the past, this strength came at significant computational expense, but in the new HMNER3 project, HMNER is now essentially **as fast as** BLAST.

As part of this evolution in the HMNER software, we are committed to making the software available to as many scientists as possible. Earlier releases of HMNER were restricted to command line use. To make the software more accessible to the wide scientific community, we now provide **servers** that allow **sequence searches** to be performed interactively via the **Web**.

The current version is **HMNER 3.0** (28 March 2010) and can be **downloaded** from the software section of the site. Previous versions of the HMNER software can be obtained from the **archive** section.

If you have used the HMNER website, please consider citing the following reference that describes this work:

HMNER web server: interactive sequence similarity searching
R.D. Finn, J. Clements, S.R. Eddy
Nucleic Acids Research (2011) Web Server Issue 39:W29-W37. [PDF](#)

Download HMNER

Get the latest version

v3.0

[Release notes](#) (28 March 2010)



[Alternative Download Options](#)
[Source](#)

Search

Perform an interactive search now.
[Search](#)

Comments or questions on the site? Send a mail to hmner@janelia.hhmi.org
Howard Hughes Medical Institute

Follow @hmm3r



HMMER

biosequence analysis using profile hidden Markov models

[Home](#) [Search](#) [Results](#) [Software](#) [Help](#) [About](#)

[phmmer](#) [hmmscan](#) [hmmsearch](#)



protein sequence vs protein sequence database

[Advanced](#)

Paste in your sequence or use the [example](#)

```
>sp|Q14807|KIF22_HUMAN
MAAGGSTQRRREMAAASAAAISGAGRCRLSKIGATRRPPPARVRVAVRLRPFVDGTAGA
SDPPCVRGMDSCSLEIANWRNHQETLKYQFDAFYGERSTQQDIYAGSVQPILRHLEGN
ASVLAYGPTGAGKTHMLGSPEQPGVIPRALMDLLQLTREEGAEGRPWALSVTMSYLEIY
QEKVLDLLDPASGDLVIREDCRGNILIPGLSQKPISSFADFERHFLPASRNRTVGATRLN
QRSSRSHAVLLVKVDQERERLAPFRQREGKLYLIDLAGESEDNRRTGNKGLRLKESGAIN
TFVLGKVVDAQNQGLPRVPYRDSKLRLLQDSLGGSAHSILIANIAPERRFYLDTVSALN
FAARSKEVINRPFTNESLQPHALGPVKLSQKELLGPPEAKRARGPEEEEEIGSPEPMAAPA
SASQKLSPLQKLSMDPAMLERLLSLDRLLASQGSQGAPLLSTPKRERMVLMKTVEEKDL
EIERLTKQKELEAKMLAQKAEKENHCPTMLRPLSHRTVTGAKPLKAVVMPLQLIQEQ
AASPNAEIHILKNKGRKRKLESLEDALEPEEKAEDCWELQISPELLAHGRQKILDLLNEGS
ARDLRSIORICPKKAOLIVGWRELFHGPESOVEDLERVEGITGKOMESELKANILGLAAGO
```

Submit

[Reset](#)

Comments or questions on the site? Send a mail to hmmer@janelia.hhmi.org
Howard Hughes Medical Institute

[Follow](#) @hmm3r



HMNER

biosequence analysis using profile hidden Markov models



Home Search **Results** Software Help About

phmmer

[Search Again](#)

Score Taxonomy Domain Download

Pfam Domains



[Show hit details](#)

Distribution of Significant Hits



« First « Previous **Page 1** of 51 Next » Last »

Query Matches (5100)

[Customize](#)

