

# Foundations of Natural Language Processing

## Lecture 13

### Heads, Dependency parsing

Shay Cohen

(slides from Henry Thompson, Alex Lascarides, Henry Thompson, Nathan Schneider and Sharon Goldwat

3 March 2020

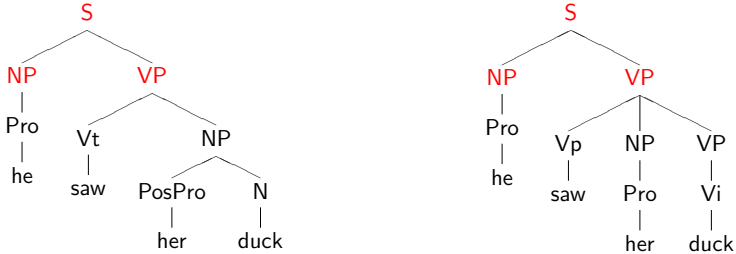


- Provide metrics for evaluating a parser
- Return to the problem of PCFGs
- Suggest a fix
- This fix leads to an approach without constituent structure!

Dependency parsing

### Evaluating parse accuracy

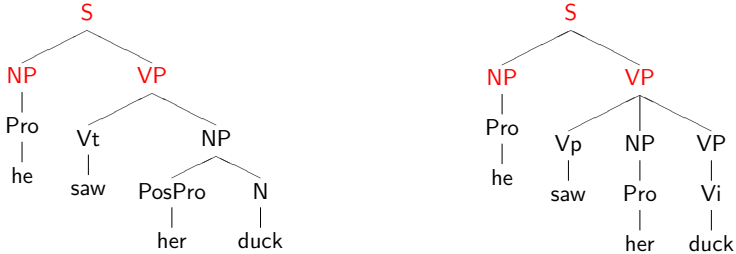
Compare gold standard tree (left) to parser output (right):



- Output constituent is counted **correct** if there is a gold constituent that spans the same sentence positions.
- Harsher measure: also require the constituent labels to match.
- **Pre-terminals** (lexical categories) don't count as constituents.

### Evaluating parse accuracy

Compare gold standard tree (left) to parser output (right):



- **Precision:** (# correct constituents)/(# in parser output) = 3/5
- **Recall:** (# correct constituents)/(# in gold standard) = 3/4
- **F-score:** balances precision/recall:  $2pr/(p+r)$

## Parsing accuracies

F-scores for parsing on WSJ corpus:

- vanilla PCFG: < 80%<sup>1</sup>
- lexicalizing + cat-splitting: 89.5% (Charniak, 2000)
- Best current parsers get about 92%
- Numbers get better if we look at top 5 or top 10

However, results on other corpora and other languages are considerably lower. Definitely not a solved problem!

<sup>1</sup>Charniak (1996) reports 81% but using gold POS tags as input.

## Recall Problem with Vanilla PCFGs No lexical dependencies

Replacing one word with another with the same POS will never result in a different parsing decision, even though it should!

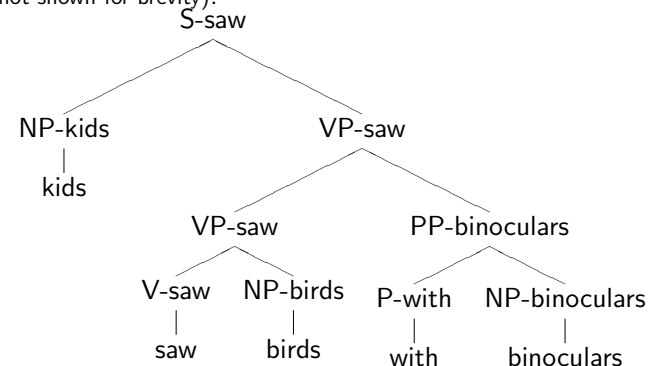
- kids saw birds with fish vs.  
kids saw birds with binoculars
- She stood by the door covered in tears vs.  
She stood by the door covered in ivy
- stray cats and dogs vs.  
Siamese cats and dogs

## Summary

- Probabilistic models of syntax can help disambiguation and speed in broad-coverage parsing.
  - by computing the probabilities of each tree or sub-tree as the product of the rules in it, and choosing the best option(s).
- Treebanks provide training data for estimating rule probabilities.
- However, to do well, we need to be clever:
  - Standard categories in the treebank don't capture some important facts about language.
  - By creating more detailed categories, we can encode more information within the PCFG framework.

