

# Hidden Markov Models

**Andrew W. Moore**

**Carnegie Mellon University**

[www.cs.cmu.edu/~awm/tutorials](http://www.cs.cmu.edu/~awm/tutorials)

with extensions from

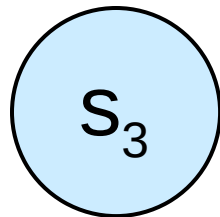
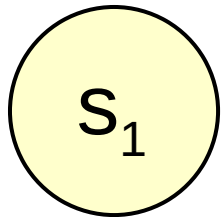
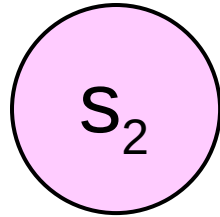
**Norbert Fuhr**

**University of Duisburg-Essen**

# A Markov System

Has  $N$  states, called  $s_1, s_2 \dots s_N$

There are discrete timesteps,  
 $t=0, t=1, \dots$



$N = 3$

$t=0$

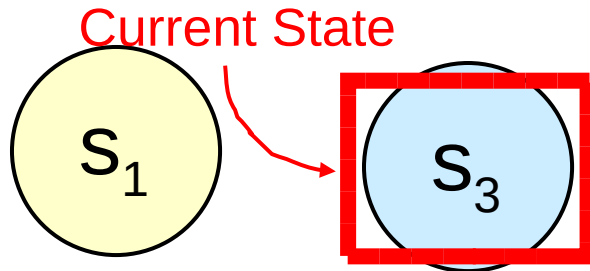
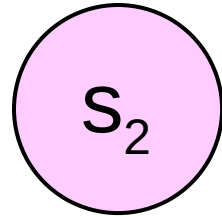
# A Markov System

Has  $N$  states, called  $s_1, s_2 \dots s_N$

There are discrete timesteps,  
 $t=0, t=1, \dots$

On the  $t$ 'th timestep the system is  
in exactly one of the available  
states. Call it  $q_t$

Note:  $q_t \in \{s_1, s_2 \dots s_N\}$



$N = 3$

$t=0$

$q_t = q_0 = s_3$

# A Markov System

Has  $N$  states, called  $s_1, s_2 \dots s_N$

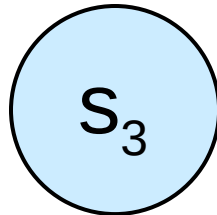
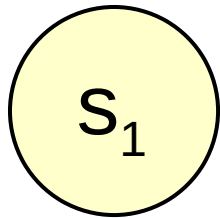
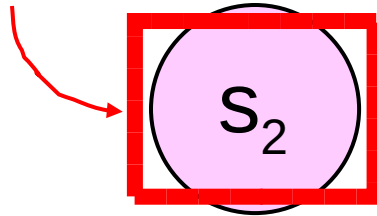
There are discrete timesteps,  
 $t=0, t=1, \dots$

On the  $t$ 'th timestep the system is  
in exactly one of the available  
states. Call it  $q_t$

Note:  $q_t \in \{s_1, s_2 \dots s_N\}$

Between each timestep, the next  
state is chosen randomly.

Current State



$N = 3$

$t=1$

$q_t = q_1 = s_2$

# A Markov System

Has  $N$  states, called  $s_1, s_2 \dots s_N$

There are discrete timesteps,  
 $t=0, t=1, \dots$

On the  $t$ 'th timestep the system is  
in exactly one of the available  
states. Call it  $q_t$

Note:  $q_t \in \{s_1, s_2 \dots s_N\}$

Between each timestep, the next  
state is chosen randomly.

The current state determines the  
probability distribution for the  
next state.

$$P(q_{t+1}=s_1|q_t=s_2) = 1/2$$

$$P(q_{t+1}=s_2|q_t=s_2) = 1/2$$

$$P(q_{t+1}=s_3|q_t=s_2) = 0$$

$$P(q_{t+1}=s_1|q_t=s_1) = 0$$

$$P(q_{t+1}=s_2|q_t=s_1) = 0$$

$$P(q_{t+1}=s_3|q_t=s_1) = 1$$

$s_2$

$s_1$

$s_3$

$N = 3$

$t=1$

$q_t=q_1=s_2$

$$P(q_{t+1}=s_1|q_t=s_3) = 1/3$$

$$P(q_{t+1}=s_2|q_t=s_3) = 2/3$$

$$P(q_{t+1}=s_3|q_t=s_3) = 0$$

# A Markov System

Has  $N$  states, called  $s_1, s_2 \dots s_N$

There are discrete timesteps,  $t=0, t=1, \dots$

On the  $t$ 'th timestep the system is in exactly one of the available states. Call it  $q_t$

Note:  $q_t \in \{s_1, s_2 \dots s_N\}$

Between each timestep, the next state is chosen randomly.

The current state determines the probability distribution for the next state.

$$P(q_{t+1}=s_1|q_t=s_2) = 1/2$$

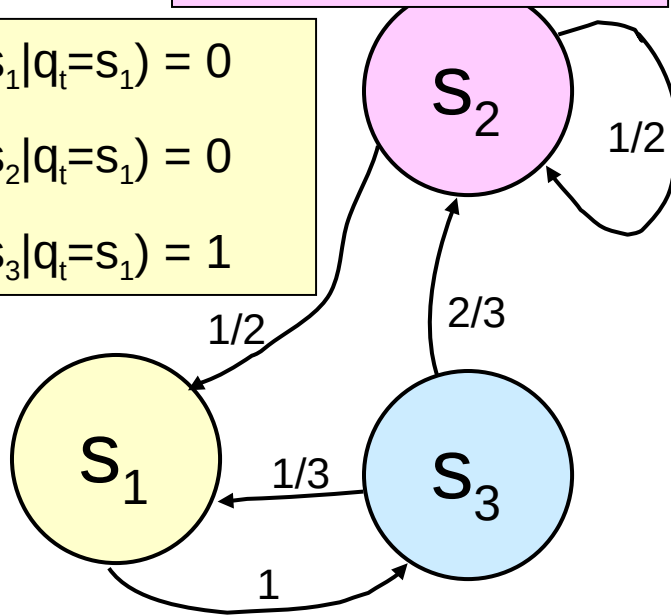
$$P(q_{t+1}=s_2|q_t=s_2) = 1/2$$

$$P(q_{t+1}=s_3|q_t=s_2) = 0$$

$$P(q_{t+1}=s_1|q_t=s_1) = 0$$

$$P(q_{t+1}=s_2|q_t=s_1) = 0$$

$$P(q_{t+1}=s_3|q_t=s_1) = 1$$



$N = 3$

$t=1$

$q_t=q_1=s_2$

$$P(q_{t+1}=s_1|q_t=s_3) = 1/3$$

$$P(q_{t+1}=s_2|q_t=s_3) = 2/3$$

$$P(q_{t+1}=s_3|q_t=s_3) = 0$$

Often notated with arcs  
between states

# Markov Property

$$P(q_{t+1}=s_1|q_t=s_2) = 1/2$$

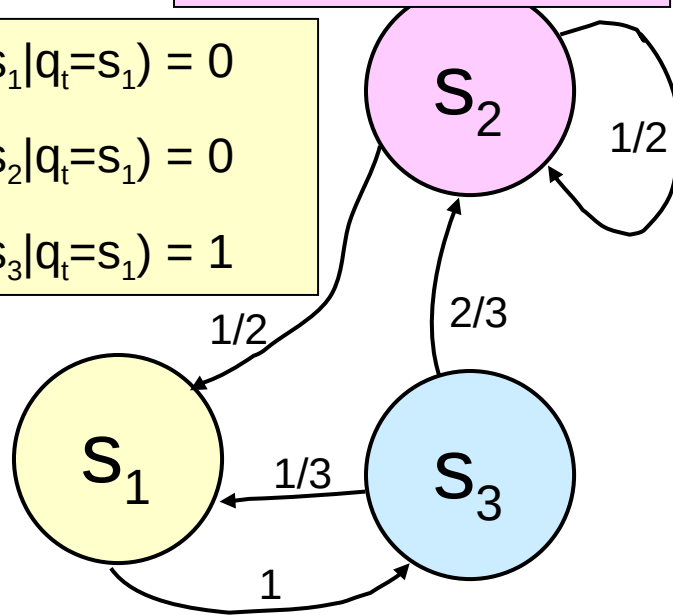
$$P(q_{t+1}=s_2|q_t=s_2) = 1/2$$

$$P(q_{t+1}=s_3|q_t=s_2) = 0$$

$$P(q_{t+1}=s_1|q_t=s_1) = 0$$

$$P(q_{t+1}=s_2|q_t=s_1) = 0$$

$$P(q_{t+1}=s_3|q_t=s_1) = 1$$



$N = 3$

$t=1$

$q_t = q_1 = s_2$

$$P(q_{t+1}=s_1|q_t=s_3) = 1/3$$

$$P(q_{t+1}=s_2|q_t=s_3) = 2/3$$

$$P(q_{t+1}=s_3|q_t=s_3) = 0$$

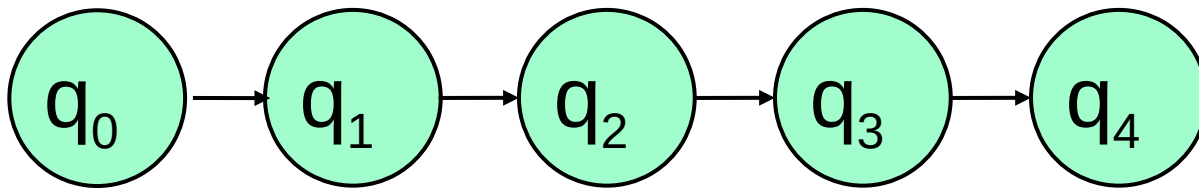
$q_{t+1}$  is conditionally independent of  $\{ q_{t-1}, q_{t-2}, \dots, q_1, q_0 \}$  given  $q_t$ .

In other words:

$$P(q_{t+1} = s_j | q_t = s_i) =$$

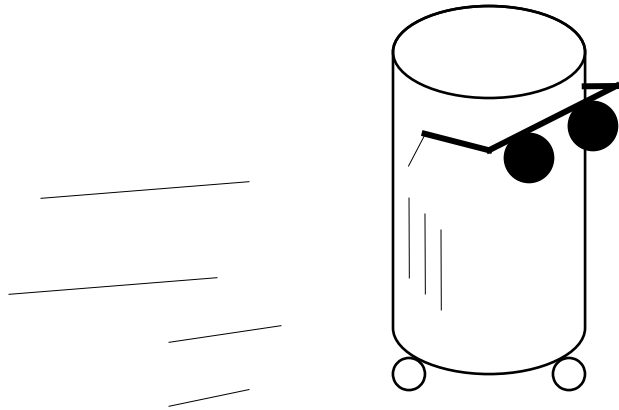
$$P(q_{t+1} = s_j | q_t = s_i, \text{any earlier history})$$

# Markov Property: Representation

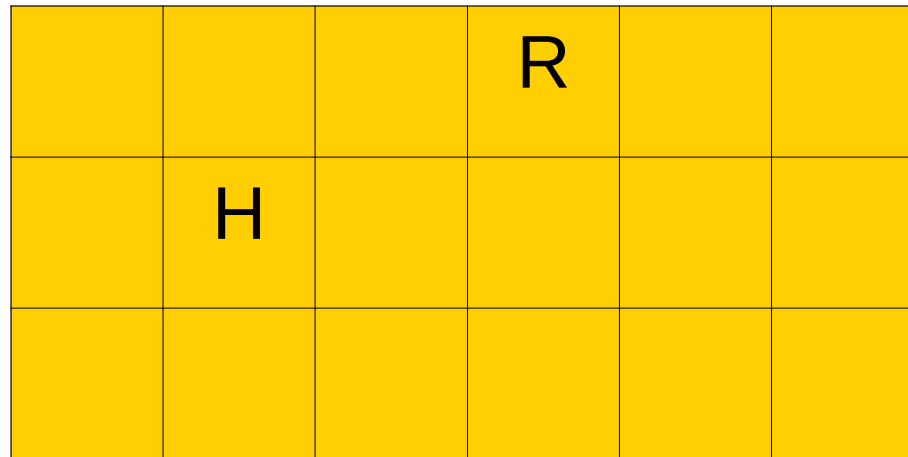




# A Blind Robot



A human and a robot wander around randomly on a grid...



**STATE  $q$  =**

Location of Robot,  
Location of Human

Note:  $N$  (num. states) =  $18^2$   
 $18 = 324$

# Dynamics of System

$q_0 =$

					R
H					

Each timestep the human moves randomly to an adjacent cell. And Robot also moves randomly to an adjacent cell.

## Typical Questions:

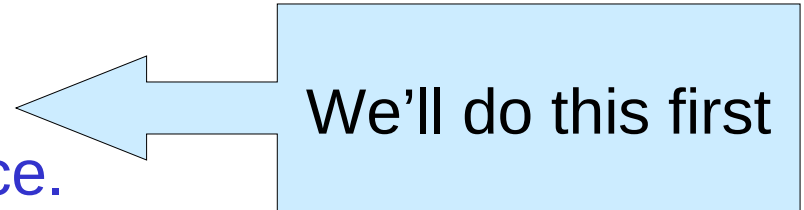
- “What’s the expected time until the human is crushed like a bug?”
- “What’s the probability that the robot will hit the left wall before it hits the human?”
- “What’s the probability Robot crushes human on next time step?”

# Example Question

“It’s currently time  $t$ , and human remains uncrushed. What’s the probability of crushing occurring at time  $t + 1$  ?”

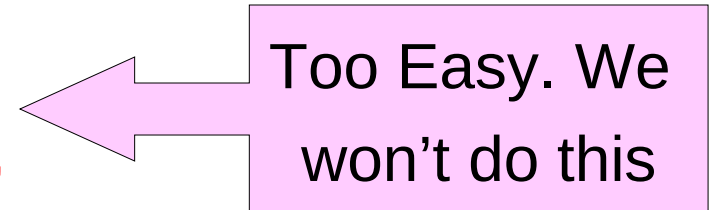
If robot is blind:

We can compute this in advance.



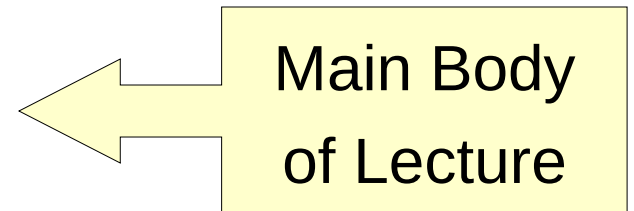
If robot is omnipotent:

(I.E. If robot knows state at time  $t$ ),  
can compute directly.



If robot has some sensors, but  
incomplete state information ...

Hidden Markov Models are  
applicable!



