# CS 364 Computational Biology

Sara Mathieson

Haverford College

# Outline

- HMM example in population genetics

- Recap Viterbi Algorithm

- Forward-Backward Algorithm

- Posterior Decoding

Notes:
- Today and Tuesday: HMMs
- Next Thursday: review
- Please fill out exam/generative AI poll!

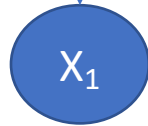# Hidden Markov models

Lunch on Monday

$X_1$

Lunch on Tuesday

$X_1$ $X_2$

Observations $\quad$ $X_1$ $\quad$ $X_2$ $\quad$ $X_3$ $\quad$ $X_4$ $\quad$ ......... $\quad$ $X_L$

Mood: "Busy" or "Relaxed"

Hidden state $Z_1$

Observations $X_1$ $X_2$ $X_3$ $X_4$ ......... $X_L$

Hidden state: $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_L$

Observations: $X_1$ $X_2$ $X_3$ $X_4$ ......... $X_L$

Markov property: hidden state i depends on state i-1 but *only* on state i-1

Hidden state

$Z_1$ → $Z_2$ → $Z_3$ → $Z_4$ ⤏ $Z_L$

Observations

$X_1$   $X_2$   $X_3$   $X_4$   ………   $X_L$

Markov property: hidden state i depends on state i-1 but *only* on state i-1
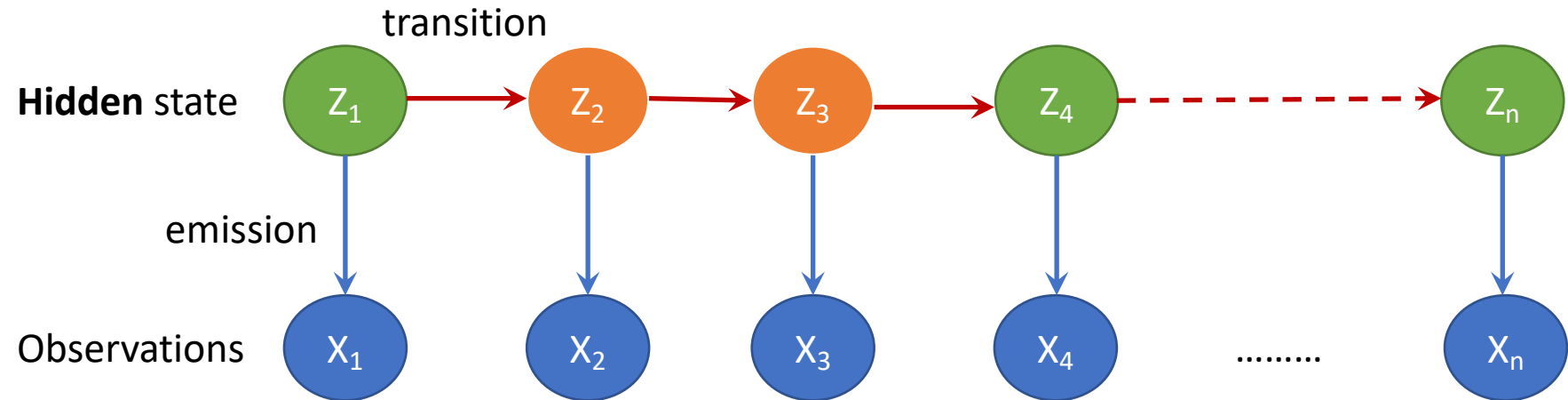


**Hidden Markov model**:

We have a series of **observations** that depend on some underlying **hidden state**.

Parameters: for each value of the hidden state, what's the probability of each possible observations [i.e. what is $P(X_i|Z_i)$?]. This is the **emission probability**

If we are in hidden state $Z_i$, what is the probability that we are in each of the possible hidden states in $Z_{i+1}$. [i.e. what is $P(Z_{i+1}|Z_i)$?]. This is the **transition probability**
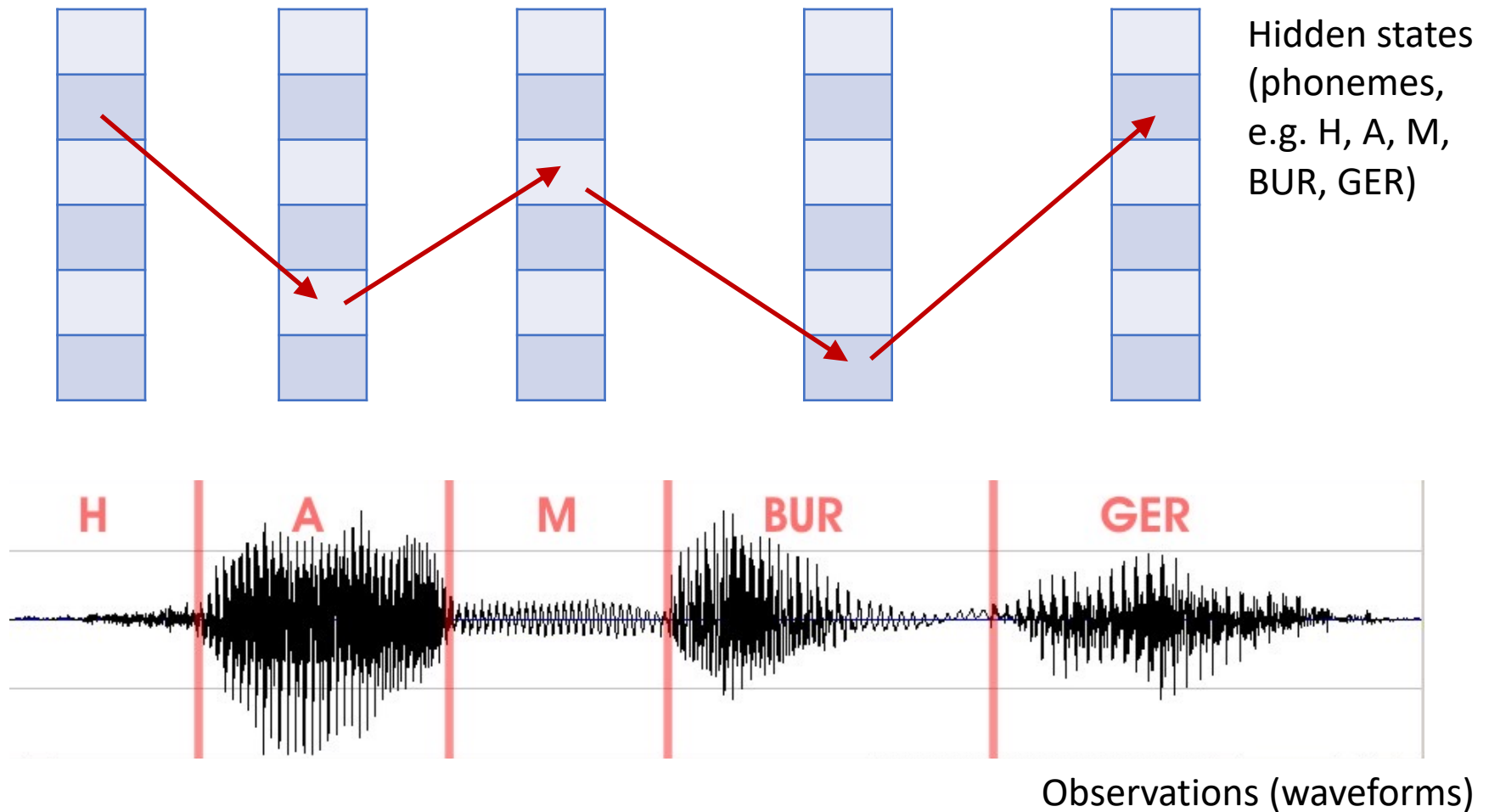
**Markov** property: hidden state i depends on state i-1 but *only* on state i-1



**What can we learn**

1) Given the observations, what were the hidden states? i.e. given the record of what I ate for lunch over a week, what was my mood on each day?

2) What are the parameters. e.g. if I am busy one day, what is the probability that I am relaxed the next day (transition)? If I am relaxed today, what is the probability that I eat out (emission)?

# Example: Speech recognition



Hidden states (phonemes, e.g. H, A, M, BUR, GER)
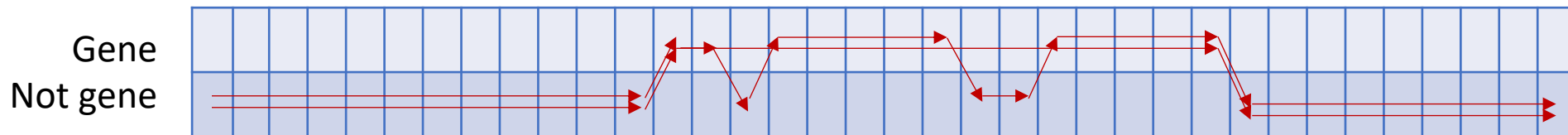
Observations (waveforms)

# Example: sequences (gene finding)

Where is the gene?

Coding sequences have higher CG content than noncoding regions

ACATCAGCTACGAT CGATGCGGGCATGAGTCCC ATATATTTAG

Gene
Not gene



Transition probabilities

|  | Gene | Not |
|---|---|---|
| **Gene** | 0.9 | 0.1 |
| **Not** | 0.05 | 0.95 |

Emission probabilities

|  | A | C | G | T |
|---|---|---|---|---|
| **Gene** | 0.25 | 0.25 | 0.25 | 0.25 |
| **Not** | 0.35 | 0.15 | 0.15 | 0.35 |

(e.g. average gene length is 10
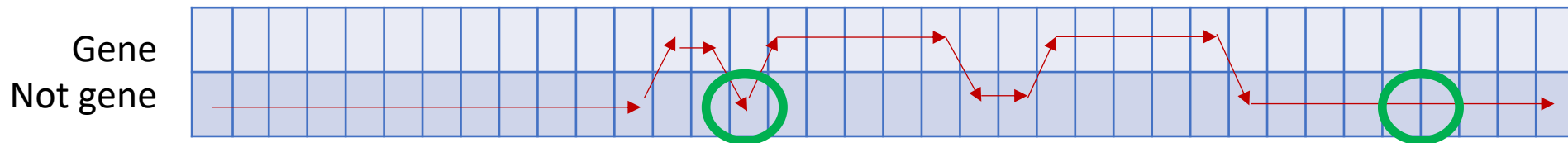and 1/20 of the genome is genic)

(e.g. GC content is 50% in genes and 30% outside)

# Example: sequences (gene finding)

Where is the gene?

