

# Hidden Markov Models

Computational Linguistics

Alexander Koller

5 November 2019

# Let's play a game

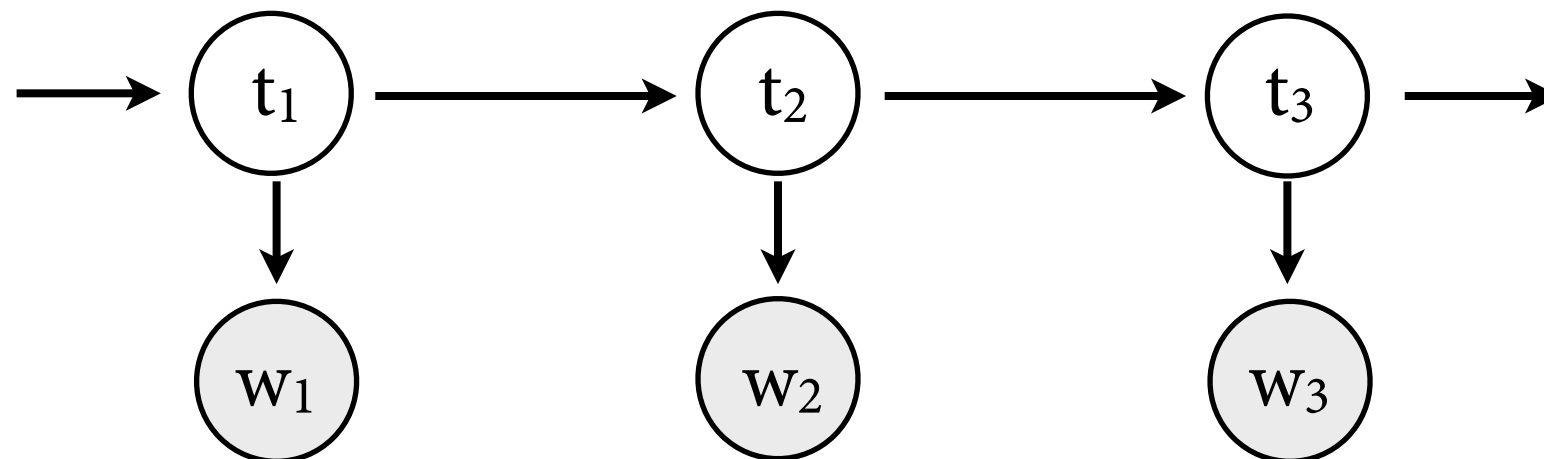
- I will write a sequence of part-of-speech tags and of words on the board.
- You take turns in giving me POS tags and words, and I will write them down.

# Penn Treebank POS tags

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &amp;</i>
CD	Cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential ‘there’	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	“	Left quote	<i>( ‘ or “)</i>
POS	Possessive ending	<i>’s</i>	”	Right quote	<i>( ’ or ”)</i>
PP	Personal pronoun	<i>I, you, he</i>	(	Left parenthesis	<i>( [ , { , &lt;)</i>
PP\$	Possessive pronoun	<i>your, one’s</i>	)	Right parenthesis	<i>( [ , }, &gt;)</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>( . ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>( : ; ... – -)</i>
RP	Particle	<i>up, off</i>			

# Hidden Markov Models

- Last week's generative story: generate words at random from n-gram model  $P(w_n \mid w_1, \dots, w_{n-1})$ .
- Replace with new generative story:
  - ▶ Language is generated by a two-step process.
  - ▶ First, generate sequence of hidden POS tags  $t_1, \dots, t_T$  tag by tag, left to right from bigram model  $P(t_i \mid t_{i-1})$ .
  - ▶ Independently, generate an observable word  $w_i$  from each  $t_i$ , at random from model  $P(w_i \mid t_i)$ .



# 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	
The	representative	put	chairs	on	the	table.	$p_1$

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

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

# Question 2: Tagging (aka Decoding)

- Given an HMM and an observed string  $w_1, \dots, w_T$ , what is the most likely sequence of hidden tags  $t_1, \dots, t_T$ ?
- We can compute  $\arg \max_{t_1, \dots, t_T} P(t_1, \dots, t_T \mid w_1, \dots, w_T)$  efficiently with the *Viterbi algorithm*.

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

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

# Question 2: Tagging (aka Decoding)

- Given an HMM and an observed string  $w_1, \dots, w_T$ , what is the most likely sequence of hidden tags  $t_1, \dots, t_T$ ?
- We can compute  $\arg \max_{t_1, \dots, t_T} P(t_1, \dots, t_T \mid w_1, \dots, w_T)$  efficiently with the *Viterbi algorithm*.

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

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

# 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.
- Can do it with maximum likelihood estimation plus smoothing.

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.						

→ next time



# 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.
- Do it efficiently with the *forward-backward algorithm* (an instance of *Expectation Maximization*).

The representative put chairs on the table.

Secretariat is expected to race tomorrow.

# Hidden Markov Models



- A *Hidden Markov Model* is 5-tuple consisting of

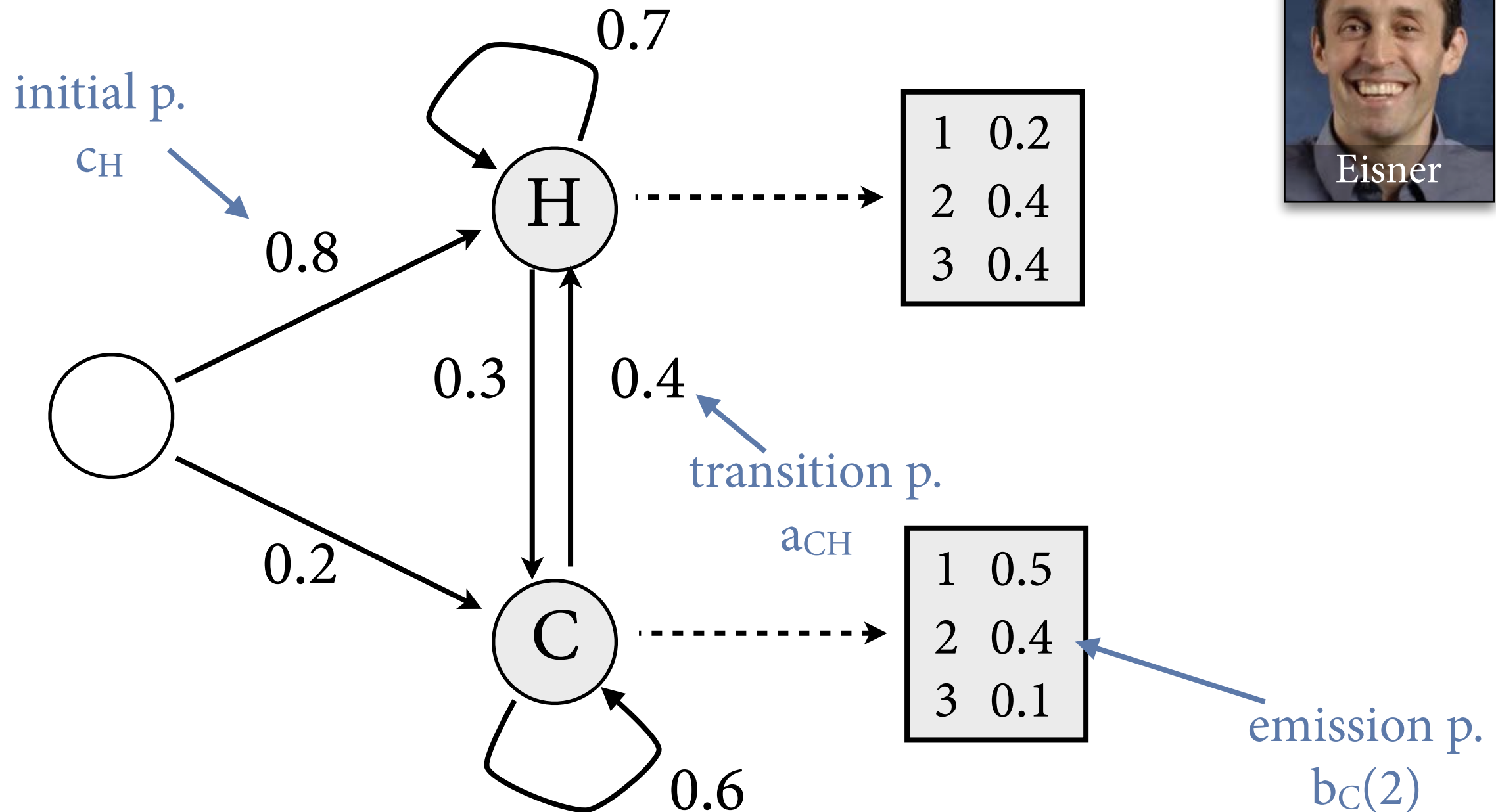
- ▶ finite set  $Q = \{q_1, \dots, q_N\}$  of *states* (= POS tags)
- ▶ finite set  $O$  of possible *observations* (= words)
- ▶ *transition probabilities*  $a_{ij} = P(X_{t+1} = q_j \mid X_t = q_i)$
- ▶ *initial probabilities*  $c_i = P(X_1 = q_i)$
- ▶ *emission probabilities*  $b_i(o) = P(Y_t = o \mid X_t = q_i)$

$$\sum_{j=1}^N a_{ij} = 1$$
$$\sum_{i=1}^N c_i = 1$$
$$\sum_{o \in O} b_i(o) = 1$$

- The HMM describes two coupled random processes:

- ▶ event  $X_t = q_i$ : At time  $t$ , HMM is in state  $q_i$ .
- ▶ event  $Y_t = o$ : At time  $t$ , HMM emits observation  $o$ .

# Example: Eisner's Ice Cream



States represent weather on a given day: Hot, Cold  
Outputs represent number of ice creams Jason eats that day

# HMMs as joint models of $x, y$

- Coupled random processes of HMM directly give us model for *joint* probability  $P(x, y)$  where
  - ▶  $y = y_1 \dots y_T$  sequence of observations,  $y_t \in O$
  - ▶  $x = x_1 \dots x_T$  sequence of hidden states,  $x_t \in Q$
- Defined as follows:

$$\begin{aligned} P(x, y) &= P(x) \cdot P(y \mid x) \\ &= P(X_1 = x_1) \cdot \prod_{t=2}^T P(X_t = x_t \mid X_1 = x_1, \dots, X_{t-1} = x_{t-1}) \\ &\quad \cdot \prod_{t=1}^T P(Y_t = y_t \mid Y_1 = y_1, \dots, Y_{t-1} = y_{t-1}, x) \\ &= P(X_1 = x_1) \cdot \prod_{t=2}^T P(X_t = x_t \mid X_{t-1} = x_{t-1}) \cdot \prod_{t=1}^T P(Y_t = y_t \mid X_t = x_t) \\ &= c_{x_1} \cdot \prod_{t=2}^T a_{x_{t-1}x_t} \cdot \prod_{t=1}^T b_{x_t}(y_t) \end{aligned}$$

# Question 2: Tagging

- Given observations  $y_1, \dots, y_T$  (# ice creams), what is the most probable sequence  $x_1, \dots, x_T$  of hidden states (temperatures)?

