

CS639 Deep Learning for NLP

Sequence Labeling I

Junjie Hu



Slides adapted from Luke, Yulia, Bob
<https://junjiehu.github.io/cs639-spring26/>

Outline

- Sequence Labeling
- **(Generative) Hidden Markov Model**
- **(Discriminative) Conditional Random Field** (next lecture)

Sequence labeling problems

- Map **a sequence of words** to **a sequence of labels**
 - Part-of-speech tagging (Church, 1988; Brants, 2000)
 - Named entity recognition (Bikel et al., 1990)
 - Text chunking and shallow parsing (Ramshaw and Marcus, 1995)
 - Word alignment of parallel text (Vogel et al., 1996)
 - Compression (Conroy and O'Leary, 2001)
 - Acoustic models, discourse segmentation, etc.

Syntax: Part-of-Speech tagging

- **Open classes** allow new members through borrowing (e.g., the noun *cafe*) and derivation (e.g., the adjective *bounteous* from the noun *bounty*). A word can have new POS tag memberships.
 - Nouns
 - Verbs
 - Adjectives
 - Adverbs
- **Closed classes** of words do not allow new members and usually involve grammatical rather than lexical words. Namely, each word has a closed, fixed POS tag membership. Reasonably easy to enumerate.
 - Prepositions
 - Determiners
 - Pronouns
 - Conjunctions
 - Auxiliary verbs

PART OF SPEECH

WORDS

DT VBZ DT JJ NN
This is a simple sentence

Part of Speech Tagging

- Penn treebank tagset (Marcus et al., 1993)

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	<i>'s</i>	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRPS	possess. pronoun	<i>your, one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WPS	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative adverb	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	<i>\$</i>
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	<i>#</i>
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &</i>	"	left quote	<i>' or "</i>
LS	list item marker	<i>1, 2, One</i>	TO	"to"	<i>to</i>	"	right quote	<i>' or "</i>
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(left paren	<i>[, (, {, <</i>
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>)	right paren	<i>],), }, ></i>
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	<i>,</i>
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	<i>. ! ?</i>
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	<i>: ; ... - -</i>

POS tagging (Example)

- System outputs:
 - The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.
 - There/EX are/VBP 70/CD children/NNS there/RB
 - Preliminary/JJ findings/NNS were/VBD reported/VBN in/IN today/NN 's/POS New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.

Universal Dependencies for All Languages

Universal Dependencies

Universal Dependencies (UD) is a framework for consistent annotation of grammar (parts of speech, morphological features, and syntactic dependencies) across different human languages. UD is an open community effort with over 300 contributors producing more than 150 treebanks in 90 languages. If you're new to UD, you should start by reading the first part of the Short Introduction and then browsing the annotation guidelines.

- [Short introduction to UD](#)
- [UD annotation guidelines](#)
- More information on UD:
 - [How to contribute to UD](#)
 - [Tools for working with UD](#)
 - [Discussion on UD](#)
 - [UD-related events](#)
- Query UD treebanks online:
 - [SETS treebank search](#) maintained by the University of Turku
 - [PML Tree Query](#) maintained by the Charles University in Prague
 - [Kontext](#) maintained by the Charles University in Prague
 - [Grew-match](#) maintained by Inria in Nancy
 - [INESS](#) maintained by the University of Bergen
- [Download UD treebanks](#)

Open class words	Closed class words	Other
ADJ	ADP	PUNCT
ADV	AUX	SYM
INTJ	CCONJ	X
NOUN	DET	
PROPN	NUM	
VERB	PART	
	PRON	
	SCONJ	

Why POS tagging?

- Goal: resolve ambiguities
- Text-to-speech
 - Words w/ slightly different pronunciations denoting different POS, e.g., record/N → /'rekərd/, record/V → /rə'kôrd/
- Lemmatization
 - saw/V → see, saw/N → saw
- Preprocessing for harder disambiguation problems
 - Syntactic parsing
 - Semantic parsing

Sequence labeling as text classification

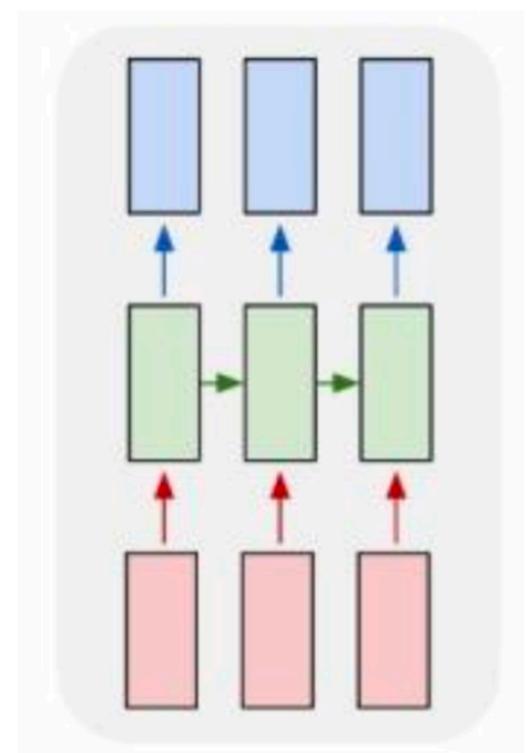
- **Generative** Model: Learn joint probability $P(X, Y)$
 - Hidden Markov Models

$$\hat{Y} = \arg \max_{\substack{y_1 \cdots y_n \\ \forall y_i \in \mathcal{C}}} P(x_1 \cdots x_n, y_1 \cdots y_n)$$

