

Hidden Markov Models and Their Applications in Bioinformatics

Instructor: Connie Li

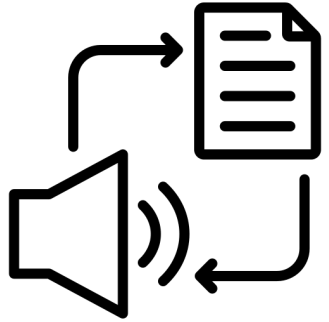
Given at the University of Calgary

November 14, 2024

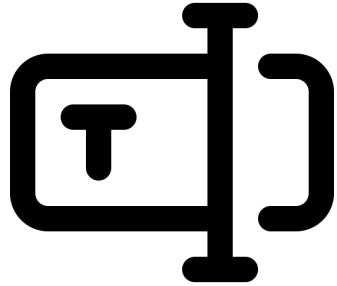
Outline: From theory to application

1. Introducing hidden Markov models
2. Formal definitions
3. Solving and inference
4. HMMs in Bioinformatics
5. Wrap up

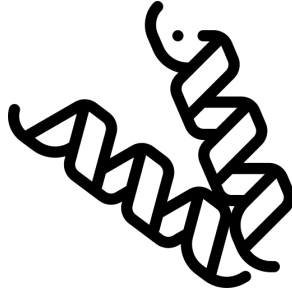
Why do we care about HMMs?



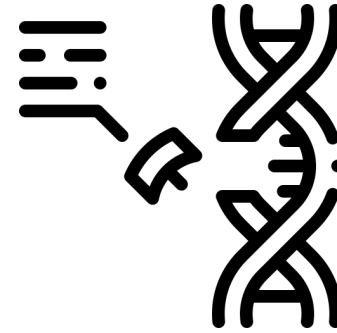
Speech to text



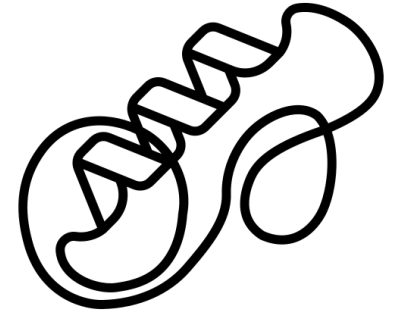
Predictive text



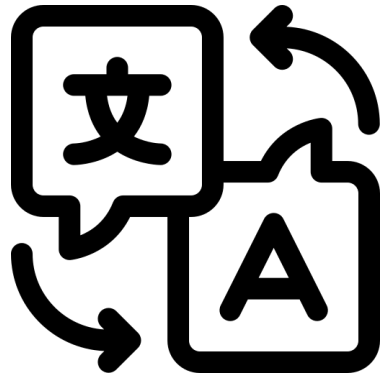
Protein family ID



Gene prediction



Protein folding



Translation



Financial modeling

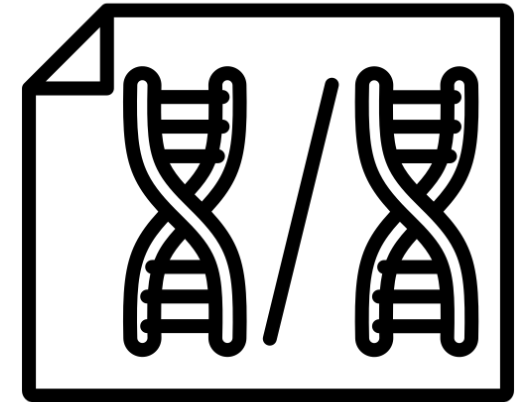
**Hidden Markov
models**



Music recognition



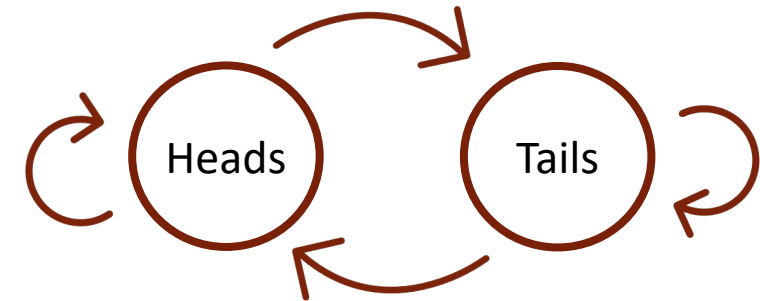
Copy number calling



Sequence alignment

Markov models

- A family of models for stochastic processes
- **Key 1:** Represent a system as a set of **states** and **transitions** between the states
 - The simplest Markov model is the Markov chain
- **Key 2:** The **Markov property** says the future state only needs depends on the present state
 - “The future only depends on today, ignore the past



Распространение закона больших чисел на величины, зависящие друг от друга.

Закон больших чисел, в силу которого, с вероятностью сколь угодно близкою к достоверности, можно утверждать, что среднее арифметическое из нескольких величин, при достаточно большом числе этих величин, будет произвольно мало отличаться от средней арифметической из их математических ожиданий, выведенъ Чебыше-

“Extension of the law of large numbers to quantities that depend on each other”

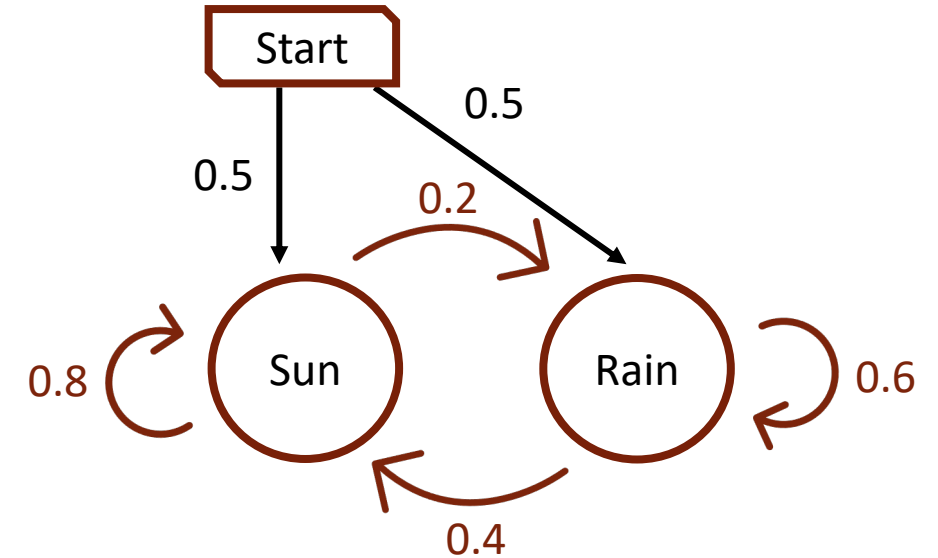


А. А. Марков (1886).

Andrey Markov*

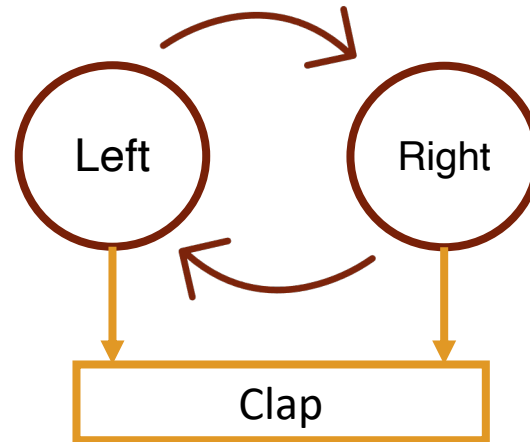
A Markov chain for the weather

- **Key 1:** Represent a system as a set of **states** and **transitions** between the states
 - We can model the weather from day to day as transitions between sunny and rainy days
- **Key 2:** The **Markov property** says the future state only depends on the present state
 - The weather tomorrow depends only on whether it's sunny or rainy today
- The **transition probabilities** describe the probability of moving from state to state
- The **initial state probabilities** describe the probability of starting in each state



Hiding a Markov chain

- An extension of Markov chains in the 1960s
- We can't see the states (“**hidden**”)
- We can see **emissions** or **observations** from the states



STATISTICAL INFERENCE FOR PROBABILISTIC FUNCTIONS OF FINITE
STATE MARKOV CHAINS

BY LEONARD E. BAUM AND TED PETRIE

Institute for Defense Analyses, Princeton, N. J.

