

Lecture 9: Formal grammars of English

Julia Hockenmaier

juliahmr@illinois.edu

3324 Siebel Center

Office Hours: Wednesday, 12:15-1:15pm

Previous key concepts

NLP tasks dealing with *words*...

- POS-tagging, morphological analysis

... require *finite-state* representations,

- Finite-State Automata and Finite-State Transducers

... the corresponding probabilistic models,

- Probabilistic FSAs and Hidden Markov Models
- Estimation: relative frequency estimation, EM algorithm

... and appropriate search algorithms

- Dynamic programming: Forward, Viterbi, Forward-Backward

The next key concepts

NLP tasks dealing with *sentences*...

- Syntactic parsing and semantic analysis

... require (at least) *context-free* representations,

- Context-free grammars, unification grammars

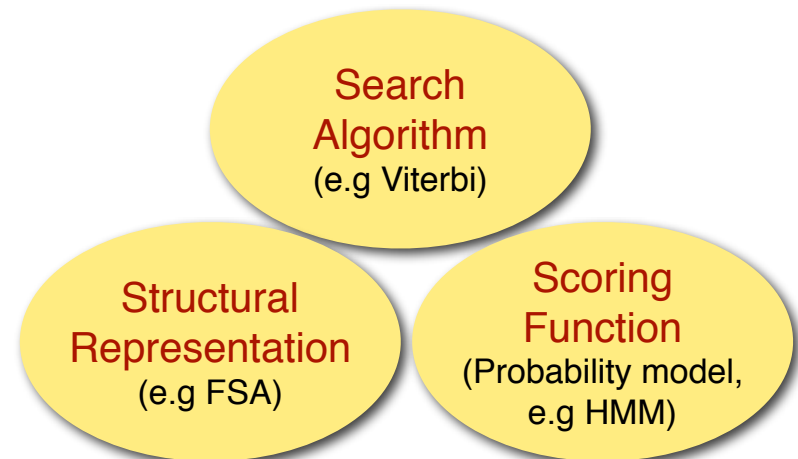
... the corresponding probabilistic models,

- Probabilistic Context-Free Grammars, Loglinear models
- Estimation: Relative Frequency estimation, EM algorithm, etc.

... and appropriate search algorithms

- Dynamic programming: chart parsing, inside-outside algorithm

Dealing with ambiguity

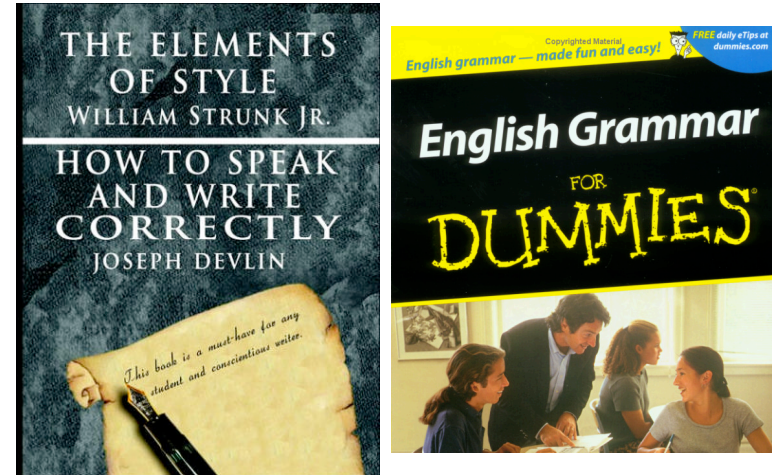


Today's lecture

Introduction to natural language syntax ('grammar'):

- Constituency and dependencies
- Context-free Grammars
- Dependency Grammars
- A simple CFG for English

What is grammar?



What is grammar?

Grammar formalisms

(= linguists' programming languages)

A precise way to define and describe the structure of sentences.

