

# Tecnologie del Linguaggio Naturale

## Parte Prima

### Lezione n. 04-1

Sintassi: la competence

15-03-2021  
relat. tra le varie parole  
└ di gruppo ← e costituenti  
└ puntuali ← e dipendente

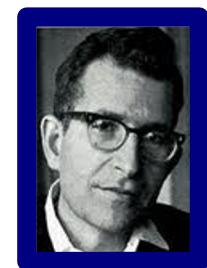
# Prologo

This is the cat that caught the rat that stole the cheese. ↗  
↓

The man who wrote the book that you told me about

*Aspects of the theory of syntax* (1957)

N. Chomsky



# Prologo

---

- English Center-Embedding

↑  
versione "centrale"



A man that a woman loves

A man that a woman that a child knows loves

A man that a woman that a child that a bird saw knows loves

A man that a woman that a child that a bird that I heard saw knows loves

# Prologo

---

E venne il macellaio,  
che uccise il toro,  
che bevve l'acqua,  
che spense il fuoco,  
che bruciò il bastone,  
che picchiò il cane,  
che morse il gatto,  
che si mangiò il topo  
che al mercato mio padre comprò.



Alla fiera dell'est (1976)

A. Branduardi

# Prologo

---

Elena e Maria odiano il latte e il caffè

rispettivamente

*importante*

# Outline

---

- 1 ○ Sintassi e grammatiche generative ↗ cfp LEFT
- 2 ○ La gerarchia di Chomsky e il linguaggio naturale  
*classe di linguaggi che contiene le lingue naturali*
- 3 ○ Mildly-Context Sensitive Languages, Tree Adjoining Grammars, Combinatory Categorial Grammars

# Outline

---

- Sintassi e grammatiche generative
- La gerarchia di Chomsky e il linguaggio naturale
- Mildly-Context Sensitive Languages, Tree Adjoining Grammars, Combinatory Categorial Grammars

# NLP pipeline

---

**Phonetics** acoustic and perceptual elements

**Phonology** inventory of basic sounds (phonemes) and basic rules for combination, e.g. vowel harmony

**Morphology** how morphemes combine to form words, relationship of phonemes to meaning

**Syntax** sentence formation, word order and the formation of constituents from word groupings

**Semantics** how do word meanings recursively compose to form sentence meanings (from syntax to logical formulas)

**Pragmatics** meaning that is not part of compositional meaning

# Syntax and Semantics

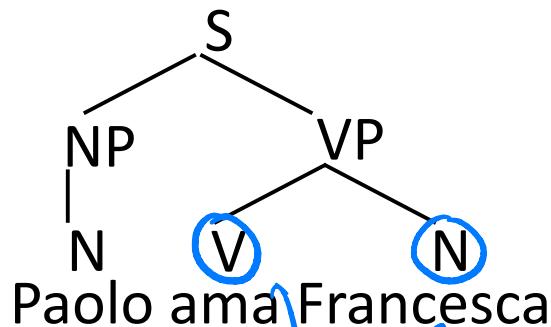
(Paolo (ama Francesca))

VS  
← la sintassi è importante!

Francesca ama Paolo

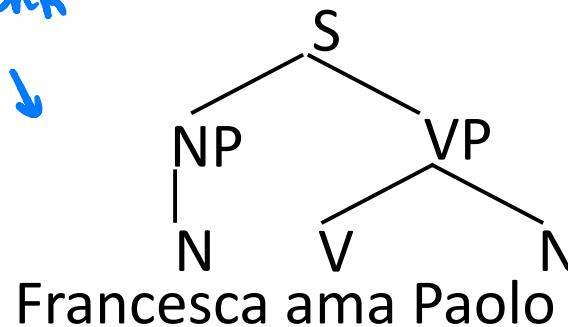
Syntactic Parsing

Syntactic Parsing



First of Species: provengono  
dal livello morfologico oppure  
dal livello sintattico stesso

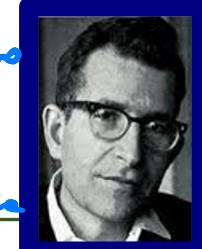
abbino a  
coaffuenti



· del livello morfologico oppure  
· del livello sintattico stesso

"la conoscenza statica che abbiamo nel cervello che riguarda la sintassi" → quella conoscenza linguistica che i nativi hanno

# Competence/Performance



↓ ↗ "la diff. tra avere una grammatica e usarla"

"Linguistic theory is concerned primarily with an ideal speaker-listener, in a completely homogeneous speech-communication, who know its (the speech community's) language perfectly and that it is unaffected by such **grammatically irrelevant conditions** as memory limitations, distractions, shifts of attention and interest, and errors (random or characteristic) in applying his knowledge of this language in actual performance."

R

competence  
"conscienza"  
vs performance  
"algoritmo"

la diff. tra teoria + pratica

# Competence/Performance in CL

---

- Competence = Grammatica Formale
- Performance = Algoritmo di Parsing

# Competence/Performance in CL

---

- Competence = Grammatica Formale
  - Performance = Algoritmo di Parsing
- > Prox. lezione: “Anatomia di un parser”

vegniamo una rapp. della  
convenzione FORNITE  
+  
+ che fanno  
+ regole

# Rewriting Systems

Emil Post , Alan Turing



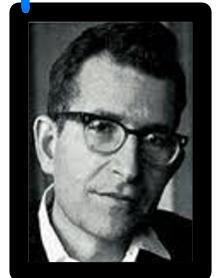
## Rewriting rule

$$\Psi \rightarrow \Theta$$



# Generative grammar

↳ è una grammatica formale: si basa sulle forme degli agg. che le compongono



$$G = (\Sigma, V, S, P)$$

$$\Sigma = \text{alphabet} \quad \xleftarrow{\text{dizionario}} \quad \text{simboli non terminali della grammatica}$$
$$V = \{A, B, \dots\}$$
$$S \in V \quad \xleftarrow{\text{simbolo di start}} \quad \text{insieme di regole di sostituzione/padroni}$$
$$P = \{\Psi \rightarrow \theta, \dots\} \quad \xleftarrow{\text{insieme delle produzioni}}$$

