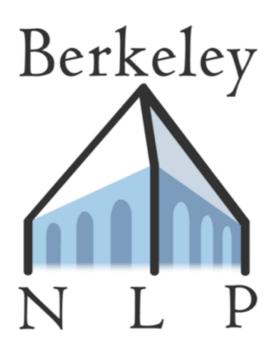
Linguistics: Syntax



EECS 183/283a: Natural Language Processing

Today



- Challenges for language modeling
 - Ambiguity
 - Grammaticality
- Formal languages
- Structures underlying natural language

Ambiguity





with

muffin

NN

chocolate

NNS

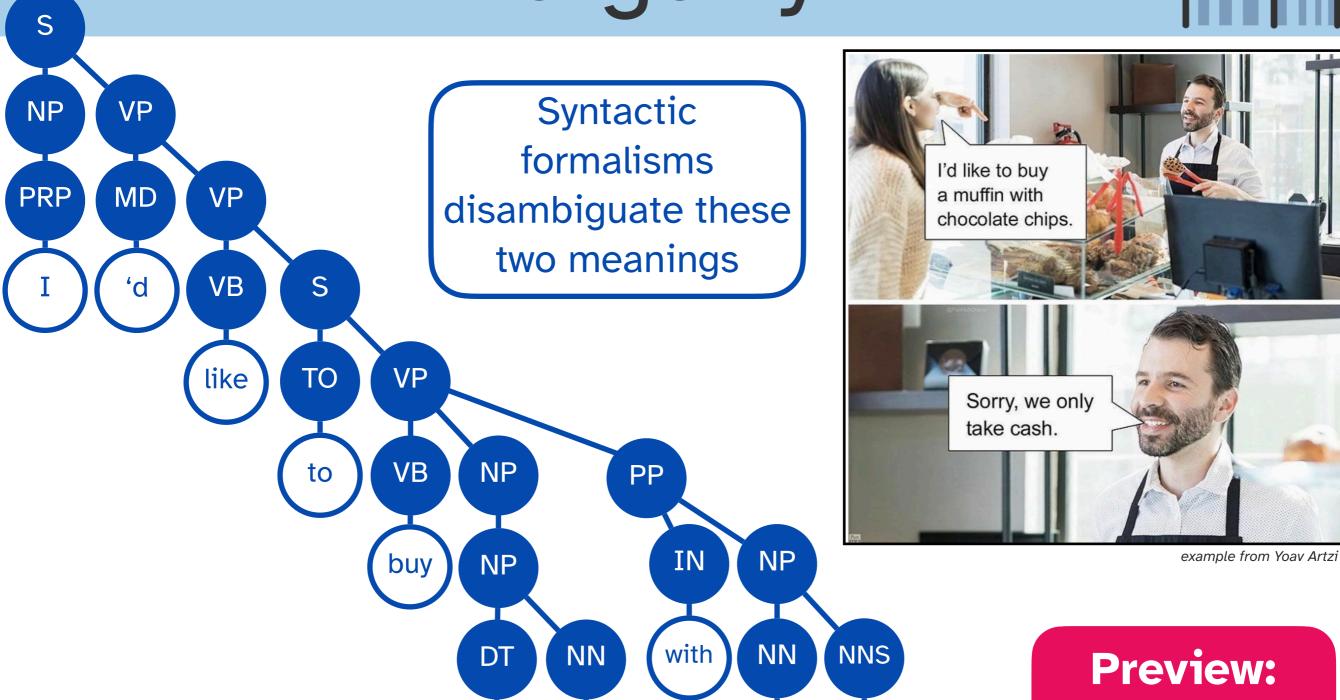
chips

Preview: constituency parsing

Ambiguity

muffin





chips

chocolate

Preview: constituency parsing



- Does a sentence sound "right" according to a native speaker?
- Grammaticality judgments are not universal!
 - No such thing as "bad grammar" disagreements over what feels grammatical or not are due to language variation
 - Grammatical sentences don't need to be meaningful

Colorless green ideas sleep furiously



- Does a sentence sound "right" according to a native speaker?
- Grammaticality judgments are not universal!
 - No such thing as "bad grammar" disagreements over what feels grammatical or not are due to language variation
 - Grammatical sentences don't need to be meaningful

Colorless green ideas sleeps furiously



- Does a sentence sound "right" according to a native speaker?
- Grammaticality judgments are not universal!
 - No such thing as "bad grammar" disagreements over what feels grammatical or not are due to language variation
 - Grammatical sentences don't need to be meaningful

Green colorless ideas sleep furiously



- Why might something sound grammatical or not?
- Our main task: we have some language L, and we want to know whether a new sentence $x \in L$
 - Should we represent L as a finite-sized set of possible sentences?
 - What about a regular expression?

Regular Expressions



- Recall: (simplified) English syllable structure (C3)V(C4), with two types of sounds:
 - Consonants C: {p,b,t,d,k,g,m,n,η,f,v,θ,ð,s,z,∫,ʒ,h,ɹ,l,j,w}
 - Vowels V: {i,u,ε,æ,ι,α,ə,υ,ɔ,aι,ου,eι}
 - Allows us to generate "unattested" but still "accepted" syllables: θ Jaim, sk Joσp, f Jiŋk
- Can syntax also be described with regular expressions?
 - How about {DET}{ADJ}*{NOUN}?

the agreement an ill proposal a tiresome merciful teacher

Regular Expressions



The cat that thinks the cow thinks the rabbits hid are wrong.

The rabbits hid.

N

Regex:

The cow thinks the rabbits hid.

VB DT N

VB

(DT N VB) *

The cat that thinks the cow thinks the rabbits hid is wrong.

VB

N

VB

DT

N

VB VB

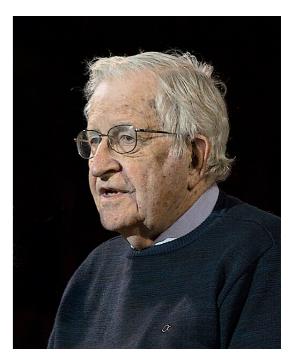
JJ

DT N IN (DT N VB) * VB JJ?

Is there another formalism that tells us whether a string is "accepted" in a language or not?



