



## **Corso di Laurea Magistrale in Ingegneria Informatica A.A. 2013-2014**

# **Linguaggi Formali e Compilatori**

*Grammatiche context-free*

**Giacomo PISCITELLI**

## Grammatiche libere da contesto (tipo 2)

Una **grammatica G libera da contesto** (non contestuale o di tipo 2) è una quadrupla  $(\Sigma, V, P, S)$  tale che:

- ✓  $\Sigma$  (alfabeto terminale) è l'insieme finito dei simboli elementari dell'alfabeto
- ✓  $V$  (alfabeto non terminale) è un insieme finito di metasimboli non terminali, alcune volte dette categorie sintattiche
- ✓  $P$  è un insieme di regole o produzioni sintattiche della forma  $A \rightarrow \alpha$
- ✓  $S \in V$  è una categoria sintattica detta assioma o simbolo distintivo della grammatica

Le regole sono coppie

$$A \rightarrow \alpha$$

dove  $A \in V$  è detto testa o **costrutto della regola sintattica** e  $\alpha$  è una stringa di simboli terminali e/o non terminali ( $\alpha \in (V \cup \Sigma)^*$ ) detta **corpo della regola sintattica**.

Se  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$  sono le produzioni di  $G$  aventi la stessa parte sinistra  $A$ , esse possono essere raggruppate usando la notazione:

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n \quad \text{oppure} \quad A \rightarrow \alpha_1 \cup \alpha_2 \cup \dots \cup \alpha_n$$

che abbrevia  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$

## Grammatiche libere da contesto (tipo 2)

Il compito principale di una grammatica è quello di generare stringhe di simboli terminali, cioè, nel caso dei LdP, frasi del linguaggio generato dalla grammatica.

Il processo di generazione delle stringhe può avvenire tramite **derivazioni** o, equivalentemente, tramite **costruzione di alberi di derivazione**. Infatti l'insieme delle stringhe definito con le derivazioni è lo stesso di quello definito con l'albero di derivazione.

Prima di esaminare in dettaglio la modalità di generazione delle stringhe tramite costruzione di alberi di derivazione, esaminiamo dapprima il processo di derivazione (o produzione) di stringhe tramite una grammatica *context free* e le caratteristiche formali del linguaggio generato da una grammatica e, in particolare, da una grammatica libera da contesto.

# Derivazioni

## Derivazioni

Si dice che la stringa  $\beta$  **diviene** la stringa  $\gamma$  (o che dalla parola  $\beta$  si può derivare in un passo la parola  $\gamma$ ) o che la stringa  $\beta$  **produce** la stringa  $\gamma$  per una grammatica  $G$  e si scrive

$$\beta \rightarrow \gamma$$

se esiste una produzione  $(A \rightarrow \alpha) \in P$  tale che  $\beta = \delta A \eta$  e  $\gamma = \delta \alpha \eta$   
con  $\beta, \gamma \in (V \cup \Sigma)^*$

$\beta \rightarrow^n \gamma$  ( $\beta$  **diviene**  $\gamma$  **in n passi**) se esiste una catena finita di stringhe (o parole)  $x_1, x_2, \dots, x_n$  tale che  $\beta = x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_{n-1} \rightarrow x_n = \gamma$

$\beta \rightarrow^* \gamma$  ( $\beta$  **diviene riflessivamente e transitivamente**  $\gamma$ ) se  $\beta \rightarrow^n \gamma$  per qualche  $n \geq 0$

$\beta \rightarrow^+ \gamma$  ( $\beta$  **diviene transitivamente**  $\gamma$ ) se  $\beta \rightarrow^n \gamma$  per qualche  $n \geq 1$

