

NLP

Introduction to NLP

251.

Probabilistic Grammars

Need for Probabilistic Parsing

- Time flies like an arrow
 - Many parses
 - Some (clearly) more likely than others
 - Need for a probabilistic ranking method

Probabilistic CFG

- Just like (deterministic) CFG, a 4-tuple (N, Σ, R, S)
 - N : non-terminal symbols
 - Σ : terminal symbols (disjoint from N)
 - R : rules $(A \rightarrow \beta) [p]$
 - β is a string from $(\Sigma \cup N)^*$
 - p is the probability $P(\beta | A)$
 - S : start symbol (from N)

Example

S -> NP VP

NP -> DT N | NP PP

PP -> PRP NP

VP -> V NP | VP PP

DT -> 'a' | 'the'

N -> 'child' | 'cake' | 'fork'

PRP -> 'with' | 'to'

V -> 'saw' | 'ate'

Example

S -> NP VP
NP -> DT N
NP -> NP PP
PP -> PRP NP
VP -> V NP
VP -> VP PP
DT -> 'a'
DT -> 'the'
N -> 'child'
N -> 'cake'
N -> 'fork'
PRP -> 'with'
PRP -> 'to'
V -> 'saw'
V -> 'ate'

Example

S -> NP VP

NP -> DT N

NP -> NP PP

PP -> PRP NP

VP -> V NP

VP -> VP PP

DT -> 'a'

DT -> 'the'

N -> 'child'

N -> 'cake'

N -> 'fork'

PRP -> 'with'

PRP -> 'to'

V -> 'saw'

V -> 'ate'

Example

S	->	NP VP	[p0=1]
NP	->	DT N	[p1]
NP	->	NP PP	[p2]
PP	->	PRP NP	[p3=1]
VP	->	V NP	[p4]
VP	->	VP PP	[p5]
DT	->	'a'	[p6]
DT	->	'the'	[p7]
N	->	'child'	[p8]
N	->	'cake'	[p9]
N	->	'fork'	[p10]
PRP	->	'with'	[p11]
PRP	->	'to'	[p12]
V	->	'saw'	[p13]
V	->	'ate'	[p14]

Probability of a Parse Tree

- The probability of a parse tree t given all n productions used to build it:

$$p(t) = \prod_{i=1}^n p(\alpha \rightarrow \beta)$$

- The most likely parse is determined as follows:

$$\operatorname{argmax}_{t \in T(s)} p(t)$$

- The probabilities are obtained using MLE from the training corpus
- The probability of a *sentence* is the *sum* of the probabilities of all of its parses

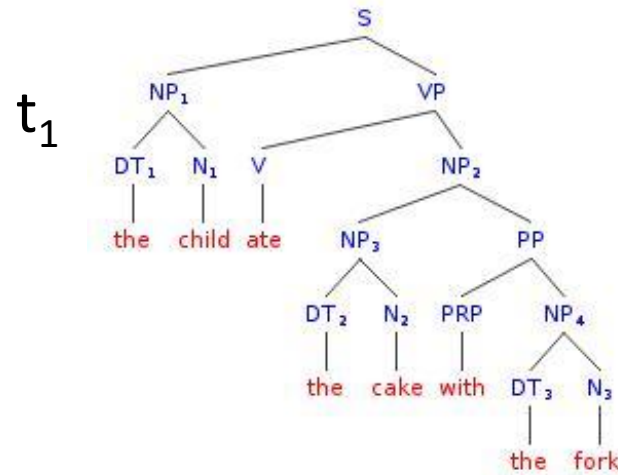
A		β	$P(\beta \mid \text{NP})$
NP	→	NP PP	0.092
NP	→	DT NN	0.087
NP	→	NN	0.047
NP	→	NNS	0.042
NP	→	DT JJ NN	0.035
NP	→	NNP	0.034
NP	→	NNP NNP	0.029
NP	→	JJ NNS	0.027
NP	→	QP -NONE-	0.018
NP	→	NP SBAR	0.017
NP	→	NP PP-LOC	0.017
NP	→	JJ NN	0.015
NP	→	DT NNS	0.014
NP	→	CD	0.014
NP	→	NN NNS	0.013
NP	→	DT NN NN	0.013
NP	→	NP CC NP	0.013

Example

S	->	NP VP	[p0=1]
NP	->	DT N	[p1]
NP	->	NP PP	[p2]
PP	->	PRP NP	[p3=1]
VP	->	V NP	[p4]
VP	->	VP PP	[p5]
DT	->	'a'	[p6]
DT	->	'the'	[p7]
N	->	'child'	[p8]
N	->	'cake'	[p9]
N	->	'fork'	[p10]
PRP	->	'with'	[p11]
PRP	->	'to'	[p12]
V	->	'saw'	[p13]
V	->	'ate'	[p14]

Example

S -> NP VP [p0=1]
NP -> DT N [p1]
NP -> NP PP [p2]
PP -> PRP NP [p3=1]
VP -> V NP [p4]
VP -> VP PP [p5]
DT -> 'a' [p6]
DT -> 'the' [p7]
N -> 'child' [p8]
N -> 'cake' [p9]
N -> 'fork' [p10]
PRP -> 'with' [p11]
PRP -> 'to' [p12]
V -> 'saw' [p13]
V -> 'ate' [p14]



Example

S -> NP VP [p0=1]

NP -> DT N [p1]

NP -> NP PP [p2]

PP -> PRP NP [p3=1]

VP -> V NP [p4]

VP -> VP PP [p5]

DT -> 'a' [p6]

DT -> 'the' [p7]

N -> 'child' [p8]

N -> 'cake' [p9]

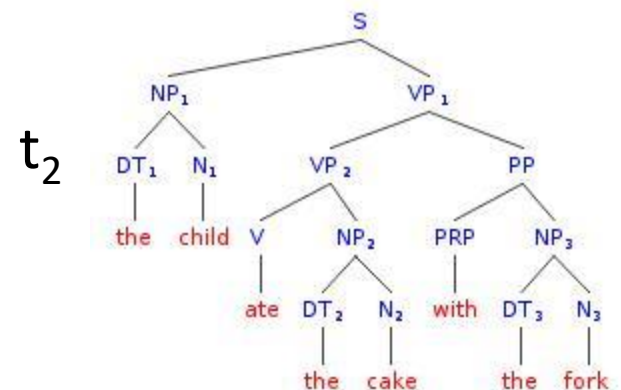
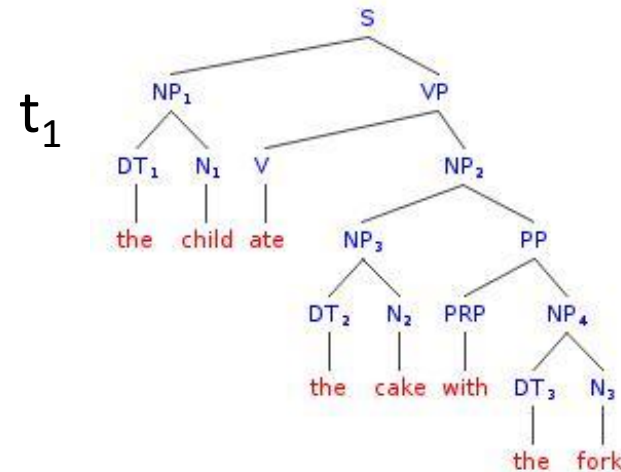
N -> 'fork' [p10]

PRP -> 'with' [p11]

PRP -> 'to' [p12]

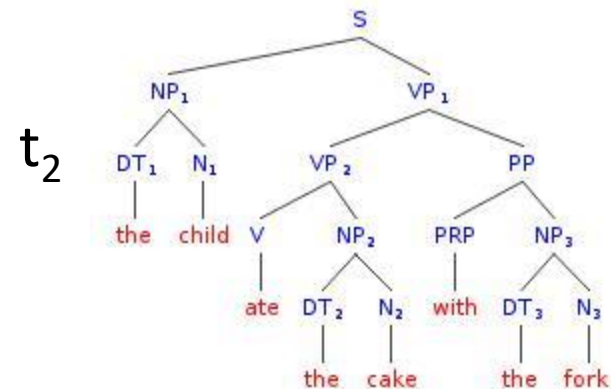
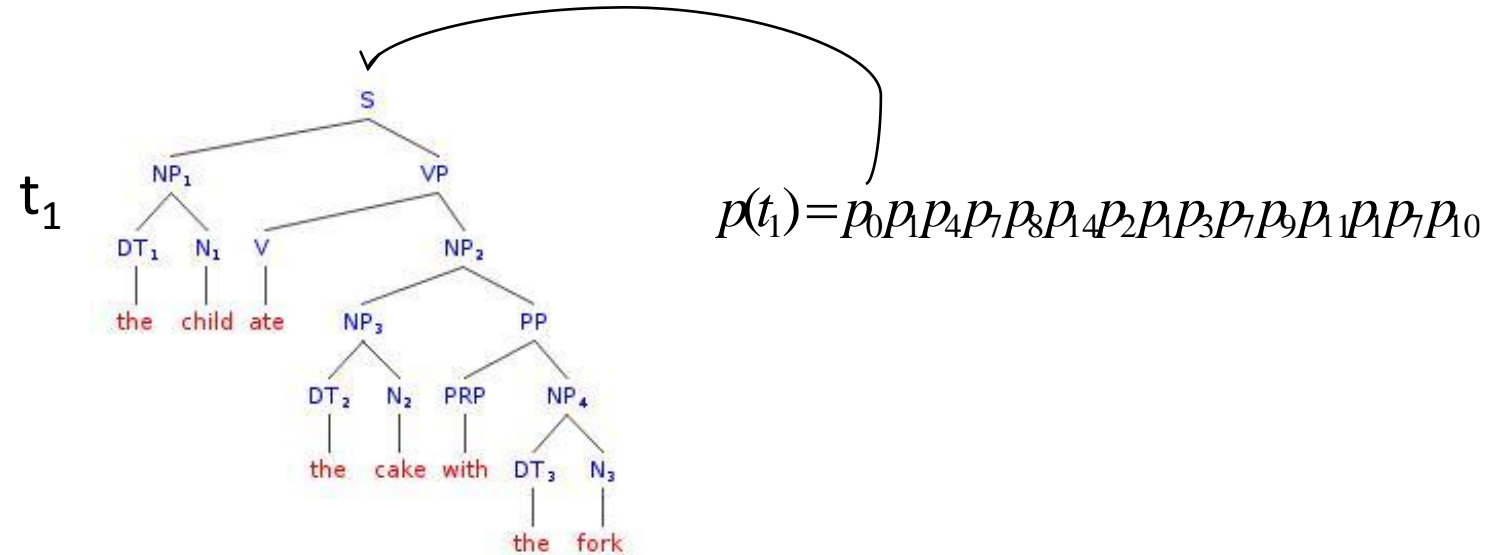
V -> 'saw' [p13]

V -> 'ate' [p14]



Example

S → NP VP [p0=1]
NP → DT N [p1]
NP → NP PP [p2]
PP → PRP NP [p3=1]
VP → V NP [p4]
VP → VP PP [p5]
DT → 'a' [p6]
DT → 'the' [p7]
N → 'child' [p8]
N → 'cake' [p9]
N → 'fork' [p10]
PRP → 'with' [p11]
PRP → 'to' [p12]
V → 'saw' [p13]
V → 'ate' [p14]



Example

S → NP VP [p0=1]

NP → DT N [p1]

NP → NP PP [p2]

PP → PRP NP [p3=1]

VP → V NP [p4]

VP → VP PP [p5]

DT → 'a' [p6]

DT → 'the' [p7]

N → 'child' [p8]

N → 'cake' [p9]

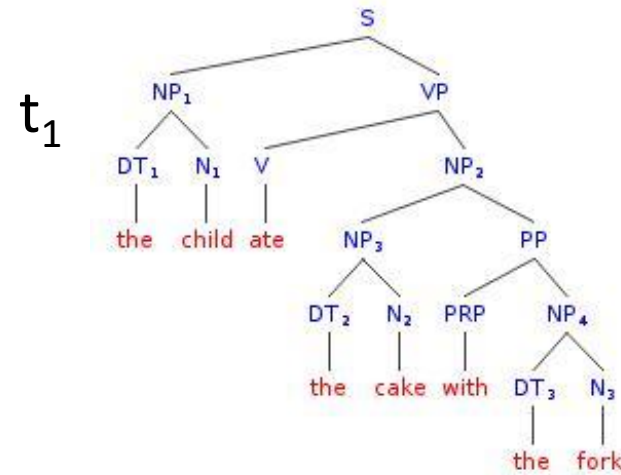
N → 'fork' [p10]

PRP → 'with' [p11]

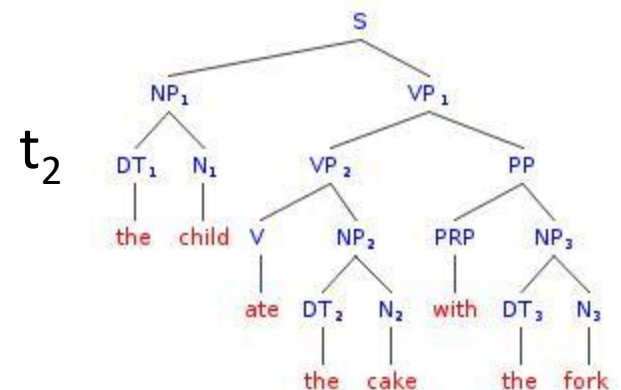
PRP → 'to' [p12]

V → 'saw' [p13]

V → 'ate' [p14]