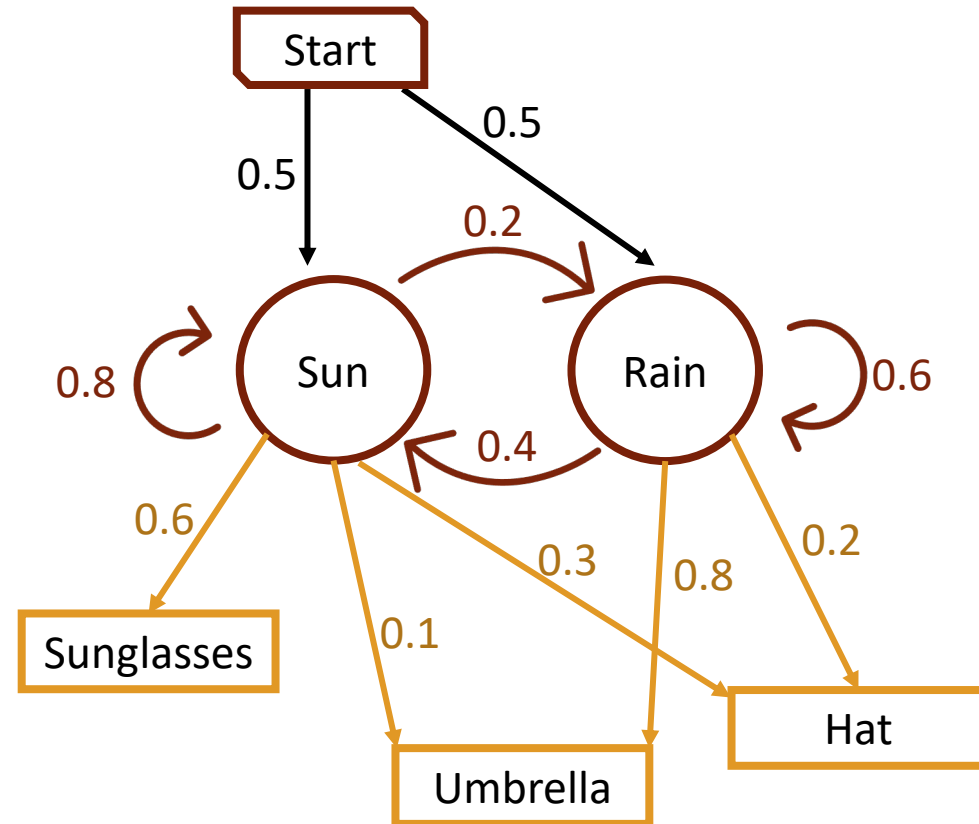
Let $\{X_t\}$ be an s state Markov process, generated by some $s \times s$ stochastic matrix $\{a_{ij}\}$ with positive entries. Let $\{Y_t\}$ be a *probabilistic function* of $\{X_t\}$, viz:

$$(0.1) \quad P\{Y_t = k \mid X_t = j, Y_{t-1}, X_{t-1}, \dots\} = b_{jk}$$

where $\{b_{jk}\}$ is an $s \times r$ matrix with positive entries and row sums $= 1$.

This paper deals with statistical estimation. We assume that the matrices $A = \{a_{ij}\}$ and $B = \{b_{jk}\}$ are unknown and we wish to recover them from an observation $\{Y_1, \dots, Y_T\}$.

How can you know the weather without seeing it?



Initial probabilities

	Sun	Rain
Start	0.5	0.5

Transition probabilities

	Sun	Rain
Sun	0.8	0.2
Rain	0.4	0.6

Emission probabilities

	Sunglasses	Umbrella	Hat
Sun	0.6	0.1	0.3
Rain	0	0.8	0.2

Sunglasses

Sunglasses

Umbrella

Sunglasses

Sunglasses

Sunglasses

Sunglasses

Formal definitions and notation

- A hidden Markov model consists of random variables q_i and y_i where

$$Q = \{q_1 \ q_2 \ \dots \ q_N\}$$

Q is a set of N hidden states q_i

$$A = a_{1,1} \ a_{1,2} \ \dots \ a_{N,N}$$

s.t. $\sum_{j=1}^N a_{i,j} = 1 \ \forall i$

with transition probabilities $a_{i,j}$ describing the probability of moving from state i to state j

$$\pi = \pi_1, \pi_2, \dots, \pi_N$$

s.t. $\sum_{j=1}^N \pi_j = 1$

and initial probabilities π_i describing the probability of starting in state i .

$$Y = \{y_1 \ y_2 \ \dots \ y_M\}$$

Y is a set of M possible emissions (or observations) y_i

$$B = b_i(y_t)$$

each with an emission probability $b_i(y_t)$ of being generated from state q_i at time t

- We call a sequence of emission observations $Y = y_1, y_2, \dots y_T$ and a path of hidden states $Q = q_1, q_2, \dots q_T$, both with length T

Key assumptions for the HMM

- Markov property:

$$P(q_t = i \mid q_1, q_2, \dots, q_{t-1}) = P(q_t = i \mid q_{t-1})$$

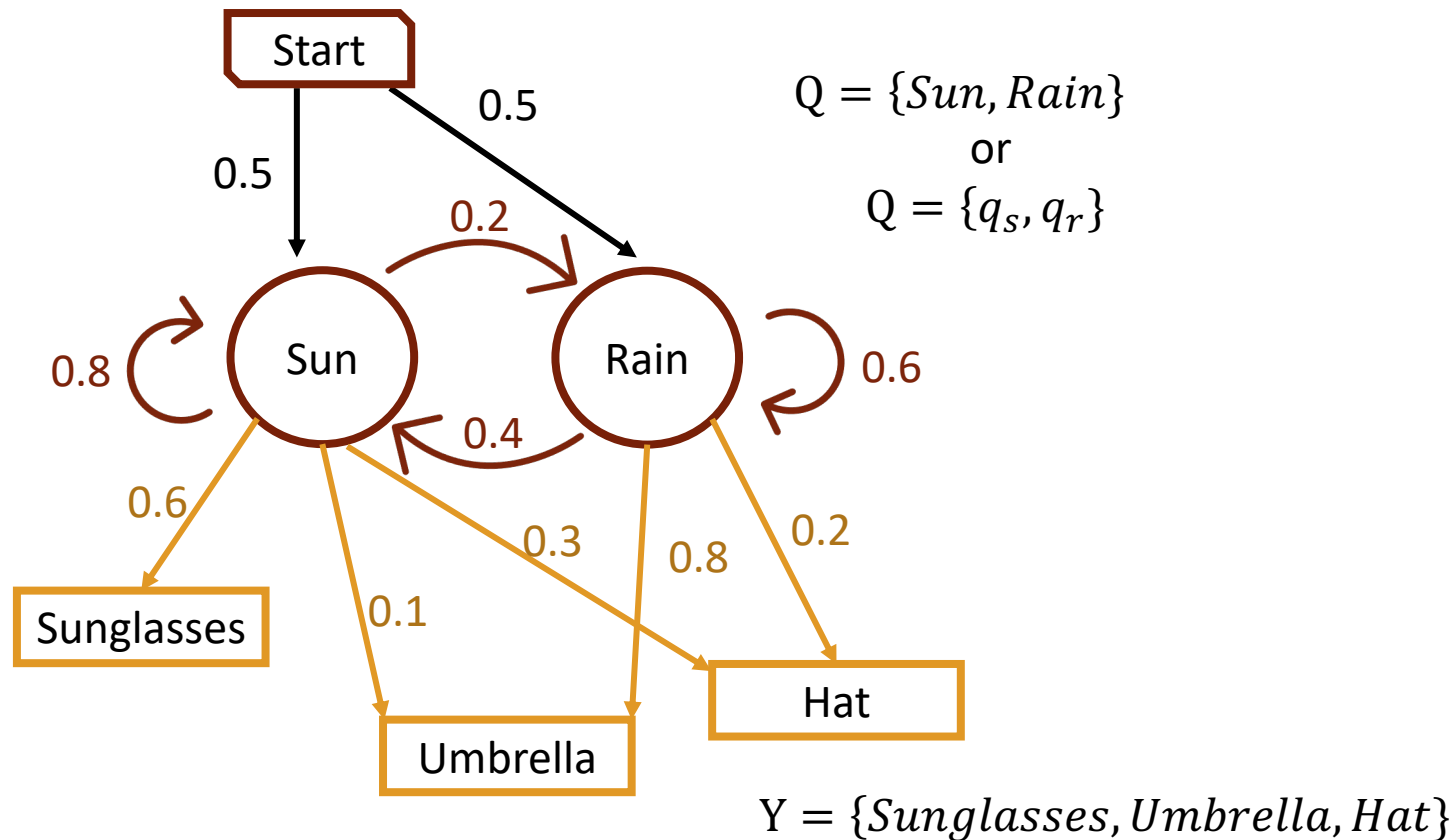
“the future only depends on the present, not the past”

- Output independence:

$$P(y_i \mid q_1, \dots, q_i, \dots, q_T, y_1, \dots, y_i, \dots, y_T) = P(y_i \mid q_i)$$

“the probability observing the emission y_i depends only on the state that produced it q_i , not any other state or emissions”

Adding labels to our weather HMM



Initial probabilities

π_i	Sun	Rain
Start	0.5	0.5

Transition probabilities

$a_{i,j}$	Sun	Rain
Sun	0.8	0.2
Rain	0.4	0.6

Emission probabilities

$b_i(y_t)$	Sunglasses	Umbrella	Hat
Sun	0.6	0.1	0.3
Rain	0	0.8	0.2

- What can we learn from our HMM?

