

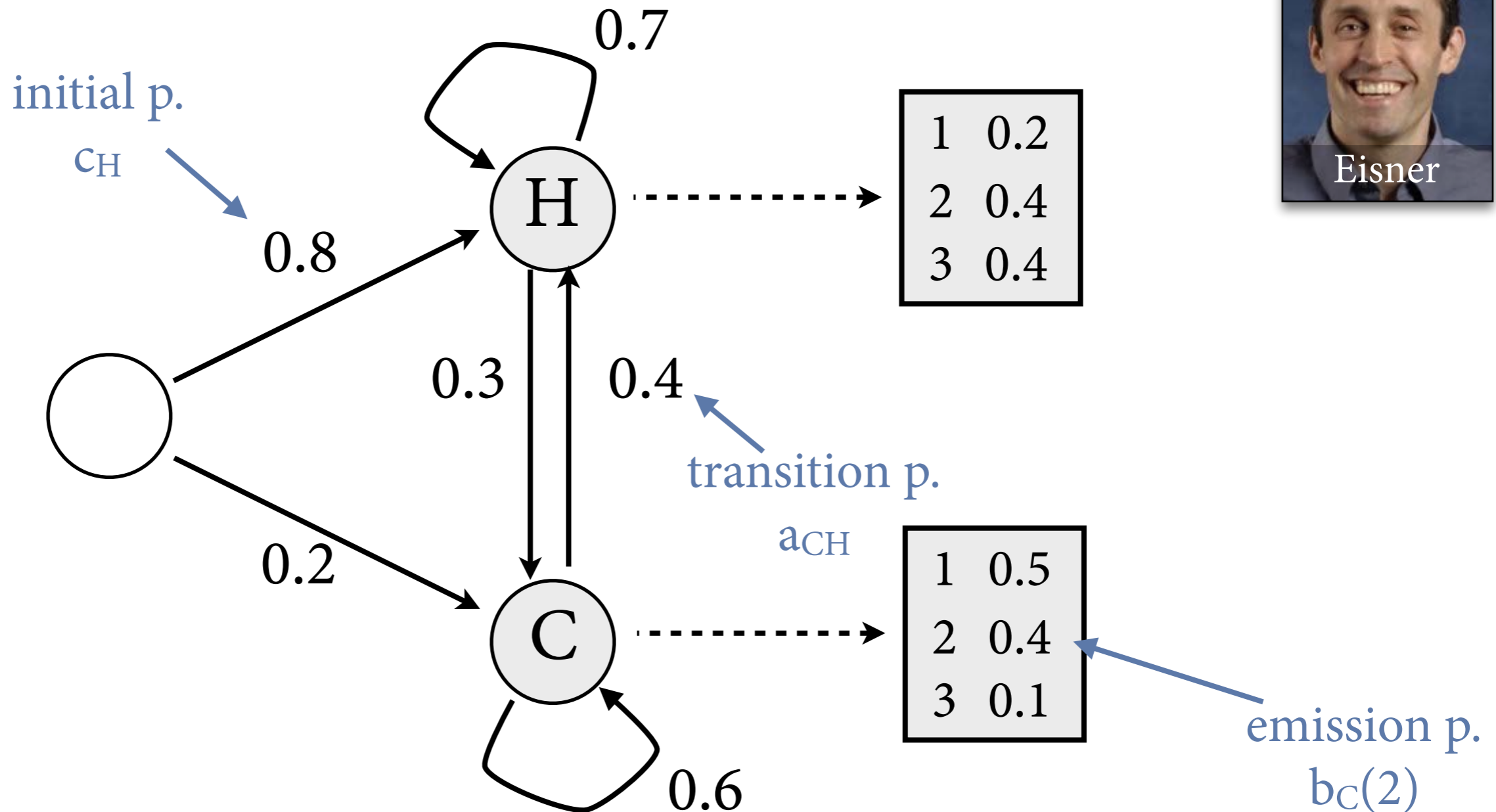
# Evaluating and Training HMMs

Computational Linguistics

Alexander Koller

9 November 2018

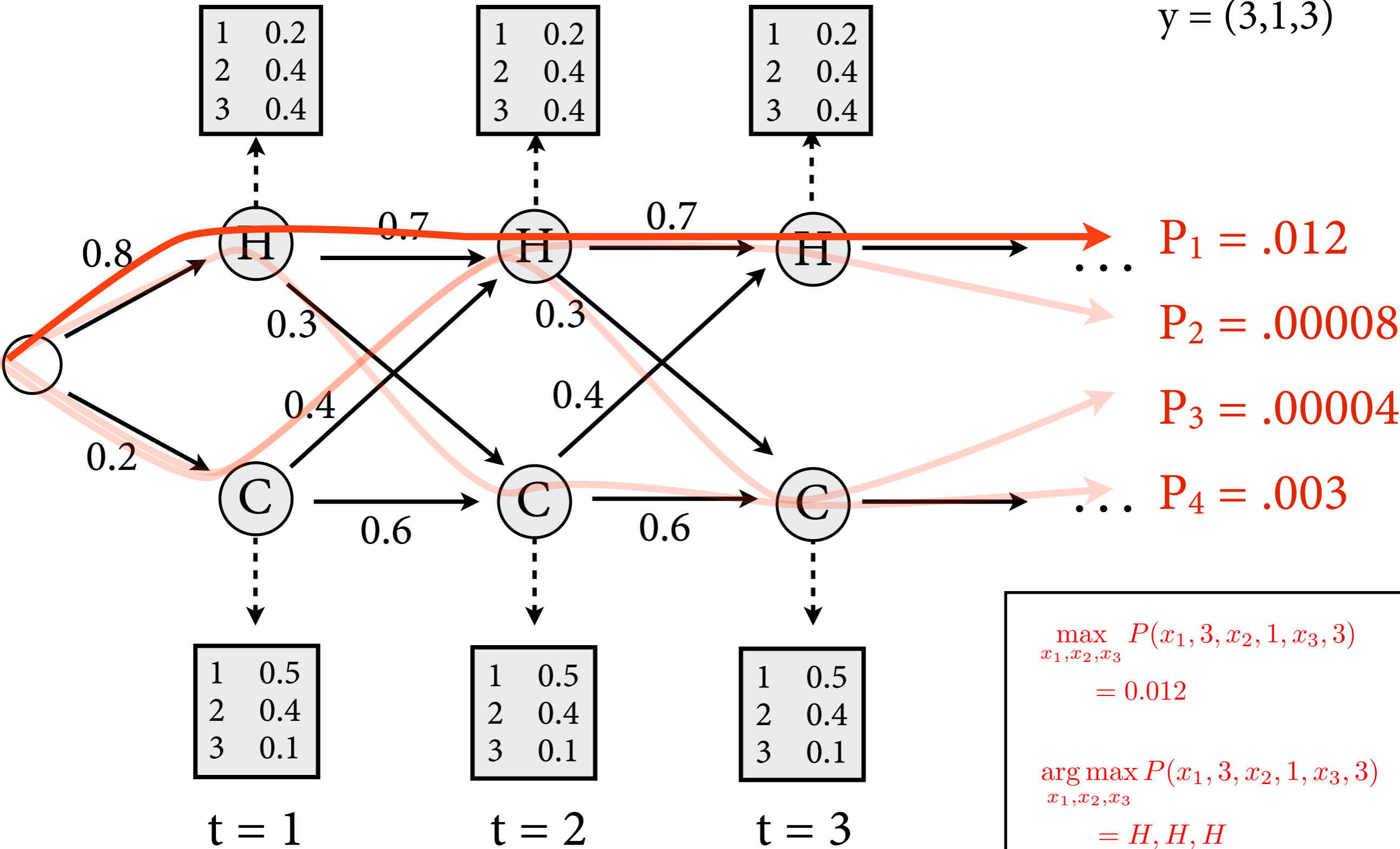
# Example HMM: Eisner's Ice Cream



States represent weather on a given day: Hot, Cold

Outputs represent number of ice creams Jason eats that day

# Ice cream trellis



1	0.2
2	0.4
3	0.4

1	0.2
2	0.4
3	0.4

1	0.2
2	0.4
3	0.4

1	0.5
2	0.4
3	0.1

1	0.5
2	0.4
3	0.1

1	0.5
2	0.4
3	0.1

$t = 1$

$t = 2$

$t = 3$

# The Viterbi Algorithm

$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

- Base case,  $t = 1$ :

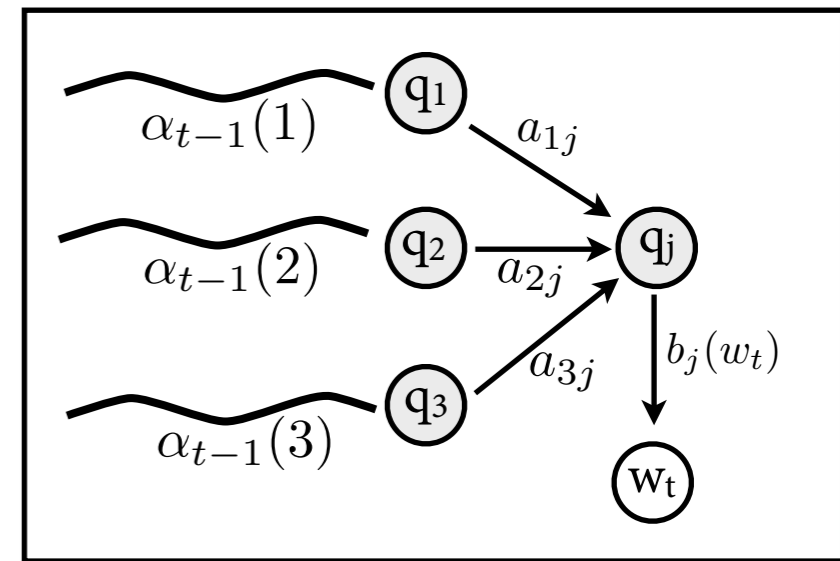
$$V_1(j) = b_j(y_1) \cdot c_j$$

- Inductive case, for  $t = 2, \dots, T$ :

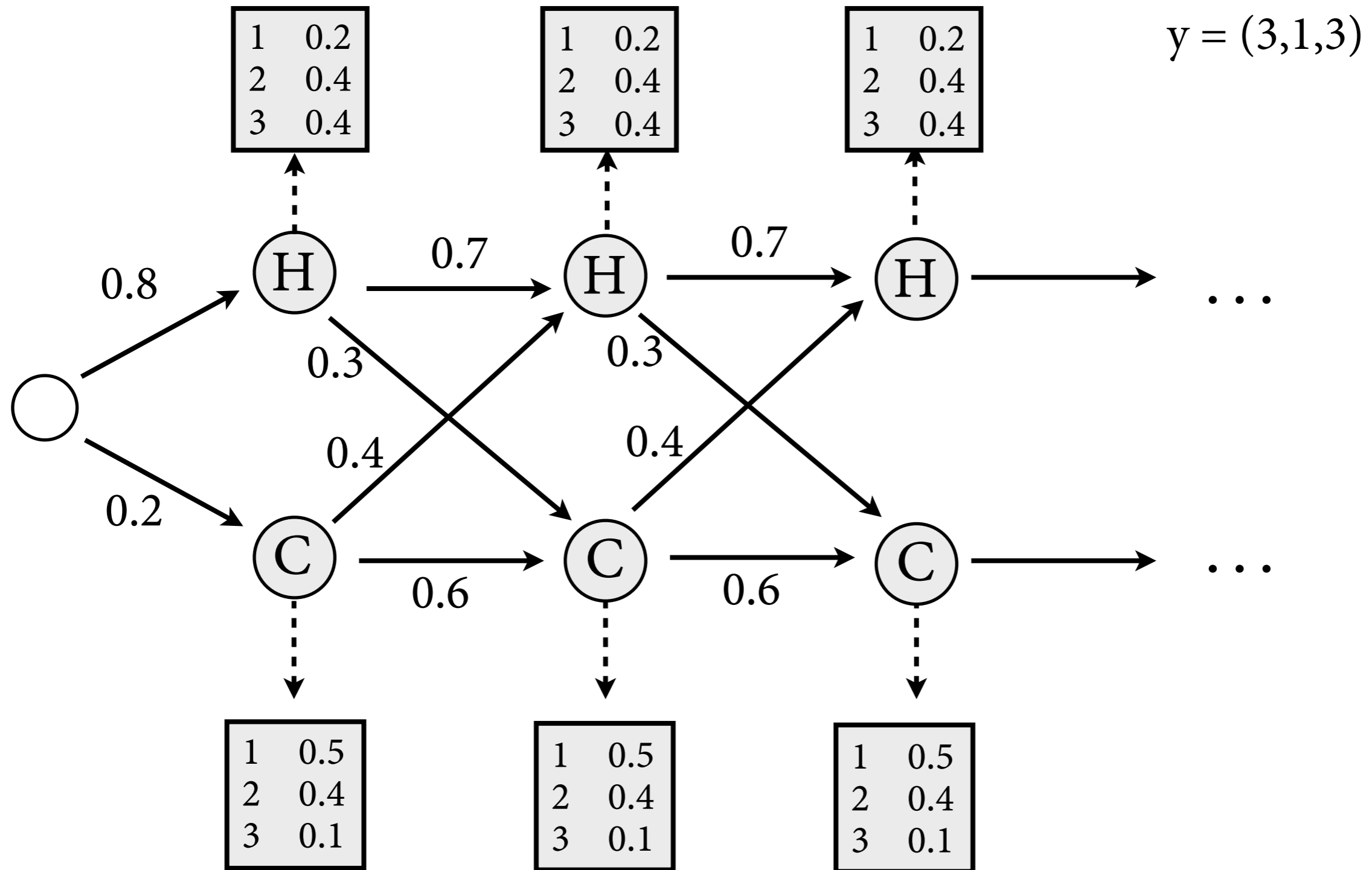
$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