## ESEMPIO

$(\{S, T\}, \{a, b\}, \{S \rightarrow aTb, T \rightarrow bSa|ab\}, S)$

✓  $aTb \rightarrow abSab$

✓  $S \rightarrow aTb \rightarrow abSab$

✓  $S \rightarrow^2 abSab$

# Linguaggio generato da una grammatica

Il **linguaggio generato da una grammatica partendo dal non terminale  $A$**  (notazione  $L_A(G)$ ) è l'insieme delle stringhe di terminali  $x \in \Sigma^*$  prodotte, a partire da  $A$ , attraverso  $n \geq 1$  passi di derivazione

$$L_A(G) \equiv \{x \in \Sigma^* \mid A \rightarrow^+ x\}$$

Il **linguaggio generato da una grammatica** (notazione  $L(G)$ ) è l'insieme delle stringhe di terminali  $x \in \Sigma^*$  prodotte, a partire dall'assioma  $S$ , attraverso  $n \geq 1$  passi di derivazione

$$L(G) \equiv \{x \in \Sigma^* \mid S \rightarrow^+ x\}$$

# Linguaggio generato da una grammatica

## ESEMPIO

$G \equiv (\{S, T\}, \{a, b\}, \{S \rightarrow aTb, T \rightarrow bSa|ab\}, S)$

✓  $T \rightarrow bSa \rightarrow baTba \rightarrow baabba$

✓  $T \rightarrow bSa \rightarrow baTba \rightarrow babSaba \rightarrow babaTbaba \rightarrow babaabbaba$

✓  $S \rightarrow aTb \rightarrow aabb$

✓  $S \rightarrow aTb \rightarrow abSab \rightarrow abaTbab \rightarrow abaabbab$

$L_T(G) \equiv \{baabba, babaabbaba, \dots\}$

$L(G) \equiv \{aabb, abaabbab, \dots\}$

## ALTRO ESEMPIO

$\Sigma = \{ (, ), id, +, *, /, - \}$

$V = \{E\}$

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow E * E$

$E \rightarrow E / E$

$E \rightarrow - E$

$E \rightarrow ( E )$

$E \rightarrow id$

# Linguaggio libero da contesto

Un linguaggio è **libero dal contesto** se esiste una grammatica libera dal contesto che lo genera.

Esempio di linguaggio context free: quello generato dalla grammatica

$$G \equiv (\{S, T\}, \{a, b\}, \{S \rightarrow aTb, T \rightarrow bSa, T \rightarrow ab\}, S)$$

Due grammatiche  $G$  e  $G'$  sono **debolmente equivalenti** se generano lo stesso linguaggio, cioè  $L(G) \equiv L(G')$

## ESEMPIO

$$G \equiv (\{S\}, \{a, b\}, \{S \rightarrow aSb|ab\}, S)$$

$$G' \equiv (\{S, A, B\}, \{a, b\}, \{S \rightarrow ASB|AB, A \rightarrow a, B \rightarrow b\}, S)$$

$$L(G) \equiv L(G') \equiv \{ab, aabb, aaabbb, \dots\} \equiv \{a^n b^n \mid n \geq 1\}$$

## Ricorsione e linguaggi finiti/infiniti

Una derivazione è **ricorsiva** se, partendo da un non terminale, produce in almeno un passo una stringa che contiene lo stesso non terminale, cioè se è della forma

$$A \rightarrow^+ \alpha A \beta$$

E' **ricorsiva sinistra** se  $\alpha = \epsilon$ ; **ricorsiva destra** se  $\beta = \epsilon$ .

Il non terminale  $A$  è detto **ricorsivo**.

