

PARSING

PCFG, Recursive Neural Network

Document classification, sentiment analysis, QA, conversation agents, summarization, translation

Discourse

Semantics

Syntax: Constituents

Syntax: Part of Speech

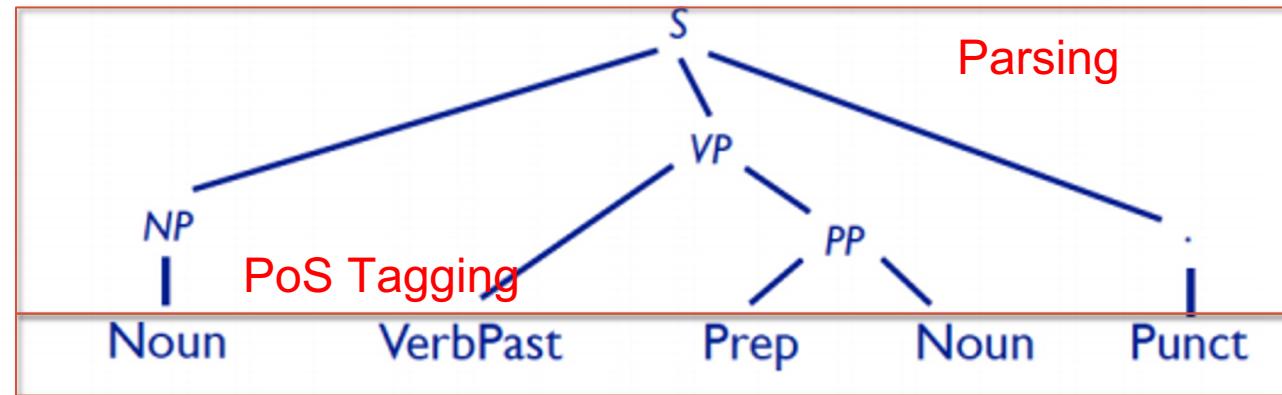
Words

Morphology

Characters

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)
SpeakerContext(s)
TemporalBefore(e, s)

Word embeddings



Alice talked to Bob.

Language modeling

talk -ed

[VerbPast]

Alice talked to Bob.

Tokenization

Discourse

Semantics

Syntax: Constituents

Syntax: Part of Speech

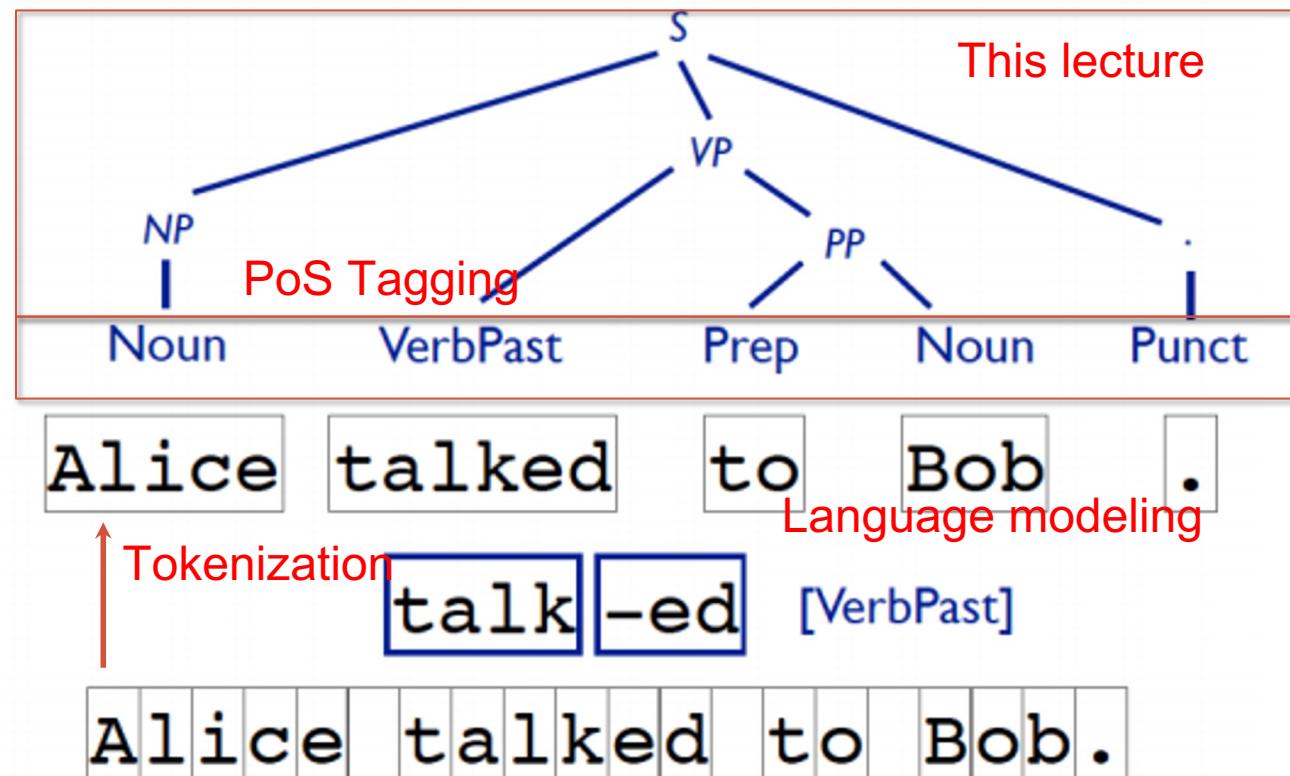
Words

Morphology

Characters

Word embeddings

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)
SpeakerContext(s)
TemporalBefore(e, s)



Use Syntax to create sentence level meanings!

Discourse

Semantics

Syntax: Constituents

Syntax: Part of Speech

Words

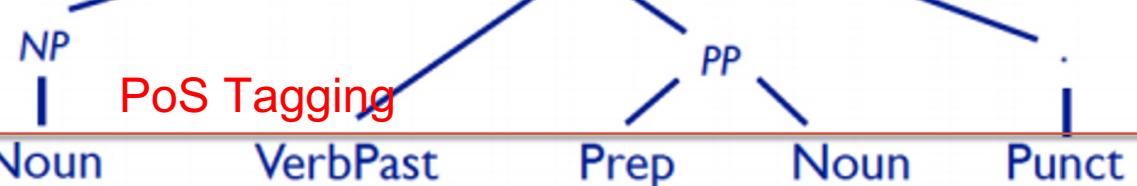
Morphology

Characters

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)
SpeakerContext(s)
TemporalBefore(e, s)

Word embeddings

Parsing



Alice talked to Bob.

Language modeling

Tokenization

talk -ed

[VerbPast]

Alice talked to Bob.

Semantic embeddings of several words

- Compositionality
- We know how to create a dense vector representation for a word
 - What about larger linguistic units? (e.g. phrase, sentence)
- We can combine smaller units into a larger unit

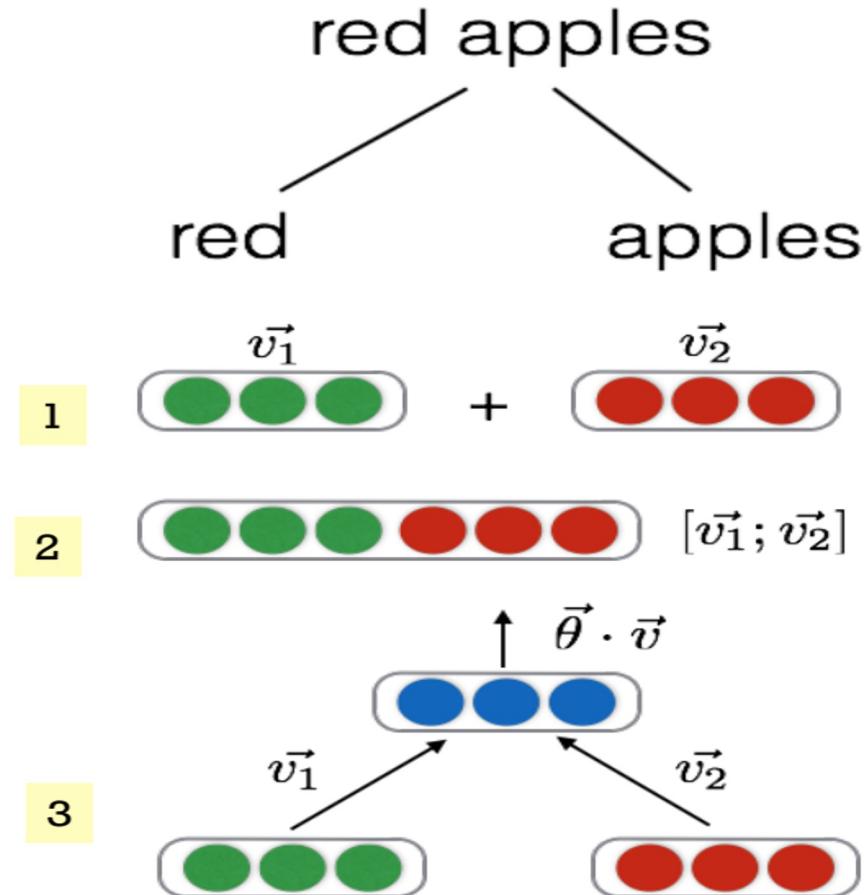
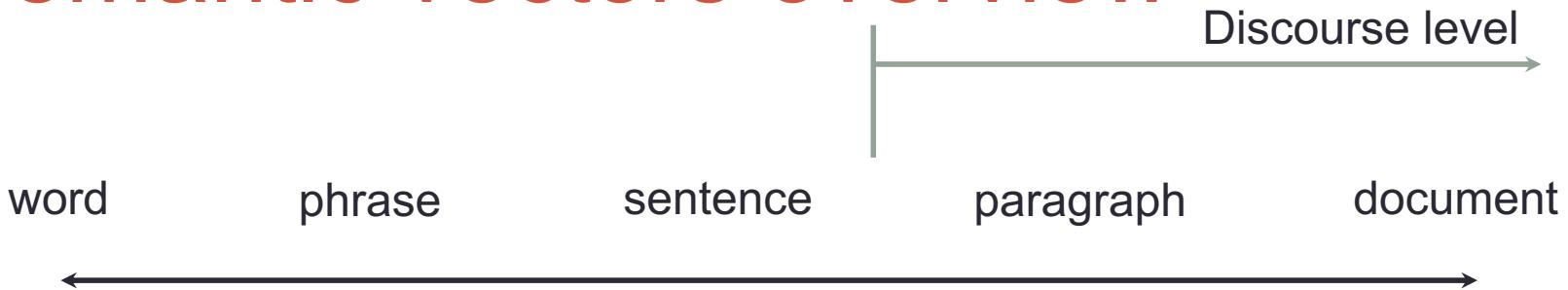


Image ref: Prof. Regina Barzilay , NLP@MIT

Semantic vectors overview



- Word vectors
 - Co-occurrence
 - PPMI
 - TFIDF
 - Word2vec
 - CBoW
 - Skip-gram

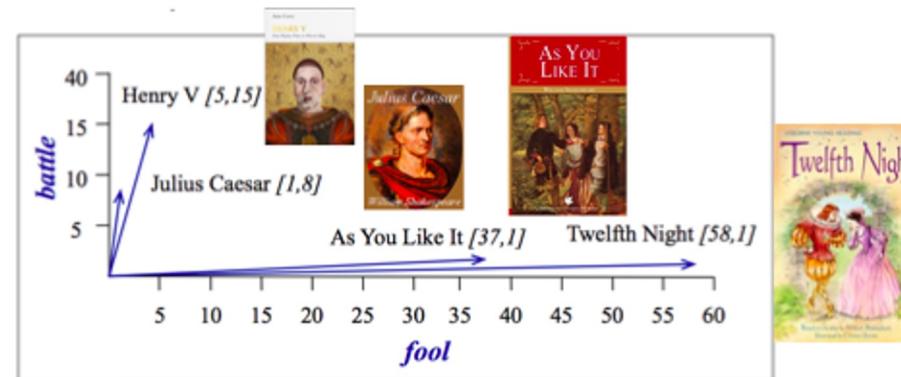
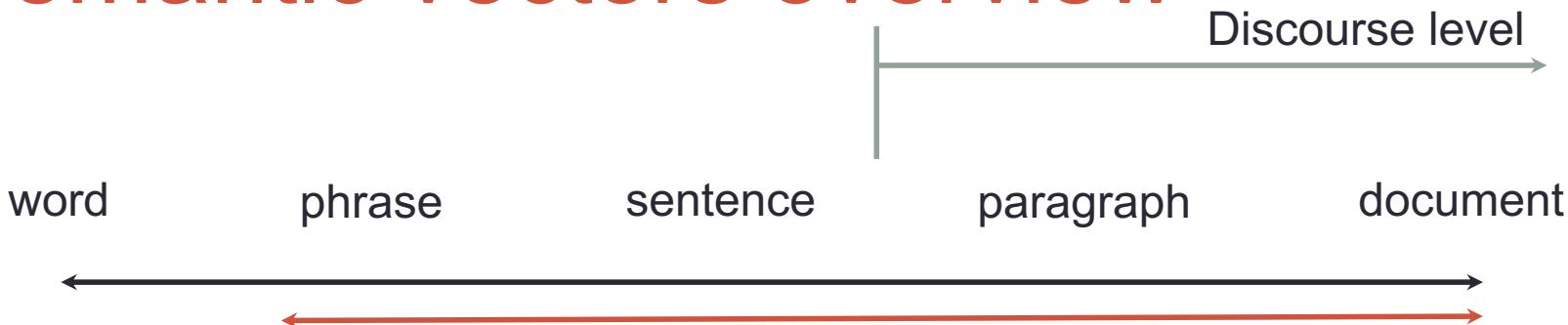


Figure 15.3 A spatial visualization of the document vectors for the four Shakespeare play documents, showing just two of the dimensions, corresponding to the words *battle* and *fool*. The comedies have high values for the *fool* dimension and low values for the *battle* dimension.

