

# Today's Lecture: HMMs

- Definitions
- Examples
- Probability calculations
  - WDAG
  - Viterbi algorithm

# HMMs: Formal Definition

- Alphabet  $\mathbf{B} = \{b\}$  of *observed symbols*
  - Set  $\mathbf{S} = \{k\}$  of *hidden states* (usually  $k = 0, 1, 2 \dots, m$ ; 0 is reserved for “begin” state, and sometimes also an “end” state)
  - (Markov chain property): prob of state occurring at given position depends only on immediately preceding state, and is given by
    - *transition probabilities* ( $a_{kl}$ ):  $a_{kl} = \text{Prob}(\text{next state is } l \mid \text{curr state is } k)$   
 $\sum_l a_{kl} = 1$ , for each  $k$ .
    - Usually, many transition probabilities are set to 0.
    - Model *topology* is the # of states, and *allowed* (i.e.  $a_{kl} \neq 0$ ) transitions.
- Sometimes omit begin state, in which case need *initiation probabilities* ( $p_k$ ) for sequence starting in a given state

*observed symbols*

A

G

C

A

T

$e_{\pi_1}(A)$

$e_{\pi_2}(G)$

$e_{\pi_3}(C)$

...

$e_{\pi_i}(A)$

...

$e_{\pi_n}(T)$

$0 \xrightarrow{a_{0\pi_1}} \pi_1 \xrightarrow{a_{\pi_1\pi_2}} \pi_2 \xrightarrow{a_{\pi_2\pi_3}} \pi_3 \xrightarrow{a_{\pi_3\pi_4}}$

$\pi_i \xrightarrow{a_{\pi_i\pi_{i+1}}}$

$\pi_n \rightarrow 0$

*unobserved states*

- Prob that symbol occurs at given sequence position depends only on hidden state at that position, and is given by

*emission probabilities:*

$$e_k(b) = \text{Prob}(\text{observed symbol is } b \mid \text{curr state is } k)$$

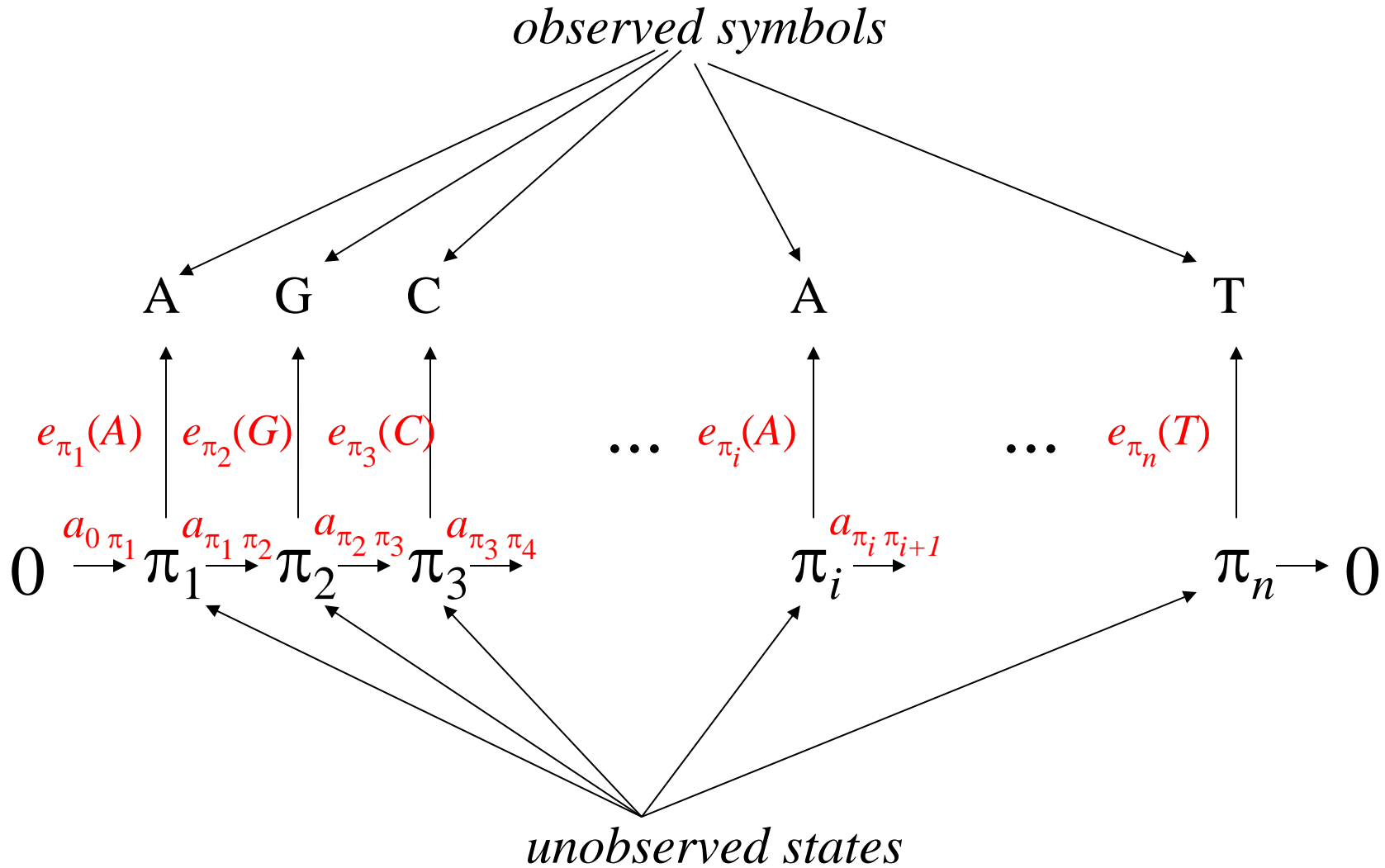
(begin and end states do not emit symbols)

