

# Natural Language Processing (7)

## Parsing (1): Constituency Parsing

Daisuke Kawahara

Department of Communications and Computer Engineering,  
Waseda University

# Lecture Plan

1. Overview of Natural Language Processing
2. Formal Language Theory
3. Word Senses and Embeddings
4. Topic Models
5. Collocations, Language Models, and Recurrent Neural Networks
6. Sequence Labeling and Morphological Analysis
7. Parsing (1)
8. Parsing (2)
9. Transfer Learning
10. Knowledge Acquisition
11. Information Retrieval, Question Answering, and Machine Translation
12. Guest Talk (1)
13. Guest Talk (2)
14. Project: Survey or Programming
15. Project Presentation

# Table of Contents

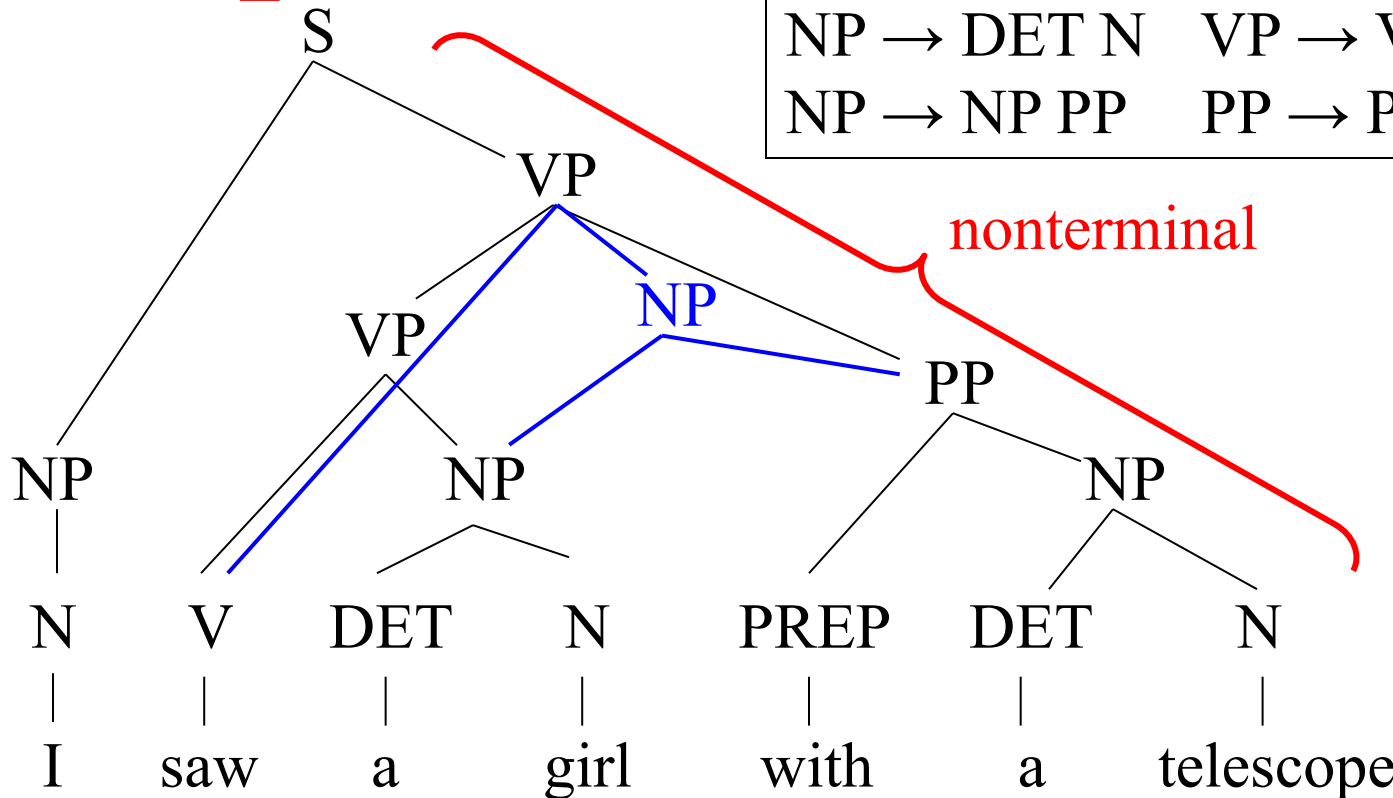
- Review of CFG and CKY parsing
- Probabilistic CFG
- Treebanks
- Extensions
  - Lexicalization
  - History
  - Nonterminal Classification
- Evaluation criteria and SOTA of English parsing

# CFG and Syntactic Structure

production rule

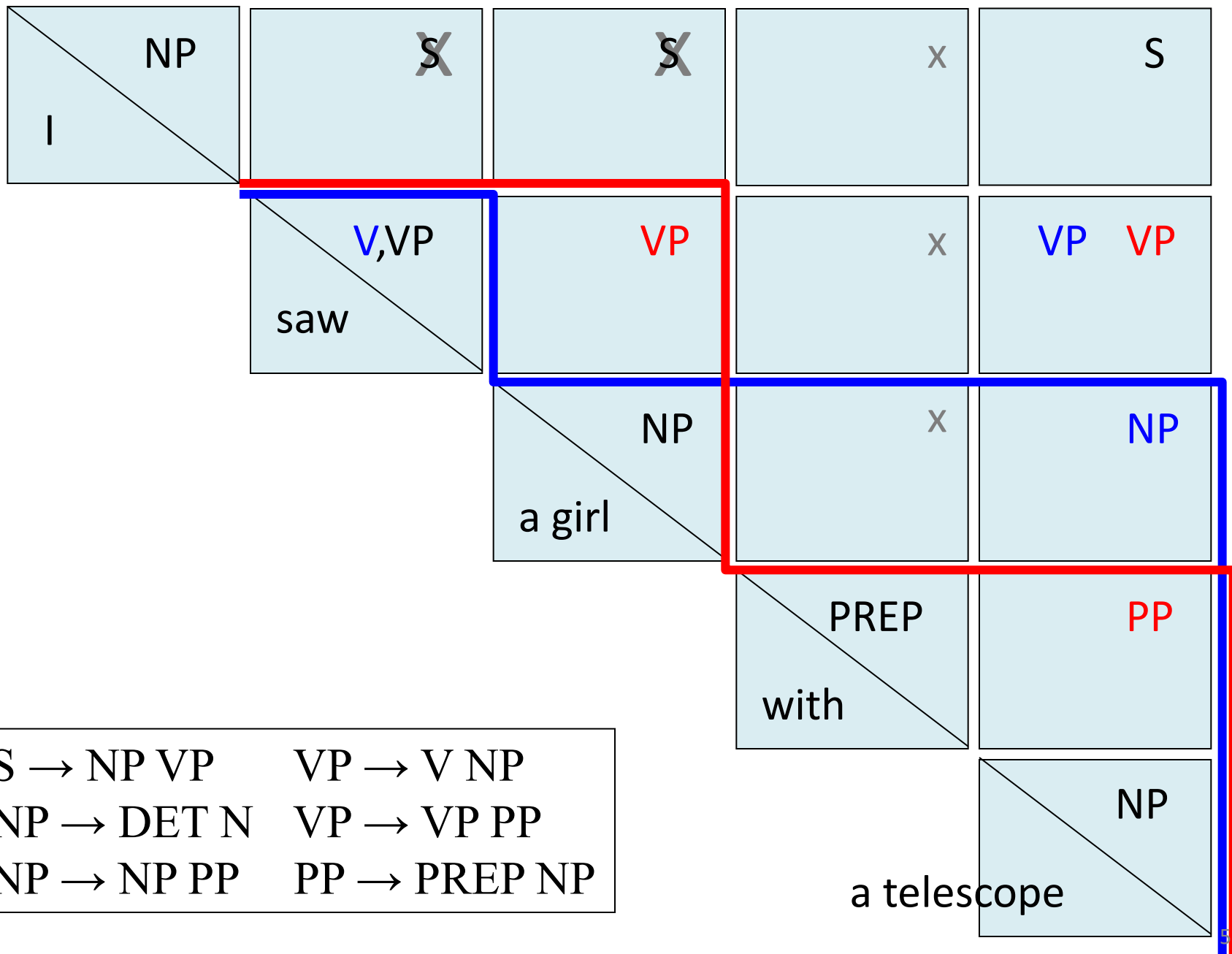
$S \rightarrow NP VP$	$VP \rightarrow V$
$NP \rightarrow N$	$VP \rightarrow V NP$
$NP \rightarrow DET N$	$VP \rightarrow VP PP$
$NP \rightarrow NP PP$	$PP \rightarrow PREP NP$

start symbol →



nonterminal

terminal →



$S \rightarrow NP VP$	$VP \rightarrow V NP$
$NP \rightarrow DET N$	$VP \rightarrow VP PP$
$NP \rightarrow NP PP$	$PP \rightarrow PREP NP$

# Probabilistic Context Free Grammar (PCFG)

- A set of terminals  $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals  $\{N^i\}, i = 1, \dots, n$
- A designated start symbol  $N^1$
- A set of rules  $\{N^i \rightarrow \zeta^j\}$
- A corresponding set of probabilities on rules

$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

# A simple PCFG

$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4
$PP \rightarrow P NP$	1.0	$NP \rightarrow \textit{astronomers}$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow \textit{ears}$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow \textit{saw}$	0.04
$P \rightarrow \textit{with}$	1.0	$NP \rightarrow \textit{stars}$	0.18
$V \rightarrow \textit{saw}$	1.0	$NP \rightarrow \textit{telescopes}$	0.1

