

CS 481

***Artificial Intelligence Language
Understanding***

February 9, 2023

Announcements / Reminders

- Please follow the Week 05 To Do List instructions
 - Quiz #04 due on Sunday (02/12/23) at 11:59 PM CST
 - PA #01 due on Monday (02/20/23) at 11:59 PM CST
-
- Exam dates:
 - Midterm: 03/02/2023 during Thursday lecture time
 - Final: 04/27/2023 during Thursday lecture time

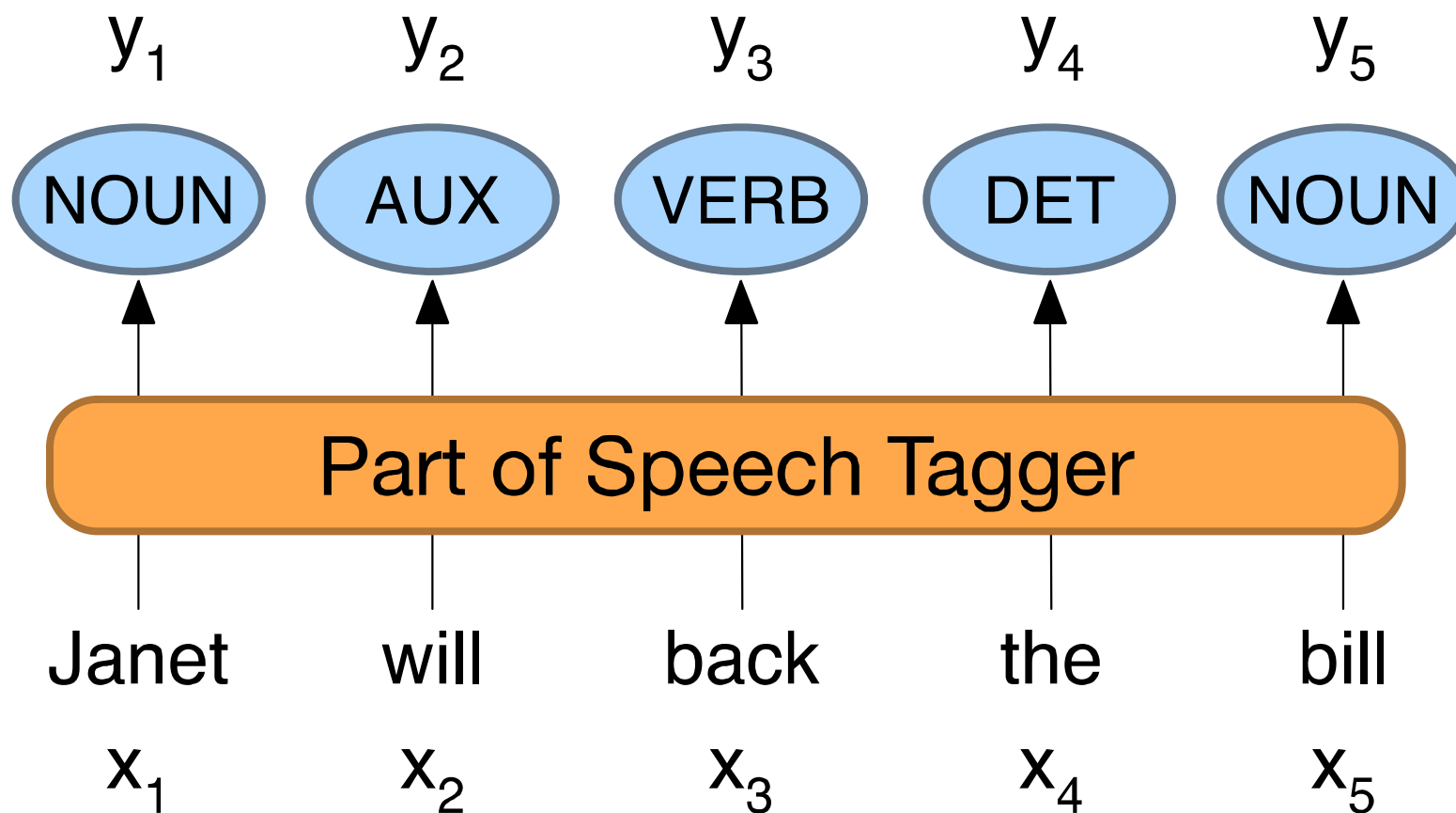
Plan for Today

- **Parts of Speech tagging - continued**

Part of Speech Tagging

- Task:

- Map sequence x_1, \dots, x_n of words to y_1, \dots, y_n of POS tags



Parts of Speech: Tagset Example

Parts of Speech in the Universal Dependencies tagset

	Tag	Description	Example
Open Class	ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by, under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, first, second</i>
	PART	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
	SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
Other	PUNCT	Punctuation	<i>; , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

Parts of Speech: Tagset Example

Penn Treebank Parts-of-speech tags:

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

Sample Tagged Sentence

There/**PRO** were/**VERB** 70/**NUM** children/**NOUN**
there/**ADV** ./**PUNC**

Preliminary/**ADJ** findings/**NOUN** were/**AUX**
reported/**VERB** in/**ADP** today/**NOUN** 's/**PART**
New/**PROPN** England/**PROPN** Journal/**PROPN**
of/**ADP** Medicine/**PROPN**

Standard POS Tagging Models

- **Supervised Machine Learning Algorithms:**
 - **Hidden Markov Models**
 - Conditional Random Fields (CRF) / Maximum Entropy Markov Models (MEMM)
 - Neural sequence models (RNNs or Transformers)
 - Large Language Models (like BERT), finetuned
- **All required a hand-labeled training set, all about equal performance (97% on English)**
 - All make use of information sources we discussed
 - Via human created features: HMMs and CRFs
 - Via representation learning: Neural LMs

Part of Speech: Conditional Probability

$$P(\textit{Category} = \textit{NOUN} \mid \textit{word} = \textit{flies}) = \frac{P(\textit{word} = \textit{flies}, \textit{Category} = \textit{NOUN})}{P(\textit{Word} = \textit{flies})}$$

where $P(\textit{Word} = \textit{flies}) > 0$

Part of Speech: Conditional Probability

$$P(C = \text{NOUN} \mid w = \text{flies}) = \frac{P(w = \text{flies}, C = \text{NOUN})}{P(w = \text{flies})}$$

where $P(w = \text{flies}) > 0$

Most Frequent Class Tagging

With:

$$P(w = \textit{flies}) = 1000 / 1\,273\,000 = 0.0008$$

$$P(w = \textit{flies}, C = \textit{NOUN}) = 400 / 1\,273\,000 = 0.0003$$

$$P(w = \textit{flies}, C = \textit{VERB}) = 600 / 1\,273\,000 = 0.0005$$

$$P(C = \textit{NOUN} \mid w = \textit{flies}) = \frac{P(w = \textit{flies}, C = \textit{NOUN})}{P(w = \textit{flies})} = \frac{0.0003}{0.0008}$$

vs.

$$P(C = \textit{VERB} \mid w = \textit{flies}) = \frac{P(w = \textit{flies}, C = \textit{VERB})}{P(w = \textit{flies})} = \frac{0.0005}{0.0008}$$

With this approach *flies* will ALWAYS be tagged as a **VERB**.

Bayes' Rule

$$P(A | B) = \frac{P(B | A) * P(A)}{P(B)}$$

POS Tagging: General Approach

Given a sequence of **words** (a “sentence”):

$$W_1, W_2, W_3, \dots, W_T$$

there is going to be a corresponding sequence of lexical categories:

$$C_1, C_2, C_3, \dots, C_T$$

What is most likely sequence of categories?

POS Tagging: General Approach

To answer this we would want to find a conditional probability:

$$P(C_1, C_2, C_3, \dots, C_T \mid w_1, w_2, w_3, \dots, w_T)$$

In other words: what is the probability of having a sequence of lexical categories

$$C_1, C_2, C_3, \dots, C_T$$

GIVEN that the sequence of words is

$$w_1, w_2, w_3, \dots, w_T ?$$

POS Tagging: General Approach

The probability we are looking for

$$P(C_1, C_2, C_3, \dots, C_T \mid w_1, w_2, w_3, \dots, w_T)$$

will require a lot of data, which we most likely won't have.

We can use Bayes' Theorem:

$$P(C_1, C_2, C_3, \dots, C_T \mid w_1, w_2, w_3, \dots, w_T) =$$

$$= \frac{P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * P(C_1, C_2, C_3, \dots, C_T)}{P(w_1, w_2, w_3, \dots, w_T)}$$

POS Tagging: General Approach

In order to find the most likely sequence:

$$C_1, C_2, C_3, \dots, C_T$$

we need to **maximize** (most likely sequence!):

$$\begin{aligned} &P(C_1, C_2, C_3, \dots, C_T \mid w_1, w_2, w_3, \dots, w_T) = \\ &= \frac{P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * (C_1, C_2, C_3, \dots, C_T)}{P(w_1, w_2, w_3, \dots, w_T)} \end{aligned}$$

POS Tagging: General Approach

Maximizing :

$$P(C_1, C_2, C_3, \dots, C_T \mid w_1, w_2, w_3, \dots, w_T)$$

in practice means **maximizing the numerator:**

$$\frac{P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * P(C_1, C_2, C_3, \dots, C_T)}{P(w_1, w_2, w_3, \dots, w_T)}$$

as denominator $P(w_1, w_2, w_3, \dots, w_T)$ will not change:

POS Tagging: General Approach

Estimating:

$$P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * P(C_1, C_2, C_3, \dots, C_T)$$

using counts once again requires a lot of data that we will likely not have.

Alternative: approximate it with N-grams (here bigrams):

$$P(C_1, C_2, C_3, \dots, C_T) = \prod_{i=1}^T P(\textcolor{red}{C}_i \mid \textit{all categories preceding } \textcolor{red}{C}_i)$$

$$P(C_1, C_2, C_3, \dots, C_T) \cong \prod_{i=1}^T P(\textcolor{red}{C}_i \mid \textcolor{red}{C}_{i-})$$

POS Tagging: General Approach

Estimating:

$$P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * P(C_1, C_2, C_3, \dots, C_T)$$

Approximate it with N-grams (here bigrams):

$$P(C_1, C_2, C_3, \dots, C_T) = \prod_{i=1}^T P(\textcolor{red}{C}_i \mid \textit{all categories preceding } \textcolor{red}{C}_i)$$

$$P(C_1, C_2, C_3, \dots, C_T) \cong \prod_{i=1}^T P(\textcolor{red}{C}_i \mid \textcolor{red}{C}_{i-1})$$

$$P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) \cong \prod_{i=1}^T P(\textcolor{blue}{w}_i \mid \textcolor{red}{C}_i)$$

POS Tagging: General Approach

With approximations:

$$P(w_1, w_2, w_3, \dots, w_T \mid C_1, C_2, C_3, \dots, C_T) * P(C_1, C_2, C_3, \dots, C_T) \cong \\ \cong \prod_{i=1}^T P(w_i \mid C_i) * P(C_i \mid C_{i-1})$$

and we want to **maximize**:

$$\prod_{i=1}^T P(w_i \mid C_i) * P(C_i \mid C_{i-1})$$

Individual probabilities can now be estimated using corpus counts!

POS Tagging: Simple Tagset

Let's assume we have a simple tagset:

- **N - NOUN**
- **V - VERB**
- **ART - ARTICLE**
- **P - PREPOSITION**

and a some synthetic corpus.

POS Tagging: General Approach

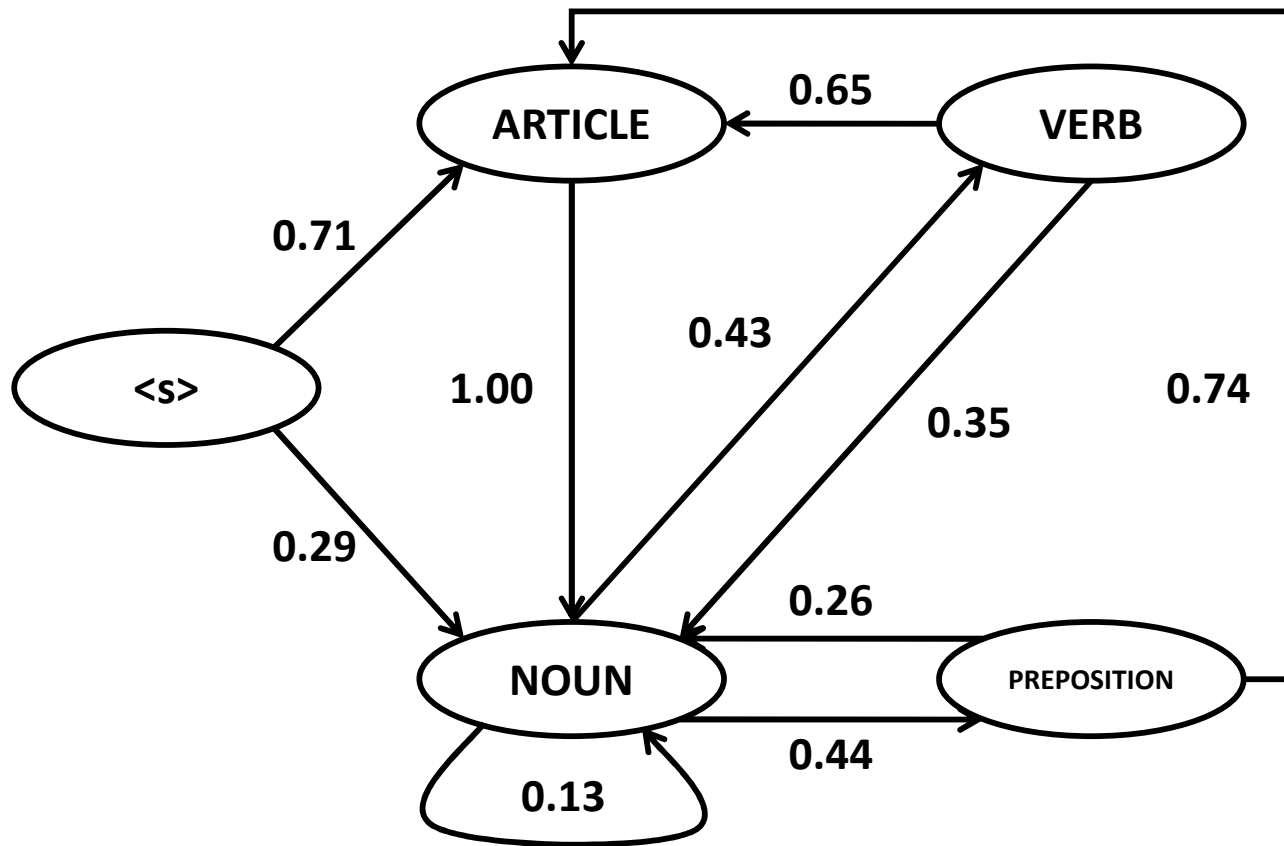
Estimations with corpus counts:

$$P(C_i = \text{VERB} \mid C_{i-1} = \text{NOUN}) = \frac{\text{Count}(\text{NOUN at position } i- \text{ and VERB at } i)}{\text{Count}(\text{NOUN at position } i-)}$$

Sample bigram probabilities from our synthetic corpus:

Category	Count at i	Pair	Count at i,i+1	P(Bigram)	Estimate
<s>	300	<s>, ARTICLE	213	P (ARTICLE <S>)	0.71
<s>	300	<s>, NOUN	87	P (NOUN <S>)	0.29
ARTICLE	558	ARTICLE, NOUN	558	P (NOUN ARTICLE)	1.00
NOUN	833	NOUN, VERB	358	P (VERB NOUN)	0.43
NOUN	833	NOUN, NOUN	108	P (NOUN NOUN)	0.13
NOUN	833	NOUN, PREPOSITION	366	P (PREPOSITION NOUN)	0.44
VERB	300	VERB, NOUN	75	P (NOUN VERB)	0.35
VERB	300	VERB, ARTICLE	194	P (ARTICLE VERB)	0.65
PREPOSITION	307	PREPOSITION, ARTICLE	226	P (ARTICLE PREPOSITION)	0.74
PREPOSITION	307	PREPOSITION, NOUN	81	P (NOUN PREPOSITION)	0.26