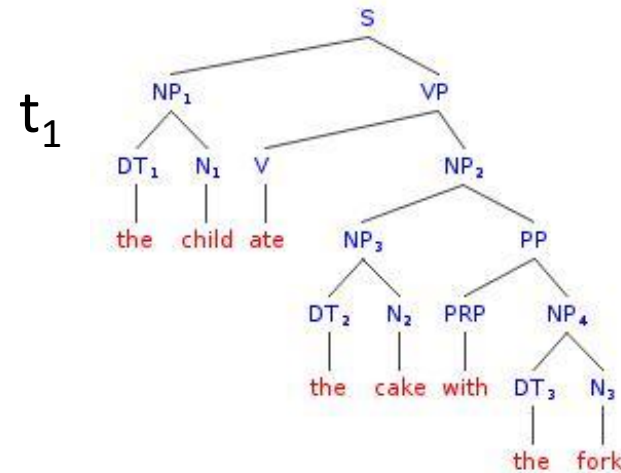
$$p(t_1) = p_0 p_1 p_4 p_7 p_8 p_{14} p_2 p_1 p_3 p_7 p_9 p_{11} p_1 p_7 p_{10}$$



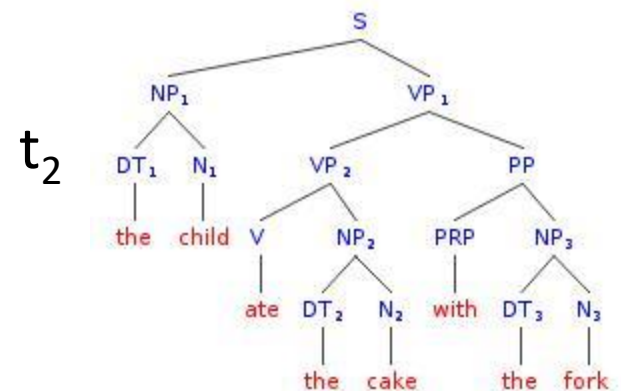
$$p(t_2) = p_0 p_1 p_5 p_7 p_8 p_4 p_3 p_{14} p_1 p_{11} p_1 p_7 p_9 p_7 p_{10}$$

Example

S → NP VP [p0=1]
 NP → DT N [p1]
 NP → NP PP [p2]
 PP → PRP NP [p3=1]
 VP → V NP [p4]
 VP → VP PP [p5]
 DT → 'a' [p6]
 DT → 'the' [p7]
 N → 'child' [p8]
 N → 'cake' [p9]
 N → 'fork' [p10]
 PRP → 'with' [p11]
 PRP → 'to' [p12]
 V → 'saw' [p13]
 V → 'ate' [p14]



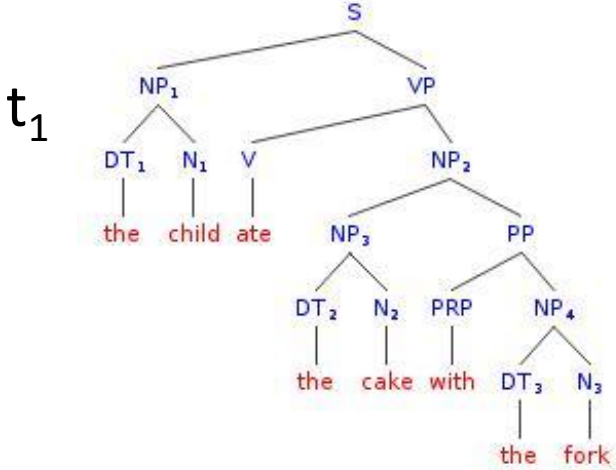
$$p(t_1) = p_0 p_1 p_4 p_7 p_8 p_{14} p_2 p_1 p_3 p_7 p_9 p_{11} p_1 p_7 p_{10}$$



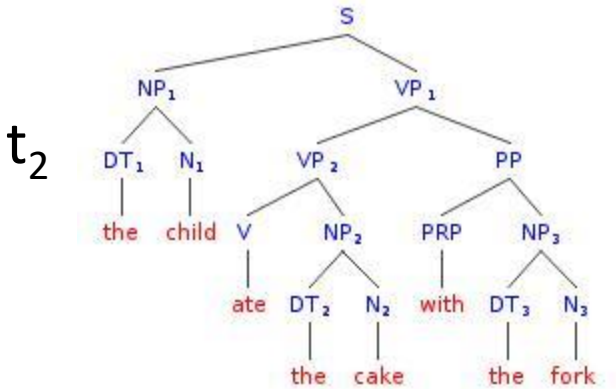
$$p(t_2) = p_0 p_1 p_5 p_7 p_8 p_4 p_3 p_{14} p_1 p_{11} p_1 p_7 p_9 p_7 p_{10}$$

Example

S	->	NP VP	[p0=1]
NP	->	DT N	[p1]
NP	->	NP PP	[p2]
PP	->	PRP NP	[p3=1]
VP	->	V NP	[p4]
VP	->	VP PP	[p5]
DT	->	'a'	[p6]
DT	->	'the'	[p7]
N	->	'child'	[p8]
N	->	'cake'	[p9]
N	->	'fork'	[p10]
PRP	->	'with'	[p11]
PRP	->	'to'	[p12]
V	->	'saw'	[p13]
V	->	'ate'	[p14]



$$p(t_1) = p_0 p_1 p_4 p_7 p_8 p_{14} p_2 p_3 p_7 p_9 p_{11} p_7 p_{10}$$



$$p(t_2) = p_0 p_1 p_5 p_7 p_8 p_4 p_3 p_{14} p_{11} p_{11} p_7 p_9 p_7 p_{10}$$

NLP

Introduction to NLP

252.

Probabilistic Parsing

Main Tasks with PCFGs

- Given a grammar G and a sentence s , let $T(s)$ be all parse trees that correspond to s
- Task 1
 - find which tree t among $T(s)$ maximizes the probability $p(t)$
- Task 2
 - find the probability of the sentence $p(s)$ as the sum of all possible tree probabilities $p(t)$

Probabilistic Parsing Methods

- Probabilistic Earley algorithm
 - Top-down parser with a dynamic programming table
- Probabilistic Cocke-Kasami-Younger (CKY) algorithm
 - Bottom-up parser with a dynamic programming table

Probabilistic Grammars

- Probabilities can be learned from a labeled corpus
 - Treebank
- Intuitive meaning
 - Parse #1 is twice as probable as parse #2
- Possible to do reranking
- Possible to combine with other stages
 - E.g., speech recognition, translation

Maximum Likelihood Estimates

- Use the parsed training set for getting the counts
 - $P_{ML}(\alpha \rightarrow \beta) = \text{Count}(\alpha \rightarrow \beta) / \text{Count}(\alpha)$
- Example:
 - $P_{ML}(S \rightarrow \text{NP VP}) = \text{Count}(S \rightarrow \text{NP VP}) / \text{Count}(S)$

Sample Probabilistic Grammar

Grammar		Lexicon
$S \rightarrow NP VP$	[.80]	$Det \rightarrow that [.10] \mid a [.30] \mid the [.60]$
$S \rightarrow Aux NP VP$	[.15]	$Noun \rightarrow book [.10] \mid flight [.30]$
$S \rightarrow VP$	[.05]	$\mid meal [.15] \mid money [.05]$
$NP \rightarrow Pronoun$	[.35]	$\mid flights [.40] \mid dinner [.10]$
$NP \rightarrow Proper-Noun$	[.30]	$Verb \rightarrow book [.30] \mid include [.30]$
$NP \rightarrow Det Nominal$	[.20]	$\mid prefer; [.40]$
$NP \rightarrow Nominal$	[.15]	$Pronoun \rightarrow I [.40] \mid she [.05]$
$Nominal \rightarrow Noun$	[.75]	$\mid me [.15] \mid you [.40]$
$Nominal \rightarrow Nominal Noun$	[.20]	$Proper-Noun \rightarrow Houston [.60]$
$Nominal \rightarrow Nominal PP$	[.05]	$\mid NWA [.40]$
$VP \rightarrow Verb$	[.35]	$Aux \rightarrow does [.60] \mid can [.40]$
$VP \rightarrow Verb NP$	[.20]	$Preposition \rightarrow from [.30] \mid to [.30]$
$VP \rightarrow Verb NP PP$	[.10]	$\mid on [.20] \mid near [.15]$
$VP \rightarrow Verb PP$	[.15]	$\mid through [.05]$
$VP \rightarrow Verb NP NP$	[.05]	
$VP \rightarrow VP PP$	[.15]	
$PP \rightarrow Preposition NP$	[1.0]	

Example from Jurafsky and Martin

Example

S	->	NP VP	[p0=1]
NP	->	DT N	[p1=.8]
NP	->	NP PP	[p2=.2]
PP	->	PRP NP	[p3=1]
VP	->	V NP	[p4=.7]
VP	->	VP PP	[p5=.3]
DT	->	'a'	[p6=.25]
DT	->	'the'	[p7=.75]
N	->	'child'	[p8=.5]
N	->	'cake'	[p9=.3]
N	->	'fork'	[p10=.2]
PRP	->	'with'	[p11=.1]
PRP	->	'to'	[p12=.9]
V	->	'saw'	[p13=.4]
V	->	'ate'	[p14=.6]

function PROBABILISTIC-CKY(*words*, *grammar*) **returns** most probable parse
and its probability

for $j \leftarrow$ **from** 1 **to** LENGTH(*words*) **do**

for all $\{ A \mid A \rightarrow \text{words}[j] \in \text{grammar} \}$

$\text{table}[j-1, j, A] \leftarrow P(A \rightarrow \text{words}[j])$

for $i \leftarrow$ **from** $j-2$ **downto** 0 **do**

for $k \leftarrow i+1$ **to** $j-1$ **do**

for all $\{ A \mid A \rightarrow BC \in \text{grammar},$

and $\text{table}[i, k, B] > 0$ **and** $\text{table}[k, j, C] > 0 \}$

if $(\text{table}[i, j, A] < P(A \rightarrow BC) \times \text{table}[i, k, B] \times \text{table}[k, j, C])$ **then**

$\text{table}[i, j, A] \leftarrow P(A \rightarrow BC) \times \text{table}[i, k, B] \times \text{table}[k, j, C]$

$\text{back}[i, j, A] \leftarrow \{k, B, C\}$

return BUILD_TREE($\text{back}[1, \text{LENGTH}(\text{words}), S]$), $\text{table}[1, \text{LENGTH}(\text{words}), S]$

the

child

ate

the

cake

with

the

fork

the

DT .75							
child							
ate							
the							
cake							
with							
the							
fork							

the

DT .75							
child	N .5						
	ate						
		the					
			cake				
				with			
					the		
						fork	

the

DT .75	NP .8						
child	N .5						
	ate						
		the					
			cake				
				with			
					the		
						fork	

the	DT .75	NP .8*.5*.75						
child	N .5							
ate								
the								
cake								
with								
the								
fork								

the	DT .75	NP .8*.5*.75						
	child	N .5						
		ate						
			the					
				cake				
					with			
						the		
							fork	

Keep only the highest score in each cell

Exercise

- Now, on your own, compute the probability of the entire sentence using Probabilistic CKY.
- Don't forget that there may be multiple parses, so you will need to add the corresponding probabilities.

Notes

- Stanford Demo
 - <http://nlp.stanford.edu:8080/parser/>
- PTB statistics
 - 50,000 sentences (40,000 training; 2,400 testing)
- PTB peculiarities
 - includes traces and other null elements
 - Flat NP structure (e.g., NP -> DT JJ JJ NNP NNS)
- Parent transformation
 - Subject NPs are more likely to have modifiers than object NPs
 - E.g., replace NP with NP^S

NLP

Introduction to NLP

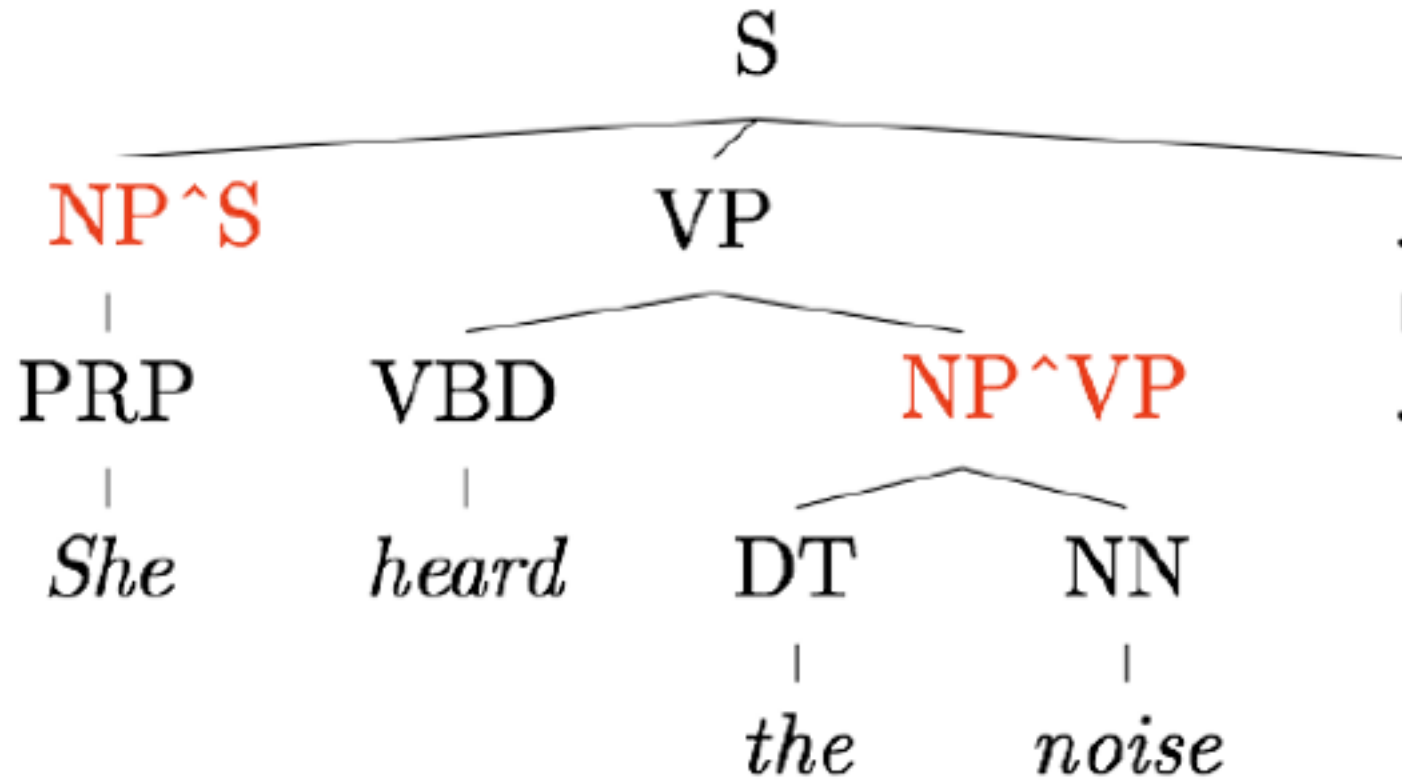
253.