- Once we have calculated all  $V$  values, we can easily calculate prob of best path:

$$\max_{x_1, \dots, x_T} P(x_1, y_1, \dots, x_T, y_T) = \max_{q \in Q} V_T(q)$$



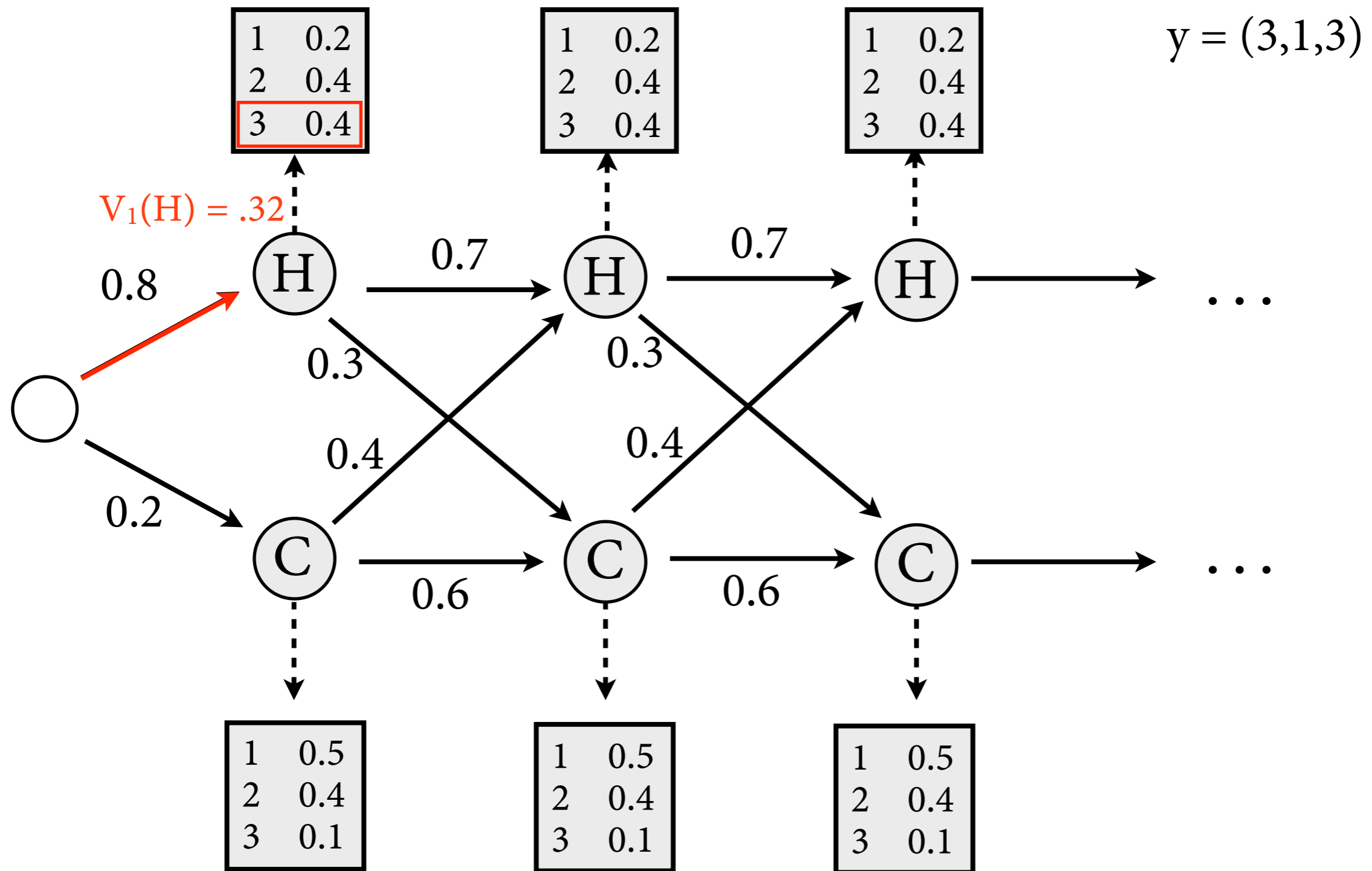
# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

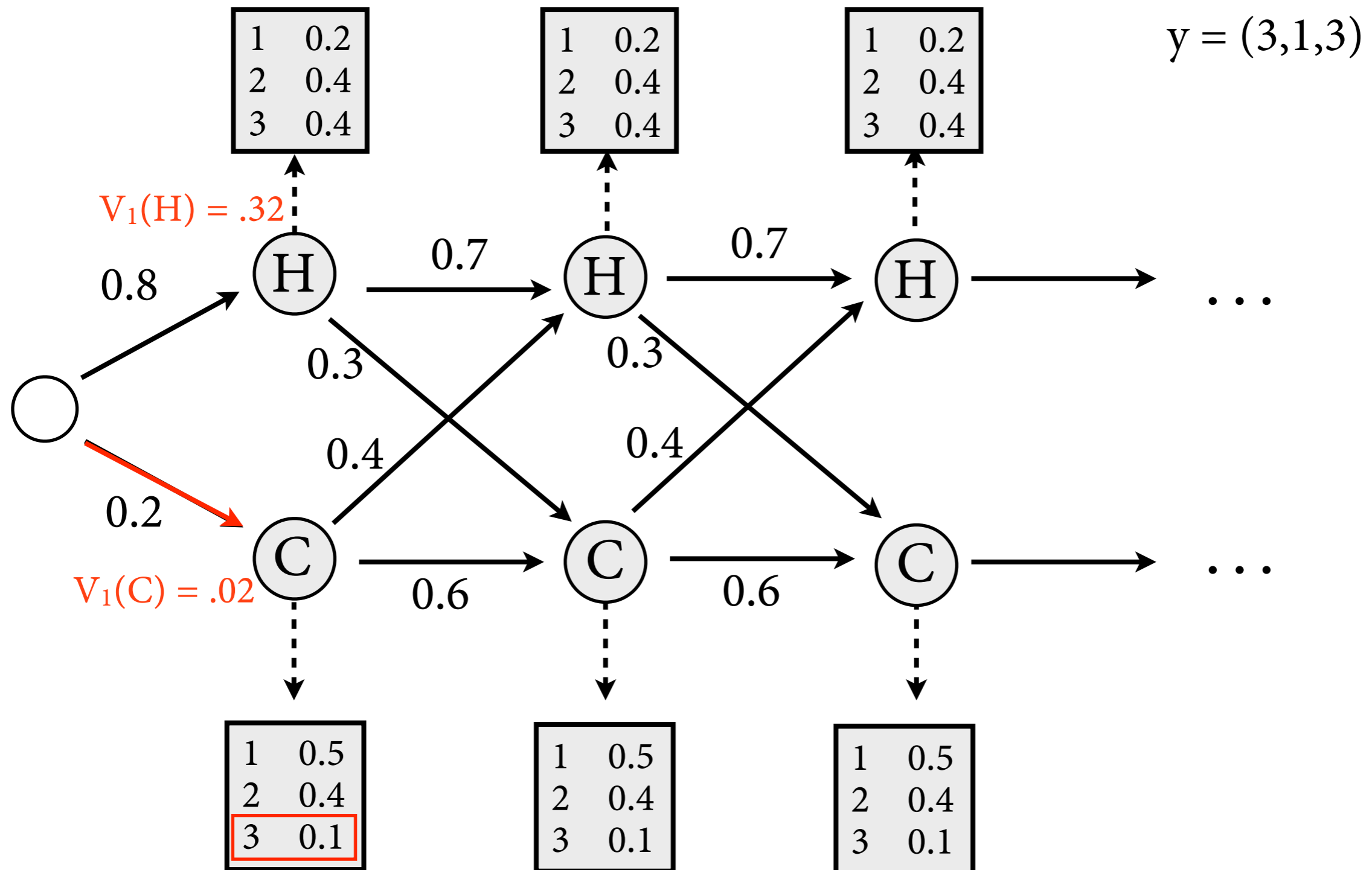
# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

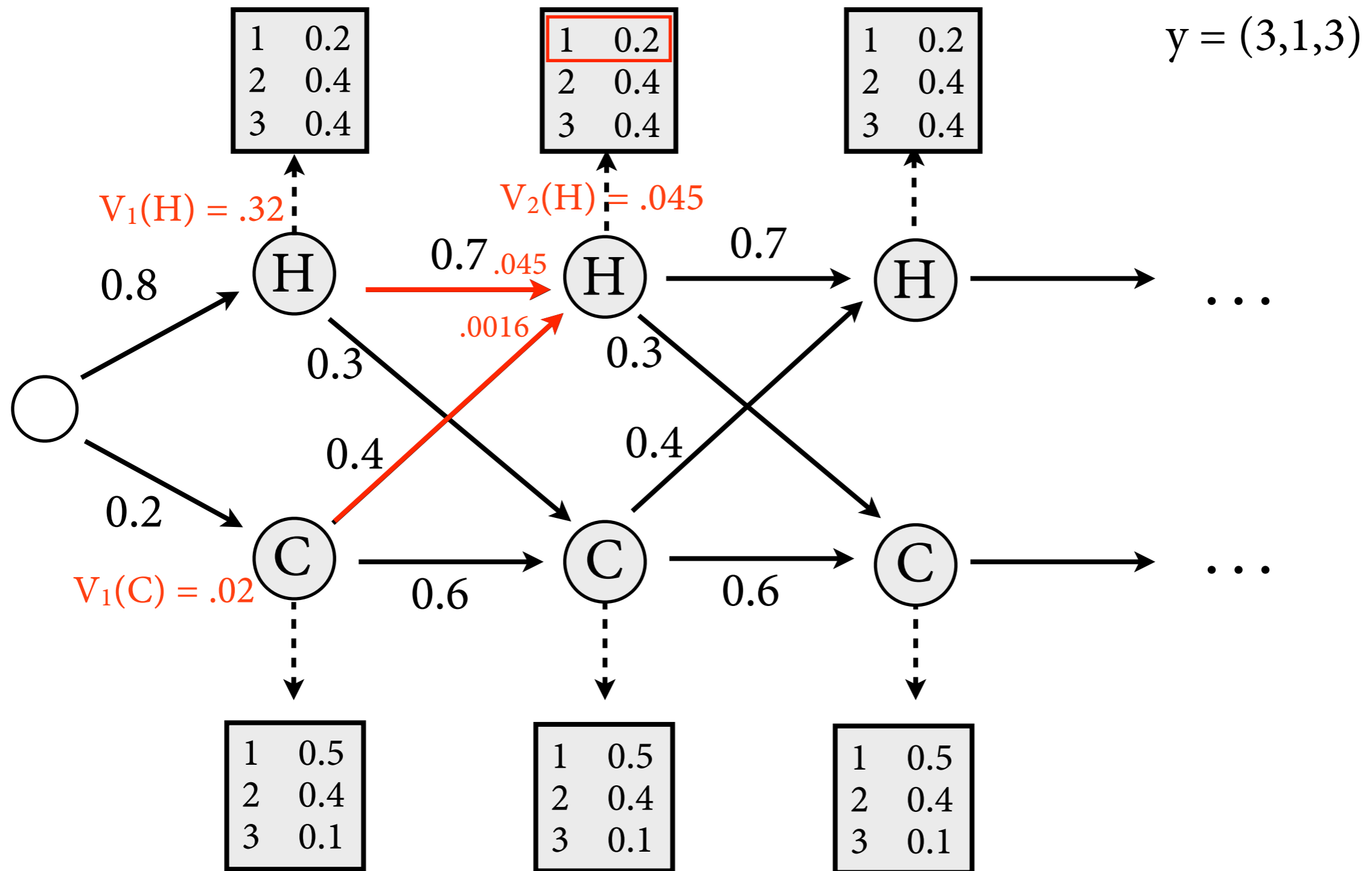
# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