$$G \equiv (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid ab\}, S)$$

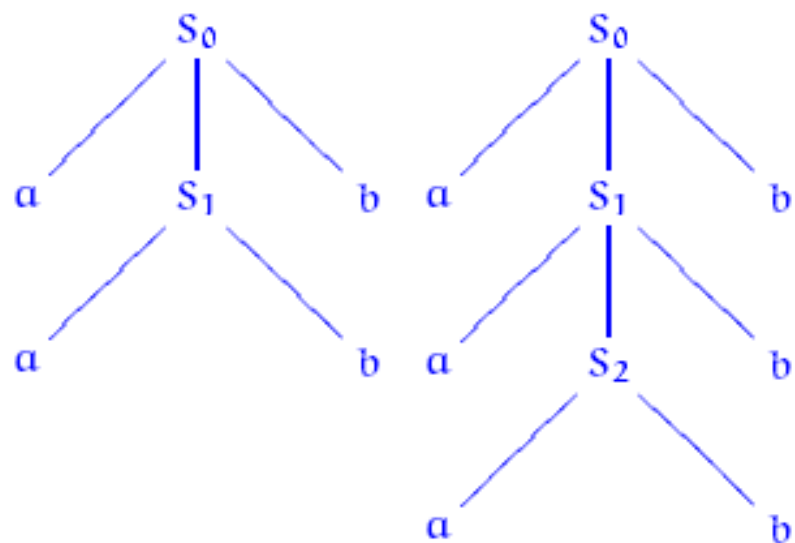
$$S \rightarrow aSb$$



# Produzioni, derivazioni e alberi sintattici

Possiamo rappresentare produzioni e derivazioni mediante alberi:

$$G \equiv (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid ab\}, S)$$



# Albero sintattico o Parse tree

Un **parse tree** o albero sintattico di una grammatica descrive graficamente come l'assioma **S** della grammatica deriva una stringa nel linguaggio da essa sotteso.

Formalmente, data una grammatica non contestuale, il suo albero sintattico avrà le seguenti caratteristiche:

- ⇒ la radice dell'albero avrà come etichetta l'**assioma**;
- ⇒ ogni nodo non terminale dell'albero corrisponderà ad **un elemento dell'alfabeto  $V$** ;
- ⇒ ogni nodo terminale (foglia) dell'albero corrisponderà ad **un elemento dell'alfabeto  $\Sigma$** ;
- ⇒ da ogni nodo non terminale dell'albero **derivano nodi terminali e/o non terminali**;
- ⇒ come caso particolare, da un nodo non terminale può derivare **un nodo terminale costituito dal simbolo vuoto  $\epsilon$** .

# Ambiguità

Una grammatica è detta **ambigua** se esistono due diverse sequenze di produzioni che generano la stessa parola:

$\text{espr} \rightarrow \text{espr} * \text{espr} \rightarrow \text{espr} * \text{espr} + \text{espr} \rightarrow 5 * 3 + 1$

$\text{espr} \rightarrow \text{espr} + \text{espr} \rightarrow \text{espr} * \text{espr} + \text{espr} \rightarrow 5 * 3 + 1$

La semantica della prima espressione vuol essere 20, quella della seconda 16.

Per risolvere il problema si definisce una precedenza tra le produzioni. Ad esempio, possiamo stabilire che una produzione

$\text{espr} \rightarrow \text{espr} + \text{espr}$

debba sempre precedere, e mai seguire, una produzione

$\text{espr} \rightarrow \text{espr} * \text{espr}$

**Lo stesso problema si presenta se usiamo due volte lo stesso operatore**, e l'operatore non è associativo (per esempio,  $<$ ).

In tal caso viene scelta una direzione per l'associatività (da sinistra a destra o viceversa).

## Ambiguità nelle istruzioni

L'ambiguità non è un problema relegato alle espressioni.

```
if ( x == 0 )  
if ( y == 0 ) printf("pippo")  
else printf("pluto");
```

Ci sono due modi di produrre il frammento precedente:

```
istr → if (espr) istr else istr → if (espr) if (espr) istr else istr  
istr → if (espr) istr → if (espr) if (espr) istr else istr
```

Nel primo caso l'**else** è associato al primo **if**: se **x** non è zero, verrà stampato **pluto**, altrimenti se **y** è zero verrà stampato **pippo**.

Nel secondo caso l'**else** è associato al secondo **if**: se **x** non è zero, viene esaminato **y**: se quest'ultimo è zero verrà stampato **pippo**, altrimenti verrà stampato **pluto**.

Le due semantiche sono chiaramente diverse (la prima stampa **pluto** in ogni caso, una volta che **x** non è zero).

Anche qui siamo di fronte a un problema di associatività, che assumiamo verso destra: l'ultimo **else** è associato all'ultimo **if**.

