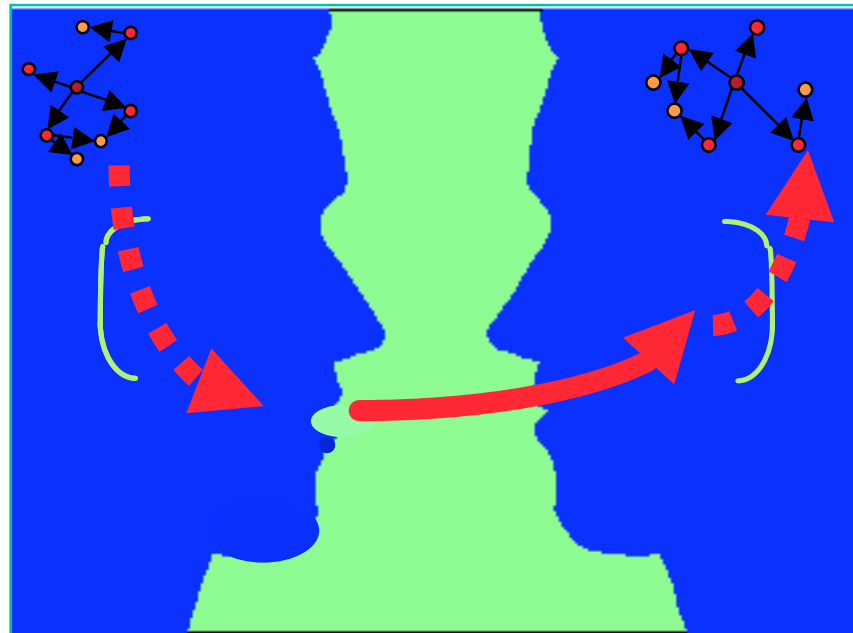

FLST

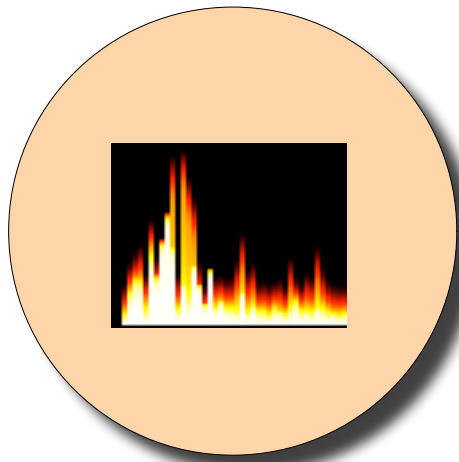
Grammars and Parsing

Hans Uszkoreit

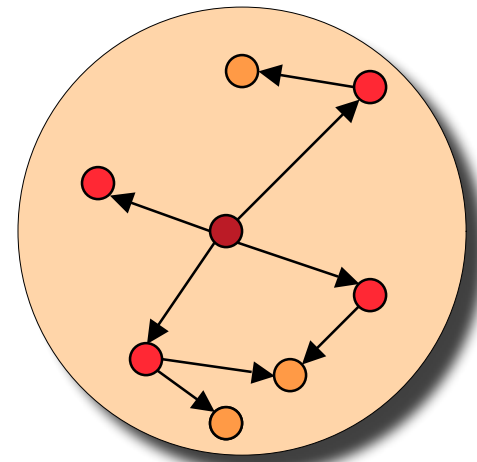
WHAT HAPPENS IN BETWEEN?



WHAT HAPPENS IN BETWEEN?

