

# Algorithms for NLP



## Part of Speech, NER, CRF

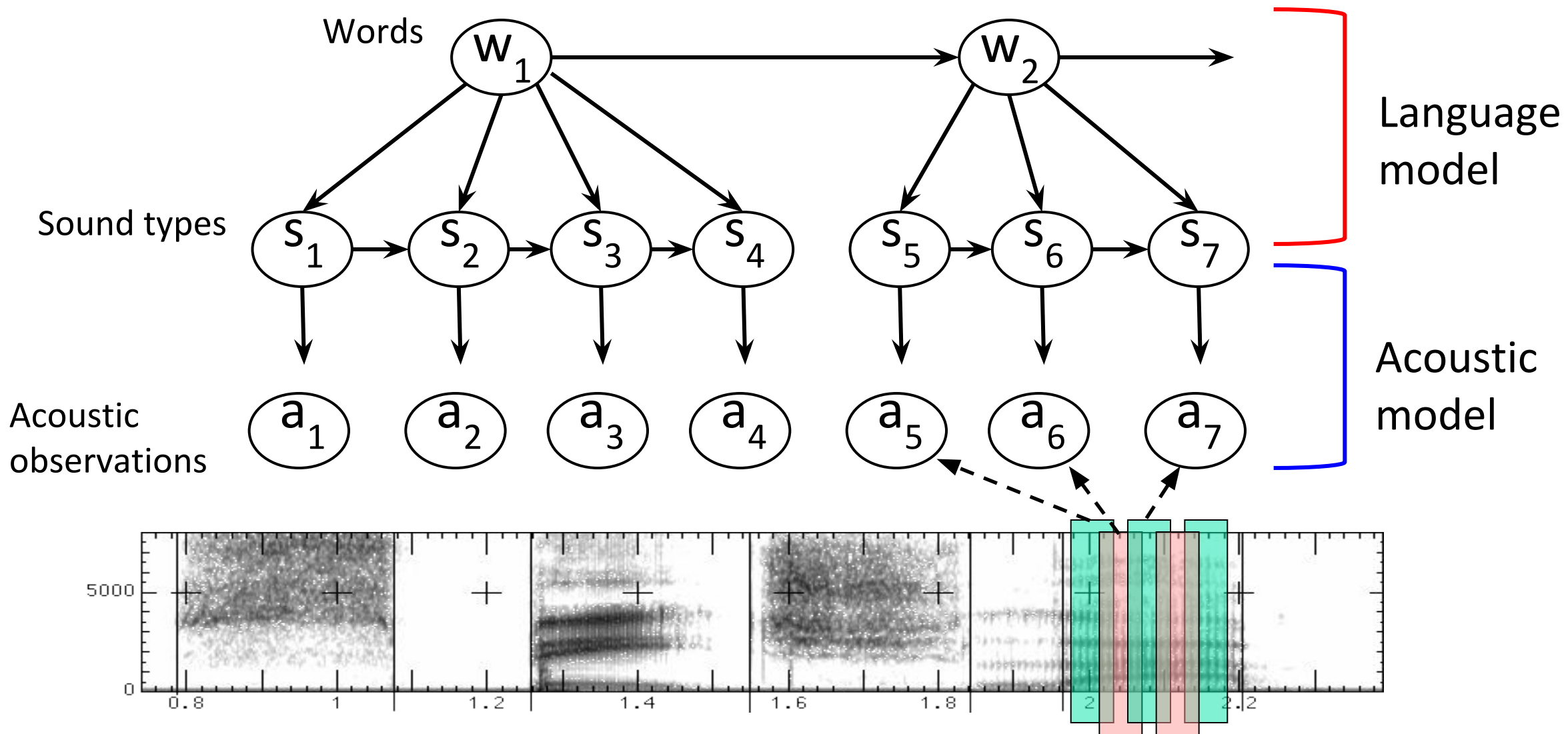
Aldrian Obaja Muis – CMU

Slides adapted from: Dan Klein – UC Berkeley

Taylor Berg-Kirkpatrick, Yulia Tsvetkov – CMU



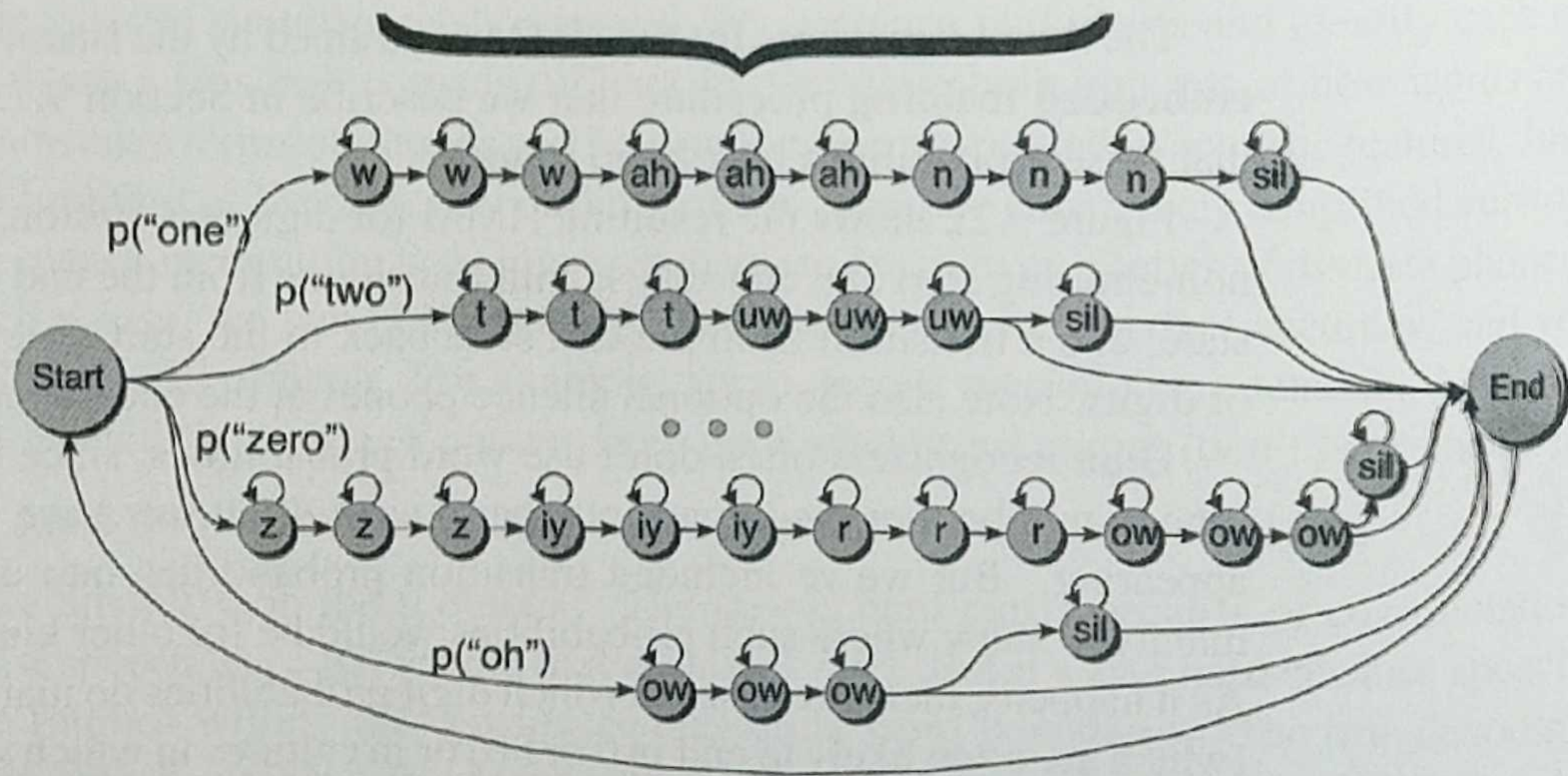
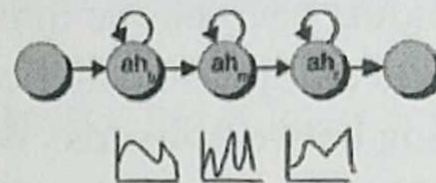
# Speech Model



### Lexicon

one	w ah n
two	t uw
three	th r iy
four	f ao r
five	f ay v
six	s ih k s
seven	s eh v ax n
eight	ey t
nine	n ay n
zero	z iy r ow
oh	ow

### Phone HMM

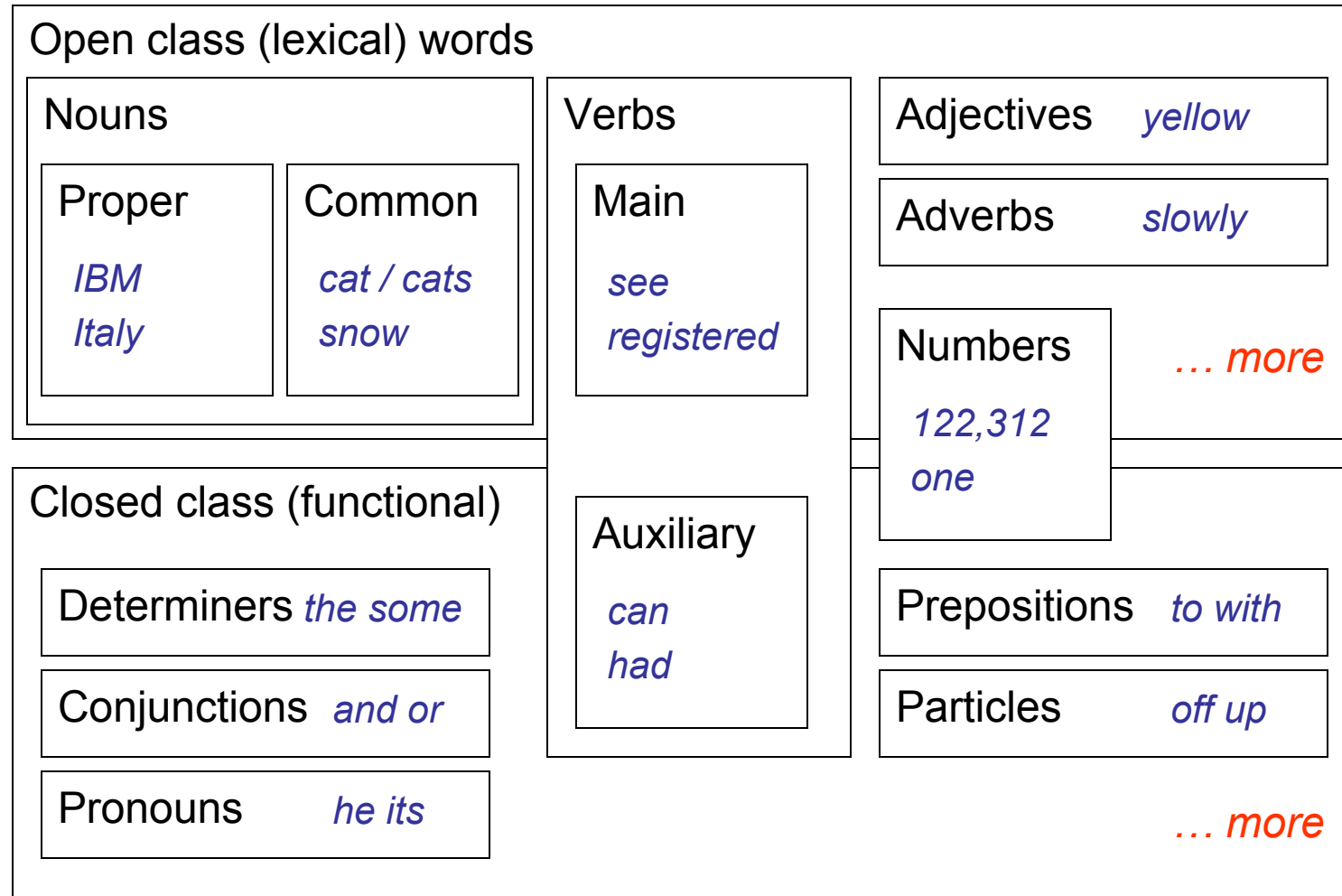


# Parts of Speech



# Parts-of-Speech (English)

- One basic kind of linguistic structure: syntactic word classes





# English Penn Treebank Part-of-Speech

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			



# Part-of-Speech in Many Languages

Language	Source	# Tags
Arabic	PADT/CoNLL07 (Hajič et al., 2004)	21
Basque	Basque3LB/CoNLL07 (Aduriz et al., 2003)	64
Bulgarian	BTB/CoNLL06 (Simov et al., 2002)	54
Catalan	CESS-ECE/CoNLL07 (Martí et al., 2007)	54
Chinese	Penn Chinese Treebank 6.0 (Palmer et al., 2007)	34
Chinese	Sinica/CoNLL07 (Chen et al., 2003)	294
Czech	PDT/CoNLL07 (Böhmová et al., 2003)	63
Danish	DDT/CoNLL06 (Kromann et al., 2003)	25
Dutch	Alpino/CoNLL06 (Van der Beek et al., 2002)	12
English	Penn Treebank (Marcus et al., 1993)	45
French	French Treebank (Abeillé et al., 2003)	30
German	Tiger/CoNLL06 (Brants et al., 2002)	54
German	Negra (Skut et al., 1997)	54
Greek	GDT/CoNLL07 (Prokopidis et al., 2005)	38
Hungarian	Szeged/CoNLL07 (Csendes et al., 2005)	43
Italian	ISST/CoNLL07 (Montemagni et al., 2003)	28
Japanese	Verbmobil/CoNLL06 (Kawata and Bartels, 2000)	80
Japanese	Kyoto4.0 (Kurohashi and Nagao, 1997)	42
Korean	Sejong ( <a href="http://www.sejong.or.kr">http://www.sejong.or.kr</a> )	187
Portuguese	Floresta Sintá(c)tica/CoNLL06 (Afonso et al., 2002)	22
Russian	SynTagRus-RNC (Boguslavsky et al., 2002)	11
Slovene	SDT/CoNLL06 (Džeroski et al., 2006)	29
Spanish	Ancora-Cast3LB/CoNLL06 (Civit and Martí, 2004)	47
Swedish	Talbanken05/CoNLL06 (Nivre et al., 2006)	41
Turkish	METU-Sabancı/CoNLL07 (Ofłazer et al., 2003)	31



# Part-of-Speech Ambiguity

- Words can have multiple parts of speech

VBD                      VB  
VBN    VBZ            VBP            VBZ  
NNP    NNS            NN            NNS    CD            NN

Fed raises interest rates 0.5 percent

Mrs./NNP Shaefer/NNP never/RB got/VBD **around/RP** to/TO joining/VBG

All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB **around/IN** the/DT corner/NN

Chateau/NNP Petrus/NNP costs/VBZ **around/RB** 250/CD

- Two basic sources of constraint:
  - Grammatical environment
  - Identity of the current word
- Many more possible features:
  - Suffixes, capitalization, name databases (gazetteers), etc...