## A way to fix PCFGs: lexicalization

Create new categories, this time by adding the **lexical head** of the phrase (note: N level under NPs not shown for brevity):



- Now consider:  
VP-saw → VP-saw PP-fish vs. VP-saw → VP-saw PP-binoculars

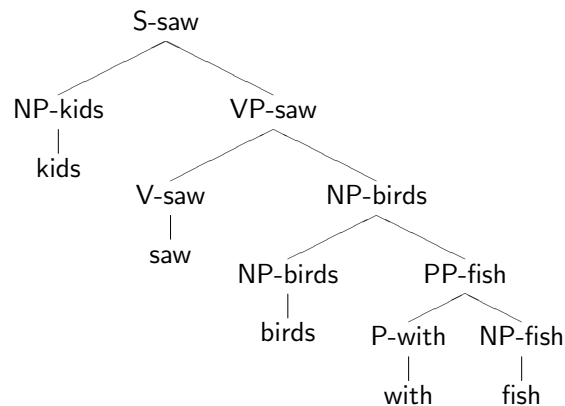
## Practical issues

- Identifying the head of every rule is not always straightforward
  - (more on this below)
- All this category-splitting makes the grammar much more **specific** (good!)
- But leads to huge grammar blowup and very sparse data (bad!)
- Lots of effort on how to balance these two issues.
  - Complex smoothing schemes (similar to N-gram interpolation/backoff).
  - More recently, increasing emphasis on automatically learned subcategories.
- But do we really need phrase structure in the first place? Not always!
- Today: Syntax (and parsing) without constituent structure.

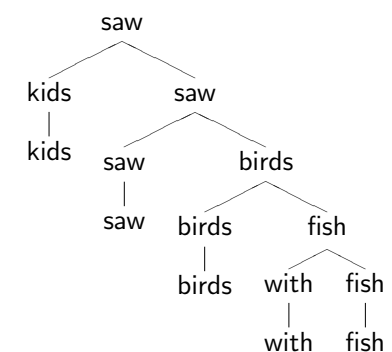
## Outline

1. Dependencies: what/why
2. Transforming constituency → dependency parse
3. Direct dependency parsing
  - Transition-based (shift-reduce)
  - Graph-based

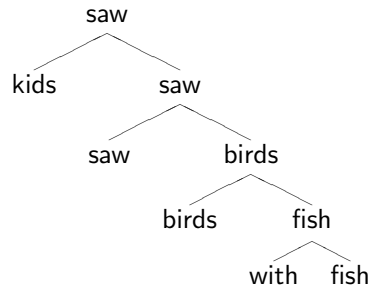
## Lexicalized Constituency Parse



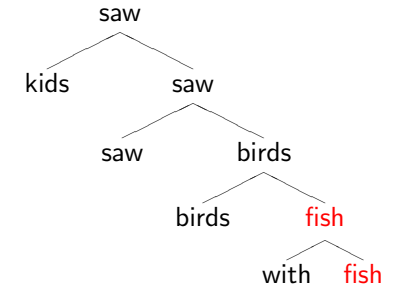
## . . . remove the phrasal categories. . .



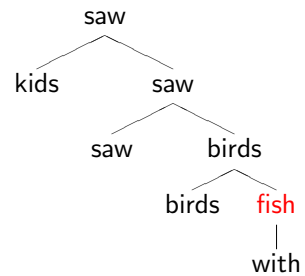
... remove the (duplicated) terminals...