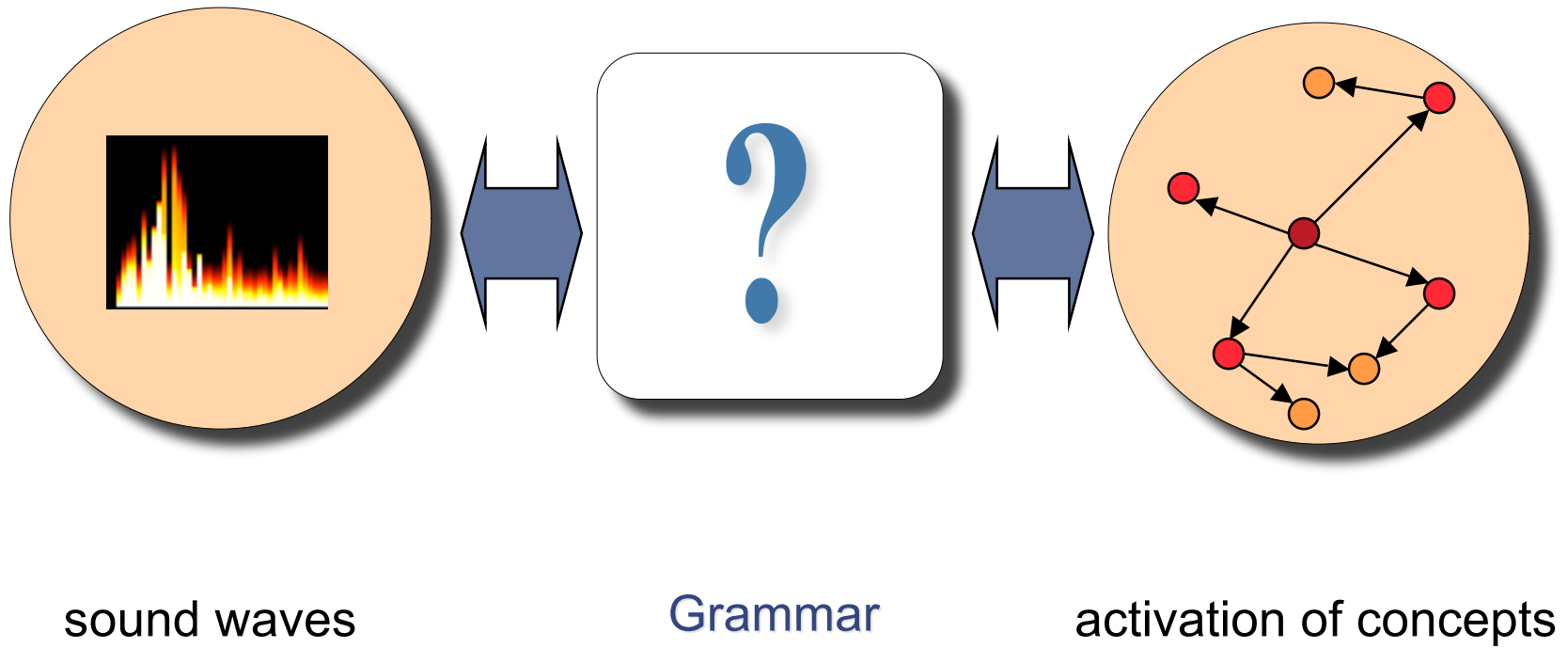


sound waves

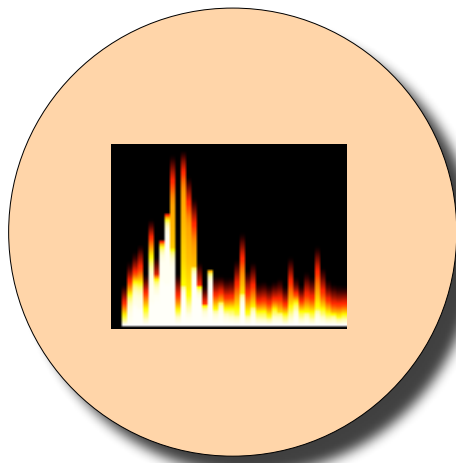


activation of concepts

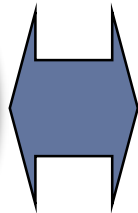
WHAT HAPPENS IN BETWEEN?



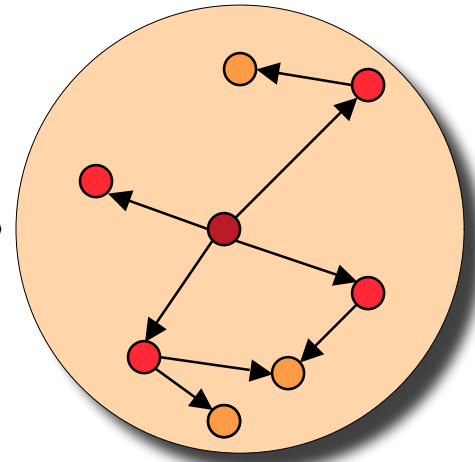
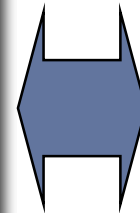
WHAT HAPPENS IN BETWEEN?