# Generative Grammars and Natural Languages

---

- ! • Generative Grammars models the natural languages as a formal languages
- The derivation tree can model the syntactic structure of the sentence

# Context-Free Grammars

$$G = (\Sigma, V, S, P)$$

- Constituency
- Grammatical relations
- Subcategorization

↑  
chomsky

$$A \rightarrow \beta$$

lo simbolo non terminale

lo posso usare per associare  
i nodi ai vari contenuti

x dare i tipi o  
alle parentesi

# Constituency



## Context-Free Grammars Constituency

Hypothesis:

Constituent  $\stackrel{\text{ASSOCIATION}}{<= >}$  non terminal symbols V

# Toy Grammar

$$G_4 = (\Sigma_4, \{S, NP, VP, V_1, V_2\}, S, P_4)$$

$$\Sigma_4 = \{I, Anna, John, Harry, saw, see, swimming\}$$

$$P_4 = \{S \rightarrow NP\ VP, VP \rightarrow V_1\ S, VP \rightarrow V_2,$$

$$NP \rightarrow I \mid John \mid Harry \mid Anna,$$

$$V_1 \rightarrow saw \mid see, V_2 \rightarrow swimming\}$$

... sono prime (non c'è ancora stato  
legato un significato)

# Toy Grammar

---

$S \rightarrow NP\ VP$

$S$

$VP \rightarrow V_1\ S$

$VP \rightarrow V_2$

$NP \rightarrow I | John | Harry | Anna$

$V_1 \rightarrow saw | see$

$V_2 \rightarrow swimming$

$S$

# Toy Grammar

→  $S \rightarrow NP\ VP$

$S \Rightarrow NP\ VP$

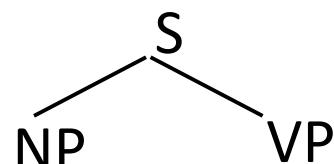
$VP \rightarrow V_1\ S$

$VP \rightarrow V_2$

$NP \rightarrow I | John | Harry | Anna$

$V_1 \rightarrow saw | see$

$V_2 \rightarrow swimming$



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

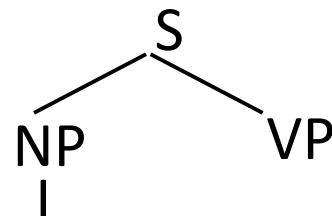
$VP \rightarrow V_2$

→  $NP \rightarrow I | John | Harry | Anna$

$V_1 \rightarrow saw | see$

$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP$



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

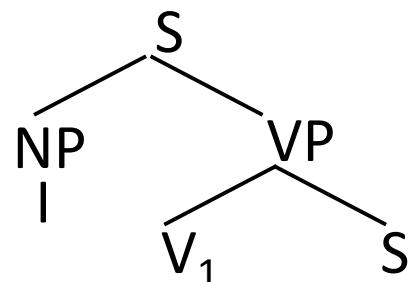
$VP \rightarrow V_2$

$NP \rightarrow I \mid John \mid Harry \mid Anna$

$V_1 \rightarrow saw \mid see$

$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1\ S$



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

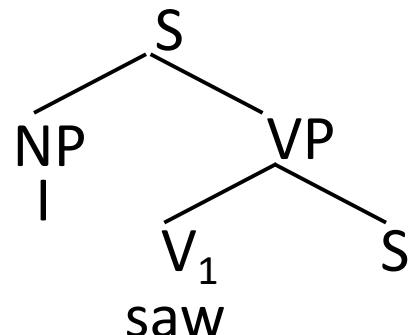
$VP \rightarrow V_2$

$NP \rightarrow I | John | Harry | Anna$

$V_1 \rightarrow saw | see$

$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1\ S \Rightarrow I\ saw\ S$



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

$VP \rightarrow V_2$

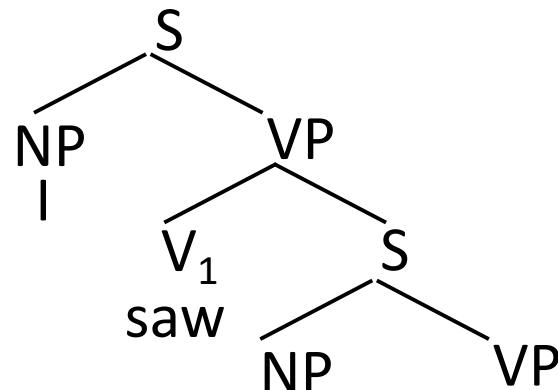
$NP \rightarrow I \mid John \mid Harry \mid Anna$