# What is $P(q_t = s)$ ? Too Slow

Step 1: Work out how to compute  $P(Q)$  for any path  $Q$   
 $= q_0 q_1 q_2 q_3 \dots q_t$

Given we know the start state  $q_0$

$$\begin{aligned} P(q_0 q_1 \dots q_t) &= P(q_0 q_1 \dots q_{t-1}) P(q_t | q_0 q_1 \dots q_{t-1}) \\ &= P(q_0 q_1 \dots q_{t-1}) P(q_t | q_{t-1}) \quad \text{WHY?} \\ &= P(q_1 | q_0) P(q_2 | q_1) \dots P(q_t | q_{t-1}) \end{aligned}$$

Step 2: Use this knowledge to get  $P(q_t = s)$

$$P(q_t = s) = \sum_{Q \in \text{Paths of length } t \text{ that end in } s} P(Q)$$

Computation is exponential in  $t$

# What is $P(q_t = s)$ ? Clever Answer

- For each state  $s_i$ , define
$$p_t(i) = \text{Prob. state is } s_i \text{ at time } t$$
$$= P(q_t = s_i)$$
- Easy to do inductive definition

$$\forall i \quad p_0(i) =$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

# What is $P(q_t = s)$ ? Clever answer

- For each state  $s_i$ , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

# What is $P(q_t = s)$ ? Clever answer

- For each state  $s_i$ , define

$$\begin{aligned} p_t(i) &= \text{Prob. state is } s_i \text{ at time } t \\ &= P(q_t = s_i) \end{aligned}$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \forall j \quad p_{t+1}(j) &= P(q_{t+1} = s_j) = \\ &= \sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) = \end{aligned}$$

# What is $P(q_t = s)$ ? Clever answer

- For each state  $s_i$ , define

$$p_t(i) = \text{Prob. state is } s_i \text{ at time } t \\ = P(q_t = s_i)$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j | q_t = s_i) P(q_t = s_i) = \sum_{i=1}^N a_{ij} p_t(i)$$

Remember,

$$a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$$



# What is $P(q_t = s)$ ? Clever answer

- For each state  $s_i$ , define  
 $p_t(i) = \text{Prob. state is } s_i \text{ at time } t$   
 $= P(q_t = s_i)$

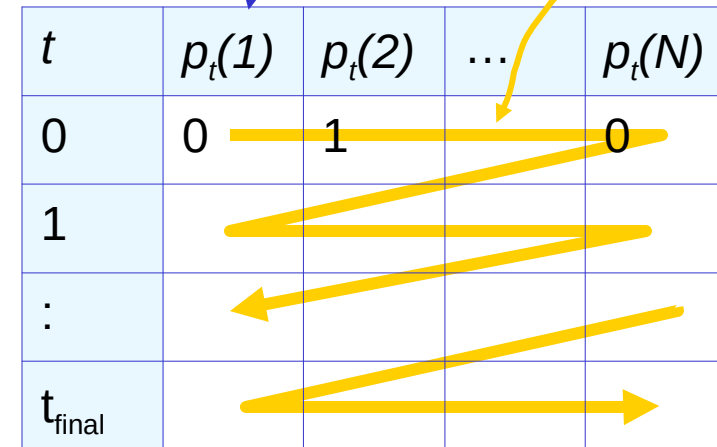
- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) = \sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j | q_t = s_i) P(q_t = s_i) = \sum_{i=1}^N a_{ij} p_t(i)$$

- Computation is simple.
- Just fill in **this** table in **this** order:



$t$	$p_t(1)$	$p_t(2)$	...	$p_t(N)$
0	0	1		0
1				
:				
$t_{\text{final}}$				

# What is $P(q_t = s)$ ? Clever answer

- For each state  $s_i$ , define

$$p_t(i) = \text{Prob. state is } s_i \text{ at time } t \\ = P(q_t = s_i)$$

- Easy to do inductive definition

$$\forall i \quad p_0(i) = \begin{cases} 1 & \text{if } s_i \text{ is the start state} \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j \wedge q_t = s_i) =$$

$$\sum_{i=1}^N P(q_{t+1} = s_j | q_t = s_i) P(q_t = s_i) = \sum_{i=1}^N a_{ij} p_t(i)$$

- Cost of computing  $P_t(i)$  for all states  $S_i$  is now  $O(t N^2)$
- The stupid way was  $O(N^t)$
- This was a simple example
- It was meant to warm you up to this trick, called **Dynamic Programming**, because HMMs do many tricks like this.

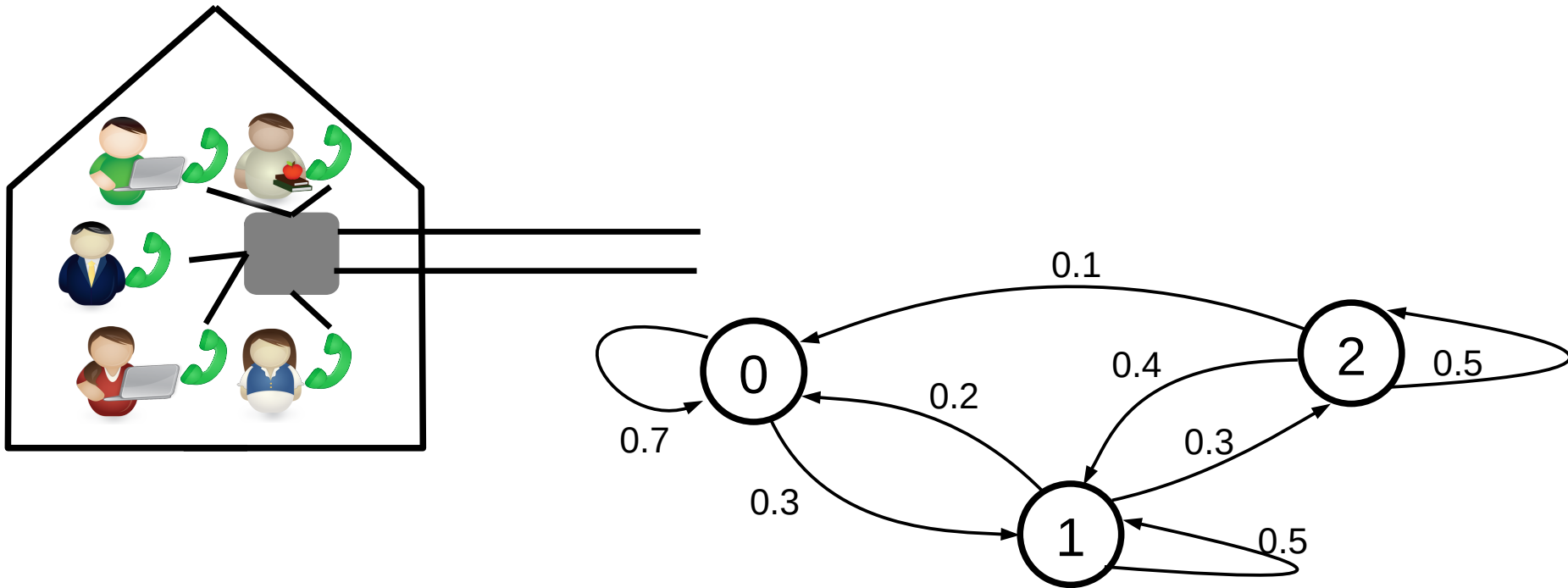
# What is $P(q_t = s)$ ?

- $p_t(i) = \text{Prob. state is } s_i \text{ at time } t = P(q_t = s_i)$
- $\mathbf{p}_t = (p_t(1), \dots, p_t(N))$ : vector of state probabilities
- $\mathbf{A} = (a_{ij})$ : Matrix of transition probabilities

$$\forall j \quad p_{t+1}(j) = P(q_{t+1} = s_j) = \sum_{i=1}^N a_{ij} p_t(i)$$

$$\mathbf{p}_{t+1} = \mathbf{p}_t \mathbf{A} = \mathbf{p}_0 \mathbf{A}^{t+1}$$

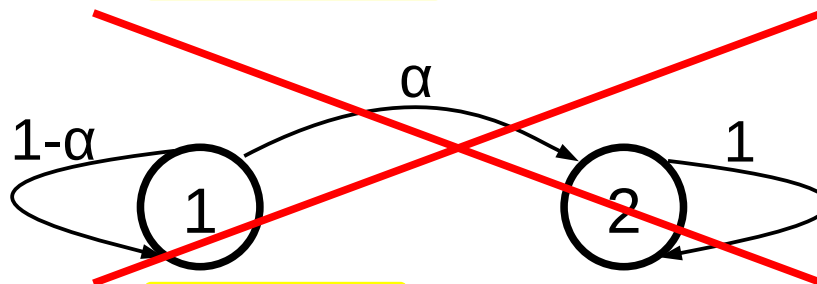
# Steady state probability



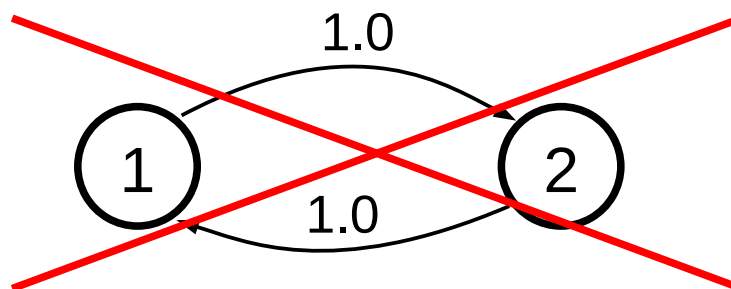
- What is the „average“ distribution over the states?
- -> steady state probability

# Steady State in MMs

- There exists a steady state, if
  - Markov model is **irreducible**, and



- Markov model is **aperiodic**

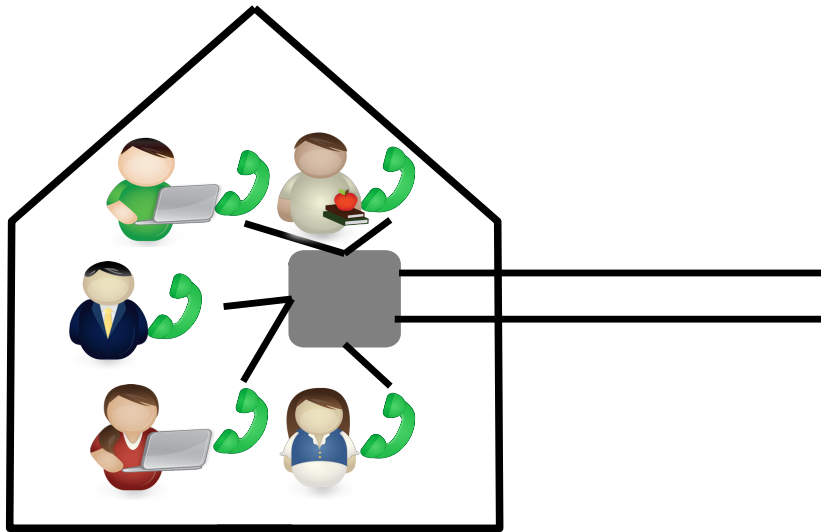


- Then there exists a steady state distribution  $\mathbf{p}_t$  that is independent of  $t$ , which we denote as  **$\pi$**

# Computing the Steady State Distribution

- $\pi = \pi A$ 
  - $\rightarrow \pi$  is left eigenvector of  $A$
  - $\rightarrow$  solve linear equation system
- Alternative method:
  - Start with arbitrary vector  $\pi_0$
  - compute iteratively  $\pi_1 = \pi_0 A$ ,  $\pi_2 = \pi_1 A$ , ...  
$$\pi_t = \pi_0 A^t$$
  - $\rightarrow$  „power method“

# Application of Power Method



$$\begin{pmatrix} 0.7 & 0.3 & 0 \\ 0.2 & 0.5 & 0.3 \\ 0.1 & 0.4 & 0.5 \end{pmatrix}$$

	$\pi_0$	$\pi_1$	$\pi_2$
$\vec{\pi}$	1.00	0.00	0.00
$\vec{\pi}.A$	0.70	0.30	0.00
$\vec{\pi}.A^2$	0.55	0.36	0.09
$\vec{\pi}.A^3$	0.47	0.38	0.15
$\vec{\pi}.A^4$	0.42	0.39	0.19
$\vec{\pi}.A^5$	0.39	0.40	0.21
$\vec{\pi}.A^6$	0.37	0.40	0.23
$\vec{\pi}.A^7$	0.36	0.40	0.23
$\vec{\pi}.A^8$	0.36	0.40	0.24
$\vec{\pi}.A^9$	0.36	0.40	0.24
$\vec{\pi}.A^{10}$	0.35	0.40	0.24
$\vec{\pi}.A^{11}$	0.35	0.41	0.24
$\vec{\pi}.A^{12}$	0.35	0.41	0.24

# Markov Modelling of Book Search

UNIVERSITÄT  
DUISBURG  
ESSEN

Open-Minded

The screenshot displays the ezDL (visitor) web application interface. The main search area on the left shows the search terms "racial discrimination" and a list of results. The results are sorted by Relevance and include book covers, titles, authors, and review counts. The top result is "American Ethnicity: The Dynamics and Consequences of Discrimination" (1997) by Jonathan H. Turner, published by McGraw-Hill College, with a rating of 3 stars and 3 reviews. Other results include "Canada's Economic Apartheid: The Social Exclusion of Racialized Groups in the New Century" (2005) by Grace-Edward Galabuzi, "Minority Voices: Linking Personal Ethnic History and the Sociological Imagination" (2004) by John P. Myers, "Only Words" (1996) by Catharine A. MacKinnon, and "Divided by Faith: Evangelical Religion and the Problem of Race in America" (2001) by Christian Smith.

The right sidebar contains a "Task" section with a "Task 1" description: "You are at an early stage of working on an assignment, and have decided to start exploring the literature in order to get an overview of your topic. Your initial idea has led to the following research need:". Below this is a "Query History" section showing two queries: "Text: racial discrimination (479) 1 min ago" and "Text: racial discrimination (0) 1 min ago".

The bottom section shows the details for the selected book, "American Ethnicity: The Dynamics and Consequences of Discrimination" (1997) by Jonathan H. Turner. The details include the price (\$54.00), publisher (Mcgraw-Hill College), ISBN (007000627X), number of pages (336), edition (2nd), Dewey (305.800973), height (75 mm), length (950 mm), width (650 mm), and weight (100 g). There are also links for "Details" and "Reviews (3)".

The bottom status bar indicates "Search finished" and "End Task".



# Eye Tracking: Areas of Interest

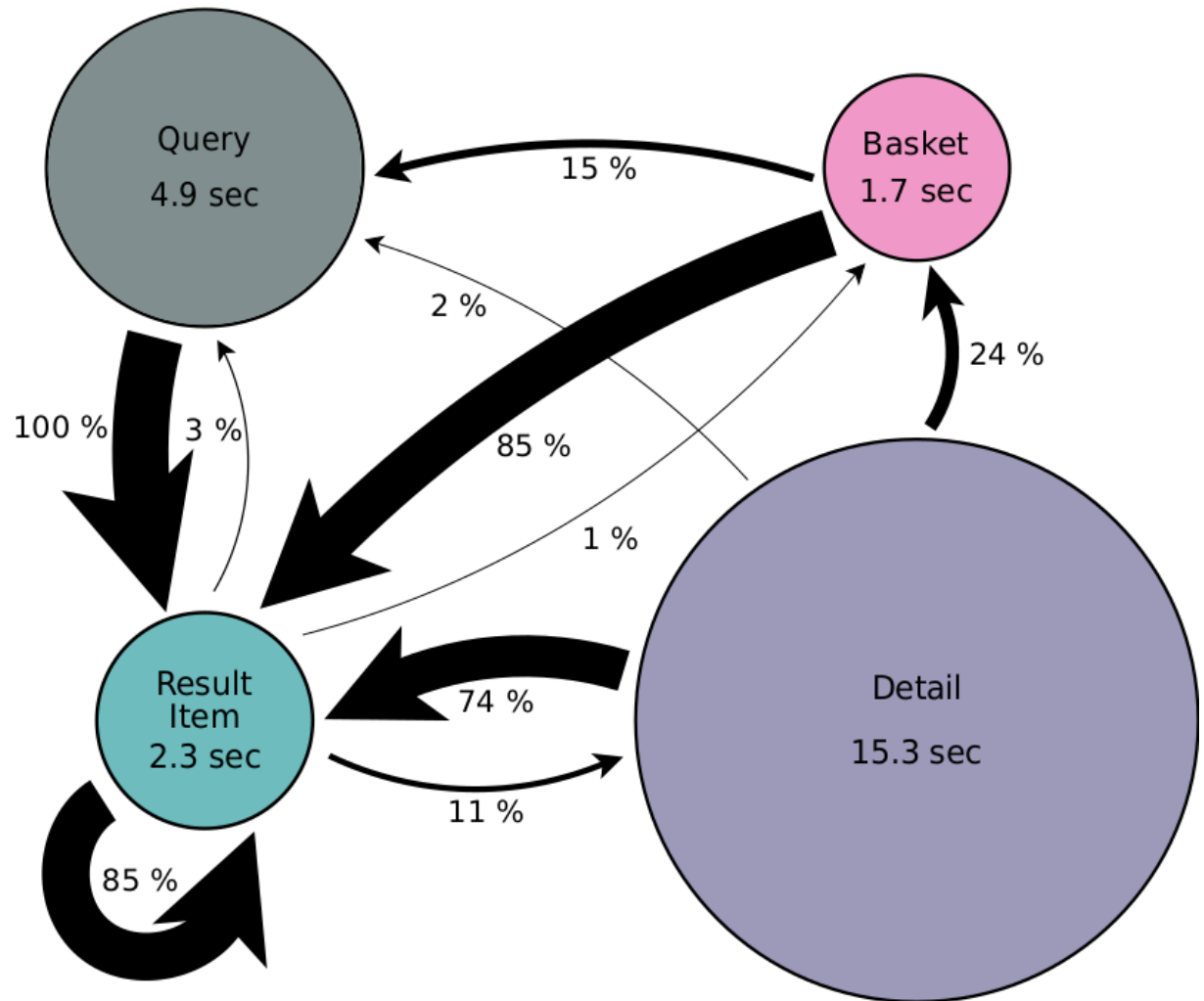
The screenshot displays a web application interface for a search tool, likely a library or research database. The interface is divided into several main sections:

- QueryView:** Located at the top left, it contains a search bar with the text "racial discrimination" and a checkbox for "search also in reviews". Below the search bar are fields for "Advanced" search criteria: "Title" (harry potter and the half-blood prince), "Author" (rowling), "Year" (2000-2005), "Abstract" (harry potter london), and "Reviews" (good introduction). A "Search" button is at the bottom right of this section.
- BasketView:** A pink rectangular area on the right side, labeled "BasketView". It contains a "Remove Selected Item" button and a "Remove All" button.
- Task:** A panel on the right side, labeled "Task". It contains a "Your current task:" section with text: "have decided to start exploring the literature in order to get an overview of your topic. Your initial idea has led to the following research need: Find highly acclaimed novels that treats issues related to racial discrimination." Below this is a "Query History" section with a list item: "Text: racial discrimination (479) 0 min ago".
- Results:** A section at the bottom left, labeled "Results". It shows a list of search results for "Text: racial discrimination". The results are displayed in a list format with columns for "ResultView", "ResultPanel", and "ResultList\_index". The results include:
  - Racial Discrimination (Issues) (Issues) (2006) by Craig Donnellan
  - Measuring Racial Discrimination (2004) by National Research Council
  - Issues on Trial - Racial Discrimination (2006) by Mitchell Young
  - Racial Slurs and Discrimination: Is It Real or Is It Imaginary (2008) by Melba Thomas
- DetailView:** A large rectangular area on the right side, labeled "DetailView". It contains a "Details" tab and a "Reviews" tab. The "Details" tab is currently selected, showing a "Click on result item to show its details." instruction.