Coding sequences have higher CG content than noncoding regions

ACATCAGCTACGATCGATGCGGGCATGAGTCCCATATATTTAG

Gene
Not gene

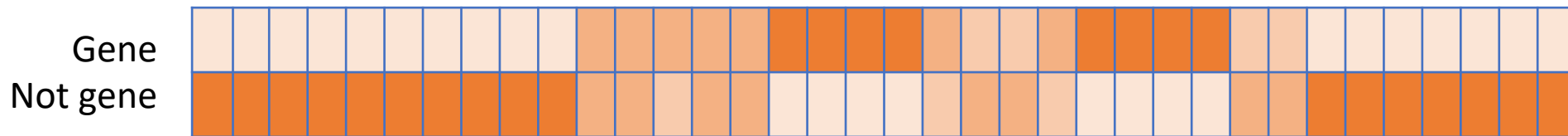The most likely path says that this nucleotide is not in a gene

But surely it is more likely to be in a gene than this one?

# Example: sequences (gene finding)

Where is the gene?

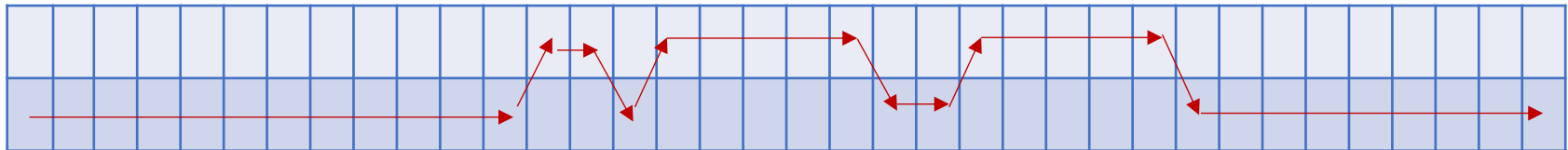Coding sequences have higher CG content than noncoding regions
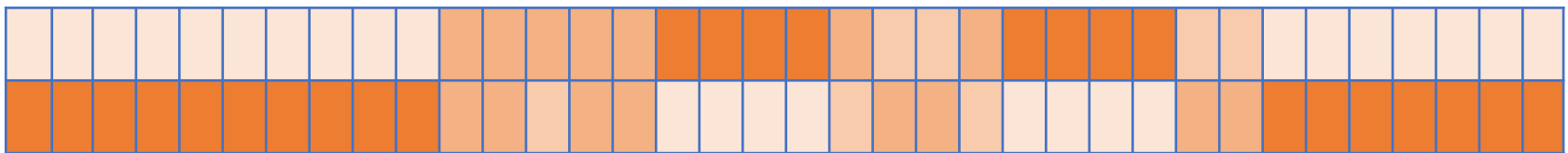
ACATCAGCTACGATCGATGCGGGCATGAGTCCCATATATTTAG

Gene
Not gene

Find probability of each hidden state, rather than most likely path

# Two ways of learning the hidden states
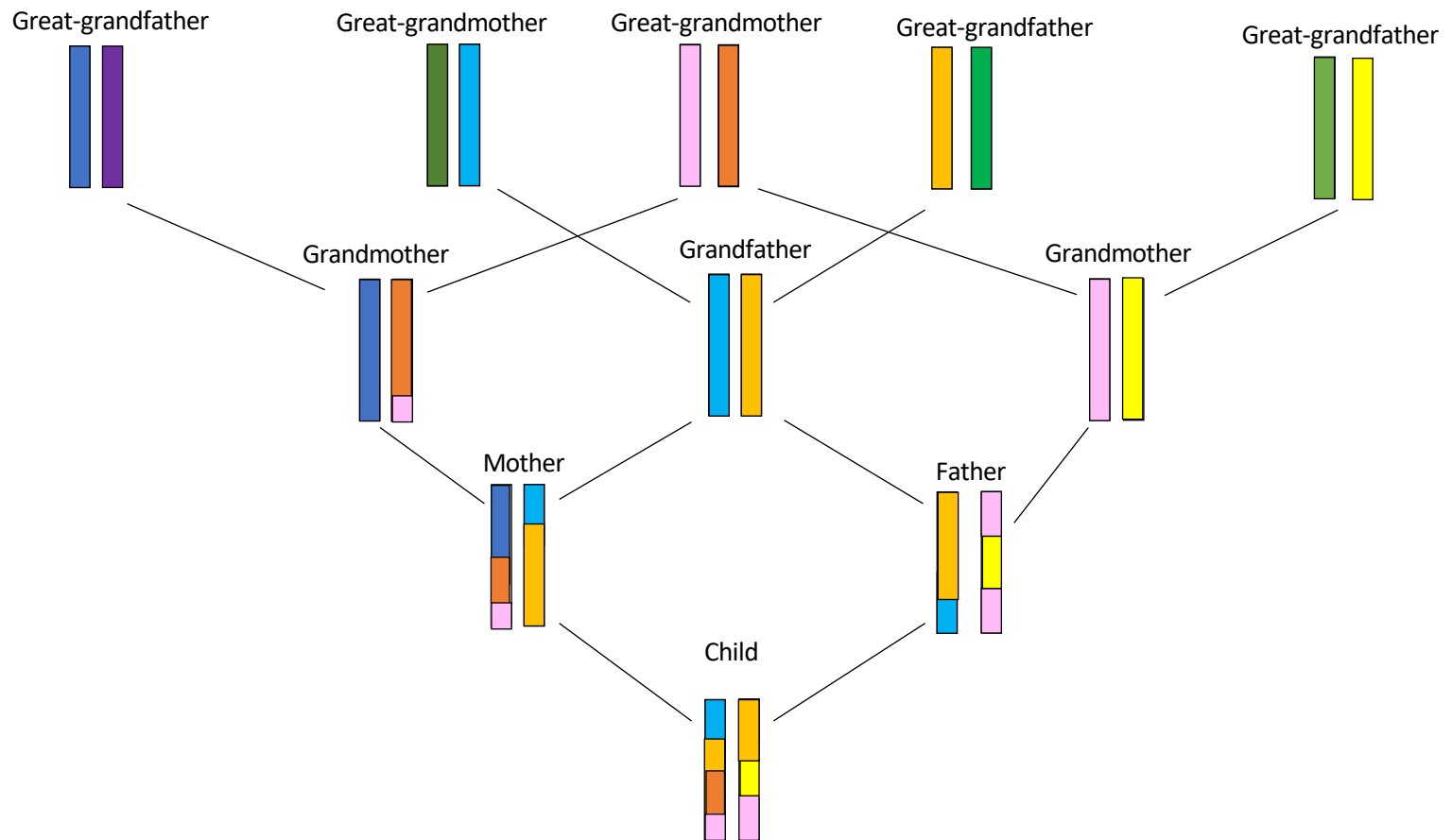
a) Find the most likely path: Viterbi algorithm



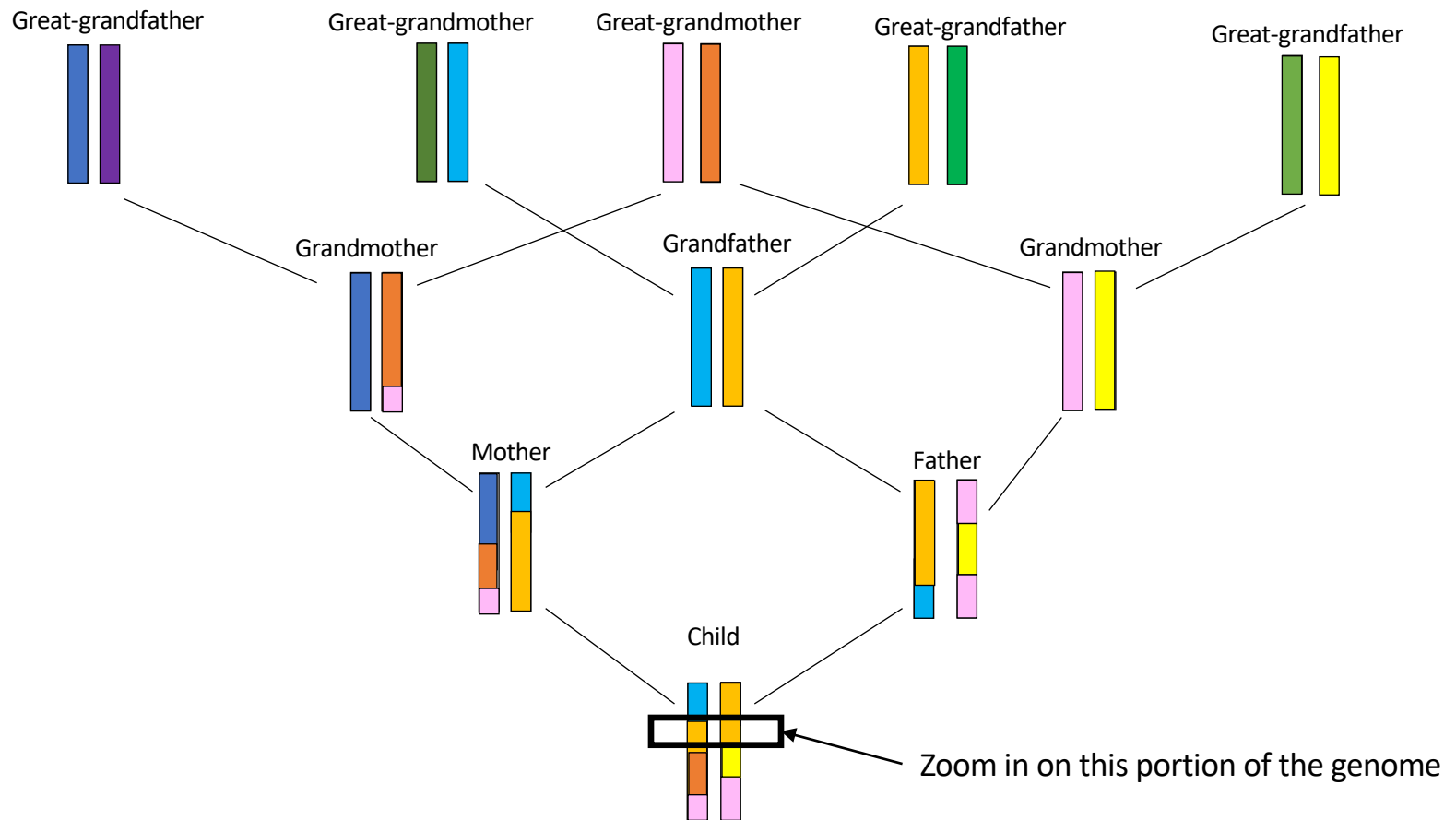b) Find posterior probability of each state: Forward-backward algorithm