$V_1 \rightarrow saw \mid see$

$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1\ S \Rightarrow$

$I\ saw\ S \Rightarrow I\ saw\ NP\ VP$



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

$VP \rightarrow V_2$

$NP \rightarrow I | John | Harry | Anna$

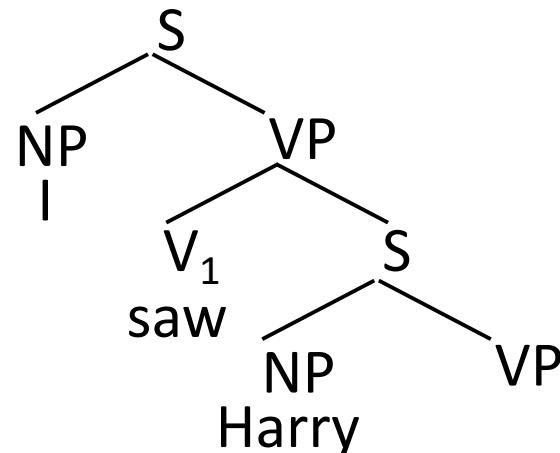
$V_1 \rightarrow saw | see$

$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1\ S \Rightarrow$

$I\ saw\ S \Rightarrow I\ saw\ NP\ VP \Rightarrow$

**I saw Harry VP**



# Toy Grammar

---

$S \rightarrow NP\ VP$

$VP \rightarrow V_1\ S$

$VP \rightarrow V_2$

$NP \rightarrow I \mid John \mid Harry \mid Anna$

$V_1 \rightarrow saw \mid see$

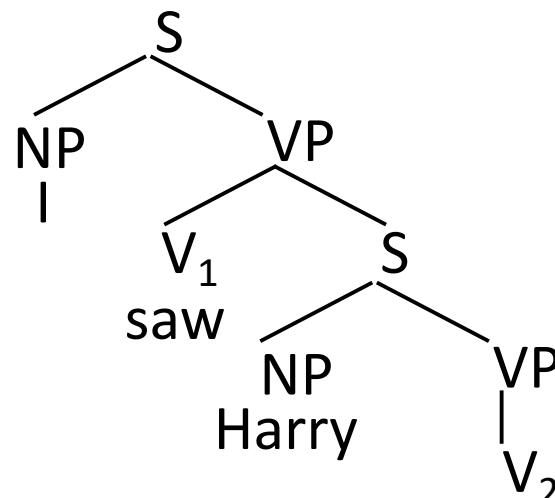
$V_2 \rightarrow swimming$

$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1\ S \Rightarrow$

$I\ saw\ S \Rightarrow I\ saw\ NP\ VP \Rightarrow$

$I\ saw\ Harry\ VP \Rightarrow$

**I saw Harry V<sub>2</sub>**



una grammatica vera è nell'ordine delle 20/30 mila regole

# Toy Grammar

$$S \rightarrow NP\ VP$$

$$VP \rightarrow V_1\ S$$

$$VP \rightarrow V_2$$

$$NP \rightarrow I \mid John \mid Harry \mid Anna$$

$$V_1 \rightarrow saw \mid see$$

$$V_2 \rightarrow swimming$$

poi non posso andare  
più avanti, perché

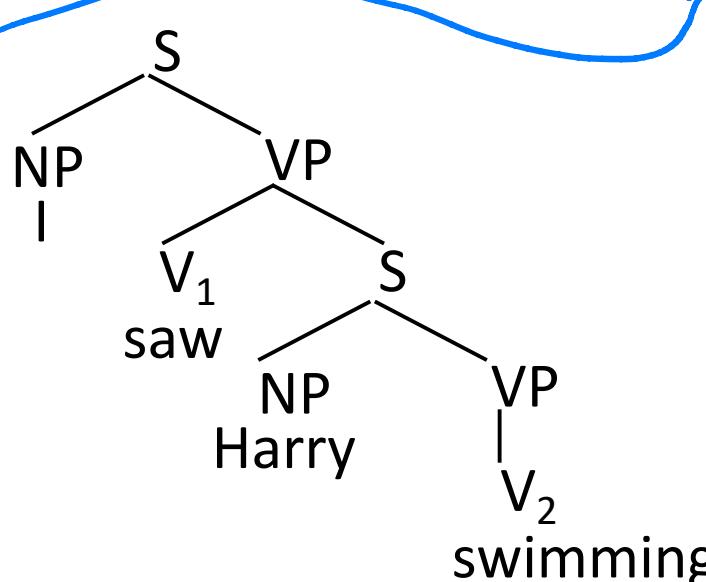
$$S \Rightarrow NP\ VP \Rightarrow I\ VP \Rightarrow I\ V_1 S \Rightarrow$$

