

SYNTAX



Matt Post
IntroHLT class
10 September 2020

*and stupor his the Fred with
pain from ease couldn't would
a set he cigarette out the that
for in wife Jones was during
caring a often drugs house but
screaming the crying at for
didn't fear worn sleep ablaze
day the from that she night*

Fred Jones was worn out from caring for his often screaming and crying wife during the day but he couldn't sleep at night for fear that she in a stupor from the drugs that didn't ease the pain would set the house ablaze with a cigarette

- 46 words, 46! permutations of those words, the vast majority of them ungrammatical and meaningless
- How is that we can
 - process and understand this sentence?
 - discriminate it from the sea of ungrammatical permutations it floats in?

Today we will cover

Linguistics

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

Computer Science

Goals for today

- After today, you should be able to
 - Give a working definition of syntax and describe how linguists think about it
 - Describe two well-known grammar formalisms and projects supporting them
 - Discuss issues related to universal language features
 - Describe the formal language hierarchy
 - Describe algorithms for parsing the two grammar formalisms

Outline

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

Linguistic fields of study

- Phonetics: sounds

Linguistic fields of study

- Phonetics: sounds
- Phonology: sound systems

Linguistic fields of study

- Phonetics: sounds
- Phonology: sound systems
- Morphology: internal word structure

Linguistic fields of study

- Phonetics: sounds
- Phonology: sound systems
- Morphology: internal word structure
- **Syntax**: external word structure (sentences)

Linguistic fields of study

- Phonetics: sounds
- Phonology: sound systems
- Morphology: internal word structure
- **Syntax**: external word structure (sentences)
- Semantics: sentence meaning

Linguistic fields of study

- Phonetics: sounds
- Phonology: sound systems
- Morphology: internal word structure
- **Syntax**: external word structure (sentences)
- Semantics: sentence meaning
- Pragmatics: contextualized meaning and communicative goals

Aside

- Much of our focus is on *written language*, but language is first and foremost spoken
- Why does this matter?

Aside

- Much of our focus is on *written language*, but language is first and foremost spoken
- Why does this matter?
- Which of these is easier for a computer to work with?

Aside

- Much of our focus is on *written language*, but language is first and foremost spoken
- Why does this matter?
- Which of these is easier for a computer to work with?
 - (written) *Dipanjan asked a question*

Aside

- Much of our focus is on *written language*, but language is first and foremost spoken
- Why does this matter?
- Which of these is easier for a computer to work with?
 - (written) *Dipanjan asked a question*
 - (spoken) *Dipanjan, uh, he, uh, um, was wondering, uh, he had a question*

Today's focus



MORGAN & CLAYPOOL PUBLISHERS

Linguistic Fundamentals for Natural Language Processing

*100 Essentials from
Morphology and Syntax*

Emily M. Bender

*SYNTHESIS LECTURES ON
HUMAN LANGUAGE TECHNOLOGIES*

Graeme Hirst, *Series Editor*

What is syntax?

- A set of constraints on the possible sentences in the language
 - *A set of constraint on the possible sentence.
 - *Dipanjan had [a] question.
 - *You are on class.
- At a coarse level, we can divide all possible sequences of words into two groups: *valid* and *invalid* (or *grammatical* and *ungrammatical*)

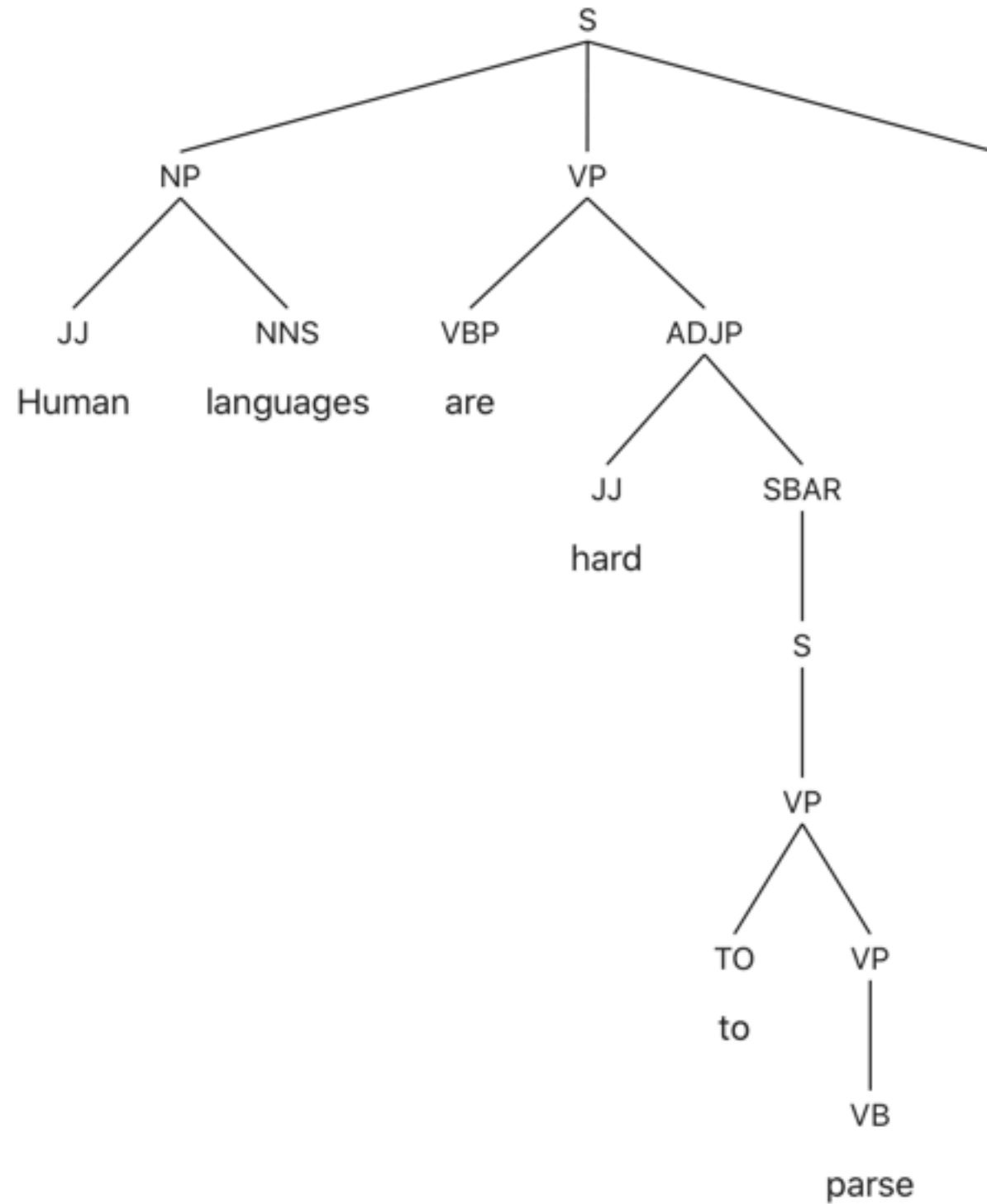
Human judgments

- How do we know what's in and out? We simply ask humans
- But how do humans know?
 - Bad idea: big lists
 - Better idea: grammars

A hierarchical view

- A grammar is a *finite set of **rules*** licensing a (possibly infinite) *number of **strings***
- e.g., some rules
 - [sentence] → [subject] [predicate]
 - [subject] → [noun phrase]
 - [noun phrase] → [determiner]? [adjective]* [noun]
 - [predicate] → [verb phrase] [adjunct]
- Rules are *phrasal or terminal*
 - Phrasal rules form **constituents** in a tree
 - Terminal rules are **parts of speech** and produce words

Example



POS Examples

- No general agreement about the exact set of parts of speech
- Penn Treebank tagset examples

POS Examples

- No general agreement about the exact set of parts of speech
- Penn Treebank tagset examples
 - nouns: NN, NNS, NNP, NNPS

POS Examples

- No general agreement about the exact set of parts of speech
- Penn Treebank tagset examples
 - nouns: NN, NNS, NNP, NNPS
 - adverbs: RB, RBR, RBS, RP

POS Examples

- No general agreement about the exact set of parts of speech
- Penn Treebank tagset examples
 - nouns: NN, NNS, NNP, NNPS
 - adverbs: RB, RBR, RBS, RP
 - verbs: VB, VBD, VBG, VBN, VBP, VBZ

POS Examples

- No general agreement about the exact set of parts of speech
- Penn Treebank tagset examples
 - nouns: NN, NNS, NNP, NNPS
 - adverbs: RB, RBR, RBS, RP
 - verbs: VB, VBD, VBG, VBN, VBP, VBZ
 - (Here, different tags are used to capture the small bit of morphology present in English)

Parts of Speech (POS)

- Three definitions of **noun**

Grammar school

(“metaphysical”)

*a person, place,
thing, or idea*

Parts of Speech (POS)

- Three definitions of **noun**

Grammar school

(“metaphysical”)

*a person, place,
thing, or idea*

Distributional

*the set of words
that have the
same distribution
as other nouns*

*{I, you, he} saw the
{bird, cat, dog}.*

Parts of Speech (POS)