- CFGs offer arbitrary expressivity through recursive structure
- You might have seen these before in programming languages classes — basically, most programming languages are CFGs
- CFGs were actually first described formally by Noam Chomsky (b. 1928), but the recursive structure of natural language was described by many linguists, including Pāṇini (~1-5th century BCE)



Σ, Wikpiedia





- Set of nonterminal symbols
- Set of terminal symbols (wordtypes)
- Set of production rules defining how nonterminal symbols could be expressed via the composition of other nonterminal and terminal symbols

Nonterminal symbols

```
DT N VB JJ IN
S NP VP SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

Nonterminal symbols

```
DT N VB JJ IN
S NP VP SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong. $\begin{tabular}{c} \begin{tabular}{c} \begin{$

Nonterminal symbols

```
DT N VB JJ IN
S NP VP SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

DT N VB DT N VB DT N VB VB

Nonterminal symbols

```
DT N VB JJ IN
S NP VP SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

VB Ν VB

Nonterminal symbols

```
N VB JJ
   VP
       SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow \{ the, a, an, ... \}
     → { cat, cow, rabbits, dogs, ... }
VB \rightarrow \{hid, is, thinks, was, are, ...\}
JJ \rightarrow \{wrong, right, blue, red, ...\}
IN \rightarrow { that, in, of, because, ... }
NP \rightarrow \{DT N, NP IN VP\}
VP \rightarrow \{VB, VB JJ, VB S\}
     \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

DT N IN VB DT N VB DT N VB VB JJ

Nonterminal symbols

DT N VB JJ IN S NP VP SBAR

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

DT N
NP

IN VB

DT N

VB DT

DT N NP

VB VB JJ

Nonterminal symbols

DT N VB JJ IN S NP VP

```
NP \rightarrow \{ DT N, NP IN VP \}
VP \rightarrow \{ VB, VB JJ, VB S \}
S \rightarrow \{ NP VP, VP \}
```



The cat that thinks the cow thinks the rabbits hid is wrong.

NP

VB

VB

VB VB

Nonterminal symbols

N VB JJ IN NP VP

Production rules

 $NP \rightarrow \{DT N, NP IN VP\}$ $VP \rightarrow \{VB, VB JJ, VB S\}$ $S \rightarrow \{NP \ VP, \ VP\}$



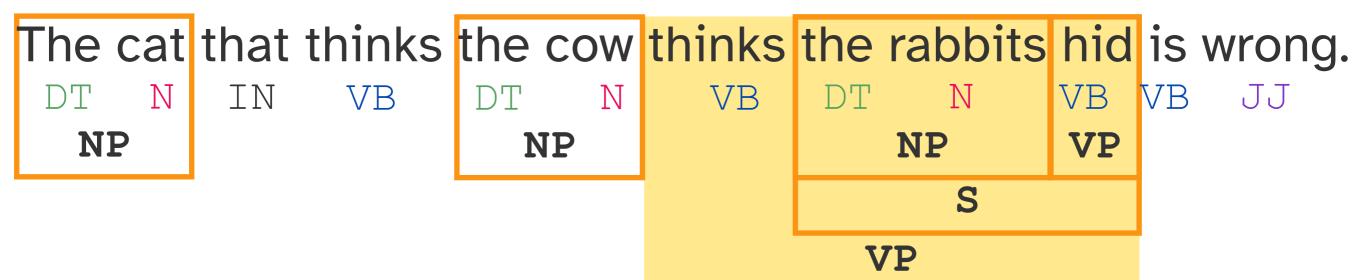
The cat that thinks the cow DT N IN VB DT N VB DT N VB VB JJ NP VP VP

Nonterminal symbols

DT N VB JJ IN S NP VP

```
NP \rightarrow \{DT N, NP IN VP\}
VP \rightarrow \{VB, VB JJ, VB S\}
S \rightarrow \{NP VP, VP\}
```



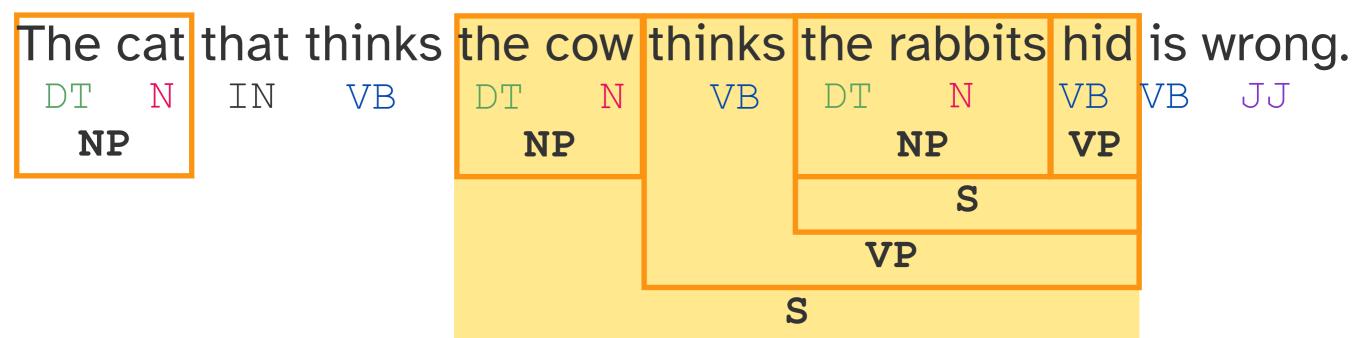


Nonterminal symbols

DT N VB JJ IN S NP VP

```
NP \rightarrow \{DT N, NP IN VP\}
VP \rightarrow \{VB, VB JJ, VB S\}
S \rightarrow \{NP VP, VP\}
```

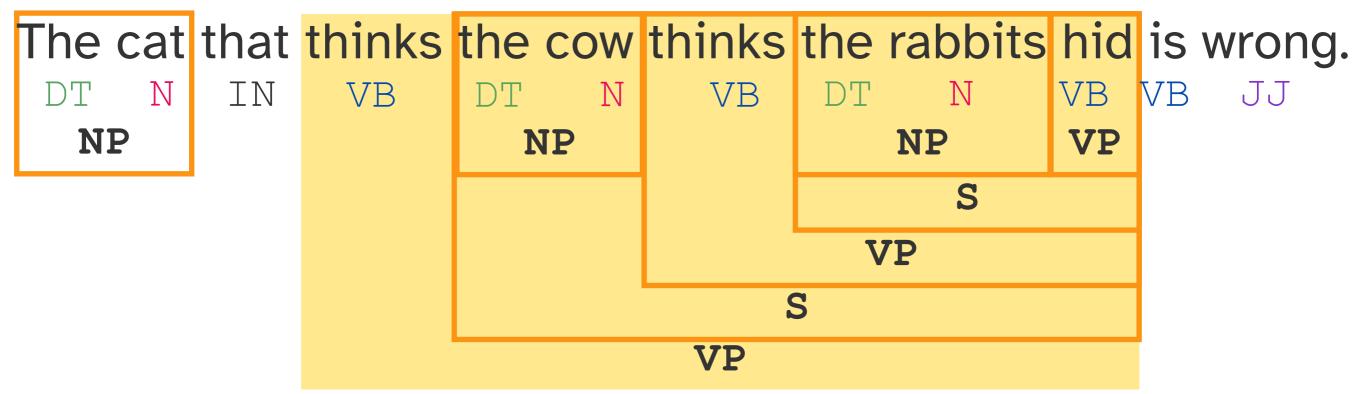




Nonterminal symbols

