



Formal Grammars for human languages

What are grammar rules?

- It is important to never split infinitives.
- Prepositions are bad to end sentences with.
- When should you use who v. whom? Whom knows?
- Follow these rules and you'll make less errors. fewer See more examples from Steven Pinker's "10 grammar rules it's OK to break"

Prescriptive grammar rules ← Important for paper writing

Descriptive grammar rules ← What you will write papers about.

A basic formalism for describing how human languages are structured is called a context free grammar or CFG.

It was formalized by Noam Chomsky in 1956. A CFG allows us to describe the structure of language in terms of recursive tree structures.





CFGs consist of a set of rules (also called productions), which have terminal symbols (words) and non-terminal symbols (parts of speech or phrase structures)

Det → a

Det → the

Noun → flight

Verb → landed

Adverb → safely

prep → in

ProperNoun → Philadelphia

Rules can also take the form

NP → Det Noun

VP → Verb PP

PP → Prep NP

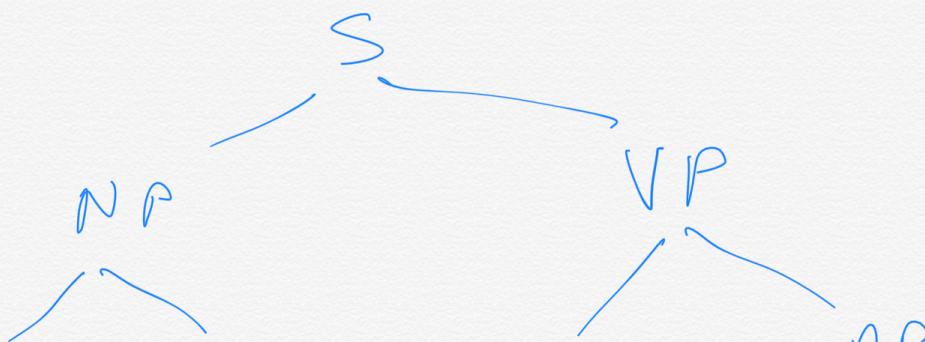
Notes



Rules can also take the form

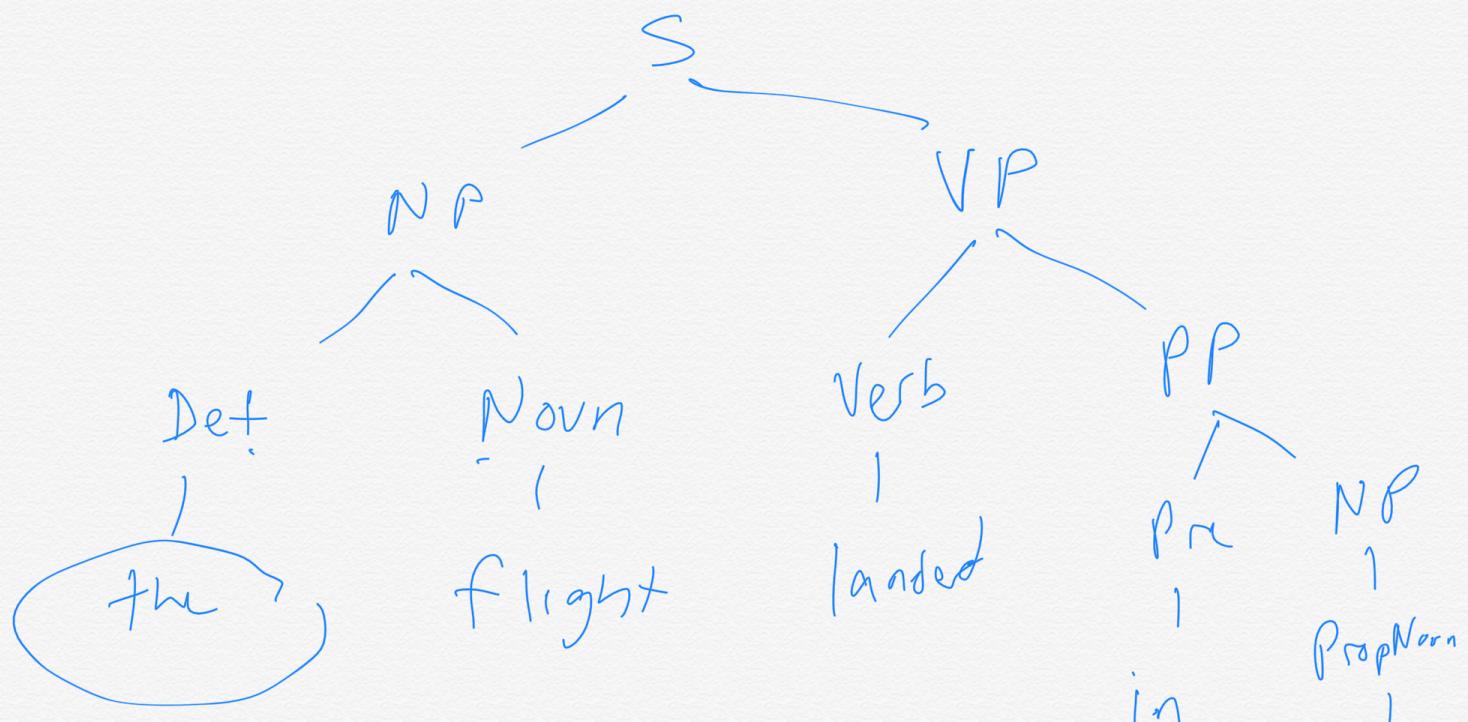
$$NP \rightarrow \text{Det } \text{Noun}$$
$$VP \rightarrow \text{Verb } PP$$
$$PP \rightarrow \text{Prep } NP$$
$$NP \rightarrow \text{Proper Noun}$$
$$VP \rightarrow VP \text{ Adverb}$$
$$VP \rightarrow \text{Verb}$$
$$S \rightarrow NP \text{ VP}$$