- Three definitions of **noun**

Grammar school

(“metaphysical”)
*a person, place,
thing, or idea*

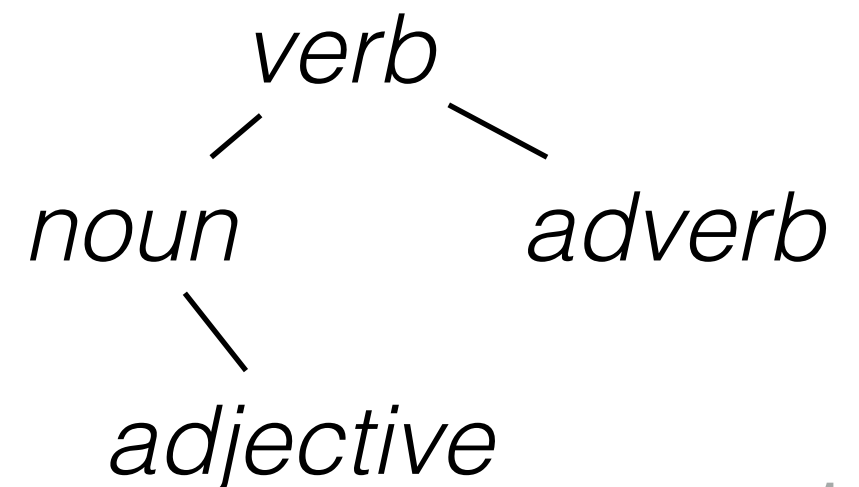
Distributional

*the set of words
that have the
same distribution
as other nouns*

*{I, you, he} saw the
{bird, cat, dog}.*

Functional

*the set of words
that serve as
arguments to
verbs*



Phrases and Constituents

- Longer sequences of words can perform the same function as individual parts of speech:
 - I saw [a_{DT} kid_N]_{NP}
 - I saw [a kid playing basketball]_{NP}
 - I saw [a kid playing basketball alone on the court]_{NP}
- This gives rise to the idea of a *phrasal constituent*, which functions as a unit in relation to the rest of the sentence

Constituent tests

- How do you know if a phrase functions as a constituent?
- A few tests
 - *Coordination*
 - Kim [read a book], [gave it to Sandy], and [left].
 - *Substitution with a word*
 - Kim read [a very interesting book about grammar].
 - Kim read [it].
 - See Bender #51

Heads, arguments, & adjuncts

- Syntax is about the relationships among words and phrases in a sentence

Heads, arguments, & adjuncts

- Syntax is about the relationships among words and phrases in a sentence
- Each constituent has its own internal structure, as well as relationship with words and constituents outside it

Heads, arguments, & adjuncts

- Syntax is about the relationships among words and phrases in a sentence
- Each constituent has its own internal structure, as well as relationship with words and constituents outside it
- Hierarchical structure among constituents
 - Top down, each constituent has a head
 - Heads have (phrasal) dependents
 - Dependents can be required (*arguments*) or optional (*adjuncts*)
 - A head word often controls the structure of its modifiers

Heads

- *Head*: “the sub-constituent which determines the internal structure and external distribution of the constituent as a whole” (Bender #52)
- Examples
 - sentence: (usually) the main verb
 - noun phrase: (usually) the main noun
 - verb phrase: (usually) the active verb

Dependents: Arguments & adjuncts

- *Dependents* of a head:
 - *Arguments*: selected/licensed by the head and **complete** the meaning
 - *Adjuncts*: not selected and **refine** the meaning

Constituent structure

- The head often constrains the internal structure of a constituent
- Examples
 - verb
 - [Kim]^{ARGUMENT} **is** [ready]^{ADJUNCT}.
 - adjective
 - Kim is [**ready**_{ADJ} [to make a pizza]_v].
 - * Kim is [**tired**_{ADJ} [to make a pizza]_v].
 - noun
 - [The [red]_{ADJ} **ball**]
 - * [The [red]_{ADJ} **ball** [the stick]_N]
 - [The [red]_{ADJ} **ball** [on top of the stick]_{PP}]

More examples

- Kim **planned** [to **give** Sandy books].
- * Kim **planned** [to **give** Sandy].
- Kim **planned** [to give books].
- * Kim **planned** [to **see** Sandy books].
- Kim [**would** [**give** Sandy books]].
- Pat [**helped** [Kim **give** Sandy books]].
- * [[**Give** Sandy books] [**surprised** Kim]].

Summary

what is
syntax?

*A finite set of rules licensing
an infinite number of strings*

*The rules specify how words and
phrases relate to one another in a
hierarchical manner*

*No one knows what the actual rules
are, but there is consensus that the
rules must exist!*

Outline

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

Treebanks

- Collections of natural text that are annotated according to a particular syntactic theory
 - Usually created by linguistic experts
 - Ideally as large as possible
 - Theories are usually coarsely divided into *constituent/phrase* or *dependency* structure

Formalisms

- **Phrase-structure** and **dependency** grammars
 - Phrase-structure: encodes the phrasal components of language
 - Dependency grammars encode the relationships between words

Penn Treebank (1993)

ABOUT
MEMBERS
COMMUNICATIONS
LANGUAGE RESOURCES ▾
Data ▾

Obtaining Data
Catalog
By Year
Top Ten Corpora
Projects
Search
Memberships
Data Scholarships
Tools >
Papers >
LR Wiki
DATA MANAGEMENT
COLLABORATIONS

Home > Language Resources > Data

Treebank-3

<i>Item Name(s):</i>	Treebank-3
<i>Author(s):</i>	Mitchell P. Marcus, Beatrice Santorini, Mary Ann Marcinkiewicz, Ann Taylor
<i>LDC Catalog No.:</i>	LDC99T42
<i>ISBN:</i>	1-58563-163-9
<i>ISLRN:</i>	141-282-691-413-2
<i>Member Year(s):</i>	1999
<i>DCMI Type(s):</i>	Text
<i>Data Source(s):</i>	telephone speech, newswire, microphone speech, transcribed speech, varied
<i>Project(s):</i>	TIDES, GALE
<i>Application(s):</i>	parsing, natural language processing, tagging
<i>Language(s):</i>	English
<i>Language ID(s):</i>	eng
<i>License(s):</i>	LDC User Agreement for Non-Members
<i>Online Documentation:</i>	LDC99T42 Documents
<i>Licensing Instructions:</i>	Subscription & Standard Members, and Non-Members
<i>Citation:</i>	Marcus, Mitchell, et al. Treebank-3 LDC99T42. Web Download. Philadelphia: Linguistic Data Consortium, 1999.
<i>Related Works:</i>	View

Introduction

This release contains the following [Treebank-2](#) Material:

- One million words of 1989 Wall Street Journal material annotated in Treebank II style.
- A small sample of ATIS-3 material annotated in Treebank II style.
- A fully tagged version of the Brown Corpus.

and the following new material:

- Switchboard tagged, dysfluency-annotated, and parsed text
- Brown parsed text

The Treebank bracketing style is designed to allow the extraction of simple predicate/argument structure. Over one million words of text are provided with this bracketing applied.

Data

The Penn Treebank

- Syntactic annotation of a million words of the 1989 Wall Street Journal, plus other corpora (released in 1993)
 - (Trivia: People often discuss “The Penn Treebank” when they mean the WSJ portion of it)

The Penn Treebank

- Syntactic annotation of a million words of the 1989 Wall Street Journal, plus other corpora (released in 1993)
 - (Trivia: People often discuss “The Penn Treebank” when they mean the WSJ portion of it)
- Contains 74 total tags: 36 parts of speech, 7 punctuation tags, and 31 phrasal constituent tags, plus some relation markings

The Penn Treebank

- Syntactic annotation of a million words of the 1989 Wall Street Journal, plus other corpora (released in 1993)
 - (Trivia: People often discuss “The Penn Treebank” when they mean the WSJ portion of it)
- Contains 74 total tags: 36 parts of speech, 7 punctuation tags, and 31 phrasal constituent tags, plus some relation markings
- Was the foundation for an entire field of research and applications for over twenty years

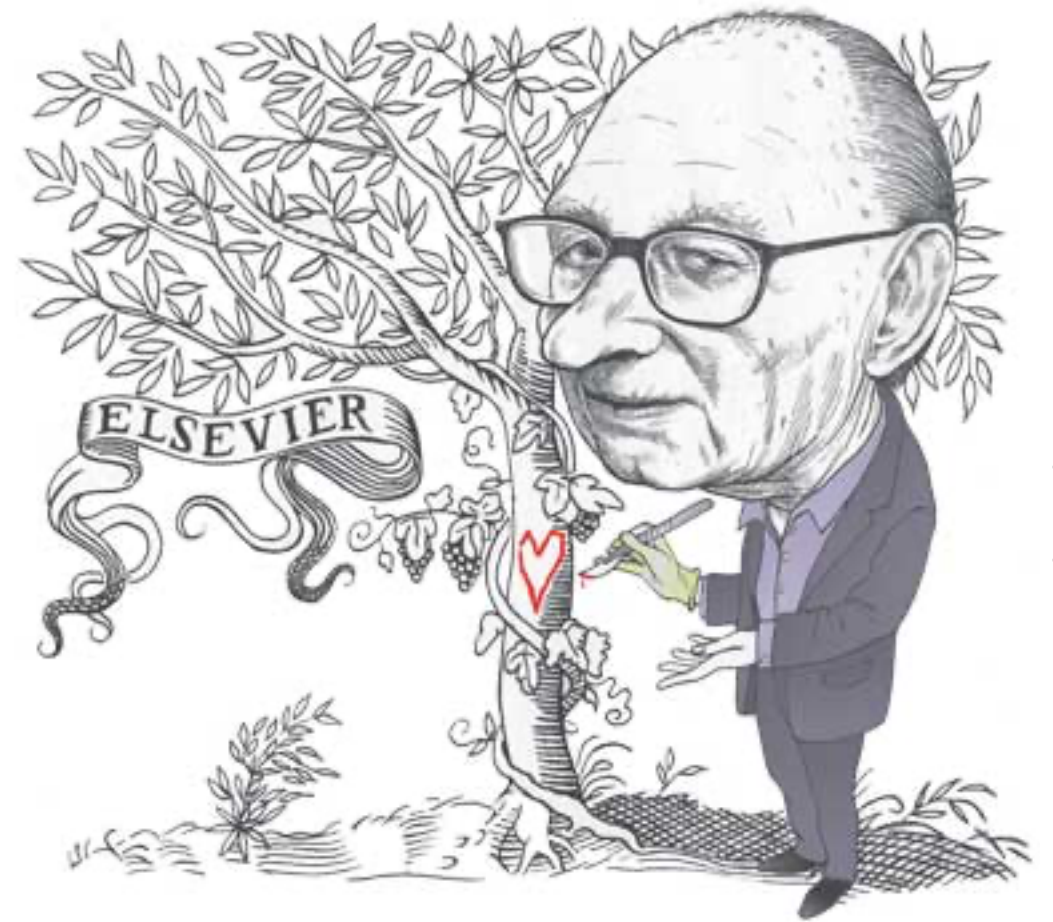
((S
 (NP-SBJ
 (NP (NNP Pierre) (NNP Vinken))
 (, ,)
 (ADJP
 (NP (CD 61) (NNS years))
 (JJ old))
 (, ,))
 (VP (MD will)
 (VP (VB join)
 (NP (DT the) (NN board))
 (PP-CLR (IN as)
 (NP (DT a) (JJ nonexecutive) (NN director)))
 (NP-TMP (NNP Nov.) (CD 29))))
 (. .)))



Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29.

((S
 (NP-SBJ
 (NP (NNP Pierre) (NNP Vinken))
 (, ,)
 (ADJP
 (NP (CD 61) (S years))
 (JJ old)
 (, ,))
 (VP (MD will ,
 (VP (VB join)
 (NP (DT the) (NN board))
 (PP-CLR (IN as)
 (NP (DT a) (JJ nonexecutive) (NN director)))
 (NP-TMP (NNP Nov.) (CD 29))))
 (. .)))

x 49,208

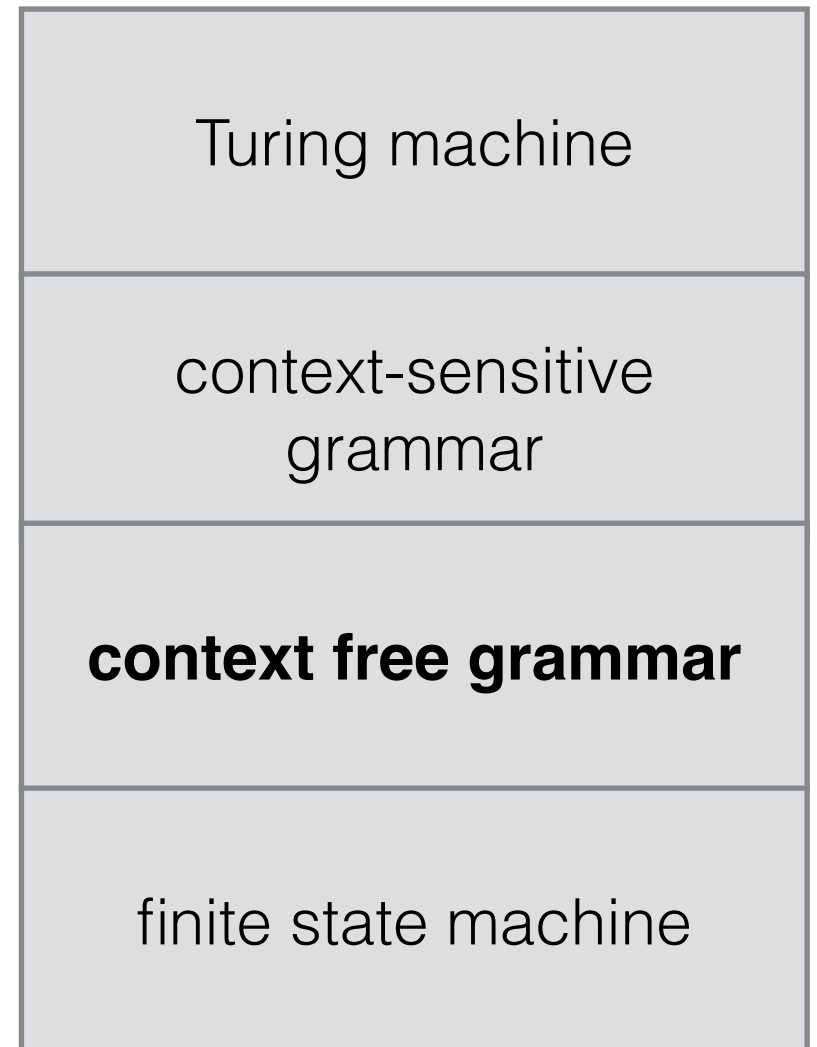


Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29.

Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone

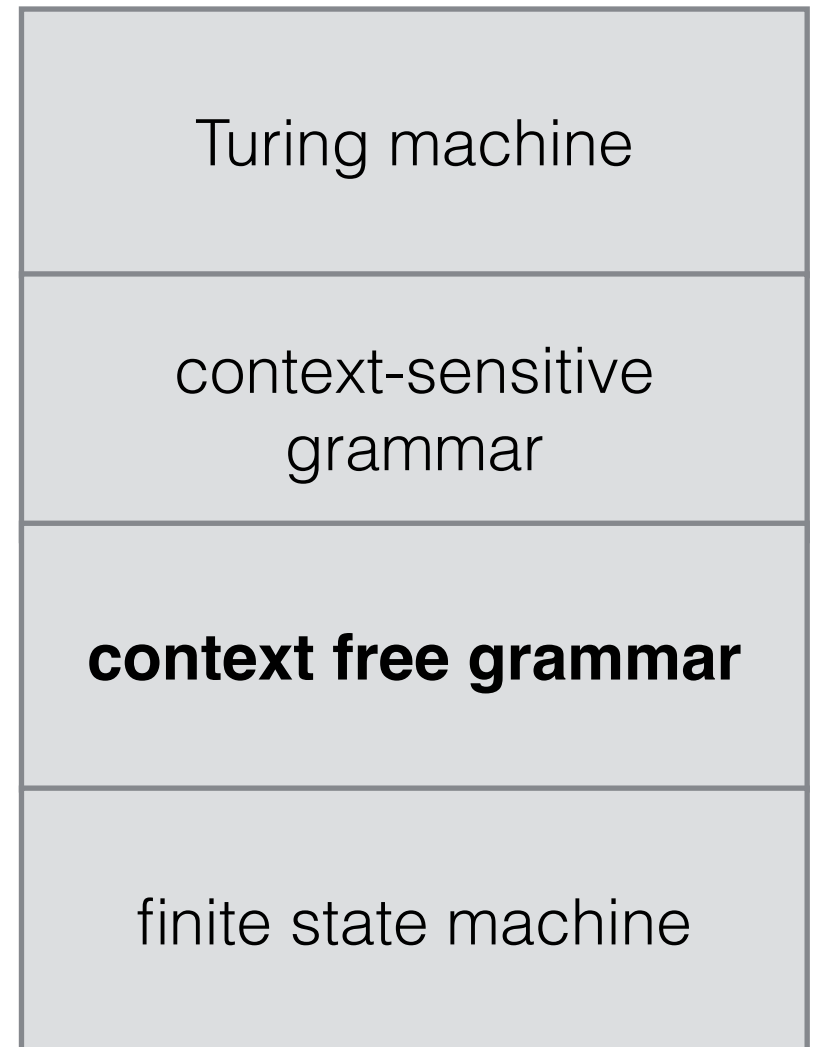
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:

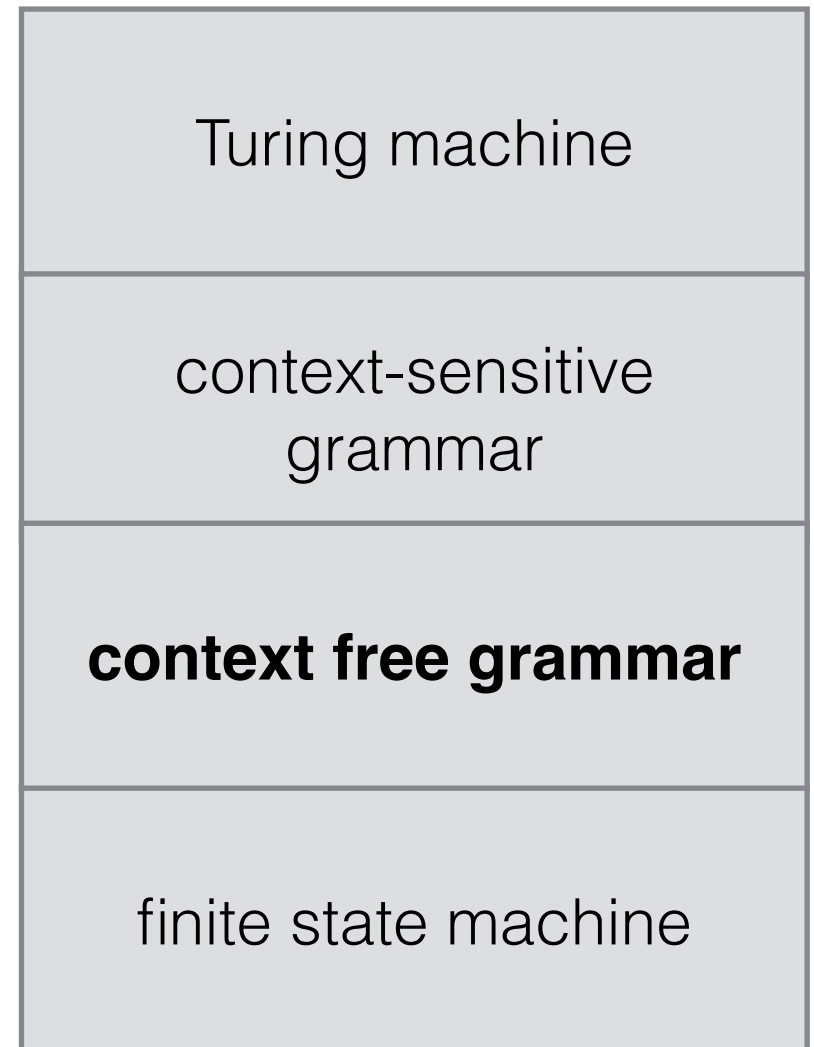
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP

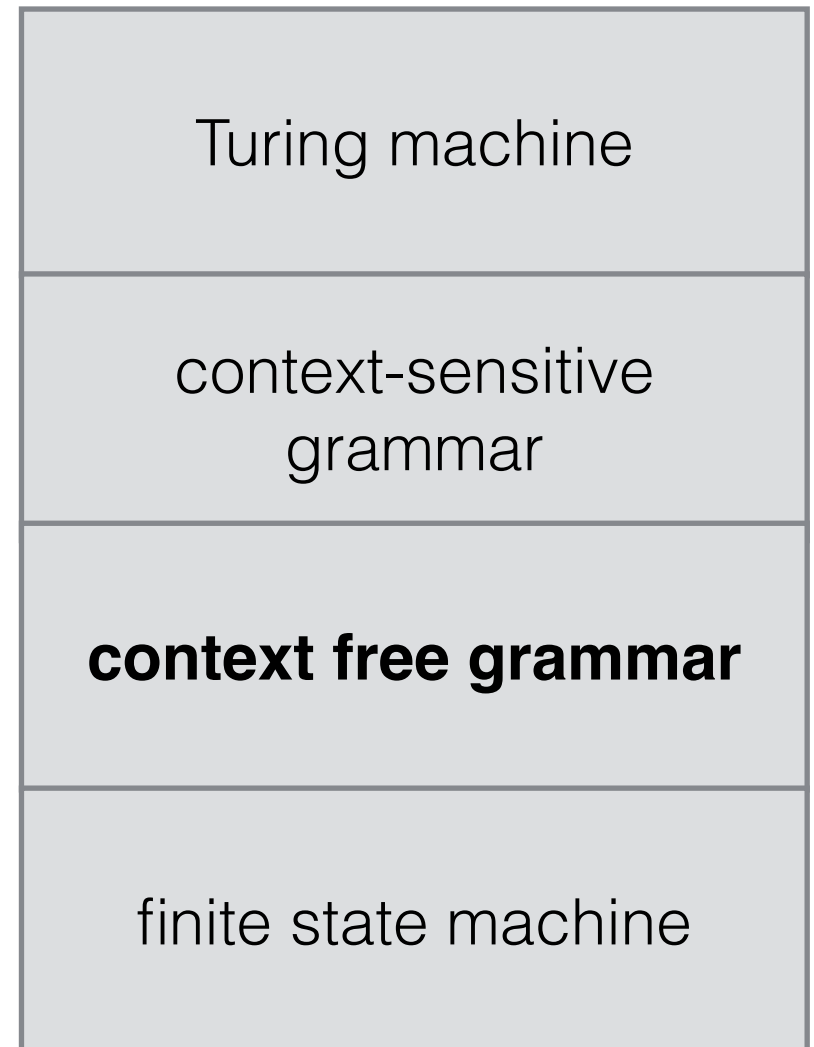
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP
 - For each leaf nonterminal:

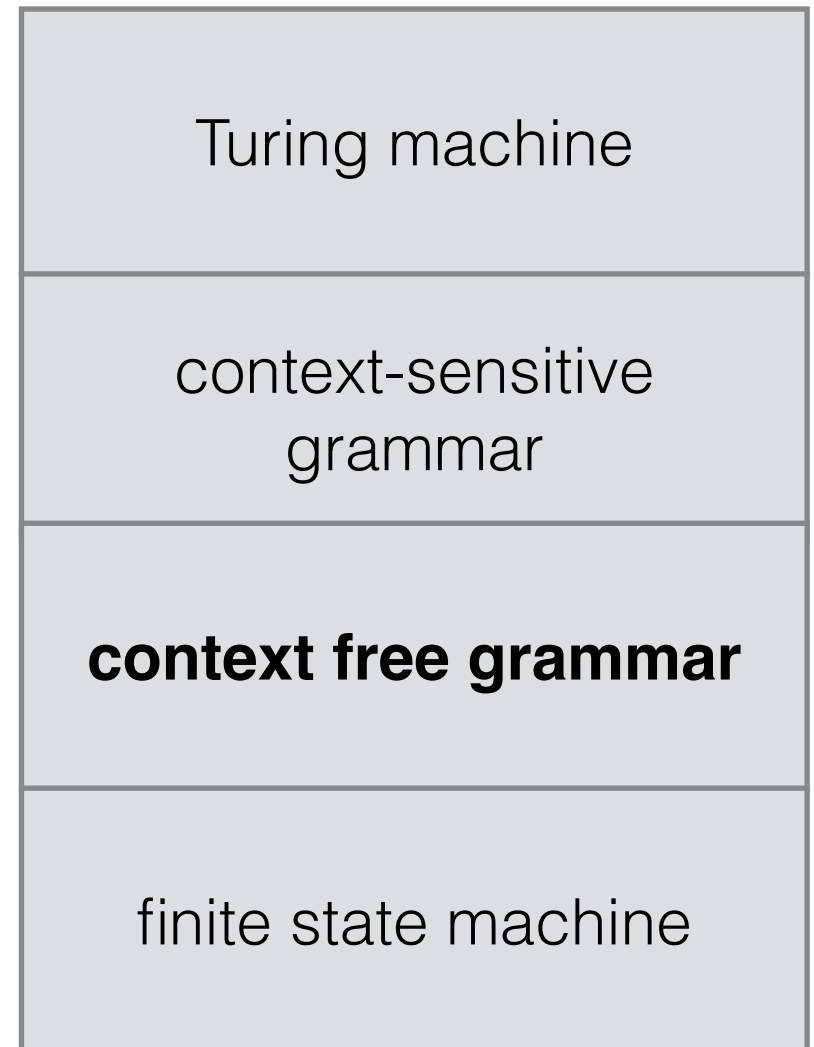
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP
 - For each leaf nonterminal:
 - Sample a rule from the set of rules for that nonterminal

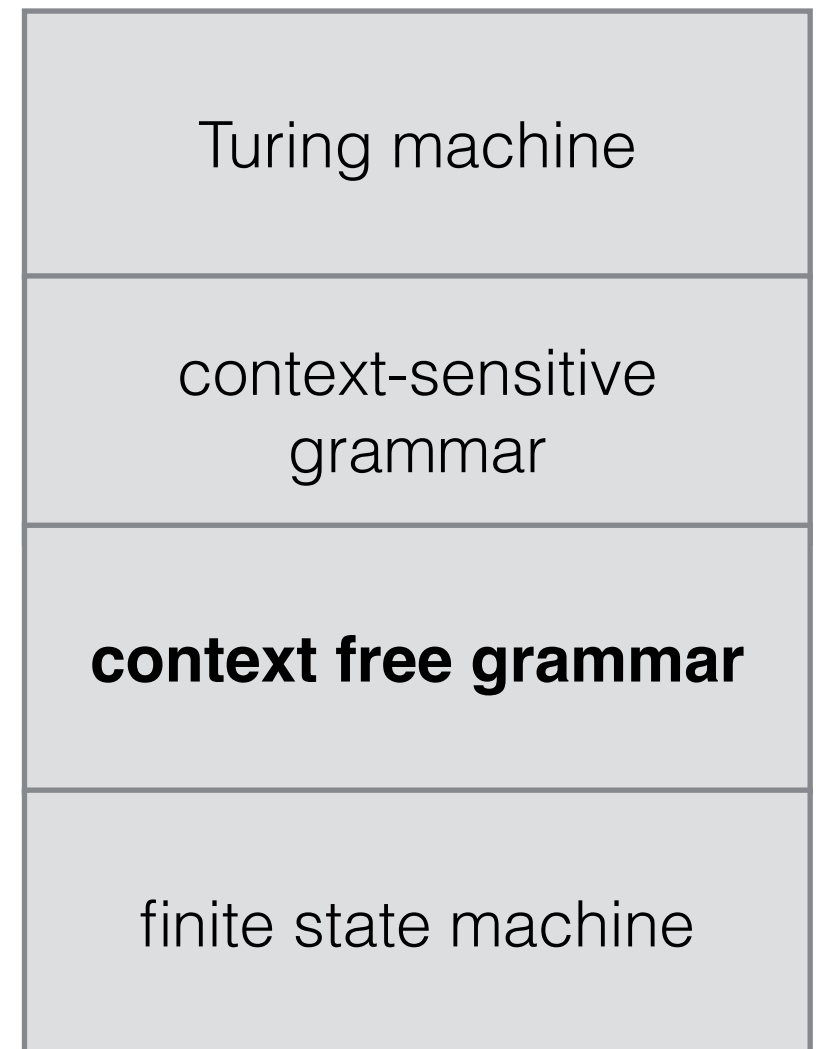
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP
 - For each leaf nonterminal:
 - Sample a rule from the set of rules for that nonterminal
 - Replace it with