sound waves

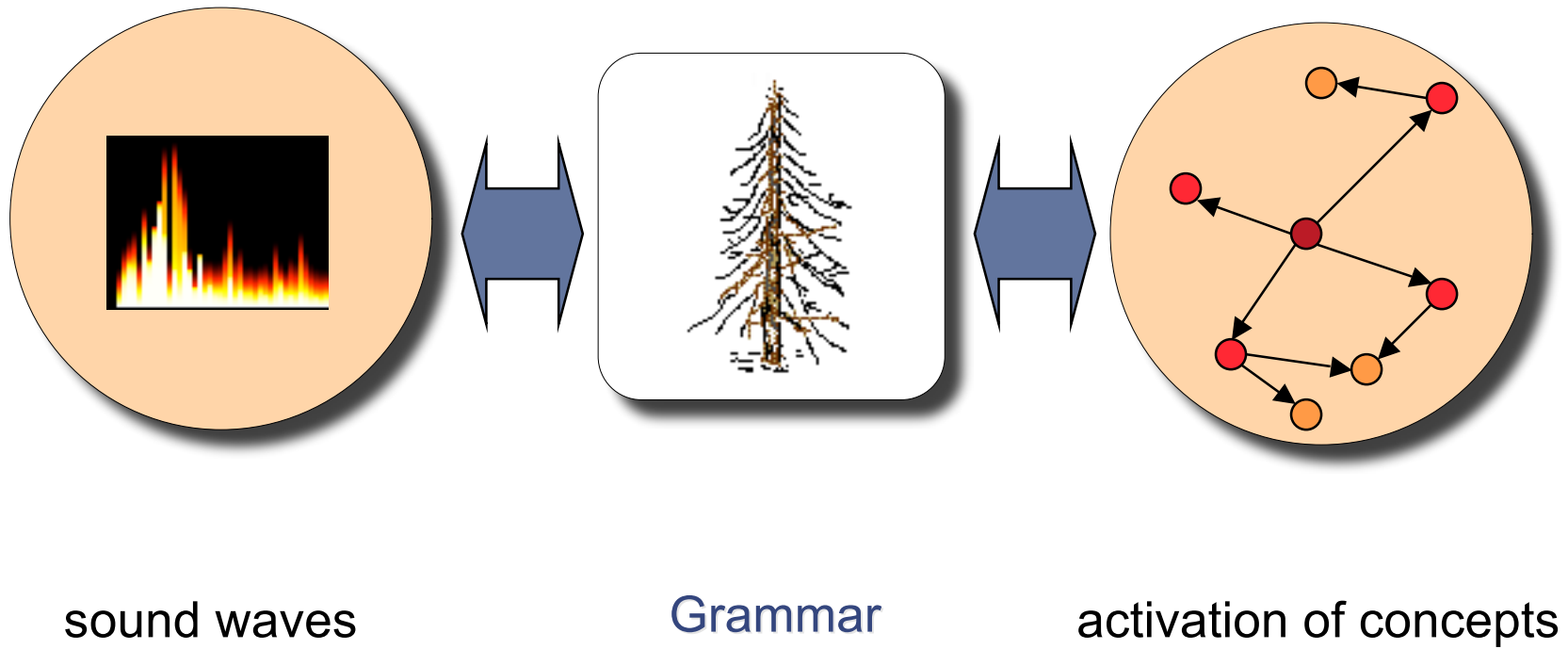


Grammar

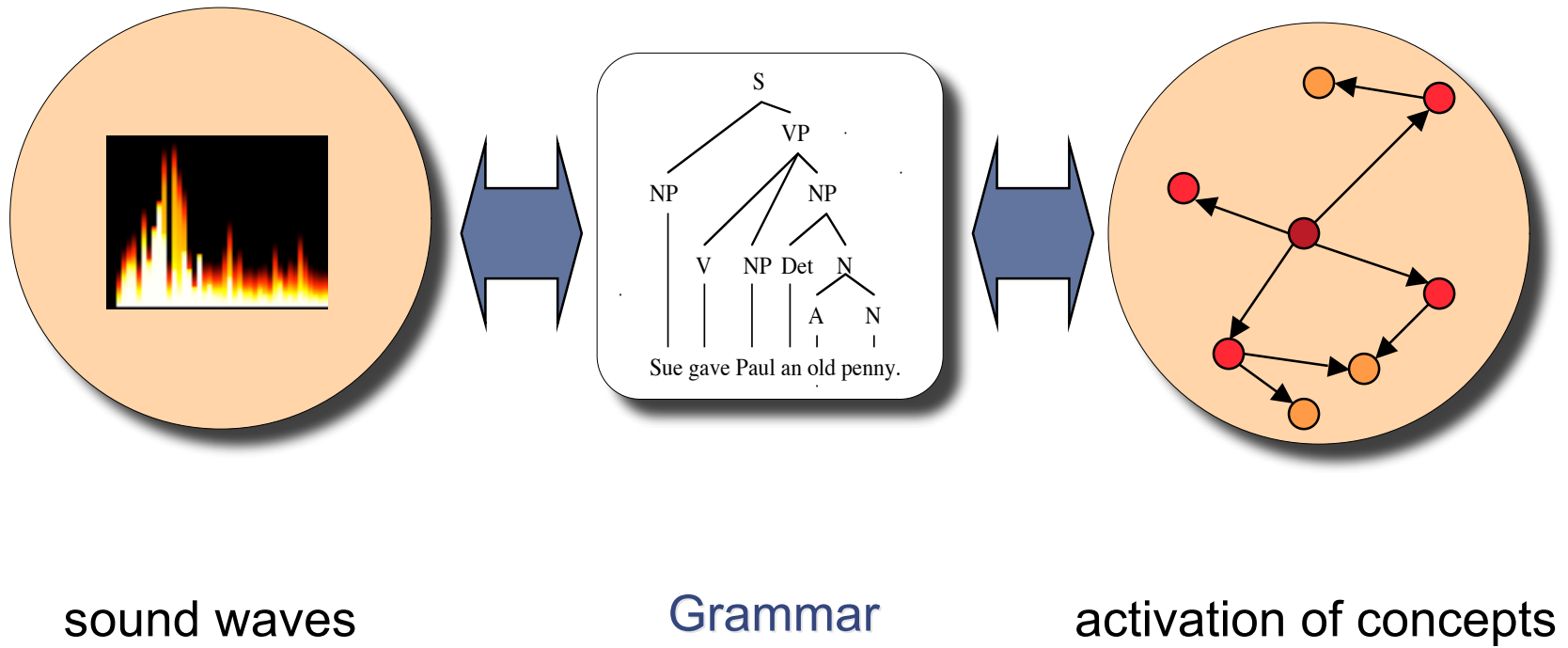