(N.B.: There are many different formalisms out there, which each define their own data structures and operations)

Specific grammars

(= linguists' programs)

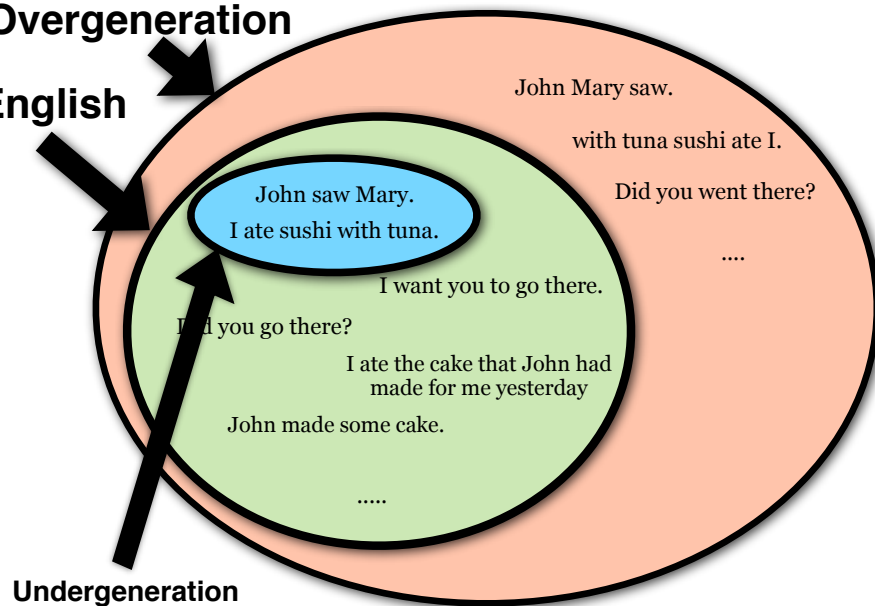
Implementations (in a particular formalism) for a particular language (English, Chinese,...)

Can we define a program that generates all English sentences?

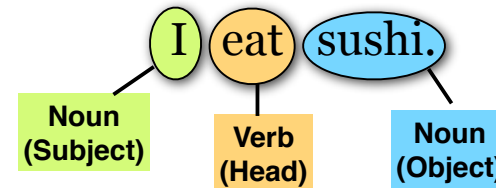
The number of sentences is infinite.
But we need our program to be finite.

Overgeneration

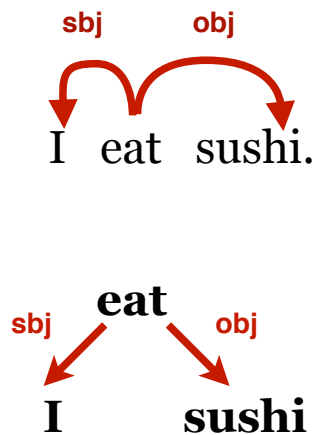
English



Basic sentence structure



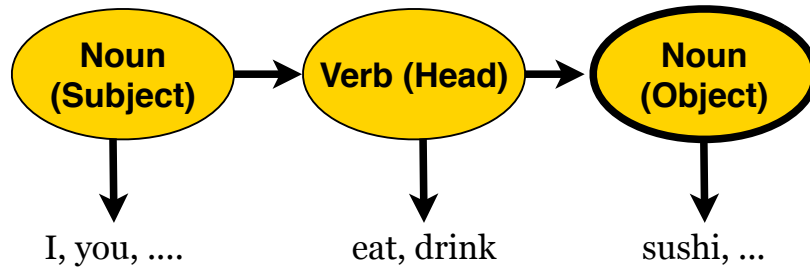
This is a dependency graph:



A finite-state-automaton (FSA)



A Hidden Markov Model (HMM)



Words take arguments

I eat sushi. ✓
I eat sushi you. ???
I sleep sushi ???
I give sushi ???
I drink sushi ?

Subcategorization:

Intransitive verbs (sleep) take only a subject.

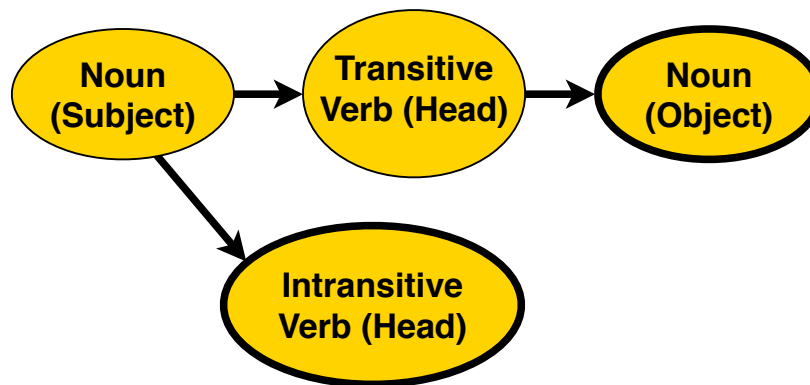
Transitive verbs (eat) take also one (direct) object.

