

CS447: Natural Language Processing

<http://courses.engr.illinois.edu/cs447>

Lecture 14:

Formal grammars

of English

Julia Hockenmaier

juliahmr@illinois.edu

3324 Siebel Center

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

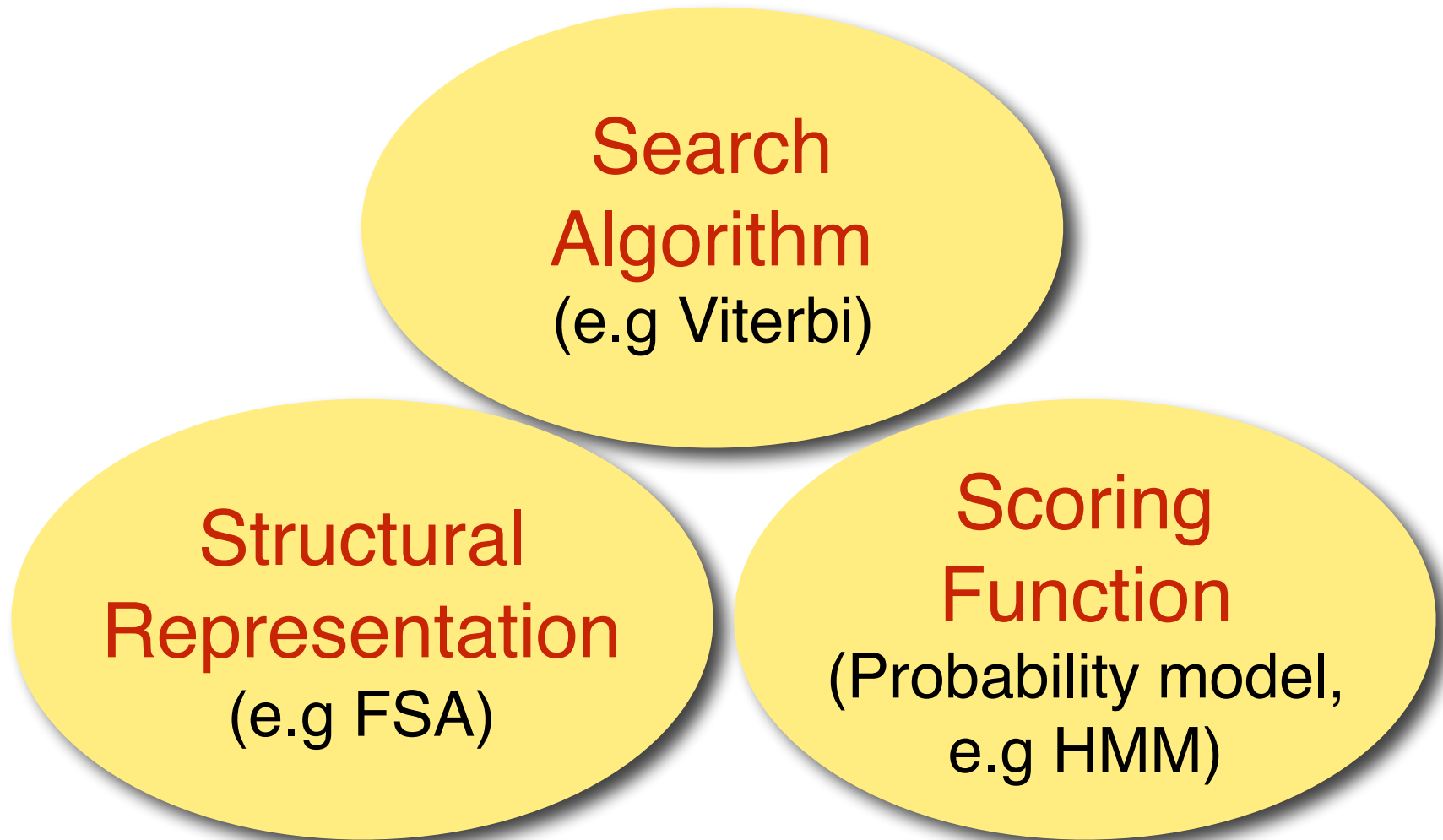
... the corresponding **probabilistic models**,

- 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'):

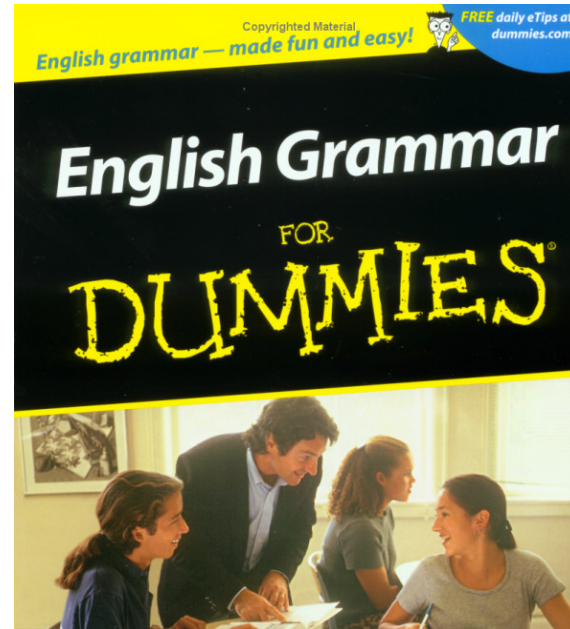
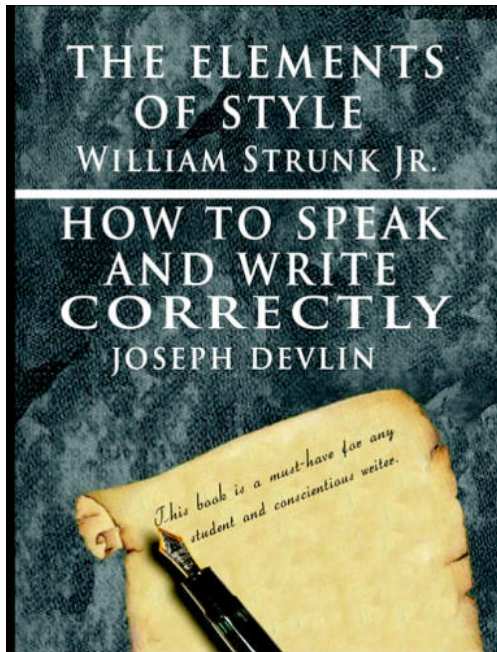
- Constituency and dependencies

- Context-free Grammars

- Dependency Grammars

- A simple CFG for English

What is grammar?