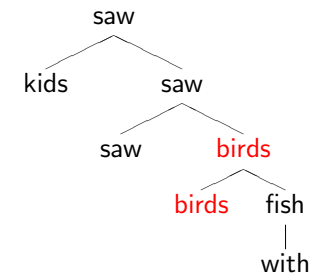
... and collapse chains of duplicates...



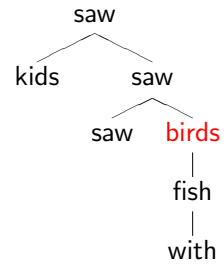
... and collapse chains of duplicates...



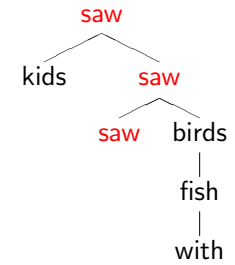
... and collapse chains of duplicates...



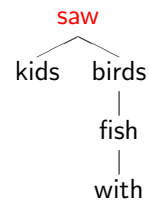
## . . . and collapse chains of duplicates. . .



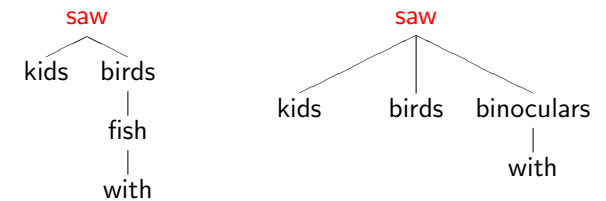
## . . . and collapse chains of duplicates. . .



## . . . and collapse chains of duplicates. . .



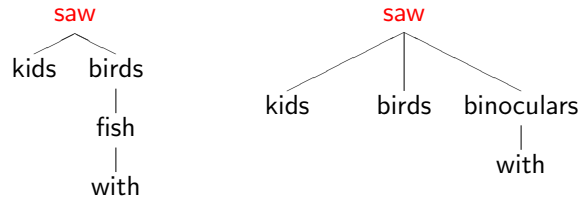
## Dependency Parse



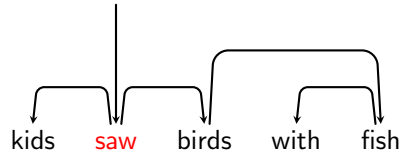
Linguists have long observed that the meanings of words within a sentence depend on one another, mostly in *asymmetric*, *binary* relations.

- Though some constructions don't cleanly fit this pattern: e.g., coordination, relative clauses.

## Dependency Parse



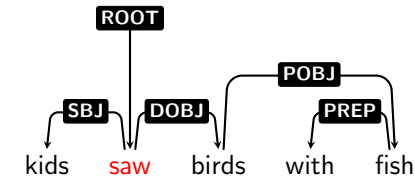
Equivalently, but showing word order (head → modifier):



Because it is a tree, every word has exactly one parent.

## Edge Labels

It is often useful to distinguish different kinds of head → modifier **relations**, by labeling edges:

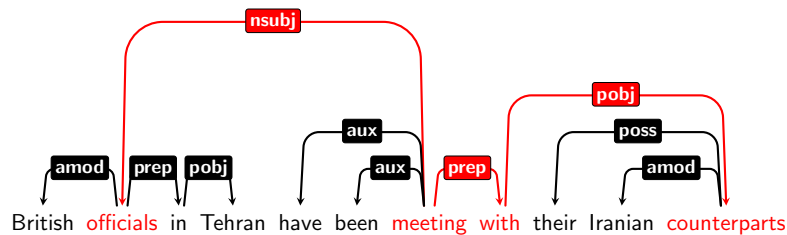


Important relations for English include **subject**, **direct object**, **determiner**, **adjective modifier**, **adverbial modifier**, etc. (Different treebanks use somewhat different label sets.)

- How would you identify the subject in a constituency parse?