- Maximum probability:

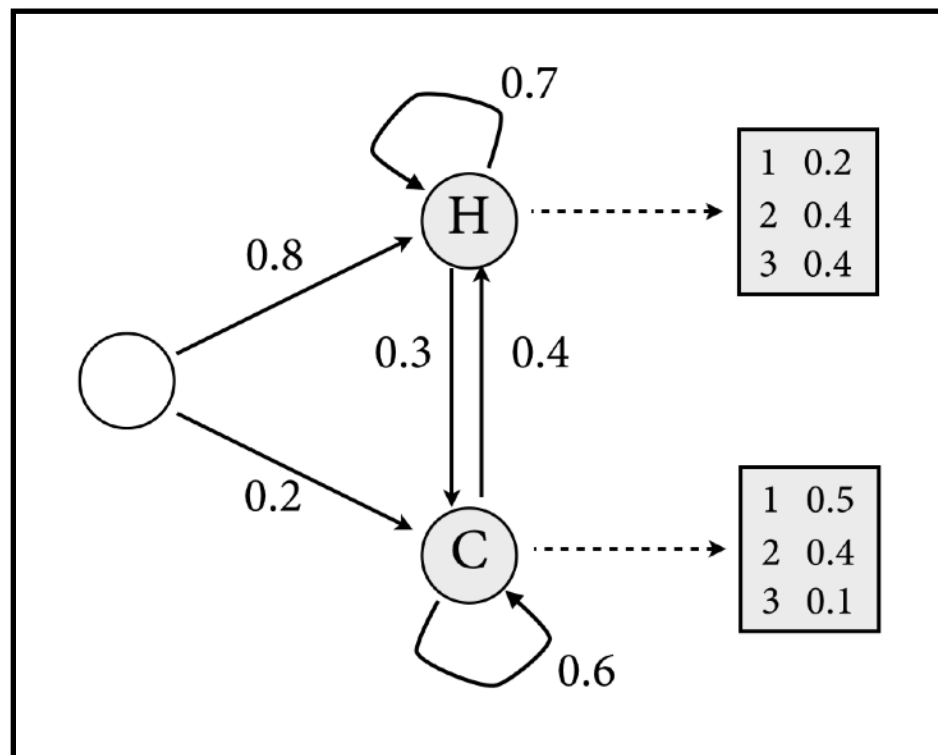
$$\max_{x_1, \dots, x_T} P(x_1, \dots, x_T \mid y_1, \dots, y_T)$$

- We are primarily interested in  $\arg \max$ :

$$\begin{aligned} & \arg \max_{x_1, \dots, x_T} P(x_1, \dots, x_T \mid y_1, \dots, y_T) \\ &= \arg \max_{x_1, \dots, x_T} \frac{P(x_1, \dots, x_T, y_1, \dots, y_T)}{P(y_1, \dots, y_T)} \\ &= \arg \max_{x_1, \dots, x_T} P(x_1, \dots, x_T, y_1, \dots, y_T) \end{aligned}$$

# Naive approach

- Let's say Jason ate 3, 1, 3 ice creams. What was the most likely weather on these three days?
- Compute  $\max P(x_1, 3, x_2, 1, x_3, 3)$  by maximizing over all possible state sequences.

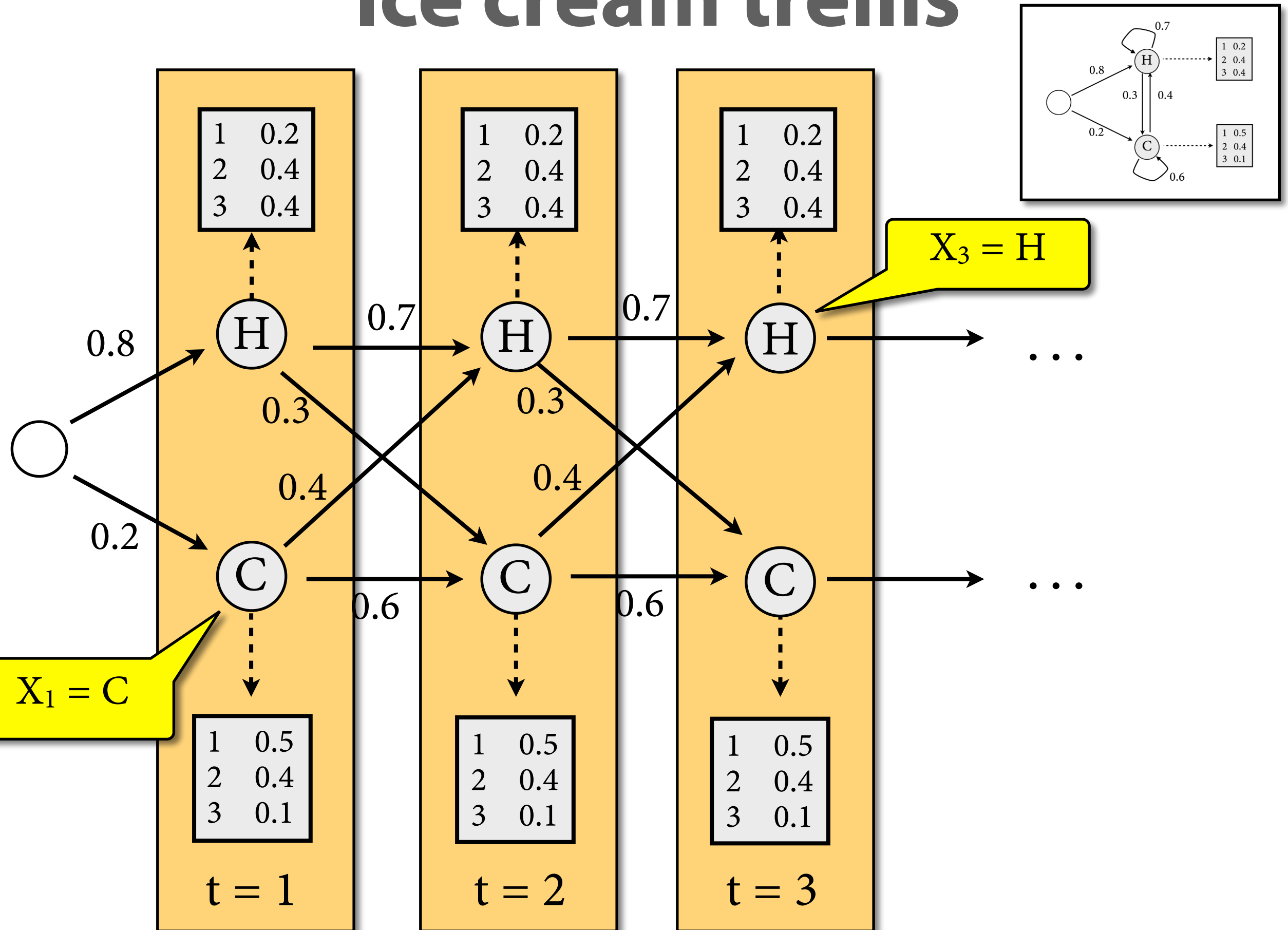


$$\begin{array}{l} P(x_1, 3, x_2, 1, x_3, 3) \text{ is max of} \\ P(H, 3, H, 1, H, 3) \quad 0.012 \\ P(H, 3, H, 1, C, 3) \quad 0.001 \\ P(H, 3, C, 1, H, 3) \quad 0.008 \\ \dots \quad \dots \\ \underline{P(C, 3, C, 1, C, 3) \quad 0.0004} \\ = 0.012 \end{array}$$

# Too expensive

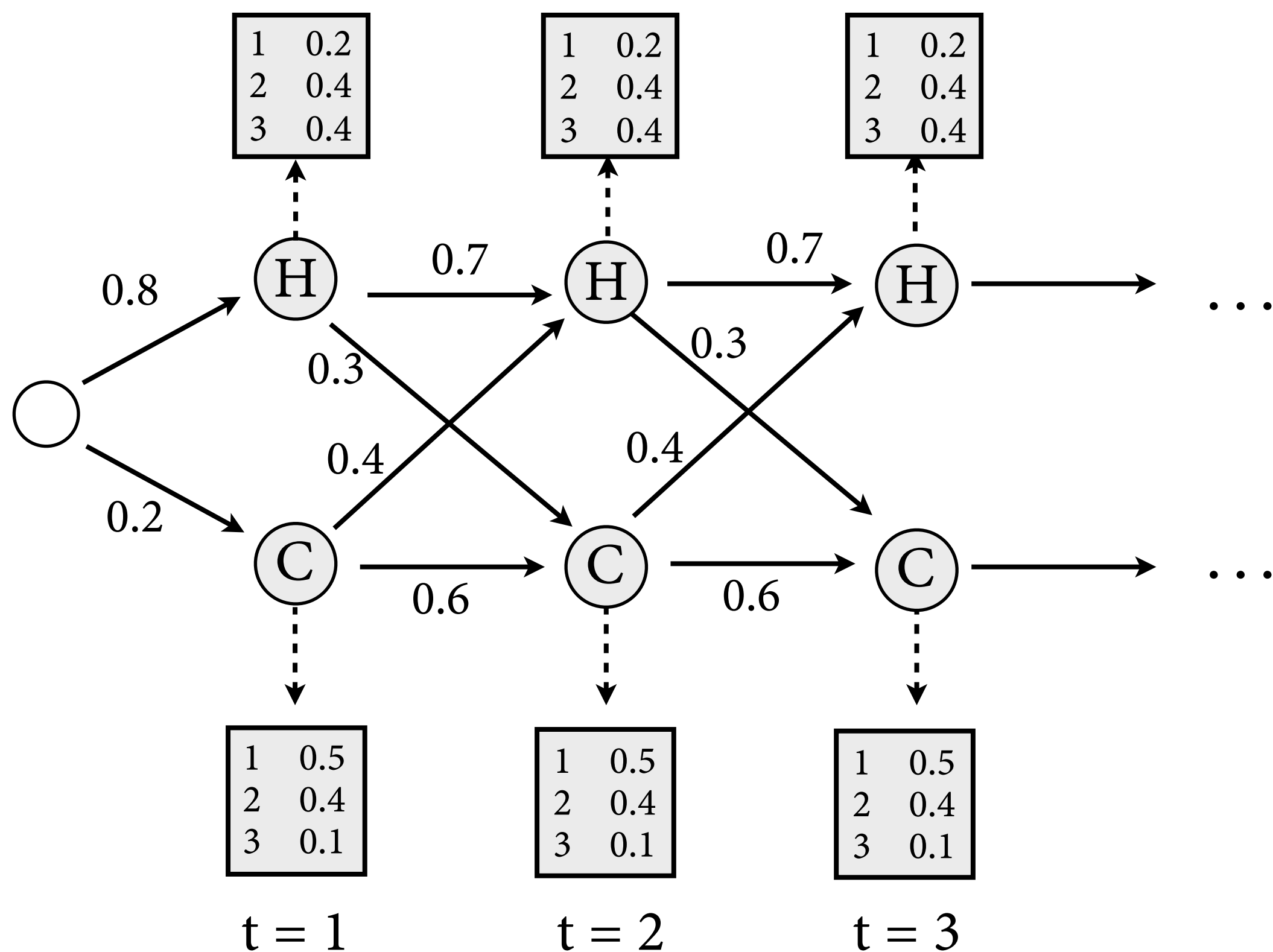
- Naive approach maximizes over exponentially large set of products. This is too slow for practical use.
- Visualize this in *trellis*: unfolding of HMM
  - ▶ one column for each time point  $t$ , represents  $X_t$
  - ▶ each column contains a copy of each state of HMM
  - ▶ edges from  $t$  to  $t+1$  = transitions of HMM
- Each path through trellis represents one state sequence.
  - ▶ So computation of  $\max P(x,y) = \max$  over all paths.