No, not
really, not in
this class

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)

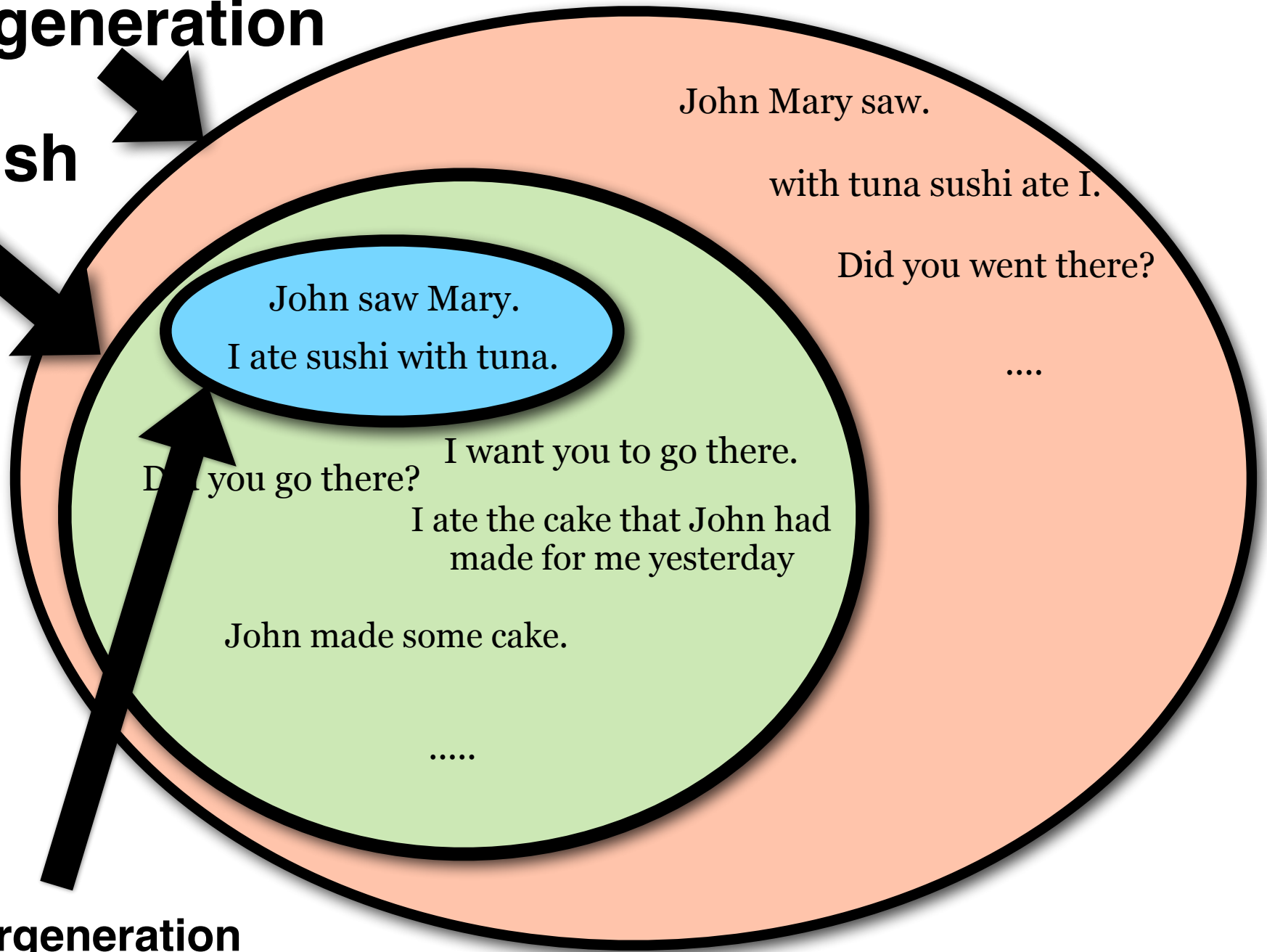
Implementations (in a particular formalism) for a particular language (English, Chinese,....)

Can we define a program that
generates all English sentences?