# Why POS Tagging?

- Useful in and of itself
  - Text-to-speech: how to pronounce “record”, “lead”, “read”?
  - Lemmatization: saw[v] → see, saw[n] → saw,
  - Quick-and-dirty NP-chunk detection: `grep {JJ | NN}* {NN | NNS}`
- Useful as a pre-processing step for parsing
  - Less tag ambiguity means fewer parses
  - However, some tag choices are better decided by parsers

DT NNP NN VBD VBN **IN** RP NN NNS  
The Georgia branch had taken **on** loan commitments ...

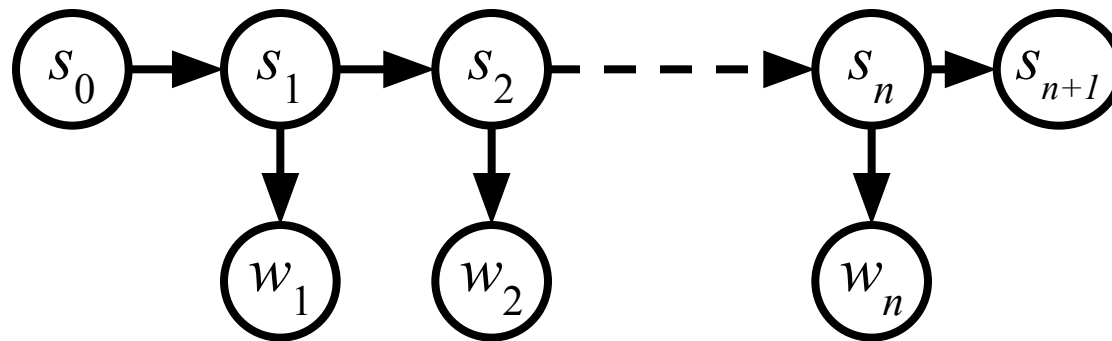
DT NN **VBN** VBD IN DT NN VBD  
The horse **raced** past the barn fell

# Part-of-Speech Tagging



# Classic Solution: HMMs

- We want a model of sequences  $s$  and observations  $w$



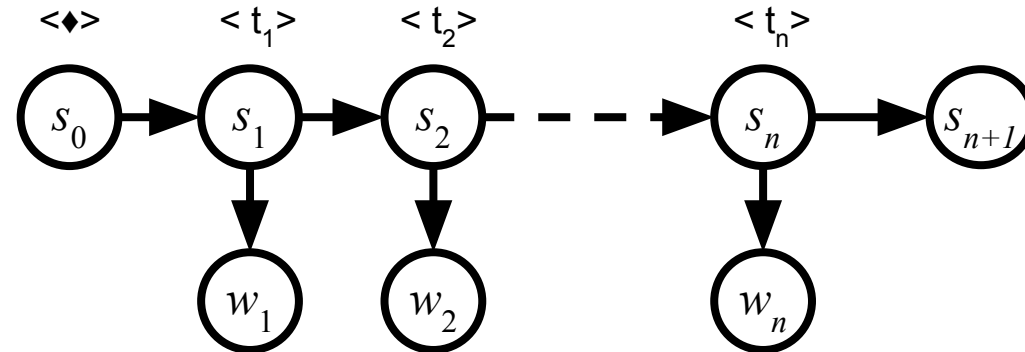
$$P(s, w) = \prod_i P(s_i | s_{i-1}) P(w_i | s_i)$$

- **Assumptions:**
  - States are tag n-grams
  - Usually add a dedicated start ( $s_0$ ) and end state ( $s_{n+1}$ )
  - 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?

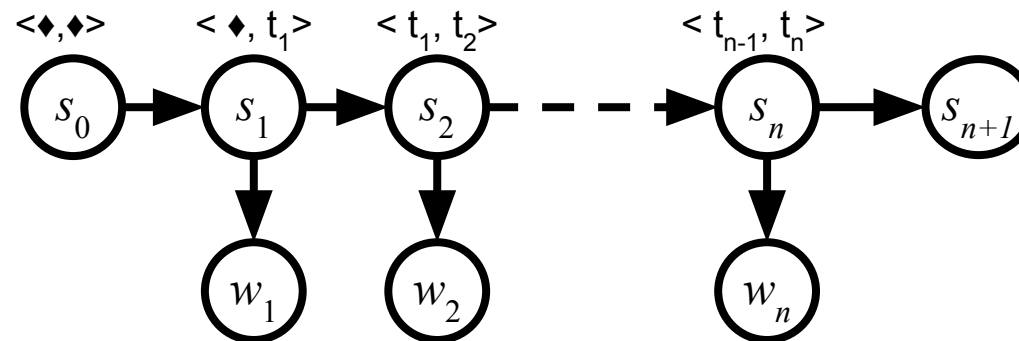


# States

- States encode what is relevant about the past
- Transitions  $P(s | s')$  encode well-formed tag sequences
  - In a bigram tagger, states = tags



- In a trigram tagger, states = tag pairs





# Estimating Transitions

---

- Use standard smoothing methods to estimate transitions:

$$P(t_i | t_{i-1}, t_{i-2}) = \lambda_2 \hat{P}(t_i | t_{i-1}, t_{i-2}) + \lambda_1 \hat{P}(t_i | t_{i-1}) + (1 - \lambda_1 - \lambda_2) \hat{P}(t_i)$$

- Can get a lot fancier (e.g. KN smoothing) or use higher orders, but in this case it doesn't buy much
- One option: encode more into the state, e.g. whether the previous word was capitalized (Brants 00) – State Splitting
- BIG IDEA: The basic approach of state-splitting / refinement turns out to be very important in a range of tasks (e.g., the states we saw in speech)



# Estimating Emissions

$$P(\mathbf{s}, \mathbf{w}) = \prod_i P(s_i | s_{i-1}) P(w_i | s_i)$$

- Emissions are trickier:

- Words we've never seen before
- Words which occur with tags we've never seen them with
- One option: break out the fancy smoothing (e.g. KN, Good-Turing)
- Issue: unknown words aren't black boxes:

343,127.23      11-year      Minteria      reintroducibly

- Basic solution: unknown words classes (affixes or shapes)

$D^+, D^+.D^+$        $D^+-x^+$        $Xx^+$        $x^+-"ly"$

- Common approach: Estimate  $P(t|w)$  and invert
- [Brants 00] used a suffix trie as its (inverted) emission model



# Disambiguation (Inference)

- Problem: find the most likely sequence under the model

$$t^* = \arg \max_t P(t|w)$$

- Given model parameters, we can score any tag sequence:  $P(t|w) \rightarrow$  likelihood

<◆,◆>	<◆,NNP>	<NNP, VBZ>	<VBZ, NN>	<NN, NNS>	<NNS, CD>	<CD, NN>	<STOP>
	NNP	VBZ	NN	NNS	CD	NN	.
	Fed	raises	interest	rates	0.5	percent	.

$P(\text{NNP}|\langle \blacklozenge, \blacklozenge \rangle) P(\text{Fed}|\text{NNP}) P(\text{VBZ}|\langle \text{NNP}, \blacklozenge \rangle) P(\text{raises}|\text{VBZ}) P(\text{NN}|\text{VBZ}, \text{NNP}) \dots$

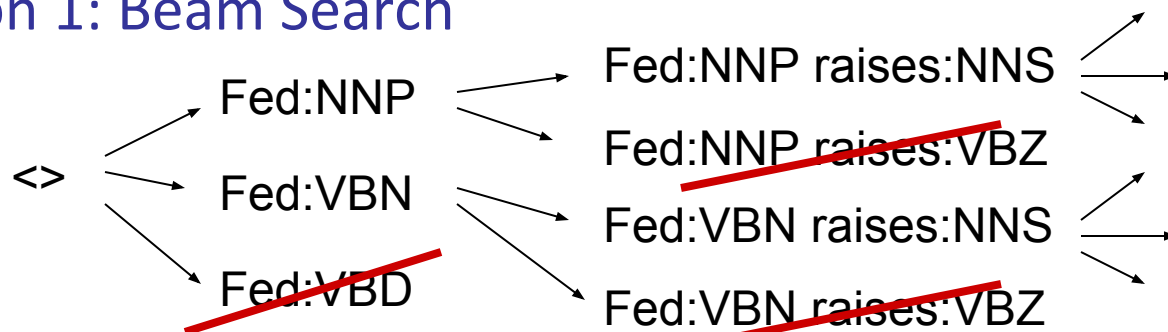
- In principle, we're done – list all possible tag sequences, score each one, pick the best one

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



# Finding the Best Trajectory

- Too many trajectories (state sequences) to list
- Option 1: Beam Search

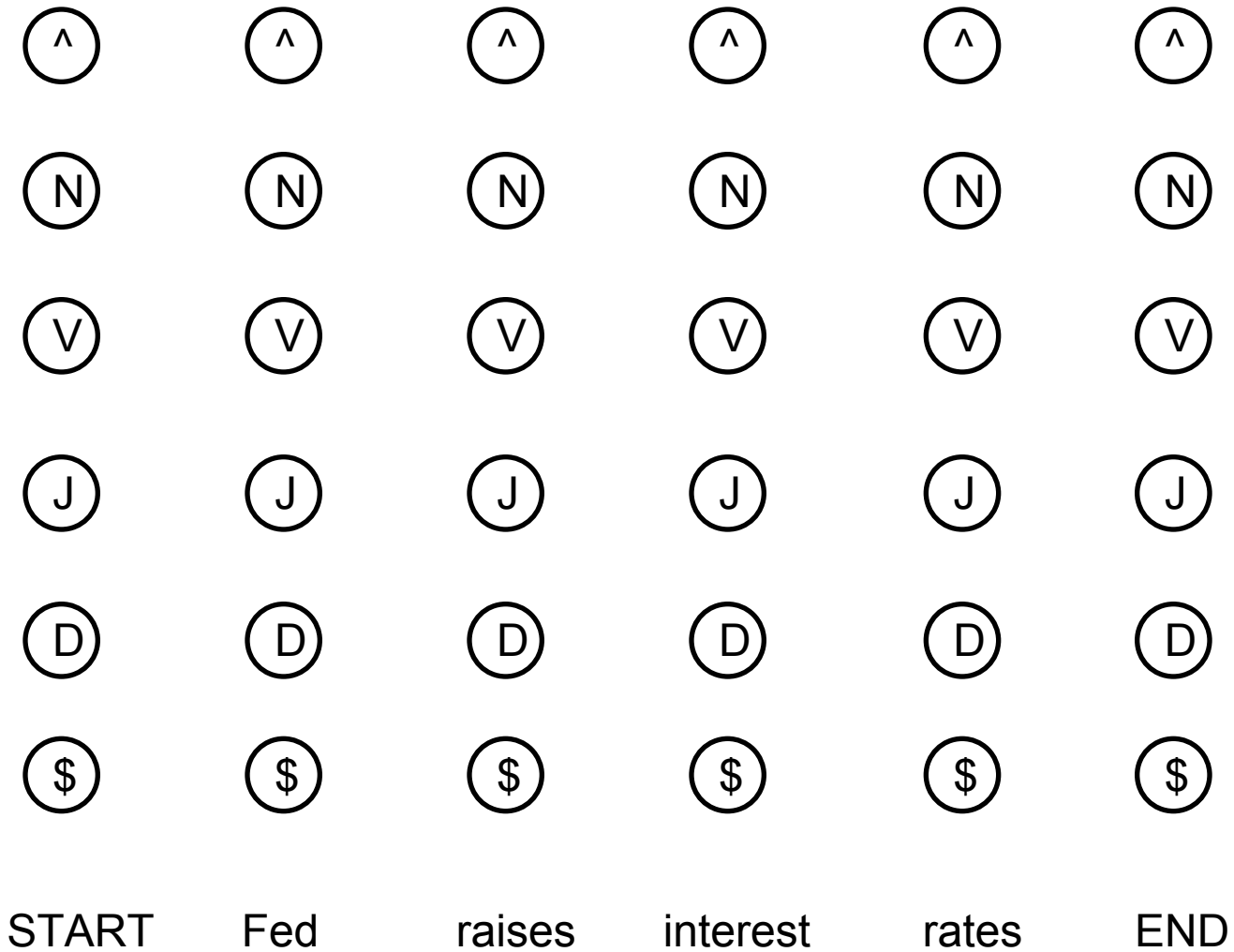


- A beam is a set of partial hypotheses
- Start with just the single empty trajectory
- At each derivation step:
  - Consider all continuations of previous hypotheses
  - Discard most, keep top k, or those within a factor of the best
- Beam search works ok in practice
  - ... but sometimes you want the optimal answer
  - ... and you need optimal answers to validate your beam search
  - ... and there's usually a better option than naïve beams



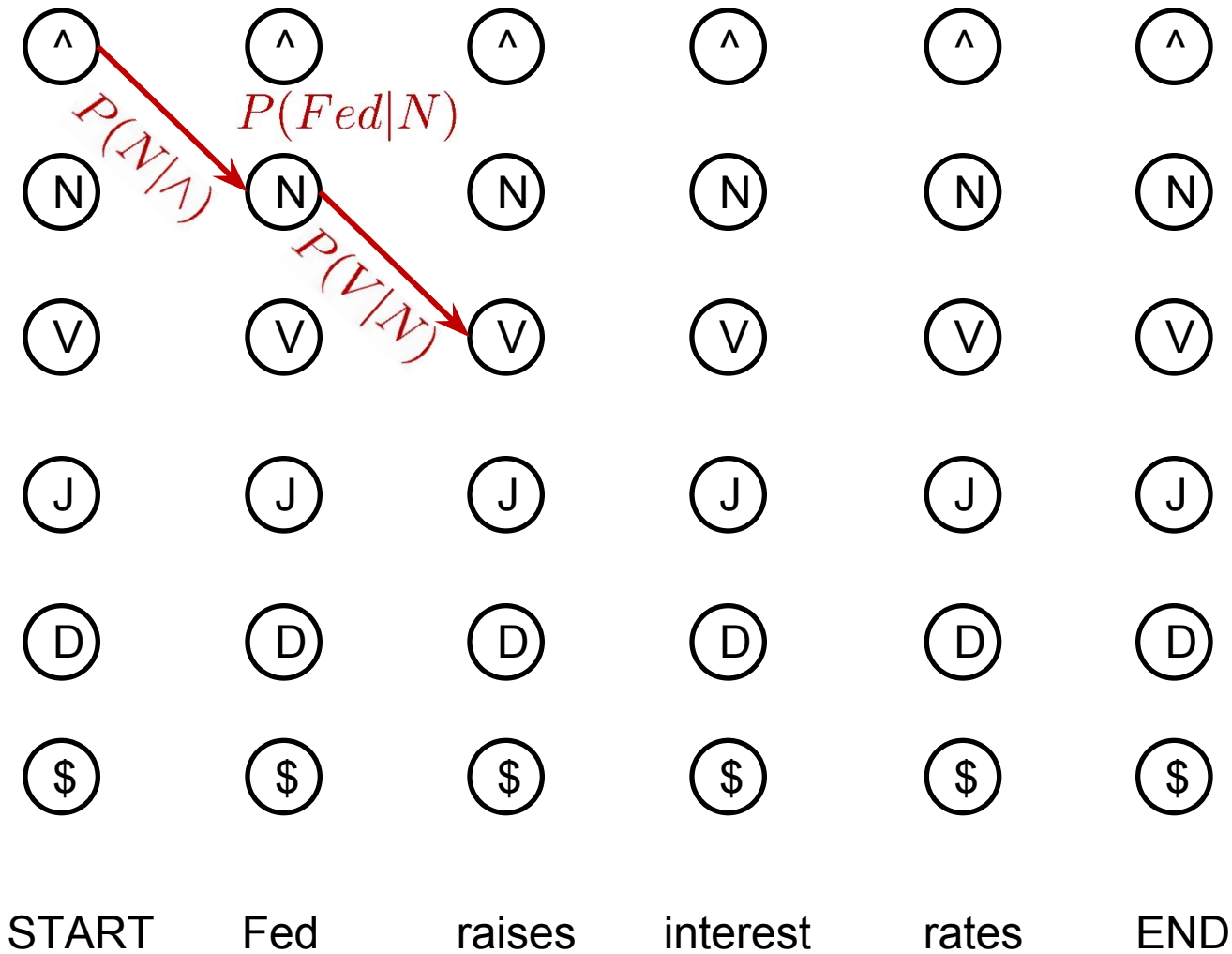


# The State Lattice / Trellis





# The State Lattice / Trellis





# The Viterbi Algorithm

- Dynamic program for computing

$$\delta_i(s) = \max_{s_0 \dots s_{i-1}} P(s_0 \dots s_{i-1} s, w_1 \dots w_{i-1})$$