- **Discriminative** Model: Learn conditional probability $P(Y|X)$
 - Conditional Random Fields
 - Neural network-based methods

$$\hat{Y} = \arg \max_Y P(Y|X)$$

- Both trained via Maximum Likelihood Estimation

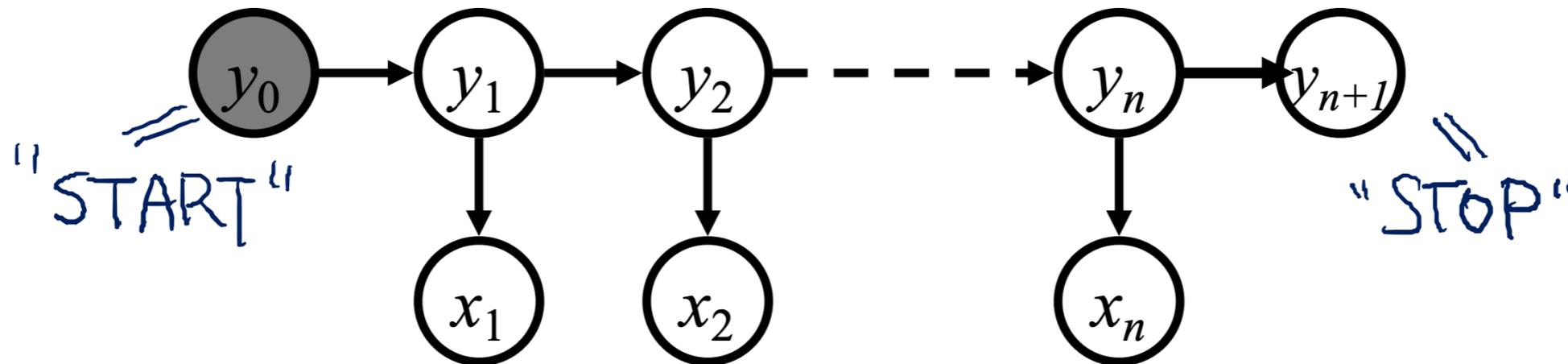


Hidden Markov Model

(Sequential Version of Naive Bayes)

Classic Solution: HMMs

- We want a model of unobservable (hidden) sequences y and observations x



$$p(x_1 \dots x_n, y_1 \dots y_{n+1}) = q(\text{STOP} | y_n) \prod_{i=1}^n q(y_i | y_{i-1}) e(x_i | y_i)$$

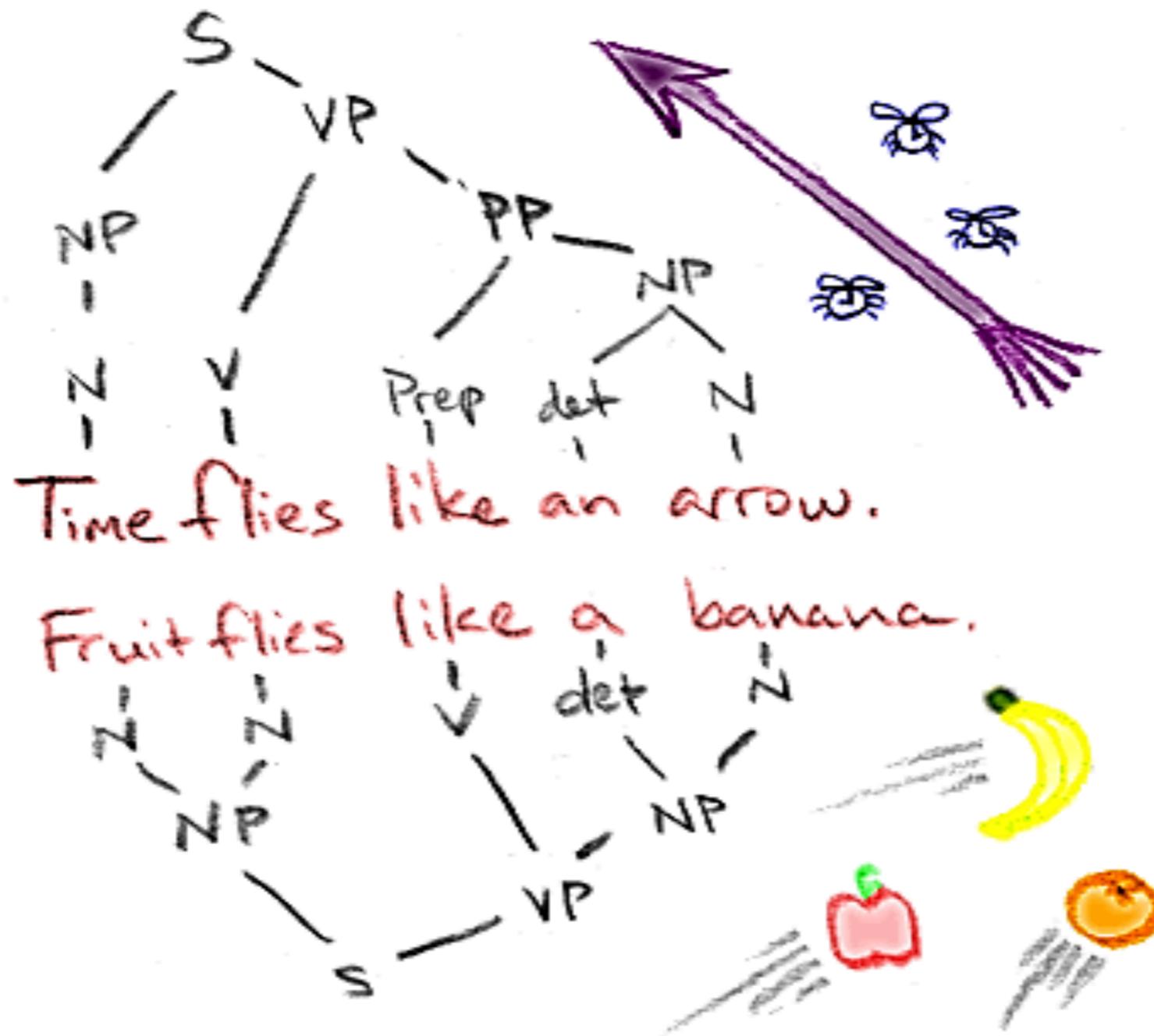
where $y_0 = \text{START}$ and we call $q(y' | y)$ the transition distribution and $e(x | y)$ the emission (or observation) distribution.