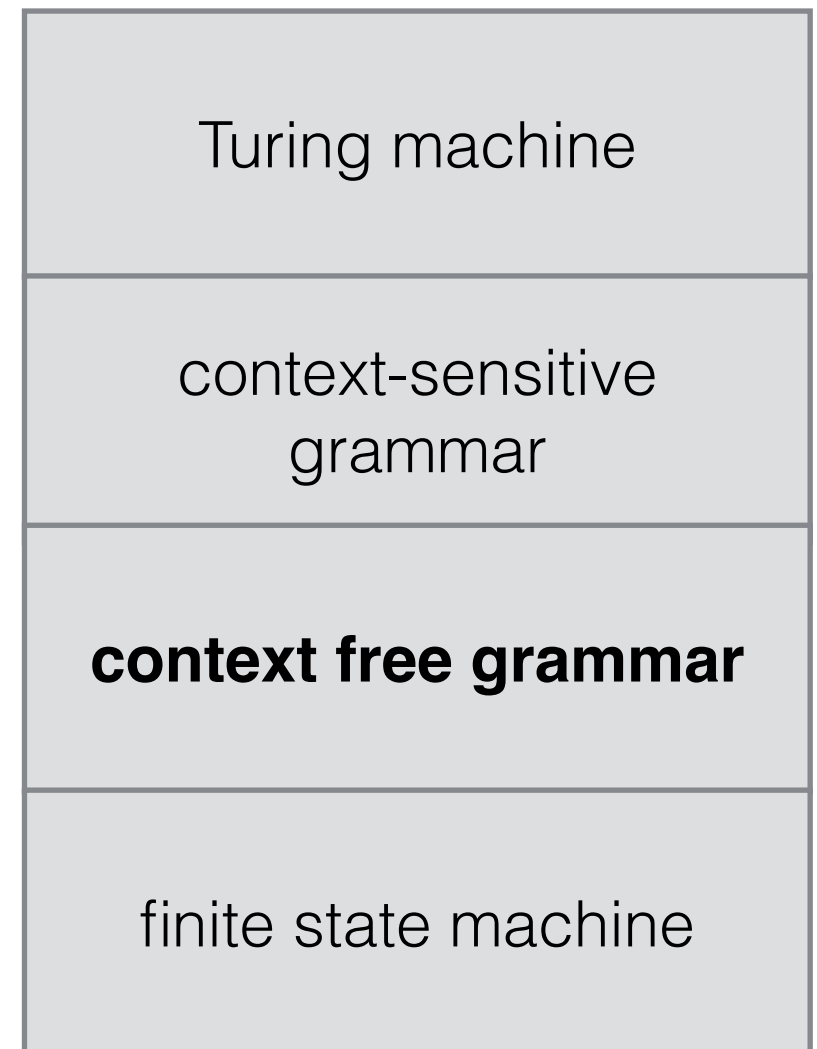
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP
 - For each leaf nonterminal:
 - Sample a rule from the set of rules for that nonterminal
 - Replace it with
 - Recurse

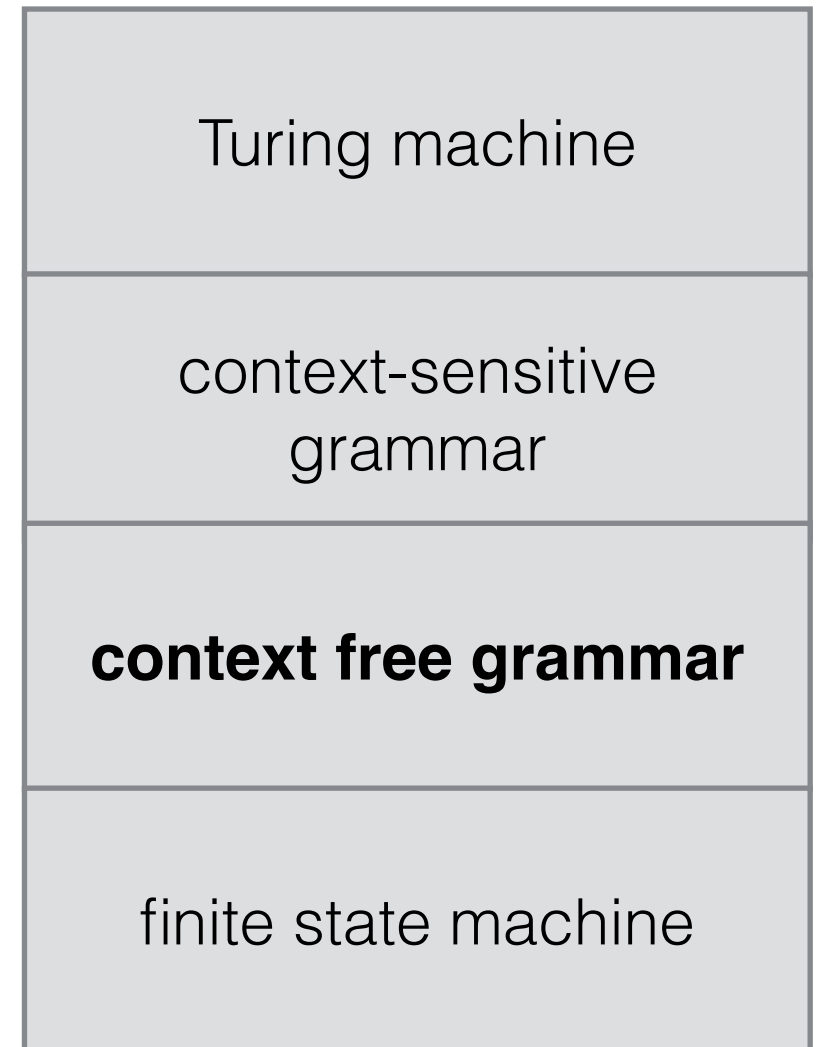
Chomsky formal language hierarchy



Context Free Grammar

- Nonterminals are rewritten based on the lefthand side alone
- Algorithm:
 - Start with TOP
 - For each leaf nonterminal:
 - Sample a rule from the set of rules for that nonterminal
 - Replace it with
 - Recurse
- Terminates when there are no more nonterminals

Chomsky formal language hierarchy



TOP

TOP \rightarrow S

TOP

S

TOP \rightarrow S

S \rightarrow VP

TOP

S

VP

TOP \rightarrow S

S \rightarrow VP

VP \rightarrow (VB \rightarrow halt) NP PP

TOP

S

VP

halt NP PP

TOP \rightarrow S

S \rightarrow VP

VP \rightarrow (VB \rightarrow halt) NP PP

NP \rightarrow (DT The)
(JJ \rightarrow market-jarring)
(CD \rightarrow 25)

TOP

S

VP

halt NP PP

halt **The market-jarring 25** PP

TOP \rightarrow S

S \rightarrow VP

VP \rightarrow (VB \rightarrow halt) NP PP

NP \rightarrow (DT The)
(JJ \rightarrow market-jarring)
(CD \rightarrow 25)

PP \rightarrow (IN \rightarrow at) NP

TOP

S

VP

halt NP PP

halt **The market-jarring 25** PP

halt The market-jarring 25 **at NP**

(NN→bond)

TOP → S

S → VP

VP → (VB→halt) NP PP

NP → (DT The)
(JJ→market-jarring)
(CD→25)

PP → (IN→at) NP

NP → (DT→the)

TOP

S

VP

halt NP PP

halt **The market-jarring 25** PP

halt The market-jarring 25 **at NP**
(NN→bond)

halt The market-jarring 25 at **the bond**

TOP → S

S → VP

VP → (VB→halt) NP PP

NP → (DT The)
(JJ→market-jarring)
(CD→25)

PP → (IN→at) NP

NP → (DT→the)

TOP

S

VP

halt NP PP

halt **The market-jarring 25** PP

halt The market-jarring 25 **at NP**
(NN→bond)

halt The market-jarring 25 **at the bond**

(TOP (S (VP (VB halt) (NP (DT The) (JJ market-jarring) (CD 25)) (PP (IN at) (NP (DT the) (NN bond))))))

TOP → S

S → VP

VP → (VB→halt) NP PP

NP → (DT The)
(JJ→market-jarring)
(CD→25)

PP → (IN→at) NP

NP → (DT→the)

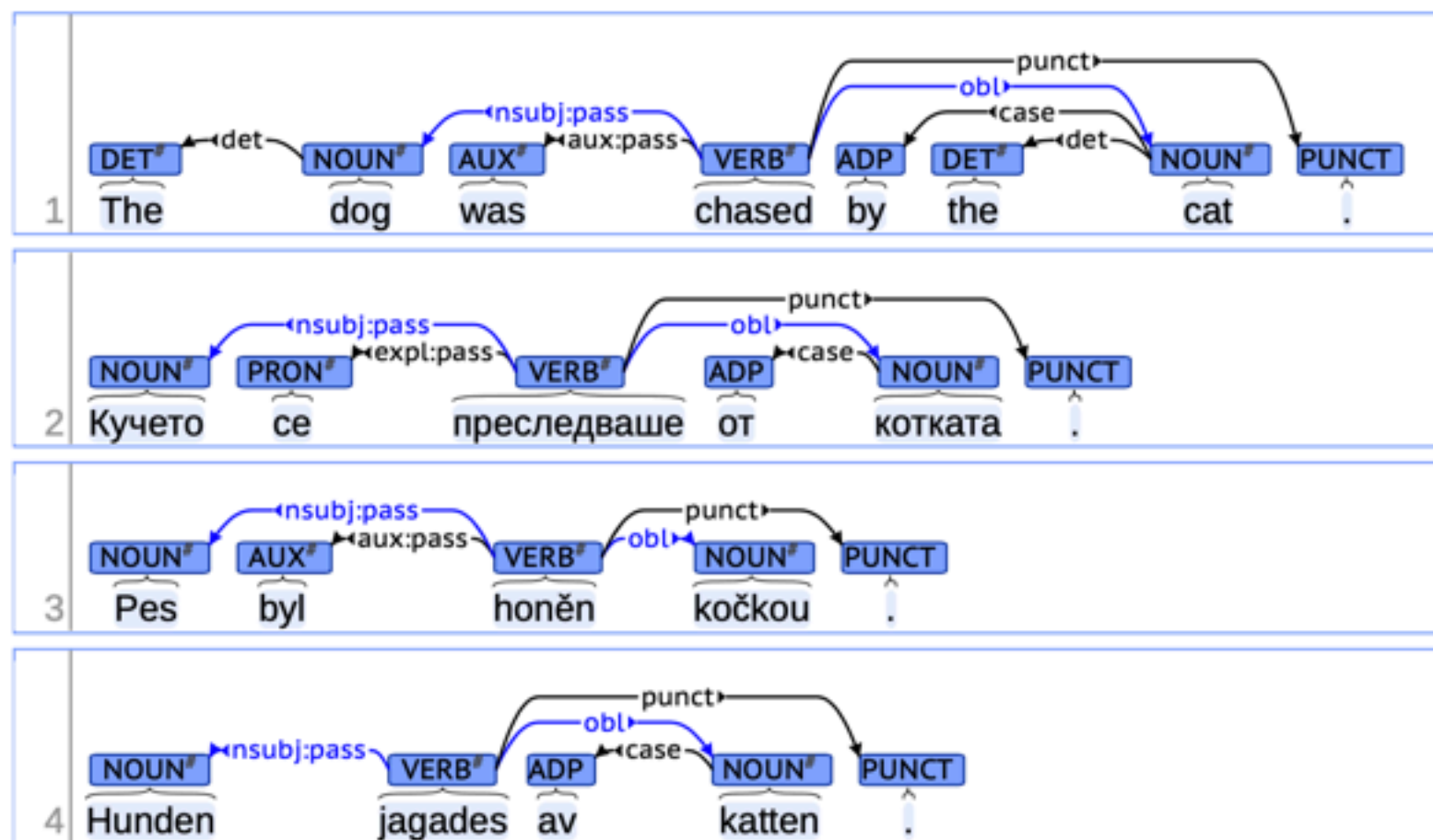
A problem with the Penn Treebank

- One language, English
 - Represents a very narrow typology (e.g., little morphology)
 - Consider the tags we looked at before
 - nouns: NN, NNS, NNP, NNPS
 - adverbs: RB, RBR, RBS, RP
 - verbs: VB, VBD, VBG, VBN, VBP, VBZ
 - How well will these generalize to other languages?