- Note that
  - there are no *direct* dependencies between observed symbols in the sequence, however
  - there are *indirect* dependencies implied by state dependencies

# Where do the parameters come from?

- Can either
  - *define* parameter values *a priori*, or
  - *estimate* them from training data (observed sequences of the type to be modelled).
- Usually one does a mixture of both –
  - model topology is defined (some transitions set to 0), but
  - remaining parameters estimated

# Hidden Markov Model



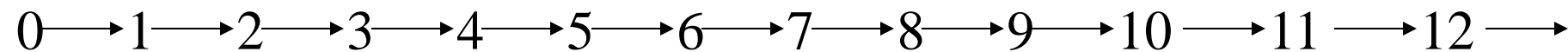
# HMM Examples

- Site models:
  - “states” correspond to positions (columns in the tables).  
state  $i$  transitions only to state  $i+1$ :
    - $a_{i,i+1} = 1$  for all  $i$ ;
    - all other  $a_{ij}$  are 0
  - emission probabilities are position-specific frequencies:  
values in frequency table columns

# Topology for Site HMM:

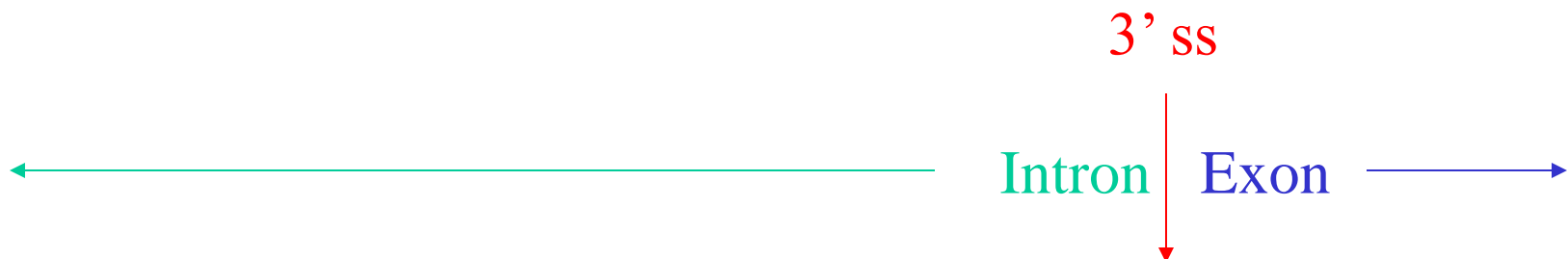
‘allowed’ transitions

(transits with non-zero prob – all are 1)





