

Analisi lessicale e sintattica

Nei programmi troviamo espressioni come questa:

$X1=34*y2+z$

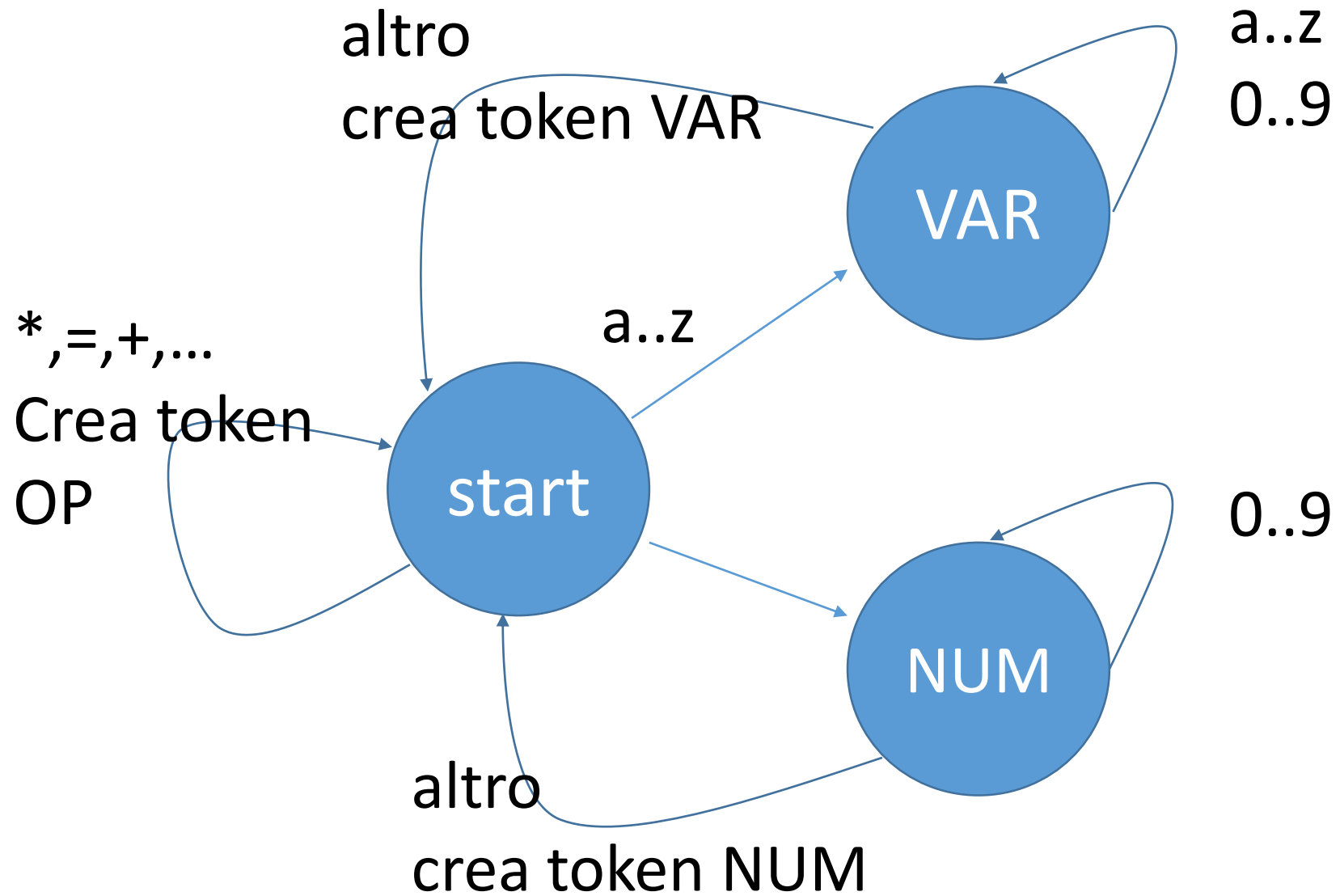
Si presentano come sequenze di caratteri

'X', '1', '=', '3', '4', '*',

Sarebbe meglio VAR(«X1»), OP('='), NUM(34), OP('*'),...

Questi oggetti si chiamano token lessicali

Questa traduzione è compito dell'analisi lessicale ed è realizzata da un'automa a stati finiti che genera output



Quindi alla seconda fase di analisi sintattica (parsing) arriva una sequenza di token

L'analisi sintattica si basa su una CFG che è fissata una volta per tutte al momento in cui viene definito il linguaggio

Ci sono diversi modi di farla.

Noi vediamo un caso semplice: analisi «recursive descent» che si può usare se la CFG è LL(1)

Esistono anche CFG LL(2), LL(3),..., LL(k)

Idea del parsing recursive descent

T1,T2,T3,.....