Lexicalized Parsing

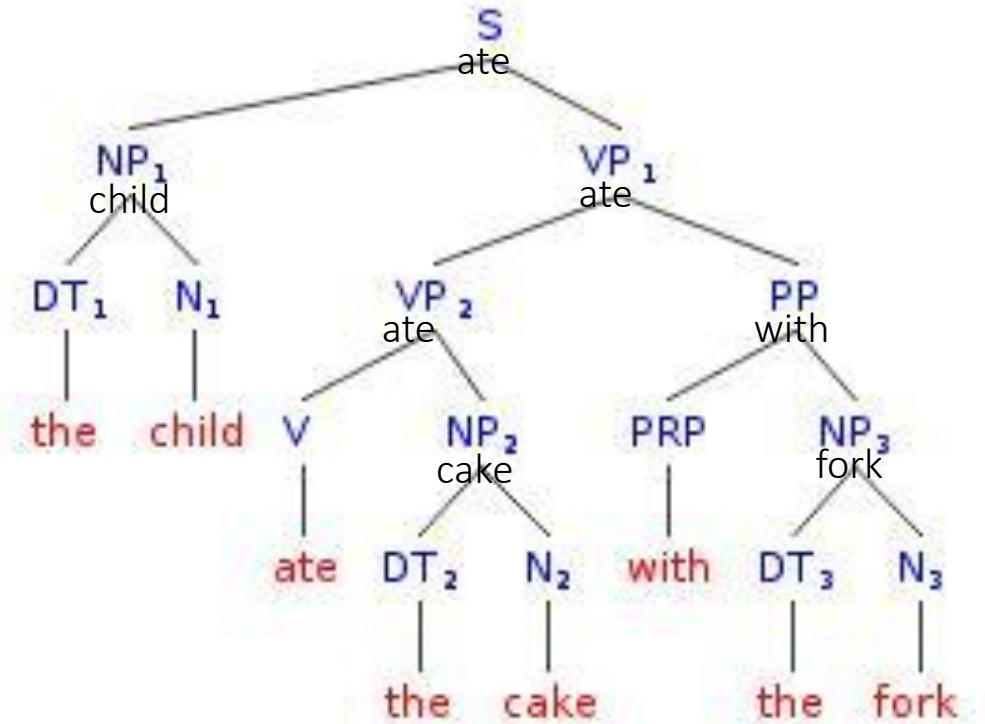
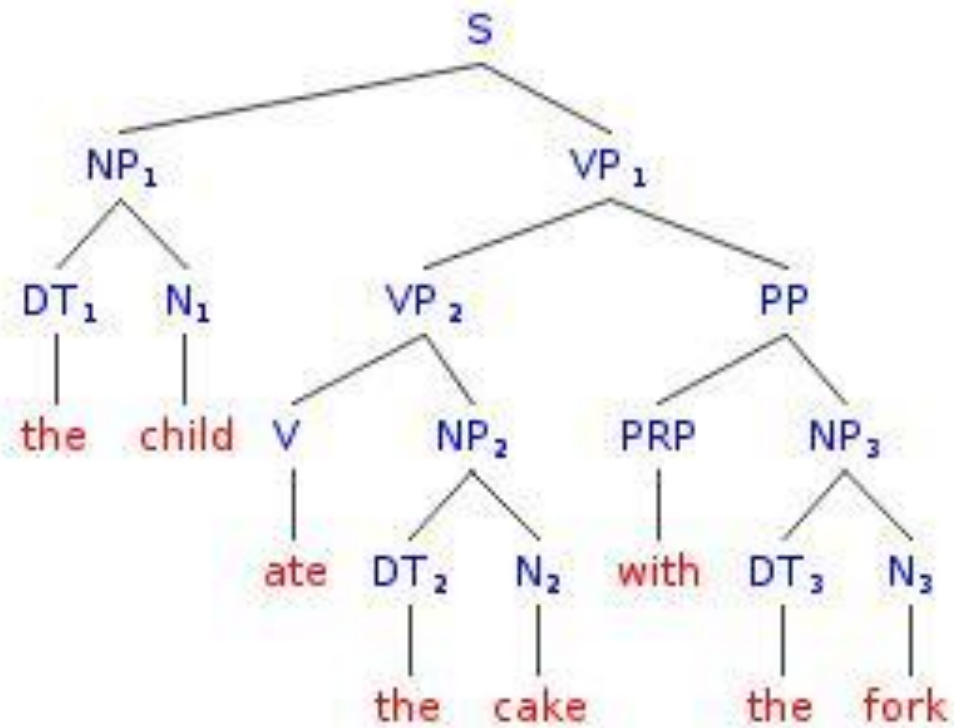
Limitations of PCFGs

- The probabilities don't depend on the specific words
 - E.g., *give* someone something (2 arguments) vs. *see* something (1 argument)
- It is not possible to disambiguate sentences based on semantic information
 - E.g., eat pizza with *pepperoni* vs. eat pizza with *fork*
- Lexicalized grammars - idea
 - Use the head of a phrase as an additional source of information
 - VP[ate] -> V[ate]
 - Fundamental idea in syntax, cf. X-bar theory, HPSG
 - Constituents receive their heads from their head child

Parent Annotation



Lexicalization



Head Extraction Example (Collins)

- NP -> DT NNP **NN** (rightmost)
- NP -> DT NN **NNP** (rightmost)
- NP -> **NP** PP (leftmost)
- NP -> DT **JJ** (rightmost)
- NP -> **DT** (rightmost leftover child)

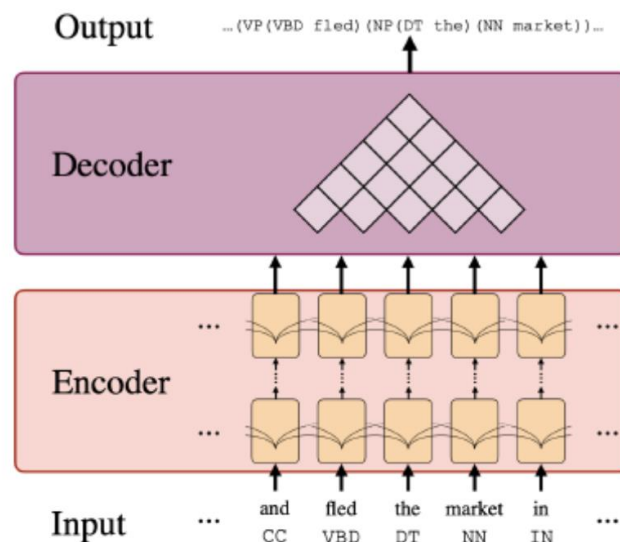
Notes

- Complexity of lexicalized parsing
 - $O(N^5g^3V^3)$, instead of $O(N^3)$ because of the lexicalization
 - N = sentence length
 - g = number of non-terminals
 - V = vocabulary size
 - Use beam search (Charniak; Collins)
- Sparse data
 - 40,000 sentences; 12,409 rules (Collins)
 - 15% of all test sentences contain a rule not seen in training (Collins)

Recent results

Labeled precision/recall on PTB-WSJ

- ▶ Vanilla PCFG: 70.6% recall, 74.8% precision
- ▶ Lexicalized PCFG: 88.1% recall, 88.3% precision
- ▶ Neuralized constituency parsing (Kitaev and Klein, 2018): 94.9% recall, 95.4% precision



Neural encoding followed by max marginal decoding: no independence assumption (read the paper)

Recent results (Label Attention Layer)

Example Input

The Label Attention Layer takes word vectors as input (red-contour matrix). In the example sentence, start and end symbols are omitted.

Select				
the				
person				
driving				

Label Attention Layer

Q is a matrix of learned query vectors. There is no more Query Matrix W^Q , and only one query vector is used per attention head. Each label is represented by one or more heads, and each head may represent one or more labels.

The query vectors q represent the attention weights from each head to dimensions of input vectors.

Computing the matrix of key vectors for the input. Each head has its own learned key matrix W^K .

The blue box outputs a vector of attention weights from each head to the words.

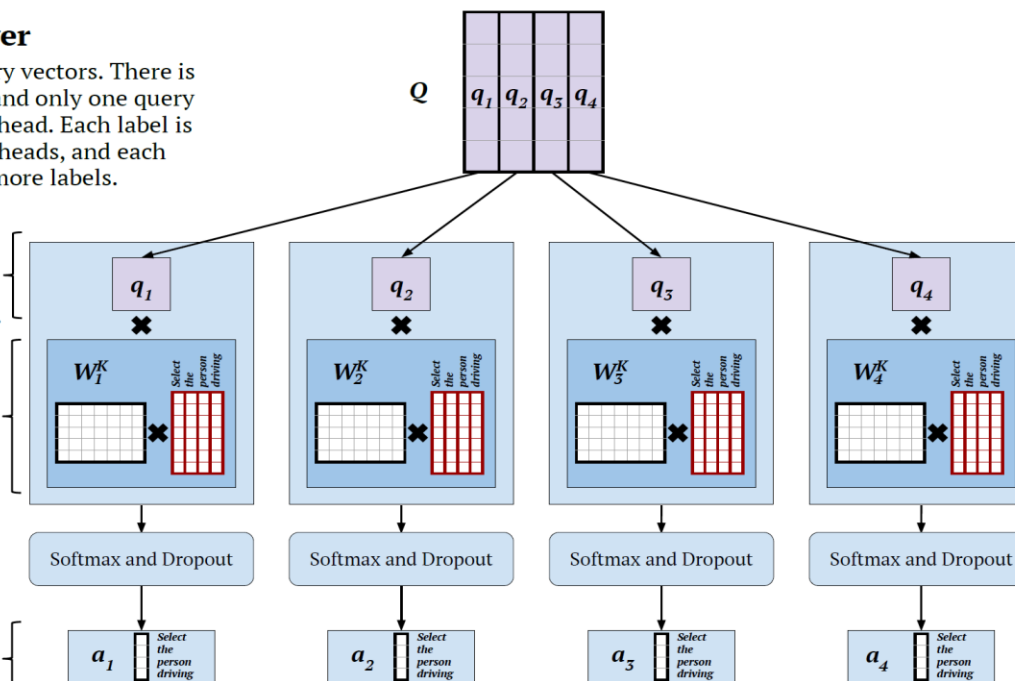


Figure 2: The architecture of the top of our proposed Label Attention Layer. In this figure, the example input sentence is “*Select the person driving*”.

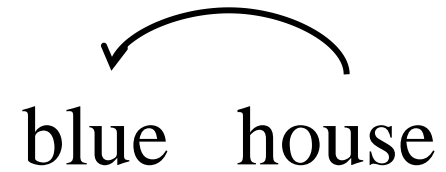
NLP

Introduction to NLP

255.

Dependency Grammars

Dependency structure



- blue
 - modifier, dependent, child, subordinate
- house
 - head, governor, parent, regent

Dependency Structure

Unionized workers are usually better paid than their non-union counterparts.

1

2

3

4

5

6

7

8

9

10

Dependency Structure



Unionized workers are usually better paid than their non-union counterparts.