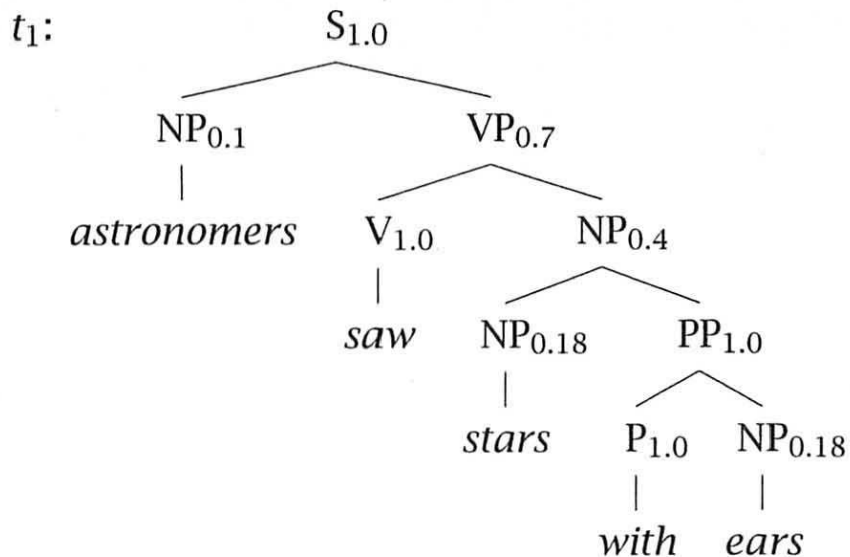
$$\sum P(NP \rightarrow *) = 1$$

Chomsky Normal Form

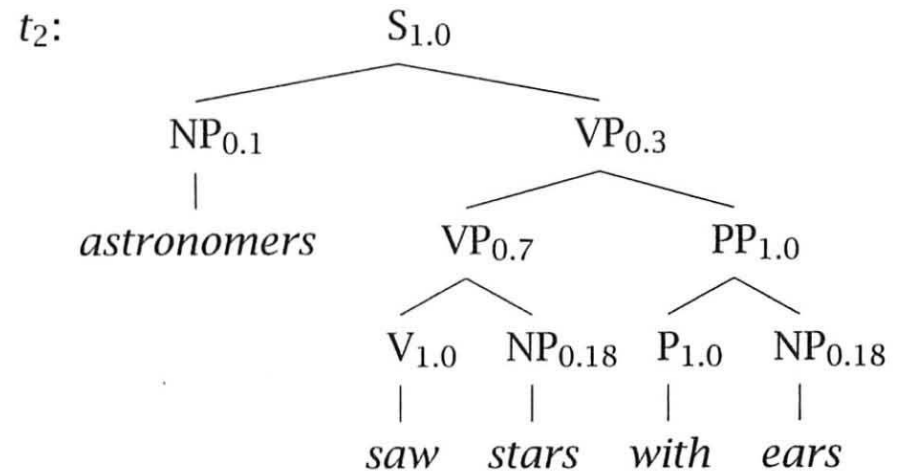
- $A \rightarrow BC$
- $A \rightarrow \alpha$

# Parse Trees

*astronomers saw stars with ears*



$$P(t_1) = 0.0009072$$



$$P(t_2) = 0.0006804$$

$$P(S) = P(t_1) + P(t_2) = 0.0015876$$



# Probabilistic Context Free Grammar (PCFG)

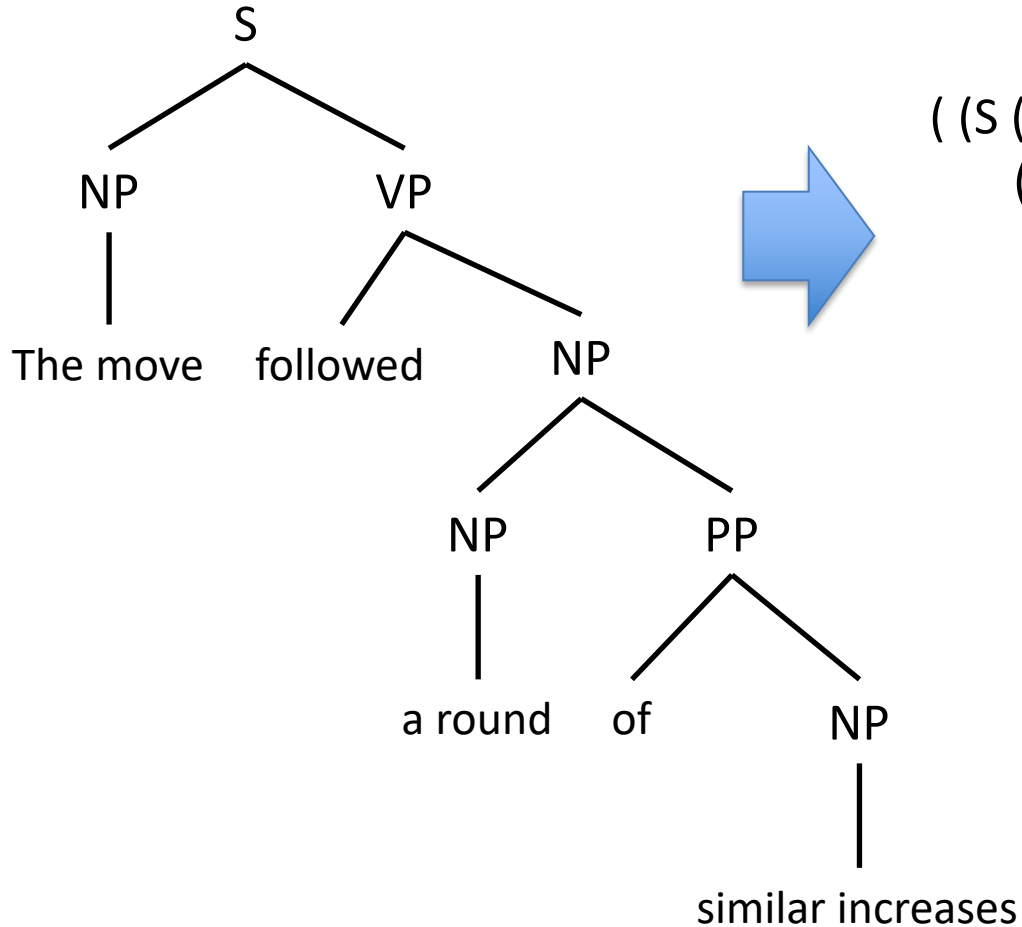
- A set of terminals  $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals  $\{N^i\}, i = 1, \dots, n$
- A designated start symbol  $N^1$
- A set of rules  $\{N^i \rightarrow \zeta^j\}$
- A corresponding set of probabilities on rules

$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

# Penn Treebank [Marcus+ 1993]

- An annotated corpus for English
  - Wall Street Journal (1 million words)
  - Brown Corpus
  - Switchboard Corpus (telephone conversation)
  - ATIS (Air Travel Information System) Corpus
- Released by LDC
- The de facto data for English parsing (training and evaluation)

# A Penn Treebank Tree



( (S (NP The move)  
(VP followed  
(NP (NP a round)  
(PP of  
(NP similar increases))))))

# Phrasal Categories of Penn Treebank

S	Simple clause (sentence)	CONJP	Multiword conjunction phrase
SBAR	S' clause with complementizer	FRAG	Fragment
SBARQ	Wh-question S' clause	INTJ	Interjection
SQ	Inverted Yes/No question S' clause	LST	List marker
SINV	Declarative inverted S' clause	NAC	Not a consistent grouping
ADJP	Adjective phrase	NX	Nominal constituent inside NP
ADVP	Adverbial phrase	PRN	Parenthetical
NP	Noun phrase	PRT	Particle
PP	Prepositional phrase	RRC	Reduced relative clause
QP	Quantifier phrase (inside NP)	UCP	Unlike coordinated phrase
VP	Verb phrase	X	Unknown or uncertain
WHNP	Wh-noun phrase	WHADJP	Wh-adjective phrase
WHPP	Wh-prepositional phrase	WHADVP	Wh-adverbial phrase

# Exercise

( (S (NP (NP (NNP Pierre) (NNP Vinken))  
    (, ,)  
    (NP (NP (CD 61) (NNS years))  
        (ADJP (JJ old))))  
    (, ,))  
(MD will)  
(VP (VB join)  
    (NP (DT the) (NN board))  
    (PP (IN as)  
        (NP (DT a) (JJ nonexecutive) (NN director)))  
    (ADVP (NP (NNP Nov.) (CD 29))))))  
(. .))

# Probabilistic Context Free Grammar (PCFG)

- A set of terminals  $\{w^k\}, k = 1, \dots, V$
- A set of nonterminals  $\{N^i\}, i = 1, \dots, n$
- A designated start symbol  $N^1$
- A set of rules  $\{N^i \rightarrow \zeta^j\}$
- A corresponding set of probabilities on rules

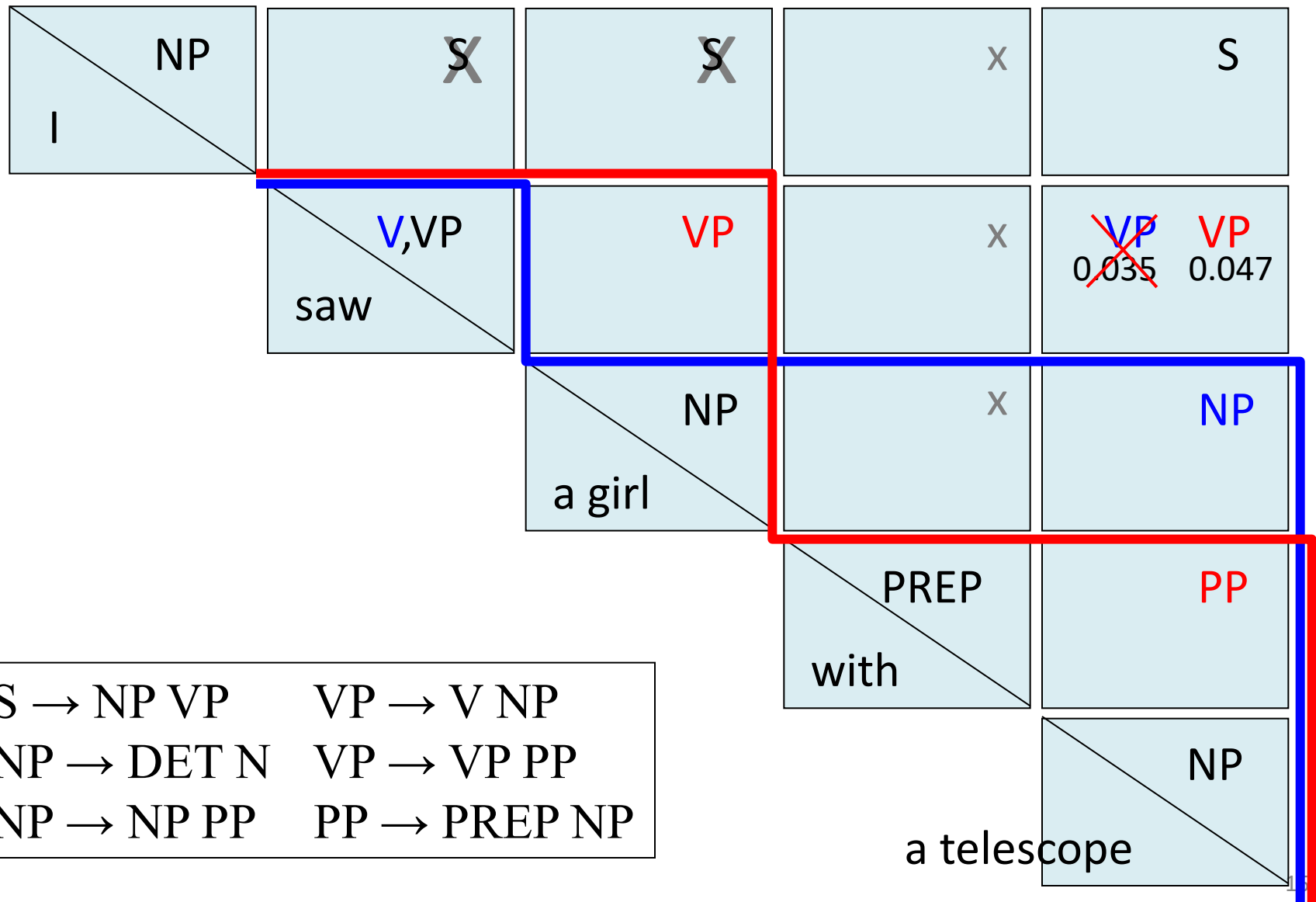
$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$