$$I\ saw\ S \Rightarrow I\ saw\ NP\ VP \Rightarrow$$

$$I\ saw\ Harry\ VP \Rightarrow$$

$$I\ saw\ Harry\ V_2 \Rightarrow$$

**I saw Harry swimming**



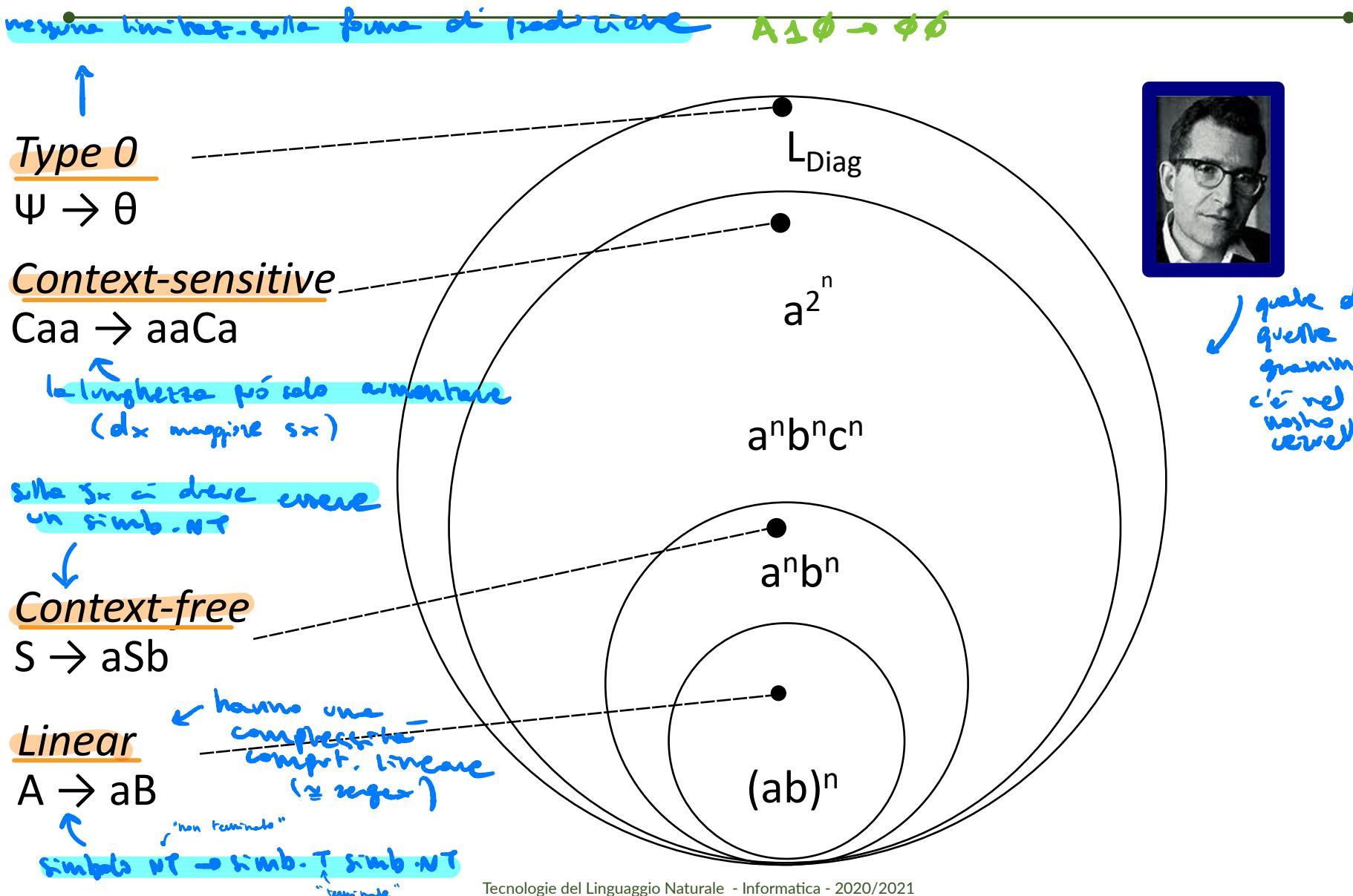
sono tutti simboli terminali

# Outline

---

- Sintassi e grammatiche generative
  - "La nostra idea è di applicare una grammatica generativa > formulare che alcune frasi nel nostro cervello sono grammaticali e altre no"
- La gerarchia di Chomsky e il linguaggio naturale
- Mildly-Context Sensitive Languages, Tree Adjoining Grammars, Combinatory Categorial Grammars

# Languages Chomsky hierarchy



# Swiss-German is NOT CF

(Shieber, 1985)

potenzialm. tutte le lingue NON sono CF



... das mer **em Hans es huus h'alfed**  
aastriiche

(ma x poco, facile  
spesso lo sono)

... that we Hans-Dat house-Acc helped paint

'... that we helped Hans paint the house'

licenziate di livello 3

... das mer **d'chind em Hans es huus l'onc**  
**h'alfe aastriiche**

rebat.  
incapaci  
che non  
possono  
generare  
da una CF

... that we the childrenAcc HansDat  
houseAcc let help paint

'... that we let the children help Hans paint  
the house'

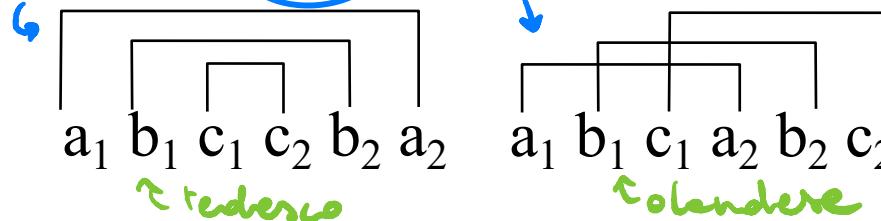
# Outline

- Sintassi e grammatiche generative
- La gerarchia di Chomsky e il linguaggio naturale
- **Mildly-Context Sensitive Languages, Tree Adjoining Grammars, Combinatory Categorial Grammars**  
~ "leggermente"

# Mildly context sensitive languages

[Joshi'85]

- Include CFG languages
- Nested and cross-serial dependencies  
☞ le grammatiche delle ling. naturali sono o nested o cross-serial



- Polynomially parsable
- Constant growth property, lineare

i ling. naturali possono essere definiti da 4 proprietà

# Constant growth property

---

- **Definition** A language  $L$  is constant growth if there is a constants  $c_0$  and a finite set of constant  $C$  such that for all  $w \in L$  where  $|w| > c_0$  there is a  $w' \in L$  such that  $|w| = |w'| + c$  for some  $c \in C$ .
- This property is the formal version of the linguistic intuition that the sentence belonging to a natural language can be built from a finite set of bounded structures using the same linear operations [Wei88].

# Languages Chomsky hierarchy

Type 0

$\Psi \rightarrow \theta$

Context-sensitive

$Caa \rightarrow aaCa$

Context-free

$S \rightarrow aSb$

Linear

$A \rightarrow aB$

$L_{\text{Diag}}$

$a^{2^n}$

$a^n b^n c^n$

$a^n b^n$

$(ab)^n$

# Languages Chomsky hierarchy

Type 0

$$\Psi \rightarrow \theta$$

Context-sensitive

$$Caa \rightarrow aaCa$$

Mildly

Context-sensitive

$$CB \rightarrow f(C,B)$$

Context-free

$$S \rightarrow aSb$$

Linear

$$A \rightarrow aB$$

$$L_{\text{Diag}}$$

$$a^{2^n}$$

$$a^n b^n c^n$$

$$a^n b^n$$

$$(ab)^n$$

in realtà è molto + vicine alle CF

↓  
posso usare gli stessi algoritmi dei lingu. di prog

# MCL $\Leftrightarrow$ TAG, HG, LIG, CCG

le nidi context-sensitivi, possono essere  
scritte attraverso grammatiche +  
semplici /  
diverse proposte

- Tree Adjoining Grammars (Joshi et al. 1975)



- Head Grammars (Pollard 1984)

- Linear Indexed Grammars (Gazdar 1985)

- Combinatory Categorial Grammars (Steedman 1985, Satta 2010)



! idea: ho sempre una grammatica generativa, ma invece che usare sfinghe, uso strutture + complete e regole di scrittura + complete