Assumptions:

- Tag/state sequence is generated by a Markov model
- Words are chosen independently, conditioned only on the tag/state
- These are totally broken assumptions: why?

Tag predictions depends on context

- Time flies like an arrow
- Fruit flies like a banana



HMM Learning and Inference

- **Learning** by **maximum likelihood estimation**: transition $q(y'|y)$ and emissions $e(x|y)$

$$p(x_1 \dots x_n, y_1 \dots y_{n+1}) = q(\text{STOP}|y_n) \prod_{i=1}^n q(y_i|y_{i-1}) e(x_i|y_i)$$

- **Inference** (linear time in sentence length!)

- Viterbi:

$$y^* = \operatorname{argmax}_{y_1 \dots y_n} p(x_1 \dots x_n, y_1 \dots y_{n+1})$$

where $y_{n+1} = \text{STOP}$

- Forward Backward:

$$p(x_1 \dots x_n, y_i) = \sum_{y_1 \dots y_{i-1}} \sum_{y_{i+1} \dots y_n} p(x_1 \dots x_n, y_1 \dots y_n)$$

Learning: Maximum Likelihood

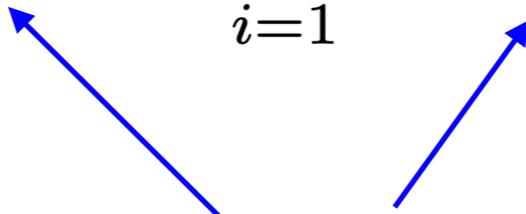
- Supervised Learning
 - Assume m fully labeled training examples:

$$\{(x^{(i)}, y^{(i)}) \mid i = 1 \cdots m\}$$

where $x^{(i)} = x_1 \cdots x_n$ and $y^{(i)} = y_1 \cdots y_n$

- What's the maximum likelihood estimate?

$$p(x_1 \cdots x_n, y_1 \cdots y_{n+1}) = q(\text{STOP} \mid y_n) \prod_{i=1}^n q(y_i \mid y_{i-1}) e(x_i \mid y_i)$$

$$q_{ML}(y_i \mid y_{i-1})$$


$$e_{ML}(x_i \mid y_i)$$


Learning: Maximum Likelihood

- MLE: counting the co-occurrence of the event

$$q_{ML}(y_i|y_{i-1}) = \frac{c(y_{i-1}, y_i)}{c(y_{i-1})} \quad e_{ML}(x|y) = \frac{c(y, x)}{c(y)}$$

- Will these estimates be high quality?
 - Which is likely to be more sparse, q or e ?
 - The emission function, because $c(y, x)$ is more likely to have sparse values.
- Can use all the same smoothing tricks we used for counting-based language models!
- Other approaches: Map low-frequency words to a small, finite set of units (e.g., prefixes, word classes), and run MLE on new sequences

Inference (Decoding)

- Problem: find the most likely (Viterbi) sequence under the model

$$y^* = \operatorname{argmax}_{y_1 \dots y_n} p(x_1 \dots x_n, y_1 \dots y_{n+1})$$