You can think of a CFG as
a device for generating a sentence
or as way to assign a structure
to an input sentence.





You can think of a CFG as a device for generating a sentence or as way to assign a structure to an input sentence.



Analysis = parsing

Philadelphia

Formal definition of a CFG.
Tuple with $\langle N, \Sigma, R, S \rangle$

out of non-terminal symbols



N





Formal definition of a CFG.
Tuple with $\langle N, \Sigma, R, S \rangle$

N - set of non-terminal symbols

Σ - set of terminal symbols
disjoint from N

R - a set of rules of the
form $A \rightarrow B$ where

$A \in N$ and B is
a string of symbols (Σ^N)*

S - start symbol ($\in N$)

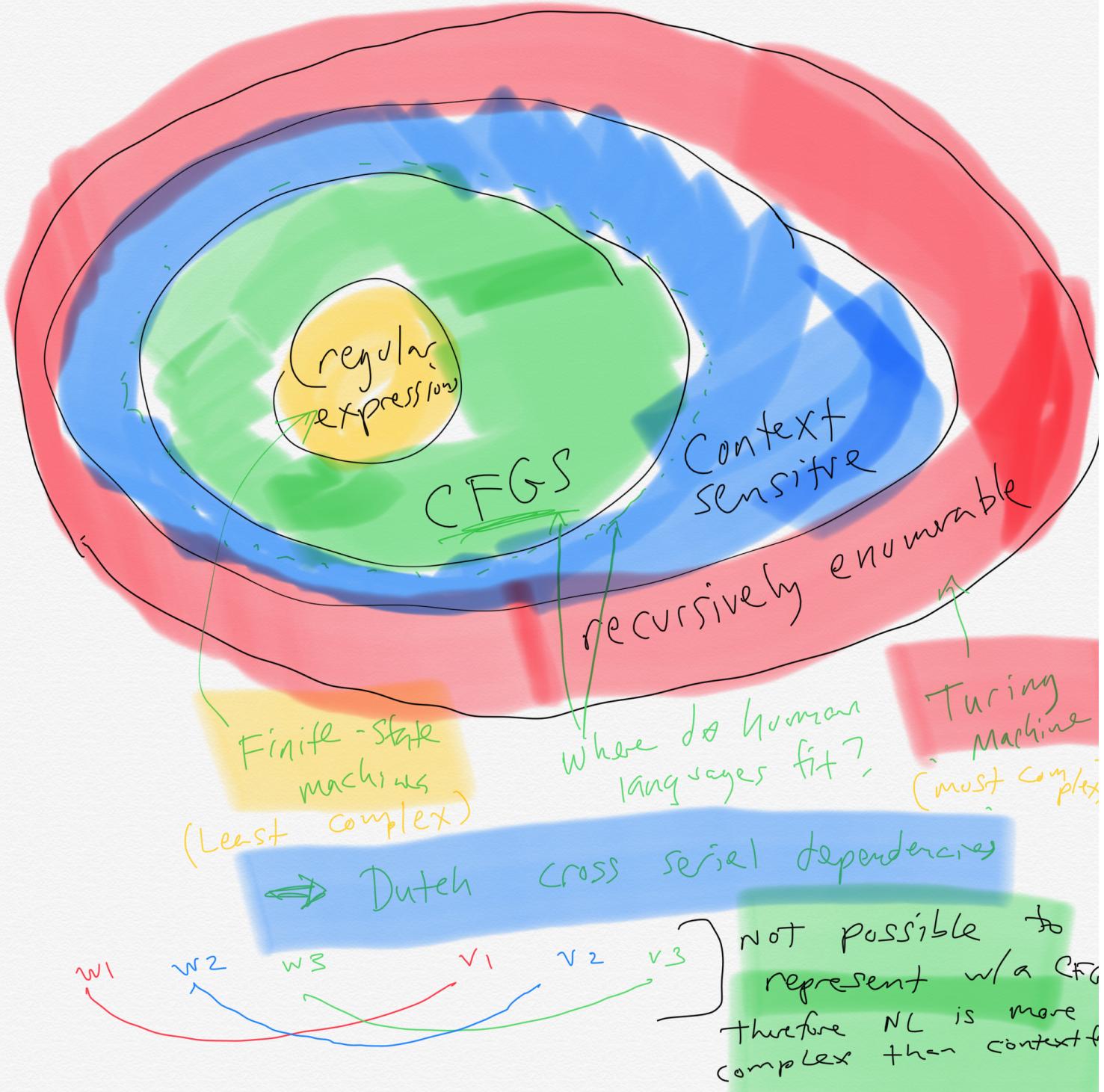
Chomsky hierarchy



< Notes



Chomsky hierarchy



Sentence level constructions in English





Sentence level constructions in English

- Declarative
- Imperative
- Yes/No questions
- Wh- questions

Subject non-subject
wh-questions wh-questions

Phrases

Noun phrases

Determiners

the cat

Posessives instead of determiners

CCB's class



Posessives instead of determiners

CCB's class

Mass nouns

{ music

pre-modifiers (numbers, adjective phrases, quantifiers)

all flights

Post-modifiers (PPs, gerundines, relative pronouns
that/who)