# Viterbi Algorithm: Example

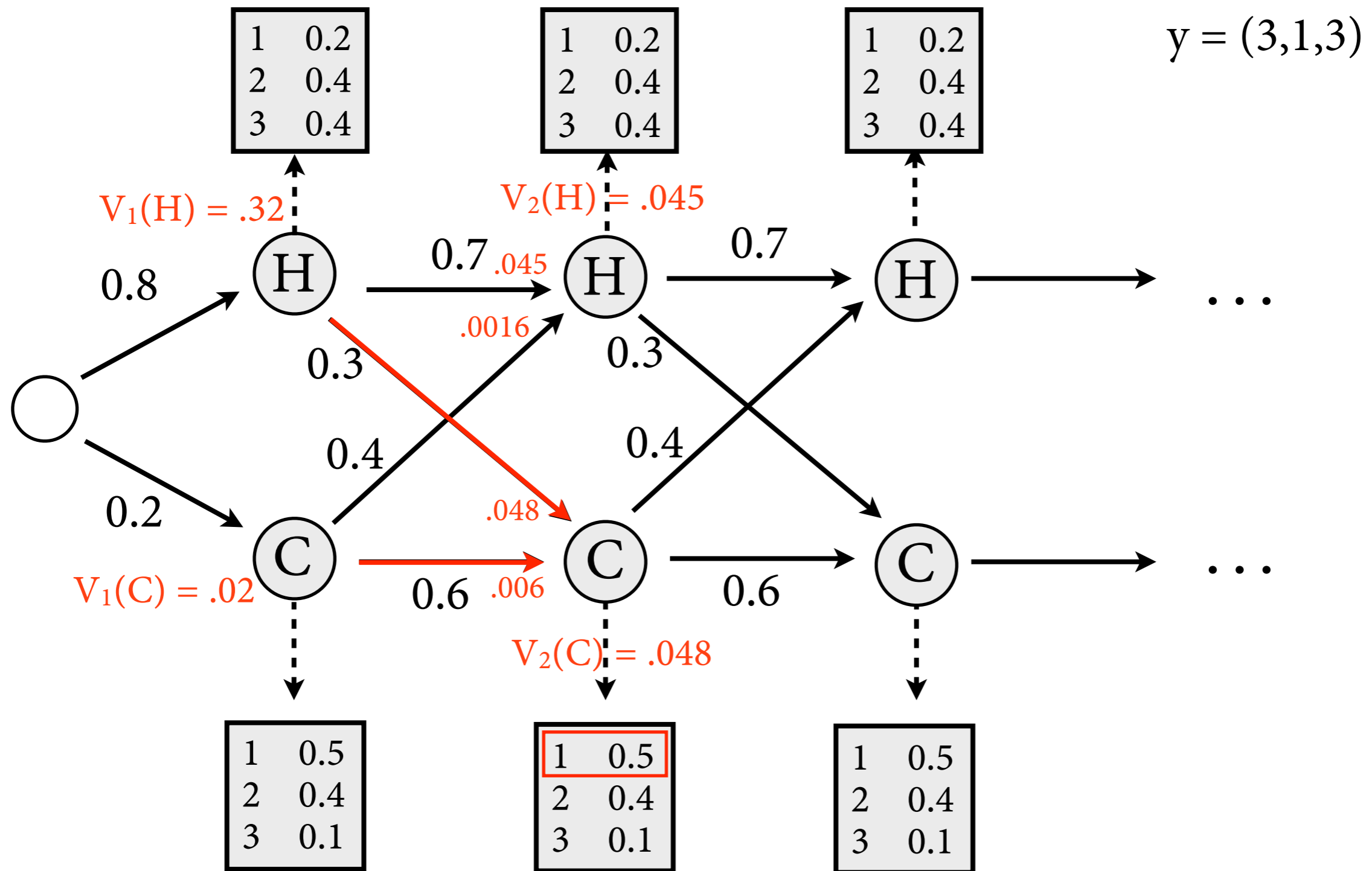


$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$



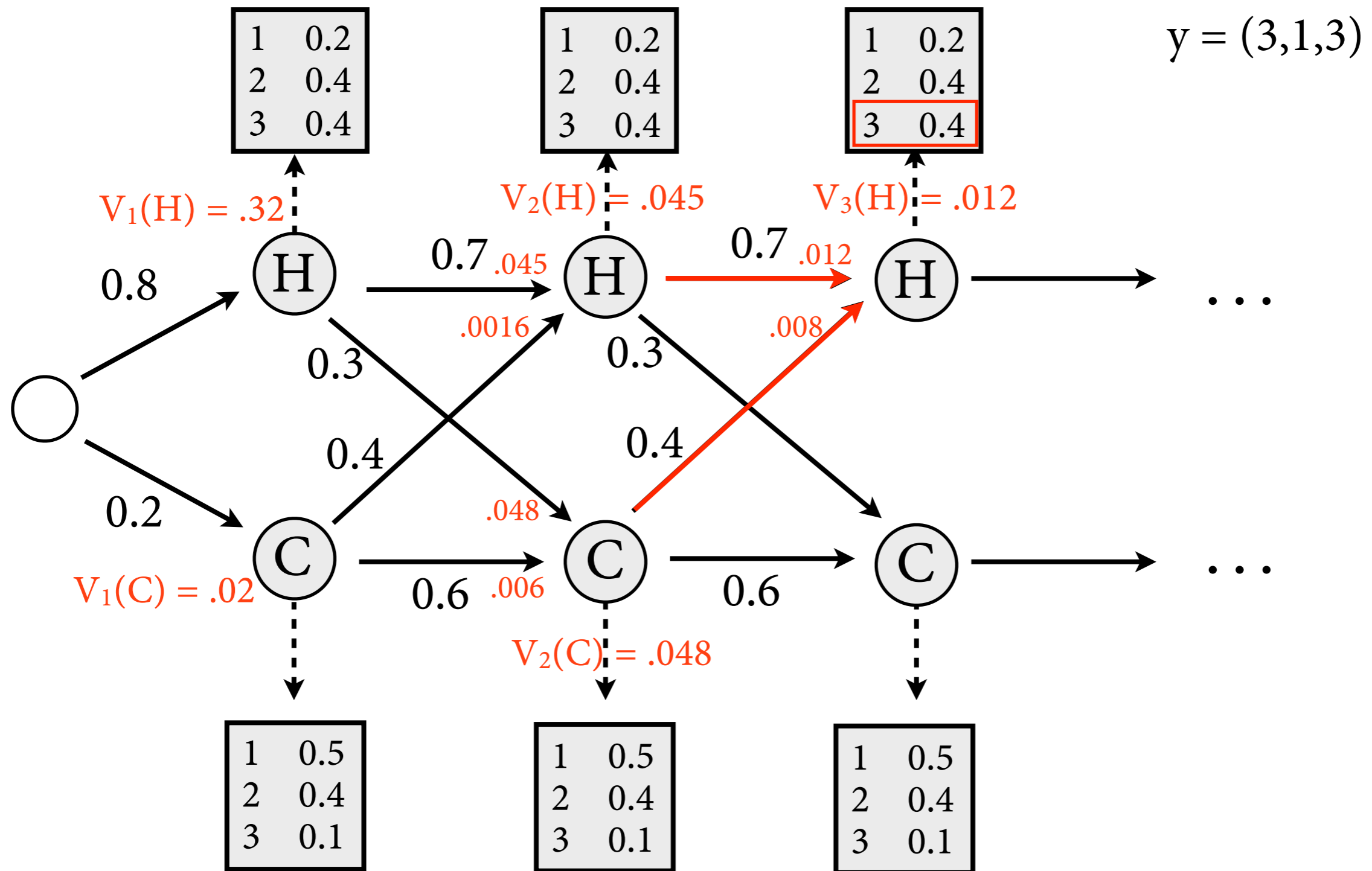
# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

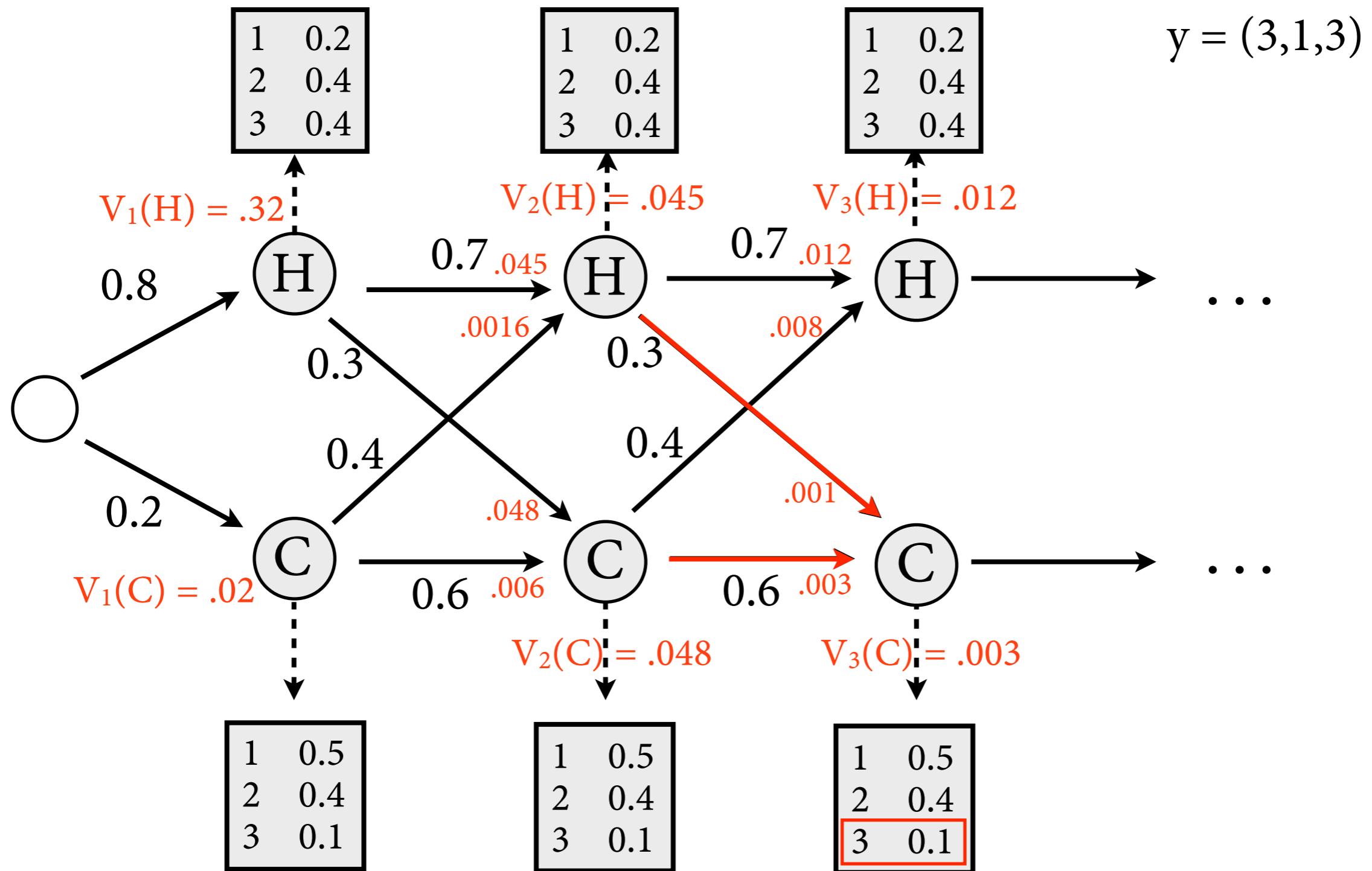
# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

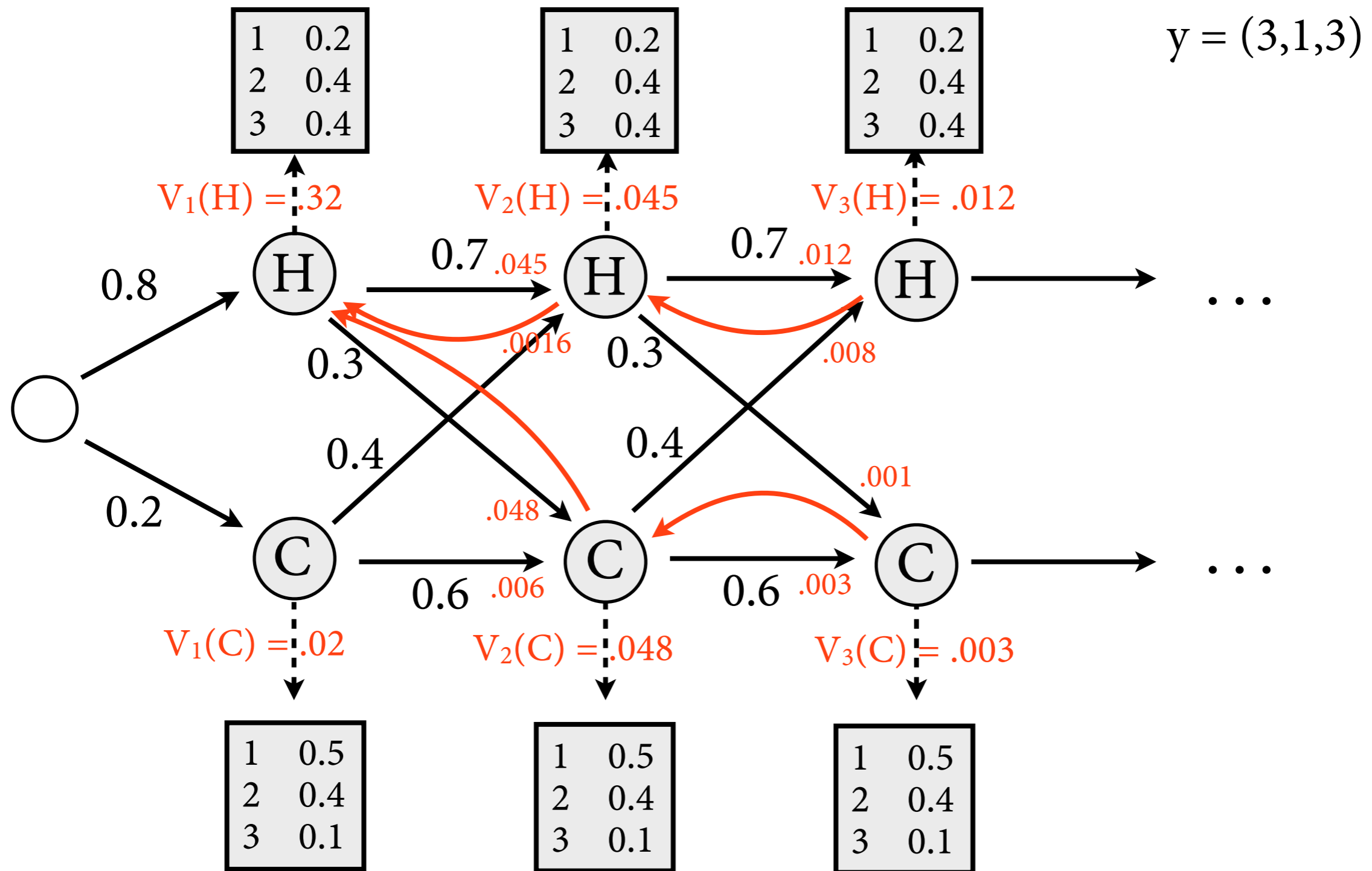
# Backpointers

- In the end, we need to reconstruct the sequence of states  $x_1, \dots, x_T$  with max probability.
- For each  $t, j$ : remember the value of  $i$  for which the maximum was achieved in *backpointer*  $bp_t(j)$ .

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

- Then just follow backpointers from right to left.

# Viterbi Algorithm: Example



$$V_t(j) = \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j)$$

$$V_t(j) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$

# Question 1: Language modeling

- Given an HMM and a string  $w_1, \dots, w_T$ , what is the likelihood  $P(w_1 \dots w_T)$ ?
- We can compute  $P(w_1 \dots w_T)$  efficiently with the *forward algorithm*.

DT	NN	VBD	NNS	IN	DT	NN	$p_1$
The	representative	put	chairs	on	the	table.	

DT	JJ	NN	VBZ	IN	DT	NN	$p_2$
The	representative	put	chairs	on	the	table.	

$$P(\text{sentence}) = p_1 + p_2$$

# The Forward Algorithm

- Key idea: *Forward probability*  $\alpha_t(j)$  that HMM outputs  $y_1, \dots, y_t$  and then ends in  $X_t = q_j$ .

$$\begin{aligned}\alpha_t(j) &= P(y_1, \dots, y_t, X_t = q_j) \\ &= \sum_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, X_1 = x_1, \dots, X_{t-1} = x_{t-1}, X_t = q_j)\end{aligned}$$

- From this, can compute easily

$$P(y_1, \dots, y_T) = \sum_{q \in Q} \alpha_T(q)$$

# The Forward Algorithm

$$\alpha_t(j) = P(y_1, \dots, y_t, X_t = q_j)$$

- Base case,  $t = 1$ :

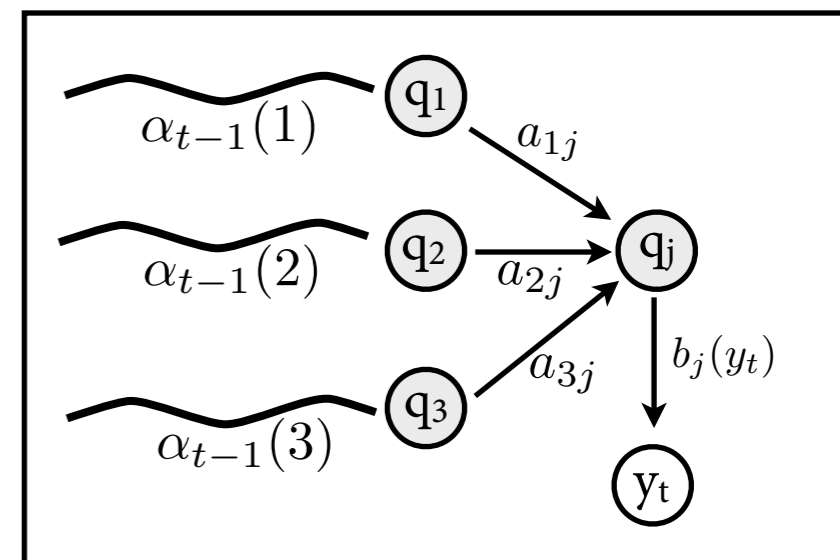
$$\alpha_1(j) = P(y_1, X_1 = q_j) = b_j(y_1) \cdot c_j$$

- Inductive case, compute for  $t = 2, \dots, T$ :

$$\alpha_t(j) = P(y_1, \dots, y_t, X_t = q_j)$$

$$= \sum_{i=1}^N P(y_1, \dots, y_{t-1}, X_{t-1} = q_i) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot P(y_t \mid X_t = q_j)$$

$$= \sum_{i=1}^N \alpha_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t)$$





# Question 3a: Supervised learning

- Given a set of POS tags and *annotated* training data  $(w_1, t_1), \dots, (w_T, t_T)$ , compute parameters for HMM that maximize likelihood of training data.

DT                  NN                  VBD      NNS      IN      DT      NN  
The representative put chairs on the table.

NNP      VBZ      VBN      TO      VB                  NR  
Secretariat is expected to race tomorrow.