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

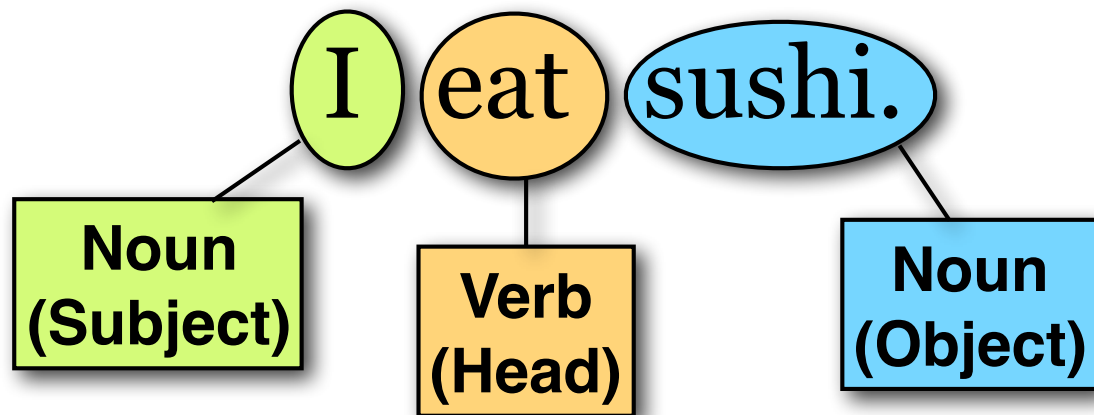
Overgeneration

English

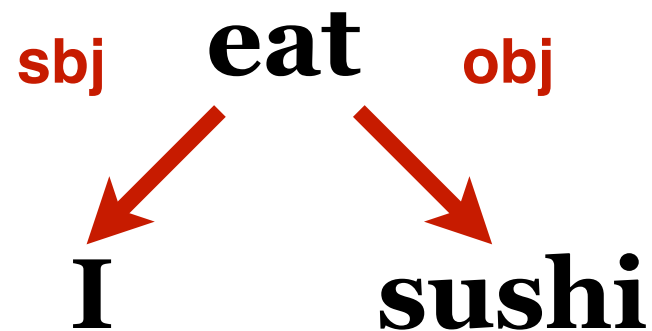
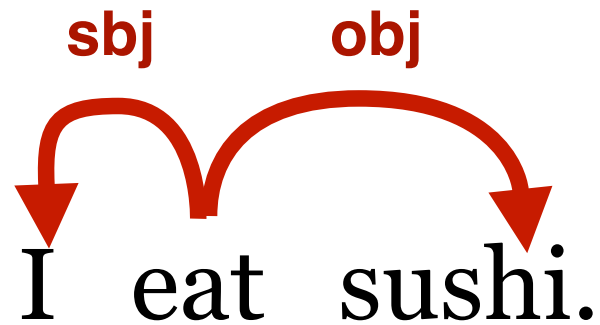


Undergeneration

Basic sentence structure



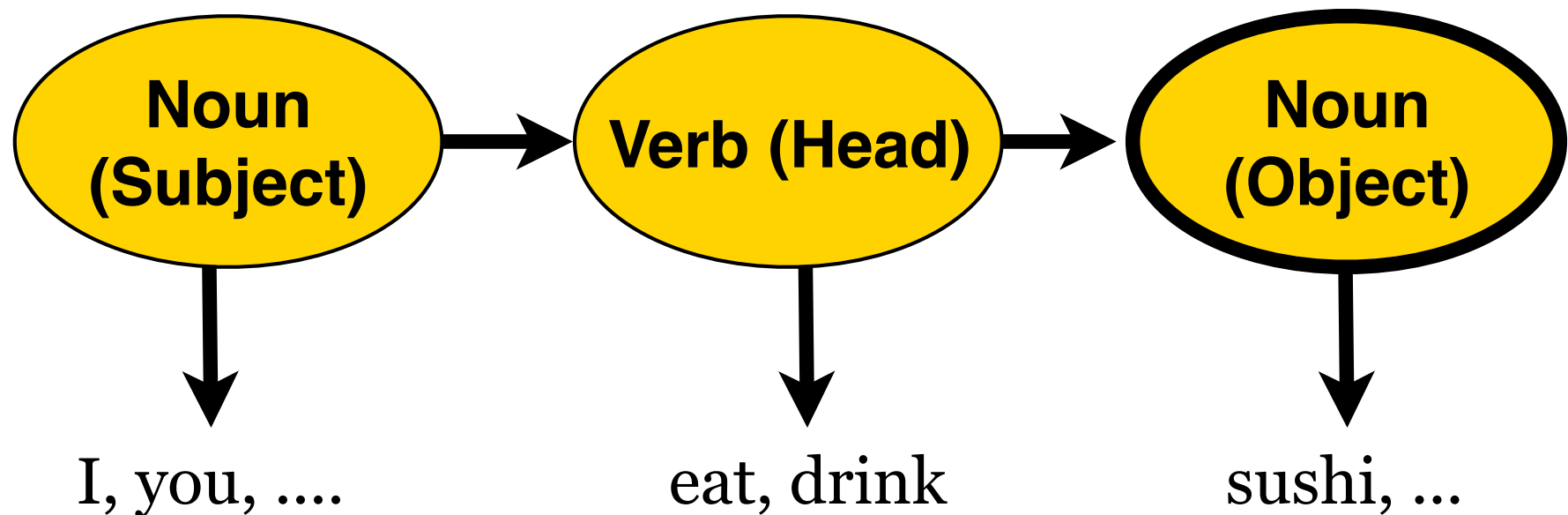
This is a dependency graph:



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

(purely syntactic: what set of arguments do words take?)

Intransitive verbs (sleep) take only a subject.