all classes holding midterms
in CIS

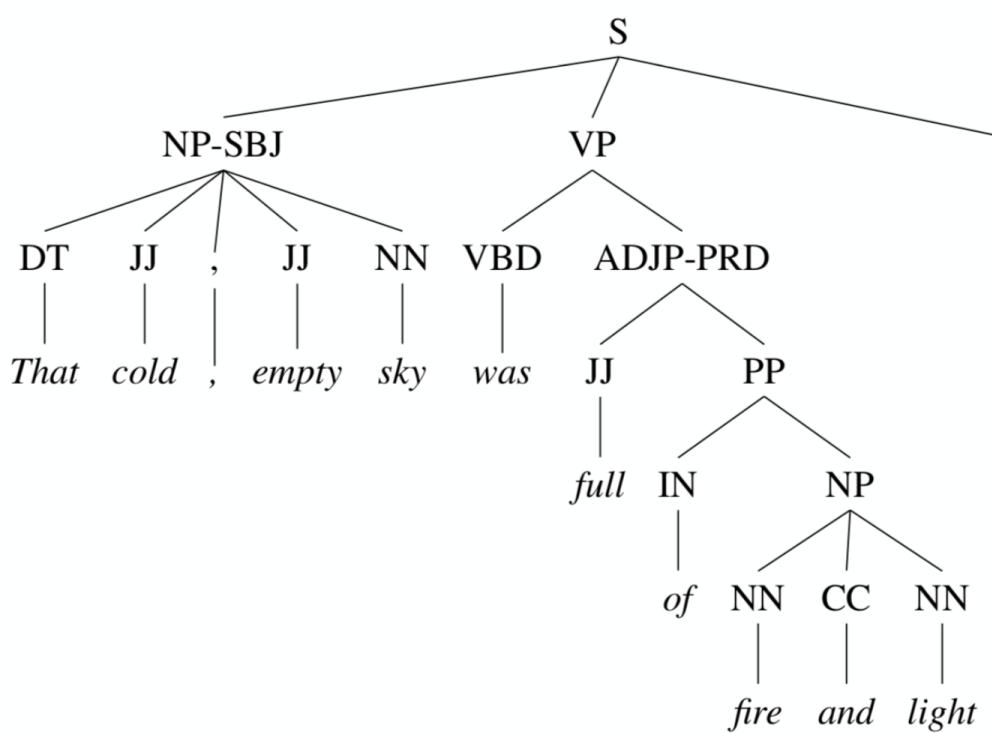
Verb phrase

arguments / valency / subcategorization
frames



Q: How do we construct a grammar for a language?
write it by hand?

Treebanks.





Grammar	Lexicon
$S \rightarrow NP VP.$	$PRP \rightarrow we he$
$S \rightarrow NP VP$	$DT \rightarrow the that those$
$S \rightarrow "S", NP VP.$	$JJ \rightarrow cold empty full$
$S \rightarrow -NONE-$	$NN \rightarrow sky fire light flight tomorrow$
$NP \rightarrow DT NN$	$NNS \rightarrow assets$
$NP \rightarrow DT NNS$	$CC \rightarrow and$
$NP \rightarrow NN CC NN$	$IN \rightarrow of at until on$
$NP \rightarrow CD RB$	$CD \rightarrow eleven$
$NP \rightarrow DT JJ, JJ NN$	$RB \rightarrow a.m.$
$NP \rightarrow PRP$	$VB \rightarrow arrive have wait$
$NP \rightarrow -NONE-$	$VBD \rightarrow was said$
$VP \rightarrow MD VP$	$VBP \rightarrow have$
$VP \rightarrow VBD ADJP$	$VBN \rightarrow collected$
$VP \rightarrow VBD S$	$MD \rightarrow should would$
$VP \rightarrow VBN PP$	$TO \rightarrow to$
$VP \rightarrow VB S$	
$VP \rightarrow VB SBAR$	
$VP \rightarrow VBP VP$	
$VP \rightarrow VBN PP$	
$VP \rightarrow TO VP$	
$SBAR \rightarrow IN S$	
$ADJP \rightarrow JJ PP$	
$PP \rightarrow IN NP$	

Figure 11.10 A sample of the CFG grammar rules and lexical entries that would be extracted from the three treebank sentences in Fig. 11.7 and Fig. 11.9.

Treebanks as grammars.