The interface also includes a "Search finished" status bar at the bottom left and an "End Task" button at the bottom right.

# Markov Modeling

- Book search task
- 12 user sessions
- Consider eye tracking + system logs



# Time to Basket (Relevance)

$$T_q = t_q + p_{qr} T_r$$

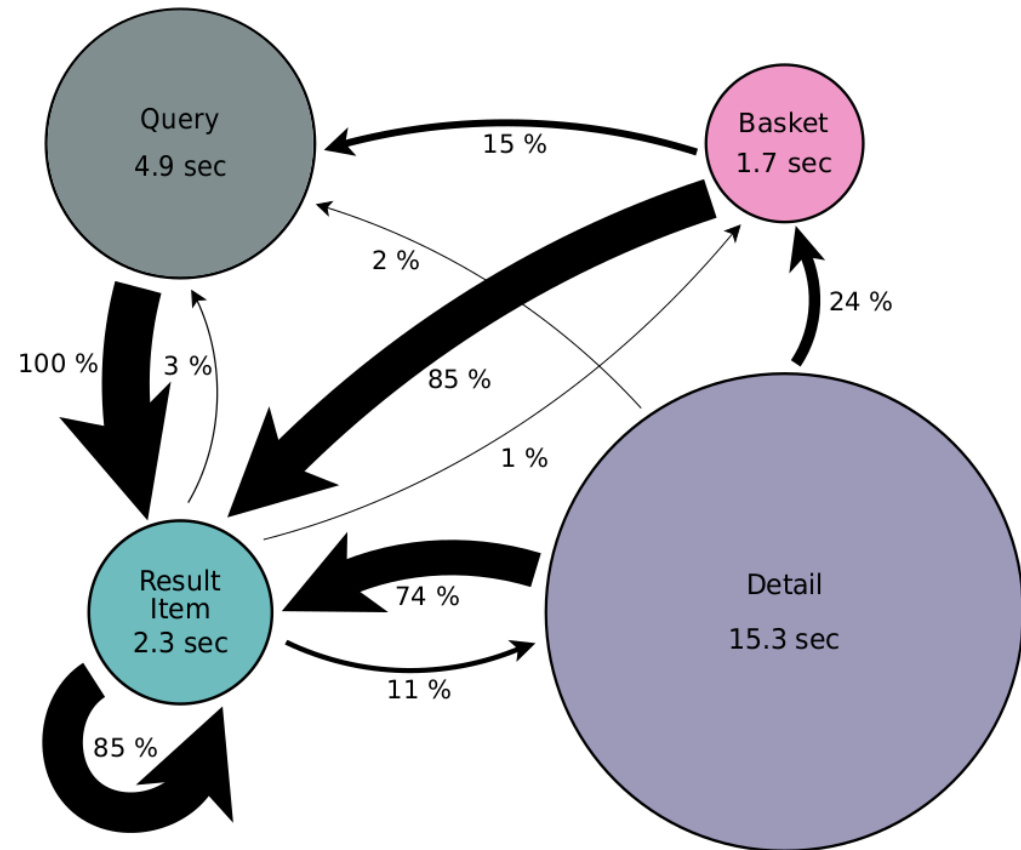
$$T_r = t_r + p_{rq} T_q + p_{rr} T_r + p_{rd} T_d$$

$$T_d = t_d + p_{dq} T_q + p_{dr} T_r$$

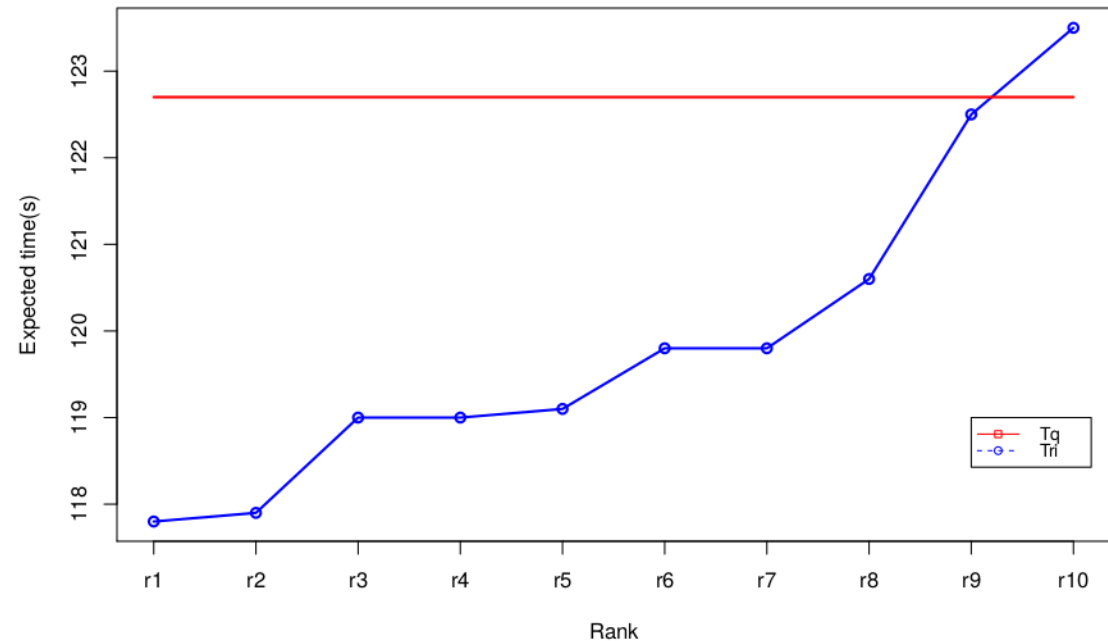
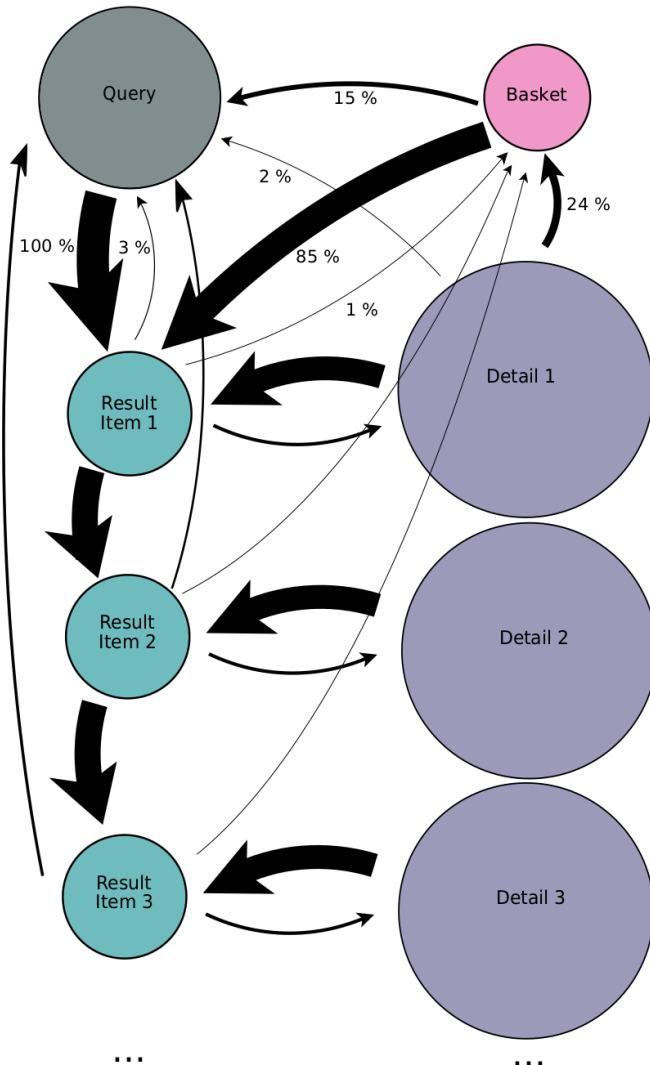
$$T_q = 122.7 \text{ s}$$

$$T_r = 117.8 \text{ s}$$

$$T_d = 104.9 \text{ s}$$

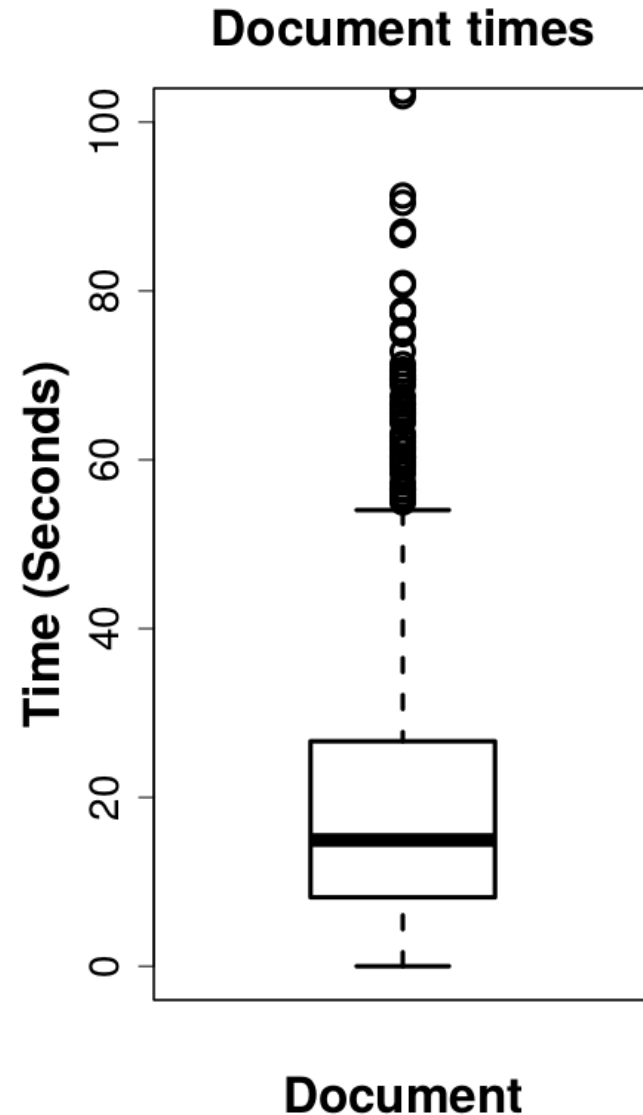


# Guidance: Modeling Retrieval Ranks



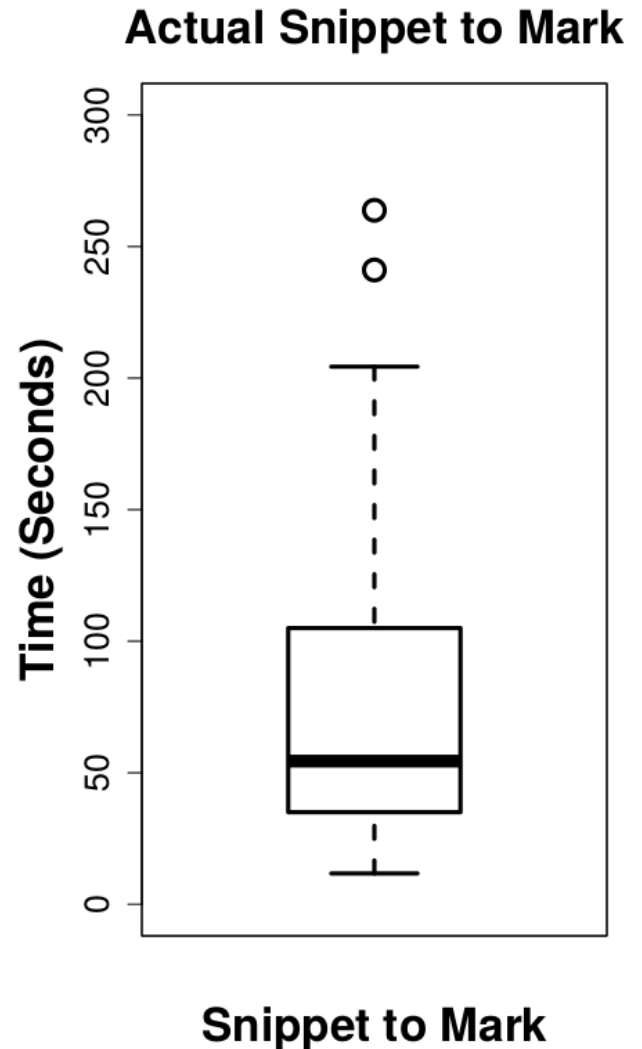
➔ After rank 9, it is better to reformulate query

- Analysis of 36 User sessions
- TREC tasks
- Find as many relevant items as possible
- Data from Maxwell & Azzopardi 2014
- Document times capped at 3.5 StdDev



