

*INF4820: Algorithms for  
Artificial Intelligence and  
Natural Language Processing*

Language Models & Hidden Markov Models

Stephan Oepen & Erik Velldal

Language Technology Group (LTG)

October 14, 2015

# Recall: By the End of the Semester . . .



. . . you should be able to determine

- ▶ which string is most likely:

- ▶ *How to recognize speech vs. How to wreck a nice beach*

# Recall: By the End of the Semester . . .



. . . you should be able to determine

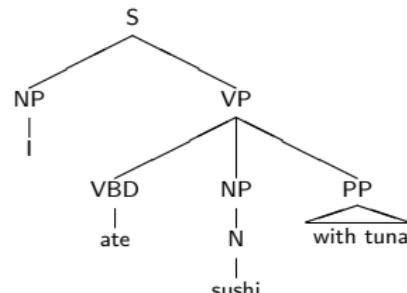
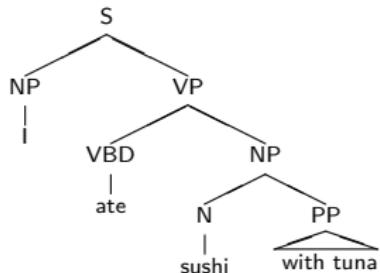
- ▶ which string is most likely:
  - ▶ *How to recognize speech* vs. *How to wreck a nice beach*
- ▶ which category sequence is most likely for *flies like an arrow*:
  - ▶ **N V D N** vs. **V P D N**

# Recall: By the End of the Semester ...



... you should be able to determine

- ▶ which string is most likely:
  - ▶ *How to recognize speech vs. How to wreck a nice beach*
- ▶ which category sequence is most likely for *flies like an arrow*:
  - ▶ **N V D N** vs. **V P D N**
- ▶ which syntactic analysis is most likely:

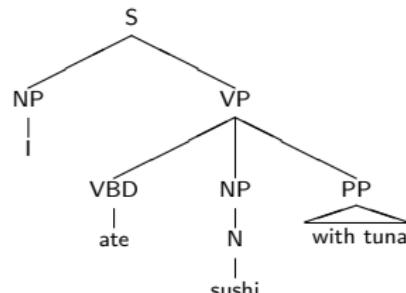
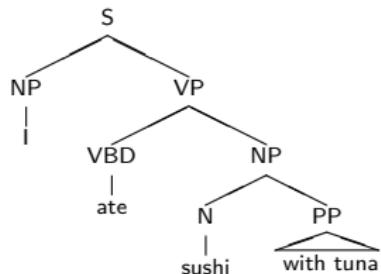


# Recall: By the End of the Semester ...



... you should be able to determine

- ▶ which string is **most likely**:
  - ▶ *How to recognize speech vs. How to wreck a nice beach*
- ▶ which category sequence is **most likely** for *flies like an arrow*:
  - ▶ **N V D N** vs. **V P D N**
- ▶ which syntactic analysis is **most likely**:



# Language Models — $N$ -Grams



A probabilistic (or stochastic) **language model**  $M$  assigns probabilities  $P_M(x)$  to all strings  $x$  in language  $L$ .

We simplify using the **Markov assumption** (limited history):

*the last  $n - 1$  elements approximate the effect of the full sequence.*

That is, instead of

- ▶  $P(w_i | w_1, \dots, w_{i-1})$

selecting an  $n$  of 3, we use

- ▶  $P(w_i | w_{i-1}, w_{i-2})$

# Language Models — $N$ -Grams



A probabilistic (or stochastic) **language model**  $M$  assigns probabilities  $P_M(x)$  to all strings  $x$  in language  $L$ .

We simplify using the **Markov assumption** (limited history):

*the last  $n - 1$  elements approximate the effect of the full sequence.*

That is, instead of

- ▶  $P(w_i|w_1, \dots w_{i-1})$

selecting an  $n$  of 3, we use

- ▶  $P(w_i|w_{i-1}, w_{i-2})$

We call these short sequences of words  $n$ -grams:

- ▶ bigrams: *I want, want to, to go, go to, to the, the beach*
- ▶ trigrams: *I want to, want to go, to go to, go to the*
- ▶ 4-grams: *I want to go, want to go to, to go to the*

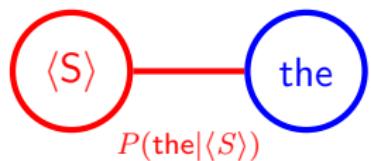
A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:

$\langle S \rangle$

A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:



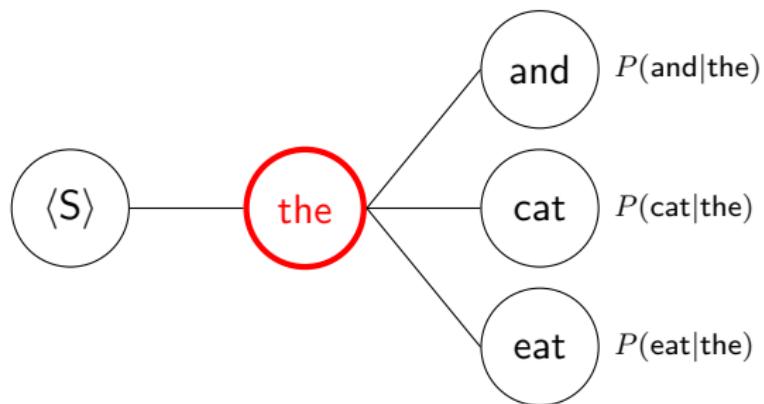
$$P(S) = P(\text{the} | \langle S \rangle)$$

# *N*-Gram Models



A generative model models a joint probability in terms of conditional probabilities.

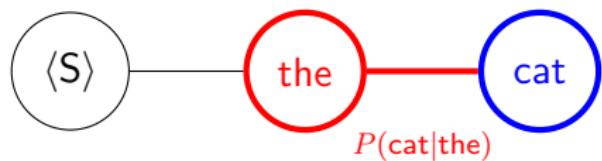
We talk about the *generative story*:



$$P(S) = P(\text{the} | \langle S \rangle)$$

A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:



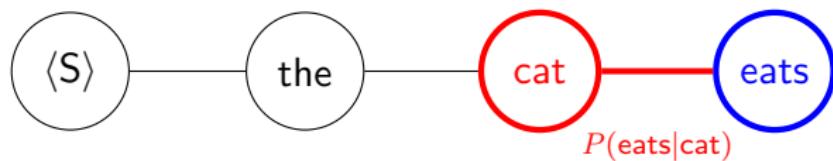
$$P(S) = P(\text{the}|\langle S \rangle) \ P(\text{cat}|\text{the})$$

# *N*-Gram Models



A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:



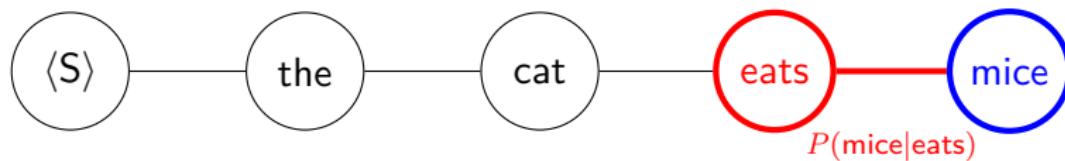
$$P(S) = P(\text{the}|\langle S \rangle) \ P(\text{cat}|\text{the}) \ P(\text{eats}|\text{cat})$$

# *N*-Gram Models



A generative model models a joint probability in terms of conditional probabilities.

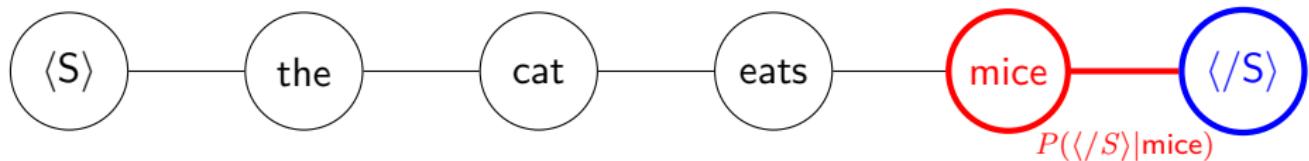
We talk about the *generative story*:



$$P(S) = P(\text{the}|\langle S \rangle) \ P(\text{cat}|\text{the}) \ P(\text{eats}|\text{cat}) \ P(\text{mice}|\text{eats})$$

A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:



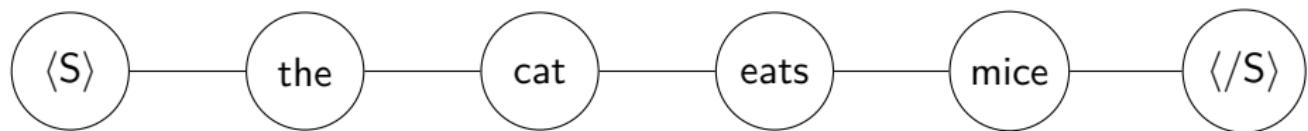
$$P(S) = P(\text{the}|\langle S \rangle) \ P(\text{cat}|\text{the}) \ P(\text{eats}|\text{cat}) \ P(\text{mice}|\text{eats}) \ P(\langle /S \rangle | \text{mice})$$

# *N*-Gram Models



A generative model models a joint probability in terms of conditional probabilities.

We talk about the *generative story*:



$$P(S) = P(\text{the}|\langle S \rangle) \ P(\text{cat}|\text{the}) \ P(\text{eats}|\text{cat}) \ P(\text{mice}|\text{eats}) \ P(\langle /S \rangle|\text{mice})$$

# *N*-Gram Models



An  $n$ -gram language model records the  $n$ -gram conditional probabilities:

$$\begin{array}{ll} P(I|\langle S \rangle) & = 0.0429 \\ P(want|I) & = 0.0111 \\ P(to|want) & = 0.4810 \\ P(go|to) & = 0.0131 \end{array} \quad \begin{array}{ll} P(to|go) & = 0.1540 \\ P(the|to) & = 0.1219 \\ P(beach|the) & = 0.0006 \end{array}$$

# *N*-Gram Models



An  $n$ -gram language model records the  $n$ -gram conditional probabilities:

$$\begin{array}{ll} P(I|\langle S \rangle) & = 0.0429 \\ P(want|I) & = 0.0111 \\ P(to|want) & = 0.4810 \\ P(go|to) & = 0.0131 \end{array} \quad \begin{array}{ll} P(to|go) & = 0.1540 \\ P(the|to) & = 0.1219 \\ P(beach|the) & = 0.0006 \end{array}$$

We calculate the probability of a sentence as (assuming bi-grams):

$$\begin{aligned} P(w_1^n) &\approx \prod_{i=1}^n P(w_i|w_{i-1}) \\ &\approx P(I|\langle S \rangle) \times P(want|I) \times P(to|want) \times P(go|to) \times P(to|go) \times \\ &\quad P(the|to) \times P(beach|the) \\ &\approx 0.0429 \times 0.0111 \times 0.4810 \times 0.0131 \times 0.1540 \times \\ &\quad 0.1219 \times 0.0006 = 3.38 \times 10^{-11} \end{aligned}$$

# Training an $N$ -Gram Model



How to estimate the probabilities of  $n$ -grams?

By counting (e.g. for trigrams):

$$P(\text{bananas} | \text{i like}) = \frac{C(\text{i like bananas})}{C(\text{i like})}$$

The probabilities are estimated using the **relative frequencies** of observed outcomes. This process is called **Maximum Likelihood Estimation (MLE)**.

# Bigram MLE Example



“I want to go to the beach”

$w_1$	$w_2$	$C(w_1 w_2)$	$C(w_1)$	$P(w_2   w_1)$
$\langle S \rangle$	I	1039	24243	0.0429
I	want	46	4131	0.0111
want	to	101	210	0.4810
to	go	128	9778	0.0131
go	to	59	383	0.1540
to	the	1192	9778	0.1219
the	beach	14	22244	0.0006

# Bigram MLE Example



“I want to go to the beach”

$w_1$	$w_2$	$C(w_1 w_2)$	$C(w_1)$	$P(w_2   w_1)$
$\langle S \rangle$	I	1039	24243	0.0429
I	want	46	4131	0.0111
want	to	101	210	0.4810
to	go	128	9778	0.0131
go	to	59	383	0.1540
to	the	1192	9778	0.1219
the	beach	14	22244	0.0006

What's the probability of *Others want to go to the beach* ?