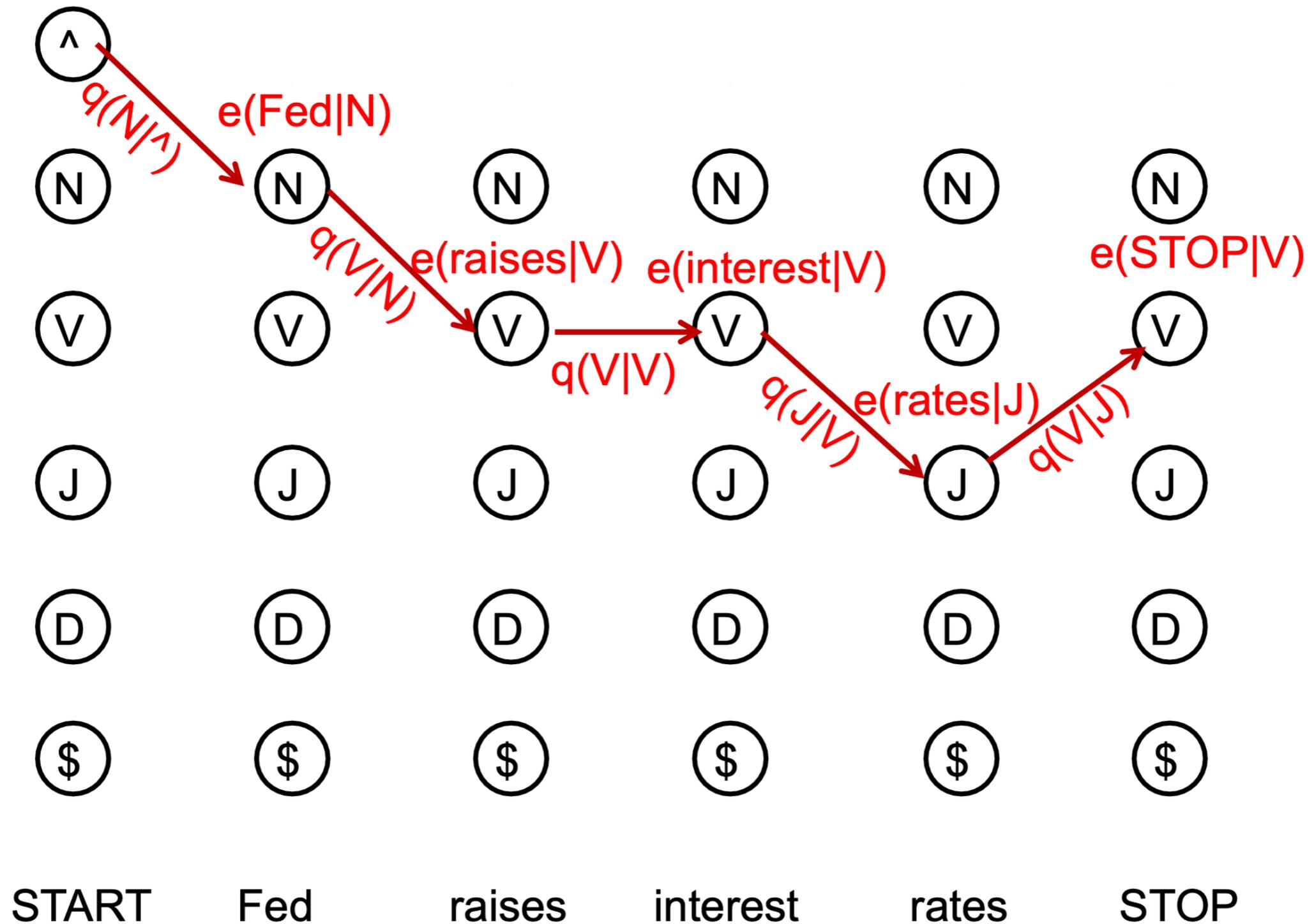
- Given model parameters, we can score any sequence pair

NNP	VBZ	NN	NNS	CD	NN	.
Fed	raises	interest	rates	0.5	percent	.

- In principle, we can list all possible tag sequences, score each one, and pick **the best one (a.k.a. the Viterbi state sequence)**

NNP	VBZ	NN	NNS	CD	NN	⇒	logP = -23
NNP	NNS	NN	NNS	CD	NN	⇒	logP = -29
NNP	VBZ	VB	NNS	CD	NN	⇒	logP = -27

The State Lattice/Trellis: Viterbi



- Brute force approach: enumerate $n^{\mathcal{K}}$ possible tag sequences

Dynamic Programming!

- Focus on max, consider special case of $n=2$
- Define $\pi(i, y_i)$ to be the max score of a sequence of length i ending in tag y_i

$$\begin{aligned} & \max_{y_1, y_2} q(STOP|y_2)q(y_2|y_1)e(x_2|y_2)q(y_1|START)e(x_1|y_1) \\ &= \max_{y_2} q(STOP|y_2)e(x_2|y_2) \max_{y_1} q(y_1|START)q(y_2|y_1)e(x_1|y_1) \\ &= \max_{y_2} q(STOP|y_2)e(x_2|y_2)\pi(2, y_2) \\ & \text{given that } \pi(2, y_2) = \max_{y_1} q(y_1|START)q(y_2|y_1)e(x_1|y_1) \end{aligned}$$

- What about the general case? (Consider $n=3$, etc...)

Dynamic Programming!

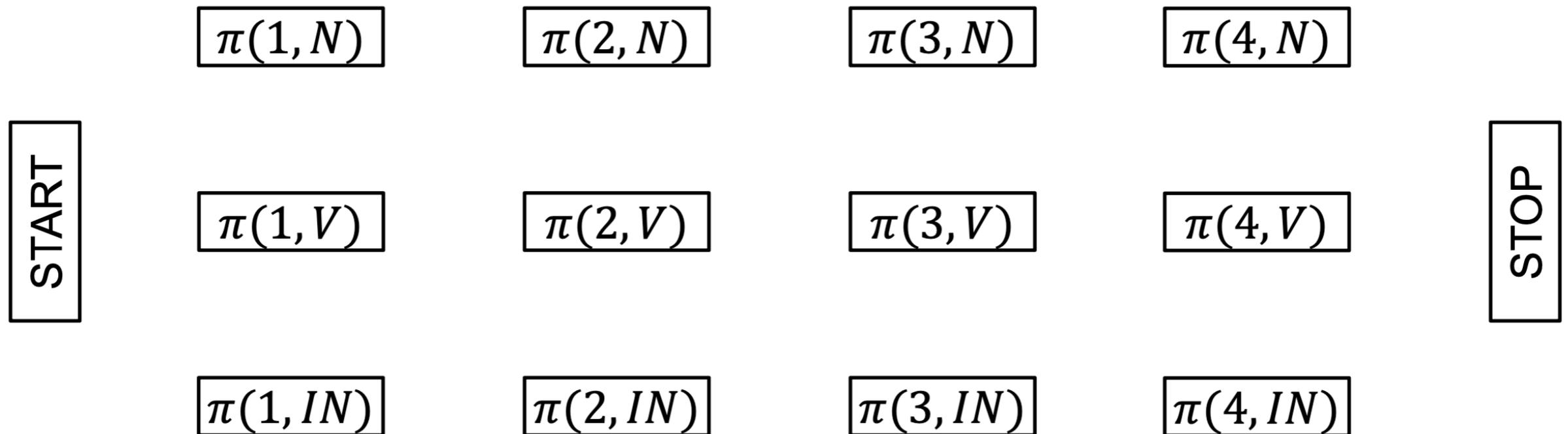
- General case
- Define $\pi(i, y_i)$ to be the max score of a sequence of length i ending in tag y_i

