

PARSER LL(1)

PARSER LL(1)

Il termine **LL(1)** ha il seguente significato:

1. La prima **L**, sta per **Left to right**, 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 guarda un solo simbolo dell'input in avanti per risolvere le scelte del parser (ci sono varianti con k simboli)

ESEMPIO

IL LINGUAGGIO DELLE PARENTESI BILANCIATE

$$S \rightarrow (S) S$$

$$S \rightarrow \epsilon$$

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)

Per comodità mettiamo il fondo della pila a destra e impiliamo a sinistra

Pila	Input	azione
\$	()\$	

PARSING LL(1)

Il parsing **inizia** inserendo il simbolo iniziale S in testa alla pila

Il parser **accetta** una stringa di input se, dopo una sequenza di azioni, la pila contiene "\$" e la stringa di input è "\$"

Ogni volta che in testa alla pila c'è un simbolo non terminale x , lo si **espande secondo una produzione $x \rightarrow y$** , che viene scelta a seconda del simbolo in testa all'input e ai valori di una tabella

Ogni volta che sulla pila c'è un simbolo terminale t , si controlla se il simbolo dell'input che si sta leggendo sia lo stesso, nel qual caso lo si elimina sia dalla pila e la testina di lettura avanza sull'input; altrimenti si è riscontrato un **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
$S\$$	$()\$$	

Pila	Input	azione
...
...
$\$$	$\$$	

Pila	Input	azione
...
$S\$$	$()\$$	$S \rightarrow (S)S$

ESEMPIO

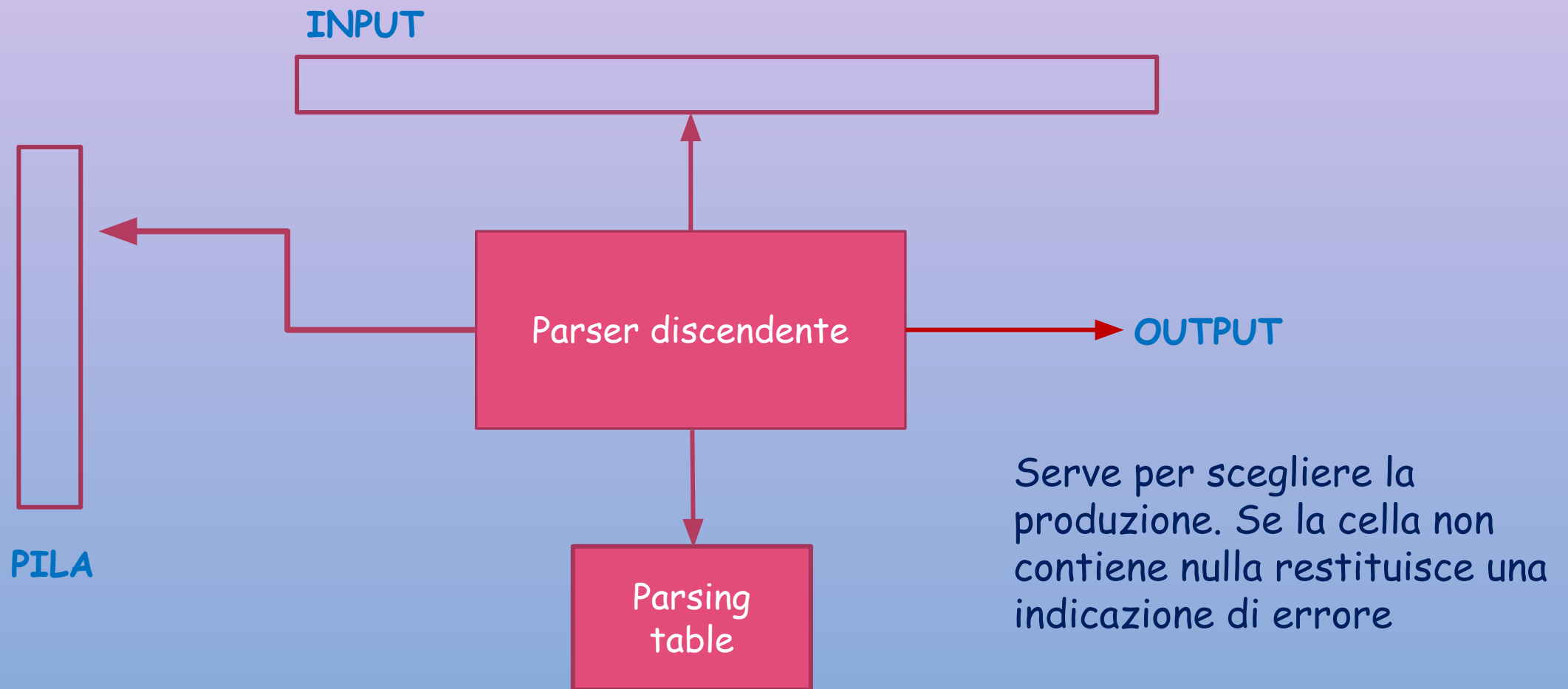
IL LINGUAGGIO DELLE PARENTESI BILANCIATE