# HMM for *C. elegans* 3' Splice Sites



|   |      |      |      |      |      |      |      |      |      |      |      |      |
|---|------|------|------|------|------|------|------|------|------|------|------|------|
| A | 3276 | 3516 | 2313 | 476  | 67   | 757  | 240  | 8192 | 0    | 3359 | 2401 | 2514 |
| C | 970  | 648  | 664  | 236  | 129  | 1109 | 6830 | 0    | 0    | 1277 | 1533 | 1847 |
| G | 593  | 575  | 516  | 144  | 39   | 595  | 12   | 0    | 8192 | 2539 | 1301 | 1567 |
| T | 3353 | 3453 | 4699 | 7336 | 7957 | 5731 | 1110 | 0    | 0    | 1017 | 2957 | 2264 |

**CONSENSUS** W W W T T t C A G r w w

|                        |   |       |       |       |       |       |       |       |       |       |       |       |       |
|------------------------|---|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Emission probabilities | A | 0.400 | 0.429 | 0.282 | 0.058 | 0.008 | 0.092 | 0.029 | 1.000 | 0.000 | 0.410 | 0.293 | 0.307 |
|                        | C | 0.118 | 0.079 | 0.081 | 0.029 | 0.016 | 0.135 | 0.834 | 0.000 | 0.000 | 0.156 | 0.187 | 0.225 |
|                        | G | 0.072 | 0.070 | 0.063 | 0.018 | 0.005 | 0.073 | 0.001 | 0.000 | 1.000 | 0.310 | 0.159 | 0.191 |
|                        | T | 0.409 | 0.422 | 0.574 | 0.896 | 0.971 | 0.700 | 0.135 | 0.000 | 0.000 | 0.124 | 0.361 | 0.276 |

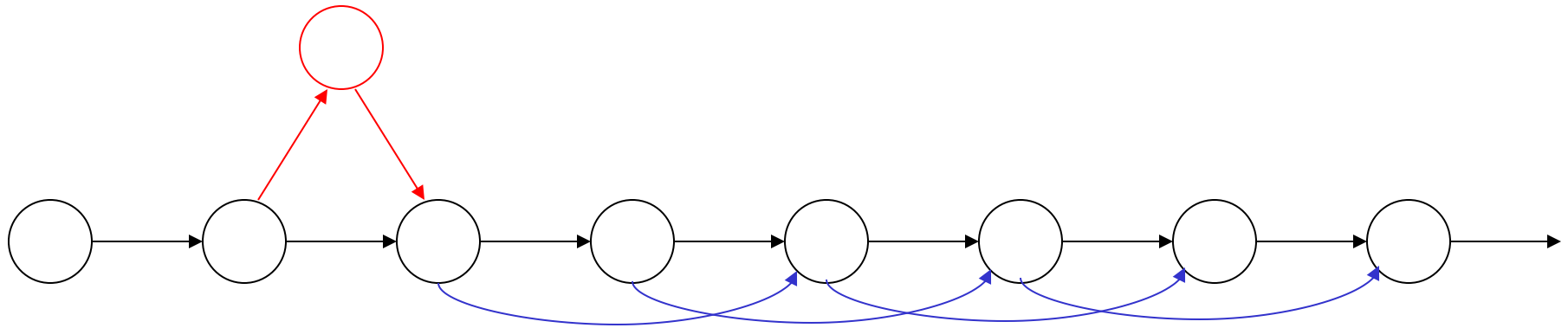
0 → 1 → 2 → 3 → 4 → 5 → 6 → 7 → 8 → 9 → 10 → 11 → 12

‘hidden’ states

- Can expand model to allow omission of nuc at some positions by including other (downstream) transitions (or via “silent states”)
- Can allow insertions by including additional states.
- transition probabilities no longer necessarily 1 or 0