$$\begin{aligned}\pi(i, y_i) &= \max_{y_1 \dots y_{i-1}} p(x_1 \dots x_i, y_1 \dots y_i) \\ &= \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \max_{y_1 \dots y_{i-2}} p(x_1 \dots x_{i-1}, y_1 \dots y_{i-1}) \\ &= \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})\end{aligned}$$

- We now have an efficient algorithm. Start with $i=0$ and work your way to the end of the sentence!

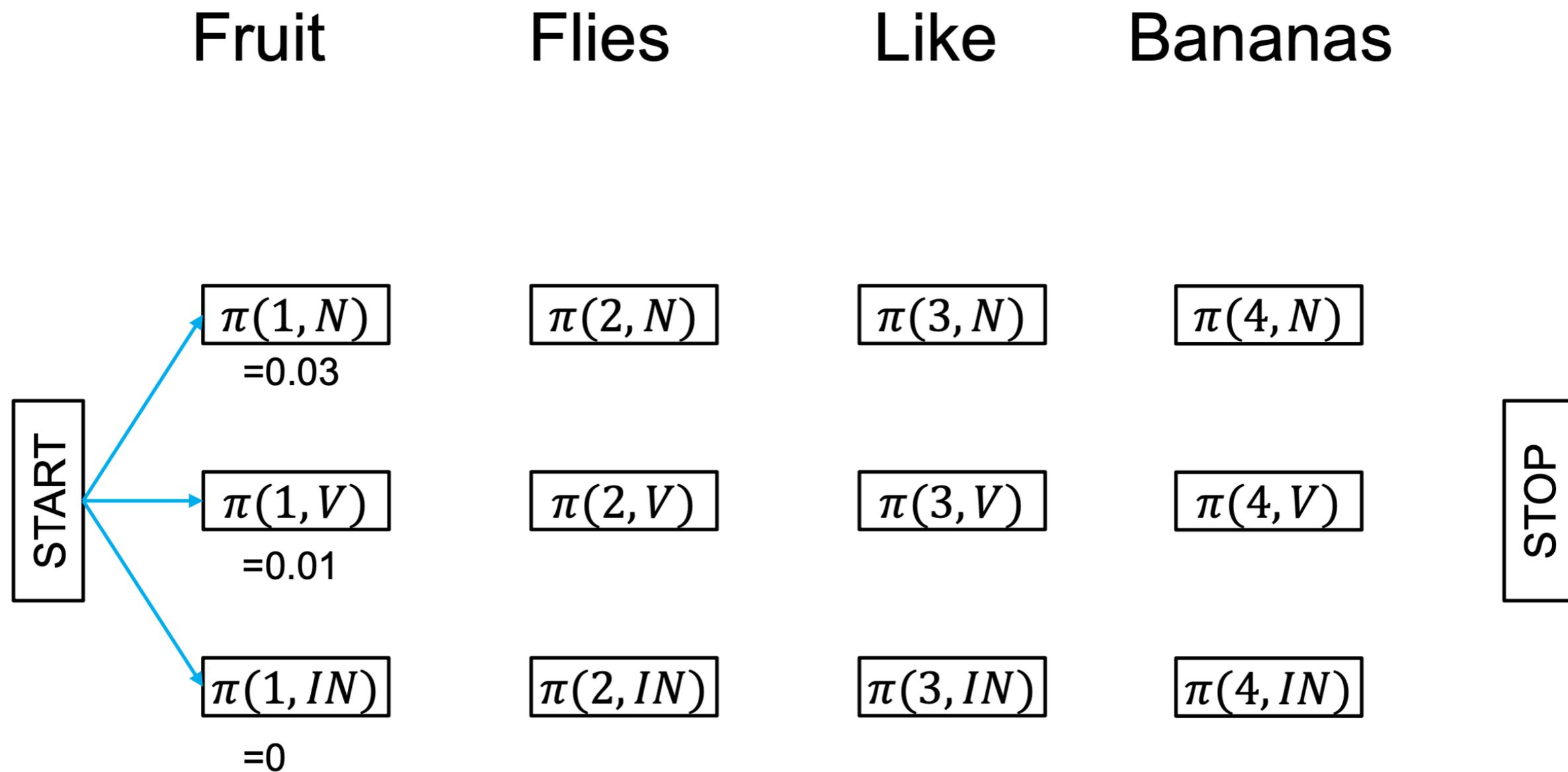
Viterbi (Example)

Fruit Flies Like Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

Viterbi (Example)