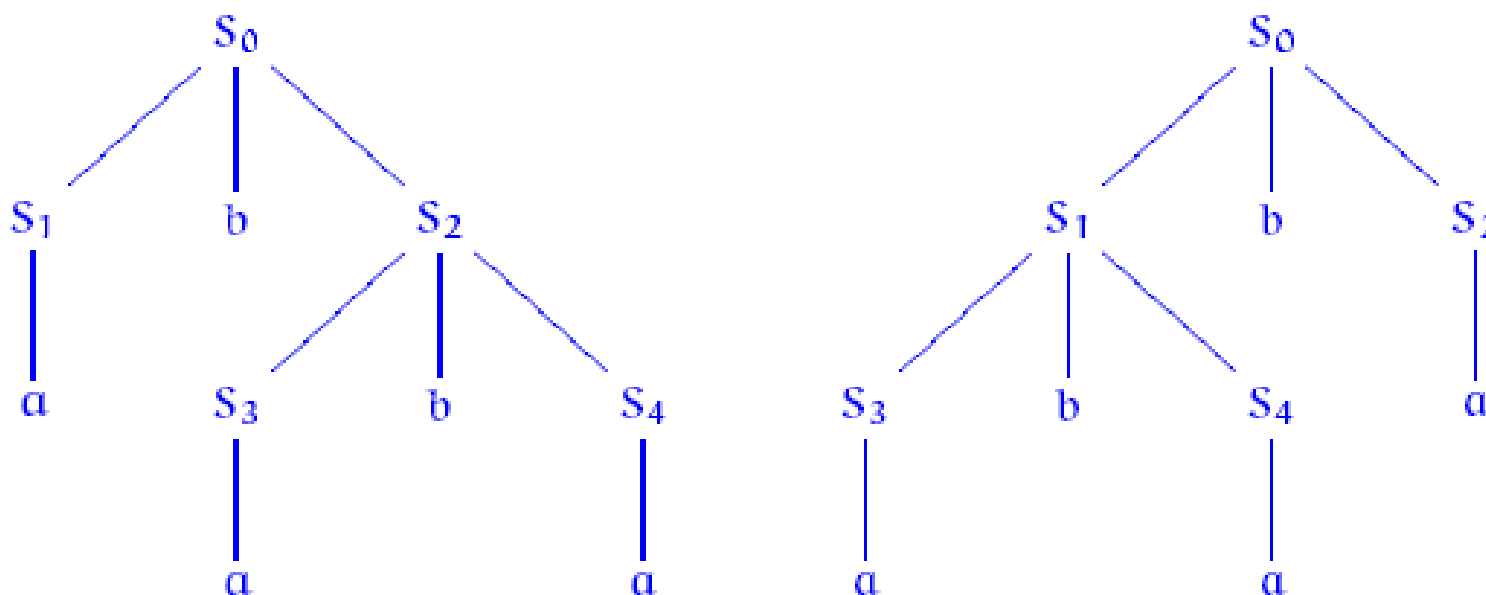
# Precedenza degli operatori

Operatori					Associatività	
( )	[ ]	.	-	>	da sinistra a destra	
++	--	!	sizeof	(tipo)	da destra a sinistra	
+	-	(unari)				
&	*	(puntatori)				
*	/	%			da sinistra a destra	
+	-				da sinistra a destra	
<	>	<=	>=		da sinistra a destra	
==	!=				da sinistra a destra	
&&					da sinistra a destra	
					da sinistra a destra	
?:					da destra a sinistra	
=	+=	-=	*=	/=	%=	da destra a sinistra

# Alberi sintattici e ambiguità

Una **frase** è **ambigua** se essa è generata dalla grammatica con due alberi sintattici distinti. In tal caso anche la **grammatica** è detta **ambigua**.

$$G = (\{S\}, \{a, b\}, \{S \rightarrow SbS|a\}, S)$$



Il grado di ambiguità di una frase è il numero dei suoi alberi sintattici distinti.  
Non è decidibile se una grammatica è ambigua.

