

# *Linguaggi Formali e Traduttori*

**Parsificazione:** grammatiche LL(1), algoritmo iterativo e algoritmo discendente ricorsivo

**Traduzione:** definizioni guidate dalla sintassi, schemi di traduzione, traduzione deterministica top-down

## Analizzatori discendenti

Costruiscono un albero di parsificazione iniziando dalla radice e creando i nodi in preordine.

Equivalentemente la parsificazione top-down può essere vista come ricerca di una derivazione leftmost per la stringa in input.

## Analizzatori ascendenti

Costruiscono l'albero dalle foglie alla radice.

L'analisi ascendente può essere vista come ricerca di una derivazione rightmost della stringa in ordine contrario.

La parsificazione nella forma più generale può richiedere il backtracking per trovare la produzione corretta da applicare durante l'esame di una stringa. In molte aree dell'informatica i linguaggi sono progettati in modo da essere analizzabili deterministicamente.

Le grammatiche che permettono analisi sintattica top-down deterministica o "*parsing predittivo discendente*" sono chiamate **LL(k)**, quelle che permettono analisi sintattica bottom-up deterministica o "*parsing predittivo ascendente*" sono chiamate **LR(k)**, dove k è il numero di simboli necessari per individuare la produzione senza ambiguità.

Automa con “*look ahead*” di lunghezza 1.

1. determinare la produzione da applicare quando un nonterminale si trova in cima alla pila.
2. se in cima alla pila si trova un terminale basta verificare la corrispondenza di questo col simbolo sul nastro di input e avanzare la testina di lettura.
3. in mancanza di questa corrispondenza la stringa in input deve essere rifiutata.

Quali sono le condizioni cui deve soddisfare la grammatica per costruire un *analizzatore deterministico discendente* usando 1 simbolo in input? Cioè quando una grammatica è LL(1)?

si associa ad ogni produzione l’«**insieme guida**».

# Parsing look ahead 1: insiemi FIRST e FOLLOW

$$\text{FIRST}(\alpha) = \{a \mid \alpha \Rightarrow^* a\beta\} \cup \{\varepsilon \mid \text{se } \alpha \Rightarrow^* \varepsilon\}$$

$F(\alpha)$  soddisfa queste regole (ricorsive):

$$\begin{array}{l} 1. F(\varepsilon) = \{\varepsilon\} \\ 2. F(a\beta) = \{a\} \\ 3. F(A\beta) = \begin{cases} F(A) & \text{se } \varepsilon \notin F(A) \\ (F(A) - \{\varepsilon\}) \cup F(\beta) & \text{se } \varepsilon \in F(A) \end{cases} \end{array}$$

$$\text{dove } F(A) = \bigcup_{i=1}^K F(\gamma_i) \quad \text{Se } A \rightarrow \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_k$$

$$\text{FOLLOW}(A) = \{a \mid S \Rightarrow^* \alpha A a \beta\} \cup \{\$ \mid \text{se } S \Rightarrow^* \alpha A\}$$

$\text{FW}(A)$  soddisfa la seguente equazione:

$$\text{FW}(A) = \{F(\beta) - \{\varepsilon\} \mid B \rightarrow \alpha A \beta \in P\}$$

$$\{\text{FW}(B) \mid B \neq A \text{ e } B \rightarrow \alpha A \beta \in P \ \& \ \varepsilon \in F(\beta)\}$$

$$\{\$ \mid A \text{ è lo start symbol di } G\}$$

$$Gui(A \rightarrow \alpha) = \begin{cases} F(\alpha) & \text{se } \varepsilon \notin F(\alpha) \\ (F(\alpha) - \{\varepsilon\}) \cup FW(A) & \text{se } \varepsilon \in F(\alpha) \end{cases}$$

Una grammatica è LL(1) se per ogni non terminale A e per ogni coppia di produzioni  $A \rightarrow \alpha$  e  $A \rightarrow \beta$ , gli insiemi guida sono disgiunti:

Eliminazione delle ricorsioni sinistre immediate

$$\begin{array}{l} A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_k \\ A \rightarrow \beta_1 | \beta_2 | \dots | \beta_h \end{array} \quad \begin{array}{l} k \geq 1, \alpha_i \neq \varepsilon \\ h \geq 1 \end{array}$$



$$\begin{array}{l} A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_h A' \\ A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_k A' | \varepsilon \end{array}$$

Algoritmo che implementa l'automa a pila look ahead 1.