- The score of a best path up to position  $i$  ending in state  $s$

$$\delta_0(s) = \begin{cases} 1 & \text{if } s = \langle \bullet, \bullet \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$\delta_i(s) = \max_{s'} P(s | s') P(w | s') \delta_{i-1}(s')$$

- Also can store a backtrace

$$\psi_i(s) = \arg \max_{s'} P(s | s') P(w | s') \delta_{i-1}(s')$$

- Memoized solution
- Iterative solution



# So How Well Does It Work?

- Choose the most common tag
    - 90.3% with a bad unknown word model
    - 93.7% with a good one
  - TnT (Brants, 2000):
    - A carefully smoothed trigram tagger
    - Suffix trees for emissions
    - 96.7% on WSJ text (SOTA is 97+%)
  - Noise in the data
    - Many errors in the training and test corpora
- DT NN IN NN VBD NNS VBD  
 The average of interbank offered rates plummeted ...
- Probably about 2% guaranteed error from noise (on this data)

JJ JJ NN  
 chief executive officer

NN JJ NN  
 chief executive officer

JJ NN NN  
 chief executive officer

NN NN NN  
 chief executive officer



# Overview: Accuracies

- Roadmap of (known / unknown) accuracies:

- Most freq tag: ~90% / ~50%
- Trigram HMM: ~95% / ~55%
- TnT (HMM++): 96.5% / 85.9%
- Maxent  $P(t|w)$ : 93.7% / 82.6%
- MEMM tagger: 96.9% / 86.9%
- State-of-the-art: 97+% / 92+%
- Upper bound: ~98%

Most errors are  
on unknown  
words



# Common Errors

- Common errors [from Toutanova & Manning 00]

	JJ	NN	NNP	NNPS	RB	RP	IN	VB	VBD	VCN	VBP	Total
JJ	0	177	56	0	61	2	5	10	15	108	0	488
NN	244	0	103	0	12	1	1	29	5	6	19	525
NNP	107	106	0	132	5	0	7	5	1	2	0	427
NNPS	1	0	110	0	0	0	0	0	0	0	0	142
RB	72	21	7	0	0	16	138	1	0	0	0	295
RP	0	0	0	0	39	0	65	0	0	0	0	104
IN	11	0	1	0	169	103	0	1	0	0	0	323
VB	17	64	9	0	2	0	1	0	4	7	85	189
VBD	10	5	3	0	0	0	0	3	0	143	2	166
VCN	101	3	3	0	0	0	0	3	108	0	1	221
VBP	5	34	3	1	1	0	2	49	6	3	0	104
Total	626	536	348	144	317	122	279	102	140	269	108	3651

NN/JJ NN/JJ NN

US Navy officer

VBD RP/IN DT NN

got around the corner

PRP RB VBD/VBN NNS IN ...

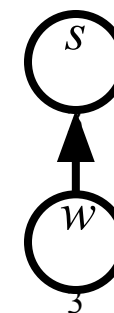
They recently sold shares at ...

# Richer Features



# Better Features

- Can do surprisingly well just looking at a word by itself:
  - Word                    the: the → DT
  - Lowercased word       Importantly: importantly → RB
  - Prefixes                unfathomable: un- → JJ
  - Suffixes                Surprisingly: -ly → RB
  - Capitalization         Meridian: CAP → NNP
  - Word shapes            35-year: d-x → JJ
- Then build a maxent (or whatever) model to predict tag
- Maxent  $P(t|w)$ :    93.7% / 82.6%







# Why Linear Context is Useful

---

- Lots of rich local information!

RB  
PRP VBD IN RB IN PRP VBD .  
They left as soon as he arrived .

- We could fix this with a feature that looked at the next word

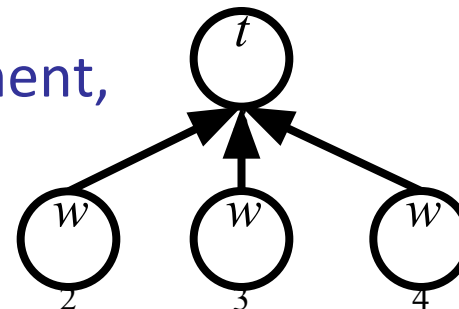
JJ  
NNP NNS VBD VBN .  
Intrinsic flaws remained undetected .

- We could fix this by linking capitalized words to their lowercase versions
- Solution: discriminative sequence models (MEMMs, CRFs)
- Reality check:
  - Taggers are already pretty good on newswire text...
  - What the world needs is taggers that work on other text!



# Sequence-Free Tagging?

- What about looking at a word and its environment, but no sequence information?



- Add in previous / next word the \_\_\_
  - Previous / next word shapes X \_\_\_ X
  - Occurrence pattern features [X: x X occurs]
  - Crude entity detection \_\_\_ ..... (Inc. | Co.)
  - Phrasal verb in sentence? put ..... \_\_\_
  - Conjunctions of these things
- 
- All features but no sequence: 96.6% / 86.8%
  - Uses lots of features: > 200K
  - Why isn't this the standard approach?



# Named Entity Recognition

- Other sequence tasks use similar models
- Example: name entity recognition (NER)

PER PER O O O O O ORG O O O O O LOC LOC O

Tim Boon has signed a contract extension with Leicestershire which will keep him at Grace Road .

## Local Context

	Prev	Cur	Next
State	Other	???	???
Word	at	Grace	Road
Tag	IN	NNP	NNP
Sig	x	Xx	Xx



# MEMM Taggers

---

- Idea: left-to-right local decisions, condition on previous tags and also entire input

$$P(\mathbf{t}|\mathbf{w}) = \prod_i P_{ME}(t_i|\mathbf{w}, t_{i-1}, t_{i-2})$$

- Train  $P(t_i|\mathbf{w}, t_{i-1}, t_{i-2})$  as a normal maxent model, then use to score sequences
  - This is referred to as an MEMM tagger [Ratnaparkhi 96]
  - Beam search effective! (Why?)
  - What about beam size 1?
- 
- Issues with local normalization, called “Label Bias Problem” (cf. Lafferty et al. 2001 and many others)
  - Should no longer be used in practice (but neural models often do this)



# NER Features

Because of regularization term, the more common prefixes can have larger weights even though entire-word features are more specific.

## Feature Weights

Feature Type	Feature	PERS	LOC
Previous word	<i>at</i>	-0.73	0.94
Current word	<i>Grace</i>	0.03	0.00
Beginning bigram	<i>&lt;G</i>	0.45	-0.04
Current POS tag	<i>NNP</i>	0.47	0.45
Prev and cur tags	<i>IN NNP</i>	-0.10	0.14
Previous state	<i>Other</i>	-0.70	-0.92
Current signature	<i>Xx</i>	0.80	0.46
Prev state, cur sig	<i>O-Xx</i>	0.68	0.37
Prev-cur-next sig	<i>x-Xx-Xx</i>	-0.69	0.37
P. state - p-cur sig	<i>O-x-Xx</i>	-0.20	0.82
...			
<b>Total:</b>		<b>-0.58</b>	<b>2.68</b>

## Local Context

	Prev	Cur	Next
State	<i>Other</i>	<i>???</i>	<i>???</i>
Word	<i>at</i>	<i>Grace</i>	<i>Road</i>
Tag	<i>IN</i>	<i>NNP</i>	<i>NNP</i>
Sig	<i>x</i>	<i>Xx</i>	<i>Xx</i>



# Decoding

- Decoding MEMM taggers:

- Just like decoding HMMs, different local scores
- Viterbi, beam search, posterior decoding

- Viterbi algorithm (HMMs):

$$\delta_i(s) = \arg \max_{s'} P(s|s')P(w_{i-1}|s')\delta_{i-1}(s')$$

- Viterbi algorithm (MEMMs):

$$\delta_i(s) = \arg \max_{s'} P(s|s', \mathbf{w})\delta_{i-1}(s')$$

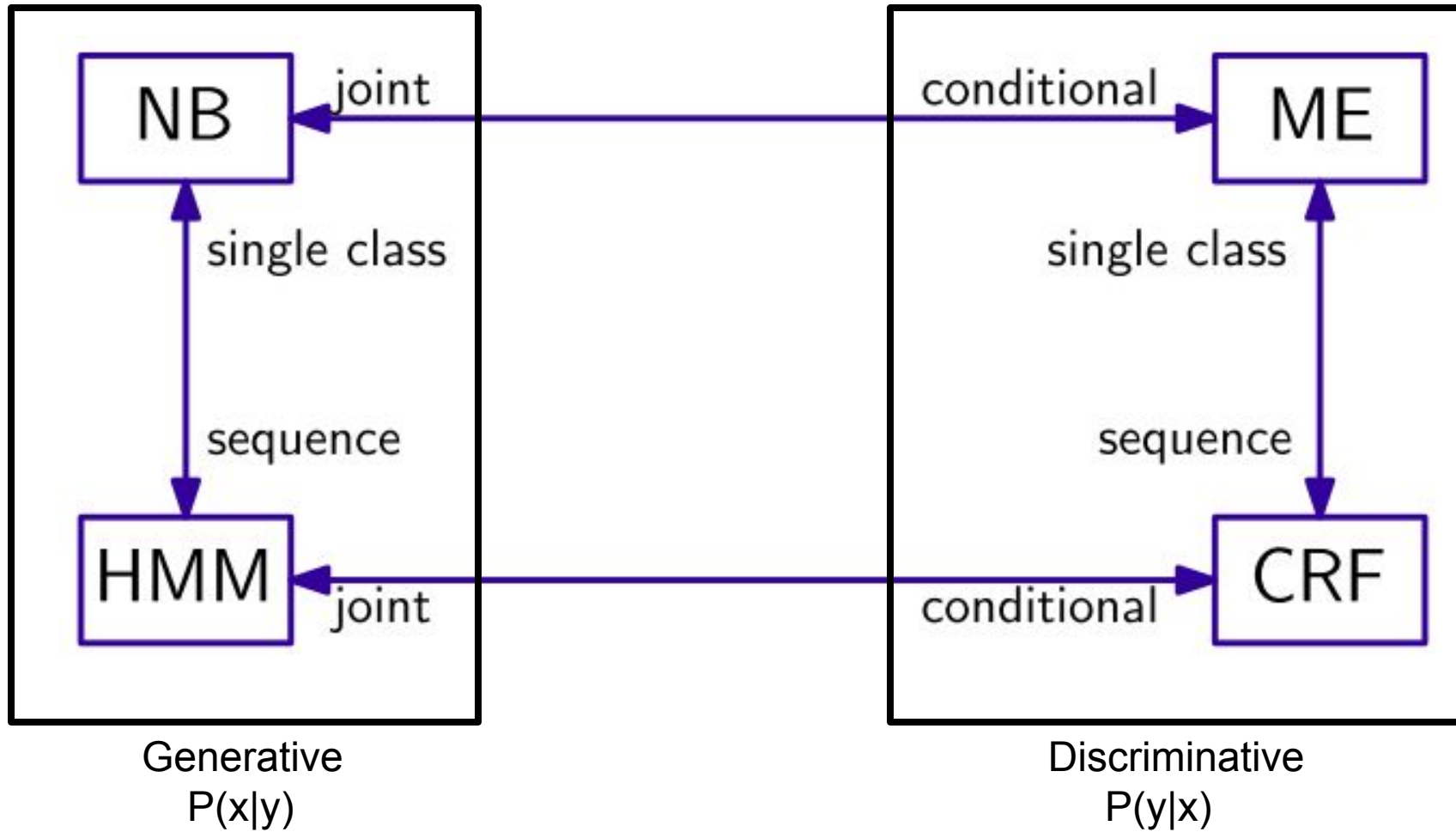
- General:

$$\delta_i(s) = \arg \max_{s'} \phi_i(s', s)\delta_{i-1}(s')$$

# Conditional Random Fields (and Friends)



# HMM and Other Probabilistic Models







# Generative vs Discriminative

---

- **Generative Models ( $P(x|y)$ ):**

- Try to approximate the process that generates the observation
  - E.g., the topic of a document determines the distribution of the words
- Turns out this can be used for classification as well by inverting the probability
- Examples: Naïve Bayes, HMM, GMM, VAE
- Focus: modeling data generation process

- **Discriminative Models ( $P(y|x)$ ):**

- If we are interested just in the tags/classes, try to model it directly
  - E.g., the words present in a document distinguish (discriminate) between different topics
- Examples: Perceptron, MEMM, SVM, CRF
- Focus: creating decision boundaries



# Perceptron Taggers

- Linear models:

$$\text{score}(\mathbf{t}|\mathbf{w}) = \lambda^\top f(\mathbf{t}, \mathbf{w})$$

- ... that decompose along the sequence

$$= \lambda^\top \sum_i f(t_i, t_{i-1}, \mathbf{w}, i)$$

- ... allow us to predict with the Viterbi algorithm

$$\mathbf{t}^* = \arg \max_{\mathbf{t}} \text{score}(\mathbf{t}|\mathbf{w})$$

- ... which means we can train with the perceptron algorithm (or related updates, like MIRA)