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

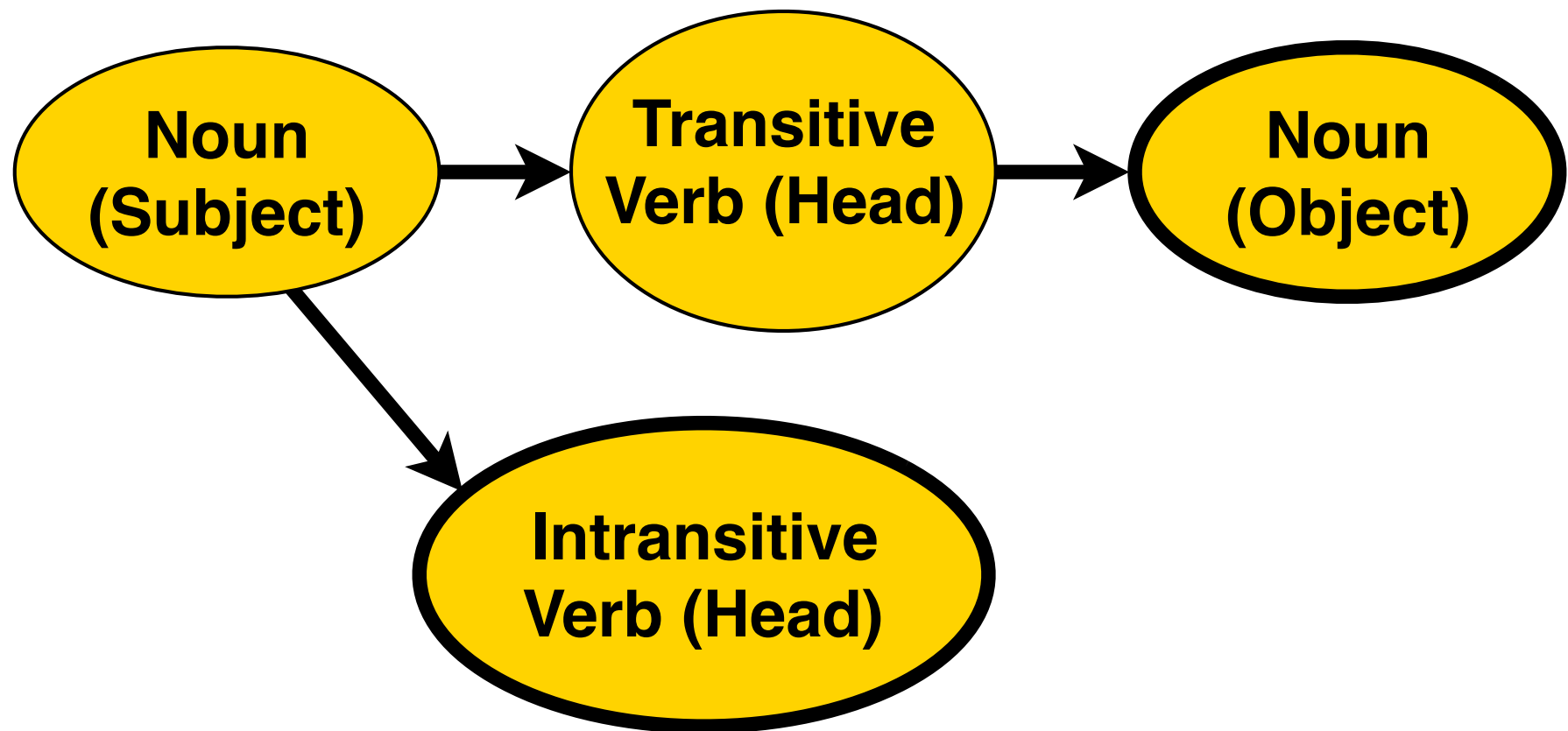
Ditransitive verbs (give) take also one (indirect) object.

Selectional preferences

(semantic: what types of arguments do words tend to take)

The object of eat should be edible.

A better FSA



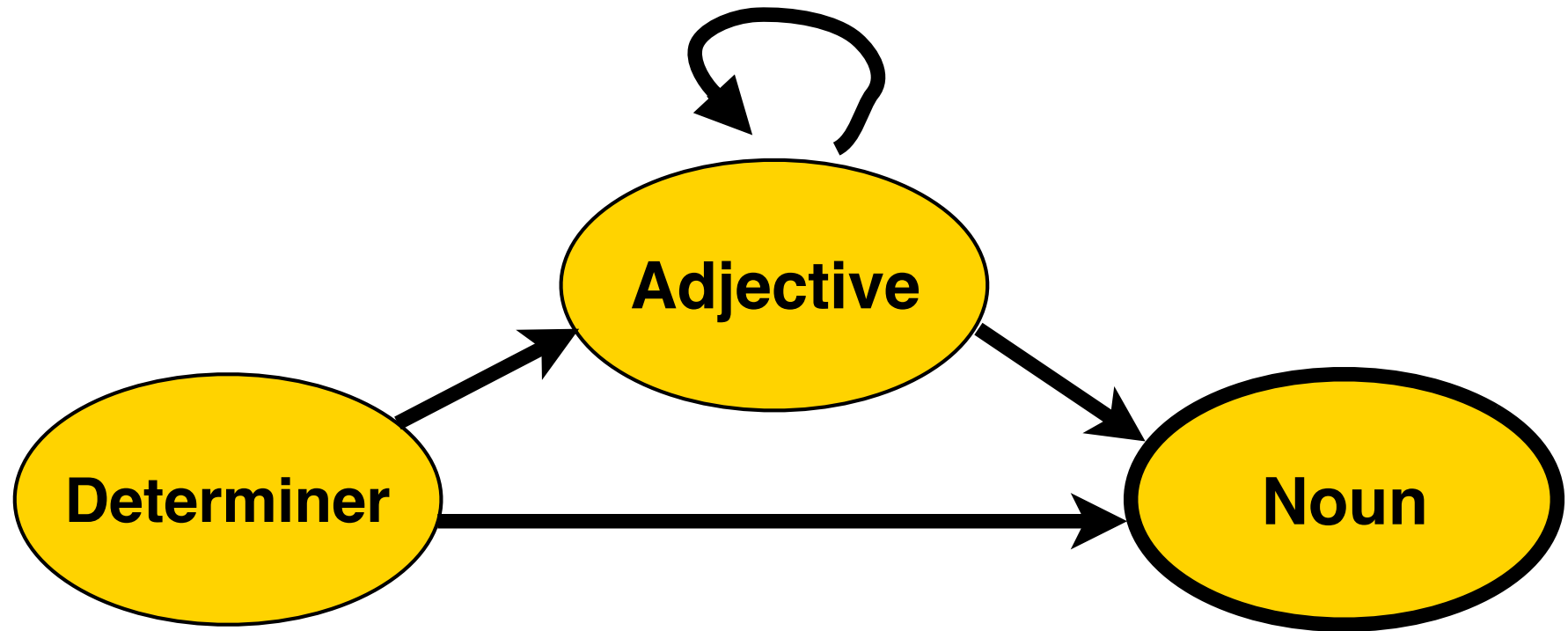
Language is recursive

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

Adjectives can **modify** nouns.