1

2

3

4

5

6

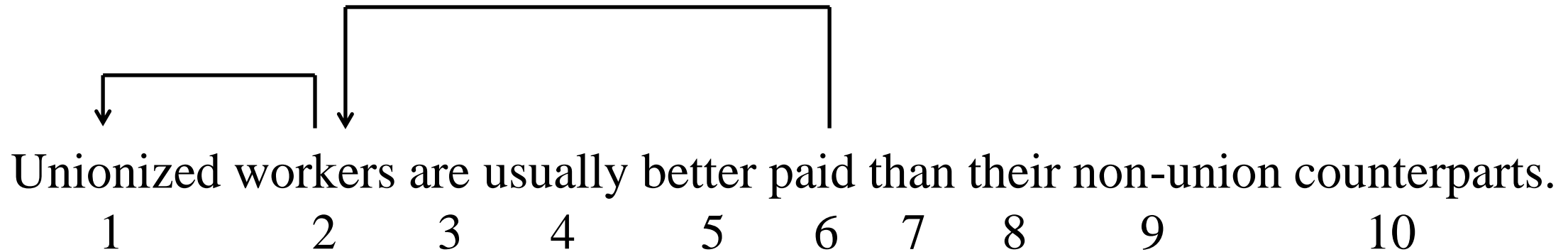
7

8

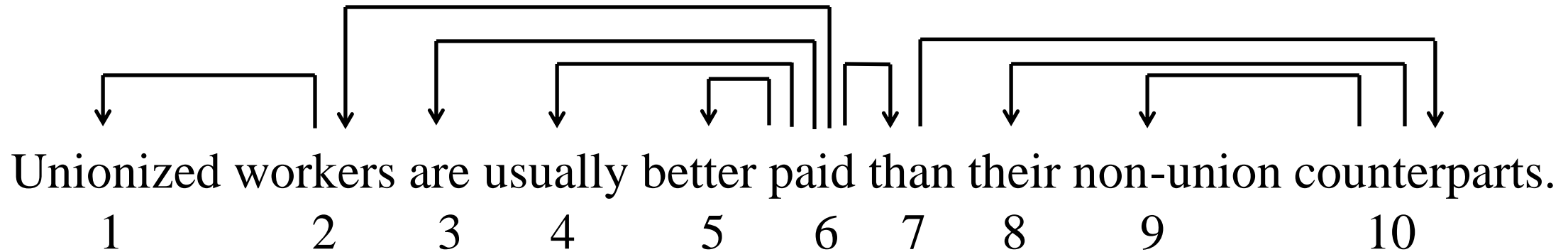
9

10

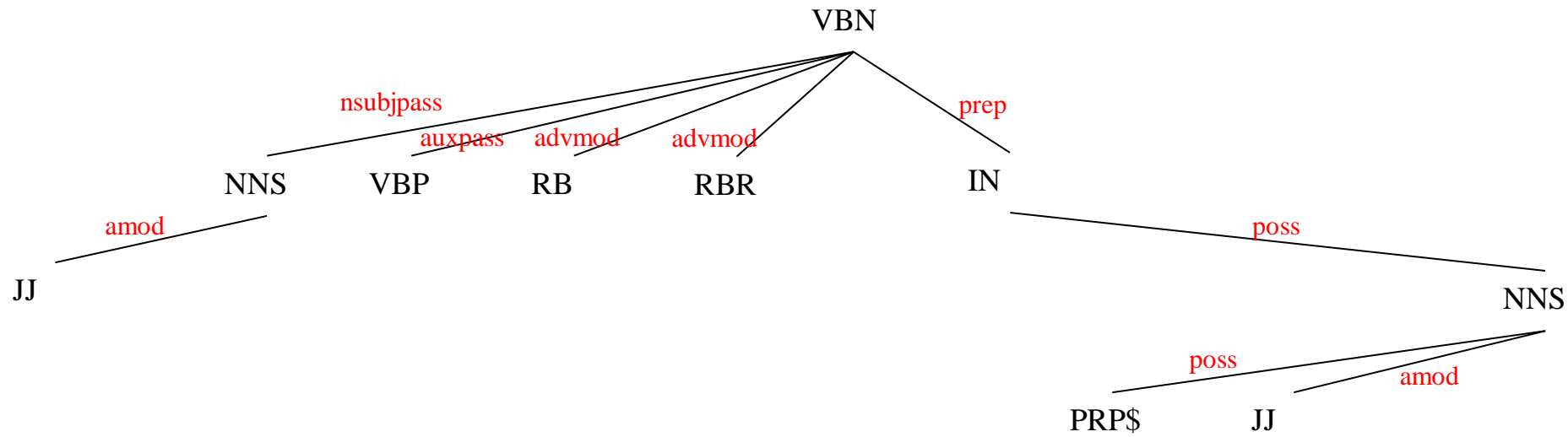
Dependency Structure



Dependency Structure



Other notations



Unionized workers are usually better paid than their non-union counterparts.

1

2

3

4

5

6

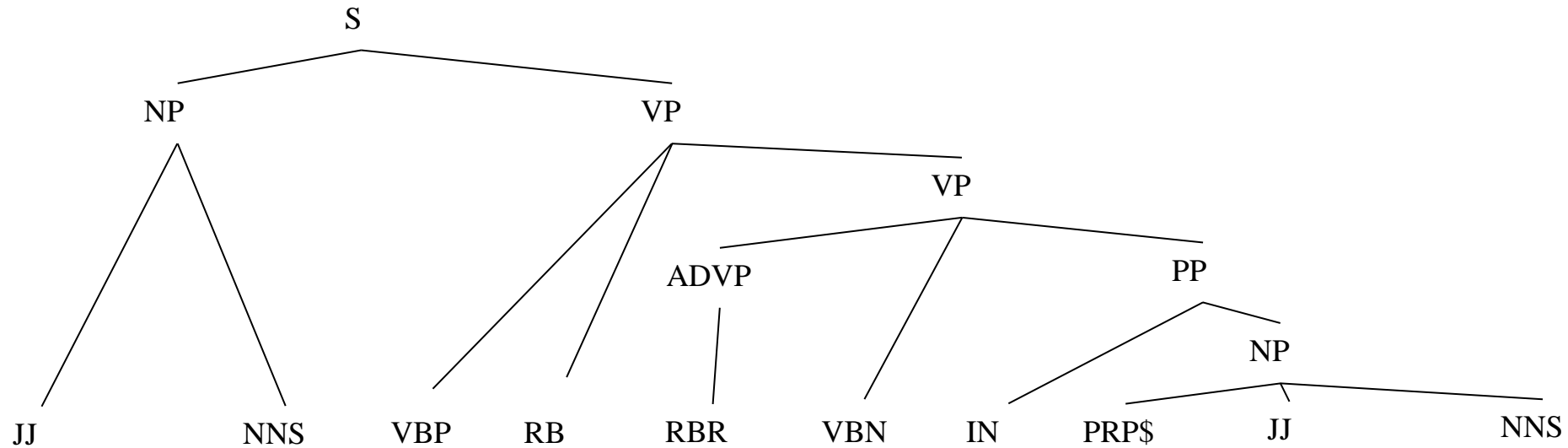
7

8

9

10

Phrase Structure



Unionized workers are usually better paid than their non-union counterparts.

1

2

3

4

5

6

7

8

9

10

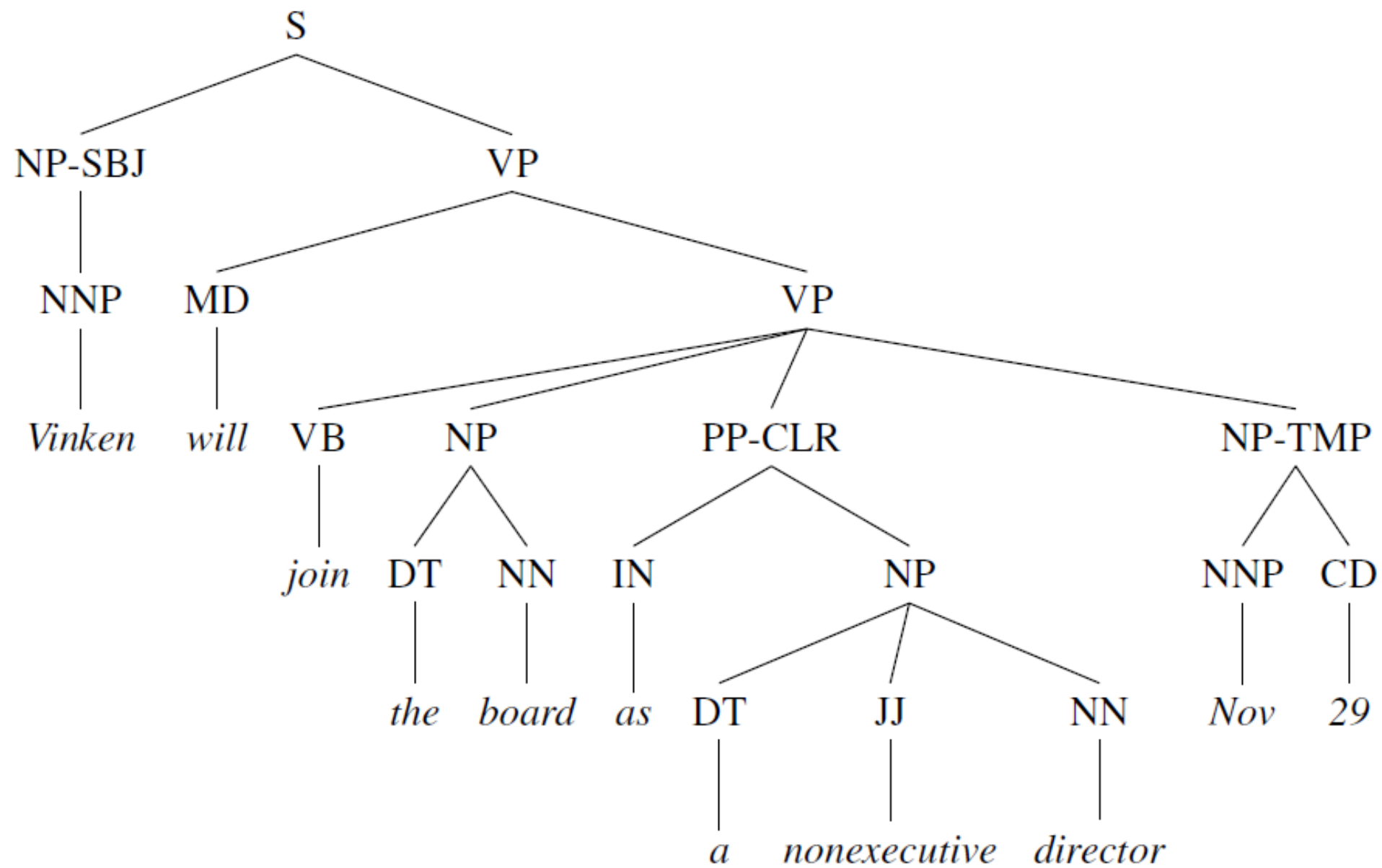
Dependency grammars

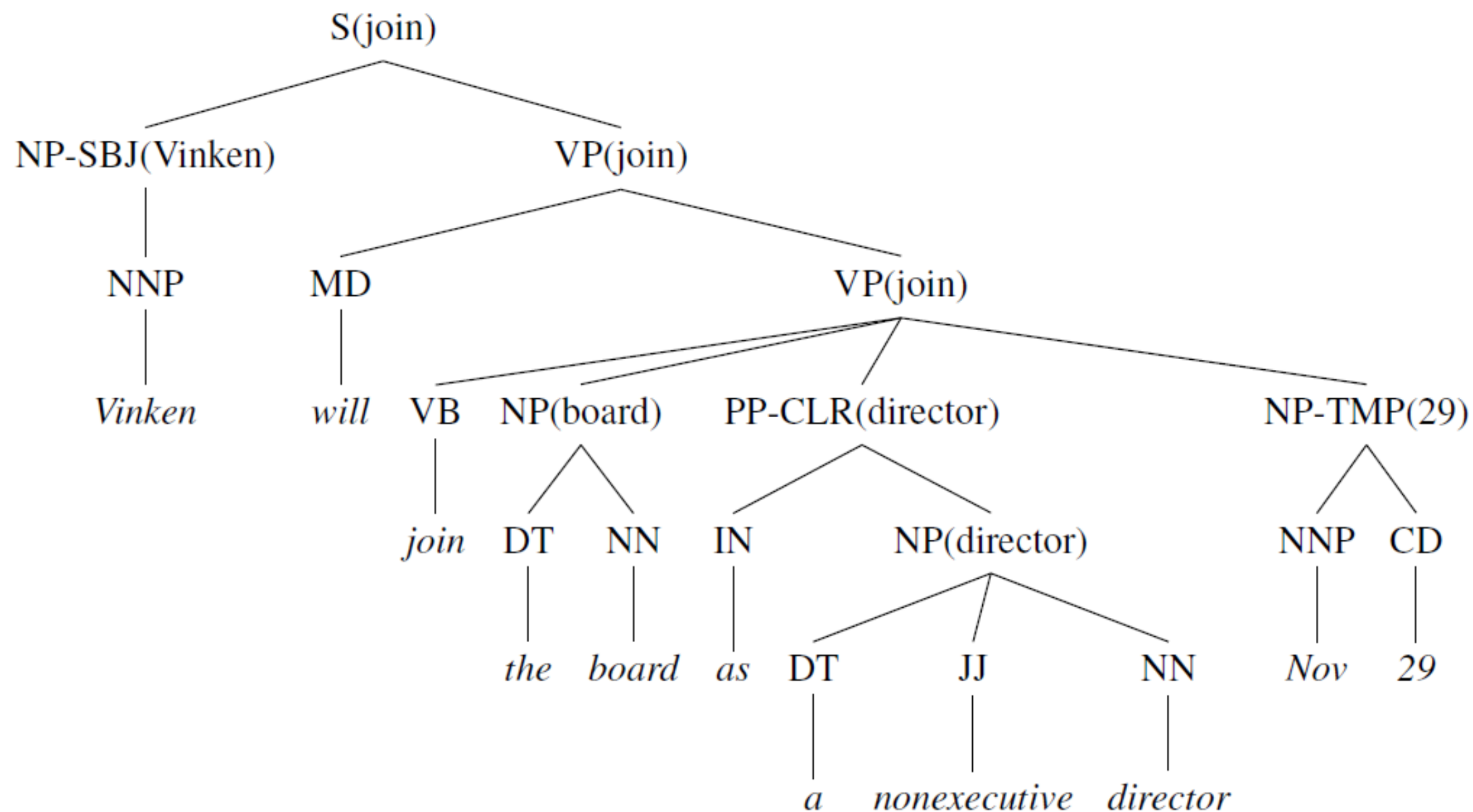
- Characteristics

- Lexical/syntactic dependencies between words
- The top-level predicate of a sentence is the root
- Simpler to parse than context-free grammars
- Particularly useful for free word order languages
- Older idea compared to constituent grammars (as far back as Pāṇini (5th century BCE))

How to identify the heads

- H=head, M=modifier
 - H determines the syntactic category of the construct
 - H determines the semantic category of the construct
 - H is required; M may be skipped
 - Fixed linear position of M with respect to H