Hidden Markov Model



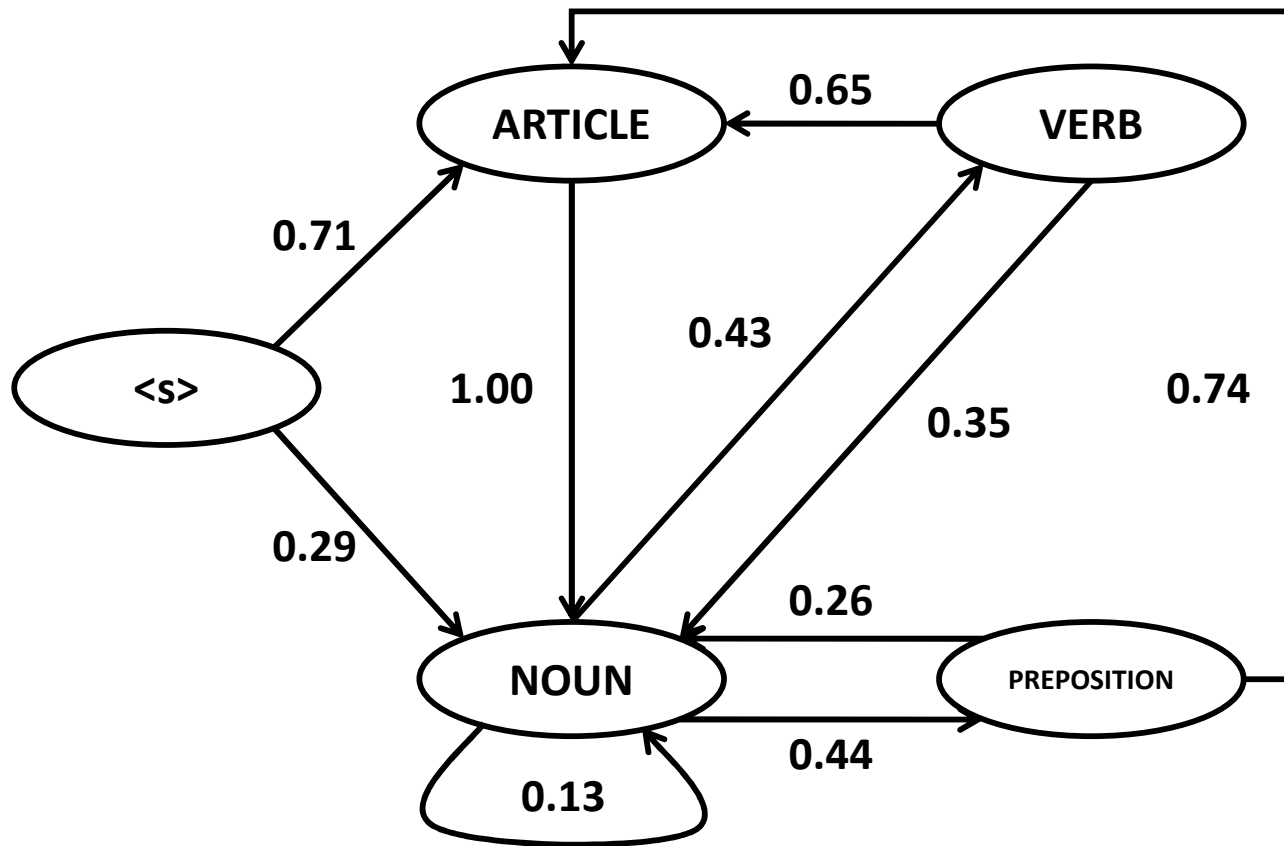
P(Bigram)	Estimate
P (ARTICLE $\langle s \rangle$)	0.71
P (NOUN $\langle s \rangle$)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.44
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

Consider a following sequence of categories (tags):

$\langle s \rangle$, ARTICLE, NOUN, VERB, NOUN

What's the probability of its occurrence in our synthetic corpus?

Hidden Markov Model (HMM)

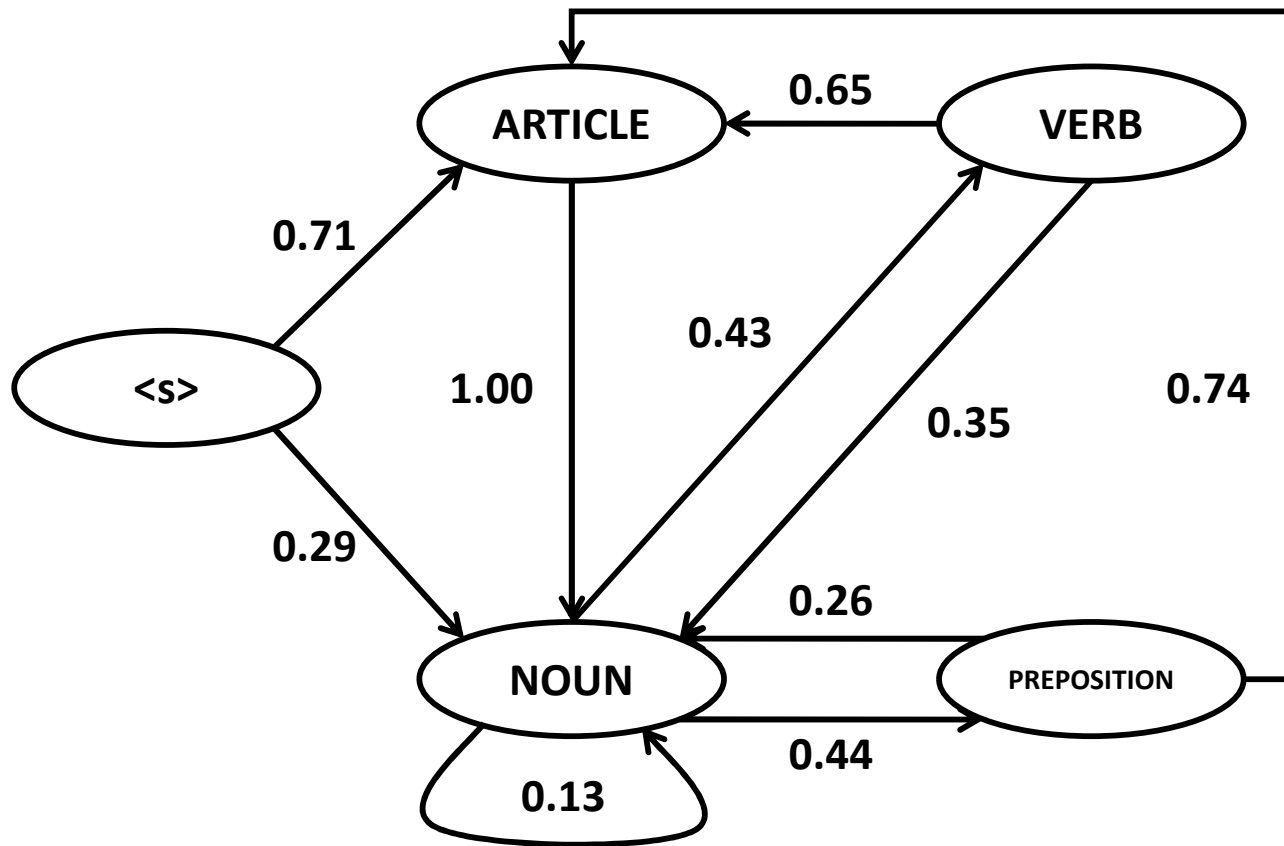


P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.44
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

The word “Hidden” in Hidden Markov Model means that for a specific sequence (of words) it is unclear what state the model is in.

The word *flies* could be generated from state NOUN and state VERB.

Hidden Markov Model

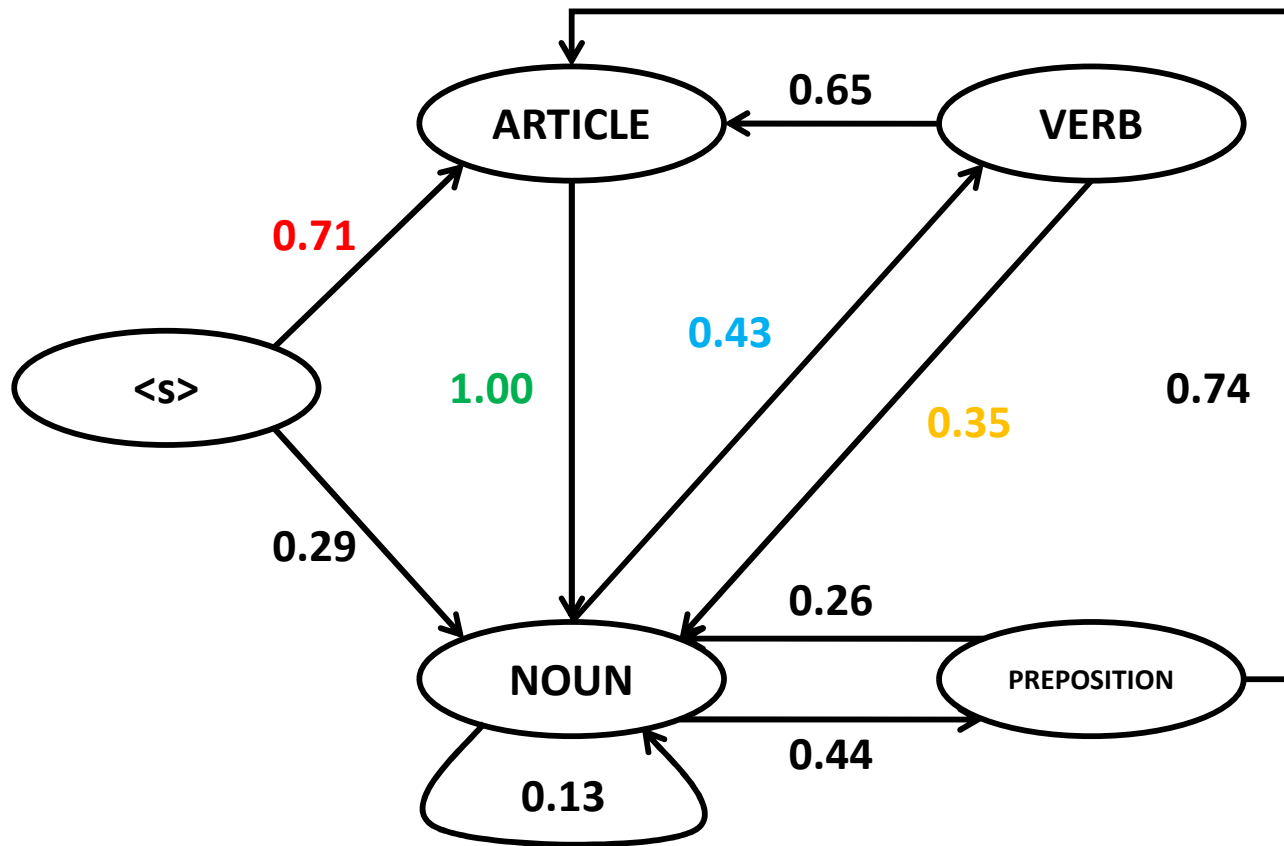


P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.26
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

Probability of occurrence of a sequence of categories (tags):

$$P(C_1, C_2, C_3, \dots, C_T) \cong \prod_{i=1}^T P(C_i | C_{i-1})$$

Hidden Markov Model



P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.44
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

Probability of occurrence of a sequence of categories (tags):

$$\begin{aligned}
 &P(<s>, \text{ARTICLE}, \text{NOUN}, \text{VERB}, \text{NOUN}) = \\
 &\cong P(\text{ART} | <s>) * P(\text{N} | \text{ART}) * P(\text{V} | \text{N}) * P(\text{N} | \text{V}) = 0.71 * 1.00 * 0.43 * 0.35 = 0.107
 \end{aligned}$$

Synthetic Corpus: Word/Tag Counts

Summary of selected word counts in the synthetic corpus:

Word/Tag	N	V	ART	P	TOTAL
<i>flies</i>	21	23	0	0	44
<i>fruit</i>	49	5	1	0	55
<i>like</i>	10	30	0	21	61
<i>a</i>	1	0	201	0	202
<i>the</i>	1	0	300	2	303
<i>flower</i>	53	15	0	0	68
<i>flowers</i>	42	16	0	0	58
<i>birds</i>	64	1	0	0	65
others	592	210	56	284	1142
TOTAL	833	300	558	307	1998

From the table we can calculate lexical generation probabilities $P(w|C)$ estimates:

$$P(\textit{the}|\text{ART}) = 300/558 = 0.54$$

$$P(a|\text{ART}) = 201/558 = 0.36$$

$$P(\textit{flies}|\text{N}) = 21/833 = 0.025$$

$$P(a|\text{N}) = 1/833 = 0.001$$

$$P(\textit{flies}|\text{V}) = 23/300 = 0.076$$

$$P(\textit{flower}|\text{N}) = 53/833 = 0.063$$

$$P(\textit{like}|\text{V}) = 30/300 = 0.1$$

$$P(\textit{flower}|\text{V}) = 15/300 = 0.05$$

$$P(\textit{like}|\text{P}) = 21/307 = 0.068$$

$$P(\textit{like}|\text{N}) = 10/833 = 0.012$$

Example

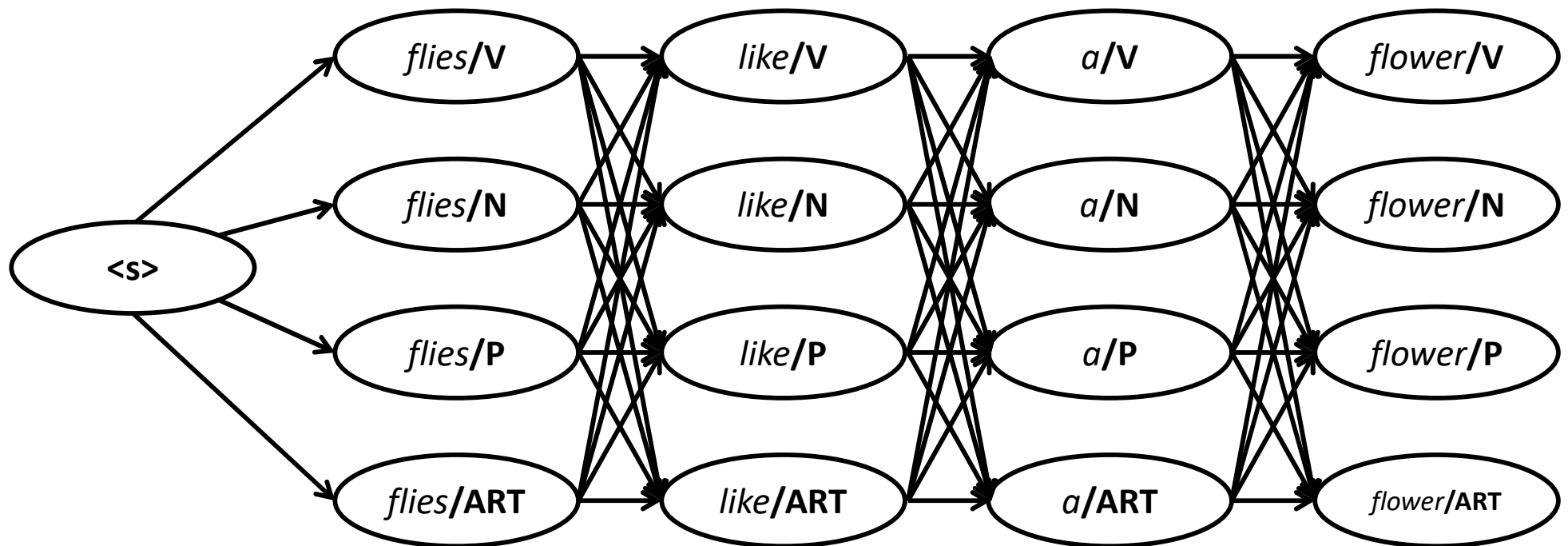
Given our synthetic corpus, what is the most like sequence of categories (tags) corresponding to a sentence:

Flies like a flower

We need to **maximize**:

$$\begin{aligned} P(w_1, w_2, w_3, \dots, w_T \mid c_1, c_2, c_3, \dots, c_T) * P(c_1, c_2, c_3, \dots, c_T) &\cong \\ &\cong \prod_{i=1}^T P(w_i \mid c_i) * P(c_i \mid c_{i-1}) \end{aligned}$$