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

Another FSA



Recursion can be more complex

the ball

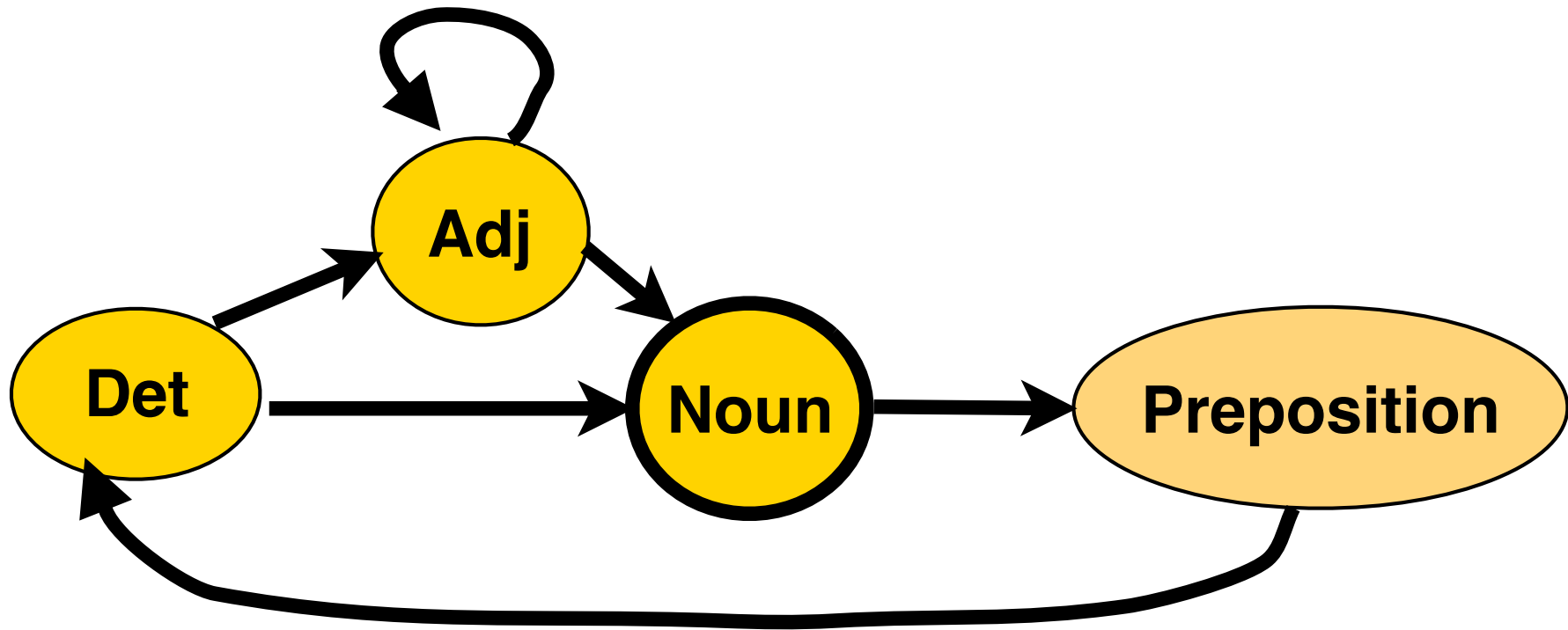
the ball in the garden

the ball in the garden behind the house

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

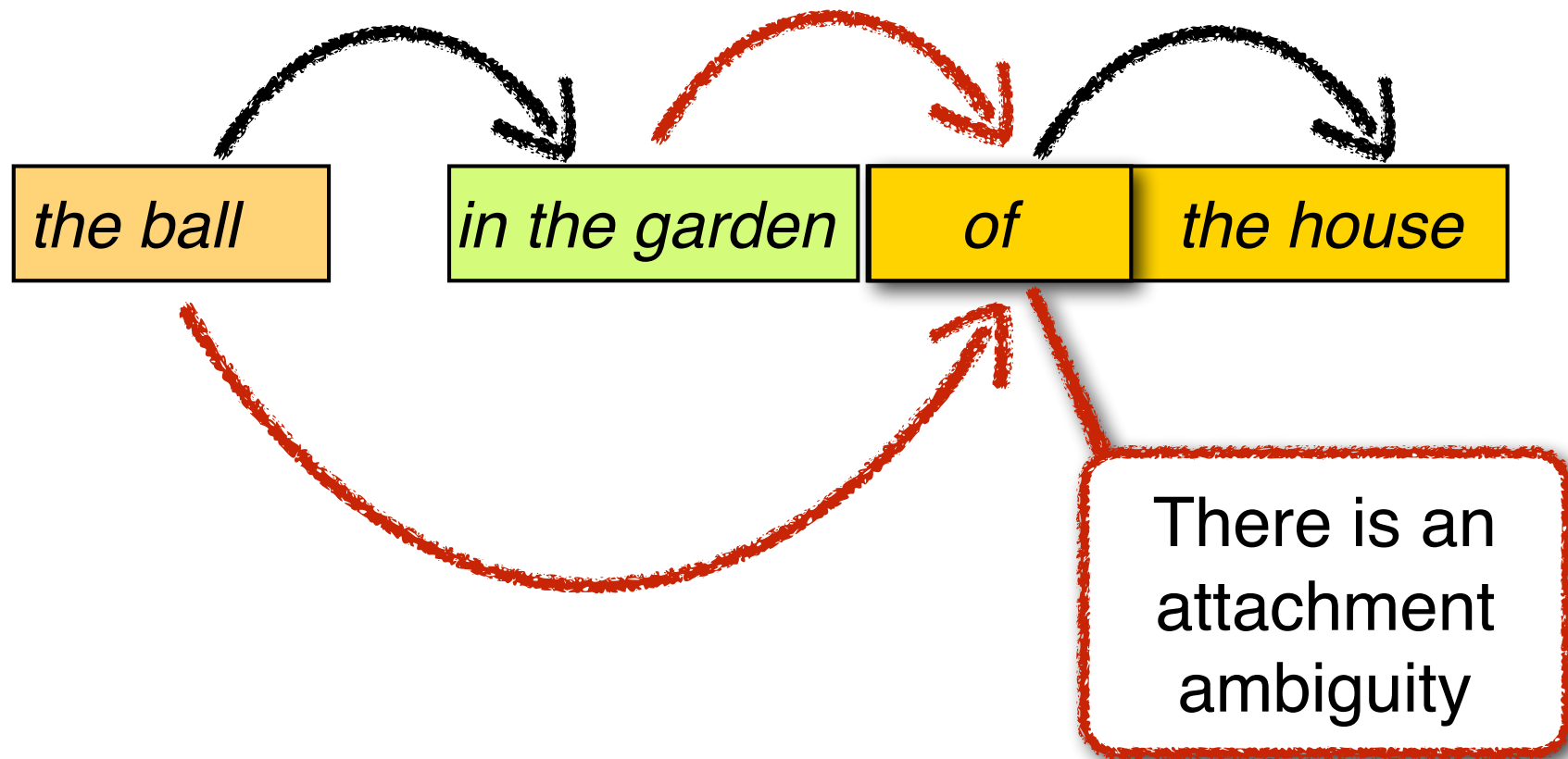
....

