

CS447: Natural Language Processing

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

Lecture 20: Expressive Grammars

Julia Hockenmaier

juliahmr@illinois.edu

3324 Siebel Center

Grammars in NLP: what and why

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,.....)

(NB: any practical parser will need to also have a model/scoring function to identify which grammatical analysis should be assigned to a given sentence)

Why study grammar?

Linguistic questions:

What kind of constructions occur in natural language(s)?

Formal questions:

Can we define formalisms that allow us to characterize which strings belong to a language?

Those formalisms have appropriate weak generative capacity

Can we define formalisms that allow us to map sentences to their appropriate structures?

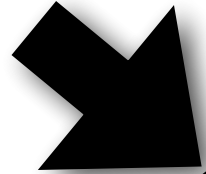
Those formalisms have appropriate strong generative capacity

Practical applications (Syntactic/Semantic Parsing):

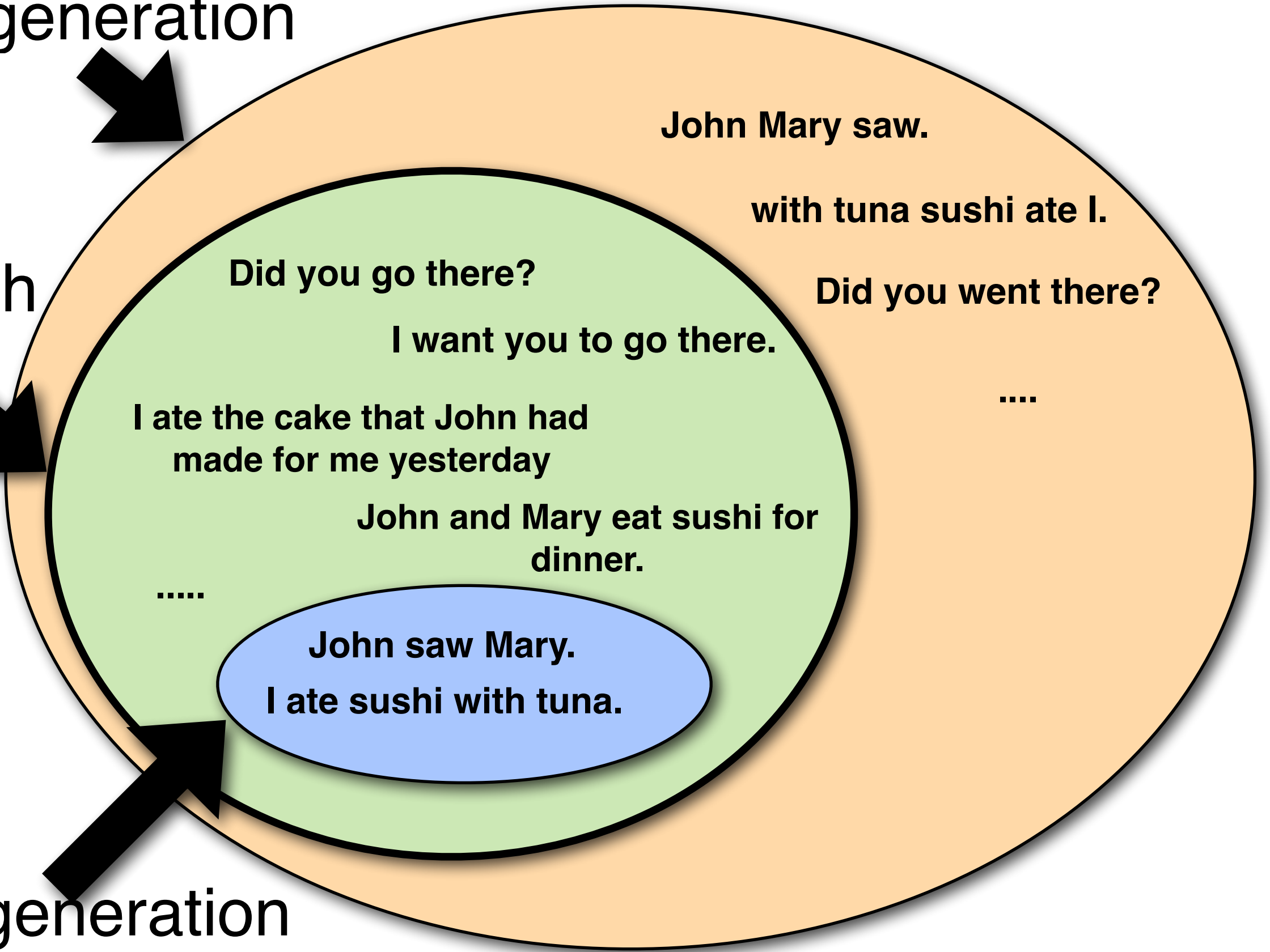
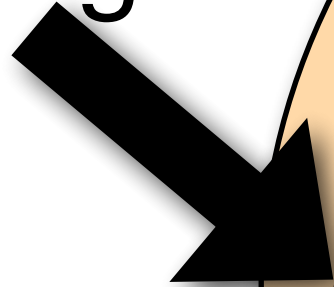
Can we identify the grammatical structure of sentences?

Can we translate sentences to appropriate meaning representations?

Overgeneration



English



John Mary saw.

with tuna sushi ate I.

Did you went there?

Did you go there?

I want you to go there.

I ate the cake that John had made for me yesterday

John and Mary eat sushi for dinner.

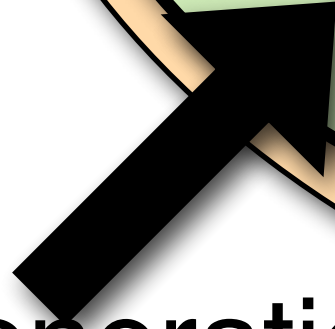
.....

John saw Mary.

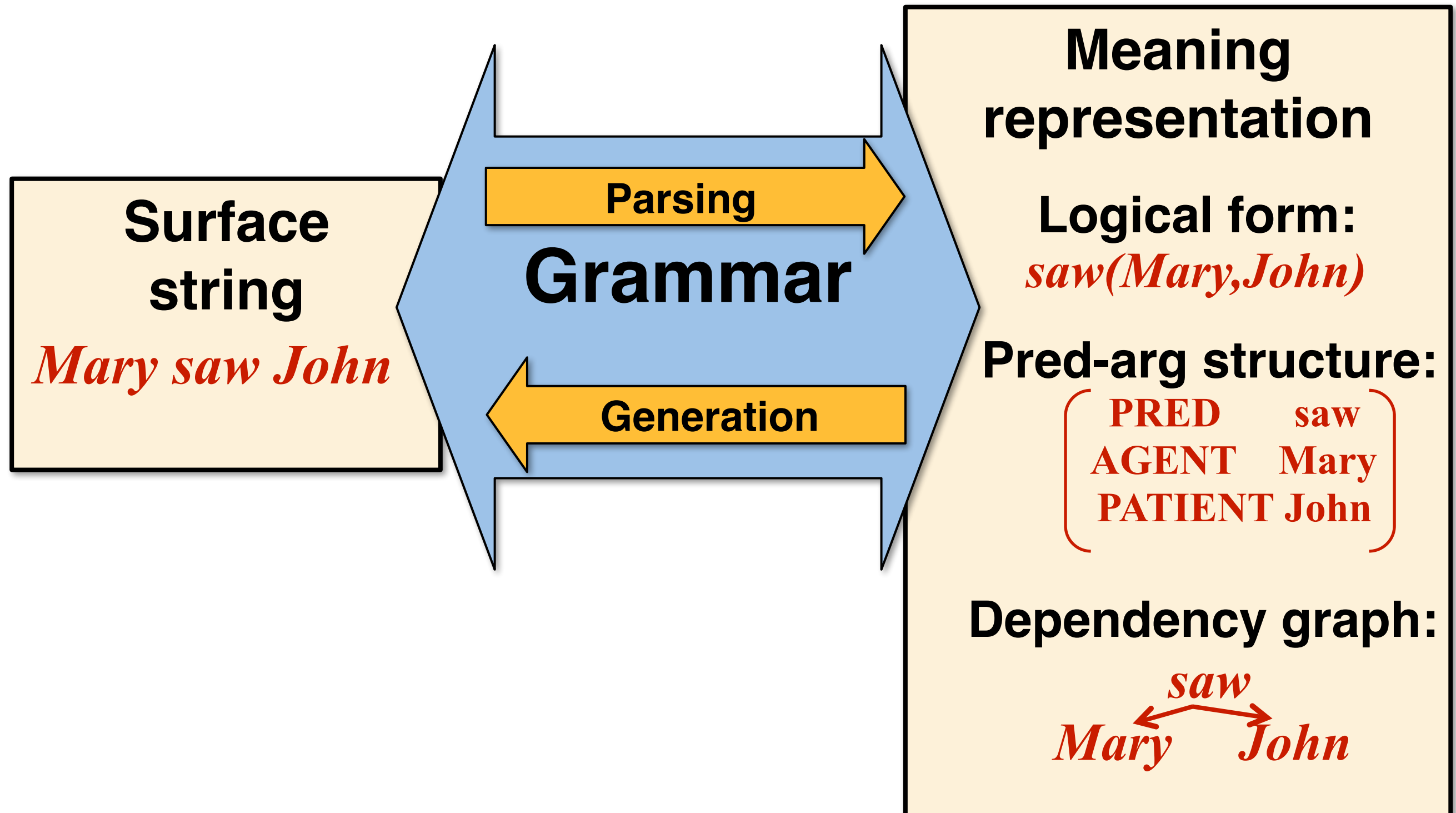
I ate sushi with tuna.

....

Undergeneration



Syntax as an interface to semantics



Grammar formalisms

Formalisms provide a **language** in which linguistic theories can be expressed and implemented

Formalisms define **elementary objects** (trees, strings, feature structures) and **recursive operations** which generate complex objects from simple objects.

Formalisms may impose **constraints** (e.g. on the kinds of dependencies they can capture)

What makes a formalism “expressive”?

“Expressive” formalisms are richer than context-free grammars.

Different formalisms use different mechanisms, data structures and operations to go beyond CFGs

Examples of expressive grammar formalisms

Tree-adjoining Grammar (TAG):

Fragments of phrase-structure trees

Lexical-functional Grammar (LFG):

Annotated phrase-structure trees (c-structure)
linked to feature structures (f-structure)

Combinatory Categorical Grammar (CCG):

Syntactic categories paired with meaning representations

Head-Driven Phrase Structure Grammar (HPSG):

Complex feature structures (Attribute-value matrices)

Why go beyond CFGs?

The dependencies so far:

Arguments:

Verbs take arguments: subject, object, complements, ...