## Dependency Paths

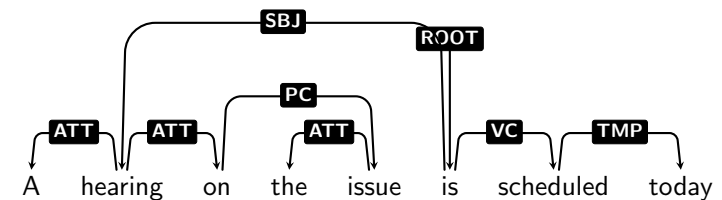
For **information extraction** tasks involving real-world relationships between entities, chains of dependencies can provide good features:



(example from Brendan O'Connor)

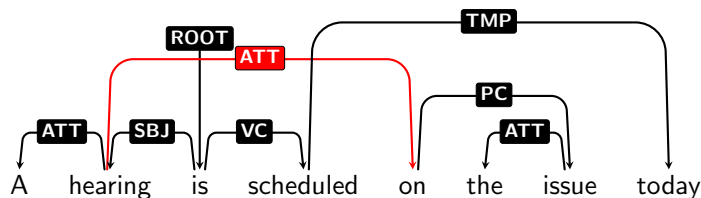
## Projectivity

- A sentence's dependency parse is said to be **projective** if every subtree (node and all its descendants) occupies a *contiguous span* of the sentence.
- = The dependency parse can be drawn on top of the sentence without any crossing edges.



# Nonprojectivity

- Other sentences are **nonprojective**:



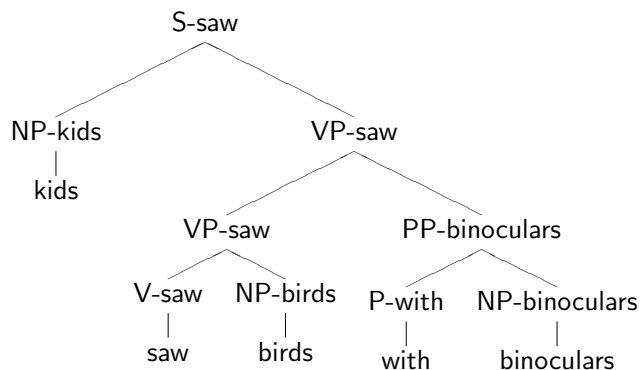
- Nonprojectivity is rare in English, but quite common in many languages.

# Outline

1. Dependencies: what/why
2. **Transforming constituency** → **dependency parse**
3. Direct dependency parsing
  - Transition-based (shift-reduce)
  - Graph-based

## Constituency Tree → Dependency Tree

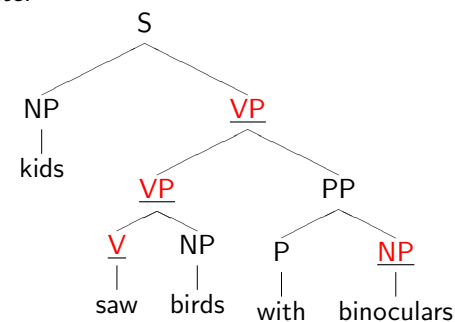
We saw how the **lexical head** of the phrase can be used to collapse down to a dependency tree:



- But how can we find each phrase's head in the first place?

## Head Rules

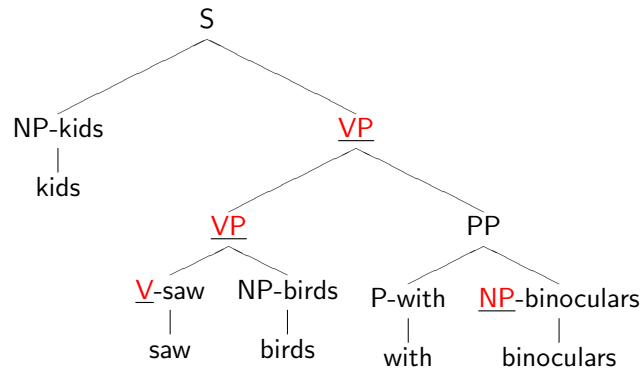
The standard solution is to use **head rules**: for every non-unary (P)CFG production, designate one RHS nonterminal as containing the head.  $S \rightarrow NP \underline{VP}$ ,  $VP \rightarrow \underline{V} PP$ ,  $PP \rightarrow P \underline{NP}$ , etc.