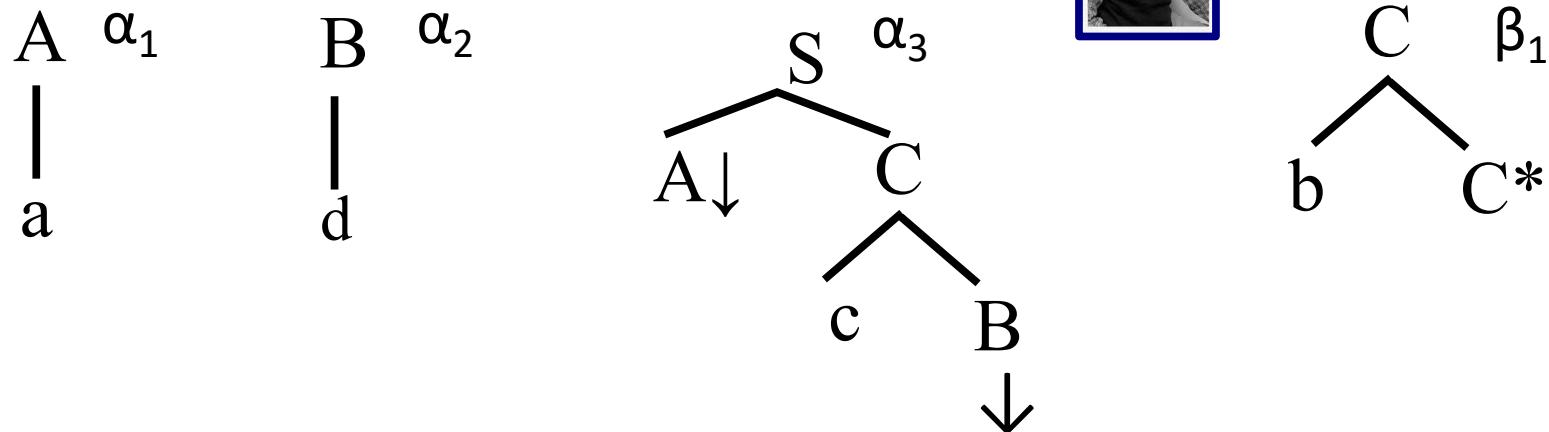
-> elementary structures



-> combination rules

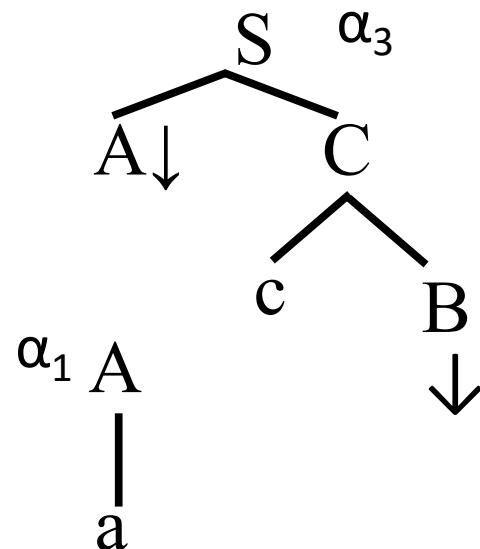
# Tree Adjoining Grammars

R  
idea: invece di  
vere stringhe,  
uso alberi



Elementary structures = multilevel trees

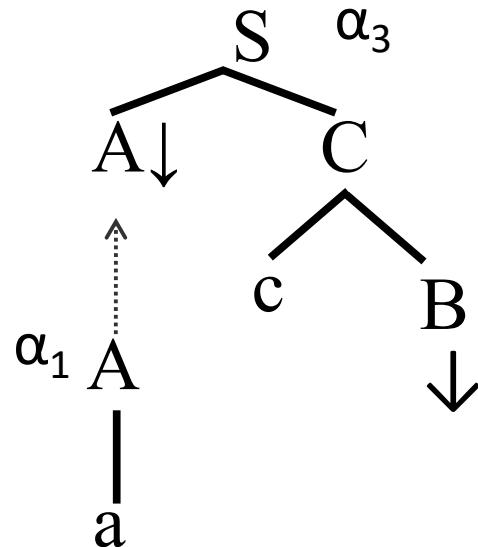
# Tree Adjoining Grammars



idea: invece che 1 singola generativa, 10 regole + complete

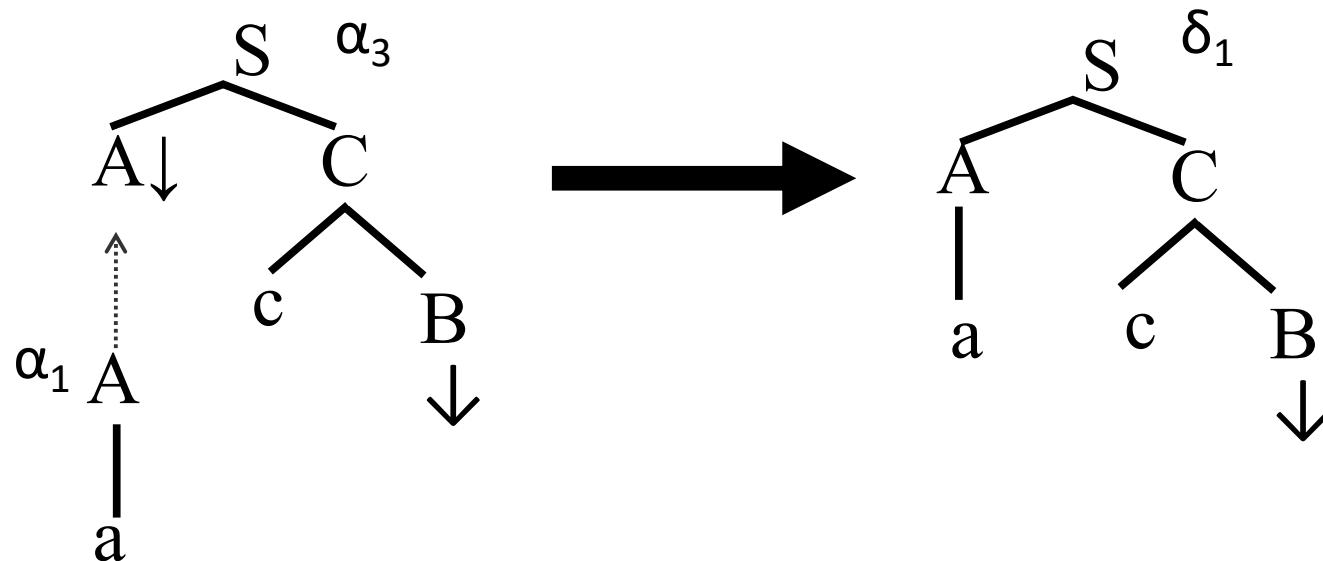
TAG operations: 1) substitution

# Tree Adjoining Grammars



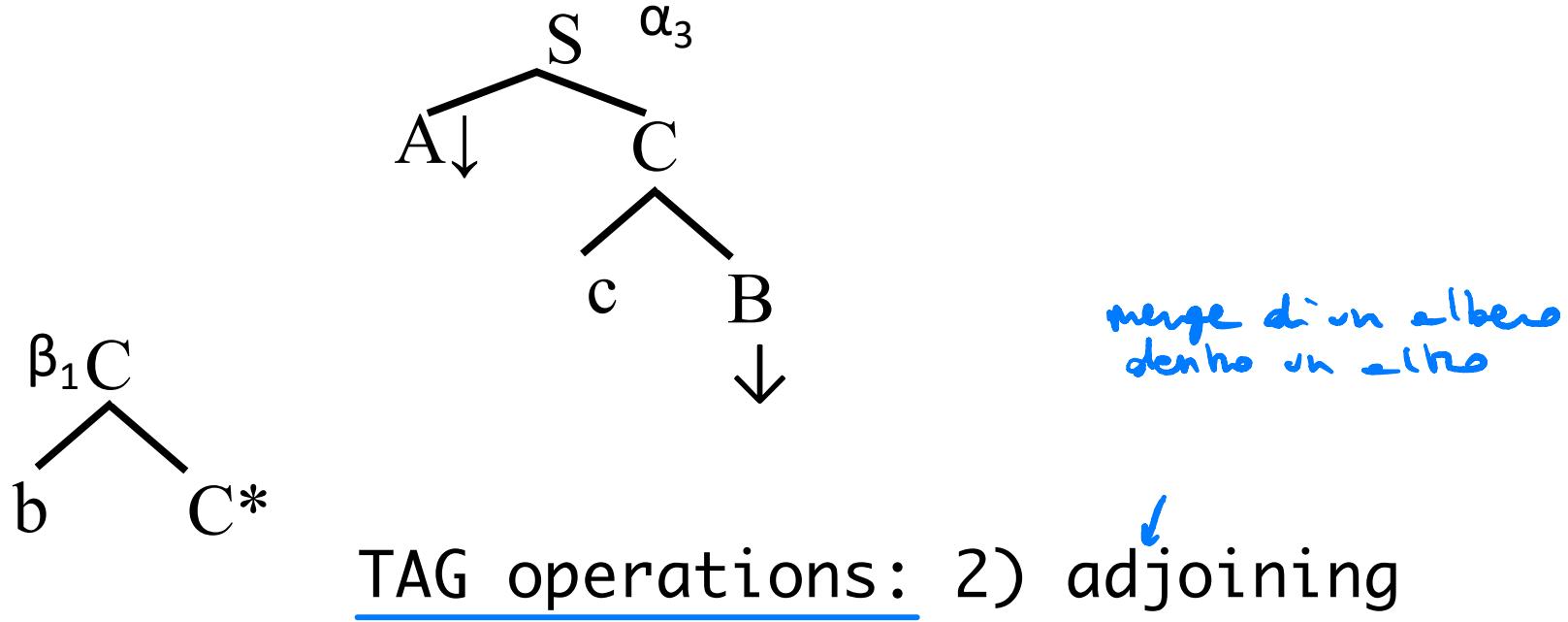