Note that the Viterbi path is NOT the same as taking the maximum cells from the Forward-Backward algorithm ("posterior decoding"). The Viterbi path must be a path that is actually possible, but the maximum posterior might not be (e.g. because it implies impossible transitions).
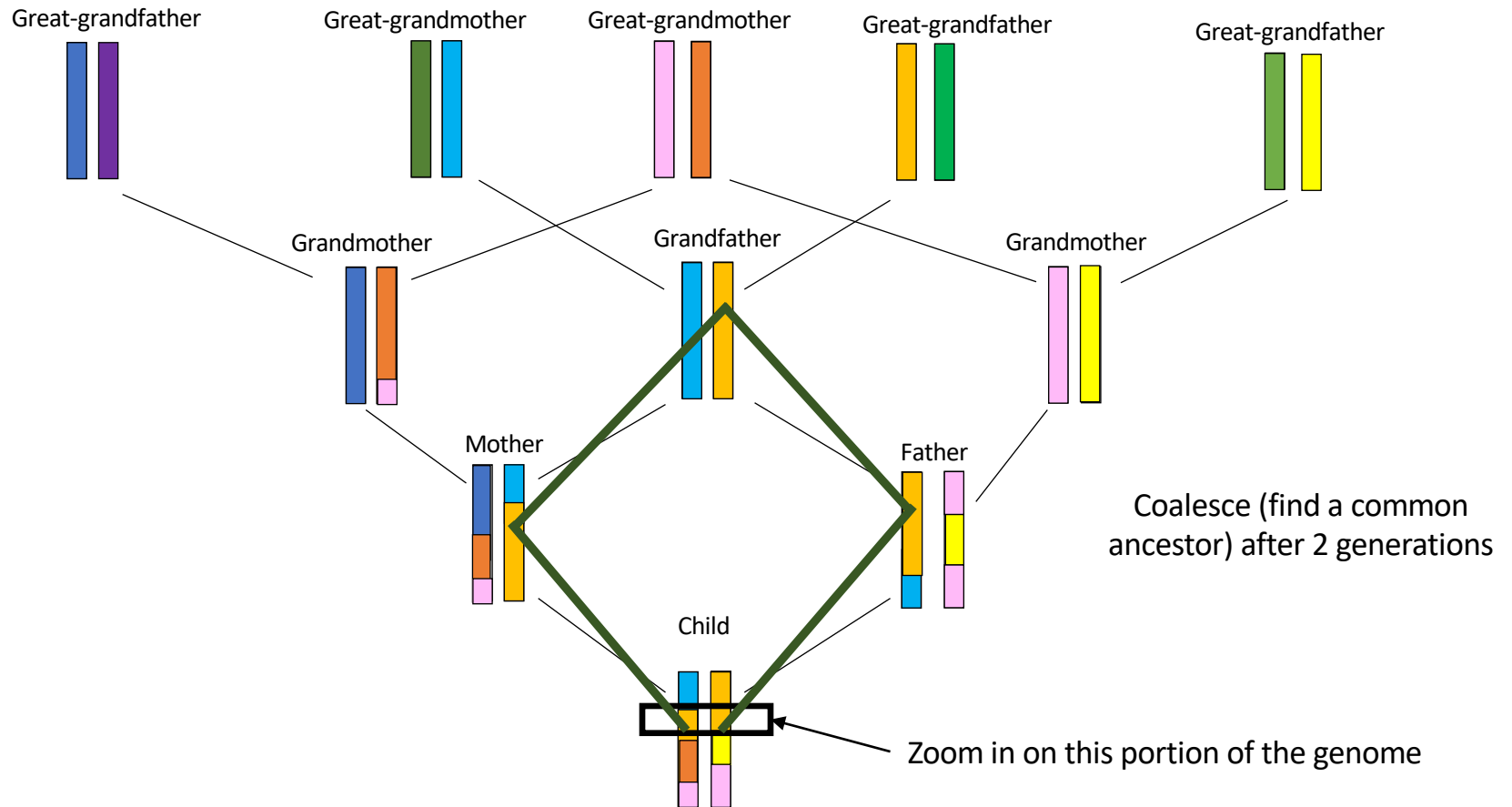
# HMM example from population genetics
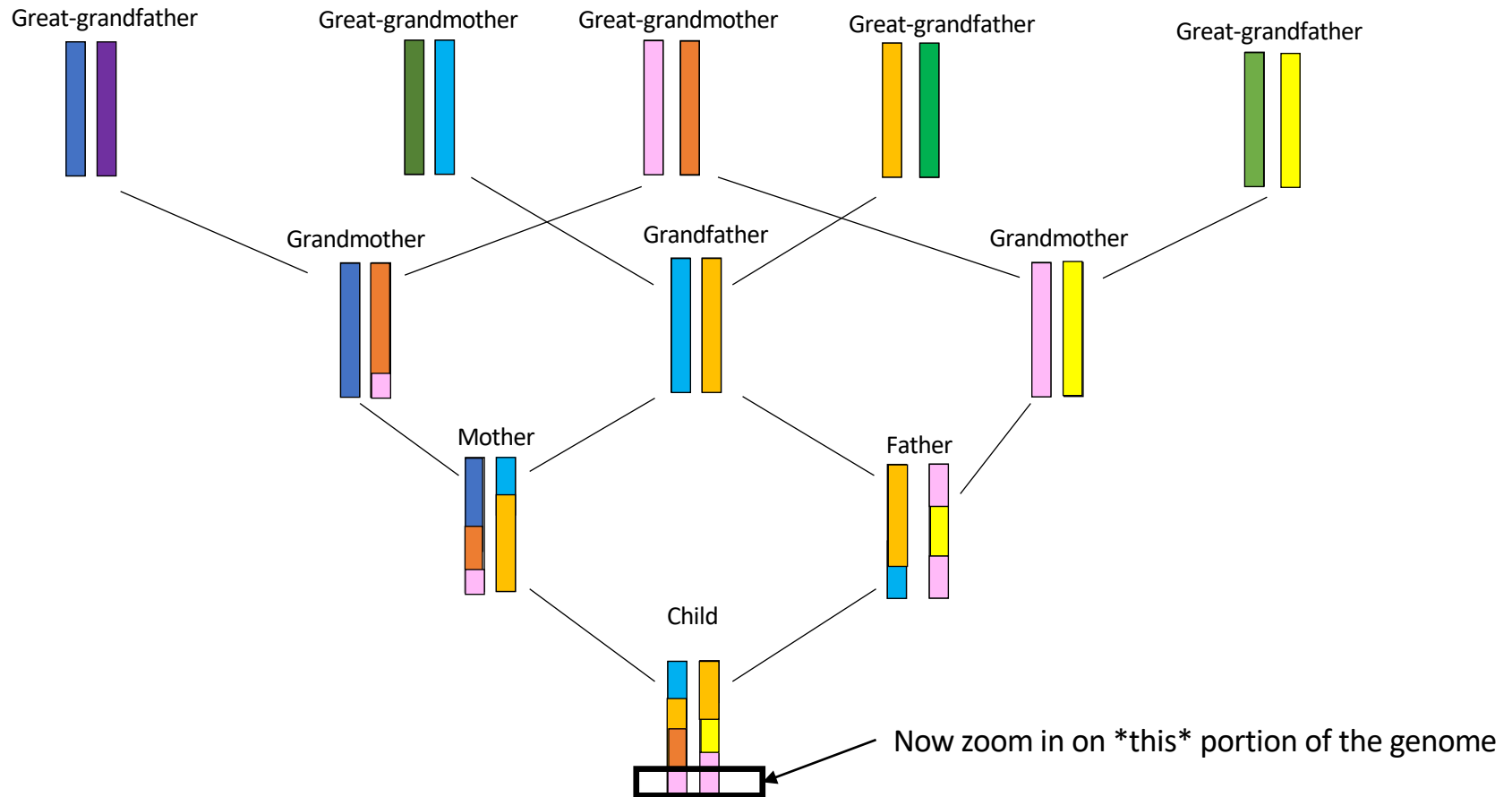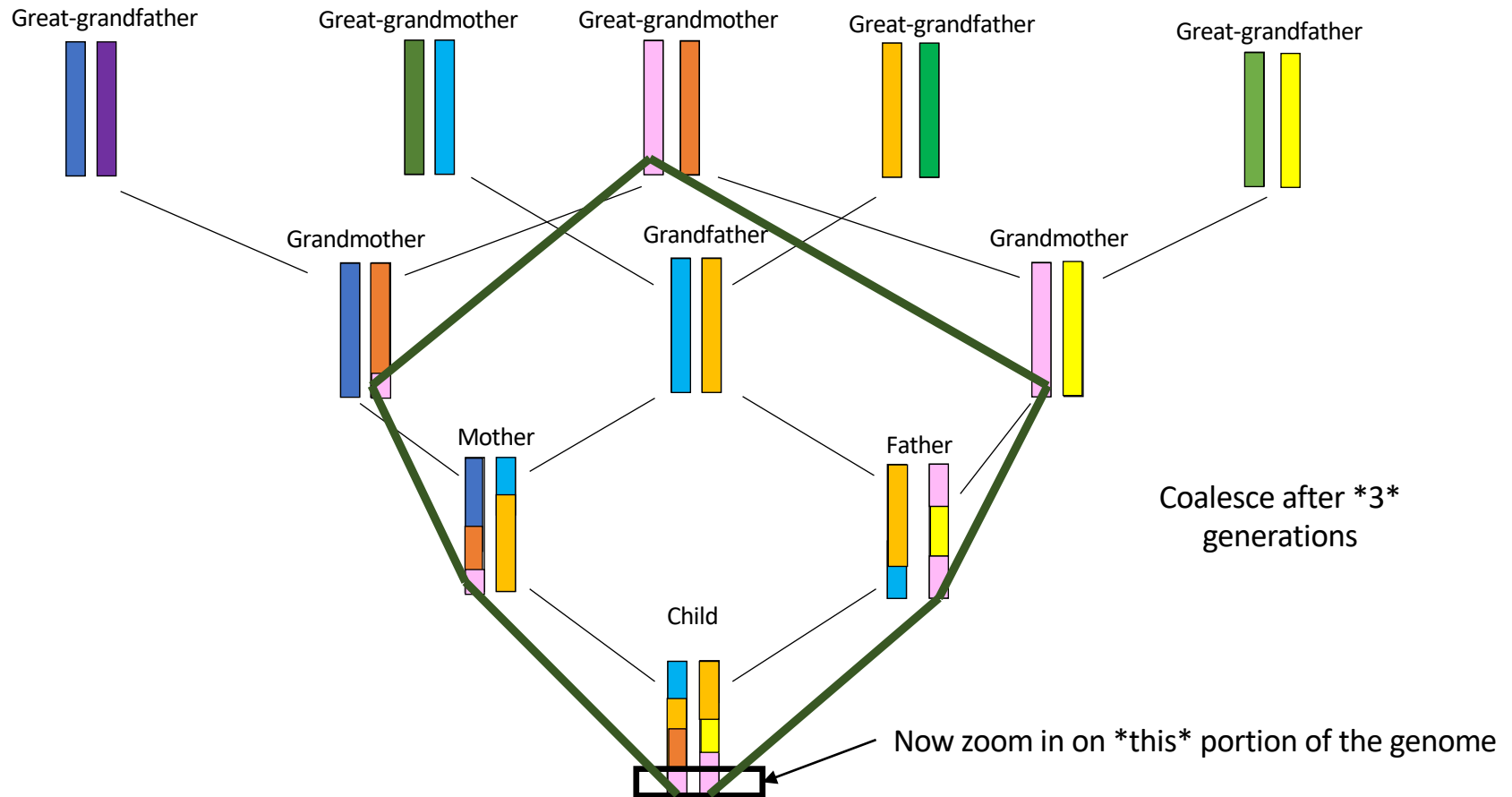
# Recombination over time

