

Dependency Parsing

Natural Language Processing

(based on revision of Chris Manning Lectures)



Suggested Readings

1. <https://web.stanford.edu/~jurafsky/slp3/12.pdf>
(context-free grammar)
2. <https://web.stanford.edu/~jurafsky/slp3/13.pdf>
(constituency grammar)
3. <https://web.stanford.edu/~jurafsky/slp3/14.pdf>
(dependency parsing)



Two views of linguistic structure:
Context-free grammar and
Dependency grammar



1. Two views of linguistic structure: Context-free grammar

Phrase structure organizes words into nested constituents

Starting unit: words

the, cat, cuddly, by, door
 det noun adj prep noun

Words combine into phrases

The cuddly cat, by the door
 np: noun phrase pp: prepositional phrase

Phrases can combine into bigger phrases

The cuddly cat by the door
 np np pp np

Noun → flights | breeze | trip | morning
 Verb → is | prefer | like | need | want | fly
 Adjective → cheapest | non-stop | first | latest
 | other | direct
 Pronoun → me | I | you | it
 Proper-Noun → Alaska | Baltimore | Los Angeles
 | Chicago | United | American
 Determiner → the | a | an | this | these | that
 Preposition → from | to | on | near
 Conjunction → and | or | but

Figure 12.2 The lexicon for \mathcal{L}_0 .

We called the word on the right (e.g., the) **terminal symbols**, and the grammar rules as **lexicon**



1. Two views of linguistic structure: Context-free grammar

NP → Det N

e.g., the cat

NP → Det (Adj) N

e.g., the large cat

PP → P NP

e.g., by the door

NP → Det (Adj) N (PP)

e.g., the large cat by the door

NP → Det (Adj)* N (PP)

e.g., the large cute furry cat by the door

VP → V PP

e.g., talk to the cat

S → NP VP

e.g., the cat walked behind the dog

More grammar rules! As much as we want....



1. Two views of linguistic structure: Context-free grammar

Using context-free grammar to parse gives us a **parse tree**

Grammar Rules	Examples
$S \rightarrow NP VP$	I + want a morning flight
$NP \rightarrow$ <i>Pronoun</i>	I
<i>Proper-Noun</i>	Los Angeles
<i>Det Nominal</i>	a + flight
$Nominal \rightarrow$ <i>Nominal Noun</i>	morning + flight
<i>Noun</i>	flights
$VP \rightarrow$ <i>Verb</i>	do
<i>Verb NP</i>	want + a flight
<i>Verb NP PP</i>	leave + Boston + in the morning
<i>Verb PP</i>	leaving + on Thursday
$PP \rightarrow$ <i>Preposition NP</i>	from + Los Angeles

Figure 12.3 The grammar for \mathcal{L}_0 , with example phrases for each rule.