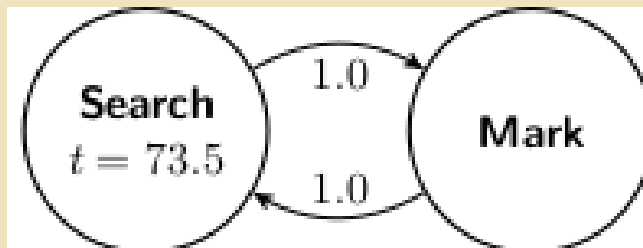
# Search Time Prediction

- Regard time from first snippet (after query or mark) to mark

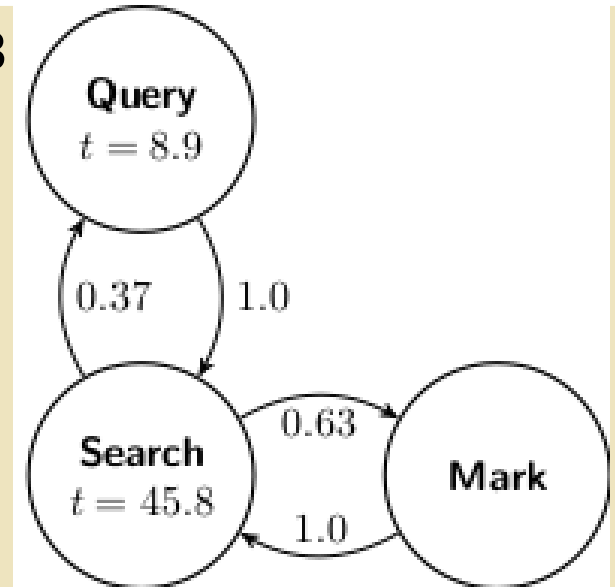


# Model Variants

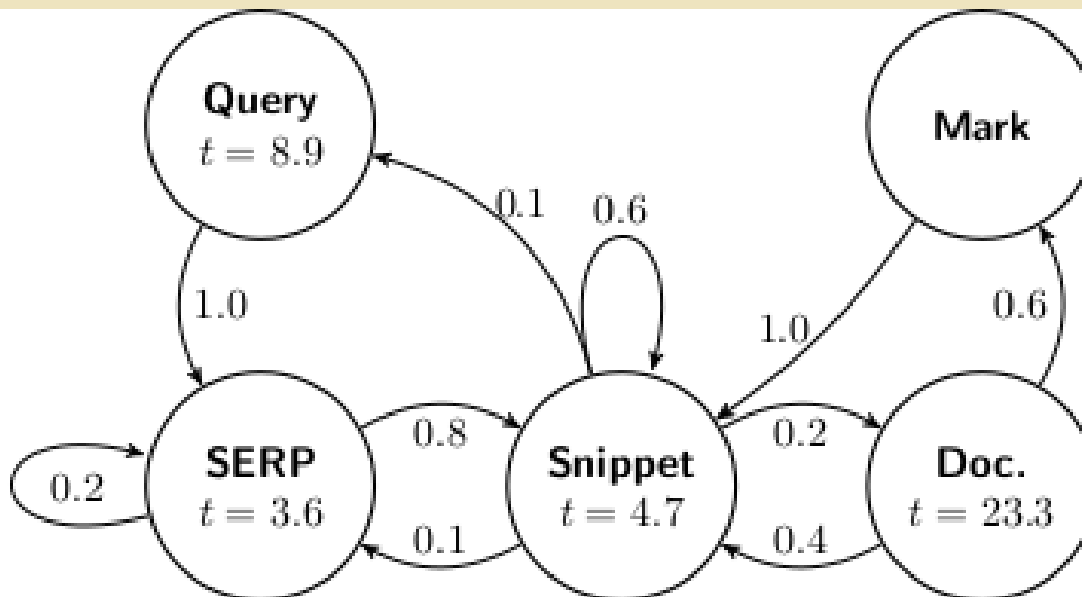
**MM2**



**MM3**

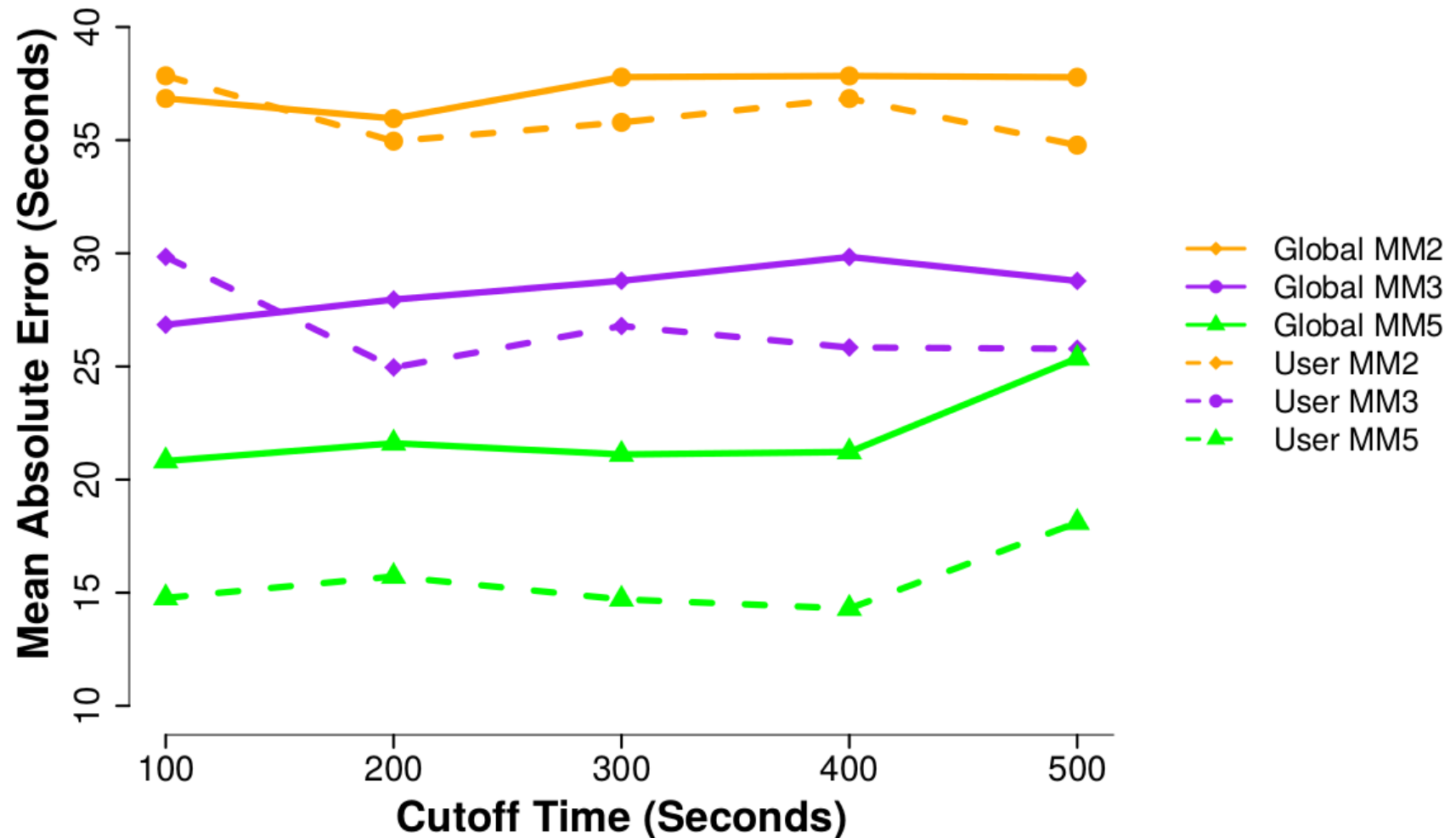


**MM5**



# Quality of Predictions

**Mean Error vs. Cutoff Times: MM2, MM3 and MM5**



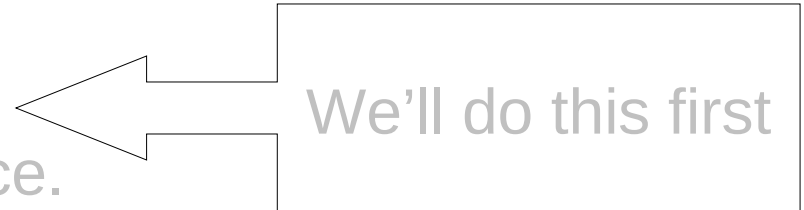


# Hidden State

“It’s currently time  $t$ , and human remains uncrushed. What’s the probability of crushing occurring at time  $t + 1$  ?”

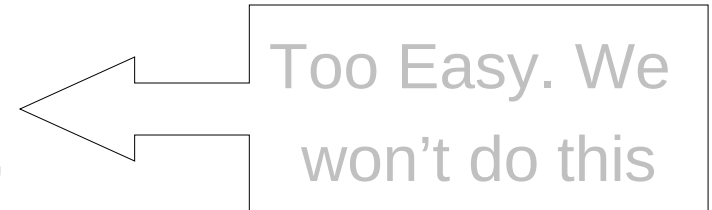
If robot is blind:

We can compute this in advance.



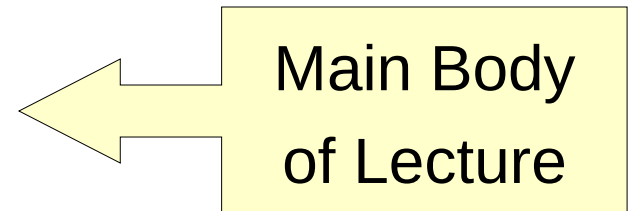
If robot is omnipotent:

(I.E. If robot knows state at time  $t$ ),  
can compute directly.



If robot has some sensors, but  
incomplete state information ...

Hidden Markov Models are  
applicable!

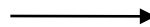


# Hidden State

- The previous example tried to estimate  $P(q_t = s_i)$  unconditionally (using no observed evidence).
- Suppose we can observe something that's affected by the true state.
- Example: Proximity sensors (tell us the contents of the 8 adjacent squares)

			$R_0$		
		H			

True state  $q_t$



W	W	W
	Ⓡ	
H		

W  
denotes  
"WALL"

What the robot sees:  
Observation  $O_t$

# Noisy Hidden State

- Example: Noisy proximity sensors. (unreliably tell us the contents of the 8 adjacent squares)

			$R_0$		
		H			

True state  $q_t$



W	W	W
	Ⓡ	
H		

W  
denotes  
"WALL"

Uncorrupted Observation



W		W
	Ⓡ	W
H	H	

What the robot sees:  
Observation  $O_t$

# Noisy Hidden State

- Example: Noisy Proximity sensors. (unreliably tell us the contents of the 8 adjacent squares)

			$R_0$		2
		H			

True state  $q_t$

$O_t$  is noisily determined depending on the current state.

Assume that  $O_t$  is conditionally independent of  $\{q_{t-1}, q_{t-2}, \dots, q_1, q_0, O_{t-1}, O_{t-2}, \dots, O_1, O_0\}$  given  $q_t$ .

In other words:

$$P(O_t = X | q_t = s_i) =$$

$$P(O_t = X | q_t = s_i, \text{any earlier history})$$



W	W	W
	Ⓡ	
H		

W  
denotes  
"WALL"

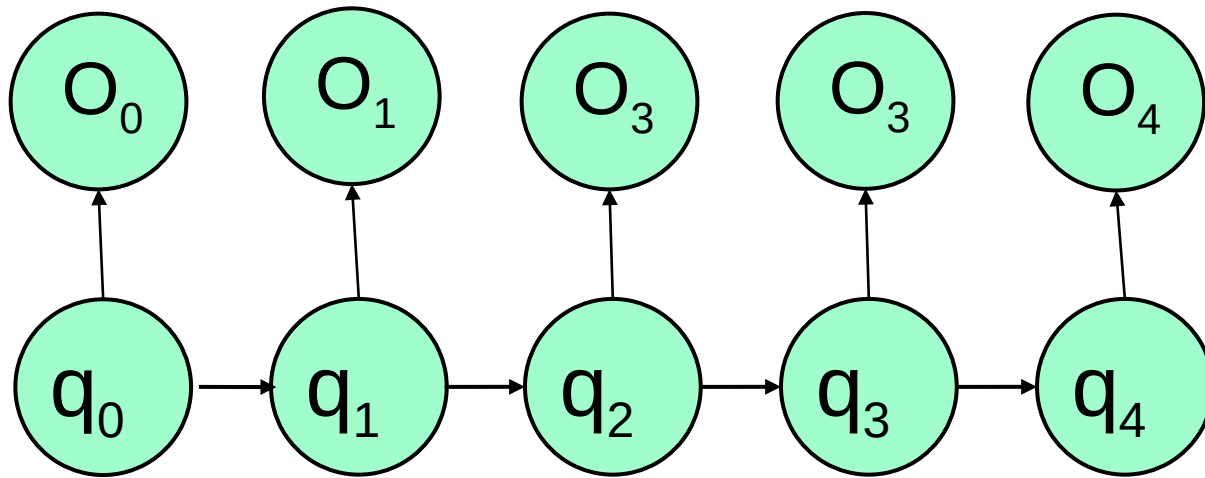
Uncorrupted Observation



W		W
	Ⓡ	W
H	H	

What the robot sees:  
Observation  $O_t$

# Noisy Hidden State: Representation



# Hidden Markov Models

Our robot with noisy sensors is a good example of an HMM

- Question 1: State Estimation

What is  $P(q_T = S_i \mid O_1 O_2 \dots O_T)$

It will turn out that a new cute D.P. trick will get this for us.

- Question 2: Most Probable Path

Given  $O_1 O_2 \dots O_T$ , what is the most probable path that I took?

And what is that probability?

Yet another famous D.P. trick, the VITERBI algorithm, gets this.

- Question 3: Learning HMMs:

Given  $O_1 O_2 \dots O_T$ , what is the maximum likelihood HMM that could have produced this string of observations?

Very very useful. Uses the E.M. Algorithm

# Are H.M.M.s Useful?

You bet !!

- Robot planning + sensing when there's uncertainty
- Speech Recognition/Understanding  
Phonemes → Words, Signal → phonemes
- Consumer decision modeling
- Economics & Finance.
- Many others ...

# HMM Notation (from Rabiner's Survey)

The states are labeled  $S_1 S_2 \dots S_N$

\*L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," Proc. of the IEEE, Vol.77, No.2, pp.257--286, 1989.

For a particular trial....

Let  $T$  be the number of observations

$T$  is also the number of states passed through

$O = O_1 O_2 \dots O_T$  is the sequence of observations

$Q = q_1 q_2 \dots q_T$  is the notation for a path of states

$\lambda = \langle N, M, \{\pi_i\}, \{a_{ij}\}, \{b_i(j)\} \rangle$  is the specification of an  
HMM



# HMM Formal Definition

An HMM,  $\lambda$ , is a 5-tuple consisting of

- $N$  the number of states
- $M$  the number of possible observations
- $\{\pi_1, \pi_2, \dots, \pi_N\}$  The starting state probabilities

$$P(q_0 = S_i) = \pi_i$$

This is new. In our previous example, start state was deterministic

- $$\begin{matrix} a_{11} & a_{12} & \dots & a_{1N} \\ a_{21} & a_{22} & \dots & a_{2N} \\ \vdots & \vdots & & \vdots \\ a_{N1} & a_{N2} & \dots & a_{NN} \end{matrix}$$

The state transition probabilities

$$P(q_{t+1}=S_j \mid q_t=S_i)=a_{ij}$$

- $$\begin{matrix} b_1(1) & b_1(2) & \dots & b_1(M) \\ b_2(1) & b_2(2) & \dots & b_2(M) \\ \vdots & \vdots & & \vdots \\ b_N(1) & b_N(2) & \dots & b_N(M) \end{matrix}$$

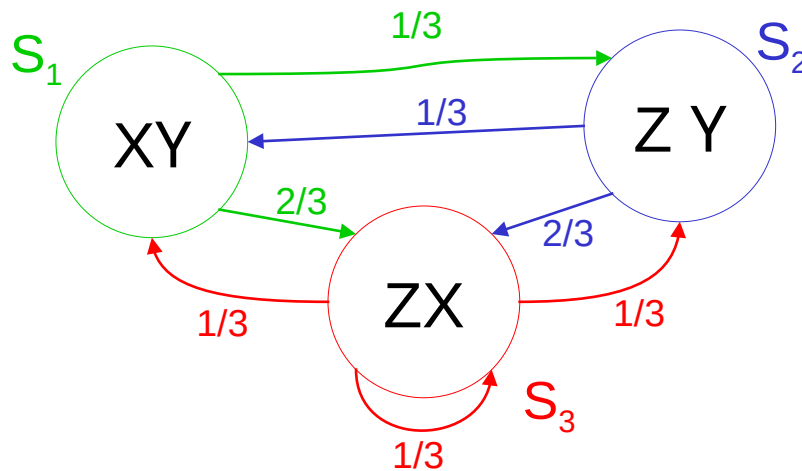
The observation probabilities

$$P(O_t=k \mid q_t=S_i)=b_i(k)$$

# Here's an HMM

Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

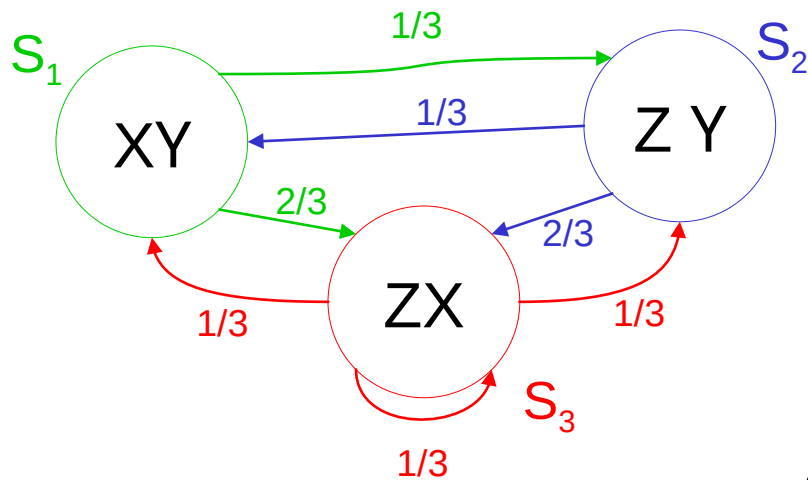
$$a_{33} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

# Here's an HMM



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

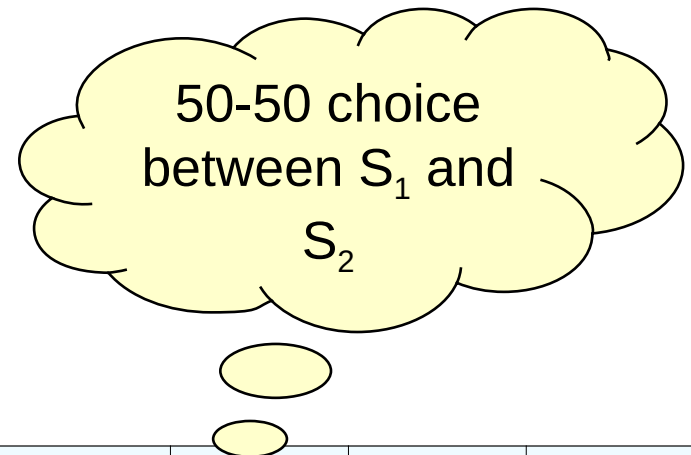
$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