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

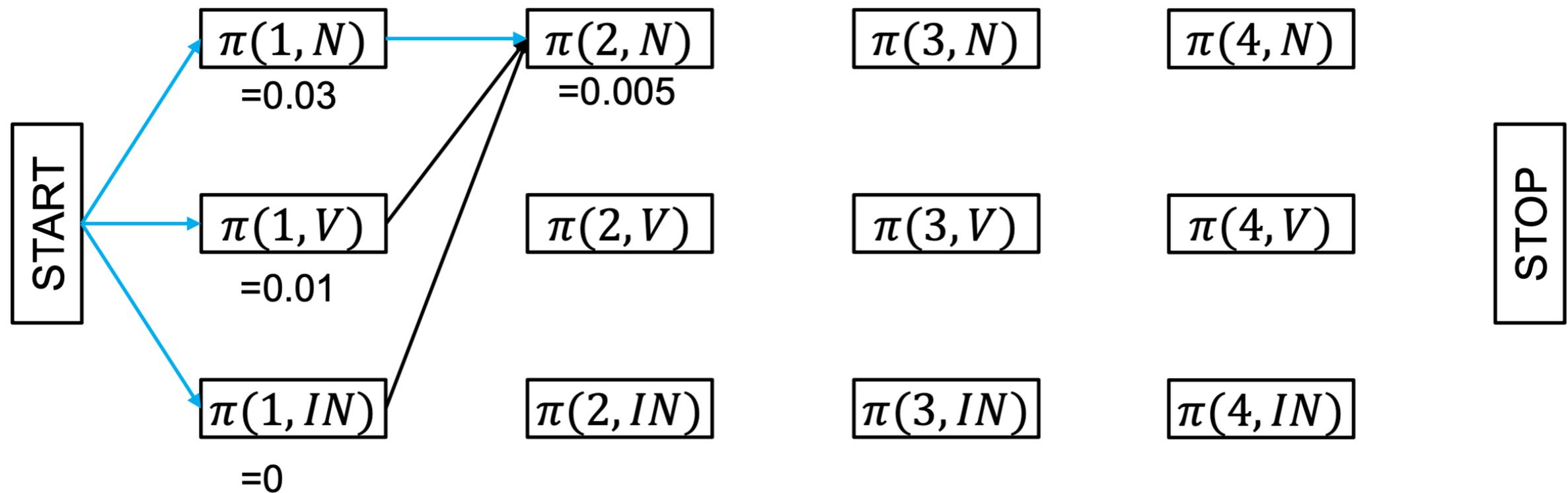
Viterbi (Example)

Fruit

Flies

Like

Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

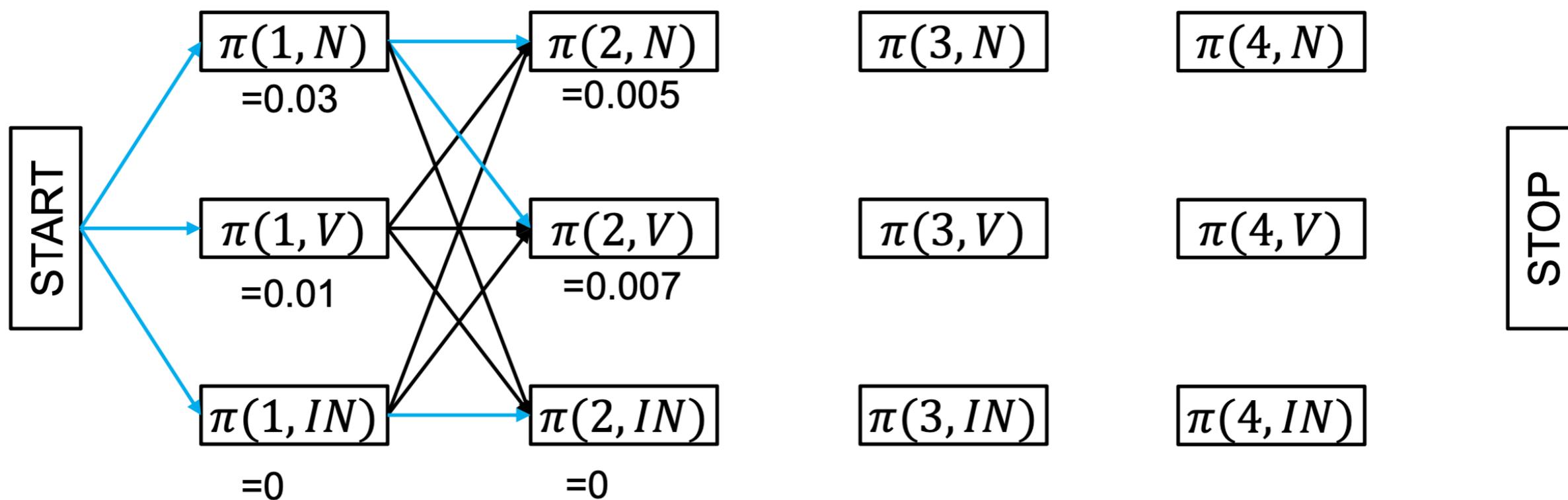
Viterbi (Example)