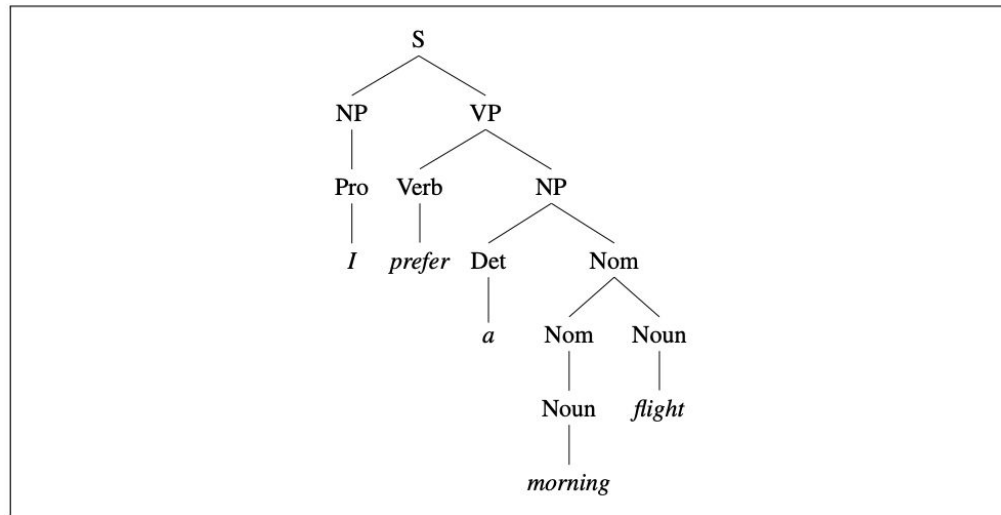


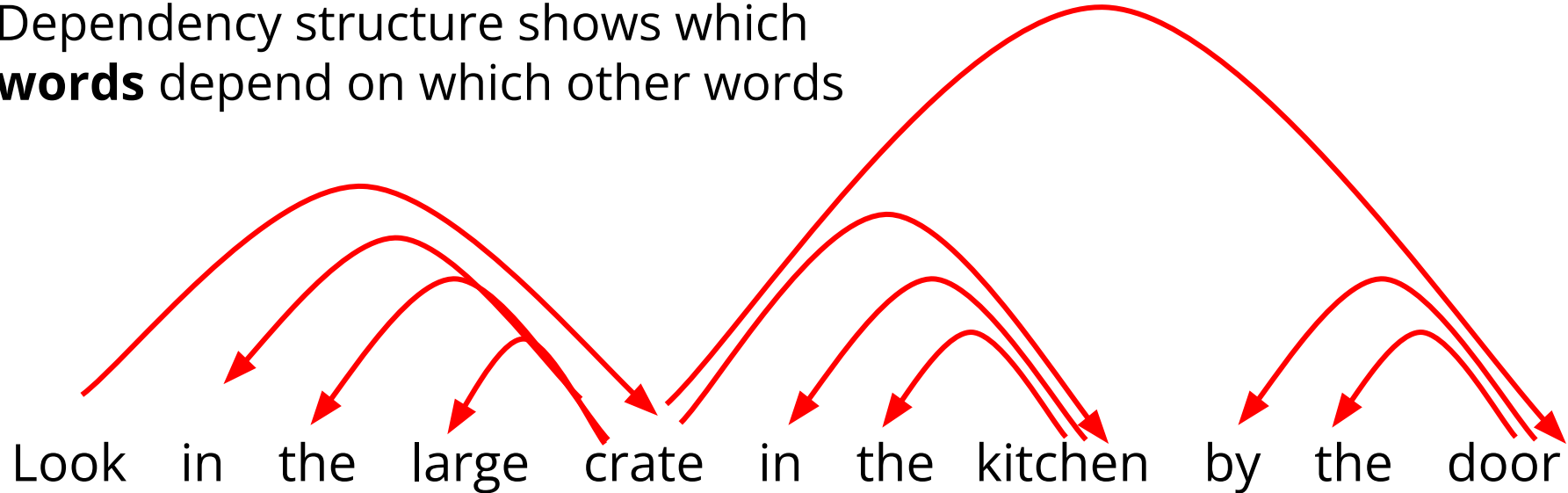
Figure 12.4 The parse tree for “I prefer a morning flight” according to grammar \mathcal{L}_0 .

CFG is useful for **grammar checking** (check whether the sentence deviates from the lexicon), **entities extraction** (check common pattern), **sentence classification** (e.g., question has certain pattern)



2. Two views of linguistic structure: Dependency grammar

Dependency structure shows which **words** depend on which other words



E.g., crate "depends on" the word "large". Here "crate" (start of the arrow) is called "head" and "large" (where the arrow point to) is called the "dependent"



2. Two views of linguistic structure: Dependency grammar

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 Some of the Universal Dependency relations (de Marneffe et al., 2014).