# Ice cream trellis



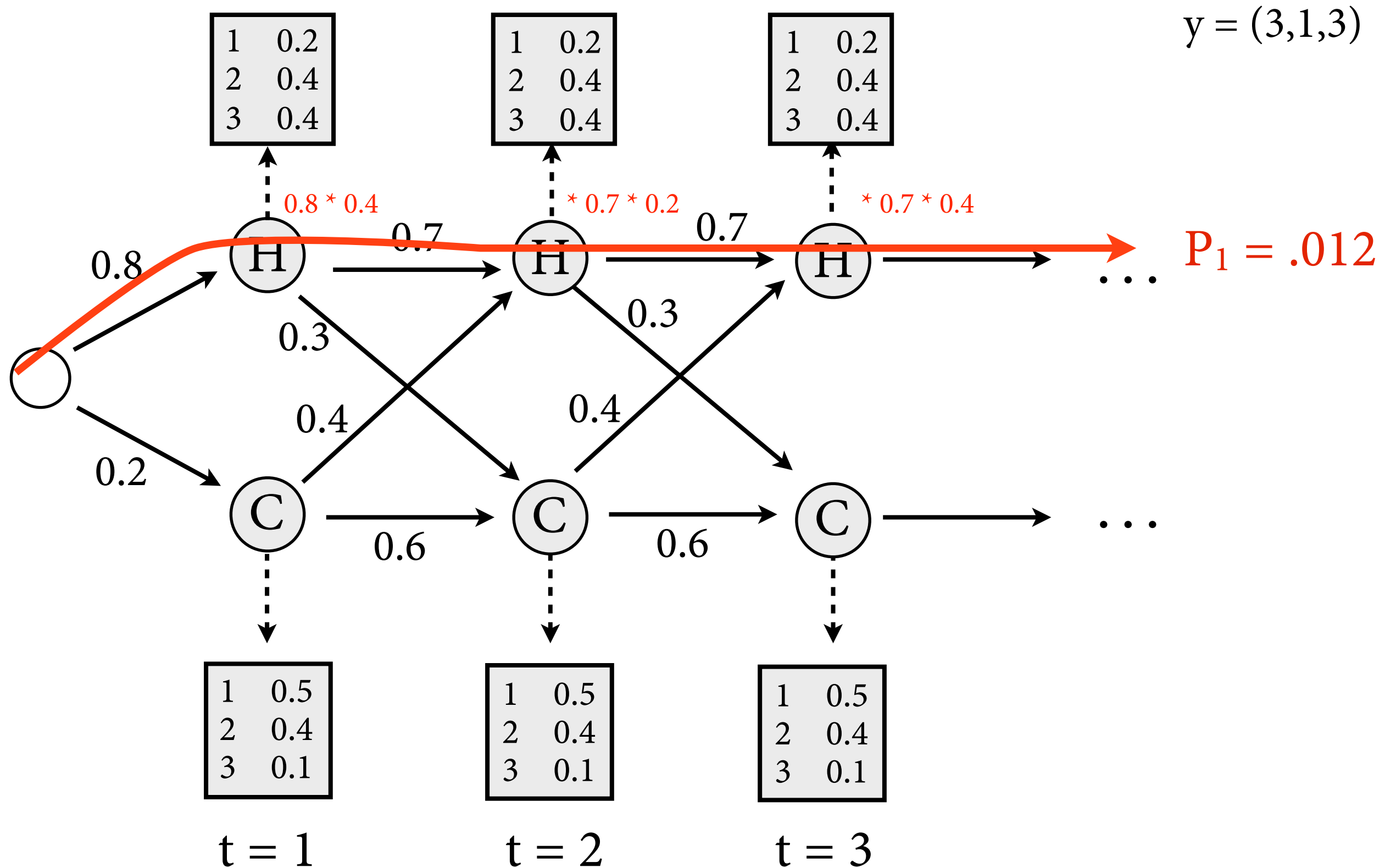


# Ice cream trellis

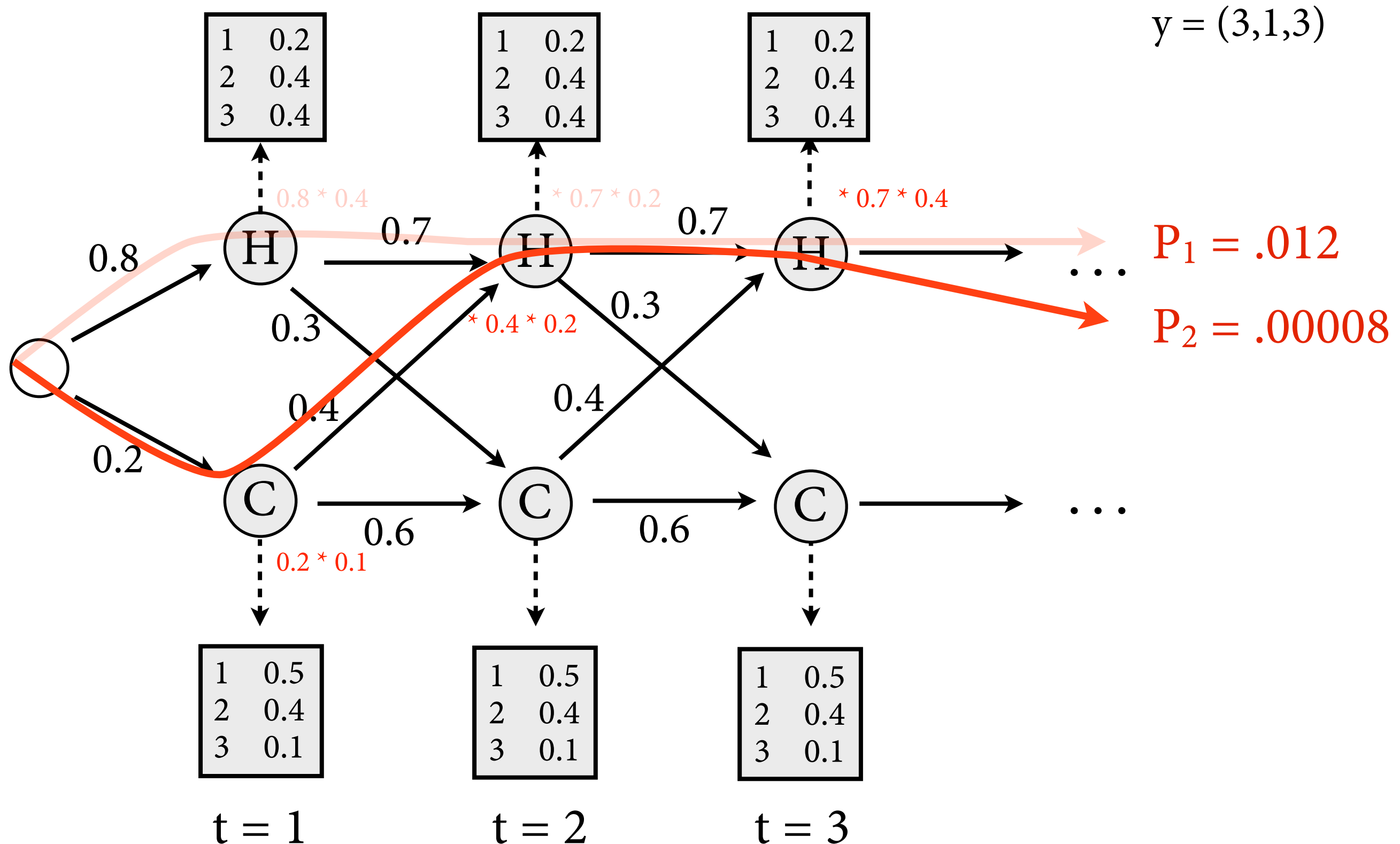


$$y = (3, 1, 3)$$

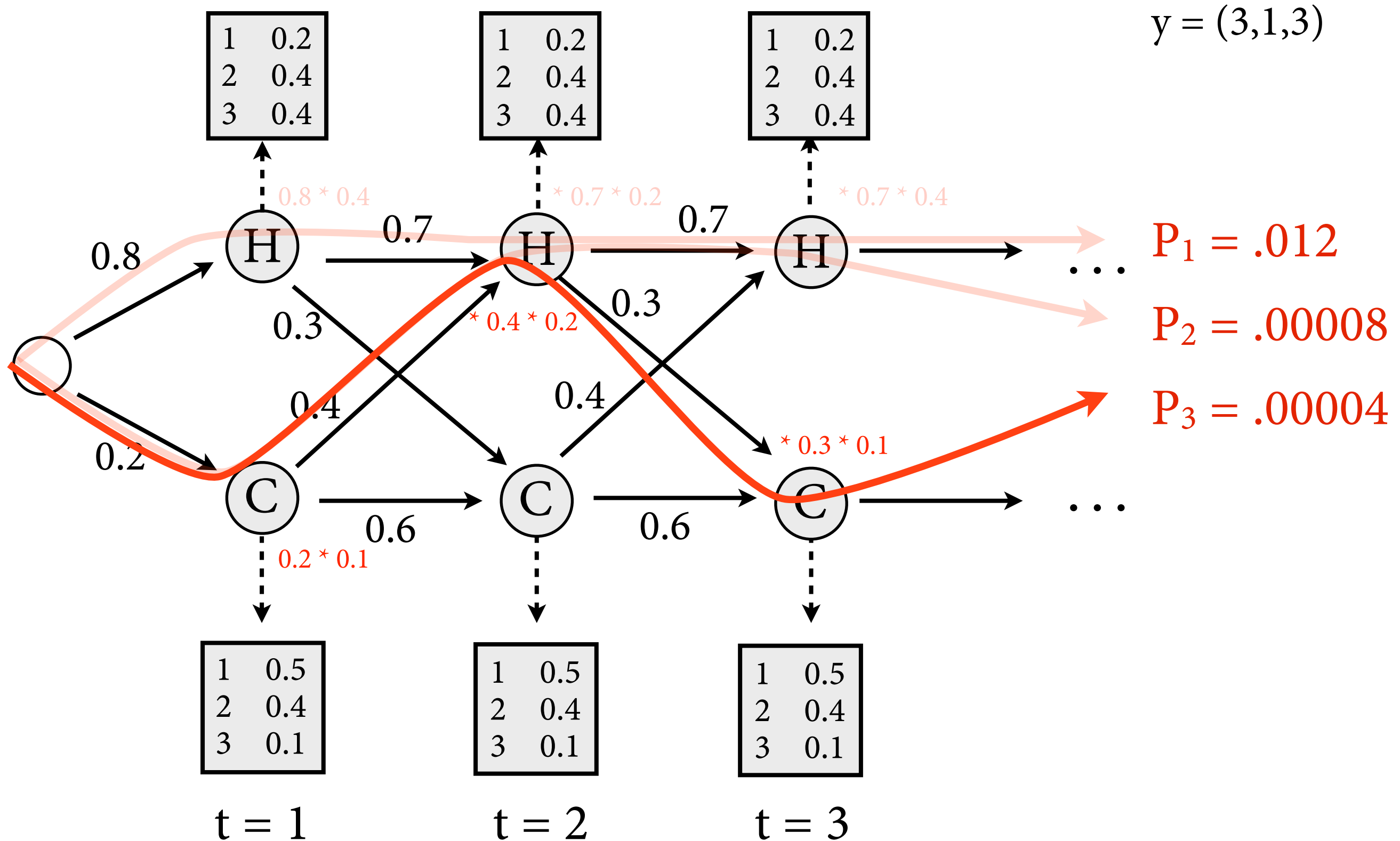
# Ice cream trellis



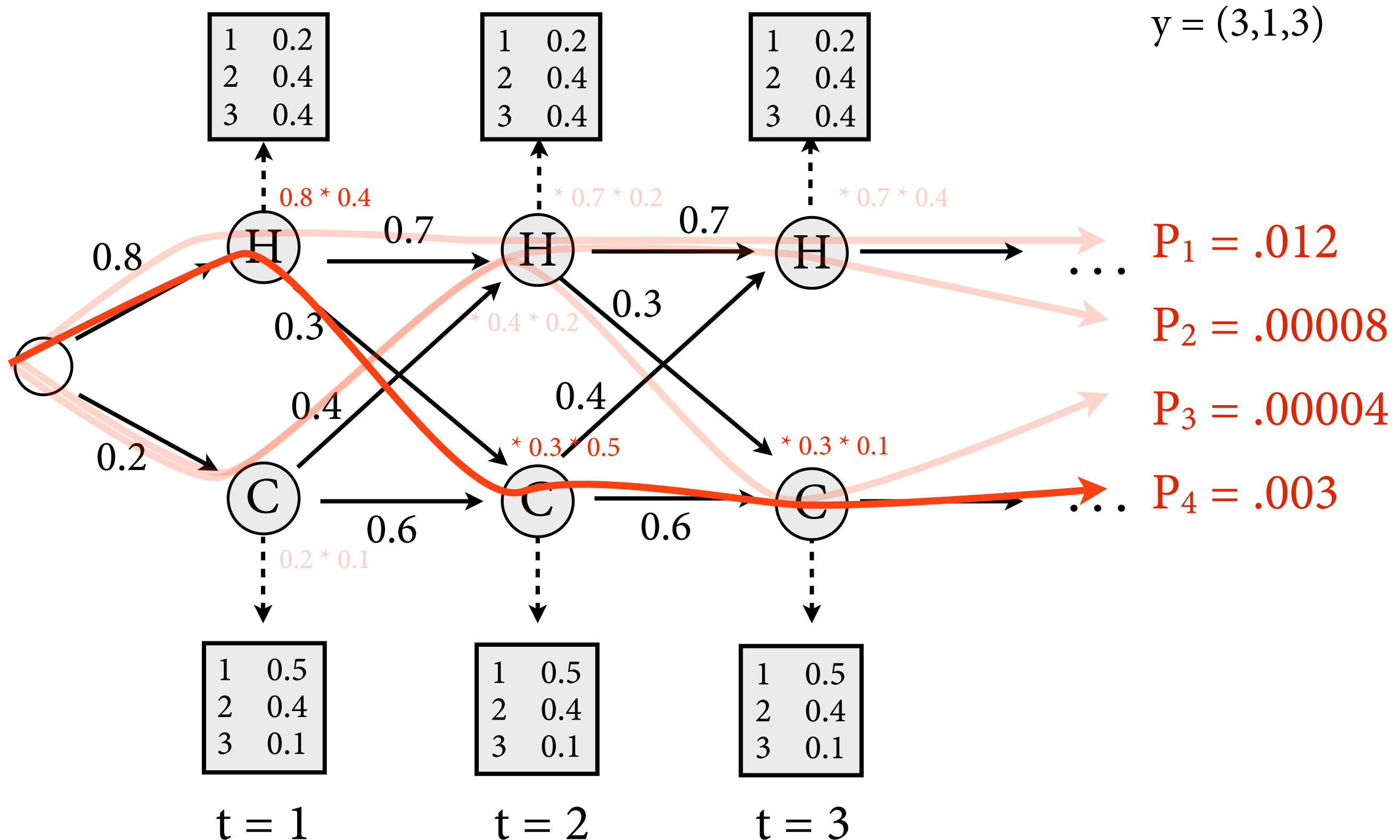
# Ice cream trellis



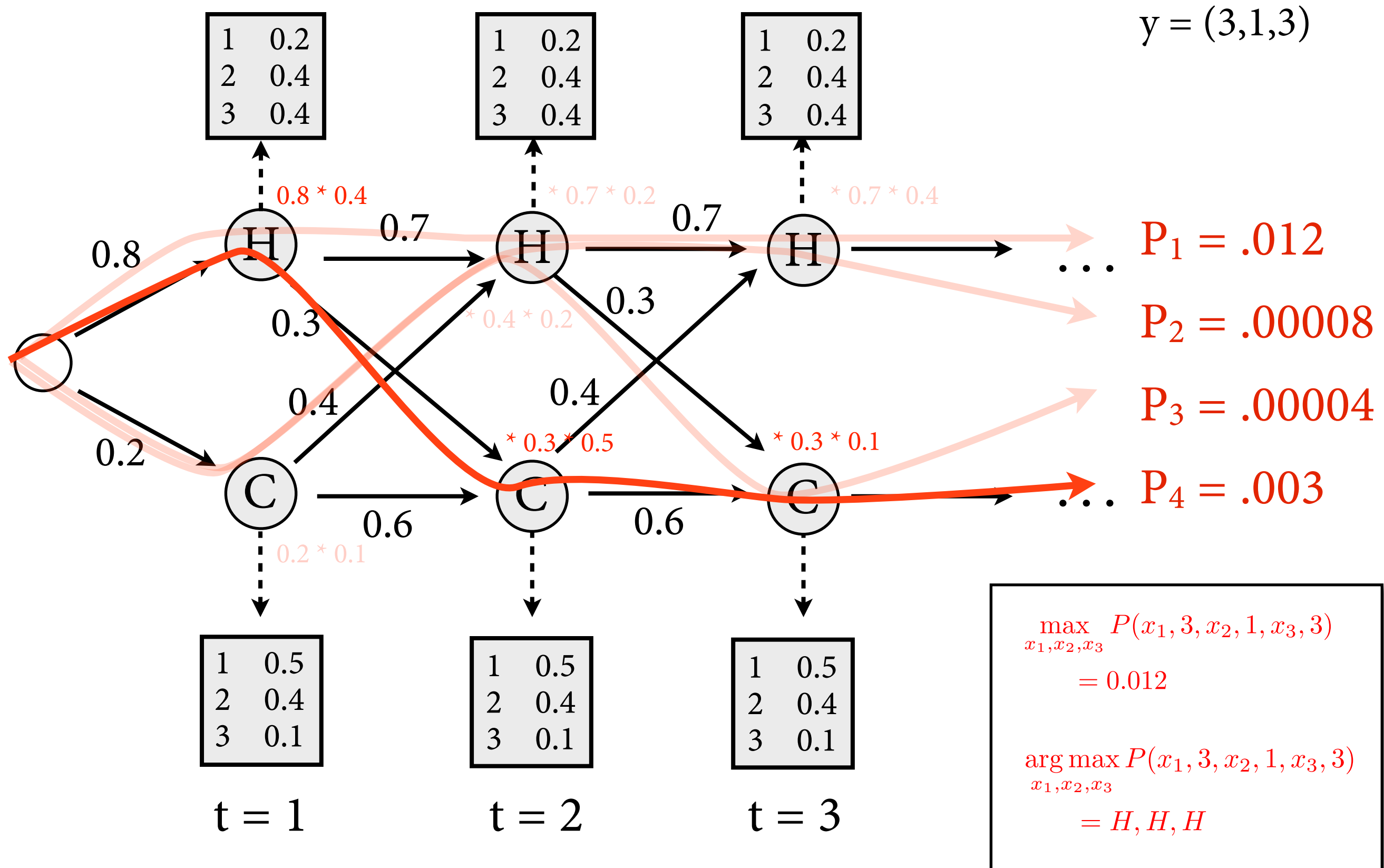
# Ice cream trellis



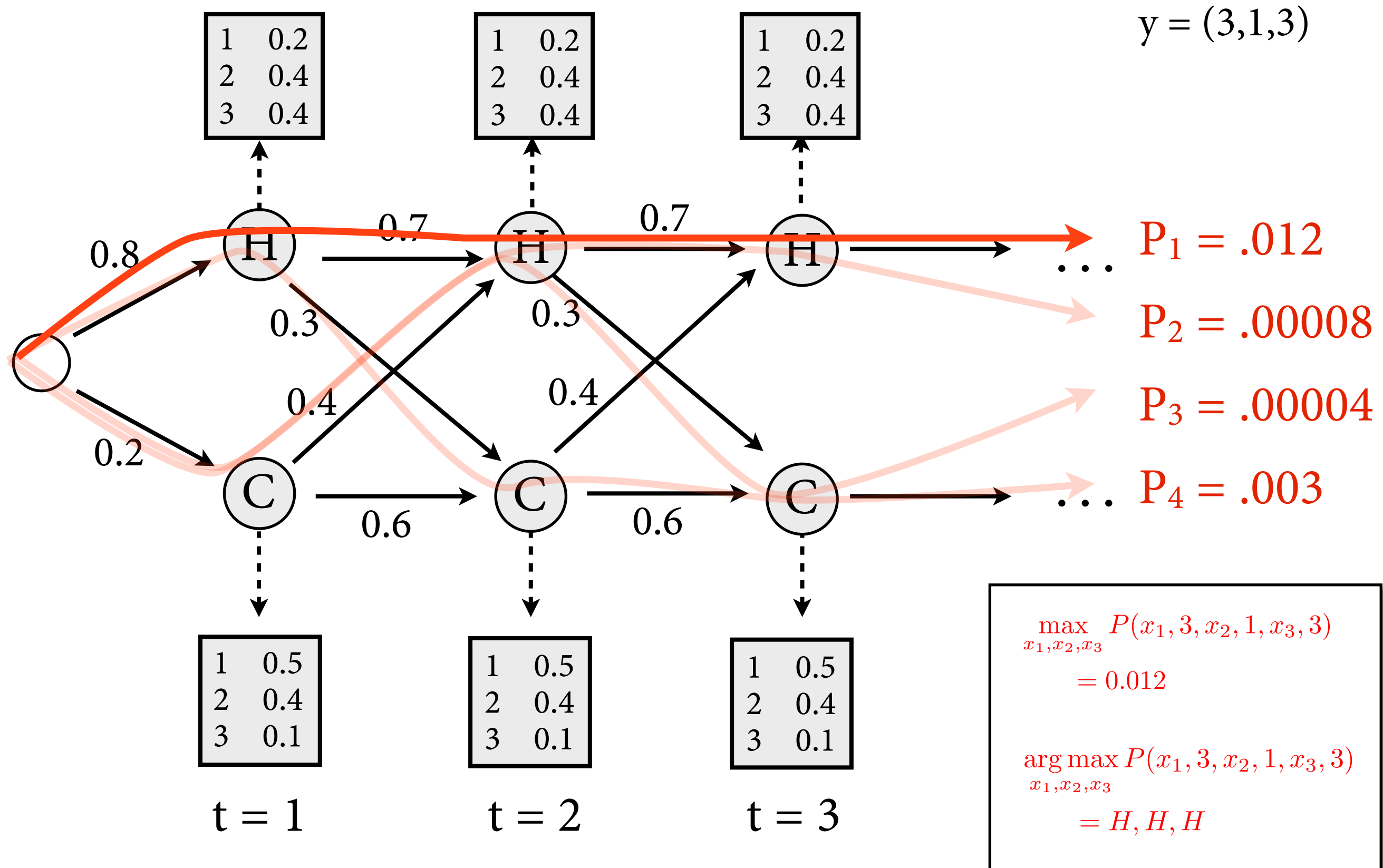
# Ice cream trellis



# Ice cream trellis



# Ice cream trellis



# Best paths to trellis nodes

- We are trying to compute the best path through the trellis:  $\arg \max_{x_1, \dots, x_T} P(X_1 = x_1, \dots, X_T = x_T, Y_1 = y_1, \dots, Y_T = y_T)$

- More general concept: *best path* through the trellis that ends in a certain node  $X_t = q$ :