Fruit

Flies

Like

Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

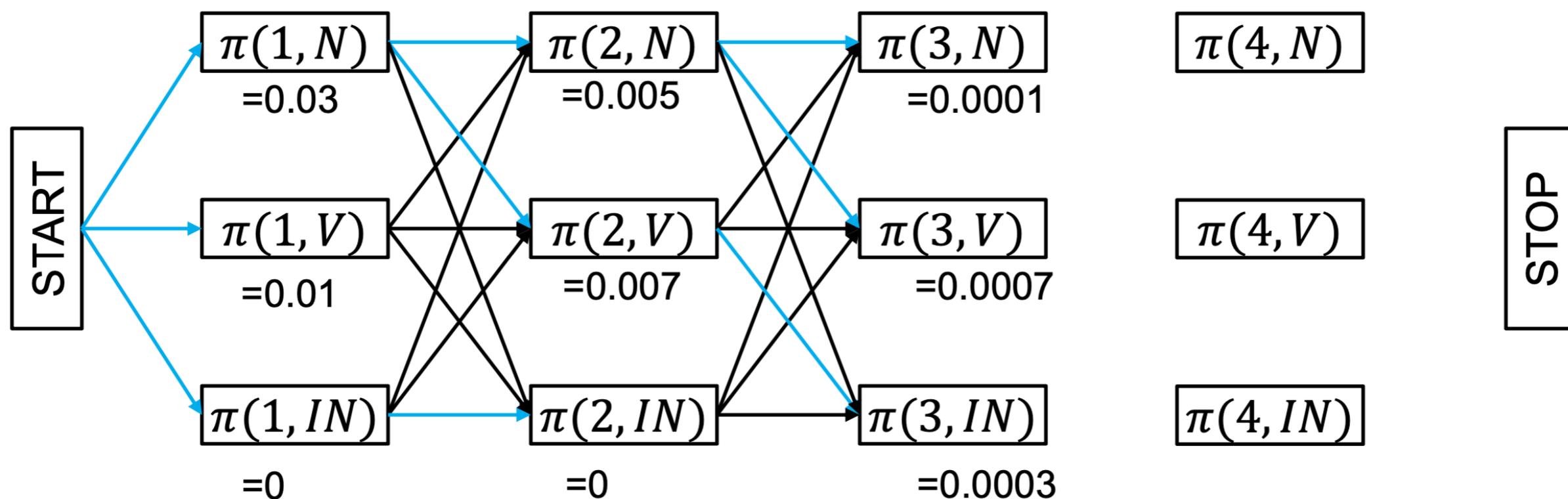
Viterbi (Example)

Fruit

Flies

Like

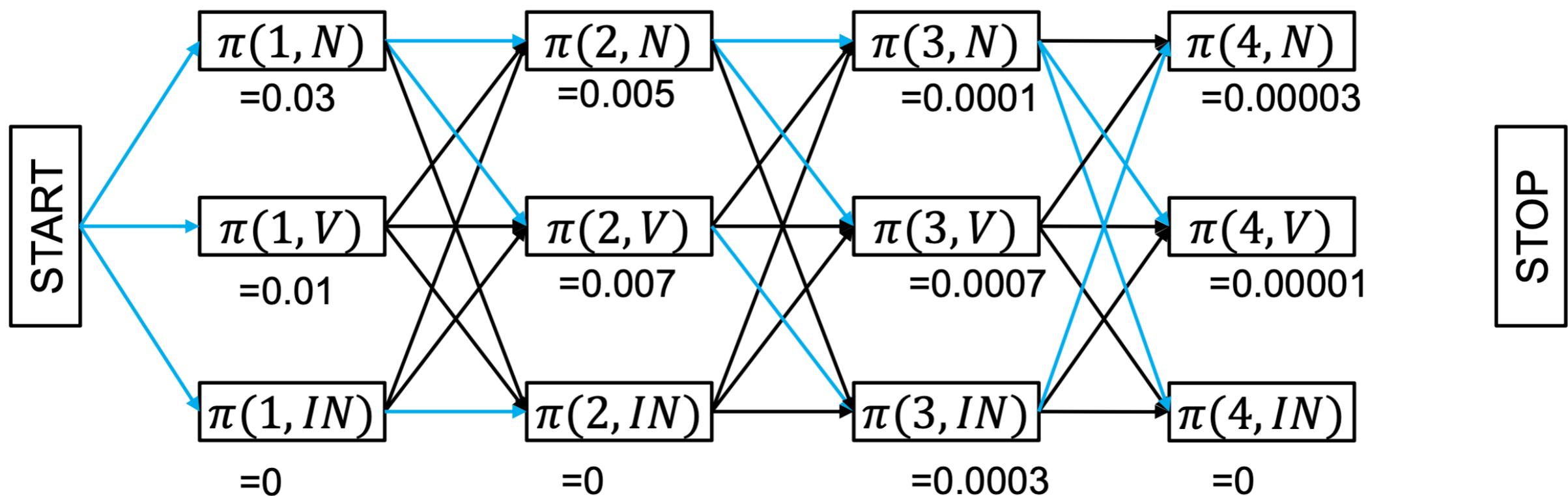
Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

Viterbi (Example)