Heads subcategorize for their arguments

Adjuncts/Modifiers:

Adjectives modify nouns, adverbs modify VPs or adjectives,

PPs modify NPs or VPs

Modifiers subcategorize for the head

Typically, these are *local* dependencies: they can be expressed *within individual CFG rules*

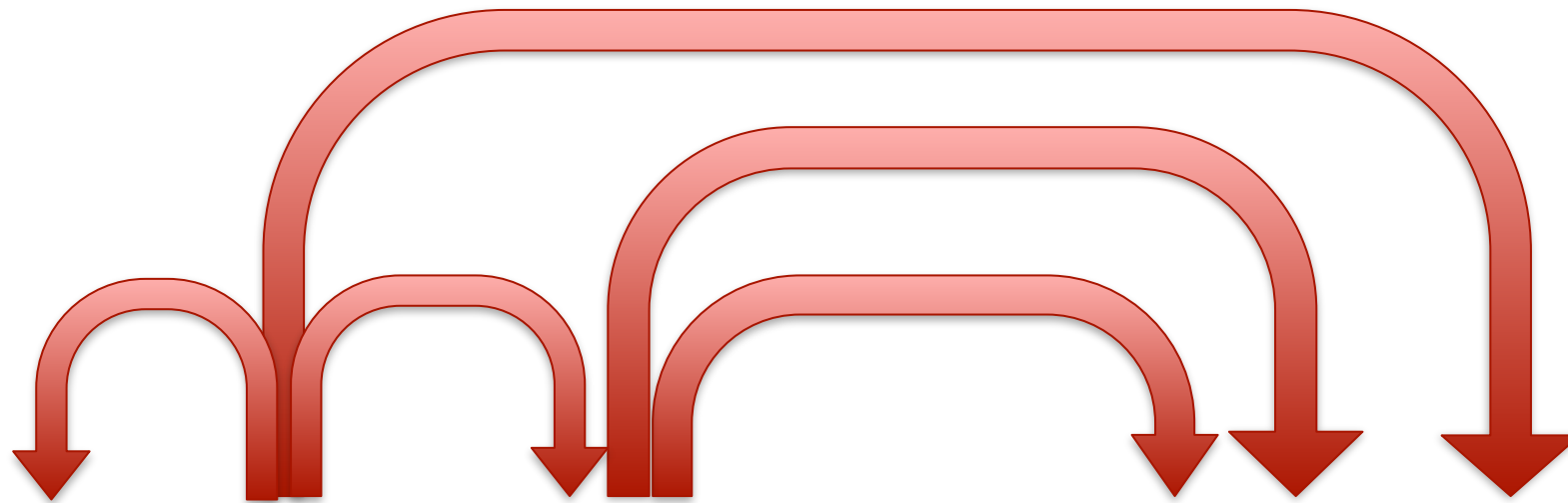


Context-free grammars

CFGs capture only **nested** dependencies

The dependency graph is a **tree**

The dependencies **do not cross**

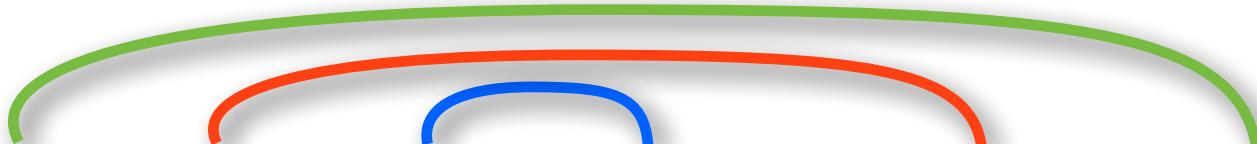


German: center embedding

...daß ich [Hans schwimmen] sah
...that I Hans swim saw
...that I saw [Hans swim]

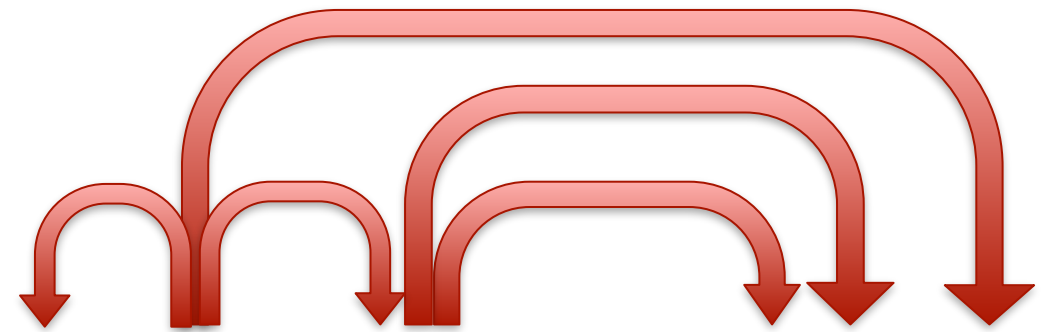
...daß ich [Maria [Hans schwimmen] helfen] sah
...that I Maria Hans swim help saw
...that I saw [Mary help [Hans swim]]

...daß ich [Anna [Maria [Hans schwimmen] helfen] lassen] sah
...that I Anna Maria Hans swim help let saw
...that I saw [Anna let [Mary help [Hans swim]]]

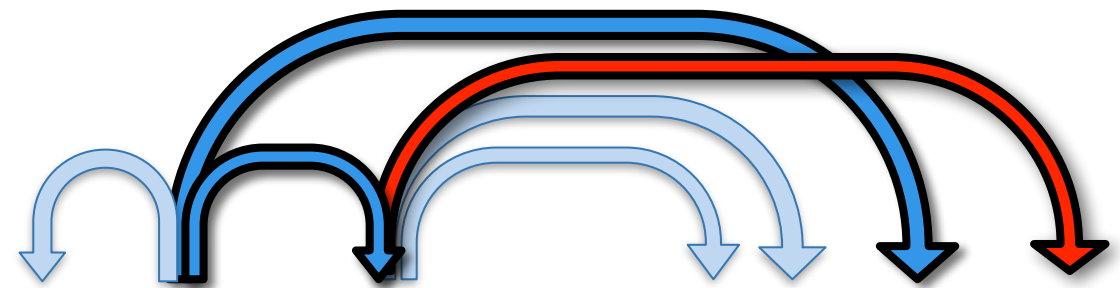
The diagram consists of three nested, upward-curving arcs above the German sentence. The innermost arc is blue and connects 'Hans' and 'schwimmen'. The middle arc is red and connects 'Maria' and 'helfen'. The outermost arc is green and connects 'Anna' and 'lassen'. These arcs illustrate the center-embedding structure of the sentence.

Dependency structures in general

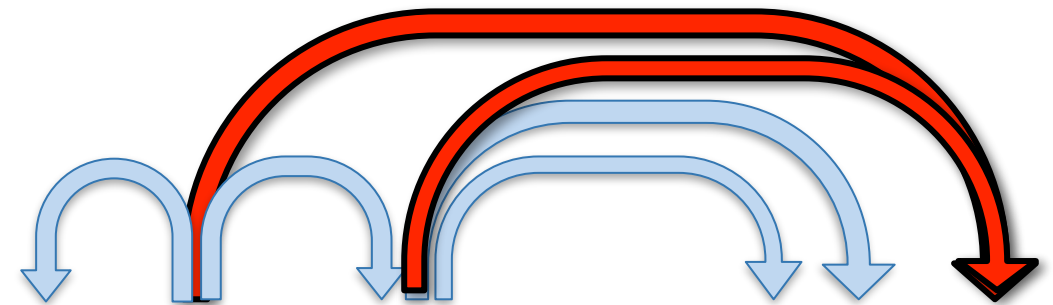
Nested (projective)
dependency trees
(CFGs)



Non-projective
dependency trees

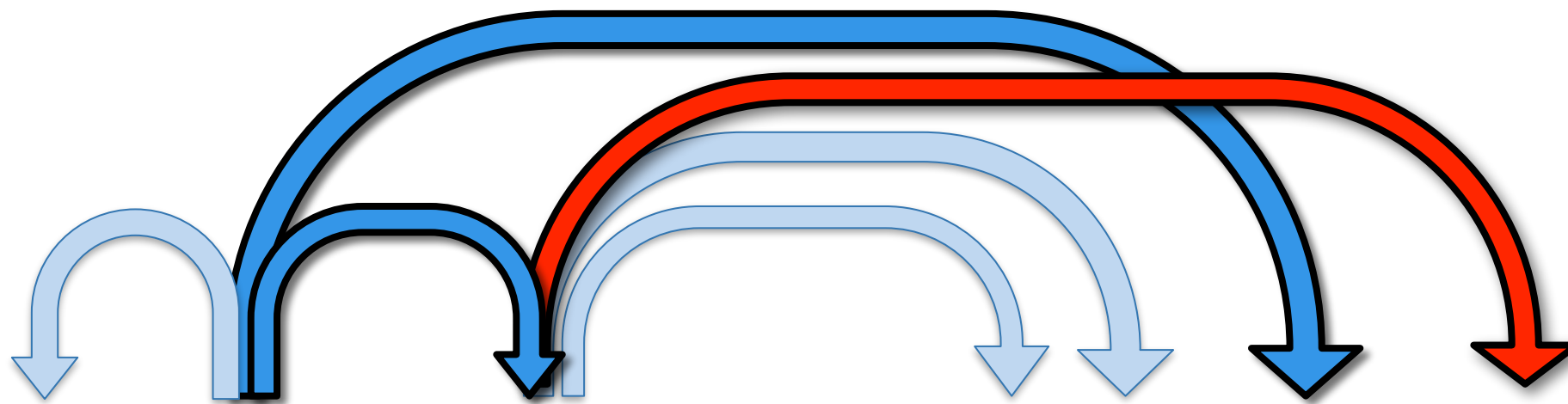


Non-local dependency
graphs



Beyond CFGs: Nonprojective dependencies

Dependencies form a **tree with crossing branches**



Dutch: Cross-Serial Dependencies

...dat ik Hans zag zwemmen

...that I Hans saw swim

...that I saw [Hans swim]

...dat ik Maria Hans zag helpen zwemmen

...that I Maria Hans saw help swim

...that I saw [Mary help [Hans swim]]

...dat ik Anna Maria Hans zag laten helpen zwemmen

...that I Anna Maria Hans saw let help swim

...that I saw [Anna let [Mary help [Hans swim]]]

Such **cross-serial** dependencies require
mildly context-sensitive grammars

Other crossing (non-projective) dependencies

(Non-local) scrambling: In a sentence with multiple verbs, the argument of a verb appears in a different clause from that which contains the verb (arises in languages with freer word order than English)

Die Pizza hat Klaus versprochen zu bringen

The pizza has Klaus promised to bring

Klaus has promised to bring the pizza

Extraposition: Here, a modifier of the subject NP is moved to the end of the sentence

The guy is coming who is wearing a hat

Compare with the non-extraposed variant

The [guy [who is wearing a hat]] is coming

Topicalization: Here, the argument of the embedded verb is moved to the front of the sentence.

Cheeseburgers, I [thought [he likes]]

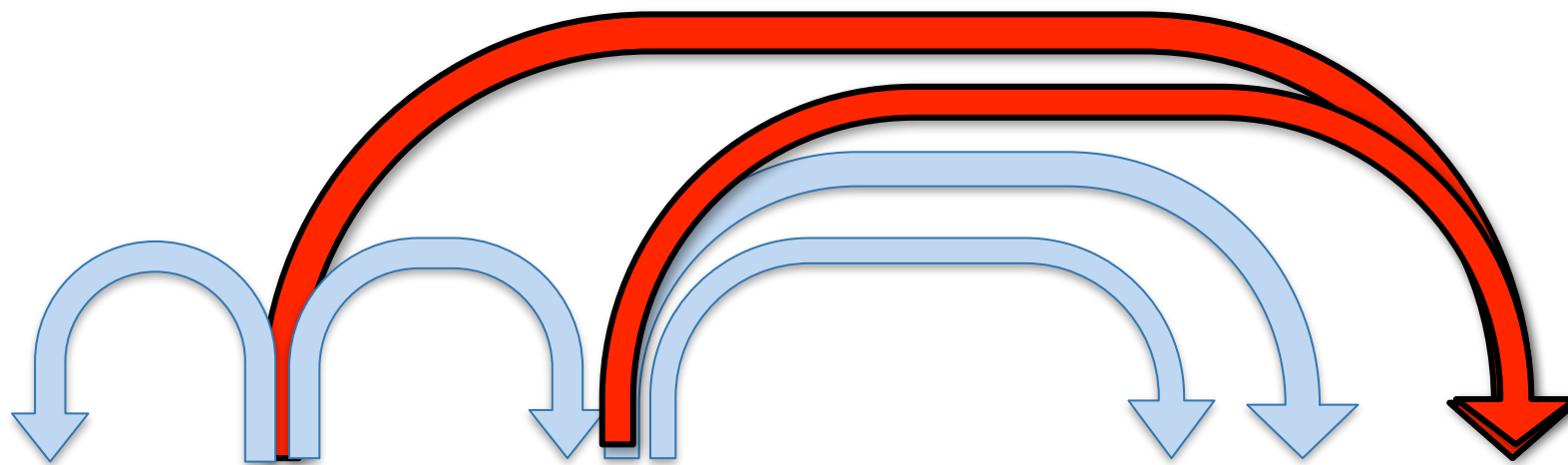
Beyond CFGs: Nonlocal dependencies

Dependencies form a **DAG**

(a node may have **multiple incoming edges**)