# Maximum likelihood training

- Estimate bigram model for state sequence:

$$a_{ij} = \frac{C(X_t = q_i, X_{t+1} = q_j)}{C(X_t = q_i)} \quad c_j = \frac{\# \text{ sentences with } X_1 = q_j}{\# \text{ sentences}}$$

- ML estimate for emission probabilities:

$$b_i(o) = \frac{C(X_t = q_i, Y_t = o)}{C(X_t = q_i)}$$

- Apply smoothing as you would for ordinary n-gram models (increase all counts C by one).

# Evaluation

- How do you know how well your tagger works?
- Run it on *test data* and evaluate *accuracy*.
  - ▶ Test data: Really important to evaluate on unseen sentences to get a fair picture of how well tagger generalizes.
  - ▶ Accuracy: Measure percentage of correctly predicted POS tags.

# Evaluation on test data

DT          NN          VBD   NNS   IN   DT   NN  
The representative put chairs on the table.

NNP   VBZ   VBN   TO   VB          NR  
Secretariat is expected to race tomorrow.

Training corpus (annotated)

*Training*

Trained system  
(e.g. HMM)

Training

# Evaluation on test data

DT NN VBD NNS IN DT NN  
The representative put chairs on the table.

NNP VBZ VBN TO VB NR  
Secretariat is expected to race tomorrow.

Training corpus (annotated)

*Training*

Trained system  
(e.g. HMM)

Training

NNP VBZ NNP  
John loves Mary.

Test corpus (annotated)

Evaluation

# Evaluation on test data

DT      NN      VBD   NNS   IN   DT   NN  
The representative put chairs on the table.

NNP   VBZ   VBN   TO   VB      NR  
Secretariat is expected to race tomorrow.

Training corpus (annotated)

*Training*

Trained system  
(e.g. HMM)

Training

NNP      VBZ      NNP  
John   loves   Mary.

Test corpus (annotated)

Test corpus  
(without annotations)

John   loves   Mary.

Evaluation

# Evaluation on test data

DT      NN      VBD   NNS   IN   DT   NN  
The representative put chairs on the table.

NNP   VBZ   VBN   TO   VB      NR  
Secretariat is expected to race tomorrow.

Training corpus (annotated)

*Training*

Trained system  
(e.g. HMM)

Training

NNP      VBZ      NNP  
John   loves   Mary.

Test corpus (annotated)

Test corpus  
(without annotations)

John   loves   Mary.

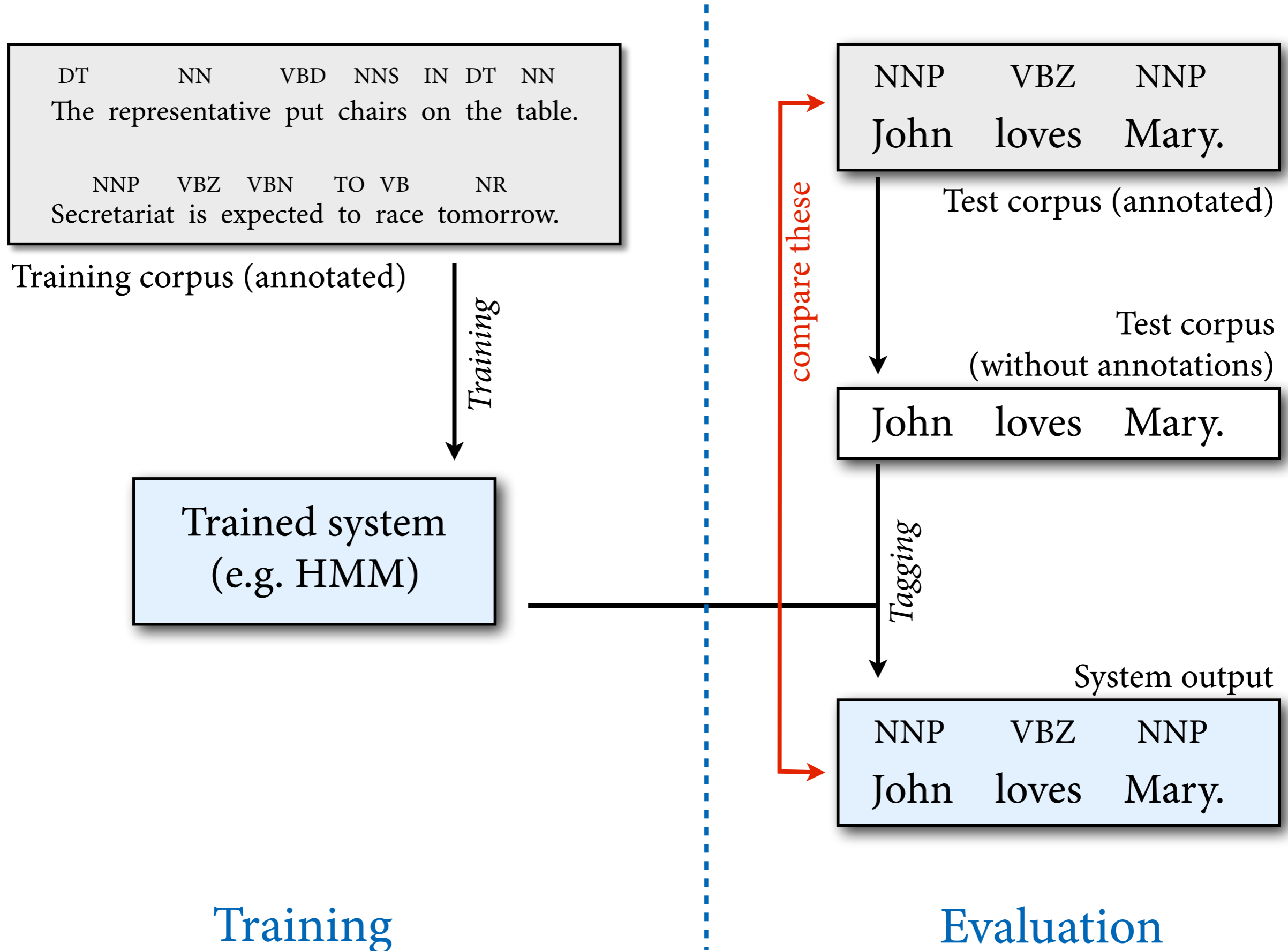
*Tagging*

System output

NNP      VBZ      NNP  
John   loves   Mary.

Evaluation

# Evaluation on test data





# Question 3b: Unsupervised learning

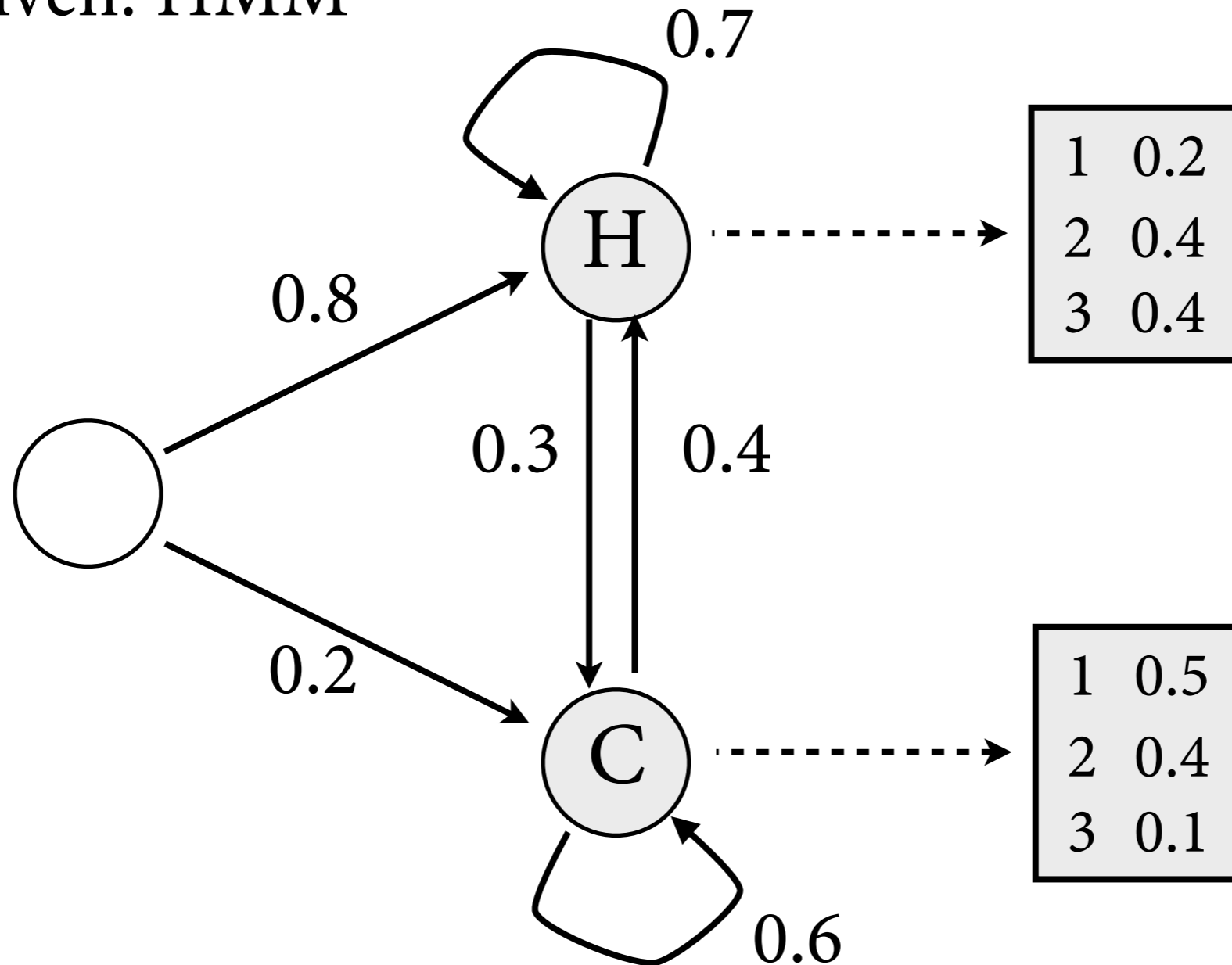
- Given a set of POS tags and *unannotated* training data  $w_1, \dots, w_T$ , compute parameters for HMM that maximize likelihood of training data.
- Relevant because annotated data is expensive to obtain, but raw text is really cheap.

The representative put chairs on the table.

Secretariat is expected to race today.