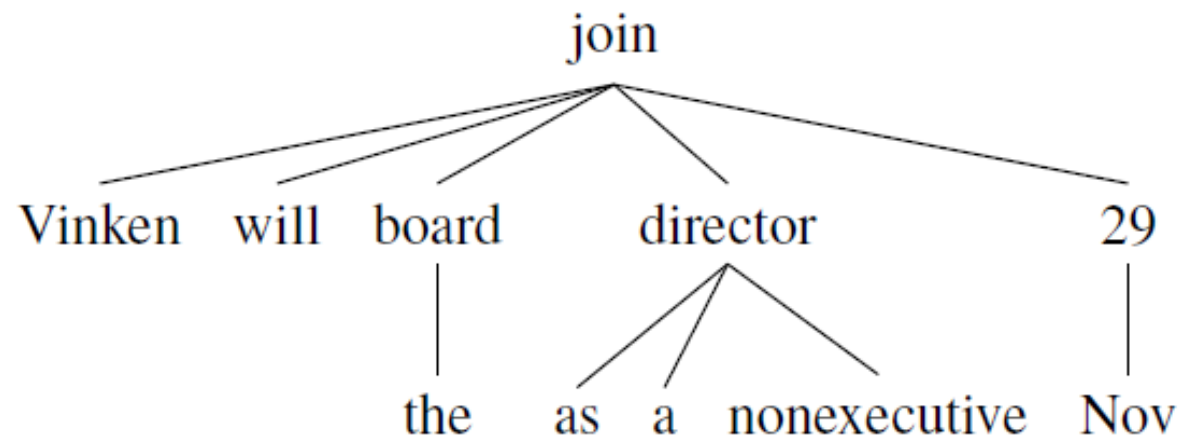
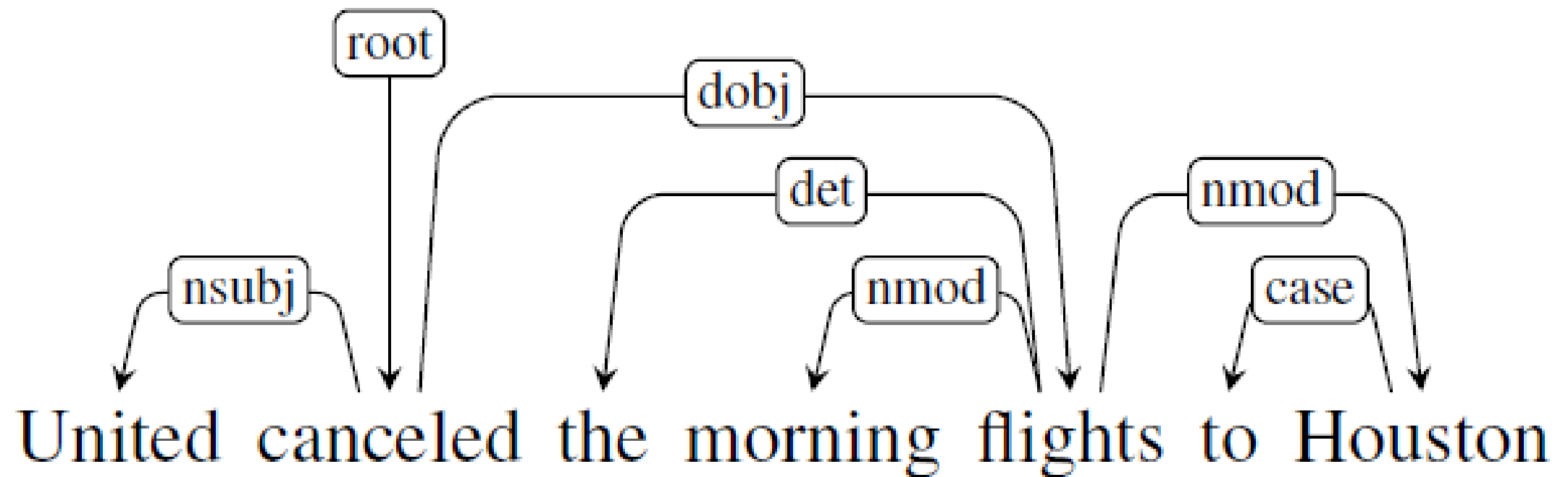


Figure 14.4 A phrase-structure tree from the *Wall Street Journal* component of the Penn Treebank 3.

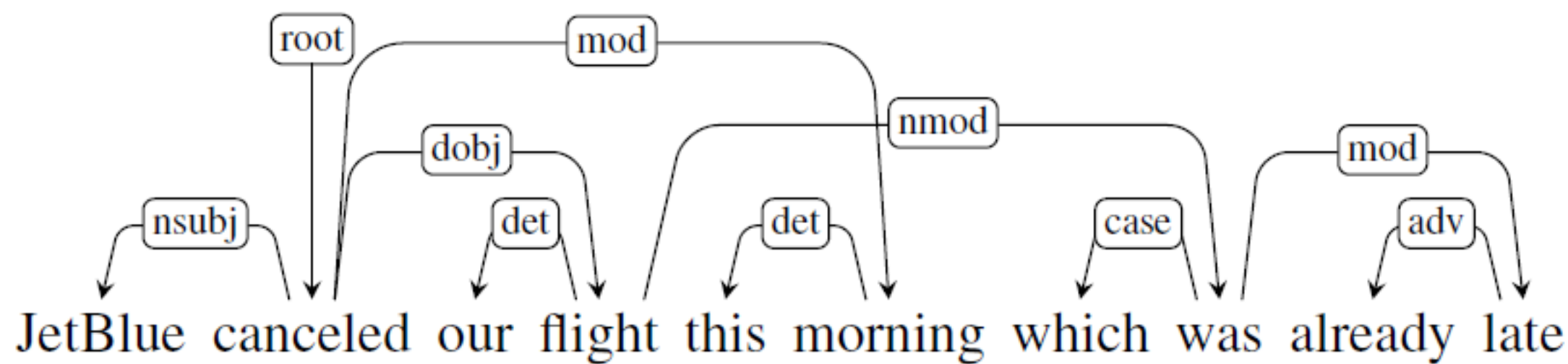
Clausal Argument Relations	Description
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
Nominal Modifier Relations	Description
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
Other Notable Relations	Description
CONJ	Conjunct
CC	Coordinating conjunction

Figure 14.2 Selected dependency relations from the Universal Dependency set. (de Marneffe et al., 2014)



Relation	Examples with <i>head</i> and dependent
NSUBJ	United <i>canceled</i> the flight.
DOBJ	United <i>diverted</i> the flight to Reno. We <i>booked</i> her the first flight to Miami.
IOBJ	We <i>booked</i> her the flight to Miami.
NMOD	We took the morning <i>flight</i> .
AMOD	Book the cheapest <i>flight</i> .
NUMMOD	Before the storm JetBlue canceled 1000 <i>flights</i> .
APPOS	<i>United</i> , a unit of UAL, matched the fares.
DET	The <i>flight</i> was canceled. Which <i>flight</i> was delayed?
CONJ	We <i>flew</i> to Denver and drove to Steamboat.
CC	We flew to Denver and <i>drove</i> to Steamboat.
CASE	Book the flight through <i>Houston</i> .

Figure 14.3 Examples of core Universal Dependency relations.

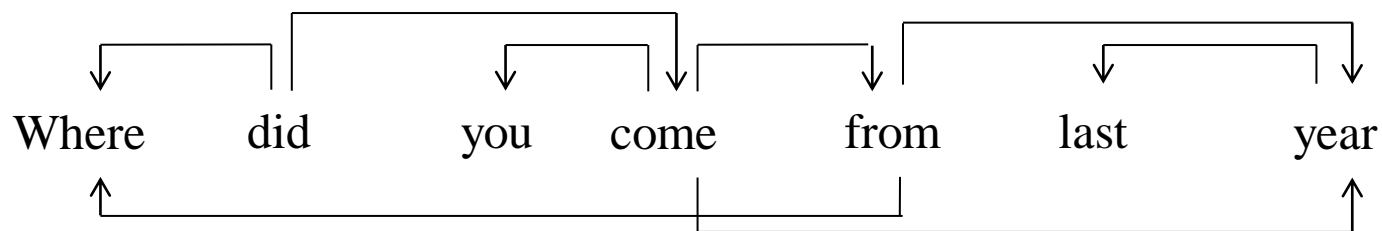


Non-Projectivity

- Rare in English
- Topicalization
 - Cats, I like a lot.
- Extraposition
 - The pizza is ready with pepperoni.

Non-projectivity

Output of (the non-projective) MSTParser



1	Where	Where	WRB	WRB	-	2	SBJ	-	-
2	did	did	VBD	VBD	-	0	ROOT	-	-
3	you	you	PRP	PRP	-	4	SBJ	-	-
4	come	come	VBP	VBP	-	2	VC	-	-
5	from	from	IN	IN	-	4	DIR	-	-
6	last	last	JJ	JJ	-	7	NMOD	-	-
7	year	year	NN	NN	-	5	PMOD	-	-
8	?	?	.	.	-	2	P	-	-

Output of Stanford parser

advmod(come-4, Where-1)
aux(come-4, did-2)
nsubj(come-4, you-3)
root(ROOT-0, come-4)
prep(come-4, from-5)
amod(year-7, last-6)
pobj(from-5, year-7)

Notes

- How to extend a projective method for non-projective parses
 - Use a SWAP operator (Nivre 2009)
- Not clear what to do with conjunctions
 - “cats, dogs, and hamsters”
 - Options: “cats” or “and”

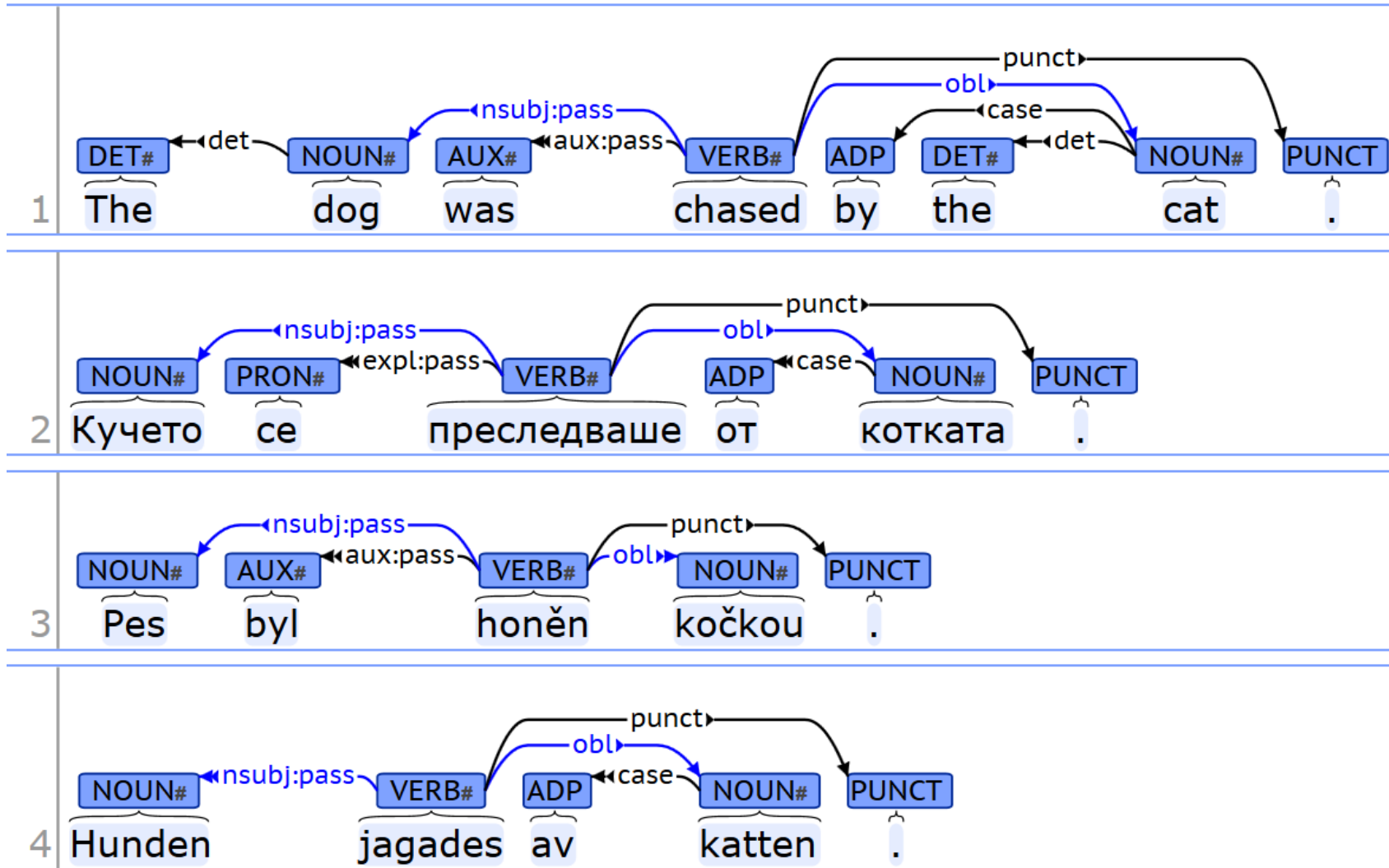
Rate of Non-Projectivity

	#T	#S	#T/#S	%NST	%NPR	%NPS	IR
Arabic	54	1.5	37.2	8.8	0.4	11.2	Yes
Bulgarian	190	12.8	14.8	14.4	0.4	5.4	No
Chinese	337	57	5.9	0.8	0.0	0.0	No
Czech	1249	72.7	17.2	14.9	1.9	23.2	Yes
Danish	94	5.2	18.2	13.9	1.0	15.6	No
Dutch	195	13.3	14.6	11.3	5.4	36.4	No
German	700	39.2	17.8	11.5	2.3	27.8	No
Japanese	151	17	8.9	11.6	1.1	5.3	No
Portuguese	207	9.1	22.8	14.2	1.3	18.9	Yes
Slovene	29	1.5	18.7	17.3	1.9	22.2	Yes
Spanish	89	3.3	27	12.6	0.1	1.7	No
Swedish	191	11	17.3	11.0	1.0	9.8	No
Turkish	58	5	11.5	33.1	1.5	11.6	No

Table 1: Treebank information; #T = number of tokens * 1000, #S = number of sentences * 1000, #T/#S = tokens per sentence, %NST = % of non-scoring tokens, %NPR = % of non-projective relations, %NPS = % of non-projective sentences, IR = has informative root labels

[CoNLL-X data: Hall and Nilsson 2006]

<https://universaldependencies.org/>



NLP

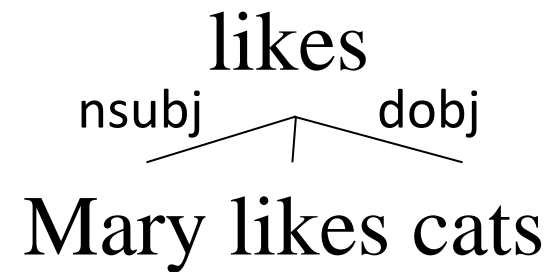
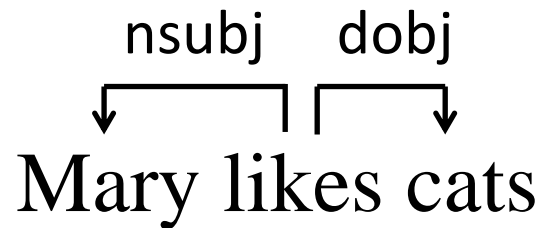
Introduction to NLP

256.