Arise in the following constructions:

- **Control** (*He has **promised** me to **go***), **raising** (*He **seems** to **go***)
- **Wh-movement** (*the **man** who you **saw** yesterday **is** here again*),
- **Non-constituent** coordination
(right-node raising, gapping, argument-cluster coordination)



Wh-Extraction (e.g. in English)

Relative clauses:

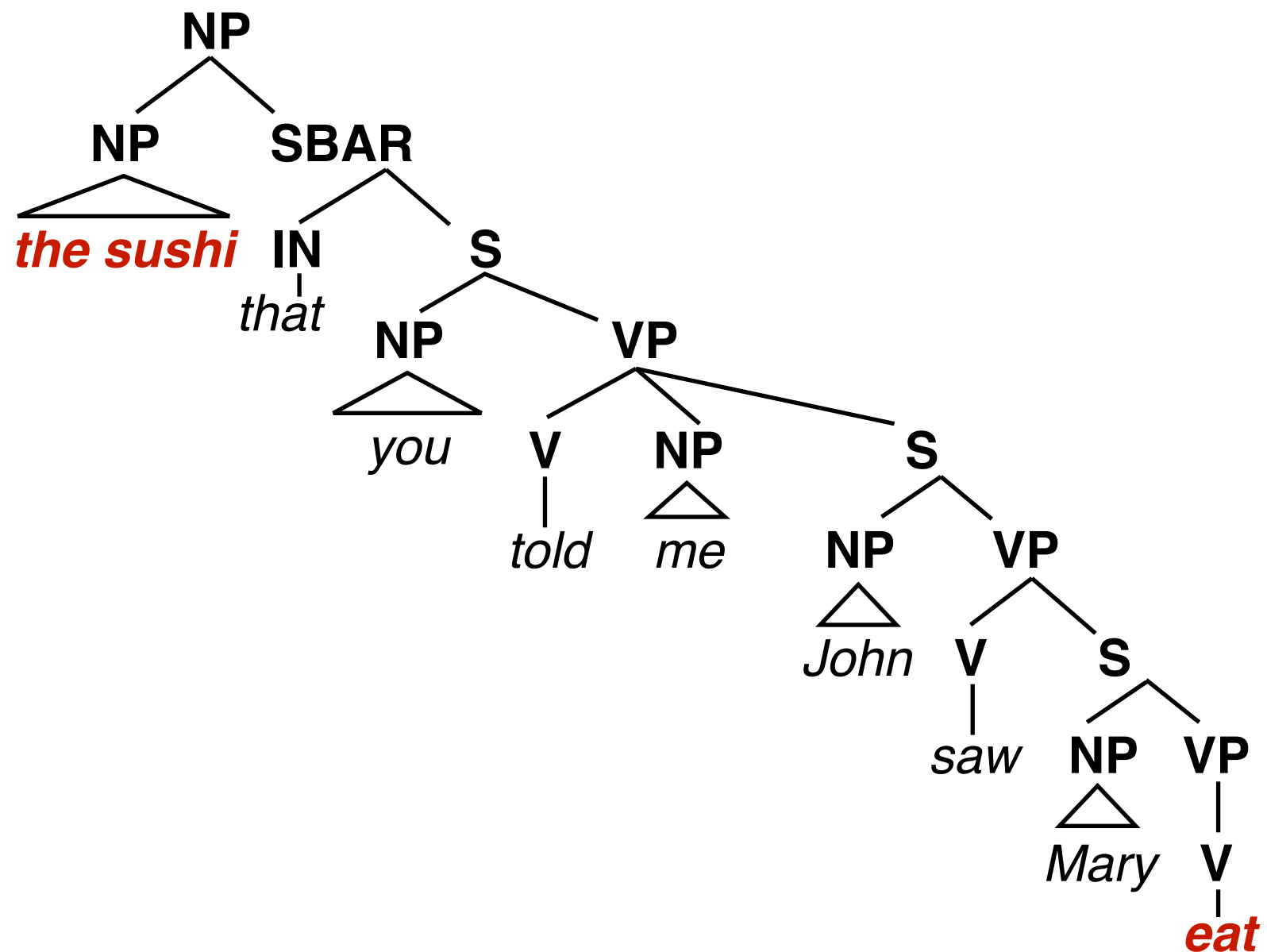
*the **sushi** that [you told me [John saw [**Mary eat**]]]*'

Wh-Questions:

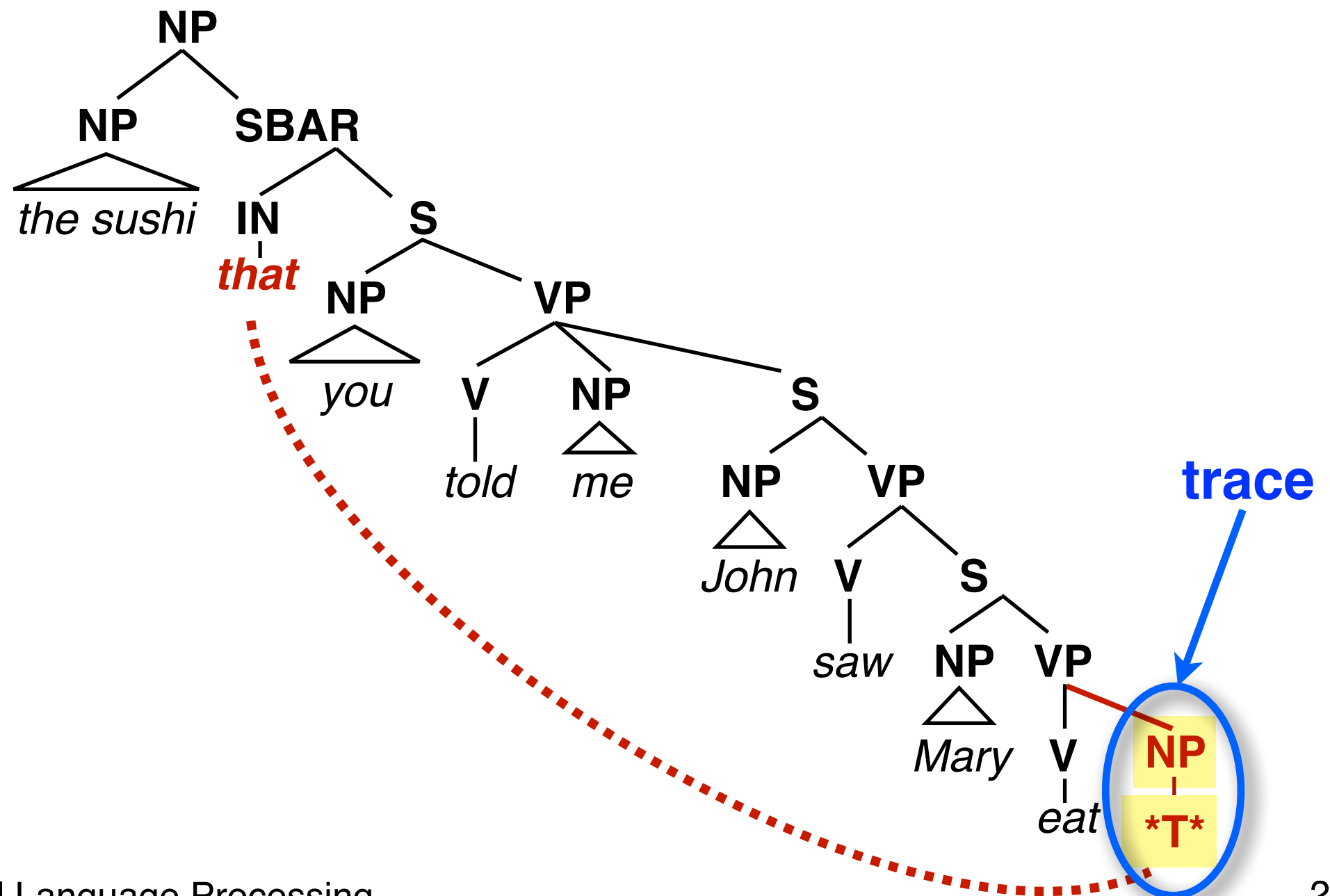
*'**what** [did you tell me [John saw [**Mary eat**]]]?'*

Wh-questions (what, who, ...) and relative clauses contain so-called *unbounded* nonlocal dependencies because the verb that subcategorizes for the moved NP may be arbitrarily deeply embedded in the tree. Linguists call this phenomenon **wh-extraction** (wh-movement).

As a phrase structure tree:



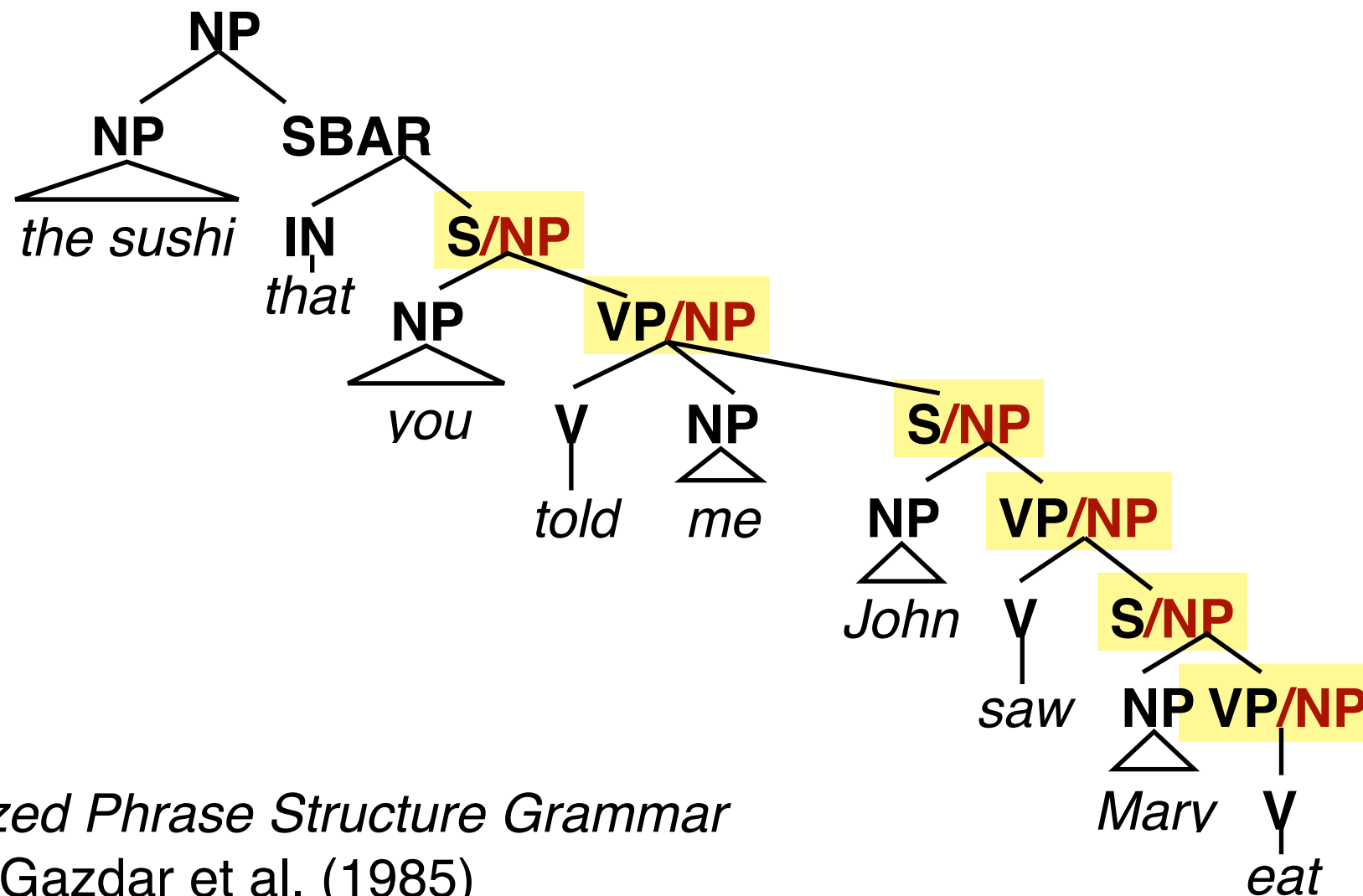
The trace analysis of wh-extraction



Slash categories for wh-extraction

Because only one element can be extracted, we can use **slash categories**.