- Doc vectors
 - Term-document
 - Bag of words model
- Recurrent networks

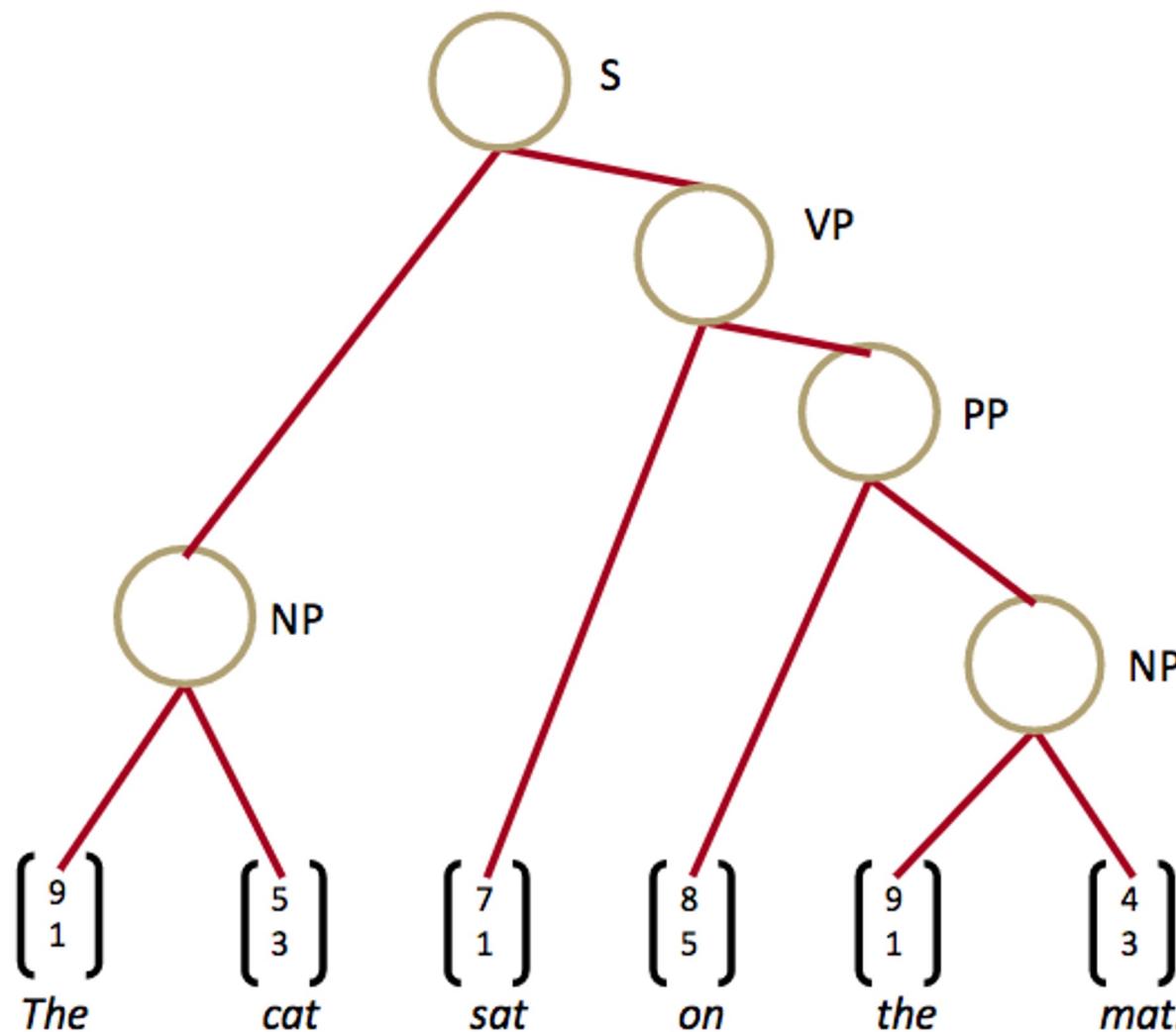
Semantic vectors overview



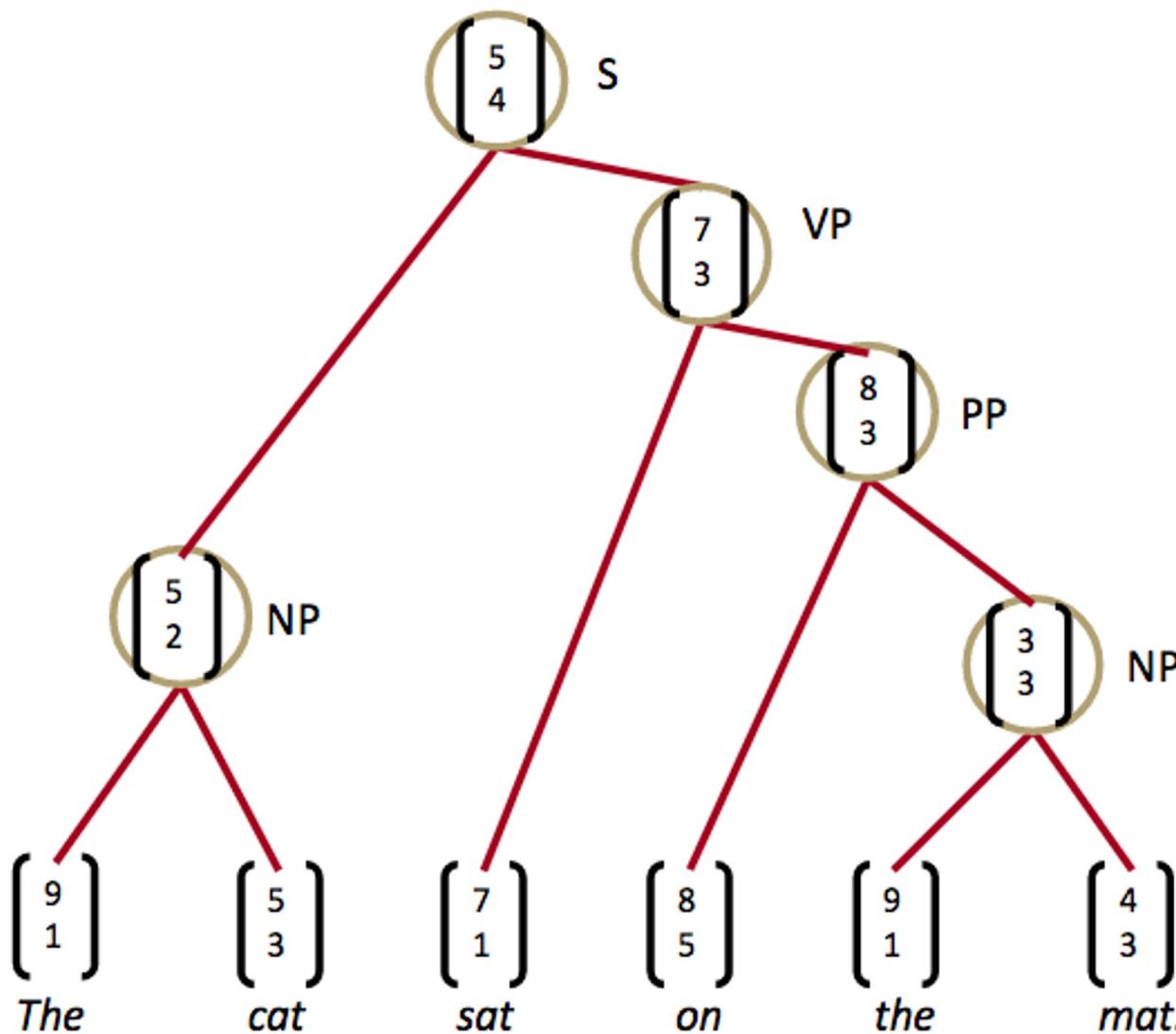
How do we represent things in this level? Without ignoring word order (bag of words)

- Word vectors
 - Co-occurrence
 - PPMI
 - TFIDF
 - Word2vec
 - CBoW
 - Skip-gram
- Doc vectors
 - Term-document
 - Bag of words model
 - Doc2Vec
 - LDA2Vec
 - Recurrent networks

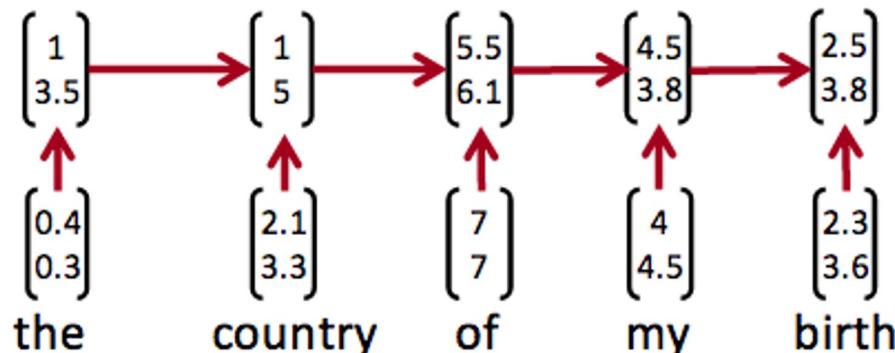
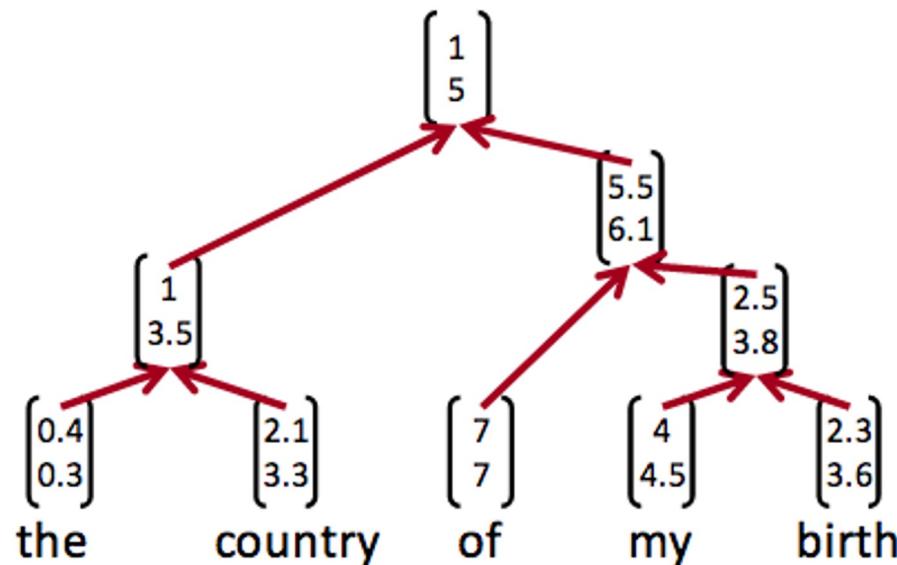
Parsing and representation big picture



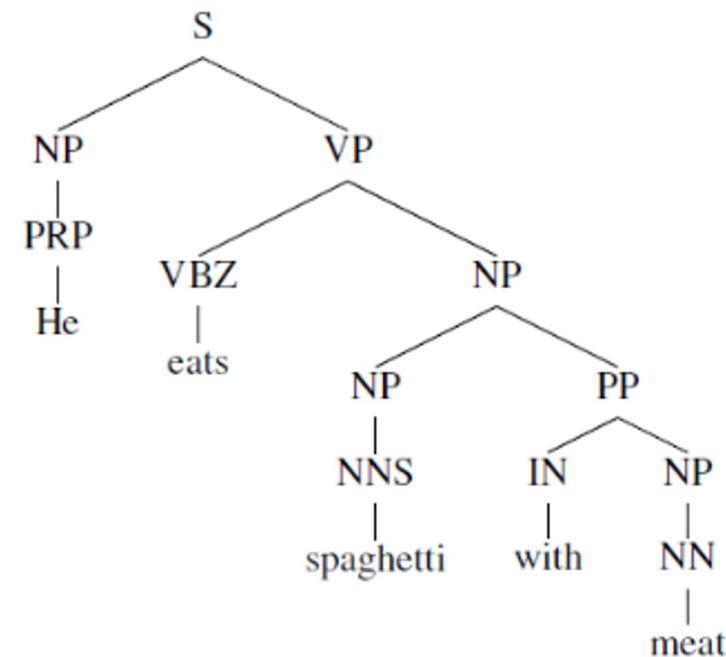
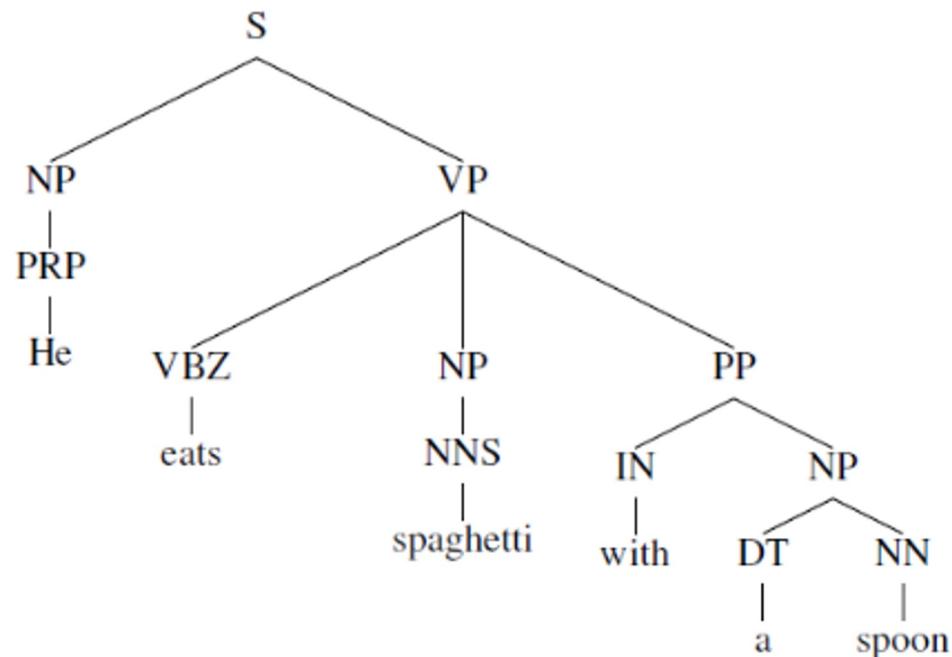
Parsing and representation big picture



Recursive vs Recurrent representation



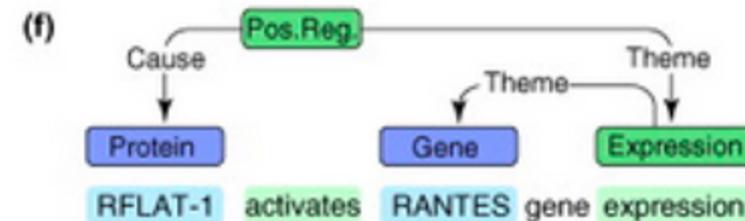
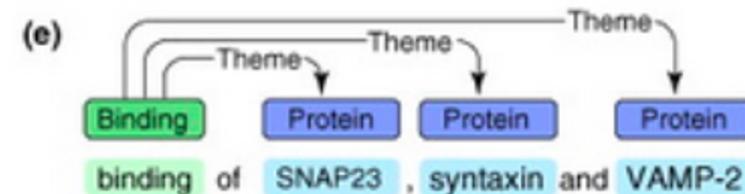
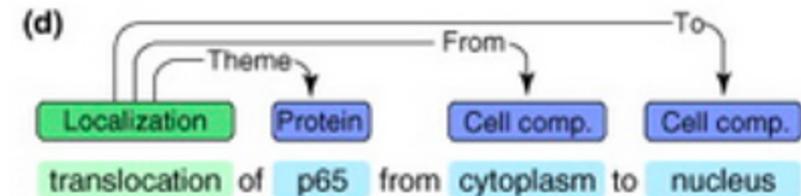
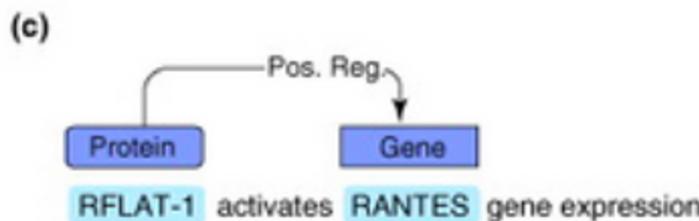
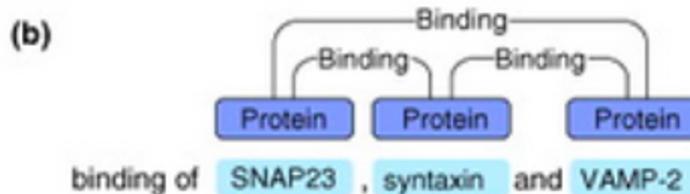
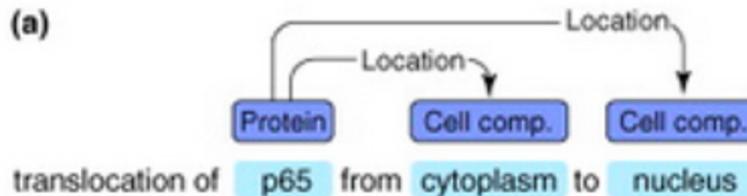
Cases where recursive might be better?



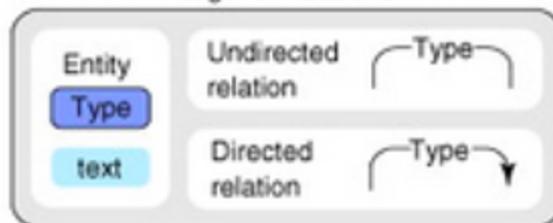
But recurrent structures might be able to learn this too...

Application of parse trees: info extraction

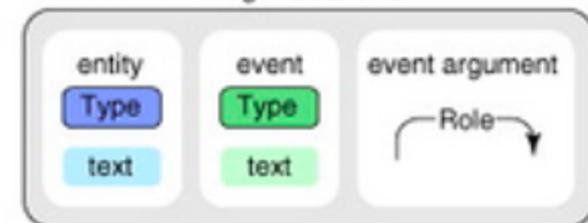
parse tree + rule based extraction = more robust



Legend : relations



Legend : events



Application of parse trees: info extraction

parse tree + rule based extraction = more robust

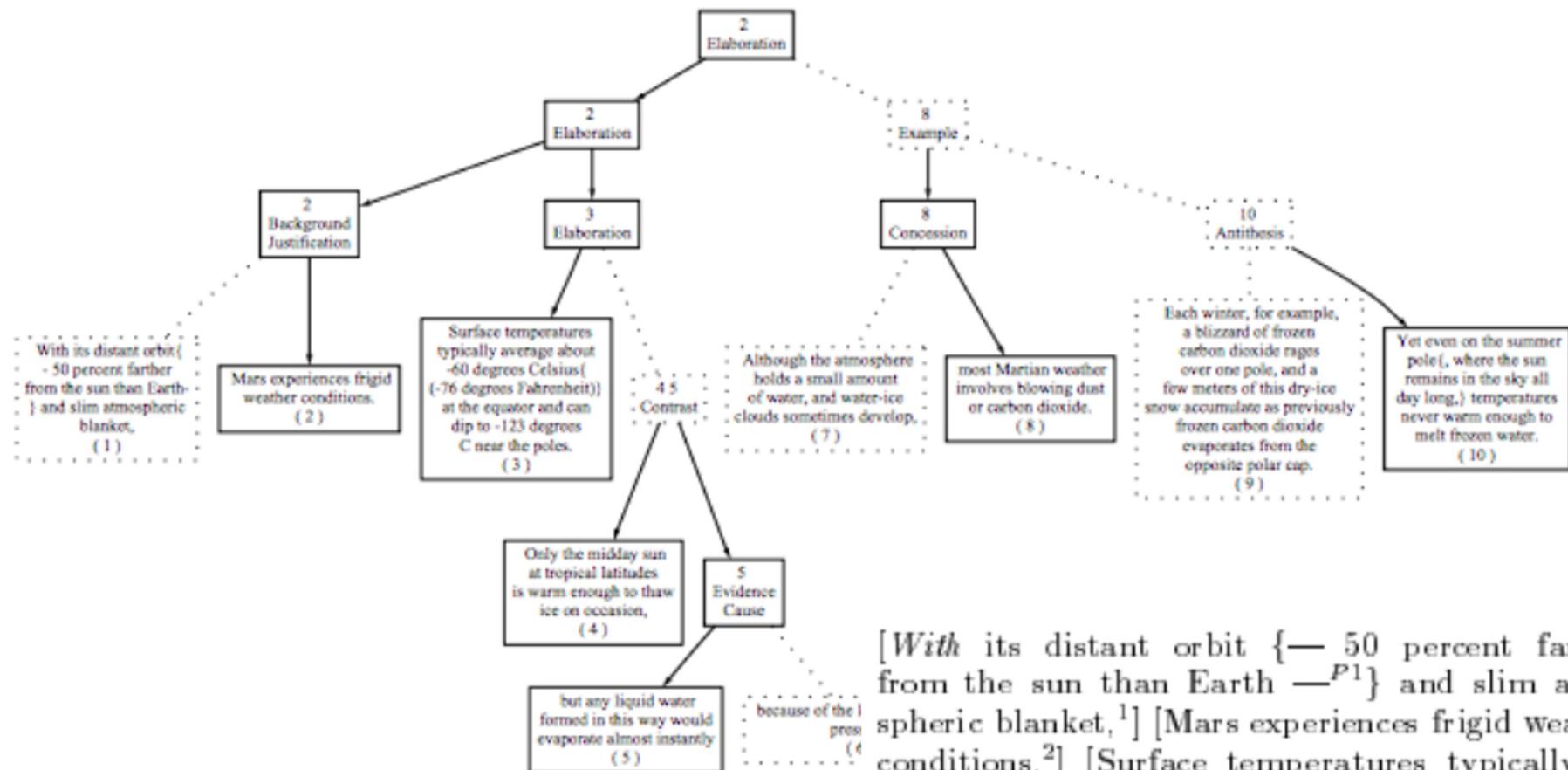
EX ขอ สั่ง พิซซ่า หน้า艰辛ยอี้น

Rule: สั่ง + XXX

ขอ สั่ง พิซซ่า ให้ มา สั่ง ที่ บ้าน ขอ หน้า艰辛ยอี้น นะ ครับ

The rule will break

A side note: discourse parsing



[With its distant orbit {— 50 percent farther from the sun than Earth ^{P1}} and slim atmospheric blanket,¹] [Mars experiences frigid weather conditions.²] [Surface temperatures typically average about —60 degrees Celsius (—76 degrees Fahrenheit) at the equator and can dip to —123 degrees C near the poles.³] [Only the midday sun at tropical latitudes is warm enough to thaw ice

Overview

- Types of grammars
 - Context Free Grammar
 - Probabilistic Context Free Grammar
 - CYK parser
 - Dependency Grammar
 - Transition-based parsing
 - Recursive neural networks

Constituents

- Groups of words behaving as a single units

- Ex: Noun phrase

Harry the Horse

The reason he comes into the house

They

A high-class spot such as Mindy's

- Checking for constituents

- See if they can appear in similar syntactic environments

They sit...