# Perceptron Algorithm

---

## Algorithm 1 Multi-class Perceptron algorithm

---

**Require:** Training data:  $\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|T|}$

- 1:  $\mathbf{w} = \mathbf{0}$  ▷ or small random vectors
- 2: **for**  $n := 1$  **to**  $N$  **do** ▷ number of iterations
- 3:     **for**  $t := 1$  **to**  $|T|$  **do** ▷ number of training instances
- 4:          $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y})$  ▷ best output
- 5:         **if**  $\mathbf{y}' \neq \mathbf{y}_t$  **then** ▷ if incorrect
- 6:              $\mathbf{w} = \mathbf{w} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t + \mathbf{y}')$  ▷ update
- 7:         **end if**
- 8:     **end for**
- 9: **end for**

---



# Transformation-Based Learning

---

- [Brill 95] presents a *transformation-based* tagger
  - Label the training set with most frequent tags

DT MD VBD VBD .  
The can was rusted .
  - Add transformation rules which reduce training mistakes
    - MD → NN : DT \_\_
    - VBD → VBN : VBD \_\_ .
  - Stop when no transformations do sufficient good
  - Does this remind anyone of anything?
- ~~Probably the most widely used tagger (esp. outside NLP)~~
- ... but definitely not the most accurate: 96.6% / 82.0 %



# Learned Transformations

- What gets learned? [from Brill 95]

#	Change Tag		Condition
	From	To	
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

#	Change Tag		Condition
	From	To	
1	NN	NNS	Has suffix <b>-s</b>
2	NN	CD	Has character <b>.</b>
3	NN	JJ	Has character <b>-</b>
4	NN	VBN	Has suffix <b>-ed</b>
5	NN	VBG	Has suffix <b>-ing</b>
6	??	RB	Has suffix <b>-ly</b>
7	??	JJ	Adding suffix <b>-ly</b> results in a word.
8	NN	CD	The word <b>\$</b> can appear to the left.
9	NN	JJ	Has suffix <b>-al</b>
10	NN	VB	The word <b>would</b> can appear to the left.
11	NN	CD	Has character <b>0</b>
12	NN	JJ	The word <b>be</b> can appear to the left.
13	NNS	JJ	Has suffix <b>-us</b>
14	NNS	VBZ	The word <b>it</b> can appear to the left.
15	NN	JJ	Has suffix <b>-ble</b>
16	NN	JJ	Has suffix <b>-ic</b>
17	NN	CD	Has character <b>1</b>
18	NNS	NN	Has suffix <b>-ss</b>
19	??	JJ	Deleting the prefix <b>un-</b> results in a word
20	NN	JJ	Has suffix <b>-ive</b>



# Maximum Entropy++

---

- Remember: maximum entropy objective

$$L(\mathbf{w}) = \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

- Problem: lots of features allow perfect fit to training set
- Regularization (compare to smoothing)

$$\max_{\mathbf{w}} \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right) - k \|\mathbf{w}\|^2$$



# Derivative for Maximum Entropy

$$L(\mathbf{w}) = -k\|\mathbf{w}\|^2 + \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial w_n} = -2kw_n + \sum_i \left( \mathbf{f}_i(\mathbf{y}^i)_n - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}_i(\mathbf{y})_n \right)$$

Big weights are bad

Total count of feature n  
in correct candidates

Expected count of  
feature n in predicted  
candidates



# Global Discriminative Taggers

---

- Newer, higher-powered discriminative sequence models
  - CRFs (also perceptrons, M3Ns)
  - Do not decompose training into independent local regions
  - Can be deathly slow to train – require repeated inference on training set
- Differences tend not to be too important for POS tagging
- Differences more substantial on other sequence tasks
- However: one issue worth knowing about in local models
  - “Label bias” and other explaining away effects
  - MEMM taggers’ local scores can be near one without having both good “transitions” and “emissions”
  - This means that often evidence doesn’t flow properly
  - Why isn’t this a big deal for POS tagging?
  - Also: in decoding, condition on predicted, not gold, histories





# Conditional Random Fields (CRFs)

---

- CRF: Make a maxent model over entire taggings

- MEMM

$$P(\mathbf{t}|\mathbf{w}) = \prod_i \frac{1}{Z(i)} \exp \left( \lambda^\top f(t_i, t_{i-1}, \mathbf{w}, i) \right)$$

- CRF

$$\begin{aligned} P(\mathbf{t}|\mathbf{w}) &= \frac{1}{Z(\mathbf{w})} \exp \left( \lambda^\top f(\mathbf{t}, \mathbf{w}) \right) \\ &= \frac{1}{Z(\mathbf{w})} \exp \left( \lambda^\top \sum_i f(t_i, t_{i-1}, \mathbf{w}, i) \right) \\ &= \frac{1}{Z(\mathbf{w})} \prod_i \phi_i(t_i, t_{i-1}) \end{aligned}$$



# Conditional Random Fields (CRFs)

- Like any maxent model, derivative is:

$$\frac{\partial L(\lambda)}{\partial \lambda} = \sum_k \left( \mathbf{f}_k(\mathbf{t}^k) - \sum_{\mathbf{t}} P(\mathbf{t}|\mathbf{w}_k) \mathbf{f}_k(\mathbf{t}) \right)$$

- So all we need is to be able to compute the expectation of each feature (for example the number of times the label pair *DT-NN* occurs, or the number of times *NN-interest* occurs) **under the model distribution**
- Critical quantity: counts of posterior marginals:

$$\text{count}(w, s) = \sum_{i:w_i=w} P(t_i = s|\mathbf{w})$$

$$\text{count}(s \rightarrow s') = \sum_i P(t_{i-1} = s, t_i = s'|\mathbf{w})$$

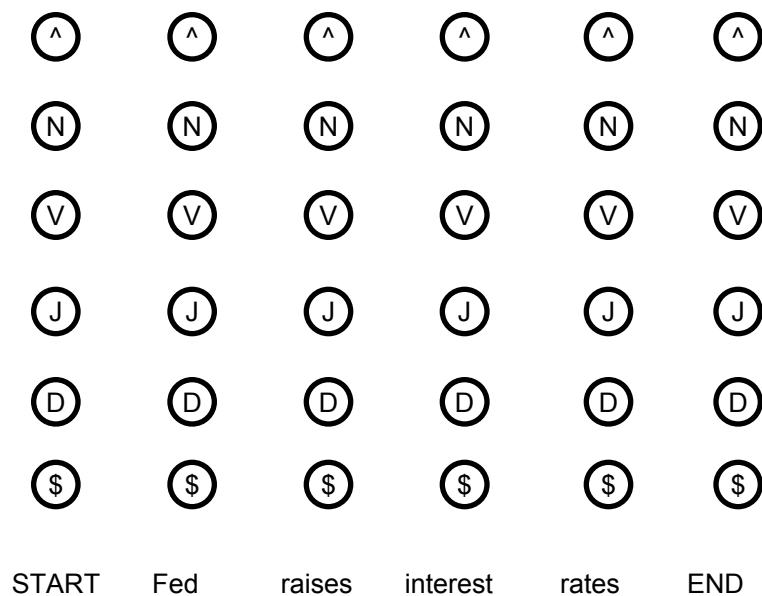


# Computing Posterior Marginals

- How many (expected) times is word  $w$  tagged with  $s$ ?

$$\text{count}(w, s) = \sum_{i:w_i=w} P(t_i = s | \mathbf{w})$$

- How to compute that marginal?



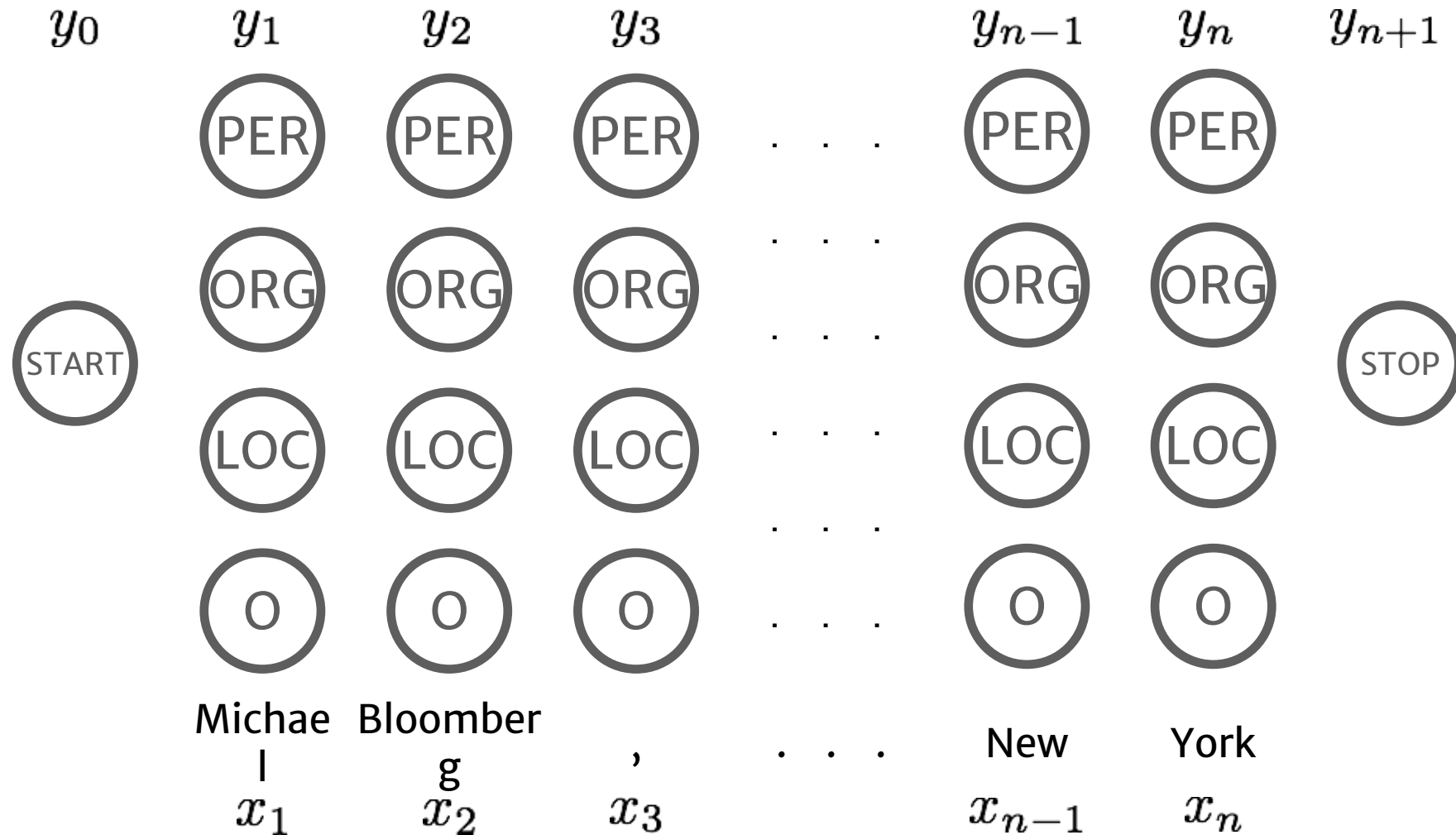
$$\alpha_i(s) = \sum_{s'} \phi_i(s', s) \alpha_{i-1}(s')$$

$$\beta_i(s) = \sum_{s'} \phi_{i+1}(s, s') \beta_{i+1}(s')$$

$$P(t_i = s | \mathbf{w}) = \frac{\alpha_i(s) \beta_i(s)}{\alpha_N(\text{END})}$$