activation of concepts

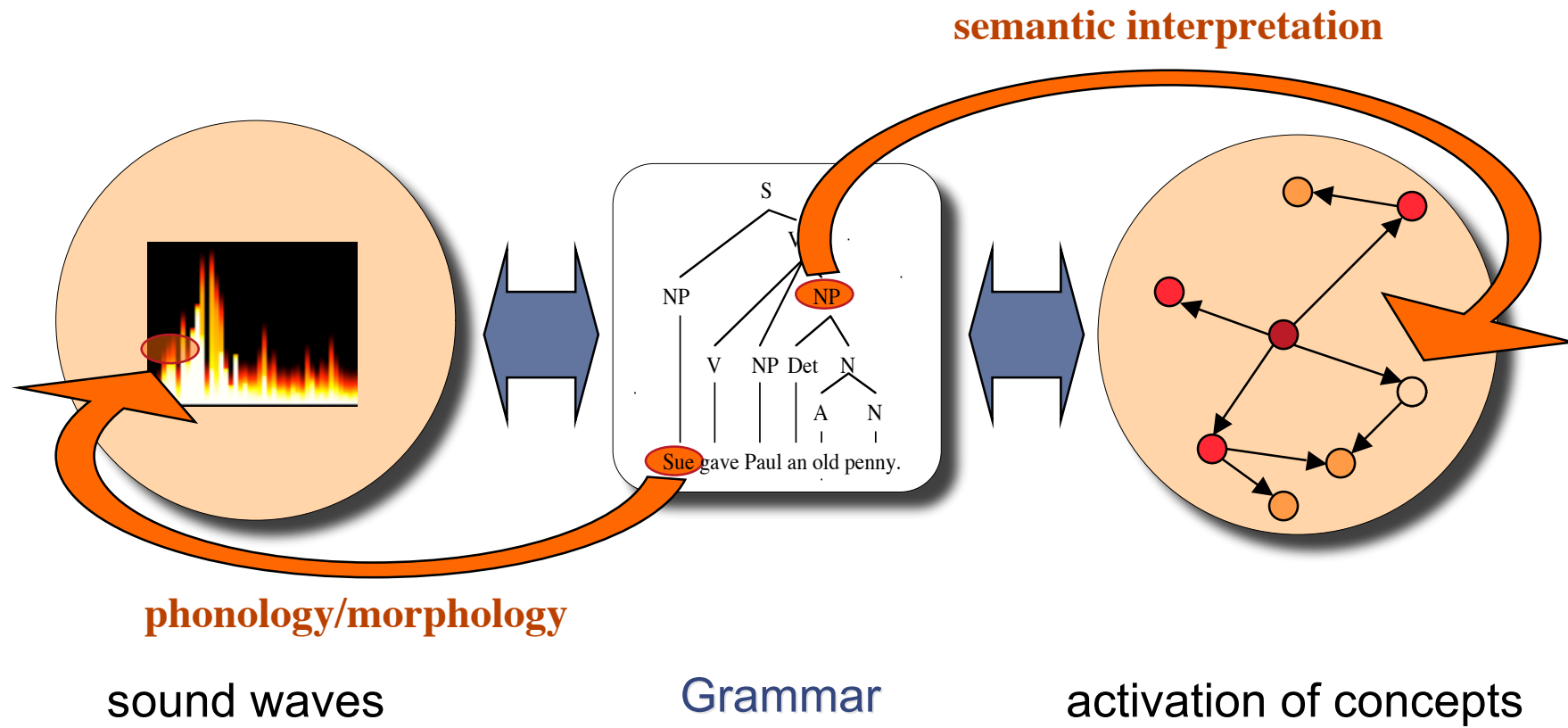
WHAT HAPPENS IN BETWEEN?