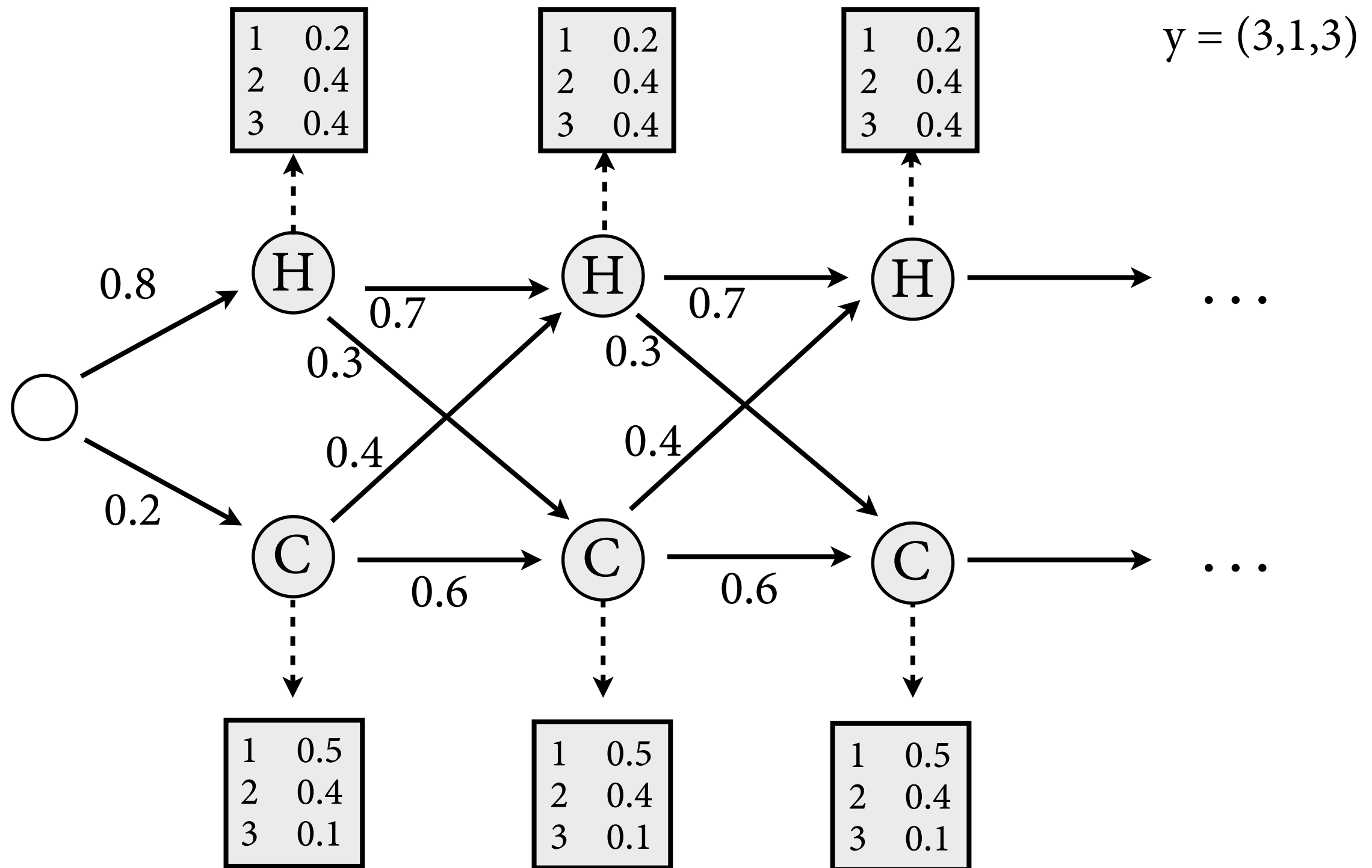
$$\arg \max_{x_1, \dots, x_{t-1}} P(X_1 = x_1, \dots, X_{t-1} = x_{t-1}, X_t = q, Y_1 = y_1, \dots, Y_t = y_t)$$

- Probability of that best path (abbreviate as  $V_t(q)$ ):

$$V_t(q) = \max_{x_1, \dots, x_T} P(X_1 = x_1, \dots, X_T = x_T, Y_1 = y_1, \dots, Y_T = y_T)$$