Target	Description	Species	E-value	Alignments (show all)
123979736	kinesin family member 22	synthetic construct	0.0e+00	show
6453818	kinesin-like protein KIF22	Homo sapiens	0.0e+00	show
30584615	Homo sapiens kinesin-like 4	synthetic construct	0.0e+00	show
123994513	kinesin family member 22	synthetic construct	0.0e+00	show
189053342	unnamed protein product	Homo sapiens	0.0e+00	show
62898423	kinesin family member 22 variant	Homo sapiens	0.0e+00	show
332845643	PREDICTED: kinesin family member 22 isoform 2	Pan troglodytes	0.0e+00	show
75062021	RecName: Full=Kinesin-like protein KIF22	Pongo abelii	0.0e+00	show
332266048	PREDICTED: kinesin-like protein KIF22-like isoform 1	Nomascus leucogenys	0.0e+00	show
297283748	PREDICTED: hypothetical protein LOC706401 isoform 3	Macaca mulatta	0.0e+00	show
296219941	PREDICTED: LOW QUALITY PROTEIN: kinesin-like protein KIF22-like	Callithrix jacchus	0.0e+00	show
296196456	PREDICTED: kinesin-like protein KIF22-like	Callithrix jacchus	0.0e+00	show
335284407	PREDICTED: kinesin-like protein KIF22-like	Sus scrofa	0.0e+00	show
221046166	unnamed protein product	Homo sapiens	0.0e+00	show
221045488	unnamed protein product	Homo sapiens	0.0e+00	show

[Score](#) [Taxonomy](#) **[Domain](#)** [Download](#)

Query



Jump to the exact match for your query architecture

« First « Previous **Page 1** of 7 Next » Last »

Domain Architectures

- | | | |
|--------------------------|--|-----------------------------|
| 3624
SEQUENCES | with domain architecture: Kinesin , <i>example:148685550</i> | View Scores |
| Show All | | |
| 126
SEQUENCES | with domain architecture: Kinesin, FHA , <i>example:157125836</i> | View Scores |
| Show All | | |
| 101
SEQUENCES | with domain architecture: Kinesin, Kinesin , <i>example:296088325</i> | View Scores |
| Show All | | |
| 80
SEQUENCES | with domain architecture: Kinesin, FHA, KIF1B, DUF3694, PH , <i>example:118101106</i> | View Scores |
| Show All | | |
| 69
SEQUENCES | with domain architecture: HHH_3 , <i>example:337289058</i> | View Scores |
| Show All | | |
| 62
SEQUENCES | with domain architecture: CH, Kinesin , <i>example:224061629</i> | View Scores |
| Show All | | |
| 60
SEQUENCES | Exact match with query architecture: Kinesin, HHH_3 , <i>example:332266048</i> | View Scores |
| Show All | | |



HMMER

biosequence analysis using profile hidden Markov models

[Home](#)
[Search](#)
[Results](#)
[Software](#)
[Help](#)
[About](#)

phmmer



[Score](#)
[Taxonomy](#)
[Domain](#)
[Download](#)

- **Job:** 9924F9AC-FEB5-11E0-A304-2B0C998A7913
- **Started:** 2011-10-24 23:01:15
- **Algorithm:** phmmer
- **HMMER Options:** -E 1 --domE 1 --incE 0.01 --incdomE 0.03 --mx BLOSUM62 --pextend 0.4 --popen 0.02 --seqdb nr

▼ Format

FASTA

Download the significant hits from your search as a gzipped FASTA file.



Full length FASTA

A gzipped file containing the full length sequences for significant search hits.



Aligned FASTA

A gzipped file containing aligned significant search hits in FASTA format.



STOCKHOLM

Download an alignment of significant hits as a gzipped STOCKHOLM file.



Text

A plain text file containing the hit alignments and scores.



XML

An XML file formatted for machine parsing of the data.



JSON

All the results information encoded as a single json string.



HMM

Profile HMM downloads are not available.

Download

[Reset](#)

Superfamily 1.75

HMM library and genome assignments server

 Search SUPERFAMILY

Home

SEARCH

[Keyword search](#)

[Sequence search](#)

BROWSE

Organisms

└─ [Taxonomy](#)

└─ [Statistics](#)

SCOP

└─ [Hierarchy](#)

Ontologies

└─ [GO](#)

└─ [EC](#)

└─ [Phenotype](#)

TOOLS

[Compare genomes](#)

[Phylogenetic trees](#)

[Web services](#)

[Downloads](#)

ABOUT

[Description](#)

[Publications](#)

[Documentation](#)

SUPERFAMILY Description

SUPERFAMILY is a database of structural and functional annotation for all proteins and genomes.