The reason he comes into the house is...

Context-Free Grammar (CFG)

- A grammar specifies what kind of parse tree can be generated.
- CFG or Phrase-Structure Grammars assumes the grammar is context-free
 - Most forms of natural language are context-free
 - Thai and English are typically CFG languages
 - Used in many programming languages
- CFG is based on constituent structures

A context-free grammar example

A CFG is defined by
 $G = (N, S, \Sigma, R)$

- $N = \{S, NP, VP, PP, DT, Vi, Vt, NN, IN\}$ Nonterminals
- $S = S$ Starting symbol
- $\Sigma = \{\text{sleeps, saw, man, woman, telescope, the, with, in}\}$

S	\Rightarrow	NP	VP
VP	\Rightarrow	Vi	
VP	\Rightarrow	Vt	NP
VP	\Rightarrow	VP	PP
NP	\Rightarrow	DT	NN
NP	\Rightarrow	NP	PP
PP	\Rightarrow	IN	NP

Vi	\Rightarrow	sleeps
Vt	\Rightarrow	saw
NN	\Rightarrow	man
NN	\Rightarrow	woman
NN	\Rightarrow	telescope
DT	\Rightarrow	the
IN	\Rightarrow	with
IN	\Rightarrow	in

Terminals

Production rules
 $X \rightarrow Y_1 \dots Y_n$
 Y_i can be terminal or nonterminal
The rule only relies on X (no context)
(vs context-sensitive)

Note: S =sentence, VP =verb phrase, NP =noun phrase, PP =prepositional phrase, DT =determiner, Vi =intransitive verb, Vt =transitive verb, NN =noun, IN =preposition

Generation using left-most derivation strategy

- Recursive strategy

```
left_most_derivation(A):
```

```
    if terminal(A):
```

```
        return A      Return if terminal
```

```
    rule = choose(rules(A))
```

```
    rhs = right_hand_side(rule)
```

```
    return concatenate([left_most_derivation(A') for A' in rhs])
```

- To generate we call `left_most_derivation(S)` (top down)
- A string belongs to the language of a CFG if there exist a sequence of left-most derivation that can generate the string

$$L = \{s \in \Sigma^* | s = \text{left_most_derivation}(S)\}$$

$R =$	$S \Rightarrow NP VP$
	$VP \Rightarrow Vi$
	$VP \Rightarrow Vt NP$
	$VP \Rightarrow VP PP$
	$NP \Rightarrow DT NN$
	$NP \Rightarrow NP PP$
	$PP \Rightarrow IN NP$

$Vi \Rightarrow$	sleeps
$Vt \Rightarrow$	saw
$NN \Rightarrow$	man
$NN \Rightarrow$	woman
$NN \Rightarrow$	telescope
$DT \Rightarrow$	the
$IN \Rightarrow$	with
$IN \Rightarrow$	in

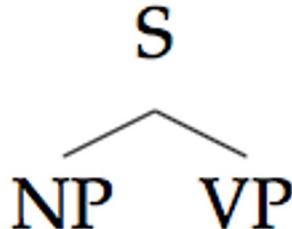
Example

- The woman saw the man with the telescope

$$R = \begin{array}{|c|c|c|} \hline S & \Rightarrow & NP \quad VP \\ \hline VP & \Rightarrow & Vi \\ \hline VP & \Rightarrow & Vt \quad NP \\ \hline VP & \Rightarrow & VP \quad PP \\ \hline NP & \Rightarrow & DT \quad NN \\ \hline NP & \Rightarrow & NP \quad PP \\ \hline PP & \Rightarrow & IN \quad NP \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline Vi & \Rightarrow & sleeps \\ \hline Vt & \Rightarrow & saw \\ \hline NN & \Rightarrow & man \\ \hline NN & \Rightarrow & woman \\ \hline NN & \Rightarrow & telescope \\ \hline DT & \Rightarrow & the \\ \hline IN & \Rightarrow & with \\ \hline IN & \Rightarrow & in \\ \hline \end{array}$$

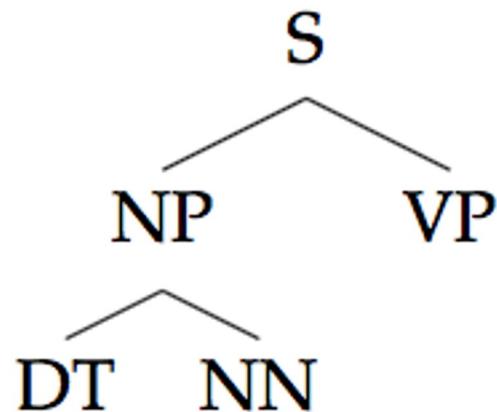
Example

- The woman saw the man with the telescope


$$R = \boxed{\begin{array}{lll} S & \Rightarrow & NP \quad VP \\ \hline VP & \Rightarrow & Vi \\ VP & \Rightarrow & Vt \quad NP \\ VP & \Rightarrow & VP \quad PP \\ \hline NP & \Rightarrow & DT \quad NN \\ NP & \Rightarrow & NP \quad PP \\ \hline PP & \Rightarrow & IN \quad NP \end{array}} \quad \boxed{\begin{array}{lll} Vi & \Rightarrow & sleeps \\ Vt & \Rightarrow & saw \\ \hline NN & \Rightarrow & man \\ NN & \Rightarrow & woman \\ NN & \Rightarrow & telescope \\ \hline DT & \Rightarrow & the \\ IN & \Rightarrow & with \\ IN & \Rightarrow & in \end{array}}$$

Example

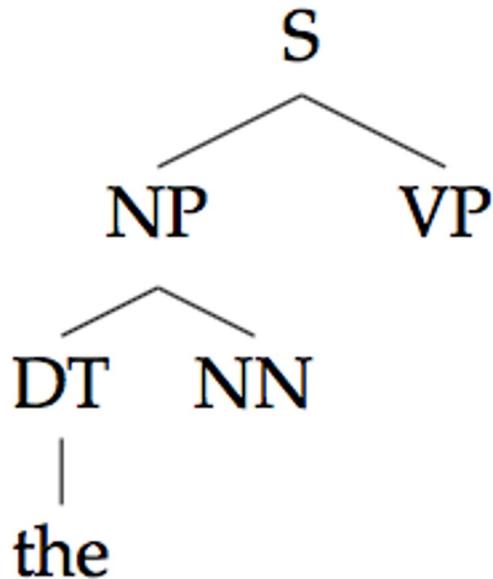
- The woman saw the man with the telescope



$R =$	$S \Rightarrow NP \quad VP$	$Vi \Rightarrow \text{sleeps}$
	$VP \Rightarrow Vi$	$Vt \Rightarrow \text{saw}$
	$VP \Rightarrow Vt \quad NP$	$NN \Rightarrow \text{man}$
	$VP \Rightarrow VP \quad PP$	$NN \Rightarrow \text{woman}$
	$NP \Rightarrow DT \quad NN$	$NN \Rightarrow \text{telescope}$
	$NP \Rightarrow NP \quad PP$	$DT \Rightarrow \text{the}$
	$PP \Rightarrow IN \quad NP$	$IN \Rightarrow \text{with}$
		$IN \Rightarrow \text{in}$

Example

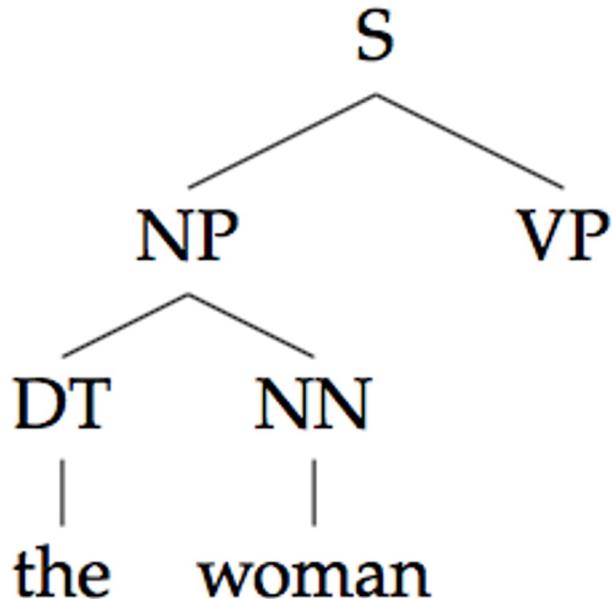
- The woman saw the man with the telescope



$R =$	$S \Rightarrow NP \quad VP$	$Vi \Rightarrow \text{sleeps}$
	$VP \Rightarrow Vi$	$Vt \Rightarrow \text{saw}$
	$VP \Rightarrow Vt \quad NP$	$NN \Rightarrow \text{man}$
	$VP \Rightarrow VP \quad PP$	$NN \Rightarrow \text{woman}$
	$NP \Rightarrow DT \quad NN$	$NN \Rightarrow \text{telescope}$
	$NP \Rightarrow NP \quad PP$	$DT \Rightarrow \text{the}$
	$PP \Rightarrow IN \quad NP$	$IN \Rightarrow \text{with}$
		$IN \Rightarrow \text{in}$

Example

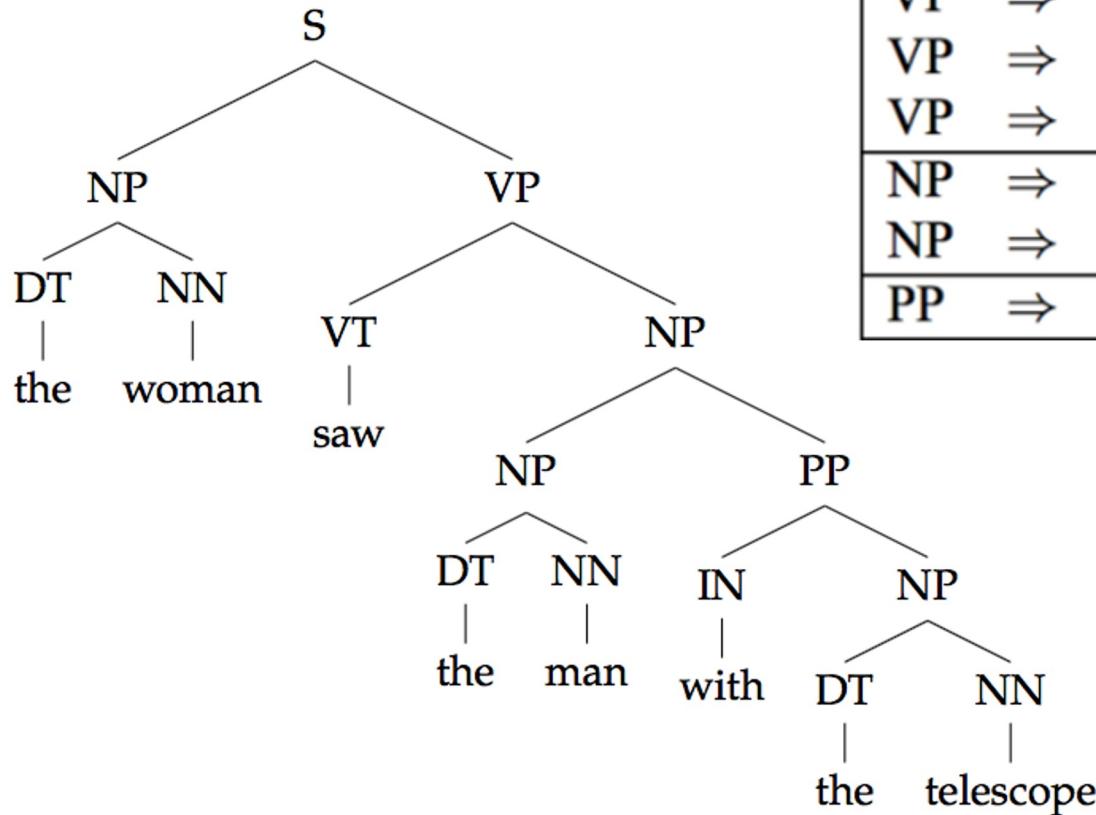
- The woman saw the man with the telescope


$$R = \begin{array}{l|ll} S & \Rightarrow & NP \quad VP \\ \hline VP & \Rightarrow & Vi \\ VP & \Rightarrow & Vt \quad NP \\ VP & \Rightarrow & VP \quad PP \\ \hline NP & \Rightarrow & DT \quad NN \\ NP & \Rightarrow & NP \quad PP \\ \hline PP & \Rightarrow & IN \quad NP \end{array}$$

Vi	\Rightarrow	sleeps
Vt	\Rightarrow	saw
NN	\Rightarrow	man
NN	\Rightarrow	woman
NN	\Rightarrow	telescope
DT	\Rightarrow	the
IN	\Rightarrow	with
IN	\Rightarrow	in

Example

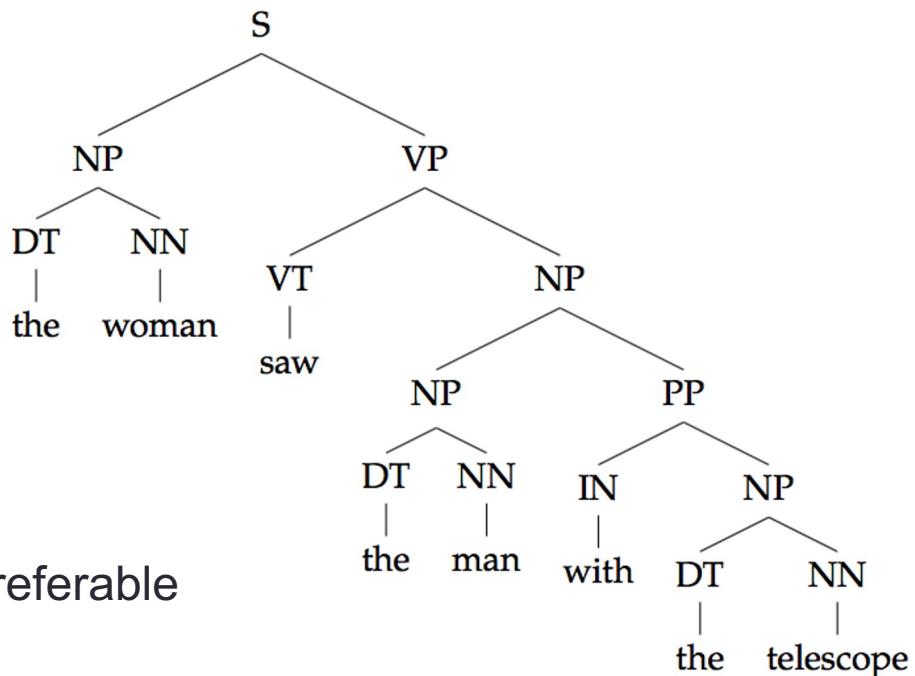
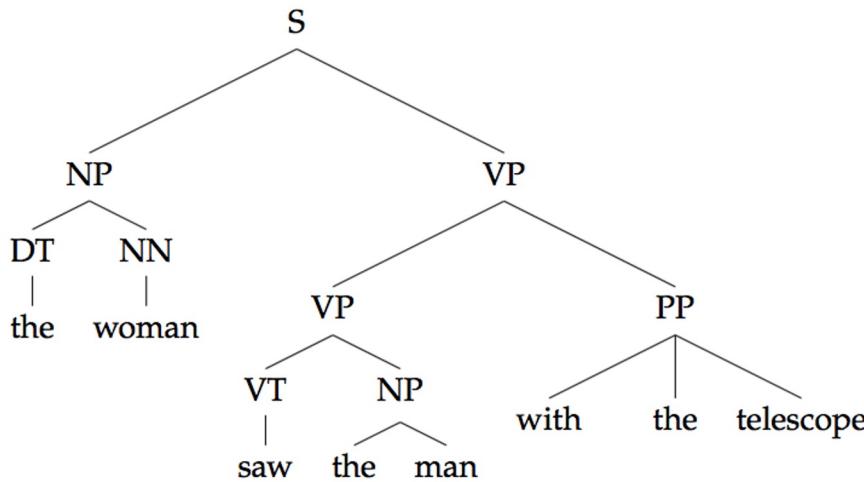
- The woman saw the man with the telescope



$R =$	$S \Rightarrow NP \quad VP$	$Vi \Rightarrow sleeps$
	$VP \Rightarrow Vi$	$Vt \Rightarrow saw$
	$VP \Rightarrow Vt \quad NP$	$NN \Rightarrow man$
	$VP \Rightarrow VP \quad PP$	$NN \Rightarrow woman$
	$NP \Rightarrow DT \quad NN$	$NN \Rightarrow telescope$
	$NP \Rightarrow NP \quad PP$	$DT \Rightarrow the$
	$PP \Rightarrow IN \quad NP$	$IN \Rightarrow with$
		$IN \Rightarrow in$

Ambiguities

- There can be multiple derivations for the same string
- These sentences are ambiguous as each parse represents a different meaning



Need a way to say which one is more preferable
Add probabilities!