# Problems with MLE of $N$ -Grams



- ▶ Data sparseness: many perfectly acceptable  $n$ -grams will not be observed
- ▶ Zero counts will result in an estimated probability of 0

# Problems with MLE of $N$ -Grams



- ▶ Data sparseness: many perfectly acceptable  $n$ -grams will not be observed
- ▶ Zero counts will result in an estimated probability of 0
- ▶ Remedy—reassign some of the probability mass of frequent events to less frequent (or unseen) events.
- ▶ Known as **smoothing** or **discounting**
- ▶ The simplest approach is **Laplace** ('add-one') smoothing:

$$P_L(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n) + 1}{C(w_{n-1}) + V}$$

# Bigram MLE Example with Laplace Smoothing



“Others want to go to the beach”

$w_1$	$w_2$	$C(w_1 w_2)$	$C(w_1)$	$P(w_2 w_1)$
$\langle S \rangle$	I	1039	24243	0.0429
$\langle S \rangle$	Others	17	24243	0.0007
I	want	46	4131	0.0111
Others	want	0	4131	0
want	to	101	210	0.4810
to	go	128	9778	0.0131
go	to	59	383	0.1540
to	the	1192	9778	0.1219
the	beach	14	22244	0.0006

# Bigram MLE Example with Laplace Smoothing



“Others want to go to the beach”

$w_1$	$w_2$	$C(w_1 w_2)$	$C(w_1)$	$P(w_2   w_1)$
$\langle S \rangle$	I	1039	24243	0.0429
$\langle S \rangle$	Others	17	24243	0.0007
I	want	46	4131	0.0111
Others	want	0	4131	0
want	to	101	210	0.4810
to	go	128	9778	0.0131
go	to	59	383	0.1540
to	the	1192	9778	0.1219
the	beach	14	22244	0.0006

# Bigram MLE Example with Laplace Smoothing



“Others want to go to the beach”

$w_1$	$w_2$	$C(w_1 w_2)$	$C(w_1)$	$P(w_2 w_1)$	$P_L(w_2 w_1)$
$\langle S \rangle$	I	1039	24243	0.0429	0.01934
$\langle S \rangle$	Others	17	24243	0.0007	0.00033
I	want	46	4131	0.0111	0.00140
Others	want	0	4131	0	0.00003
want	to	101	210	0.4810	0.00343
to	go	128	9778	0.0131	0.00328
go	to	59	383	0.1540	0.00201
to	the	1192	9778	0.1219	0.03035
the	beach	14	22244	0.0006	0.00029

$$P_L(w_n|w_{n-1}) = \frac{C(w_{n-1} w_n) + 1}{C(w_{n-1}) + 29534}$$

- ▶ The likelihood of the next word depends on its context.
- ▶ We can calculate this using the chain rule:

$$P(w_1^N) = \prod_{i=1}^N P(w_i | w_1^{i-1})$$

- ▶ In an  $n$ -gram model, we approximate this with a Markov chain:

$$P(w_1^N) \approx \prod_{i=1}^N P(w_i | w_{i-n+1}^{i-1})$$

- ▶ We use Maximum Likelihood Estimation to estimate the conditional probabilities.
- ▶ Smoothing techniques are used to avoid zero probabilities.

# Parts of Speech



- ▶ Known by a variety of names: part-of-speech, POS, lexical categories, word classes, morpho-syntactic classes, . . .
- ▶ ‘Traditionally’ defined semantically (e.g. “nouns are naming words”), but (arguably) more accurately by their distributional properties.

# Parts of Speech



- ▶ Known by a variety of names: part-of-speech, POS, lexical categories, word classes, morpho-syntactic classes, . . .
- ▶ ‘Traditionally’ defined semantically (e.g. “nouns are naming words”), but (arguably) more accurately by their distributional properties.
- ▶ Open-classes
  - ▶ New words created/updated/deleted all the time

# Parts of Speech



- ▶ Known by a variety of names: part-of-speech, POS, lexical categories, word classes, morpho-syntactic classes, . . .
- ▶ ‘Traditionally’ defined semantically (e.g. “nouns are naming words”), but (arguably) more accurately by their distributional properties.
  
- ▶ Open-classes
  - ▶ New words created/updated/deleted all the time
- ▶ Closed-classes
  - ▶ Smaller classes, relatively static membership
  - ▶ Usually function words

# Open Class Words



- ▶ Nouns: dog, Oslo, scissors, snow, people, truth, cups
  - ▶ proper or common; countable or uncountable; plural or singular; masculine, feminine, or neuter; ...

# Open Class Words



- ▶ Nouns: dog, Oslo, scissors, snow, people, truth, cups
  - ▶ proper or common; countable or uncountable; plural or singular; masculine, feminine, or neuter; ...
- ▶ Verbs: fly, rained, having, ate, seen
  - ▶ transitive, intransitive, ditransitive; past, present, passive; stative or dynamic; plural or singular; ...

# Open Class Words



- ▶ Nouns: dog, Oslo, scissors, snow, people, truth, cups
  - ▶ proper or common; countable or uncountable; plural or singular; masculine, feminine, or neuter; ...
- ▶ Verbs: fly, rained, having, ate, seen
  - ▶ transitive, intransitive, ditransitive; past, present, passive; stative or dynamic; plural or singular; ...
- ▶ Adjectives: good, smaller, unique, fastest, best, unhappy
  - ▶ comparative or superlative; predicative or attributive; intersective, subsective, or scopal; ...

# Open Class Words



- ▶ Nouns: dog, Oslo, scissors, snow, people, truth, cups
  - ▶ proper or common; countable or uncountable; plural or singular; masculine, feminine, or neuter; ...
- ▶ Verbs: fly, rained, having, ate, seen
  - ▶ transitive, intransitive, ditransitive; past, present, passive; stative or dynamic; plural or singular; ...
- ▶ Adjectives: good, smaller, unique, fastest, best, unhappy
  - ▶ comparative or superlative; predicative or attributive; intersective, subsective, or scopal; ...
- ▶ Adverbs: again, somewhat, slowly, yesterday, aloud
  - ▶ intersective; scopal; discourse; degree; temporal; directional; comparative or superlative; ...

# Closed Class Words



- ▶ Prepositions: on, under, from, at, near, over, ...
- ▶ Determiners: a, an, the, that, ...
- ▶ Pronouns: she, who, I, others, ...
- ▶ Conjunctions: and, but, or, when, ...
- ▶ Auxiliary verbs: can, may, should, must, ...
- ▶ Interjections, particles, numerals, negatives, politeness markers, greetings, existential there ...

(Examples from Jurafsky & Martin, 2008)

# POS Tagging



The (automatic) assignment of POS tags to word sequences

- ▶ non-trivial where words are ambiguous: *fly* (v) vs. *fly* (n)
- ▶ choice of the correct tag is *context-dependent*
- ▶ useful in pre-processing for parsing, etc; but also directly for text-to-speech synthesis: **content** (n) vs. **content** (adj)

# POS Tagging



The (automatic) assignment of POS tags to word sequences

- ▶ non-trivial where words are ambiguous: *fly* (v) vs. *fly* (n)
- ▶ choice of the correct tag is *context-dependent*
- ▶ useful in pre-processing for parsing, etc; but also directly for text-to-speech synthesis: **content** (n) vs. **content** (adj)
- ▶ difficulty and usefulness can depend on the *tagset*
  - ▶ English
    - ▶ Penn Treebank (PTB)—45 tags: NNS, NN, NNP, JJ, JJR, JJS
    - <http://bulba.sdsu.edu/jeanette/thesis/PennTags.html>
  - ▶ Norwegian
    - ▶ Oslo-Bergen Tagset—multi-part: <subst appell fem be ent>
    - <http://tekstlab.uio.no/obt-ny/english/tags.html>

# Labeled Sequences



- We are interested in the probability of sequences like:

flies	like	the	wind	or	flies	like	the	wind
NN S	VB	DT	NN		VBZ	P	DT	NN

- We are interested in the probability of sequences like:

flies	like	the	wind	or	flies	like	the	wind
NNS	VB	DT	NN		VBZ	P	DT	NN

- In normal text, we see the words, but not the tags.
- Consider the POS tags to be underlying skeleton of the sentence, unseen but influencing the sentence shape.
- A structure like this, consisting of a **hidden** state sequence, and a related **observation** sequence can be modelled as a *Hidden Markov Model*.

# Hidden Markov Models



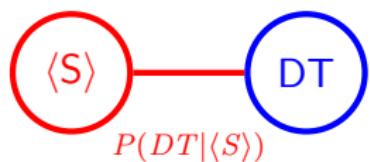
The generative story:

 $\langle S \rangle$

# Hidden Markov Models



The generative story:

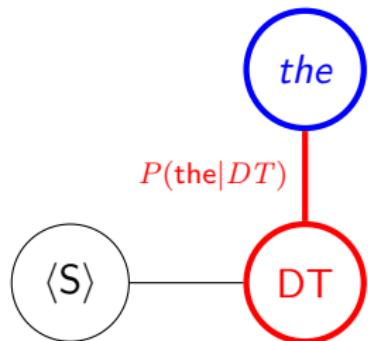


$$P(S, O) = P(\text{DT} | \langle S \rangle)$$

# Hidden Markov Models



The generative story:

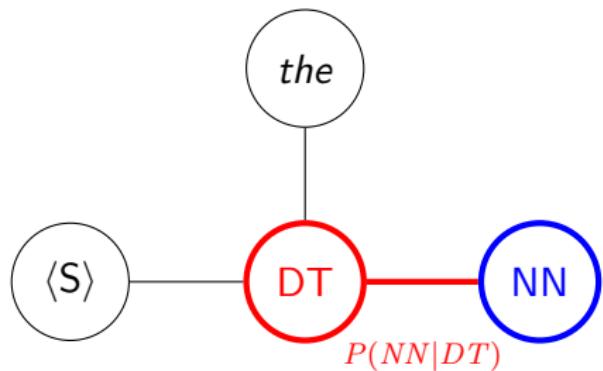


$$P(S, O) = P(\text{DT} | \langle S \rangle) P(\text{the} | \text{DT})$$

# Hidden Markov Models



The generative story:

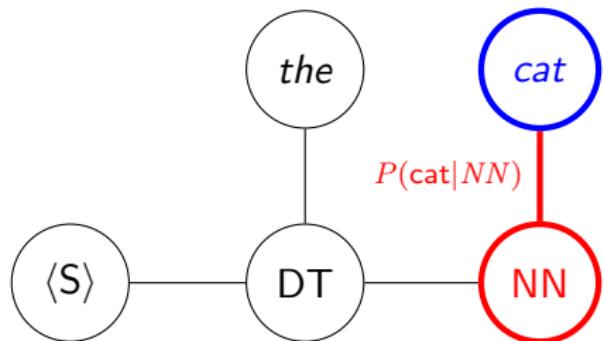


$$P(S, O) = P(\text{DT} | \langle S \rangle) P(\text{the} | \text{DT}) P(\text{NN} | \text{DT})$$

# Hidden Markov Models



The generative story:

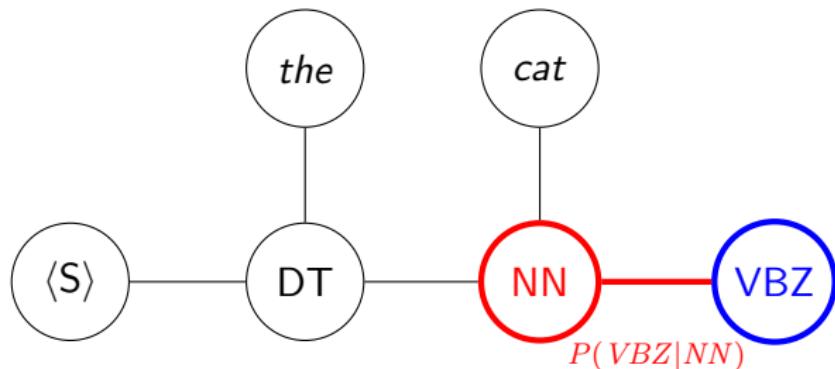


$$P(S, O) = P(\text{DT} | \langle S \rangle) P(\text{the} | \text{DT}) P(\text{NN} | \text{DT}) P(\text{cat} | \text{NN})$$

# Hidden Markov Models



The generative story:

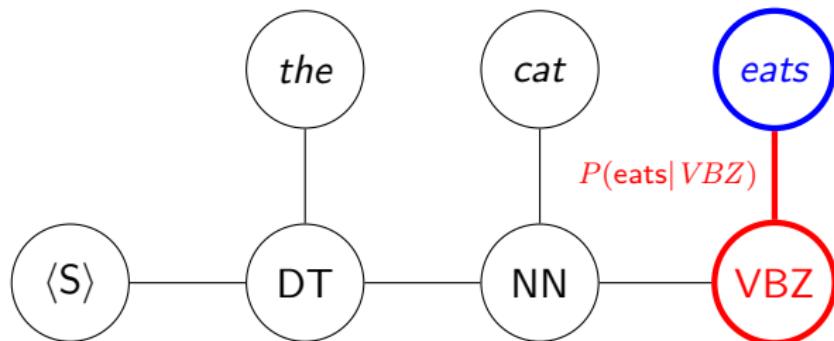


$$P(S, O) = P(\text{DT} | \langle S \rangle) P(\text{the} | \text{DT}) P(\text{NN} | \text{DT}) P(\text{cat} | \text{NN}) \\ P(\text{VBZ} | \text{NN})$$

# Hidden Markov Models



The generative story:

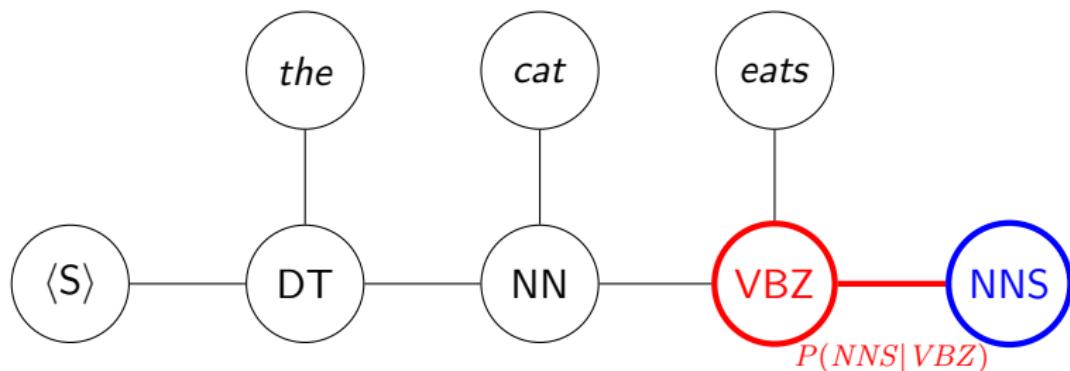


$$P(S, O) = P(\text{DT}|\langle S \rangle) P(\text{the}|\text{DT}) P(\text{NN}|\text{DT}) P(\text{cat}|\text{NN}) \\ P(\text{VBZ}|\text{NN}) P(\text{eats}|\text{VBZ})$$

# Hidden Markov Models



The generative story:

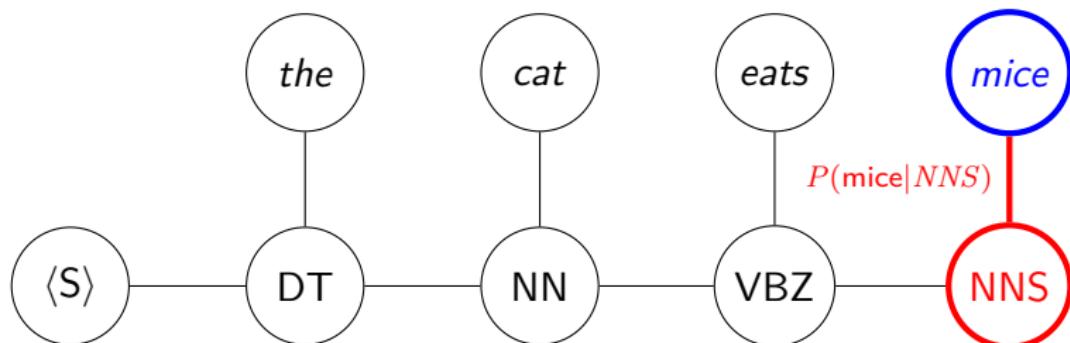


$$P(S, O) = P(\text{DT} | \langle S \rangle) P(\text{the} | \text{DT}) P(\text{NN} | \text{DT}) P(\text{cat} | \text{NN}) \\ P(\text{VBZ} | \text{NN}) P(\text{eats} | \text{VBZ}) P(\text{NNS} | \text{VBZ})$$

# Hidden Markov Models



The generative story:

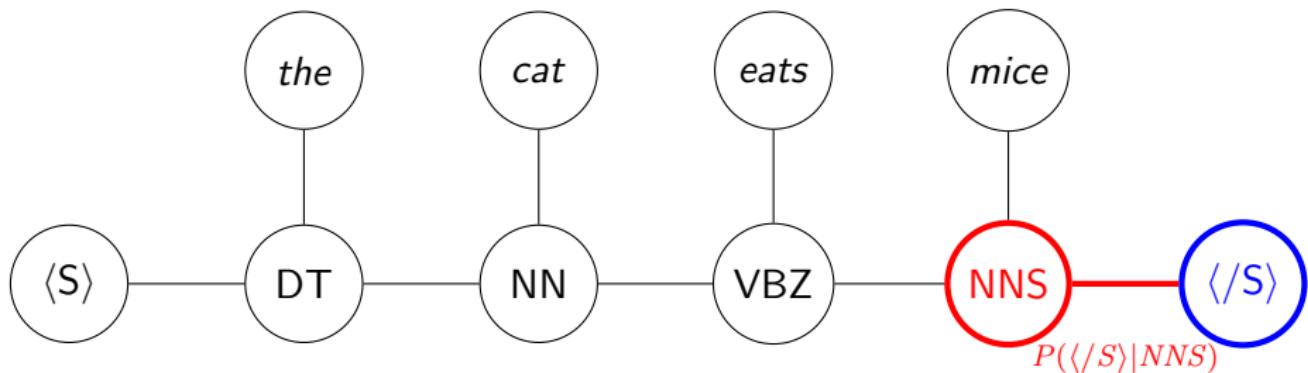


$$P(S, O) = P(\text{DT}|\langle S \rangle) P(\text{the}|DT) P(\text{NN}|DT) P(\text{cat}|NN) \\ P(\text{VBZ}|NN) P(\text{eats}|VBZ) P(\text{NNS}|VBZ) P(\text{mice}|NNS)$$

# Hidden Markov Models



The generative story:

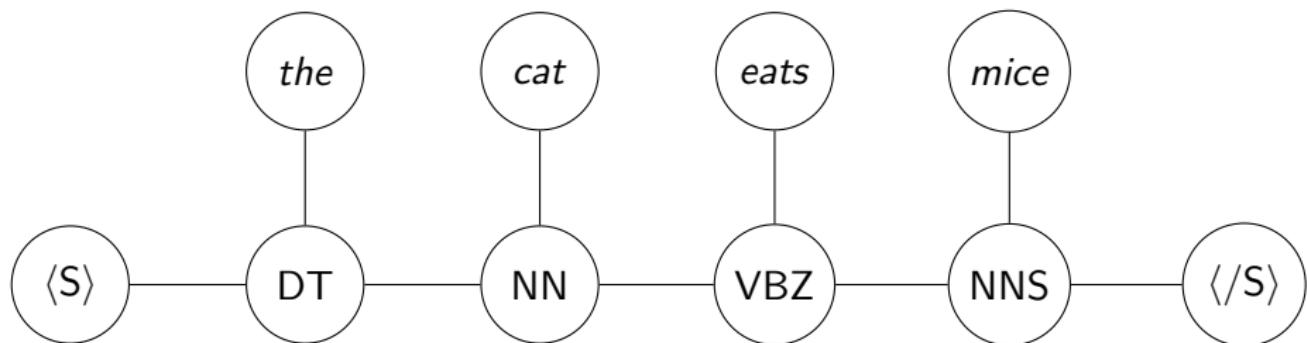


$$\begin{aligned}P(S, O) = & P(\text{ DT}|<S>) \ P(\text{the}|DT) \ P(\text{NN}|DT) \ P(\text{cat}|NN) \\& P(\text{VBZ}|NN) \ P(\text{eats}|VBZ) \ P(\text{NNS}|VBZ) \ P(\text{mice}|NNS) \\& P(</S>|NNS)\end{aligned}$$

# Hidden Markov Models



The generative story:



$$\begin{aligned} P(S, O) = & P(\text{ DT}|<\text{S}>) \ P(\text{the}|DT) \ P(\text{NN}|DT) \ P(\text{cat}|NN) \\ & P(\text{VBZ}|NN) \ P(\text{eats}|VBZ) \ P(\text{NNS}|VBZ) \ P(\text{mice}|NNS) \\ & P(</\text{S}>|\text{NNS}) \end{aligned}$$

# Hidden Markov Models



For a bi-gram HMM, with observations  $O_1^N$ :

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i) \quad \text{where} \quad s_0 = \langle S \rangle, \quad s_{N+1} = \langle /S \rangle$$

# Hidden Markov Models



For a bi-gram HMM, with observations  $O_1^N$ :

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i) \quad \text{where} \quad s_0 = \langle S \rangle, \quad s_{N+1} = \langle /S \rangle$$

- ▶ The **transition probabilities** model the probabilities of moving from state to state.

For a bi-gram HMM, with observations  $O_1^N$ :

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i) \quad \text{where} \quad s_0 = \langle S \rangle, \quad s_{N+1} = \langle /S \rangle$$

- ▶ The transition probabilities model the probabilities of moving from state to state.
- ▶ The **emission probabilities** model the probability that a state *emits* a particular observation.

# Using HMMs



The HMM models the process of generating the labeled sequence. We can use this model for a number of tasks:

- ▶  $P(S, O)$  given  $S$  and  $O$
- ▶  $P(O)$  given  $O$
- ▶  $S$  that maximizes  $P(S|O)$  given  $O$
- ▶  $P(s_x|O)$  given  $O$
- ▶ We can learn the model parameters, given labeled observations.



The HMM models the process of generating the labeled sequence. We can use this model for a number of tasks:

- ▶  $P(S, O)$  given  $S$  and  $O$
- ▶  $P(O)$  given  $O$
- ▶  $S$  that maximizes  $P(S|O)$  given  $O$
- ▶  $P(s_x|O)$  given  $O$
- ▶ We can learn the model parameters, given labeled observations.

# Using HMMs



The HMM models the process of generating the labeled sequence. We can use this model for a number of tasks:

- ▶  $P(S, O)$  given  $S$  and  $O$
- ▶  $P(O)$  given  $O$
- ▶  $S$  that maximizes  $P(S|O)$  given  $O$
- ▶  $P(s_x|O)$  given  $O$
- ▶ We can learn the model parameters, given labeled observations.

# Using HMMs



The HMM models the process of generating the labeled sequence. We can use this model for a number of tasks:

- ▶  $P(S, O)$  given  $S$  and  $O$
- ▶  $P(O)$  given  $O$
- ▶  $S$  that maximizes  $P(S|O)$  given  $O$
- ▶  $P(s_x|O)$  given  $O$
- ▶ We can learn the model parameters, given labeled observations.

Our observations will be **words** ( $w_i$ ), and our states PoS **tags** ( $t_i$ ).



As so often in NLP, we learn an HMM from labeled data:

## Transition Probabilities

Based on a training corpus of previously tagged text, with tags as our states, the MLE can be computed from the counts of observed tags:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$



As so often in NLP, we learn an HMM from labeled data:

## Transition Probabilities

Based on a training corpus of previously tagged text, with tags as our states, the MLE can be computed from the counts of observed tags:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

## Emission Probabilities

Computed from relative frequencies in the same way, with the words as observations:

$$P(w_i | t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

# Implementation Considerations



$$\begin{aligned} P(S, O) &= P(s_1|\langle S \rangle)P(o_1|s_1)P(s_2|s_1)P(o_2|s_2)P(s_3|s_2)P(o_3|s_3)\dots \\ &= 0.0429 \times 0.0031 \times 0.0044 \times 0.0001 \times 0.0072 \times \dots \end{aligned}$$

# Implementation Considerations



$$\begin{aligned} P(S, O) &= P(s_1|\langle S \rangle)P(o_1|s_1)P(s_2|s_1)P(o_2|s_2)P(s_3|s_2)P(o_3|s_3)\dots \\ &= 0.0429 \times 0.0031 \times 0.0044 \times 0.0001 \times 0.0072 \times \dots \end{aligned}$$

- ▶ Multiplying many small probabilities → risk of numeric underflow

# Implementation Considerations



$$\begin{aligned}P(S, O) &= P(s_1|\langle S \rangle)P(o_1|s_1)P(s_2|s_1)P(o_2|s_2)P(s_3|s_2)P(o_3|s_3)\dots \\&= 0.0429 \times 0.0031 \times 0.0044 \times 0.0001 \times 0.0072 \times \dots\end{aligned}$$

- ▶ Multiplying many small probabilities → risk of numeric underflow
- ▶ Solution: work in  $\log(\text{arithmic})$  space:
  - ▶  $\log(AB) = \log(A) + \log(B)$
  - ▶ hence  $P(A)P(B) = \exp(\log(A) + \log(B))$
  - ▶  $\log(P(S, O)) = -1.368 + -2.509 + -2.357 + -4 + -2.143 + \dots$

# Implementation Considerations



$$\begin{aligned} P(S, O) &= P(s_1|\langle S \rangle)P(o_1|s_1)P(s_2|s_1)P(o_2|s_2)P(s_3|s_2)P(o_3|s_3)\dots \\ &= 0.0429 \times 0.0031 \times 0.0044 \times 0.0001 \times 0.0072 \times \dots \end{aligned}$$

- ▶ Multiplying many small probabilities → risk of numeric underflow
- ▶ Solution: work in  $\log(\text{arithmic})$  space:
  - ▶  $\log(AB) = \log(A) + \log(B)$
  - ▶ hence  $P(A)P(B) = \exp(\log(A) + \log(B))$
  - ▶  $\log(P(S, O)) = -1.368 + -2.509 + -2.357 + -4 + -2.143 + \dots$

Still, the issues related to MLE that we discussed for  $n$ -gram models also apply here ...