# Ambiguità

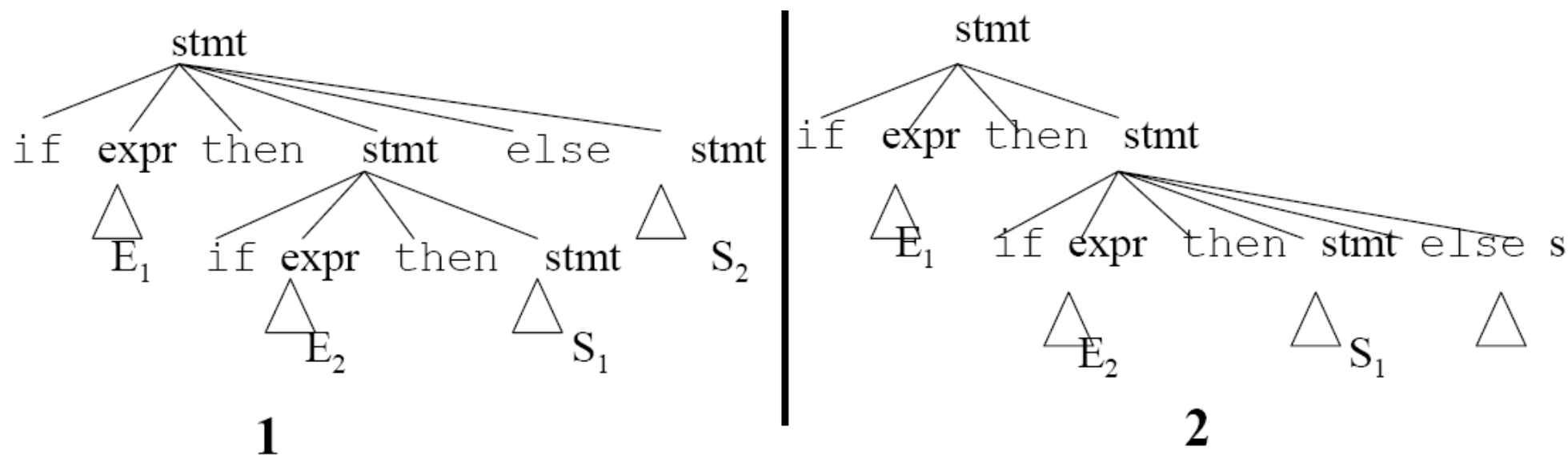
Per poter costruire un parser, la grammatica non deve essere ambigua.

Le ambiguità nella grammatica devono essere eliminate durante il progetto del compilatore.

Grammatica non ambigua → unica selezione nell'albero sintattico per una frase.

`stmt` → `if expr then stmt` | `if expr then stmt else stmt` | `otherstmts`

`if E1 then if E2 then S1 else S2`



# Ambiguità

- Noi assumiamo il secondo albero sintattico (**else** corrisponde all'**if** più vicino).
- Di conseguenza dobbiamo eliminare l'ambiguità con tale obiettivo.
- La grammatica non-ambigua sarà:

`stmt → matchedstmt | unmatchedstmt`

`matchedstmt → if expr then matchedstmt else matchedstmt | otherstmts`

`unmatchedstmt → if expr then stmt |  
if expr then matchedstmt else unmatchedstmt`



## Ambiguità – Precedenza degli operatori

Grammatiche ambigue possono essere rese non-ambigue in accordo con le precedenze degli operatori e con le regole di associatività degli operatori.

L'**associatività** esprime a quali operandi si applica un operatore: nella stragrande maggioranza dei linguaggi di programmazione, i 4 operatori aritmetici (+, -, \*, /) associano **verso sinistra**, nel senso che un operando con lo stesso operatore aritmetico alla sua sinistra e alla sua destra, è associato con l'operatore a sinistra. L'operatore di esponenziazione (^) associa **verso destra**.

In presenza di operatori diversi, la **precedenza** regola l'ordine di applicazione degli operatori.

precedenze:

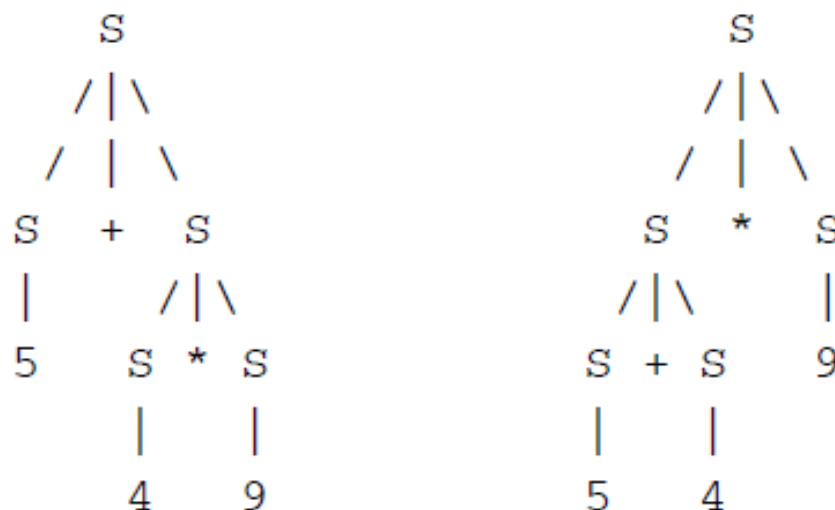
- ^ (right to left)
- \* / (left to right)
- + - (left to right)

Esempio: grammatica G per espressioni aritmetiche:

$$S \rightarrow S + S \mid S - S \mid S * S \mid S / S \mid ( 0 \mid 1 \mid \dots \mid 9 )^+$$

la stringa  $s = 5 + 4 * 9$  ammette due alberi sintattici distinti.

## Ambiguità – Precedenza degli operatori



Una grammatica ambigua NON è adatta all'analisi sintattica (e traduzione automatica):

- la struttura sintattica di una frase non è unica
- l'analizzatore sintattico dovrebbe restituire tutti gli alberi sintattici corrispondenti alla stringa (oppure uno a caso)
- ambiguità sintattica → ambiguità semantica

In presenza di grammatiche ambigue non è possibile eseguire l'analisi sintattica (e la traduzione) in modo deterministico: di conseguenza, tali grammatiche sono poco utilizzate nelle applicazioni informatiche.

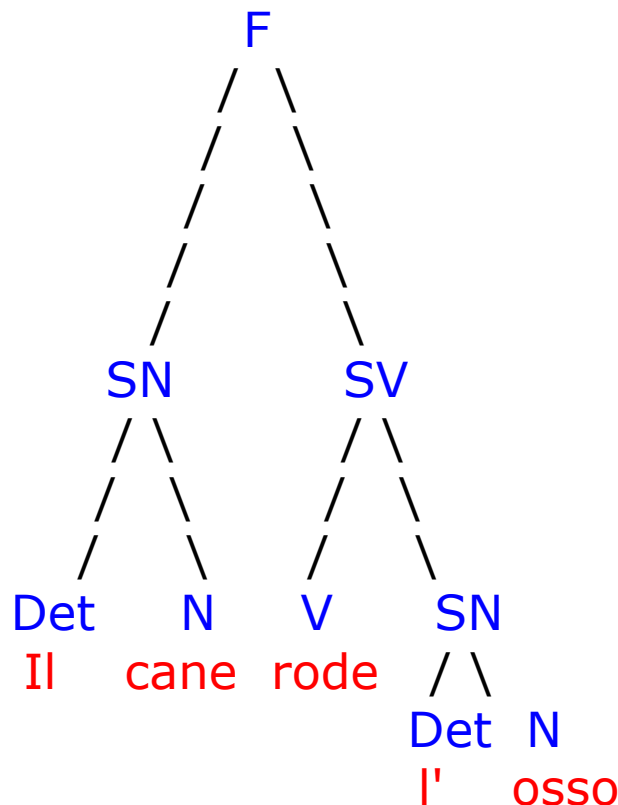
# Produzioni, derivazioni e alberi sintattici nei linguaggi naturali

Il lavoro dei linguisti nell'ambito della grammatica generativa vede spesso tale albero di derivazioni come un oggetto di studio fondamentale.

Secondo questo punto di vista, una frase non è semplicemente una stringa di parole, ma piuttosto un albero con rami subordinati e sopraordinati connessi a dei nodi.

Sostanzialmente, il modello ad albero funziona in qualche modo come il seguente esempio (si è considerata una frase semplice), nel quale **F** è una frase (*sentence*), **Det** è un articolo determinativo (*determiner*), **N** è un sostantivo (*noun*), **V** è un verbo (*verb*), **SN** è un sintagma nominale (*noun phrase*) e **SV** è sintagma verbale (*verb phrase*).