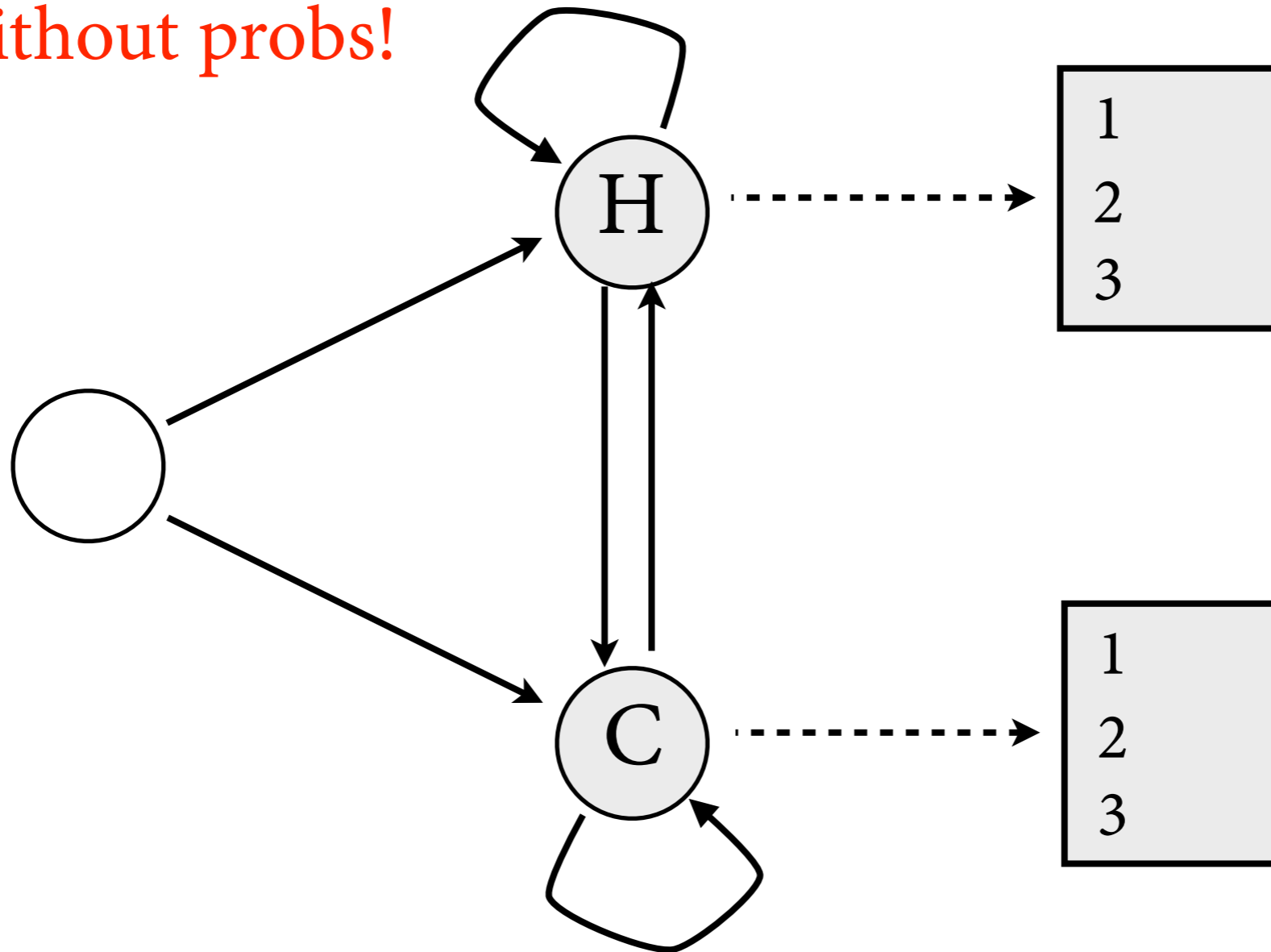
# The setup

Given: HMM



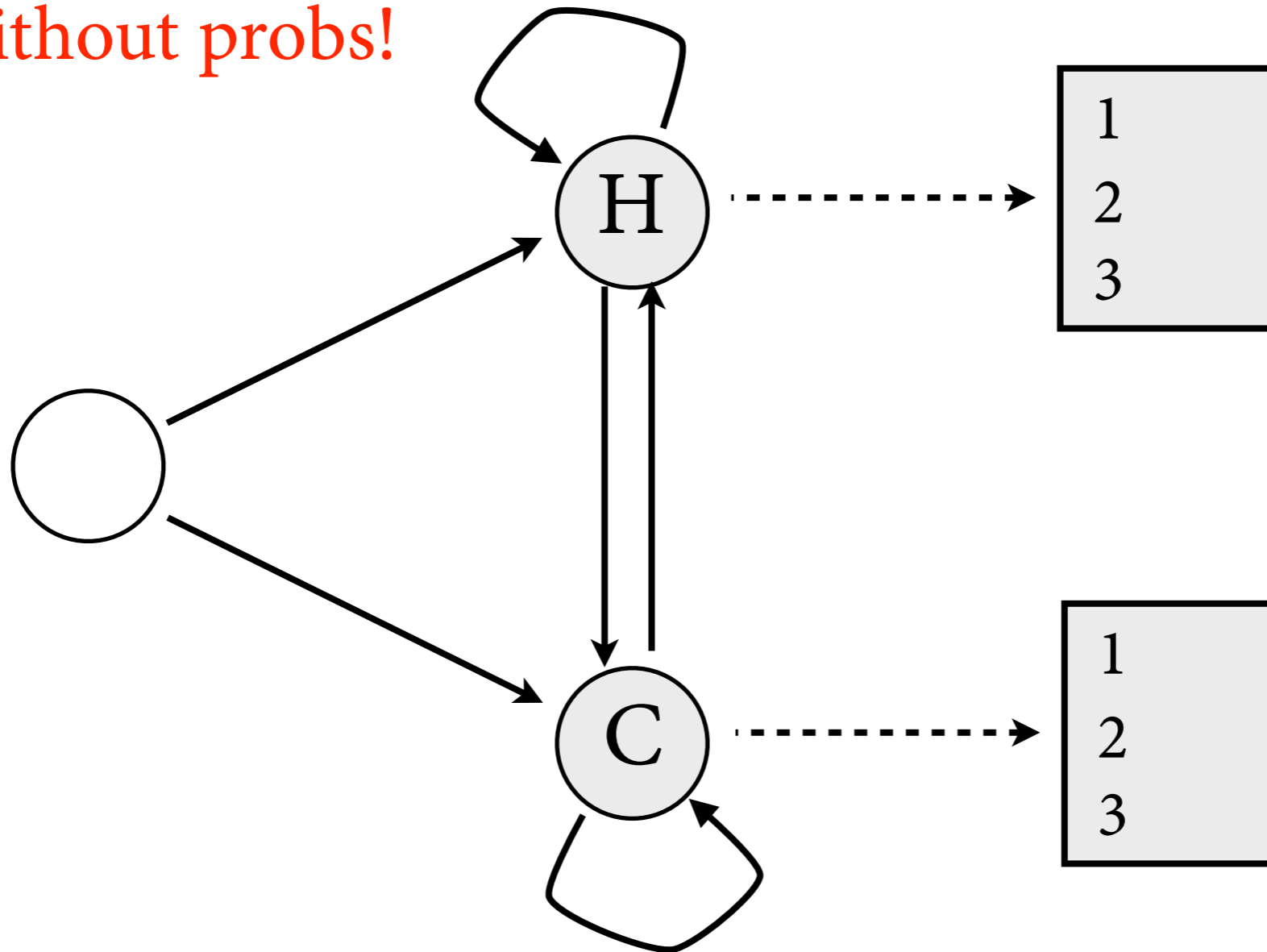
# The setup

Given: HMM  
without probs!



# The setup

Given: HMM  
without probs!



Observations: 2, 3, 3, 2, 3, 2, 3, 2, 2, 3, 1, 3, 3, ...

# The setup

- If we had counts of state transitions in corpus, we could simply use ML estimation.

$$a_{ij} = \frac{C(q_i \rightarrow q_j)}{C(q_i \rightarrow \bullet)}$$

- Idea: replace actual counts by *estimated* counts.

$$a_{ij} \approx \frac{\hat{C}(q_i \rightarrow q_j)}{\hat{C}(q_i \rightarrow \bullet)}$$

- How can we estimate counts?

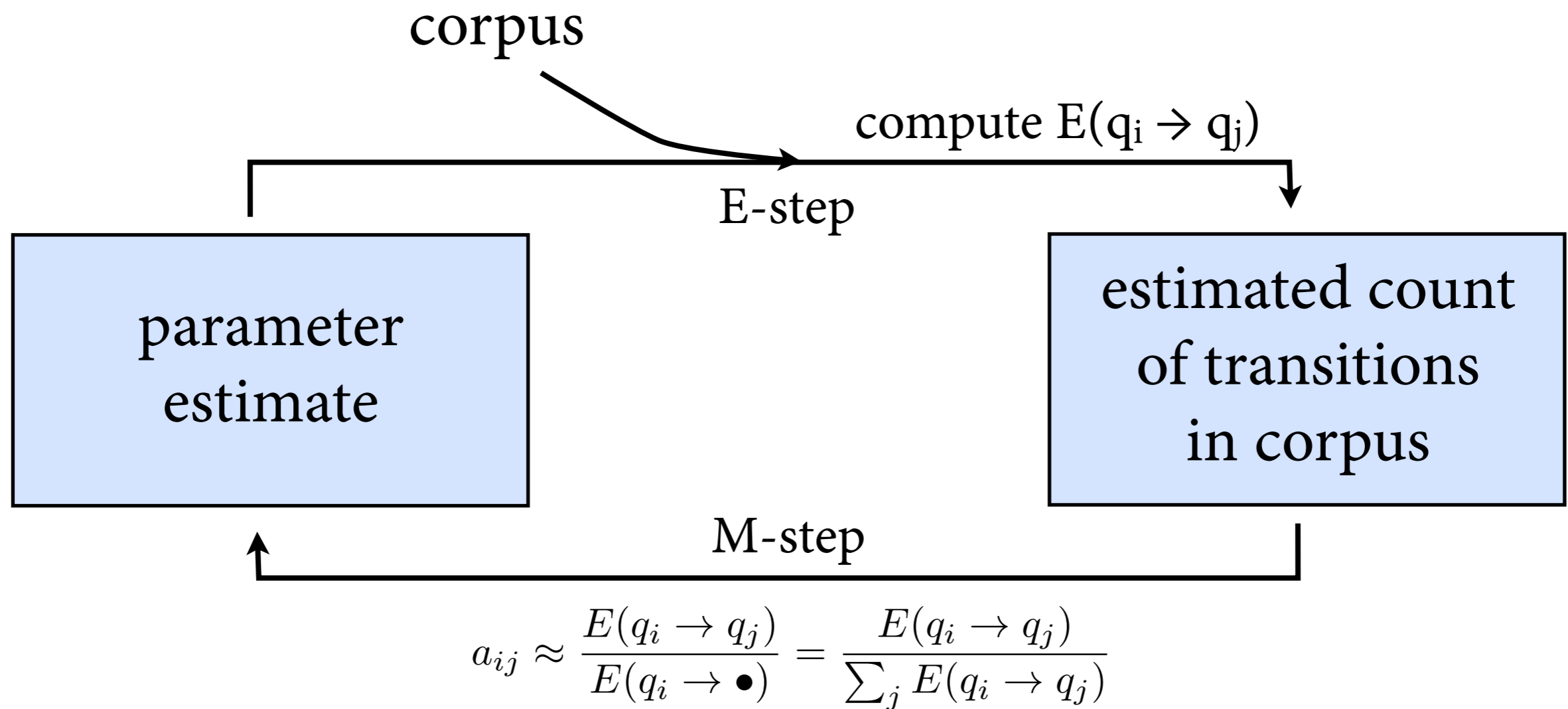
# Estimated counts

Observations $y$	3	1	3	$C(H \rightarrow H)$	$P(x   y)$	
Hidden states $x$	H	H	H	2	*	0.408
	H	H	C	+ 1	*	0.034
	H	C	H	+ 0	*	0.272
		...				
	C	C	C	+ 0	*	0.136
				<hr/>		
				<b>0.864</b>		

$$\hat{C}(q_i \rightarrow q_j) = E(q_i \rightarrow q_j) = \sum_x \hat{P}(x | y) \cdot C(q_i \rightarrow q_j \text{ in } x)$$

$$= \sum_{t=1}^{M-1} \hat{P}(X_t = q_i, X_{t+1} = q_j | y)$$

# Expectation Maximization



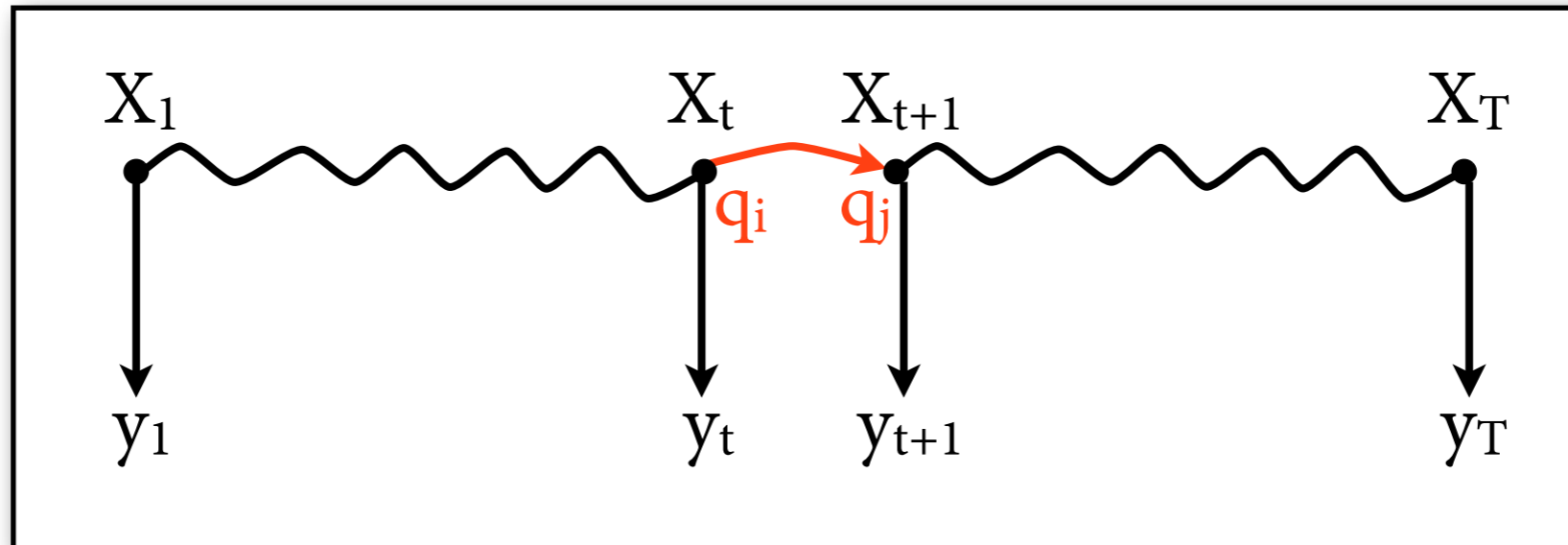
# Plan for computing E

$$E(q_i \rightarrow q_j) = \sum_{t=1}^{M-1} \hat{P}(X_t = q_i, X_{t+1} = q_j | y)$$

- How can we compute  $\hat{P}$  efficiently?  
Challenge: It is conditioned on  $y$ .
- We compute  $\xi_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j | y)$   
$$= \frac{\hat{P}(X_t = q_i, X_{t+1} = q_j, y)}{\hat{P}(y)}$$
- Do it in two steps:
  - ▶ compute  $\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$
  - ▶ compute  $P(y)$

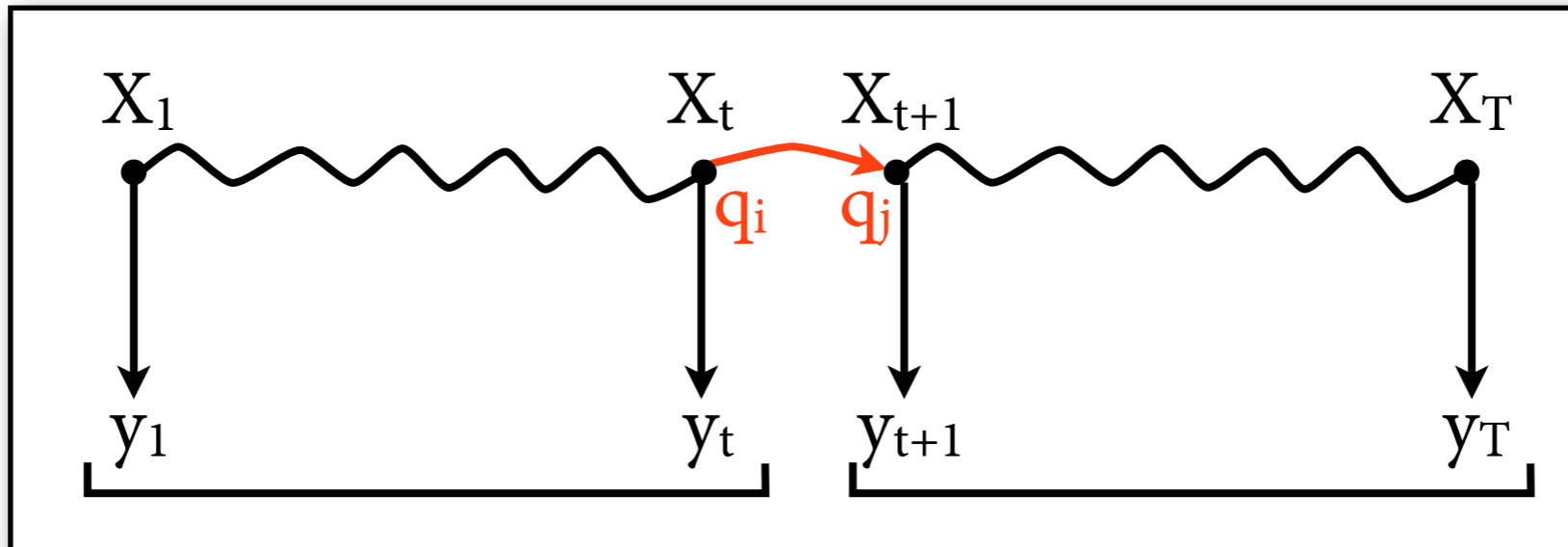


$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$



$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

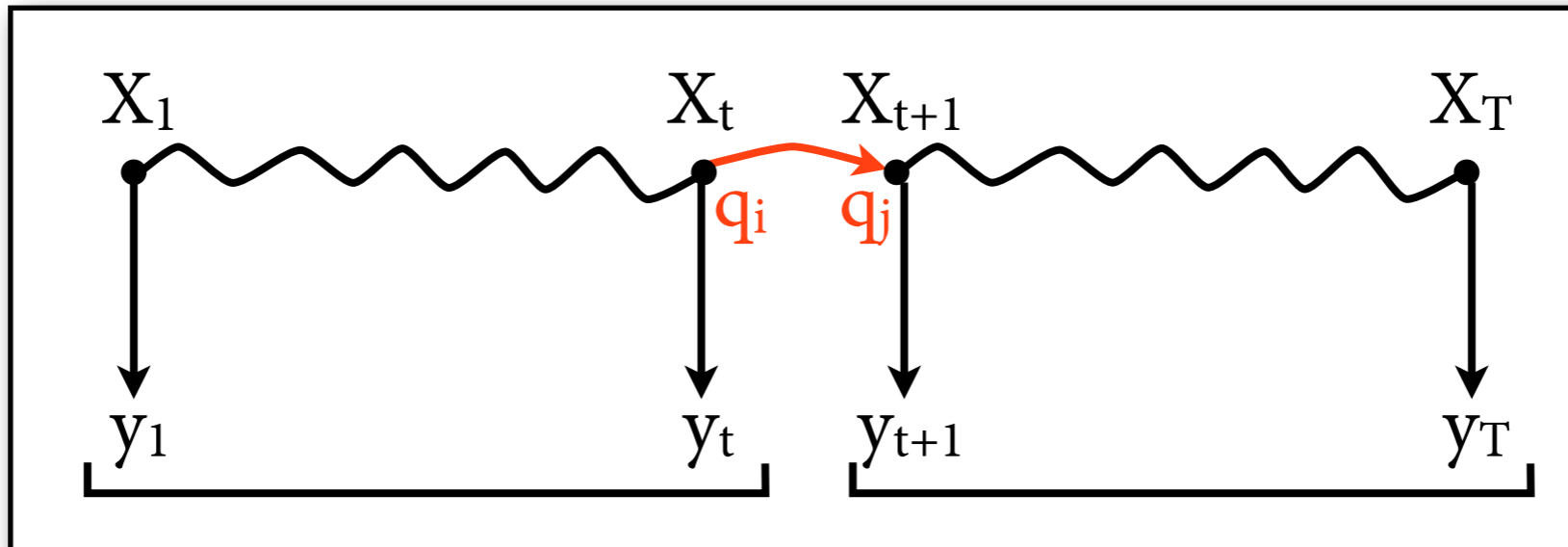
$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$