What is the probability of...

- Observing a particular sequence of emissions?

sunglasses, hat, **umbrella**

- First, think of one possible hidden state path

sun, **sun**, rain

$$\begin{aligned} P(\text{SHU} | \text{sun sun rain}) \\ &= \pi_{\text{sun}} b_{\text{sun}}(S) * a_{\text{sun}, \text{sun}} b_{\text{sun}}(H) * a_{\text{sun}, \text{rain}} b_{\text{rain}}(U) \\ &= (0.5 * 0.6) * (0.8 * 0.3) * (0.2 * 0.8) \\ &= 0.01152 \end{aligned}$$

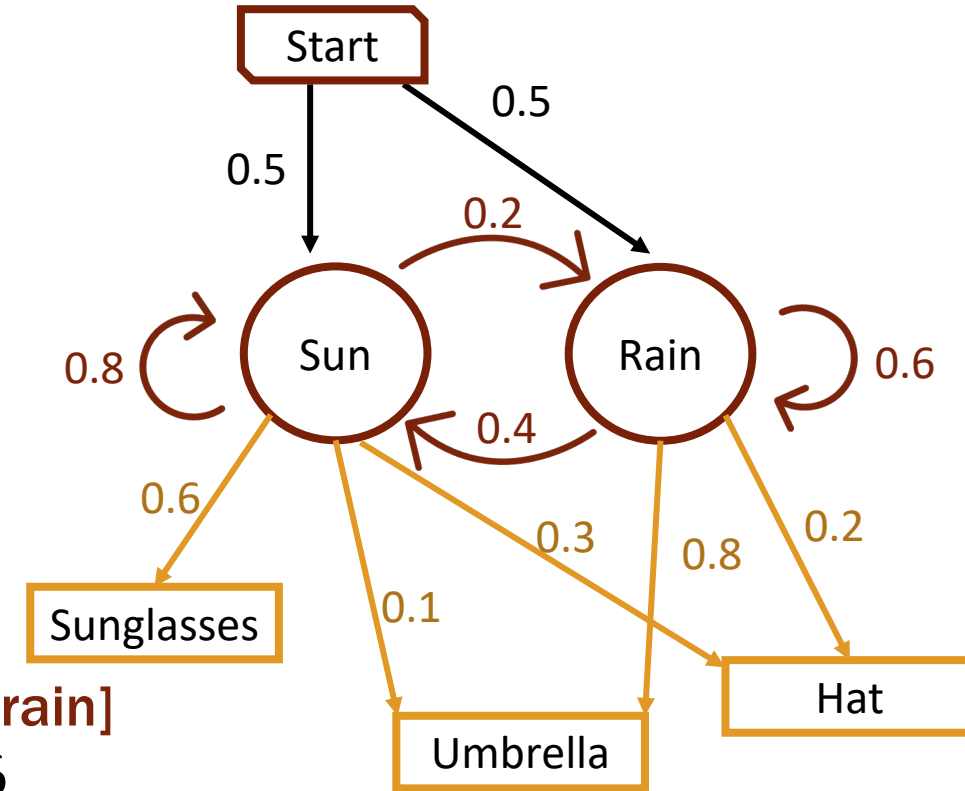
- How many possible hidden state paths are there?

[sun, sun, rain] [sun, sun, sun] [sun, rain, sun] [sun, rain, rain]

$$\begin{aligned} P(\text{SHU}) &= 0.01152 + 0.00576 + 0.00048 + 0.00576 \\ &= 0.02352 \end{aligned}$$

- What about the paths starting with rain?

- Notice: You can also identify the most likely path of hidden states



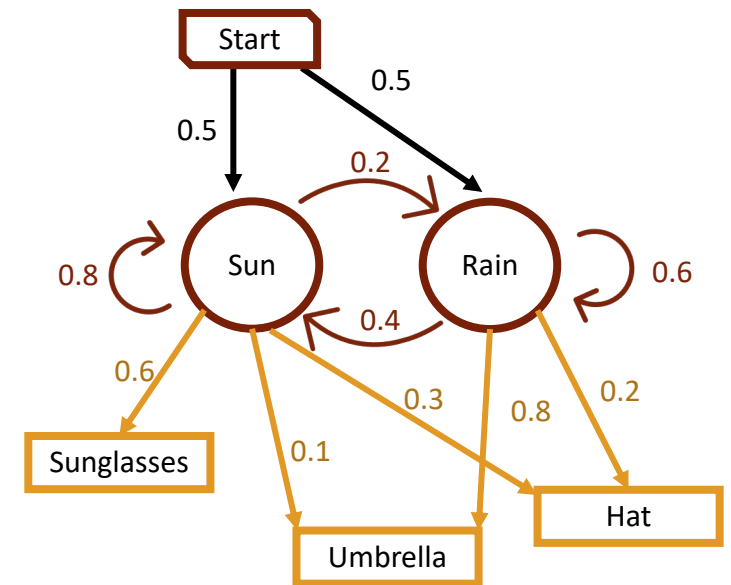
Four typical HMM inference problems

- Given the parameters of an HMM and an observed sequence of T emissions,
 - what is the probability of that observed sequence?
(i. **Scoring**)
 - what is the most likely hidden states path?
(ii. **Decoding**)
 - what is the distribution of hidden states at time k ?
(iii. **Filtering** when $k = T$; iv. **Smoothing** when $k < T$)
- These problems become increasingly complex with increasing N , M , and T
 - How complex?
- Usually not possible to compute by traversing all the paths

Some insights about state sequences

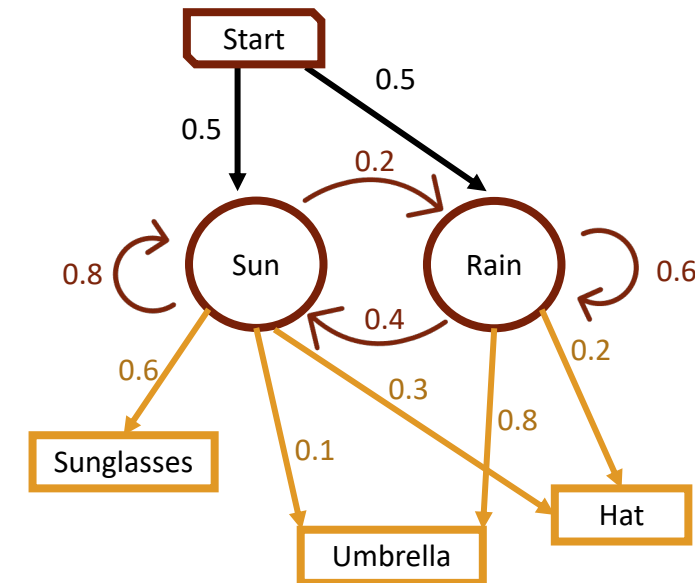
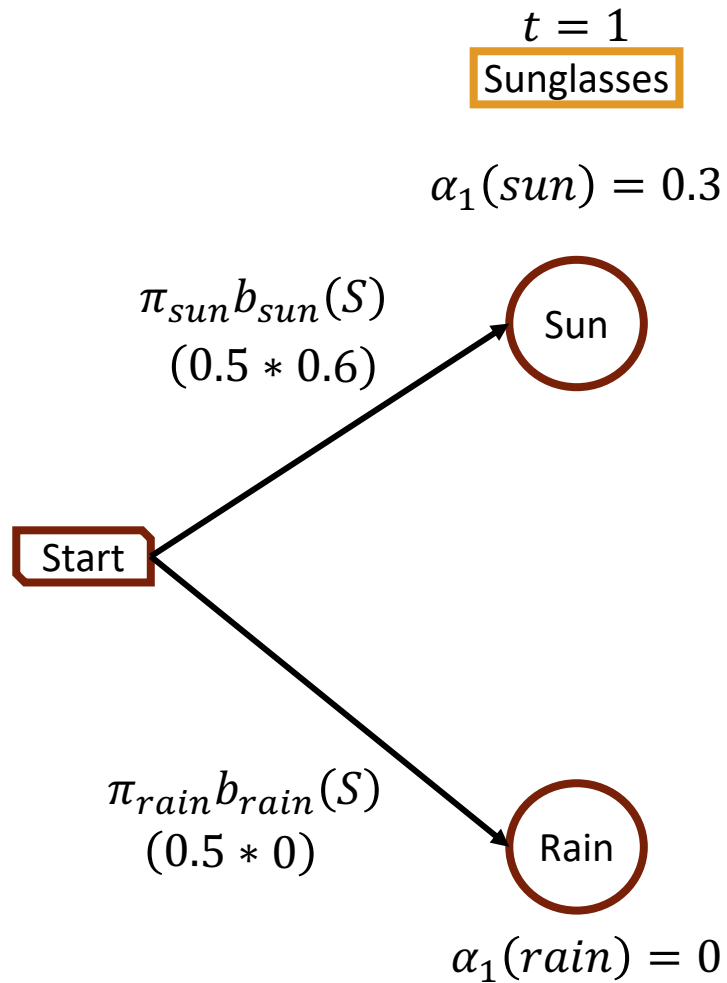
- Given emission sequence $Y = \text{sunglasses, hat, umbrella}$,
All possible paths are:

$t = 1$ $y_1 = \text{sunglasses}$	$t = 2$ $y_2 = \text{hat}$	$t = 3$ $y_3 = \text{umbrella}$
sun	sun	rain
sun	sun	sun
sun	rain	sun
sun	rain	rain
rain	sun	rain
rain	sun	sun
rain	rain	sun
rain	rain	rain

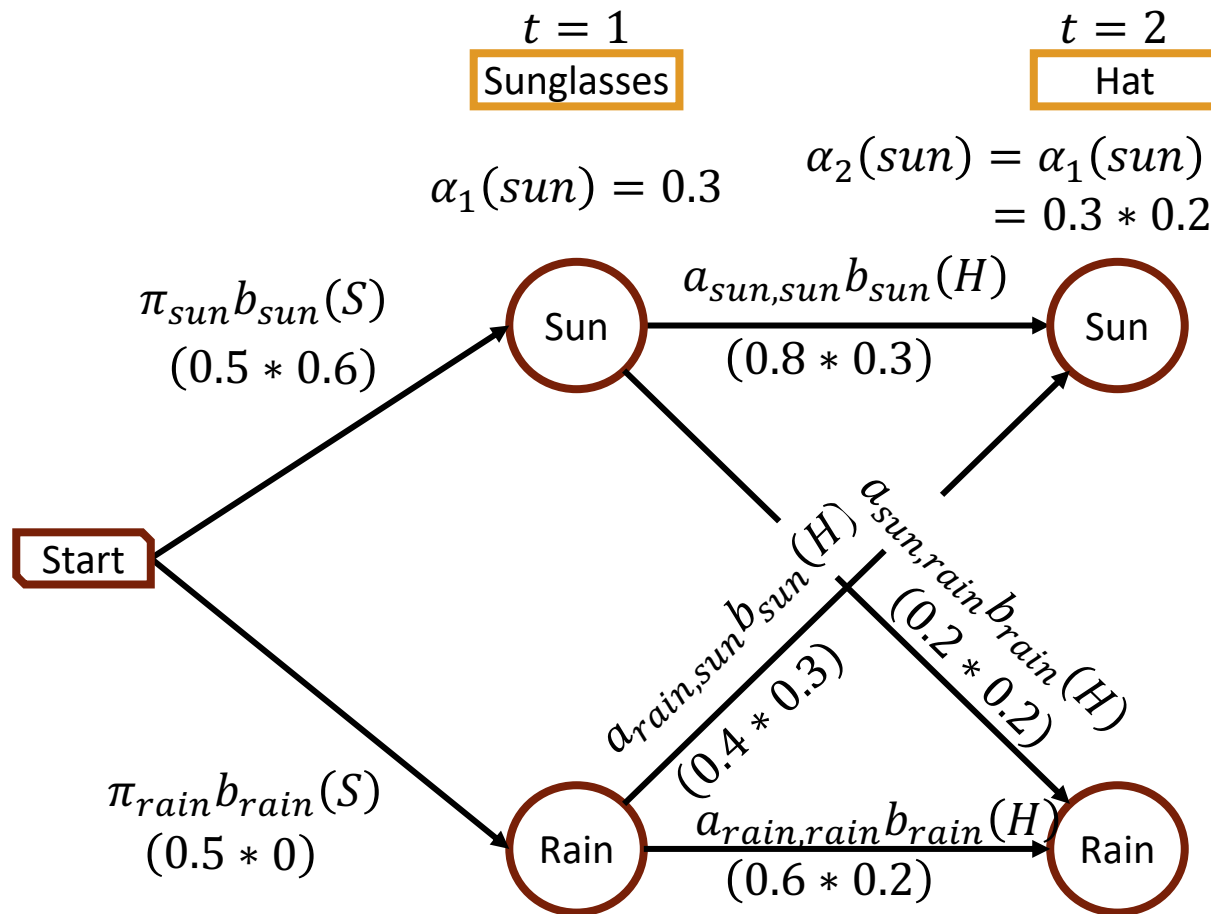


- Note: at $t = 1$, there are two unique paths
at $t = 2$, there are four unique paths
at $t = 3$, there are eight unique paths
- At $t = T$, there are N^T unique paths, but we can save time by re-using information

Another way to solve the scoring problem: the trellis



The forward trellis for scoring

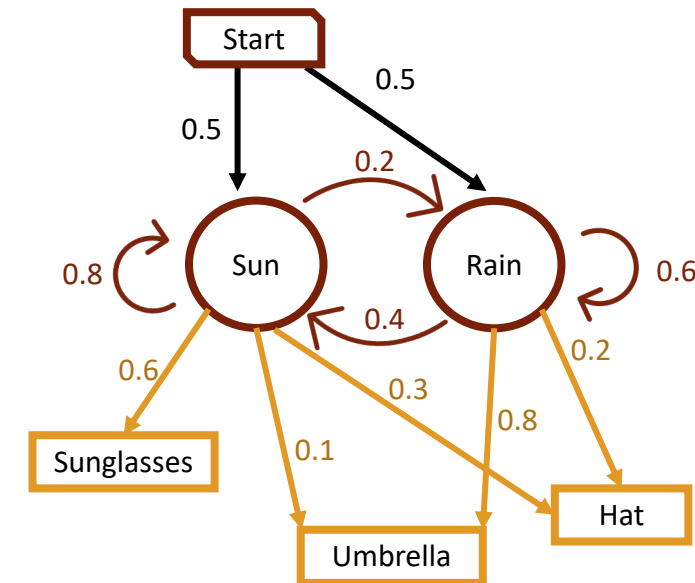


$$\alpha_1(sun) = 0.3$$

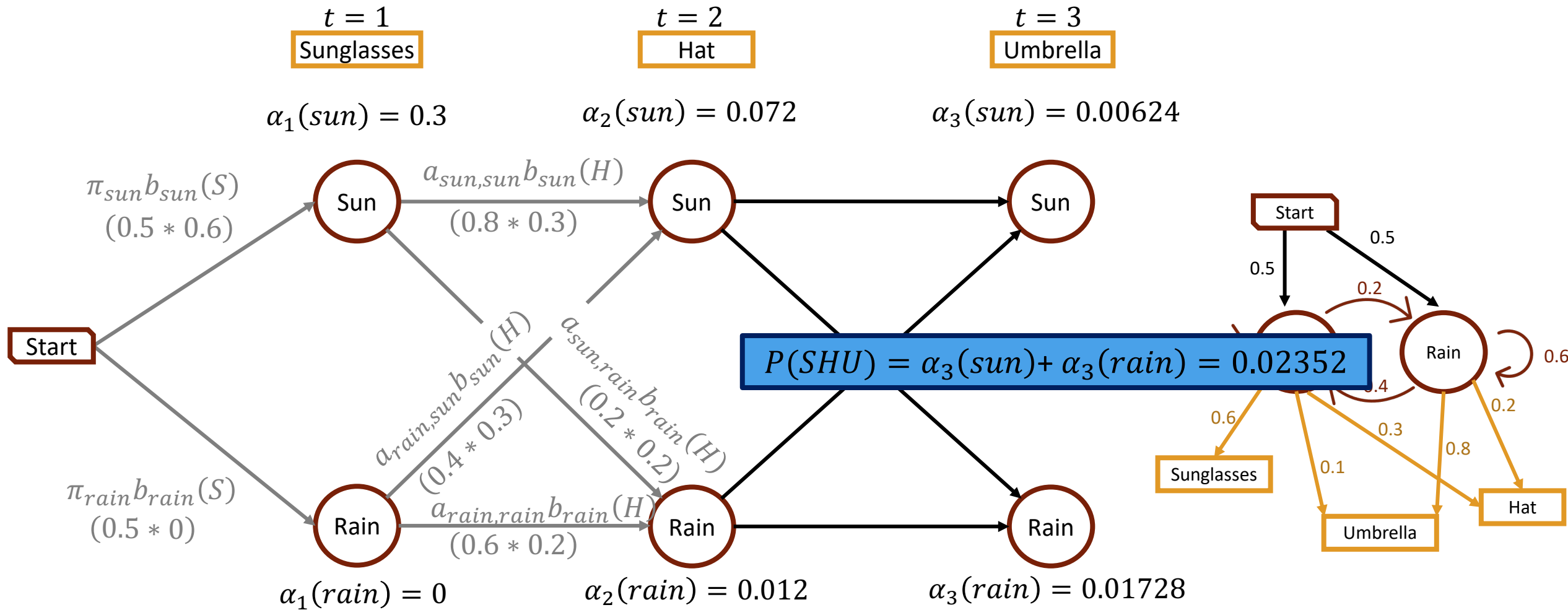
$$\alpha_2(sun) = \alpha_1(sun) * a_{sun,sun} b_{sun}(H) + \alpha_1(rain) a_{rain,sun} b_{sun}(H) = 0.3 * 0.24 + 0 * 0.12 = 0.072$$