# Produzioni, derivazioni e alberi sintattici nei linguaggi naturali



Questo diagramma ad albero è anche chiamato un **indicatore sintagmatico** (*phrase marker*).

Esso può essere rappresentato più convenientemente in forma testuale, sebbene il risultato sia meno facilmente leggibile.

# Produzioni, derivazioni e alberi sintattici nei linguaggi naturali

In forma testuale la suddetta frase sarebbe resa nel modo seguente:

[F [SN [Det Il ] [N cane ] ] [SV [v rode ] [SN [Det l' ] [N osso ] ] ] ]

Comunque Chomsky ha affermato che anche le grammatiche a struttura sintagmatica (*phrase structure grammars*) sono inadeguate per descrivere i linguaggi naturali. Per far fronte a questa necessità, egli formulò allora il sistema più complesso della **grammatica trasformativa**.

## Forma normale di Chomsky

Data una grammatica, si possono costruire delle **grammatiche equivalenti** che abbiano particolari forme delle produzioni, dette **forme normali**.

Una grammatica è in **forma normale di Chomsky** se le sue regole sono caratterizzate dalle seguenti proprietà:

1. **le regole sono omogenee binarie**: cioè esse sono della forma

$$A \rightarrow BC \quad \text{con } B, C \in V$$

la stringa è perciò formata da due non terminali;

2. **le regole sono terminali unitarie**: cioè esse sono della forma

$$A \rightarrow a \quad \text{con } a \in \Sigma$$

la stringa è formata da un terminale;

3. la regola  $S \rightarrow \epsilon$  si applica solo se il linguaggio contiene  $\epsilon$ .

La forma normale di Chomsky è caratterizzata da un albero sintattico binario.

Data una grammatica qualsiasi, si può dimostrare che esiste un semplice algoritmo per costruire una grammatica equivalente in forma normale di Chomsky.

## Parser top-down e bottom-up

L'analizzatore sintattico (parser) può lavorare, come vedremo, in una varietà di modi.

I più comuni sono il metodo **top-down** (nel quale l'ordine con cui il parse tree viene esplorato procede dalla radice dell'albero sintattico fino alle foglie) e il metodo **bottom-up** (nel quale l'ordine procede dalle foglie alla radice) – ma in quasi tutti esso processa l'input da sinistra a destra, un simbolo alla volta, applicando una *left-most derivation*.

Si dice, in sostanza, che le grammatiche si prestano per lo più al *L-parsing*, cioè **le frasi possono essere analizzate da sinistra a destra**.

# Derivazioni Left-Most e Right-Most

Si consideri la seguente grammatica:

$$E \Rightarrow -E \mid (E) \mid E+E \mid E * E \mid \text{id}$$

avente come simboli terminali  $(, ), +, -, *, \text{id}$  e come simbolo non terminale  $E$ , che è anche il simbolo iniziale.

## Left-Most Derivation

Una “left-most derivation” è una derivazione nella quale durante ogni passo solo il non-terminale più a sinistra è sostituito.

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\text{id}+E) \Rightarrow -(\text{id}+\text{id})$$

## Right-Most Derivation

Una “right-most derivation” è una derivazione nella quale durante ogni passo solo il non-terminale più a destra è sostituito.

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+\text{id}) \Rightarrow -(\text{id}+\text{id})$$



## Ricorsione sinistra (Left Recursion)

Una grammatica è "**left recursive**" se ha un non terminale  $A$  tale che

$$A \Rightarrow A\alpha \text{ per qualche stringa } \alpha$$

Le tecniche di parsing Top-down non possono gestire grammatiche left-recursive, perché ne potrebbe derivare un *loop* infinito.

Una grammatica left-recursive deve essere convertita in una equivalente non left-recursive.

La ricorsione sinistra può comparire in un singolo passo della derivazione (*immediate left-recursion*), o può comparire in più che un passo.