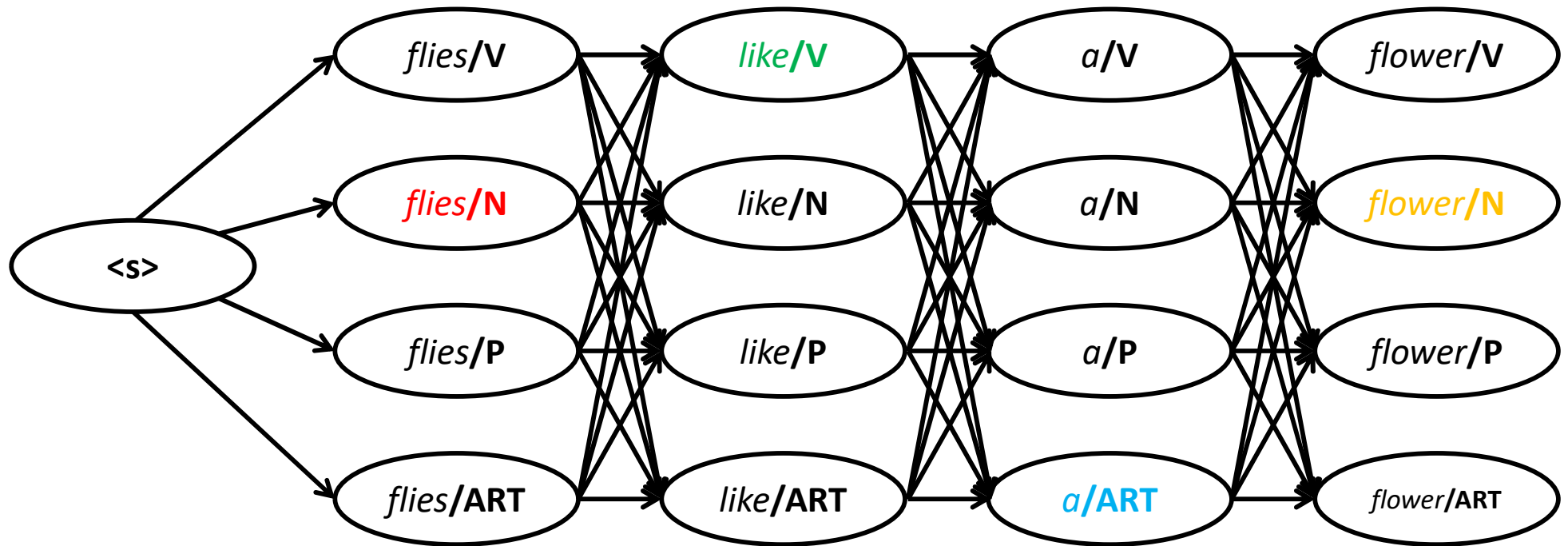
Example: All Possible Sequences



Every sequence can be assigned a probability:

$$P(w_1, w_2, w_3, \dots, w_T \mid c_1, c_2, c_3, \dots, c_T) \cong \prod_{i=1}^T P(w_i \mid c_i)$$

Example: All Possible Sequences



Every sequence can be assigned a probability:

$$\prod_{i=1}^T P(w_i | C_i) = P(\textit{flies}|\textit{N}) * P(\textit{like}|\textit{V}) * P(\textit{a}|\textit{ART}) * P(\textit{flower}|\textit{N})$$

Synthetic Corpus: Word/Tag Counts

Summary of selected word counts in the synthetic corpus:

Word/Tag	N	V	ART	P	TOTAL
<i>flies</i>	21	23	0	0	44
<i>fruit</i>	49	5	1	0	55
<i>like</i>	10	30	0	21	61
<i>a</i>	1	0	201	0	202
<i>the</i>	1	0	300	2	303
<i>flower</i>	53	15	0	0	68
<i>flowers</i>	42	16	0	0	58
<i>birds</i>	64	1	0	0	65
others	592	210	56	284	1142
TOTAL	833	300	558	307	1998

From the table we can calculate lexical generation probabilities $P(w|C)$ estimates:

$$P(\textit{the}|\text{ART}) = 300/558 = 0.54$$

$$P(a|\text{ART}) = 201/558 = 0.36$$

$$P(\textit{flies}|\text{N}) = 21/833 = 0.025$$

$$P(a|\text{N}) = 1/833 = 0.001$$

$$P(\textit{flies}|\text{V}) = 23/300 = 0.076$$

$$P(\textit{flower}|\text{N}) = 53/833 = 0.063$$

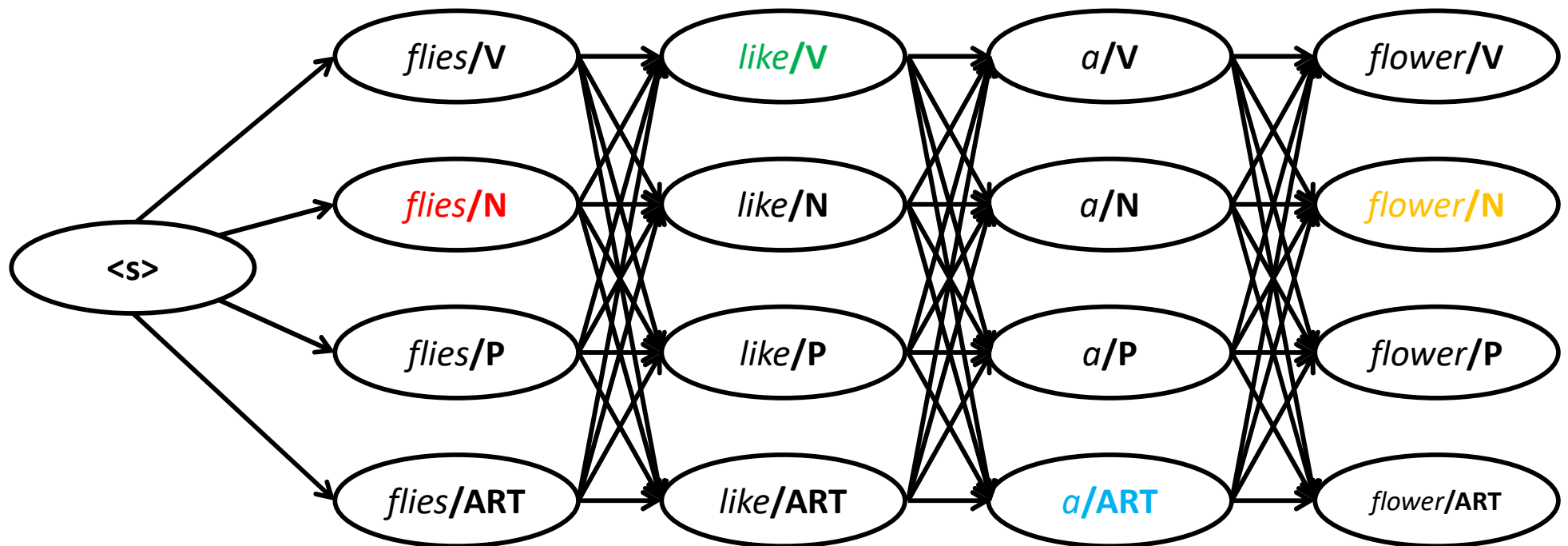
$$P(\textit{like}|\text{V}) = 30/300 = 0.1$$

$$P(\textit{flower}|\text{V}) = 15/300 = 0.05$$

$$P(\textit{like}|\text{P}) = 21/307 = 0.068$$

$$P(\textit{like}|\text{N}) = 10/833 = 0.012$$

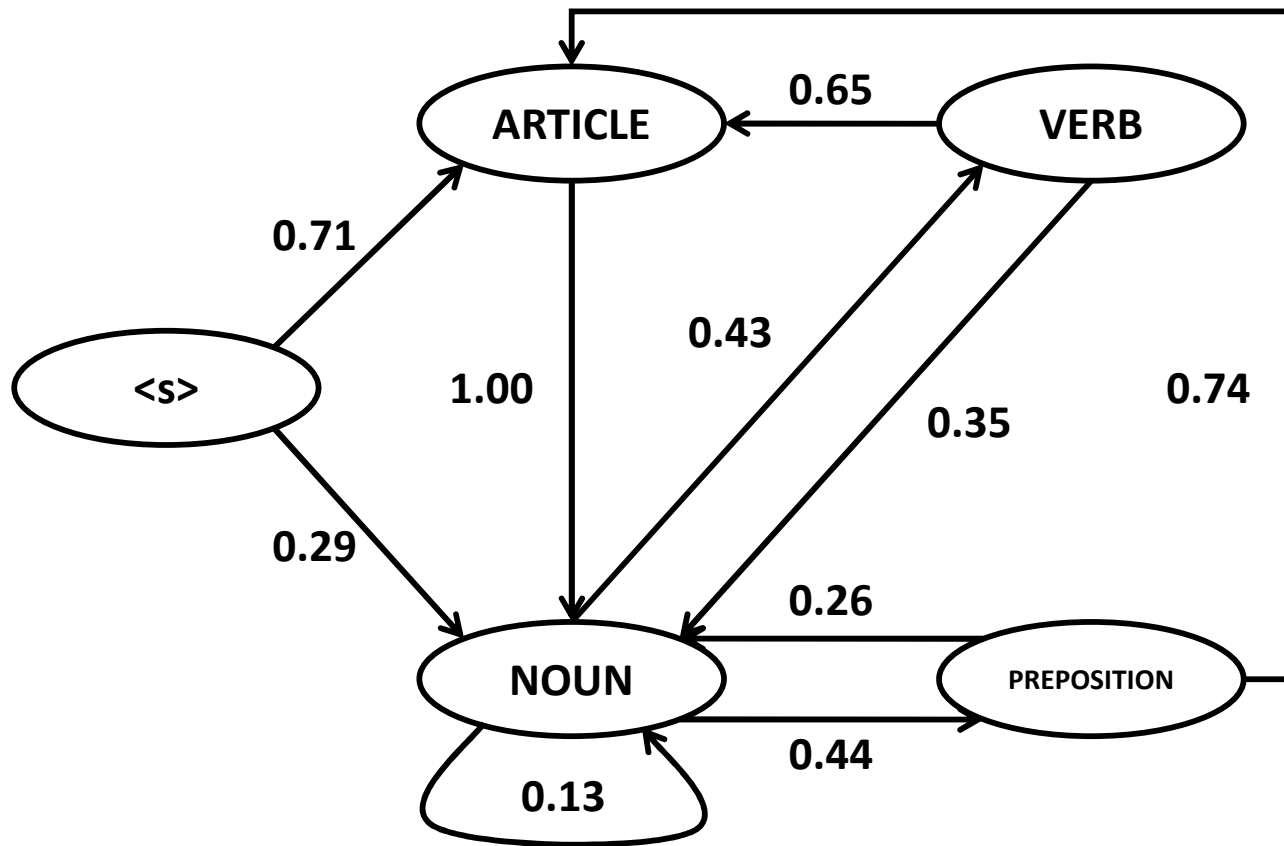
Example: All Possible Sequences



Every sequence can be assigned a probability:

$$\prod_{i=1}^T P(w_i | C_i) = 0.025 * 0.1 * 0.36 * 0.063 = 5.4 * 10^{-5}$$

Hidden Markov Model

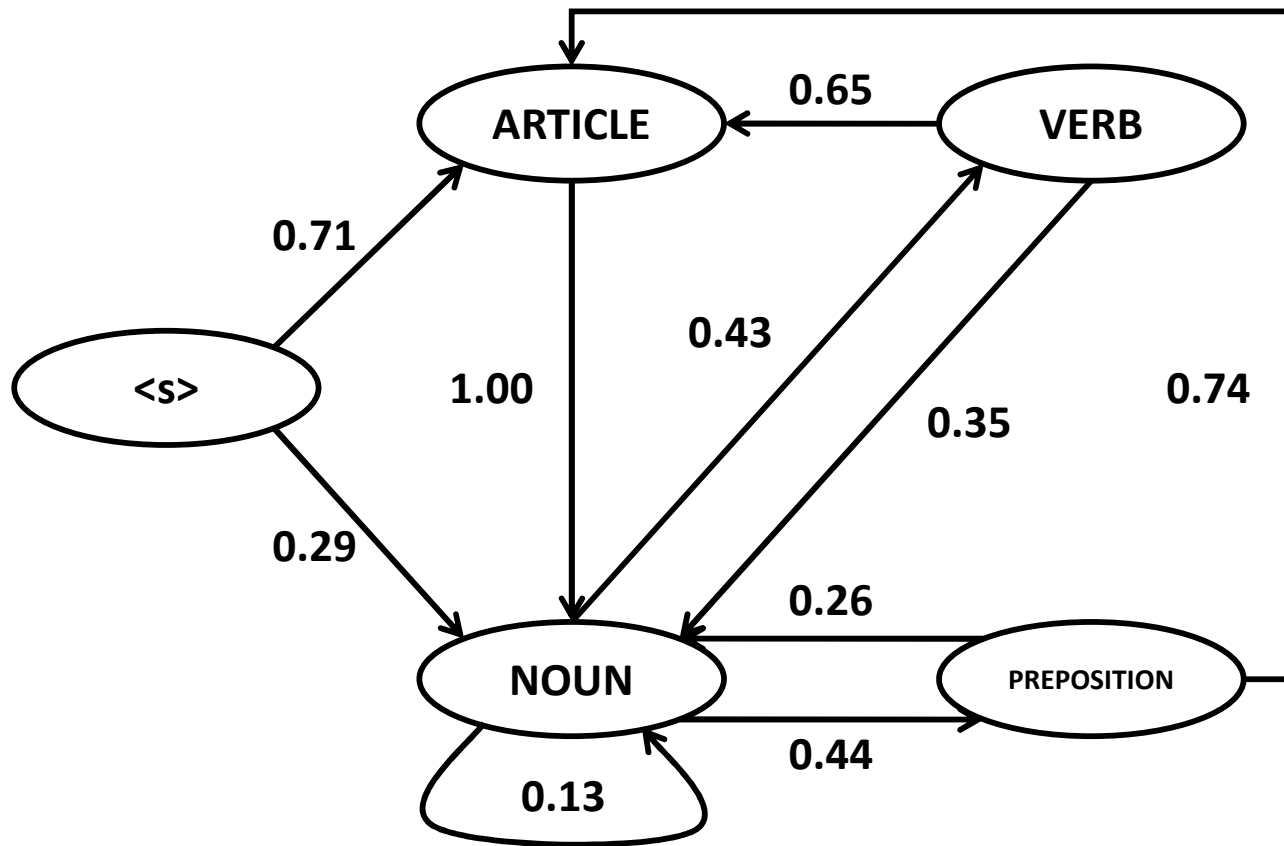


P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.26
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

