

# Parser LL(1)

Lunedì 12 Novembre

# Parser LL(1)

Il termine LL(1) ha il seguente significato:

- 1. la prima L, significa che l'input è analizzato da sinistra verso destra
- 2. la seconda L, significa il parser costruisce una derivazione leftmost per la stringa di input
- 3. il numero 1, significa che l'algoritmo utilizza soltanto un solo simbolo dell'input per risolvere le scelte del parser (ci sono varianti con k simboli)

# Esempio

Il linguaggio delle parentesi bilanciate

$S \rightarrow ( S ) S$

$S \rightarrow \varepsilon$

e vediamo come opera il parser **LL(1)**

per riconoscere la stringa "( )"

Il parser consiste di una **pila**, che contiene inizialmente il simbolo "\$" (fondo della pila), ed un **input**, la cui fine è marcata dal simbolo "\$"

(EOF generato dallo scanner)

pila	input	azione
\$	( ) \$	

-il parsing inizia inserendo il simbolo iniziale in testa alla pila

pila	input	azione
\$ S	( ) \$	

- il parser **accetta** una stringa di input se, dopo una sequenza di azioni, la pila contiene “\$” e la stringa di input è “\$”

pila	input	azione
...	... ..	
\$	\$	accept

- ogni volta che in testa alla pila c'è un simbolo non terminale X, lo si espande secondo una produzione  $X \rightarrow \gamma$ , che viene scelta a seconda del simbolo in testa all'input e ai valori di una tabella (la parte destra della produzione viene invertita sulla pila)

pila	input	azione
\$ S	( ) \$	$S \rightarrow ( S ) S$

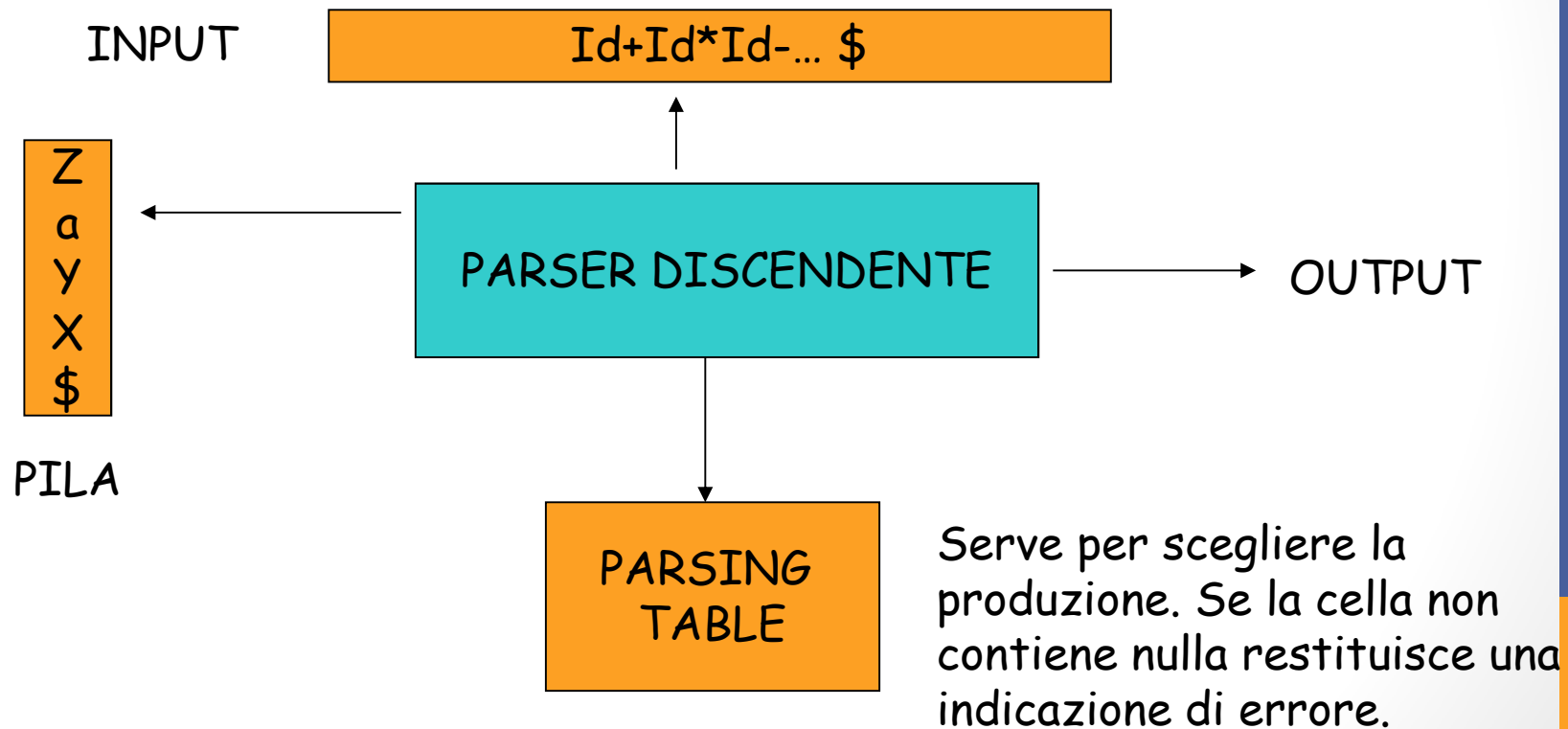
-ogni volta che sulla pila c'è un simbolo terminale **t**, si controlla che in testa all'input ci sia anche lo stesso simbolo, nel qual caso lo si elimina sia dalla pila che dall'input; altrimenti è **errore**