Prepositional phrase attachment ambiguity

San Jose cops kill man with knife

Text Paper Translate Listen

sub_j obl:oblique obj

BBC Sign in News Sport Weather Shop Reel Travel

NEWS

Home Video World US & Canada UK Business Tech Science Stories

Science & Environment

Scientists count whales from space

By Jonathan Amos
BBC Science Correspondent

San Jose cops kill man with knife

Text Paper Translate Listen

sub_j nmod: noun modifier obj

BBC Sign in News Sport Weather Shop Reel Travel

NEWS

Home Video World US & Canada UK Business Tech Science Stories

Science & Environment

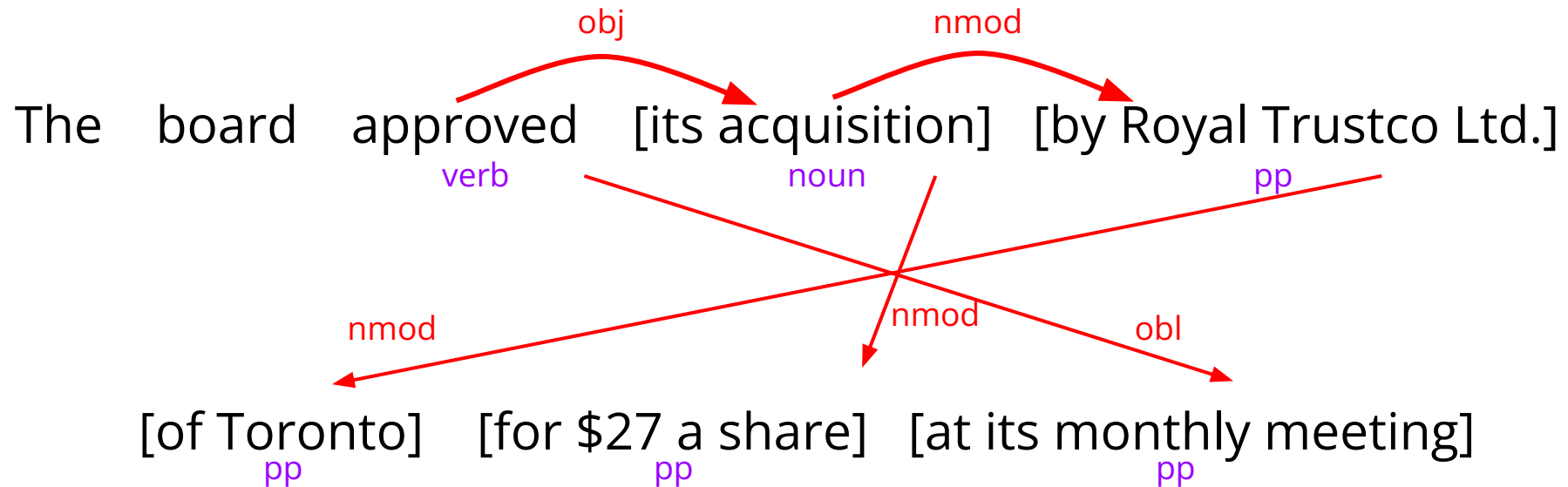
Scientists count whales from space

By Jonathan Amos
BBC Science Correspondent

<https://universaldependencies.org/u/dep/>



How to link?

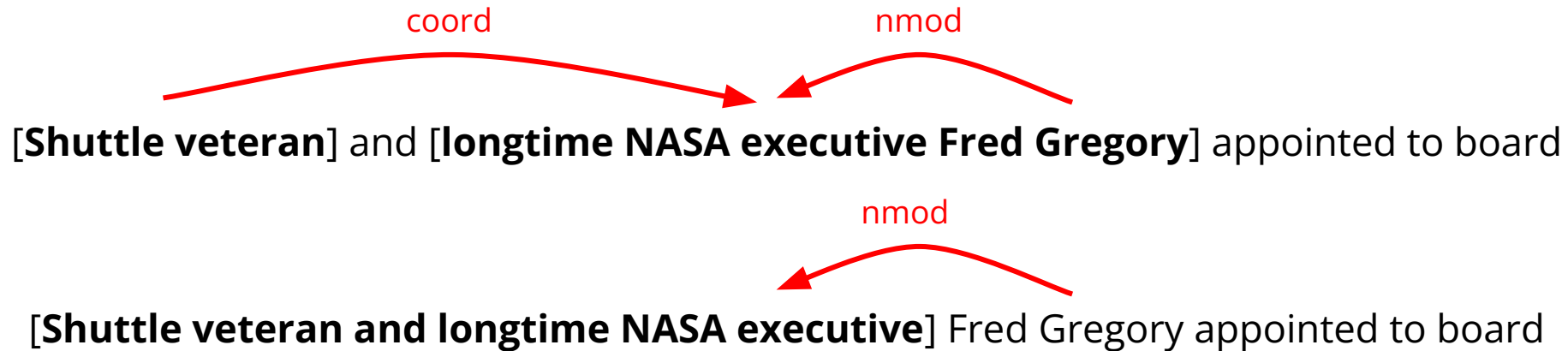


The number of possible parses grows exponentially w.r.t. to the number of prepositional phrases (here $n = 4$) (i.e., Catalan numbers)