TAG operations: 1) substitution

# Tree Adjoining Grammars

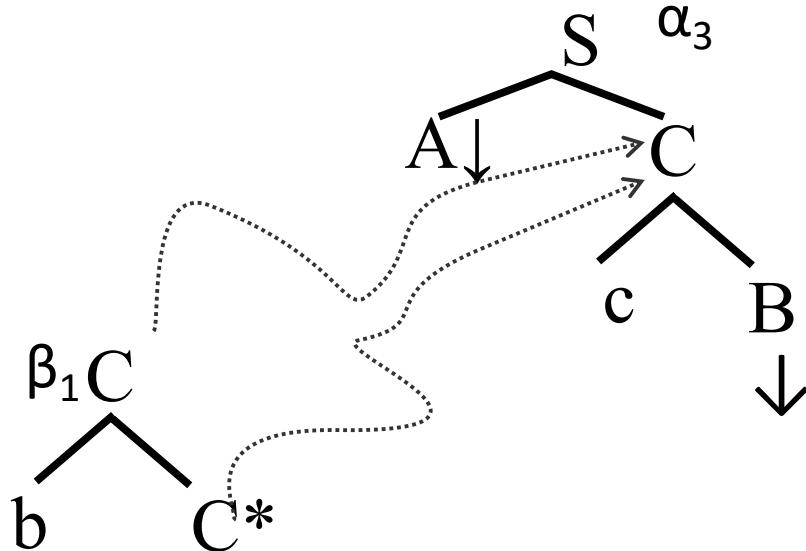


TAG operations: 1) substitution

# Tree Adjoining Grammars

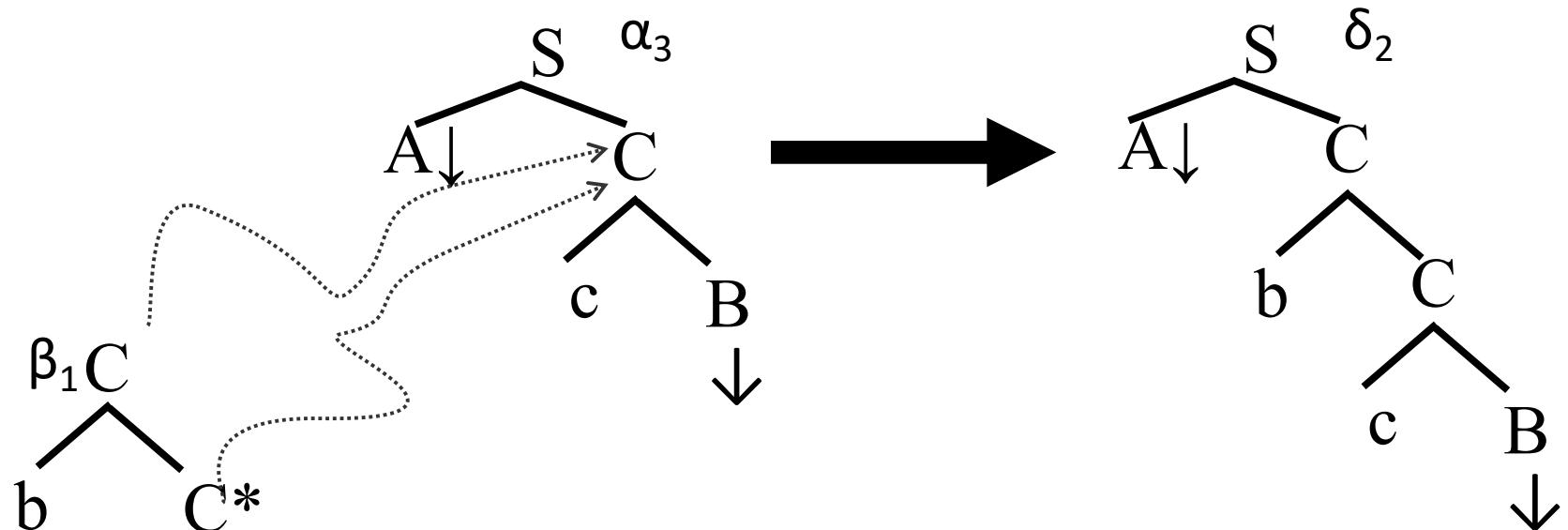


# Tree Adjoining Grammars



TAG operations: 2) adjoining

# Tree Adjoining Grammars



TAG operations: 2) adjoining



struttura  
dati

regola di  
sintassi

CF standard

stringa

Sostituzione

Tree Adjoining

alberi

sostituzione  
+  
joining

perché fatto?: per riuscire a modellare veri fenomeni linguistici  
consente  