Fruit Flies Like Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

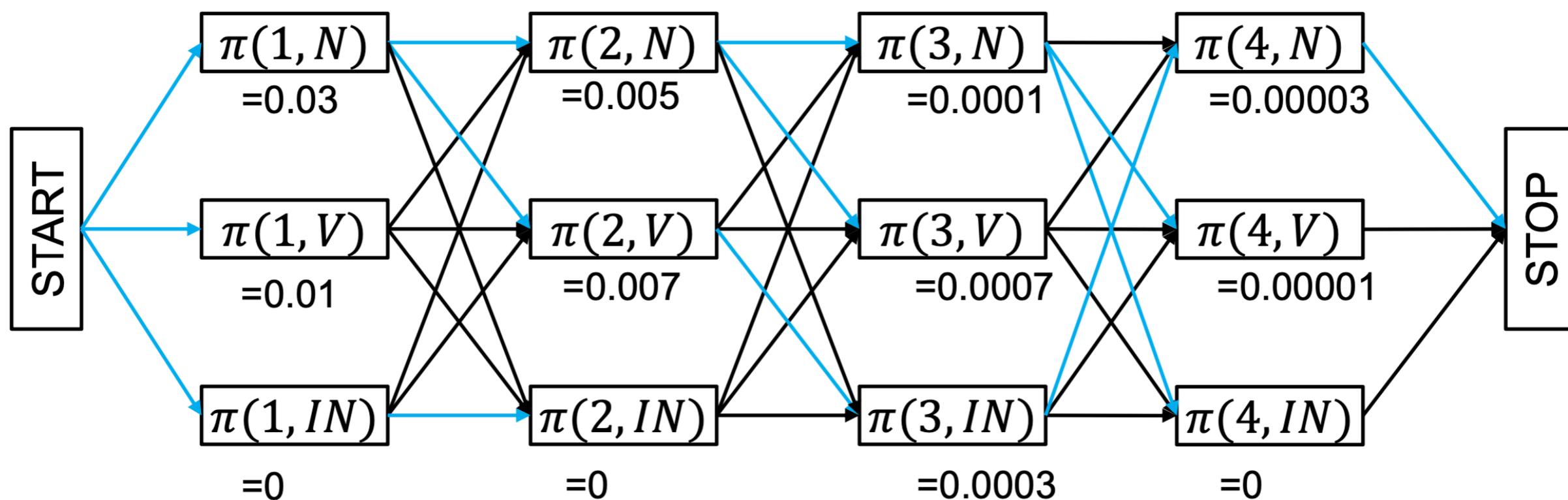
Viterbi (Example)

Fruit

Flies

Like

Bananas



$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

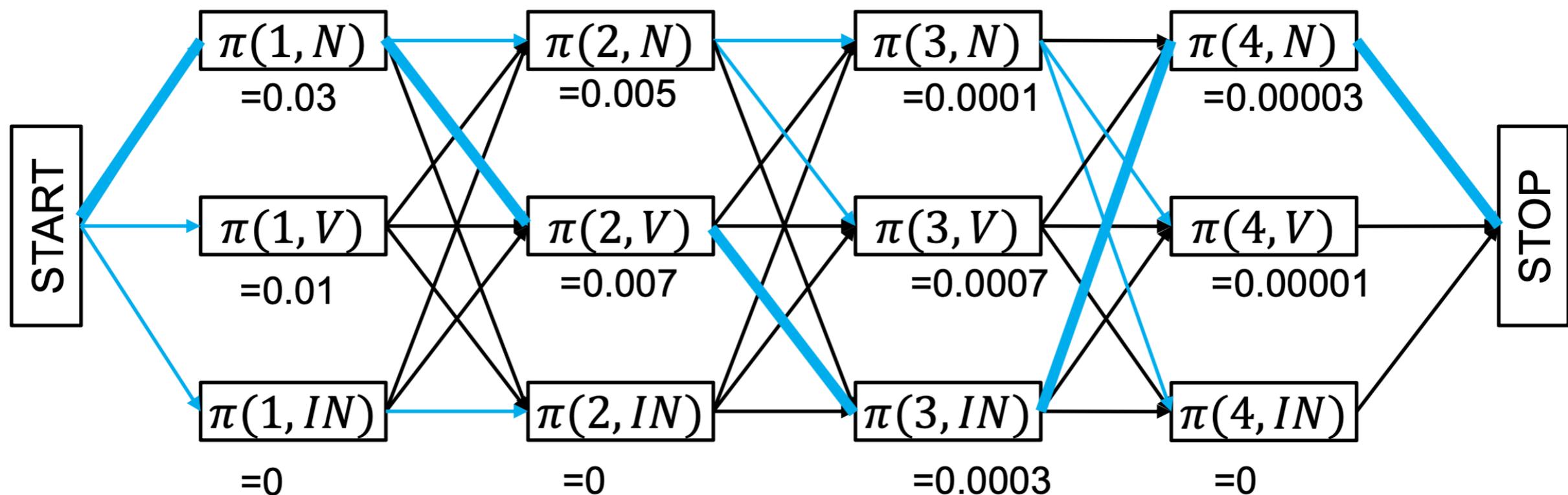
Viterbi (Example)

Fruit

Flies

Like

Bananas



$$bp(i, y_i) = \arg \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

Why is this not a greedy algorithm? Why does this find max P(.)?

Viterbi Algorithm

- Dynamic programming (for all i)

$$\pi(i, y_i) = \max_{y_1 \dots y_{i-1}} p(x_1 \dots x_i, y_1 \dots y_i)$$

- Iterative computation

$$\pi(0, y_0) = \begin{cases} 1 & \text{if } y_0 == \textit{START} \\ 0 & \text{otherwise} \end{cases}$$

For $i = 1 \dots n$:

$$\pi(i, y_i) = \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i - 1, y_{i-1})$$

- Store back pointers:

$$bp(i, y_i) = \arg \max_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i - 1, y_{i-1})$$

- What is the final solution? $bp(n + 1, \textit{STOP})$

Viterbi Algorithm: Time complexity

- Linear in sentence length n
- Polynomial in the number of possible tags \mathcal{K}

$$\pi(i, y_i) = \max_{\underline{y_{i-1}}} e(x_i | y_i) q(y_i | y_{i-1}) \pi(i-1, y_{i-1})$$

iterate over all possible tags

- Specifically:

$O(n|\mathcal{K}|)$ entries in $\pi(i, y_i)$

$O(|\mathcal{K}|)$ time to compute each $\pi(i, y_i)$

- Total runtime:

$$O(n|\mathcal{K}|^2)$$

Questions?