The SUPERFAMILY annotation is based on a collection of **hidden Markov models**, which represent structural protein domains at the [SCOP](#) superfamily level. A superfamily groups together domains which have an evolutionary relationship. The annotation is produced by scanning protein sequences from over **1,700 completely sequenced genomes** against the hidden Markov models.

For each **protein** you can:

- Submit sequences for [SCOP classification](#)
- View domain organisation, sequence alignments and protein sequence details

For each **genome** you can:

- Examine superfamily assignments, phylogenetic trees, domain organisation lists and networks
- Check for over- and under-represented superfamilies within a genome

For each **superfamily** you can:

- Inspect SCOP classification, functional annotation, Gene Ontology annotation, InterPro abstract and genome assignments
- Explore taxonomic distribution of a superfamily across the tree of life

All annotation, models and the database dump are freely available for [download](#) to everyone.

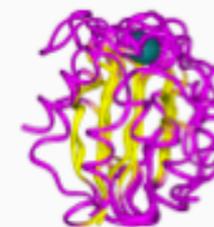
[Description cont.](#)

Jump to [[SUPERFAMILY description](#) · [Recent news](#)]

Major Features

Superfamily 1.75

HMM library and genome assignments server


 Search SUPERFAMILY

[Home](#) > [Keyword search](#) > kinesin

SEARCH

[Keyword search](#)

[Sequence search](#)

BROWSE

Organisms

[Taxonomy](#)

[Statistics](#)

SCOP

[Hierarchy](#)

Ontologies

[GO](#)

[EC](#)

[Phenotype](#)

TOOLS

[Compare genomes](#)

[Phylogenetic trees](#)

[Web services](#)

[Downloads](#)

ABOUT

[Description](#)

[Publications](#)

[Documentation](#)

Keyword Search Results

Another Search:

Results 1-3 of 3 for **kinesin**.

1.

SCOP classification

Class : [Alpha and beta proteins \(a/b\)](#)

Fold : [P-loop containing nucleoside triphosphate hydrolases](#)

Superfamily : [P-loop containing nucleoside triphosphate hydrolases](#)

Family : [Motor proteins](#)

Protein : [Kinesin](#)

Protein : [Kinesin motor Ncd \(non-claret disjunctional\)](#)

Protein : [Kinesin heavy chain-like protein](#)

Superfamily

[Alignments](#)

[Genome assignments](#)

[Taxonomic distribution](#)

[Domain combinations](#)

Family

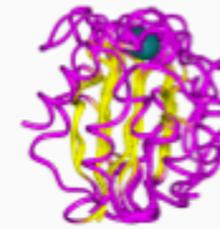
[Alignments](#)

[Genome assignments](#)

[Taxonomic distribution](#)

Superfamily 1.75

HMM library and genome assignments server


 Search SUPERFAMILY

[Home](#) > [SCOP hierarchy](#) > Motor proteins family

SEARCH

[Keyword search](#)
[Sequence search](#)

[Structural Classification](#)

[Genome Assignments](#)

[Sequence Alignments](#)

[Taxonomic Distribution](#)

BROWSE

Organisms
 - [Taxonomy](#)
 - [Statistics](#)
 SCOP
 - [Hierarchy](#)
 Ontologies
 - [GO](#)
 - [EC](#)
 - [Phenotype](#)

Motor proteins family

SCOP classification

Root: [SCOP hierarchy in SUPERFAMILY \[SCOP_0\]](#) (11)
 Class: [Alpha and beta proteins \(a/b\) \[SCOP_51349\]](#) (147)
 Fold: [P-loop containing nucleoside triphosphate hydrolases \[SCOP_52539\]](#)
 Superfamily: [P-loop containing nucleoside triphosphate hydrolases \[SCOP_52540\]](#) (24)
Family: [Motor proteins \[SCOP_52641\]](#) (4)