*Input:* 1) Stringa da parsificare, seguita da un simbolo di fine stringa: \$

2) Tabella M che memorizza, per ogni coppia  
< variabile A, simbolo in input a >  
la produzione  $A \rightarrow \alpha$  se  $a \in \text{Gui}(A \rightarrow \alpha)$   
' ' se  $a \notin \text{Gui}(A \rightarrow \alpha)$ .

*Output:* Stringa accettata o segnalazione di errore;  
Elenco delle produzioni usate per generare la stringa.

# Analizzatore iterativo

```
AUTOMA( )
    cc ← PROSS( )          /*1^ simbolo della stringa in input */
    X ← top (STACK)
    while (not_empty (STACK))
        if (X ∈ Σ and X = cc)
            cc ← PROSS( )
            pop (STACK)
        elseif (X ∈ V and M[X, cc] = X → α and α ≠ cc.β )
            pop (STACK)
            push (α, STACK)
            stampa (X → α)
        elseif (X ∈ V and M[X, cc] = X → α and α = cc.β )
            cc ← PROSS( )
            pop (STACK)
            push (β, STACK)
            stampa (X → α)
        else ERRORE (...)
        X ← top (STACK)

    if (cc = '$') stampa ("stringa accettata")
    else ERRORE (...)
```

Ad ogni variabile  $A$  con produzioni

$A \rightarrow \alpha_1$

$A \rightarrow \alpha_2$

...

$A \rightarrow \alpha_k$

si associa una funzione:

function  $A( )$

if ( $cc \in \text{Gui}(A \rightarrow \alpha_1)$ )  $\text{body}(\alpha_1)$

elseif ( $cc \in \text{Gui}(A \rightarrow \alpha_2)$ )  $\text{body}(\alpha_2)$

....

elseif ( $cc \in \text{Gui}(A \rightarrow \alpha_k)$ )  $\text{body}(\alpha_k)$

else **ERRORE** (...)

# Analizzatore a discesa ricorsiva

Se  $\alpha = \varepsilon$ ,  $\text{body}(\varepsilon) = \underline{\text{do nothing}}$

Se  $\alpha = X_1 \dots X_m$ ,  $\text{body}(X_1 \dots X_m)$  è così definito:

$\text{body}(X_1 \dots X_m) = \text{act}(X_1) \text{act}(X_2) \dots \text{act}(X_m)$

$$\text{act}(X) = \begin{cases} X() & \text{se } X \in V \\ \text{cc} \leftarrow \text{PROSS}() & \text{se } X = \text{cc} \\ \text{ERRORE}(\dots) & \text{altrimenti} \end{cases}$$

```
main discesa_ricorsiva( )
```

```
    cc ← PROSS( )
```

```
    S( )
```

```
    if (cc = '$') stampa "stringa accettata"
```

```
    else ERRORE(...)
```

# Definizioni guidate dalla sintassi: SDD

Tecniche di traduzione guidate dalle grammatiche libere.

- si associa informazione a un costrutto di un linguaggio associando “attributi” al simbolo della grammatica che rappresenta il costrutto.
- in una *definizione guidata dalla sintassi* (syntax-directed definition) i valori degli attributi sono calcolati da “regole semantiche” associate alle produzioni della grammatica.

La traduzione è il risultato della valutazione delle regole semantiche.

Le definizioni guidate dalla sintassi, quando i valori degli attributi sono definiti solo in funzione di costanti e dei valori di altri attributi, vengono anche chiamate grammatiche ad attributi.

Gli attributi per i simboli terminali hanno i valori lessicali forniti dall'analizzatore lessicale. Nelle SDD non vi sono regole semantiche per calcolare i valori degli attributi per i simboli terminali.