Yet another FSA

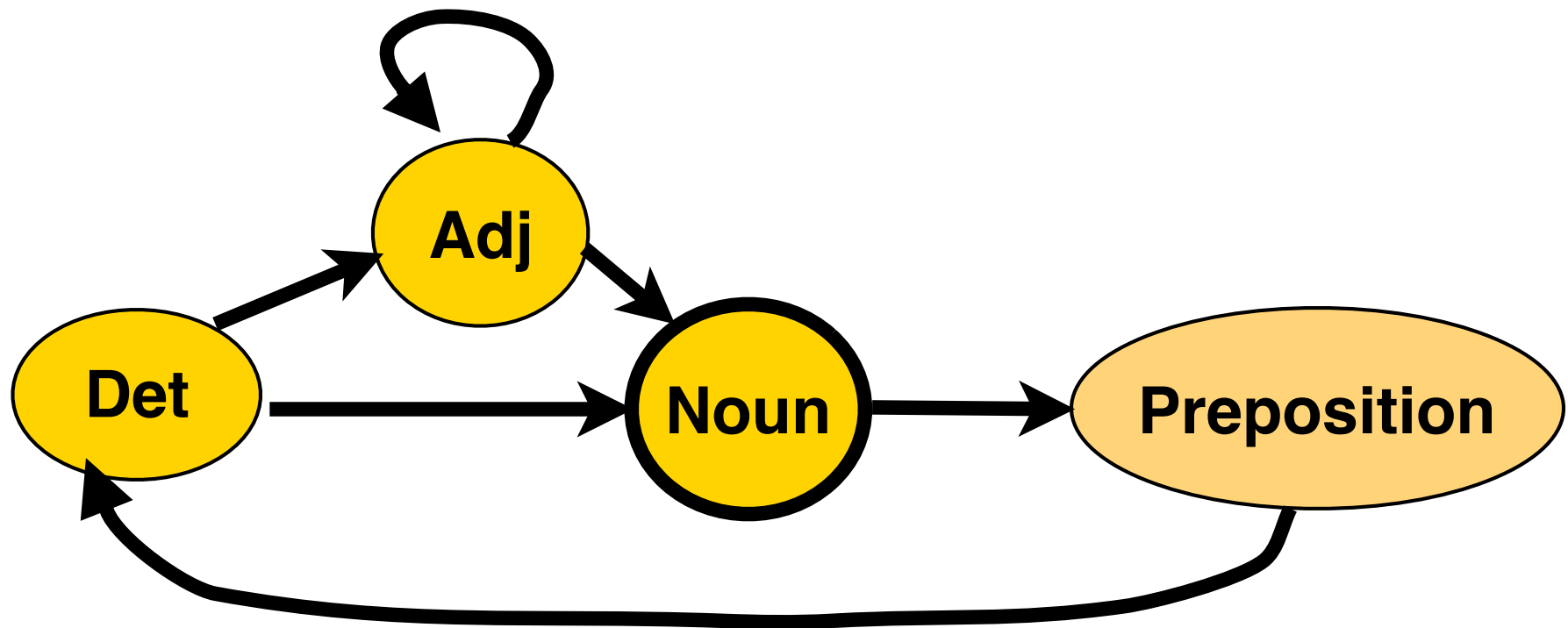


So, why do we need anything
beyond regular (finite-state) grammars?

What does this mean?



FSAs do not generate hierarchical structure

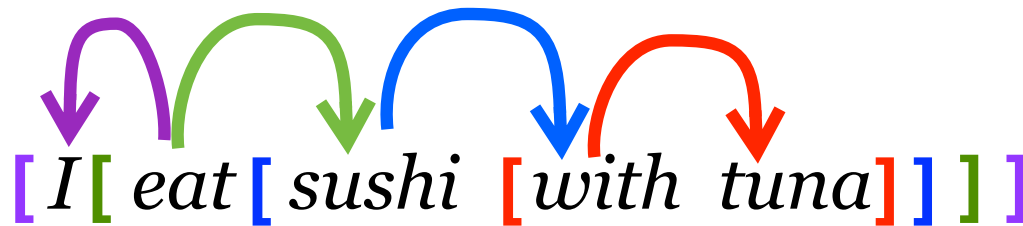


What is the structure of a sentence?

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 $\langle \mathbf{N}, \mathbf{\Sigma}, \mathbf{R}, S \rangle$ consisting of:

A set of nonterminals \mathbf{N}

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

A set of terminals $\mathbf{\Sigma}$

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

A set of rules \mathbf{R}

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

A start symbol $S \in \mathbf{N}$

An example

DT \rightarrow {the, a}

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

P \rightarrow {in, behind, with}

NP \rightarrow DT N

NP \rightarrow NP PP

PP \rightarrow P NP

N: noun

P: preposition

NP: “noun phrase”

PP: “prepositional phrase”

CFGs define parse trees

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

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

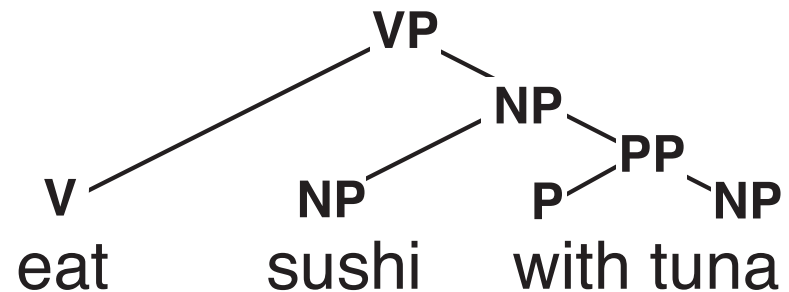
$V \rightarrow \{\text{eat}\}$

$NP \rightarrow N$

$NP \rightarrow NP \ PP$

$PP \rightarrow P \ NP$

$VP \rightarrow V \ NP$



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** ate the corn.