in maniera molto + semplice

# TAG and MCSL

- TAG properly contains all context-free languages (finitely ambiguous). Theorem (Schabes 1990)
- TAG is polynomially parsable}:  $O(n^6)$ 
  - Embedded Push Down Automata, CKY (Vijay-Shanker 1987)
  - Left-to-right parser (Schabes 1990)

# TAG and MCSL

---

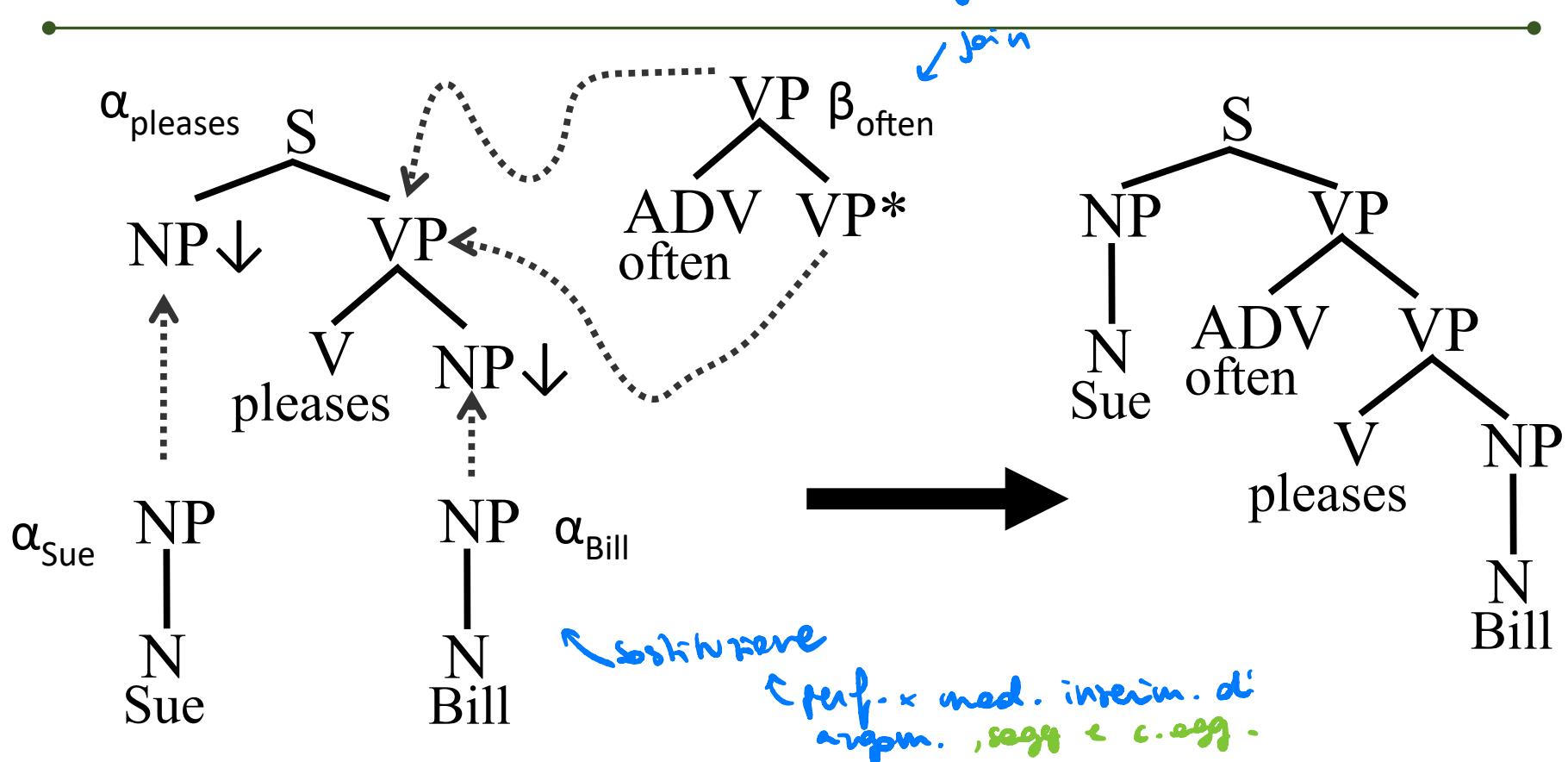
- TAG captures only certain types of dependencies
  - Cross-serial dependencies: verb-raised analysis (Kroch  
Santorini 1991)
  - No mix-languages
- TAG has the constant--growth property: (Weir  
1988)