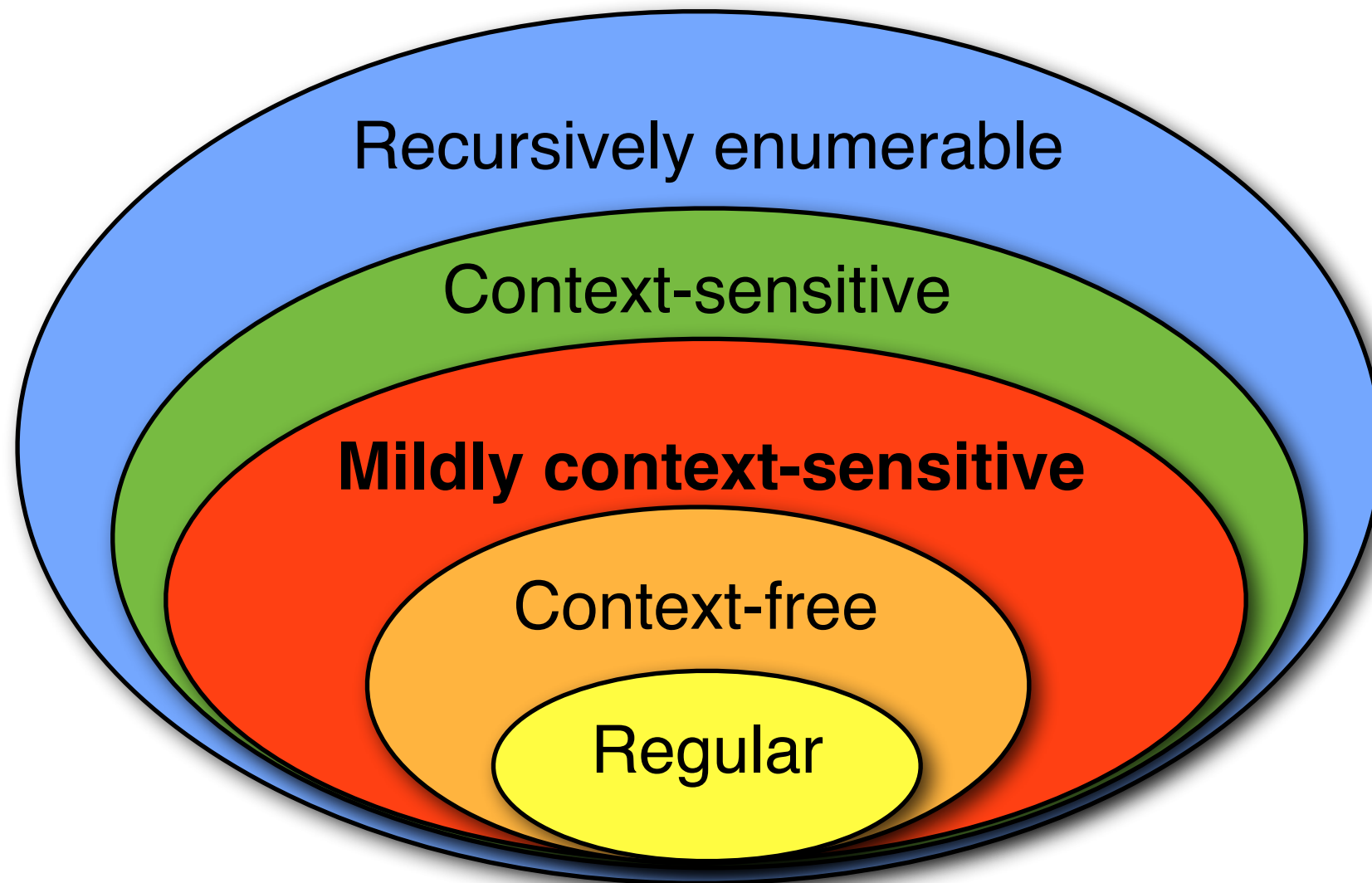
This is still a CFG: the set of nonterminals is finite.



Generalized Phrase Structure Grammar
(GPSG), Gazdar et al. (1985)

Two mildly context-sensitive formalisms: TAG and CCG

The Chomsky Hierarchy



Mildly context-sensitive grammars

Contain all context-free grammars/languages

Can be parsed in polynomial time (TAG/CCG: $O(n^6)$)

(*Strong* generative capacity) capture certain kinds of dependencies: **nested** (like CFGs) and **cross-serial** (like the Dutch example), but not the MIX language:

MIX: the set of strings $w \in \{a, b, c\}^*$ that contain equal numbers of *as*, *bs* and *cs*

Have the **constant growth** property:

the length of strings grows in a linear way

The power-of-2 language $\{a^{2^n}\}$ does not have the constant growth property.

TAG and CCG are lexicalized formalisms

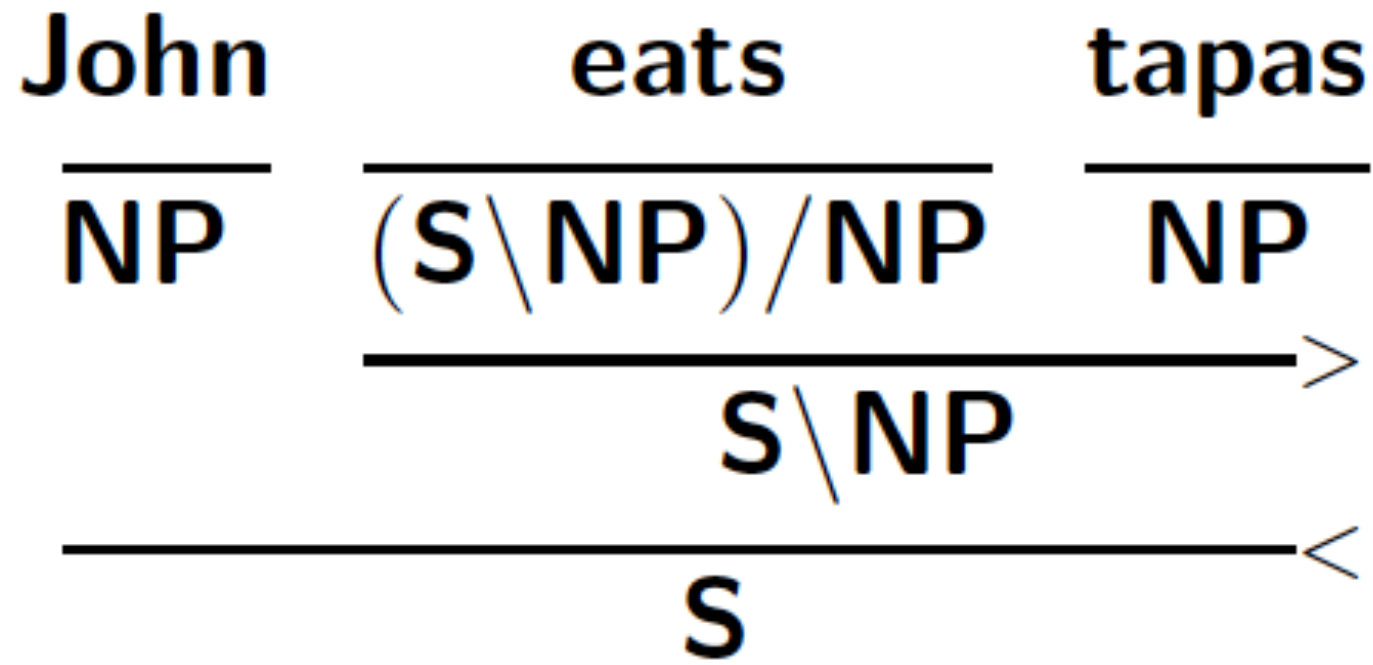
The lexicon:

- pairs words with elementary objects
- specifies all language-specific information (e.g. subcategorization information)

The grammatical operations:

- are universal
- define (and impose constraints on) recursion.

A (C)CG derivation



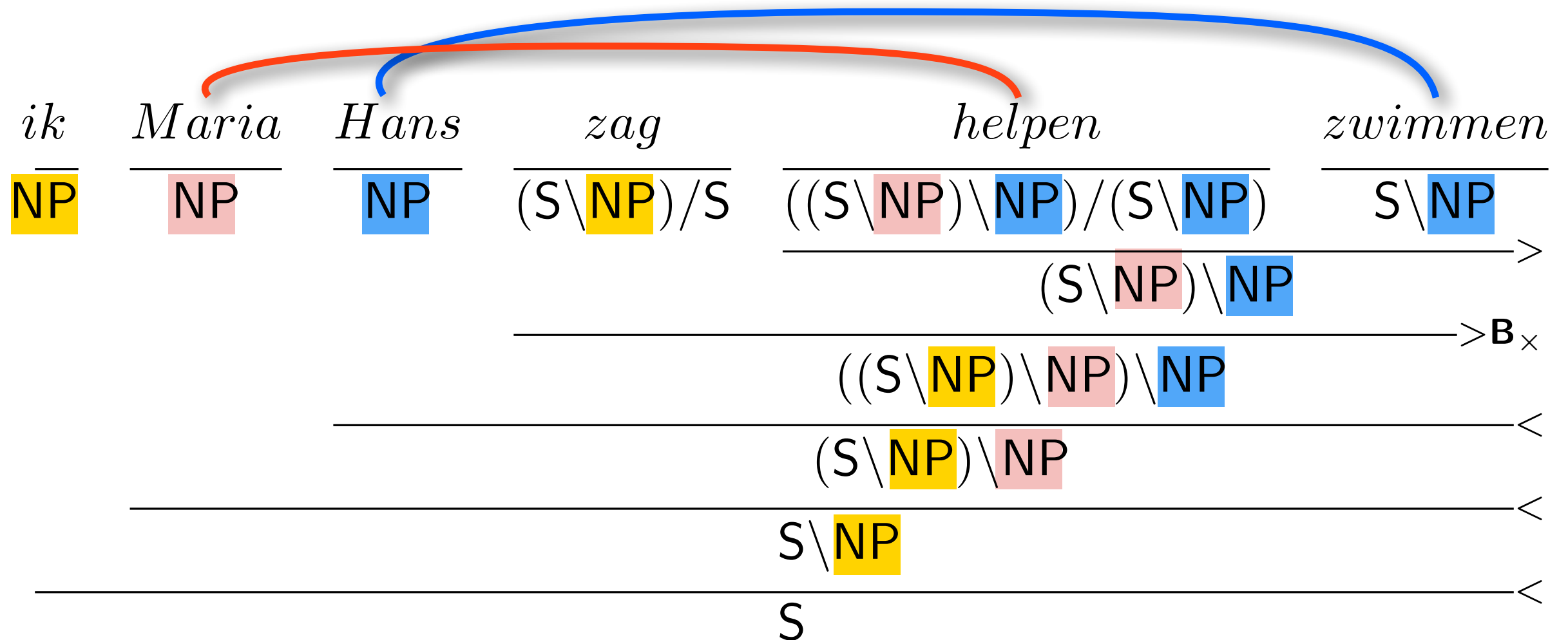
CCG categories are defined recursively:

- Categories are atomic (S, NP) or complex (S\NP, (S\NP)/NP)
- Complex categories (X/Y or X\Y) are functions:
X/Y combines with an adjacent argument to its right of category Y to return a result of category X.

Function categories can be composed, giving more expressive power than CFGs

More on CCG in one of our Semantics lectures!