Maximum Likelihood Estimation using a treebank

$$\hat{P}(N^i \rightarrow \zeta^j) = \frac{C(N^i \rightarrow \zeta^j)}{C(N^i)}$$

# A Dynamic Programming Algorithm



# Extensions

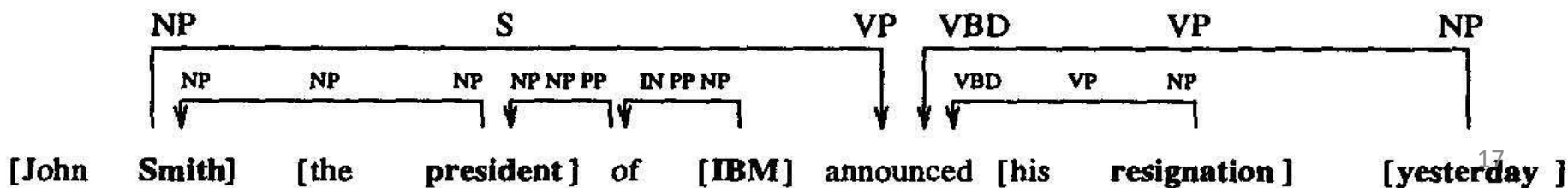
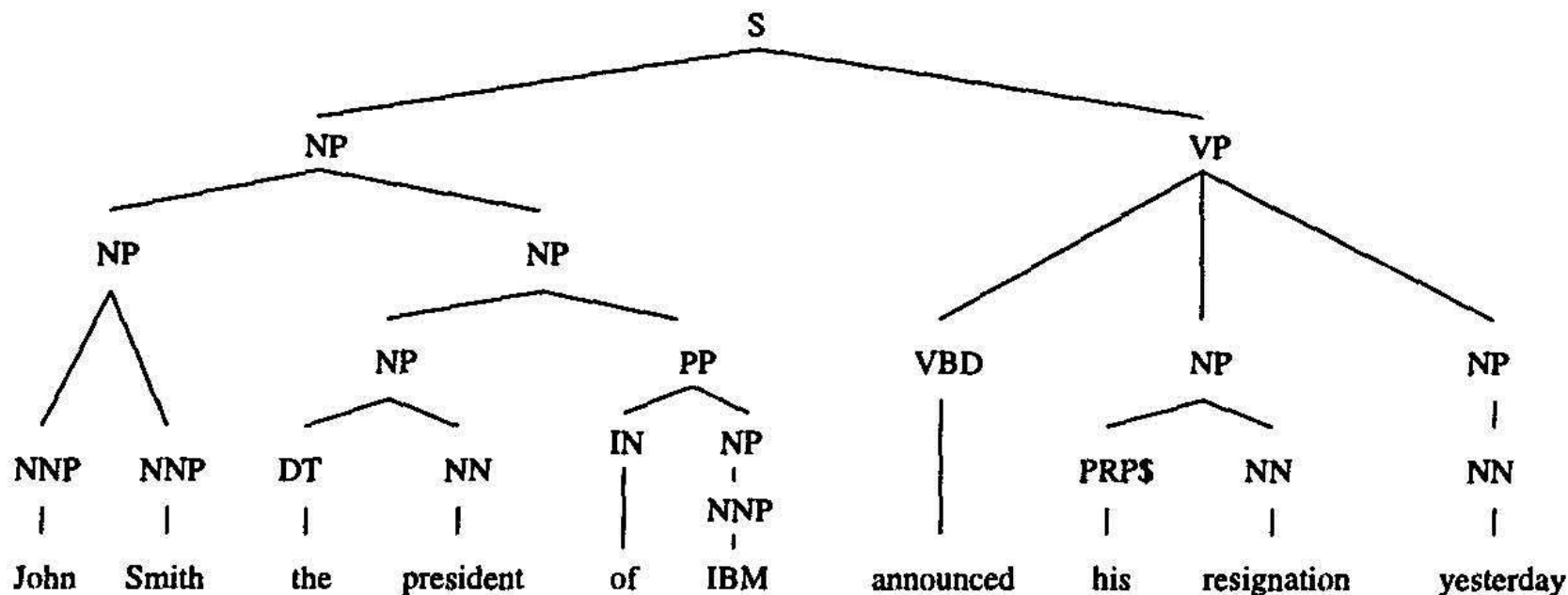
- Lexicalization
  - Use lexical information
- History
  - PCFG: too strong independence assumption
  - Use history
- Nonterminal Classification



# Lexicalization

A New Statistical Parser Based on Bigram Lexical Dependencies

[Collins 1996]



# Lexicalization

A New Statistical Parser Based on Bigram Lexical Dependencies

[Collins 1996]

$$T_{best} = \arg \max_T P(T | S) = \arg \max_T P(B | S) \times P(D | S, B)$$

Base NP Model

Dependency Model

$B = \{ [\text{John Smith}], [\text{the president}], [\text{IBM}], [\text{his resignation}], [\text{yesterday}] \}$

$D = \{$


$\begin{array}{ccc} \text{NP} & \text{S} & \text{VP} \\ \hline \downarrow & & \downarrow \end{array}$	$\begin{array}{ccc} \text{NP} & \text{NP} & \text{NP} \\ \hline \downarrow & & \downarrow \end{array}$	$\begin{array}{ccc} \text{NP} & \text{NP} & \text{PP} \\ \hline \downarrow & & \downarrow \end{array}$	$\begin{array}{ccc} \text{IN} & \text{PP} & \text{NP} \\ \hline \downarrow & & \downarrow \end{array}$	$\begin{array}{ccc} \text{VBD} & \text{VP} & \text{NP} \\ \hline \downarrow & & \downarrow \end{array}$
Smith announced	Smith president	president of	of IBM	announced resignation

$\}$

$\begin{array}{ccc} \text{VBD} & \text{VP} & \text{NP} \\ \hline \downarrow & & \downarrow \end{array}$   
 announced yesterday }

# Base NP Model

$$P(B | S) = \prod_{i=2..n} \hat{P}(G_i | w_{i-1}, t_{i-1}, w_i, t_i, c_i)$$

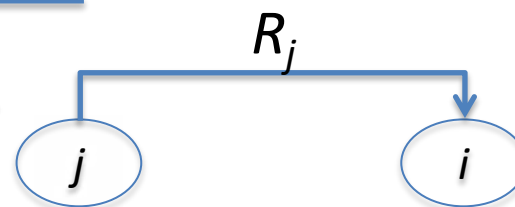


S(tart)
C(ontinue)
E(nd)
B(etween)
N(ull)

S John C Smith B the C president E of S IBM E announced ...

# Dependency Model

$$P(D | S, B) = \prod_{j=1}^m P(\underbrace{AF(j) = (i, R_j)}_{\text{Diagram}} | S, B)$$



e.g.  $AF(1) = (5, \langle \text{NP}, \text{S}, \text{VP} \rangle)$

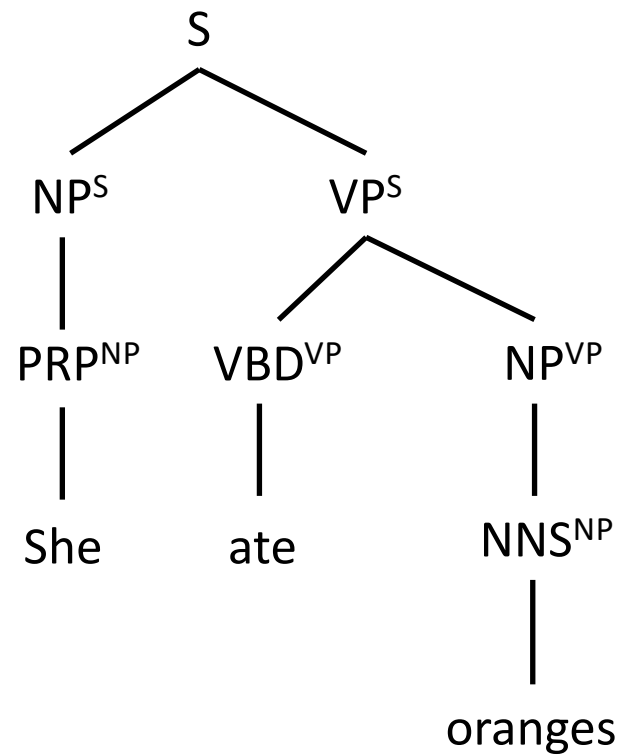
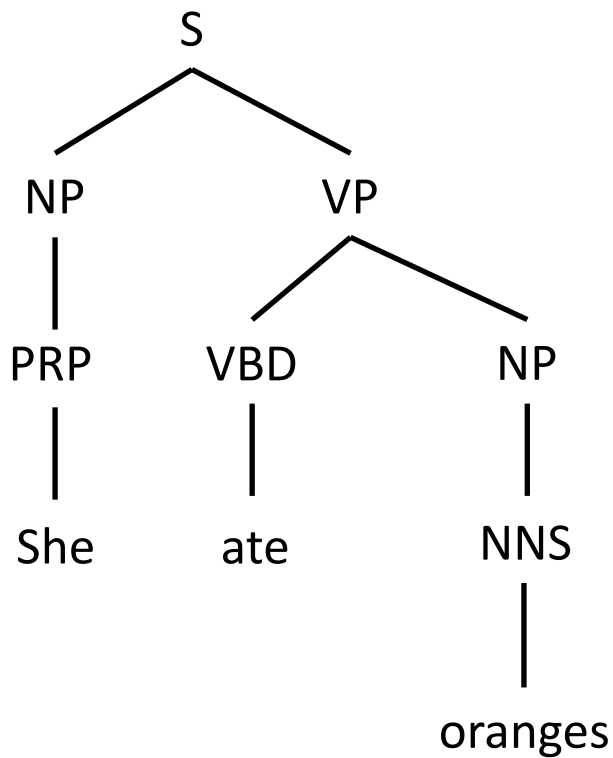
$$= \prod_{j=1}^m \frac{F(R_j | \langle w_j, t_j \rangle, \langle w_i, t_i \rangle)}{\sum_{k=1 \dots m, k \neq j, p \in P} F(p | \langle w_j, t_j \rangle, \langle w_k, t_k \rangle)}$$

$$F(R | \langle a, b \rangle, \langle c, d \rangle) = C(R, \langle a, b \rangle, \langle c, d \rangle) / C(\langle a, b \rangle, \langle c, d \rangle)$$

e.g. 
$$\frac{C(\langle \text{NP}, \text{S}, \text{VP} \rangle, \langle \text{Smith}, \text{NNP} \rangle, \langle \text{announced}, \text{VBD} \rangle)}{C(\langle \text{Smith}, \text{NNP} \rangle, \langle \text{announced}, \text{VBD} \rangle)}$$