<https://universaldependencies.org/u/dep/>



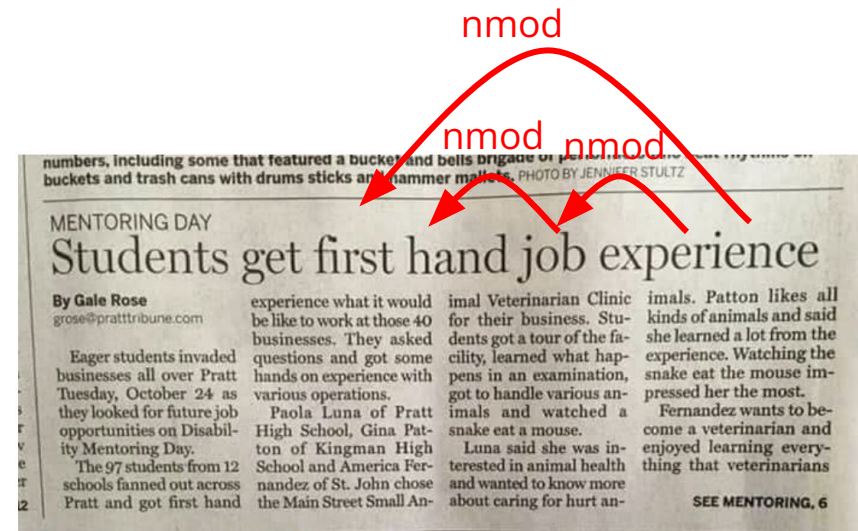
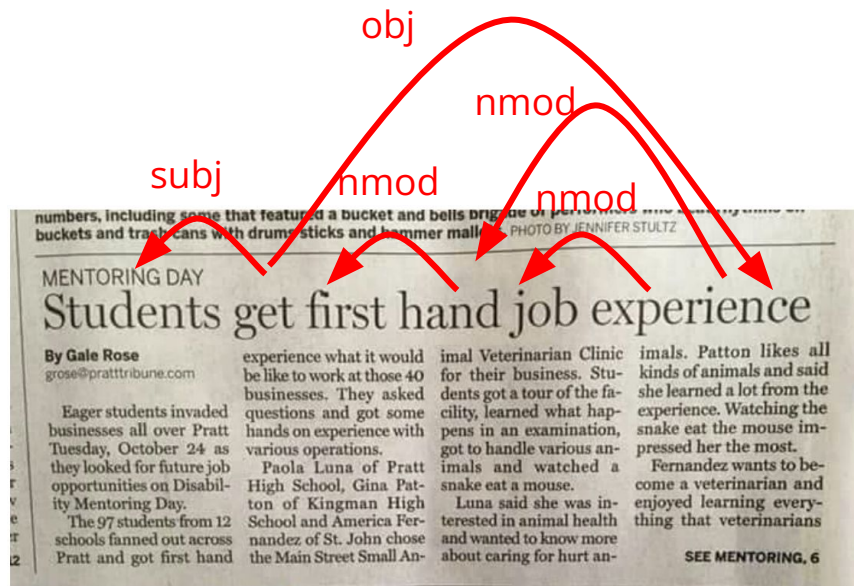
Coordination scope ambiguity