$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid y_1, \dots, y_t, X_t = q_i) \\ \cdot \hat{P}(y_{t+2}, \dots, y_T \mid y_1, \dots, y_{t+1}, X_t = q_i, X_{t+1} = q_j)$$

$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

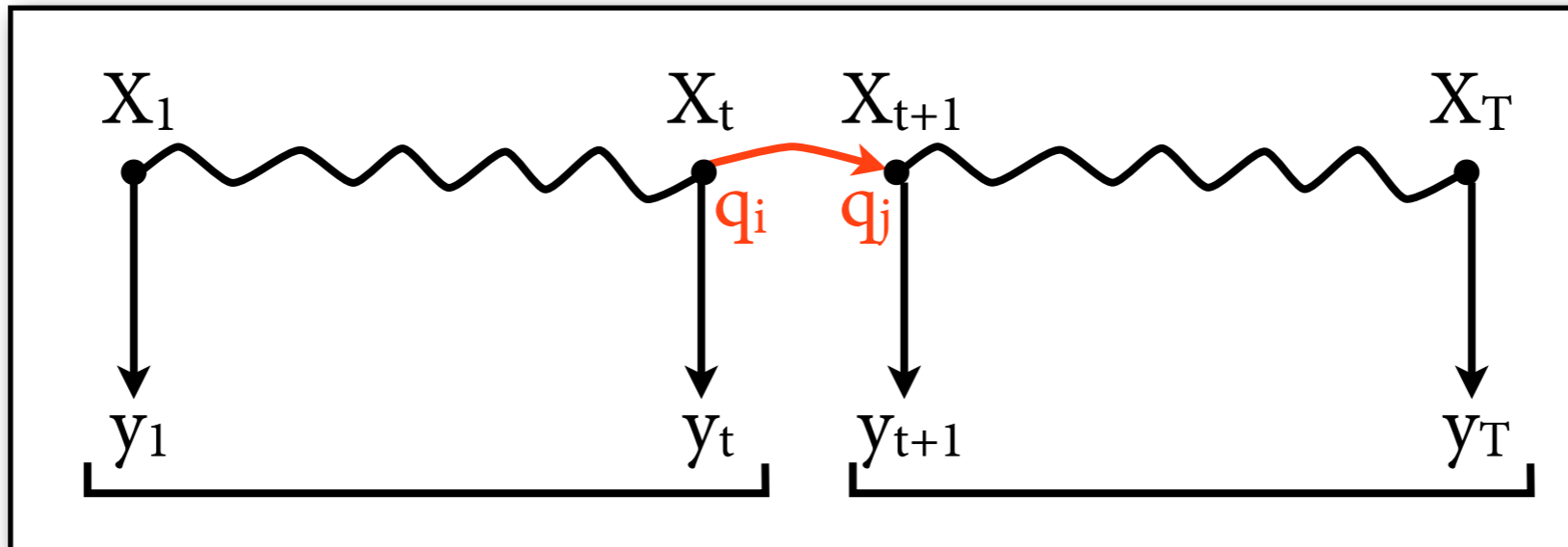


$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid y_1, \dots, y_t, X_t = q_i) \\ \cdot \hat{P}(y_{t+2}, \dots, y_T \mid y_1, \dots, y_{t+1}, X_t = q_i, X_{t+1} = q_j)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid X_t = q_i) \cdot \hat{P}(y_{t+2}, \dots, y_T \mid X_{t+1} = q_j)$$

$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

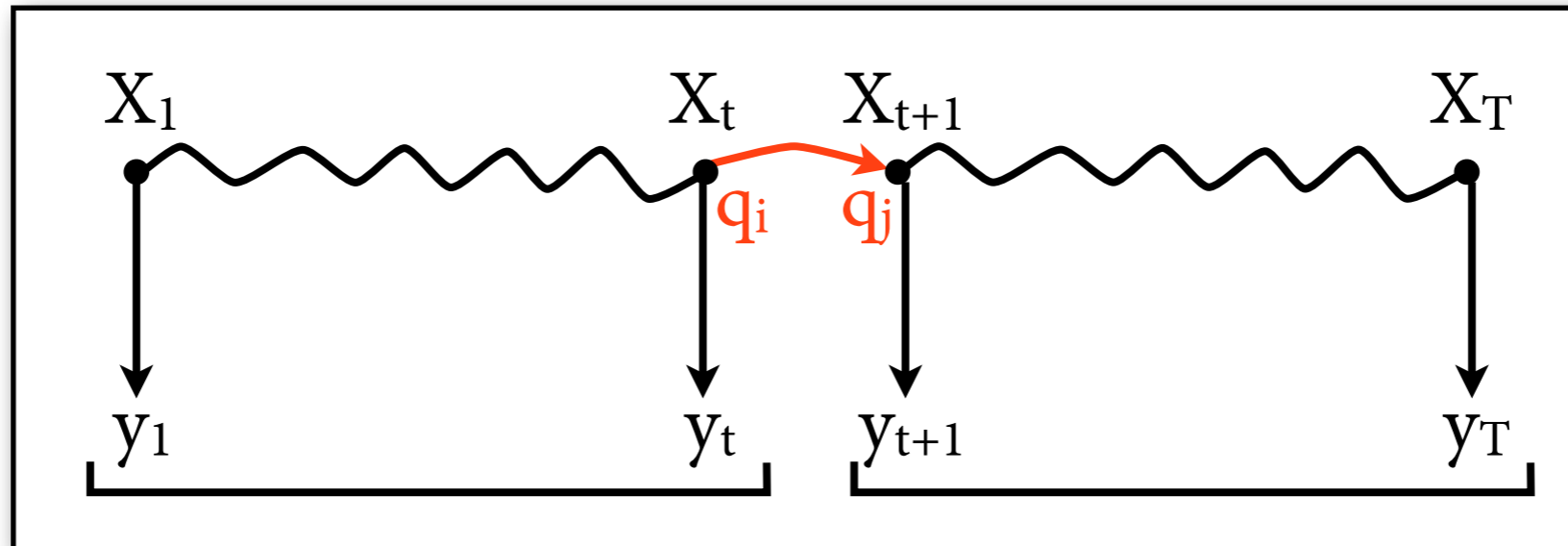


$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid y_1, \dots, y_t, X_t = q_i) \\ \cdot \hat{P}(y_{t+2}, \dots, y_T \mid y_1, \dots, y_{t+1}, X_t = q_i, X_{t+1} = q_j)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid X_t = q_i) \cdot \hat{P}(y_{t+2}, \dots, y_T \mid X_{t+1} = q_j) \\ \cdot a_{ij} \cdot b_j(w_{t+1}) \cdot$$

$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$



$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid y_1, \dots, y_t, X_t = q_i) \\ \cdot \hat{P}(y_{t+2}, \dots, y_T \mid y_1, \dots, y_{t+1}, X_t = q_i, X_{t+1} = q_j)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid X_t = q_i) \cdot \hat{P}(y_{t+2}, \dots, y_T \mid X_{t+1} = q_j)$$

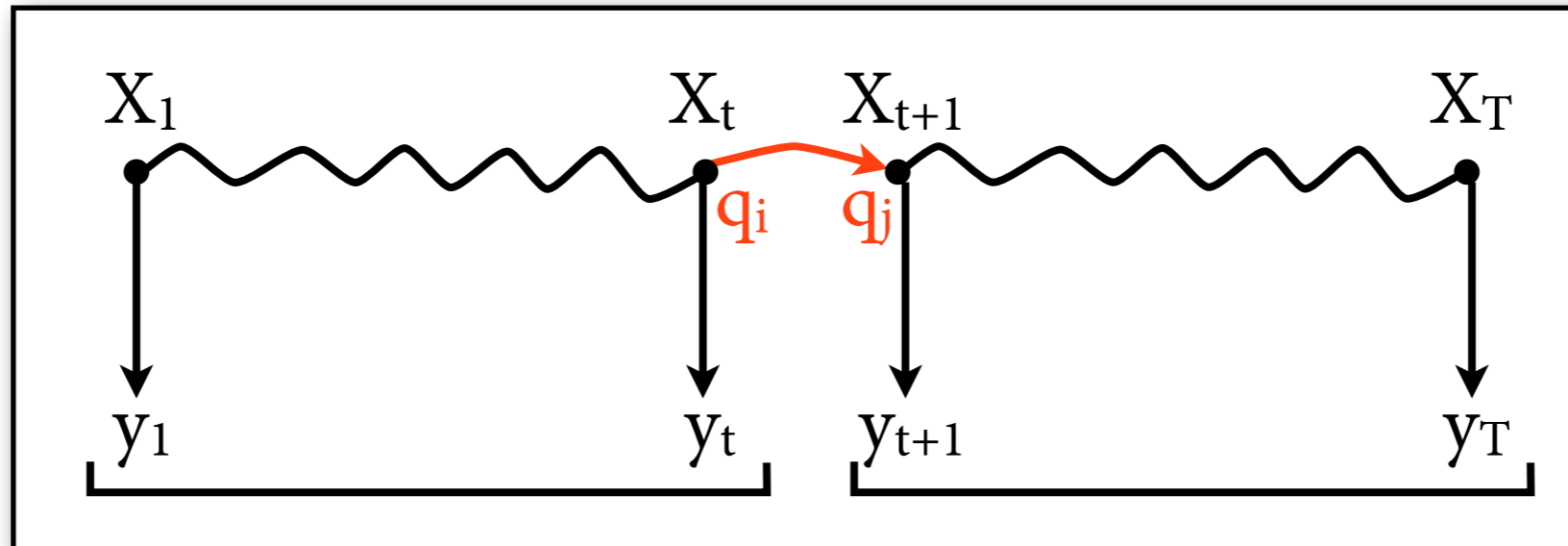
$$= \alpha_t(i)$$

$$\cdot a_{ij} \cdot b_j(w_{t+1}) \cdot$$

forward prob:

$$\alpha_t(i) = P(y_1, \dots, y_t, X_t = q_i)$$

$$\xi'_t(i, j) = \hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$



$$\hat{P}(X_t = q_i, X_{t+1} = q_j, y)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid y_1, \dots, y_t, X_t = q_i) \\ \cdot \hat{P}(y_{t+2}, \dots, y_T \mid y_1, \dots, y_{t+1}, X_t = q_i, X_{t+1} = q_j)$$

$$= \hat{P}(y_1, \dots, y_t, X_t = q_i) \cdot \hat{P}(y_{t+1}, X_{t+1} = q_j \mid X_t = q_i) \cdot \hat{P}(y_{t+2}, \dots, y_T \mid X_{t+1} = q_j)$$

$$= \alpha_t(i)$$

$$\cdot a_{ij} \cdot b_j(w_{t+1}) \cdot$$

$$\beta_{t+1}(j)$$

forward prob:

$$\alpha_t(i) = P(y_1, \dots, y_t, X_t = q_i)$$

backward prob:

$$\beta_t(i) = P(y_{t+1}, \dots, y_T \mid X_t = q_i)$$

# Backward probabilities

$$\beta_t(i) = P(y_{t+1}, \dots, y_T \mid X_t = q_i)$$

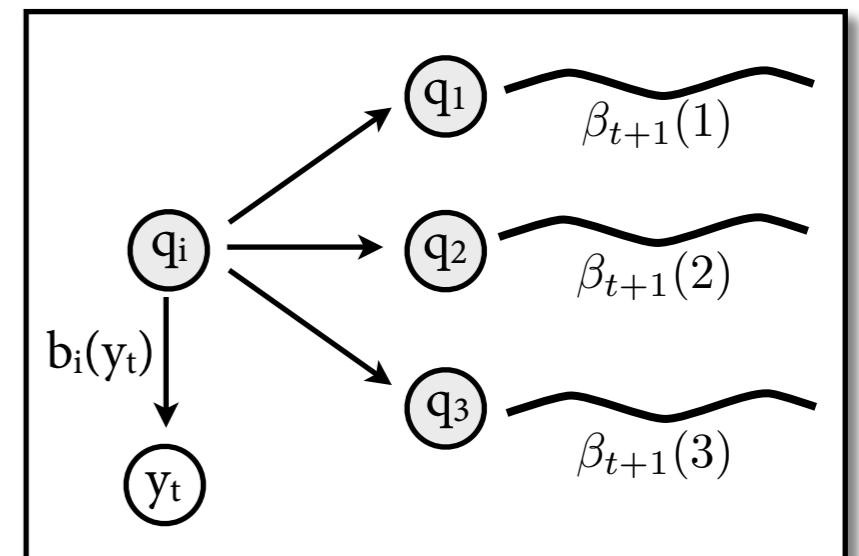
- Base case,  $t = T$ :

$$\beta_T(i) = 1 \text{ for all } i^*$$

- Inductive case, compute for  $t = T-1, \dots, 1$ :

$$\beta_t(i) = \sum_{j=1}^N a_{ij} \cdot b_j(y_{t+1}) \cdot \beta_{t+1}(j)$$

- Exact mirror image of forward.