Start randomly in state 1 or 2

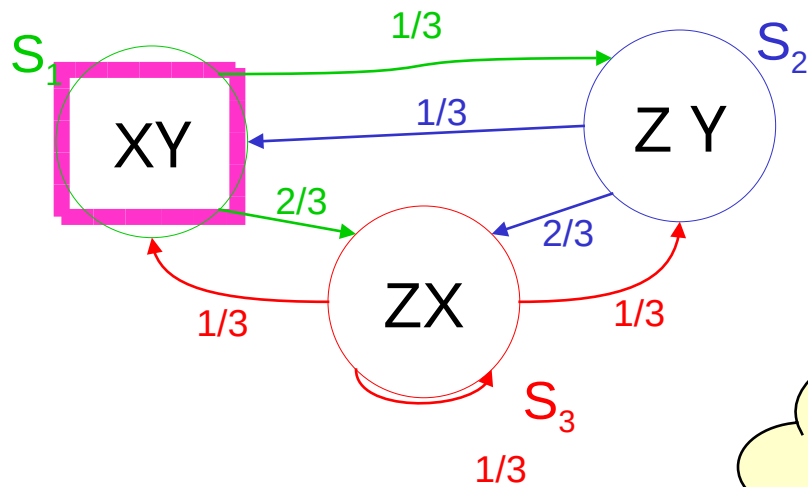
Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:



$q_0 =$	<u>—</u>	$O_0 =$	—
$q_1 =$	—	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

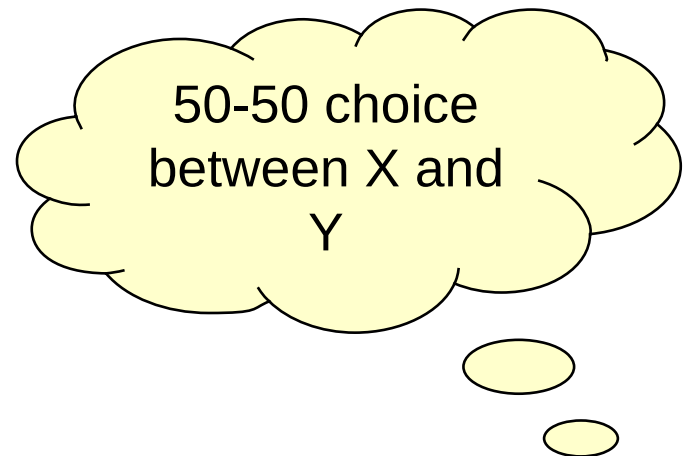
# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

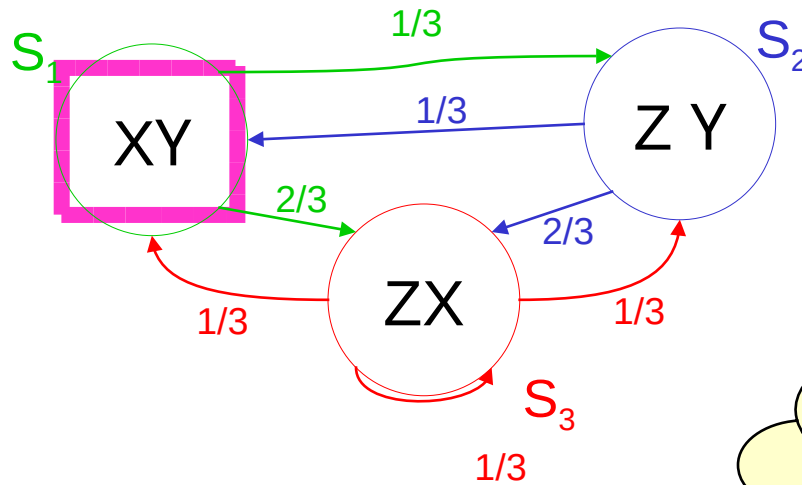
$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	$S_1$	$O_0 =$	<u>   </u>
$q_1 =$	—	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

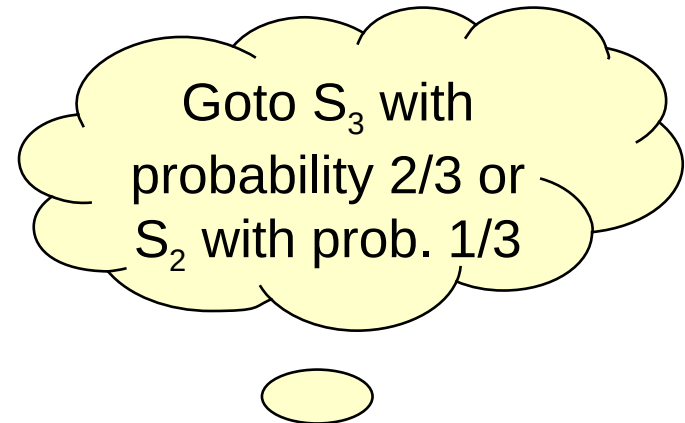
# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

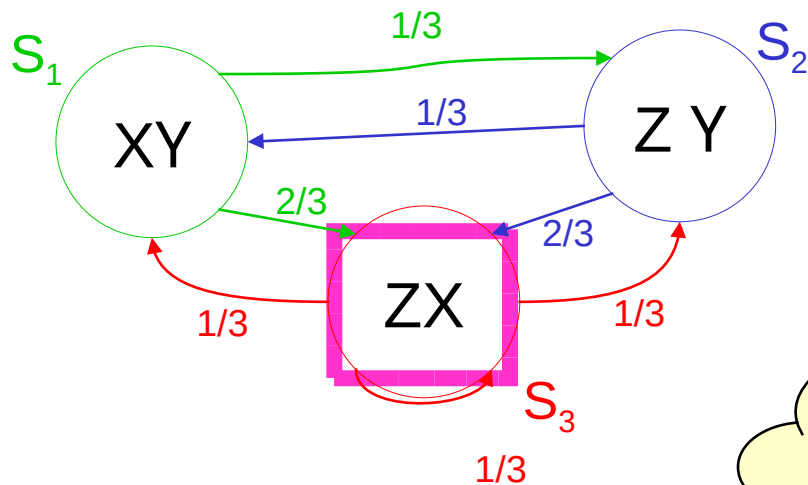
$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	$S_1$	$O_0 =$	X
$q_1 =$	—	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

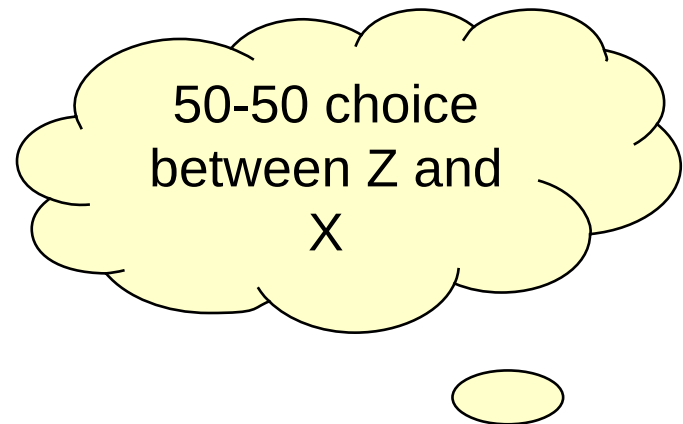
# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

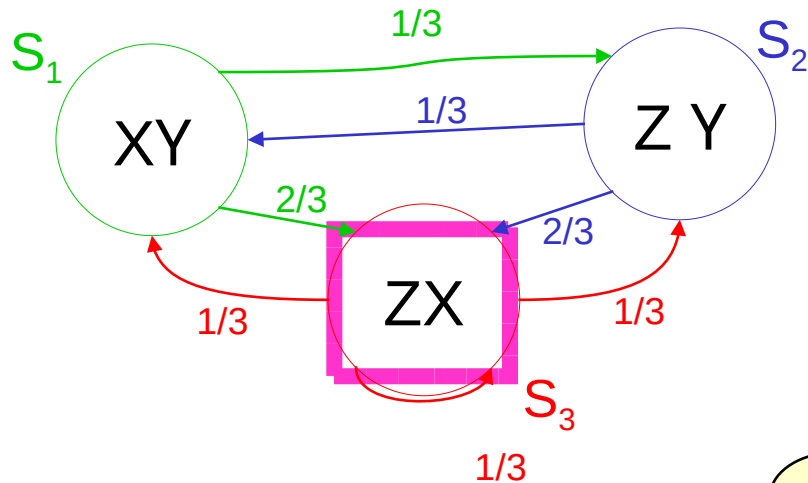
$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	$S_1$	$O_0 =$	X
$q_1 =$	$S_3$	$O_1 =$	—
$q_2 =$	—	$O_2 =$	—

# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:

$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

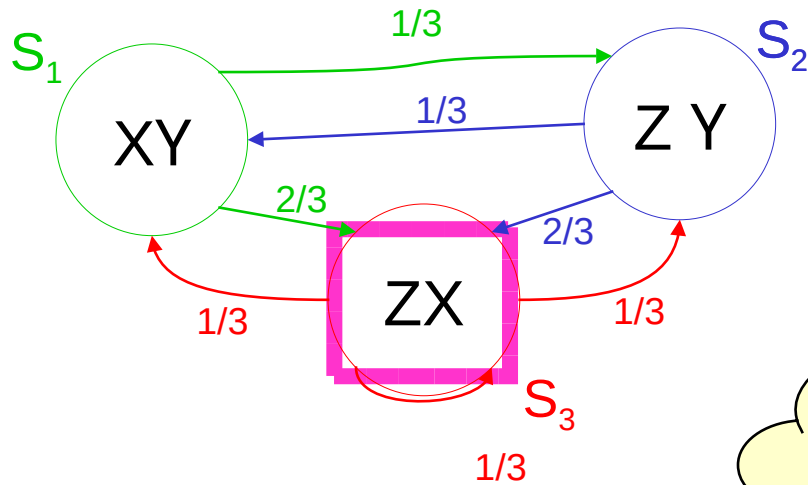
$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

Each of the three next states is equally likely

$q_0 =$	$S_1$	$O_0 =$	X
$q_1 =$	$S_3$	$O_1 =$	X
$q_2 =$	—	$O_2 =$	—

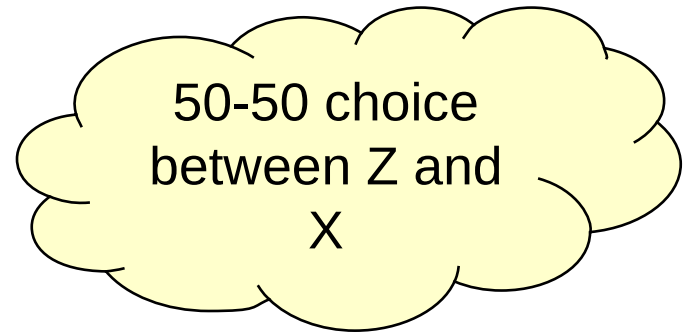
# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:



$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

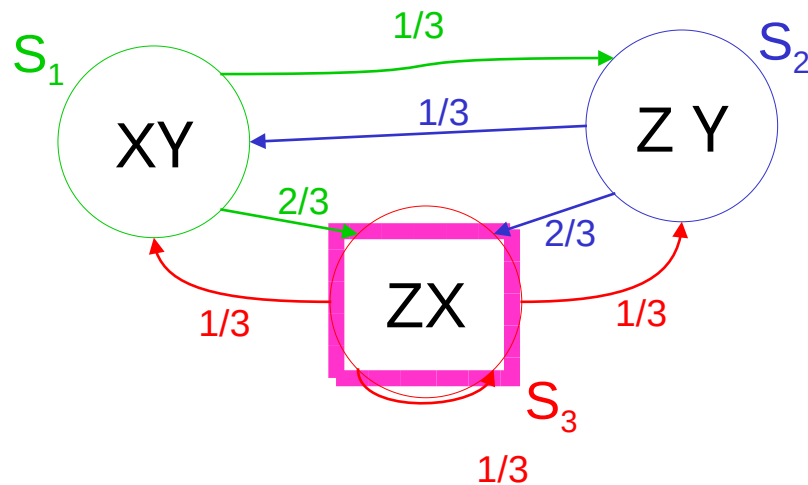
$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	$S_1$	$O_0 =$	X
$q_1 =$	$S_3$	$O_1 =$	X
$q_2 =$	$S_3$	$O_2 =$	—



# Here's an HMM



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:

$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

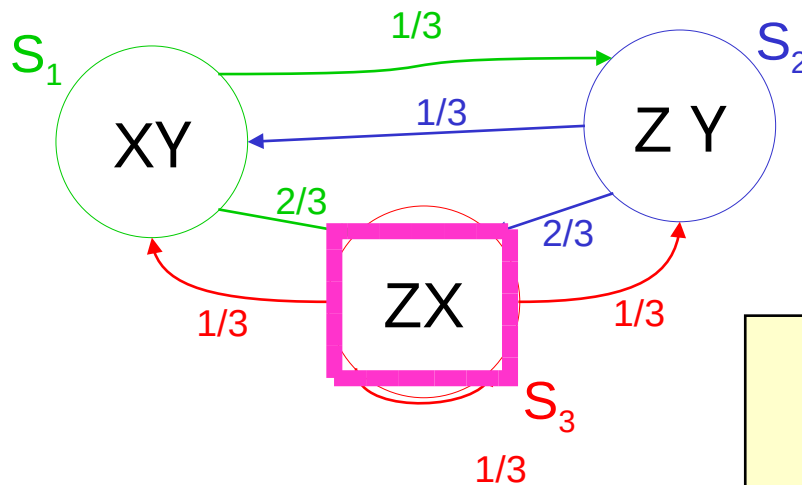
$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	$S_1$	$O_0 =$	X
$q_1 =$	$S_3$	$O_1 =$	X
$q_2 =$	$S_3$	$O_2 =$	Z

# State Estimation



Start randomly in state 1 or 2

Choose one of the output symbols in each state at random.

Let's generate a sequence of observations:

This is what the observer has to work with...

$$N = 3$$

$$M = 3$$

$$\pi_1 = 1/2$$

$$\pi_2 = 1/2$$

$$\pi_3 = 0$$

$$a_{11} = 0$$

$$a_{12} = 1/3$$

$$a_{13} = 2/3$$

$$a_{21} = 1/3$$

$$a_{22} = 0$$

$$a_{23} = 2/3$$

$$a_{31} = 1/3$$

$$a_{32} = 1/3$$

$$a_{33} = 1/3$$

$$b_1(X) = 1/2 \quad b_1(Y) = 1/2 \quad b_1(Z) = 0$$

$$b_2(X) = 0 \quad b_2(Y) = 1/2 \quad b_2(Z) = 1/2$$

$$b_3(X) = 1/2 \quad b_3(Y) = 0 \quad b_3(Z) = 1/2$$

$q_0 =$	?	$O_0 =$	X
$q_1 =$	?	$O_1 =$	X
$q_2 =$	?	$O_2 =$	Z

# Prob. of a series of observations

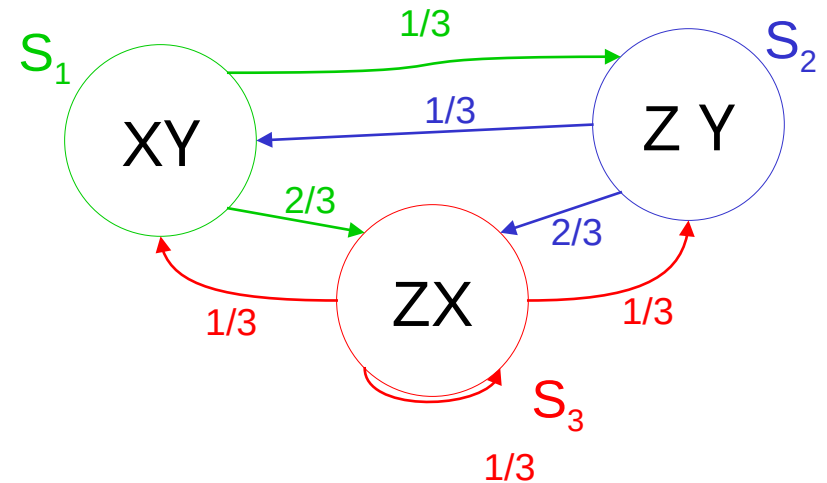
What is  $P(\mathbf{O}) = P(O_1 O_2 O_3)$   
 $= P(O_1=X \wedge O_2=X \wedge O_3=Z)$ ?