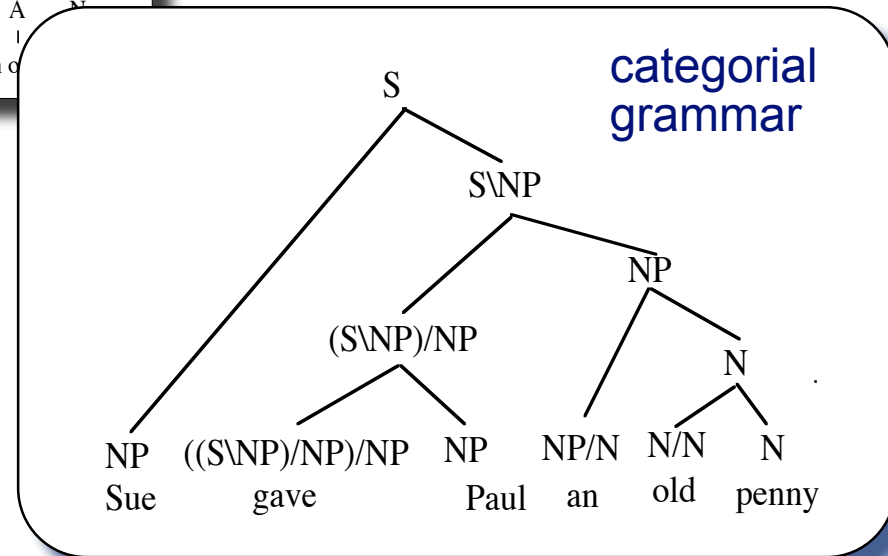
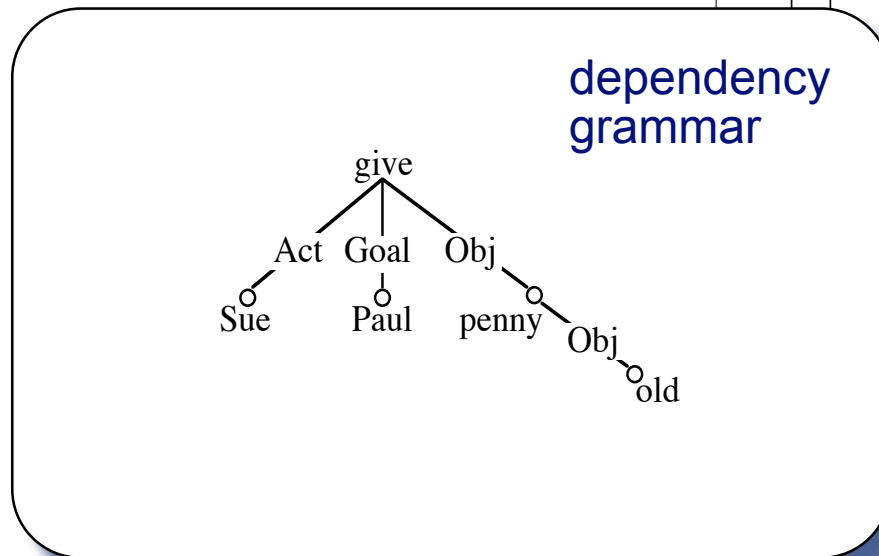
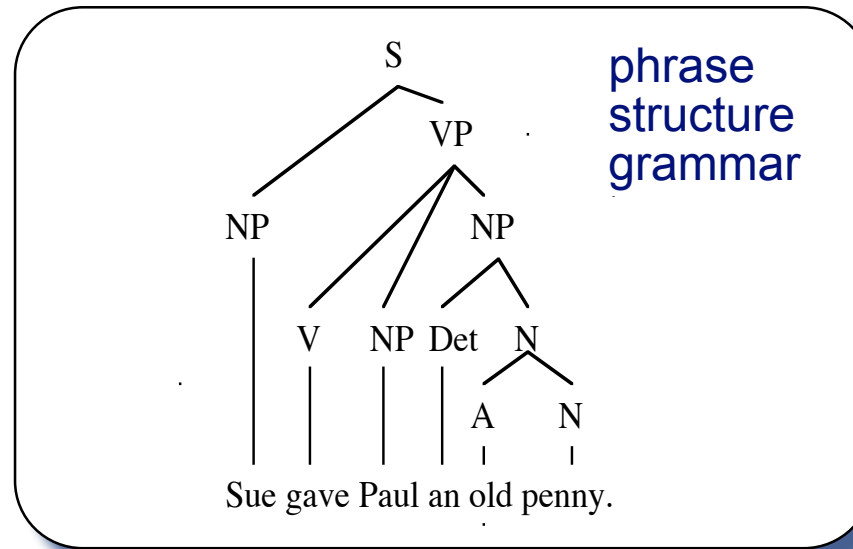
WHAT HAPPENS IN BETWEEN?

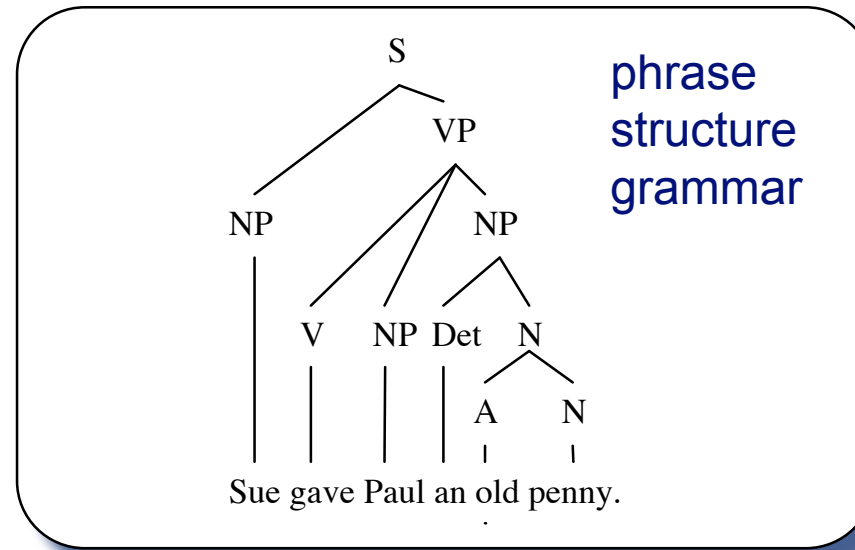


WHAT HAPPENS IN BETWEEN?