Ditransitive verbs (give) take also one (indirect) object.

Selectional preferences:

The object of eat should be edible.

A better FSA



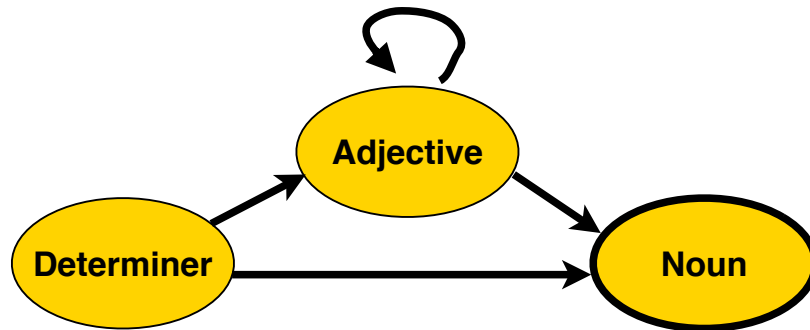
Language is recursive

the ball
*the **big** ball*
*the **big, red** ball*
*the **big, red, heavy** ball*
....

Adjectives can **modify** nouns.

The **number of modifiers/adjunctions** a word can have is (in theory) **unlimited**.

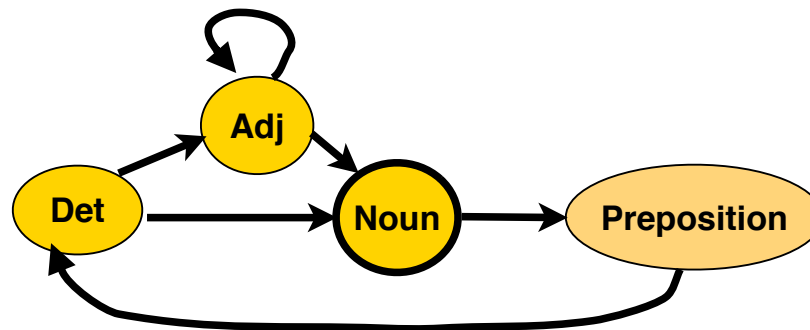
Another FSA



Recursion can be more complex

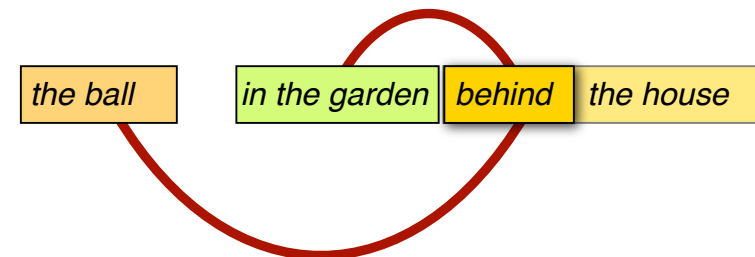
the ball
the ball in the garden
the ball in the garden behind the house
the ball in the garden behind the house next to the school
....

Yet another FSA

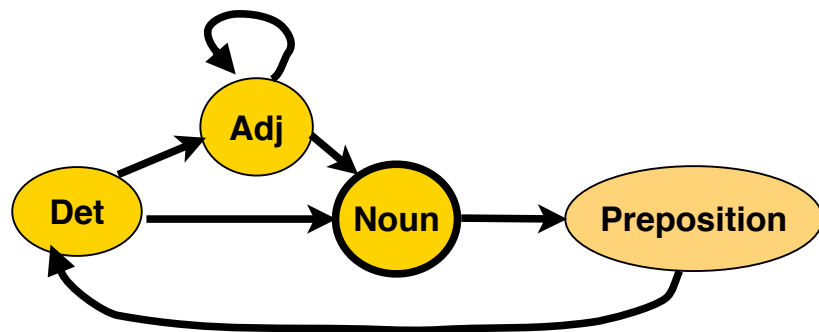


So, why do we need anything
beyond *regular grammars*?

What does this *mean*?