# Recombination over time



Great-grandfather

Great-grandmother

Great-grandmother

Great-grandfather

Great-grandfather

Grandmother

Grandfather

Grandmother

Mother

Father

Child

Zoom in on this portion of the genome

# Recombination over time

# Recombination over time



Now zoom in on *this* portion of the genome

# Recombination over time

Great-grandfather

Great-grandmother

Great-grandmother

Great-grandfather

Great-grandfather

Grandmother

Grandfather

Grandmother

Mother

Father

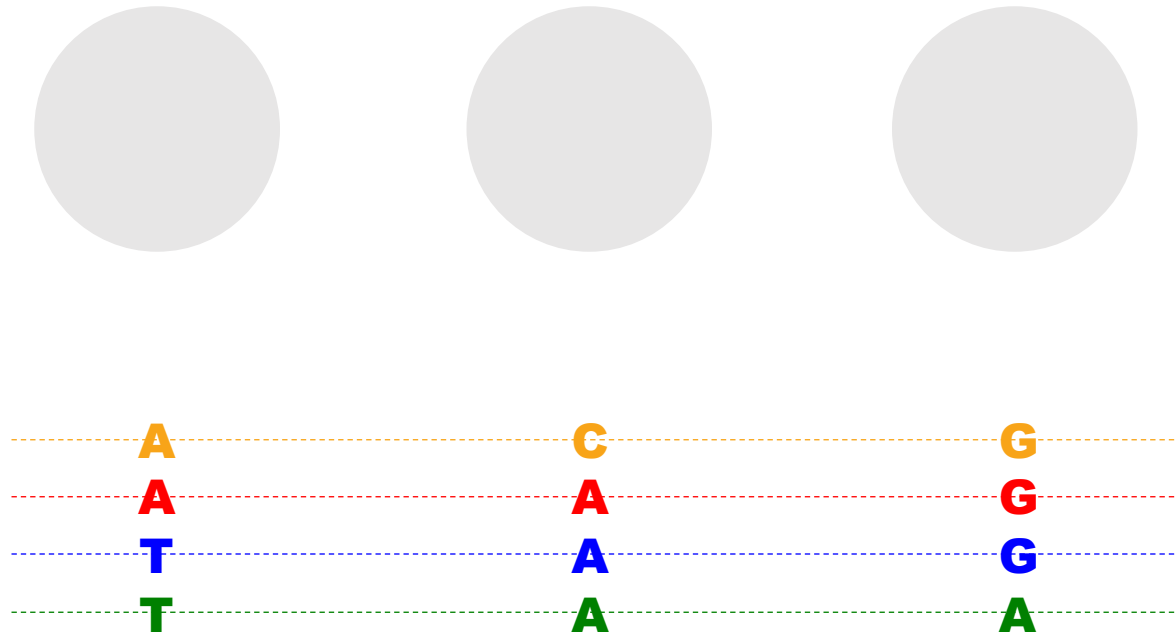Coalesce after *3* generations
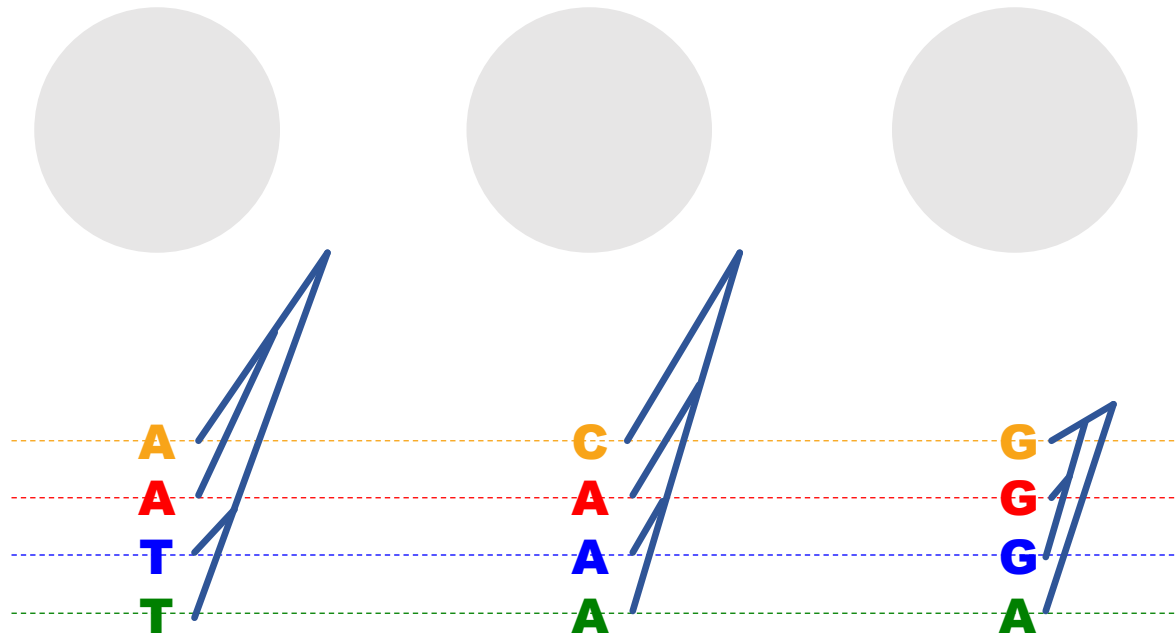
Child

Now zoom in on *this* portion of the genome

# How could we encode this as an HMM?

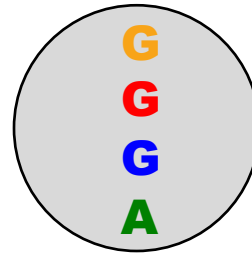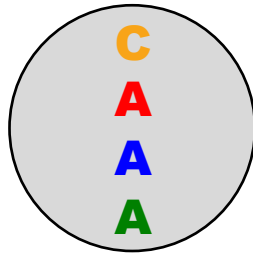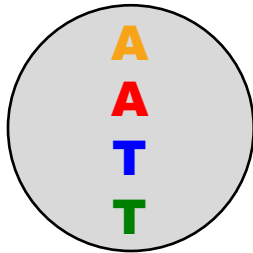- Take-home message: the tree changes across the genome! Both topology (for n > 2) and branch lengths
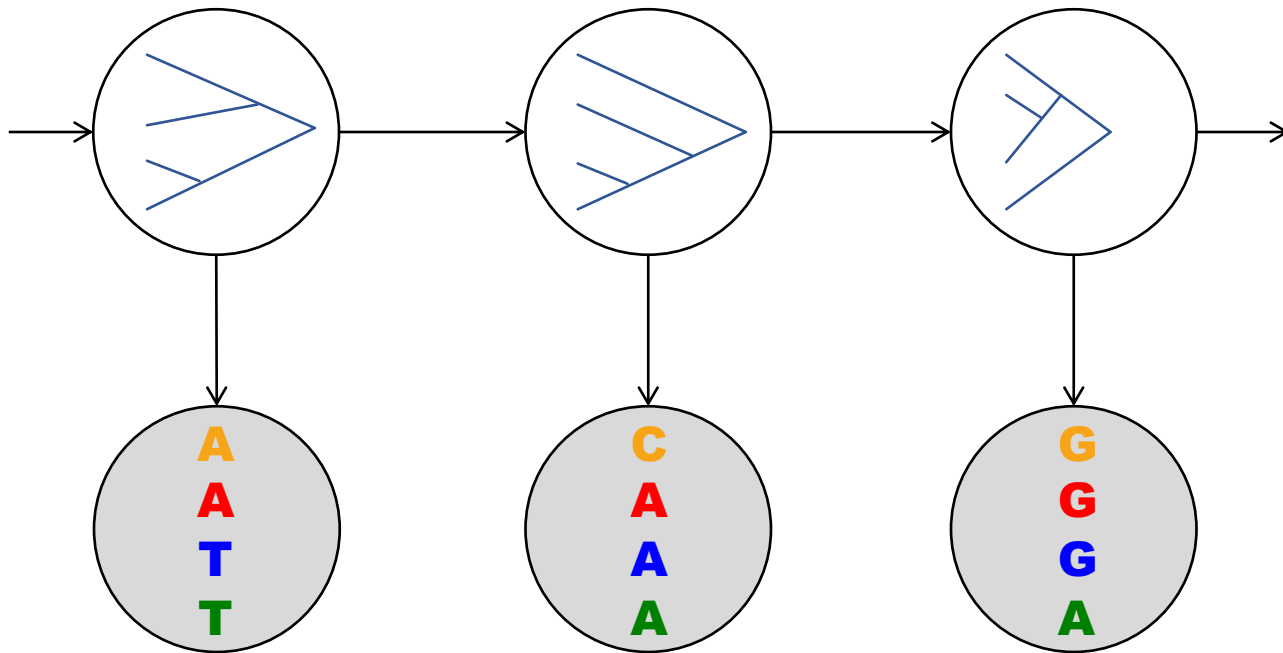
# Sequence data at many sites

# Tree changes along the genome!

# HMM observations: sequence data

# HMM hidden states: the tree

# Number of possible trees grows exponentially…



One person, two chromosomes!