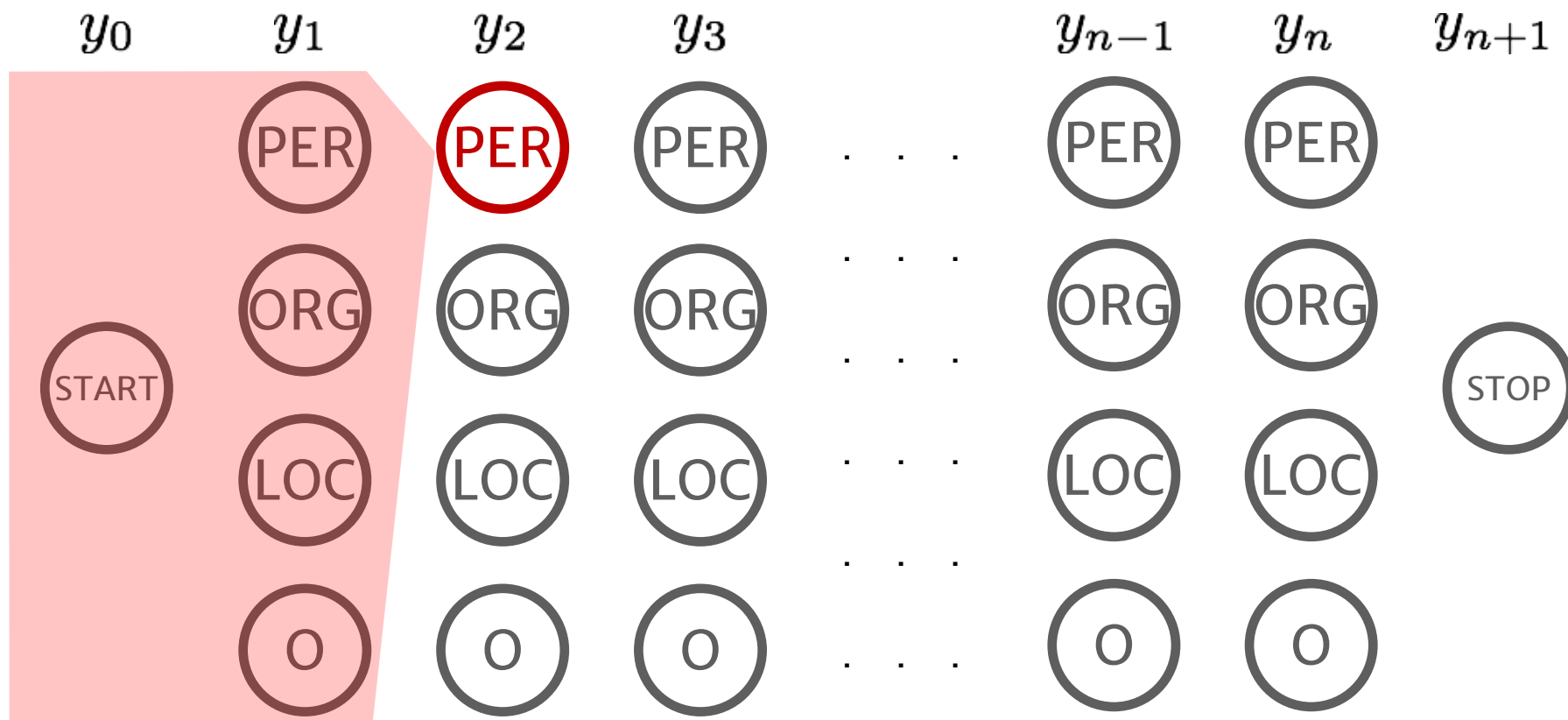


# Dynamic programming





# Forward pass

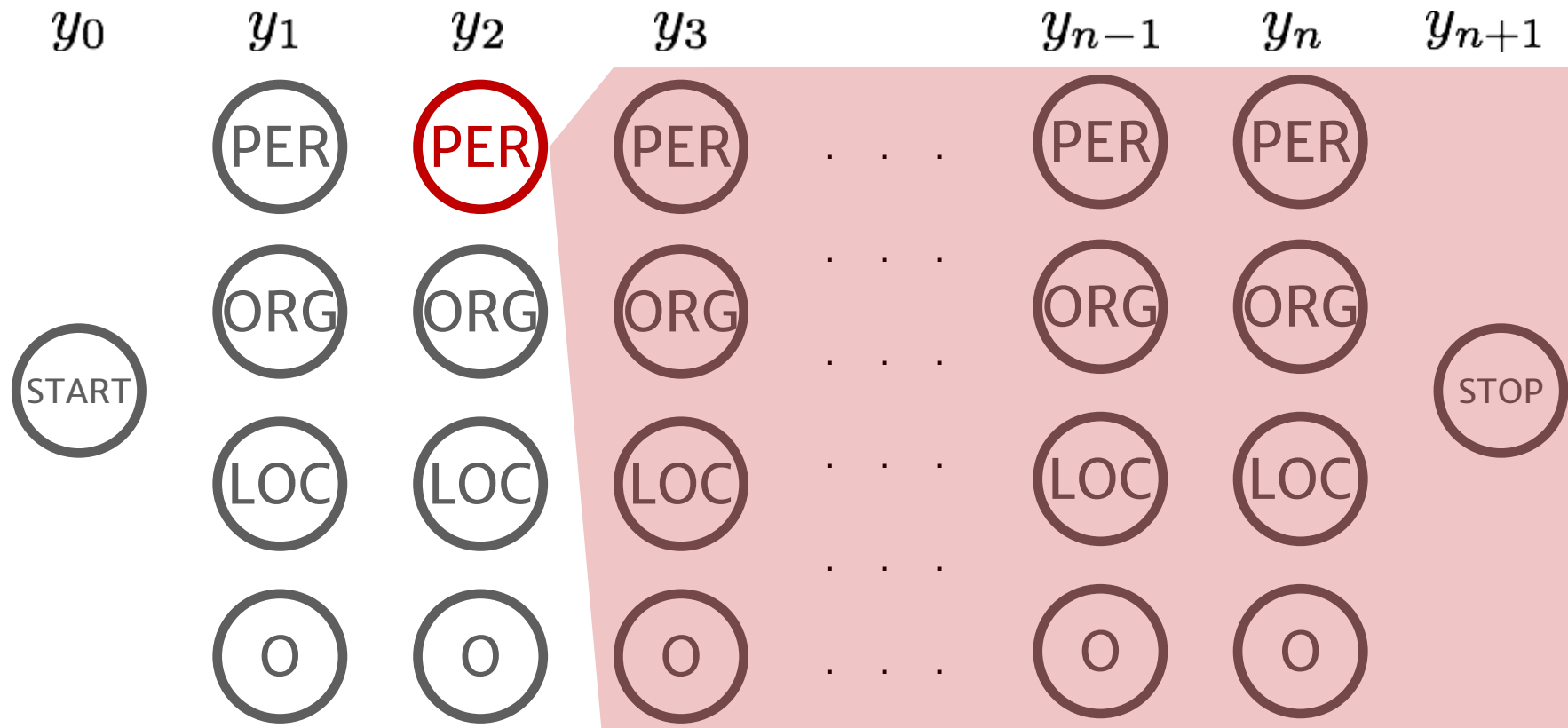


$$\alpha_t(s) = \sum_{s'} \exp(w^\top f(x, y_t = s, y_{t-1} = s')) \alpha_{t-1}(s')$$

$$\alpha_0(s) = \mathbb{1}[s = \text{START}]$$



# Backward pass

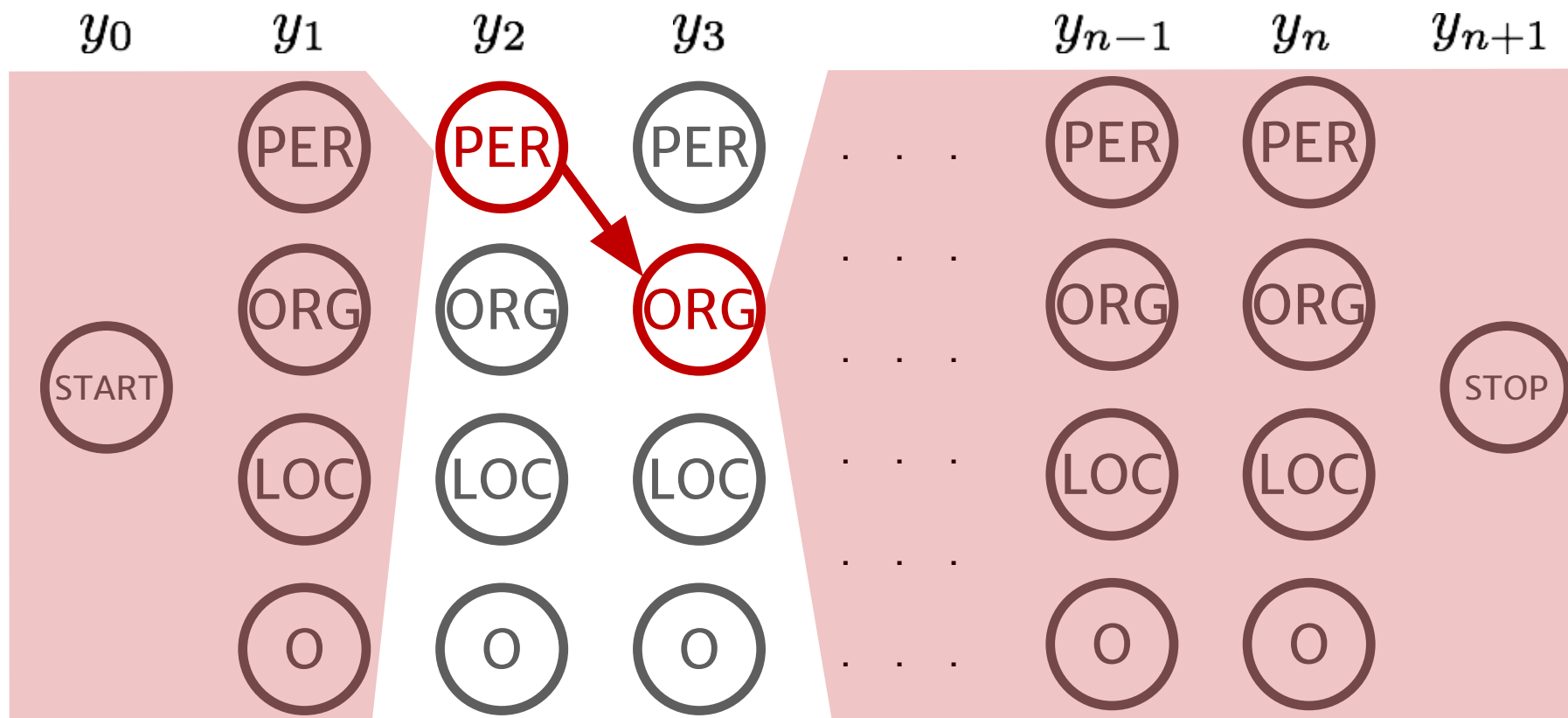


$$\beta_t(s) = \sum_{s'} \exp(w^\top f(x, y_t = s, y_{t+1} = s')) \beta_{t+1}(s')$$

$$\beta_{n+1}(s) = \mathbb{1}[s = \text{STOP}]$$



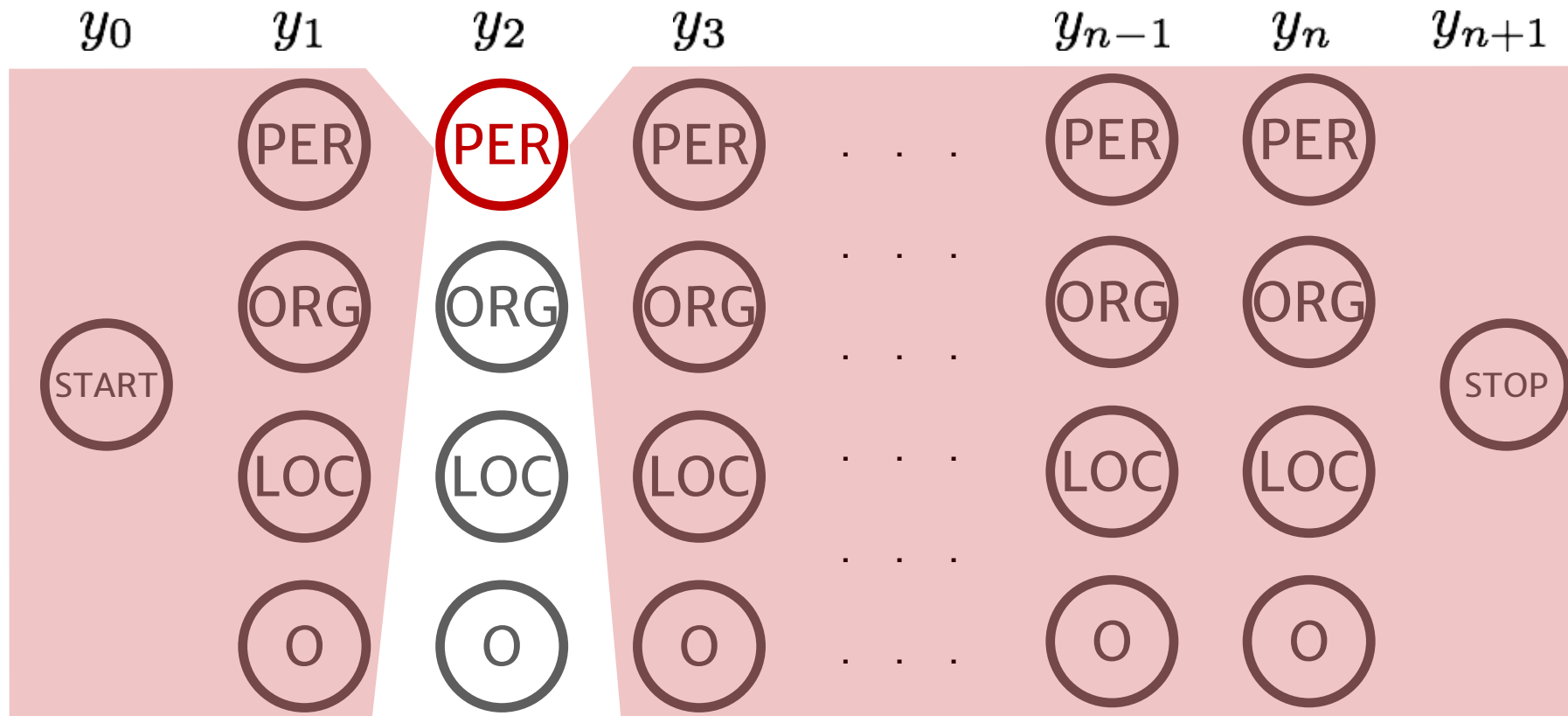
# Computing marginals



$$P(y_t = s, y_{t-1} = s' | x, w) = \frac{\alpha_{t-1}(s') \exp(w^\top f(x, y_t = s, y_{t-1} = s')) \beta_t(s)}{\alpha_{n+1}(\text{STOP})}$$



# Computing marginals

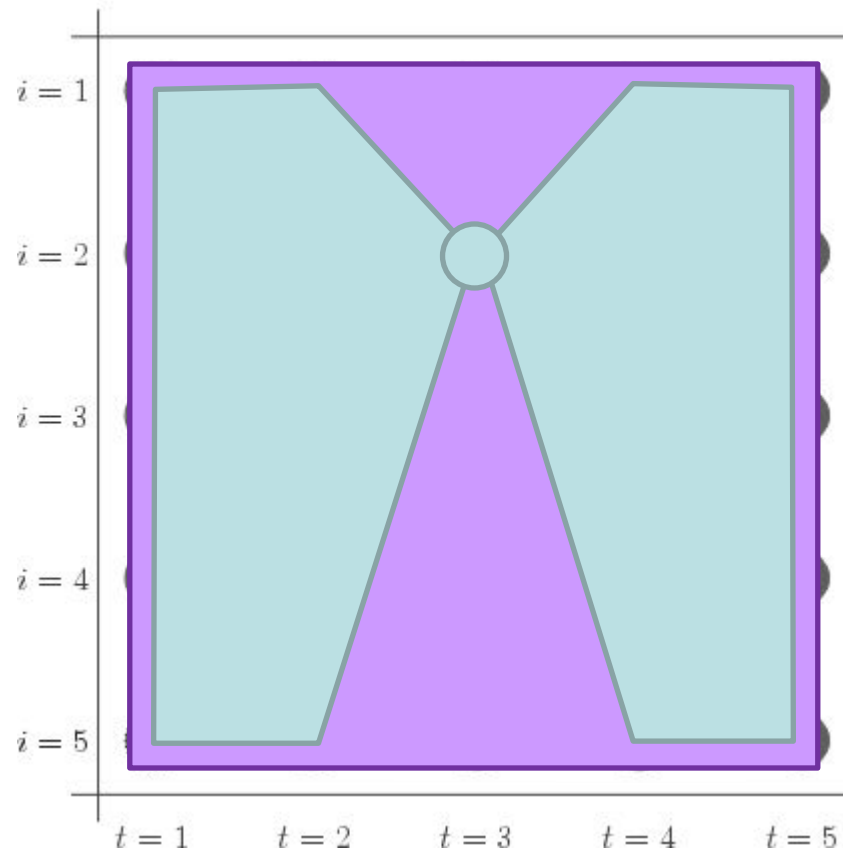


$$P(y_t = s | x, w) = \frac{\alpha_t(s)\beta_t(s)}{\alpha_{n+1}(\text{STOP})}$$





# Forward-Backward

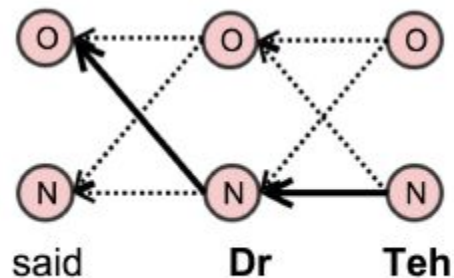


$$P(s_t|x) = \frac{P(s_t, x)}{P(x)}$$

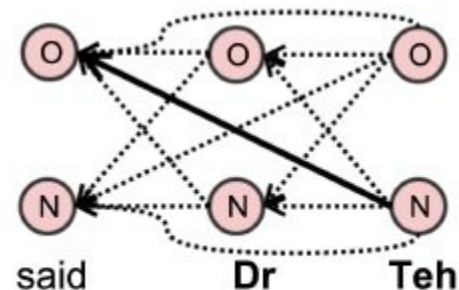
= sum of all paths through s at t  
sum of all possible paths