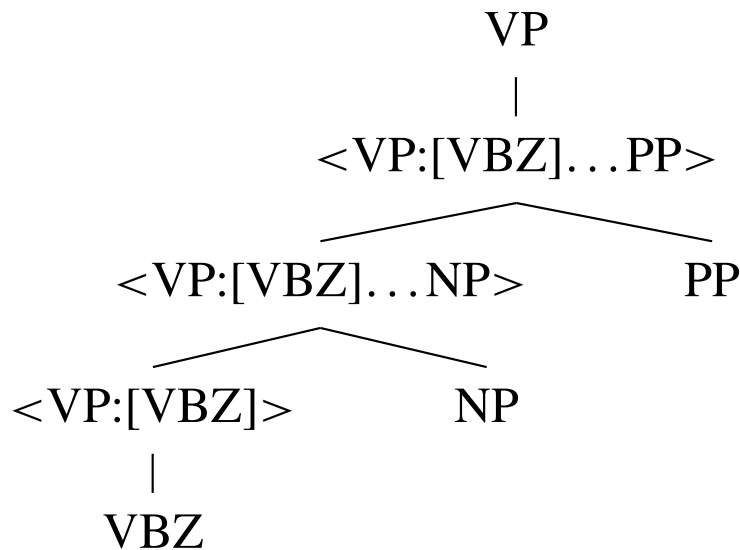
# History: Parent Annotation

[Johnson 1998]



# Vertical/Horizontal Markovization

[Klein & Manning 2003]



$v=1, h=1$  markovization of “VP  $\rightarrow$  VBZ NP PP”

F1 and grammar size:

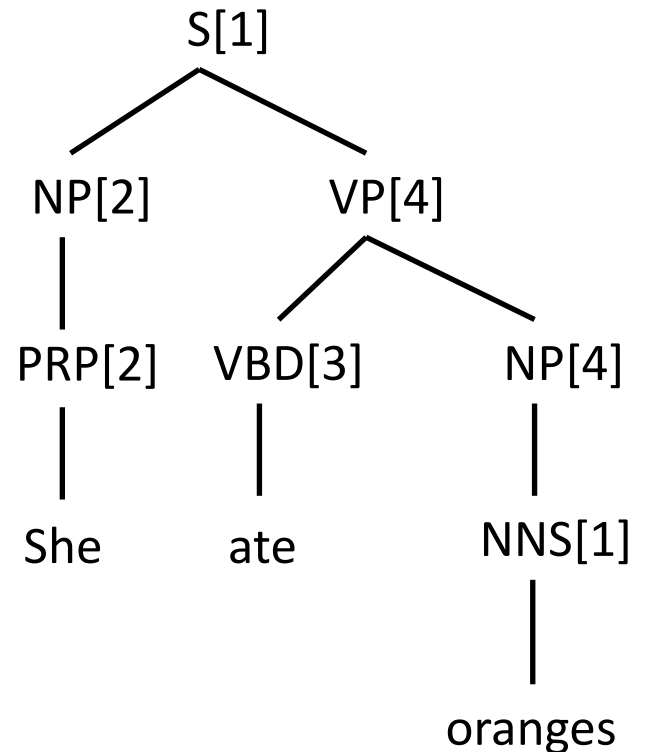
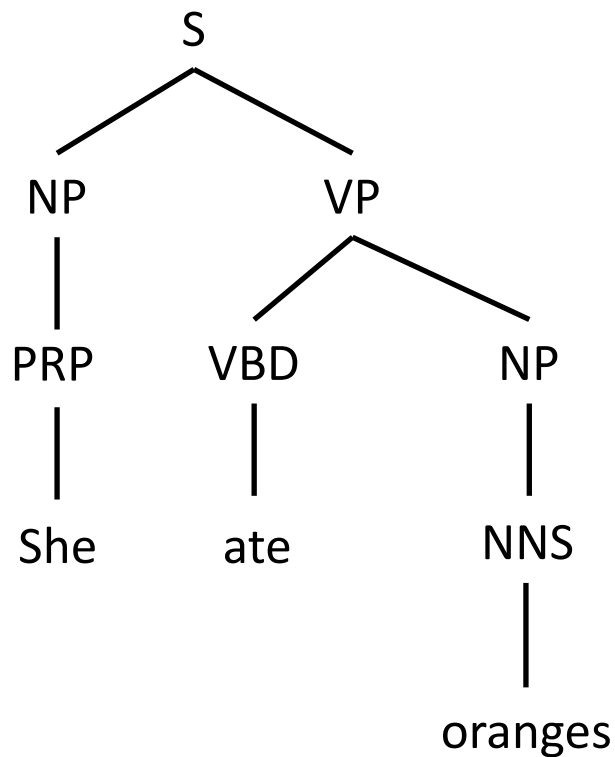
Vertical Order		Horizontal Markov Order				
		$h = 0$	$h = 1$	$h \leq 2$	$h = 2$	$h = \infty$
$v = 1$	No annotation	71.27 (854)	72.5 (3119)	73.46 (3863)	72.96 (6207)	72.62 (9657)
$v \leq 2$	Sel. Parents	74.75 (2285)	77.42 (6564)	77.77 (7619)	77.50 (11398)	76.91 (14247)
$v = 2$	All Parents	74.68 (2984)	77.42 (7312)	77.81 (8367)	77.50 (12132)	76.81 (14666)
$v \leq 3$	Sel. GParents	76.50 (4943)	78.59 (12374)	79.07 (13627)	78.97 (19545)	78.54 (20123)
$v = 3$	All GParents	76.74 (7797)	79.18 (15740)	79.74 (16994)	79.07 (22886)	78.72 (22002)

PCFG

Parent Annotation

# Nonterminal Classification

[Matsuzaki+ 2005] [Petrov+ 2006]



VBZ			
VBZ-0	gives	sells	takes
VBZ-1	comes	goes	works
VBZ-2	includes	owns	is
VBZ-3	puts	provides	takes
VBZ-4	says	adds	Says
VBZ-5	believes	means	thinks
VBZ-6	expects	makes	calls
VBZ-7	plans	expects	wants
VBZ-8	is	's	gets
VBZ-9	's	is	remains
VBZ-10	has	's	is
VBZ-11	does	Is	Does

DT			
DT-0	the	The	a
DT-1	A	An	Another
DT-2	The	No	This
DT-3	The	Some	These
DT-4	all	those	some
DT-5	some	these	both
DT-6	That	This	each
DT-7	this	that	each
DT-8	the	The	a
DT-9	no	any	some
DT-10	an	a	the
DT-11	a	this	the

IN			
IN-0	In	With	After
IN-1	In	For	At
IN-2	in	for	on
IN-3	of	for	on
IN-4	from	on	with
IN-5	at	for	by
IN-6	by	in	with
IN-7	for	with	on
IN-8	If	While	As
IN-9	because	if	while
IN-10	whether	if	That
IN-11	that	like	whether
IN-12	about	over	between
IN-13	as	de	Up
IN-14	than	ago	until
IN-15	out	up	down

NNP			
NNP-0	Jr.	Goldman	INC.
NNP-1	Bush	Noriega	Peters
NNP-2	J.	E.	L.
NNP-3	York	Francisco	Street
NNP-4	Inc	Exchange	Co
NNP-5	Inc.	Corp.	Co.
NNP-6	Stock	Exchange	York
NNP-7	Corp.	Inc.	Group
NNP-8	Congress	Japan	IBM
NNP-9	Friday	September	August
NNP-10	Shearson	D.	Ford
NNP-11	U.S.	Treasury	Senate
NNP-12	John	Robert	James
NNP-13	Mr.	Ms.	President
NNP-14	Oct.	Nov.	Sept.
NNP-15	New	San	Wall

CD			
CD-0	1	50	100
CD-1	8.50	15	1.2
CD-2	8	10	20
CD-3	1	30	31
CD-4	1989	1990	1988
CD-5	1988	1987	1990
CD-6	two	three	five
CD-7	one	One	Three
CD-8	12	34	14
CD-9	78	58	34
CD-10	one	two	three
CD-11	million	billion	trillion

RB			
RB-0	recently	previously	still
RB-1	here	back	now
RB-2	very	highly	relatively
RB-3	so	too	as
RB-4	also	now	still
RB-5	however	Now	However
RB-6	much	far	enough
RB-7	even	well	then
RB-8	as	about	nearly
RB-9	only	just	almost
RB-10	ago	earlier	later
RB-11	rather	instead	because
RB-12	back	close	ahead
RB-13	up	down	off
RB-14	not	Not	maybe
RB-15	n't	not	also