# Ice Cream and Global Warming



## Missing records of weather in Baltimore for Summer 2007

- ▶ Jason likes to eat ice cream.
- ▶ He records his daily ice cream consumption in his diary.
- ▶ The number of ice creams he ate was influenced, but not entirely determined by the weather.
- ▶ Today's weather is partially predictable from yesterday's.



## Missing records of weather in Baltimore for Summer 2007

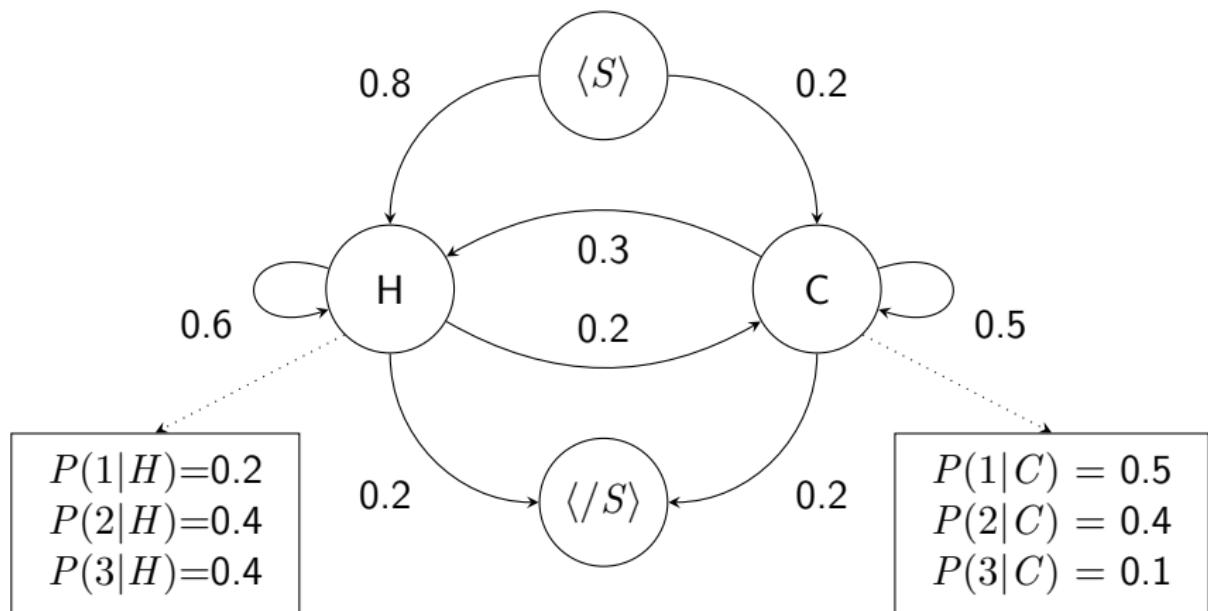
- ▶ Jason likes to eat ice cream.
- ▶ He records his daily ice cream consumption in his diary.
- ▶ The number of ice creams he ate was influenced, but not entirely determined by the weather.
- ▶ Today's weather is partially predictable from yesterday's.

## A Hidden Markov Model

with:

- ▶ Hidden states:  $\{H, C\}$  (plus pseudo-states  $\langle S \rangle$  and  $\langle /S \rangle$ )
- ▶ Observations:  $\{1, 2, 3\}$

# Ice Cream and Global Warming



# Using HMMs

The HMM models the process of generating the labeled sequence. We can use this model for a number of tasks:

- ▶  $P(S, O)$  given  $S$  and  $O$
- ▶  $P(O)$  given  $O$
- ▶  $S$  that maximizes  $P(S|O)$  given  $O$
- ▶  $P(s_x|O)$  given  $O$
- ▶ We can also learn the model parameters, given a set of observations.

# Part-of-Speech Tagging

We want to find the tag sequence, given a word sequence. With tags as our states and words as our observations, we know:

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i)$$

We want:  $P(S|O)$

# Part-of-Speech Tagging

We want to find the tag sequence, given a word sequence. With tags as our states and words as our observations, we know:

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i)$$

We want:  $P(S|O) = \frac{P(S, O)}{P(O)}$

# Part-of-Speech Tagging

We want to find the tag sequence, given a word sequence. With tags as our states and words as our observations, we know:

$$P(S, O) = \prod_{i=1}^{N+1} P(s_i | s_{i-1}) P(o_i | s_i)$$

We want:  $P(S|O) = \frac{P(S, O)}{P(O)}$

Actually, we want the state sequence  $\hat{S}$  that maximizes  $P(S|O)$ :

$$\hat{S} = \arg \max_S \frac{P(S, O)}{P(O)}$$

Since  $P(O)$  always is the same, we can drop the denominator.

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle | H) = 0.2 & P(\langle /S \rangle | C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

if  $O = 3 \ 1 \ 3$

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle | H) = 0.2 & P(\langle /S \rangle | C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

if  $O = 3 \ 1 \ 3$

$P(H \langle S \rangle) = 0.8$	$P(C \langle S \rangle) = 0.2$	$\langle S \rangle$	H	H	H	$\langle /S \rangle$
$P(H H) = 0.6$	$P(C H) = 0.2$					
$P(H C) = 0.3$	$P(C C) = 0.5$					
$P(\langle /S \rangle H) = 0.2$	$P(\langle /S \rangle C) = 0.2$					
$P(1 H) = 0.2$	$P(1 C) = 0.5$					
$P(2 H) = 0.4$	$P(2 C) = 0.4$					
$P(3 H) = 0.4$	$P(3 C) = 0.1$					

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

if  $O = 3\ 1\ 3$

$P(H \langle S \rangle) = 0.8$	$P(C \langle S \rangle) = 0.2$	$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$P(H H) = 0.6$	$P(C H) = 0.2$						
$P(H C) = 0.3$	$P(C C) = 0.5$						
$P(\langle /S \rangle H) = 0.2$	$P(\langle /S \rangle C) = 0.2$						
$P(1 H) = 0.2$	$P(1 C) = 0.5$						
$P(2 H) = 0.4$	$P(2 C) = 0.4$						
$P(3 H) = 0.4$	$P(3 C) = 0.1$						

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

if  $O = 3 \ 1 \ 3$

$P(H \langle S \rangle) = 0.8$	$P(C \langle S \rangle) = 0.2$	$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$P(H H) = 0.6$	$P(C H) = 0.2$	$\langle S \rangle$	H	H	C	$\langle /S \rangle$	
$P(H C) = 0.3$	$P(C C) = 0.5$						
$P(\langle /S \rangle H) = 0.2$	$P(\langle /S \rangle C) = 0.2$						
$P(1 H) = 0.2$	$P(1 C) = 0.5$						
$P(2 H) = 0.4$	$P(2 C) = 0.4$						
$P(3 H) = 0.4$	$P(3 C) = 0.1$						

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle | H) = 0.2 & P(\langle /S \rangle | C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3 1 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle | H) = 0.2 & P(\langle /S \rangle | C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3 1 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3 1 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3 1 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200
$\langle S \rangle$	C	H	H	$\langle /S \rangle$	0.0000576

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3\ 1\ 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200
$\langle S \rangle$	C	H	H	$\langle /S \rangle$	0.0000576
$\langle S \rangle$	C	H	C	$\langle /S \rangle$	0.0000048

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \end{array}$$

if  $O = 3 \ 1 \ 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200
$\langle S \rangle$	C	H	H	$\langle /S \rangle$	0.0000576
$\langle S \rangle$	C	H	C	$\langle /S \rangle$	0.0000048
$\langle S \rangle$	C	C	H	$\langle /S \rangle$	0.0001200

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3\ 1\ 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200
$\langle S \rangle$	C	H	H	$\langle /S \rangle$	0.0000576
$\langle S \rangle$	C	H	C	$\langle /S \rangle$	0.0000048
$\langle S \rangle$	C	C	H	$\langle /S \rangle$	0.0001200
$\langle S \rangle$	C	C	C	$\langle /S \rangle$	0.0000500

# Decoding

## Task

Find the most likely state sequence  $\hat{S}$ , given an observation sequence  $O$ .

HMM

$$\begin{array}{ll} P(H|\langle S \rangle) = 0.8 & P(C|\langle S \rangle) = 0.2 \\ P(H|H) = 0.6 & P(C|H) = 0.2 \\ P(H|C) = 0.3 & P(C|C) = 0.5 \\ P(\langle /S \rangle|H) = 0.2 & P(\langle /S \rangle|C) = 0.2 \\ \\ P(1|H) = 0.2 & P(1|C) = 0.5 \\ P(2|H) = 0.4 & P(2|C) = 0.4 \\ P(3|H) = 0.4 & P(3|C) = 0.1 \end{array}$$

if  $O = 3\ 1\ 3$

$\langle S \rangle$	H	H	H	$\langle /S \rangle$	0.0018432
$\langle S \rangle$	H	H	C	$\langle /S \rangle$	0.0001536
$\langle S \rangle$	H	C	H	$\langle /S \rangle$	0.0007680
$\langle S \rangle$	H	C	C	$\langle /S \rangle$	0.0003200
$\langle S \rangle$	C	H	H	$\langle /S \rangle$	0.0000576
$\langle S \rangle$	C	H	C	$\langle /S \rangle$	0.0000048
$\langle S \rangle$	C	C	H	$\langle /S \rangle$	0.0001200
$\langle S \rangle$	C	C	C	$\langle /S \rangle$	0.0000500

# Dynamic Programming

For (only) two states and a (short) observation sequence of length three, comparing all possible sequences may be workable, but . . .

# Dynamic Programming

For (only) two states and a (short) observation sequence of length three, comparing all possible sequences may be workable, but ...

- ▶ for  $N$  observations and  $L$  states, there are  $L^N$  sequences;
- ▶ we end up doing the same partial calculations over and over again.

# Dynamic Programming

For (only) two states and a (short) observation sequence of length three, comparing all possible sequences may be workable, but ...

- ▶ for  $N$  observations and  $L$  states, there are  $L^N$  sequences;
- ▶ we end up doing the same partial calculations over and over again.

## Dynamic Programming:

- ▶ records sub-problem solutions for further re-use
- ▶ useful when a complex problem can be described recursively
- ▶ examples: Dijkstra's shortest path, minimum edit distance, longest common subsequence, Viterbi algorithm

# Dynamic Programming

For (only) two states and a (short) observation sequence of length three, comparing all possible sequences may be workable, but ...

- ▶ for  $N$  observations and  $L$  states, there are  $L^N$  sequences;
- ▶ we end up doing the same partial calculations over and over again.

## Dynamic Programming:

- ▶ records sub-problem solutions for further re-use
- ▶ useful when a complex problem can be described recursively
- ▶ examples: Dijkstra's shortest path, minimum edit distance, longest common subsequence, **Viterbi algorithm**

# Viterbi Algorithm

Recall our problem:

$$\text{maximize } P(s_1 \dots s_n | o_1 \dots o_n) = P(s_1 | s_0)P(o_1 | s_1)P(s_2 | s_1)P(o_2 | s_2) \dots$$

Our recursive sub-problem:

$$v_i(x) = \max_{k=1}^L [v_{i-1}(k) \cdot P(x|k) \cdot P(o_i|x)]$$

The variable  $v_i(x)$  represents the maximum probability that the  $i$ -th state is  $x$ , given that we have seen  $O_1^i$ .