Dependency Parsing

Classic Techniques

- Dynamic programming
 - CKY – similar to lexicalized PCFG, cubic complexity (Eisner 96)

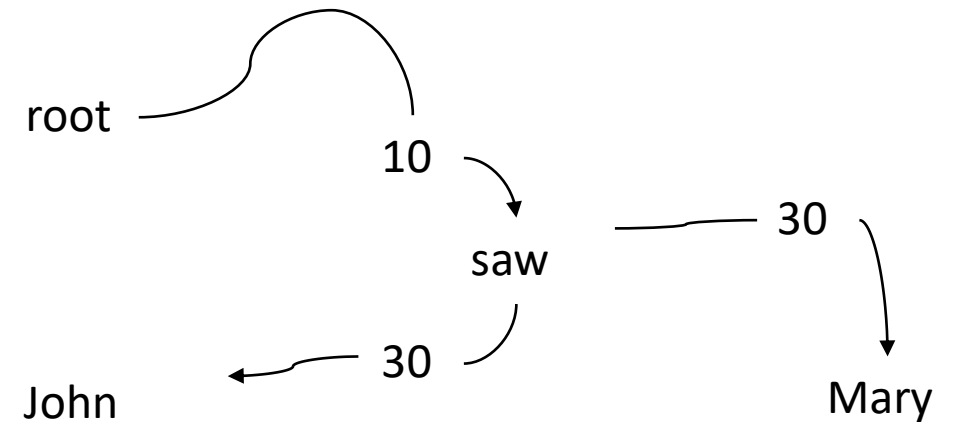
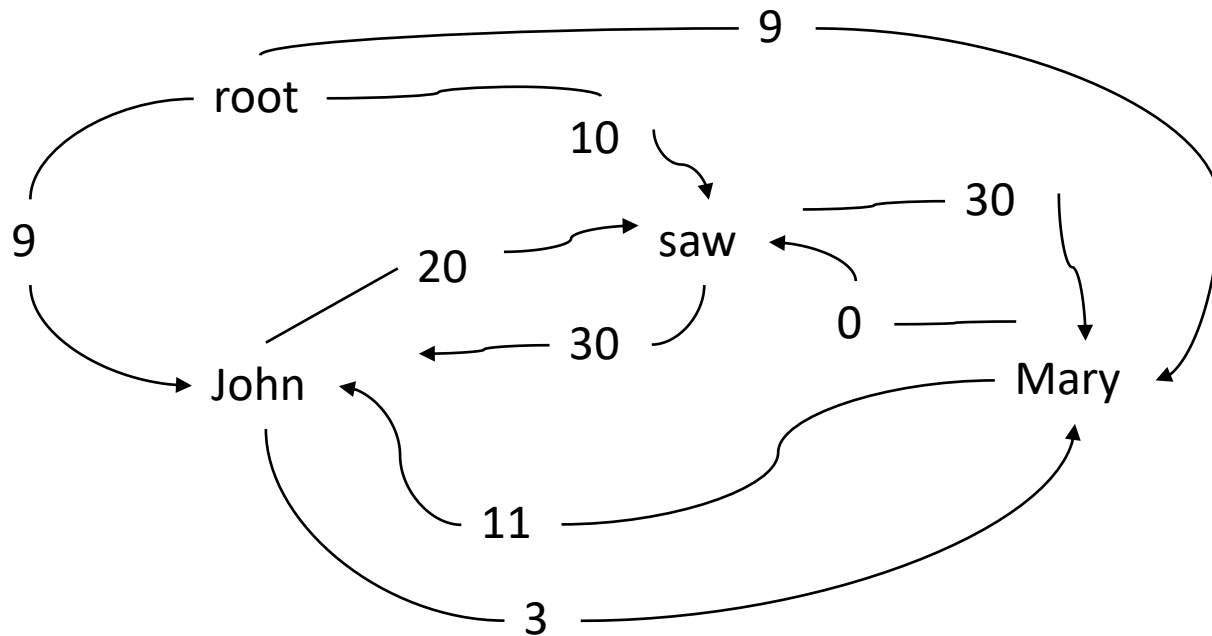


Graph-based Dependency Parsing

- McDonald et al. 2005
- Dependency parsing is equivalent to search for a maximum spanning tree (MST) in a directed graph.
- Efficient algorithm for finding MST for directed graphs
 - Chu and Liu (1965) and Edmonds (1967) give an.

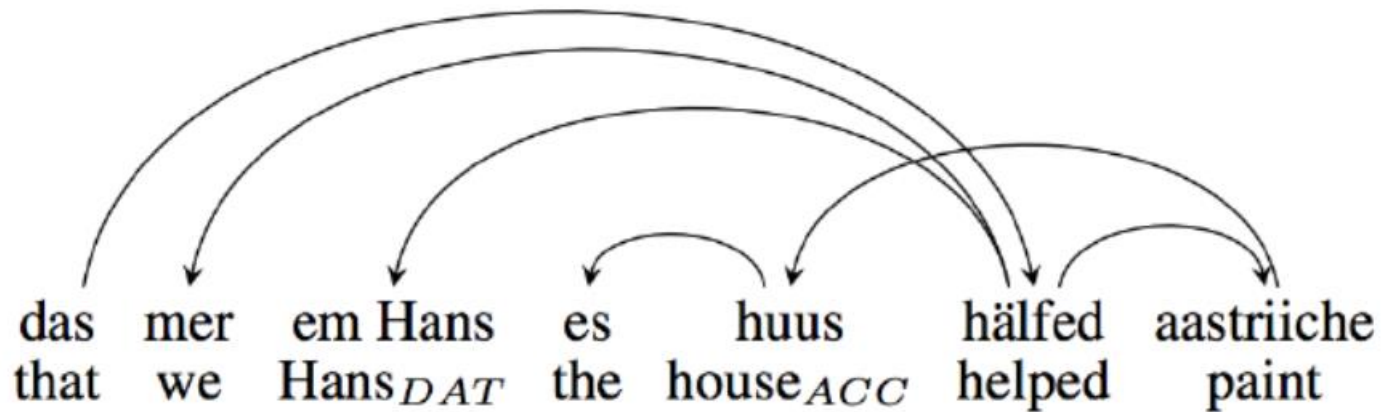
MST Parser example

- Consider the sentence “John saw Mary”
- Recursively remove cycle
- The Chu-Liu-Edmonds algorithm gives the MST on the right hand side (right). This is in general a non-projective tree.



Notes

- Complexity
 - Interestingly, MST is $O(n^2)$, compared with $O(n^3)$ for Eisner, even though MST is non-projective.
- Example of a highly non-projective language
 - Swiss German



[example from Pitler 2013]

NLP

Introduction to NLP

257.

Transition-based Dependency Parsing

Transition-Based Parsing

- Similar to shift-reduce
- Produces a single (projective) tree
- Data structures
 - Stack of partially processed (unattached) words
 - Input buffer
 - Set of dependency arcs
 - Attach the word on the top of the stack to the word at the current position in the buffer (or in the other direction)

Transition-Based Parsing

- Initial configuration
 - Stack (including the root token w_0)
 - Buffer (sentence)
 - Arcs (empty)
- Goal configuration
 - Stack (empty)
 - Buffer (empty)
 - Arcs (complete tree)

Example

- “Book me the morning flight”

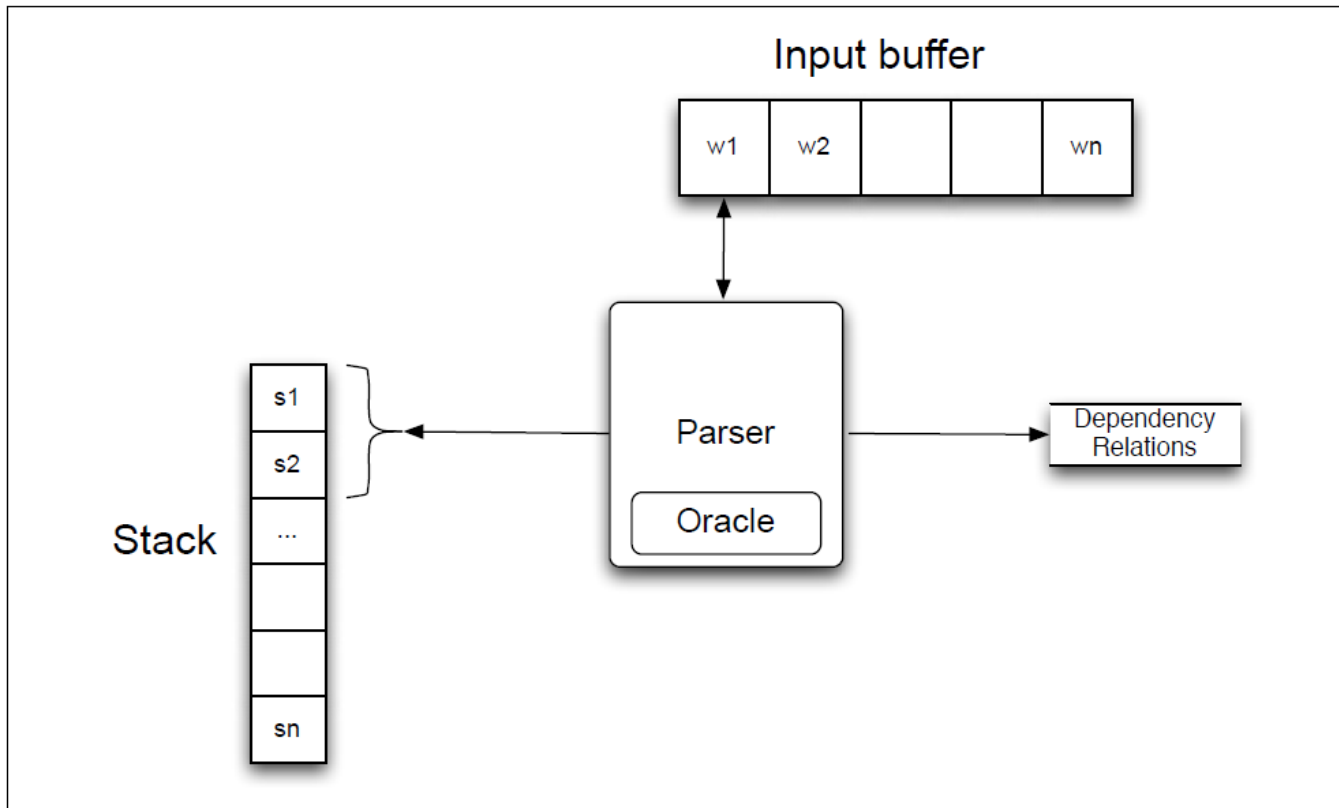


Figure 14.5 Basic transition-based parser. The parser examines the top two elements of the stack and selects an action based on consulting an oracle that examines the current configuration.