```
NP \rightarrow \{DT N, NP IN VP\}
VP \rightarrow \{VB, VB JJ, VB S\}
S \rightarrow \{NP VP, VP\}
```

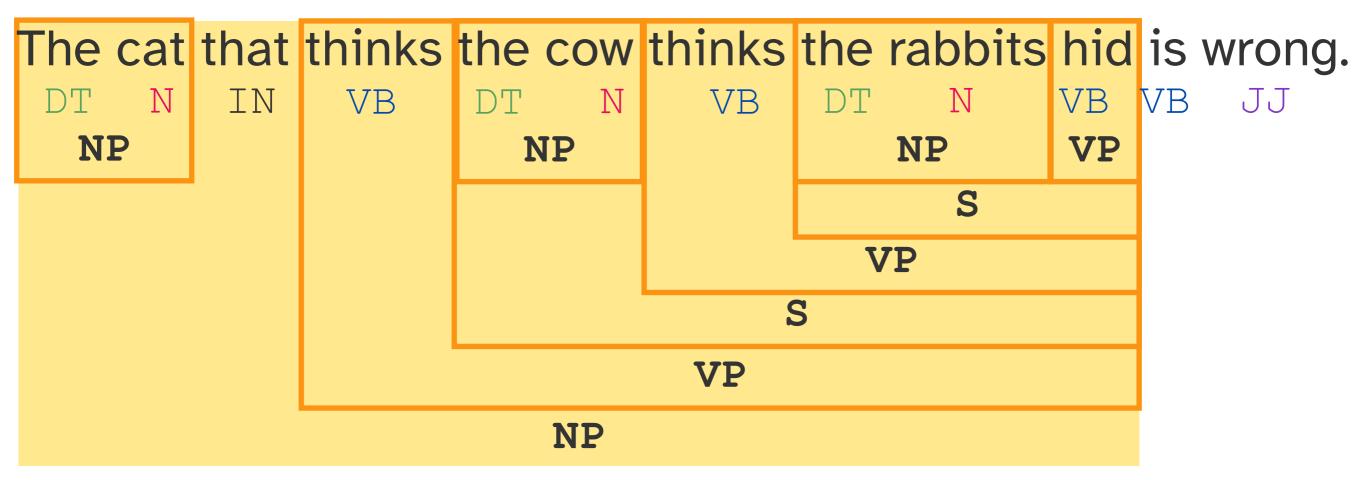




Nonterminal symbols

$$NP \rightarrow \{DT N, NP IN VP\}$$
 $VP \rightarrow \{VB, VB JJ, VB S\}$
 $S \rightarrow \{NP VP, VP\}$

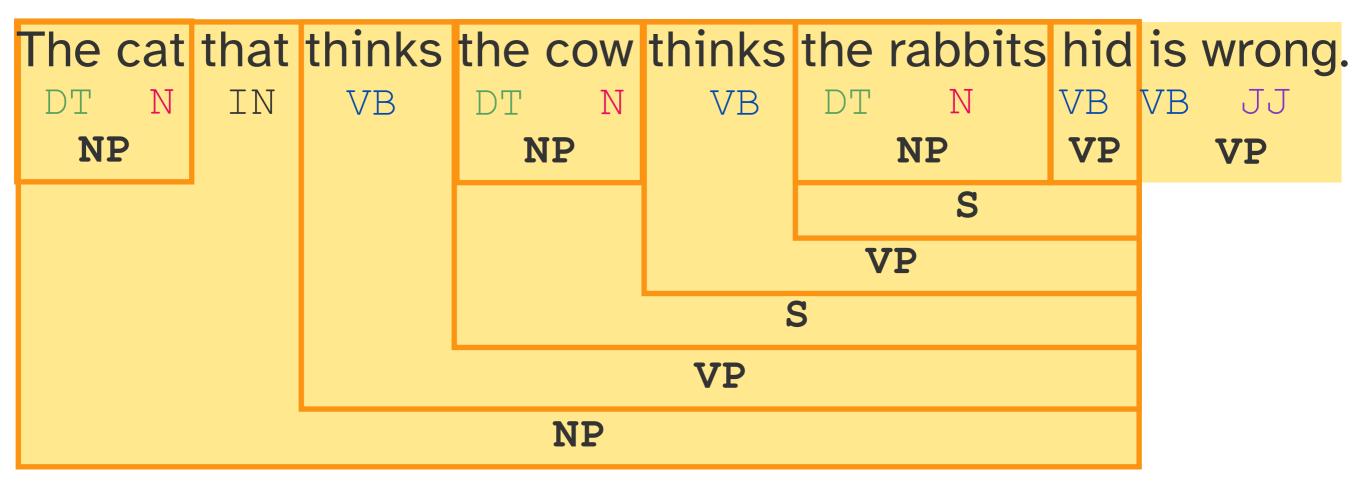




Nonterminal symbols

$$NP \rightarrow \{DT N, NP IN VP\}$$
 $VP \rightarrow \{VB, VB JJ, VB S\}$
 $S \rightarrow \{NP VP, VP\}$

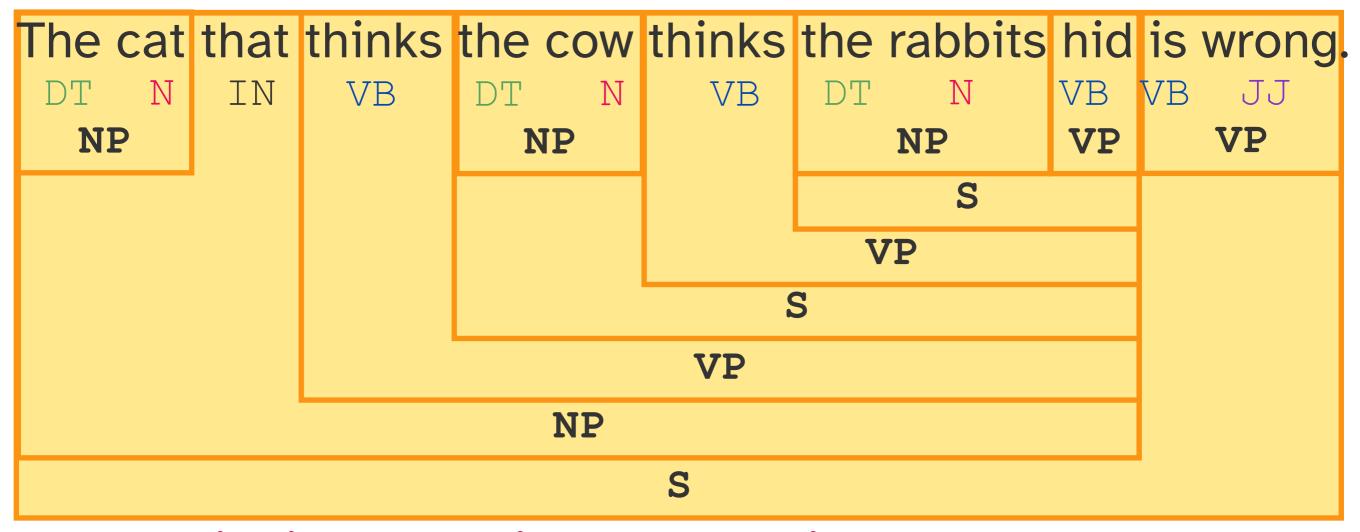




Nonterminal symbols