# Viterbi Algorithm

Recall our problem:

$$\text{maximize } P(s_1 \dots s_n | o_1 \dots o_n) = P(s_1 | s_0) P(o_1 | s_1) P(s_2 | s_1) P(o_2 | s_2) \dots$$

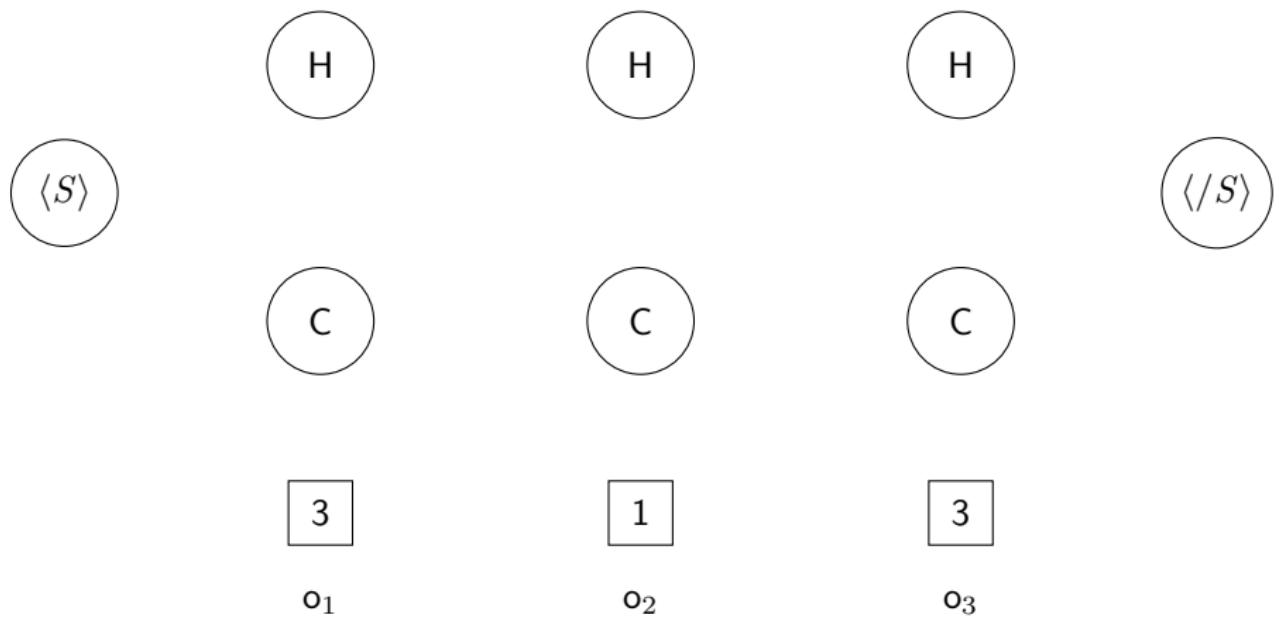
Our recursive sub-problem:

$$v_i(x) = \max_{k=1}^L [v_{i-1}(k) \cdot P(x|k) \cdot P(o_i|x)]$$

The variable  $v_i(x)$  represents the maximum probability that the  $i$ -th state is  $x$ , given that we have seen  $O_1^i$ .

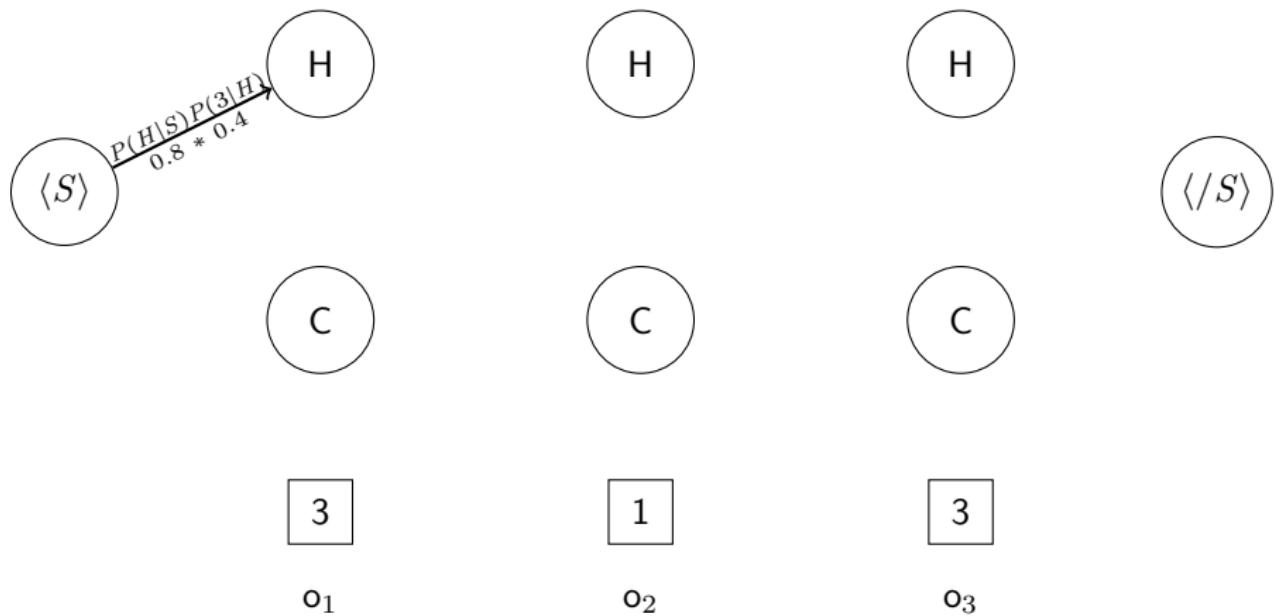
At each step, we record backpointers showing which previous state led to the maximum probability.

# An Example of the Viterbi Algorithmn

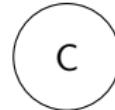
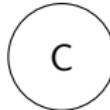
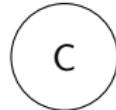
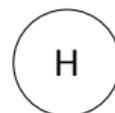
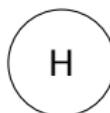
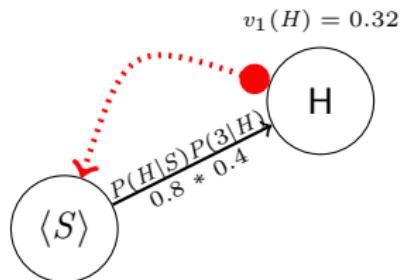


# An Example of the Viterbi Algorithmn

$$v_1(H) = 0.32$$



# An Example of the Viterbi Algorithmn



3

$o_1$

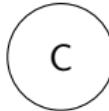
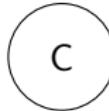
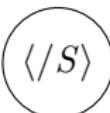
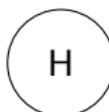
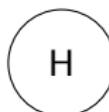
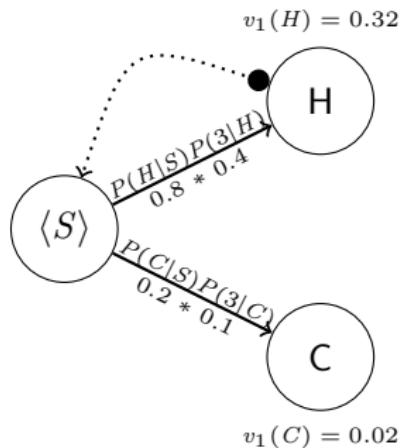
1