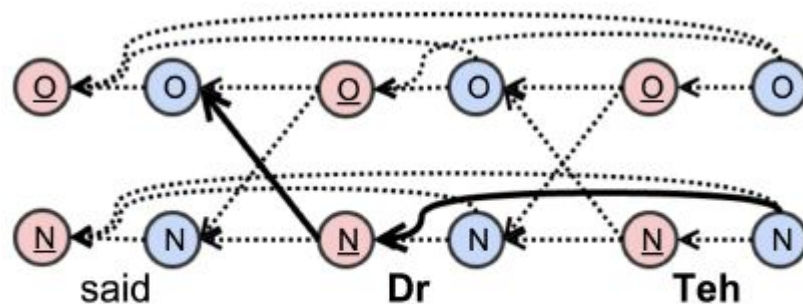
# Variants of “All possible paths”



Linear CRF



Semi-Markov CRF

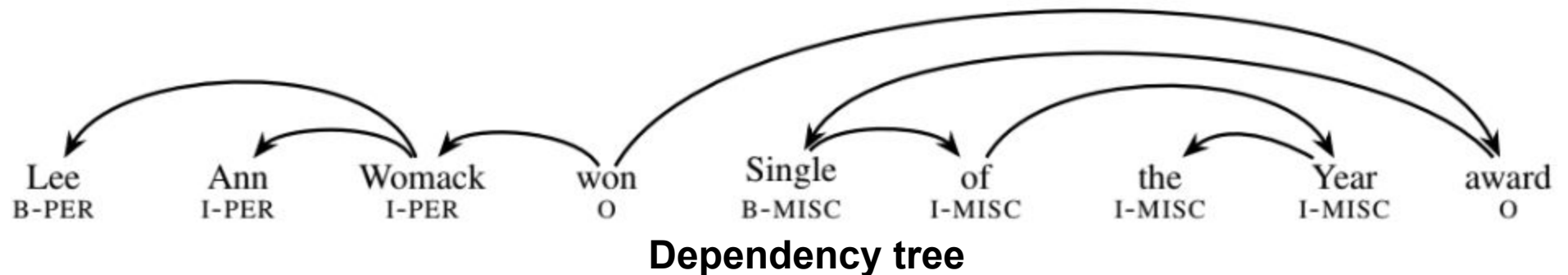
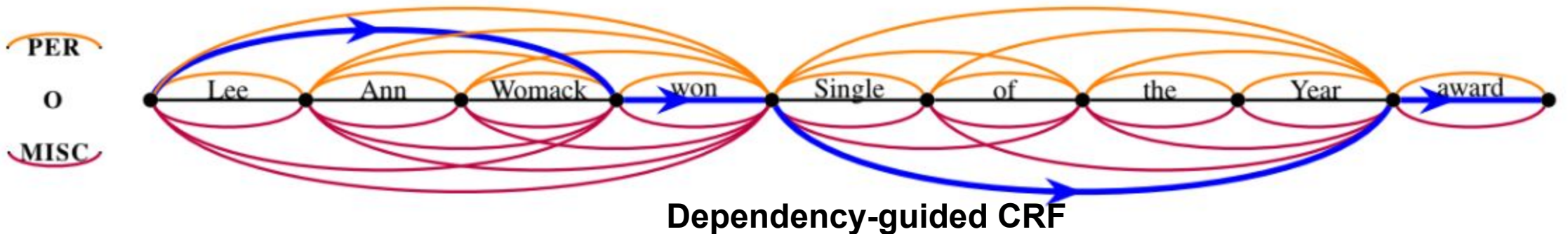
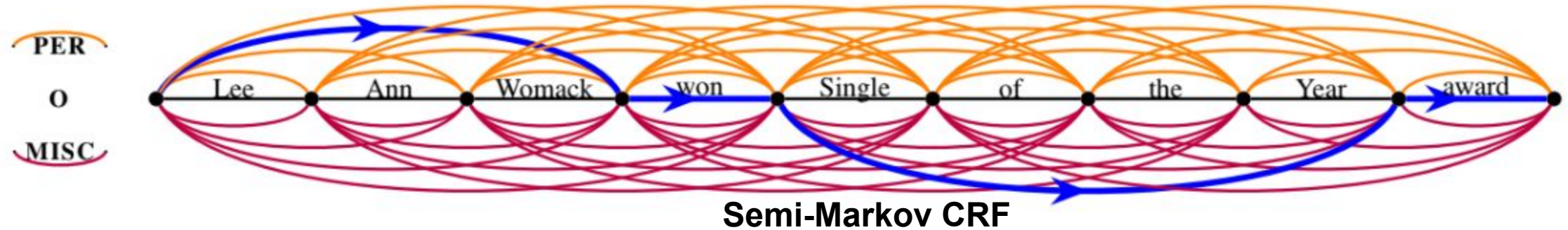


Weak Semi-Markov CRF

[Muis and Lu. 2016. Weak Semi-Markov CRFs for Noun Phrase Chunking in Informal Text. NAACL]



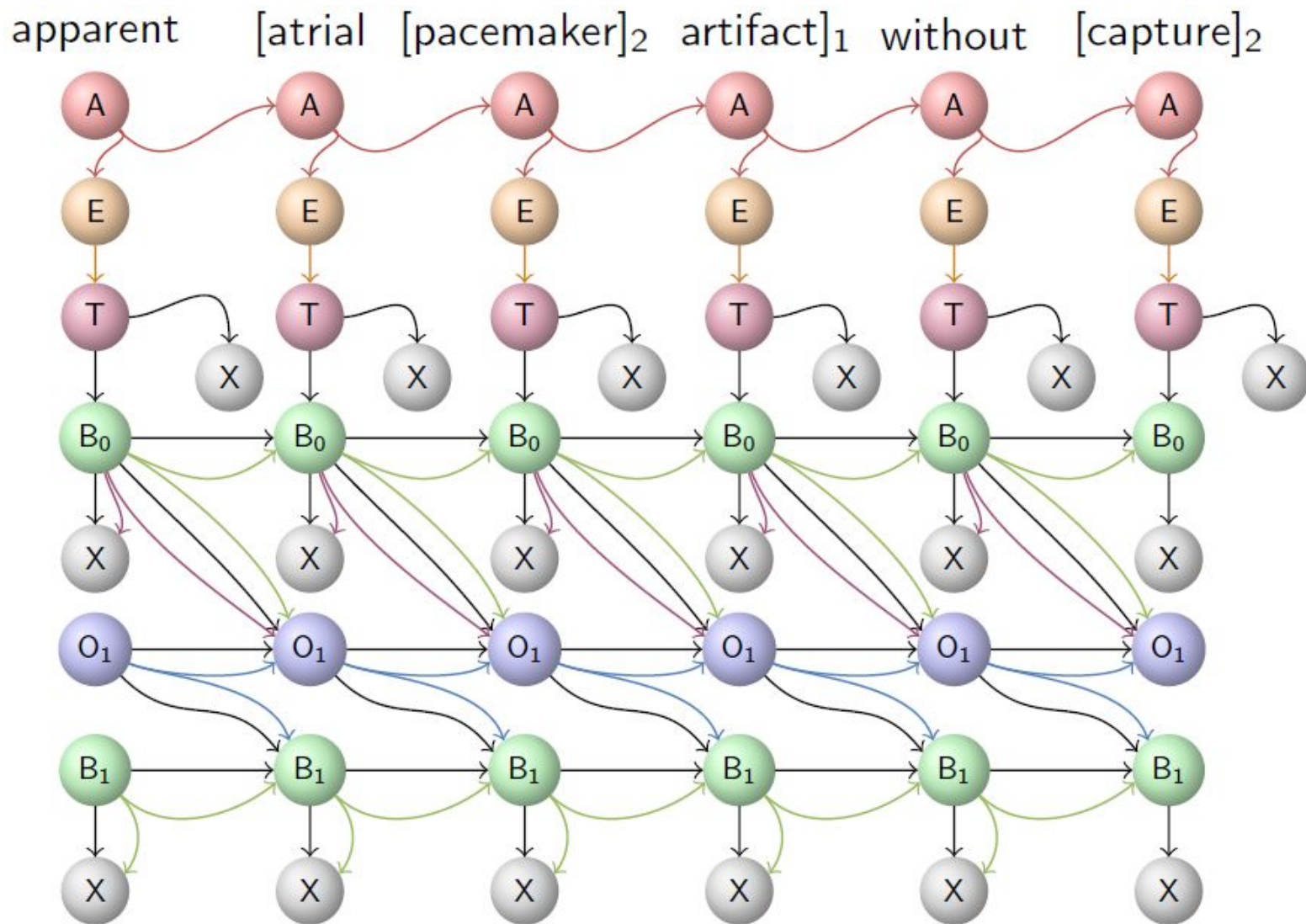
# Variants of “All possible paths”



[Jie, Muis and Lu. 2017. Efficient Dependency-guided Named Entity Recognition. AAAI]



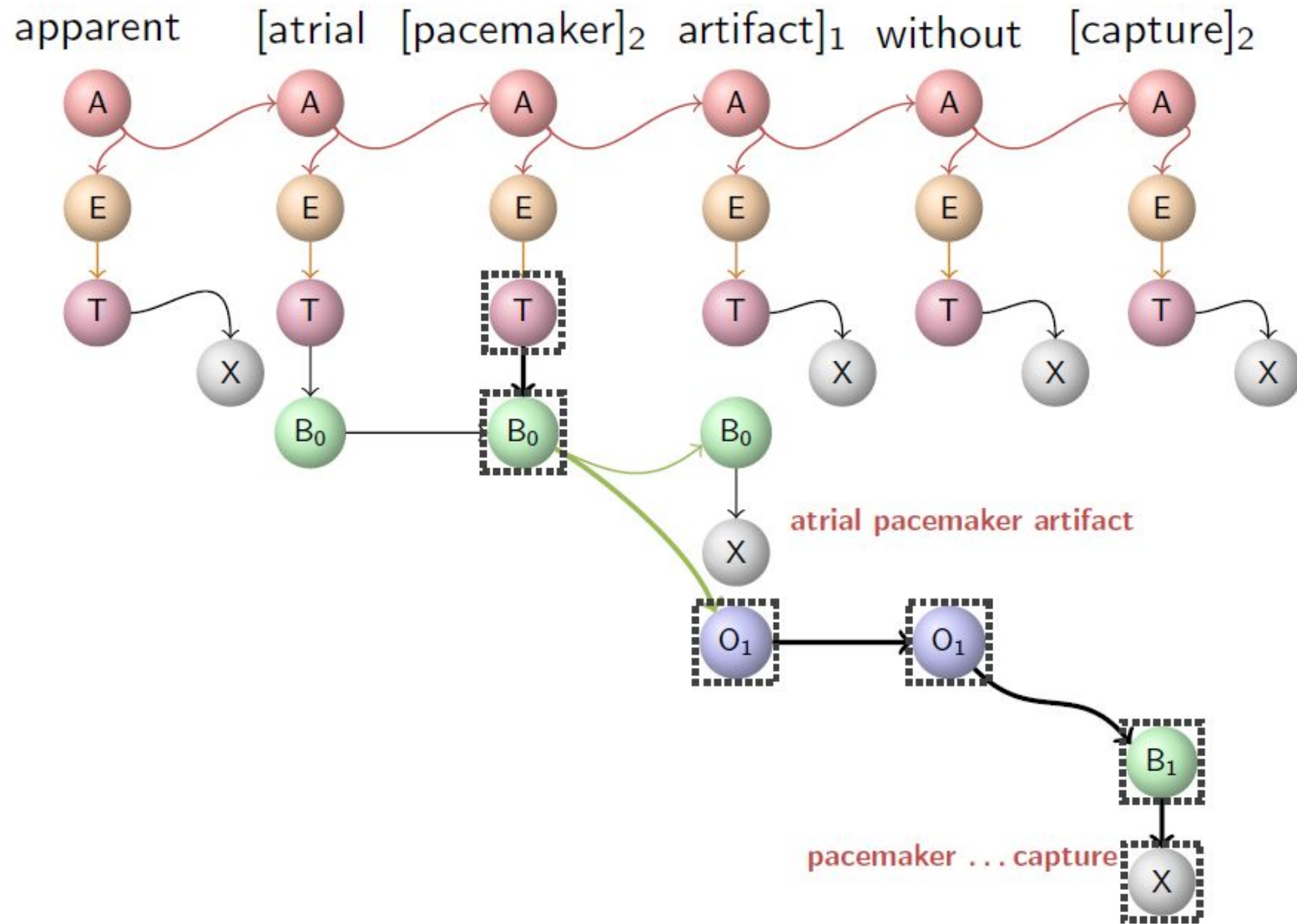
# Variants of “All possible paths”



[Muis and Lu. 2016. Learning to Recognize Discontiguous Entities. EMNLP]



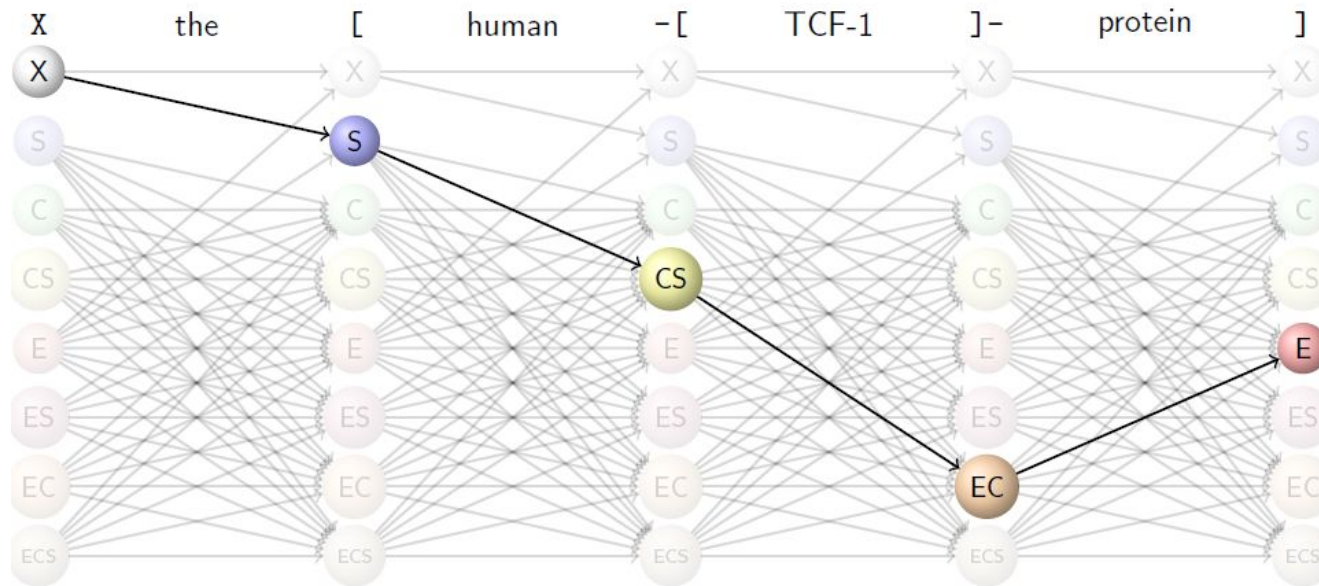
# Variants of “All possible paths”