Eliminazione di una *Immediate Left-Recursion*

$$A \rightarrow A\alpha \mid \beta \quad \text{dove } \beta \neq A\gamma \quad \text{e} \quad \gamma \in (V \cup T)^*$$

$\Downarrow$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \varepsilon \quad \text{grammatica equivalente}$$

La ricorsione sinistra di una generica grammatica libera da contesto può essere sistematicamente eliminata attraverso l'algoritmo 4.8 riportato nel testo consigliato.

# Sintassi

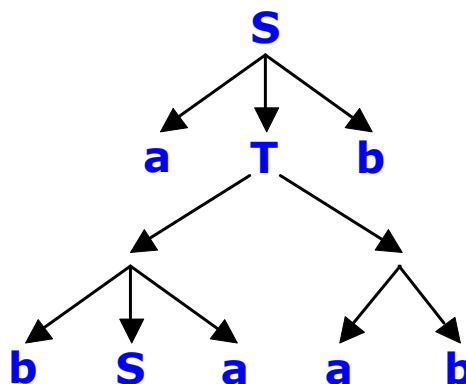
## In conclusione

La sintassi è costituita da un insieme di regole che definiscono le frasi formalmente corrette e allo stesso tempo permettono di assegnare ad esse una struttura che rappresenta graficamente il processo di derivazione (albero sintattico).

Considerata la seguente grammatica:

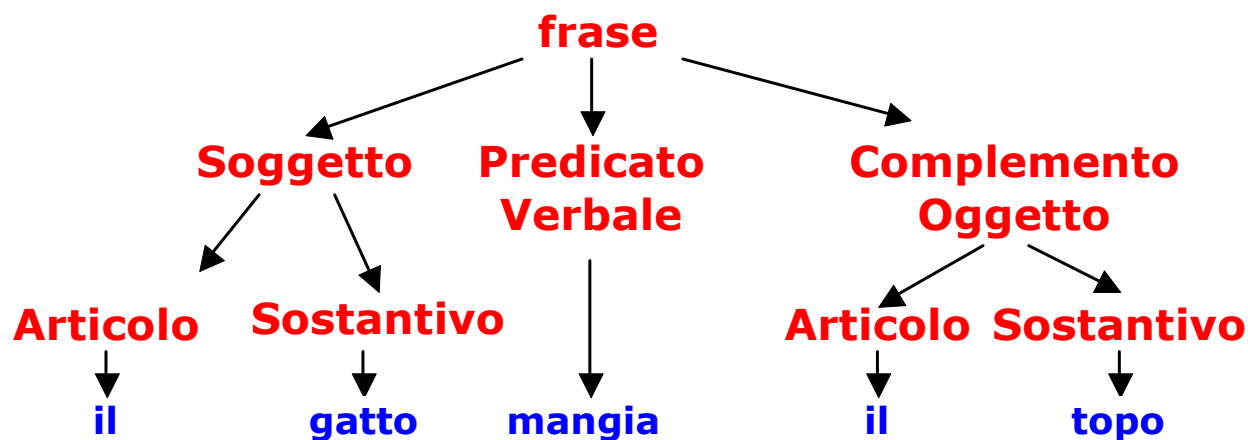
$$G \equiv (\{S, T\}, \{a, b\}, \{S \rightarrow aTb, T \rightarrow bSa, T \rightarrow ab\}, S)$$

l'insieme delle sue regole può essere rappresentato graficamente tramite il seguente albero:



# Sintassi context-free

Se, analogamente, si vuole rappresentare graficamente la struttura di una **frase** del linguaggio italiano, che ha come costituenti un **Soggetto**, un **Predicato Verbale** e un **Complemento Oggetto**, ecc. , si potrà ricorrere al seguente albero sintattico:



## Sintassi context-free

Se, ancora, si vuole rappresentare la struttura di un programma in **un linguaggio di programmazione**, si dovrà costruire, come in seguito esamineremo, un albero sintattico dal quale si evinca che un programma ha come costituenti **le parti dichiarative e quelle esecutive**.

Le parti dichiarative definiscono i dati usati dal programma. Le parti esecutive si articolano nelle istruzioni, che possono essere di vari tipi. I costituenti del livello più basso sono gli elementi lessicali già considerati, che dalla sintassi sono visti come atomi indecomponibili. Infatti la loro definizione spetta al livello lessicale.