For any sequence of categories (tags), their probability is:

$$P(C_1, C_2, C_3, \dots, C_T) \cong \prod_{i=1}^T P(C_i | C_{i-1})$$

Hidden Markov Model

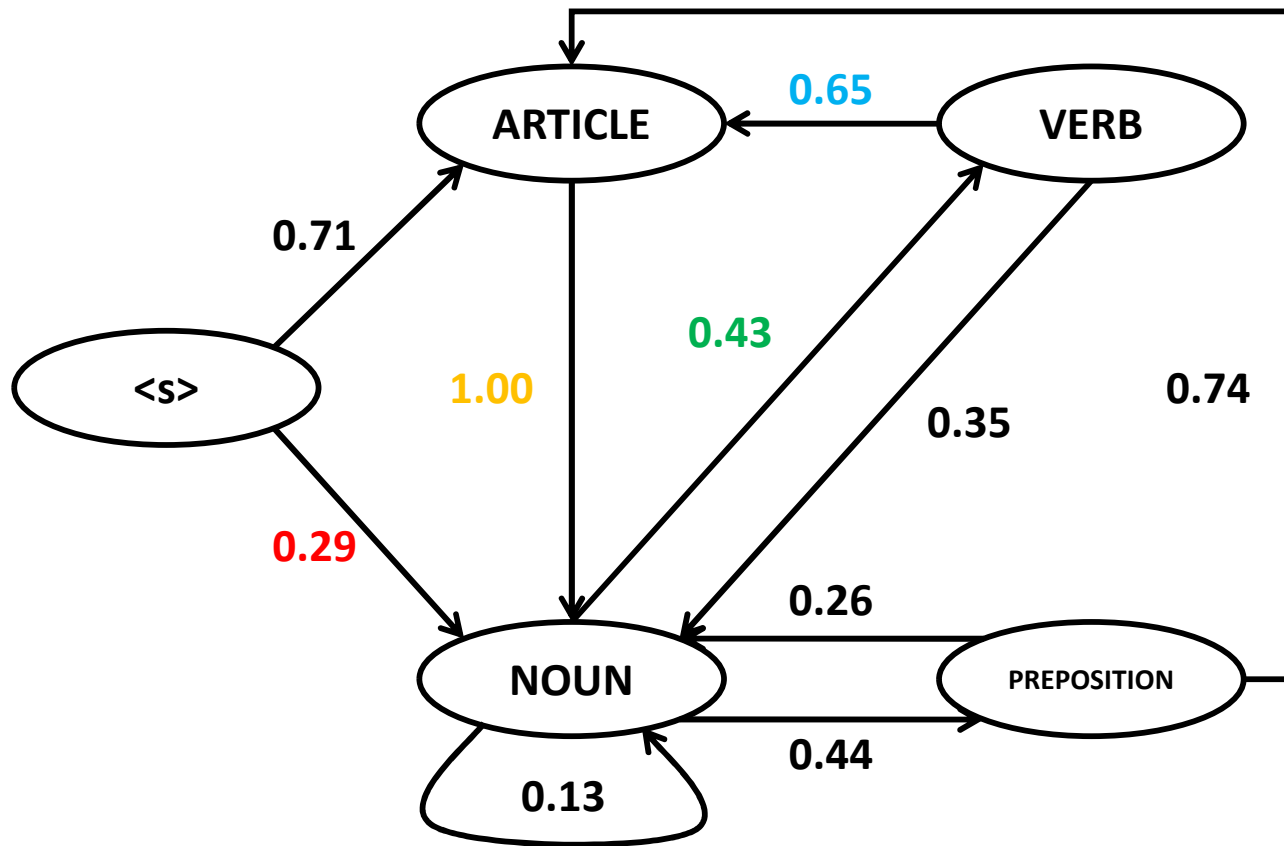


P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.26
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

For any sequence of categories (tags), their probability is:

$$\prod_{i=1}^T P(C_i | C_{i-1}) = P(N | <s>) * (V|N) * (ART|V) * (N|ART)$$

Hidden Markov Model



P(Bigram)	Estimate
P (ARTICLE <S>)	0.71
P (NOUN <S>)	0.29
P (NOUN ARTICLE)	1.00
P (VERB NOUN)	0.43
P (NOUN NOUN)	0.13
P (PREPOSITION NOUN)	0.44
P (NOUN VERB)	0.35
P (ARTICLE VERB)	0.65
P (ARTICLE PREPOSITION)	0.74
P (NOUN PREPOSITION)	0.26

For any sequence of categories (tags), their probability is:

$$\prod_{i=1}^T P(C_i | C_{i-1}) = 0.29 * 0.43 * 0.65 * 1.00 = 0.081$$

Example

Given our synthetic corpus, what is the most like sequence of categories (tags) corresponding to a sentence:

Flies like a flower

For example:

$$P(\textit{Flies, like, a, flower} \mid N, V, ART, N) * P(N, V, ART, N)$$

$$\cong 5.4 * 10^{-5} * 0.081$$

POS Tagging: Simple Tagset

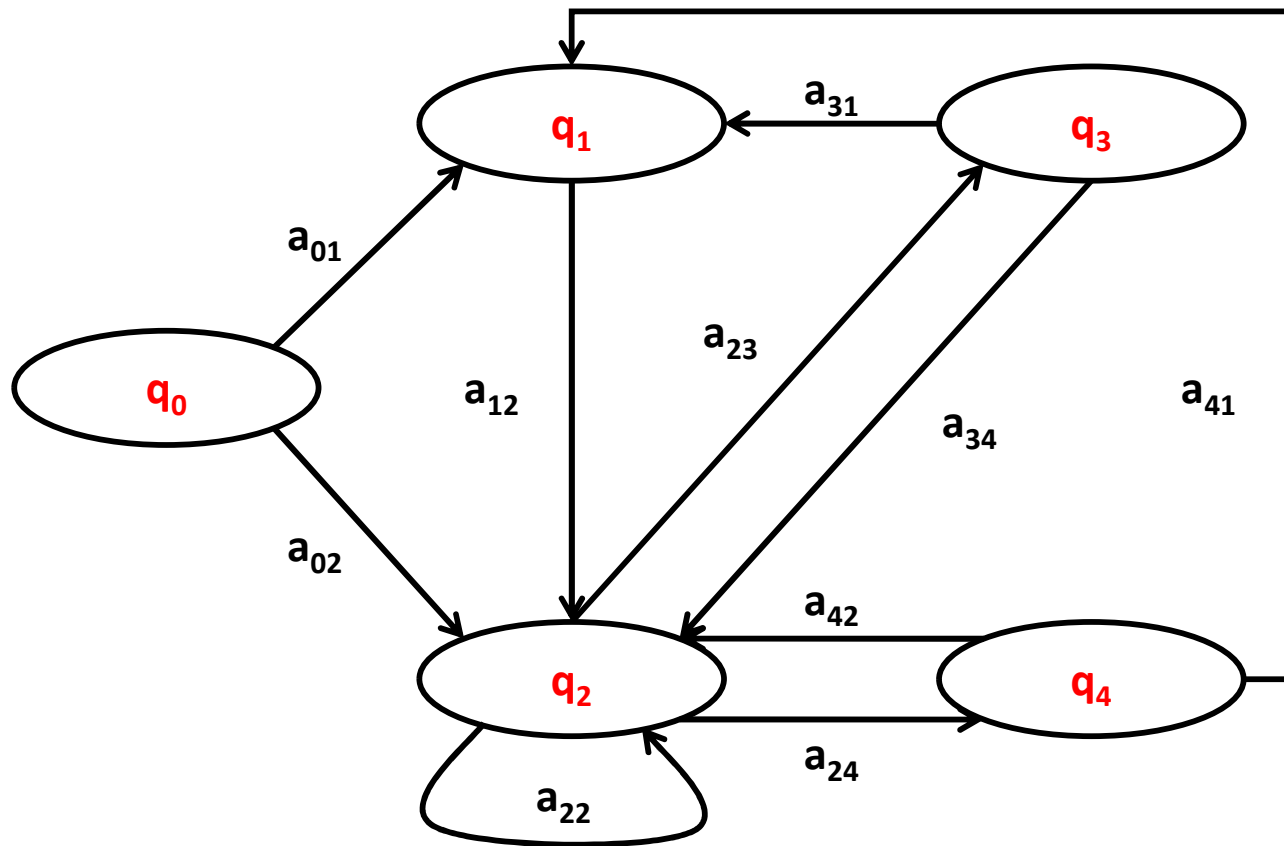
Let's assume we have a simple tagset:

- N - NOUN
- V - VERB
- ART - ARTICLE
- P - PREPOSITION

and some synthetic corpus. Ex. sentence:

Flies like a flower

Hidden Markov Model



HMMs are specified with:

- A set of N states:
 $Q = \{q_1, q_2, \dots, q_N\}$
- A **transition probability** matrix A , where each a_{ij} represents the probability of moving from state q_i to state q_j
- A sequence of observations O :
 $O = o_1, o_2, \dots, o_T$
- A sequence of observation likelihoods (**emission probabilities**): probability of observation o_T being generated by a state q_i
 $B = b_i(o_t)$
- Special <s> and end (final) states

q_0 and q_F

Transition probability matrix A						
	q_0	q_1	q_2	q_3	q_4	Notes
q_0	a_{00}	a_{01}	a_{02}	a_{03}	a_{04}	row sum = 1
q_1	a_{10}	a_{11}	a_{12}	a_{13}	a_{14}	row sum = 1
q_2	a_{20}	a_{21}	a_{22}	a_{23}	a_{24}	row sum = 1
q_3	a_{30}	a_{31}	a_{32}	a_{33}	a_{34}	row sum = 1
q_4	a_{40}	a_{41}	a_{42}	a_{43}	a_{44}	row sum = 1

POS Tagging: General Approach

Estimations with corpus counts:

$$P(C_i = \text{VERB} \mid C_{i-1} = \text{NOUN}) = \frac{\text{Count}(\text{NOUN at position } i-1 \text{ and VERB at } i)}{\text{Count}(\text{NOUN at position } i-1)}$$

Sample bigram probabilities from our synthetic corpus:

Category	Count at i	Pair	Count at i,i+1	P(Bigram)	Estimate
<s>	300	<s>, ARTICLE	213	P (ARTICLE <S>)	0.71
<s>	300	<s>, NOUN	87	P (NOUN <S>)	0.29
ARTICLE	558	ARTICLE, NOUN	558	P (NOUN ARTICLE)	1.00
NOUN	833	NOUN, VERB	358	P (VERB NOUN)	0.43
NOUN	833	NOUN, NOUN	108	P (NOUN NOUN)	0.13
NOUN	833	NOUN, PREPOSITION	366	P (PREPOSITION NOUN)	0.44
VERB	300	VERB, NOUN	75	P (NOUN VERB)	0.35
VERB	300	VERB, ARTICLE	194	P (ARTICLE VERB)	0.65
PREPOSITION	307	PREPOSITION, ARTICLE	226	P (ARTICLE PREPOSITION)	0.74
PREPOSITION	307	PREPOSITION, NOUN	81	P (NOUN PREPOSITION)	0.26

Transition
probabilities

Synthetic Corpus: Word/Tag Counts

Summary of selected word counts in the synthetic corpus:

Word/Tag	N	V	ART	P	TOTAL
<i>flies</i>	21	23	0	0	44
<i>fruit</i>	49	5	1	0	55
<i>like</i>	10	30	0	21	61
<i>a</i>	1	0	201	0	202
<i>the</i>	1	0	300	2	303
<i>flower</i>	53	15	0	0	68
<i>flowers</i>	42	16	0	0	58
<i>birds</i>	64	1	0	0	65
others	592	210	56	284	1142
TOTAL	833	300	558	307	1998

From the table we can calculate lexical generation probabilities $P(w|C)$ estimates:

$$P(\textit{the}|\text{ART}) = 300/558 = 0.54$$

$$P(a|\text{ART}) = 201/558 = 0.36$$

$$P(\textit{flies}|\text{N}) = 21/833 = 0.025$$

$$P(a|\text{N}) = 1/833 = 0.001$$

$$P(\textit{flies}|\text{V}) = 23/300 = 0.076$$

$$P(\textit{flower}|\text{N}) = 53/833 = 0.063$$

$$P(\textit{like}|\text{V}) = 30/300 = 0.1$$

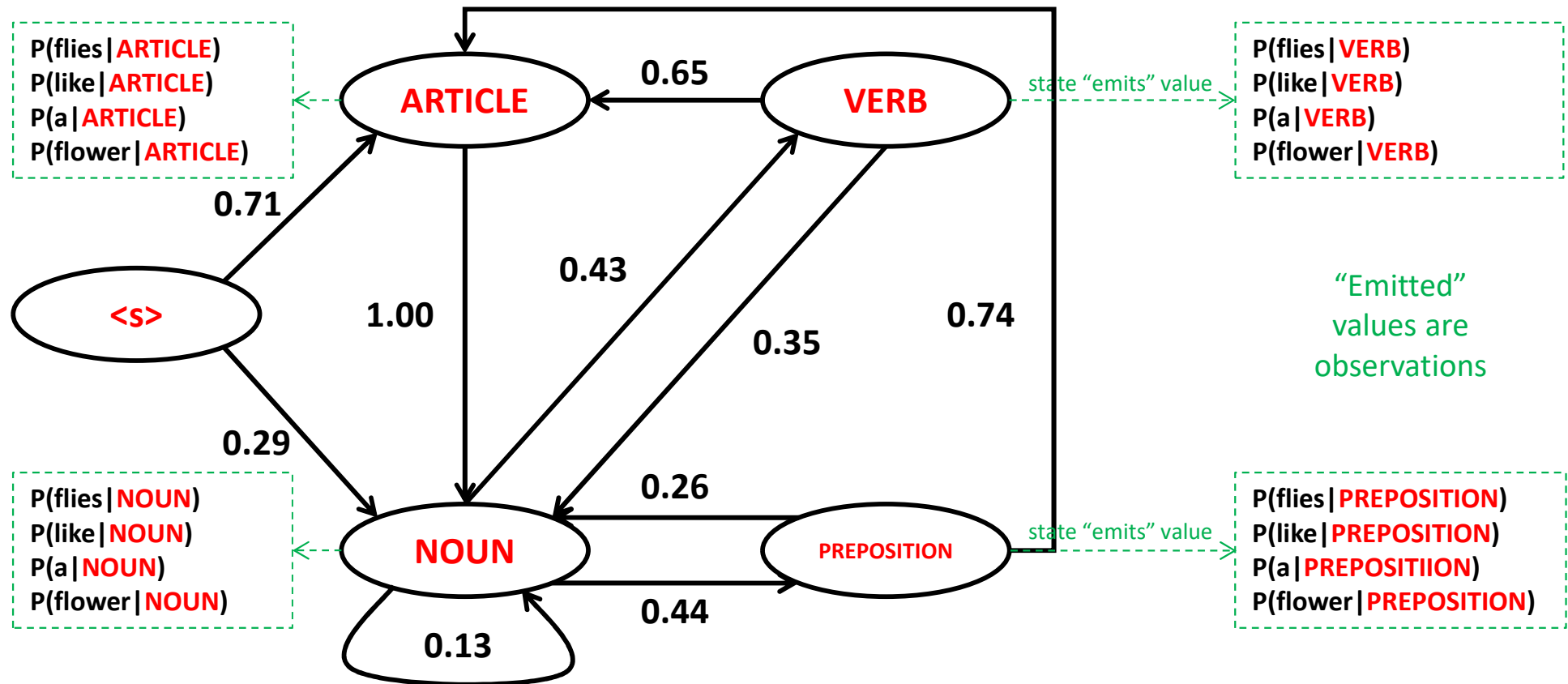
$$P(\textit{flower}|\text{V}) = 15/300 = 0.05$$

$$P(\textit{like}|\text{P}) = 21/307 = 0.068$$

$$P(\textit{like}|\text{N}) = 10/833 = 0.012$$

Emission
probabilities

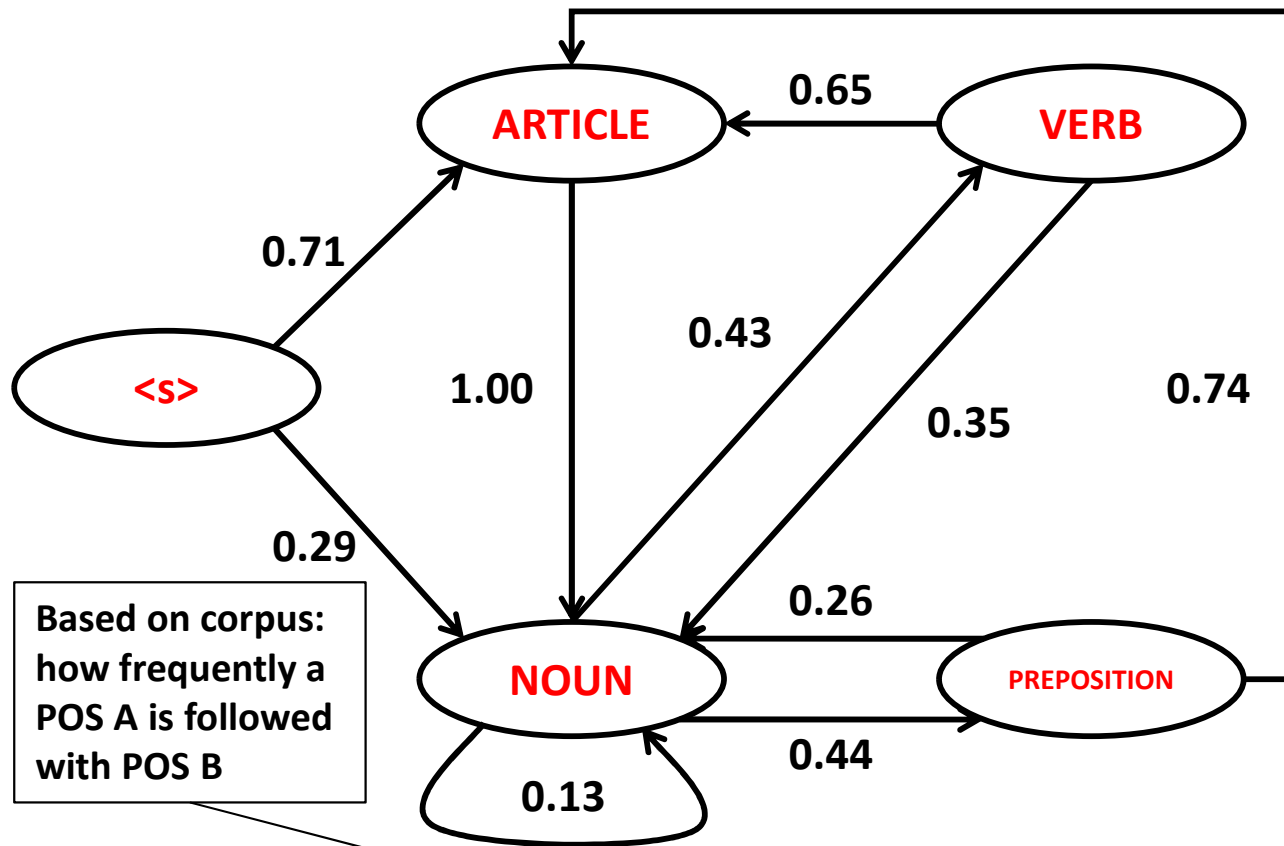
Hidden Markov Model



Transition probability matrix					
	<s>	ARTICLE	NOUN	VERB	PREPOSITION
<s>	0.00	0.71	0.29	0.00	0.00
ARTICLE	0.00	0.00	1.00	0.00	0.00
NOUN	0.00	0.00	0.13	0.43	0.44
VERB	0.00	0.65	0.35	0.00	0.00
PREPOSITION	0.00	0.74	0.26	0.00	0.00