JJS			
JJS-0	largest	latest	biggest
JJS-1	least	best	worst
JJS-2	most	Most	least

PRP			
PRP-0	It	He	I
PRP-1	it	he	they
PRP-2	it	them	him

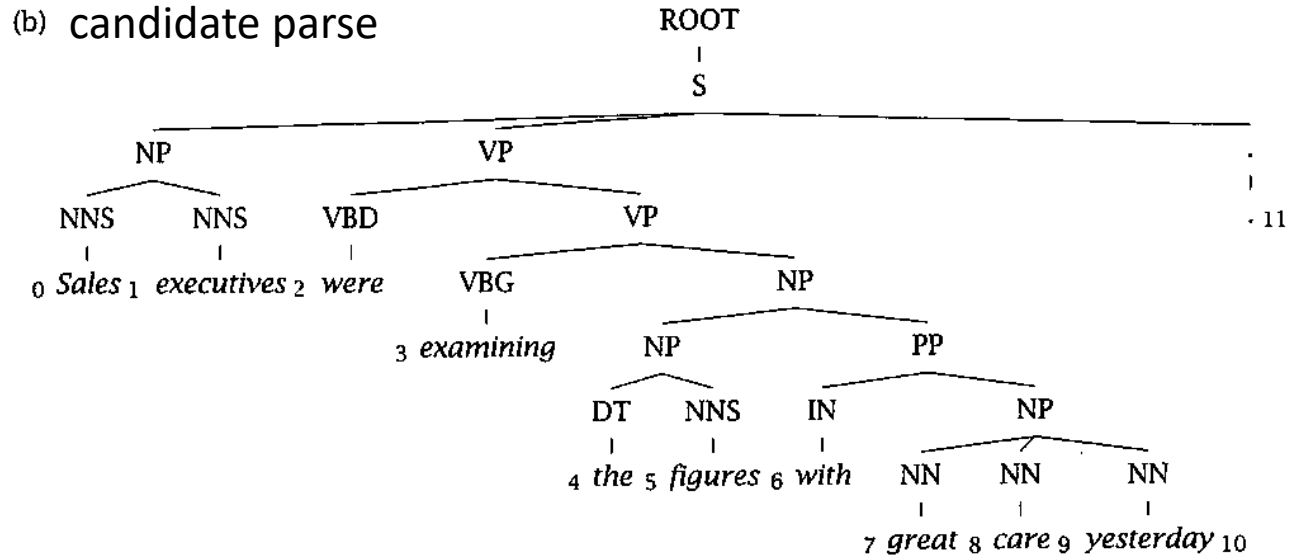
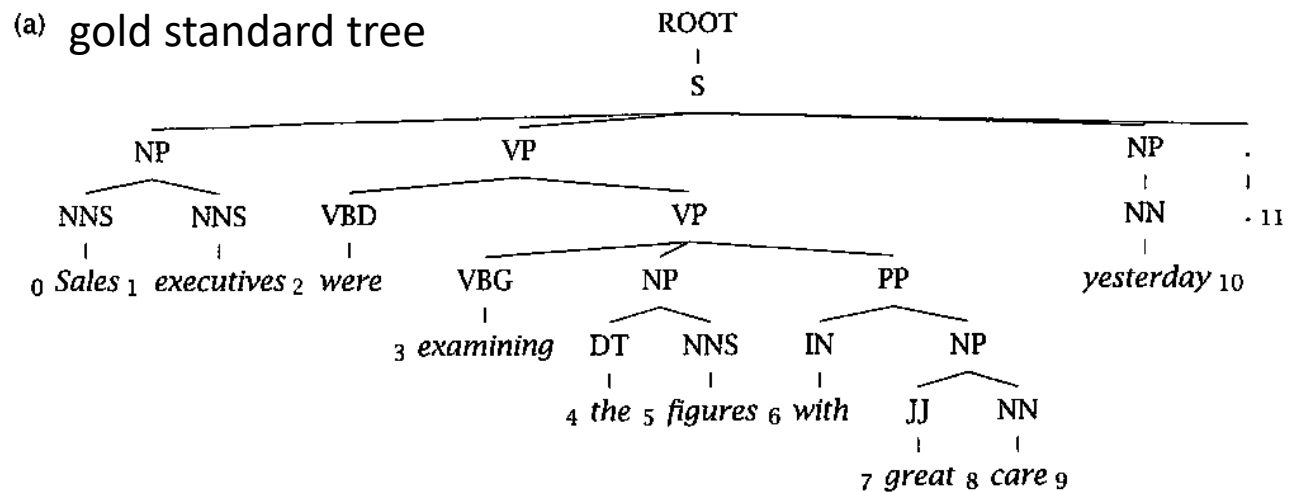
RBR			
RBR-0	further	lower	higher
RBR-1	more	less	More
RBR-2	earlier	Earlier	later

[Petrov+ 2006]



# Evaluation

- Precision
- Recall
- Labeled Precision
- Labeled Recall



- (c) Brackets in gold standard tree (a.):  
 S-(0:11), NP-(0:2), VP-(2:9), VP-(3:9), NP-(4:6), PP-(6:9), NP-(7,9), \*NP-(9:10)
- (d) Brackets in candidate parse (b.):  
 S-(0:11), NP-(0:2), VP-(2:10), VP-(3:10), NP-(4:10), NP-(4:6), PP-(6:10), NP-(7,10)
- (e) Precision:  $3/8 = 37.5\%$  Crossing Brackets: 0  
 Recall:  $3/8 = 37.5\%$  Crossing Accuracy: 100%  
 Labeled Precision:  $3/8 = 37.5\%$  Tagging Accuracy:  $10/11 = 90.9\%$   
 Labeled Recall:  $3/8 = 37.5\%$

# Performance on English (supervised)

	LP	LR	F1		LP	LR	F1
[Magerman 1995]	84.0	84.3	84.2	[Socher+ 2013]			90.4
[Charniak 1997]	86.7	86.6	86.7	[Watanabe+ 2015]			90.7
[Collins 1997]	87.5	88.1	87.8	[Mi&Huang 2015]	90.7	90.9	90.8
[Charniak 2000]	89.6	89.5	89.6	[Cross&Huang 2016]	90.5	92.1	91.3
[Petrov&Klein 2007]	90.2	89.9	90.1	[Dyer+ 2016]			91.7
[Carreras+ 2008]	90.7	91.4	91.1	[Stern+ 2017]	90.6	93.0	91.8
[Shindo+ 2012]			91.1	[Stern+ 2017]	92.6	92.6	92.6
[Zhu+ 2013]	90.2	90.7	90.4	[Gaddy+ 2018]			92.1
				[Kitaev&Klein 2018]	93.2	93.9	93.6

# Performance on English (semi-supervised, reranking, etc.)

	F1	Note
[Charniak&Johnson 2005]	91.0	reranking
[McClosky+ 2006]	92.1	self-training
[Shindo+ 2012]	92.4	ensemble
[Vinyals+ 2015]	92.8	tri-training
[Dyer+ 2016]	92.4	reranking
[Choe&Charniak 2016]	93.8	tri-training
[Kuncoro+ 2017]	93.6	reranking
[Liu&Zhang 2017]	94.2	reranking/tri-training
[Fried+ 2017]	94.7	ensemble/reranking

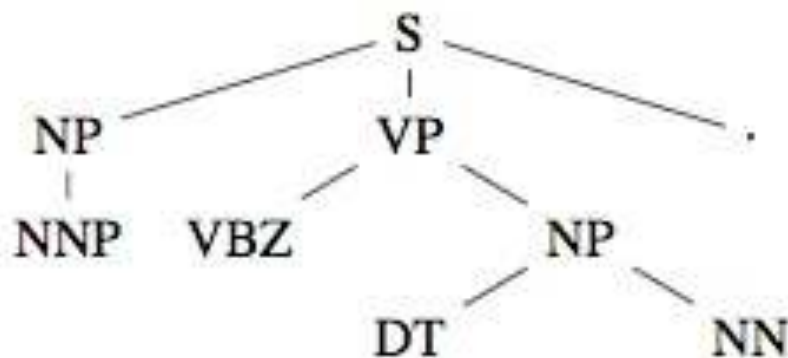
# Performance on English (pre-trained)

	F1	Note
[Kitaev&Klein 2018]	95.1	pre-train (ELMo)
[Kitaev+ 2019]	95.6	pre-train (BERT)
[Zhou&Zhao 2019]	96.3	pre-train (XLNet)
[Yang&Deng 2020]	96.3	pre-train (XLNet)
[Mrini+ 2020]	96.4	pre-train (XLNet)

# Serialization [Vinyals+ 2015]

John has a dog .

→



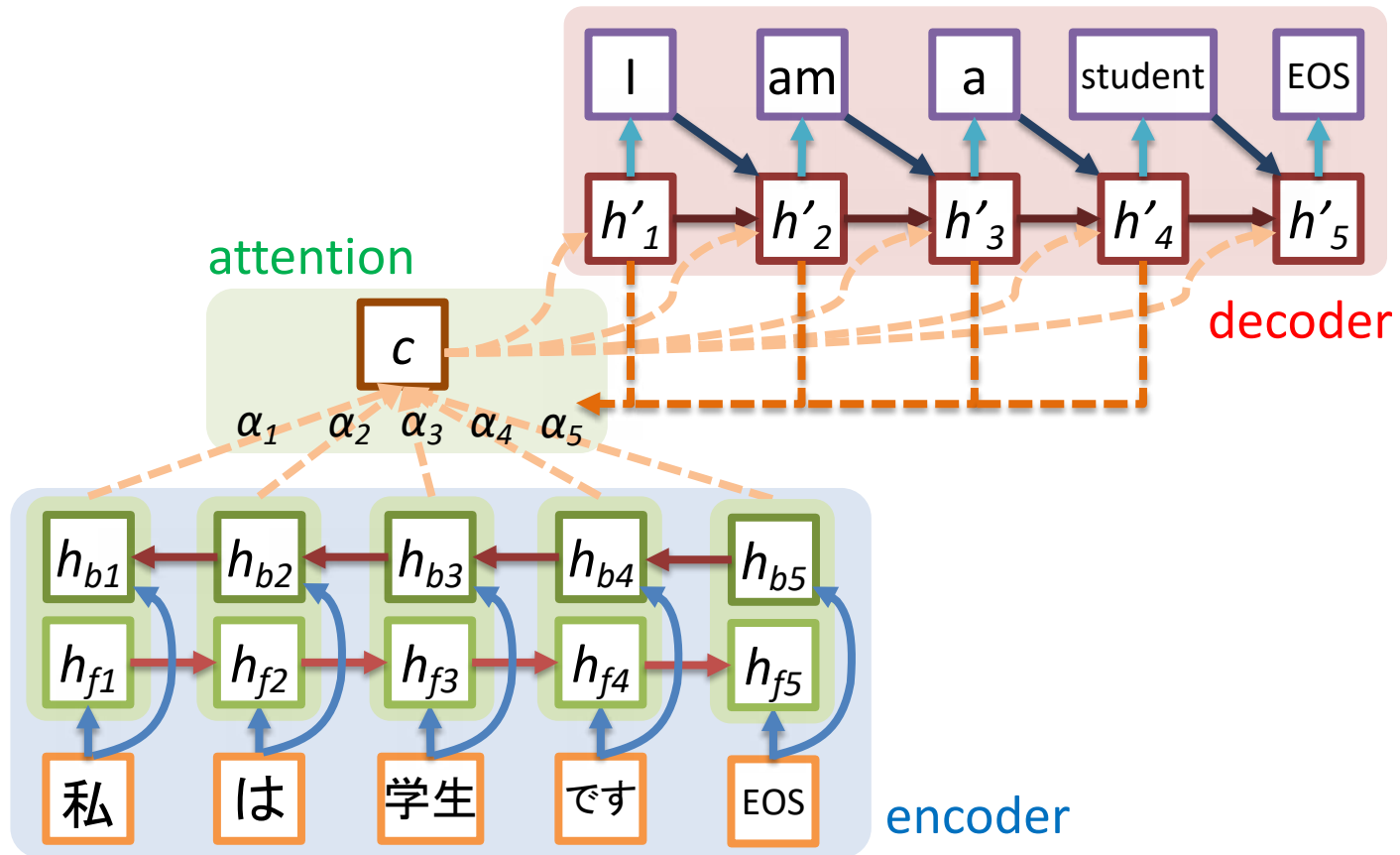
John has a dog .