$$\alpha_1(rain) = 0$$

$$\alpha_2(rain) = \alpha_1(sun) * a_{sun,rain} b_{rain}(H) + \alpha_1(rain) a_{rain,rain} b_{rain}(H) = 0.3 * 0.04 + 0 * 0.12 = 0.012$$



The forward trellis for scoring



Dynamic programming for HMMs

- Our trellis approach motivates the **Forward Algorithm**, which computes the **forward probabilities** $\alpha_t(j)$ at each timepoint t
- At each timepoint, we only used information from the previous timepoint
 - Excepting at the first timepoint, where we calculated $\alpha_1(sun)$ and $\alpha_1(rain)$,
 - We calculated $\alpha_2(sun)$ using information from $\alpha_1(sun)$ and $\alpha_1(rain)$
 - We calculated $\alpha_3(sun)$ using information from $\alpha_2(sun)$ and $\alpha_2(rain)$
- This is a **dynamic programming** approach using a **recursive** relationships
 - “Divide and conquer”
 - **Recursion** involves breaking up a problem into simpler versions of itself
 - We require a **base case**, **recursive step**, and a **termination step**

The Forward Algorithm

Scoring problem: Given an HMM $\Theta = (A, B)$ and an observation sequence Y , what is the likelihood $P(Y|\Theta)$?

1. Base case:

$$\alpha_1(j) = \pi_j b_j(y_1); \quad 1 \leq j \leq N$$

The previous forward probability

2. Recursion:

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(y_t); \quad 1 \leq j \leq N, 1 \leq t \leq T$$

The transition probability from q_i to q_j

3. Termination:

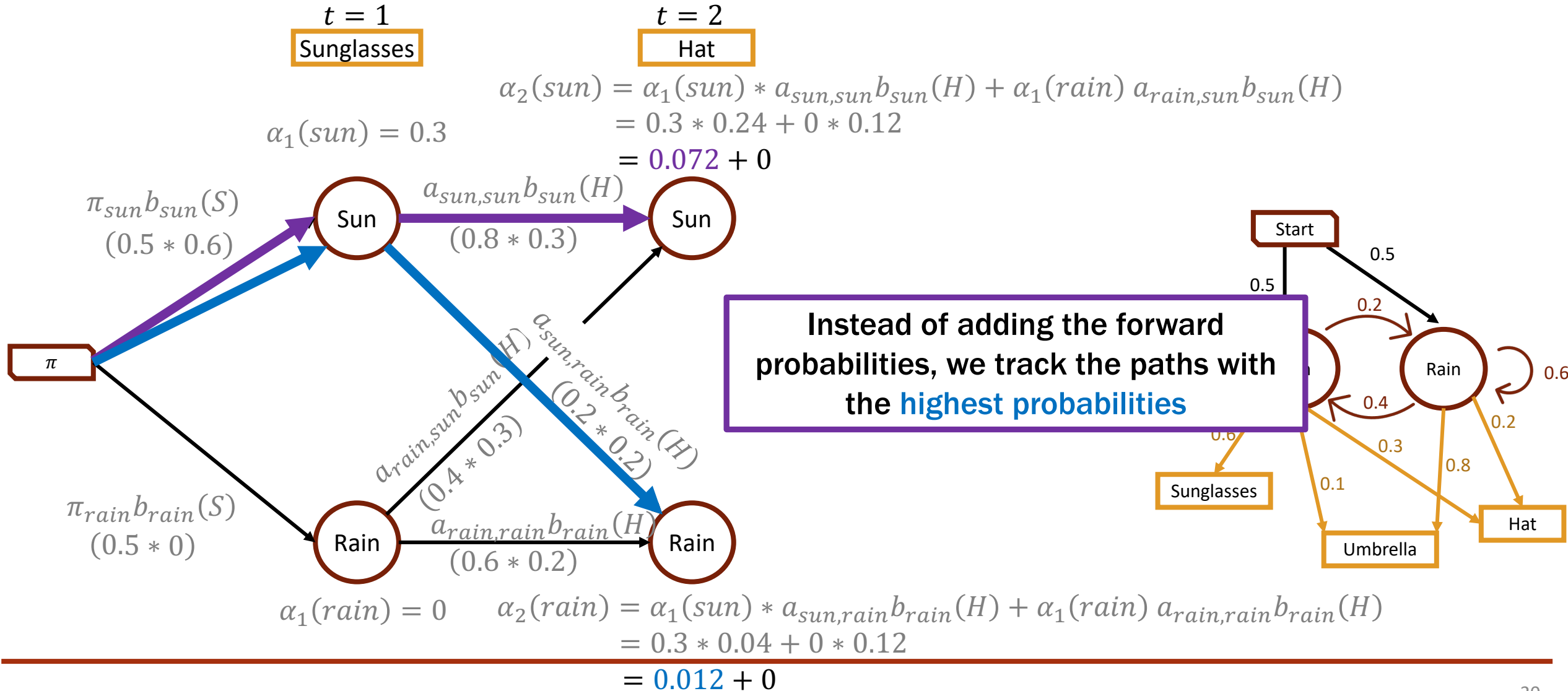
$$P(Y|\Theta) = \sum_{i=1}^N \alpha_T(i)$$

The emission probability of observation y_t given state j

Notes on the Forward Algorithm

- The Forward Algorithm is much more efficient than traversing all paths N^T
 - $O(N^2T)$ time
- Formal derivation relies on liberal application of chain rule and Markov property (see additional materials on git)
- We saw how to use it to **score** a sequence of emissions
- The Forward Algorithm is also useful in **filtering**
 - What is the most likely hidden state at the end of a sequence of emissions?
- A small adjustment allows us to identify the most likely underlying hidden state sequence

The trellis again, but we don't take the sum



The Viterbi Algorithm

Decoding problem: Given an HMM $\Theta = (A, B)$ and an observation sequence Y , what is the most likely sequence of hidden states?

Forward algorithm

1. Base case:

$$v_1(j) = \pi_j b_j(y_1); \quad 1 \leq j \leq N$$
$$bt_1(j) = 0; \quad 1 \leq j \leq N$$

$$\alpha_1(j) = \pi_j b_j(y_1); \quad 1 \leq j \leq N$$

2. Recursion:

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(y_t); \quad 1 \leq j \leq N, 1 \leq t \leq T$$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(y_t); \quad 1 \leq j \leq N, 1 \leq t \leq T$$

$$bt_t(j) = \operatorname{argmax}_{i=1}^N v_{t-1}(i) a_{ij} b_j(y_t); \quad 1 \leq j \leq N, 1 \leq t \leq T$$