# Insertions & Deletions in Site Model

insertion state



other transitions correspond  
to deletions

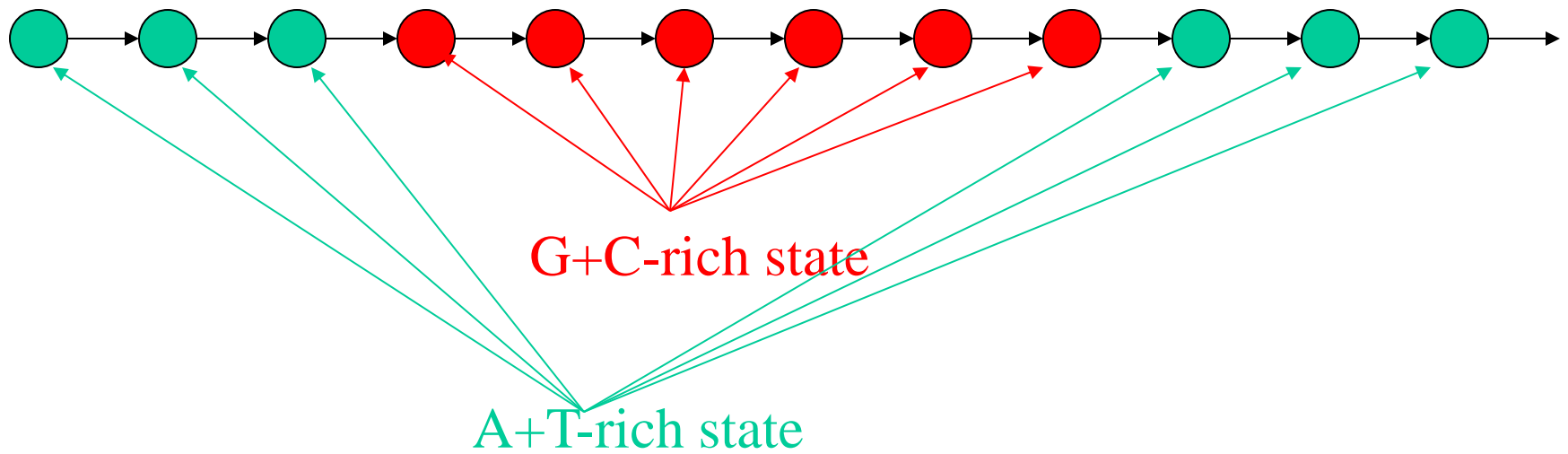
# Examples (cont'd) – 1-state HMMs

- single state, emitting residues with specified freqs:  
= ‘background’ model

# Examples (cont'd) – 2-state HMMs

- if  $a_{11}$  and  $a_{22}$  are small (close to 0), and  $a_{12}$  and  $a_{21}$  are large (close to 1), then get (nearly) periodic model with period 2; e.g.
  - dinucleotide repeat in DNA, or
  - (some) beta strands in proteins.
- if  $a_{11}$  and  $a_{22}$  large, and  $a_{12}$  and  $a_{21}$  small, then get models of alternating regions of different compositions (specified by emission probabilities), e.g.
  - higher vs. lower G+C content regions (RNA genes in thermophilic bacteria); or
  - hydrophobic vs. hydrophilic regions of proteins (e.g. transmembrane domains).