[Example from Jurafsky and Martin]

MaltParser (Nivre 2008)

- The reduce operations combine an element from the stack and one from the buffer
- Arc-standard parser
 - The actions are shift, left-arc, right-arc
- Arc-eager parser
 - The actions are shift, reduce, left-arc, right-arc

(Arc-Eager) MaltParser Actions

Shift	$\frac{[\dots]_S \quad [w_i, \dots]_Q}{[\dots, w_i]_S \quad [\dots]_Q}$
Reduce	$\frac{[\dots, w_i]_S \quad [\dots]_Q \quad \exists w_k : w_k \rightarrow w_i}{[\dots]_S \quad [\dots]_Q}$
Left-Arc _r	$\frac{[\dots, w_i]_S \quad [w_j, \dots]_Q \quad \neg \exists w_k : w_k \rightarrow w_i}{[\dots]_S \quad [w_j, \dots]_Q \quad w_i \xleftarrow{r} w_j}$
Right-Arc _r	$\frac{[\dots, w_i]_S \quad [w_j, \dots]_Q \quad \neg \exists w_k : w_k \rightarrow w_j}{[\dots, w_i, w_j]_S \quad [\dots]_Q \quad w_i \xrightarrow{r} w_j}$

[Example from Nivre and Kuebler]

Transition Configuration

	([ROOT],	[Economic, ..., .],	\emptyset)
SH \Rightarrow	([ROOT, Economic],	[news, ..., .],	\emptyset)
LA _{ATT} \Rightarrow	([ROOT],	[news, ..., .],	$A_1 = \{(news, ATT, Economic)\}$)
SH \Rightarrow	([ROOT, news],	[had, ..., .],	A_1)
LA _{SBJ} \Rightarrow	([ROOT],	[had, ..., .],	$A_2 = A_1 \cup \{(had, SBJ, news)\}$)
RA _{PRED} \Rightarrow	([ROOT, had],	[little, ..., .],	$A_3 = A_2 \cup \{(root, PRED, had)\}$)
SH \Rightarrow	([ROOT, had, little],	[effect, ..., .],	A_3)
LA _{ATT} \Rightarrow	([ROOT, had],	[effect, ..., .],	$A_4 = A_3 \cup \{(effect, ATT, little)\}$)
RA _{OBJ} \Rightarrow	([ROOT, had, effect],	[on, ..., .],	$A_5 = A_4 \cup \{(had, OBJ, effect)\}$)
RA _{ATT} \Rightarrow	([ROOT, ... on],	[financial, markets, .],	$A_6 = A_5 \cup \{(effect, ATT, on)\}$)
SH \Rightarrow	([ROOT, ..., financial],	[markets, .],	A_6)
LA _{ATT} \Rightarrow	([ROOT, ... on],	[markets, .],	$A_7 = A_6 \cup \{(markets, ATT, financial)\}$)
RA _{PC} \Rightarrow	([ROOT, ..., markets],	[.],	$A_8 = A_7 \cup \{(on, PC, markets)\}$)
RE \Rightarrow	([ROOT, ..., on],	[.],	A_8)
RE \Rightarrow	([ROOT, ..., effect],	[.],	A_8)
RE \Rightarrow	([ROOT, had],	[.],	A_8)
RA _{PU} \Rightarrow	([ROOT, ..., .],	[.],	$A_9 = A_8 \cup \{(had, PU, .)\}$)

Figure 3.7: Arc-eager transition sequence for the English sentence in figure 1.1 (LA_r = LEFT-ARC_r, RA_r = RIGHT-ARC_r, RE = REDUCE, SH = SHIFT).

[Example from Kuebler, McDonald, Nivre]

Example

- Example: “People want to be free”

- [ROOT] [People, want, to, be, free] \emptyset
- Shift [ROOT, People] [want, to, be, free]
- LA_{nsubj} [ROOT] [want, to, be, free] $A_1 = \{nsubj(want, people)\}$
- RA_{root} [ROOT, want] [to, be, free] $A_2 = A_1 \cup \{root(ROOT, want)\}$

- Characteristics

- Approximate the oracle with a classifier: $o(c) = \operatorname{argmax}_t \mathbf{w} \cdot \mathbf{f}(c, t)$
- There is no search in the greedy version (although beam search also works)
- The final list of arcs is returned as the dependency tree
- Trained on a dependency treebank

- Very fast method

Feature Model

Table 3.1: Feature model for transition-based parsing.

f_i	Address	Attribute
1	STK[0]	FORM
2	BUF[0]	FORM
3	BUF[1]	FORM
4	LDEP(STK[0])	DEPREL
5	RDEP(STK[0])	DEPREL
6	LDEP(BUF[0])	DEPREL
7	RDEP(BUF[0])	DEPREL

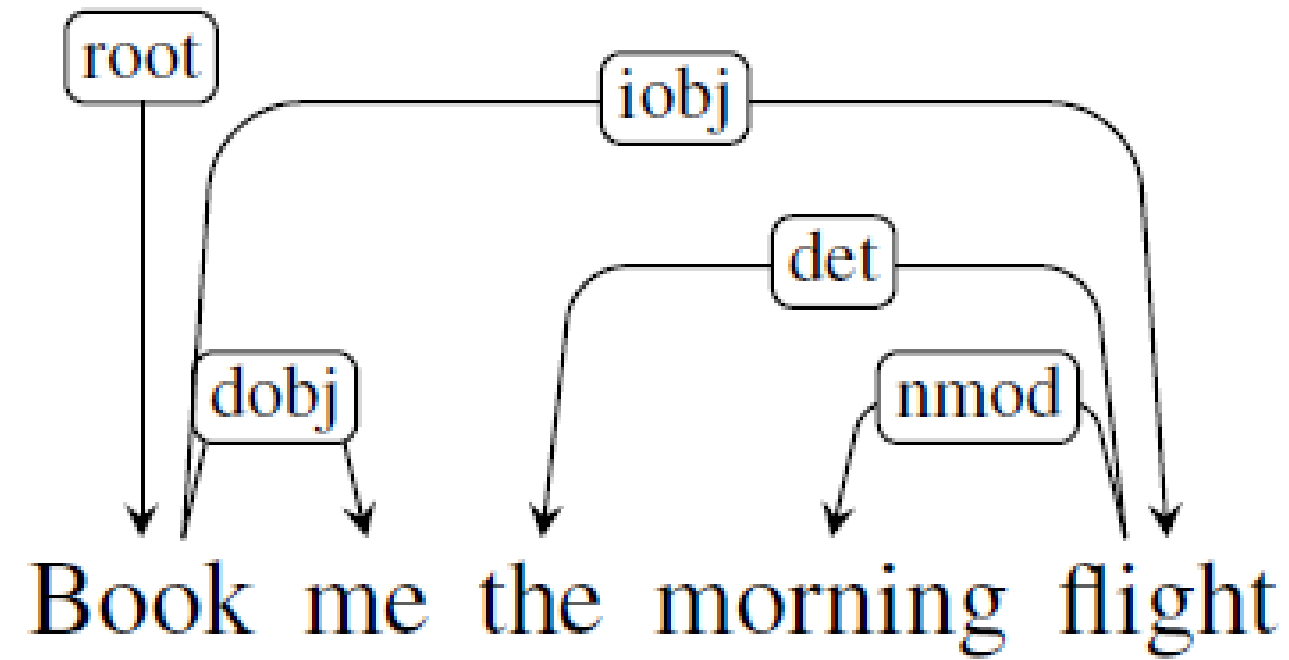
[Example from Kuebler, McDonald, Nivre]

Feature Vectors

$f(c_0)$	=	(ROOT	Economic	news	NULL	NULL	NULL	NULL)
$f(c_1)$	=	(Economic	news	had	NULL	NULL	NULL	NULL)
$f(c_2)$	=	(ROOT	news	had	NULL	NULL	ATT	NULL)
$f(c_3)$	=	(news	had	little	ATT	NULL	NULL	NULL)
$f(c_4)$	=	(ROOT	had	little	NULL	NULL	SBJ	NULL)
$f(c_5)$	=	(had	little	effect	SBJ	NULL	NULL	NULL)
$f(c_6)$	=	(little	effect	on	NULL	NULL	NULL	NULL)
$f(c_7)$	=	(had	effect	on	SBJ	NULL	ATT	NULL)
$f(c_8)$	=	(effect	on	financial	ATT	NULL	NULL	NULL)
$f(c_9)$	=	(on	financial	markets	NULL	NULL	NULL	NULL)
$f(c_{10})$	=	(financial	markets	.	NULL	NULL	NULL	NULL)
$f(c_{11})$	=	(on	markets	.	NULL	NULL	ATT	NULL)
$f(c_{12})$	=	(effect	on	.	ATT	NULL	NULL	ATT)
$f(c_{13})$	=	(had	effect	.	SBJ	NULL	ATT	ATT)
$f(c_{14})$	=	(ROOT	had	.	NULL	NULL	SBJ	OBJ)
$f(c_{15})$	=	(had	.	NULL	SBJ	OBJ	NULL	NULL)
$f(c_{16})$	=	(ROOT	had	NULL	NULL	NULL	SBJ	PU)
$f(c_{17})$	=	(NULL	ROOT	NULL	NULL	NULL	NULL	PRED)
$f(c_{18})$	=	(ROOT	NULL	NULL	NULL	PRED	NULL	NULL)

Figure 3.5: Feature vectors for the configurations in figure 3.2.

[Example from Kuebler, McDonald, Nivre]



Introduction to NLP

258.

Evaluation of Dependency Parsing

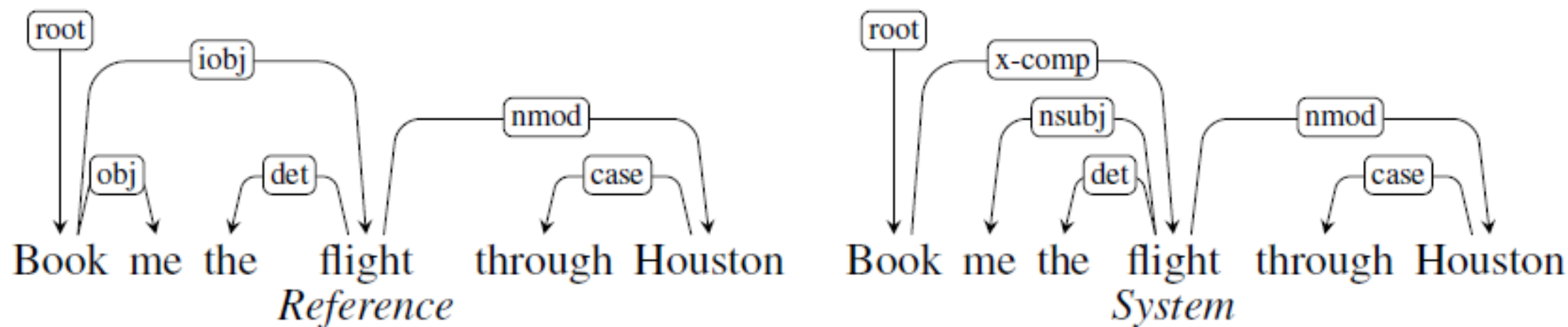


Figure 14.15 Reference and system parses for *Book me the flight through Houston*, resulting in an LAS of 3/6 and an UAS of 4/6.

Evaluation of Dependency Parsing

- Attachment Score (Buchholz & Marsi 2006)
 - # correct deps/# deps (attached to the right head)
 - Unlabeled dependency accuracy (UAS)
 - Labeled dependency accuracy (LAS)

1	Unionized	Unionized	VBN	VBN	–	2	NMOD	–	–
2	workers	workers	NNS	NNS	–	3	SBJ	–	–
3	are	are	VBP	VBP	–	0	ROOT	–	–
4	usually	usually	RB	RB	–	3	TMP	–	–
5	better	better	RBR	RBR	–	4	ADV	–	–
6	paid	paid	VBN	VBN	–	5	AMOD	–	–
7	than	than	IN	IN	–	5	AMOD	–	–
8	their	their	PRP\$	PRP\$	–	10	NMOD	–	–
9	non-union	non-union	JJ	JJ	–	10	NMOD	–	–
10	counterparts	counterparts	NNS	NNS	–	7	PMOD	–	–

External Links

- <http://ilk.uvt.nl/conll/>
 - CONLL-X Shared task
- <http://ufal.mff.cuni.cz/pdt2.0/>
 - Prague Dependency Treebank
- <http://nextens.uvt.nl/depparse-wiki/SharedTaskWebsite>
- <http://nextens.uvt.nl/depparse-wiki/DataOverview>
- <http://maltparser.org/>
 - Joakim Nivre's Maltparser
- <http://www.cs.ualberta.ca/~lindek/minipar.htm>
 - Dekang Lin's Minipar
- <http://www.link.cs.cmu.edu/link/>
 - Daniel Sleator and Davy Temperley's Link parser

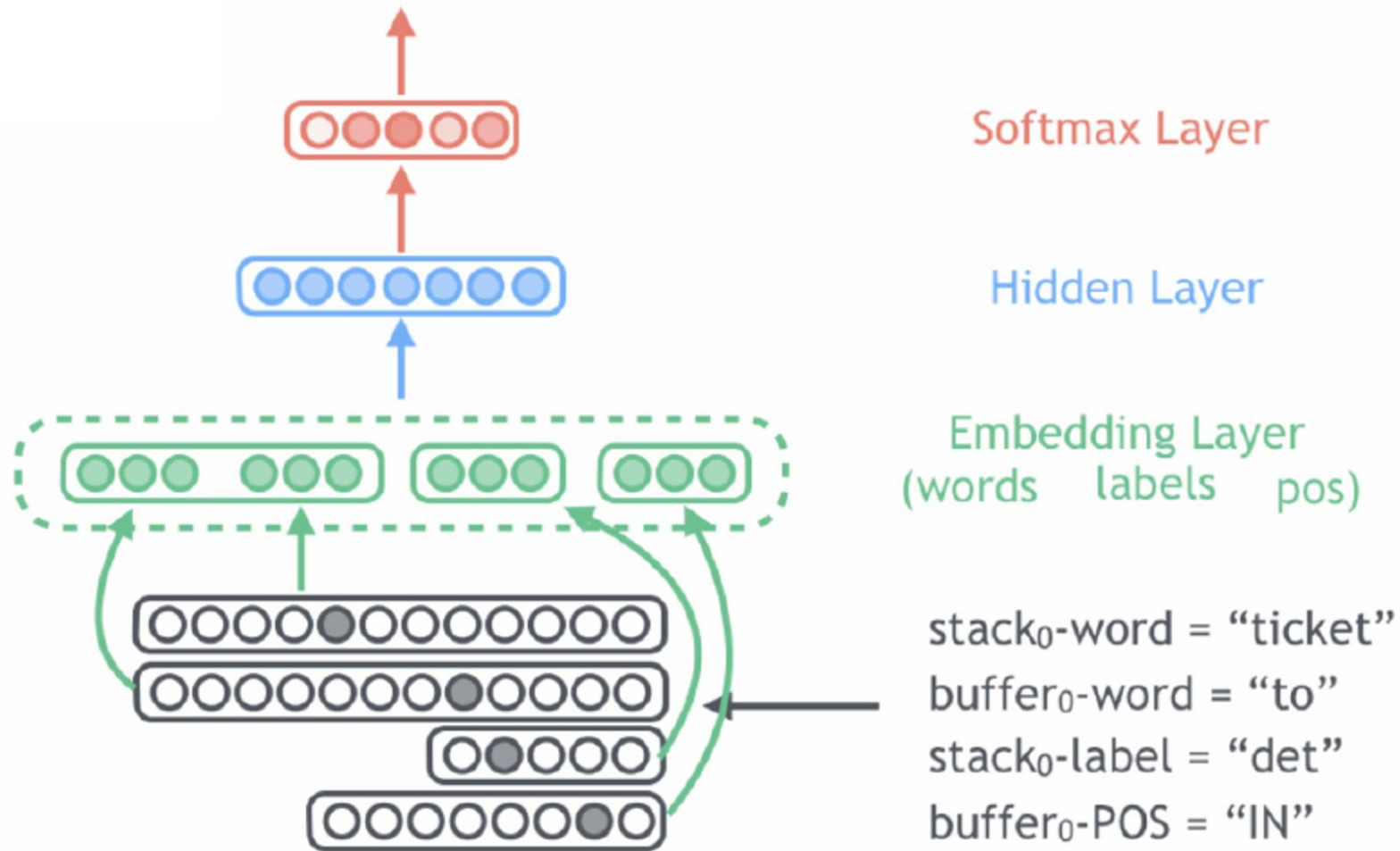
Notes

- The original versions of MSTParser and MaltParser from 2007 achieve about 81% accuracy
 - Highest in Japanese (91-92%)
 - Lowest in Arabic and Turkish (63-67%)
- Non-projective parsing is harder than projective parsing