Dependency Treebanks (2012)

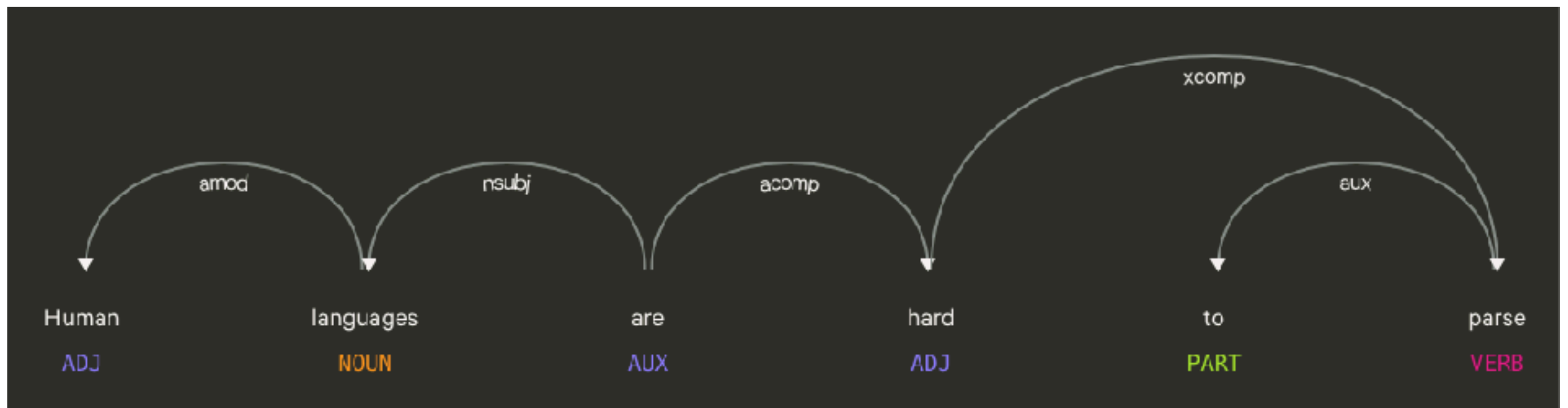
Universal Dependencies

- Dependency trees annotated across languages in a consistent manner



Example

- Instead of encoding phrase structure, it encodes dependencies between words
- Often more directly encodes information we care about (i.e., *who did what to whom*)



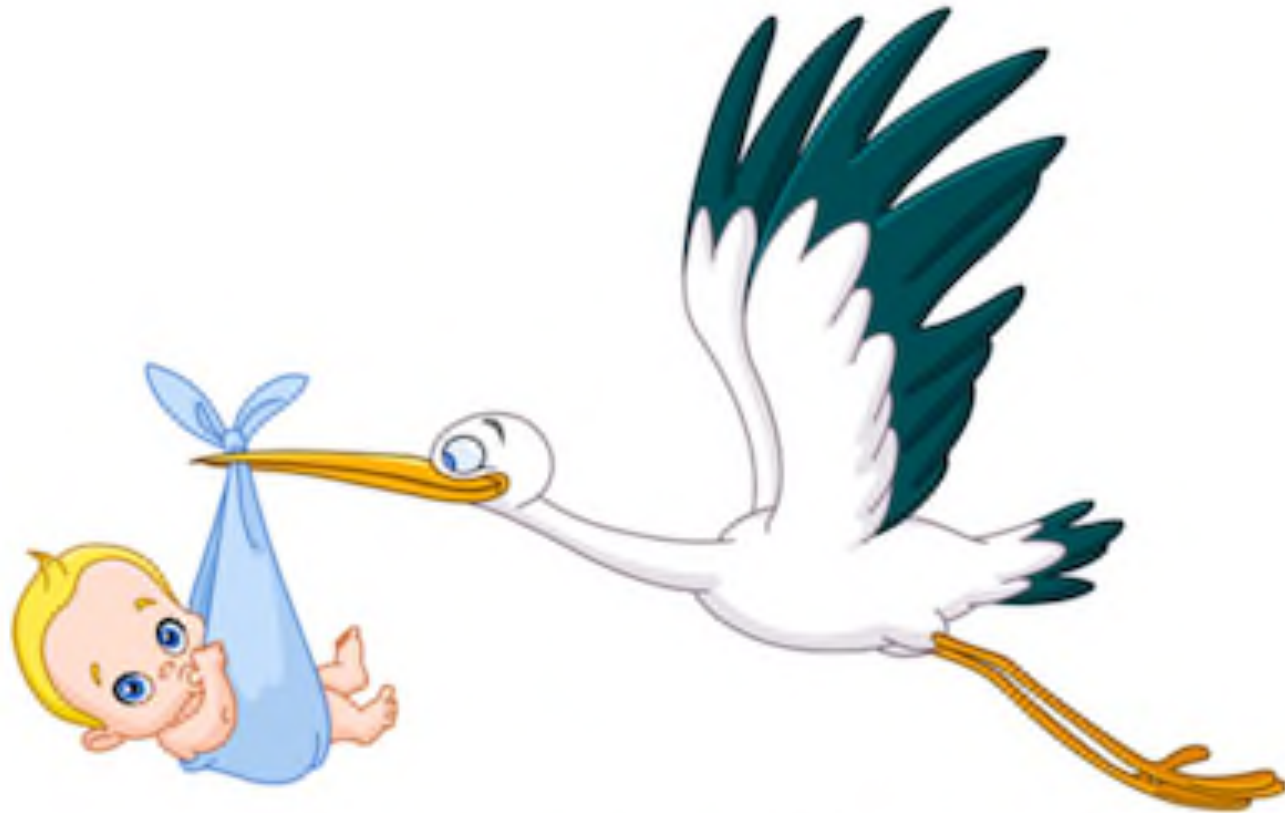
Guiding principles

- Works for individual languages
- Suitable across languages
- Easy to use when annotating
- Easy to parse quickly
- Understandable to laypeople
- Usable by downstream tasks

Universal Dependencies

- Parts of speech
 - open class
 - ADJ, ADV, INTJ, NOUN, PROPN, VERB
 - closed class
 - ADP, AUX, CCONJ, DET, NUM, PART, PRON, SCONJ
 - other
 - PUNCT, SYM, X

Where do grammars come from?



Where do grammars come from?

- Treebanks!
- Given a treebank, and a formalism, we can learn statistics by counting over the annotated instances

Probabilities

- For example, a context-free grammar

– $S \rightarrow NP , NP VP .$ [0.002]

– $NP \rightarrow NNP NNP$ [0.037]

– $, \rightarrow ,$ [0.999]

– $NP \rightarrow *$ [X]

– $VP \rightarrow VB NP$ [0.057]

– $NP \rightarrow PRP\$ NN$ [0.008]

– $. \rightarrow .$ [0.987]

- Probabilities given as $P(X) = \sum_{X' \in N} \frac{P(X)}{P(X')}$

Summary

Grammars are learned from Treebanks

where do
grammars
come from?

Treebanks are annotated according to a particular theory or formalism

Outline

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

Formal Language Theory

- Consider the claims underlying our grammar-based view of language
 1. Sentences are either in or out of a language
 2. Sentences have an invisible hidden structure

Formal Language Theory

- Consider the claims underlying our grammar-based view of language
 1. Sentences are either in or out of a language
 2. Sentences have an invisible hidden structure
- We can generalize this discussion to make a connection between natural and other kinds of languages

Formal Language Theory

- Consider the claims underlying our grammar-based view of language
 1. Sentences are either in or out of a language
 2. Sentences have an invisible hidden structure
- We can generalize this discussion to make a connection between natural and other kinds of languages
- Consider, for example, *computer programs*
 - They either compile or don't compile
 - Their structure determines their interpretation

Formal Language Theory

- Generalization: define a **language** to be a set of strings under some alphabet, Σ
 - e.g., the set of valid English sentences (where the “alphabet” is English words), or the set of valid Python programs

Formal Language Theory

- Generalization: define a **language** to be a set of strings under some alphabet, Σ
 - e.g., the set of valid English sentences (where the “alphabet” is English words), or the set of valid Python programs
- Formal Language Theory provides a common framework for studying properties of these languages, e.g.,
 - Is this file a valid C++ program? A valid Czech sentence?
 - What is the structure?
 - How hard / time-consuming is it to answer these questions?

The Chomsky Hierarchy

- Definitions: given
 - an alphabet (Σ),
 - terminal symbols, e.g., $a \in \Sigma$
 - nonterminal symbols, e.g., $\{S, N, A, B\}$
 - α, β, γ , strings of terminals and/or nonterminals

Type	Rules	Name	Recognized by
3	$A \rightarrow aB$	Regular	Regular expressions
2	$A \rightarrow \alpha$	Context-free	Pushdown automata
1	$\alpha A \beta \rightarrow \alpha \gamma \beta$	Context-sensitive	Linear-bounded Turing machine
0	$\alpha A \beta \rightarrow \gamma$	Recursively enumerable	Turing Machines

Problems

- What is the value?