# Definizioni guidate dalla sintassi: SDD

Due tipi di attributi:

- **sintetizzati**: un attributo sintetizzato per una variabile  $A$  in un nodo  $n$  dell'albero di parsificazione è definito da una regola semantica associata alla produzione in  $n$  e il suo valore è calcolato solo in termini dei valori degli attributi nei nodi figli di  $n$  e in  $n$  stesso.  
( $A$  è il simbolo a sinistra nella produzione, cioè la testa).
- **ereditati**: un attributo ereditato per una variabile  $A$  in un nodo  $n$  dell'albero di parsificazione è definito da una regola semantica associata alla produzione nel nodo padre di  $n$  e il suo valore è calcolato solo in termini dei valori degli attributi del padre di  $n$ , di  $n$  stesso e dei suoi fratelli.  
( $A$  è un simbolo nel corpo della produzione, cioè al membro destro).
- *albero annotato*
- attributi sintetizzati valutati in ordine *bottom-up*
- negli SDD che hanno sia attributi sintetizzati, sia ereditati, non vi è garanzia che vi sia almeno un ordine in cui valutare gli attributi
- *grafo delle dipendenze*: Se il grafo presenta un ciclo non è possibile valutare gli attributi. L'ordine di valutazione deve rispettare un ordinamento topologico dei vertici del grafo.

# Schemi di traduzione: SDT

Gli **schemi di traduzione (SDT)** sono un'utile notazione per specificare la traduzione durante la parsificazione.

Uno schema di traduzione è una definizione guidata dalla sintassi in cui le azioni semantiche, racchiuse tra parentesi graffe, sono inserite nei membri destri delle produzioni, in posizione tale che il valore di un attributo sia disponibile quando un'azione fa ad esso riferimento.

Gli schemi di traduzione impongono un ordine di valutazione da sinistra a destra e permettono che nelle azioni semantiche siano contenuti frammenti di programma e in generale side-effect che non influiscano sulla valutazione degli attributi.

# Definizioni S-attribuite e L-attribuite

Una SDD è **S-attribuita** se tutti gli attributi sono sintetizzati

Una SDD è **L-attribuita** se gli attributi sono:

- sintetizzati e/o
- ereditati, che soddisfano il seguente vincolo:

per ogni produzione  $A \rightarrow X_1 X_2 \dots X_n$ , ogni attributo ereditato di  $X_j$  dipende solo da:

- attributi ereditati o sintetizzati dei simboli  $X_1 X_2 \dots X_{j-1}$  a sinistra di  $X_j$  nella produzione;
- attributi ereditati di  $A$ ;
- altri attributi di  $X_j$ , purchè non vi siano cicli nel grafo delle dipendenze formati dagli attributi di questa occorrenza di  $X_j$ .

Trasformazione di SDD S-attribuite e L-attribuite in Schemi di Traduzione (con azioni specificate al momento giusto)

Tipicamente gli SDT sono implementati durante la parsificazione, in particolare in due casi:

- SDD L-attribuita e grammatica LL-parsificabile
- SDD S-attribuita e grammatica LR-parsificabile

# Valutazione top-down di SDD L-attribuite

main *traduzione\_discesa\_ricorsiva*( )

var risultato

cc ← PROSS

risultato ← S( )

if (cc = '\$') stampa (“stringa corretta,  
la sua traduzione è:” risultato)

else ERRORE(...)

function A( $e_1, \dots, e_n$ )

var  $s_1, \dots, s_m, X_{1-a_1}, \dots, X_{1-a_k}, \dots, X_{h-a_1}, \dots, X_{h-a_r}$

if (cc ∈ Gui(A →  $\alpha_1$ )) body'( $\alpha_1$ )

elseif (cc ∈ Gui(A →  $\alpha_2$ )) body'( $\alpha_2$ )

....

elseif (cc ∈ Gui(A →  $\alpha_k$ )) body'( $\alpha_k$ )

else ERRORE(...)

return (< $s_1, \dots, s_m$ >)

Codice per le parti destre delle produzioni (body'( $\alpha_i$ )):

- Per ogni terminale i valori degli attributi vengono assegnati alle corrispondenti variabili e l'esame passa al simbolo successivo.
- Per ogni non terminale B si genera un'assegnazione  $c \leftarrow B(b_1, \dots, b_n)$ , che è una chiamata alla funzione associata a B. Poiché la SDD è L-attribuita, gli argomenti (attributi) saranno già stati calcolati e memorizzati nelle variabili locali.
- Le azioni semantiche vengono ricopiate dopo aver sostituito i riferimenti agli attributi con le variabili corrispondenti.