$o_2$

3

$o_3$

# An Example of the Viterbi Algorithmn



3

$o_1$

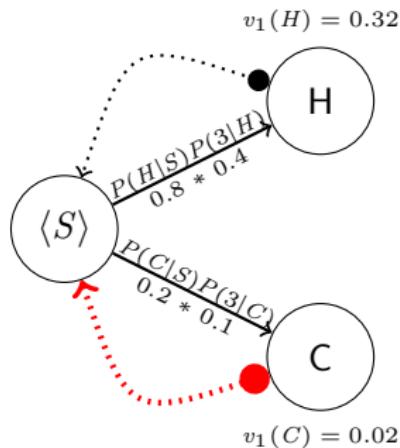
1

$o_2$

3

$o_3$

# An Example of the Viterbi Algorithmn



3

$o_1$

1

$o_2$

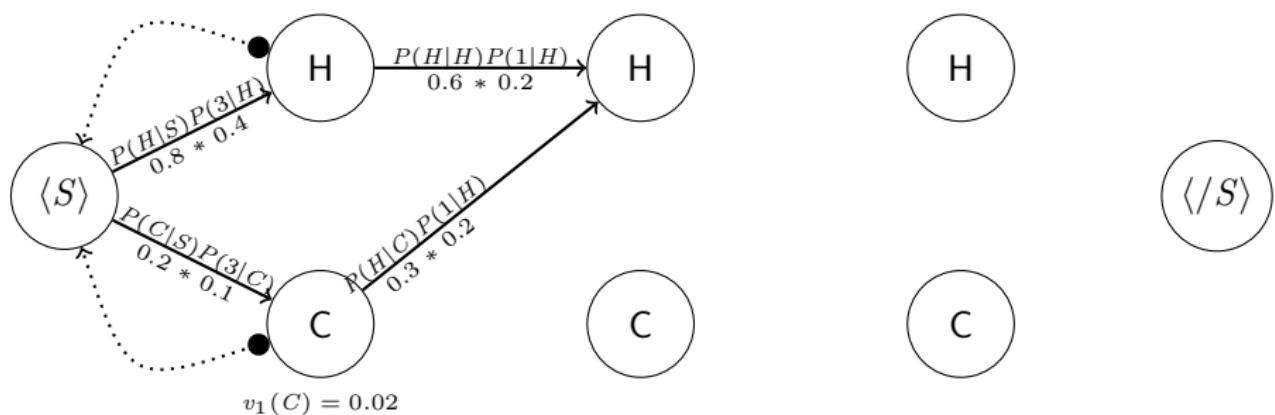
3

$o_3$

# An Example of the Viterbi Algorithmn

$$v_2(H) = \max(.32 * .12, .02 * .06)$$

$$= .0384$$



3

o<sub>1</sub>

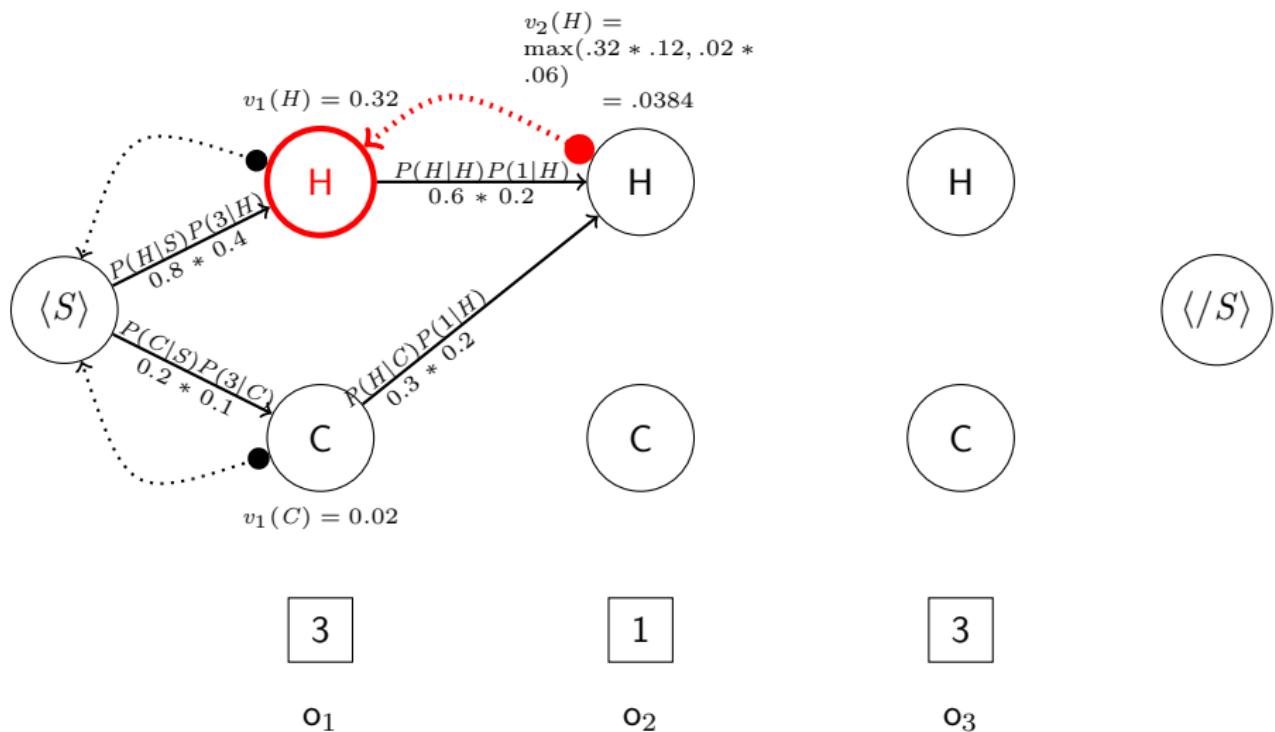
1

o<sub>2</sub>

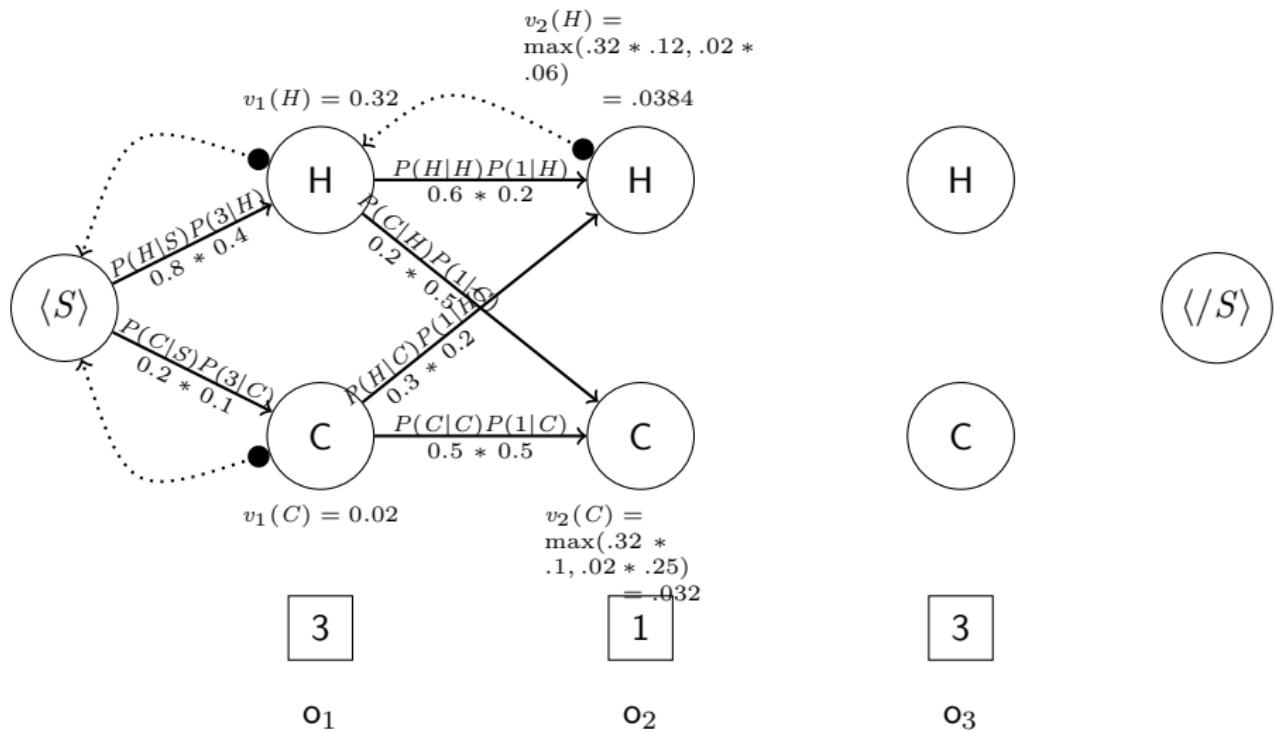
3

o<sub>3</sub>

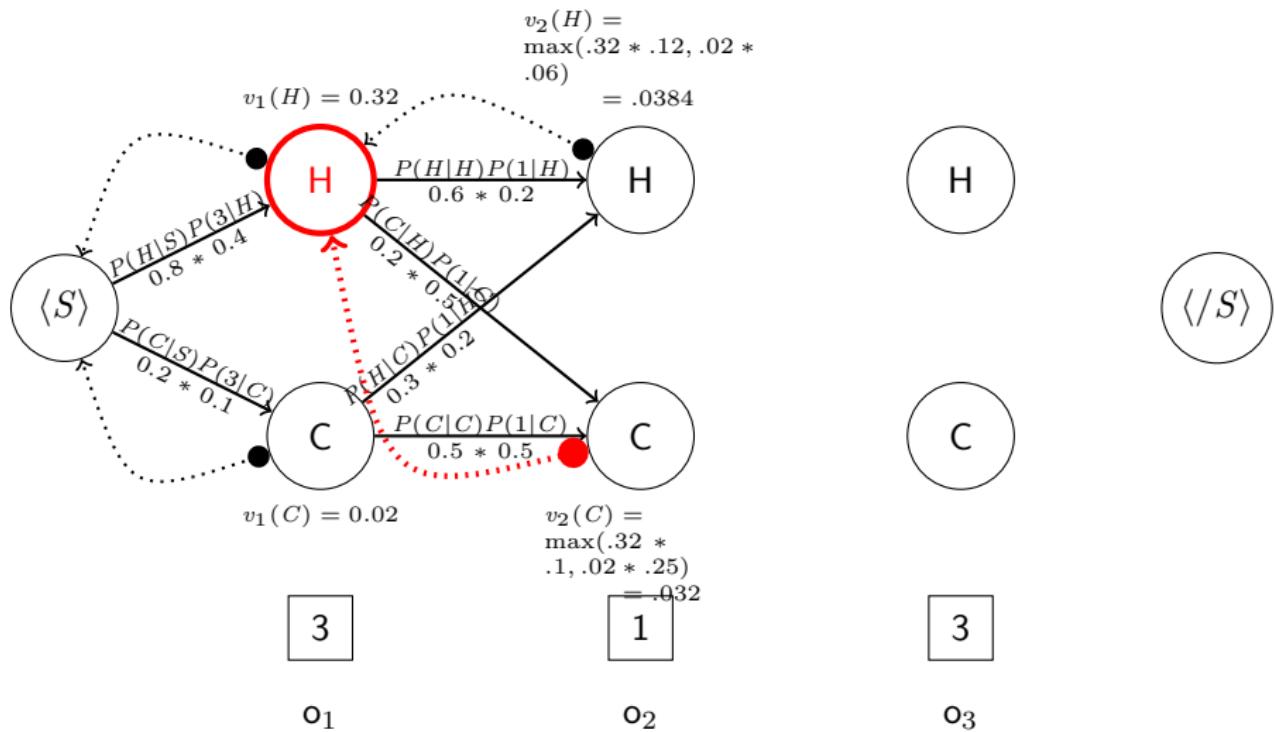
# An Example of the Viterbi Algorithmn



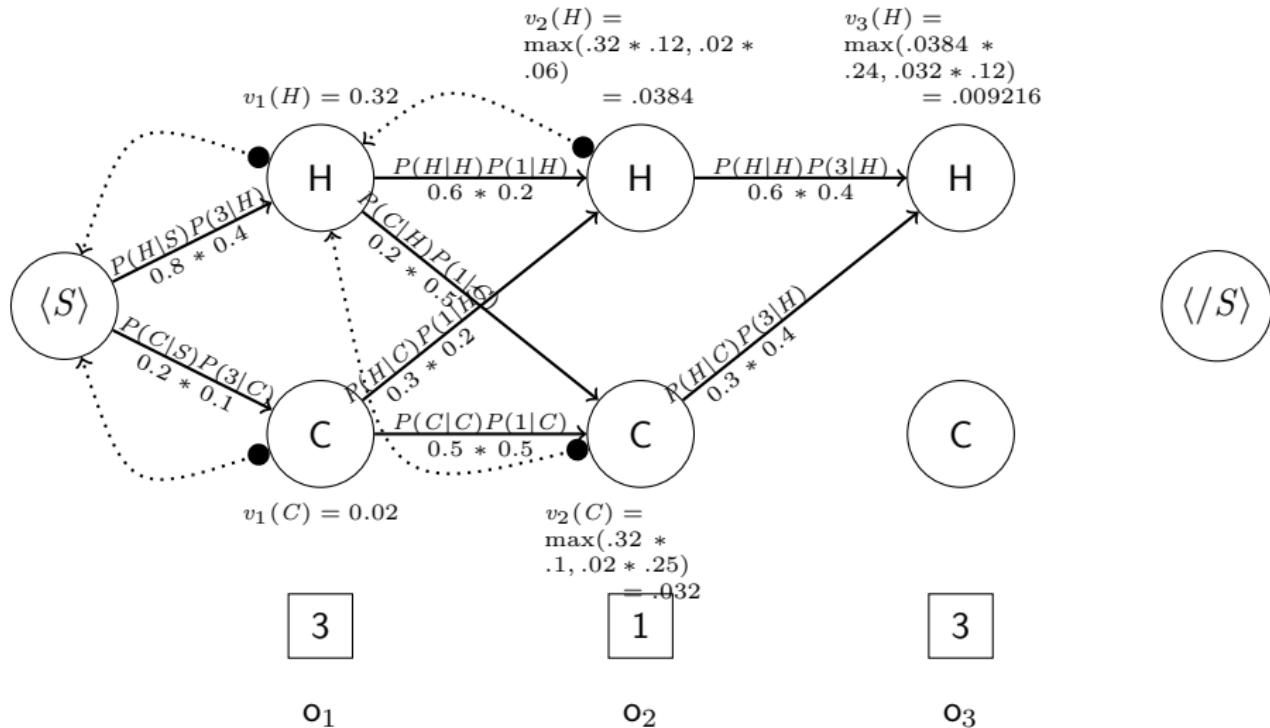
# An Example of the Viterbi Algorithmn



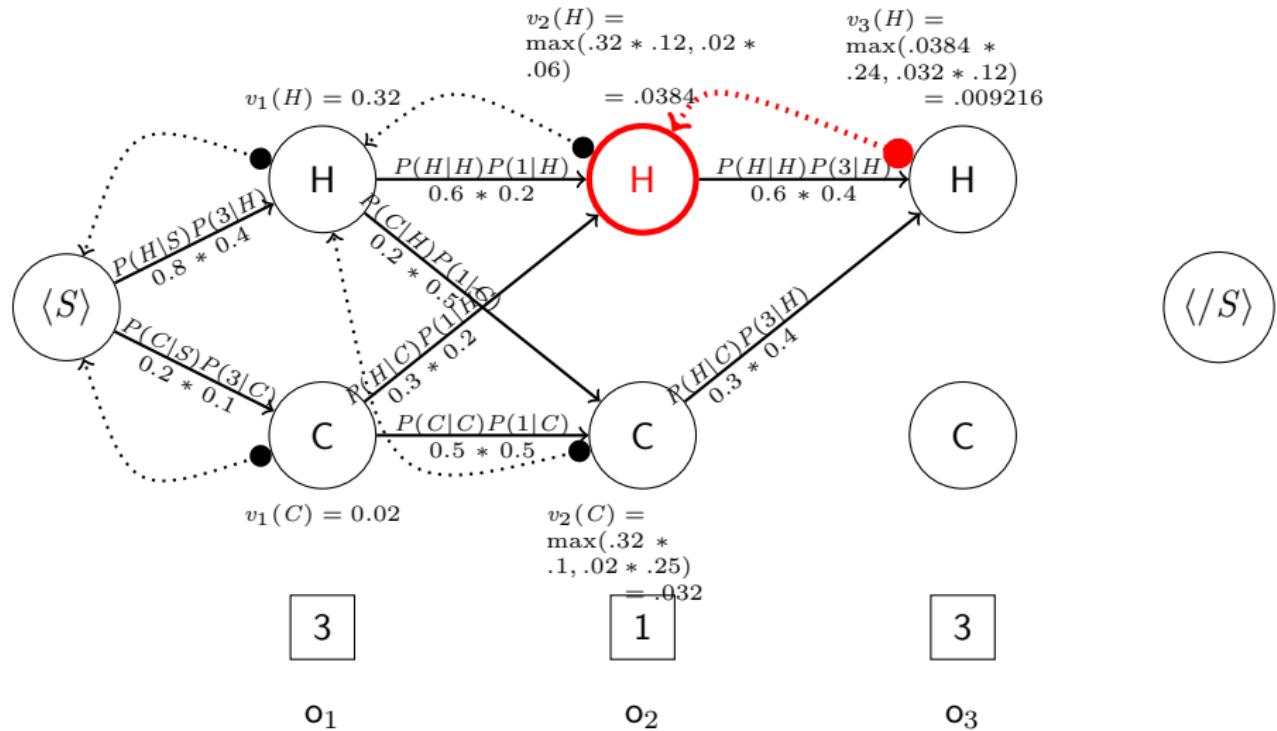
# An Example of the Viterbi Algorithm



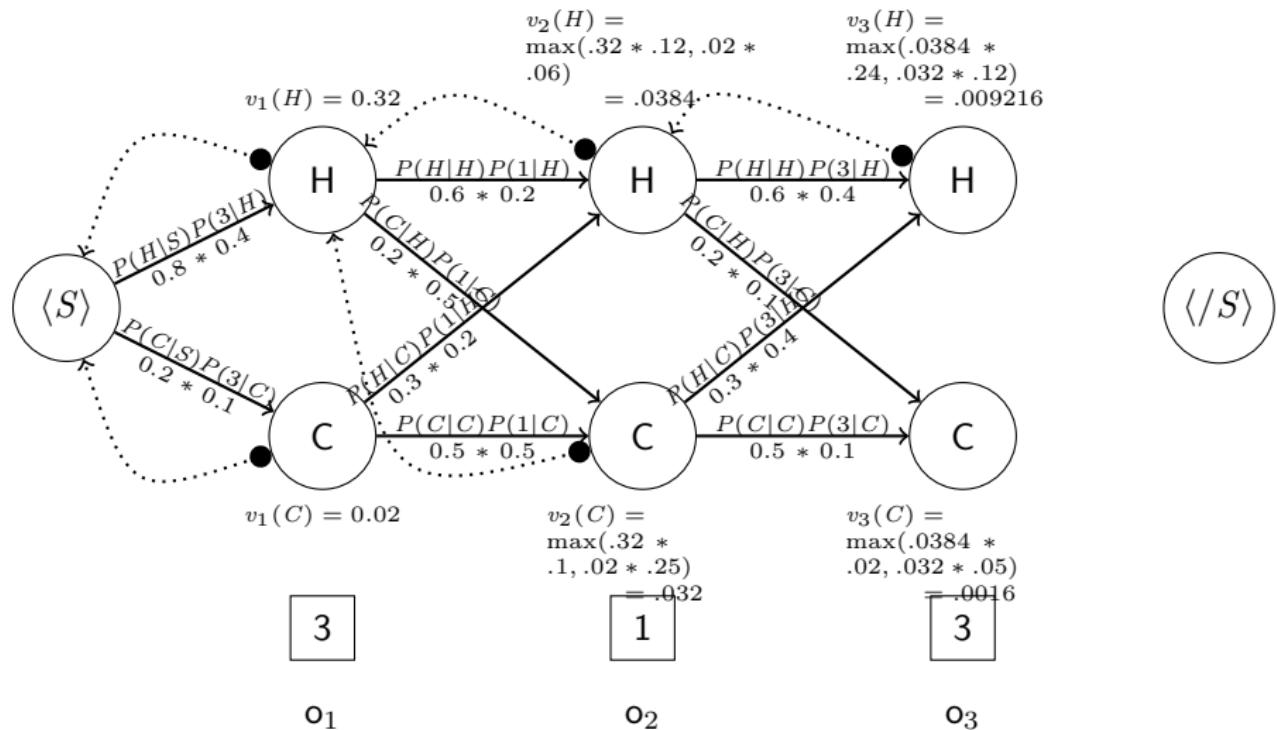
# An Example of the Viterbi Algorithmn



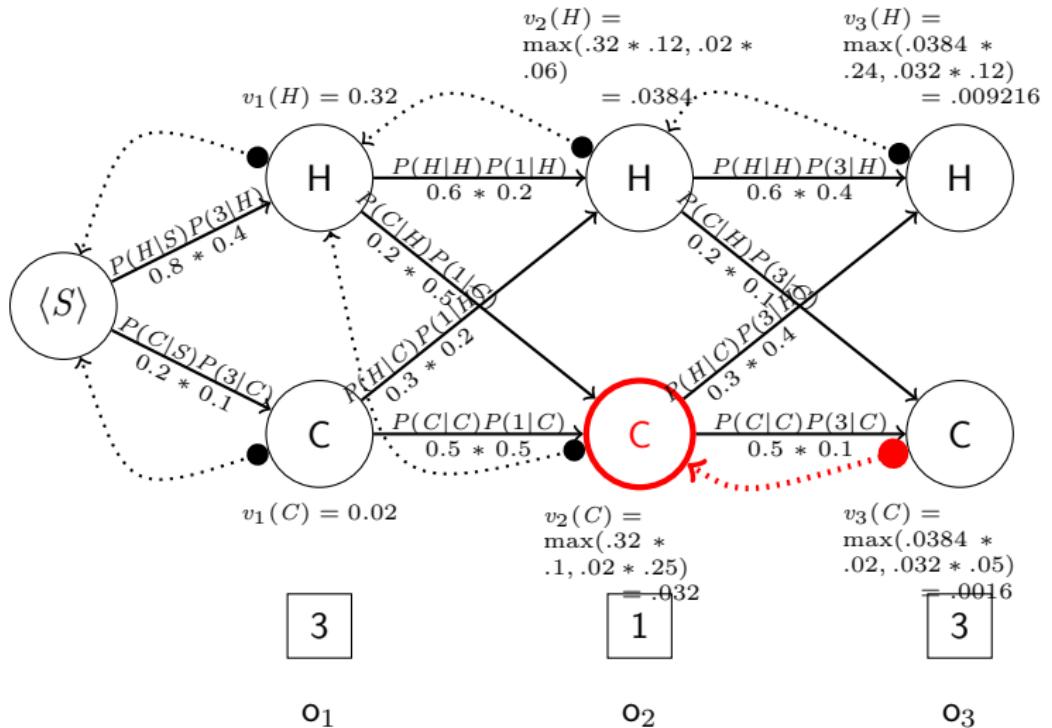
# An Example of the Viterbi Algorithm



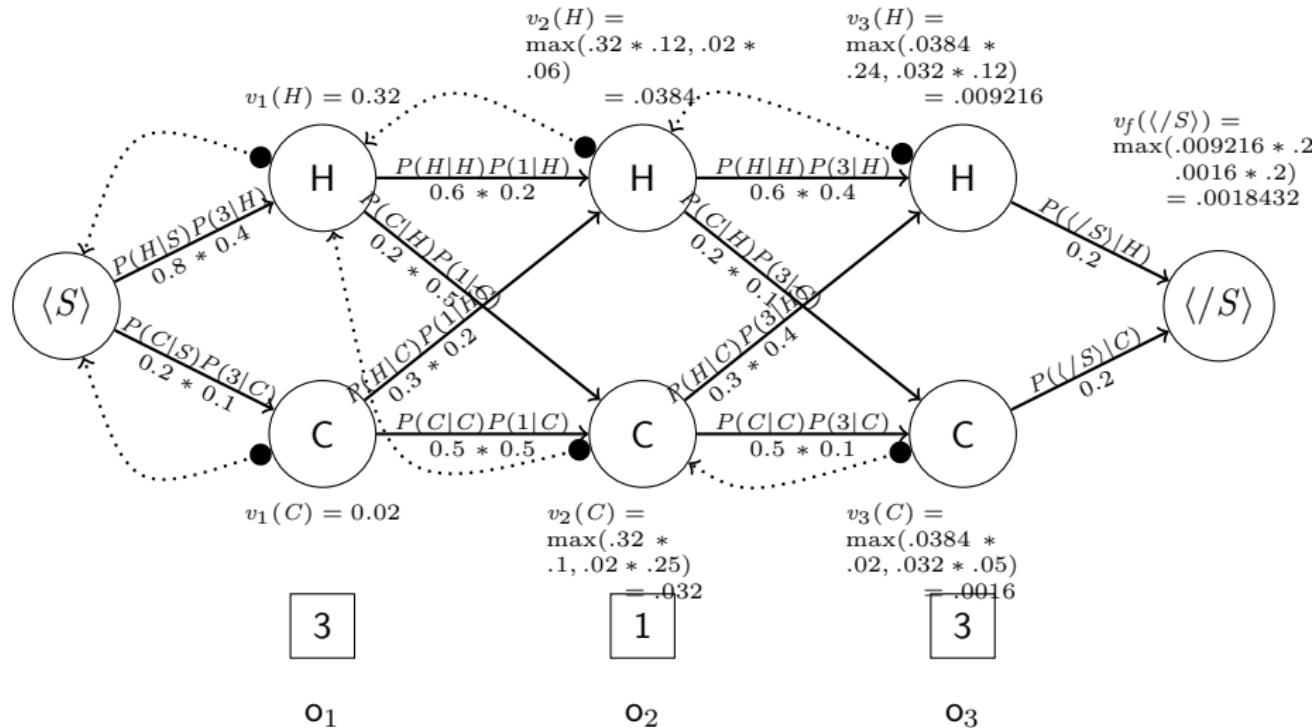
# An Example of the Viterbi Algorithmn



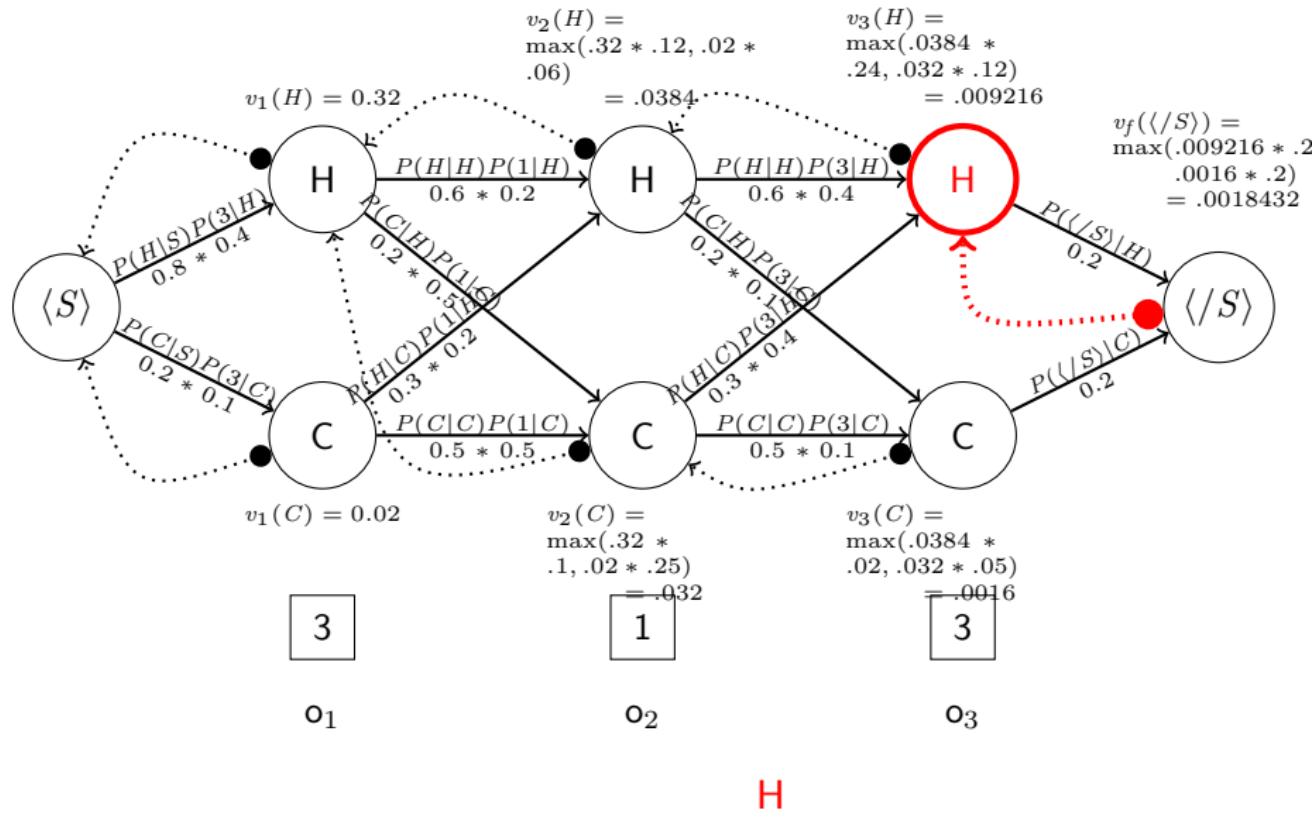
# An Example of the Viterbi Algorithm



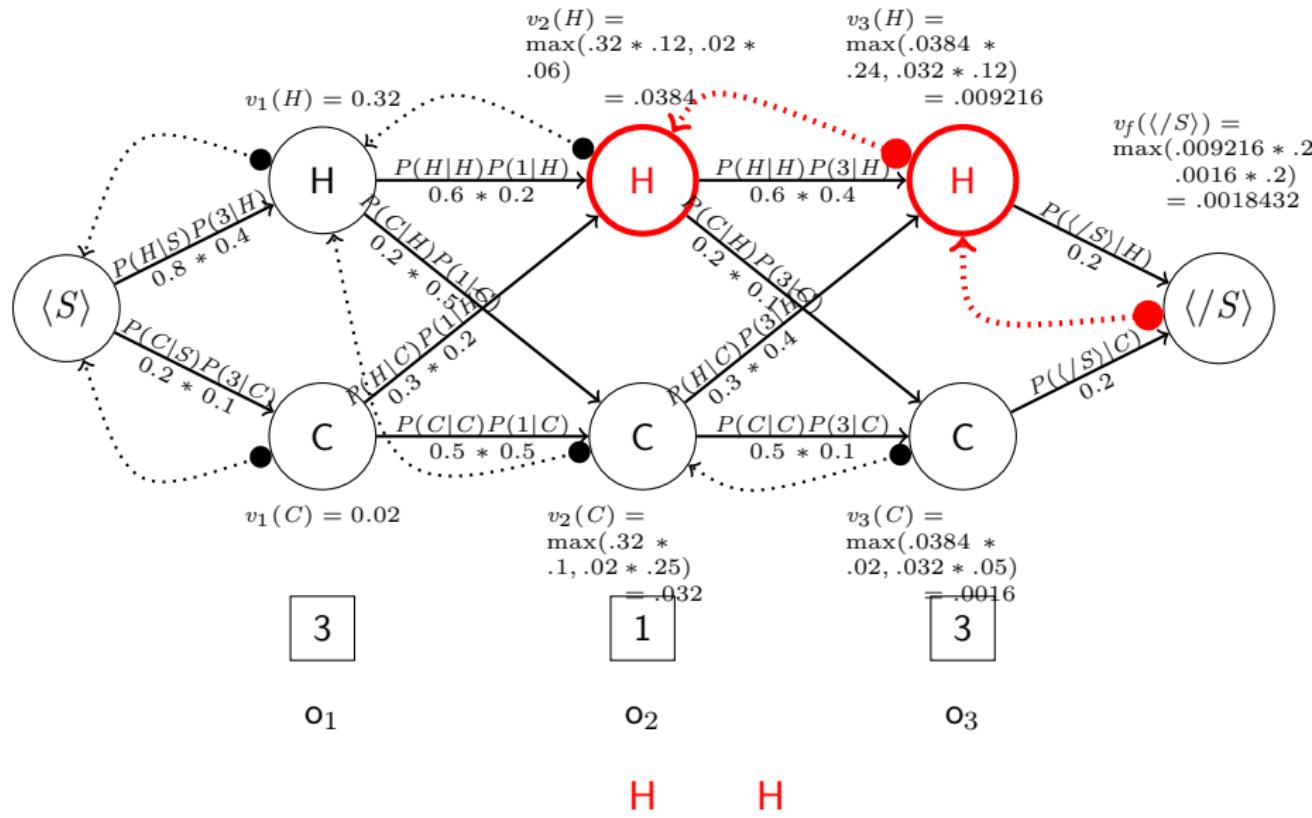