....

CFGs and center embedding

Formally, these sentences are all grammatical, because they can be generated by the CFG that is required for the first sentence:

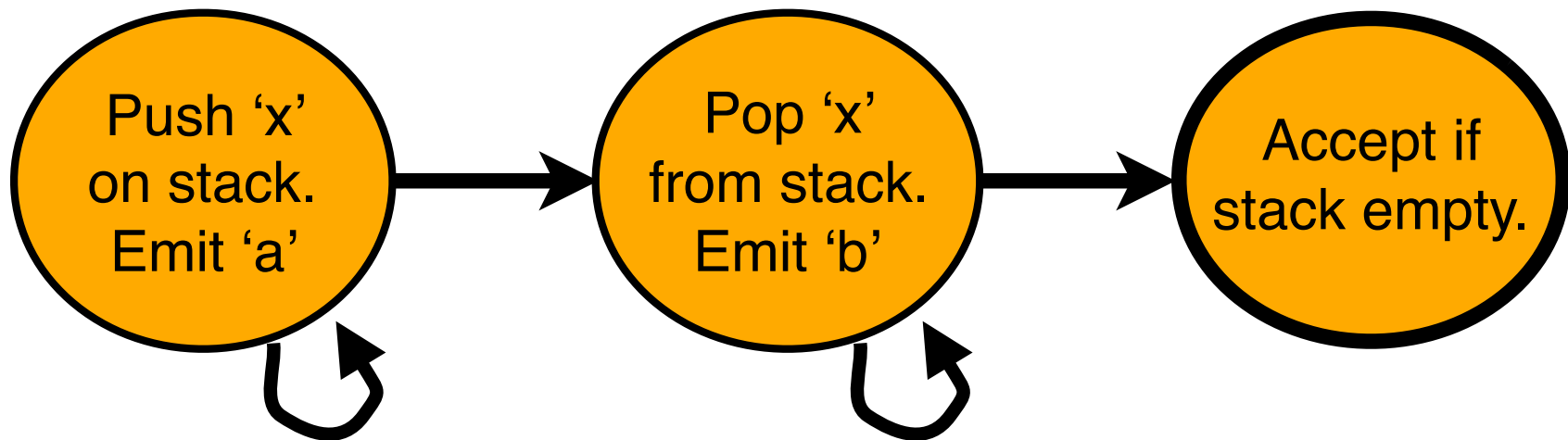
$$\begin{array}{ll} S & \rightarrow NP \quad VP \\ NP & \rightarrow NP \quad RelClause \\ RelClause & \rightarrow that \quad NP \quad ate \end{array}$$

Problem: CFGs are not able to capture **bounded recursion**. ('only embed one or two relative clauses').

To deal with this discrepancy between what the model predicts to be grammatical, and what humans consider grammatical, 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 $a^n b^n$

Action

1. Push x on stack. Emit a.
2. Push x on stack. Emit a.
3. Push x on stack. Emit a.
4. Push x on stack. Emit a.
5. Pop x off stack. Emit b.
6. Pop x off stack. Emit b.
7. Pop x off stack. Emit b.
8. Pop x off stack. Emit b

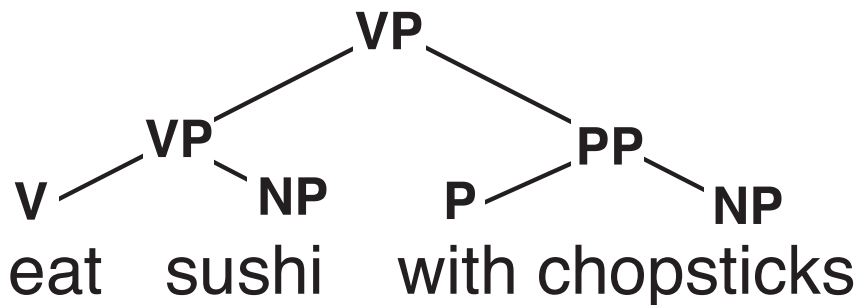
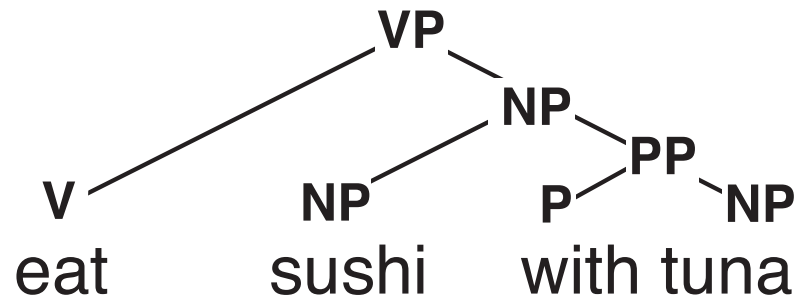
Stack String