```
NP \rightarrow \{DT N, NP IN VP\}
VP \rightarrow \{VB, VB JJ, VB S\}
S \rightarrow \{NP VP, VP\}
```





NP is singular, so its corresponding VP should be too

Nonterminal symbols

$$NP \rightarrow \{DT N, NP IN VP\}$$
 $VP \rightarrow \{VB, VB JJ, VB S\}$
 $S \rightarrow \{NP VP, VP\}$

Generating from a Grammar



Nonterminal symbols

```
DT N VB JJ IN
S NP VP SBAR
```

Terminal symbols (vocabulary) are, cat, cow, hid, is, rabbits, that, the, thinks, wrong, ...

```
DT \rightarrow {the, a, an, ...}

N \rightarrow {cat, cow, rabbits, dogs, ...}

VB \rightarrow {hid, is, thinks, was, are, ...}

JJ \rightarrow {wrong, right, blue, red, ...}

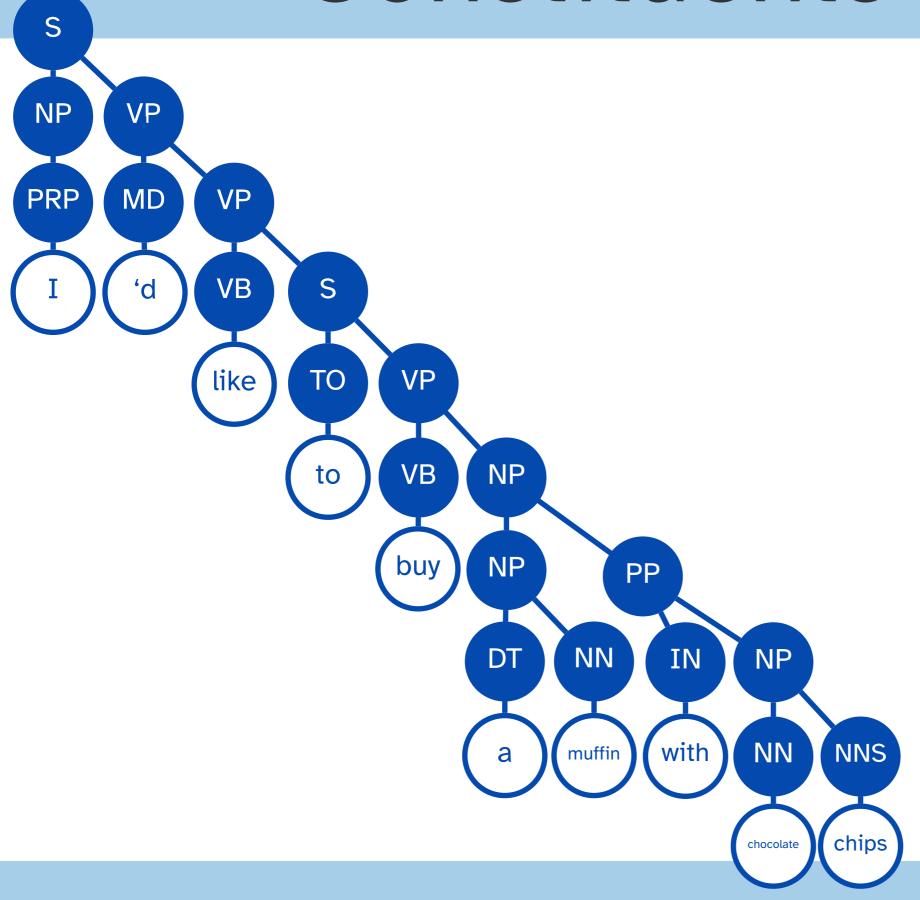
IN \rightarrow {that, in, of, because, ...}