## Programma

$P ::= SL EOF$

## Istruzioni

- $S \rightarrow ID := E$
- $S \rightarrow \text{print } (E)$
- $S \rightarrow \text{read } (ID)$
- $S \rightarrow \text{while } (B) S$
- $S \rightarrow \text{if } (B) S$
- $S \rightarrow \text{if } (B) S \text{ else } S$
- $S \rightarrow \{SL\}$
- $SL \rightarrow SL ; S$
- $SL \rightarrow S$

## Espressioni aritmetiche

- $E \rightarrow T E'$
- $E' \rightarrow + T E'$
- $E' \rightarrow - T E'$
- $E' \rightarrow \varepsilon$
- $T \rightarrow F T'$
- $T' \rightarrow * F T'$
- $T' \rightarrow / F T'$
- $T' \rightarrow \varepsilon$
- $F \rightarrow (E)$
- $F \rightarrow \text{NUM}$
- $F \rightarrow \text{ID}$

## Espressioni booleane

- $B \rightarrow E \text{ RELOP } E$

# Dal linguaggio P al bytecode: espressioni aritmetiche

Per le espressioni aritmetiche usiamo un attributo sintetizzato *.code* il cui valore è la traduzione dell'espressione nel bytecode.

$$E \rightarrow T E' \{E.code = T.code \parallel E'.code\}$$
$$E' \rightarrow + T E'_1 \{E'.code = T.code \parallel \text{"iadd"} \parallel E'_1.code\}$$
$$E' \rightarrow - T E'_1 \{E'.code = T.code \parallel \text{"isub"} \parallel E'_1.code\}$$
$$E' \rightarrow \varepsilon \{E'.code = \text{" "}\}$$
$$T \rightarrow F T' \{T.code = F.code \parallel T'.code\}$$
$$T' \rightarrow * F T'_1 \{T'.code = F.code \parallel \text{"imul"} \parallel T'_1.code\}$$
$$T' \rightarrow / F T'_1 \{T'.code = F.code \parallel \text{"idiv"} \parallel T'_1.code\}$$
$$T' \rightarrow \varepsilon \{T'.code = \text{" "}\}$$
$$F \rightarrow (E) \{F.code = E.code\}$$
$$F \rightarrow \text{NUM} \{F.code = \text{ldc}(\text{NUM.val})\}$$
$$F \rightarrow \text{ID} \{F.code = \text{iload}(\text{addr}(\text{ID.lessema}))\}$$

Per le espressioni booleane usiamo tre attributi:

- a) *.code* sintetizzato per la traduzione in bytecode
- b) *.true* e *.false*, attributi ereditati che servono per definire le label delle prime istruzioni che traducono gli statement da eseguire nei casi B vero e B falso rispettivamente.

$$B \rightarrow E_1 == E_2 \quad \{B.code = E_1.code \parallel E_2.code \parallel \\ \parallel \text{'if\_icmpeq'} B.true \parallel \text{'goto'} B.false\}$$

# Dal linguaggio P al bytecode: SDT istruzioni

$P \rightarrow \{SL.next = newlabel( )\} SL EOF \{P.code = SL.code \parallel label(SL.next) \parallel 'stop'\}$

$S \rightarrow ID := E \{S.code = E.code \parallel istore(addr(id.lessema))\}$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = S.next\}$   
B)  $\{S_1.next = S.next\} S_1 \{S.code = B.code \parallel label(B.true) \parallel S_1.code\}$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = newlabel( )\}$   
B)  $\{S_1.next = S.next\} S_1$  else  
 $\{S_2.next = S.next\} S_2 \{S.code = B.code \parallel label(B.true) \parallel S_1.code \parallel$   
 $\parallel 'goto' S_1.next) \parallel label(B.false) \parallel S_2.code\}$

$S \rightarrow while ( \{begin = newlabel( ) ; B.true = newlabel( ) ; B.false = S.next\}$   
B)  $\{S_1.next = begin\}$   
 $S_1 \{S.code = label(begin) \parallel B.code \parallel label(B.true) \parallel S_1.code \parallel 'goto S_1.next\}$

$S \rightarrow \{ \{SL.next = S.next\} SL \} \{S.code = SL.code\}$

$SL \rightarrow \{SL_1.next = newlabel( )\} SL_1 ; \{S.next = SL.next\}$   
S  $\{SL.code = SL_1.code \parallel label(SL_1.next) \parallel S.code\}$

$SL \rightarrow \{S.next = SL.next\} S \{SL.code = S.code\}$

Sia per le espressioni aritmetiche, sia, e soprattutto, per gli statement, i valori degli attributi *.code* sono stringhe piuttosto lunghe.