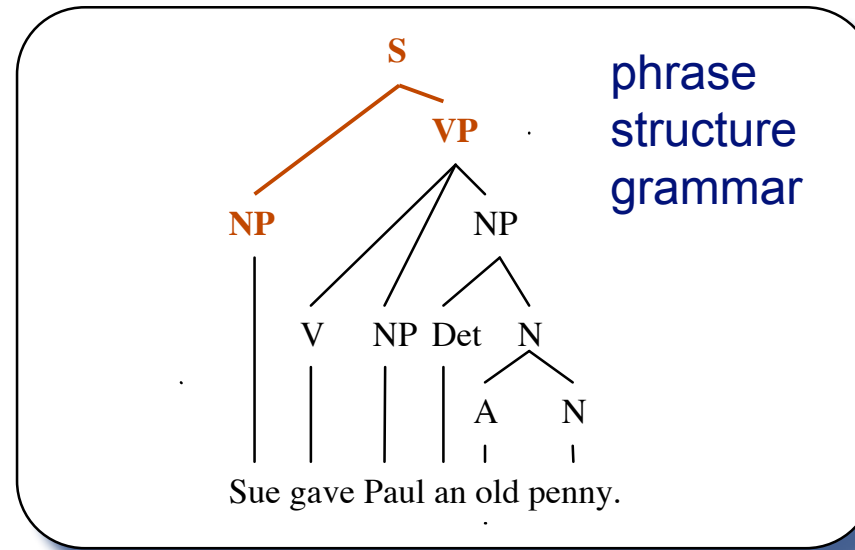


THREE TRADITIONS

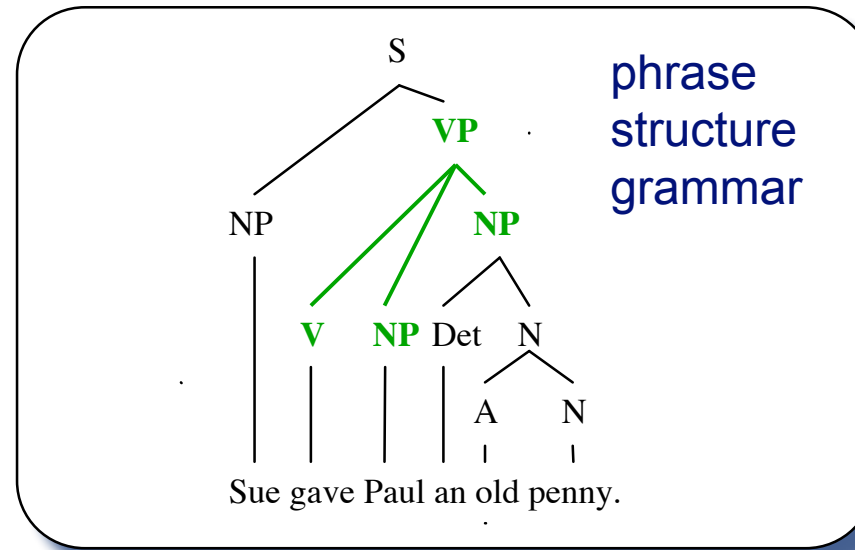




$S \rightarrow NP VP$

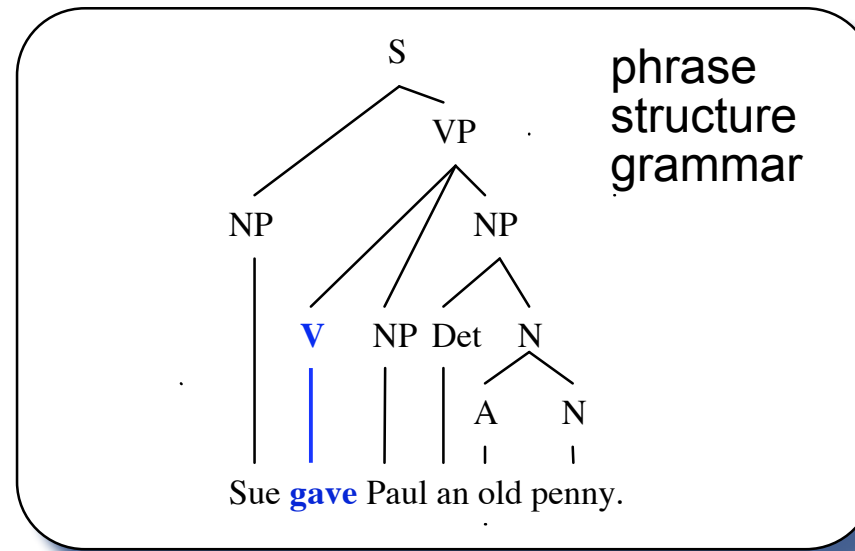


$S \rightarrow NP VP$



$S \rightarrow NP VP$
 $VP \rightarrow V NP NP$

Grammar



$S \rightarrow NP VP$
 $VP \rightarrow V NP NP$

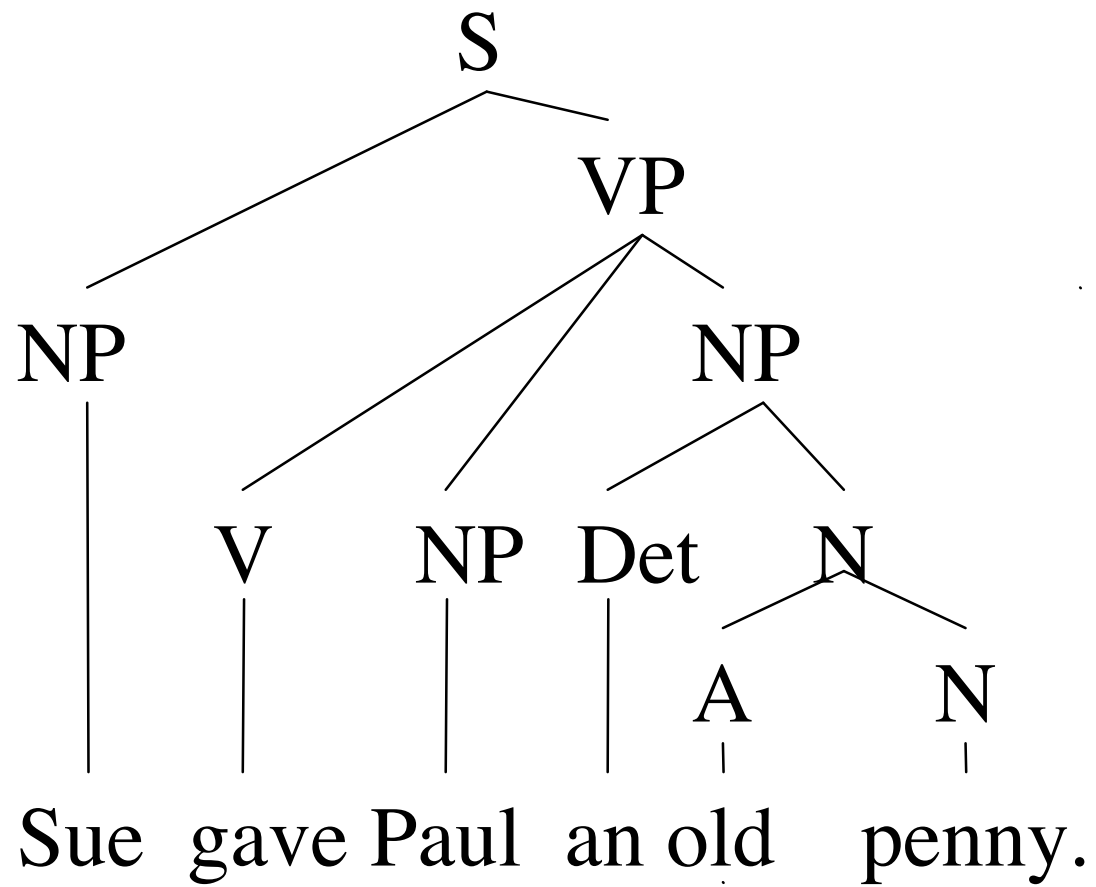
$V \rightarrow \text{gave}$

symbol conventions

	single symbols	strings
non-terminal	A, B, C, \dots	\dots, X, Y, Z
terminals	a, b, c, \dots	\dots, x, y, z
unspecified symbols	$\alpha, \beta, \gamma, \dots$	$\dots, \phi, \chi, \psi, \omega$
start symbol	S	
empty string	ϵ	
integers	$\dots, i, j, k, l, m, n, \dots$	

Why Syntax

- Einen Hund hat dieser Mann gebissen.
- Ein Hund hat diesen Mann gebissen.
- This man has bitten a dog.
- A dog this man has bitten.
- A dog has bitten this man.
- Peter promised Paul, to process the files.
- Peter persuaded Paul, to process the files.



Formal Grammar

A language over an alphabet (vocabulary) Σ is a subset of Σ^* .

A formal grammar G_L for a language L is a quadruple $(V_N, V_T, \{S\}, P)$.