NP \rightarrow {DT N, NP IN VP}

VP \rightarrow {VB, VB JJ, VB S}

S \rightarrow {NP VP, VP}
```

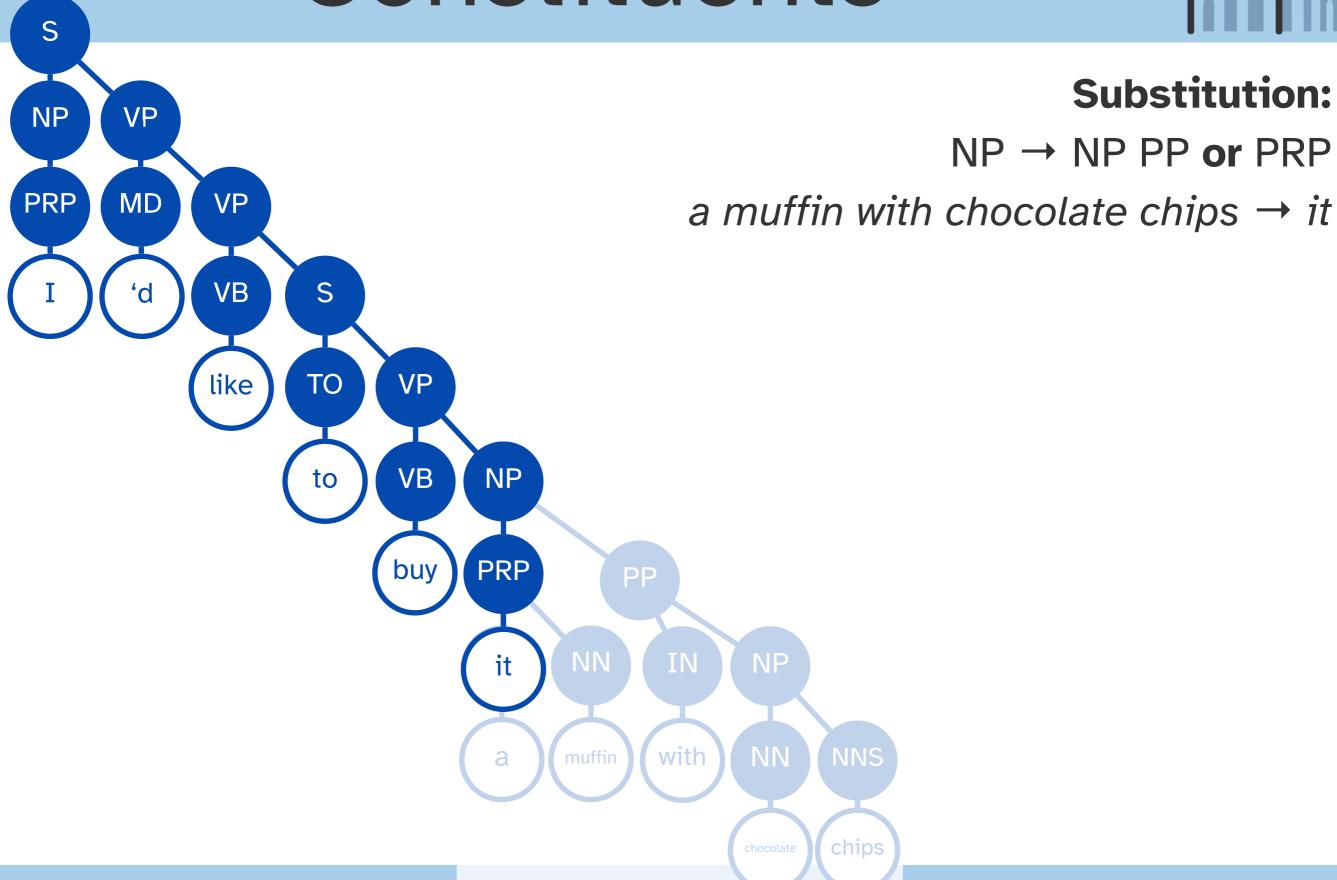




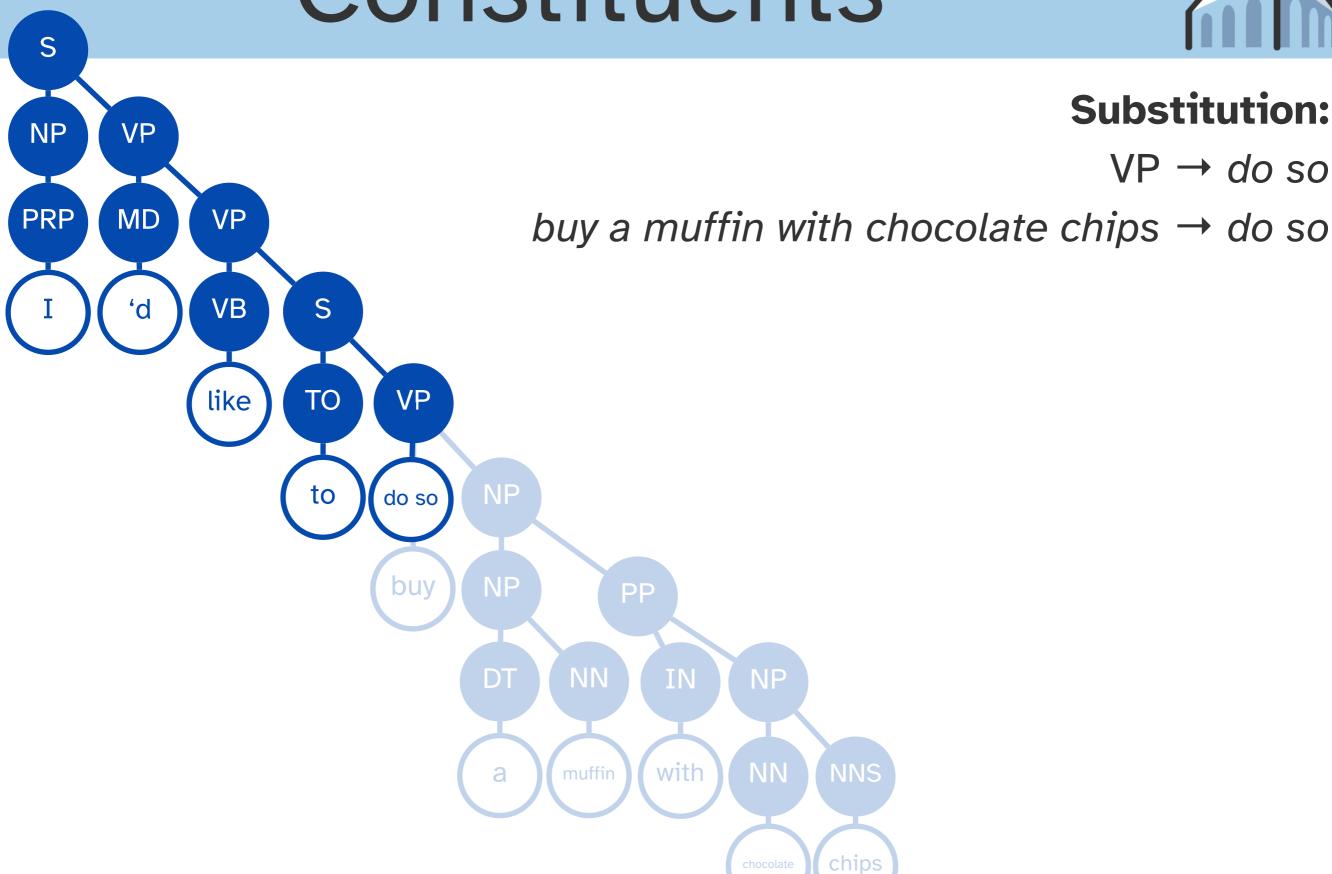


Substitution:

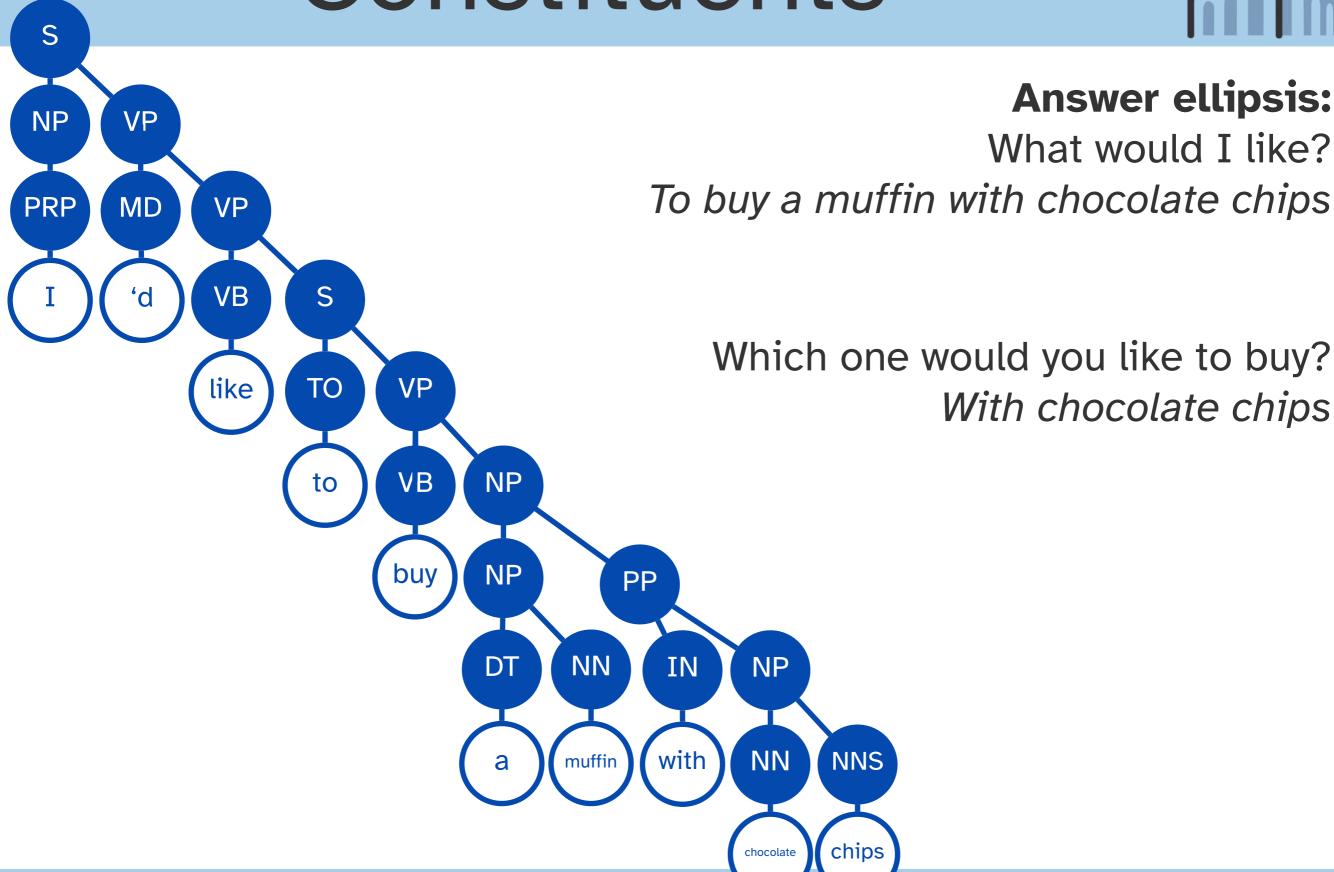
 $NP \rightarrow NP PP or PRP$





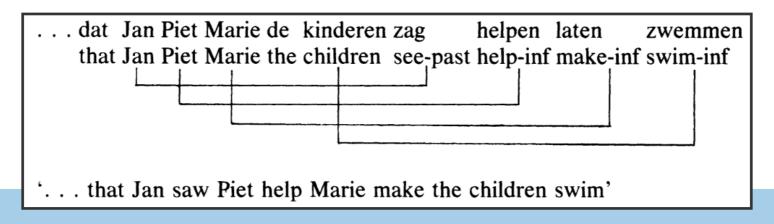








- CFG is often called phrase structure or constituency grammar
- Each production rule describes a constituent