Probabilistic Context-Free Grammar

- Production rules now have **probabilities**

S	\Rightarrow	NP	VP	1.0
VP	\Rightarrow	Vi		0.4
VP	\Rightarrow	Vt	NP	0.4
VP	\Rightarrow	VP	PP	0.2
NP	\Rightarrow	DT	NN	0.3
NP	\Rightarrow	NP	PP	0.7
PP	\Rightarrow	P	NP	1.0

Vi	\Rightarrow	sleeps	1.0
Vt	\Rightarrow	saw	1.0
NN	\Rightarrow	man	0.7
NN	\Rightarrow	woman	0.2
NN	\Rightarrow	telescope	0.1
DT	\Rightarrow	the	1.0
IN	\Rightarrow	with	0.5
IN	\Rightarrow	in	0.5

The probability of a (sentence, parse tree) pair is

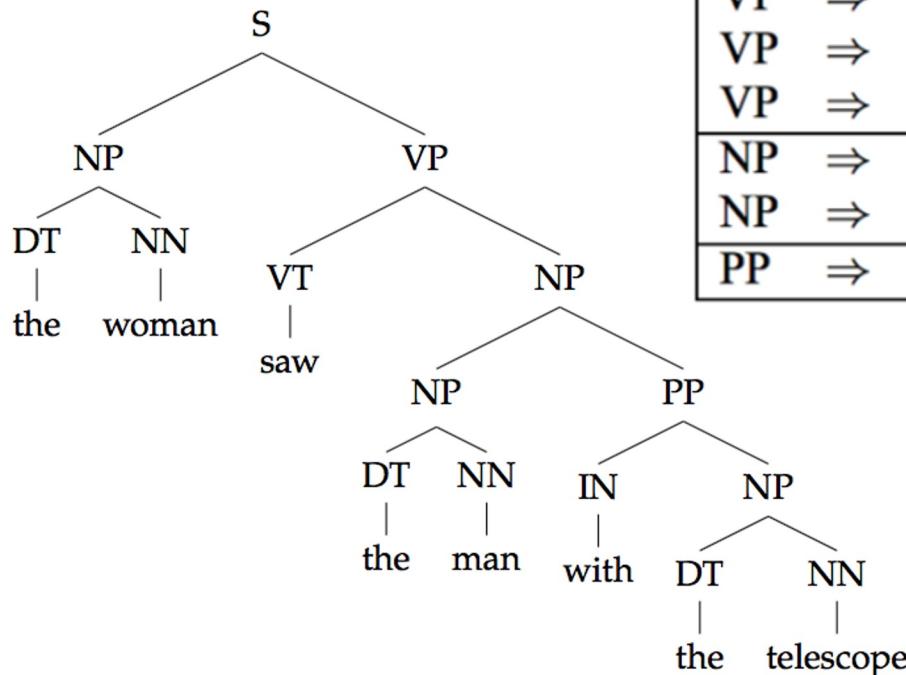
$$p(S, T) = \prod_{i=1}^m p(\alpha_i \rightarrow \beta_i | \alpha_i)$$

m is the number of transitions

$$P(NN \rightarrow \text{man} | NN) = 0.7$$

P(S, T)

- $P(S, T) = [1.0 * 0.3 * 1.0 * 0.2] * [0.4 * 1.0 * 0.7 * (0.3 * 1.0 * 0.7) * (1.0 * 0.5 * 0.3 * 1.0 * 0.1)]$



S	\Rightarrow	NP	VP	1.0
VP	\Rightarrow	Vi		0.4
VP	\Rightarrow	Vt	NP	0.4
VP	\Rightarrow	VP	PP	0.2
NP	\Rightarrow	DT	NN	0.3
NP	\Rightarrow	NP	PP	0.7
PP	\Rightarrow	P	NP	1.0

Vi	\Rightarrow	sleeps	1.0
Vt	\Rightarrow	saw	1.0
NN	\Rightarrow	man	0.7
NN	\Rightarrow	woman	0.2
NN	\Rightarrow	telescope	0.1
DT	\Rightarrow	the	1.0
IN	\Rightarrow	with	0.5
IN	\Rightarrow	in	0.5

Estimating transition probabilities

- Counts from training set!

$$P(S \rightarrow NP\ VP|S) = \frac{count(S \rightarrow NP\ VP)}{count(S)}$$

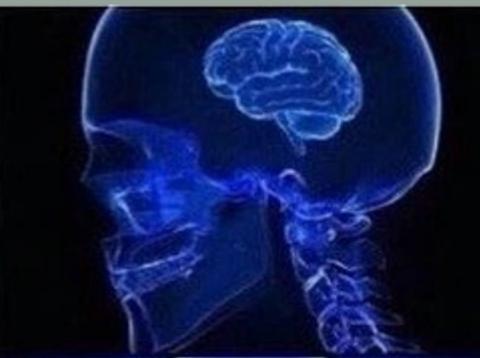
PCFG related questions

- What is the most likely parse? (**parsing task**)
 - $\operatorname{argmax}_T P(T, S)$
- What is the probability of the sentence, $P(S)$? (**Language modeling task**)
 - $P(S) = \sum_T P(T, S)$

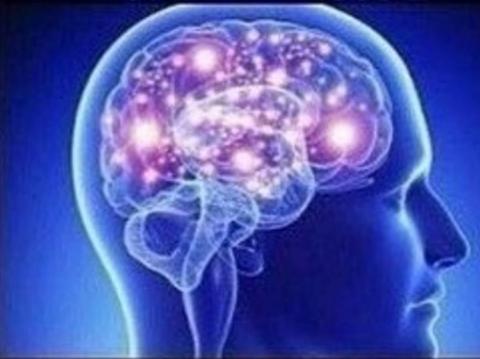
How?

Dynamic programming
(CYK algorithm –
Cocke–Younger–Kasami algorithm)

BRUTE FORCE



MAXIMAL MATCHING



VITERBI



CYK



Chomsky Normal Form

- CYK can be used if the CFG is in Chomsky Normal Form
- A CFG is in Chomsky Normal Form if each rule either converts to **two nonterminals** or a **single terminal**.

$$X \rightarrow Y_1 Y_2$$

Uppercase letters mean nonterminal
Lowercase letters mean terminal

$$X \rightarrow y$$

- Any CFG can be converted to CNF

- For example, NP \rightarrow DT, ADJ, NN
 - Turns into two rules
 - NP \rightarrow DT, ADJP
 - ADJP \rightarrow ADJ, NN

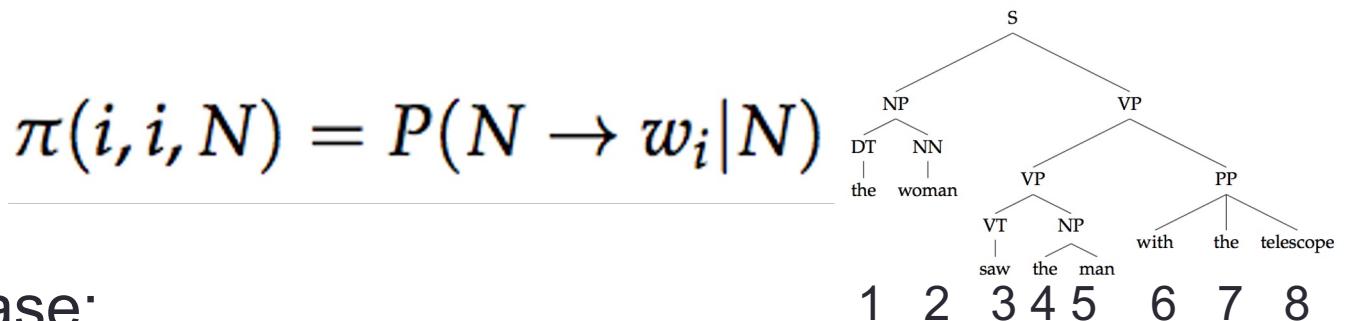
$R =$	$S \Rightarrow NP VP$	$Vi \Rightarrow sleeps$	Production rules
	$VP \Rightarrow Vi$	$Vt \Rightarrow saw$	$X \rightarrow Y_1 \dots Y_n$
	$VP \Rightarrow Vt NP$	$NN \Rightarrow man$	Y_i can be terminal or nonterminal
	$VP \Rightarrow VP PP$	$NN \Rightarrow woman$	The rule only relies on
	$NP \Rightarrow DT NN$	$NN \Rightarrow telescope$	X (no context) (vs context-sensitive)
	$NP \Rightarrow NP PP$	$DT \Rightarrow the$	
	$PP \Rightarrow IN NP$	$IN \Rightarrow with$	
		$IN \Rightarrow in$	

CYK algorithm for parsing

3 dimensional

- $\pi(i, j, N)$ probability that words i to j can be generated by nonterminal N

- Base case: $\pi(i, i, N) = P(N \rightarrow w_i | N)$



$$\pi(i, j, N) = \max_{k, P, Q} Pr(N \rightarrow P \mid Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

where $k \in \{i, \dots, j - 1\}$, $P \in \mathcal{N}$, and $Q \in \mathcal{N}$.

Saves the rule that gave max probability to backtrack

CYK algorithm for language modeling

- $\pi(i, j, N)$ probability that words i to j can be generated by nonterminal N
- Base case: $\pi(i, i, N) = P(N \rightarrow w_i | N)$

- Inductive case:

$$\begin{aligned}\pi(i, j, N) &= \max_{k, P, Q} Pr(N \rightarrow P \ Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q) \\ \text{argmax}_T P(T, S) \quad \text{vs} \quad P(S) &= \sum_T P(T, S)\end{aligned}$$

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P \ Q | N) \cdot \pi(\underline{i}, \underline{k}, P) \cdot \pi(\underline{k + 1}, \underline{j}, Q)$$

where $k \in \{i, \dots, j - 1\}$, $P \in \mathcal{N}$, and $Q \in \mathcal{N}$.

Language modeling example CYK with PCFG

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N)$$

A

a			
b			
c			

B

a			
b			
c			

Base case

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, i, N) = P(N \rightarrow w_i | N)$$

A

a	0.2		
b		0	
c			0

B

a	0		
b		0.1	
c			0.2

1st step

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

A

B

a	0.2		
b		0	
c			0

Only consider i, j where $i < j$

a	0		
b		0.1	
c			0.2

1st step

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

A

a	0.2	0.016	
b		0	
c			0

$$\begin{aligned}
 \pi[1, 2, A] &= P(A \rightarrow AA) \cdot P(A \rightarrow a | A) \cdot P(A \rightarrow b | A) \\
 &\quad + P(A \rightarrow AB) \cdot P(A \rightarrow a | A) \cdot P(B \rightarrow b | B) \\
 &\quad + P(A \rightarrow BA) \cdot P(B \rightarrow a | B) \cdot P(A \rightarrow b | A) \\
 &\quad + P(A \rightarrow BB) \cdot P(B \rightarrow a | B) \cdot P(B \rightarrow b | B) \\
 &= 0 + P(A \rightarrow AB) \cdot \pi[1, 1, A] \cdot \pi[2, 2, B] + 0 + 0 \\
 &= 0.8 \cdot 0.2 \cdot 0.1 = 0.016
 \end{aligned}$$

1st step

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

$$\begin{aligned}\pi(1, 2, B) &= P(B \rightarrow BB)P(B \rightarrow a | B)P(B \rightarrow b | B) \\ &= P(B \rightarrow BB)\pi(1, 1, B)\pi(2, 2, B) \\ &= 0.7 * 0 * 0.1 = 0\end{aligned}$$

B

a	0	0	
b		0.1	
c			0.2

2nd step

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

A

B

a	0.2	0.016	
b		0	0
c		0	

a	0	0	
b		0.1	0.014
c			0.2

$$\pi(2,3,A) = P(A \rightarrow A B) \pi(2,2,A) \pi(3,3,B) = 0.8 * 0 * 0.2 = 0$$

2nd step

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

A

B

a	0.2	0.016	
b		0	0
c		0	

a	0	0	
b		0.1	0.014
c			0.2

$$\pi(2, 3, B) = P(B \rightarrow B | B) \pi(2, 2, B) \pi(3, 3, B) + 0 = 0.7 * 0.1 * 0.2 = 0.014$$

Finish

- $N = \{A, B\}$
- $\Sigma = \{a, b, c\}$
- $S = \{A\}$

Find $P(\text{"abc"})$

$A \rightarrow AB$	0.8
$A \rightarrow a$	0.2
$B \rightarrow BB$	0.7
$B \rightarrow b$	0.1
$B \rightarrow c$	0.2

$$\pi(i, j, N) = \sum_{k, P, Q} P(N \rightarrow P | Q | N) \cdot \pi(i, k, P) \cdot \pi(k + 1, j, Q)$$

A

B

a	0.2	0.016	0.0048
b		0	0
c		0	0

a	0	0	0
b		0.1	0.014
c			0.2

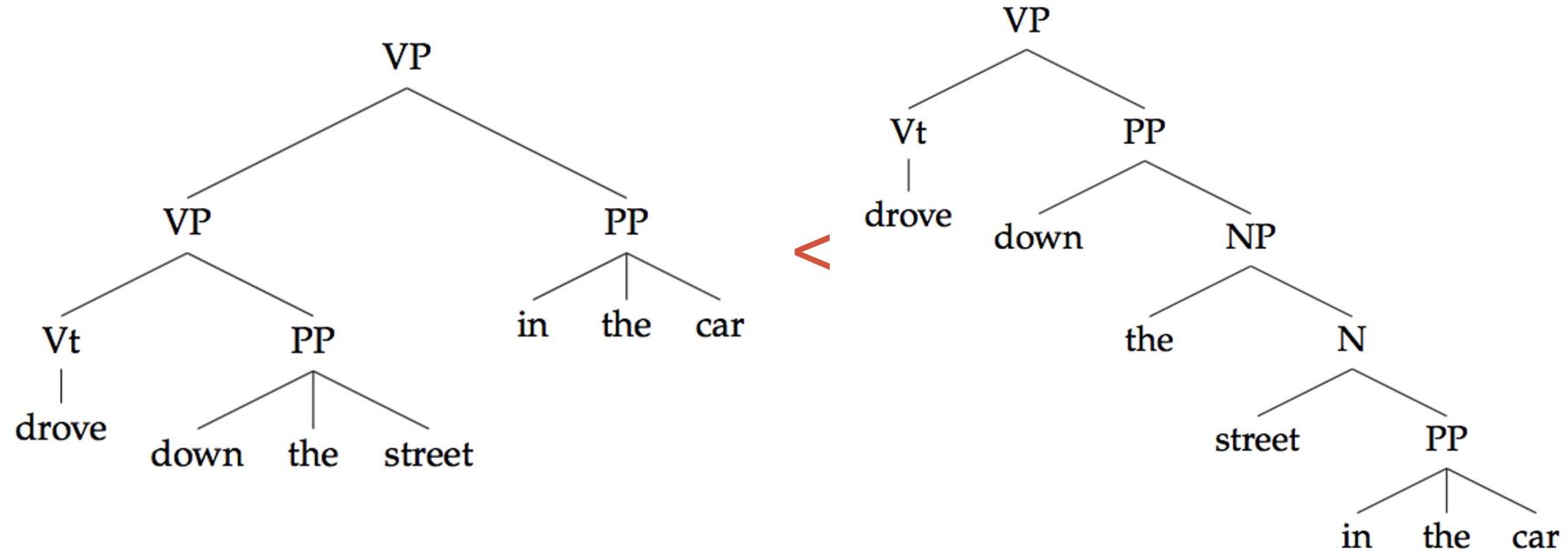
$$\begin{aligned} \pi(1, 3, A) &= P(A \rightarrow A | B) \pi(1, 2, A) \pi(3, 3, B) + P(A \rightarrow A | B) \pi(1, 1, A) \pi(2, 3, B) \\ &= 0.8 * 0.016 * 0.2 + 0.8 * 0.2 * 0.014 = 0.0048 \end{aligned}$$