# An Example of the Viterbi Algorithmmn



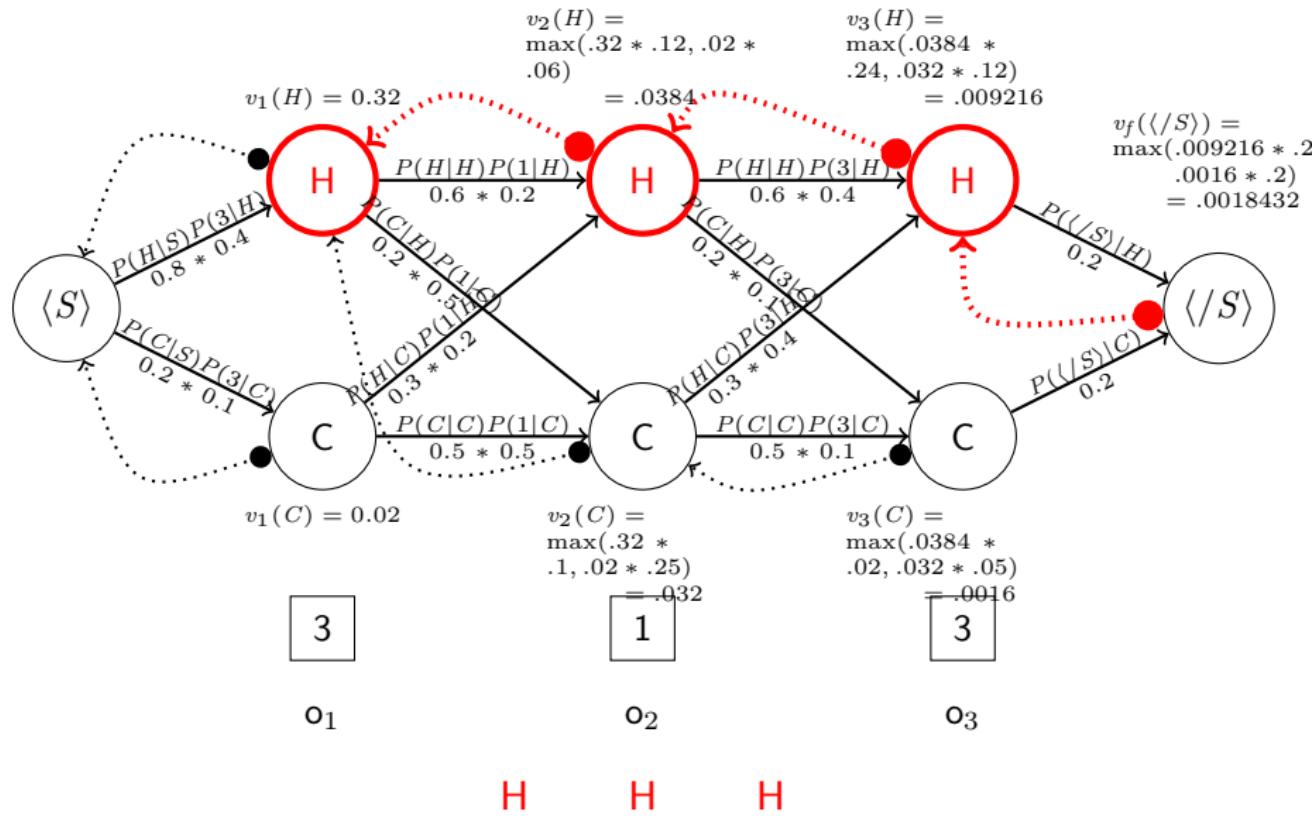
# An Example of the Viterbi Algorithmmn



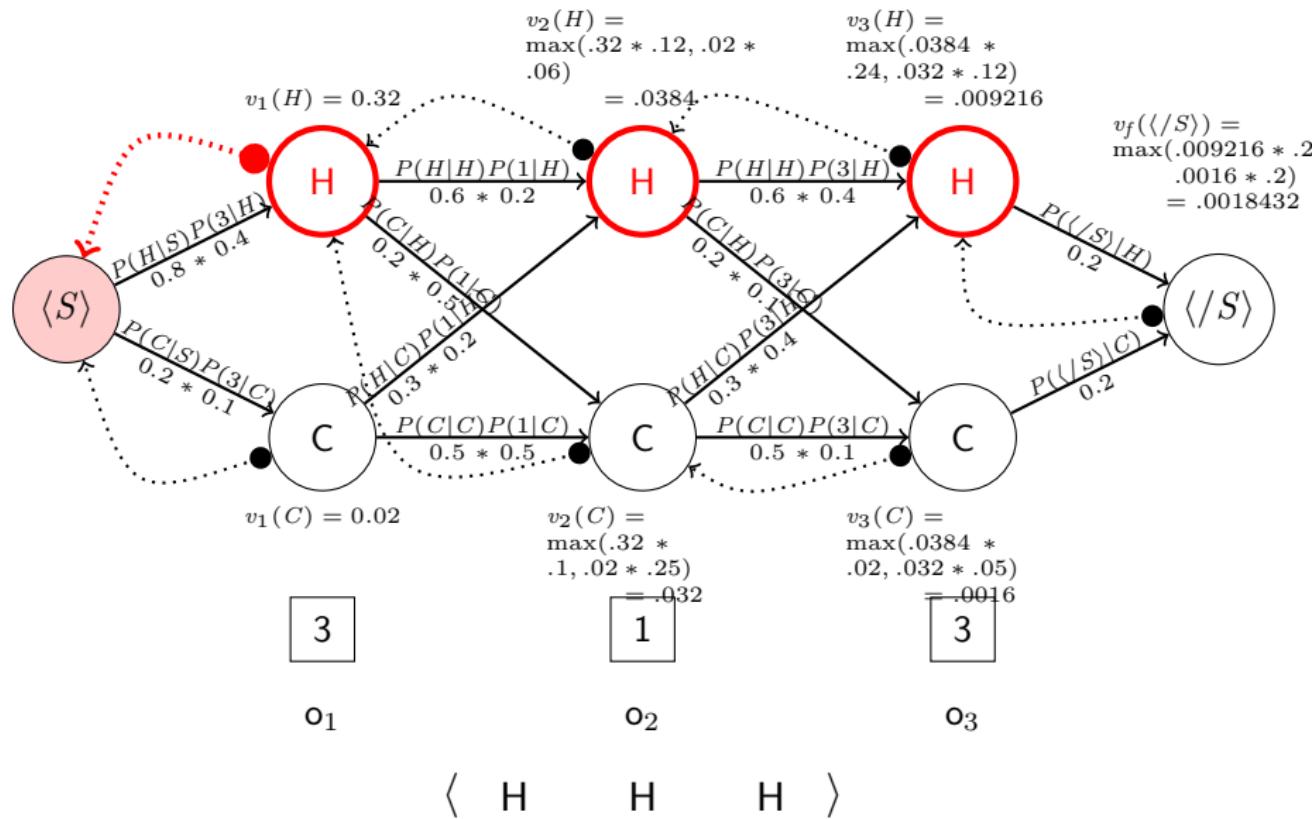
# An Example of the Viterbi Algorithmmn



# An Example of the Viterbi Algorithmmn



# An Example of the Viterbi Algorithmmn



# Pseudocode for the Viterbi Algorithm

**Input:** *observations* of length  $N$ , state set of size  $L$

**Output:** *best-path*

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, s] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, s] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

## Diversion: Complexity and O(N)

Big-O notation describes the complexity of an algorithm.

- ▶ it describes the worst-case *order of growth* in terms of the size of the input
- ▶ only the largest order term is represented
- ▶ constant factors are ignored
- ▶ determined by looking at loops in the code

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$L$

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$L$

L

N

31

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$L$

L

N

L

31

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$L$

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$$L + L^2N$$

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$$L + L^2N + N$$

# Pseudocode for the Viterbi Algorithm

**Input:** observations of length  $N$ , state set of length  $L$

**Output:** best-path

create a path probability matrix  $viterbi[N, L + 1]$

create a path backpointer matrix  $backpointer[N, L + 1]$

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[1, s] \leftarrow trans(\langle S \rangle, s) \times emit(o_1, s)$

$backpointer[1, s] \leftarrow 0$

**end**

**foreach** time step  $i$  from 2 to  $N$  **do**

**foreach** state  $s$  from 1 to  $L$  **do**

$viterbi[i, s] \leftarrow \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s) \times emit(o_i, s)$

$backpointer[i, s] \leftarrow \arg \max_{s'=1}^L viterbi[i - 1, s'] \times trans(s', s)$

**end**

**end**

$viterbi[N, L + 1] \leftarrow \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

$backpointer[N, L + 1] \leftarrow \arg \max_{s=1}^L viterbi[N, 1] \times trans(s, \langle /S \rangle)$

**return** the path by following backpointers from  $backpointer[N, L + 1]$

$$O(L^2 N)$$

L

N

L

N