- Constituent constructions are independent of one another (this is why the grammar is <u>context-free</u>)
 - Augmenting a CFG with agreement (e.g., distinguishing plural vs. singular NPs and plural vs. singular VPs plural) means it is no longer context-free
 - Also, some languages aren't even context-free, not even considering agreement:



Probabilistic CFG



Production rules

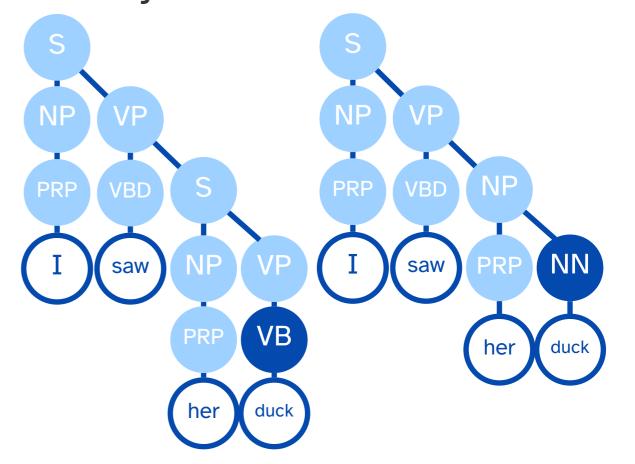
augmented with probabilities