Dutch cross-serial dependencies



Tree-Adjoining Grammar

(Lexicalized) Tree-Adjoining Grammar

TAG is a tree-rewriting formalism:

TAG defines operations (**substitution**, **adjunction**) on trees.

The **elementary objects** in TAG are trees (not strings)

TAG is lexicalized:

Each elementary tree is **anchored** to a lexical item (word)

“Extended domain of locality”:

The elementary tree contains all arguments of the anchor.

TAG requires a linguistic theory which specifies the shape of these elementary trees.

TAG is mildly context-sensitive:

can capture Dutch cross-serial dependencies

but is still efficiently parseable

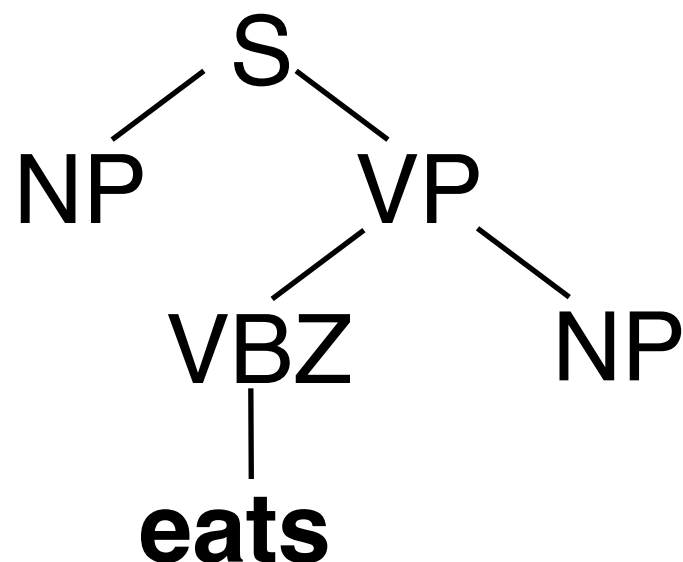
AK Joshi and Y Schabes (1996)
Tree Adjoining Grammars.
In G. Rosenberg and A. Salomaa,
Eds., Handbook of Formal
Languages

Extended domain of locality

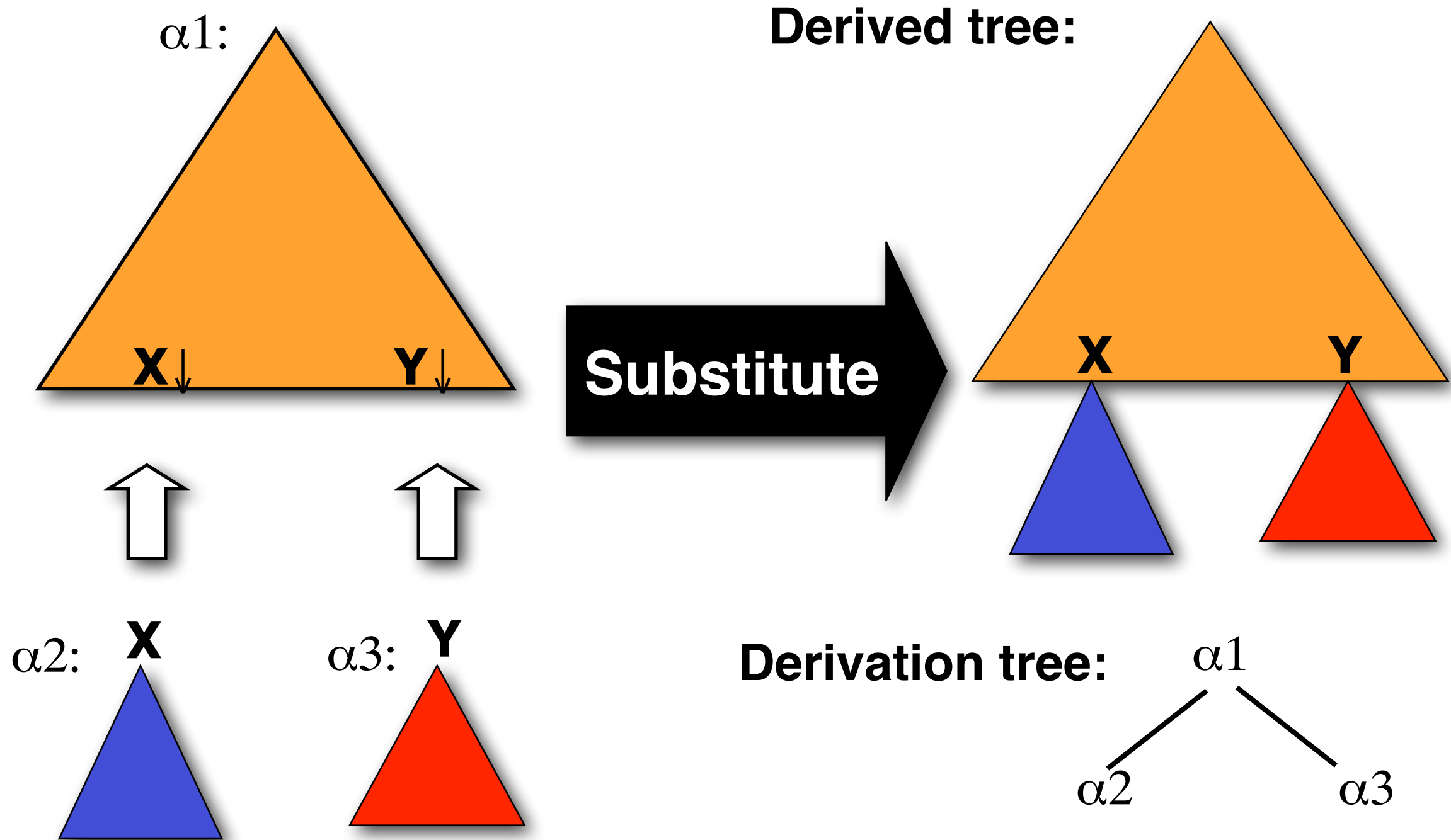
We want to capture **all arguments of a word** in a **single elementary object**.

We also want to retain certain syntactic structures (e.g. VPs).

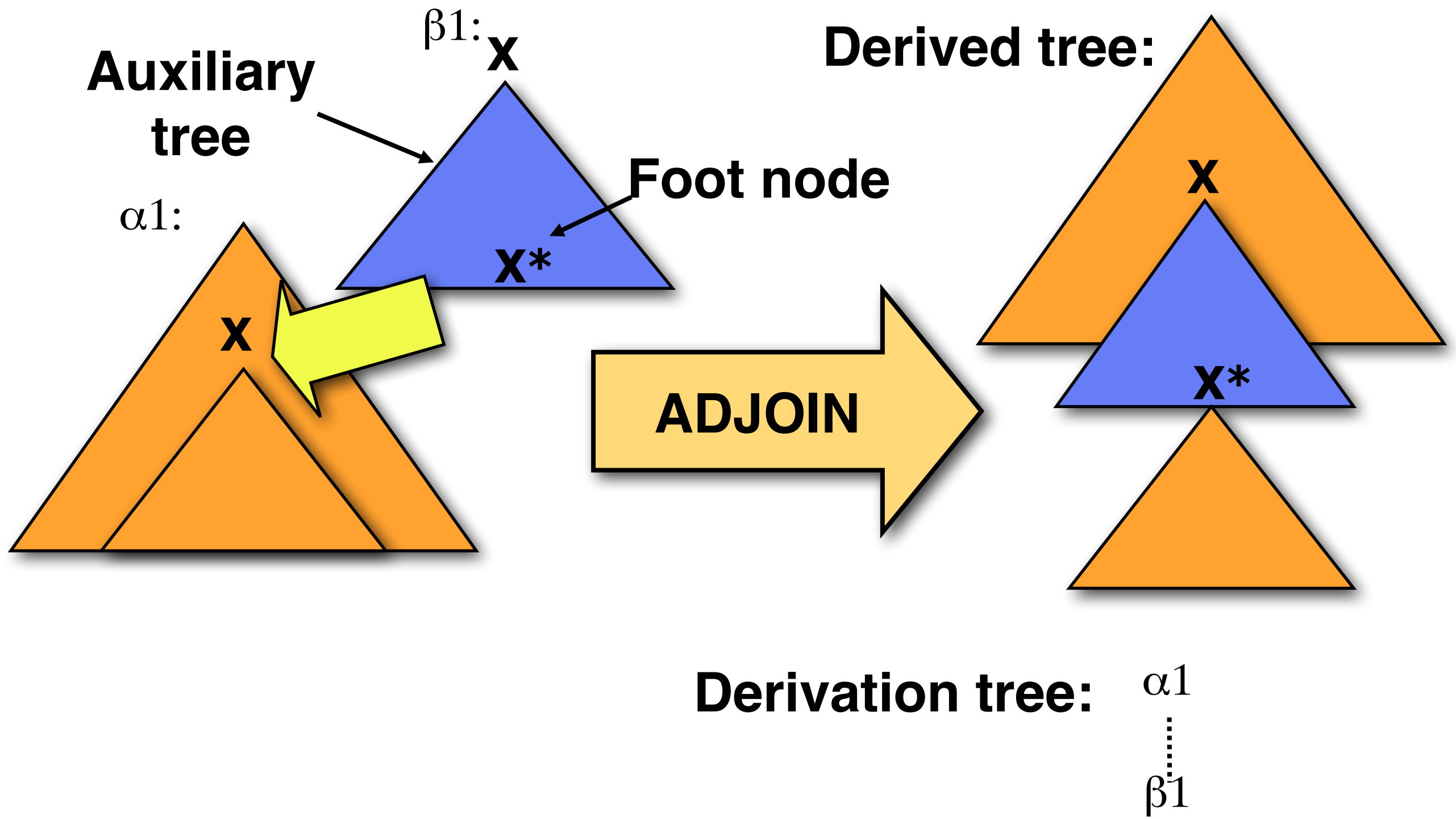
Our elementary objects are tree fragments:



TAG substitution (arguments)

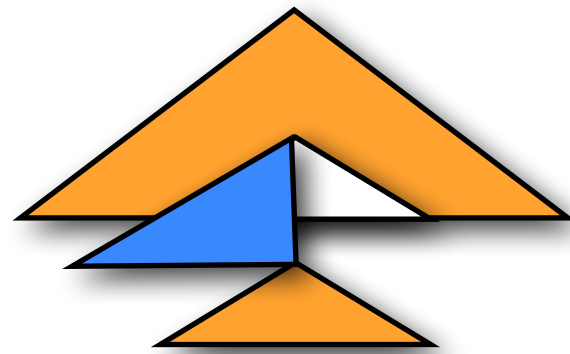


TAG adjunction

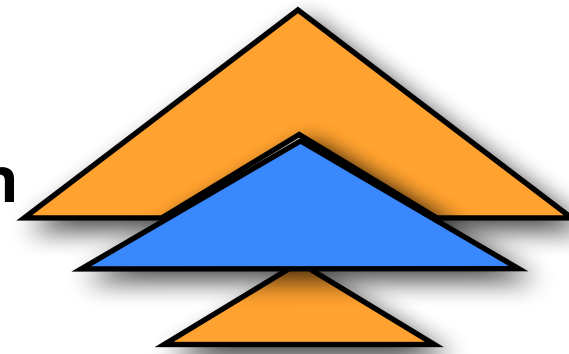


The effect of adjunction

**TIG:
sister
adjunction**



**TAG:
wrapping
adjunction**



No adjunction: TSG (Tree substitution grammar)

TSG is context-free

Sister adjunction: TIG (Tree insertion grammar)

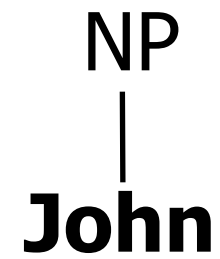
TIG is also context-free, but has a linguistically more adequate treatment of modifiers

Wrapping adjunction: TAG (Tree-adjoining grammar)