Emission probability matrix				
	flies	like	a	flower
<s>	0.000	0.000	0.000	0.000
ARTICLE	0.000	0.000	0.360	0.000
NOUN	0.025	0.012	0.001	0.063
VERB	0.076	0.100	0.000	0.050
PREPOSITION	0.000	0.068	0.000	0.000

Hidden Markov Model



Based on corpus:
how frequently a
POS A is followed
with POS B

Based on corpus:
how frequently a
word X is tagged
with Y

Transition probability matrix

	<s>	ARTICLE	NOUN	VERB	PREPOSITION
<s>	0.00	0.71	0.29	0.00	0.00
ARTICLE	0.00	0.00	1.00	0.00	0.00
NOUN	0.00	0.00	0.13	0.43	0.44
VERB	0.00	0.65	0.35	0.00	0.00
PREPOSITION	0.00	0.74	0.26	0.00	0.00

Emission probability matrix

	flies	like	a	flower
<s>	0.000	0.000	0.000	0.000
ARTICLE	0.000	0.000	0.360	0.000
NOUN	0.025	0.012	0.001	0.063
VERB	0.076	0.100	0.000	0.050
PREPOSITION	0.000	0.068	0.000	0.000

Example

Given our synthetic corpus, what is the most like sequence of categories (tags) corresponding to a sentence:

Flies like a flower

We need to **maximize**:

$$\begin{aligned} P(w_1, w_2, w_3, \dots, w_T \mid c_1, c_2, c_3, \dots, c_T) * P(c_1, c_2, c_3, \dots, c_T) &\cong \\ &\cong \prod_{i=1}^T P(w_i \mid c_i) * P(c_i \mid c_{i-1}) \end{aligned}$$

Example

Given our synthetic corpus, what is the most like sequence of categories (tags) corresponding to a sentence:

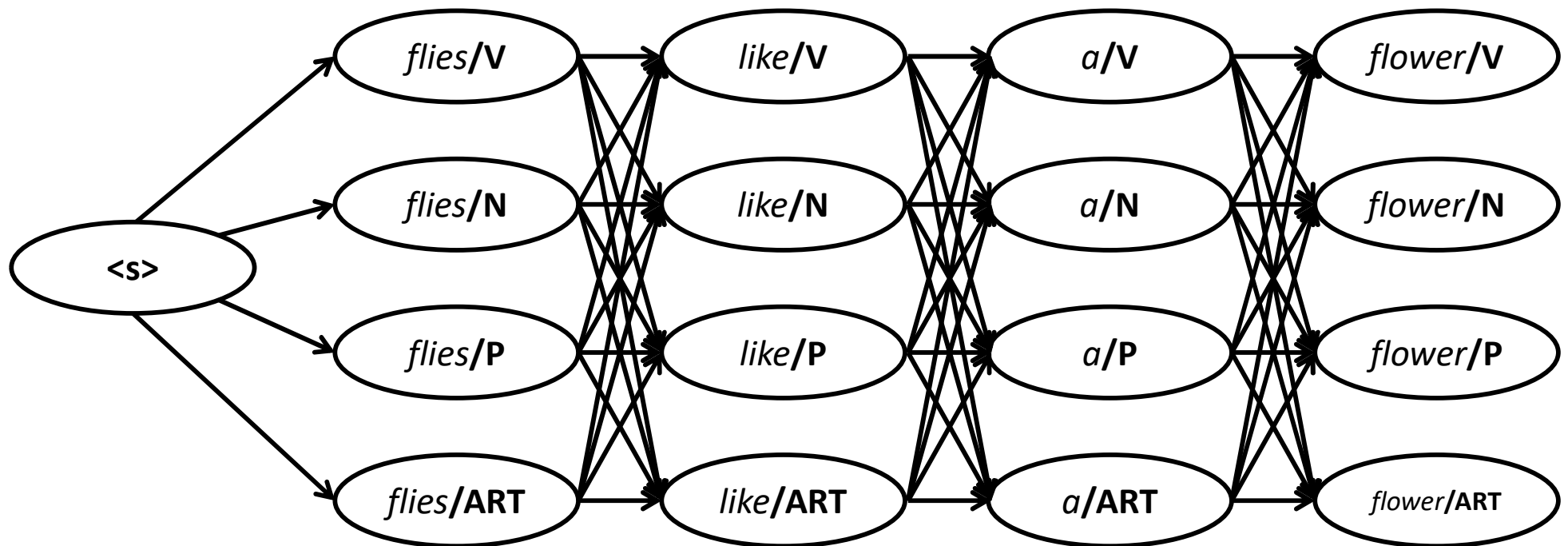
Flies like a flower

For example:

$$P(\textit{Flies, like, a, flower} \mid N, V, ART, N) * P(N, V, ART, N)$$

$$\cong 5.4 * 10^{-5} * 0.081$$

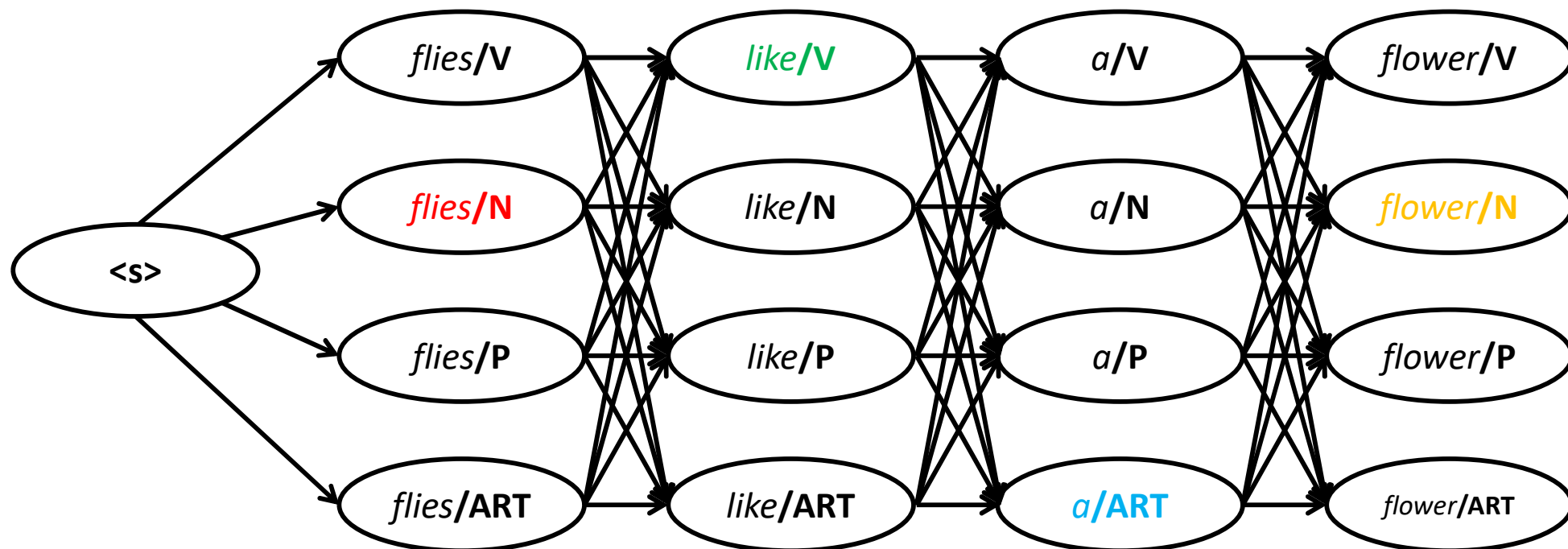
Example: All Possible Sequences



Every sequence can be assigned a probability:

$$P(w_1, w_2, w_3, \dots, w_T \mid c_1, c_2, c_3, \dots, c_T) \cong \prod_{i=1}^T P(w_i \mid c_i)$$

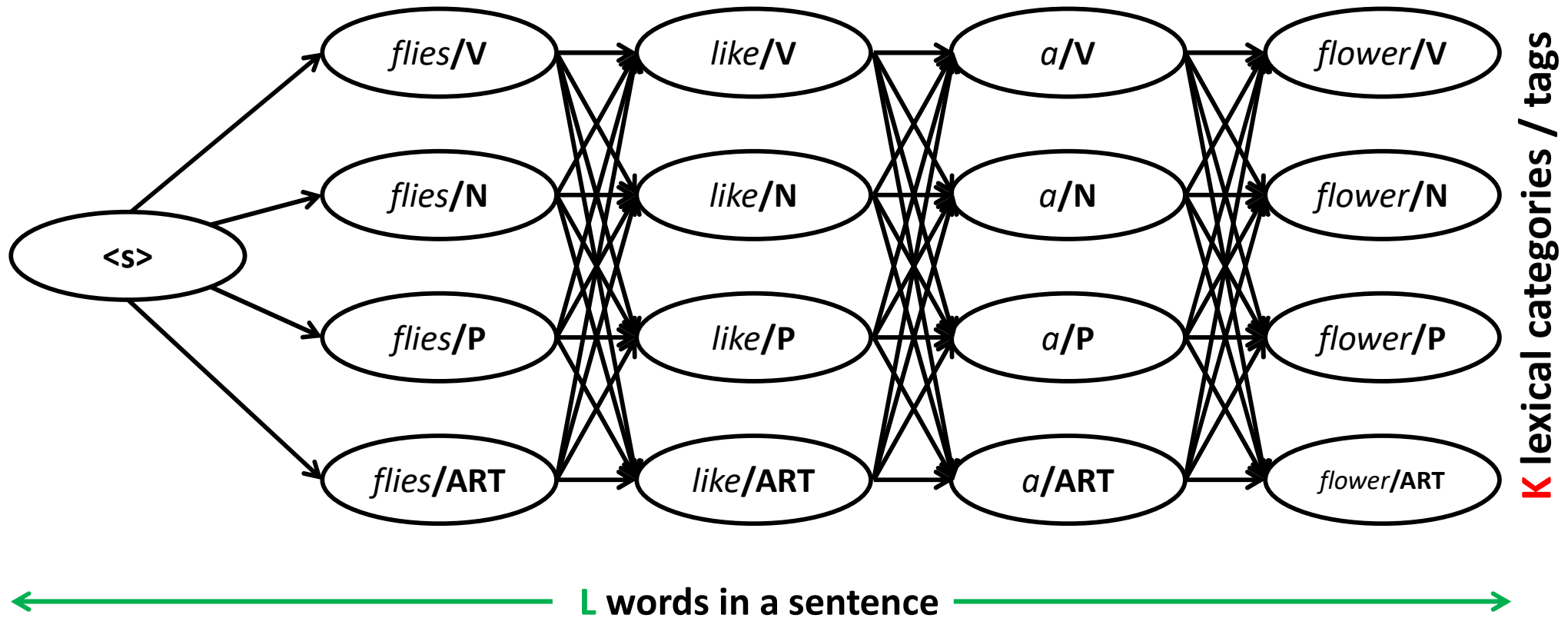
Example: Best Option



Best option will be:

$$\prod_{i=1}^T P(w_i | C_i) = P(\text{flies}|N) * P(\text{like}|V) * P(a|ART) * P(\text{flower}|N)$$