<https://universaldependencies.org/u/dep/>



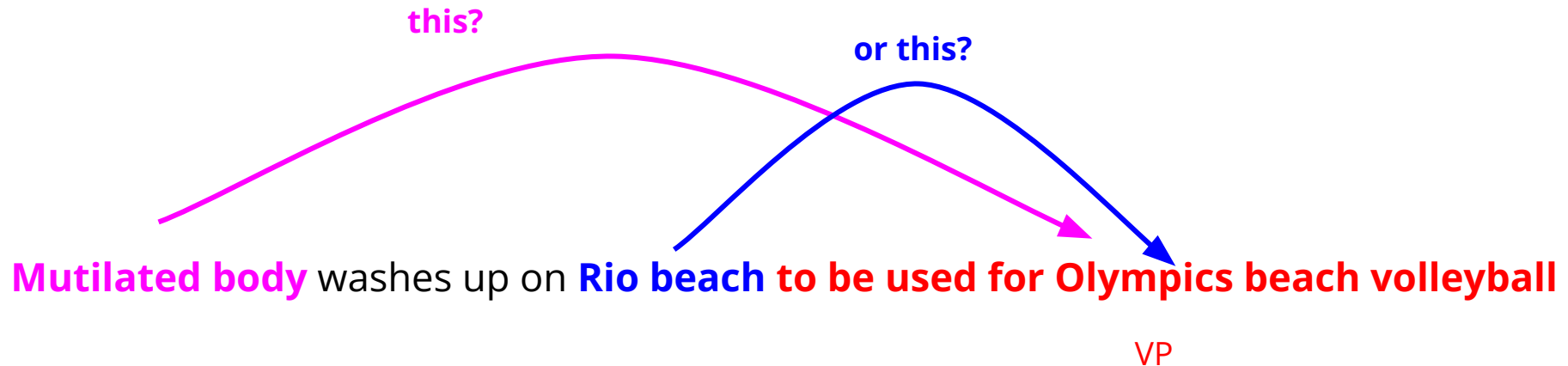
Adjectival/Adverbial Modifier ambiguity



<https://universaldependencies.org/u/dep/>



Verb phrase (VP) attachment ambiguity



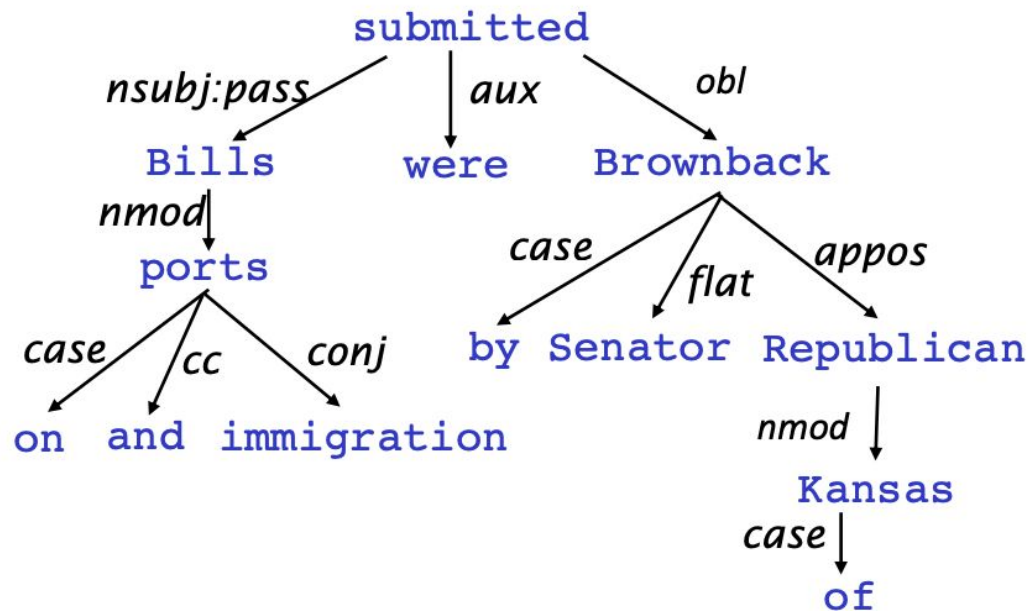
<https://universaldependencies.org/u/dep/>



Dependency as tree

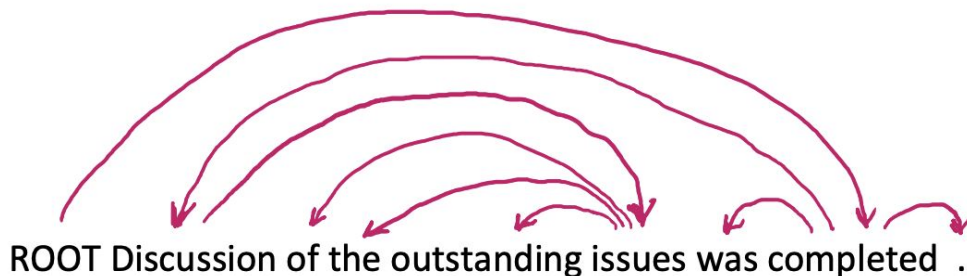
Dependencies can be redrawn as tree (motivated by computer science)

Dependencies usually form a **connected, acyclic, single-root** graph.



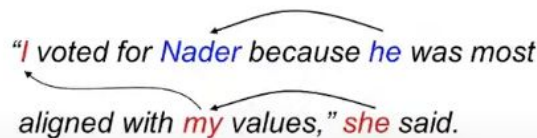
Dependency as tree

- Some people draw the arrows one way; some the other way!
 - **Tesniere (1959)** had them point from head to dependent - we follow that convention
- We usually add a **fake ROOT** so every word is a dependent of precisely 1 other node



Why we need to learn sentence structure?

- Humans need to work out what *depends* what
 - Thus, similarly, a model needs to understand sentence structure in order to interpret language correctly
 - E.g., **conference resolution**



"I voted for Nader because he was most aligned with my values," she said.

The diagram shows two sentences with dependency arcs. The first sentence is "I voted for Nader because he was most" and the second is "aligned with my values," she said. Arcs connect "I" to "aligned", "Nader" to "my", "because" to "values", and "he" to "said".