$(5 + 7) * 11$

- Who did what to whom?

Him the Almighty hurled

Dipanjan taught Johnmark

*If we have a grammar, we can answer these with **parsing***₄₆

Parsing

- If the grammar has certain properties (Type 2 or 3), we can efficiently answer two questions with a **parser**
 - Is the sentence in the language of the parser?
 - What is the structure above that sentence?

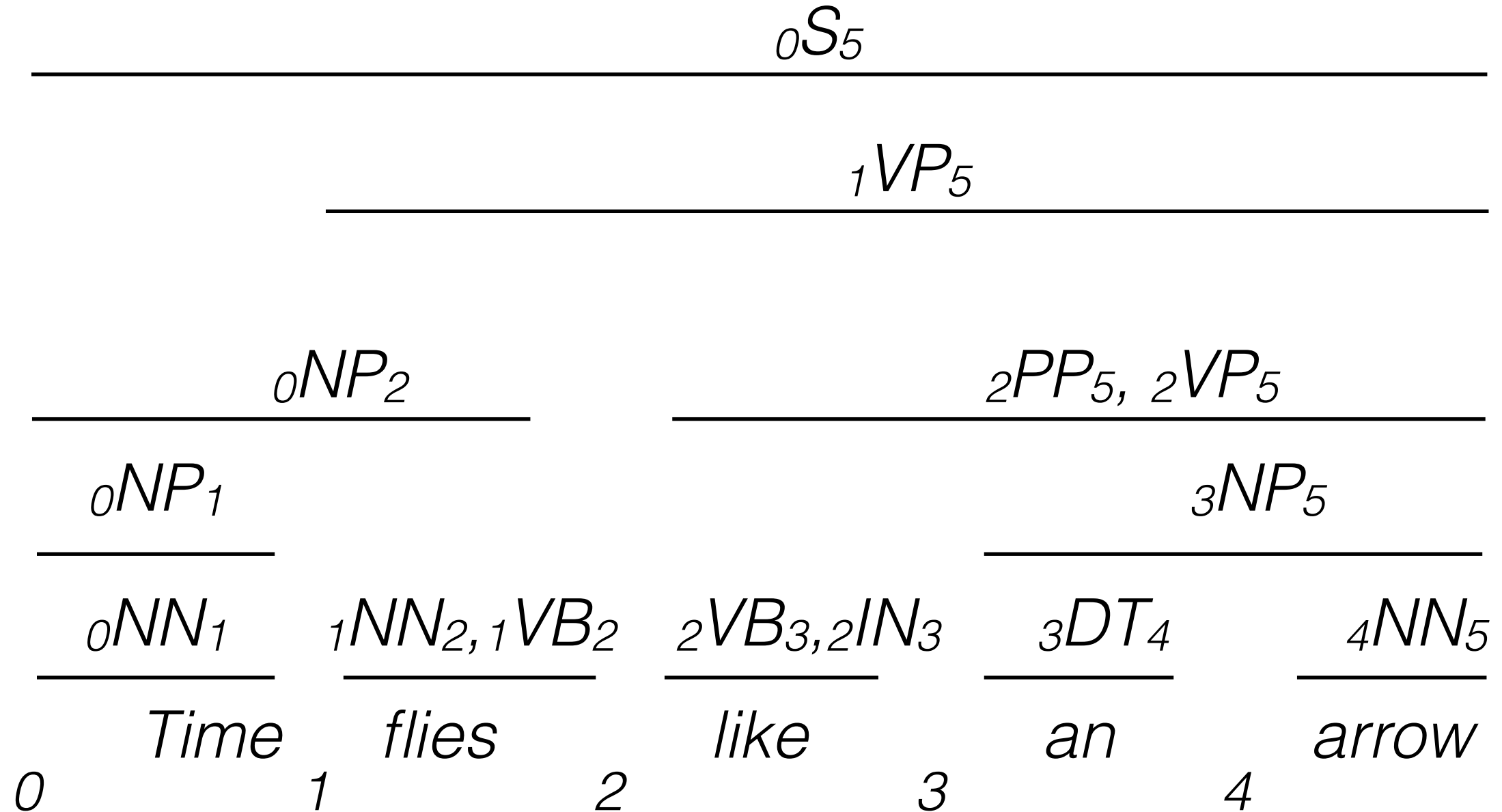
Algorithms

- The **CKY algorithm** for parsing with constituency grammars
- **Transition-based** parsing with dependency grammars

Chart parsing for constituency grammars

- Maintains a chart of nonterminals spanning words, e.g.,
 - NP over words 1..4 and 2..5
 - VP over words 4..6 and 4..8
 - etc

Chart parsing for constituency grammars



CKY algorithm

- How do we produce this chart? Cocke-Younger-Kasami (CYK/CKY)
- Basic idea is to apply rules in a bottom-up fashion, applying all rules, and (recursively) building larger constituents from smaller ones
- Input: sentence of length N
for width in $2..N$
 for begin i in $1..{N - width}$
 $j = i + width$
 for split k in ${i + 1}..{j - 1}$
 for all rules $A \rightarrow B C$
 create iA_j if iB_k and kC_j

CKY algorithm

0 *Time* 1 *flies* 2 *like* 3 *an* 4 *arrow* 5 52

CKY algorithm

$\frac{NN}{Time}$ $\frac{NN,VB}{flies}$ $\frac{VB,IN}{like}$ $\frac{DT}{an}$ $\frac{NN}{arrow}$