```
\begin{array}{lll} \text{DT} & \to \ ^{5}the, ^{2}a, ^{2}an, \dots \ ) \\ \text{N} & \to \ ^{6}cat, ^{2}cow, ^{2}rabbits, ^{0}cogs, \dots \ ) \\ \text{VB} & \to \ ^{6}hid, ^{1}is, ^{1}thinks, ^{1}was, ^{3}are, \dots \ ) \\ \text{JJ} & \to \ ^{6}that, ^{1}in, ^{1}of, ^{1}blue, ^{1}red, \dots \ ) \\ \text{IN} & \to \ ^{6}that, ^{1}in, ^{1}of, ^{1}because, \dots \ ) \\ \text{NP} & \to \ ^{6}that, ^{1}in, ^{1}of, ^{1}because, \dots \ ) \\ \text{NP} & \to \ ^{6}that, ^{1}ve, ^{1}v
```

p(rule|nonterminal)

learn from data

syntactic parsing: given a PCFG, which rule applications generated our sentence? Why is this hard?



Combinatory Categorial Grammar (CCG)



- Another way of representing a constituency grammar: bottom-up
- Elements of a CCG:
 - Lexical items (wordtypes)
 - Each paired with a syntactic type (≈ nonterminal or composition thereof)

$\mathrm{the}:NP/N$	$\mathrm{dog}:N$	$\mathrm{John}:NP$	$\mathrm{bit}: (S\backslash NP)/NP$
If a Noun appears to the right, then it creates a Noun	Noun	Noun Phrase	If a Noun Phrase appears to the right, then it creates an element with the type