→

(S (NP NNP )<sub>NP</sub> (VP VBZ (NP DT NN )<sub>NP</sub> )<sub>VP</sub> . )<sub>S</sub>

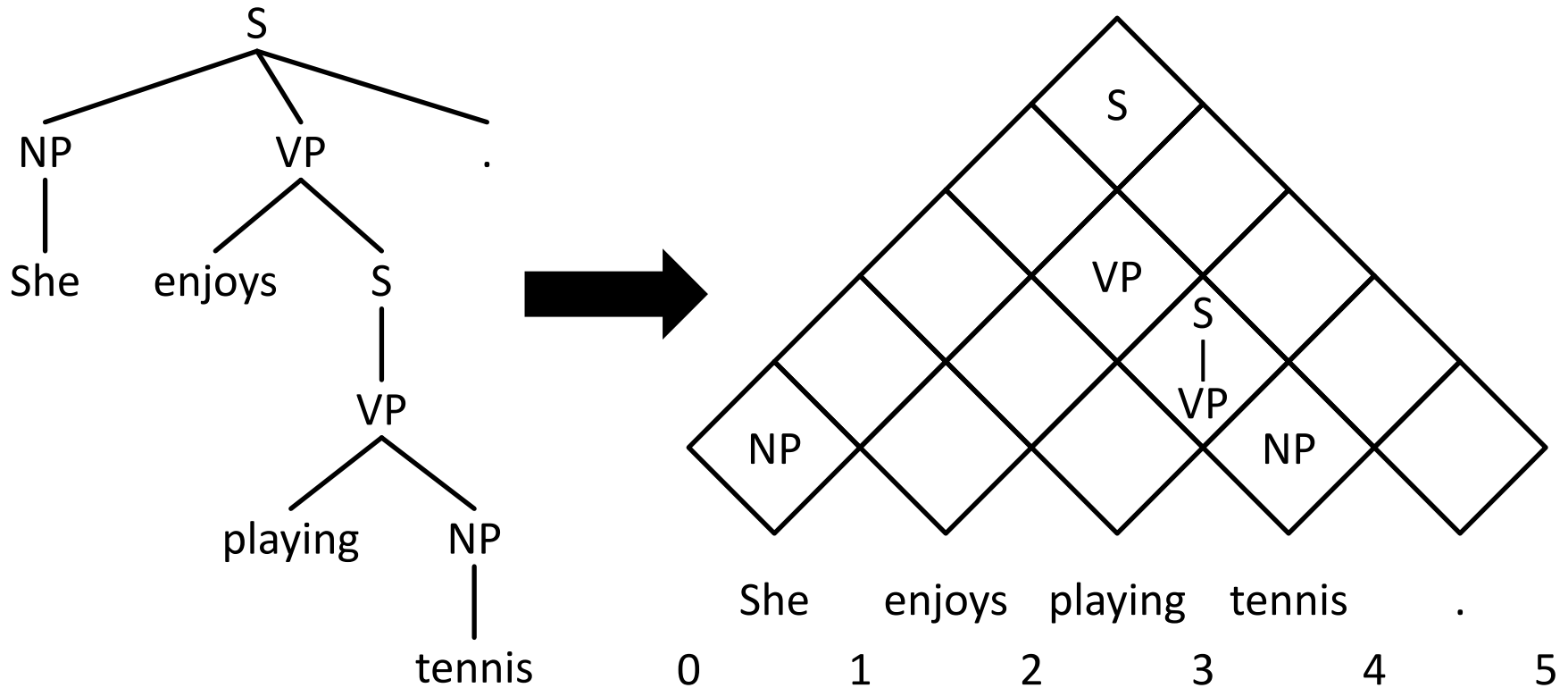
# Attention-based Neural Machine Translation

[Bahdanau+ 2014]

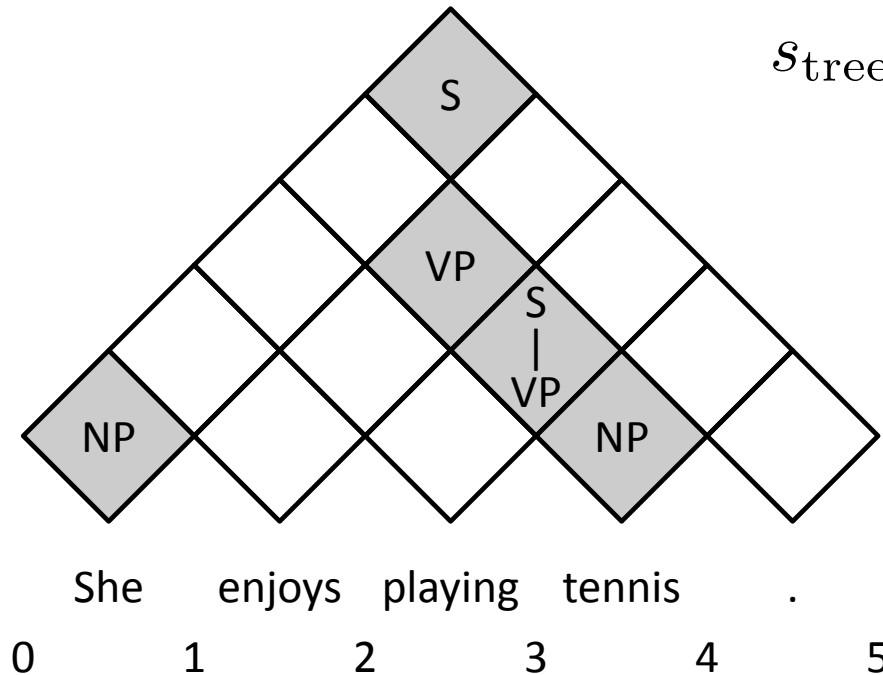


# Neural Span Classification

[Stern+ 2017]



# Neural Span Classification

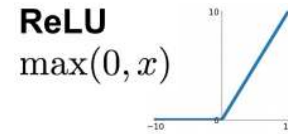


$$\begin{aligned}
 s_{\text{tree}}(T) &= \sum_{(\ell, (i, j)) \in T} s(i, j, \ell) \\
 &= s(0, 5, S) \\
 &\quad + s(0, 1, \text{NP}) \\
 &\quad + s(1, 4, \text{VP}) \\
 &\quad + s(2, 4, \text{S-VP}) \\
 &\quad + s(3, 4, \text{NP})
 \end{aligned}$$

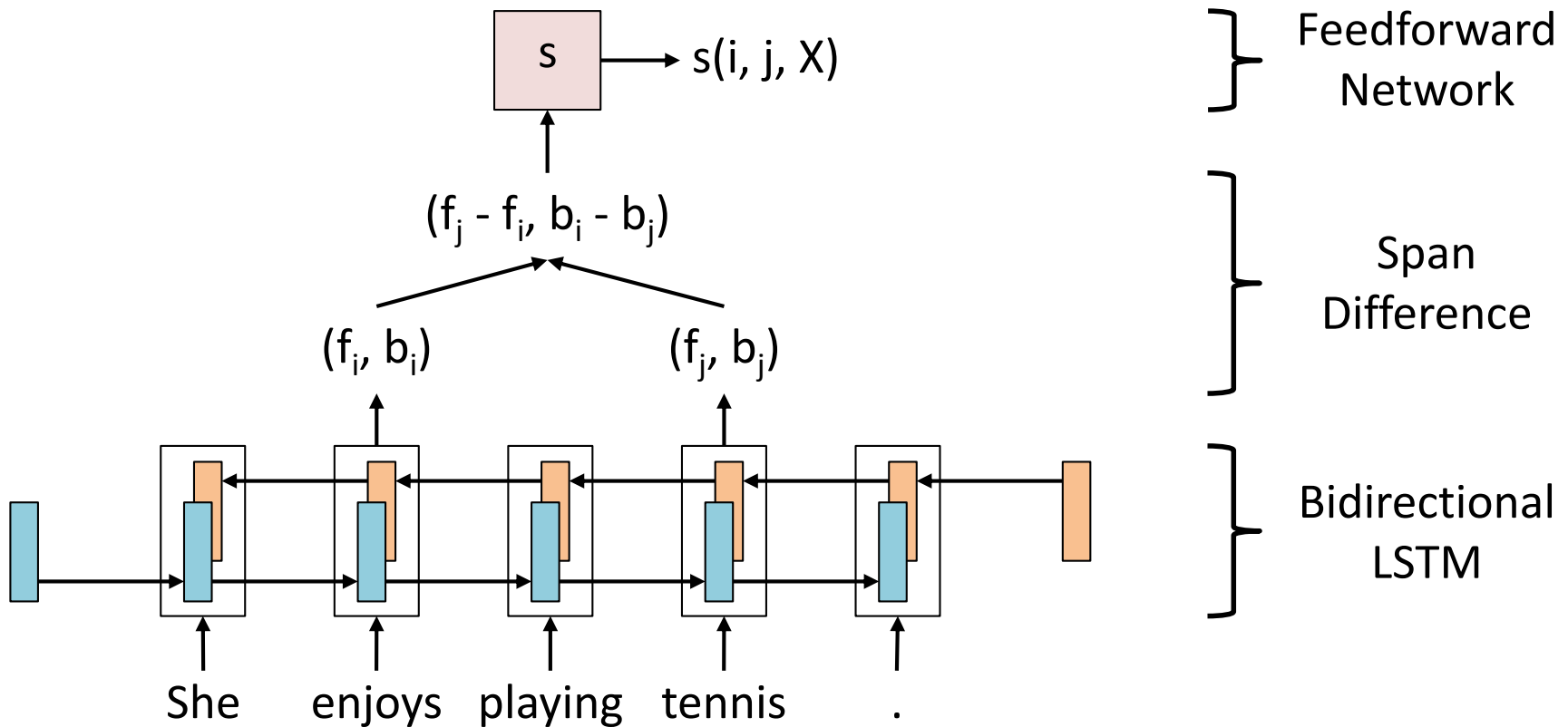
$$s_{\text{best}}(i, j) = \overbrace{\max_{\ell} [s(i, j, \ell)]}^{\text{Pick best label}} + \underbrace{\max_k [s_{\text{best}}(i, k) + s_{\text{best}}(k, j)]}_{\text{Pick best split point}}$$