0 1 2 3 4 5

CKY algorithm

$NP \rightarrow NN$

$NP \rightarrow DT NN$

NN

NN, VB

VB, IN

DT

NN

Time

flies

like

an

arrow

0

1

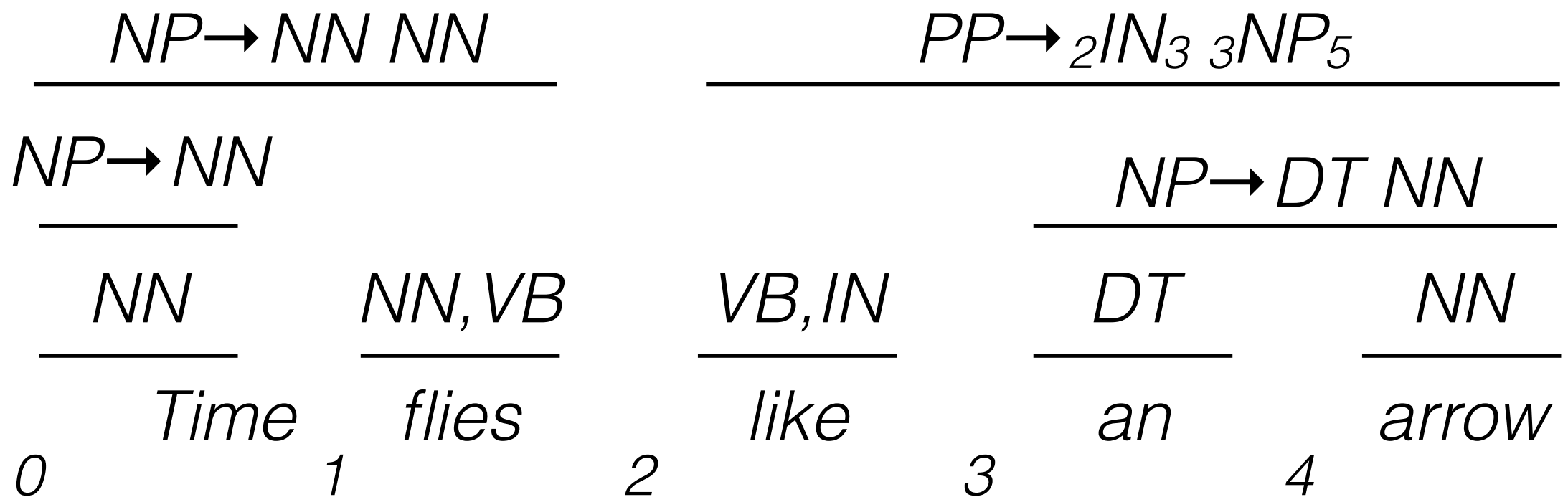
2

3

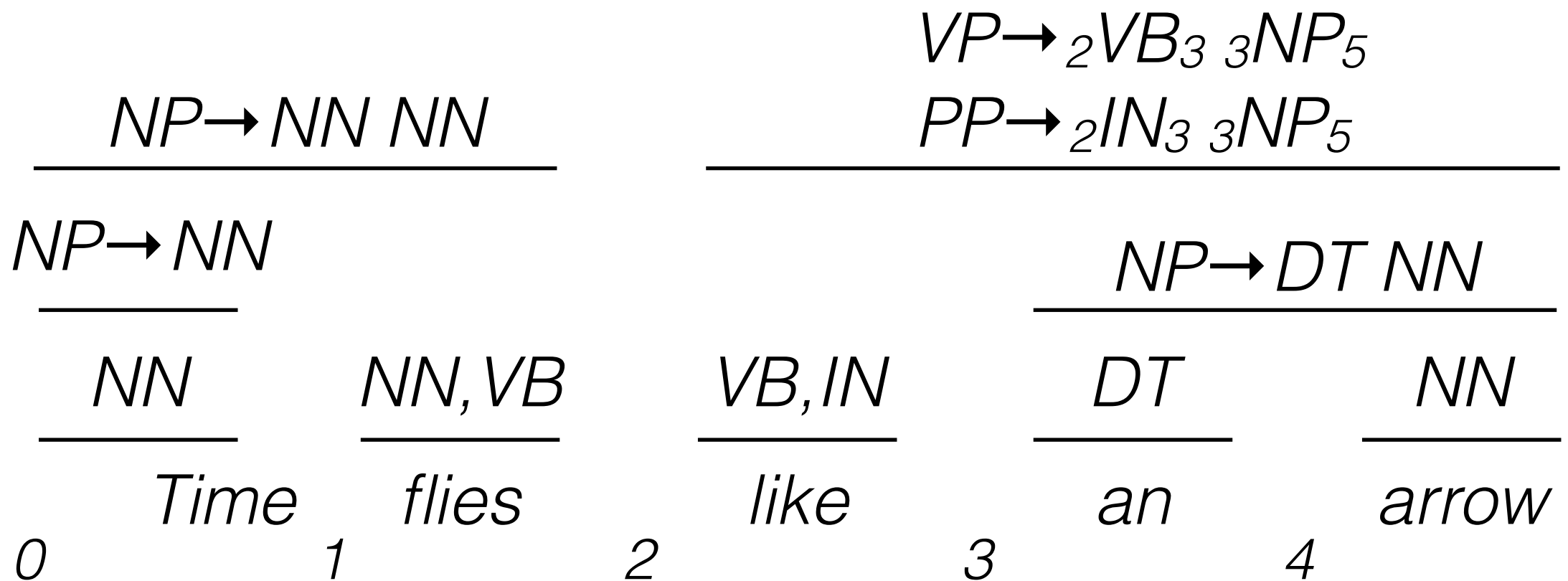
4

5

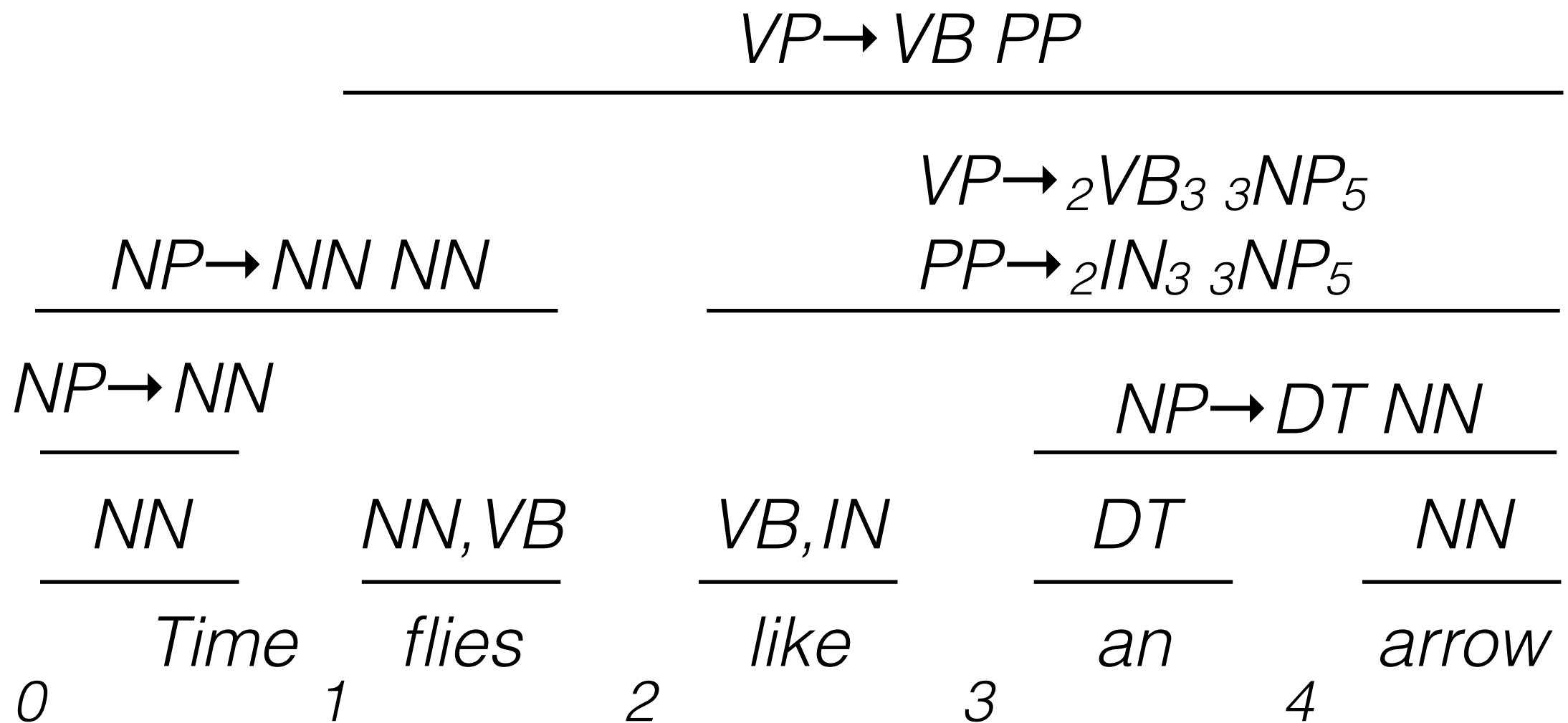
CKY algorithm



CKY algorithm



CKY algorithm



CKY algorithm

$S \rightarrow {}_0NP_1 {}_1VP_5$

$VP \rightarrow VB PP$

$VP \rightarrow {}_2VB_3 {}_3NP_5$

$PP \rightarrow {}_2IN_3 {}_3NP_5$

$NP \rightarrow NN NN$

$NP \rightarrow NN$

$NP \rightarrow DT NN$

NN

NN, VB

VB, IN

DT

NN

Time

flies

like

an

arrow

0

1

2

3

4

5

CKY algorithm

$S \rightarrow {}_0NP_2 {}_2VP_5$

$S \rightarrow {}_0NP_1 {}_1VP_5$

$VP \rightarrow VB PP$

$VP \rightarrow {}_2VB_3 {}_3NP_5$

$PP \rightarrow {}_2IN_3 {}_3NP_5$

$NP \rightarrow NN NN$

$NP \rightarrow NN$

$NP \rightarrow DT NN$

NN

NN, VB

VB, IN

DT

NN

Time

flies

like

an

arrow

0

1

2

3

4

5

CKY algorithm

- Termination: is there a chart entry at $0S_N$?
 - ✓ string is in the language
 - Obtain the structure by following backpointers
 - Not covered: adding probabilities to rules to resolve ambiguities

Dependency parsing

- The situation is different in many ways
 - We're no longer building labeled constituents
 - Instead, we're searching for word dependencies

Dependency parsing

- The situation is different in many ways
 - We're no longer building labeled constituents
 - Instead, we're searching for word dependencies
- This is accomplished by a stack-based **transition parser**
 - Repeatedly (a) shift a word onto the stack or (b) create a LEFT or RIGHT dependency from the top two words

ROOT human languages are hard to parse

```
graph LR; ROOT --> are; human --> languages; are --> hard; hard --> to; to --> parse;
```

step

stack

words

action

relation

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	
8	[are,hard,to,parse]	[]	LEFTARC	to←parse

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	
8	[are,hard,to,parse]	[]	LEFTARC	to←parse
9	[are,hard,parse]	[]	RIGHTARC	hard→parse

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	
8	[are,hard,to,parse]	[]	LEFTARC	to←parse
9	[are,hard,parse]	[]	RIGHTARC	hard→parse
10	[are,hard]	[]	RIGHTARC	are→hard

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	
8	[are,hard,to,parse]	[]	LEFTARC	to←parse
9	[are,hard,parse]	[]	RIGHTARC	hard→parse
10	[are,hard]	[]	RIGHTARC	are→hard
11	[are]	[]	RIGHTARC	ROOT→are

ROOT human languages are hard to parse

step	stack	words	action	relation
0	[]	[human,langs,are,hard,to,parse]	SHIFT	
1	[human]	[langs,are,hard,to,parse]	SHIFT	
2	[human,langs]	[are,hard,to,parse]	LEFTARC	human←langs
3	[langs]	[are,hard,to,parse]	SHIFT	
4	[langs,are]	[hard,to,parse]	LEFTARC	langs←are
5	[are]	[hard,to,parse]	SHIFT	
6	[are,hard]	[to,parse]	SHIFT	
7	[are,hard,to]	[parse]	SHIFT	
8	[are,hard,to,parse]	[]	LEFTARC	to←parse
9	[are,hard,parse]	[]	RIGHTARC	hard→parse
10	[are,hard]	[]	RIGHTARC	are→hard
11	[are]	[]	RIGHTARC	ROOT→are
12	[]	[]	DONE	

Unanswered questions

- How do we score rules (for constituency parsing) and actions and relations (for dependency parsing)?
 - Probabilities can be read from Treebanks
 - Actions can be informed by feature selection
- How do we know the right path to take?
 - We can try multiple paths using beam search
 - We get lots of savings via dynamic programming

Summary

For context-free grammars, the (weighted) CKY algorithm can be used to find the most probable (maximum a posteriori) tree given a certain grammar

For dependency grammars, the most popular approach is a variation of transition-based parsers

how can a computer find a sentence's structure?

Resources

- Demos:
 - AllenNLP: <https://demo.allennlp.org>
 - Berkeley Neural Parser: <https://parser.kitaev.io>
 - Spacy dependency parser: <https://explosion.ai/demos/displacy>

Outline

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

Outline

what is
syntax?

where do
grammars
come from?

how can a
computer find
a sentence's
structure?

the study of the
internal structure of
sentences (in natural
and synthetic
languages)

Outline

what is
syntax?

the study of the
internal structure of
sentences (in natural
and synthetic
languages)

where do
grammars
come from?

they are created by
linguists, usually
under particular
grammatical theories

how can a
computer find
a sentence's
structure?

Outline

what is
syntax?

the study of the
internal structure of
sentences (in natural
and synthetic
languages)

where do
grammars
come from?

they are created by
linguists, usually
under particular
grammatical theories

how can a
computer find
a sentence's
structure?

train a grammar from
a treebank and then
apply that grammar
to new sentences
using parsing
algorithms