x	a
xx	aa
xxx	aaa
xxxx	aaaa
xxx	aaab
xx	aaabb
x	aaabbb
	aaabbbb

Defining grammars for natural language

Two ways to represent structure

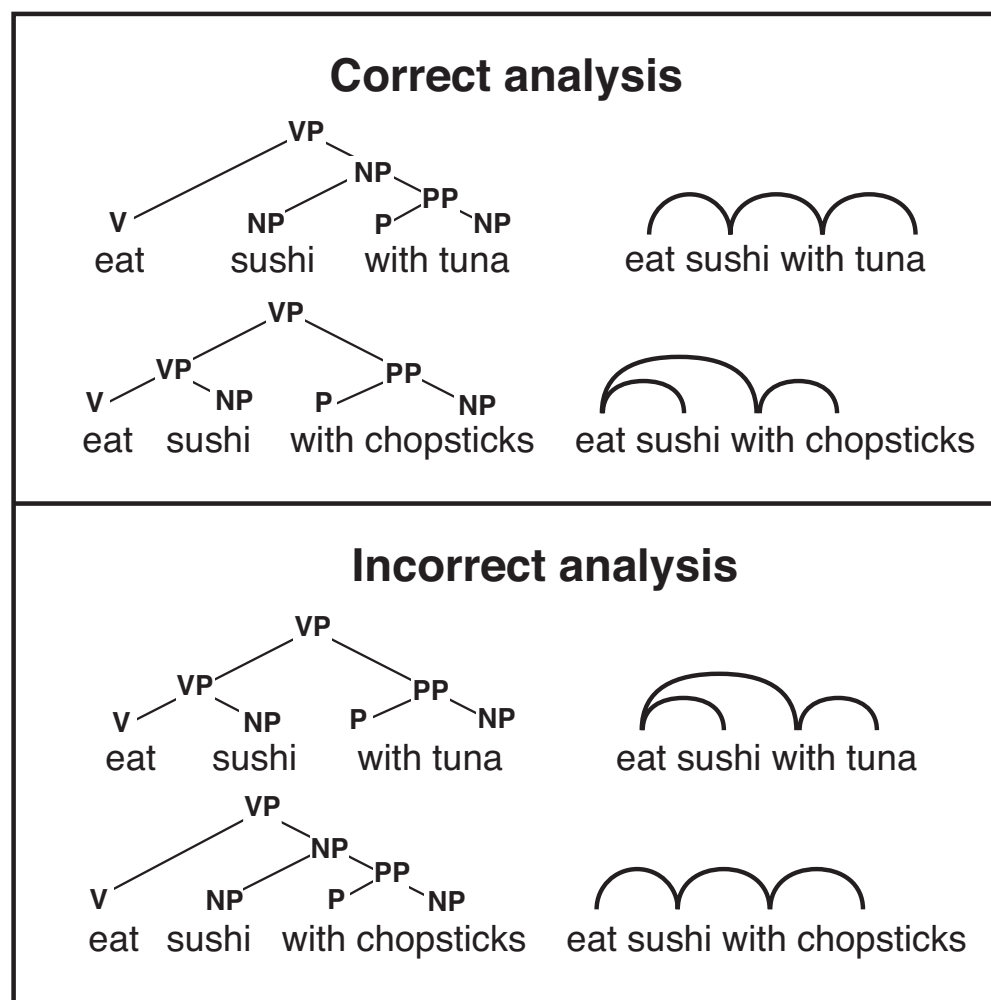
Phrase structure trees



Dependency trees



Structure (syntax) corresponds to meaning (semantics)



Dependency grammar

DGs describe the structure of sentences as a directed acyclic graph.

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

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

Typically, the graph is assumed to be a **tree**.

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:

No object: *[John] likes. No subject: *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 \rightarrow Pronoun

NP \rightarrow ProperName

NP \rightarrow Det Noun

Det \rightarrow {a, the, every}

Pronoun \rightarrow {he, she,...}

ProperName \rightarrow {John, Mary,...}

Noun \rightarrow AdjP Noun

Noun \rightarrow N

NP \rightarrow NP PP

NP \rightarrow NP RelClause

Adjective phrases (AdjP) and prepositional phrases (PP)

AdjP \rightarrow Adj

AdjP \rightarrow Adv AdjP

Adj \rightarrow {big, small, red,...}

Adv \rightarrow {very, really,...}

PP \rightarrow P NP

P \rightarrow {with, in, above,...}

The verb phrase (VP)

He [eats].

He [eats sushi].

He [gives John sushi].

He [eats sushi with chopsticks].

$VP \rightarrow V$

$VP \rightarrow V\ NP$

$VP \rightarrow V\ NP\ PP$

$VP \rightarrow VP\ PP$

$V \rightarrow \{\text{eats, sleeps gives,...}\}$

Capturing subcategorization

He [eats]. ✓

He [eats sushi]. ✓

He [gives John sushi]. ✓

He [eats sushi with chopsticks]. ✓

*He [eats John sushi]. ???

$VP \rightarrow V_{\text{intrans}}$

$VP \rightarrow V_{\text{trans}} NP$

$VP \rightarrow V_{\text{ditrans}} NP NP$

$VP \rightarrow VP PP$

$V_{\text{intrans}} \rightarrow \{\text{eats, sleeps}\}$

$V_{\text{trans}} \rightarrow \{\text{eats}\}$

$V_{\text{trans}} \rightarrow \{\text{gives}\}$

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].

$VP \rightarrow Vcomp S$

$Vcomp \rightarrow \{\text{says, think, believes}\}$

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]].

$VP \rightarrow V_{\text{have}} VP_{\text{pastPart}}$

$VP \rightarrow V_{\text{be}} VP_{\text{pass}}$

$VP_{\text{pastPart}} \rightarrow V_{\text{pastPart}} NP$

$VP_{\text{pass}} \rightarrow V_{\text{pastPart}} PP$

$V_{\text{have}} \rightarrow \{\text{has}\}$

$V_{\text{pastPart}} \rightarrow \{\text{eaten, seen}\}$

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]

$S \rightarrow S \text{ conj } S$

$NP \rightarrow NP \text{ conj } NP$

$VP \rightarrow VP \text{ conj } VP$

He says [he eats sushi].

$VP \rightarrow V_{\text{comp}} S$

$V_{\text{comp}} \rightarrow \{\text{says, think, believes}\}$

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 \rightarrow Aux NP VP_{inf}

YesNoQ \rightarrow 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