# Lexicalized Tree Adjoining Grammars

---

- Extended domain of locality
- Recursion Factorization by adjoining operation
- Lexicalization

# LTAG



structures = multilevel trees

modellare proprietà verbali transittive

operations = substitution, adjoining

*title al solo*

*sono anch'esse midly blable*

# Combinatory Categorial Grammar

- Generative -> top-down
- Categoriali ->bottom-up

The Syntactic Process, *Mark Steedman*

Categories, Lexicon, Rules

# Combinatory Categorial Grammar

---

- Generative -> top-down
- Categorial -> bottom-up

*el g shwthra dati*

Category

- >  $A \subseteq C$ , where A is a given set of atomic elements
- >  $(X/Y), (X\backslash Y) \in C$ , if  $X, Y \in C$

- Lexicon
  - Paolo : NP
  - Francesca : NP
  - amare :  
 $(S \setminus NP) / NP$
- Rules
  - $X/Y \ Y \Rightarrow X$   
->
  - $Y \ X \setminus Y \Rightarrow X$   
-<
  - ...

Q

invece delle reg. di scrittura, ho  
regole di combinazione

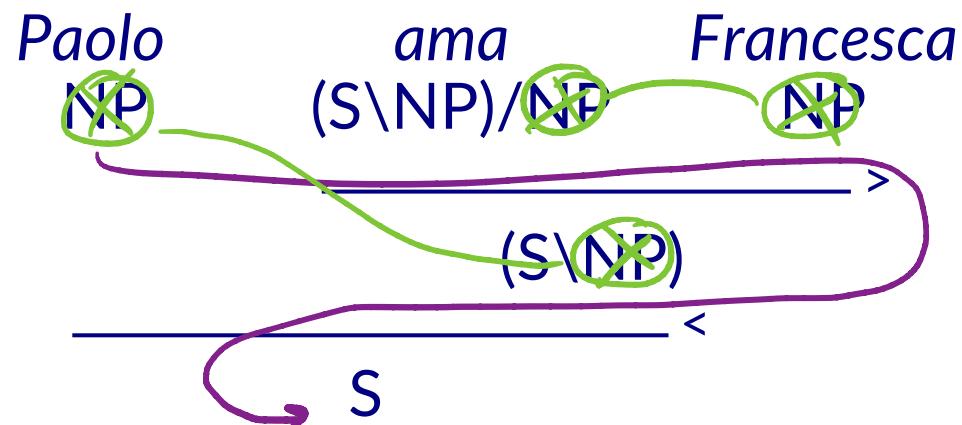
# CCG

- Lexicon

- Paolo : NP
- Francesca : NP
- amare :  
 $(S \setminus NP) / NP$

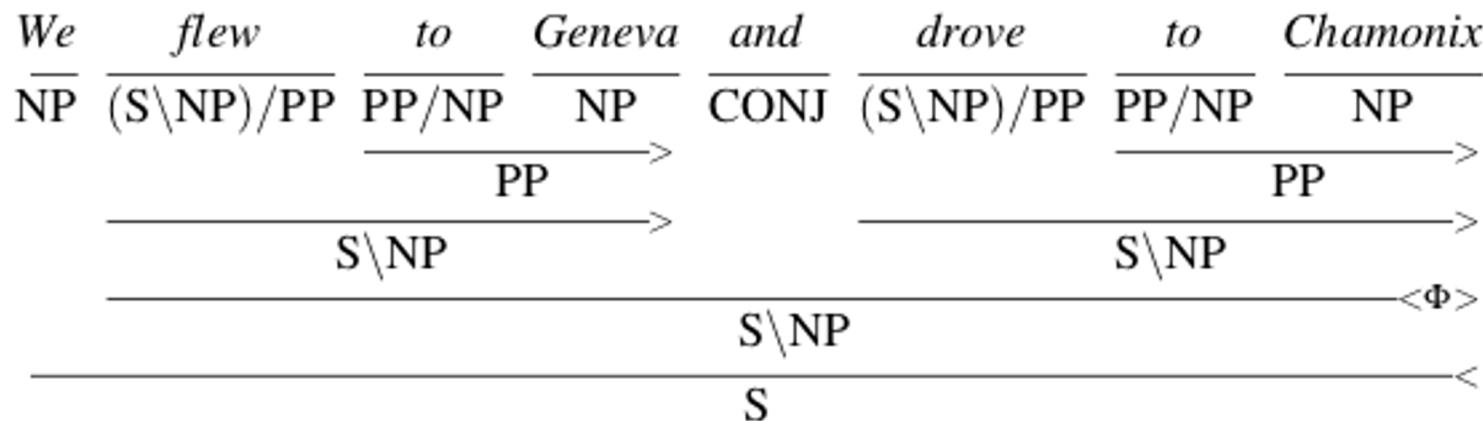
- Rules

- $X/Y \quad Y \Rightarrow X$   
->
- $Y \quad X \setminus Y \Rightarrow X$   
-<
- ...



**CCG** → sono utili in certi particolari

- Coordinazione:  $X \text{ CONJ } X \Rightarrow X$



R ↓ nei saggi sono questo

ho 2 strade:

1. il NL  $\approx$  CF (xe la diff. è poco) ↗
2. usare una TAF o una CCT + altre + modellare una midly bubble  
(e non fare l'approssimazione  $NL \approx CF$ )