# Neural Span Classification



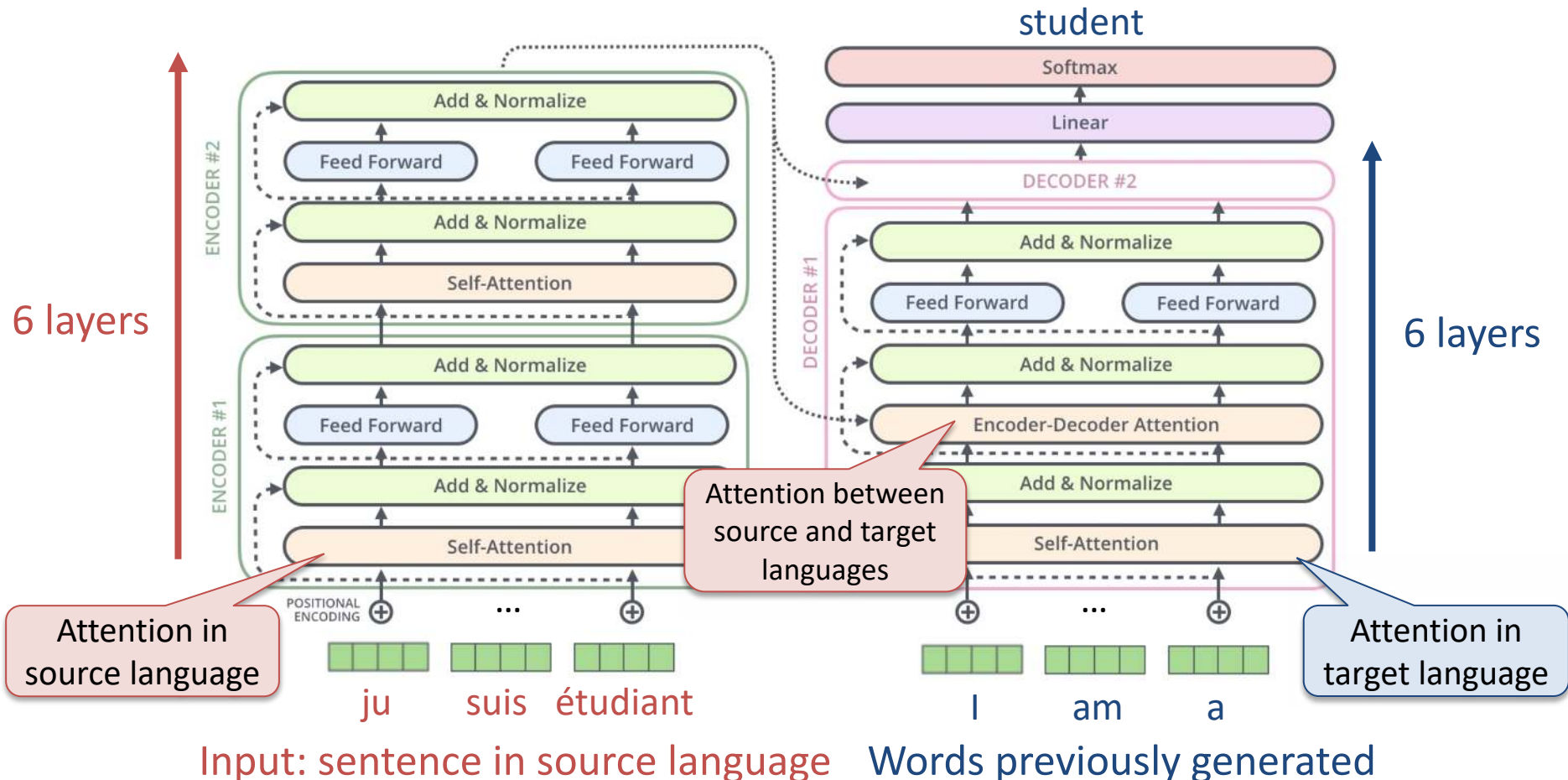
$$\text{FeedForward}(x) = W_2 \text{relu}(W_1 x + b_1) + b_2$$



# Transformer: “Attention is All You Need”

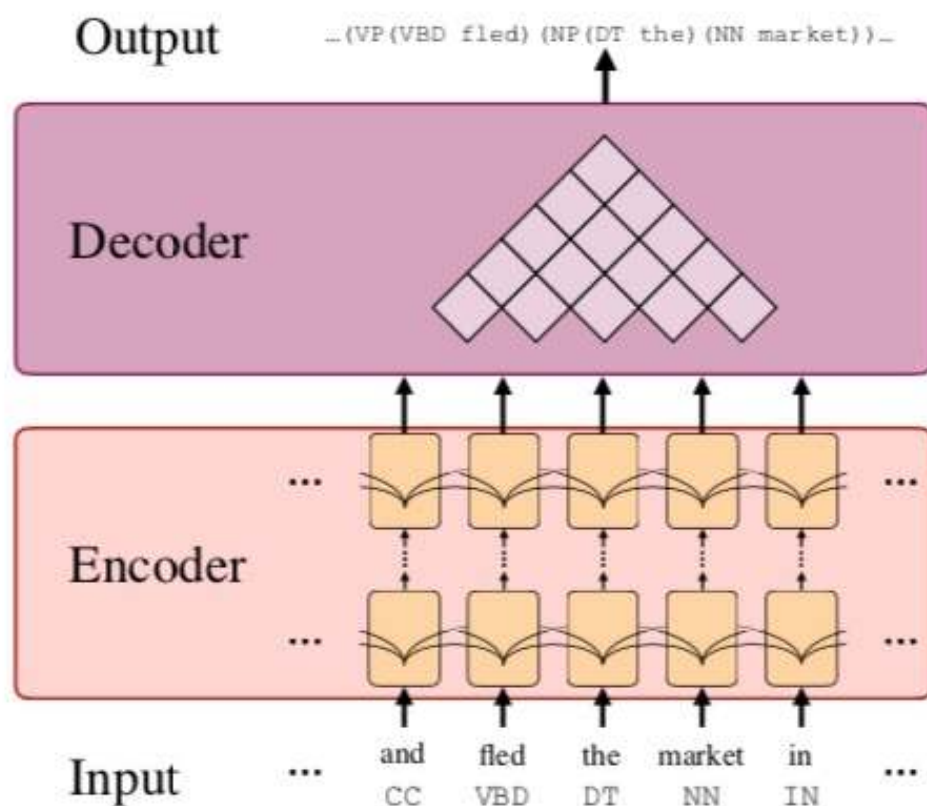
[Vaswani+ 2017]

Output: next word in target language

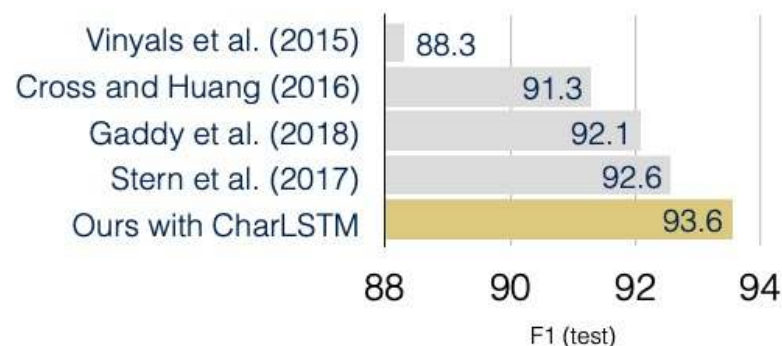


# Using Self Attention

[Kitaev&Klein 2018]



## Single Model, WSJ Only

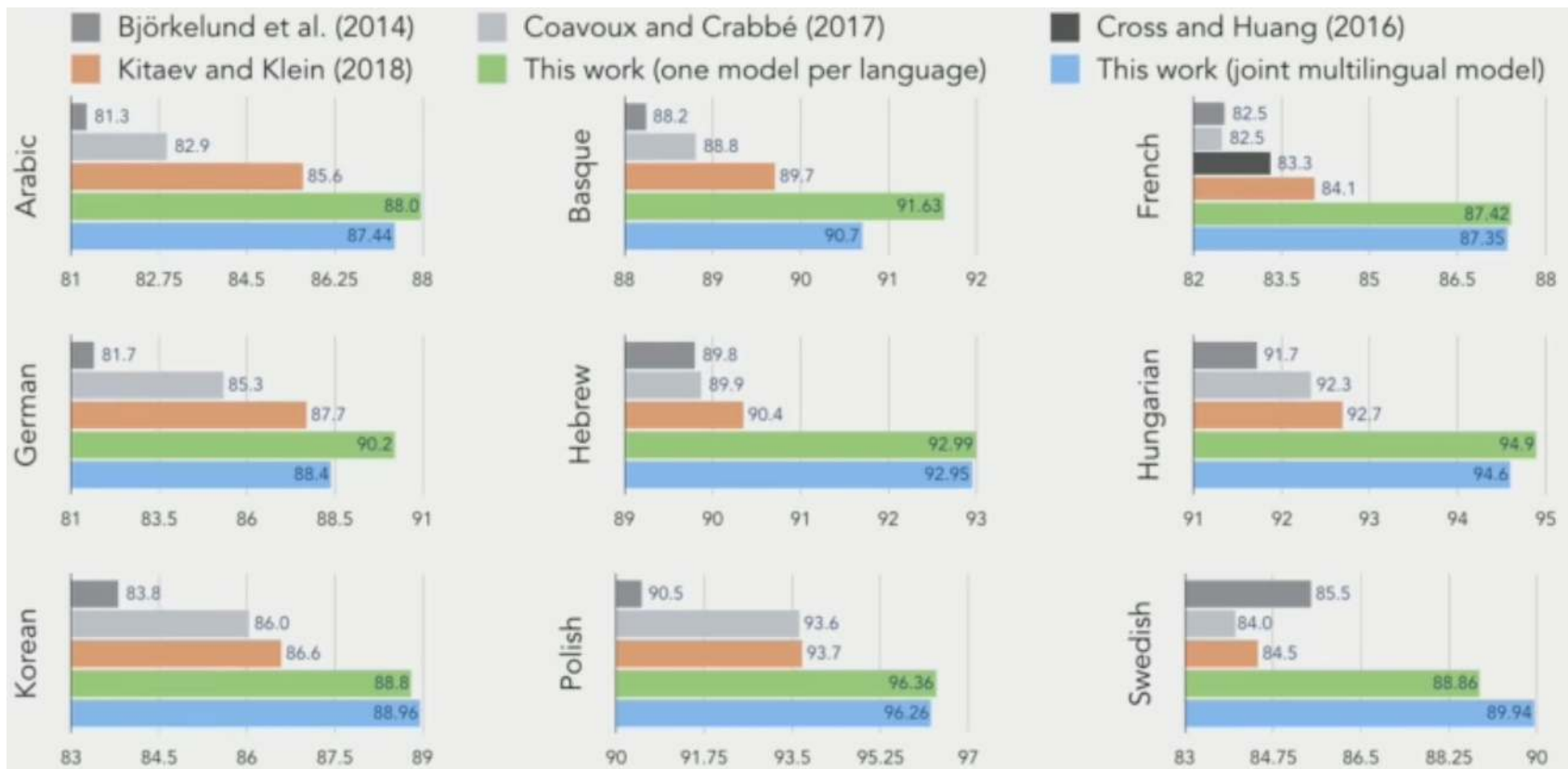
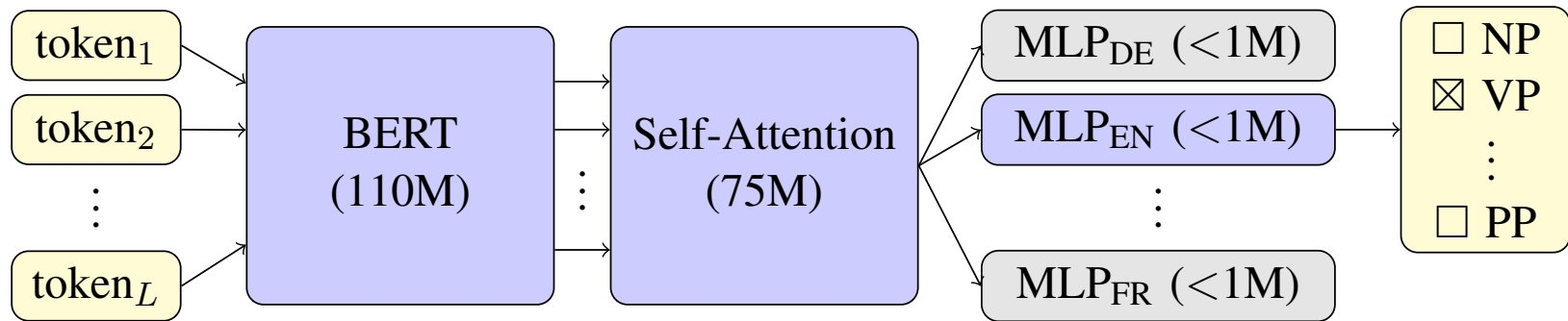