**A**
**T**

**A**
**A**

**G**
**G**

Now the hidden state becomes the *time* of coalescence

# PSMC: pairwise sequentially Markovian coalescent

- The distribution of pairwise coalescence times should be **exponential** with parameter 1



Image: wikipedia

- If this differs from the exponential distribution, there were probably **population size changes**

- If all coalescence times are very recent, small population size

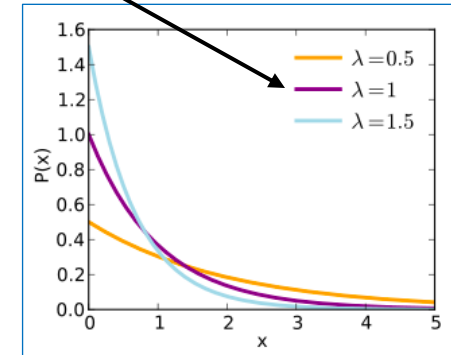- If all coalescence times are very ancient, large population size
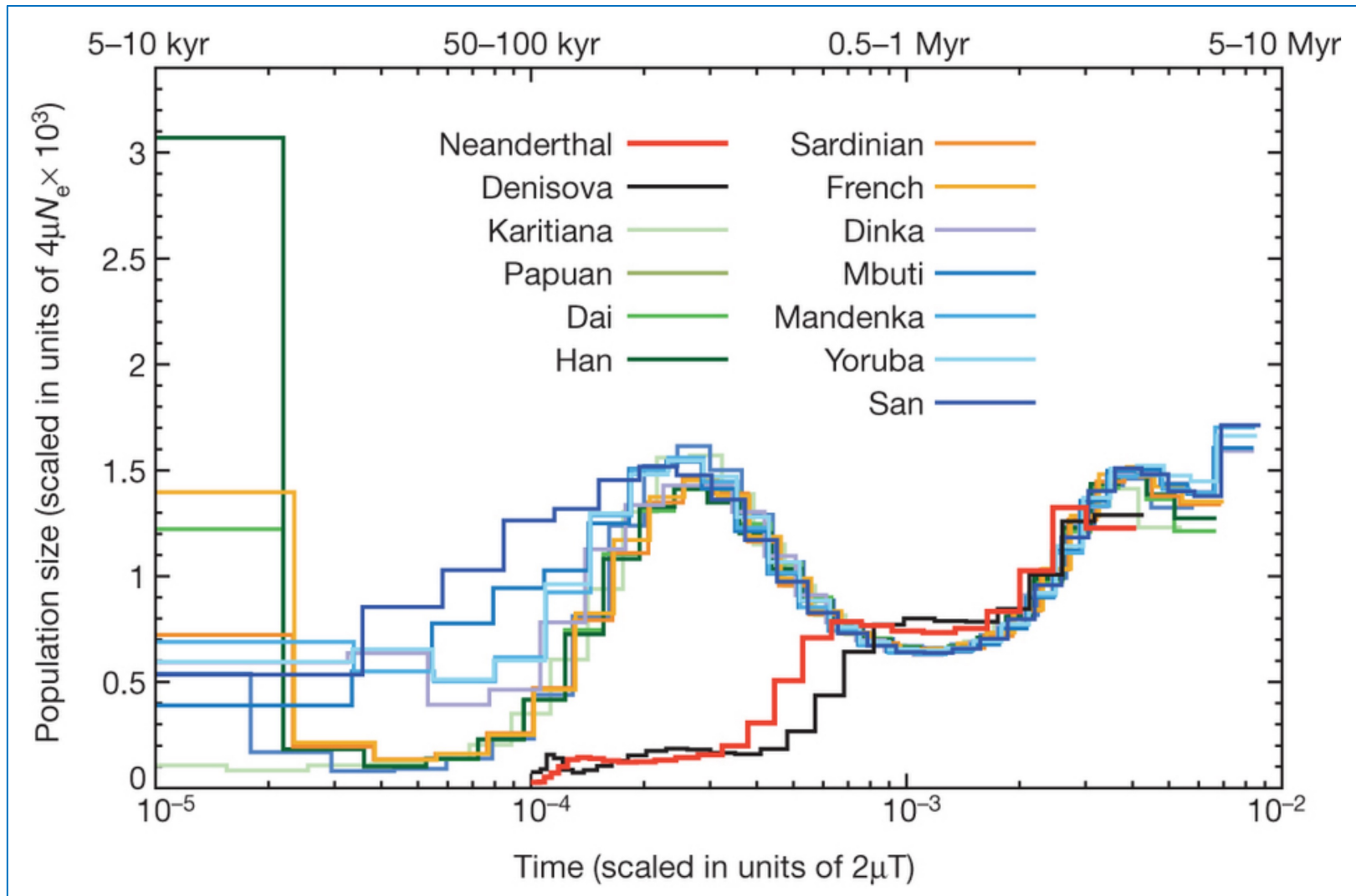
# PSMC: an HMM for two sequences

"The complete genome sequence of a Neanderthal from the Altai Mountains", Prufer et al (2014)

# HMM definition

- Transition probabilities:

(*K* x *K* matrix)

$$a_{kl} = P(z_i = l | z_{i-1} = k)$$

$z_{i-1}$      $z_i$

$k \longrightarrow l$

# HMM definition

- Transition probabilities:

($K$ x $K$ matrix)

$$a_{kl} = P(z_i = l | z_{i-1} = k)$$

$z_{i-1}$    $z_i$

$k \rightarrow l$

- Emission probabilities:

($K$ x $B$ matrix)

$$e_k(b) = P(x_i = b | z_i = k)$$

$z_i$

$k$

$b$

$x_i$

# HMM definition

- Transition probabilities:

(*K* x *K* matrix)

$$a_{kl} = P(z_i = l | z_{i-1} = k)$$

$z_{i-1}$ → $z_i$

$k$ → $l$

- Emission probabilities:

(*K* x *B* matrix)

$$e_k(b) = P(x_i = b | z_i = k)$$

$z_i$

$k$

$b$

$x_i$

- A way to deal with the first state:

$z_0$   $z_1$   $z_1$

$0$ →     $k$

$x_1$   $x_1$

$\pi_k = p(z_1 = k)$

$Z_0$ = Special start state with no emission

$\pi_k$ = probability of starting in state $k$

(*K* x *1* vector)

Do this way for lab!

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities ($\boldsymbol{a}$ and $\boldsymbol{e}$ matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $\boldsymbol{z}^*$

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, ..., x_L)$ and transition/emission probabilities ($\boldsymbol{a}$ and $\boldsymbol{e}$ matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $\boldsymbol{z}^*$

- **Initialization:** create a $K$ x $L$ matrix, this will be our dynamic programming (DP) table V

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, …, x_L)$ and transition/emission probabilities (***a*** and ***e*** matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence ***z***$^*$

- **Initialization:** create a *K* x *L* matrix, this will be our dynamic programming (DP) table

$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities ($a$ and $e$ matrices)
- **Output:** most probable (i.e. most likely) hidden state sequence $z^*$

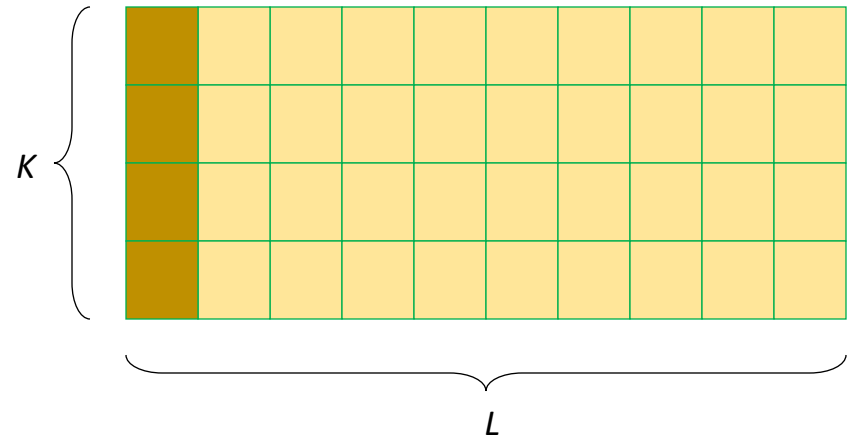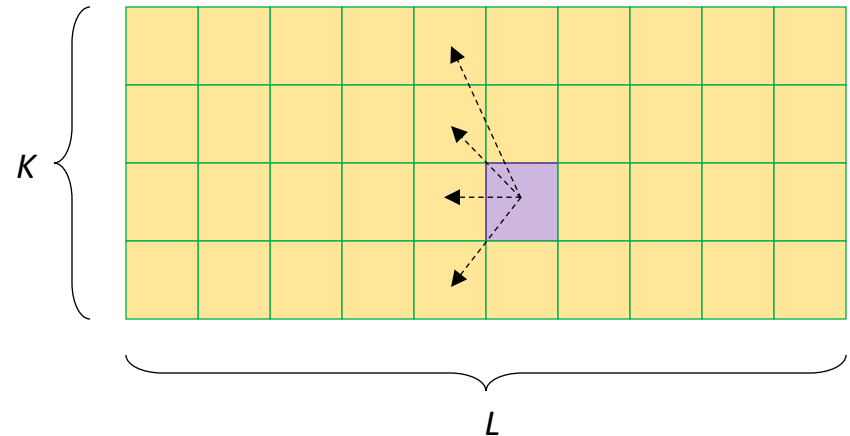- **Initialization:** create a $K \times L$ matrix, this will be our dynamic programming (DP) table