[Muis and Lu. 2016. Learning to Recognize Discontiguous Entities. EMNLP]



# Variants of “All possible paths”



Linear-CRF:

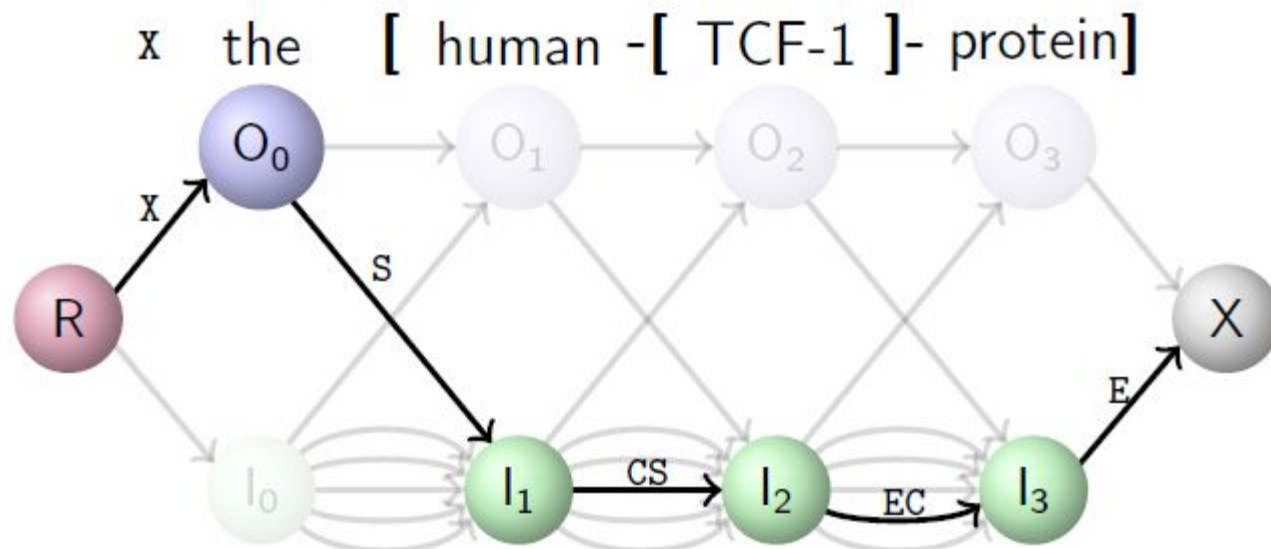
8 states per word

$8 \times 8 = 64$  edges between words

Multigraph CRF

2 states

8 edges between words



[Muis and Lu. 2017. Labeling Gaps between Words: Recognizing Overlapping Mentions with Mention Separators. EMNLP]



# Variants on Objective Functions

$$\begin{aligned} \text{Structured Perceptron: } & \min_{\mathbf{w}} \sum_{(x,y) \in \mathcal{T}} \left( \underbrace{-\mathbf{w} \cdot \Phi(x, y)}_{\text{Gold Score}} + \underbrace{\max_{\hat{y} \in \mathcal{Y}}}_{\text{Selector}} \left( \underbrace{\mathbf{w} \cdot \Phi(x, \hat{y})}_{\text{Other Score}} \right) \right) \underbrace{\hspace{10em}}_{\text{Cost}} \\ \text{CRF: } & \min_{\mathbf{w}} \sum_{(x,y) \in \mathcal{T}} \left( -\mathbf{w} \cdot \Phi(x, y) + \text{softmax}_{\hat{y} \in \mathcal{Y}} \left( \mathbf{w} \cdot \Phi(x, \hat{y}) \right) \right) \\ \text{Max-Margin (SSVM): } & \min_{\mathbf{w}} \sum_{(x,y) \in \mathcal{T}} \left( -\mathbf{w} \cdot \Phi(x, y) + \max_{\hat{y} \in \mathcal{Y}} \left( \mathbf{w} \cdot \Phi(x, \hat{y}) + \text{cost}(y, \hat{y}) \right) \right) \\ \text{Softmax-Margin: } & \min_{\mathbf{w}} \sum_{(x,y) \in \mathcal{T}} \left( -\mathbf{w} \cdot \Phi(x, y) + \text{softmax}_{\hat{y} \in \mathcal{Y}} \left( \mathbf{w} \cdot \Phi(x, \hat{y}) + \text{cost}(y, \hat{y}) \right) \right) \end{aligned}$$

[see also Gimpel and Smith. 2010.

Softmax-Margin Training for Structured Log-linear Models]



# EngCG Tagger

- English constraint grammar tagger
  - [Tapanainen and Voutilainen 94]
  - Something else you should know about
  - Hand-written and knowledge driven
  - “Don’t guess if you know” (general point about modeling more structure!)
  - Tag set doesn’t make all of the hard distinctions as the standard tag set (e.g. JJ/NN)
  - They get stellar accuracies: 99% on *their* tag set
  - Linguistic representation matters...
  - ... but it’s easier to win when you make up the rules

```
walk
walk <SV> <SVO> V SUBJUNCTIVE VFIN
walk <SV> <SVO> V IMP VFIN
walk <SV> <SVO> V INF
walk <SV> <SVO> V PRES -SG3 VFIN
walk N NOM SG
```

```
walk V-SUBJUNCTIVE V-IMP V-INF
V-PRES-BASE N-NOM-SG
```





# Domain Effects

---

- Accuracies degrade outside of domain
  - Up to triple error rate
  - Usually make the most errors on the things you care about in the domain (e.g. protein names)
- Open questions
  - How to effectively exploit unlabeled data from a new domain/language (what could we gain?) [Daume III. 2007 (FEDA), Kim et al. 2016 (Neural FEDA), Chen et al. 2016 (Deep Averaging Networks), Muis et al. 2018 (Distant Supervision)]
  - How to best incorporate domain lexica in a principled way (e.g. UMLS specialist lexicon, ontologies)

# Unsupervised Tagging



# Unsupervised Tagging?

---

- AKA part-of-speech induction
- Task:
  - Raw sentences in
  - Tagged sentences out
- Obvious thing to do:
  - Start with a (mostly) uniform HMM
  - Run EM (Expectation-Maximization)
  - Inspect results



# EM for HMMs: Process

---

- Alternate between recomputing distributions over hidden variables (the tags) and re-estimating parameters
- Crucial step: we want to tally up how many (fractional) counts of each kind of transition and emission we have under current params:

$$\text{count}(w, s) = \sum_{i:w_i=w} P(t_i = s | \mathbf{w})$$

$$\text{count}(s \rightarrow s') = \sum_i P(t_{i-1} = s, t_i = s' | \mathbf{w})$$

- Same quantities we needed to train a CRF!



# EM for HMMs: Quantities

---

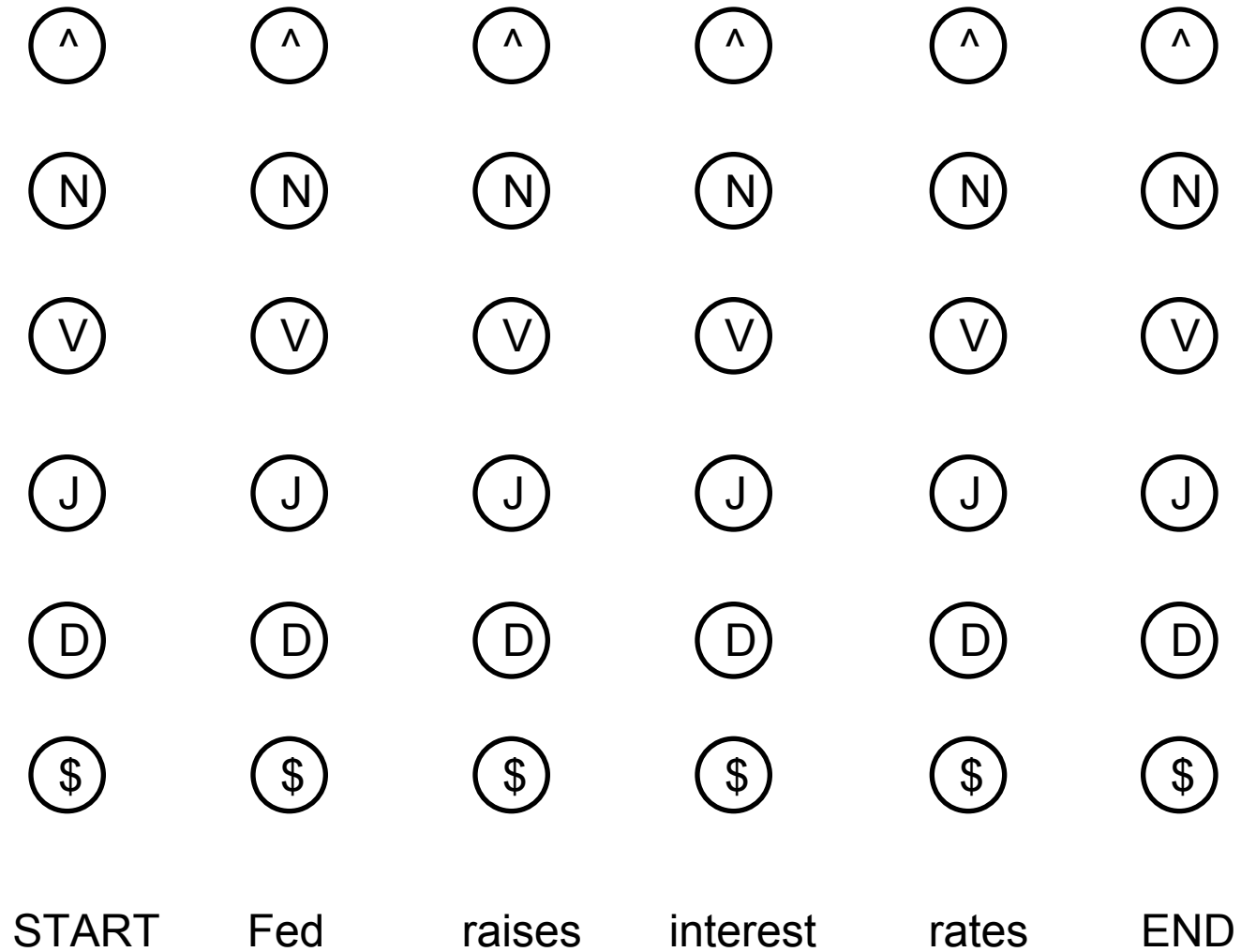
- Total path values (correspond to probabilities here):

$$\begin{aligned}\alpha_i(s) &= P(w_0 \dots w_i, s_i) \\ &= \sum_{s_{i-1}} P(s_i | s_{i-1}) P(w_i | s_i) \alpha_{i-1}(s_{i-1})\end{aligned}$$

$$\begin{aligned}\beta_i(s) &= P(w_{i+1} \dots w_n | s_i) \\ &= \sum_{s_{i+1}} P(s_{i+1} | s_i) P(w_{i+1} | s_{i+1}) \beta_{i+1}(s_{i+1})\end{aligned}$$



# The State Lattice / Trellis





# EM for HMMs: Process

---

- From these quantities, can compute expected transitions:

$$\text{count}(s \rightarrow s') = \frac{\sum_i \alpha_i(s) P(s'|s) P(w_i|s) \beta_{i+1}(s')}{P(\mathbf{w})}$$

- And emissions:

$$\text{count}(w, s) = \frac{\sum_{i:w_i=w} \alpha_i(s) \beta_{i+1}(s)}{P(\mathbf{w})}$$



# Meriardo: Setup

---

- Some (discouraging) experiments [Meriardo 94]
- Setup:
  - You know the set of allowable tags for each word
  - Fix  $k$  training examples to their true labels
    - Learn  $P(w|t)$  on these examples
    - Learn  $P(t|t_{-1}, t_{-2})$  on these examples
  - On  $n$  examples, re-estimate with EM
- Note: we know allowed tags but not frequencies





# Meritaldo: Results

Number of tagged sentences used for the initial model							
	0	100	2000	5000	10000	20000	all
Iter	Correct tags (% words) after ML on 1M words						
0	77.0	90.0	95.4	96.2	96.6	96.9	97.0
1	80.5	92.6	95.8	96.3	96.6	96.7	96.8
2	81.8	93.0	95.7	96.1	96.3	96.4	96.4
3	83.0	93.1	95.4	95.8	96.1	96.2	96.2
4	84.0	93.0	95.2	95.5	95.8	96.0	96.0
5	84.8	92.9	95.1	95.4	95.6	95.8	95.8
6	85.3	92.8	94.9	95.2	95.5	95.6	95.7
7	85.8	92.8	94.7	95.1	95.3	95.5	95.5
8	86.1	92.7	94.6	95.0	95.2	95.4	95.4
9	86.3	92.6	94.5	94.9	95.1	95.3	95.3
10	86.6	92.6	94.4	94.8	95.0	95.2	95.2

“There is no data like more data” – Mercer, 1985



# Distributional Clustering

◆ *the president said that the downturn was over* ◆

<i>president</i>	<i>the __ of</i>
<i>president</i>	<i>the __ said</i>
<i>governor</i>	<i>the __ of</i>
<i>governor</i>	<i>the __ appointed</i>
<i>said</i>	<i>sources __</i> ◆
<i>said</i>	<i>president __ that</i>
<i>reported</i>	<i>sources __</i> ◆

*president  
governor*

*said  
reported*

*the  
a*

[Finch and Chater 92, Shuetze 93, many others]



# Distributional Clustering

---

- Three main variants on the same idea:
  - Pairwise similarities and heuristic clustering
    - E.g. [Finch and Chater 92]
    - Produces dendrograms
  - Vector space methods
    - E.g. [Shuetze 93]
    - Models of ambiguity
  - Probabilistic methods
    - Various formulations, e.g. [Lee and Pereira 99]
  - Basis of (neural) word embedding