Slow, stupid way:

$$\begin{aligned} P(O) &= \sum_{Q \in \text{Paths of length 3}} P(O \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(O|Q)P(Q) \end{aligned}$$

How do we compute  $P(Q)$  for an arbitrary path  $Q$ ?

How do we compute  $P(O|Q)$  for an arbitrary path  $Q$ ?



# Prob. of a series of observations

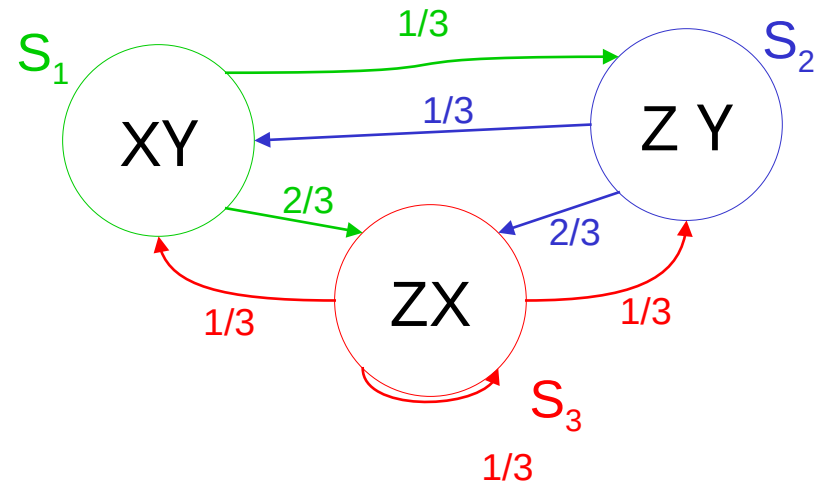
What is  $P(\mathbf{O}) = P(O_1 O_2 O_3) =$   
 $P(O=X \wedge O_2=X \wedge O_3=Z)$ ?

Slow, stupid way:

$$\begin{aligned} P(O) &= \sum_{Q \in \text{Paths of length 3}} P(O \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(O|Q)P(Q) \end{aligned}$$

How do we compute  $P(Q)$  for  
an arbitrary path  $Q$ ?

How do we compute  $P(O|Q)$   
for an arbitrary path  $Q$ ?



$$P(Q) = P(q_1, q_2, q_3)$$

$$= P(q_1) P(q_2, q_3 | q_1) \text{ (chain rule)}$$

$$= P(q_1) P(q_2 | q_1) P(q_3 | q_2, q_1) \text{ (chain)}$$

$$= P(q_1) P(q_2 | q_1) P(q_3 | q_2) \text{ (why?)}$$

Example in the case  $Q = S_1 S_3 S_3$ :

$$= 1/2 * 2/3 * 1/3 = 1/9$$

# Prob. of a series of observations

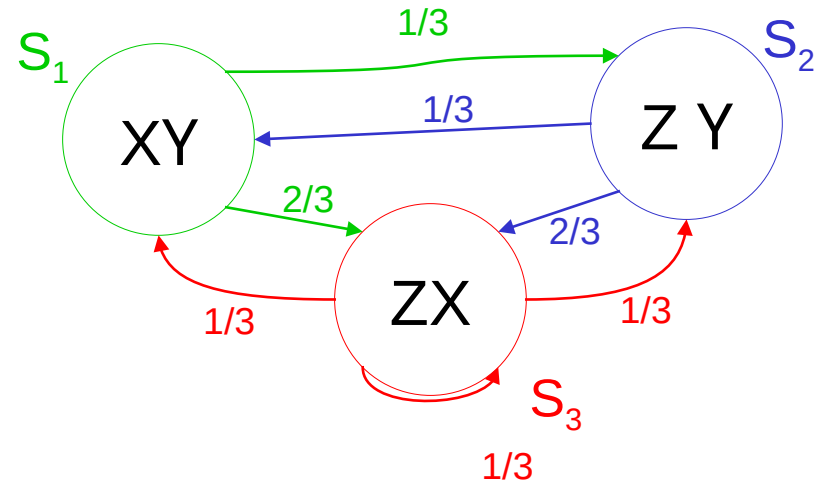
What is  $P(\mathbf{O}) = P(O_1 O_2 O_3) =$   
 $P(O=X \wedge O_2=X \wedge O_3=Z)$ ?

Slow, stupid way:

$$\begin{aligned} P(O) &= \sum_{Q \in \text{Paths of length 3}} P(O \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(O|Q)P(Q) \end{aligned}$$

How do we compute  $P(Q)$  for  
an arbitrary path  $Q$ ?

How do we compute  $P(O|Q)$   
for an arbitrary path  $Q$ ?



$P(O|Q)$

$$= P(O_1 O_2 O_3 | q_1 q_2 q_3)$$

$$= P(O_1 | q_1) P(O_2 | q_2) P(O_3 | q_3) \text{ (why?)}$$

Example in the case  $Q = S_1 S_3 S_3$ :

$$= P(X | S_1) P(X | S_3) P(Z | S_3) =$$

$$= 1/2 * 1/2 * 1/2 = 1/8$$

# Prob. of a series of observations

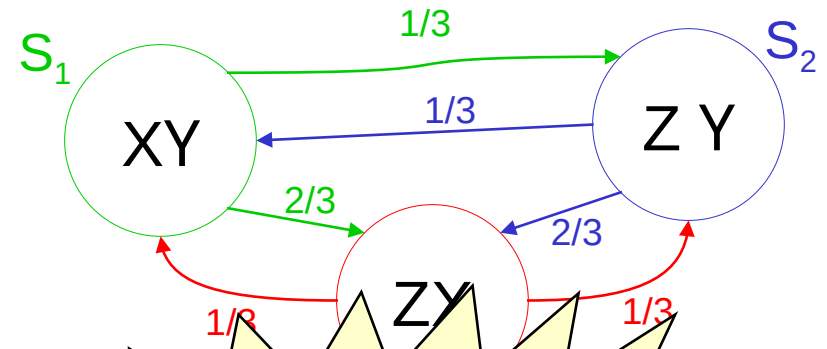
What is  $P(\mathbf{O}) = P(O_1 O_2 O_3) =$   
 $P(O=X \wedge O_2=X \wedge O_3=Z)$ ?

Slow, stupid way:

$$\begin{aligned} P(\mathbf{O}) &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O} \wedge Q) \\ &= \sum_{Q \in \text{Paths of length 3}} P(\mathbf{O}|Q)P(Q) \end{aligned}$$

How do we compute  $P(Q)$  for  
an arbitrary path  $Q$ ?

How do we compute  $P(\mathbf{O}|Q)$   
for an arbitrary path  $Q$ ?



$P(\mathbf{O})$  would need 27  $P(Q)$   
computations and 27  $P(\mathbf{O}|Q)$   
computations

A sequence of 20 observations would need  $3^{20} =$   
3.5 billion computations and 3.5 billion  $P(\mathbf{O}|Q)$   
computations

So let's be smarter...

# The Probability of a given series of observations, non-exponential-cost-style

Given observations  $O_1 O_2 \dots O_T$

Define

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i \mid \lambda) \quad \text{where } 1 \leq t \leq T$$

$\alpha_t(i)$  = Probability that, in a random trial,

- We'd have seen the first  $t$  observations
- We'd have ended up in  $S_i$  as the  $t$ 'th state visited.

In our example, what is  $\alpha_2(3)$  ?

# $\alpha_t(i)$ : easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i \mid \lambda)$$

( $\alpha_t(i)$  can be defined stupidly by considering all paths length “t”. How?)

$$\begin{aligned}\alpha_1(i) &= P(O_1 \wedge q_1 = S_i) \\ &= P(q_1 = S_i) P(O_1 \mid q_1 = S_i) \\ &= \pi_i b_i(O_1)\end{aligned}$$

$$\begin{aligned}\alpha_{t+1}(j) &= P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j) \\ &= \end{aligned}$$



# $\alpha_t(i)$ : easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i \mid \lambda)$$

$$\alpha_1(i) = P(O_1 \wedge q_1 = S_i)$$

$$= P(q_1 = S_i) P(O_1 \mid q_1 = S_i)$$

$$= \pi_i b_i(O_1)$$

$$\alpha_{t+1}(j) = P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j)$$

$$= \sum_{i=1}^N P(O_1 O_2 \dots O_t \wedge q_t = S_i \wedge O_{t+1} \wedge q_{t+1} = S_j)$$

# $\alpha_t(i)$ : easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i \mid \lambda)$$

$$\alpha_1(i) = P(O_1 \wedge q_1 = S_i)$$

$$= P(q_1 = S_i) P(O_1 \mid q_1 = S_i)$$

$$= \pi_i b_i(O_1)$$

$$\alpha_{t+1}(j) = P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j)$$

$$= \sum_{i=1}^N P(O_1 O_2 \dots O_t \wedge q_t = S_i \wedge O_{t+1} \wedge q_{t+1} = S_j)$$

$$= \sum_{i=1}^N P(O_{t+1}, q_{t+1} = S_j \mid O_1 O_2 \dots O_t \wedge q_t = S_i) P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

$$= \sum_i P(O_{t+1}, q_{t+1} = S_j \mid q_t = S_i) \alpha_t(i)$$

# $\alpha_t(i)$ : easy to define recursively

$$\alpha_t(i) = P(O_1 O_2 \dots O_T \wedge q_t = S_i \mid \lambda)$$

$$\alpha_1(i) = P(O_1 \wedge q_1 = S_i)$$

$$= P(q_1 = S_i) P(O_1 \mid q_1 = S_i)$$

$$= \pi_i b_i(O_1)$$

$$\alpha_{t+1}(j) = P(O_1 O_2 \dots O_t O_{t+1} \wedge q_{t+1} = S_j)$$

$$= \sum_{i=1}^N P(O_1 O_2 \dots O_t \wedge q_t = S_i \wedge O_{t+1} \wedge q_{t+1} = S_j)$$

$$= \sum_{i=1}^N P(O_{t+1}, q_{t+1} = S_j \mid O_1 O_2 \dots O_t \wedge q_t = S_i) P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

$$= \sum_i P(O_{t+1}, q_{t+1} = S_j \mid q_t = S_i) \alpha_t(i)$$

$$= \sum_i P(q_{t+1} = S_j \mid q_t = S_i) P(O_{t+1} \mid q_{t+1} = S_j) \alpha_t(i)$$

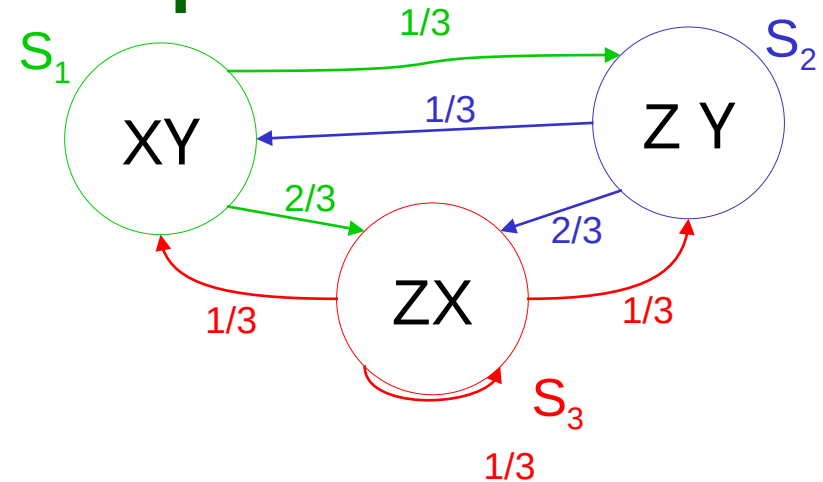
$$= \sum_i a_{ij} b_j(O_{t+1}) \alpha_t(i)$$

# in our example

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i | \lambda)$$

$$\alpha_1(i) = b_i(O_1) \pi_i$$

$$\alpha_{t+1}(j) = \sum_i a_{ij} b_j(O_{t+1}) \alpha_t(i)$$



WE SAW  $O_1 O_2 O_3 = X X Z$

$$\alpha_1(1) = \frac{1}{4}$$

$$\alpha_1(2) = 0$$

$$\alpha_1(3) = 0$$

$$\alpha_2(1) = 0$$

$$\alpha_2(2) = 0$$

$$\alpha_2(3) = \frac{1}{12}$$

$$\alpha_3(1) = 0$$

$$\alpha_3(2) = \frac{1}{72}$$

$$\alpha_3(3) = \frac{1}{72}$$

# Easy Question

We can cheaply compute

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

(How) can we cheaply compute

$$P(O_1 O_2 \dots O_t) \quad ?$$

(How) can we cheaply compute

$$P(q_t = S_i | O_1 O_2 \dots O_t)$$

# Easy Question

We can cheaply compute

$$\alpha_t(i) = P(O_1 O_2 \dots O_t \wedge q_t = S_i)$$

(How) can we cheaply compute

$$P(O_1 O_2 \dots O_t) \quad ? \quad \sum_{i=1}^N \alpha_t(i)$$

(How) can we cheaply compute

$$P(q_t = S_i | O_1 O_2 \dots O_t) \quad \frac{\alpha_t(i)}{\sum_{j=1}^N \alpha_t(j)}$$

# Most probable path given observations

What's most probable path given  $O_1 O_2 \dots O_T$ , i.e.

What is  $\operatorname{argmax}_Q P(Q|O_1 O_2 \dots O_T)$ ?

Slow, stupid answer:

$$\begin{aligned} & \operatorname{argmax}_Q P(Q|O_1 O_2 \dots O_T) \\ &= \operatorname{argmax}_Q \frac{P(O_1 O_2 \dots O_T|Q) P(Q)}{P(O_1 O_2 \dots O_T)} \\ &= \operatorname{argmax}_Q P(O_1 O_2 \dots O_T|Q) P(Q) \end{aligned}$$

# Efficient MPP computation

We're going to compute the following variables:

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 \dots O_t)$$

= The Probability of the path of Length  $t-1$  with the maximum chance of doing all these things:

...OCCURRING

and

...ENDING UP IN STATE  $S_i$

and

...PRODUCING OUTPUT  $O_1 \dots O_t$

DEFINE:  $mpp_t(i)$  = that path

So:  $\delta_t(i) = \text{Prob}(mpp_t(i))$



# The Viterbi Algorithm

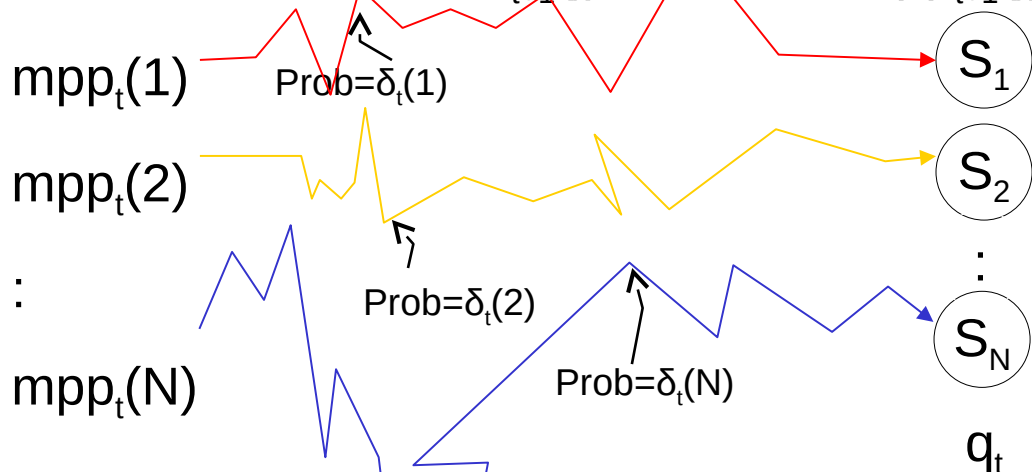
$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

$$mpp_t(i) = \arg \max_{q_1 q_2 \dots q_{t-1}} P(q_1 q_2 \dots q_{t-1} \wedge q_t = S_i \wedge O_1 O_2 \dots O_t)$$

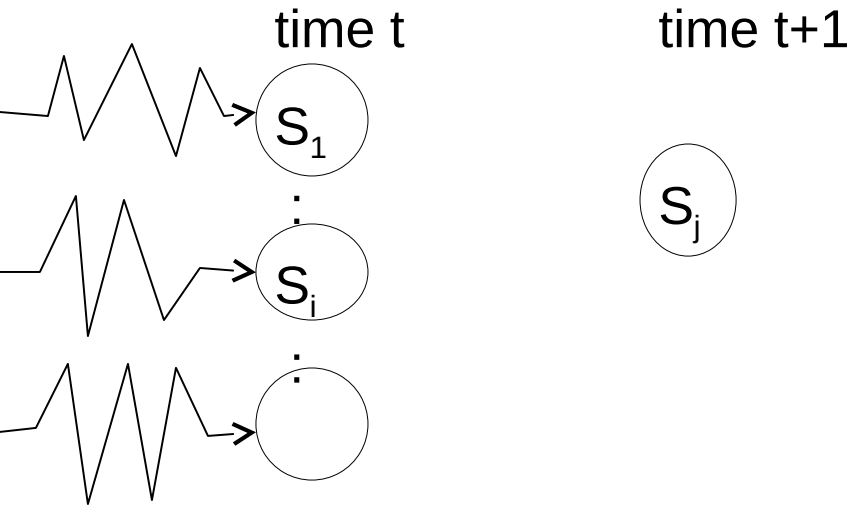
$$\begin{aligned} \delta_1(i) &= \max P(q_1 = S_i \wedge O_1) \\ &= P(q_1 = S_i) P(O_1 | q_1 = S_i) \\ &= \pi_i b_i(O_1) \end{aligned}$$

Now, suppose we have all the  $\delta_t(i)$ 's and  $mpp_t(i)$ 's for all  $i$ .