FSAs do not generate hierarchical structure

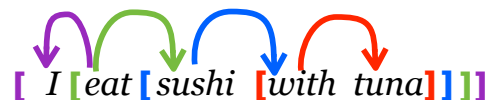


What is the structure of a sentence?

Sentence structure is hierarchical:

A sentence consists of **words** (*I, eat, sushi, with, tuna*)
..which form phrases or **constituents**: “*sushi with tuna*”

Sentence structure defines dependencies between words or phrases:



Strong vs. weak generative capacity

Formal language theory:

- defines language as string sets
- is only concerned with generating these strings (**weak generative capacity**)

Formal/Theoretical syntax (in linguistics):

- defines language as sets of strings with (hidden) structure
- is also concerned with generating the right structures (**strong generative capacity**)

Context-free grammars (CFGs) capture recursion

Language has complex **constituents**

(“the garden behind the house”)

Syntactically, these constituents behave just like simple ones.

(“behind the house” can always be omitted)

CFGs define **nonterminal categories** to capture equivalent constituents.

Context-free grammars

A CFG is a 4-tuple $\langle N, \Sigma, R, S \rangle$ consisting of:

A set of nonterminals N

(e.g. $N = \{S, NP, VP, PP, Noun, Verb, \dots\}$)

A set of terminals Σ

(e.g. $\Sigma = \{I, you, he, eat, drink, sushi, ball, \}$)

A set of rules R

$R \subseteq \{A \rightarrow \beta \text{ with left-hand-side (LHS) } A \in N$
and right-hand-side (RHS) $\beta \in (N \cup \Sigma)^* \}$

A start symbol $S \in N$

An example

DT $\rightarrow \{\text{the, a}\}$

N $\rightarrow \{\text{ball, garden, house, sushi}\}$

P $\rightarrow \{\text{in, behind, with}\}$

NP $\rightarrow \text{DT N}$

NP $\rightarrow \text{NP PP}$

PP $\rightarrow \text{P NP}$

N: noun

P: preposition

NP: “noun phrase”

PP: “prepositional phrase”

CFGs define parse trees

N $\rightarrow \{\text{sushi, tuna}\}$

P $\rightarrow \{\text{with}\}$

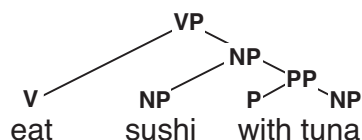
V $\rightarrow \{\text{eat}\}$

NP $\rightarrow \text{N}$

NP $\rightarrow \text{NP PP}$

PP $\rightarrow \text{P NP}$

VP $\rightarrow \text{V NP}$



CFGs and center embedding

The mouse ate the corn.

The mouse **that the snake ate** ate the corn.

The mouse **that the snake that the hawk ate** ate the corn.

....

These sentences are all grammatical.

They can be generated by a CFG:

S $\rightarrow \text{NP VP}$

NP $\rightarrow \text{NP RelClause}$

RelClause $\rightarrow \text{that NP ate}$

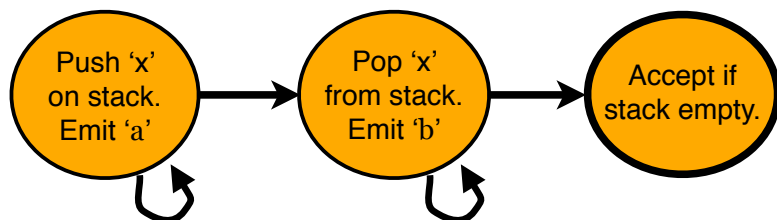
Linguists distinguish between a speaker's

- **competence** (grammatical knowledge) **and**

- **performance** (processing and memory limitations)

CFGs are equivalent to Pushdown automata (PDAs)

PDAs are FSAs with an additional stack:
Emit a symbol *and* push/pop a symbol from the stack



This is equivalent to the following CFG:

$$\begin{array}{ll} S \rightarrow a X b & S \rightarrow a b \\ X \rightarrow a X b & X \rightarrow a b \end{array}$$