$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

- **Recursion:**

$$V_k(i) = e_k(x_i) \cdot \max_l \left\{ V_l(i-1) \cdot a_{lk} \right\}$$

$K$

$L$

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities (**a** and **e** matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $\boldsymbol{z}^*$

- **Initialization:** create a $K$ x $L$ matrix, this will be our dynamic programming (DP) table

$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

- **Recursion:**

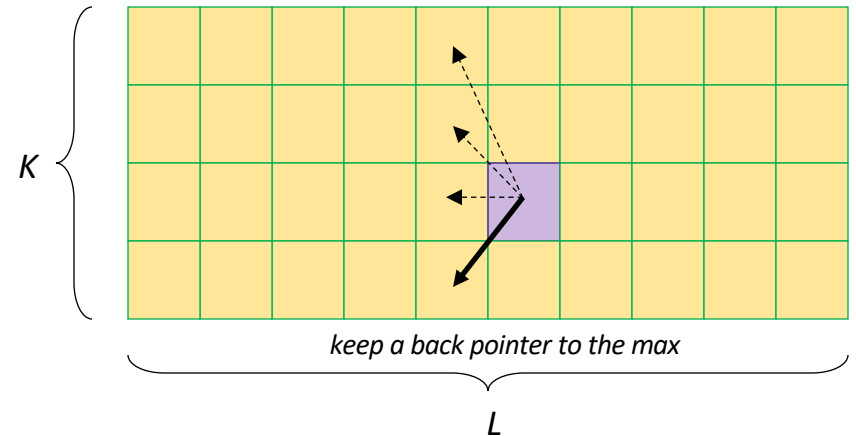$$V_k(i) = e_k(x_i) \cdot \max_l \left\{ V_l(i-1) \cdot a_{lk} \right\}$$



$K$

*keep a back pointer to the max*

$L$

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities ($a$ and $e$ matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $z^*$

- **Initialization:** create a $K$ x $L$ matrix, this will be our dynamic programming (DP) table

$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

- **Recursion:**

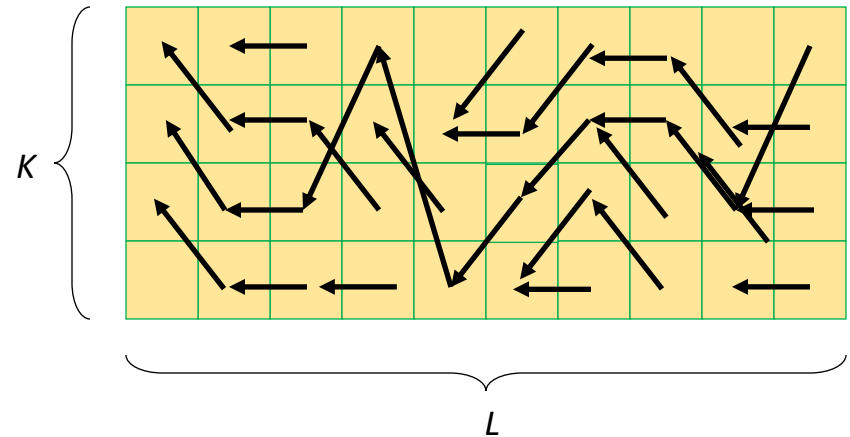$$V_k(i) = e_k(x_i) \cdot \max_l \left\{ V_l(i-1) \cdot a_{lk} \right\}$$

$K$

$L$

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, …, x_L)$ and transition/emission probabilities ($a$ and $e$ matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $z^*$

- **Initialization:** create a $K \times L$ matrix, this will be our dynamic programming (DP) table

$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

- **Recursion:**

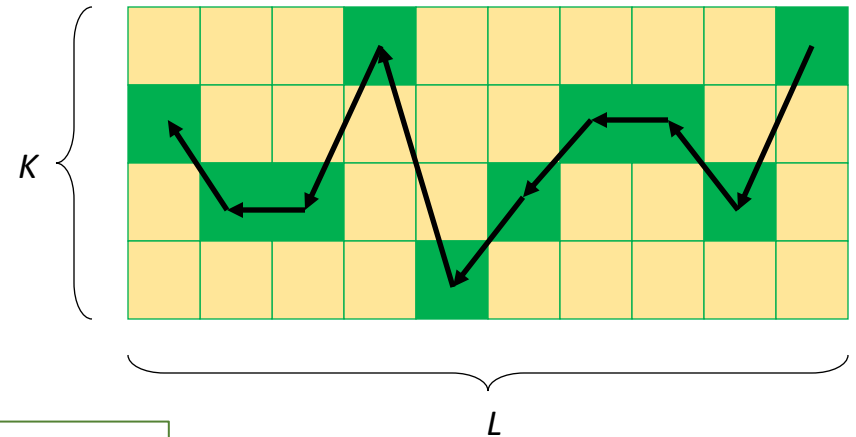$$V_k(i) = e_k(x_i) \cdot \max_l \left\{ V_l(i-1) \cdot a_{lk} \right\}$$

- **Termination and traceback:**

$$P(\vec{x}, \vec{z}^*) = \max_k \left\{ V_k(L) \right\}$$

# Viterbi Algorithm

- **Input:** observed sequence $(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities (**a** and **e** matrices)

- **Output:** most probable (i.e. most likely) hidden state sequence $\boldsymbol{z}^*$

- **Initialization:** create a $K \times L$ matrix, this will be our dynamic programming (DP) table
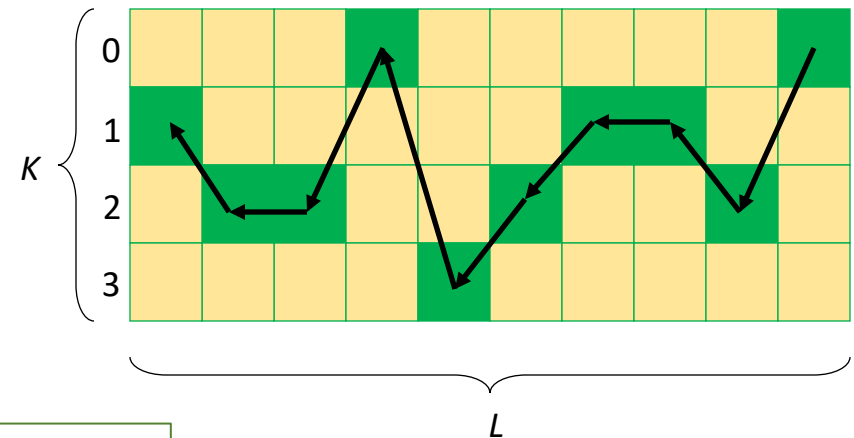
$$V_k(1) = \pi_k e_k(x_1)$$

*(Note: there are lots of ways to initialize, this avoids a special start state.)*

- **Recursion:**

$$V_k(i) = e_k(x_i) \cdot \max_l \left\{ V_l(i-1) \cdot a_{lk} \right\}$$
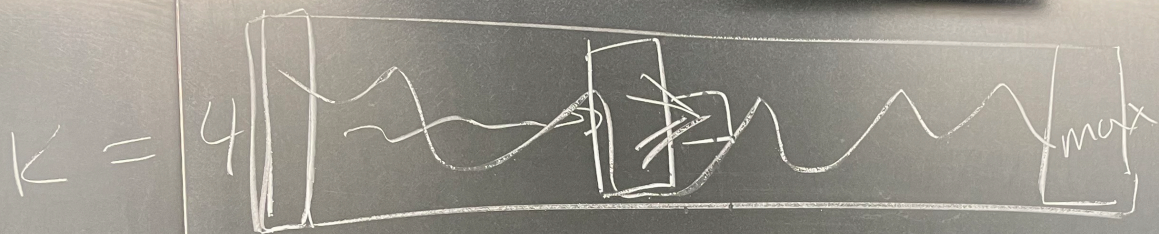


- **Termination and traceback:**

$$P(\vec{x}, \vec{z}^*) = \max_k \left\{ V_k(L) \right\}$$

$\boldsymbol{z}^* = (1,2,2,0,3,2,1,1,2,0)$

K states, L observations: what is the complexity of the Viterbi algorithm?

# Viterbi algorithm worksheet
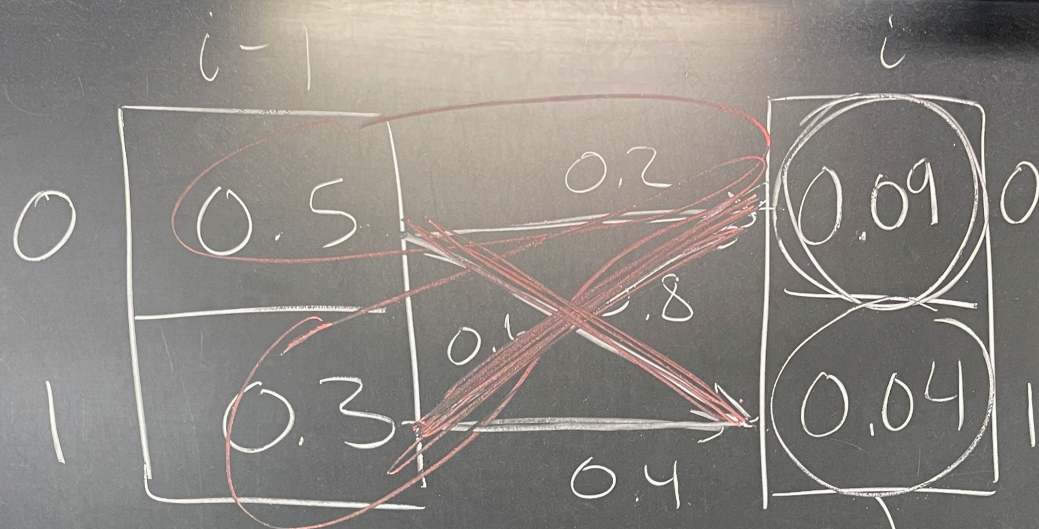
$k = 4$



$$L = 100,000$$

① $\quad K \cdot K \cdot K \cdots K = K^L$

$\quad\quad$ i=1 $\quad$ i=2 $\quad$ i=3 $\quad\quad\quad$ i=L

② $\quad$ at least $\cancel{O(KL)}$

$$O(LK^2)$$

③

|  | $i-1$ | | | $i$ | |
|---|---|---|---|---|---|
| 0 | 0.5 | 0.2 | | 0.09 | 0 |
| 1 | 0.3 | 0.1 0.8 | | 0.04 | 1 |
| | | 0.4 | | | |

(A)

State seq

| 1 | 0 |
|---|---|

$(z_1, z_2) = \vec{z}^*$

$$V_0(i) = e_0(A) \cdot \max\{V_0(i-1)a_{00}, V_1(i-1)a_{10}\}$$

$$= 0.5 \cdot \max\{0.5 \cdot 0.2, \boxed{0.3 \cdot 0.6}\}$$

$$\underbrace{\phantom{0.5 \cdot 0.2}}_{0.1} \qquad \underbrace{\phantom{0.3 \cdot 0.6}}_{0.18}$$

$$= \boxed{0.09}$$

$$V_1(i) = 0.1 \cdot \max\{\boxed{0.5 \cdot 0.8}, 0.3 \cdot 0.4\}$$

$$= 0.1 \cdot 0.4$$

$$= \boxed{0.04}$$

⑤ log space!