PCFG weakness

- Lack of sensitivity to lexical info (does not consider semantics)
- Lack of sensitivity to structural frequency

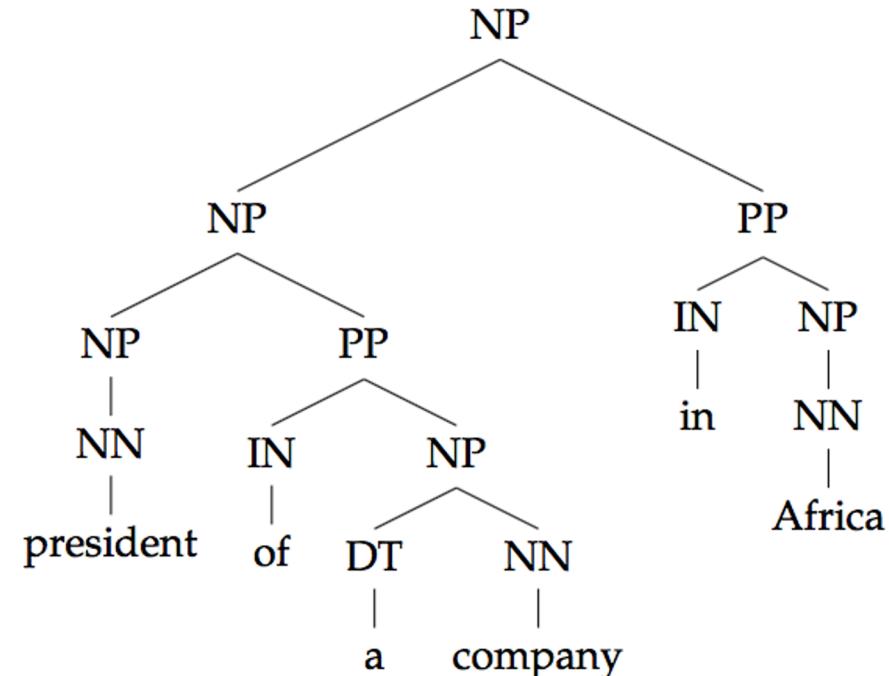
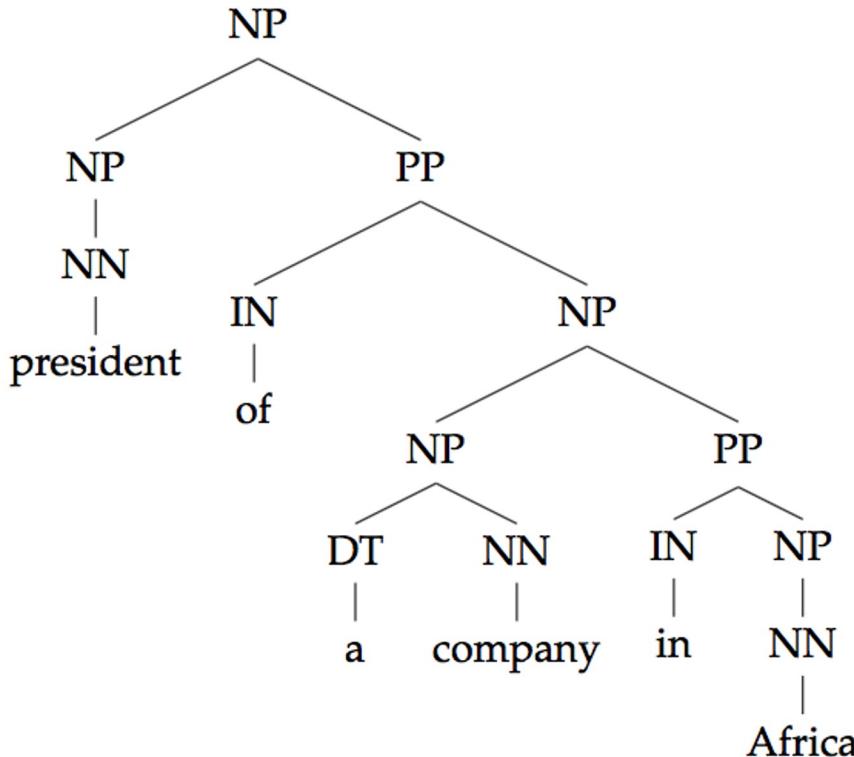
Lack of lexical info

The probabilities only see terminals and expansions



The street in the car

Lack of sensitivity to structural frequency



Both trees have same expansion and so same probability

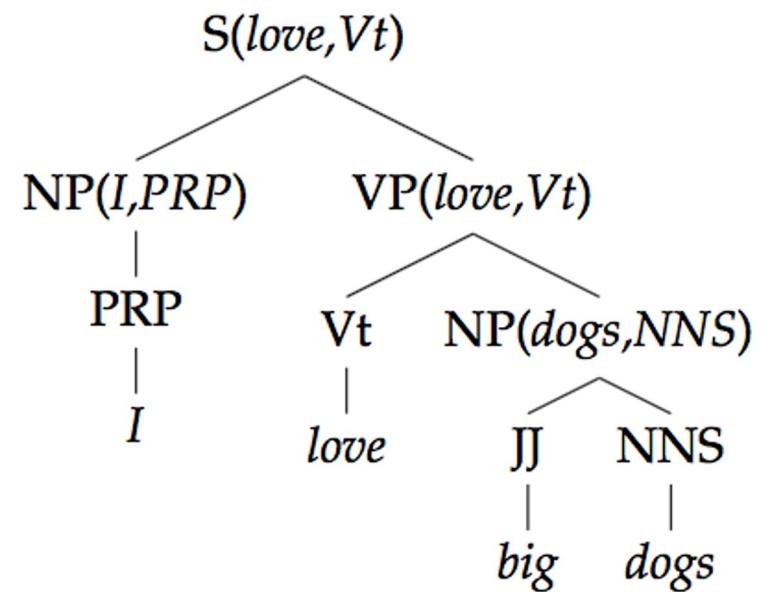
Turns out, left hand side structure appears more frequently in the training set

Methods to improve PCFG

- Add lexical information to the tree
 - Label each node with associate words
 - Lexicalized trees
 - Latent variable grammars
 - Look at bigger portions of the tree at a time

S	⇒	NP	VP	1.0
VP	⇒	Vi		0.4
VP	⇒	Vt	NP	0.4
VP	⇒	VP	PP	0.2
NP	⇒	DT	NN	0.3
NP	⇒	NP	PP	0.7
PP	⇒	P	NP	1.0

Vi	⇒	sleeps	1.0
Vt	⇒	saw	1.0
NN	⇒	man	0.7
NN	⇒	woman	0.2
NN	⇒	telescope	0.1
DT	⇒	the	1.0
IN	⇒	with	0.5
IN	⇒	in	0.5



Lexicalized tree

Overview

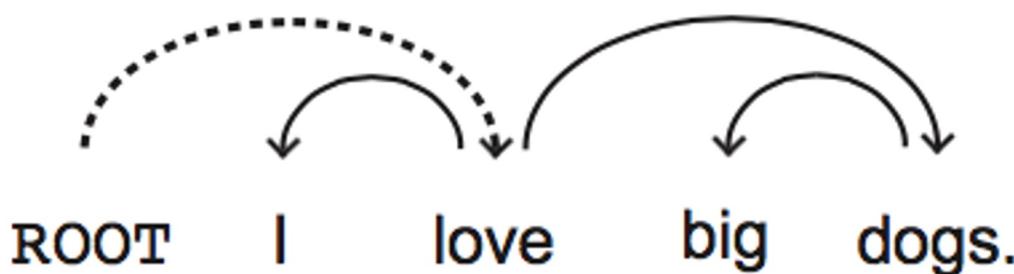
- Types of grammars
 - Context Free Grammar
 - Probabilistic Context Free Grammar
 - CYK parser
 - Dependency Grammar
 - Transition-based parsing
 - Recursive neural networks

Dependency grammar

- CFG is based on constituency relation
- In dependency grammar the structure is composed of lexical items (words) linked by edges to form a tree
- Assumptions
 - Each words in a sentence is related or modifies another word
 - All words have a direct or indirect relation to the main verb

Example

- Add ROOT node as the root of the tree
- The main verb always point to ROOT



- Each arc can have a category for the relationship.
- Each word can have a PoS label

A \rightarrow B means
A governs B or
B depends on A
A is the head of B

Constituency structures vs dependency structures

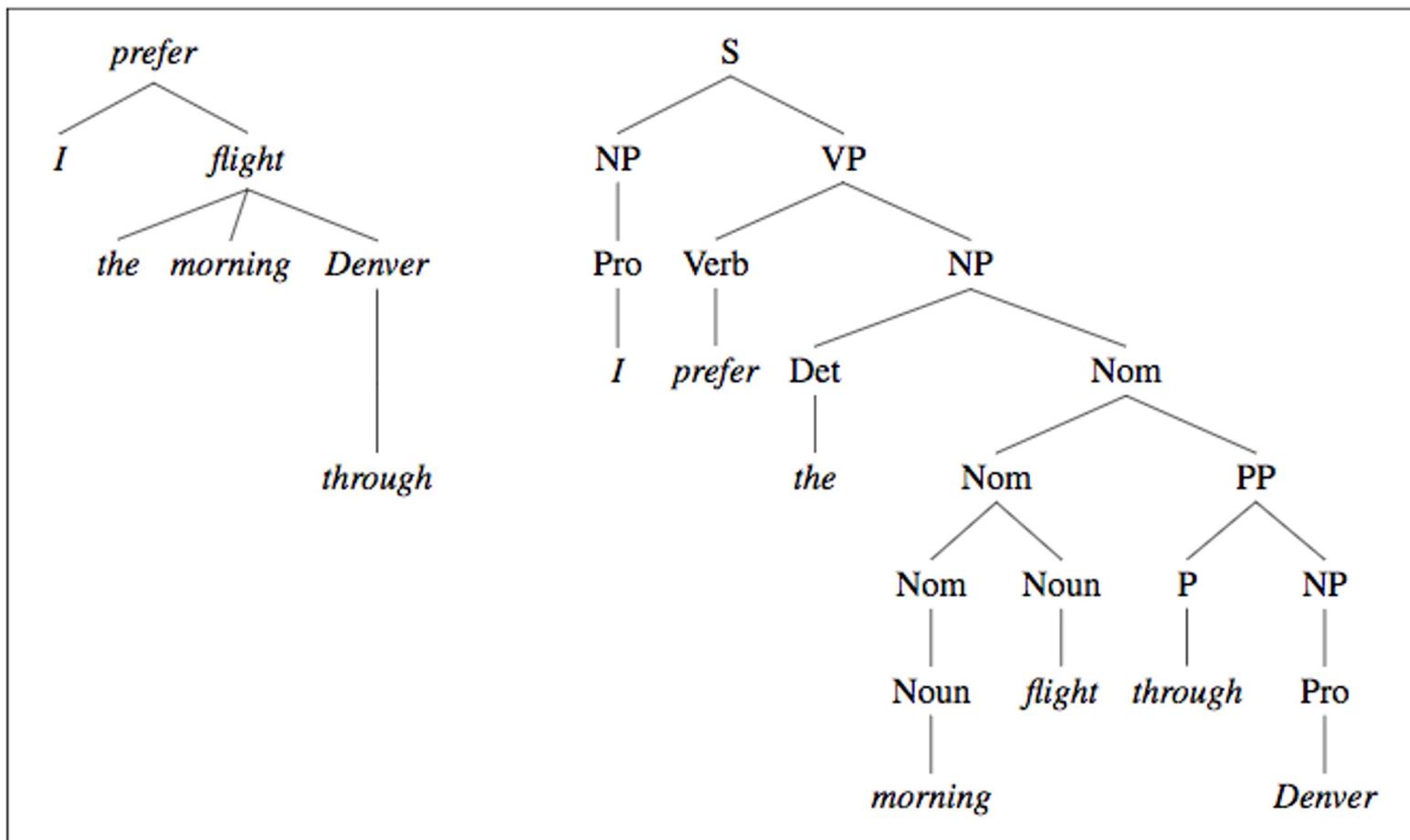


Figure 14.1 A dependency-style parse alongside the corresponding constituent-based analysis for *I prefer the morning flight through Denver.*

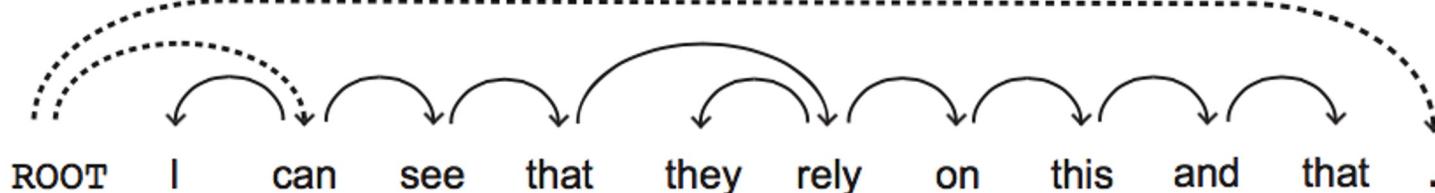
Constituency structures vs dependency structures

- Constituency structures use more nodes to represent sentences at different levels.
- Constituency structures explicitly label non-terminal nodes (NP vs VP)
- Constituency structures encode more info than dependency structures
- You can convert constituency structures to dependency structures
 - Dependency parsers trained on this is usually better than dependency parsers trained on original dependency structures

Criteria for heads (basics)

- Head H. Dependent D
- D modifies H
 - Big (D) dogs (H), willow (D) tree (H)
- H can often replace D
 - I love big (D) dogs (H) -> I love dogs
- H is obligatory while D sometimes is optional
- H determines whether D is obligatory
 - Sarah sneezed (H) vs George kicks (H) the chair (D)
- More criteria! **Mostly depends on corpus standard**

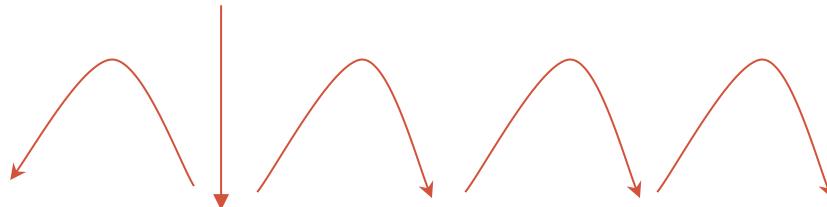
Ex



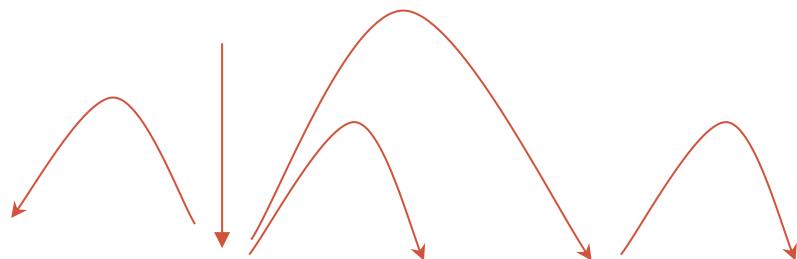
Dependency and meaning

- Scientists study whales from space
- Scientists study whales from space

Dependency and meaning



- Scientists study whales from space



- Scientists study whales from space

Sample Thai sentences

- วัน หนึ่ง เกิด เหตุ ที่ ไม่ คาดคิด ขึ้น มา
- การ ที่ เรา จะ รู้ ถึง สิ่ง เหล่านั้น ได้ ก็ ต้อง อาศัย ปัจจัย หลาย ๆ อย่าง ประกอบ เข้า ด้วยกัน

Note these criteria depends on the corpus convention

Sample Thai sentences

- วัน หนึ่ง เกิด เหตุ ที่ ไม่ คาดคิด ขึ้น มา

Note these criteria depends on the corpus convention

Sample Thai sentences

• การ ที่ เราก็ จะ รู้ ถึง สิ่ง เหล่านั้น ได้ ก็ ต้อง อาศัย

ปัจจัย หลาย ๆ อย่าง ประกอบ เข้า ด้วยกัน

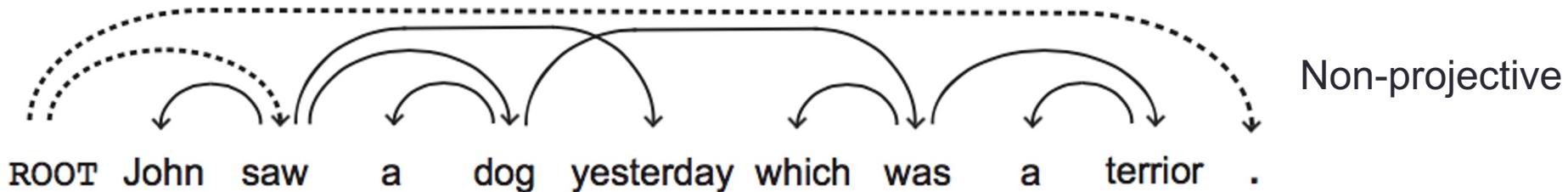
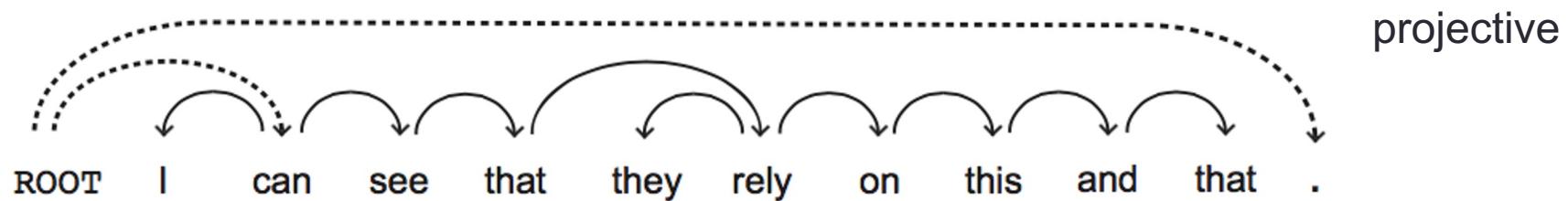
Note these criteria depends on the corpus convention

Dependency graph requirements

- Syntactic structure is complete (connectedness, spanning)
- Hierarchical (acyclic)
- Every word has a single head

Projectivity

- A dependency graph is projective if the arcs do not cross



English and Thai are mostly projective.