S

T1 deve indicare con quale produzione espandere S

Supponiamo sia $S \rightarrow AB$, allora

T2 deve indicare con quale produzione espandere A

Passeremo a B quando la procedura ha consumato i token che corrispondono all'albero generato da A

LL(2) = possiamo guardare 2 token

LL(k) ne possiamo guardare k

Questo parsing è top-down : l'albero di derivazione è «costruito» dalla radice alla frontiera

Esistono grammatiche che non sono LL(1) e nemmeno LL(k)

Ci sono metodi di parsing bottom-up, cioè dalla frontiera alla radice. Sono quelli utilizzabili per una classe più ampia di CFG.

Sarebbe comodo se ogni produzione:

$A \rightarrow T B \dots$, dove T sta per Token

Ma in generale non è così.

Però possiamo calcolarci

$\text{First}(A) = \{ T_1 \mid T_1 \dots T_k \text{ in } L(A) \}$

$\text{First}(A)$ può contenere ε

Consideriamo questa CFG per le espressioni

1. $E \rightarrow E + T \mid T$

2. $T \rightarrow T * F \mid F$

3. $F \rightarrow \text{NUM} \mid \text{VAR} \mid (E)$

Si inizia con $\text{First}(E) = \text{First}(T) = \text{First}(F) = \emptyset$

Poi si aggiunge $($, NUM e VAR

Quindi ora possiamo provare a costruire un parser recursive descent:

Con E se il primo token è NUM posso usare 1.1, ma anche 1.2 !! Non va

Non può funzionare perché la grammatica è ricorsiva a sinistra:

$E \rightarrow E+T \mid T$ e quindi guardando il primo token non potrò mai sapere se si applica $E+T$ o T

Anche E genera $T+T\dots$. E quindi non è distinguibile da T

Dobbiamo cambiare la grammatica eliminando la ricorsione a sinistra

$E \rightarrow T E'$

$E' \rightarrow + T E' \mid \varepsilon$

$T \rightarrow F T'$

$T' \rightarrow * F T' \mid \varepsilon$

$F \rightarrow \text{NUM} \mid \text{VAR} \mid (E)$

1. $E \rightarrow T E'$

2. $E' \rightarrow + T E' \mid \varepsilon$

3. $T \rightarrow F T'$

4. $T' \rightarrow * F T' \mid \varepsilon$

5. $F \rightarrow \text{NUM} \mid \text{VAR} \mid (E)$

First	E	E'	T	T'	F
0	0	0	0	0	0
1		+ ε		* ε	(NUM VAR
2			(NUM VAR		
3	(NUM VAR				

	(VAR	NUM	*	+)
E	1	1	1			
E'					2.1	
T	3	3	3			
T'				4.1		
F	5.3	5.2	5.1			

Non sappiamo come usare le ε -produzioni

Dobbiamo sapere per ogni variabile quali terminali possono venire subito dopo

$\text{Follow}(A) = \{T \mid S \Rightarrow^* \alpha A \beta \text{ e } \text{First}(\beta) \text{ contiene } T\}$

1. $E \rightarrow T E'$

2. $E' \rightarrow + T E' \mid \varepsilon$

3. $T \rightarrow F T'$

4. $T' \rightarrow * F T' \mid \varepsilon$

5. $F \rightarrow \text{NUM} \mid \text{VAR} \mid (E)$

Follow	E	E'	T	T'	F
0	\$	0	0	0	0
1)	\$	+		*
2)	\$	+	
3)	\$	+
4)	\$
5)

1. $E \rightarrow T E'$

2. $E' \rightarrow + T E' \mid \varepsilon$

3. $T \rightarrow F T'$

4. $T' \rightarrow * F T' \mid \varepsilon$

5. $F \rightarrow \text{NUM} \mid \text{VAR} \mid (E)$

$\text{Follow}(E) = \{\$,)\}$

$\text{Follow}(E') = \{\$,)\}$

$\text{Follow}(T) = \{\$,), +\}$

$\text{Follow}(T') = \{\$,), +\}$

$\text{Follow}(F) = \{\$,), +, *\}$

	(VAR	NUM	*	+)	\$
E	1	1	1				
E'					2.1	2.2	2.2
T	3	3	3				
T'				4.1	4.2	4.2	4.2
F	5.3	5.2	5.1				

	(VAR	NUM	*	+)	\$
E	1	1	1				

Token* parseE(Token* T)

{

 switch (T[0])

 case '(', VAR, NUM :

 Token * k1=parseT(T+1); return parseE'(k1);

 default: throw errore

}

	(VAR	NUM	*	+)	\$
T'				4.1	4.2	4.2	4.2

```

Token* parseT'(Token * T)
{
    switch(T[0])
    case '*' : Token* k1=parseF(T+1);
                return parseT'(k1);
    case '+', ')', '$' : return T;

}

```