④ i+1

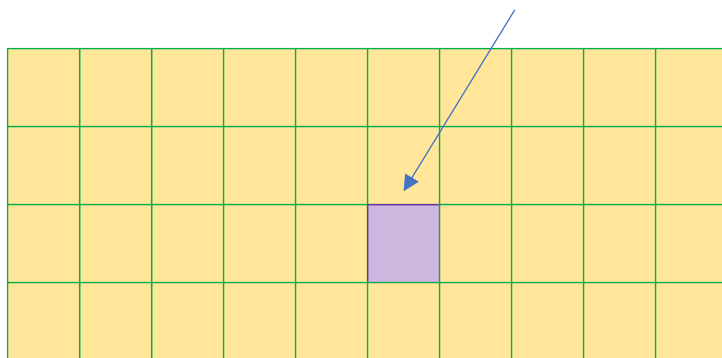| 0.0024 |
|--------|
| 0.0216 |

state seq

101

# Forward-backward Algorithm

- **Input:** observed sequence $x=(x_1, x_2, \ldots, x_L)$ and transition/emission probabilities (**a** and **e** matrices)

- **Output:** posterior probability of being in each hidden state at each time point $P(Z_i=k|x)$

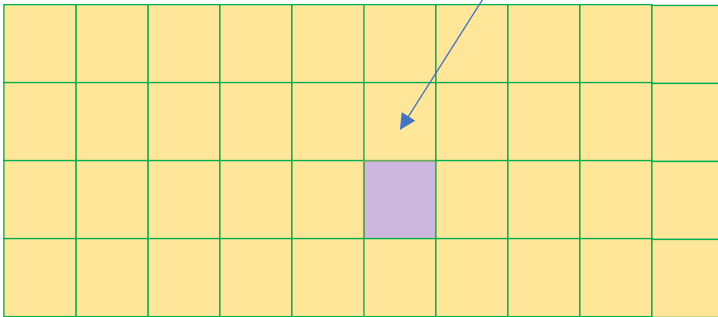What is the probability of being in this state, conditional on the observations?



Want to compute $P(Z_i=k|x_1, x_2 \ldots x_L) = P(Z_i=k, x_1, x_2 \ldots x_L) / P(x_1, x_2 \ldots x_L)$
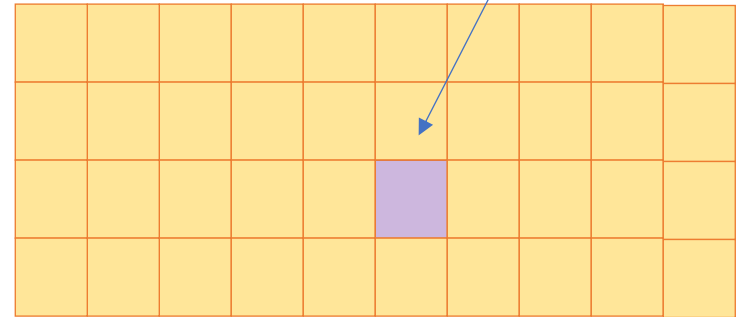
# Forward-backward Algorithm

Forward matrix $f_k(i)$

Backward matrix $b_k(i)$

What is the probability of being in this state, conditional on all the earlier observations

What is the probability of all the later observations, conditional on being in this state



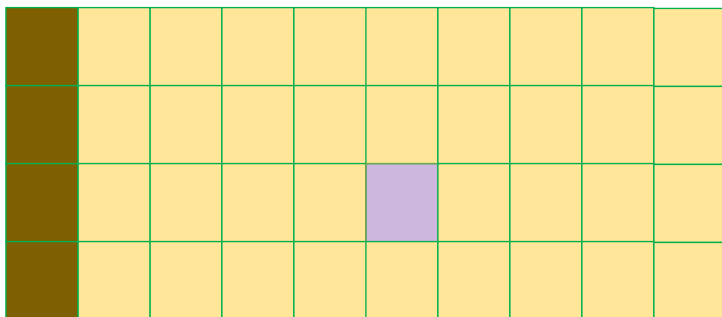Want to compute $P(Z_i=k|x_1,x_2...x_L) = P(Z_i=k, x_1,x_2...x_L) / P(x_1,x_2...x_L)$

But we can break this up: $P(Z_i=k, x_1,x_2...x_L) = P(Z_i=k, x_1,x_2...x_i) \, P(x_{i+1},x_2...x_L|Z_i=k)$

Probability of being in state k at time i, given the sequence up to time i: Forward algorithm

Probability of the sequence after time i, given that you are in state k at time i: Backward algorithm
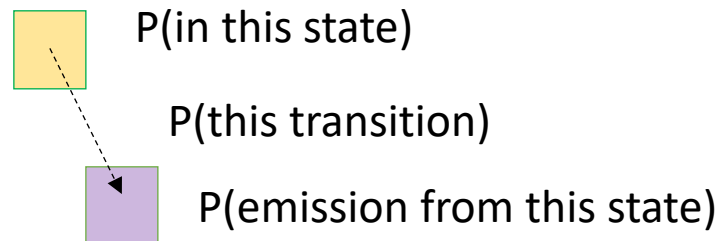
# Forward Algorithm
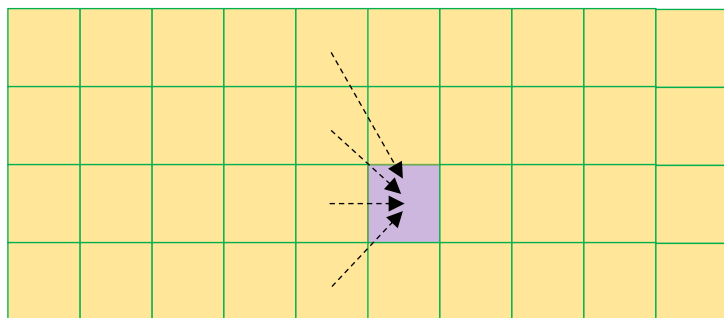
1. $f_k(1) = \pi_k e_k(x_1)$

$P(Z_i=k, x_1,x_2...x_i )$

Probability of being in state k at time i, given the sequence up to time i: Forward algorithm

# Forward Algorithm

1. $f_k(1) = \pi_k e_k(x_1)$

2. $f_k(i) = \Sigma_{j=1}^{K} f_j(i-1) a_{jk} e_k(x_i)$

P(in this state)

P(this transition)

P(emission from this state)

$P(Z_i = k, x_1, x_2 \ldots x_i)$

Probability of being in state k at time i, given the sequence up to time i: Forward algorithm

# Forward Backward Algorithm

## Goal

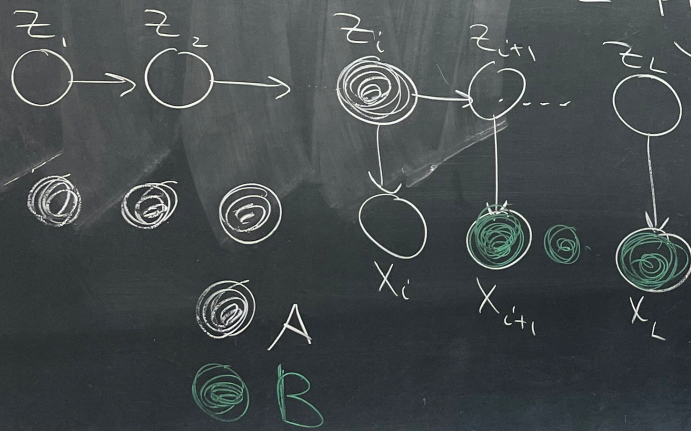$$P(z_i = k \mid \vec{x}) = \frac{P(\vec{x}, z_i = k)}{P(\vec{x})}$$

(posterior)

forward DP

$$P(\vec{x}, z_i = k) = P(\underbrace{x_1, x_2 \cdots x_i}_{\text{up to } i}, z_i = k, \underbrace{x_{i+1}, x_{i+2} \cdots x_L}_{\text{after } i})$$

$$\overset{\text{BAYES}}{=} \underbrace{P(x_1, x_2 \cdots x_i, z_i = k)}_{A} \cdot \underbrace{P(x_{i+1} \cdots x_L \mid z_i = k, \underbrace{\cancel{x_1, x_2 \cdots x_i}}_{A})}_{B}$$

forward
probabilities

backward
probabilities

$z_1 \rightarrow z_2 \rightarrow \cdots z_i \rightarrow z_{i+1} \cdots z_L$

$x_i \quad x_{i+1} \quad x_L$

A

B

# Forward Algorithm

$f_k(i)$ = prob of observing $x_1, \ldots x_i$ & ending in state $k$

$$= P(x_1, \ldots x_i, z_i = k)$$

recursion: $f_k(i) = e_k(x_i) \sum_{\ell} f_\ell(i-1) \, a_{\ell k}$

all possible prev states

$\underbrace{\quad}_{\ell}$ integrating over all possibilities

initialization (same as viterbi)

$$f_k(1) = \pi_k e_k(x_1)$$
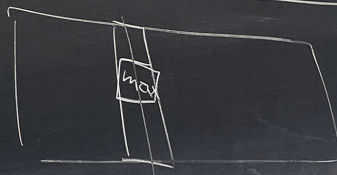
termination $\quad P(\vec{x}) = \sum_k f_k(L)$