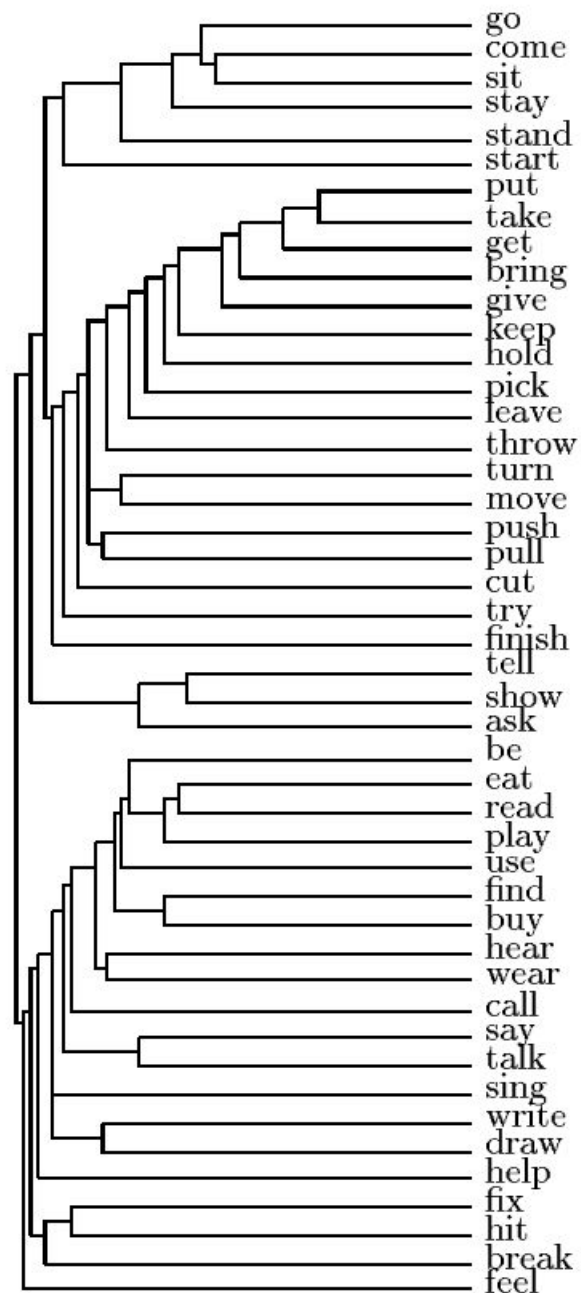
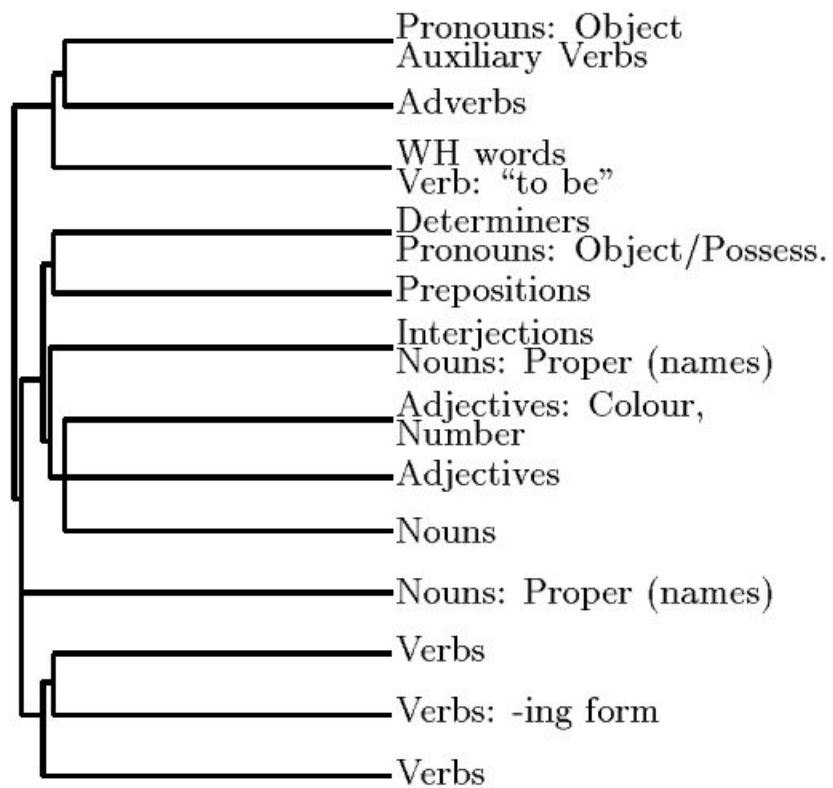


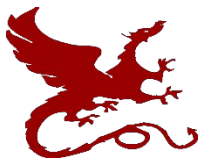
# Nearest Neighbors

word	nearest neighbors
accompanied	submitted banned financed developed authorized headed canceled awarded barred
almost	virtually merely formally fully quite officially just nearly only less
causing	reflecting forcing providing creating producing becoming carrying particularly
classes	elections courses payments losses computers performances violations levels pictures
directors	professionals investigations materials competitors agreements papers transactions
goal	mood roof eye image tool song pool scene gap voice
japanese	chinese iraqi american western arab foreign european federal soviet indian
represent	reveal attend deliver reflect choose contain impose manage establish retain
think	believe wish know realize wonder assume feel say mean bet
york	angeles francisco sox rouge kong diego zone vegas inning layer
on	through in at over into with from for by across
must	might would could cannot will should can may does helps
they	we you i he she nobody who it everybody there

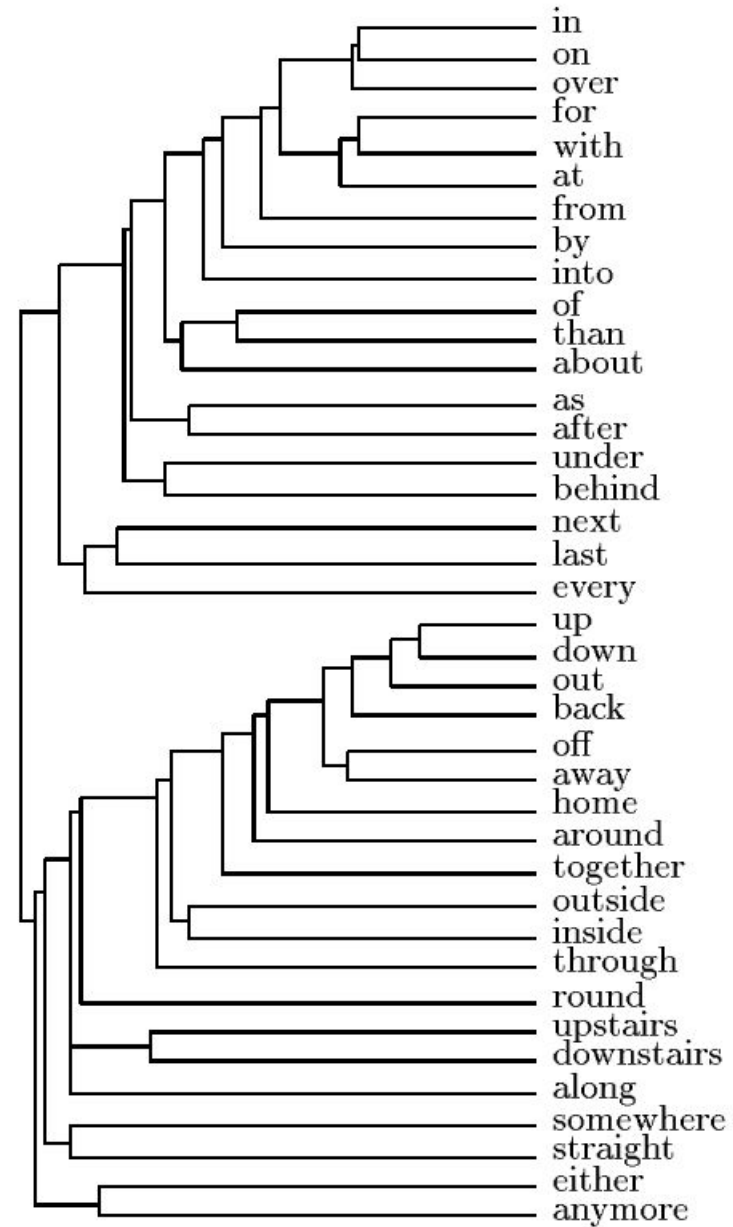
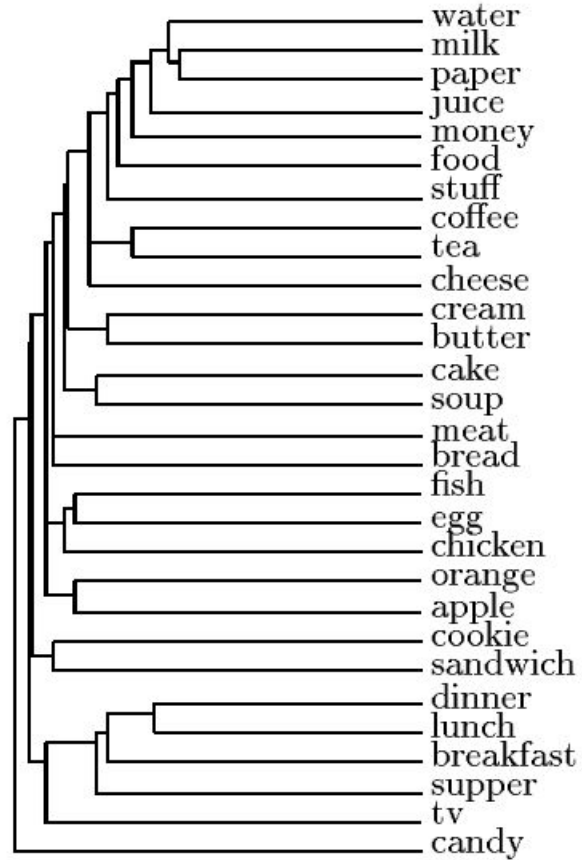


# Dendrograms





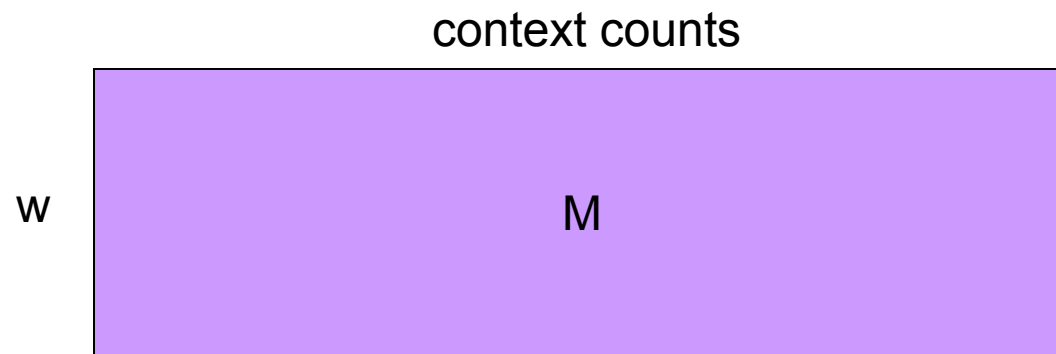
# Dendrograms



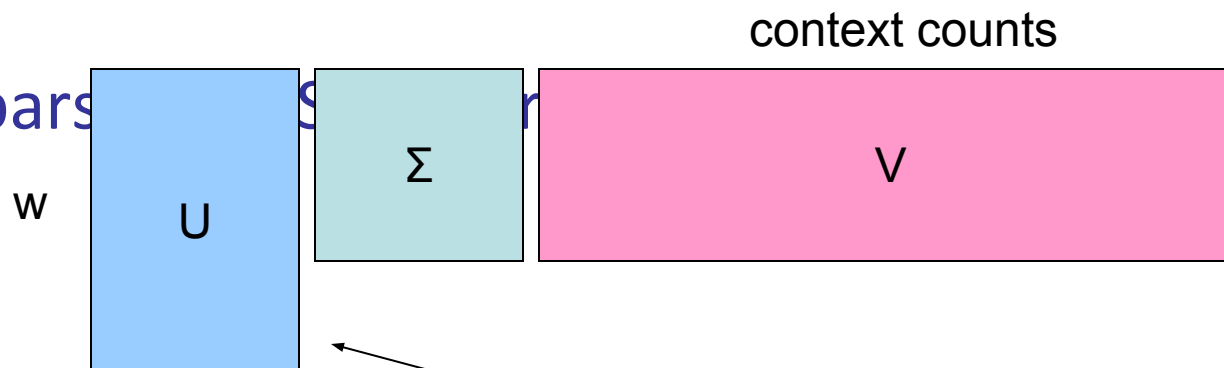


# Vector Space Version

- [Shuetze 93] clusters words as points in  $R^n$



- Vectors too sparse

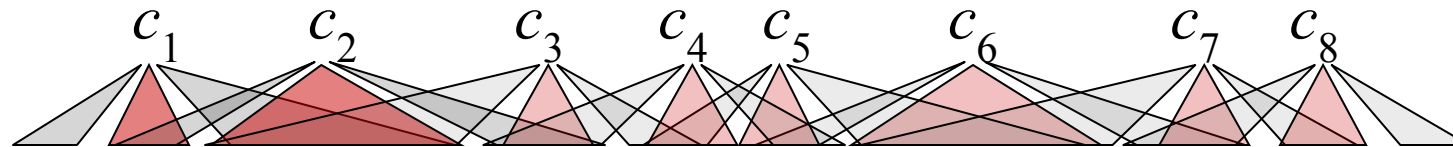


Cluster these 50-200 dim vectors instead.

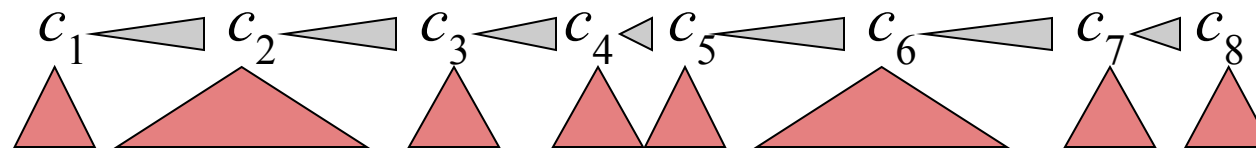


# A Probabilistic Version?

$$P(S, C) = \prod_i P(c_i)P(w_i | c_i)P(w_{i-1}, w_{i+1} | c_i)$$



◆ *the president said that the downturn was over* ◆



◆ *the president said that the downturn was over* ◆





# What Else?

---

- Various newer ideas:

- Context distributional clustering [Clark 00]
- Morphology-driven models [Clark 03]
- Contrastive estimation [Smith and Eisner 05]
- Feature-rich induction [Haghighi and Klein 06]

- Also:

- What about ambiguous words?
- Using wider context signatures has been used for learning synonyms (what's wrong with this approach?)
- Can extend these ideas for grammar induction (later)