2. Termination:

$$P^* = \max_{i=1}^N v_T(i)$$

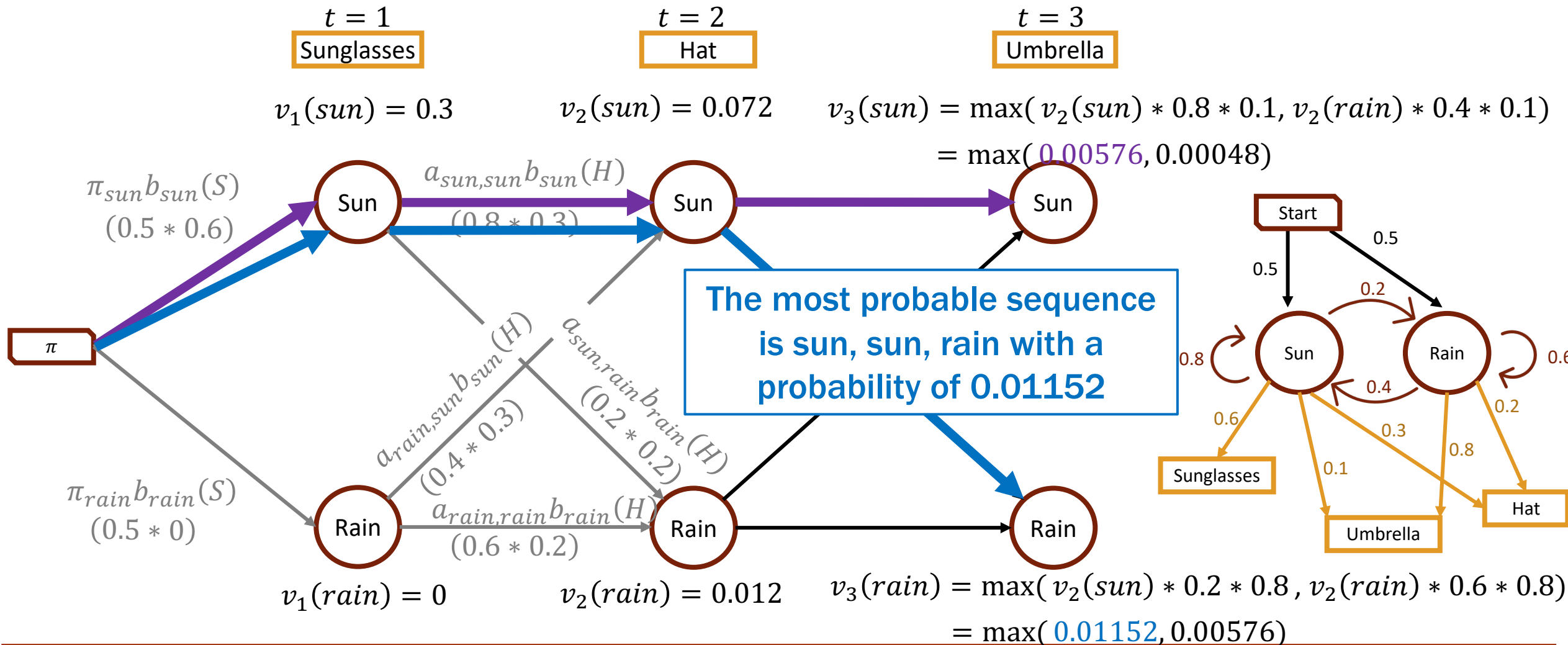
Probability of the best path

$$P(Y|\Theta) = \sum_{i=1}^N \alpha_T(i)$$

$$q_T^* = \operatorname{argmax}_{i=1}^N v_T(i)$$

Backtrace to the best path

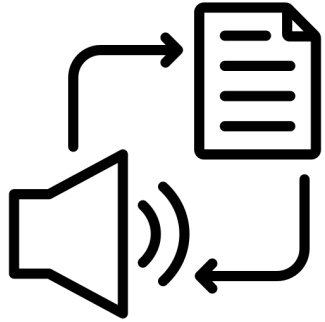
The Viterbi trellis for decoding



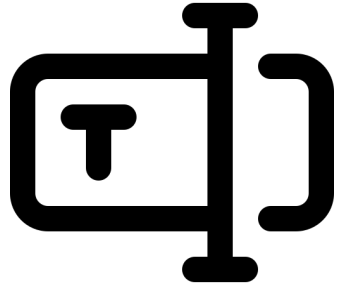
Four typical HMM inference problems and one more

- Given the parameters of an HMM and an observed sequence of emissions,
 - what is the probability of that observed sequence?
 - i. **Scoring** → **Forward Algorithm, Backward Algorithm**
 - what is the most likely sequence of hidden states?
 - ii. **Decoding** → **Viterbi Algorithm**
 - what is the distribution of hidden states at time k ?
 - iii. **Filtering** when $t = T$ → **Forward Algorithm**
 - iv. **Smoothing** when $t < T$ → **Forward-Backward Algorithm**
- How do we parameterize an HMM to begin with?
 - Given all hidden states and emission sequences, how do we compute π, A and B ?
 - v. **Parameterization** → **Baum-Welch Algorithm**
 - Baum-Welch is an iterative method that uses the forward-backward algorithm

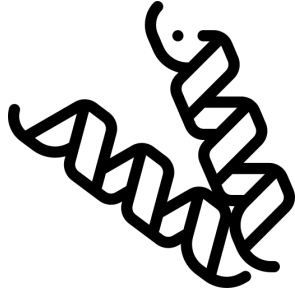
Why do we care about HMMs?



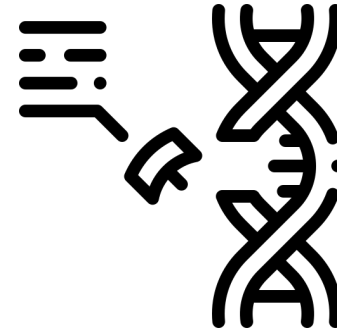
Speech to text



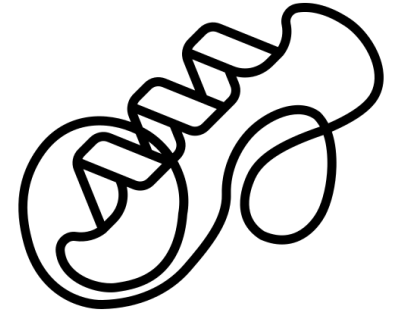
Predictive text



Protein family ID

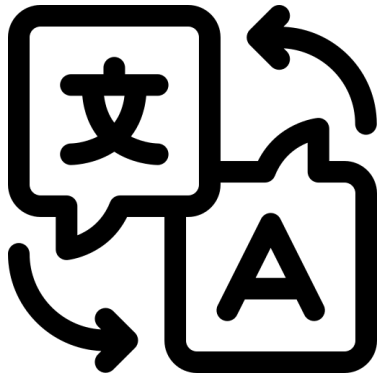


Gene prediction



Protein folding

Any system that involves change between states
Where we observe data emitted from those states