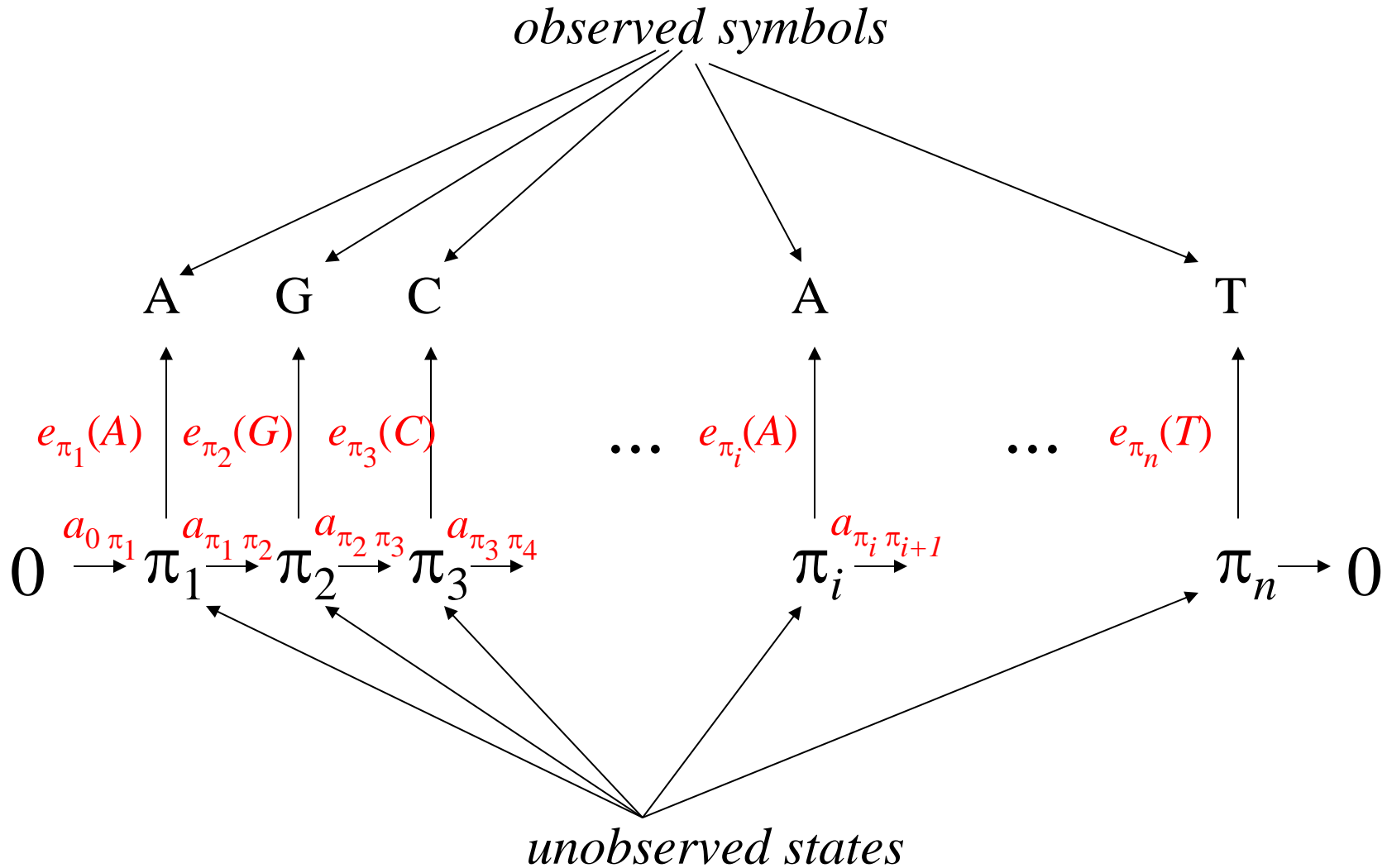
A A T G C C T G G A T A



# 2-state HMMs

- Can find most probable state decomposition (‘Viterbi path’) consistent with observed sequence
- Advantages over linked-list dynamic programming method (lecture 3) for finding high-scoring segments:
  - That method assumes you *know* appropriate parameters to find targeted regions; HMM method can *estimate* parameters.
  - HMM (easily) finds **multiple** segments
  - HMM can attach *probabilities* to alternative decompositions
  - HMM generalization to *> 2 types* of segments is easy – just allow more states!
- Disadvantage:
  - Markov assumption on state transitions implies geometric distribution for lengths of regions -- may not be appropriate

# Hidden Markov Model





# HMM Probabilities of Sequences

- Prob of **sequence of states**  $\pi_1\pi_2\pi_3 \dots \pi_n$  is  
 $a_{0\pi_1} a_{\pi_1\pi_2} a_{\pi_2\pi_3} a_{\pi_3\pi_4} \dots a_{\pi_{n-1}\pi_n}$ .
- Prob of **seq of observed symbols**  $b_1b_2b_3 \dots b_n$ ,  
*conditional on state sequence* is  
 $e_{\pi_1}(b_1)e_{\pi_2}(b_2) e_{\pi_3}(b_3) \dots e_{\pi_n}(b_n)$
- **Joint probability** =  $a_{0\pi_1} \prod_{i=1}^n a_{\pi_i\pi_{i+1}} e_{\pi_i}(b_i)$   
(define  $a_{\pi_n\pi_{n+1}}$  to be 1)
- (Unconditional) prob of observed sequence  
= **sum (of joint probs)** over all possible state paths
  - not practical to compute directly, by ‘brute force’! We will use dynamic programming.