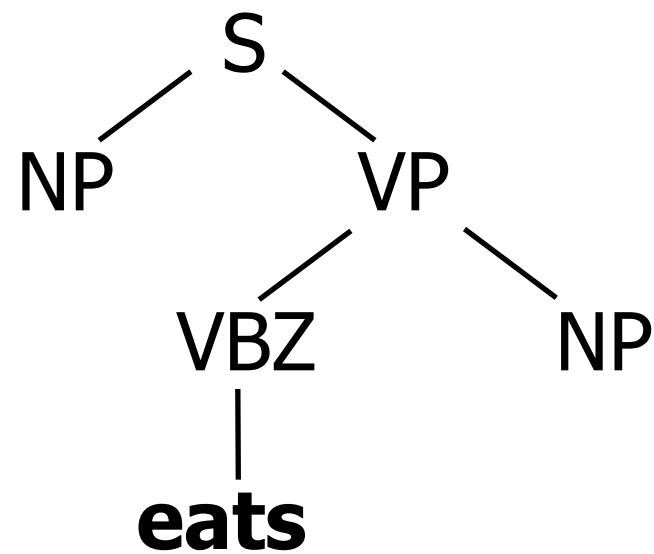
TAG is mildly context-sensitive

A small TAG lexicon

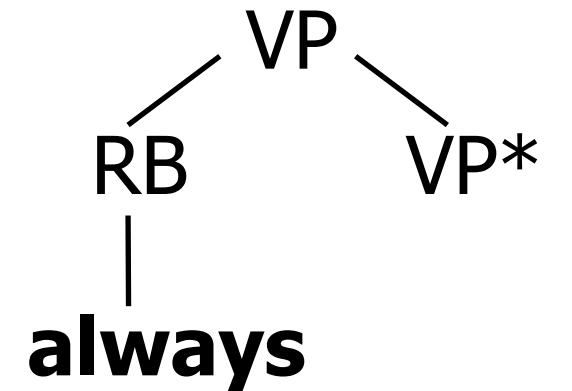
α_2 :



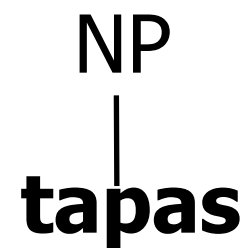
α_1 :



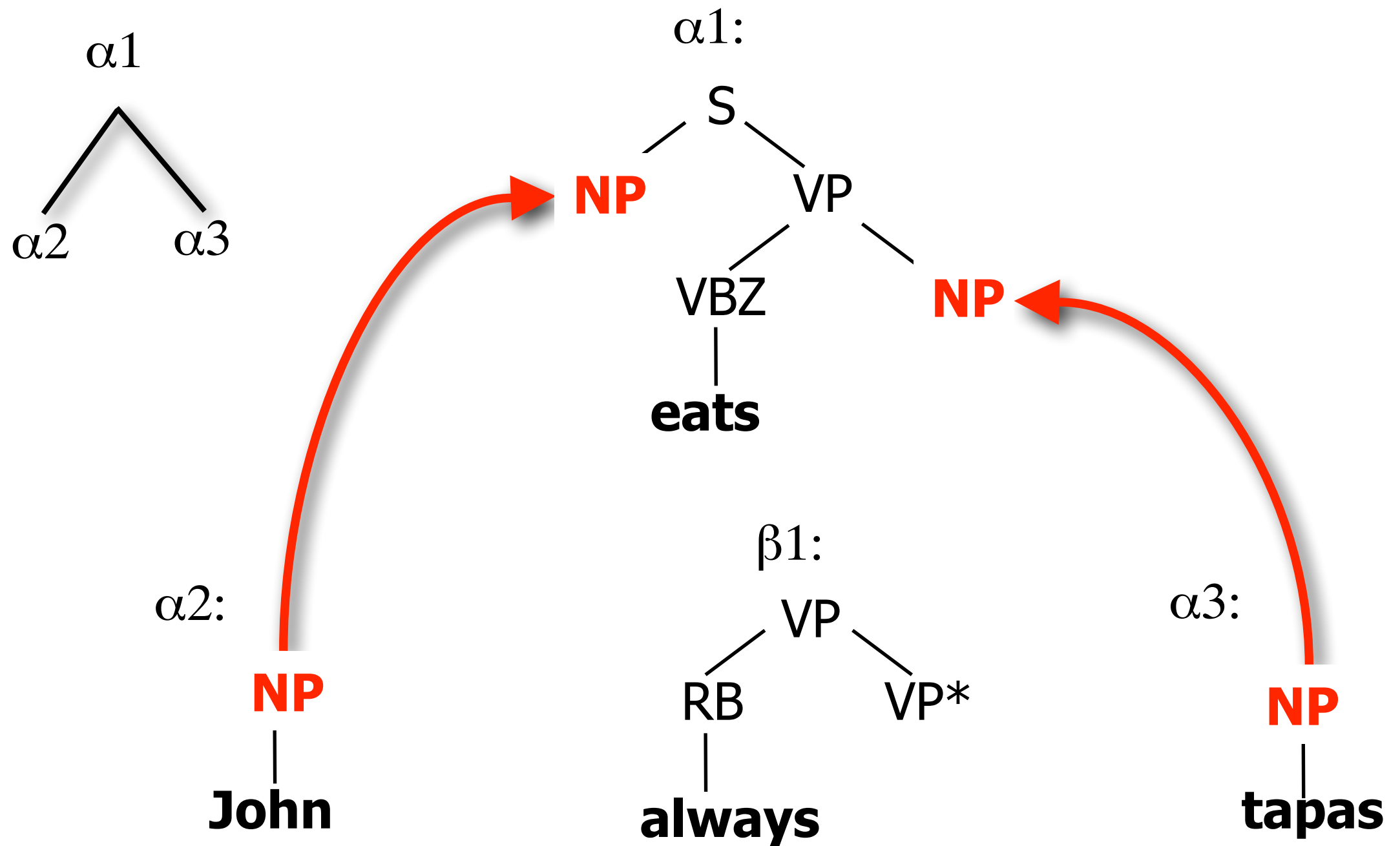
β_1 :



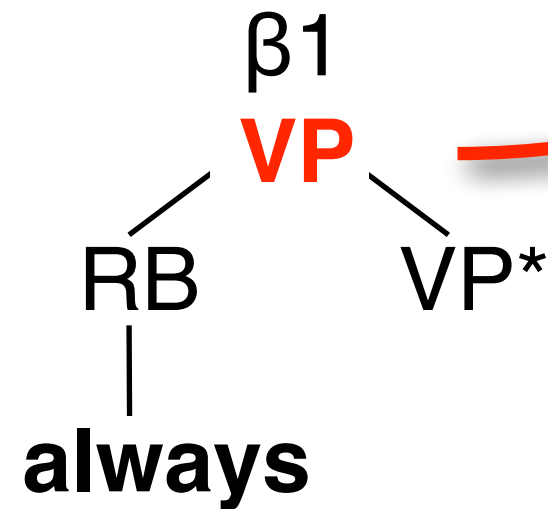
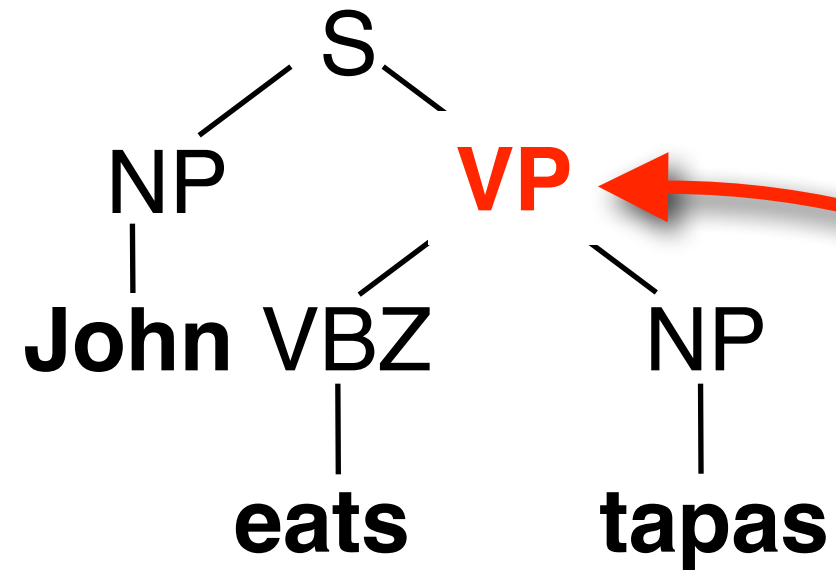
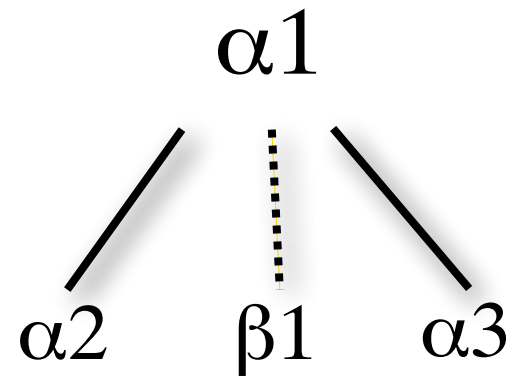
α_3 :



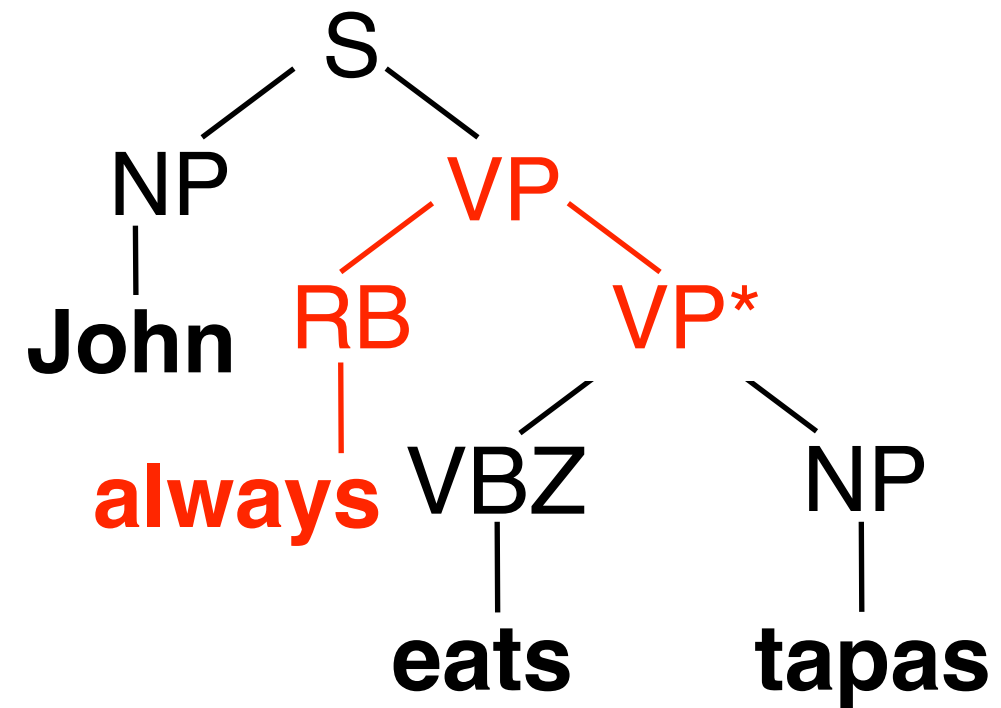
A TAG derivation



A TAG derivation

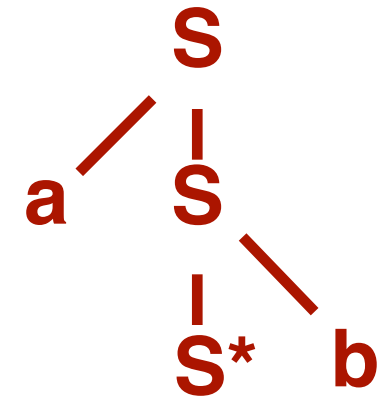
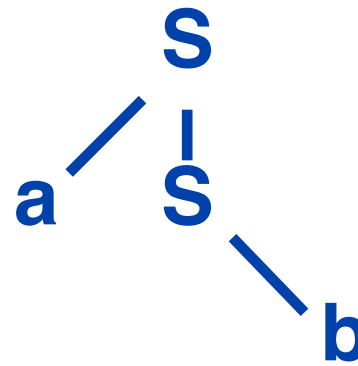