per costruire un parser **LL(1)**, bisogna costruire una tabella – la **tabella LL(1)** – che determina la regola da usare per l'espansione, dati il simbolo non-terminale e il carattere in input

pila	input	azione
\$ S	( ) \$	$S \rightarrow ( S ) S$
\$ S ) S (	( ) \$	match
\$ S ) S	) \$	$S \rightarrow \varepsilon$
\$ S )	) \$	match
\$ S	\$	$S \rightarrow \varepsilon$
\$	\$	accept

Se l'input è generato dalla grammatica questo parsing fornisce una derivazione leftmost, altrimenti produce un'indicazione d'errore.

# I Parser LL(1) sono parser discendenti non ricorsivi



# Costruzione della parsing table di un LL(1)

1. la tabella ha simboli non-terminali come righe, e simboli terminali come colonne
  2. per ogni regola  $X \rightarrow \gamma$  di  $G$ , si inserisce tale regola nella casella  $(X, t)$ , per ogni  $t$  tale che  $\gamma \Rightarrow^* t \beta$
  3. per ogni regola  $X \rightarrow \gamma$  di  $G$ , per cui  $\gamma \Rightarrow^* \epsilon$ , si inserisce nella casella  $(X, t)$  tale regola, per ogni  $t$  tale che  $S \Rightarrow^* \beta X t \alpha$
- (Bisogna conoscere questi simboli  $t$  in grado di far effettuare la scelta della produzione  
le regole 2 e 3 sono difficili da implementare: gli algoritmi per risolverle sono discussi in seguito vedi **FIRST** e **FOLLOW**)

**per esempio:** la tabella per la grammatica

$S \rightarrow (S)S, S \rightarrow \epsilon$		
(	)	\$
S	$S \rightarrow (S)S$	$S \rightarrow \epsilon$
	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$

# First

□ E' una funzione che consente di costruire le entries della tabella, quando possibile.

**FIRST( $\gamma$ )**: insieme dei simboli **terminali** che si trovano all'inizio delle stringhe derivate da  $\gamma$  ( $\gamma$  è una stringa di simboli terminali e non)



# First

1. se  $X$  è un terminale, allora  $\text{FIRST}(X) = \{ X \}$ ,
2. se  $X \rightarrow \varepsilon$  appartiene alla grammatica, allora aggiungi  $\varepsilon$  a  $\text{FIRST}(X)$ ,
3. se  $X \rightarrow Y_1 Y_2 \dots Y_k$  appartiene alla grammatica, allora:
  - se  $a \in \text{FIRST}(Y_i)$  per qualche  $i$  ed  $\varepsilon$  sta in  $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \dots, \text{FIRST}(Y_{i-1})$ , aggiungi  $a$  in  $\text{FIRST}(X)$  ;
  - se tutti gli insiemi  $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \dots, \text{FIRST}(Y_k)$ , contengono  $\varepsilon$ , aggiungi  $\varepsilon$  a  $\text{FIRST}(X)$ ;
4. per definire l'insieme **FIRST**( $\gamma$ ), dove  $\gamma = X_1 X_2 \dots X_k$  (una stringa di terminali e non), si procede reiterando le regole seguenti:
  - aggiungi **FIRST**( $X_1$ )  $\setminus \{ \varepsilon \}$  a **FIRST**( $\gamma$ )
  - se per qualche  $i < k$ , tutti gli insiemi **FIRST**( $X_1$ ),  $\dots$ , **FIRST**( $X_i$ ) contengono  $\varepsilon$ , allora aggiungi **FIRST**( $X_{i+1}$ )  $\setminus \{ \varepsilon \}$  a **FIRST**( $\gamma$ )
  - se tutti gli insiemi **FIRST**( $X_1$ ),  $\dots$ , **FIRST**( $X_k$ ) contengono  $\varepsilon$ , allora aggiungi  $\varepsilon$  a **FIRST**( $\gamma$ )