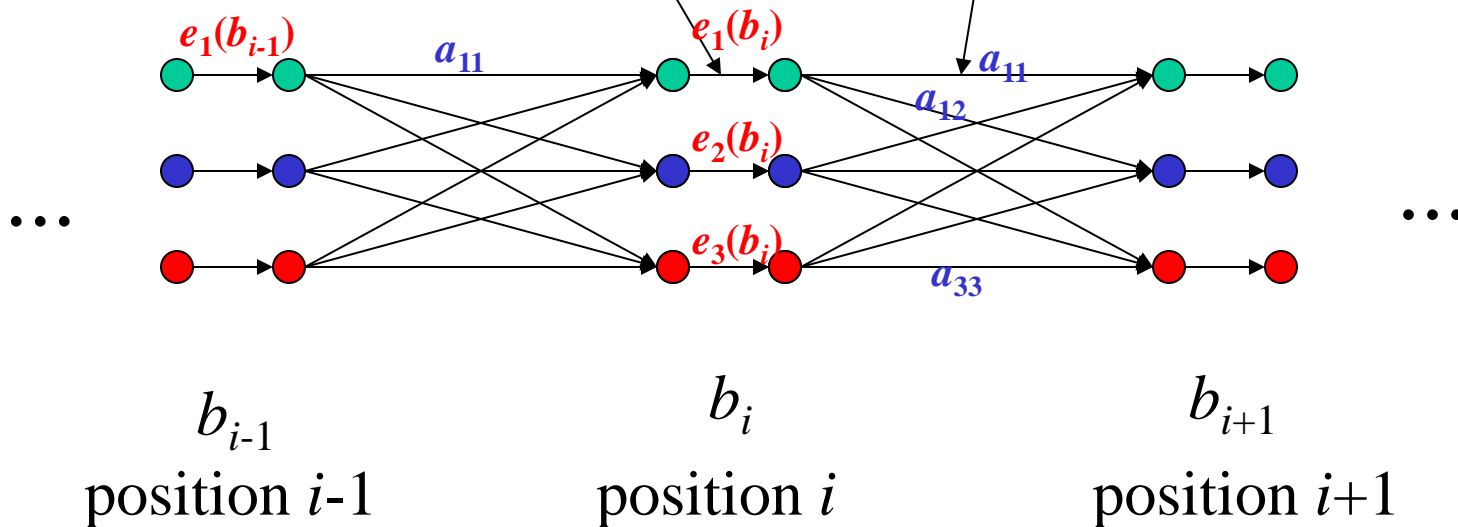
# Computing HMM Probabilities

- WDAG structure for sequence HMMs:
  - for  $i^{\text{th}}$  position in seq ( $i = 1, \dots, n$ ), have 2 nodes for each state:
    - total # nodes =  $2ns + 1$ , where  $n = \text{seq length}$ ,  $s = \# \text{ states}$
  - Pair of nodes for a given state at  $i^{\text{th}}$  position is connected by an *emission edge*
    - Weight is the emission prob for  $i^{\text{th}}$  observed residue.
    - Can omit node pair if emission prob = 0.
  - Have *transition edges* connecting (right-hand) state nodes at position  $i$  with (left-hand) state nodes at position  $i+1$ 
    - Weights are transition probs
    - Can omit edges with transition prob = 0.

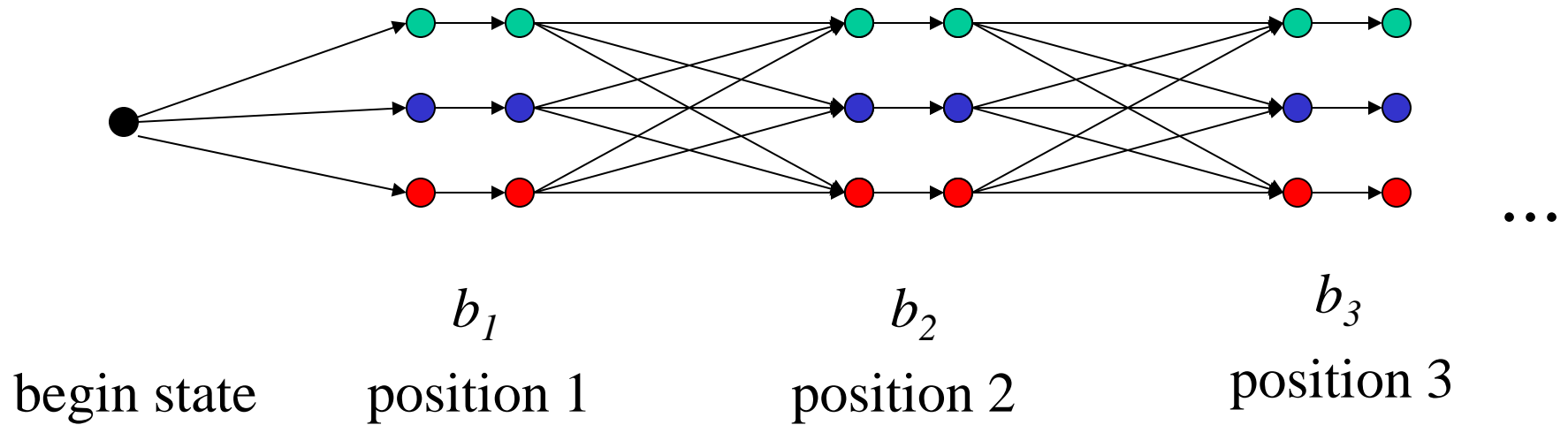
# WDAG for 3-state HMM, length $n$ sequence