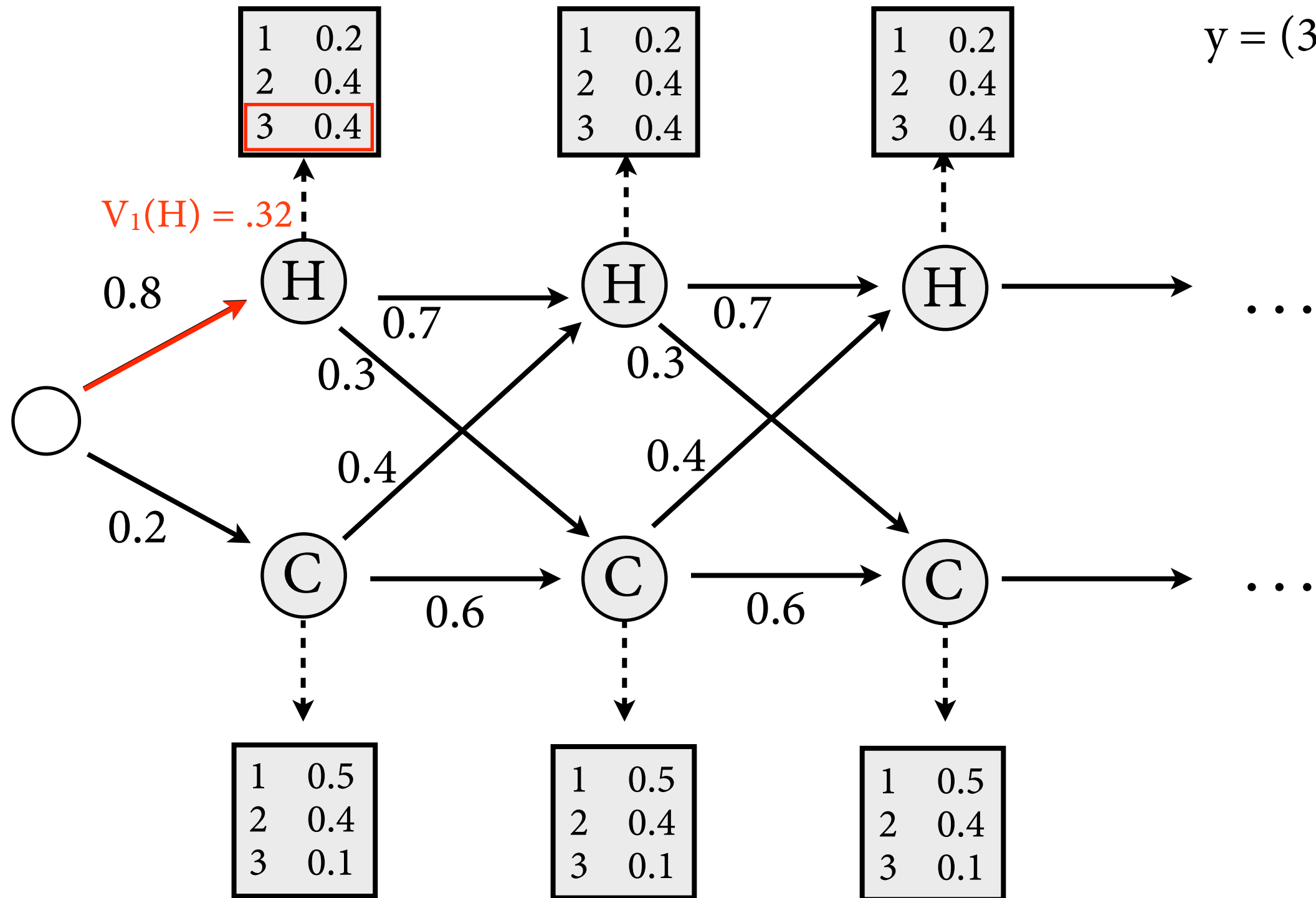
# Computing best path probs



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

# Computing best path probs

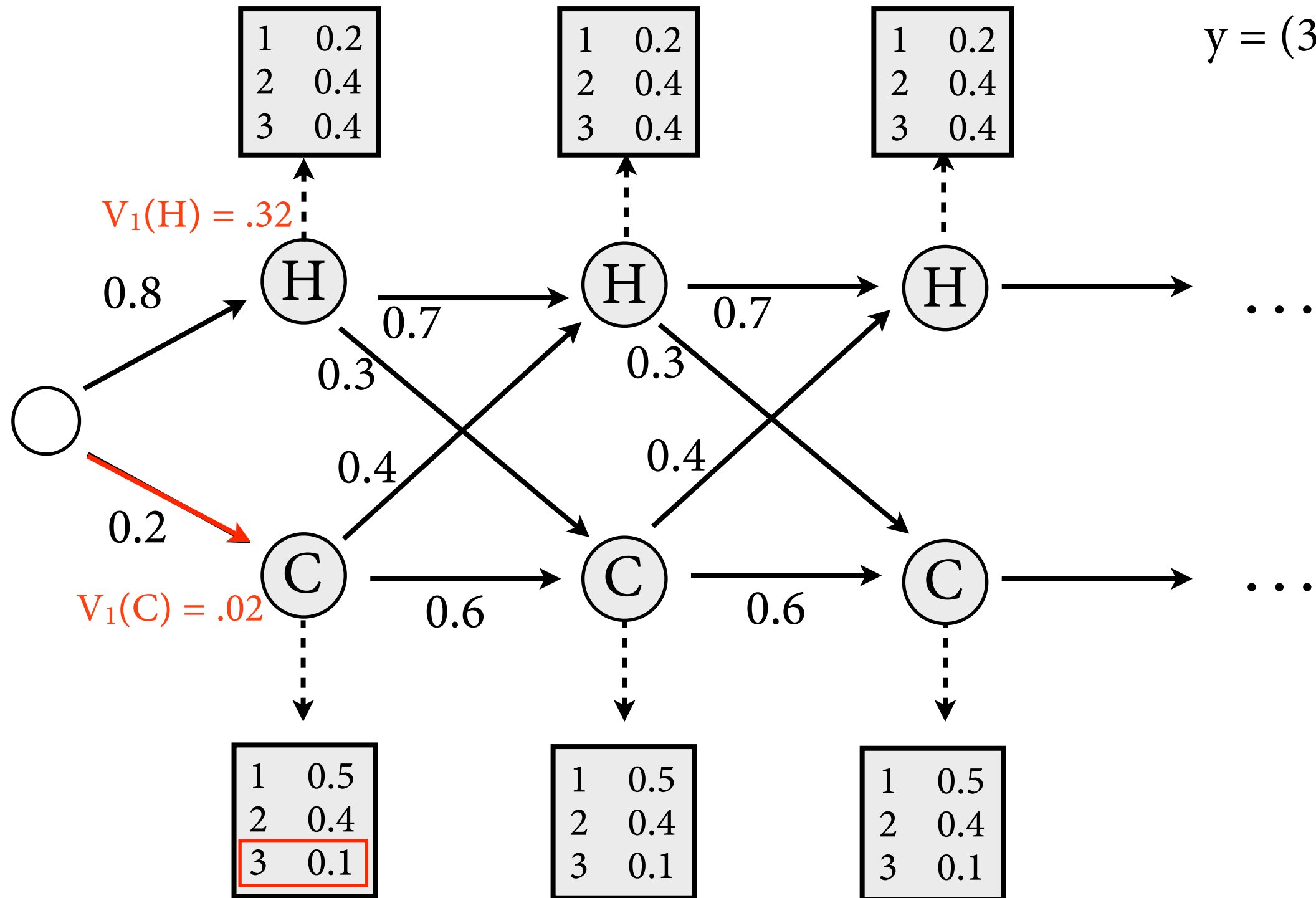
$y = (3,1,3)$



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

# Computing best path probs

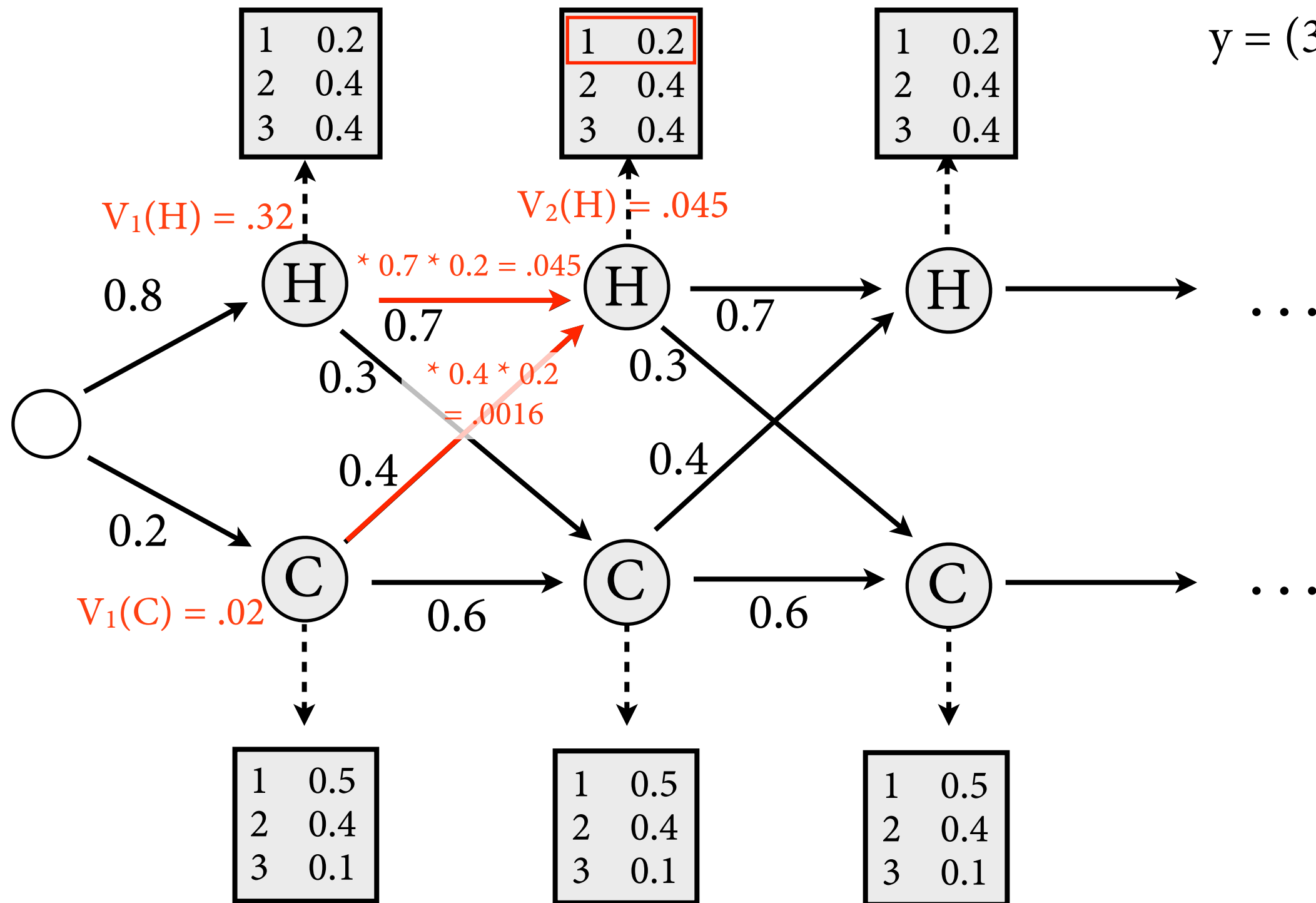
$y = (3,1,3)$



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

# Computing best path probs

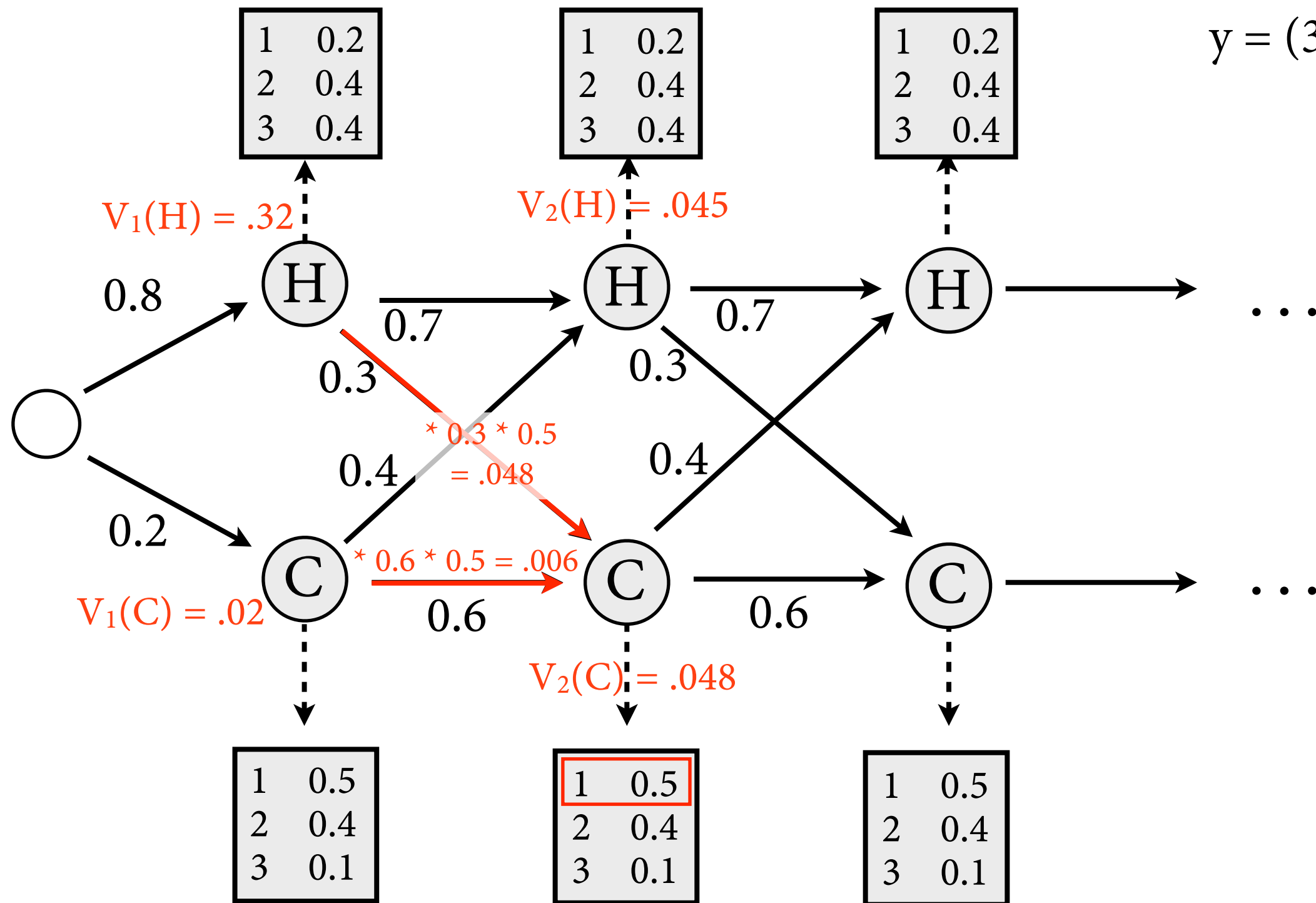
$y = (3,1,3)$



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

# Computing best path probs

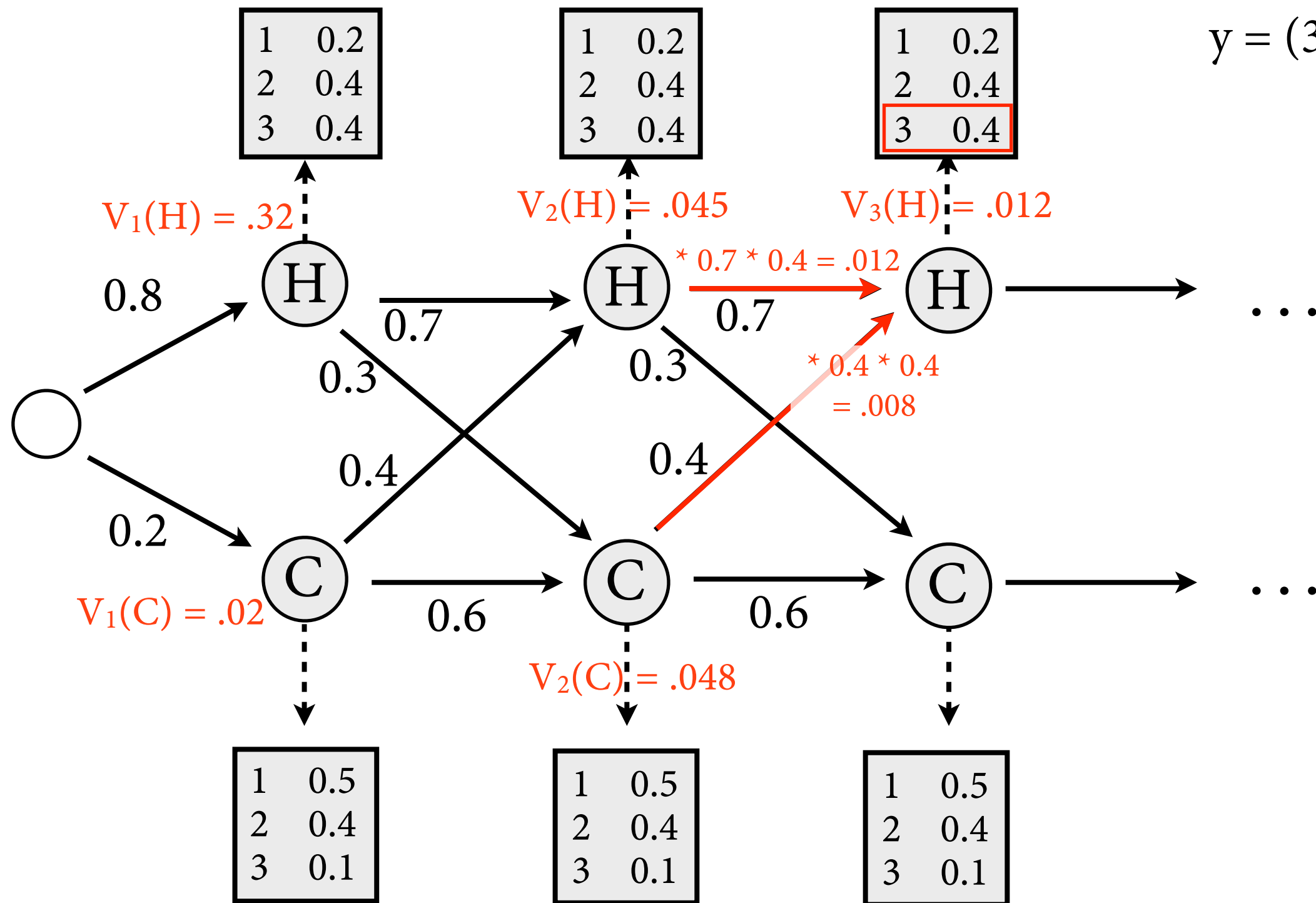
$y = (3,1,3)$



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

# Computing best path probs

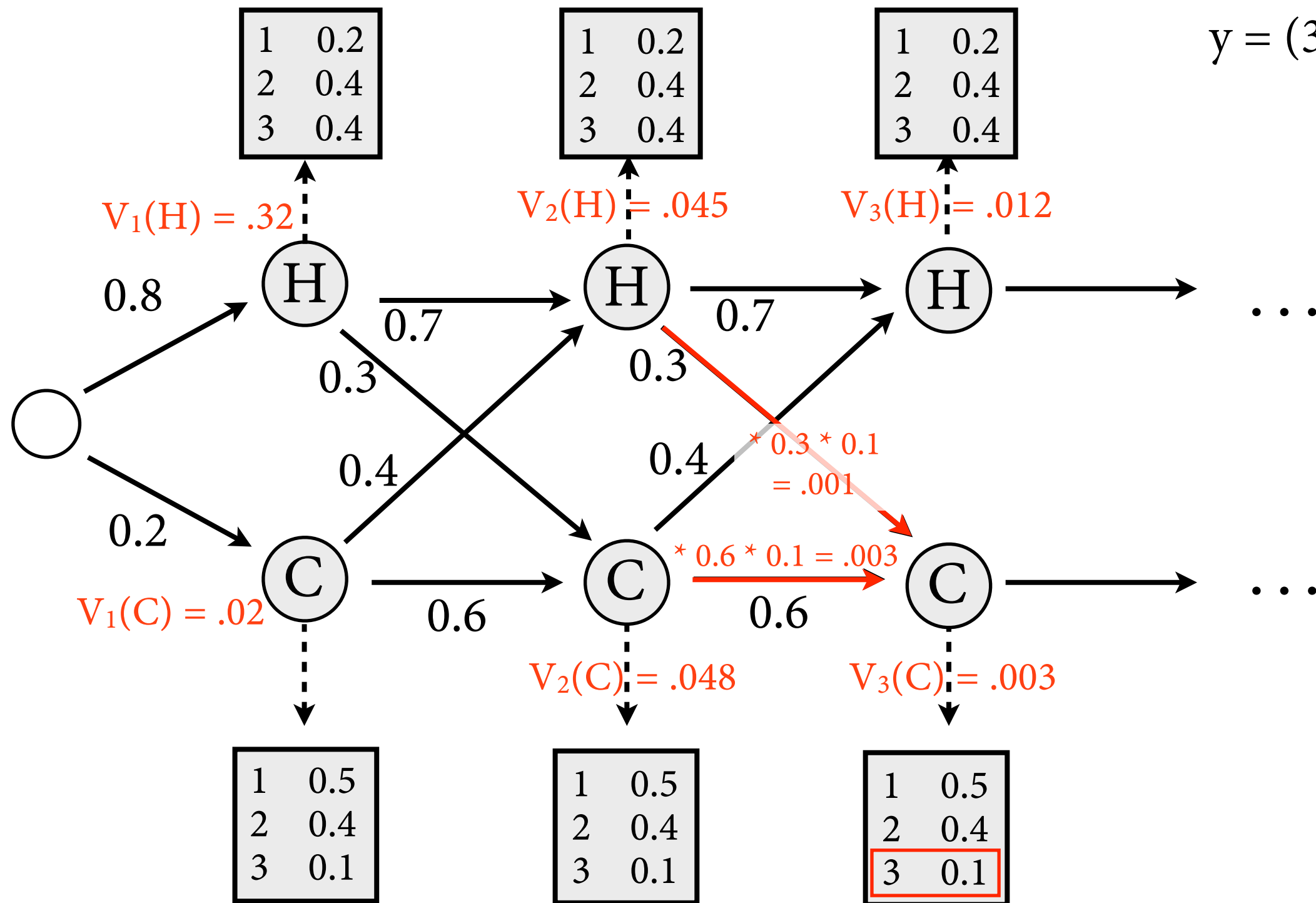
$y = (3,1,3)$



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

# Computing best path probs

$y = (3,1,3)$



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

# The Viterbi Algorithm

- Because of statistical independencies, can decompose joint probability:

$$\begin{aligned} P(x_1, y_1, \dots, x_t, y_t) &= P(y_t \mid x_1, \dots, x_t, y_1, \dots, y_{t-1}) \cdot P(x_t \mid x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}) \\ &\quad \cdot P(x_1, y_1, \dots, x_{t-1}, y_{t-1}) \\ &= P(y_t \mid x_t) \cdot P(x_t \mid x_{t-1}) \cdot P(x_1, y_1, \dots, x_{t-1}, y_{t-1}) \end{aligned}$$

- Thus, maximum has recursive structure:

$$\begin{aligned} V_t(j) &= \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j) \\ &= \max_{x_1, \dots, x_{t-1}} P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid x_{t-1}) \cdot P(y_1, \dots, y_{t-1}, x_1, \dots, x_{t-1}) \\ &= \max_{x_{t-1}} P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid x_{t-1}) \cdot \left( \max_{x_1, \dots, x_{t-2}} P(y_1, \dots, y_{t-1}, x_1, \dots, x_{t-1}) \right) \\ &= \max_i P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot \left( \max_{x_1, \dots, x_{t-2}} P(y_1, \dots, y_{t-1}, x_1, \dots, X_{t-1} = q_i) \right) \\ &= \max_i P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot V_{t-1}(i) = \max_{i=1}^N V_{t-1}(i) \cdot a_{ij} \cdot b_j(y_t) \end{aligned}$$



# The Viterbi Algorithm

- Because of statistical independencies, can decompose joint probability:

$$\begin{aligned} P(x_1, y_1, \dots, x_t, y_t) &= P(y_t \mid x_1, \dots, x_t, y_1, \dots, y_{t-1}) \cdot P(x_t \mid x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}) \\ &\quad \cdot P(x_1, y_1, \dots, x_{t-1}, y_{t-1}) \\ &= P(y_t \mid x_t) \cdot P(x_t \mid x_{t-1}) \cdot P(x_1, y_1, \dots, x_{t-1}, y_{t-1}) \end{aligned}$$

- Thus, maximum has recursive structure:

$$\begin{aligned} V_t(j) &= \max_{x_1, \dots, x_{t-1}} P(y_1, \dots, y_t, x_1, \dots, x_{t-1}, X_t = q_j) \\ &= \max_{x_1, \dots, x_{t-1}} P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid x_{t-1}) \cdot P(y_1, \dots, y_{t-1}, x_1, \dots, x_{t-1}) \\ &= \max_{x_{t-1}} P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid x_{t-1}) \cdot \left( \max_{x_1, \dots, x_{t-2}} P(y_1, \dots, y_{t-1}, x_1, \dots, x_{t-1}) \right) \\ &= \max_i P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot \left( \max_{x_1, \dots, x_{t-2}} P(y_1, \dots, y_{t-1}, x_1, \dots, X_{t-1} = q_i) \right) \\ &= \max_i P(y_t \mid X_t = q_j) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot V_{t-1}(i) \end{aligned}$$

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

# 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,  $t = 1$ :

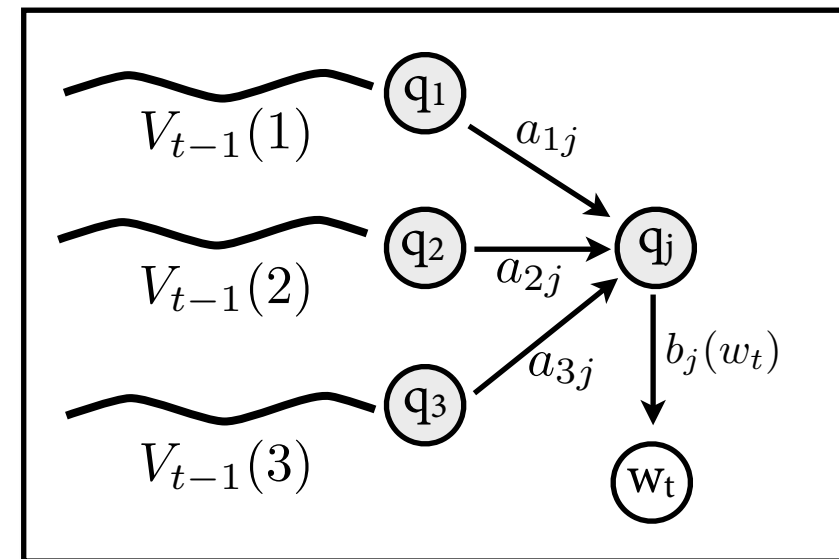
$$V_1(j) = P(y_1, X_1 = q_j) = b_j(y_1) \cdot c_j$$

- 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:

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



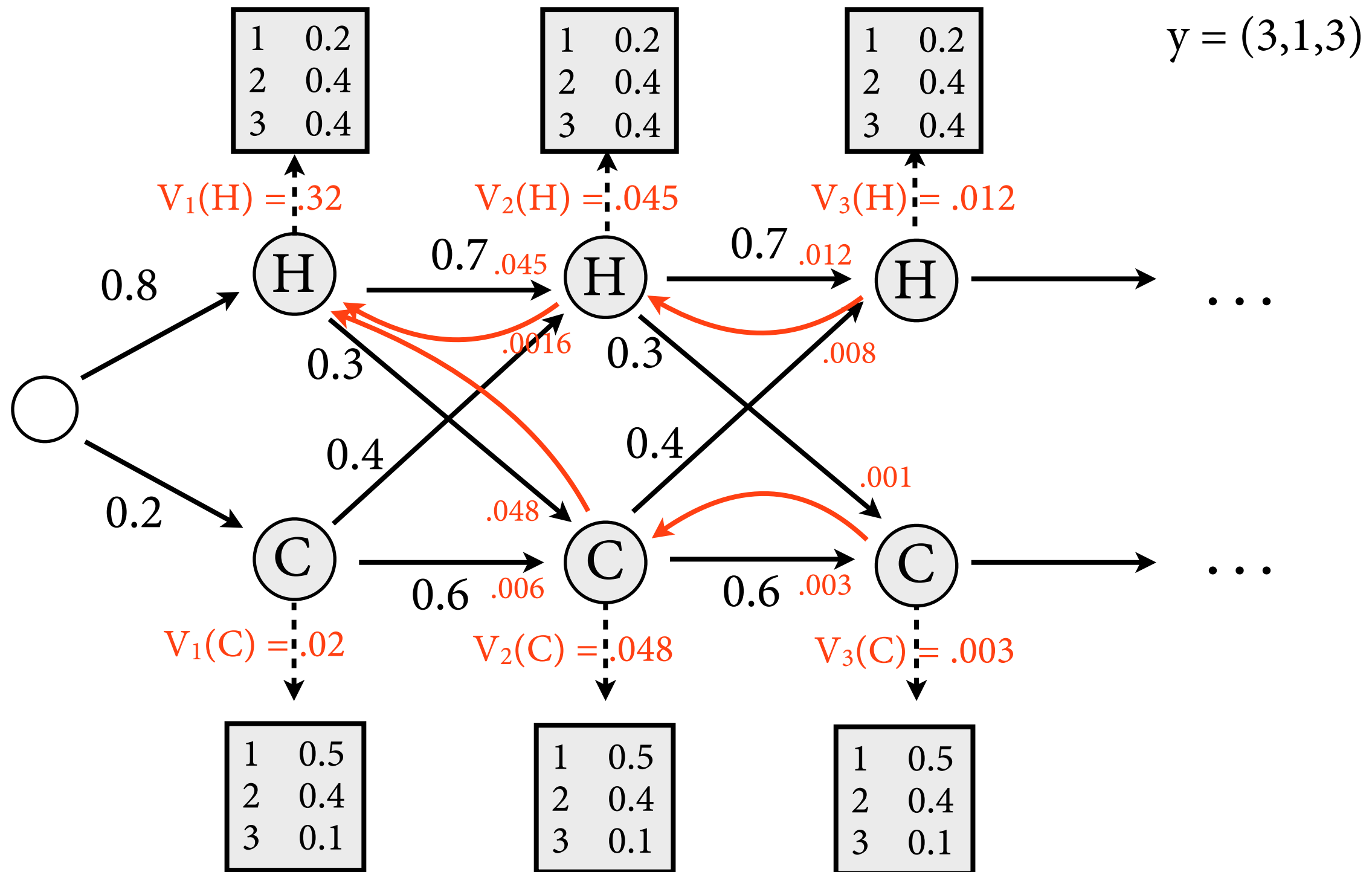
# Backpointers

- 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



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

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

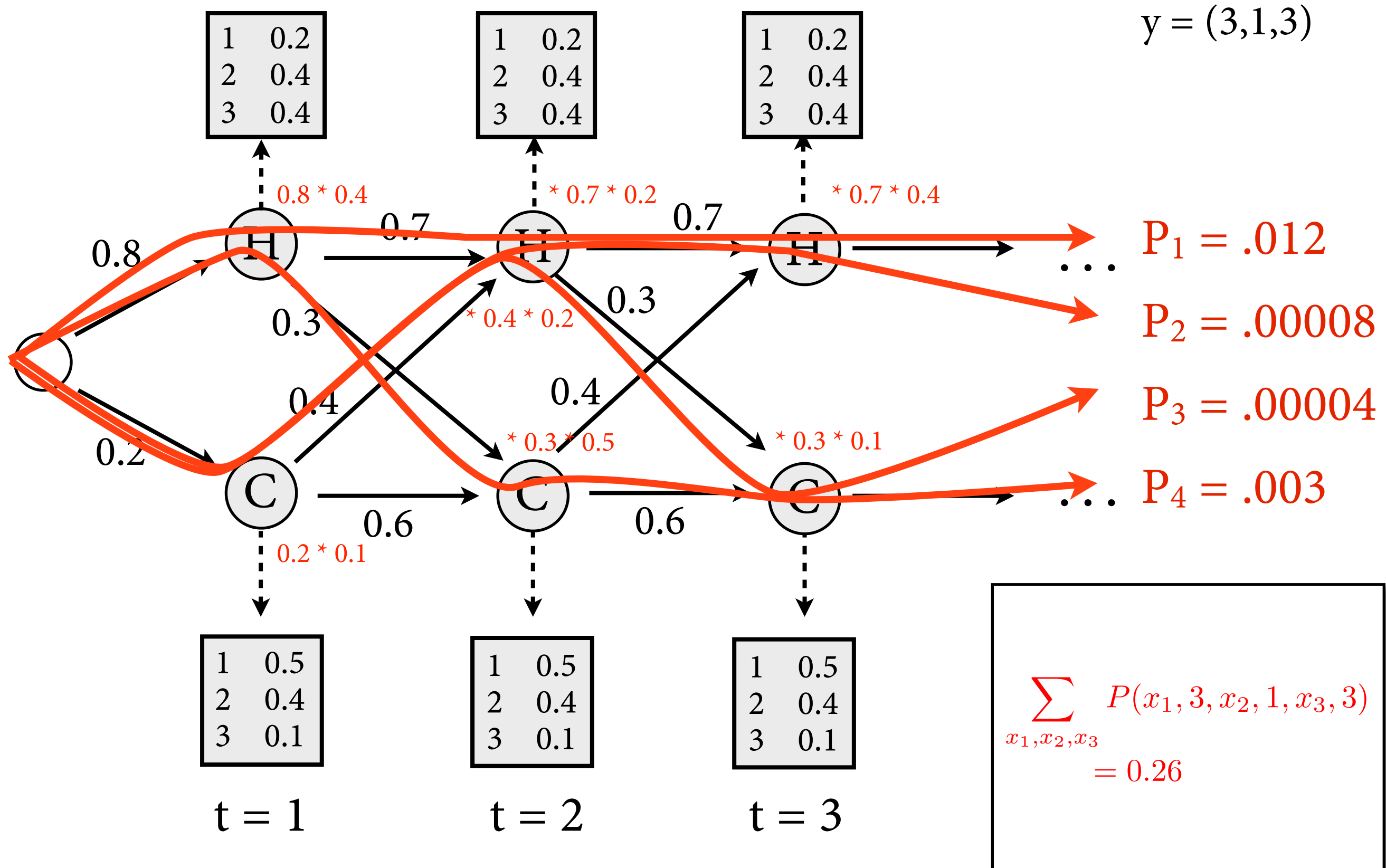
# Question 1: Likelihood, $P(y)$

- How likely is it that Jason Eisner ate 3 ice creams on day 1, 1 ice cream on day 2, 3 ice creams on day 3?
- Want to compute:  $P(3, 1, 3)$ .
- Same problem as with max:
  - ▶ Output 3, 1, 3 can be emitted by many different state sequences.
  - ▶ Obtain by marginalization:

$$P(3, 1, 3) = \sum_{x_1, x_2, x_3 \in Q} P(x_1, 3, x_2, 1, x_3, 3)$$

- ▶ Naive computation is far too slow.

# Ice cream trellis



# 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,  $t = 1$ :

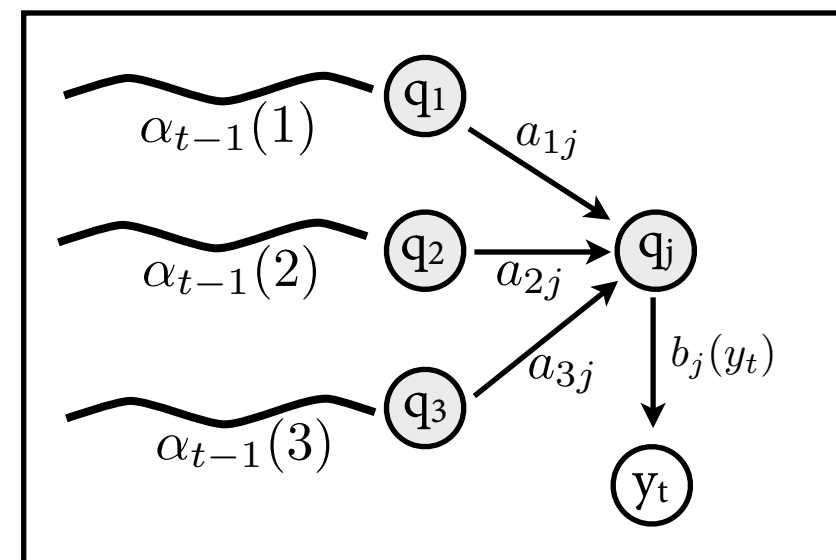
$$\alpha_1(j) = P(y_1, X_1 = q_j) = b_j(y_1) \cdot c_j$$

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

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

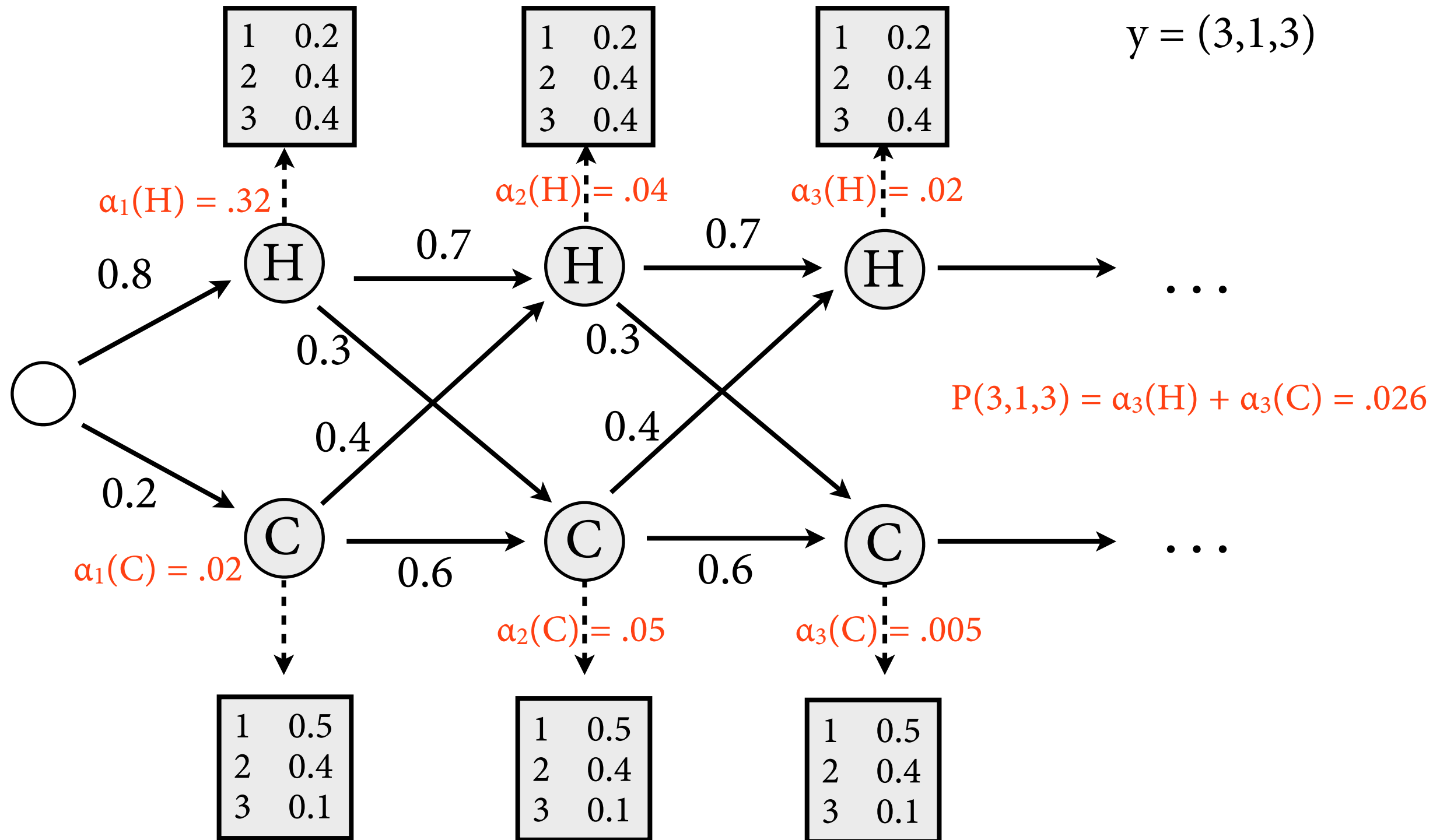
$$= \sum_{i=1}^N P(y_1, \dots, y_{t-1}, X_{t-1} = q_i) \cdot P(X_t = q_j \mid X_{t-1} = q_i) \cdot P(y_t \mid X_t = q_j)$$

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





# P(3,1,3) with Forward



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

$$\alpha_1(j) = b_j(y_1) \cdot a_{0j}$$

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

# Runtime

- Forward and Viterbi have the same runtime, dominated by inductive step:

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

- Compute  $N \cdot T$  values for  $V_t(j)$ . Each computation step requires iteration over  $N$  predecessor states.
- Total runtime is  $O(N^2 \cdot T)$ , i.e.
  - ▶ linear in sentence length
  - ▶ quadratic in size of tag set

# Summary

- Hidden Markov Models popular model for POS tagging (and other applications, see later).
- Two coupled random processes:
  - ▶ bigram model for hidden states (“Markov Chain”)
  - ▶ model for producing observable output from each state
- Efficient algorithms for common problems:
  - ▶ Likelihood computation: Forward algorithm
  - ▶ Best state sequence: Viterbi algorithm.