Some languages are more non-projective than others,
for example, German, Dutch, Czech.

When picking a parser algo, check whether it assumes projectivity

Transition-based parsing (Nivre 2007)

- Use a stack and buffer data structure and sequentially add edges
- Characteristics
 - Greedy algo. Only goes left to right. No backtracking
 - Fast $O(n)$
 - Requires projectivity
 - The algo is closely related to how humans parse sentences (left to right one word at a time instead of looking at the sentence as a whole)

Arc-standard Transition-based parsing

- A stack σ , written with top of the stack to the right
 - Starts with the ROOT symbol
 - A buffer β , written with top to the left
 - Starts with the input sentence
 - A set of dependency arcs A
 - Starts of empty
 - A set of actions
- $$\sigma = [\text{ROOT}], \beta = w_1, \dots, w_n, A = \emptyset$$
1. Shift $\sigma, w_i | \beta, A \rightarrow \sigma | w_i, \beta, A$ Move buffer to stack
 2. Left-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_j, \beta, A \cup \{r(w_j, w_i)\}$
 3. Right-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_i, \beta, A \cup \{r(w_i, w_j)\}$
- Finishes when β becomes empty

I ate fish

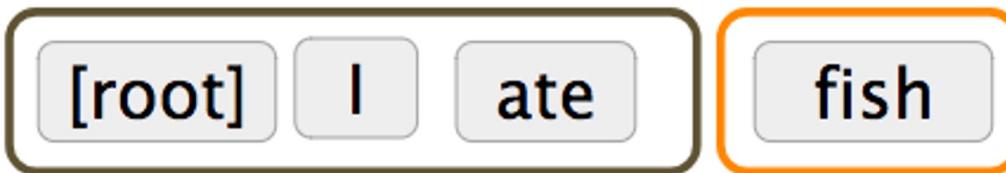
Start



Shift



Shift



Start: $\sigma = [\text{ROOT}], \beta = w_1, \dots, w_n, A = \emptyset$

1. Shift $\sigma, w_i|\beta, A \rightarrow \sigma|w_i, \beta, A$
2. Left-Arc, $\sigma|w_i|w_j, \beta, A \rightarrow \sigma|w_j, \beta, A \cup \{r(w_j, w_i)\}$
3. Right-Arc, $\sigma|w_i|w_j, \beta, A \rightarrow \sigma|w_i, \beta, A \cup \{r(w_i, w_j)\}$

Finish: $\beta = \emptyset$

I ate fish

Left Arc



Shift



Right Arc



Right Arc



Discriminative parsing

- How to choose an action?
 - Shift, left-arc, right-arc
- Each action is predicted by a discriminative classifier (SVM, logistic regression, Neural networks) over legal moves
 - Features: top two word from stack, POS, children info; first word in buffer, POS, children info; etc.
- Greedy and no beamsearch
 - But you can include beamsearch (modern parsers do)

Discriminative parsing with neural networks

Shift, left, right (in actual, left/right + type of dependency)
So $2N+1$, where N = dependency types

Softmax layer:

$$p = \text{softmax}(W_2 h)$$

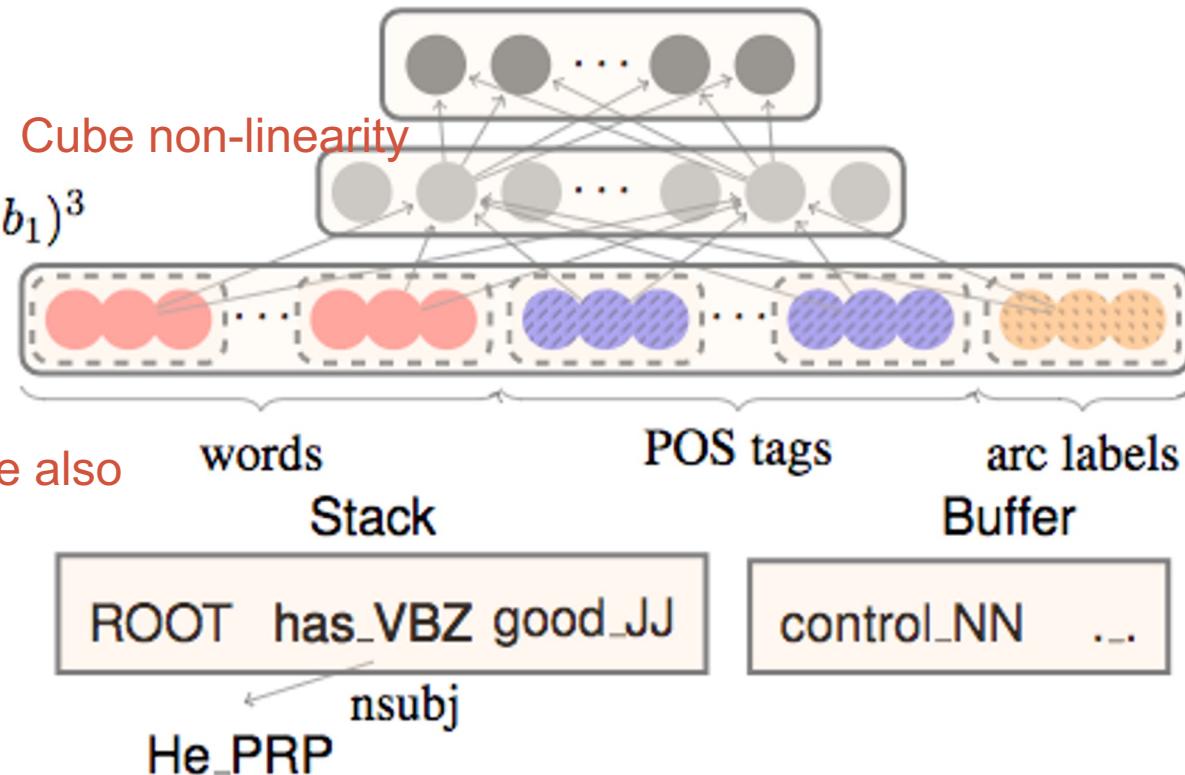
Hidden layer:

$$h = (W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1)^3$$

Input layer: $[x^w, x^t, x^l]$

POS and arc labels are also embeddings

Configuration



Improvements

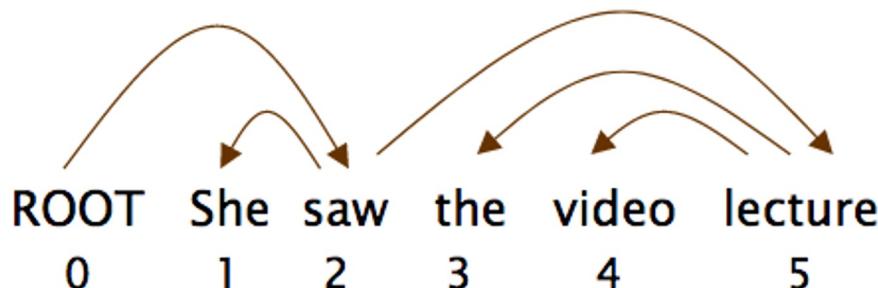
- Bigger networks
- Beam search
- Global inference over the sequence (CRF style)
- Lead to SyntaxNet and Parsey McParseFace model

<https://github.com/tensorflow/models/tree/master/research/syntaxnet>

<https://research.googleblog.com/2016/05/announcing-syntaxnet-worlds-most.html>

Parsing evaluation

- Labeled parsing accuracy (LAS)
- Unlabeled parsing accuracy (UAS)



$$\text{Acc} = \frac{\# \text{ correct deps}}{\# \text{ of deps}}$$

$$\text{UAS} = 4 / 5 = 80\%$$

$$\text{LAS} = 2 / 5 = 40\%$$

Gold

1	2	She	nsubj
2	0	saw	root
3	5	the	det
4	5	video	nn
5	2	lecture	obj

Parsed

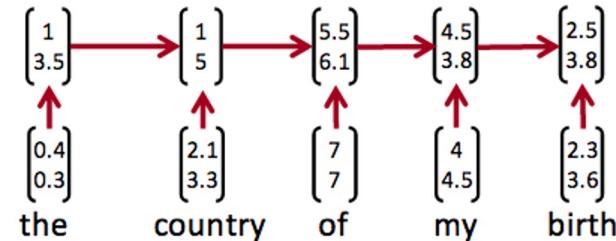
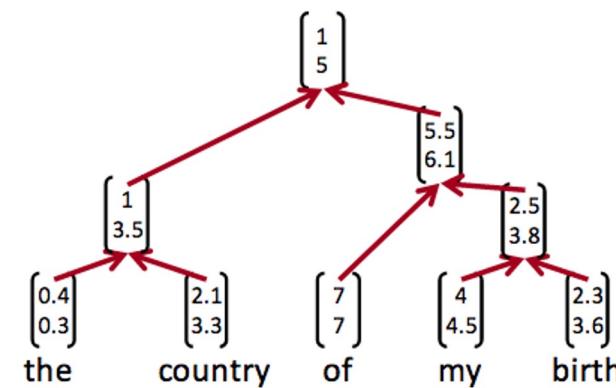
1	2	She	nsubj
2	0	saw	root
3	4	the	det
4	5	video	nsubj
5	2	lecture	ccomp

UAS of Parsey McParseFace

Model	News	Web	Questions
Martins et al. (2013)	93.10	88.23	94.21
Zhang and McDonald (2014)	93.32	88.65	93.37
Weiss et al. (2015)	93.91	89.29	94.17
Andor et al. (2016)*	94.44	90.17	95.40
Parsey McParseface	94.15	89.08	94.77

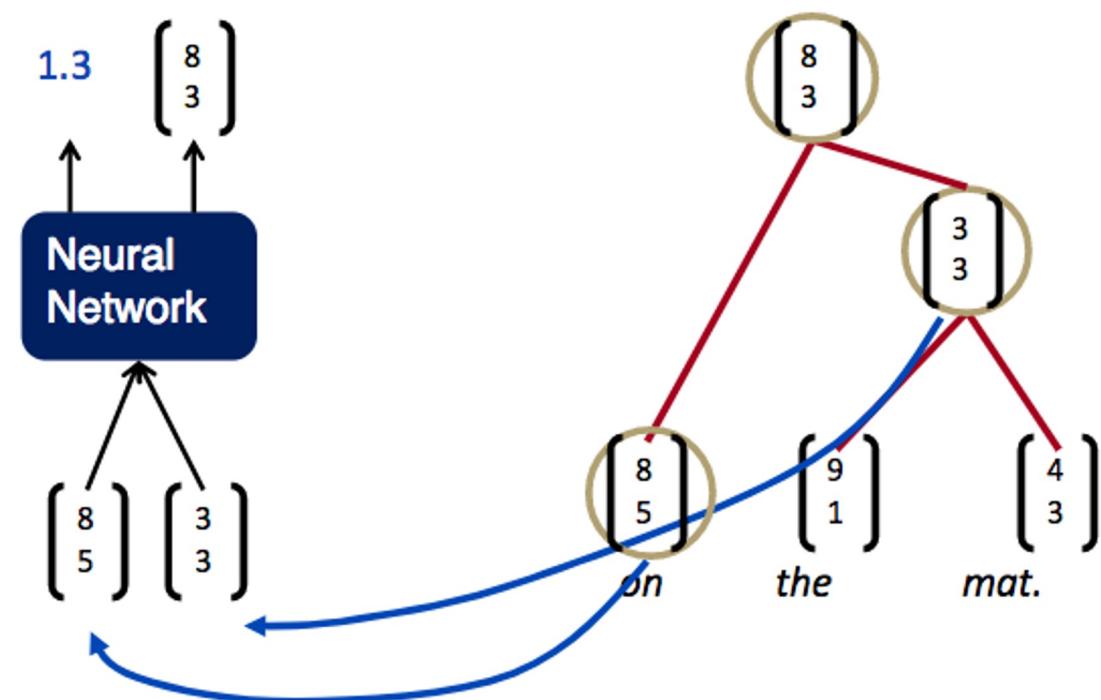
Recursive neural networks

- Not really used in parsing anymore but interesting concept
- Recursive vs Recurrent

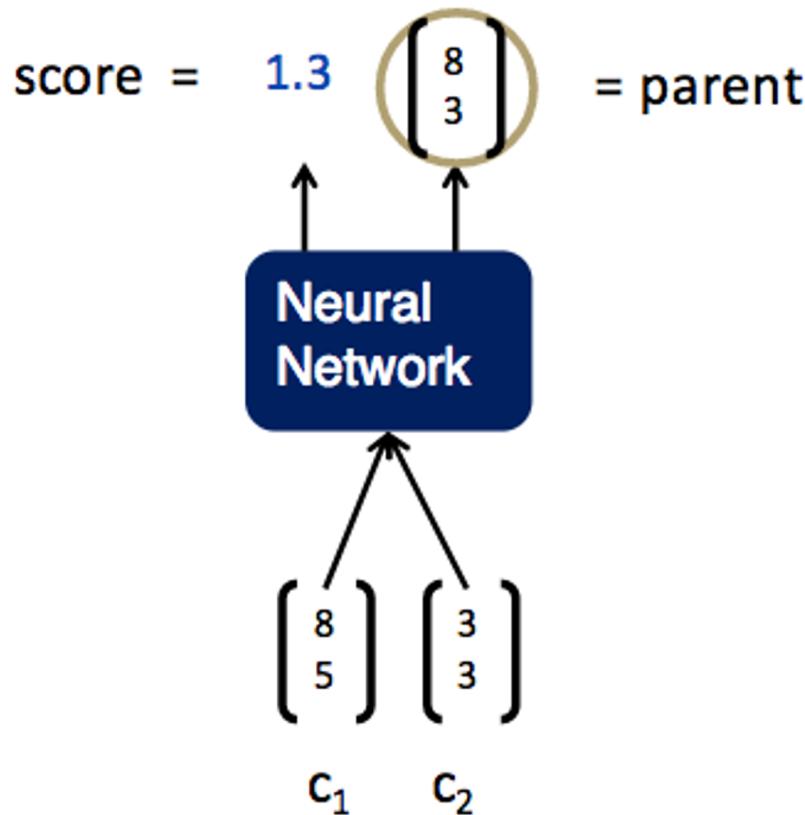


Recursive neural networks

- Concept: try different connections, see if which one gives the highest score (graph-based dependency parsers)
- Inputs: two candidate children's representations
- Output:
 - Semantic representation of the parent
 - Score of new node