weights are emission  
probabilities  $e_k(b_i)$  for  $i^{\text{th}}$   
residue  $b_i$

weights are transition  
probabilities  $a_{kl}$

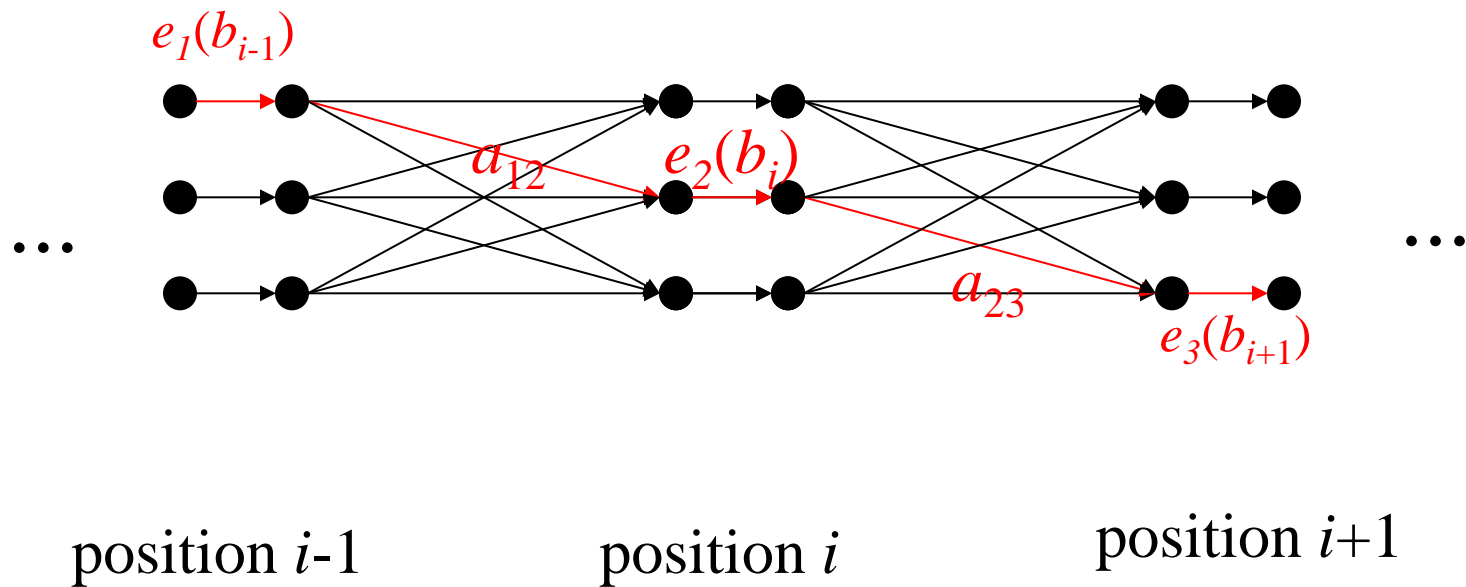


# Beginning of Graph



- ***Paths*** through graph from begin node to end node correspond to ***sequences of states***
- ***Product weight*** along path
  - = ***joint probability*** of state sequence & observed symbol sequence
- Sum of (product) path weights, over all paths,
  - = ***probability of observed sequence***
- Sum of (product) path weights over
  - all paths going through a particular node, or
  - all paths that include a particular edge,***divided by*** prob of observed sequence,
  - = ***posterior probability*** of that edge or node
- ***Highest-weight path*** = ***highest probability state sequence***

# Path Weights



- By general results on WDAGs, can use dynamic programming to find highest weight path:
  - = “**Viterbi algorithm**” to find highest probability path (most probable “parse”)
  - in this case can use log probabilities & sum weights
  - (N.B. paths are constrained to begin at the begin node!)

The Viterbi path is  
the *most probable parse!*