- Human language is **ambiguous** by nature
- Knowing these ambiguities allow us to watch out for these potential problems in our model



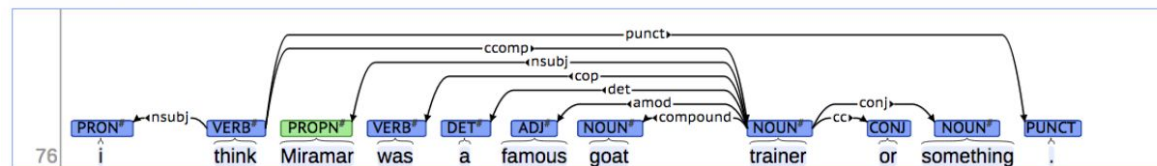
Dependencies Treebanks



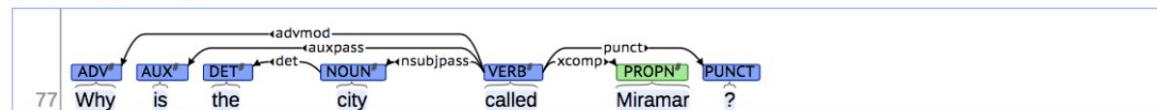
The rise of annotated data and Universal Dependencies treebank

Brown corpus (1967; PoS tagged 1979); Lancaster-IBM Treebank (starting late 1980s); Marcus et al. 1993, The Penn Treebank, Computational Linguistics; Universal Dependencies: <http://universaldependencies.org/>

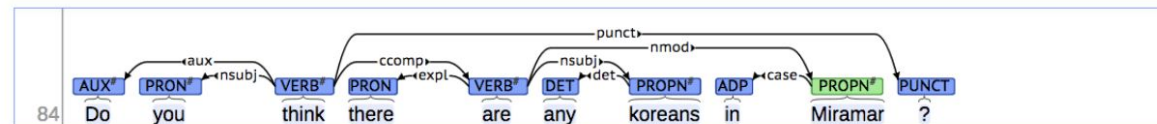
[context] [conllu]



[context] [conllu]



[context] [conllu]



- Many parsers, part-of-speech taggers, etc. can be built on it
- Broad coverage, not just a few intuitions
- Contain frequencies and distributional information
- Provide a way to evaluate existing NLP systems - turn NLP into actual science



So how do we build a parser once we got these dependencies?

Many information you can use from the treebanks

1. **Bilexical affinities**

- a. Check whether both ends of dependencies seem right

2. **Dependency distance**

- a. Check whether the distance is too far away from usual

3. **Intervening material**

- a. Because sentences are normally organized around verbs, dependencies rarely span intervening verbs
- b. Because punctuation (e.g., commas) indicates segment, it also indicates less plausible spanning intervening punctuations

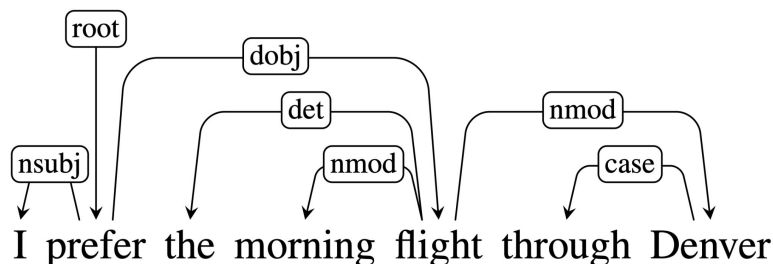
4. **Valency of heads**

- a. Given a head, what is usually its dependents and on which side?



So how do we build a parser once we got these dependencies?

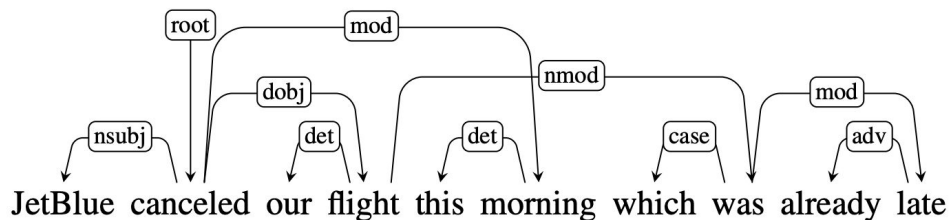
- A sentence is **parsed** by choosing for each word, what other (including ROOT) it is a dependent of
- To make the problem easier to solve, we impose some **constraints**:
 - There is a single designated root node that has no incoming arcs
 - With the exception of the root node, each vertex has exactly one incoming arc
 - There is a unique path from the root node to each vertex in V
 - Don't want cycles, e.g., $A \rightarrow B$, $B \rightarrow A$



Projectivity

Note that some dependency parsing algorithms are **projective**:

Projectivity: An arc from a head to a dependent is said to be projective **if there is a path from the head to every word that lies between the head and the dependent in the sentence**. A dependency tree is then said to be projective if all the arcs that make it up are projective.