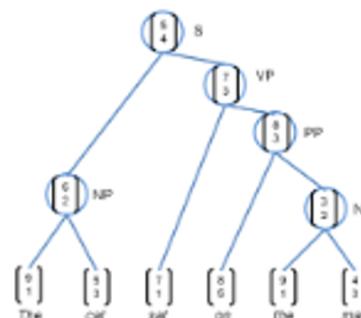
Shared recursive structure



$$\text{score} = U^\top p$$

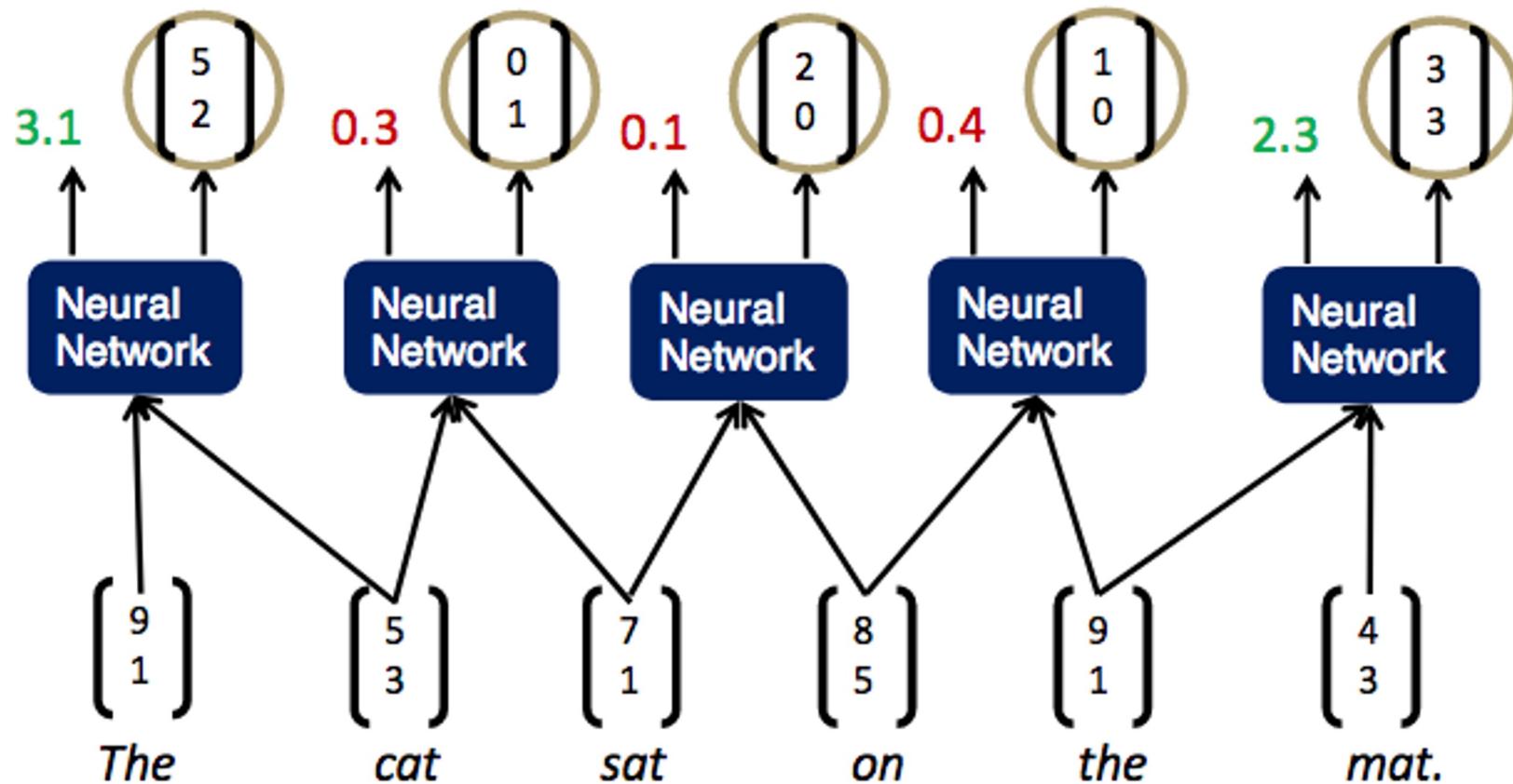
$$p = \tanh\left(w \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} + b\right),$$

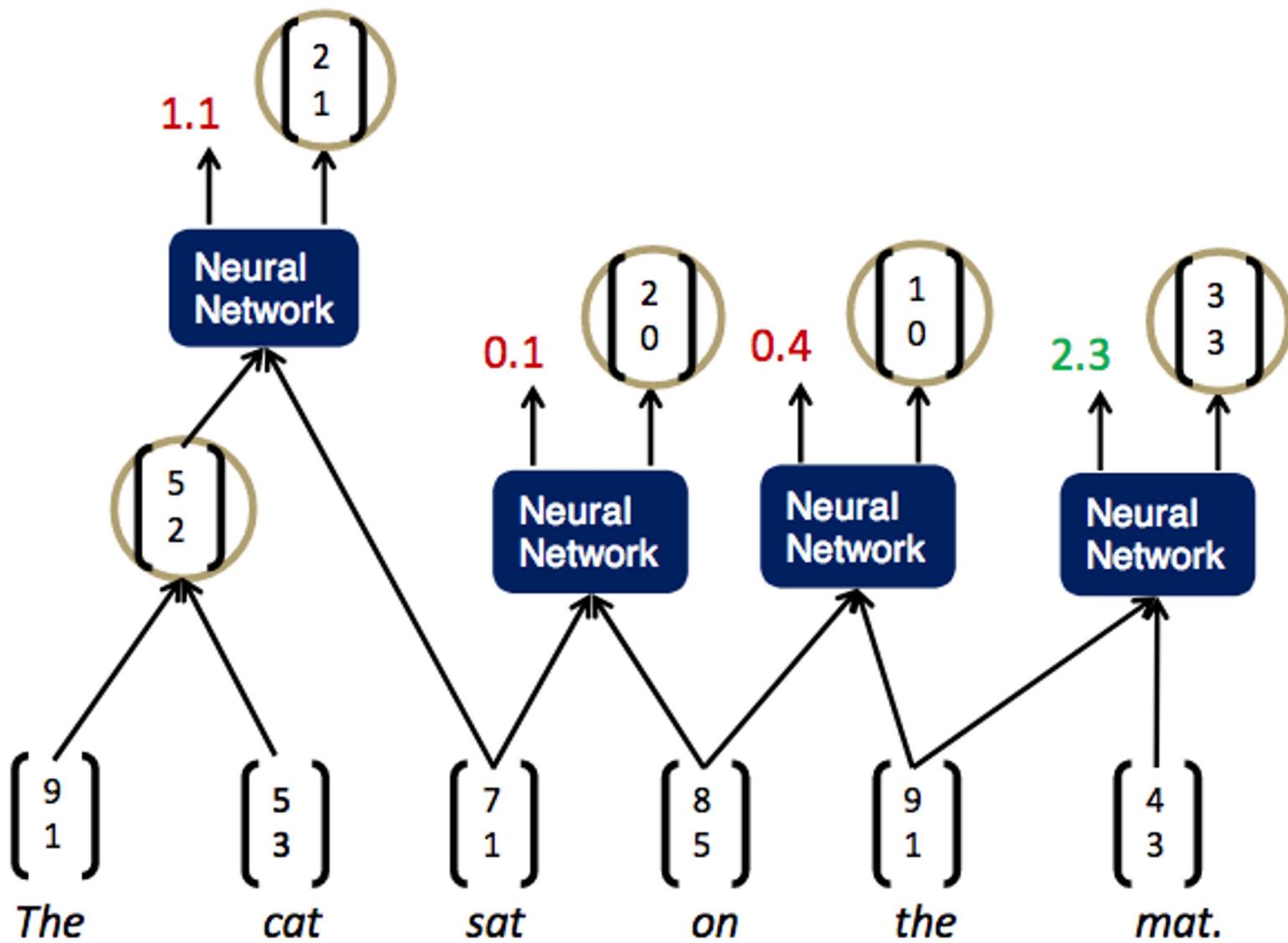
**Same W parameters at all nodes
of the tree**

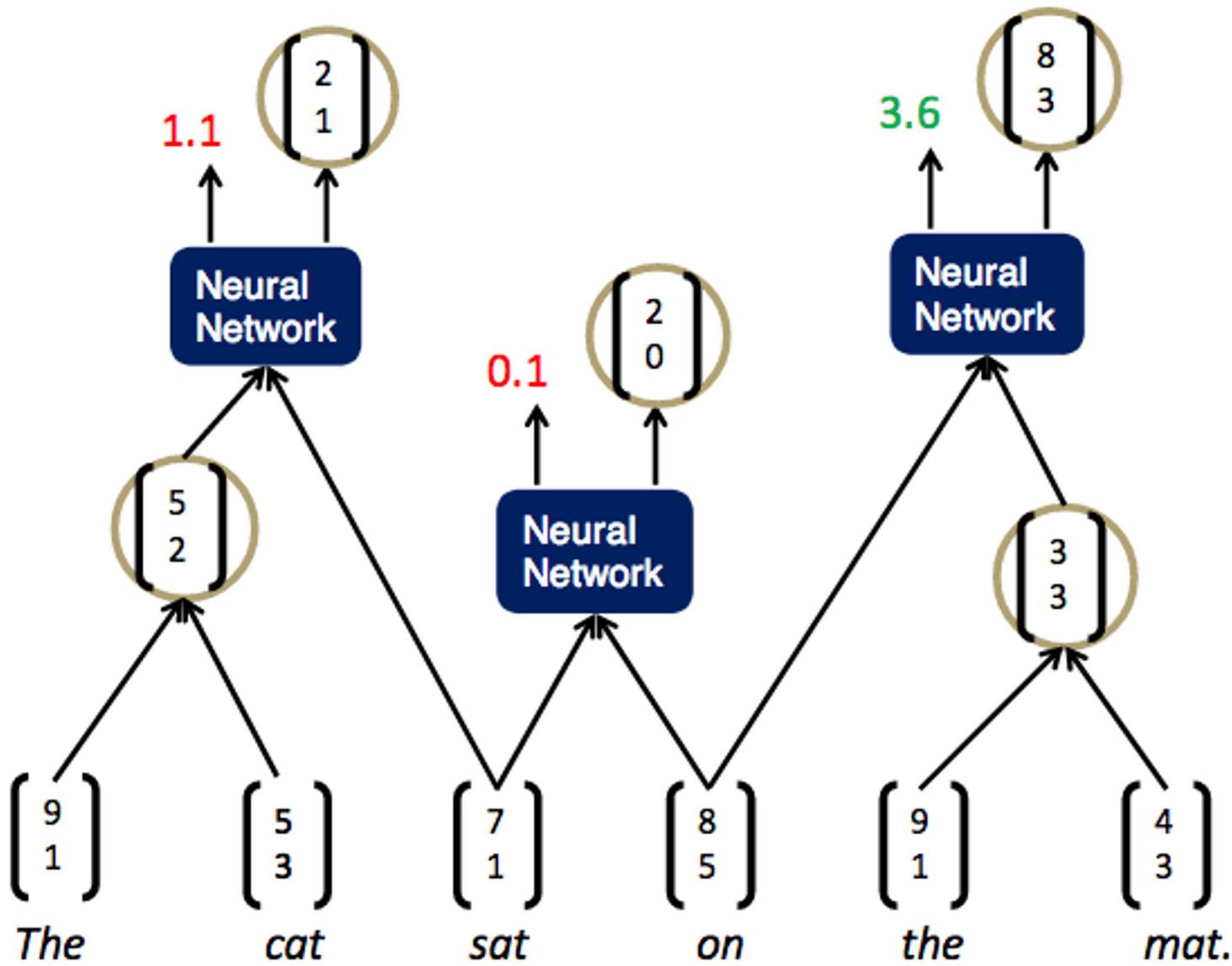


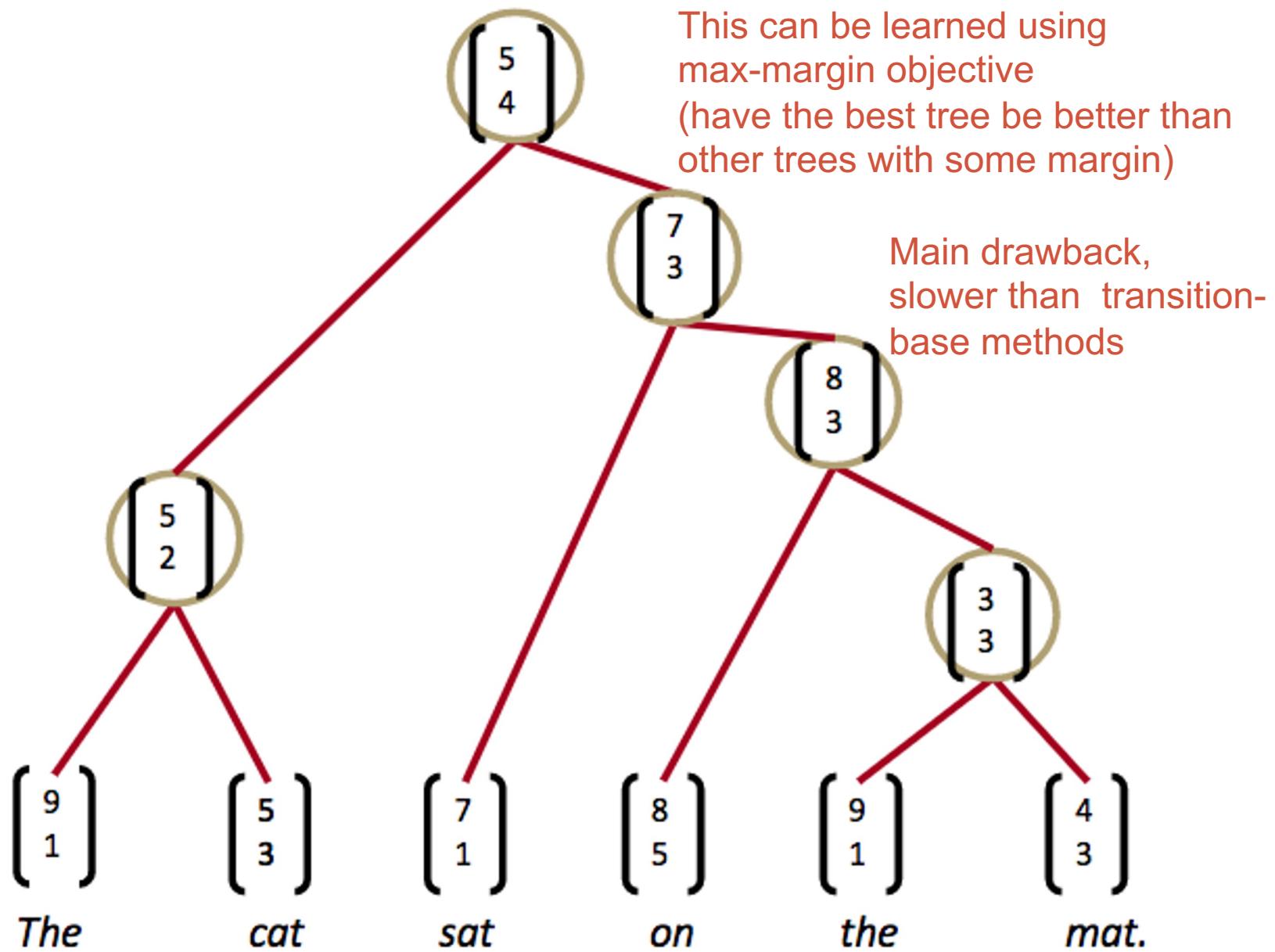
Parsing with recursive networks

Follow the highest scoring pair

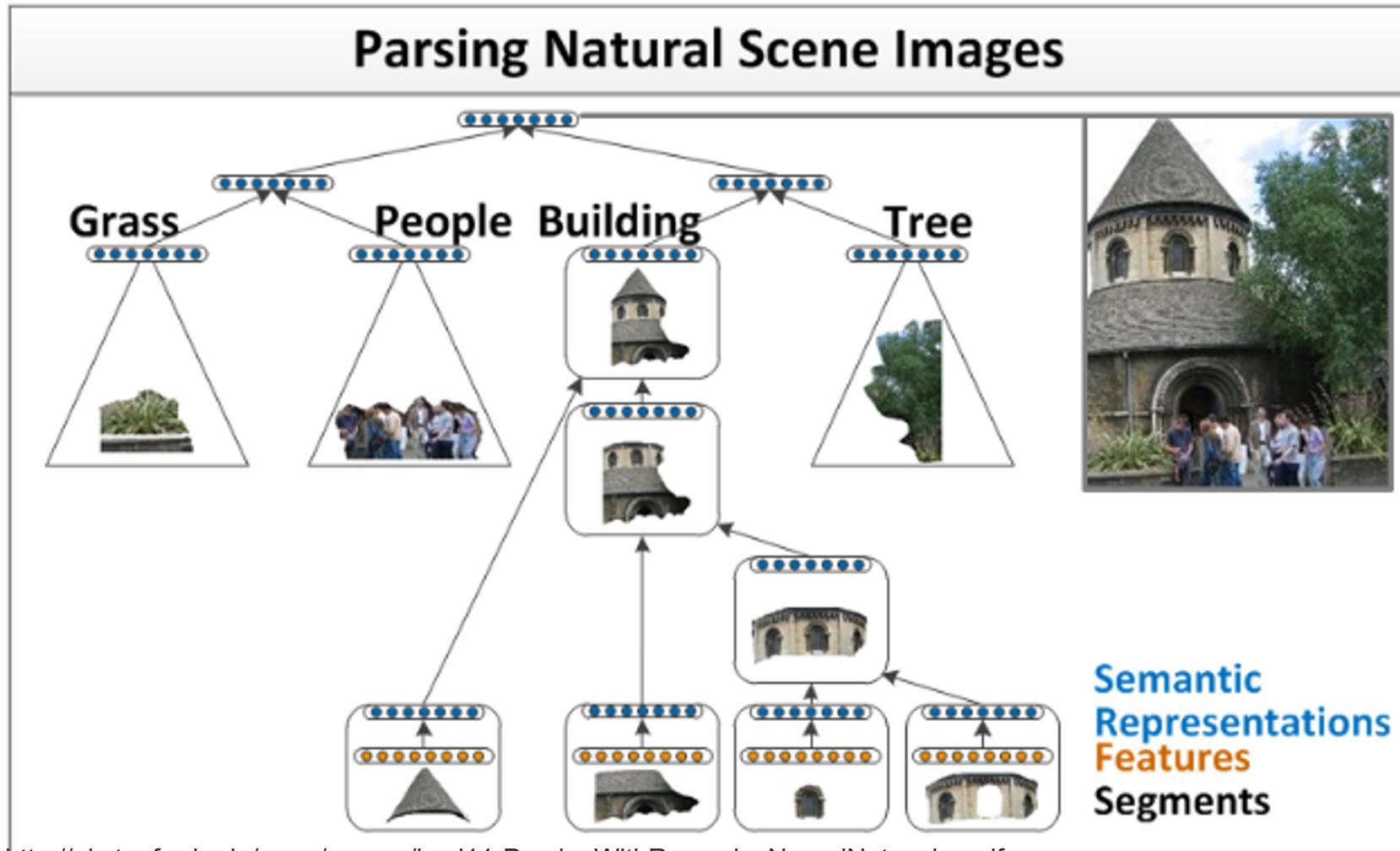








Recursive neural networks for scene parsing



Resources for parsing

- Orchid & corpuses by NECTEC
- <http://universaldependencies.org/>

▶	Flag	Language	Count	Dependencies	Properties	
▶		Russian	3	1,226K		IE, Slavic
▶		Sanskrit	1	1K		IE, Indic
▶		Serbian	1	86K		IE, Slavic
▶		Slovak	1	106K		IE, Slavic
▶		Slovenian	2	170K		IE, Slavic
▶		Spanish	3	1,004K		IE, Romance
▶		Swedish	3	195K		IE, Germanic
▶		Swedish Sign Language	1	1K		Sign Language
▶		Tamil	1	9K		Dravidian, Southern
▶		Telugu	1	6K		Dravidian, South Central
▶		Thai	1	23K		Tai-Kadai
▶		Turkish	2	74K		Turkic, Southwestern
▶		Ukrainian	1	100K		IE, Slavic
▶		Upper Sorbian	1	10K		IE, Slavic
▶		Urdu	1	138K		IE, Indic
▶		Uyghur	1	15K		Turkic, Southeastern
▶		Vietnamese	1	43K		Austro-Asiatic