Translation



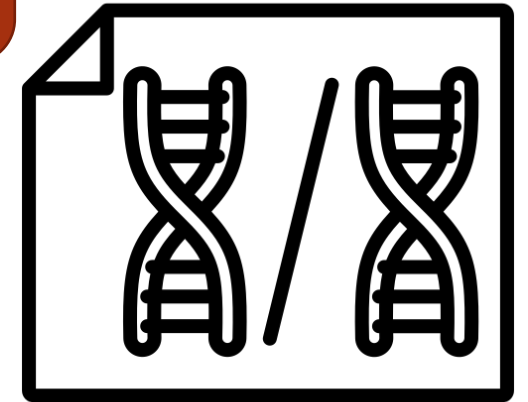
Financial modeling



Music recognition



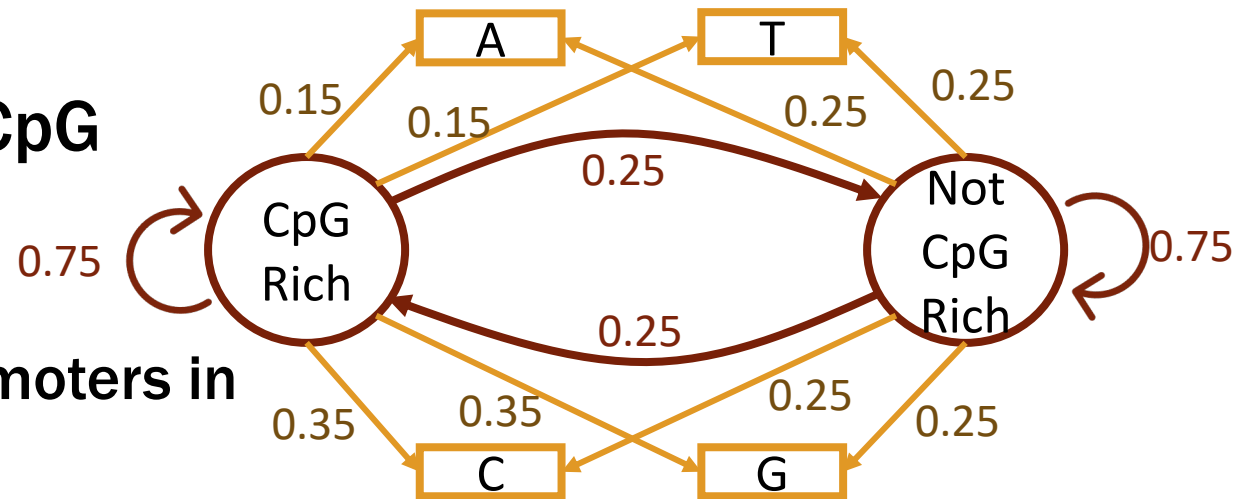
Copy number calling



Sequence alignment

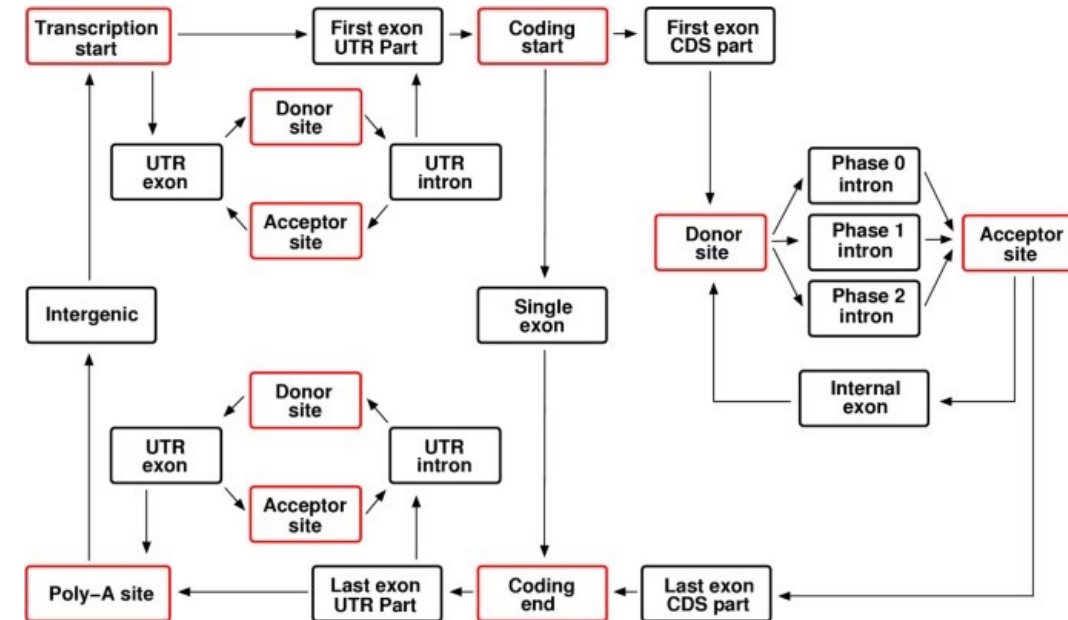
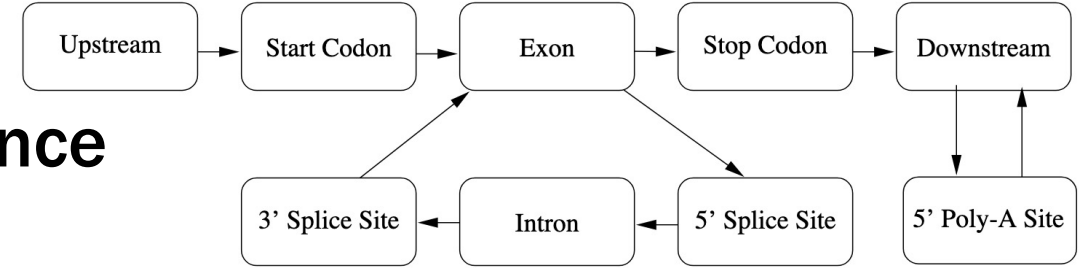
An HMM for CpG islands

- CpG sites are dinucleotide pairings of 5'-cytosine - phosphate - guanine - 3'
- CpG islands have a high frequency of CpG sites
 - Associated with gene promoters
 - Almost all housekeeping genes have promoters in CpG islands
- Finding CpG islands can help us identify genes
- Given a sequence CTATAGCGCGCATCAATGTCTTTCGCCGTATT, where are the CG-rich regions?



Gene prediction

- Our simple HMM for CpG islands encodes underlying characteristics of a DNA sequence
- We can extend the idea to genome annotation
- What characterizes a gene?
 - What is a gene?
 - Where is the start? The end?
 - What differentiates exon from intron?
- After learning the parameters, we can compute the most likely hidden state sequence and identify genic regions



Sequence alignment (especially proteins)

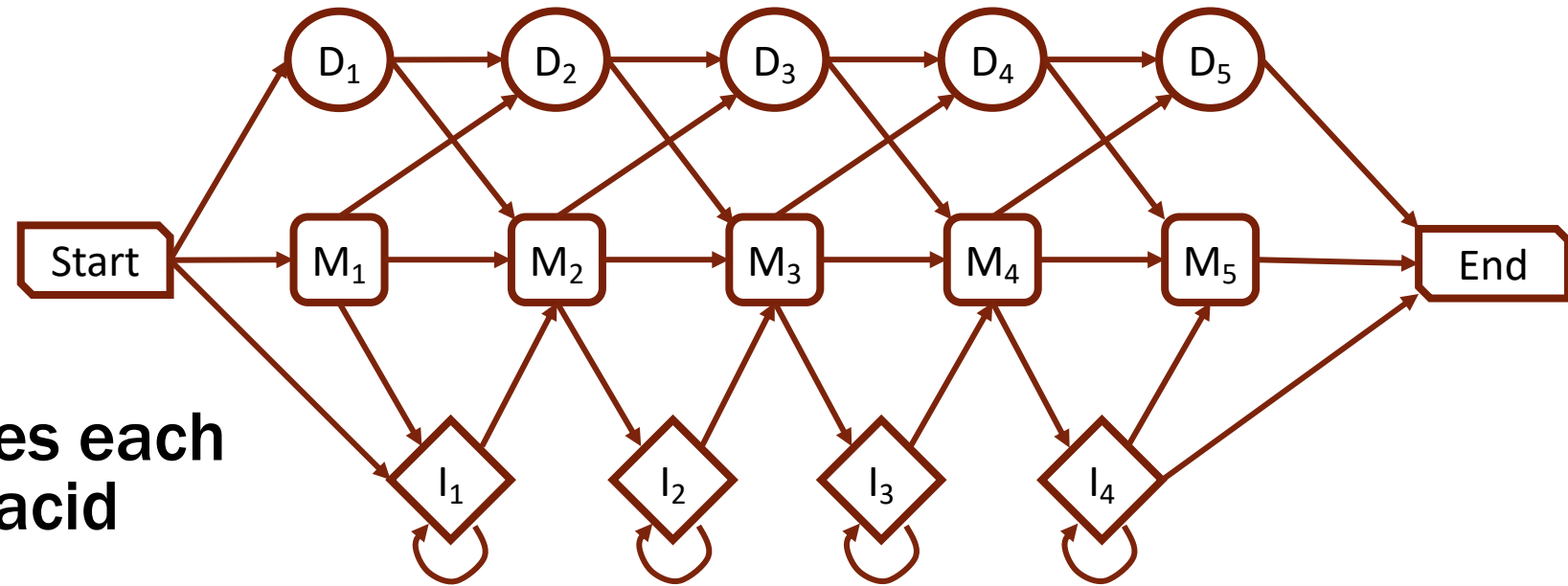
- Given a protein, identify its **protein family**
 - A group of proteins that share an evolutionary origin
 - They share similarities in function, structure, and **sequence**
- We can construct an HMM describing sequence patterns in a protein family
 - **Profile HMM**
 - Family protein sequences can be aligned or not aligned (learn by Baum-Welch EM)
- Given a new protein sequence, we can check its family membership and the most likely alignment
- Pfam is a database of profile HMMs
 - 21,979 families as of 2024

Protein sequence alignment

- Given a set of protein sequences, we can formulate an HMM

Aligned sequences

Seq1	H	N	Y	-	H	S
Seq2	H	H	Y	-	H	G
Seq3	N	H	Y	-	-	S
Seq4	T	N	Y	g	F	S
Seq5	N	G	Y	-	H	G



- Match and Insertion states each have 20 possible amino acid emissions
- Deletion states are silent
- We can use the aligned sequences as training data for parameterization

Key take homes

- The Hidden Markov model is a powerful tool for modeling systems we cannot directly observe
- We use dynamic programming to efficiently solve inference problems
 - Can we can use a trellis approach for scoring, decoding, and filtering
 - We have intuition for the Forward and Viterbi algorithms
- HMMs are common in bioinformatics
 - Genome annotation (including CpG island and gene identification)
 - Copy number calling
 - And especially sequence alignment

Additional resources

- https://github.com/conniehli/HMM_Materials
 - More math on the Forward and Backward algorithms
 - Self study exercises on the Forward-Backward and Baum-Welch algorithms
 - Exercises expanding on our protein family profile HMM
- Interested in HMMs in sequence alignment?
 - *Biological sequence analysis*, Cambridge University Press Durbin, Eddy, Krogh, & Mitchison (1998)
 - *Pfam: The protein families database in 2021*, Nucleic Acids Research, Mistry et al (2021)
- Curious about gene finding?
 - *Finding genes in DNA with a Hidden Markov Model*, Henderson, Salzberg & Fasman, Journal of Computational Biology (1997)
 - *Using database matches with HMMGene for automated gene detection in Drosophila*, Krogh, Genome Research (2000)

Thanks!