HOW TO GET  $\delta_{t+1}(j)$  and  $mpp_{t+1}(j)$ ?



# The Viterbi Algorithm

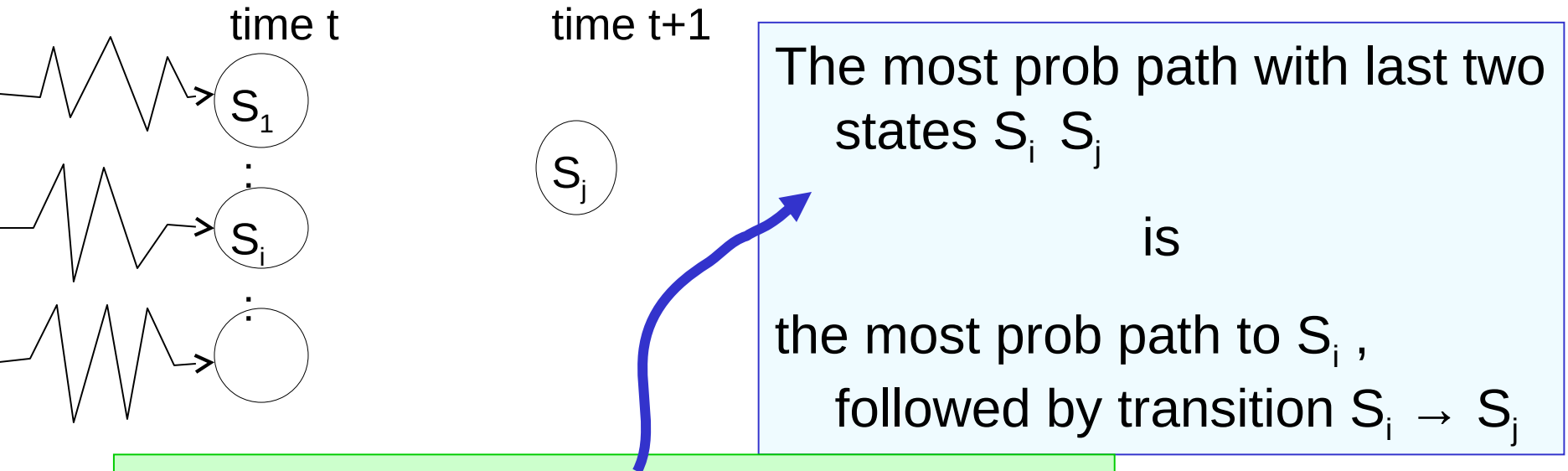


The most prob. path with last  
two states  $S_i$   $S_j$

is

the most prob path to  $S_i$  ,  
followed by transition  $S_i \rightarrow S_j$

# The Viterbi Algorithm



What is the prob of that path?

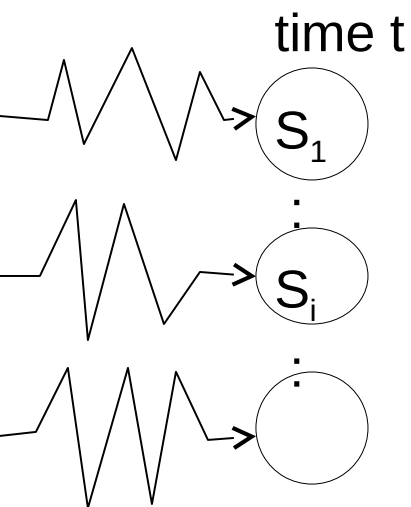
$$\delta_t(i) \times P(S_i \rightarrow S_j \wedge O_{t+1} | \lambda)$$

$$= \delta_t(i) a_{ij} b_j(O_{t+1})$$

SO The most probable path to  $S_j$  has  $S_{i^*}$  as its penultimate state

where  $i^* = \operatorname{argmax}_i \delta_t(i) a_{ij} b_j(O_{t+1})$

# The Viterbi Algorithm



time t+1



The most prob path with last two states  $S_i$   $S_j$

is

the most prob path to  $S_i$ , followed by transition  $S_i \rightarrow S_j$

What is the prob of that path?

$$\delta_t(i) \times P(S_i \rightarrow S_j | O_{t+1})$$

$$= \delta_t(i) a_{ij} b_j(O_{t+1})$$

SO The most probable

$S_{i^*}$  as its penultimate state

where  $i^* = \arg\max_i \delta_t(i) a_{ij} b_j(O_{t+1})$

Summary:

$$\left. \begin{aligned} \delta_{t+1}(j) &= \delta_t(i^*) a_{ij} b_j(O_{t+1}) \\ mpp_{t+1}(j) &= mpp_{t+1}(i^*) S_{i^*} \end{aligned} \right\} \text{ with } i^* \text{ defined to the left}$$

# What's Viterbi used for?

Classic Example

Speech recognition:

Signal  $\rightarrow$  words

HMM  $\rightarrow$  observable is signal

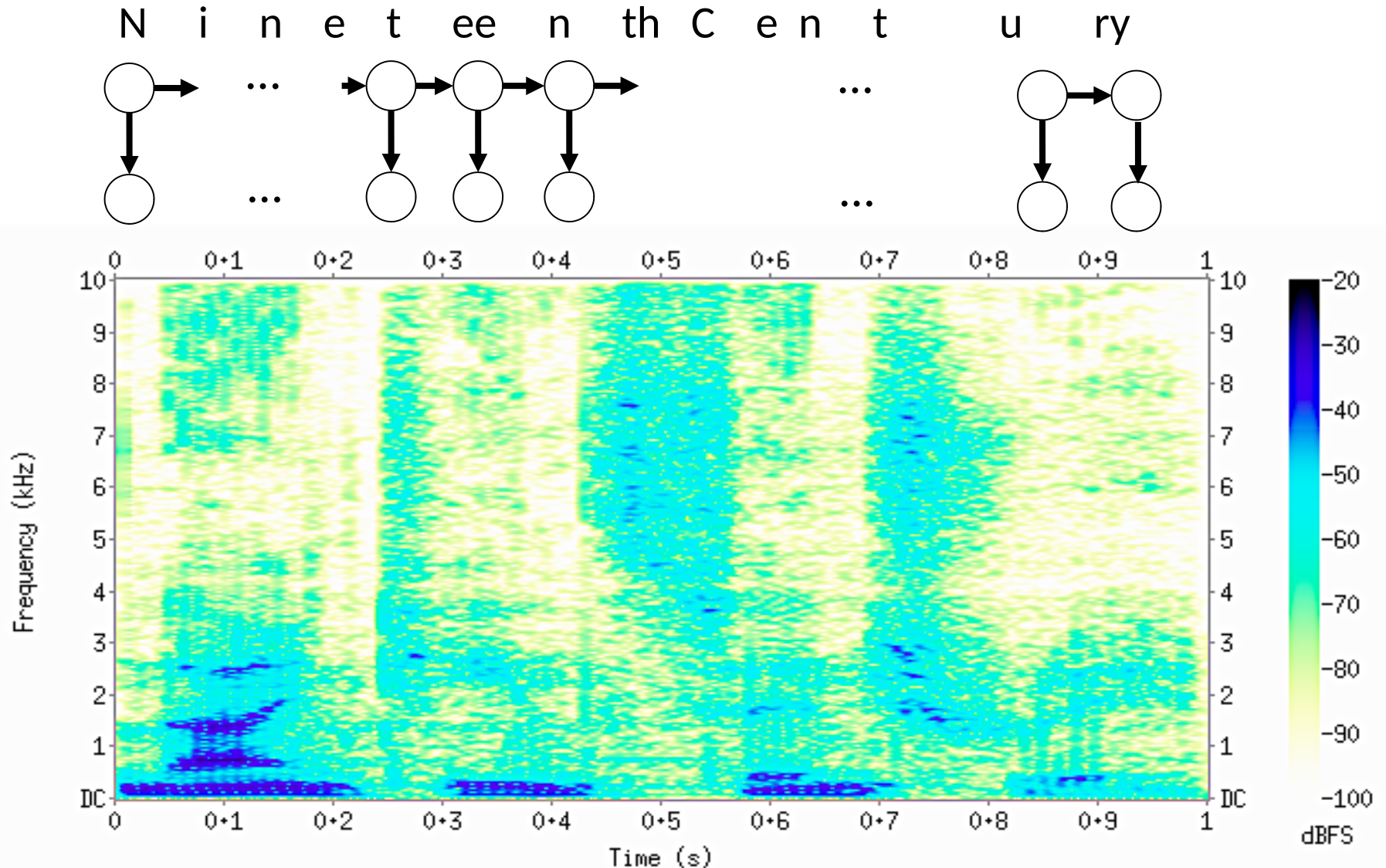
$\rightarrow$  Hidden state is part of word  
formation

What is the most probable word given this signal?

**UTTERLY GROSS SIMPLIFICATION**

In practice: many levels of inference; not  
one big jump.

# Example: HMMs for speech recognition



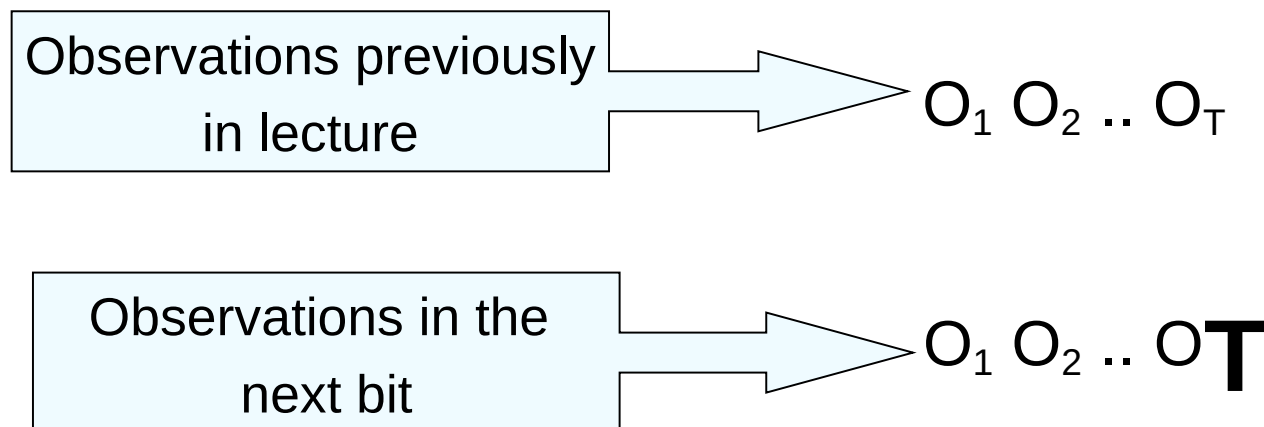
Spectrogram for audio for "nineteenth century" – From Wikipedia

# HMMs are used and useful

But how do you design an HMM?

Occasionally, (e.g. in our robot example) it is reasonable to deduce the HMM from first principles.

But usually, especially in Speech or Genetics, it is better to infer it from large amounts of data.  $O_1 O_2 \dots O_T$  with a big “T”.



# Inferring an HMM

Remember, we've been doing things like

$$P(O_1 O_2 \dots O_T \mid \lambda)$$

That “ $\lambda$ ” is the notation for our HMM parameters.

Now We have some observations and we want to estimate  $\lambda$  from them.

AS USUAL: We could use

$$\text{MAX LIKELIHOOD} \quad \lambda = \underset{\lambda}{\operatorname{argmax}} P(O_1 \dots O_T \mid \lambda)$$

BAYES

Work out  $P(\lambda \mid O_1 \dots O_T)$

and then take  $E[\lambda]$  or  $\underset{\lambda}{\operatorname{max}} P(\lambda \mid O_1 \dots O_T)$



# Max likelihood HMM estimation

Define

$$\gamma_t(i) = P(q_t = S_i \mid O_1 O_2 \dots O_T, \lambda)$$

$$\varepsilon_t(i,j) = P(q_t = S_i \wedge q_{t+1} = S_j \mid O_1 O_2 \dots O_T, \lambda)$$

$\gamma_t(i)$  and  $\varepsilon_t(i,j)$  can be computed efficiently  $\forall i,j,t$

(Details in Rabiner paper)

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{Expected number of transitions out of state } i \text{ during the path}$$

$$\sum_{t=1}^{T-1} \varepsilon_t(i,j) = \text{Expected number of transitions from state } i \text{ to state } j \text{ during the path}$$

$$\gamma_t(i) = P(q_t = S_i | O_1 O_2 \dots O_T, \lambda)$$

$$\varepsilon_t(i, j) = P(q_t = S_i \wedge q_{t+1} = S_j | O_1 O_2 \dots O_T, \lambda)$$

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions out of state } i \text{ during path}$$

$$\sum_{t=1}^{T-1} \varepsilon_t(i, j) = \text{expected number of transitions out of } i \text{ and into } j \text{ during path}$$

# HMM estimation

Notice 
$$\frac{\sum_{t=1}^{T-1} \varepsilon_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \frac{\left( \begin{array}{c} \text{expected frequency} \\ i \rightarrow j \end{array} \right)}{\left( \begin{array}{c} \text{expected frequency} \\ i \end{array} \right)}$$

= Estimate of Prob(Next state  $S_j$  | This state  $S_i$ )

We can re-estimate

$$a_{ij} \leftarrow \frac{\sum \varepsilon_t(i, j)}{\sum \gamma_t(i)}$$

We can also re-estimate

$$b_j(O_k) \leftarrow L$$

(see Rabiner)

# EM for HMMs

If we knew  $\lambda$  we could estimate EXPECTATIONS of quantities such as

Expected number of times in state  $i$

Expected number of transitions  $i \rightarrow j$

If we knew the quantities such as

Expected number of times in state  $i$

Expected number of transitions  $i \rightarrow j$

We could compute the MAX LIKELIHOOD estimate of

$$\lambda = \langle \{a_{ij}\}, \{b_i(j)\}, \pi_i \rangle$$

Roll on the EM Algorithm...

# EM 4 HMMs

1. Get your observations  $O_1 \dots O_T$
  2. Guess your first  $\lambda$  estimate  $\lambda(0)$ ,  $k=0$
  3.  $k = k+1$
  4. Given  $O_1 \dots O_T$ ,  $\lambda(k)$  compute
$$\gamma_t(i) , \varepsilon_t(i,j) \quad \forall 1 \leq t \leq T, \quad \forall 1 \leq i \leq N, \quad \forall 1 \leq j \leq N$$
  5. Compute expected freq. of state  $i$ , and expected freq.  $i \rightarrow j$
  6. Compute new estimates of  $a_{ij}$ ,  $b_j(k)$ ,  $\pi_i$  accordingly.  
Call them  $\lambda(k+1)$
  7. Goto 3, unless converged.
- **Also known (for the HMM case) as the BAUM-WELCH algorithm.**

# Bad News

- There are lots of local minima

# Good News

- The local minima are usually adequate models of the data.

# Notice

- EM does not estimate the number of states. That must be given.
- Often, HMMs are forced to have some links with zero probability. This is done by setting  $a_{ij}=0$  in initial estimate  $\lambda(0)$
- Easy extension of everything seen today: HMMs with real valued outputs

## Red Noise

Trade-off between too few states (inadequately modeling the structure in the data) and too many (fitting the noise).