V_N - non-terminal vocabulary (auxiliary vocabulary)

V_T - terminal vocabulary

$$(V_T \cap V_N = \emptyset, \quad L \subseteq V_T^*, \quad V = V_T \cup V_N)$$

$\{S\}$ - singleton with the start symbol (set of axioms)

P - set of productions, rule set

set of rules of the form $\omega_1 \phi \omega_2 \rightarrow \omega_1 \psi \omega_2$

usually written as $\phi \rightarrow \psi$

derivation

relation "follows":

If $G = (V_N, V_T, \{S\}, P)$, then ψ follows from ϕ according to G iff there are strings $\phi_1, \phi_2, \chi, \omega$, so that $\phi = \phi_1 \chi \phi_2$ und $\psi = \phi_1 \omega \phi_2$ und $\chi \rightarrow \omega \in P$.

Notation: $\phi \Rightarrow_G \psi$

derivation:

A sequence of strings $\phi_1, \phi_1, \dots, \phi_n$ is a derivation according to G iff $\phi_i \Rightarrow_G \phi_{i+1}$ for all $i, 1 \leq i \leq n$.

If there is *derivable* according to G from ϕ to ψ we can write this: $\phi \xRightarrow{*}_G \psi$

The relation *derivable* is transitive and is moreover defined to be reflexive.