A TAG derivation

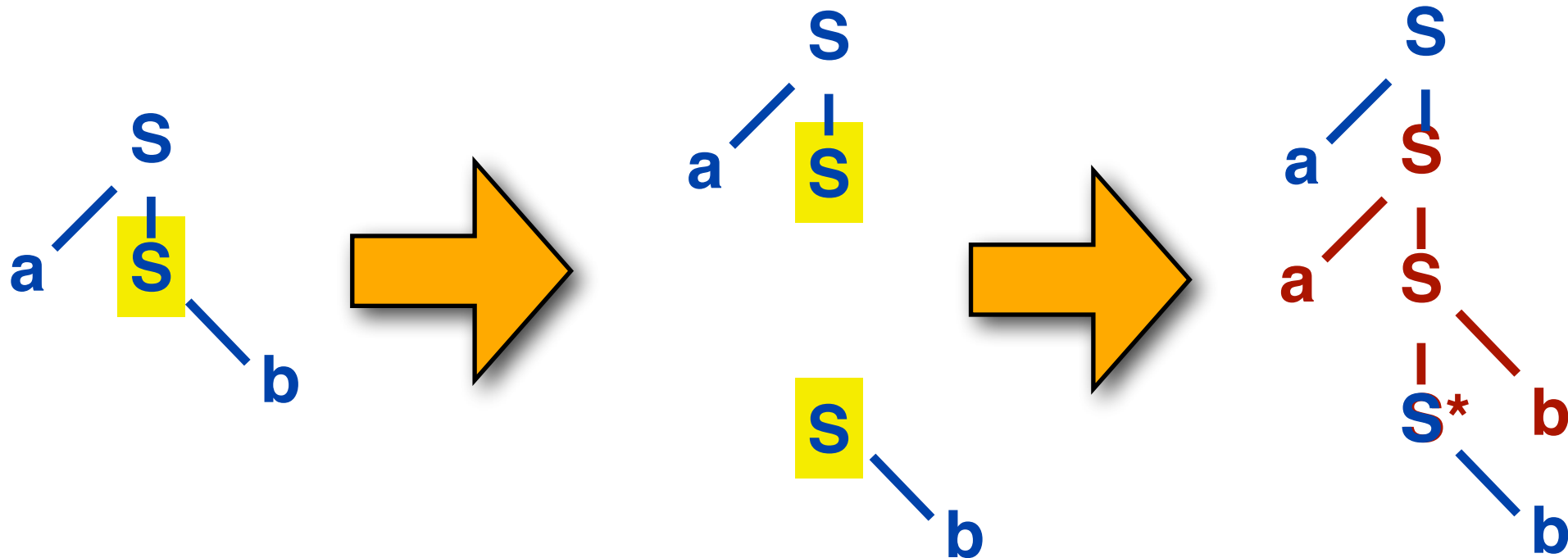


$a^n b^n$: Cross-serial dependencies

Elementary trees:



Deriving **aabb**



Feature Structure Grammars

Why feature structures

Feature structures form the basis for many grammar formalisms used in computational linguistics.

Feature structure grammars (aka **attribute-value grammars**, or **unification grammars**) can be used as

- a more compact way of representing rich CFGs
- a way to represent more expressive grammars

Simple grammars overgenerate

$S \rightarrow NP VP$
 $VP \rightarrow Verb NP$
 $NP \rightarrow Det Noun$
 $Det \rightarrow the \mid a \mid these$
 $Verb \rightarrow eat \mid eats$
 $Noun \rightarrow cake \mid cakes \mid student \mid students$

This generates ungrammatical sentences like
“these student eats a cakes”

We need to capture (number/person) agreement

Refining the nonterminals

$S \rightarrow NP_{sg} VP_{sg}$

$S \rightarrow NP_{pl} VP_{pl}$

$VP_{sg} \rightarrow VerbSg NP$

$VP_{pl} \rightarrow VerbPl NP$

$NP_{sg} \rightarrow DetSg NounSg$

$DetSg \rightarrow the \mid a$

... ..

This yields **very large grammars**.

What about person, case, ...?

Difficult to capture **generalizations**.

Subject and verb have to have number agreement

NP_{sg} , NP_{pl} and NP are three distinct nonterminals

Feature structures

Replace atomic categories with feature structures:

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{SG} \\ \text{PERS} & 3 \\ \text{CASE} & \text{NOM} \end{bmatrix}$$
$$\begin{bmatrix} \text{CAT} & \text{VP} \\ \text{NUM} & \text{SG} \\ \text{PERS} & 3 \\ \text{VFORM} & \text{FINITE} \end{bmatrix}$$

A **feature structure** is a list of **features** (= attributes), e.g. CASE, and **values** (eg NOM).

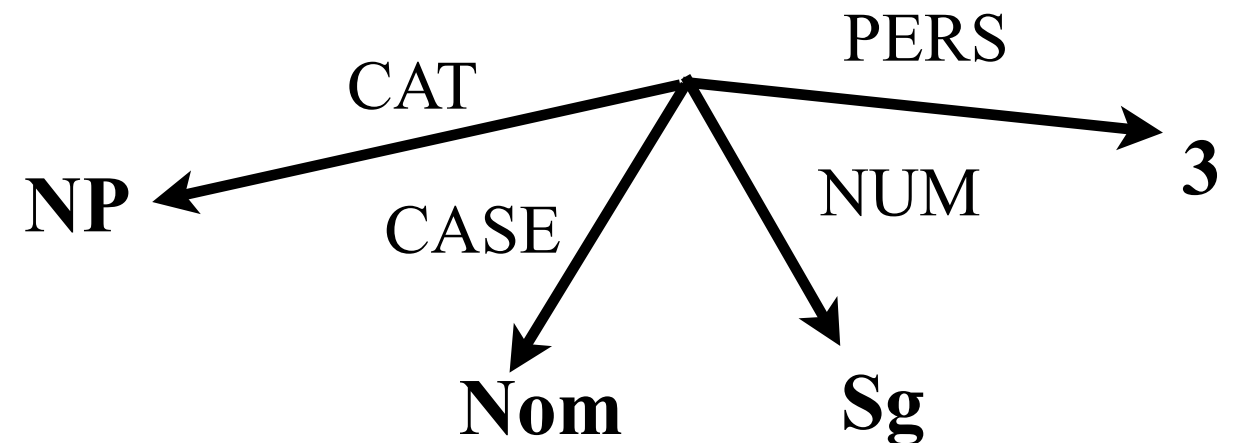
We often represent feature structures as **attribute value matrices (AVM)**

Usually, values are **typed** (to avoid CASE:SG)

Feature structures as directed graphs

| | |
|------|-----|
| CAT | NP |
| NUM | SG |
| PERS | 3 |
| CASE | NOM |

=



Complex feature structures

We distinguish between **atomic** and **complex** feature values.

A complex value is a feature structure itself.

This allows us to capture better generalizations.

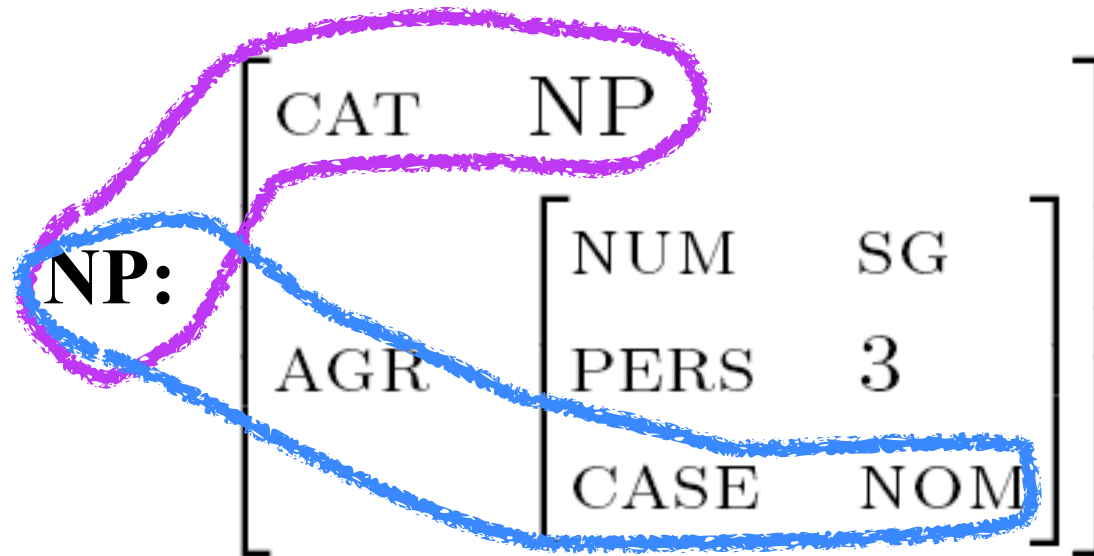
Only atomic values:

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{SG} \\ \text{PERS} & 3 \\ \text{CASE} & \text{NOM} \end{bmatrix}$$

Complex values:

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{SG} \\ \text{PERS} & 3 \\ \text{CASE} & \text{NOM} \end{bmatrix} \end{bmatrix}$$

Feature paths



A **feature path** allows us to identify particular values in a feature structure:

$\langle \mathbf{NP} \ \mathbf{CAT} \rangle = \mathbf{NP}$

$\langle \mathbf{NP} \ \mathbf{AGR} \ \mathbf{CASE} \rangle = \mathbf{NOM}$

Unification

Two **feature structures A and B unify** ($A \sqcup B$) if they can be merged into one consistent feature structure C:

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{SG} \\ \text{CASE} & \text{NOM} \end{bmatrix} \sqcup \begin{bmatrix} \text{CAT} & \text{NP} \\ \text{PERS} & 3 \end{bmatrix} = \begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{SG} \\ \text{PERS} & 3 \\ \text{CASE} & \text{NOM} \end{bmatrix}$$

Otherwise, unification **fails**:

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{SG} \\ \text{CASE} & \text{NOM} \end{bmatrix} \sqcup \begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{PL} \end{bmatrix} = \emptyset$$

PATR-II style feature structures

CFG rules are augmented with constraints:

$$\mathbf{A}_0 \rightarrow \mathbf{A}_1 \dots \mathbf{A}_n \\ \{\text{set of constraints}\}$$

There are two kinds of constraints:

Unification constraints:

$$\langle \mathbf{A}_i \text{ feature-path} \rangle = \langle \mathbf{A}_j \text{ feature-path} \rangle$$

Value constraints:

$$\langle \mathbf{A}_i \text{ feature-path} \rangle = \text{atomic value}$$

A grammar with feature structures

| | | | | |
|-------------|--|----------|---|--------------------|
| S | → NP VP | | Grammar rule | |
| | $\langle \mathbf{NP} \mathit{NUM} \rangle$ | = | $\langle \mathbf{VP} \mathit{NUM} \rangle$ | Constraints |
| | $\langle \mathbf{NP} \mathit{CASE} \rangle$ | = | <i>nom</i> | |
| NP | → DT NOUN | | Grammar rule | |
| | $\langle \mathbf{NP} \mathit{NUM} \rangle$ | = | $\langle \mathbf{NOUN} \mathit{NUM} \rangle$ | Constraints |
| | $\langle \mathbf{NP} \mathit{CASE} \rangle$ | = | $\langle \mathbf{NOUN} \mathit{CASE} \rangle$ | |
| NOUN | → <i>cake</i> | | Lexical entry | |
| | $\langle \mathbf{NOUN} \mathit{NUM} \rangle$ | = | <i>sg</i> | Constraints |

With complex feature structures

| | | | | |
|-------------|---|----------|--|--------------------|
| S | → NP VP | | Grammar rule | |
| | $\langle \mathbf{NP} \mathit{AGR} \rangle$ | = | $\langle \mathbf{VP} \mathit{AGR} \rangle$ | Constraints |
| | $\langle \mathbf{NP} \mathit{CASE} \rangle$ | = | <i>nom</i> | |
| NP | → DT NOUN | | Grammar rule | |
| | $\langle \mathbf{NP} \mathit{AGR} \rangle$ | = | $\langle \mathbf{NOUN} \mathit{AGR} \rangle$ | Constraints |
| NOUN | → <i>cake</i> | | Lexical entry | |
| | $\langle \mathbf{NOUN} \mathit{AGR} \mathit{NUM} \rangle$ | = | <i>sg</i> | Constraints |

Complex feature structures capture better generalizations (and hence require fewer constraints) — cf. the previous slide

The head feature

Instead of implicitly specifying heads for each rewrite rule, let us define a **head feature**.