Parsing corpus on LST20

A Little Surprise for NLP Community

Thai Link Parser (Prachya Boonkwan)	Thai Link Parser (Prachya Boonkwan)	Thai Link Parser (Prachya Boonkwan)
<p>C Thai Parser</p> <p>parser for Thai language constituent and dependency structures, as well as semantic roles dependency structures close to Universal Dependencies Project (>80% TEDVAL score) if more than 110,000 words and proper names Dictionary (2011 ed.) and LST20 running fast, memory-saving, suitable for low-systems, and robust to heavy loads (on web API) AI-for-Thai (coming soon) and on ithub.com/kaamanita/link-grammar</p>	<p>NEC Thai Parser</p> <p>a member of รัฐวิยากรศนกภาษาไทยในระดับประเทศไทย กรุํงรัตนวัสดุ (constituency) ความสัมพันธ์ของคำ (sy) และหน้าที่ภายในภาษาบันทึกคร่าวๆ ได้แบบ unlabeled tree edit distance (TED-ก 80%) โดยเก็บข้อมูลต้นฉบับไว้ยังรูป Universal Dependencies Project พิเศษและซื่อสัมพันธ์กับ ภาษาพจนานุกรมฉบับคิด 554 และลักษณะของ LST20 มากกว่า 110,000 คำ ใช้หน่วยความจำบ้อย เหมาะสำหรับระบบเล็ก ก่อให้รับเปิดตัวบน web API ได้มาบาน AI-for-Thai (เร็วๆ นี้) และดาวน์โหลดได้แล้วที่ ithub.com/kaamanita/link-grammar</p>	<p>NEle of Parse Tree</p> <p>a member of ตัวอย่างเช่น cost vector = (UNUSED=0 DIS= 1.00 LEN=13) +----->AV/pr----->+ <- S<->O >+ ->+ ->+ >+>PO>+ >+>PO>+ ต.ก. เป็น.v สถาบัน.n วิจัย.g เมือง.pnn ชาติ.k ใน.pan ด้าน.a ()</p> <p>รัฐ.v กpn ชาติ.k)) ก ติสือภารกิจ.j)) /github.com/kaamanita/link-grammar</p>

<https://web.facebook.com/dancearmy/posts/pfbid02GpMzUCgwqX3g87wVyioMYDGTVwnH8ZmFCVFRQKnC8erpqbbNqUps9AXFidVqA8cnl>

NECTEC Thai Parser

- Syntactic analyzer for Thai language
- Producing constituent and dependency structures, as well as general semantic roles
- Yielding dependency structures close to Universal Dependencies Project (>80% TEDEVAL score)
- Consisting of more than 110,000 words and proper names from ROYIN Dictionary (2011 ed.) and LST20
- Features: lightning fast, memory-saving, suitable for low-resourced systems, and robust to heavy loads (on web API)
- Available on AI-for-Thai (coming soon) and on <https://github.com/kaamanita/link-grammar>



- CC-BY License
- Easy-to-use interface with C/C++ and Python
- Full documentation for Thai Link Grammar coming soon

Example of Parse Tree



Linkage 2, cost vector = (UNUSED=0 DIS= 1.00 LEN=13)

(S (NP เนคเทค.ก)

ເປັນ. V

(NP สถาบัน.ก วิจัย.ก)

(PP แห่ง. pnn)

(N)

(NP ด้าน. ก อิเล็กทรอนิกส์. ง))

<https://github.com/kaamanita/link-grammar>

Computer Science > Computation and Language

[Submitted on 12 Aug 2020]

The Annotation Guideline of LST20 Corpus

Prachya Boonkwan, Vorapon Luantangsrusuk, Sitthaa Phaholphinyo, Kanyanat Kriengket, Dhanon Leenoi, Charun Phrombut, Monthika Boriboon, Krit Kosawat, Thepchai Supnithi

This report presents the annotation guideline for LST20, a large-scale corpus with multiple layers of linguistic annotation for Thai language processing. Our guideline consists of five layers of linguistic annotation: word segmentation, POS tagging, named entities, clause boundaries, and sentence boundaries. The dataset complies to the CoNLL-2003-style format for ease of use. LST20 Corpus offers five layers of linguistic annotation as aforementioned. At a large scale, it consists of 3,164,864 words, 288,020 named entities, 248,962 clauses, and 74,180 sentences, while it is annotated with 16 distinct POS tags. All 3,745 documents are also annotated with 15 news genres. Regarding its sheer size, this dataset is considered large enough for developing joint neural models for NLP. With the existence of this publicly available corpus, Thai has become a linguistically rich language for the first time.

Comments: 28 pages, 3 tables, 2 figures

Subjects: Computation and Language (cs.CL)

MSC classes: 68T50

ACM classes: I.2.7

Cite as: arXiv:2008.05055 [cs.CL]

(or arXiv:2008.05055v1 [cs.CL] for this version)

<https://doi.org/10.48550/arXiv.2008.05055>

Download:

- PDF only



Current browse context:
cs.CL

< prev | next >
new | recent | 2008

Change to browse by:
cs

References & Citations

- NASA ADS
- Google Scholar
- Semantic Scholar

DBLP – CS Bibliography
[listing](#) | [bibtex](#)

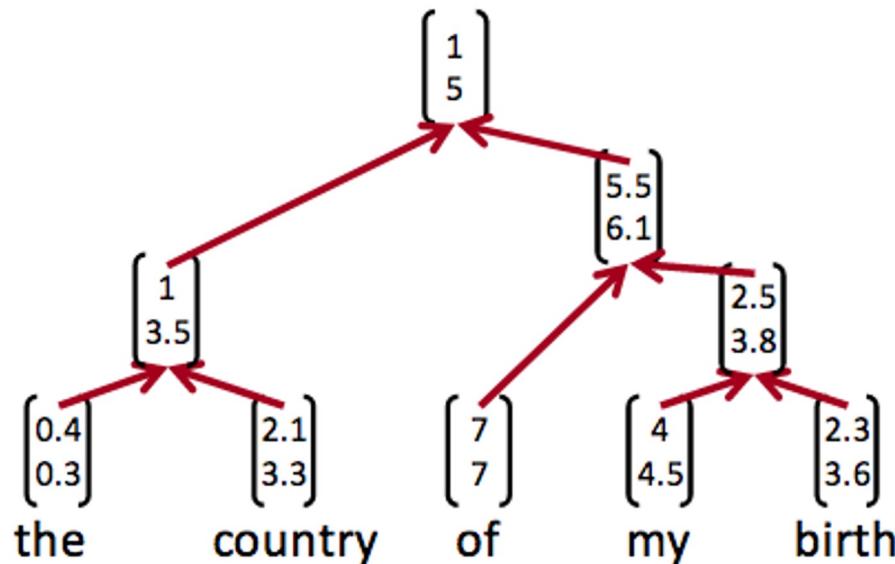
[Export Bibtex Citation](#)

Bookmark

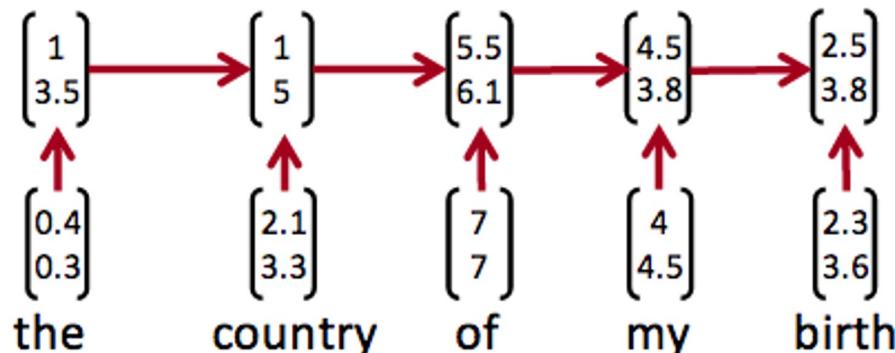


<https://arxiv.org/abs/2008.05055>

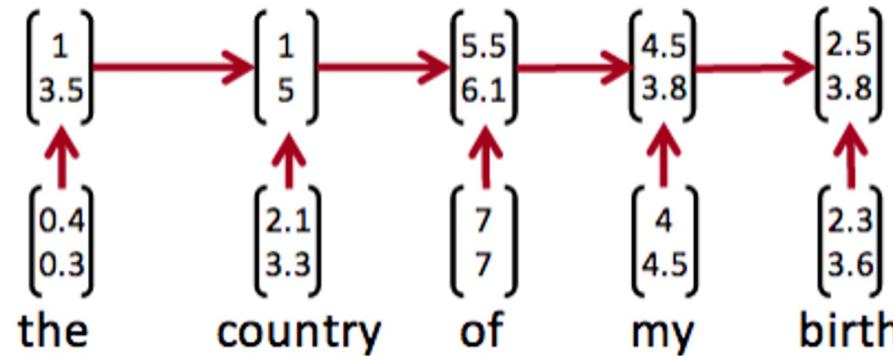
Recursive vs Recurrent representation



resurgence of recursive idea in newer BERTs



Towards unsupervised dependency parsing

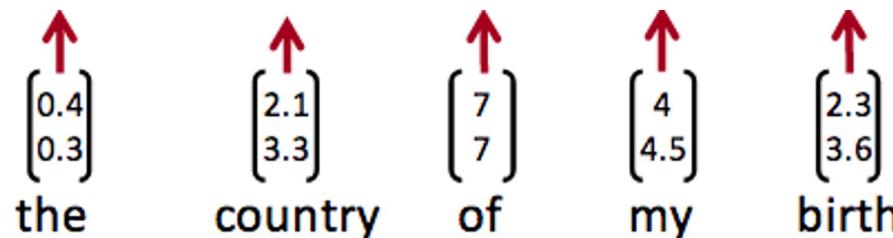


Towards unsupervised dependency parsing

Neural network that decides structure

$$\begin{bmatrix} 2.5 \\ 3.8 \end{bmatrix}$$

What is this piece?
Next lecture! (attention modeling)



A top down end-to-end approach

- Input: text
- Output: some task, sentiment analysis score
- Automatically gets parse tree (without any treebank corpus)

Example of generated parse tree



Three men drink at a reflective bar



Workers at Basking Robbins are filling orders



Three men are socializing during happy hour



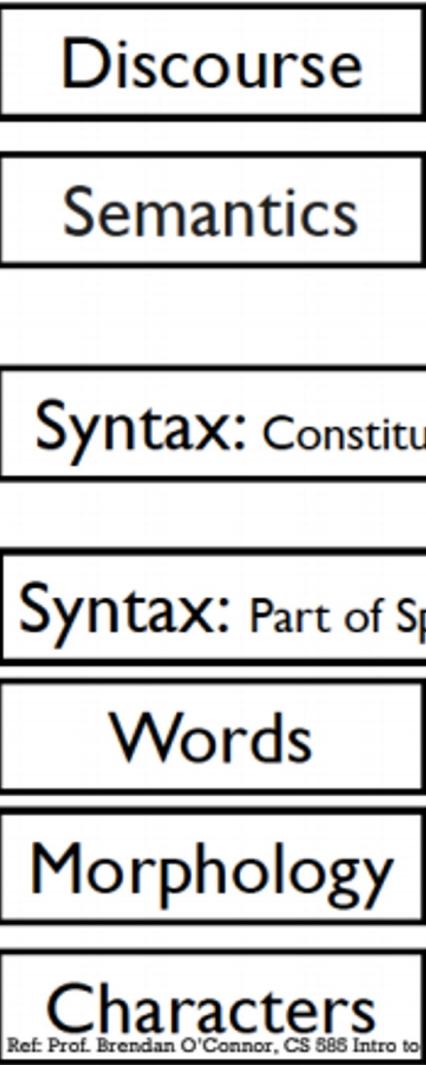
Workers filling orders at Basking Robbins

<https://arxiv.org/pdf/1705.09207.pdf>

Can also use recursive neural networks to learn unsupervised parse trees

Example <https://arxiv.org/pdf/1707.02786.pdf>

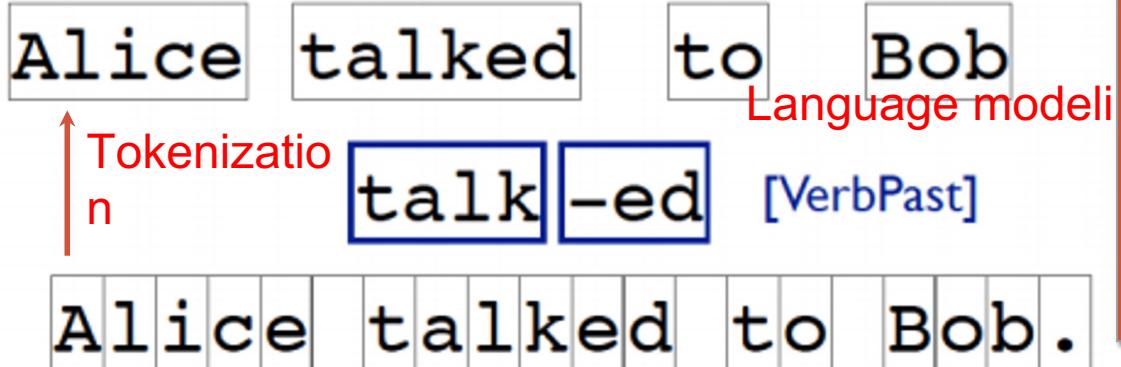
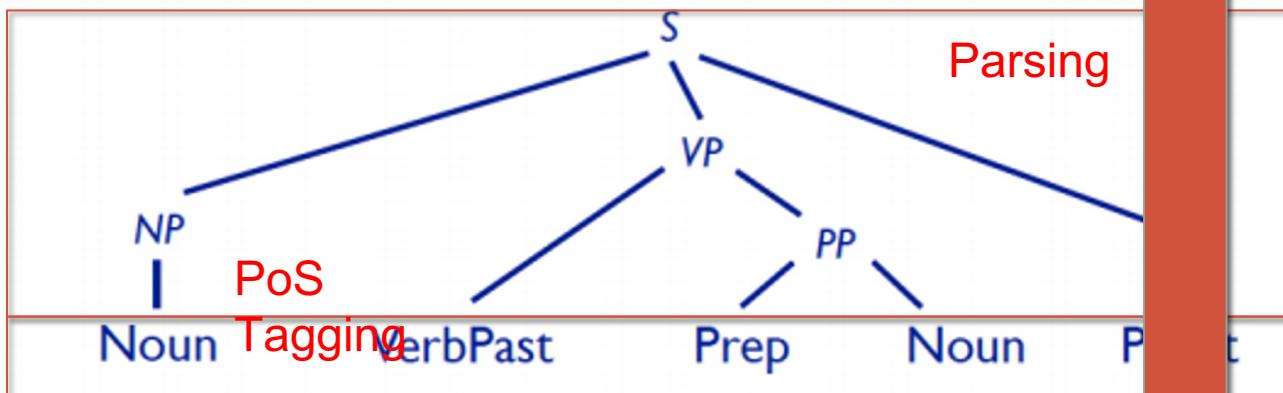
Typical flow



Word embeddings

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)

SpeakerContext(s)
TemporalBefore(e, s)



Summary

- Parsing
- Types of grammars
 - Context Free Grammar
 - Probabilistic Context Free Grammar (Constituency parsers)
 - CYK parser
 - Dependency Grammar (Dependency parsers)
 - Transition-based parsing
 - Recursive Neural networks