The generated language

The language L: A string ω is in L according to G_L iff the following three conditions are fulfilled:

1. $\omega \in V_T^*$
2. $S \xRightarrow[G]{*} \omega$
3. There is no χ , so that $\omega \xRightarrow[G]{*} \chi$ and $\omega \neq \chi$.

We say that G_L generates the language L. The language L generated by G is also written as $L(G)$.

Weak Equivalence: Two grammar G_1 and G_2 are weakly equivalent, if they generate the same language.

Types of Grammars

Type 0 (unrestricted rewriting systems):

Every formal grammar according to the definition is of type 0.

Type 1 (context sensitive grammars):

Every production is of the form $\phi A \psi \rightarrow \phi \omega \psi$, where $A \in V_N$, $\omega \neq \varepsilon$.

Type 2 (context free grammars):

Every production is of the form $A \rightarrow \omega$, where $\omega \neq \varepsilon$.

Type 3 (regular grammars):

Every production is of the form $A \rightarrow x B$ or $A \rightarrow x$, where $x \neq \varepsilon$.

context-free derivations

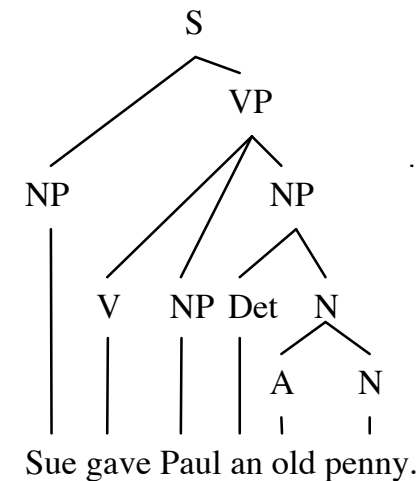
S
NP VP
DET ADJ N VP
DET ADJ N V NP
DET ADJ N V DET ADJ N
.
.
.
ein kleines Mädchen sucht ein kleines Mädchen

Trees

the notion of syntactic constituent tree

coded information

1. the hierarchic organisation of a sentence in terms of constituents
2. the assignment of each constituent to a linguistic class (category)
3. the linear sequence of the constituents



relations: immediate dominance - dominance
immediate precedence - precedence

constituent structure tree: quintuple(N, Q, D, P, L)

N - finite set of nodes

Q - finite set of labels

D - weak partial order in $N \times N$, the dominance relation
(reflexive, transitive und antisymmetric)

P - strong partial order in $N \times N$, the precedence relation
(irreflexive, transitive und asymmetric)

L - function from N into Q , the labelling function

conditions:

(single) root condition

exclusivity condition

no crossing condition / no tangling condition

root condition

There is exactly one node that dominates all other nodes of the tree.

exclusivity condition

For any two nodes x and y holds

either	$D(x,y)$ or $D(y,x)$
or	$P(x,y)$ or $P(y,x)$
but never both.	

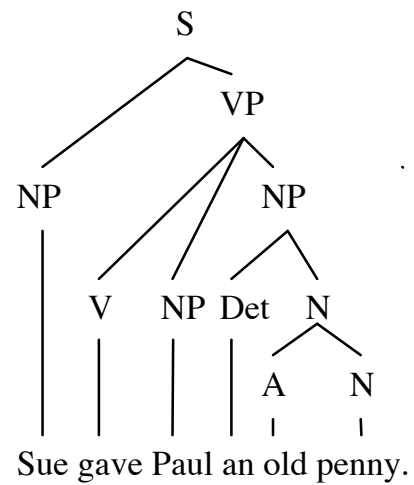
no-tangling condition

If $P(x,y)$ then for all x' dominated by x [$D(x,x')$] and for all y' dominated by y [$D(y, y')$] that x' precedes y' [$P(x', y')$]

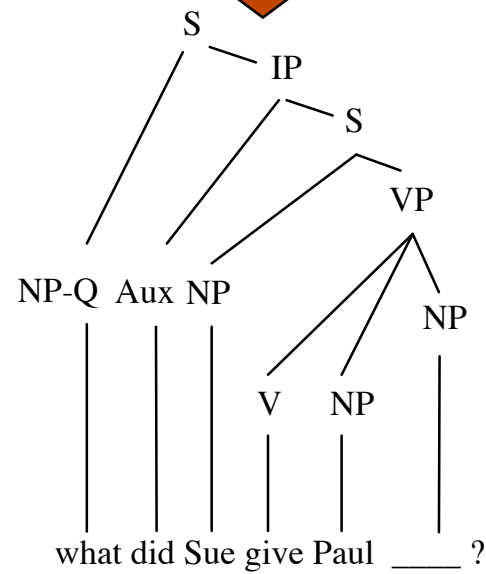
Parsing

- The syntactic analysis of strings according to a grammar we call parsing
- input: terminal string
output: structure of the sentence, i.e. all constituents usually a tree
- parsing algorithms:
top-down vs bottom-up
left-right vs right left vs birectional or island parsing
deterministic vs. non-deterministic

Grammar



transformational
grammar



Unification grammar