Se  $X \rightarrow^* \varepsilon$  allora  $\varepsilon \in \text{FIRST}(X)$  ( $X$  è detto annullabile)

# Esempio di calcolo di First

Nel caso della grammatica delle espressioni aritmetiche senza ricorsioni sinistre:

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \text{ id}$$

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

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

$$F \rightarrow (E) \text{ id}$$

$$\text{FIRST}(\text{id})=\{\text{id}\} \quad \text{FIRST}(()=\{(\}$$

$$\text{FIRST}(+)=\{+\}$$

$$\text{FIRST}(*)=\{*\}$$

$$\text{FIRST}())=\{)\}$$

$$\text{FIRST}(F)=\{\text{id}, (\}$$

$$\text{FIRST}(T)=\text{FIRST}(E)=\text{FIRST}(F)$$

$$\text{FIRST}(E')=\{+, \varepsilon\} \quad \text{FIRST}(T')=\{*, \varepsilon\}$$

## algoritmo in pseudo-C per calcolare FIRST

```
for all terminals X and  $\epsilon$  do FIRST(X) = {X} ;
for all non-terminals X do FIRST(X) = { } ;
while (there are changes to any FIRST(X)) do
  for each production  $X \rightarrow Y_1 Y_2 \dots Y_k$  do
    {
      i := 1 ; continue := true ;
      while (continue == true && i <= k) do
        {
          add FIRST( $Y_i$ ) \ { $\epsilon$ } to FIRST(X) ;
          if ( $\epsilon$  is not in FIRST( $Y_i$ )) continue := false ;
          i := i+1 ;
        }
      if (continue == true) add  $\epsilon$  to FIRST(X) ;
    }
```

# FIRST può non bastare

**Osservazione:** se la grammatica contiene due produzioni

$X \rightarrow \gamma_1$

$X \rightarrow \gamma_2$

(stesso simbolo non-terminale a sinistra, due sequenze differenti a destra)

e **FIRST**( $\gamma_1$ )  $\cap$  **FIRST**( $\gamma_2$ ) è non vuota.

Allora la grammatica non può essere analizzata col parsing top-down:  
se  $t \in \mathbf{FIRST}(\gamma_1) \cap \mathbf{FIRST}(\gamma_2)$  allora, il parsing discendente non saprà cosa fare quando il primo simbolo dell'input è  $t$

# Follow

Dato un **non terminale**  $X$ , l'insieme **FOLLOW**( $X$ ) è l'insieme dei simboli terminali, eventualmente  $\$$ , che appaiono alla destra di  $X$  in qualche forma sentenziale ed è definito come segue:

1. se  $X$  è l'assioma, allora aggiungi  $\$$  a **FOLLOW**( $X$ )
2. se c'è una produzione  $A \rightarrow \alpha X \gamma$ , allora aggiungi **FIRST**( $\gamma$ ) \setminus \{ \epsilon \} a **FOLLOW**( $X$ )
3. se c'è una produzione  $A \rightarrow \alpha X \gamma$  per cui  $\epsilon \in \mathbf{FIRST}(\gamma)$ , allora aggiungi **FOLLOW**( $A$ ) a **FOLLOW**( $X$ )

## Osservazioni:

1. quando l'assioma non compare a destra delle produzioni, il “\$” è l'unico simbolo nel suo insieme **FOLLOW**
2. l'insieme **FOLLOW** non contiene mai “ $\epsilon$ ”
3. **FOLLOW** è definito soltanto per non-terminali: potremmo generalizzare la definizione ma ciò è inutile per la tabella **LL(1)**
4. la definizione di **FOLLOW** lavora “alla destra” di una produzione, mentre quella di **FIRST** lavora “alla sinistra”: una produzione  $X \rightarrow \alpha$  non ha alcuna informazione su **FOLLOW**( $X$ ), se  $X$  non è presente in  $\alpha$

# Algoritmo in pseudo-C per calcolare FOLLOW

```
FOLLOW(Axiom) = { $ } ;  
for all (nonterminal(X) && X != Axiom) do FOLLOW(X)= {};  
while (there are changes to any FOLLOW set) do  
    for each production  $X \rightarrow Y_1 Y_2 \dots Y_k$  do  
        for each nonterminal( $Y_i$ ) do {  
            add  $\text{FIRST}(Y_{i+1} \dots Y_k) \setminus \{ \epsilon \}$  to  $\text{FOLLOW}(Y_i)$  ;  
/* Note: if  $i=k$  then  $Y_{i+1} \dots Y_k = \epsilon$  */  
            if  $\epsilon$  is in  $\text{FIRST}(Y_{i+1} \dots Y_k)$   
                add  $\text{FOLLOW}(X)$  to  $\text{FOLLOW}(Y_i)$ ;  
        }  
    }
```