## together

$$p(\vec{x}, z_i = k) = f_k(i) b_k(i)$$

$$\boxed{p(z_i = k \mid \vec{x}) = \frac{f_k(i) b_k(i)}{P(\vec{x})}}$$
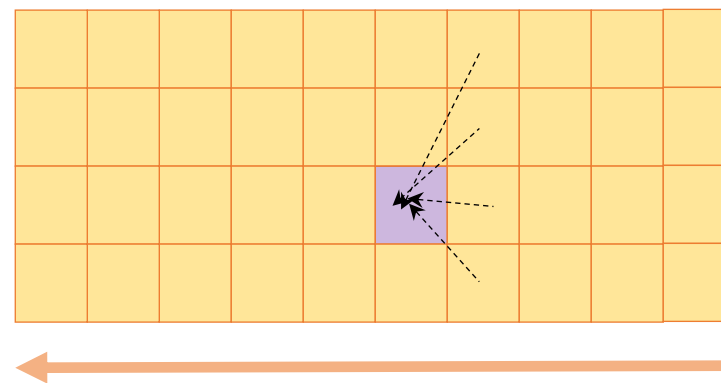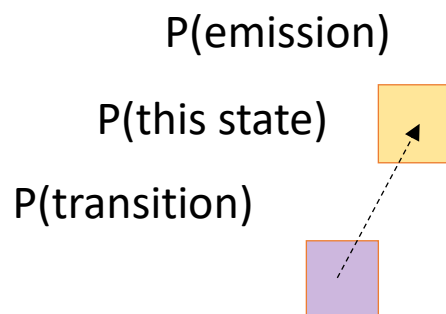
posterior prob

posterior decoding

$$\boxed{\hat{z}_i = \underset{k}{\arg\max} \, P(z_i = k \mid \vec{x})}$$

max

# Backward Algorithm

1 $b_k(K) = 1$

2 $b_k(i) = \Sigma_{j=1}^{K} b_j(i+1) a_{ij} e_j(x_{i+1})$

P(emission)

P(this state)

P(transition)

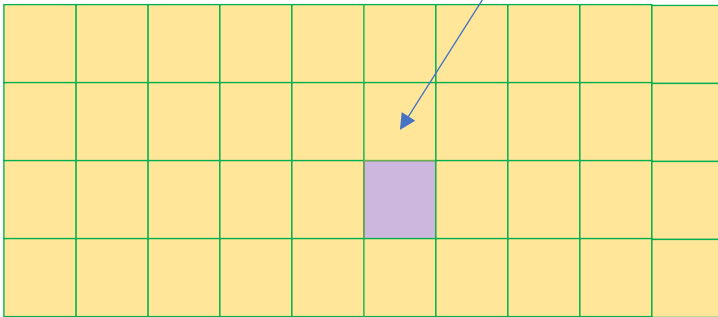$P(x_{i+1}, x_2 ... x_L | Z_i = k)$

Probability of the sequence after time i, given that you are in state k at time i: Backward algorithm
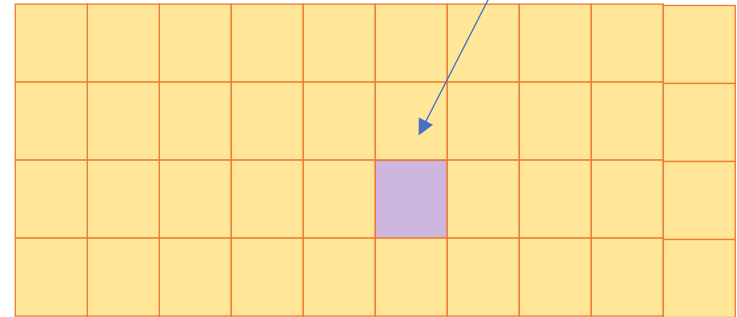
# Forward-backward Algorithm

Forward matrix $f_K(i)$

What is the probability of being in this state, conditional on all the earlier observations

Backward matrix $b_k(i)$

What is the probability of all the later observations, conditional on being in this state

Want to compute $P(Z_i=k|x_1,x_2...x_L) = P(Z_i=k, x_1,x_2...x_L) / P(x_1,x_2...x_L)$

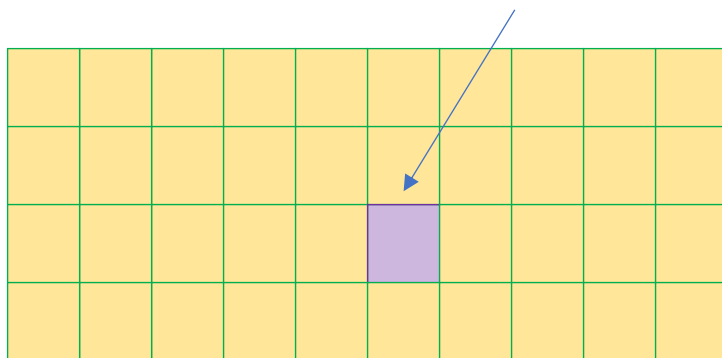But we can break this up: $P(Z_i=k, x_1,x_2...x_L) = P(Z_i=k, x_1,x_2...x_i) P(x_{i+1},x_2...x_L|Z_i=k)$

Probability of being in state k at time i, given the sequence up to time i: Forward algorithm

Probability of the sequence after time i, given that you are in state k at time i: Backward algorithm

# Forward-backward Algorithm

- **Input:** observed sequence $x=(x_1,x_2,…,x_L)$ and transition/emission probabilities ($a$ and $e$ matrices)
- **Output:** posterior probability of being in each hidden state at each time point $P(Z_i=k|x)$

What is the probability of being in this state, conditional on the observations?



Want to compute $P(Z_i=k|x_1,x_2…x_L) = P(Z_i=k, x_1,x_2…x_L) / P(x_1,x_2…x_L)$

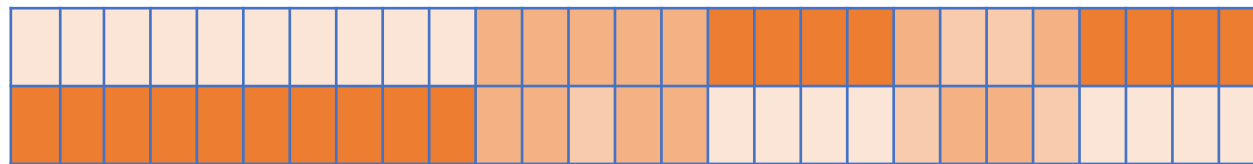What is the runtime of the forward-backward algorithm?

# Parameter estimation

So far we assumed that we know the emission and transition probabilities, and the initial state distribution. If we don't know these, can we estimate them?

Yes – using the Baum-Welch algorithm; a special case of the Expectation-Maximization (EM) algorithm:

0. Guess some parameters and run the F-B algorithm

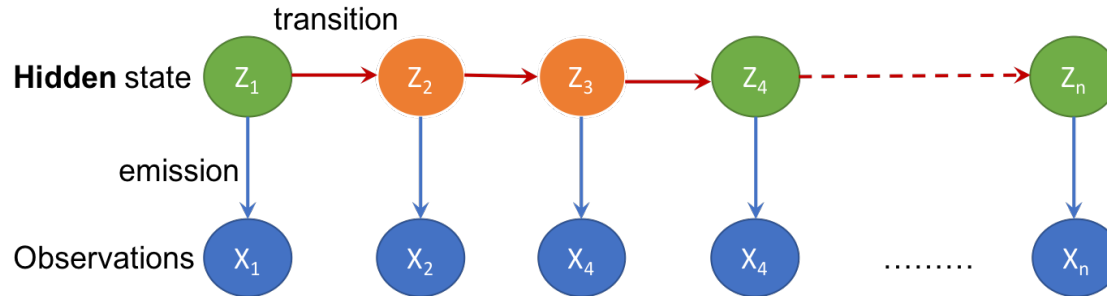1. Use estimated parameters to run the F-B algorithm

2. Use the output of the F-B algorithm to estimate the parameters

(estimate the parameters by using the empirical probabilities from the FB matrix)
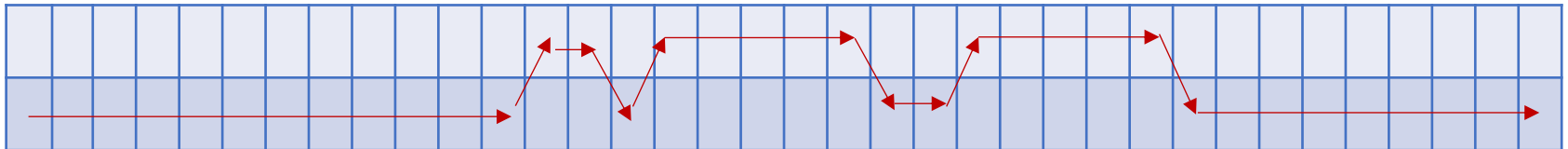
# Summary of HMMs

Structure of a Hidden Markov Model: Observations, hidden states, emissions, transitions
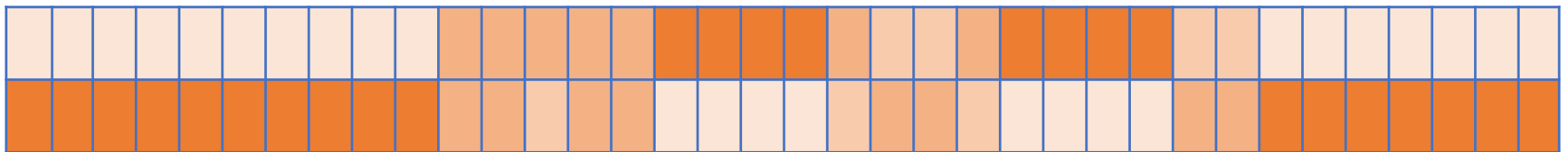


Three inference problems:

1) What is the most likely sequence of hidden states?: Viterbi algorithm



2) What is the most likely hidden state at each observation? Forward-Backward algorithm



3) What are the parameters? Baum-Welch algorithm