In this example, the arc from **flight** to its modifier was is non-projective since there is no path from **flight** to the intervening words **this** and **morning**. A dependency tree is projective if it can be drawn with **no crossing edges**. Here there is no way to link flight to its dependent was without crossing the arc that links morning to its head.



Projectivity

There are computational limitations to the most widely used families of parsing algorithms. The **transition-based approaches** we gonna discuss shortly **can only produce projective trees**, hence any sentences with non-projective structures will necessarily contain some errors. This limitation is one of the motivations for the more flexible graph-based parsing approach.



Dependency Parsing Methods



Greedy transition-based parsing [Nivre IWPT 2003]

- A simple form of greedy discriminative dependency parser
- The parser does a sequence of bottom up actions
 - Roughly like “shift” or “reduce” in a shift-reduce parser, but the “reduce” actions are specialized to create dependencies with head on left or right
- The parser has:
 - **a stack σ** , written with top to the right
 - which starts with the ROOT symbol
 - **a buffer β** , written with top to the left
 - which starts with the input sentence
 - a set of **dependency arcs A**
 - which starts off empty
 - a set of **actions**
 - shift or left arc or right arc

An efficient algorithm for projective dependency parsing, Joakim Nivre, 2003. <https://aclanthology.org/W03-3017/>



Greedy transition-based parsing

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_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_j, \beta, \boxed{AU\{r(w_j, w_i)\}}$

3. Right-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_i, \beta, \boxed{AU\{r(w_i, w_j)\}}$

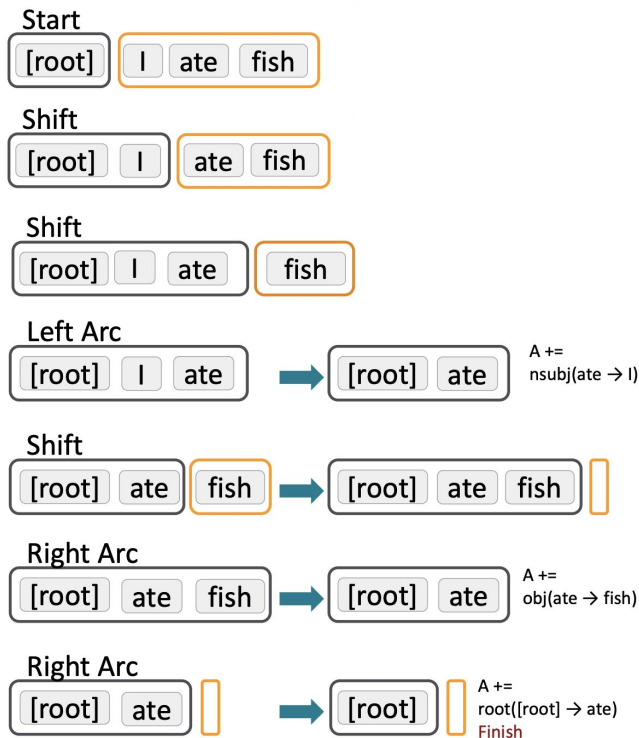
adding the
relationship

Finish: $\sigma = [w]$, $\beta = \emptyset$



Greedy transition-based parsing

- 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_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_j, \beta, A \cup \{r(w_i, w_j)\}$
 3. Right-Arc_r $\sigma | w_i | w_j, \beta, A \rightarrow \sigma | w_i, \beta, A \cup \{r(w_i, w_j)\}$
- Finish:** $\sigma = [w]$, $\beta = \emptyset$



if the only thing in the stack is ROOT, we can only SHIFT

once we shift, we can either perform dependencies by doing LEFT-ARC or RIGHT-ARC. But here, let's say my model tells me to SHIFT more

similarly, my model gonna decide whether to LEFT-ARC or RIGHT-ARC. It seems to be LEFT ARC, because "I" is the subject of "ate"

the dependent is gone. We add the dependencies to A

We SHIFT again. Nothing left on the buffer. We must do either RIGHT-ARC or LEFT-ARC. Let's say the model tells me to perform RIGHT-ARC.

the dependent is gone. We add the dependencies to A

the dependent is gone. We add the dependencies to A. We stop when the buffer is empty and the stack contains only ROOT.