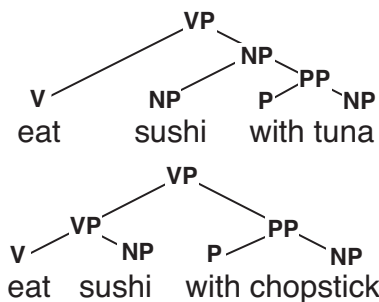
Generating $a^n b^n$

1. Push x on stack. Emit a.	Stack:	String:
2. Push x on stack. Emit a.	Stack: x	String: a
3. Push x on stack. Emit a.	Stack: xx	String: aa
4. Push x on stack. Emit a.	Stack: xxx	String: aaa
5. Pop x off stack. Emit b.	Stack: xxx	String: aaa
6. Pop x off stack. Emit b.	Stack: xx	String: aaa
7. Pop x off stack. Emit b.	Stack: x	String: aab
8. Pop x off stack. Emit b.	Stack:	String: aabb

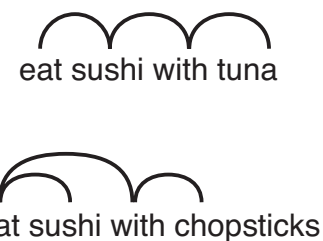
Defining grammars for natural language

Two ways to represent structure

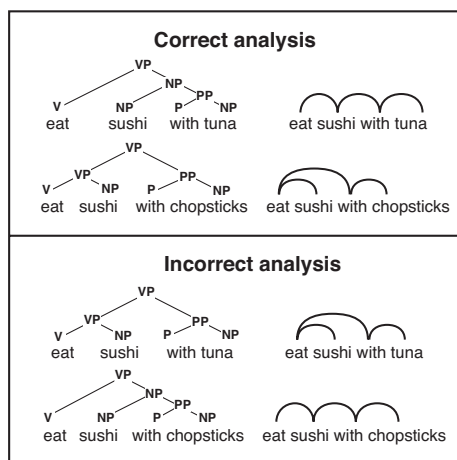
Phrase structure trees



Dependency trees



Structure (syntax) corresponds to meaning (semantics)



Dependency grammar

DGs describe the structure of sentences as a graph.

The **nodes** of the graph are the **words**

The **edges** of the graph are the **dependencies**.

Note: the relationship between DG and CFGs:

If a CFG phrase structure tree is translated into DG, the resulting dependency graph has **no crossing edges**.

Constituents: Heads and dependents

There are different kinds of constituents:

Noun phrases: the man, a girl with glasses, Illinois

Prepositional phrases: with glasses, in the garden

Verb phrases: eat sushi, sleep, sleep soundly

Every phrase has a **head**:

Noun phrases: the man, a girl with glasses, Illinois

Prepositional phrases: with glasses, in the garden

Verb phrases: eat sushi, sleep, sleep soundly

The other parts are its **dependents**.

Dependents are either **arguments** or **adjuncts**

Is string α a constituent?

He talks **[in class]**.

Substitution test:

Can α be replaced by a single word?

He talks **[there]**.

Movement test:

Can α be moved around in the sentence?

[In class], he talks.

Answer test:

Can α be the answer to a question?

Where does he talk? - **[In class]**.

Arguments are obligatory

Words subcategorize for specific sets of arguments:

Transitive verbs (sbj + obj): [John] likes [Mary]

All arguments have to be present:

*[John] likes. *likes [Mary].

No argument can be occupied multiple times:

*[John] [Peter] likes [Ann] [Mary].

Words can have multiple *subcat frames*:

Transitive *eat* (sbj + obj): [John] eats [sushi].

Intransitive *eat* (sbj): [John] eats.

Adjuncts are optional

Adverbs, PPs and adjectives can be adjuncts:

Adverbs: John runs [**fast**].

a [**very**] heavy book.

PPs: John runs [**in the gym**].

the book [**on the table**]

Adjectives: a [**heavy**] book

There can be an arbitrary number of adjuncts:

John saw Mary.

John saw Mary [**yesterday**].

John saw Mary [**yesterday**] [**in town**]

John saw Mary [**yesterday**] [**in town**] [**during lunch**]

[**Perhaps**] John saw Mary [**yesterday**] [**in town**] [**during lunch**]

A context-free grammar for a fragment of English

Noun phrases (NPs)

Simple NPs:

[**He**] sleeps. (pronoun)

[**John**] sleeps. (proper name)

[**A student**] sleeps. (determiner + noun)

Complex NPs:

[**A tall student**] sleeps. (det + adj + noun)

[**The student in the back**] sleeps. (NP + PP)

[**The student who likes MTV**] sleeps. (NP + Relative Clause)

The NP fragment

NP → Pronoun

NP → ProperName

NP → Det Noun

Det → {a, the, every}

Pronoun → {he, she,...}

ProperName → {John, Mary,...}

Noun → AdjP Noun

Noun → N

NP → NP PP

NP → NP RelClause

Adjective phrases (AdjP) and prepositional phrases (PP)

AdjP → Adj

AdjP → Adv AdjP

Adj → {big, small, red,...}

Adv → {very, really,...}

PP → P NP

P → {with, in, above,...}

The verb phrase (VP)

He [eats].