## Multi-Model / External



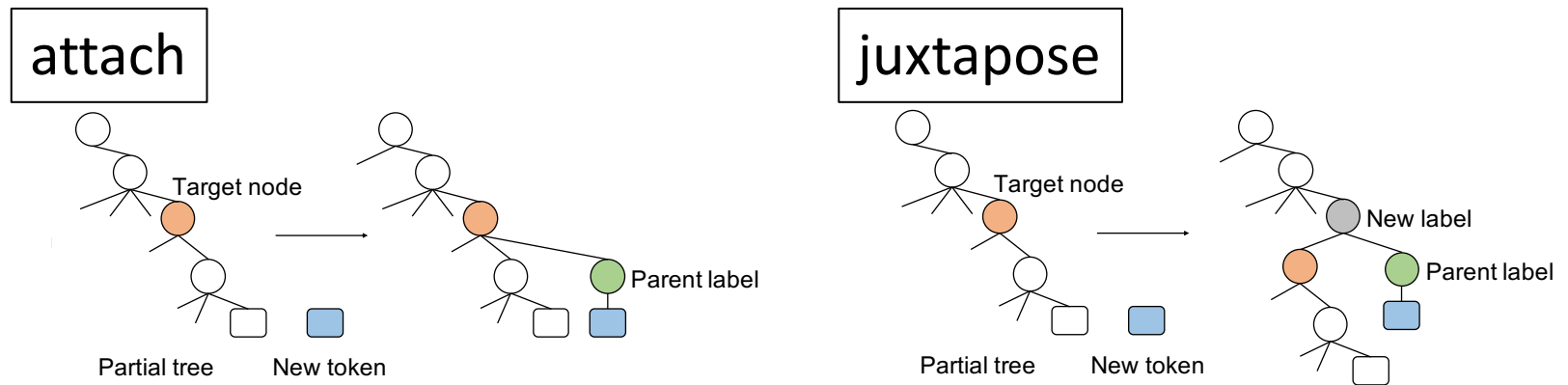
# Joint Multilingual Span Classification

[Kitaev+ 2019]

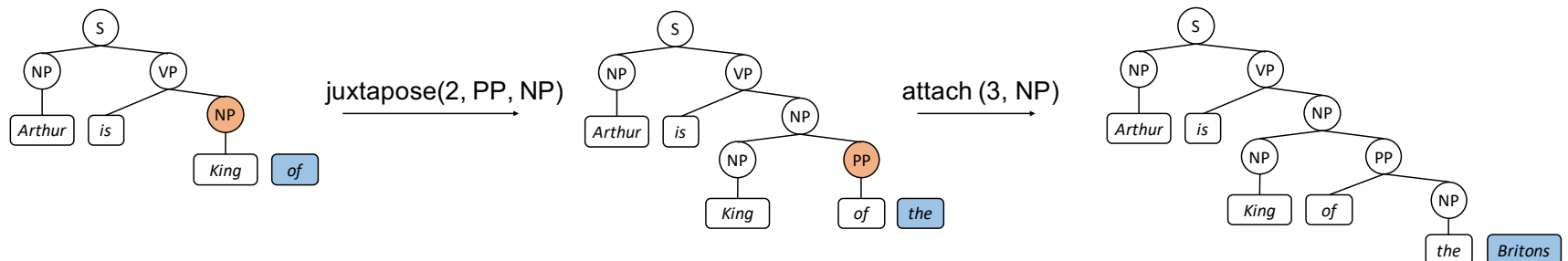


# Transition-based Constituency Parsing

- An attach-juxtapose parser [Yang&Deng 2020]

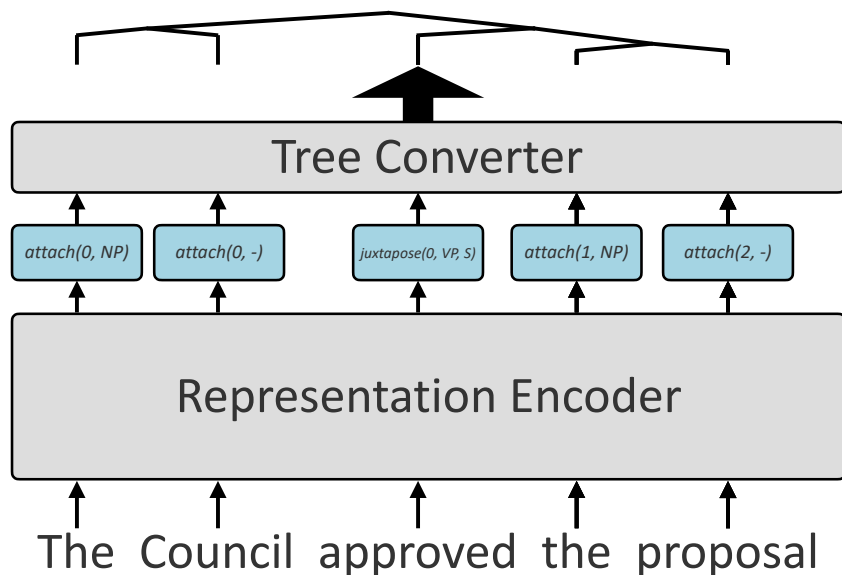


Example:

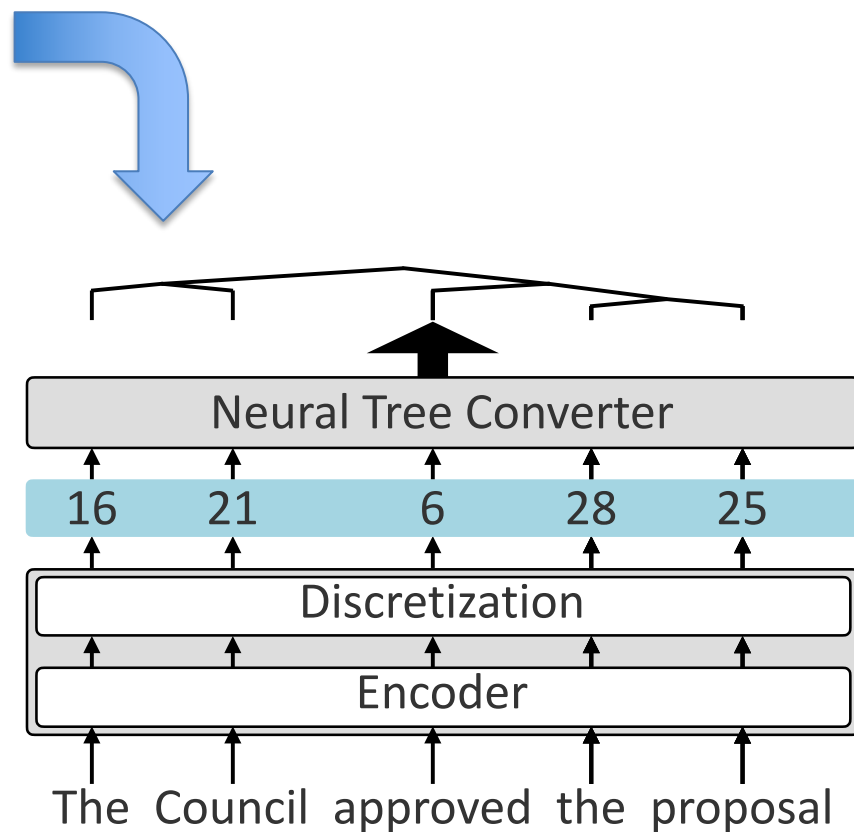


# Learning Representations for Parsing

[Kitaev+ 2022]



Representation	Encoder Type		
	Bi ( $\leftrightarrow$ )		Uni ( $\rightarrow$ )
	BERT	GPT-2	GPT-2
Span Classification (Kitaev et al., 2019)	95.59	<b>95.10<sup>†</sup></b>	93.95 <sup>†</sup>
Attach-Juxtapose (Yang and Deng, 2020)	<b>95.79</b>	94.53 <sup>†</sup>	87.66 <sup>†</sup>
Learned (This work)	95.55	–	<b>94.97</b>



# Summary

- Review of CFG and CKY parsing
- Probabilistic CFG
- Treebanks
- Extensions
  - Lexicalization
  - History
  - Nonterminal Classification
- Evaluation criteria and SOTA of English parsing