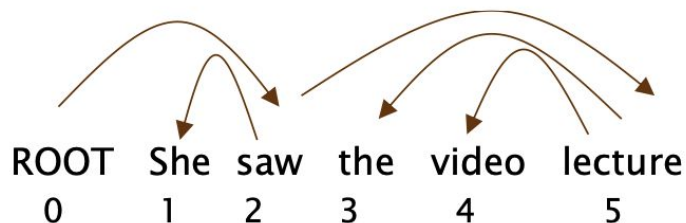
MaltParser [Nivre and Hall 2006]

- **So how to choose the right action??**
 - Answer: Stand back, I know **machine learning**!
- Each action is predicted by a **discriminative classifier** (e.g., SVM) over each legal move
 - Max of 3 untyped choices (shift, left, right); max of $|R| \times 2$ (left, right) + 1 (shift) when typed (consider also the relations)
 - **Features: top of stack word, POS; first in buffer word; etc.**
 - For each word, the features are represented by a one-hot encoding vector. Since there could be millions of features, this vector tends to be very big....(10^6)
 - There is NO search (in the simplest form)
 - But you can profitably do a beam search if you wish (slower but better)
 - You keep k good parse prefixes at each time step
 - The model's accuracy is fractionally below the state of the art in dependency parsing, but it provides very fast linear time parsing, with high accuracy – great for parsing the web

MaltParser: A Data-Driven Parser-Generator for Dependency Parsing, Nivre et al., 2006. <https://aclanthology.org/L06-1084/>



Evaluation of the parser



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

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

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

UAS (unlabeled accuracy score):

Just count how many match without considering the relations (e.g., nsubj)

LAS (labeled accuracy score): the label (e.g., nsubj) must also match

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

Neural Dependency Parsing



A neural transition-based dependency parser [Chen and Manning 2014]

- The traditional dependency parser **feature vector is big** and thus computationally costly
- Instead use **word embeddings!**
 - Concatenated along with part of speech tags (POS) and dependency labels (represented as one-hot, but they are much smaller than the full vector anyway)
- Still use the **transition-based approach**

Other research has further improved by adding deeper network, adding beam search, e.g., SyntaxNet and the Parsey McParseFace model (2016)

Parser	UAS	LAS	sent. / s
MaltParser	89.8	87.2	469
MSTParser	91.4	88.1	10
TurboParser	92.3	89.6	8
C & M 2014	92.0	89.7	654

MST and TurboParser is graph-based with higher accuracy but is much slower....

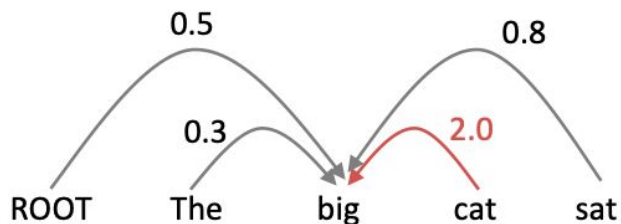
A Fast and Accurate Dependency Parser using Neural Networks, Chen and Manning, 2014. <https://aclanthology.org/D14-1082/>

Announcing Syntaxnet: The World's Most Accurate Parser Goes Open source <https://ai.googleblog.com/2016/05/announcing-syntaxnet-worlds-most.html>



A neural graph-based dependency parser [Dozat and Manning 2017]

- Can handle **non-projectivity**, unlike transition-based
- Compute a score for **every possible dependency** for each word
 - n^2 possible dependencies in a sentence of length n (this makes parsing slow)
- Repeat the same process for each other word
- Determine the optimal tree by using MST algorithm (in our Data Structure and Algorithms class)



e.g., picking the head for “big”

Method	UAS	LAS (PTB WSJ SD 3.3)
Chen & Manning 2014	92.0	89.7
Weiss et al. 2015	93.99	92.05
Andor et al. 2016	94.61	92.79
Dozat & Manning 2017	95.74	94.08

Stanford's Graph-based Neural Dependency Parser at the CoNLL 2017 Shared Task, Dozat et al., 2017. <https://aclanthology.org/K17-3002.pdf>