- Heuristics to scale this to large grammars: e.g., within an **NP**, last immediate **N** child is the head.

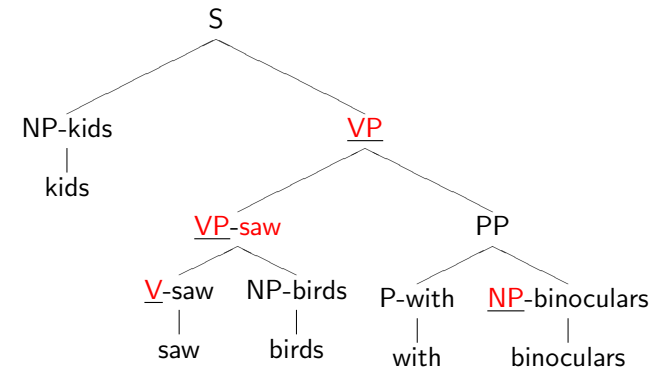
## Head Rules

Then, propagate heads up the tree:



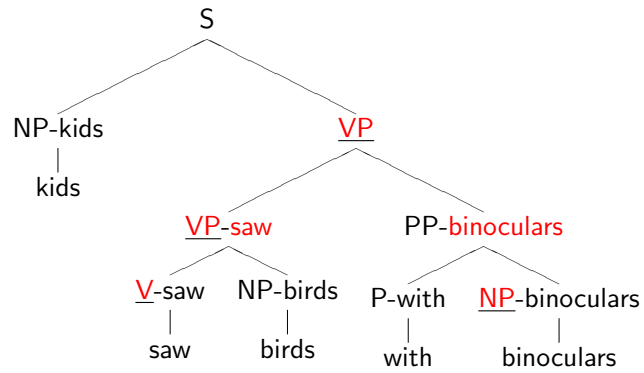
## Head Rules

Then, propagate heads up the tree:



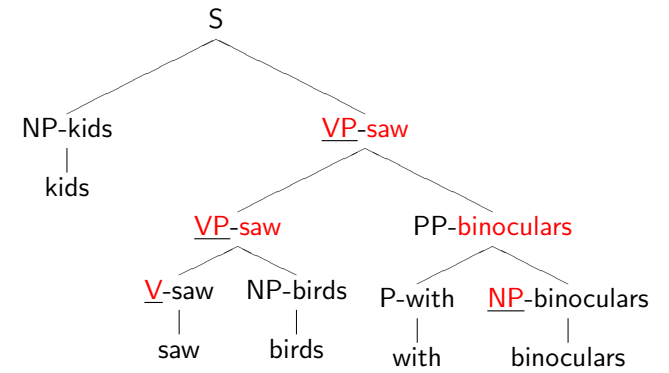
## Head Rules

Then, propagate heads up the tree:



## Head Rules

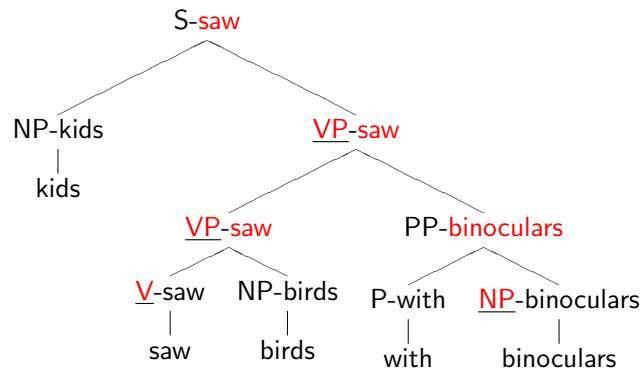
Then, propagate heads up the tree:





## Head Rules

Then, propagate heads up the tree:



## Dependency Parsing

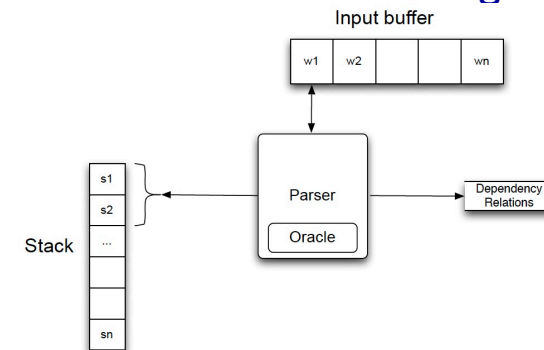
Some of the algorithms you have seen for PCFGs can be adapted to dependency parsing.

- **CKY** can be adapted, though efficiency is a concern: obvious approach is  $O(Gn^5)$ ; Eisner algorithm brings it down to  $O(Gn^3)$ 
  - N. Smith's slides explaining the Eisner algorithm: <http://courses.cs.washington.edu/courses/cse517/16wi/slides/an-dep-slides.pdf>
- **Shift-reduce**: more efficient, doesn't even require a grammar!

## Outline

1. Dependencies: what/why
2. Transforming constituency → dependency parse
3. **Direct dependency parsing**
  - Transition-based (shift-reduce)
  - Graph-based

## Transition-based Parsing: Shift Reduce Parser



3 possible actions:

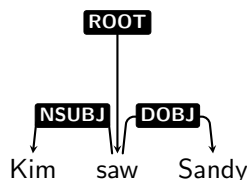
- LeftArc**: Assign head-dependent relation between  $s1$  and  $s2$ ; pop  $s2$
- RightArc**: Assign head-dependent relation between  $s2$  and  $s1$ ; pop  $s1$
- Shift**: Put  $w1$  on top of the stack.

Remember, dependency relation points *from* head *to* dependent

Both **LeftArc** and **RightArc** pop the dependent, leaving the head at the top of the stack

## Example

Step	Stack	Word List	Action	Relations
0	[root]	[Kim,saw,Sandy]		
1	[root, Kim]	[saw,Sandy]	Shift	
2	[root, Kim, saw]	[Sandy]	Shift	
3	[root, saw]	[Sandy]	LeftArc	nsubj(saw, Kim)
4	[root, saw, Sandy]	[]	Shift	
5	[root, saw]	[]	RightArc	dobj(saw, Sandy)
6	[root]	[]	RightArc	root→saw



## Graph-based Parsing

- Global algorithm: From the fully connected directed graph of all possible edges, choose the best ones that form a tree.
- **Edge-factored** models: Classifier assigns a nonnegative score to each possible edge; **maximum spanning tree** algorithm finds the spanning tree with highest total score in  $O(n^2)$  time.
  - Edge-factored assumption can be relaxed (higher-order models score larger units; more expensive).
  - Unlabeled parse  $\rightarrow$  edge-labeling classifier (pipeline).
- Pioneering work: McDonald's MSTPARSER
- Can be formulated as constraint-satisfaction with **integer linear programming** (Martins's TURBOPARSER)

## Transition-based Parsing

- Latent structure is just edges between words. Train a **classifier** as the oracle to predict next action (SHIFT, LEFTARC, or RIGHTARC), and proceed left-to-right through the sentence.  $O(n)$  time complexity!
- Only finds **projective** trees (without special extensions)
- Pioneering system: Nivre's MALTPARSER
- See <http://spark-public.s3.amazonaws.com/nlp/slides/Parsing-Dependency.pdf> (Jurafsky & Manning Coursera slides) for details and examples

## Graph-based vs. Transition-based vs. Conversion-based

- TB: Features in scoring function can look at any part of the stack; no optimality guarantees for search; linear-time; (classically) projective only
- GB: Features in scoring function limited by factorization; optimal search within that model; quadratic-time; no projectivity constraint
- CB: In terms of accuracy, sometimes best to first constituency-parse, then convert to dependencies (e.g., STANFORD PARSER).
  - Slower than direct methods.
  - And, you need a grammar and head rules.