Ad esempio, per lo statement: `while (x < 100) y := x * 2`

*S.code* = L2 iload(*addr(x)*) ldc 100 if\_icmplt L3 goto L1 L3  
    iload(*addr(x)*) ldc 2 imul istore(*addr(y)*) goto L2

1. L'attributo principale è sintetizzato
2. Le regole semantiche sono tali che:
  - a) L'attributo è ottenuto dal concatenamento di stringhe
  - b) Gli attributi dei non terminali si presentano nella regola semantica nello stesso ordine in cui i non terminali si presentano nel corpo della produzione



Traduzione on-the-fly

## Espressioni aritmetiche

$$E \rightarrow T E'$$
$$E' \rightarrow + T \{emit(iadd)\} E'_1$$
$$E' \rightarrow - T \{emit(isub)\} E'_1$$
$$E' \rightarrow \varepsilon$$
$$T \rightarrow F T'$$
$$T' \rightarrow * F \{emit(imul)\} T'_1$$
$$T' \rightarrow / T \{emit(idiv)\} T'_1$$
$$T' \rightarrow \varepsilon$$
$$F \rightarrow (E)$$
$$F \rightarrow \mathbf{NUM} \{emit(ldc(\mathbf{NUM}.val))\}$$
$$F \rightarrow \mathbf{ID} \{emit(iloadd(addr(\mathbf{ID}.lessema)))\}$$

## Espressioni booleane

$$B \rightarrow E_1 == E_2 \{emit('if_icmpeq' B.true ; emit('goto' B.false)\}$$

# Schemi di traduzione 'on-the-fly': statement

$P \rightarrow \{SL.next = newlabel( )\} SL \{emitlabel(SL.next)\} EOF \{emit('stop')\}$

$S \rightarrow ID := E \{emit(istore (addr(ID.lessema)))\}$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = S.next\}$   
B)  $\{emitlabel(B.true) ; S_1.next = S.next\} S_1$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = newlabel( )\}$   
B)  $\{emitlabel(B.true) ; S_1.next = S.next \}$   
 $S_1 \{emit('goto' S_1.next)\}$  else  $\{emitlabel(B.false), S_2.next = S.next \} S_2$

$S \rightarrow while ( \{begin=newlabel(); emitlabel(begin); B.true=newlabel( ) ; B.false=S.next\}$   
B)  $\{emitlabel(B.true); S_1.next = begin\}$   
 $S_1 \{emit('goto' S_1.next)\}$

$S \rightarrow \{ \{SL.next = S.next\} SL \}$

$SL \rightarrow \{SL_1.next = newlabel( )\} SL_1 ; \{emitlabel(SL_1.next) ; S.next = SL.next\} S$

$SL \rightarrow \{S.next = SL.next\} S$

# Schemi di traduzione 'on-the-fly': statement

$P \rightarrow \{SL.next = newlabel( )\} SL \{emitlabel(SL.next)\} EOF \{emit('stop')\}$

$S \rightarrow ID := E \{emit(istore (addr(ID.lessema)))\}$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = S.next\}$   
B)  $\{emitlabel(B.true) ; S_1.next = S.next\} S_1$

$S \rightarrow if ( \{B.true = newlabel( ) ; B.false = newlabel( )\}$   
B)  $\{emitlabel(B.true) ; S_1.next = S.next \}$   
 $S_1 \{emit('goto' S_1.next)\}$  else  $\{emitlabel(B.false), S_2.next = S.next \} S_2$

$S \rightarrow while ( \{begin=newlabel(); emitlabel(begin); B.true=newlabel( ) ; B.false=S.next\}$   
B)  $\{emitlabel(B.true); S_1.next = begin\}$   
 $S_1 \{emit('goto' S_1.next)\}$

$S \rightarrow \{ \{SL.next = S.next\} SL \}$

$SL \rightarrow \{SL_1.next = newlabel( )\} SL_1 ; \{emitlabel(SL_1.next) ; S.next = SL.next\} S$

$SL \rightarrow \{S.next = SL.next\} S$