# Esempio di calcolo di Follow

Nel caso della grammatica delle espressioni aritmetiche senza ricorsioni sinistre:

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \text{ id}$$

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

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

$$F \rightarrow (E) \text{ id}$$

$$\text{FIRST}(\text{id})=\{\text{id}\} \quad \text{FIRST}('(')=\{('\}$$

$$\text{FIRST}('+')=\{+\}$$

$$\text{FIRST}('*')=\{*\}$$

$$\text{FIRST}('')=\{')'\}$$

$$\text{FIRST}(F)=\{\text{id}, '('\}$$

$$\text{FIRST}(T)=\text{FIRST}(E)=\text{FIRST}(F)$$

$$\text{FIRST}(E')=\{+, \varepsilon\} \quad \text{FIRST}(T')=\{*, \varepsilon\}$$

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

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

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

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

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

# Grammatica LL(1)

*Una grammatica è LL(1) se e solo se per ogni produzione del tipo  $A \rightarrow \alpha | \beta$  si ha:*

- $\alpha$  e  $\beta$  non derivano stringhe che cominciano con lo stesso simbolo a.
- Al più uno tra i due può derivare la stringa vuota.
- Se  $\beta \rightarrow^* \epsilon$  allora  $\alpha$  non deriva stringhe che cominciano con terminali che stanno in FOLLOW(A). Analogamente per  $\alpha$

Equivalentemente affinché una grammatica sia LL(1) deve avvenire che per ogni coppia di produzioni  $A \rightarrow \alpha | \beta$

1. FIRST( $\alpha$ ) e FIRST( $\beta$ ) devono essere disgiunti
2. Se  $\epsilon$  è in FIRST( $\beta$ ) allora FIRST( $\alpha$ ) e FOLLOW(A) devono essere disgiunti.



# Come costruire la tabella?

1. Per ogni regola  $X \rightarrow \alpha$  di  $G$ , si inserisce nella casella  $(X, t)$  la regola  $X \rightarrow \alpha$ , per ogni  $t$  tale che  $t \in \mathbf{FIRST}(\alpha)$
2. Per ogni regola  $X \rightarrow \alpha$  di  $G$ , per cui  $\alpha \Rightarrow^* \varepsilon$  ( $\varepsilon \in \mathbf{FIRST}(\alpha)$ ), si inserisce nella casella  $(X, t)$  la regola  $X \rightarrow \alpha$ , per ogni  $t$  tale che  $t \in \mathbf{FOLLOW}(X)$ .  
Se  $\varepsilon \in \mathbf{FIRST}(\alpha)$  and  $\$ \in \mathbf{FOLLOW}(X)$ , si inserisce la regola  $X \rightarrow \alpha$  in  $(X, \$)$ .
3. Le caselle non definite definiscono un errore.

NOTA: Se  $G$  è ricorsiva sinistra o ambigua, la tabella avrà caselle con valori multipli.

# Altra definizione di Grammatica LL(1)

*Una grammatica la cui tabella **LL(1)** non contiene più di un elemento nelle caselle è detta **LL(1)***

**Osservazione:** per costruzione una grammatica **LL(1)** non è ambigua, né ricorsiva sinistra

# Algoritmo di parsing LL(1)

```
LL1_parser( stack p, input i, LL1_table M, initial S ) {  
    /* p, i sono pile, S è l'assioma */  
    int error = 0 ; p = push(p, $) ;  
    p = push(p, S) ;  
    while ( top(p) ≠ $ && top(i) ≠ $ && !error){  
        if isterminal(top(p)) {  
            if top(p) == top(i) { p = pop(p) ; i = avanza(i) ; }  
            else error = 1 ;}  
        else { /* top(p) è un non-terminale, bisogna espandere */  
            if isempty(M[top(p),top(i)]) error = 1;  
            else { /* M[top(p),top(i)] == X -> X1 ... Xn; */  
                p = pop(p) ;  
                for (j = n ; j > 0 ; j = j-1)  
                    p = push(p, Xj) ;}}  
        }  
        if (!error) accept() ;  
        else raise_error() ;  
    }  
}
```

# Complessità di calcolo del parser LL(1)

E' lineare nella lunghezza  $n$  della stringa sorgente, poiché viene consumato un carattere per volta.

**ESEMPIO:** costruire il parser **LL(1)** per la grammatica

$S \rightarrow ( S ) \quad S \rightarrow [ S ] \quad S \rightarrow \{ S \}$

$S \rightarrow \epsilon$

insiemi **FIRST**, **FOLLOW**:

	<b>FIRST</b>	<b>FOLLOW</b>
S	(,[,{,ε	),],},\$

la tabella **LL(1)**:

	(	[	{	)	]	}	\$
S	$S \rightarrow (S)$	$S \rightarrow [S]$	$S \rightarrow \{S\}$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$