PennTreebank has 4500 VP rules
and 10005 NP rules

Common modifications to CFGs

- lexicalized rules that encode the head of the phrase





Common modifications to CFGs

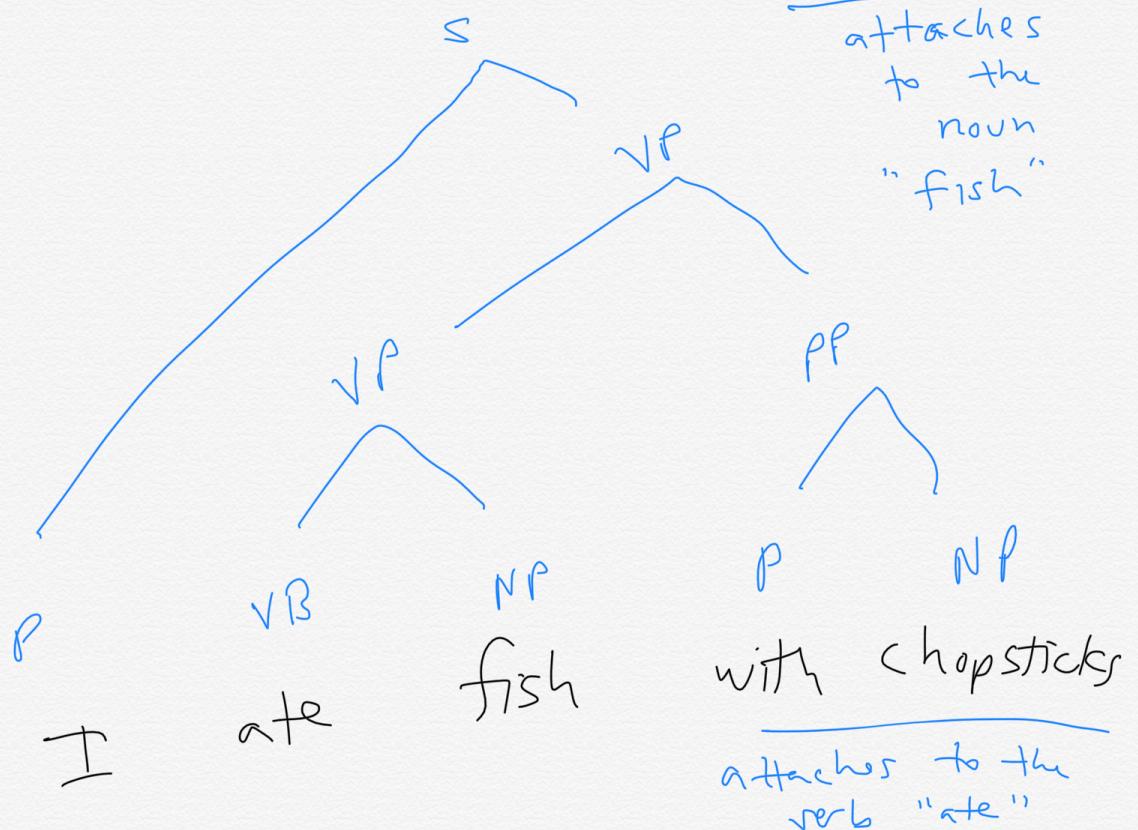
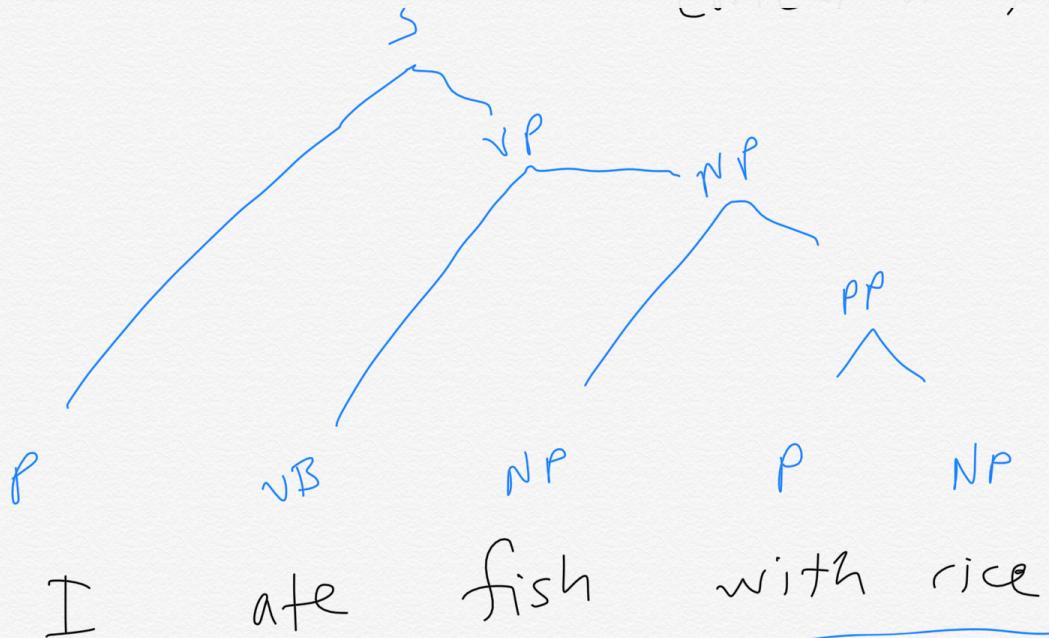
- lexicalized rules that encode the head of the phrase
- estimating probabilities for all production rules. (See ch. 13)
- Synchronous CFGs for translation

Important concepts

- Parsing: derive one or more parse trees for an input sentence.
- Language is ambiguous so this is hard to find the correct analysis.



< Notes



Different syntactic analysers encode different semantics.