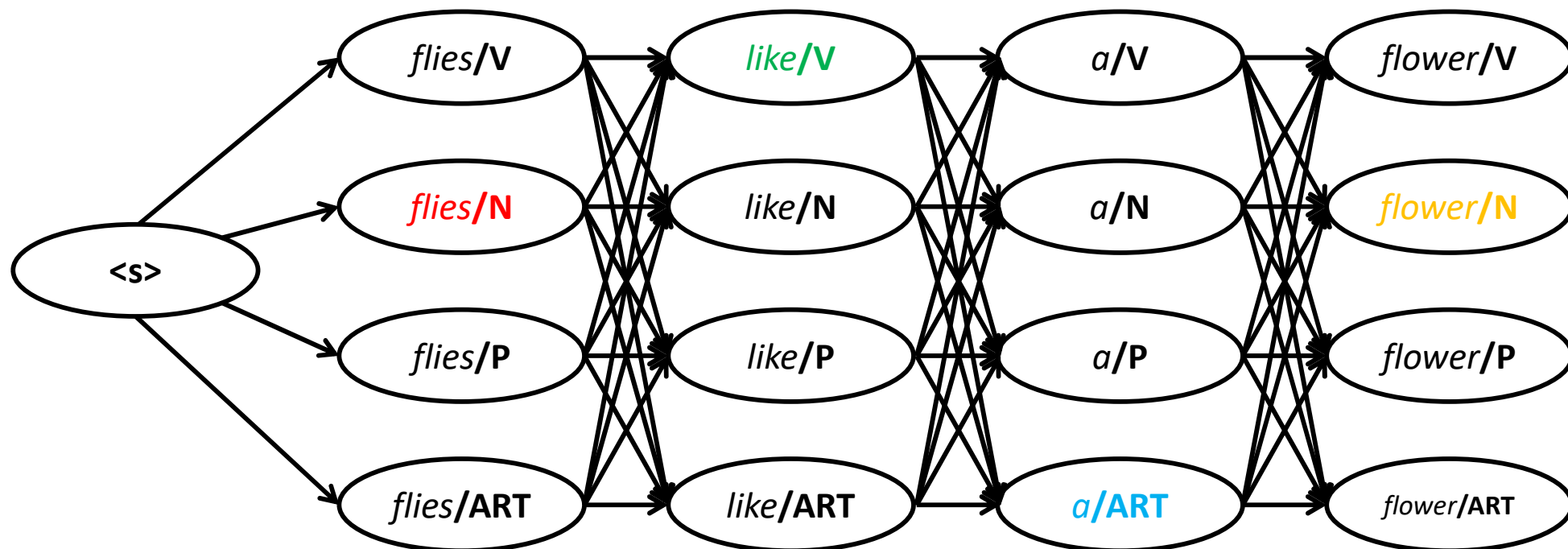
Example: All Possible Sequences



Brute force approach time complexity: $O(K^L)$

$$K = 20, L = 10 \rightarrow 20^{10} = 10240000000000$$

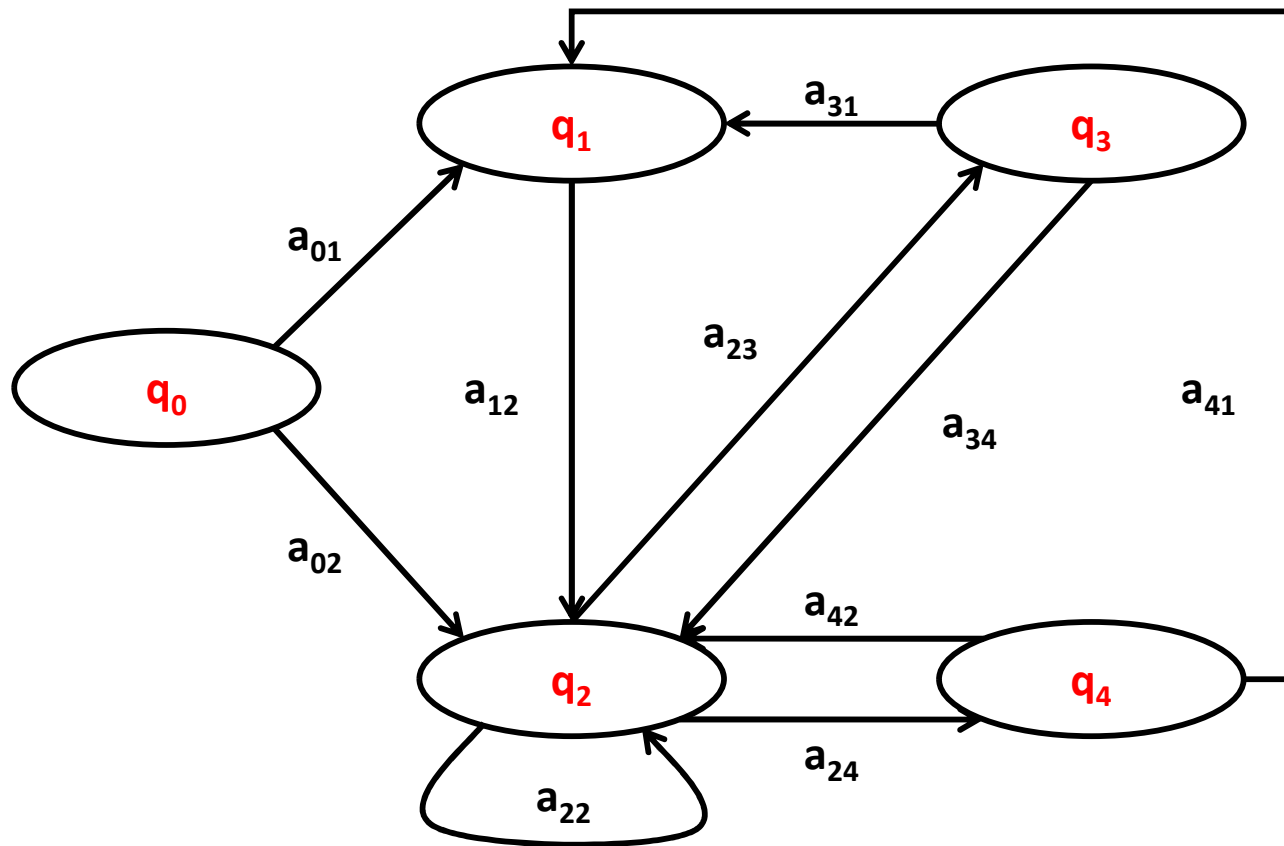
Example: Best Option



How can we efficiently find:

$$\prod_{i=1}^T P(w_i | C_i) = P(\text{flies} | N) * P(\text{like} | V) * P(a | ART) * P(\text{flower} | N)$$

Hidden Markov Model



HMMs are specified with:

- A set of N states:

$$Q = \{q_1, q_2, \dots, q_N\}$$
- A **transition probability** matrix A , where each a_{ij} represents the probability of moving from state q_i to state q_j
- A sequence of observations O :

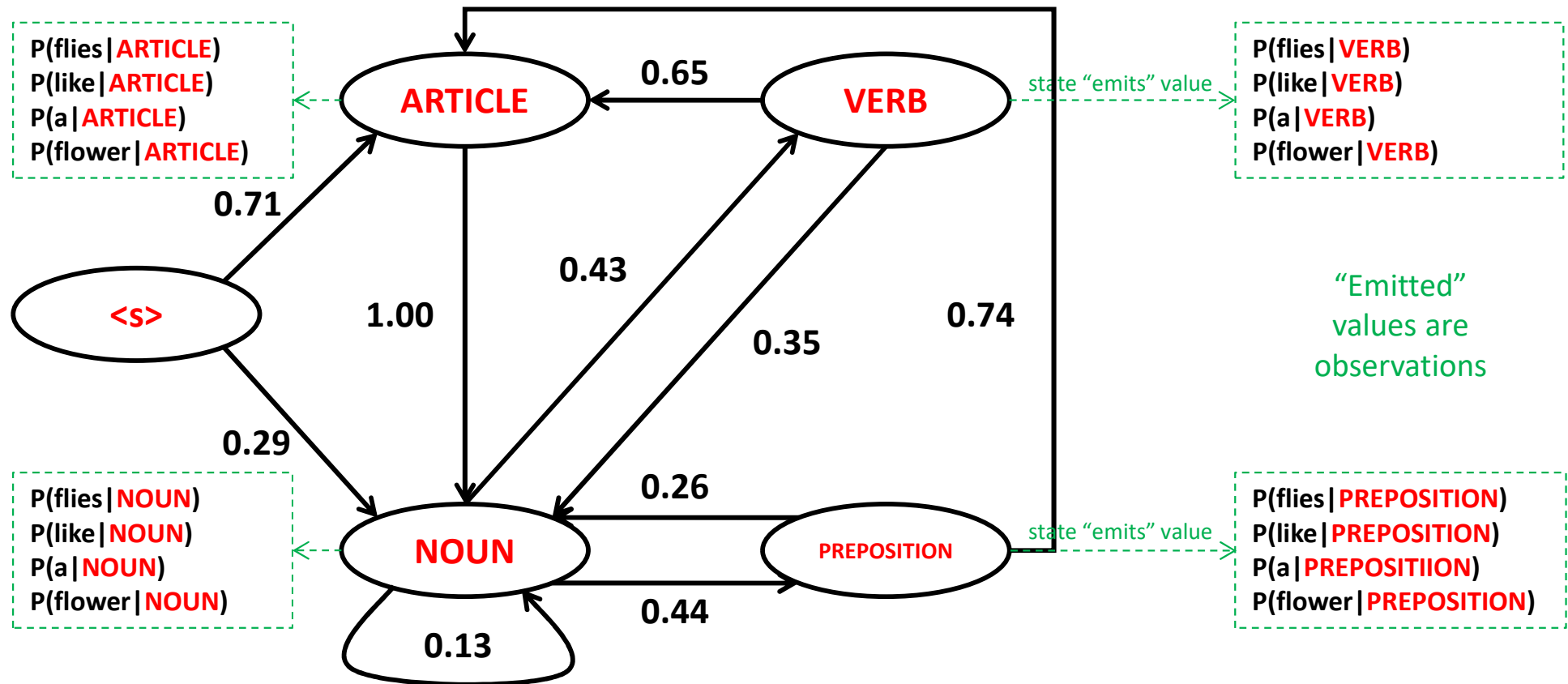
$$O = o_1, o_2, \dots, o_T$$
- A sequence of observation likelihoods (**emission probabilities**): probability of observation o_T being generated by a state q_i

$$B = b_i(o_t)$$
- Special <s> and end (final) states

q_0 and q_F

Transition probability matrix A						
	q ₀	q ₁	q ₂	q ₃	q ₄	Notes
q ₀	a ₀₀	a ₀₁	a ₀₂	a ₀₃	a ₀₄	row sum = 1
q ₁	a ₁₀	a ₁₁	a ₁₂	a ₁₃	a ₁₄	row sum = 1
q ₂	a ₂₀	a ₂₁	a ₂₂	a ₂₃	a ₂₄	row sum = 1
q ₃	a ₃₀	a ₃₁	a ₃₂	a ₃₃	a ₃₄	row sum = 1
q ₄	a ₄₀	a ₄₁	a ₄₂	a ₄₃	a ₄₄	row sum = 1

Hidden Markov Model



Transition probability matrix					
	<S>	ARTICLE	NOUN	VERB	PREPOSITION
<S>	0.00	0.71	0.29	0.00	0.00
ARTICLE	0.00	0.00	1.00	0.00	0.00
NOUN	0.00	0.00	0.13	0.43	0.44
VERB	0.00	0.65	0.35	0.00	0.00
PREPOSITION	0.00	0.74	0.26	0.00	0.00

Emission probability matrix				
	flies	like	a	flower
<S>	0.000	0.000	0.000	0.000
ARTICLE	0.000	0.000	0.360	0.000
NOUN	0.025	0.012	0.001	0.063
VERB	0.076	0.100	0.000	0.050
PREPOSITION	0.000	0.068	0.000	0.000

Hidden Markov Models: Decoding

The task of **determining which sequence of variables is the underlying source of some sequence of observations** is called the **decoding**:

Given as input an HMM $\alpha = (A, B)$ and a sequence of observations o_1, o_2, \dots, o_T find the most probable sequence of states q_1, q_2, \dots, q_T .

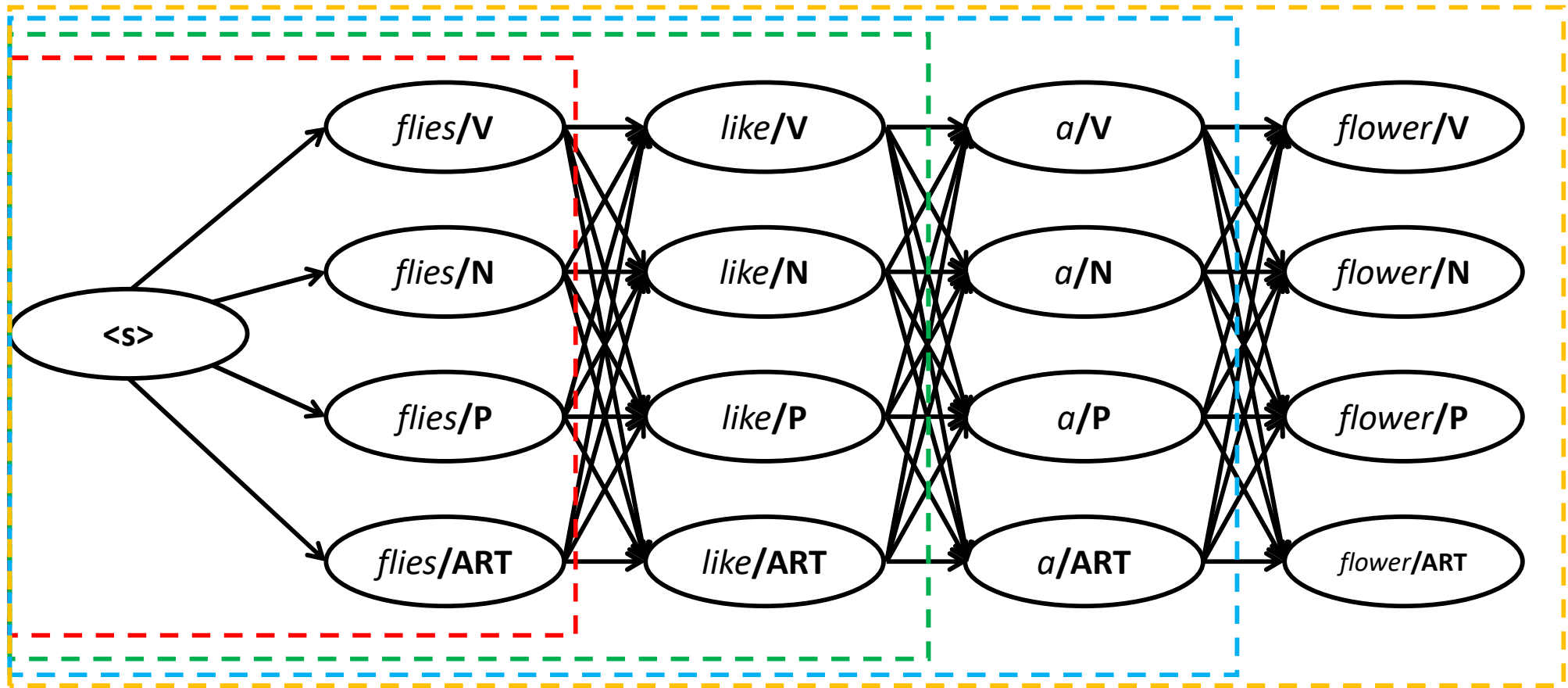
or in our case:

*Given as input an HMM $\alpha = (A, B)$ and a sequence of **words** w_1, w_2, \dots, w_T find the most probable sequence of **tags/states** c_1, c_2, \dots, c_T .*

A - transition probabilities matrix

B - emission probabilities matrix

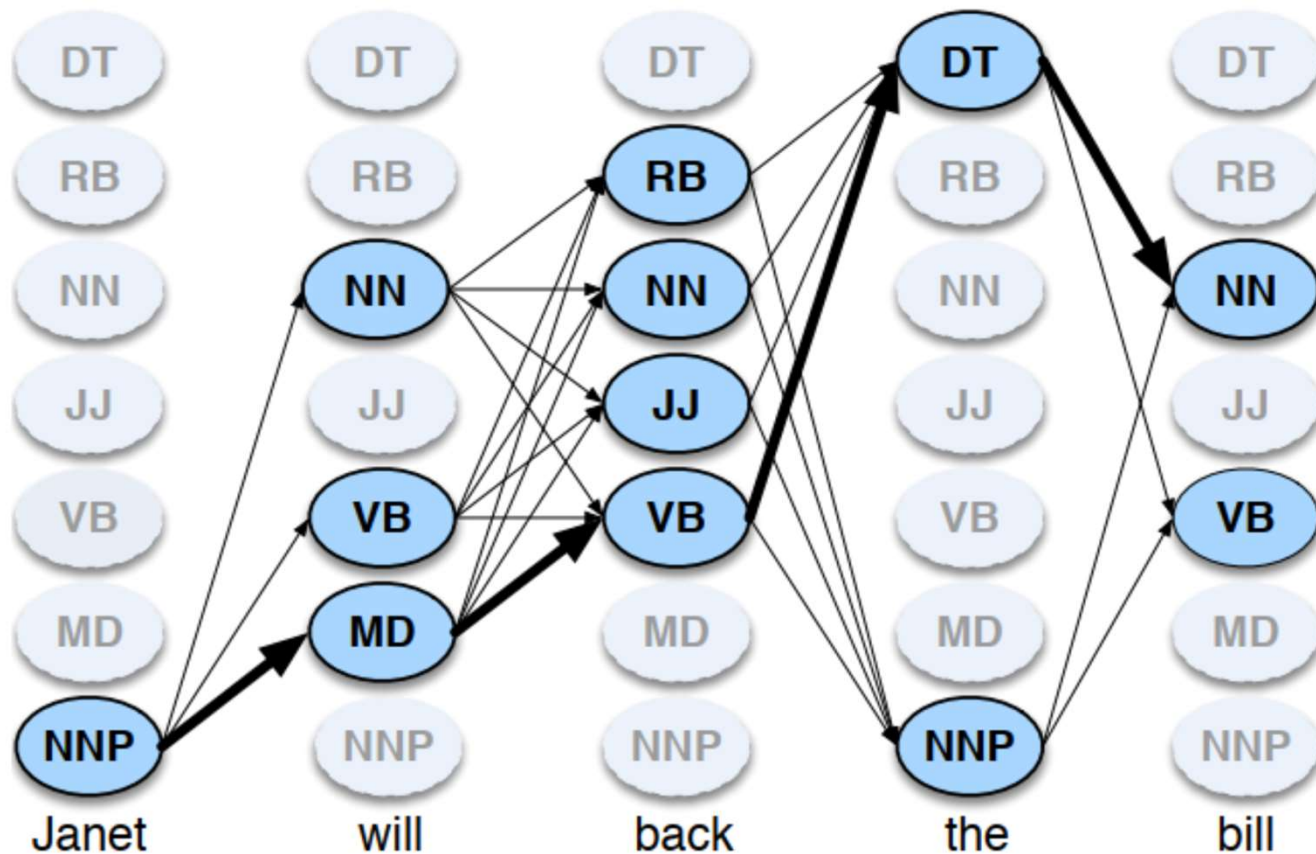
Viterbi Algorithm: the Idea



Maximizing means maximizing , , and .