\*) this is different in J&M because of  $q_F$

# Putting it all together

- Compute estimated transition counts for all  $i, j, t$ :

$$\xi_t(i, j) = \frac{\xi'_t(i, j)}{\hat{P}(y)} = \frac{\alpha_t(i) \cdot a_{ij} \cdot b_j(y_{t+1}) \cdot \beta_{t+1}(j)}{\sum_q \alpha_T(q)}$$

- Compute overall estimated transition counts:

$$E(q_i \rightarrow q_j) = \sum_{t=1}^{T-1} \xi_t(i, j)$$

- Revised estimate of transition probabilities:

$$a_{ij} \approx \frac{E(q_i \rightarrow q_j)}{E(q_i \rightarrow \bullet)}$$



# The other parameters

- Revise initial and emission probabilities using estimated counts, in completely analogous way.
- Here's what it looks like for emission prob:

$$\gamma_t(j) = P(X_t = q_j \mid y) = \frac{\hat{P}(X_t = q_j, y)}{\hat{P}(y)} = \frac{\alpha_t(j) \cdot \beta_t(j)}{\hat{P}(y)}$$

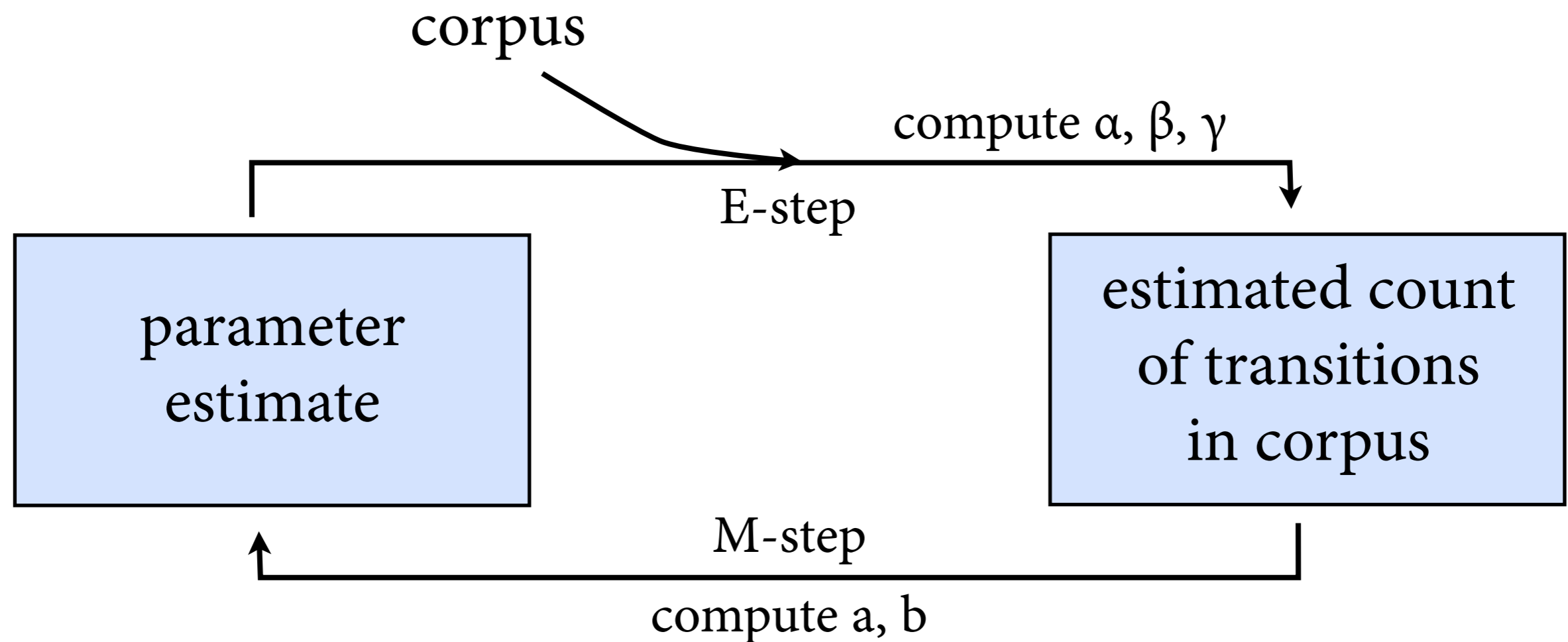
$$b_j(o) \approx \left( \sum_{\substack{t=1 \\ y_t=o}}^T \gamma_t(j) \right) / \sum_{t=1}^T \gamma_t(j)$$

estimated count of  
o emitted in state  $q_j$

estimated count of  
state  $q_j$

# Forward-Backward Algorithm

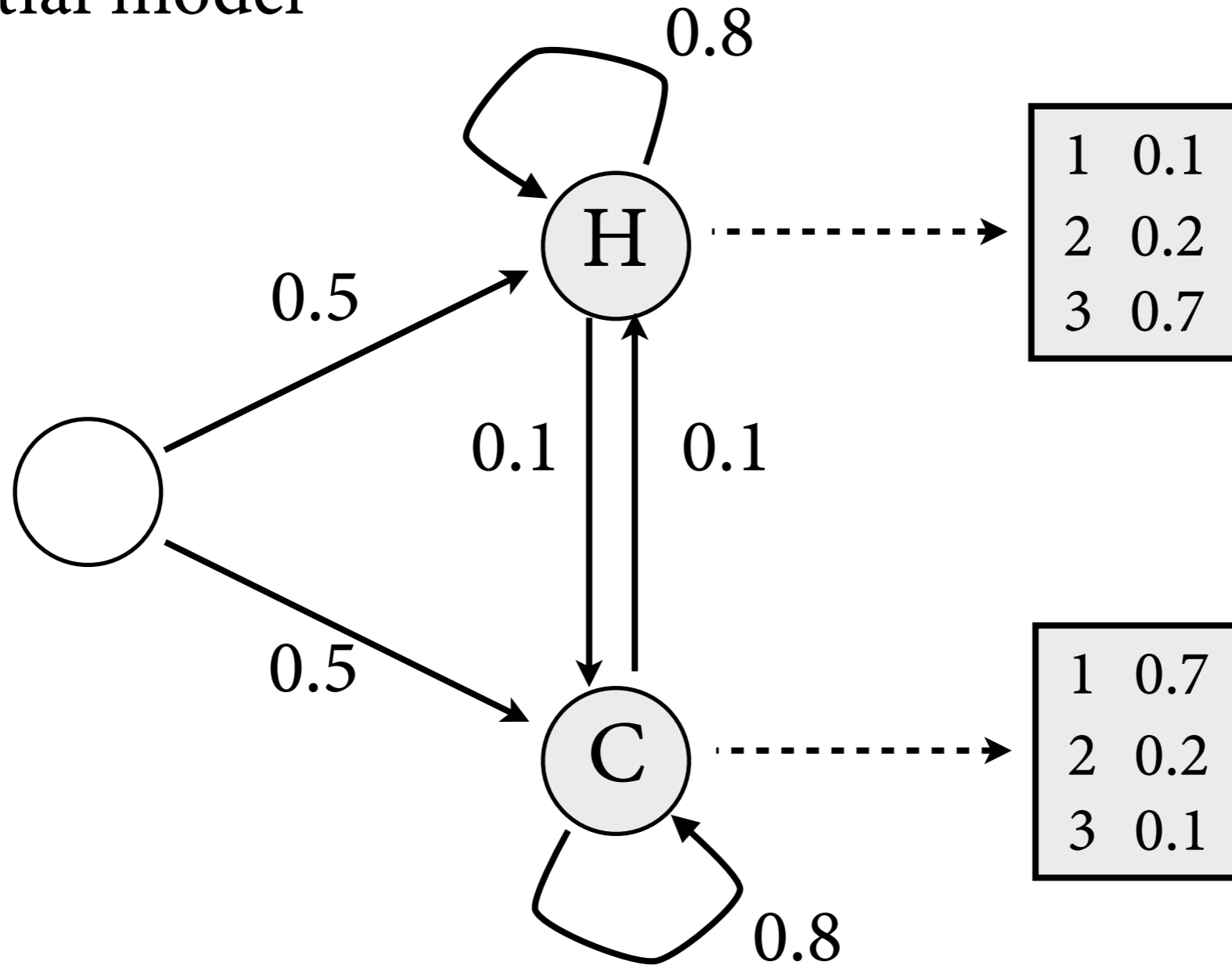
Initialization: start with some estimation of parameters.



Continue computation until parameters don't change much.

# Example

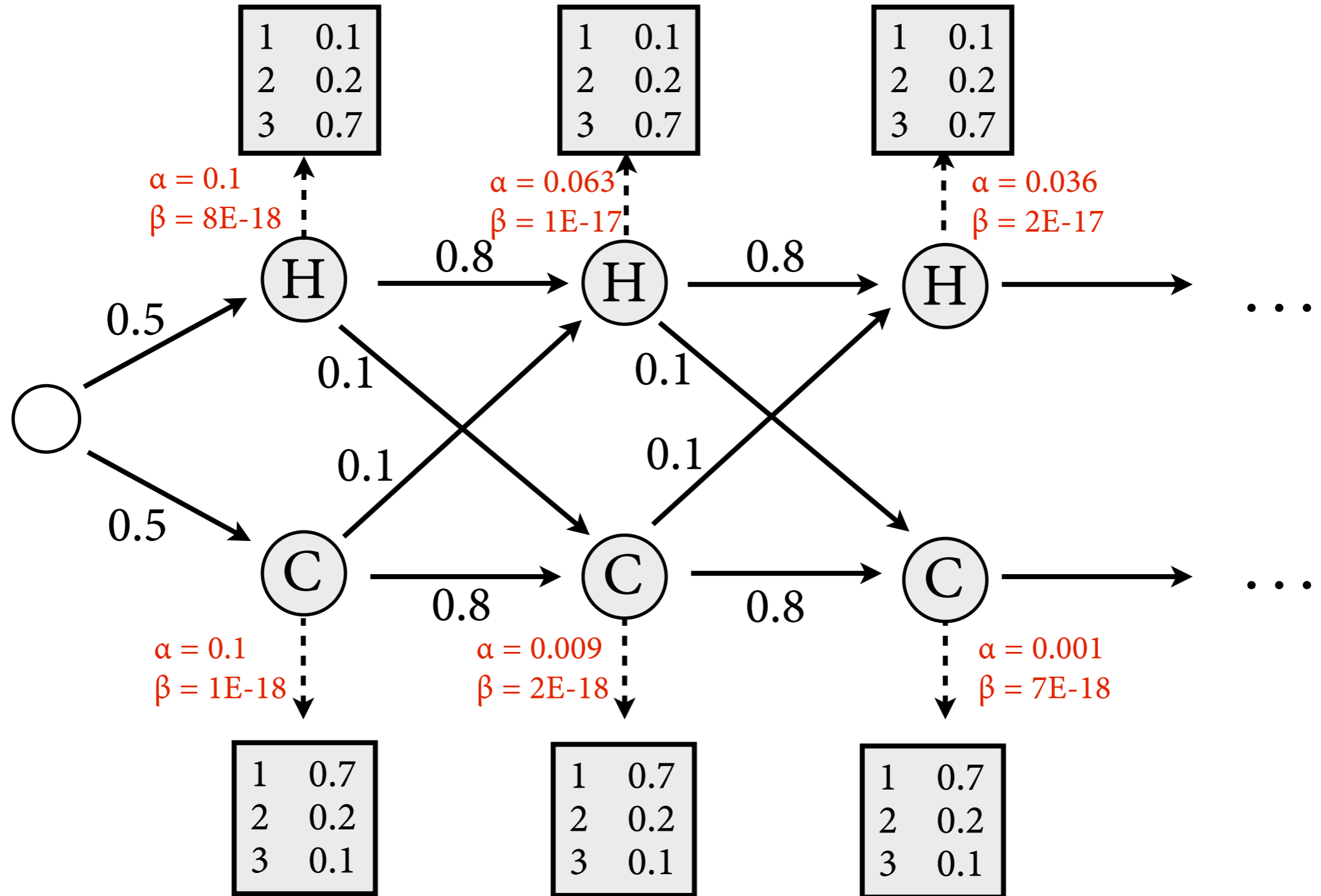
Step 1: initial model



Probs do not sum to one because Eisner's HMMs have  $q_F$ .

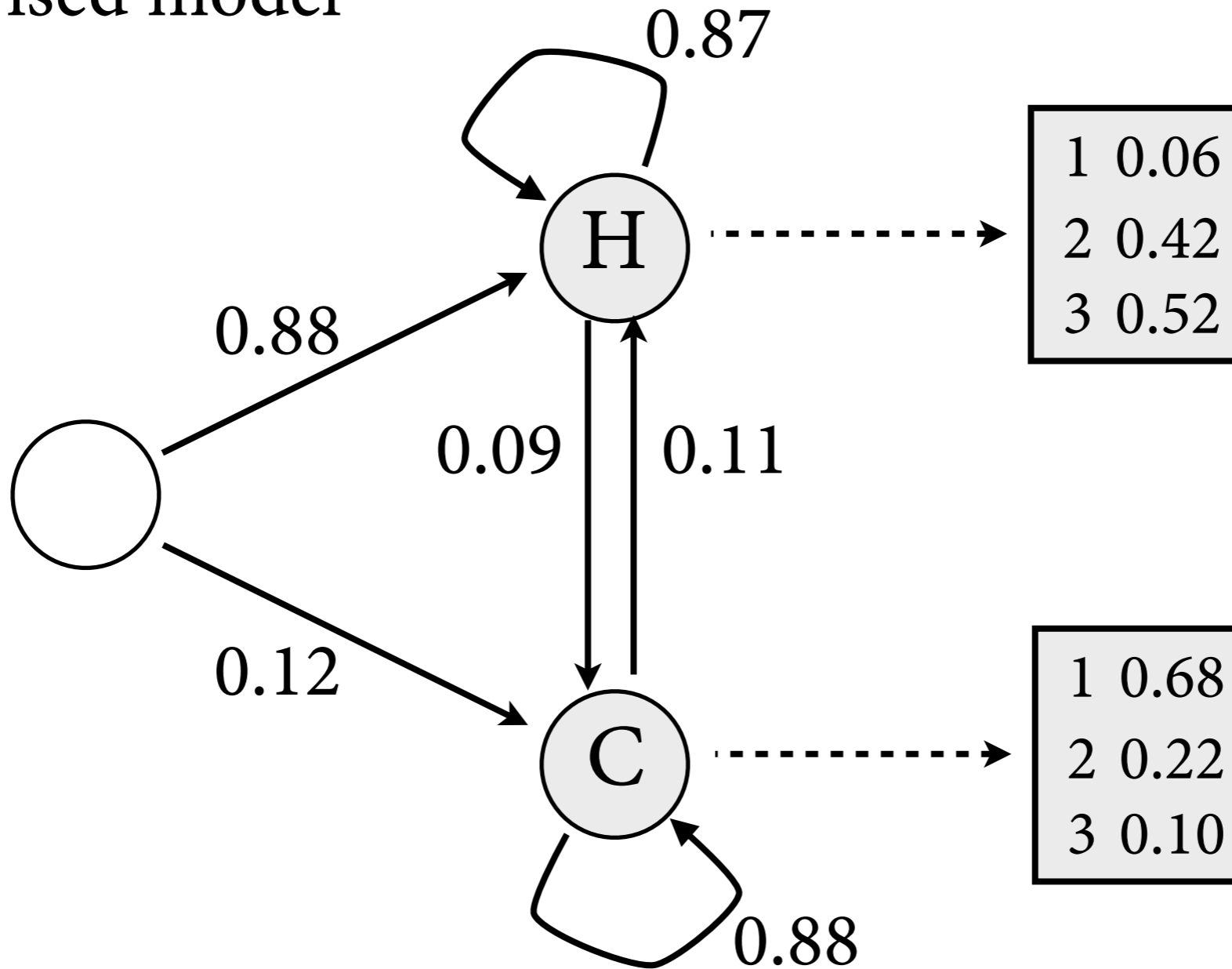
2, 3, 3, 2, 3, 2, 3, 2, 2, 3, 1, 3, 3, 1, 1, 1, 2, 1, 1, 1, 3, 1, 2, 1, 1, 1, 2, 3, 3, 2, 3, 2, 2