Introduction to NLP

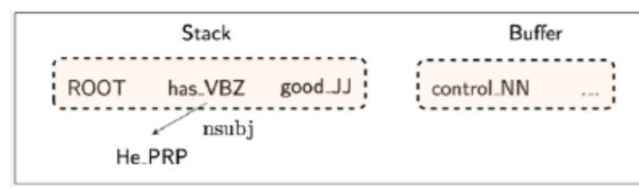
Neural Dependency Parsing

Neural dependency parsing

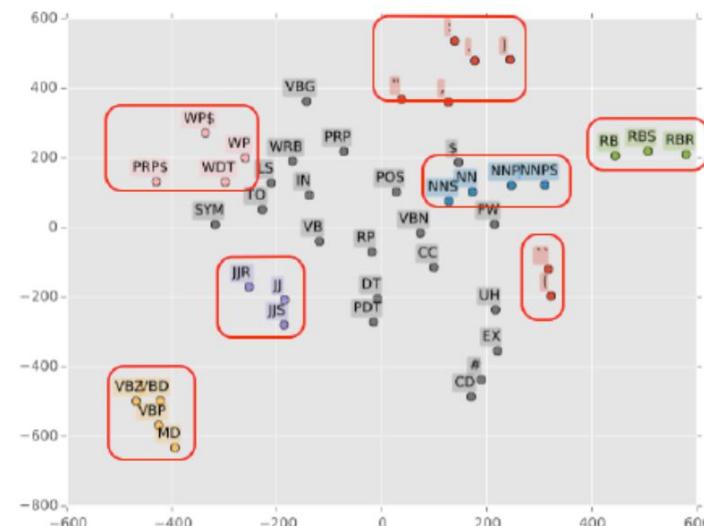


Neural dependency parsing

- Used pre-trained word embeddings
- Part-of-speech tags and dependency labels are also represented as vectors
- No feature template any more!



	word	POS	dep.
s1	good	JJ	∅
s2	has	VBZ	∅
b1	control	NN	∅
lc(s1)	→ ∅	+	∅
rc(s1)	∅	+	∅
lc(s2)	He	PRP	nsubj
rc(s2)	∅	∅	∅



- A simple feedforward NN: what is left is backpropagation!

(Chen and Manning, 2014): A Fast and Accurate Dependency Parser using Neural Networks

NLP

Introduction to NLP

263.

Tree Adjoining Grammars

The Chomsky Hierarchy

- Regular languages
 - Recognized by Finite State Automata
- Context-free languages
 - Recognized by Push Down Automata
- Context-sensitive languages
- Recursively enumerable languages

The Chomsky Hierarchy

- Regular languages
- Context-free languages
- **Mildly context-sensitive languages**
- Context-sensitive languages
- Recursively enumerable languages

Mildly Context-Sensitive Grammars

- Superset of CFG
- Polynomial parsing
 - $O(n^6)$
- Constant growth property
 - (string length grows linearly)
- Cannot handle the language of strings with the same number of as, bs, and cs.

[https://en.wikipedia.org/wiki/Mildly_context-sensitive_grammar_formalism]

[Example from Julia Hockenmaier]

Other Formalisms

- Tree Substitution Grammar (TSG)
 - Terminals generate entire tree fragments
 - TSG and CFG are formally equivalent

Mildly Context-Sensitive Grammars

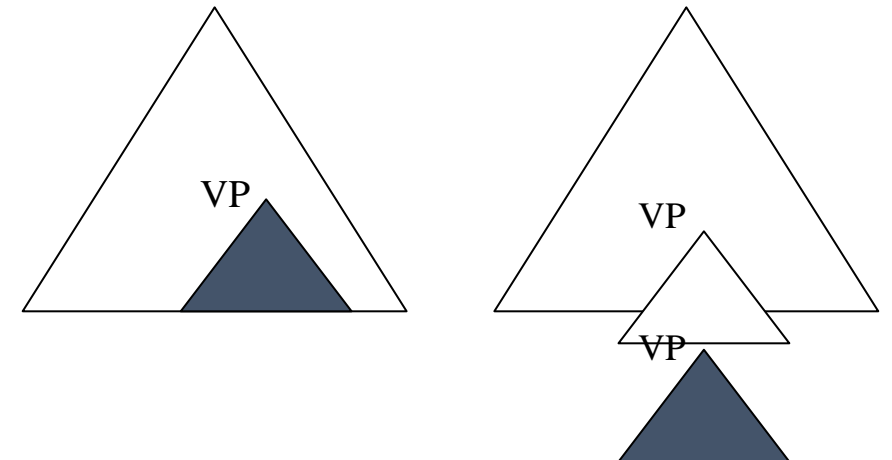
- More powerful than TSG
- Examples:
 - Tree Adjoining Grammar (TAG)
 - Combinatory Categorical Grammar (CCG)

(Tree) Operations in TAG

- Substitution
 - Insert an initial tree to the bottom of a tree
- Adjunction (not in TSG)
 - Insert an auxiliary tree fragment in the middle of a tree
 - Used for long-distance dependencies and for optional modifiers
- Features
 - Each elementary tree has features that can be associated with the top half and with the bottom half
 - Unification is needed
- (Lexicalization)
 - LTAG: each initial or auxiliary tree is labeled with a lexical item

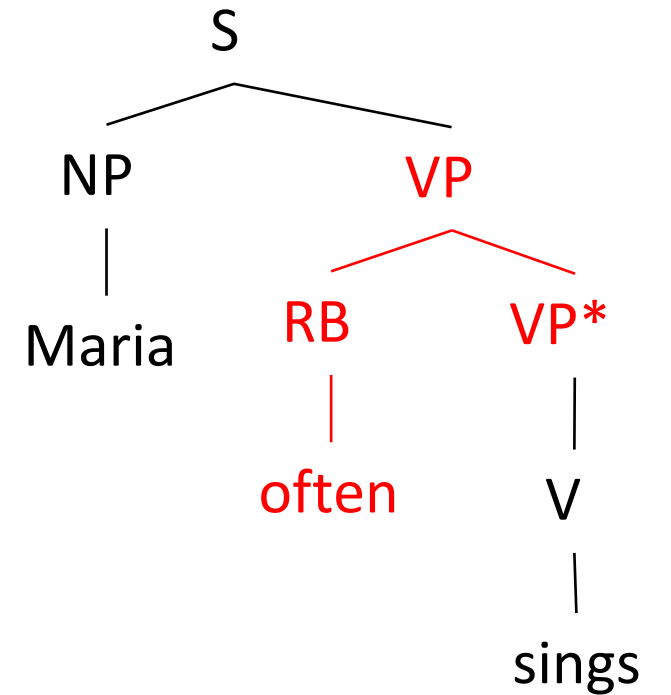
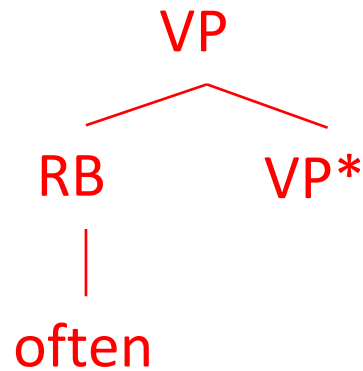
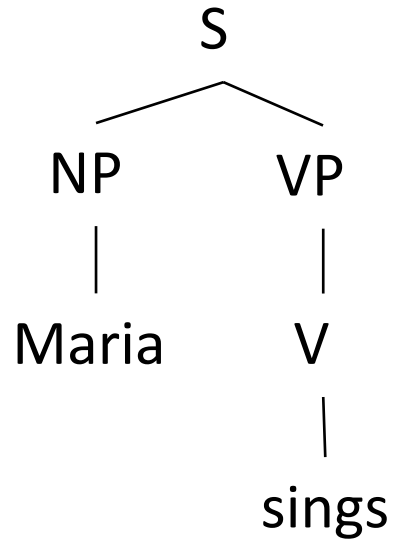
Tree Adjoining Grammar (TAG)

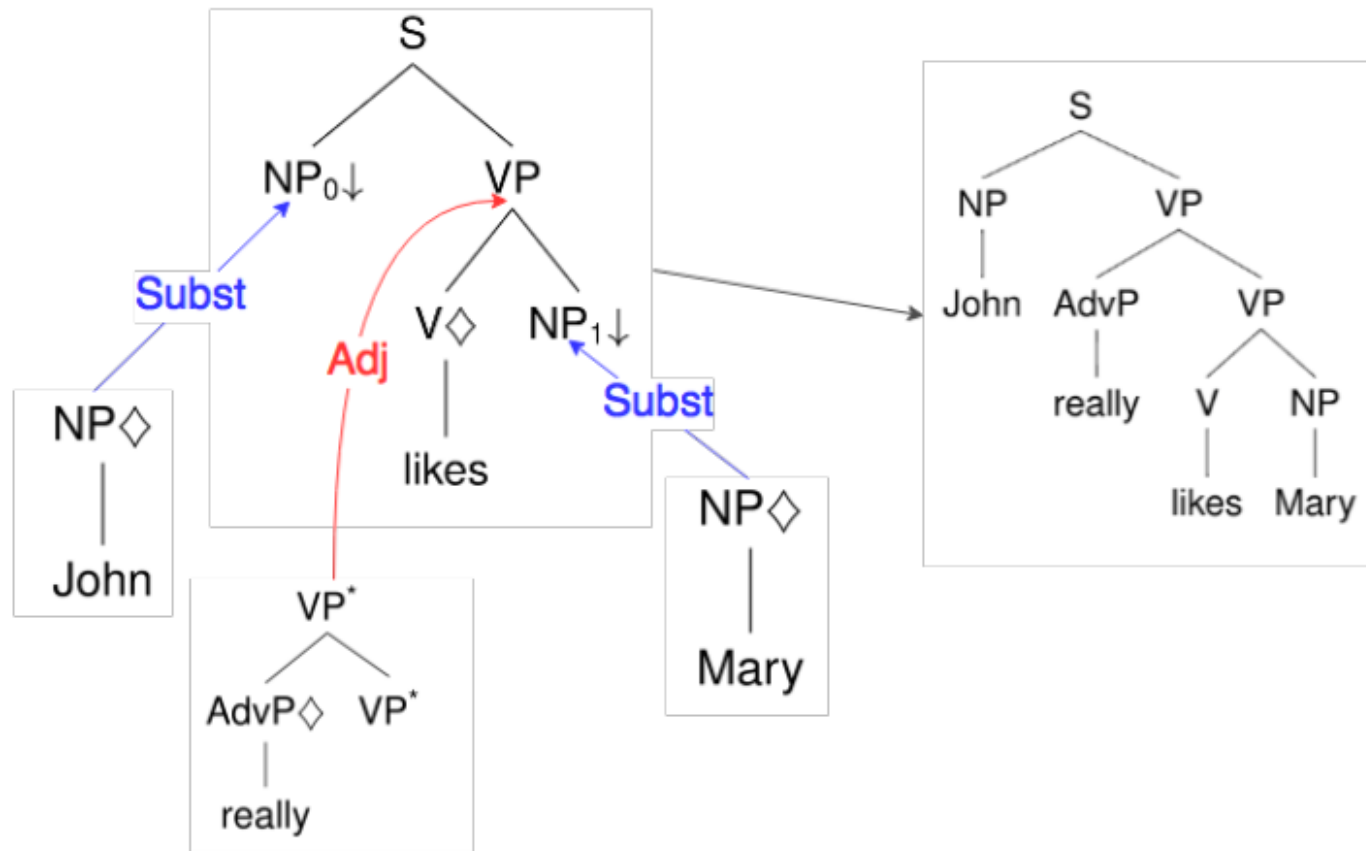
- It can generate languages like $a^n b^n c^n d^n$ or ww (cross-serial dependencies):
 - e.g., Mary gave a book and a magazine to Chen and Mike, respectively.
- Expressive power
 - TAG is formally more powerful than CFG
 - TAG is less powerful than CSG
- Card game online (*broken links*)
 - <http://www.ltaggame.com/>
 - <http://www.ltaggame.com/family.html>



TAG Example

- Maria sings
- Maria *often* sings (optional modifier)





Adjunction allows for unbounded recursion while still enforcing agreement.

John **smartly occasionally really only** likes Mary...

[example from Jungo Kasai]

NLP