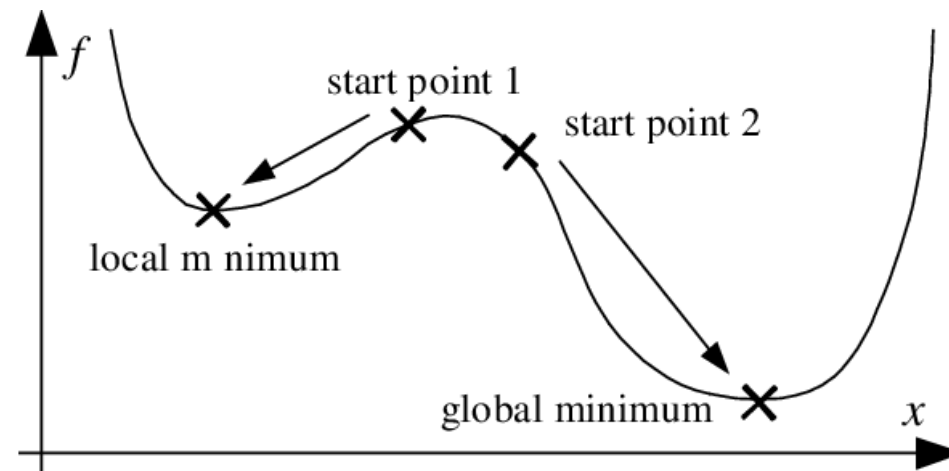
Questions?

Bonus slides

- Baum-Welch intuition
- HMMs in CNA calling

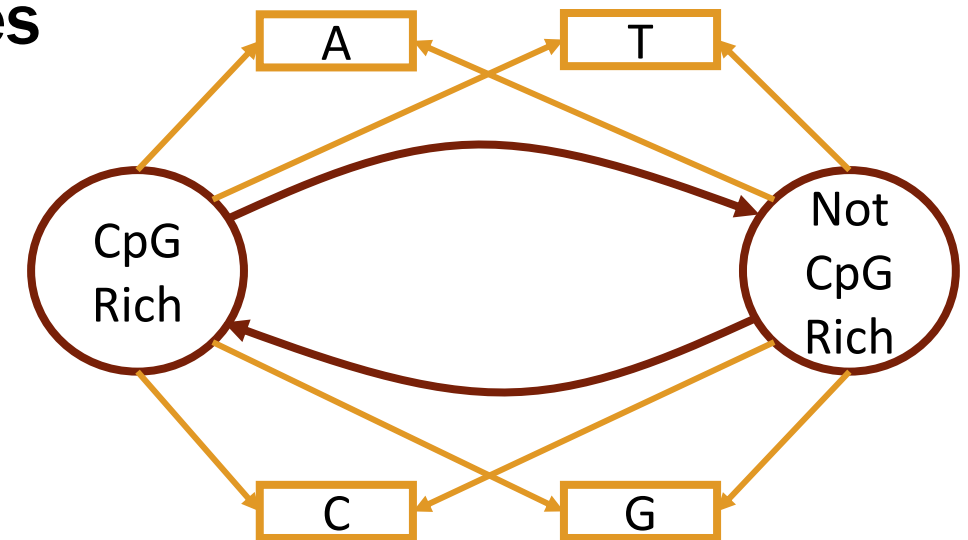
Notes on the Baum-Welch Algorithm

- The Baum-Welch is an **expectation-maximization** algorithm
 - An iterative method that converges on a local optimum
- It combines concepts from the forward and backward algorithms
- The outputs are estimates for initial probabilities π , transition probabilities A and emission probabilities B to characterize the HMM $\Theta = (A, B, \pi)$



An HMM for CpG islands

- CpG sites are dinucleotide pairings of 5'- cytosine - phosphate - guanine - 3'
- CpG islands have a high frequency of CpG sites
 - Associated with gene promoters
 - Almost all housekeeping genes have promoters in CpG islands
- Finding CpG islands can help us identify genes

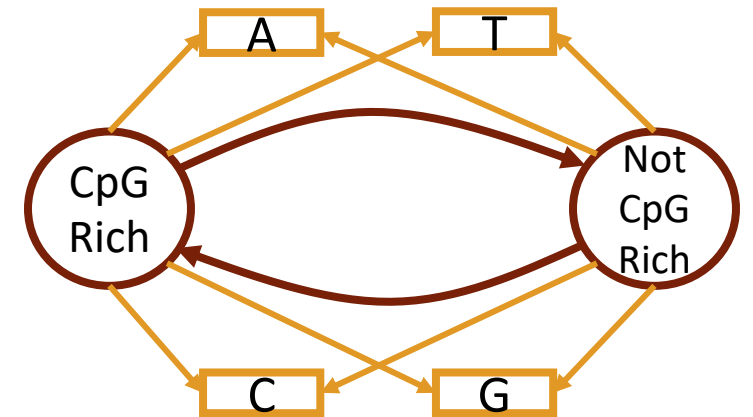


- What are reasonable guesses for (A, B, π) ?

Intuition on the Baum-Welch Algorithm

- Given a **DNA sequence** (emission sequence), we can guess at the **underlying hidden states**

A G T G G A T G C T G A C G C G C G C G C G C G C
N N N N N N N N N N N N N N N N N R R R R R R R R R R
G C G C G C G C G C A T G A T T A A G C G T A C C T C
R R R R R R R R R R R R R R R R R N N N N N N N N N N
A T C T C A C C A A T A C A T A G A G A G A G T A C A
N N N N N N N N N N N R R R R R R N N N N N N N N N N
T A T C G C G C G C G C G C G C G C G C G C G C G C G
N N N R

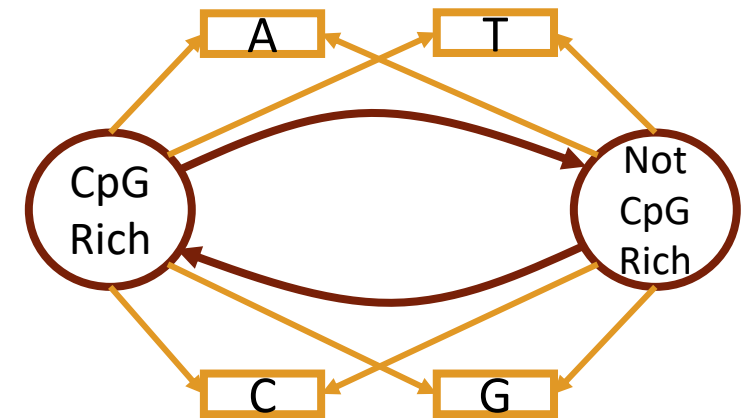


- We can estimate $(\hat{A}, \hat{B}, \hat{\pi})$ based on this guess
- This is our initialization of the algorithm

Intuition on the Baum-Welch Algorithm

1. Using estimated $(\hat{A}, \hat{B}, \hat{\pi})$, we look for the most likely hidden state path

A G T G G A T G C T G A C G C G C G C G C G C G C
N
G C G C G C G C G C A T G A T T A A G C G T A C C T C
R R R R R R R R R R R R R R R R N N N N N N N N N N N
A T C T C A C C A A T A C A T A G A G A G A G T A C A
N
T A T C G C G C G C G C G C G C G C G C G C G C G C G
N N N R



2. Based on the new estimated hidden sequence, re-estimate $(\hat{A}, \hat{B}, \hat{\pi})$
3. Repeat 1 - 2 until convergence

Detecting copy number changes

- Normal human genome is diploid
- An individual may:
 - inherit variation in the copy number of a gene (copy number variant; CNV)
 - acquire copy number changes in lifetime (copy number alteration; CNA)
 - CNAs frequently observed and implicated in cancer
- We are interested in calling copy number changes from DNA microarray and sequencing data