# E-Step



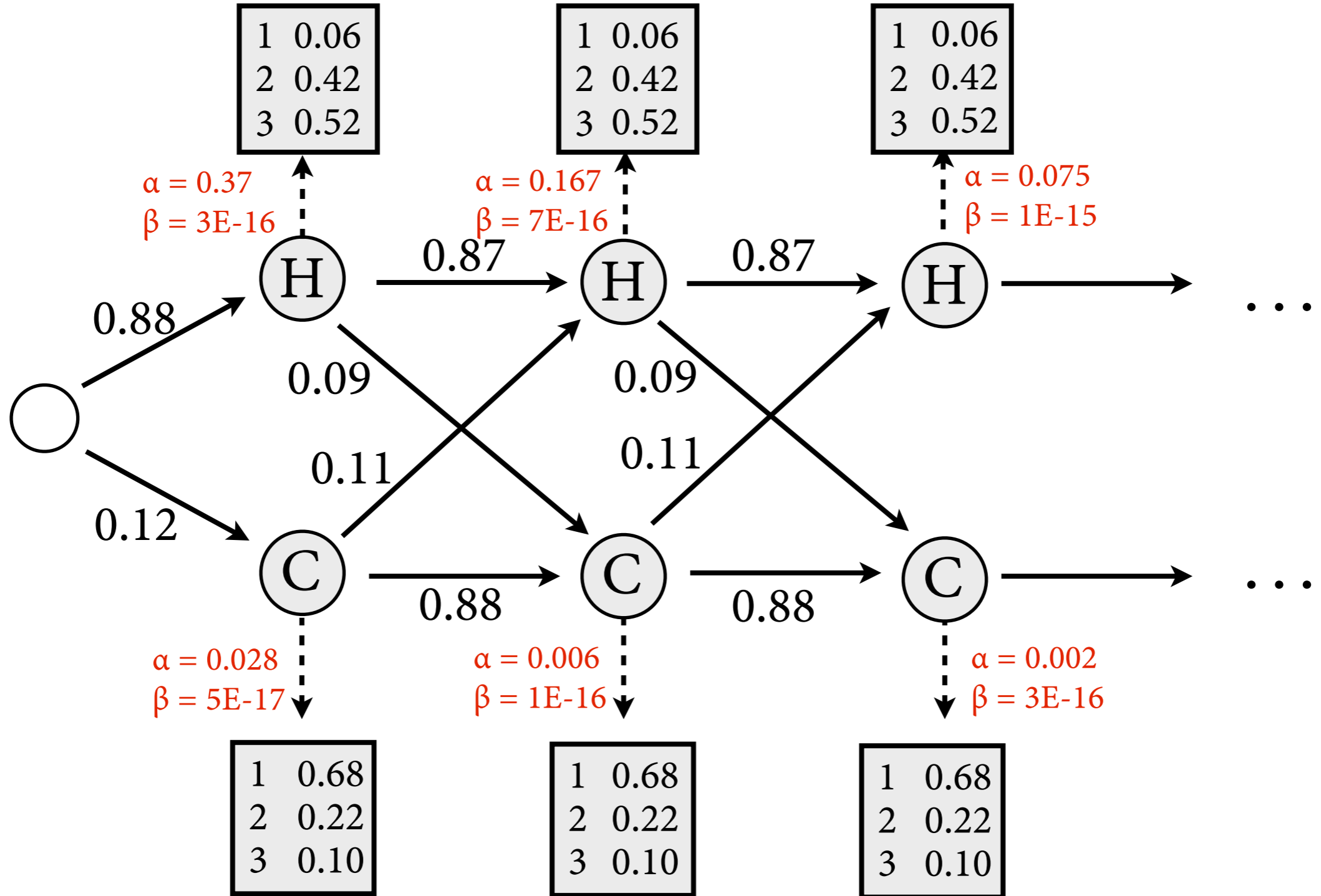
# M-Step

Step 2: revised model



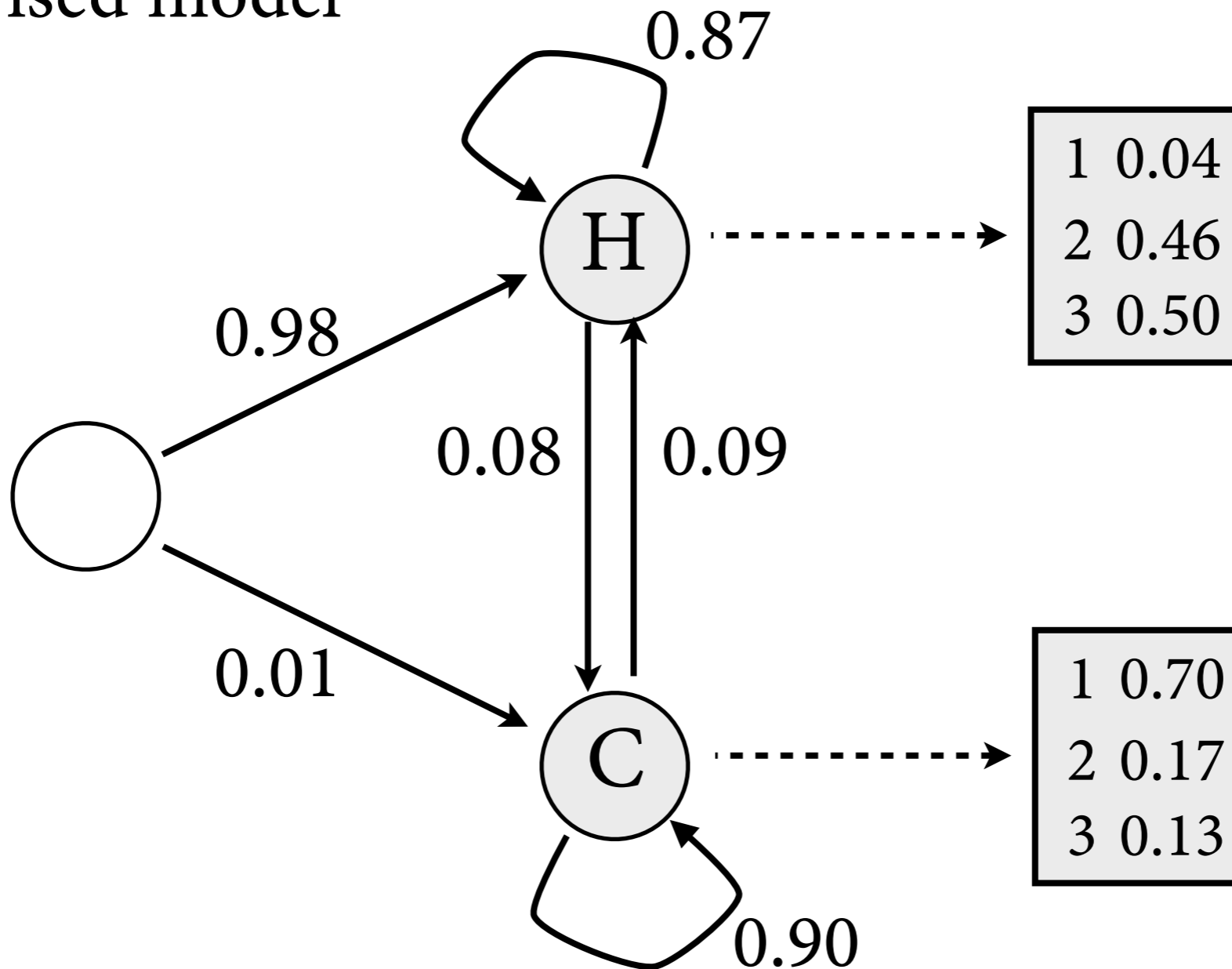
2, 3, 3, 2, 3, 2, 3, 2, 2, 3, 1, 3, 3, 1, 1, 1, 2, 1, 1, 1, 3, 1, 2, 1, 1, 1, 2, 3, 3, 2, 3, 2, 2

# E-Step



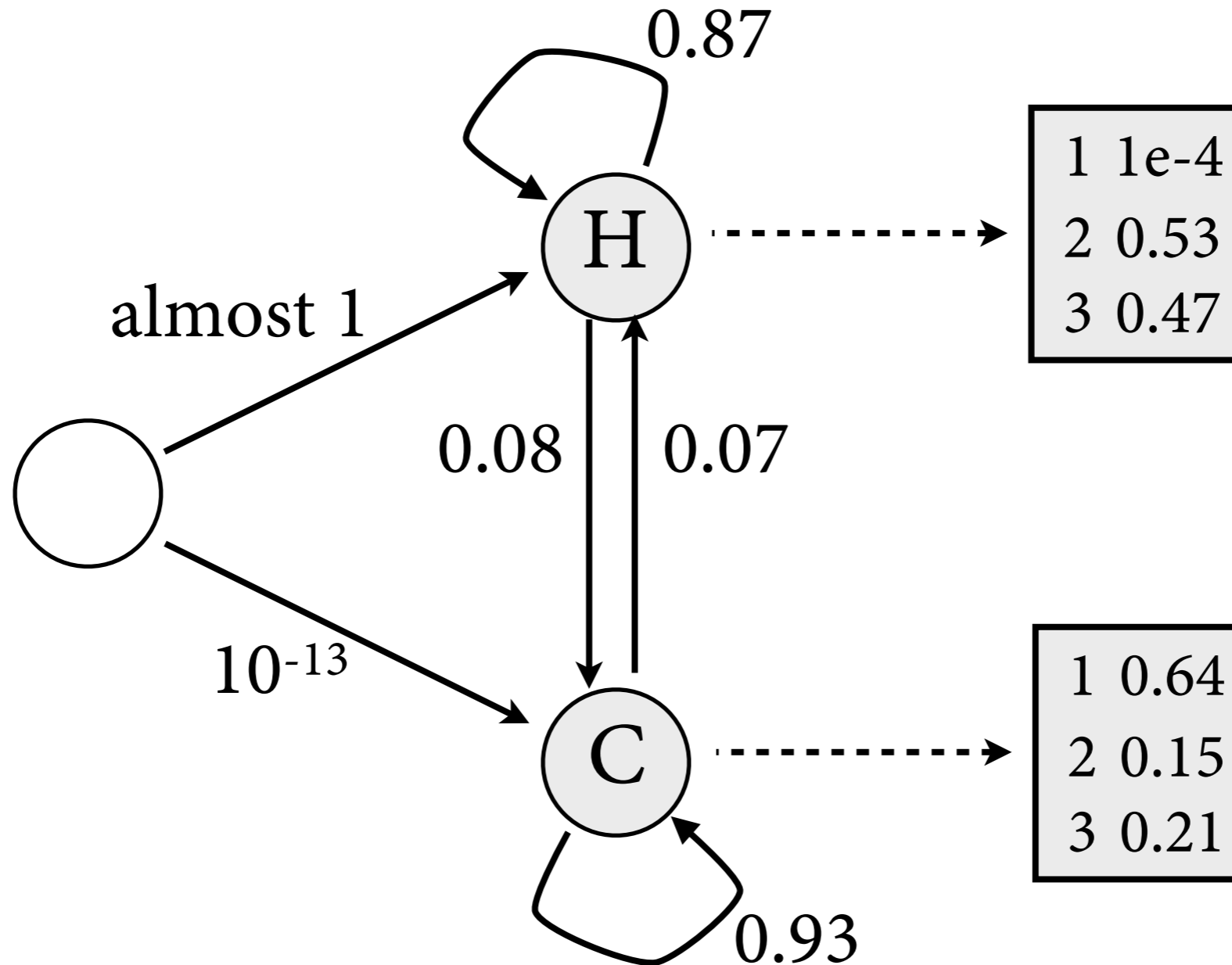
# M-Step

Step 3: revised model



2, 3, 3, 2, 3, 2, 3, 2, 2, 3, 1, 3, 3, 1, 1, 1, 2, 1, 1, 1, 3, 1, 2, 1, 1, 1, 2, 3, 3, 2, 3, 2, 2

# Result after 10 iterations



2, 3, 3, 2, 3, 2, 3, 2, 2, 3, 1, 3, 3, 1, 1, 1, 2, 1, 1, 1, 3, 1, 2, 1, 1, 1, 2, 3, 3, 2, 3, 2, 2

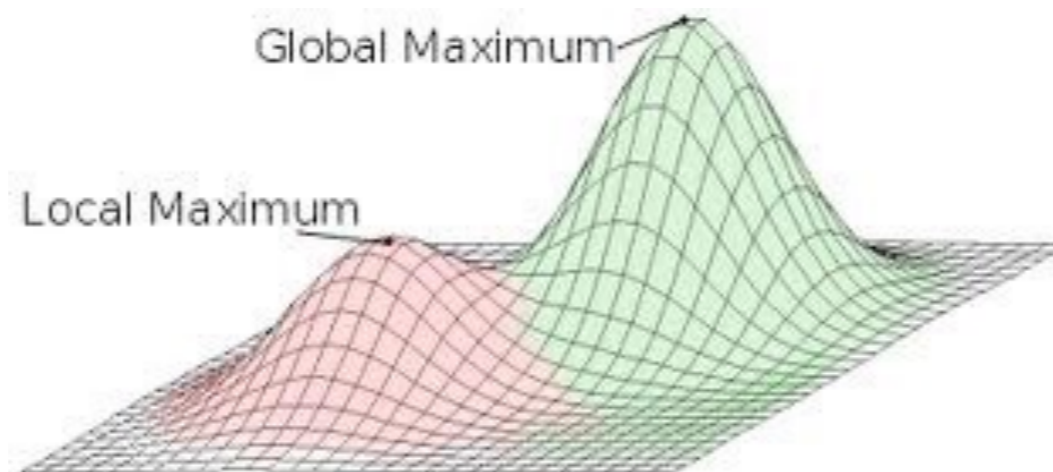


# Some remarks

- Forward-backward algorithm also called *Baum-Welch Algorithm* after inventors.
- Special case of the *expectation maximization* algorithm:
  - ▶ E-Step: Compute expected values of relevant counts based on current parameter estimate.
  - ▶ M-Step: Adjust model based on estimated counts.
- Runtime of each iteration is  $O(N^2 T)$ .  
Most of the time goes into E-step.

# Some remarks

- EM algorithm is guaranteed to improve likelihood of corpus in each iteration.
- However, can run into *local maxima*: would have to go through worse model to find globally best one.
- Extremely sensitive to initial parameter estimate. Only useful in practice if HMM structure very strongly constrained (e.g. speech recognition).



# Conclusion

- Evaluate tagger on *accuracy* on *unseen* data.
- Training algorithms for HMM estimation:
  - ▶ *supervised* training from annotated data:  
maximum likelihood
  - ▶ *unsupervised* training from unannotated data:  
forward-backward