He [eats sushi].

He [gives John sushi].

He [eats sushi with chopsticks].

VP → V

VP → V NP

VP → V NP PP

VP → VP PP

V → {eats, sleeps gives,...}

Capturing subcategorization

He [eats]. ✓

He [eats sushi]. ✓

He [gives John sushi]. ✓

He [eats sushi with chopsticks]. ✓

*He [eats John sushi]. ???

VP → V_{intrans}

VP → V_{trans} NP

VP → V_{ditrans} NP NP

VP → VP PP

V_{intrans} → {eats, sleeps}

V_{trans} → {eats}

V_{trans} → {gives}

Sentences

[He eats sushi].
[Sometimes, he eats sushi].
[In Japan, he eats sushi].

$S \rightarrow NP VP$
 $S \rightarrow AdvP S$
 $S \rightarrow PP S$

He says [he eats sushi].
 $VP \rightarrow V_{comp} S$
 $V_{comp} \rightarrow \{\text{says, think, believes}\}$

Sentences redefined

[He eats sushi]. ✓
*[I eats sushi]. ???
*[They eats sushi]. ???

$S \rightarrow NP_{3sg} VP_{3sg}$
 $S \rightarrow NP_{1sg} VP_{1sg}$
 $S \rightarrow NP_{3pl} VP_{3pl}$

We need features to capture agreement:
(number, person, case,...)

Complex VPs

In English, **simple tenses** have separate forms:

present tense: the girl eats sushi
simple past tense: the girl ate sushi

Complex tenses, progressive aspect and passive voice consist of auxiliaries and participles:

past perfect tense: the girl has eaten sushi
future perfect: the girl will have eaten sushi
passive voice: the sushi was eaten by the girl
progressive: the girl is/was/will be eating sushi

VPs redefined

He [has [eaten sushi]].
The sushi [was [eaten by him]].

$VP \rightarrow V_{have} VP_{pastPart}$
 $VP \rightarrow V_{be} VP_{pass}$
 $VP_{pastPart} \rightarrow V_{pastPart} NP$
 $VP_{pass} \rightarrow V_{pastPart} PP$
 $V_{have} \rightarrow \{\text{has}\}$
 $V_{pastPart} \rightarrow \{\text{eaten, seen}\}$

We need more nonterminals (e.g. $VP_{pastpart}$).
N.B.: We call $VP_{pastPart}$, VP_{pass} , etc. 'untensed' VPs

Coordination

[He eats sushi] and [she drinks tea]
[John] and [Mary] eat sushi.
He [eats sushi] and [drinks tea]

S → S conj S
NP → NP conj NP
VP → VP conj VP

He says [he eats sushi].
VP → V_{comp} S
V_{comp} → {says, think, believes}

Relative clauses

Relative clauses **modify a noun phrase**:
the girl [that eats sushi]

Relative clauses **lack a noun phrase**, which is understood to be filled by the NP they modify:
'the girl that eats sushi' implies 'the girl eats sushi'

There are **subject** and **object relative clauses**:
subject: 'the girl that eats sushi'
object: 'the sushi that the girl eats'

Yes/No questions

Yes/no questions consist of an auxiliary, a subject and an (untensed) verb phrase:

does she eat sushi?
have you eaten sushi?

YesNoQ → Aux NP VP_{inf}
YesNoQ → Aux NP VP_{pastPart}

Wh-questions

Subject wh-questions consist of an wh-word, an auxiliary and an (untensed) verb phrase:

Who has eaten the sushi?

Object wh-questions consist of an wh-word, an auxiliary, an NP and an (untensed) verb phrase:

What does Mary eat?

Today's key concepts

Natural language syntax

Constituents

Dependencies

Context-free grammar

Arguments and modifiers

Recursion in natural language

Today's reading

Textbook:

Jurafsky and Martin, Chapter 12, sections 1-7