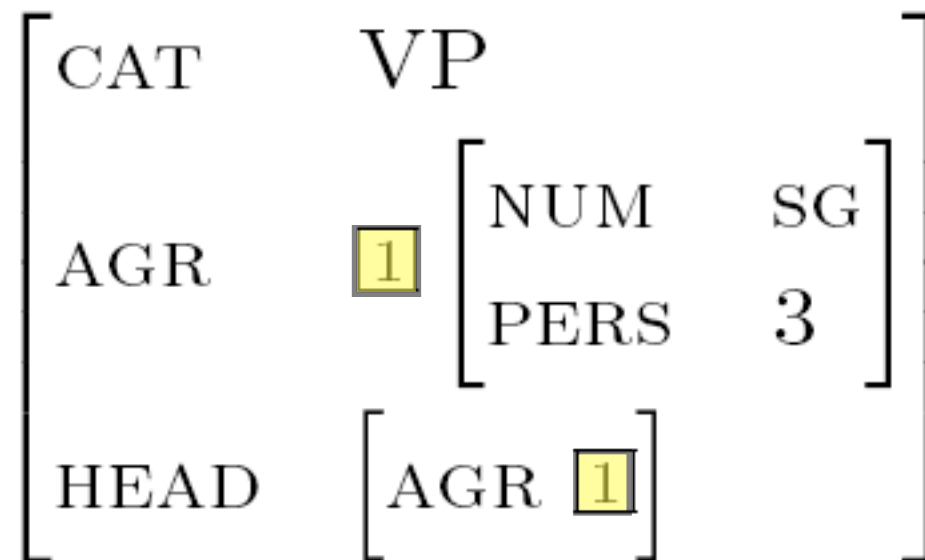
The head of a VP has the same agreement feature as the VP itself:

$$\left[\begin{array}{l} \text{CAT} \\ \text{AGR} \\ \text{HEAD} \end{array} \right. \left. \begin{array}{l} \text{VP} \\ \left[\begin{array}{ll} \text{NUM} & \text{SG} \\ \text{PERS} & 3 \end{array} \right] \\ \left[\begin{array}{l} \text{AGR} \\ \left[\begin{array}{ll} \text{NUM} & \text{SG} \\ \text{PERS} & 3 \end{array} \right] \end{array} \right] \end{array} \right]$$

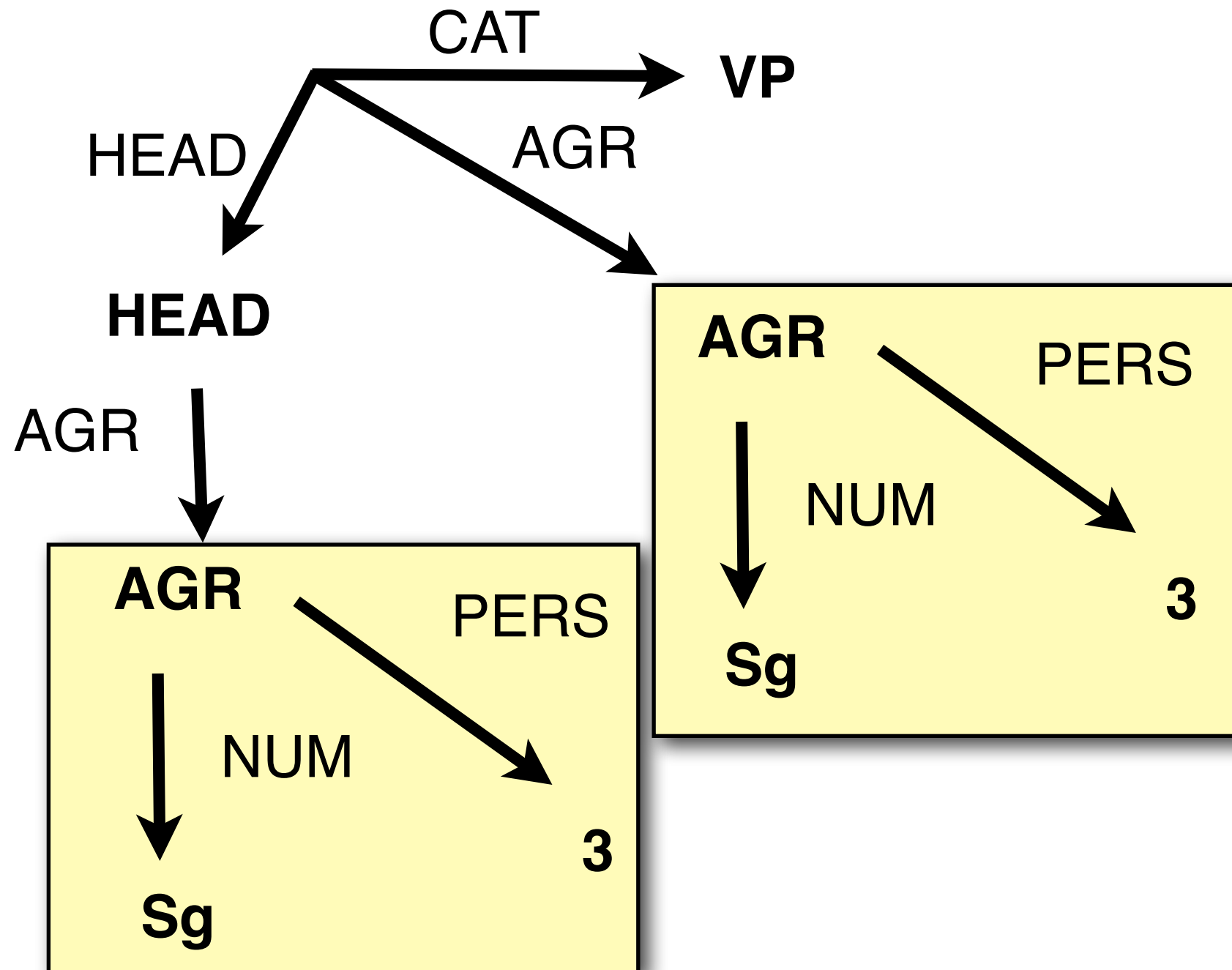
Re-entrancies

What we *really* want to say is that the agreement feature of the head is *identical* to that of the VP itself.

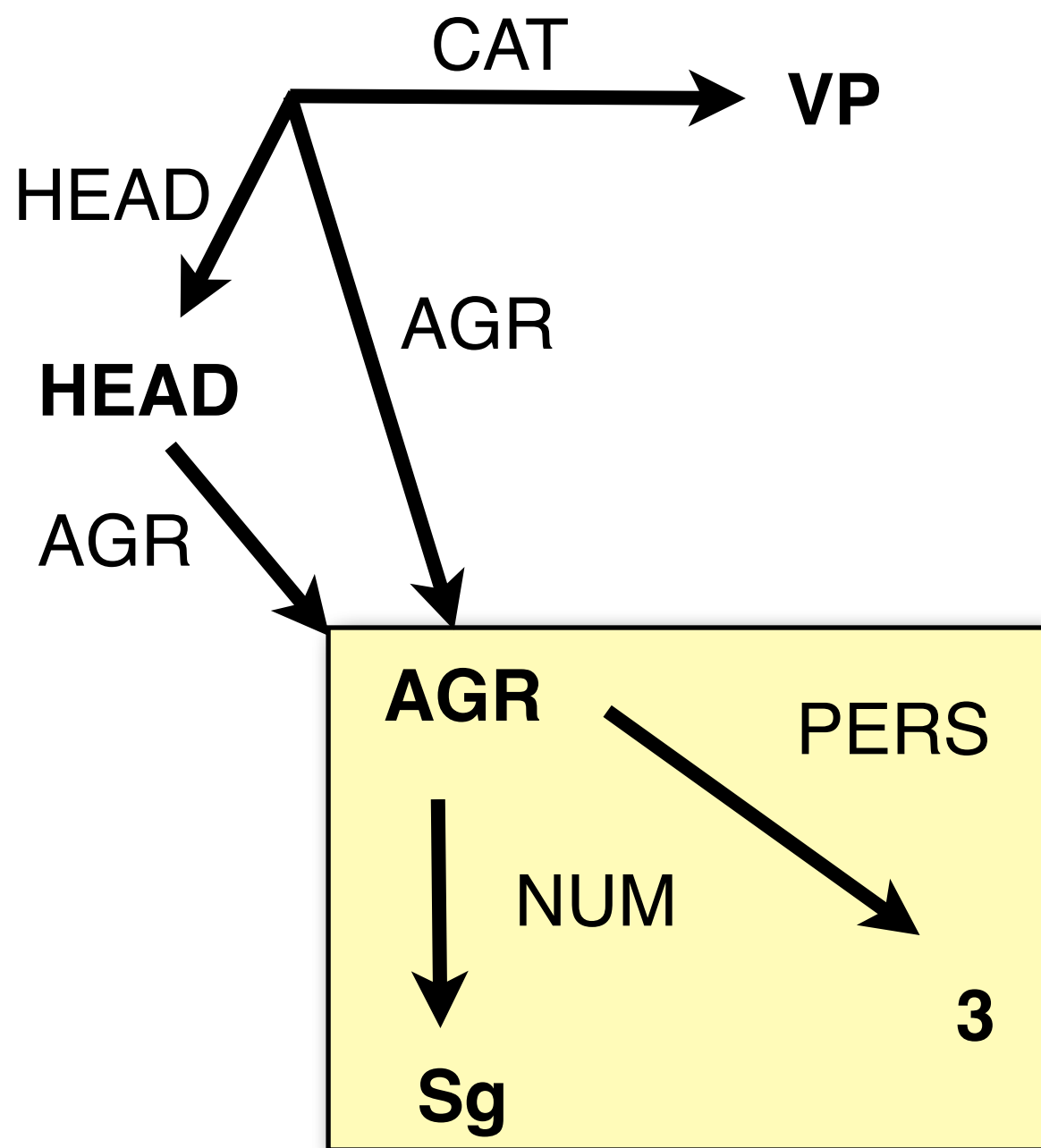
This corresponds to a **re-entrancy** in the FS (indicated via coindexation $\boxed{1}$)



Re-entrancies - not like this:



Re-entrancies - but like this:



Attribute-Value Grammars and CFGs

If every feature can only have **a finite set of values**, any attribute-value grammar can be compiled out into a (possibly huge) context-free grammar

Going beyond CFGs

The power-of-2 language: $L_2 = \{w^i \mid i \text{ is a power of } 2\}$

L_2 is a (fully) context-sensitive language.

(*Mildly* context-sensitive languages have the **constant growth property** (the length of words always increases by a constant factor c))

Here is a feature grammar which generates L_2 :

$$A \rightarrow a$$

$$\langle A \ F \rangle = 1$$

$$A \rightarrow A_1 \ A_2$$

$$\langle A \ F \rangle = \langle A_1 \rangle$$

$$\langle A \ F \rangle = \langle A_2 \rangle$$