- There are lots of

Thus #states is a **regularization parameter**.

- The local minimum data.

Blah blah blah... bias variance tradeoff...blah  
blah...cross-validation...blah blah....AIC,  
BIC....blah blah (same ol' same ol')

## Notice

- EM does not estimate the number of states. That must be given.
- Often, HMMs are forced to have some links with zero probability. This is done by setting  $a_{ij}=0$  in initial estimate  $\lambda(0)$
- Easy extension of everything seen today: HMMs with real valued outputs

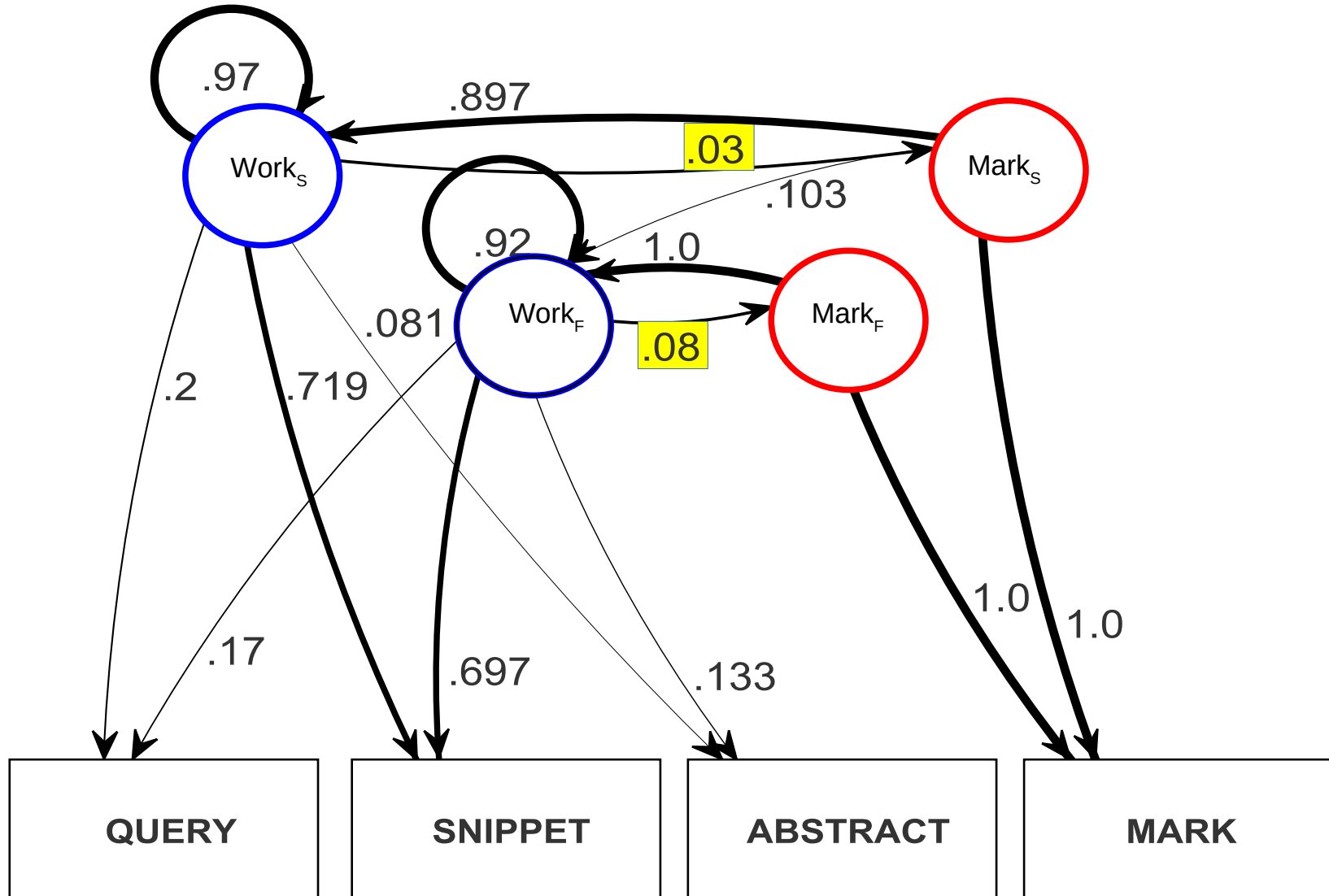
# What You Should Know

- What is a Markov model, what is an HMM?
- Which questions can be answered with a MM?
- What are the central questions for an HMM?
- What are the basic methods for addressing these problems?
- Give some examples of applications of MMs and HMMs!

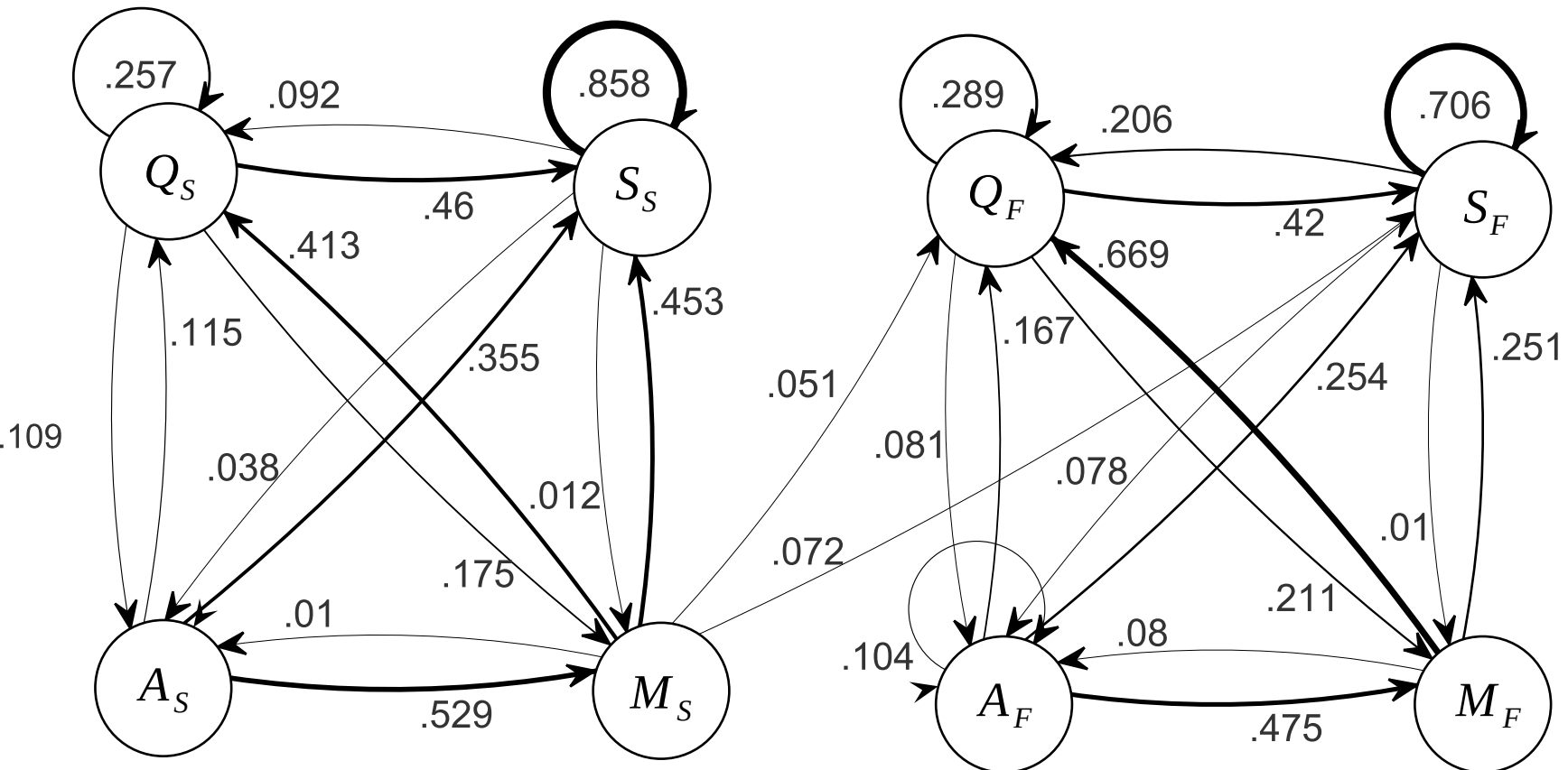
- Goal: detection of search phases
  - adaptive system behavior: ranking, search suggestions...
  - user guidance
- Data
  - Search logs from German Social Science Info Centre (GESIS)
  - Consider only sessions with  $> 3$  documents marked relevant
  - 1642 search sessions



# Simple Discrete HMM



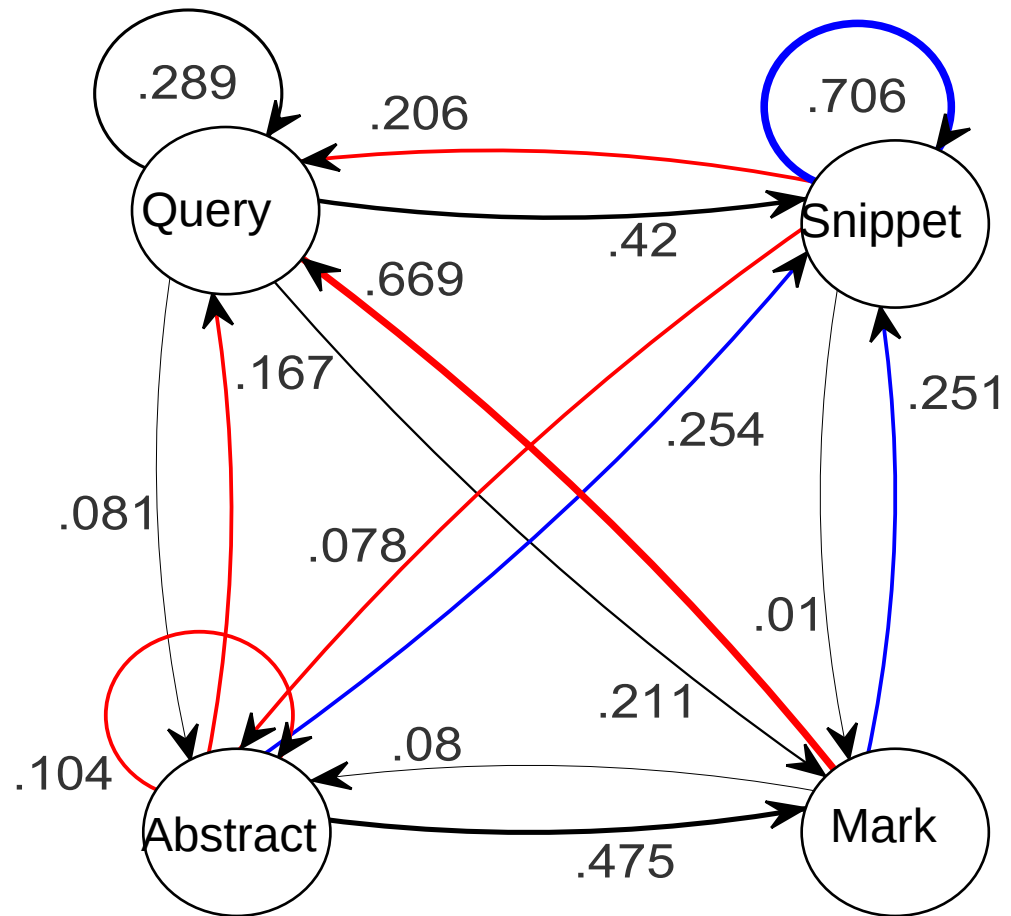
- 4 states per search phase
- Consider discrete signals + action times



# Differences Between Phases

In 2nd phase, users

- formulate more queries
- look at fewer snippets
- click more often on a snippet



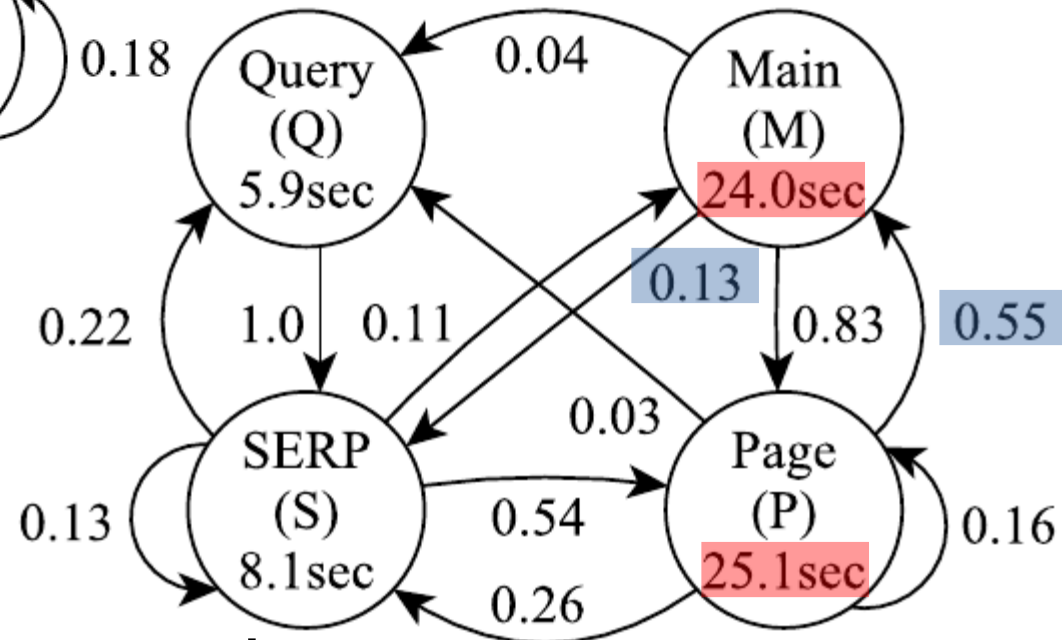
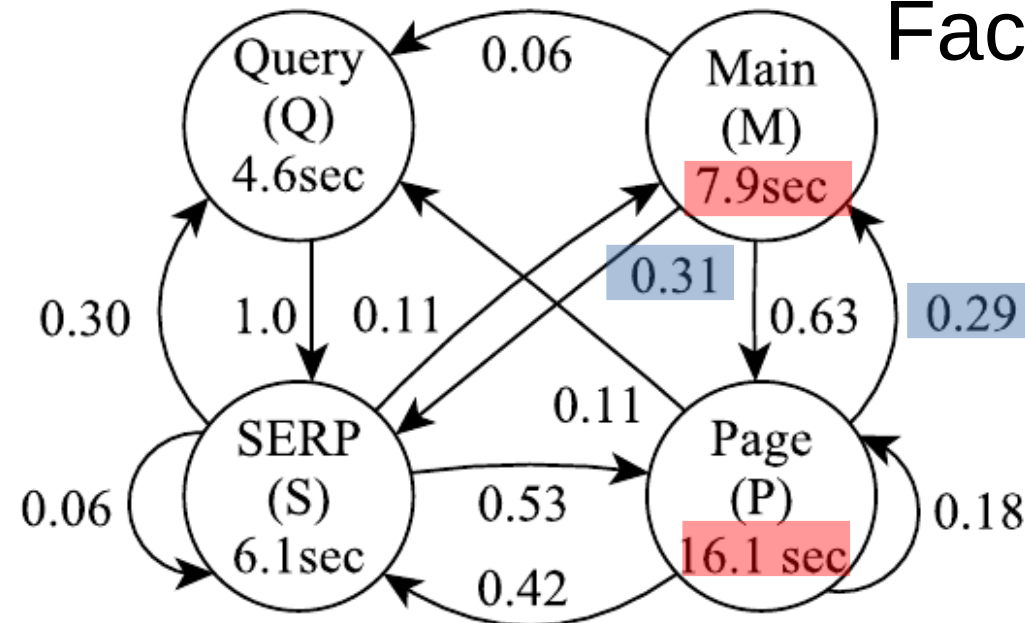
- Given: training sample with observation sequences for different classes  $c_k$
- Train separate HMM  $\lambda_k$  for each class  $c_k$
- Testing: determine most probable class  $C$  for given observation sequence  $\varepsilon_i$

$$C(\varepsilon_i) = \underset{c_k}{\operatorname{argmax}} P(\varepsilon_i | \lambda_k)$$

# Recognizing Type of Search

## Fact Finding

(Main contains the task description)



## Exploratory Search