In other words: maximize $P()$ for all “sub-sentences”.

Example: Possible Tag Sequences



Brown corpus tags:

NN - common noun
 NNP - singular proper noun
 DT - singular determiner
 VB - verb, base form
 JJ - adjective
 MD - modal auxiliary
 RB - adverb

	NP	MD	VB	JJ	NN	RB	DT		Janet	will	back	the	bill
<s>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026	NP	0.000032	0	0	0.000048	0
NP	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025	MD	0	0.308431	0	0	0
MD	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041	VB	0	0.000028	0.000672	0	0.000028
VB	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231	JJ	0	0	0.000340	0	0
JJ	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036	NN	0	0.000200	0.000223	0	0.002337
NN	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068	RB	0	0	0.010446	0	0
RB	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479	DT	0	0	0	0.506099	0
DT	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017						

Transition probabilities matrix

Emission probabilities matrix

Viterbi Algorithm: Pseudocode

function VITERBI(*observations* of len T , *state-graph* of len N) **returns** *best-path*, *path-prob*

create a path probability matrix $viterbi[N, T]$

for each state s **from** 1 **to** N **do** ; initialization step

$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$

$backpointer[s, 1] \leftarrow 0$

for each time step t **from** 2 **to** T **do** ; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$; termination step

$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$; termination step

$bestpath \leftarrow$ the path starting at state $bestpathpointer$, that follows $backpointer[]$ to states back in time

return $bestpath$, $bestpathprob$

Viterbi Algorithm: Pseudocode

function VITERBI(*observations of len T , state-graph of len N*) **returns** *best-path, path-prob*

create a path probability matrix *viterbi*[N, T]

for each state s **from** 1 **to** N **do** ; initialization step

$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$

$backpointer[s, 1] \leftarrow 0$

for each time step t **from** 2 **to** T **do** ; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

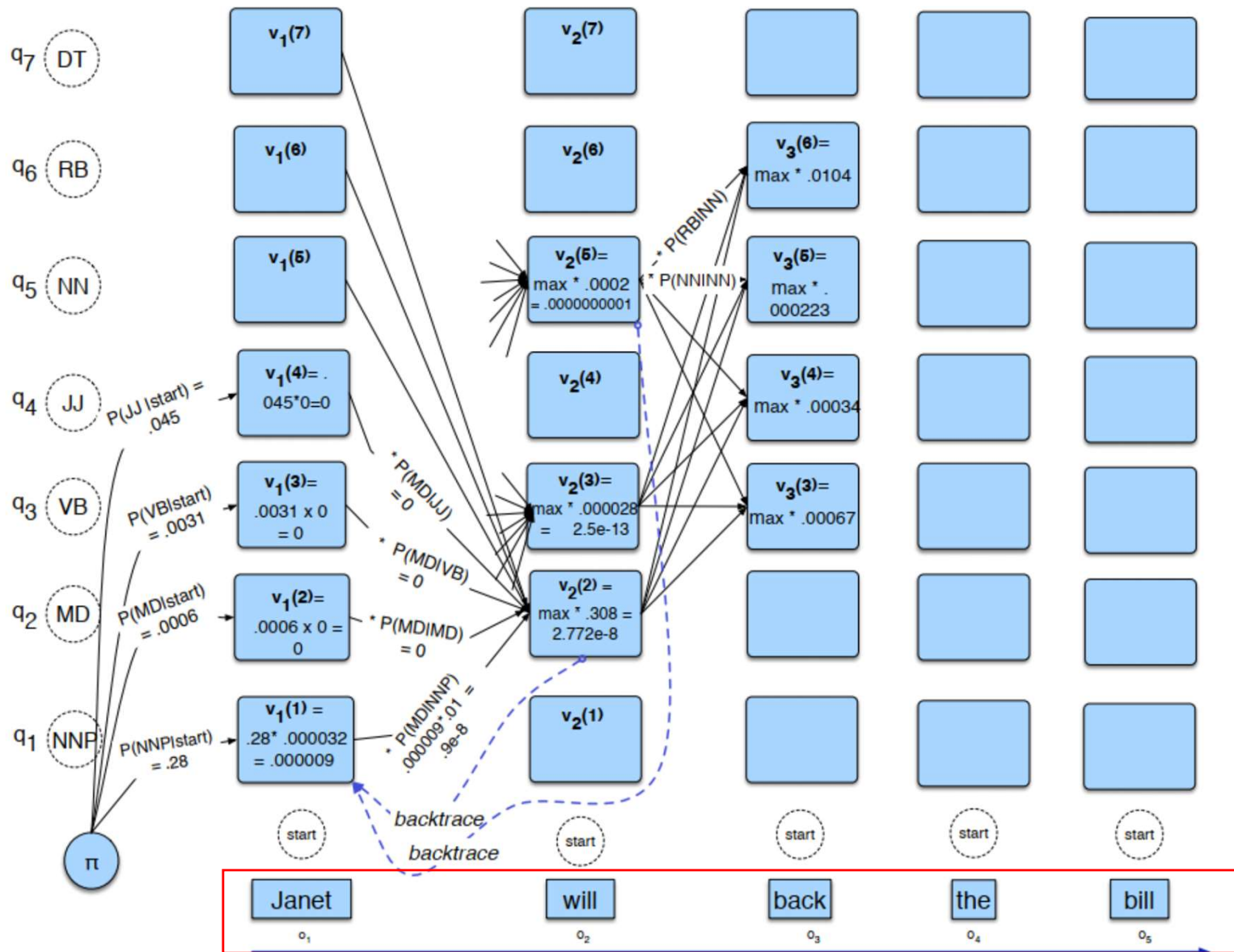
$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$; termination step

$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$; termination step

$bestpath \leftarrow$ the path starting at state $bestpathpointer$, that follows $backpointer[]$ to states back in time

return $bestpath, bestpathprob$

Viterbi Algorithm: Example



Viterbi Algorithm: Pseudocode

function VITERBI(*observations* of len T , *state-graph* of len N) **returns** *best-path*, *path-prob*

create a path probability matrix $viterbi[N, T]$

for each state s **from** 1 **to** N **do** ; initialization step

$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$

$backpointer[s, 1] \leftarrow 0$

for each time step t **from** 2 **to** T **do** ; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

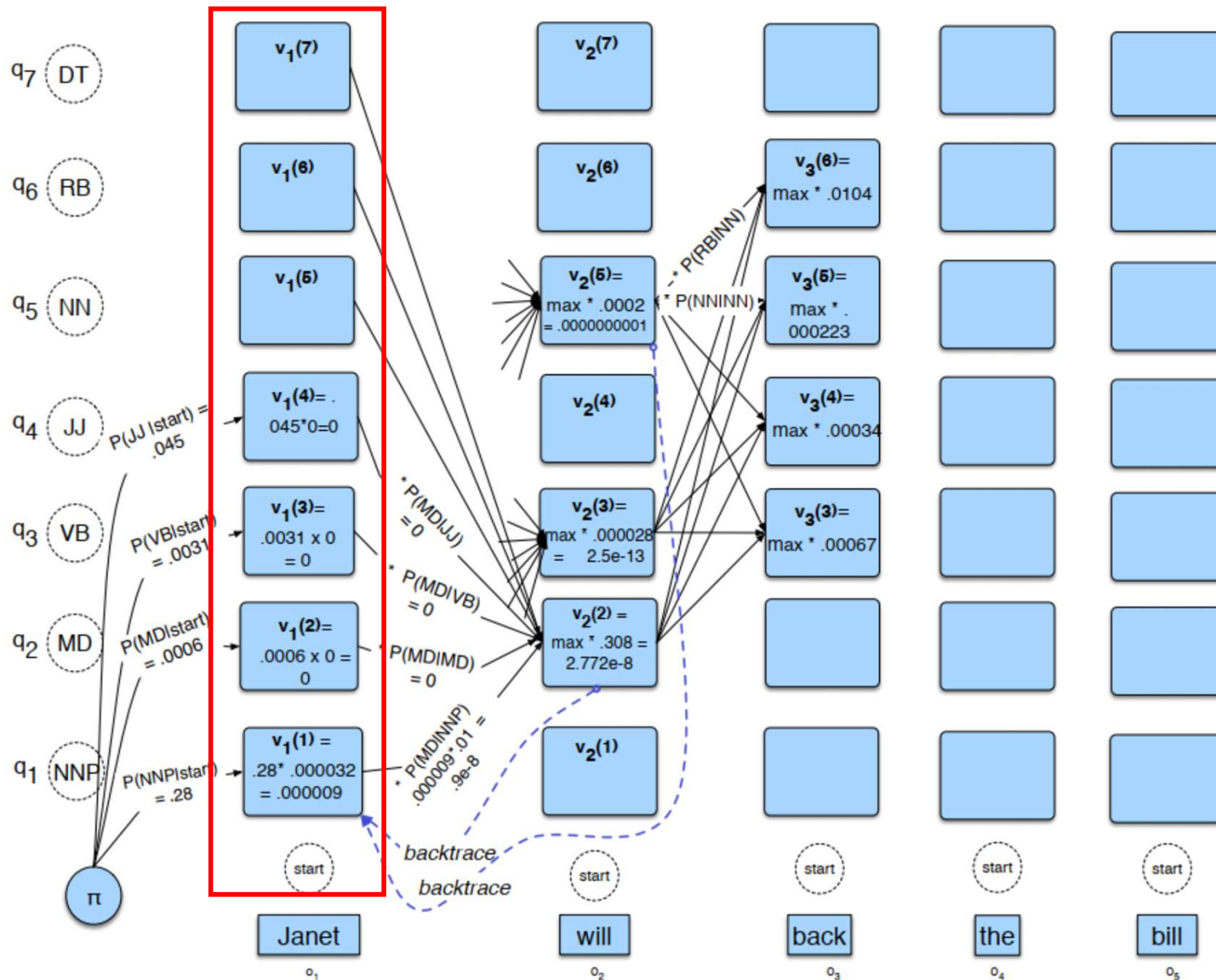
$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$; termination step

$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$; termination step

$bestpath \leftarrow$ the path starting at state $bestpathpointer$, that follows $backpointer[]$ to states back in time

return $bestpath$, $bestpathprob$

Viterbi Algorithm: Example



Initialization step:

$$v_1(1) = v_1(\text{NNP}) = P(T) * P(E) = (\text{NNP} | <s>) * P(\text{Janet} | \text{NNP})$$

$$v_1(2) = v_1(\text{MD}) = P(T) * P(E) = (\text{MD} | <s>) * P(\text{Janet} | \text{MD})$$

$$v_1(3) = v_1(\text{VB}) = P(T) * P(E) = (\text{VB} | <s>) * P(\text{Janet} | \text{VB})$$

$$v_1(4) = v_1(\text{JJ}) = P(T) * P(E) = (\text{JJ} | <s>) * P(\text{Janet} | \text{JJ})$$

$$v_1(5) = v_1(\text{NN}) = P(T) * P(E) = (\text{NN} | <s>) * P(\text{Janet} | \text{NN})$$

$$v_1(6) = v_1(\text{RB}) = P(T) * P(E) = (\text{RB} | <s>) * P(\text{Janet} | \text{RB})$$

$$v_1(7) = v_1(\text{DT}) = P(T) * P(E) = (\text{DT} | <s>) * P(\text{Janet} | \text{DT})$$

$P(T)$ - P(transition)
 $P(E)$ - P(emission)

Viterbi Algorithm: Example

Transition probabilities matrix

	NNP	MD	VB	JJ	NN	RB	DT
<s>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
NNP	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025
MD	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041
VB	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231
JJ	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036
NN	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068
RB	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479
DT	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017

Emission probabilities matrix

	Janet	will	back	the	bill
NNP	0.000032	0	0	0.000048	0
MD	0	0.308431	0	0	0
VB	0	0.000028	0.000672	0	0.000028
JJ	0	0	0.000340	0	0
NN	0	0.000200	0.000223	0	0.002337
RB	0	0	0.010446	0	0
DT	0	0	0	0.506099	0

Initialization step:

$$v_1(1) = v_1(\text{NNP}) = P(T) * P(E) = (\text{NNP} | <s>) * P(\text{Janet} | \text{NNP})$$

$$v_1(2) = v_1(\text{MD}) = P(T) * P(E) = (\text{MD} | <s>) * P(\text{Janet} | \text{MD})$$

$$v_1(3) = v_1(\text{VB}) = P(T) * P(E) = (\text{VB} | <s>) * P(\text{Janet} | \text{VB})$$

$$v_1(4) = v_1(\text{JJ}) = P(T) * P(E) = (\text{JJ} | <s>) * P(\text{Janet} | \text{JJ})$$

$$v_1(5) = v_1(\text{NN}) = P(T) * P(E) = (\text{NN} | <s>) * P(\text{Janet} | \text{NN})$$

$$v_1(6) = v_1(\text{RB}) = P(T) * P(E) = (\text{RB} | <s>) * P(\text{Janet} | \text{RB})$$

$$v_1(7) = v_1(\text{DT}) = P(T) * P(E) = (\text{DT} | <s>) * P(\text{Janet} | \text{DT})$$

$P(T)$ - P(transition)

$P(E)$ - P(emission)

Viterbi Algorithm: Pseudocode

function VITERBI(*observations* of len T , *state-graph* of len N) **returns** *best-path*, *path-prob*

create a path probability matrix $viterbi[N, T]$

for each state s **from** 1 **to** N **do** ; initialization step

$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$

$backpointer[s, 1] \leftarrow 0$

for each time step t **from** 2 **to** T **do** ; recursion step

for each state s **from** 1 **to** N **do**

$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$

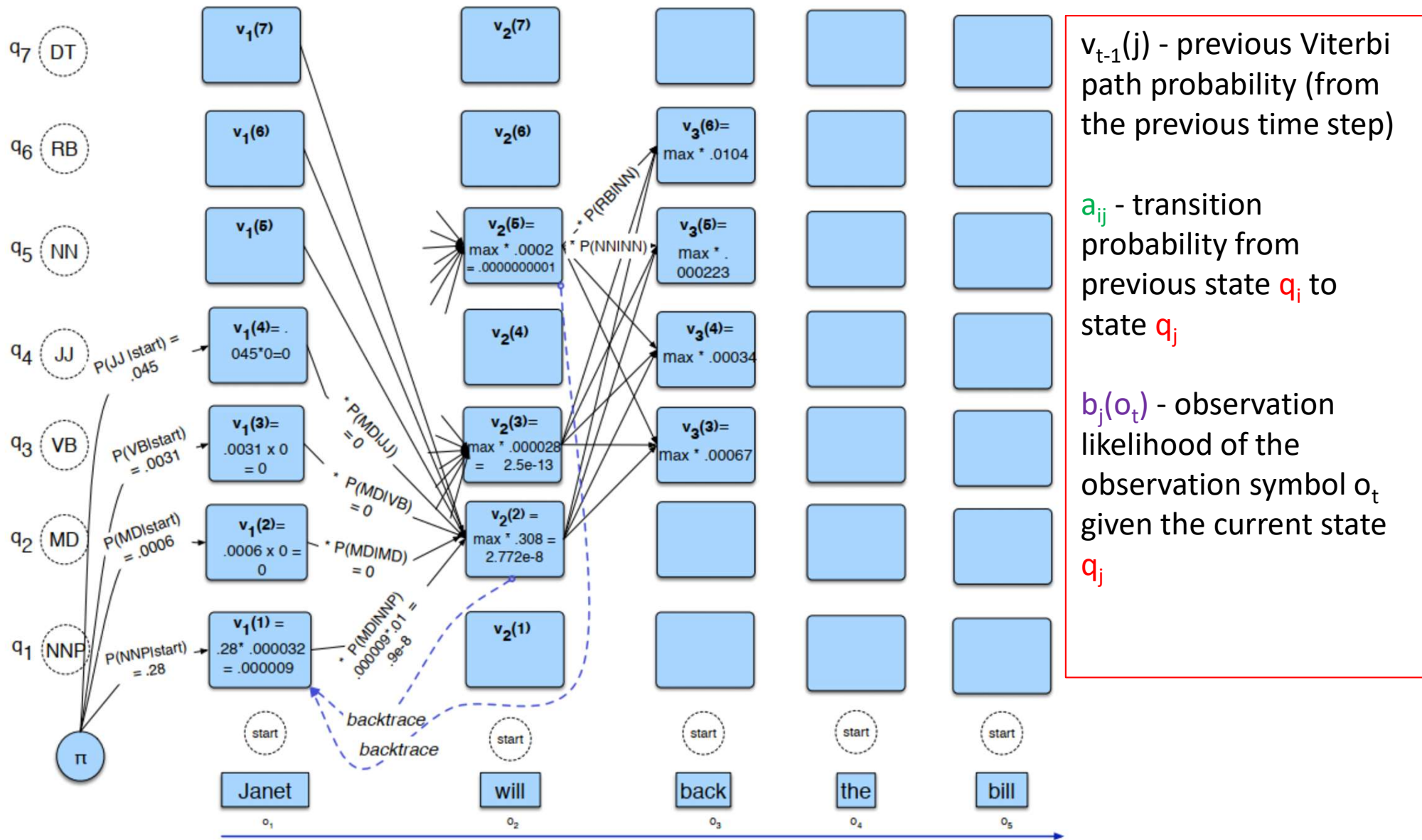
$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$; termination step