Phrase	$N \rightarrow dog$	$NP \rightarrow PRO$ $PRO \rightarrow John$	S\NP
$NP \rightarrow DT N$ $DT \rightarrow the$			If an NP appears to the left of that element, it creates a Sentence

Combinatory Categorial Grammar (CCG)

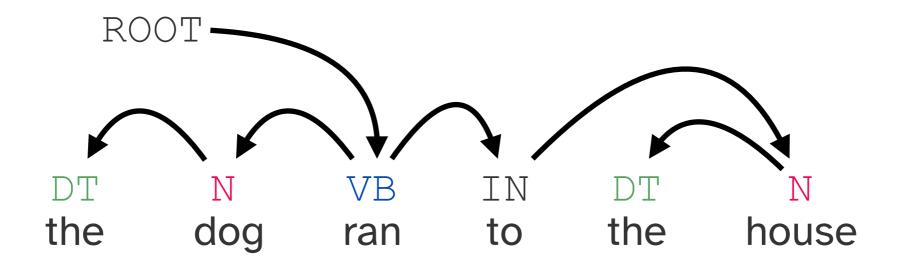


- Another way of representing a constituency grammar: bottom-up
- Elements of a CCG:
 - Lexical items (wordtypes)
 - Each paired with a syntactic type (≈ nonterminal or composition thereof)

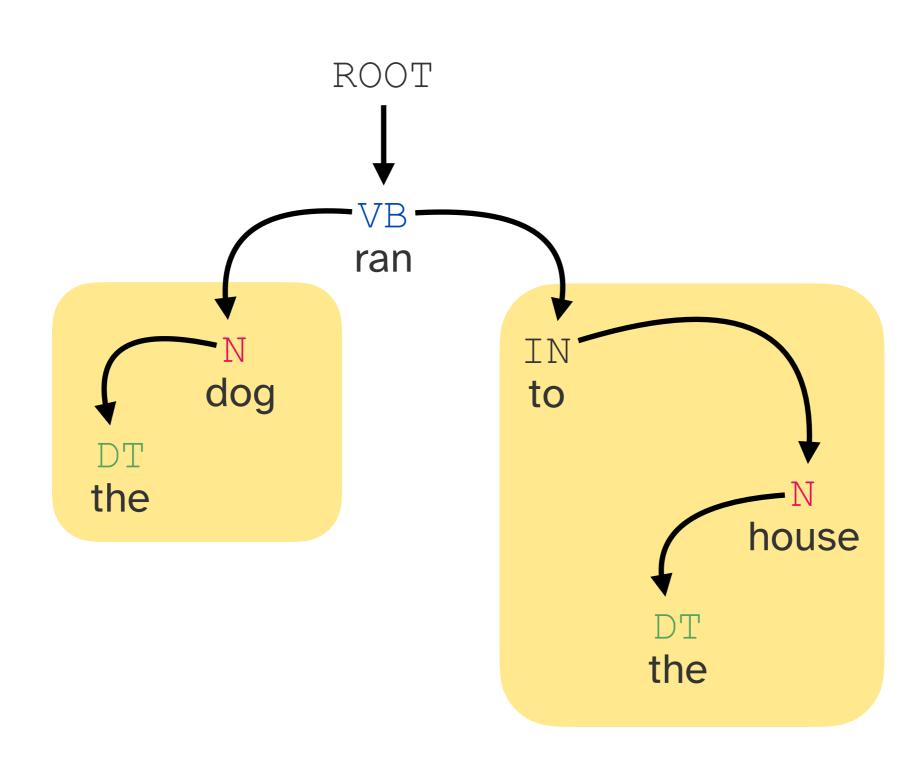
$$he : NP/N \qquad he : NP \qquad he : (S ackslash NP)/NP$$



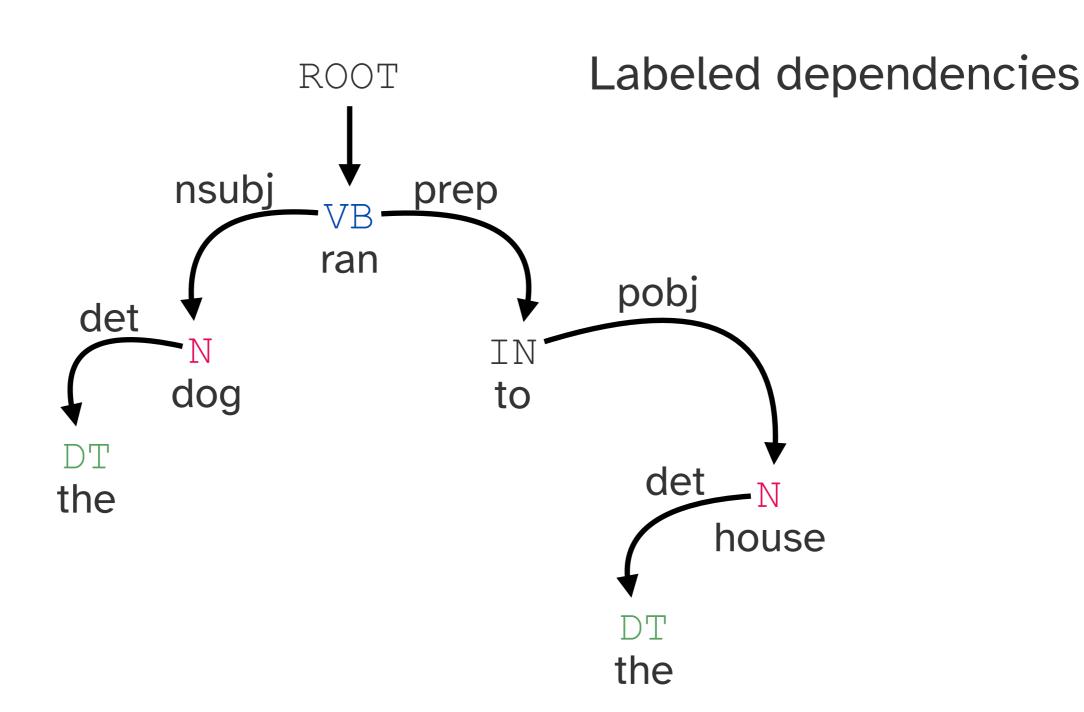
- A dependency grammar includes:
 - Terminal symbols (wordtypes)
 - Parts of speech
 - Some constraints on which parts of speech can be attached to other parts of speech



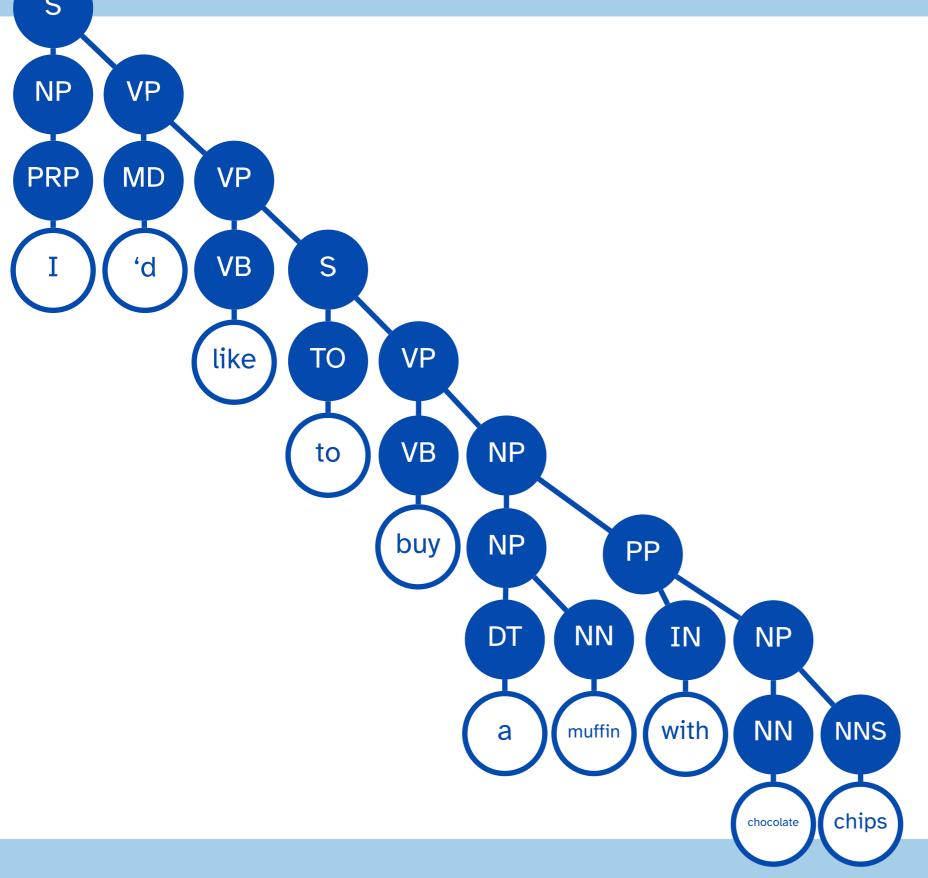




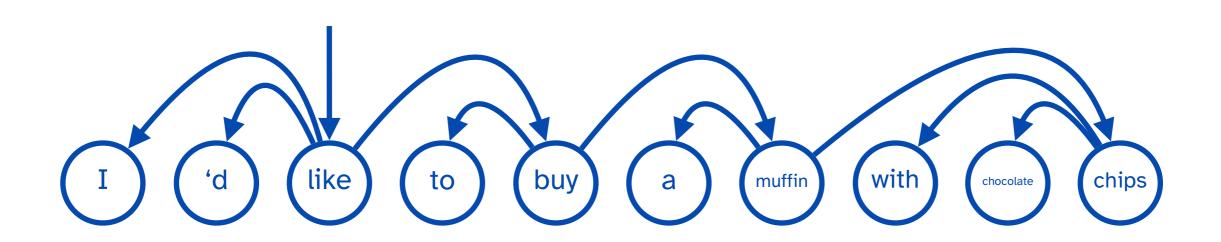




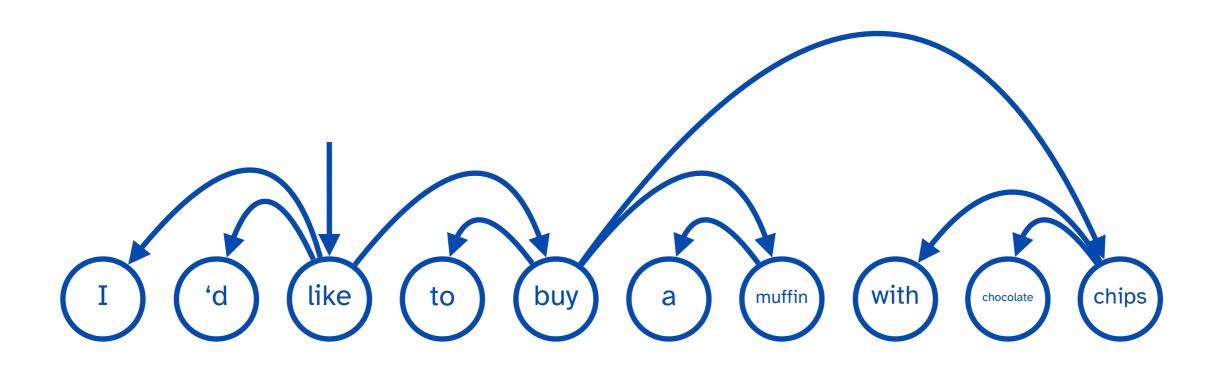












Universal Dependencies



Huge resource of grammar annotations across ~150 languages

- Parts of speech
- Morphological features
- Syntactic dependency parses

English

Bulgarian

Czech

Swedish