Family statistics

	Genomes (351)	Uniprot 2011_09	PDB chains (SCOP 1.75)
Domains	15,002	10,025	46
Proteins	14,811	9,977	46

TOOLS

[Compare genomes](#)
[Phylogenetic trees](#)
[Web services](#)
[Downloads](#)

Gene Ontology (high-coverage)

[\(show details\)](#)

ABOUT

[Description](#)
[Publications](#)
[Documentation](#)

	GO term	FDR (all)	SDFO lev
Biological Process (BP)	multicellular organismal process	0	<i>Least Inf</i>
Biological Process (BP)	biological regulation	0.03575	<i>Least Inf</i>

SEARCH

[Keyword search](#)
[Sequence search](#)

**Structural
Classification**

**Genome
Assignments**

**Sequence
Alignments**

**Domain
Combinations**

**Taxonomic
Distribution**

BROWSE

Organisms

[Taxonomy](#)

[Statistics](#)

SCOP

[Hierarchy](#)

Ontologies

[GO](#)

[EC](#)

[Phenotype](#)

TOOLS

[Compare genomes](#)

[Phylogenetic trees](#)

[Web services](#)

[Downloads](#)

ABOUT

[Description](#)

[Publications](#)

[Documentation](#)

HELP

[User support](#)

[Contact us](#)

[Email list](#)

[Sitemap](#)

P-loop containing nucleoside triphosphate hydrolases superfamily**SCOP classification**

Root: [SCOP hierarchy in SUPERFAMILY \[SCOP_0\]](#) (11)

Class: [Alpha and beta proteins \(a/b\) \[SCOP_51349\]](#) (147)

Fold: [P-loop containing nucleoside triphosphate hydrolases \[SCOP_52539\]](#)

Superfamily: [P-loop containing nucleoside triphosphate hydrolases \[SCOP_52540\]](#) (24)

Families: [Nucleotide and nucleoside kinases \[SCOP_52541\]](#) (20)

[Shikimate kinase \(AroK\) \[SCOP_52566\]](#)

[Chloramphenicol phosphotransferase \[SCOP_52569\]](#)

[Gluconate kinase \[SCOP_75195\]](#)

[Plasmid maintenance system epsilon/zeta, toxin zeta subunit \[SCOP_82395\]](#)

[Adenosine-5'phosphosulfate kinase \(APS kinase\) \[SCOP_52572\]](#)

[ATP sulfurylase C-terminal domain \[SCOP_64011\]](#)

[PAPS sulfotransferase \[SCOP_52575\]](#) (14)

[Phosphoribulokinase/pantothenate kinase \[SCOP_52584\]](#) (5)

[6-phosphofructo-2-kinase/fructose-2,6-bisphosphatase, kinase domain \[SCOP_52589\]](#)

[G proteins \[SCOP_52592\]](#) (78)

[Motor proteins \[SCOP_52641\]](#) (4)

[Nitrogenase iron protein-like \[SCOP_52652\]](#) (15)

[RecA protein-like \(ATPase-domain\) \[SCOP_52670\]](#) (17)

[Bacterial cell division inhibitor Sula \[SCOP_89678\]](#)

[ABC transporter ATPase domain-like \[SCOP_52686\]](#) (23)

[Tandem AAA-ATPase domain \[SCOP_81268\]](#) (23)

[Extended AAA-ATPase domain \[SCOP_81269\]](#) (28)

[RNA helicase \[SCOP_52724\]](#) (3)

[Helicase-like "domain" of reverse gyrase \[SCOP_69496\]](#)

[DNA helicase UvsW \[SCOP_102396\]](#)

[YjeE-like \[SCOP_75213\]](#)

[Type II thymidine kinase \[SCOP_117558\]](#)



weblogo.berkeley.edu

That's it!