$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$; termination step

$bestpath \leftarrow$ the path starting at state $bestpathpointer$, that follows $backpointer[]$ to states back in time

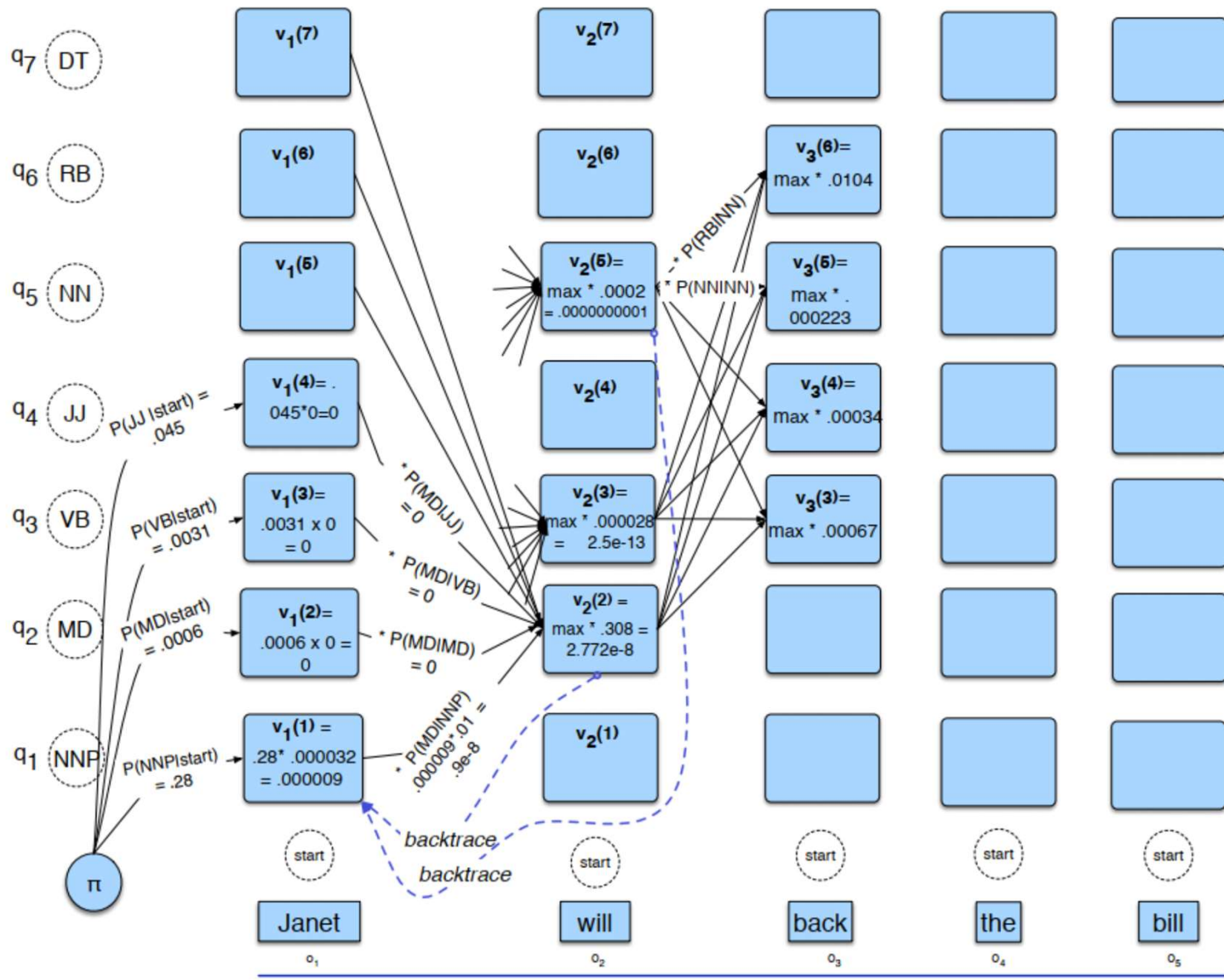
return $bestpath$, $bestpathprob$

Viterbi Algorithm: Example



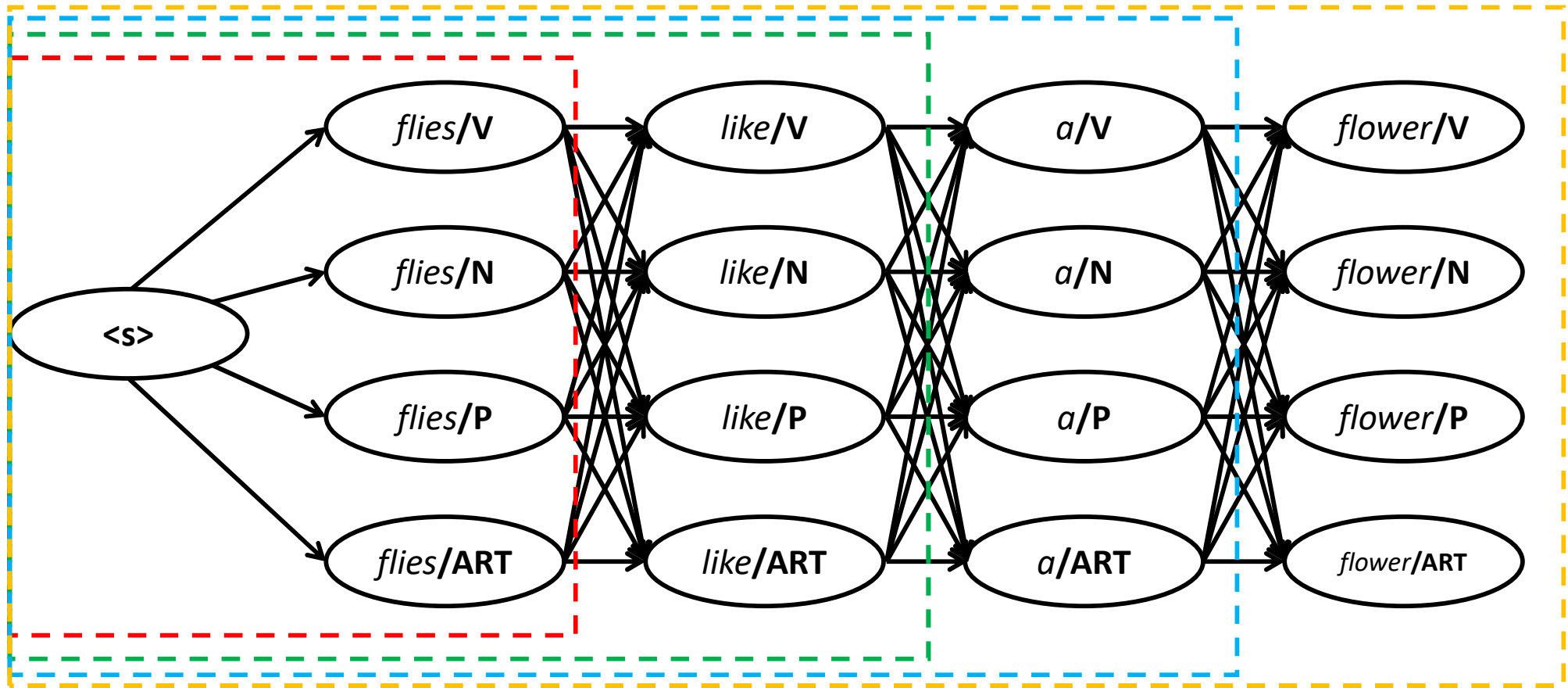
For each state q_j at time t , compute: $v_t(j) = \max_{i=1 \text{ to } N} v_{t-1}(i) * a_{ij} * b_j(o_t)$

Viterbi Algorithm: Example



For each state q_j at time t , compute: $v_t(j) = \max_{i=1t} v_{t-1}(i) * P(\text{transition}) * P(\text{emission})$

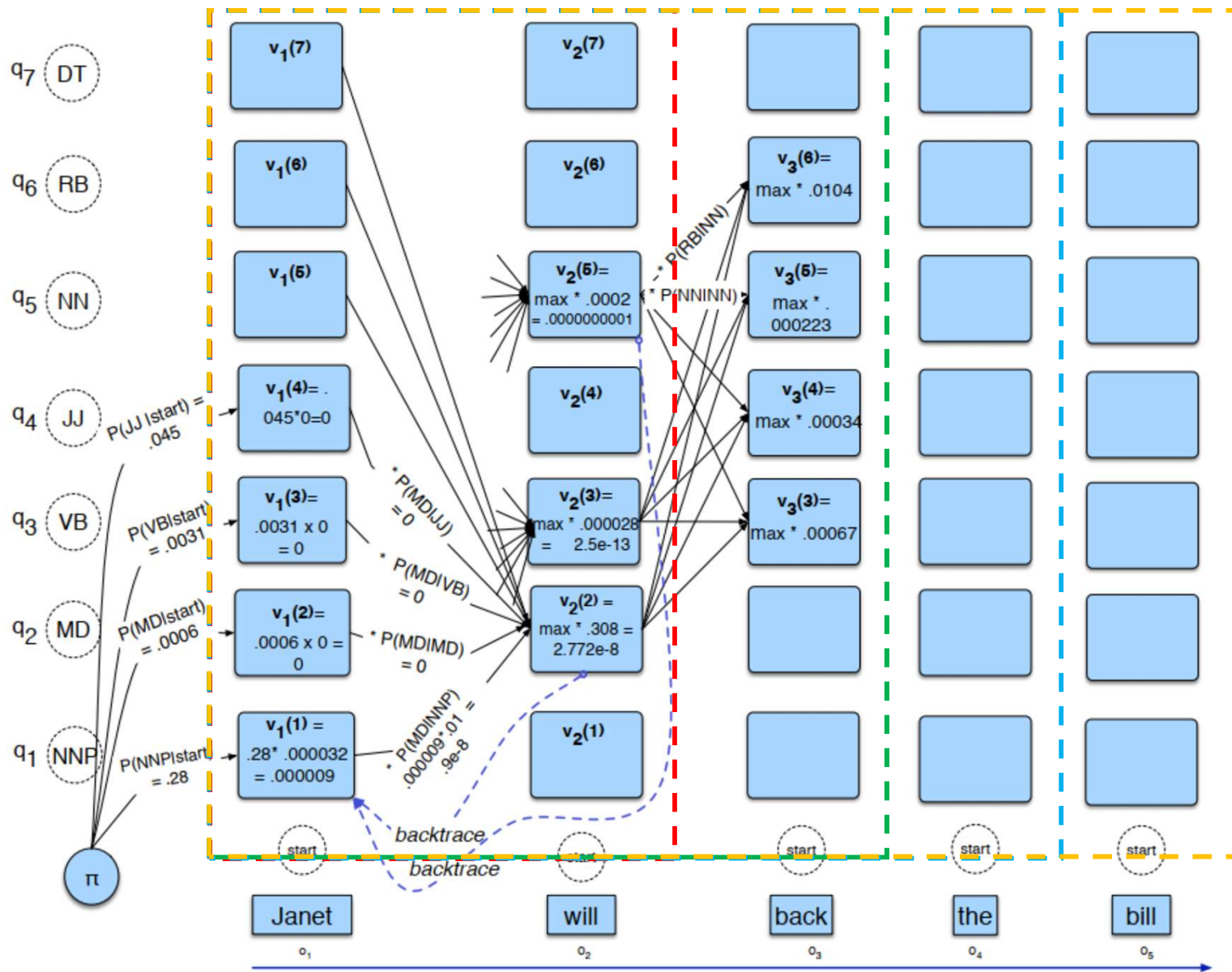
Viterbi Algorithm: the Idea



Maximizing means maximizing , , and .

In other words: maximize $P()$ for all “sub-sentences”.

Viterbi Algorithm: Example



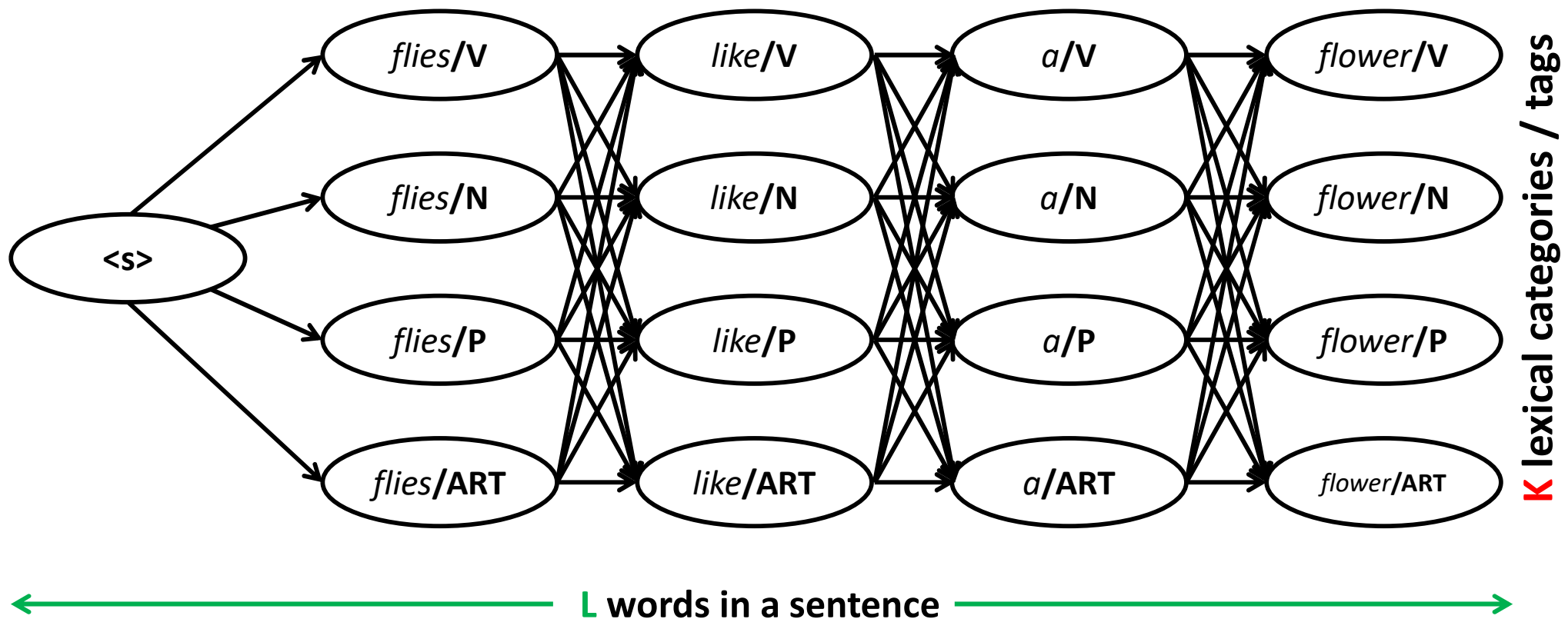
Maximizing [] means maximizing [], [], and [] .

Illinois Institute of Technology



Some paths need not to be explored!

Example: Viterbi



Viterbi algorithm approach time complexity: $O(L * K^2)$

$K = 20, L = 10 \rightarrow 10 * 20^2 = 4000$

versus brute force approach time complexity: $O(K^L)$

$K = 20, L = 10 \rightarrow 20^{10} = 10240000000000$