## Choosing a Parser: Criteria

- Target representation: constituency or dependency?
- Efficiency? In practice, both runtime and memory use.
- Incrementality: parse the whole sentence at once, or obtain partial left-to-right analyses/expectations?
- Retractable system?

## Choosing a Parser: Performance

SOTA for English constituency parsing (WSJ §23): 91%–92%  $F_1$

Parser	LR	LP	F1	#Toks/s.
harniak 2000	85.3	85.5	85.4	1.3
lein and anning 2003	85.3	85.5	85.4	1.3
etrov and lein 200	0.0	0.3	0.1	1
arreras et al 2008	0	1	1.1	
hu et al 2013	0.3	0	0	1.2
Stanford Shitka et al 201	88.8	88.8	88.8	12
all et al 201	88.8	88.8	88.8	12
<b>This work</b>	8	0	0.2	5
harniak and Johnson 2005 *	1.2	1.8	1.5	8
Socher et al 2013 *	8	1.8	8	0
hu et al 2013 *	1.1	1.5	1.3	

Table 3: Results on the English §23. All systems re-optimizing runtimes were run on the same machine. Marked as \* are reranking and semi-supervised parsers.

(Fernández-González and Martins, 2015)

## Choosing a Parser: Performance

Constituency parsing in other languages

Parser	Basque	French	German	Hebrew	Hungar.	Korean	Polish	Swedish	Avg.
Erkley	0.50	<b>80.38</b>	83.0	8	81.2	1.2	23	1	8.5
Erkley-agg			82.8	85.2	85.22	8.5	8.5	80	81.1
all et al 201	83.3	0	8.3	8.18	<b>88.25</b>	<b>80.18</b>	0	82.00	83.2
rabbe and Seddah 201	85.35	8	15	8.1	8.51	35	<b>91.60</b>	82.2	83
<b>This work</b>	<b>85.90</b>	8.5	<b>78.66</b>	<b>88.97</b>	88.1	28	1.20	<b>82.80</b>	<b>84.22</b>
orkelund et al 201	88.2	82.53	81	8.80	1.2	83.81	0.50	85.50	8.2

(Fernández-González and Martins, 2015)

## Choosing a Parser: Performance

SOTA for English dependency parsing (WSJ §23):  
93%–94% UAS, 91%–92% LAS (Zhou et al., 2015)

System	UAS	LAS	Speed
baseline greedy parser	91.47	90.43	0.001
Huang and Sagae (2010)	92.10		0.04
Zhang and Nivre (2011)	92.90	91.80	0.03
Choi and McCallum (2013)	92.96	91.93	0.009
Ma et al. (2014)	93.06		
Bohnet and Nivre (2012)†‡	93.67	92.68	0.4
Suzuki et al. (2009)†	93.79		
Koo et al. (2008)†	93.16		
Chen et al. (2014)†	93.77		
beam size			
training	decoding		
100	100	<b>93.28</b>	<b>92.35</b>
100	64	93.20	92.27
100	16	92.40	91.95

Table 5: Results on WSJ. Speed: sentences per second. †: semi-supervised learning. ‡: joint POS-tagging and dependency parsing models.

**Labelled Attachment Score (LAS)** is stricter than **Unlabelled Attachment Score (UAS)**

## Summary

- While constituency parses give hierarchically nested phrases, dependency parses represent syntax with trees whose edges connect words in the sentence. (No abstract phrase categories like NP.) Edges often labeled with relations like [subject](#).
- Head rules govern how a lexicalized constituency grammar can be extracted from a treebank, and how a constituency parse can be converted to a dependency parse.
- For English, it is often fastest and most convenient to parse directly to dependencies. Two main paradigms, graph-based and transition-based, with different kinds of models and search algorithms.
- [Google "online dependency parser"](#).  
[Try out the Stanford parser and SEMAFOR!](#)