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

Il linguaggio genera «()»?
Se l'input è generato dalla grammatica
questo parsing fornisce una
derivazione leftmost, altrimenti
produce un'indicazione d'errore.

Pila	Input	Azione
S\$	() \$	$S \rightarrow (S)S$
(S) S \$	() \$	match
S) S \$) \$	$S \rightarrow \epsilon$
) S \$) \$	match
S \$	\$	$S \rightarrow \epsilon$
\$	\$	accept

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. La tabella contiene regole di produzione
3. Ogni regola $x \rightarrow y$ di G per cui $y \Rightarrow^* t\beta$ (t appare come primo simbolo in una derivazione di y), si inserisce nella casella (x, t) (per ogni t per cui questo avviene)
4. Ogni regola $x \rightarrow y$ di G , per cui $y \Rightarrow^* \epsilon$, si inserisce nella casella (x, t) , 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**

ESEMPIO

Data la grammatica:

$S \rightarrow (S) S,$

$S \rightarrow \epsilon$

Parsing Table

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

FIRST E FOLLOW

FIRST e **FOLLOW** sono due funzioni, utilizzate sia dai parser top down che dai parser bottom-up, associate a una grammatica G .

In particolare nel parser top-down permettono di effettuare una scelta fra diverse produzioni in base al solo simbolo di input immediatamente successivo alla posizione della testina di lettura

FIRST

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

First(α): insieme dei simboli **terminali** che si trovano all'inizio delle stringhe derivabili da α (α è una stringa di simboli terminali e non).

Se $\alpha \Rightarrow^* \epsilon$ allora anche ϵ appartiene a **FIRST(α)**

Nella slide successiva l'algoritmo per la costruzione di **FIRST(α)**

FIRST

1. Se x è un terminale, allora $\text{FIRST}(x) = \{x\}$,
2. Se $x \rightarrow \epsilon$ appartiene alla grammatica, allora aggiungi ϵ 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 ϵ sta in $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \dots, \text{FIRST}(Y_{i-1})$, aggiungi a a $\text{FIRST}(X)$;
 - se tutti gli insiemi $\text{FIRST}(Y_1), \text{FIRST}(Y_2), \dots, \text{FIRST}(Y_k)$, contengono ϵ , aggiungi ϵ a $\text{FIRST}(X)$;
4. Per definire l'insieme **FIRST** (γ), 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) $\setminus \{\epsilon\}$ a **FIRST**(γ)
 - Se per qualche $i < k$, tutti gli insiemi **FIRST** (x_1), ..., **FIRST** (x_i) contengono ϵ , allora aggiungi **FIRST** (x_{i+1}) $\setminus \{\epsilon\}$ a **FIRST** (γ)
 - Se tutti gli insiemi **FIRST** (x_1), ..., **FIRST** (x_k) contengono ϵ , allora aggiungi ϵ a **FIRST**(γ)

NOTA: Se $x \rightarrow^* \epsilon$ allora $\epsilon \in \text{FIRST}(x)$ (x è detto annullabile)

ESEMPIO

Nel caso della grammatica delle espressioni aritmetiche senza ricorsioni sinistre:

$E \rightarrow E+T \mid T$
 $T \rightarrow T*F \mid F$
 $F \rightarrow (E) \mid id$

$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow (E) \mid id$

Simboli
terminali

$FIRST(id) = \{id\}$
 $FIRST(() = \{(\}$
 $FIRST(+) = \{+\}$
 $FIRST(*) = \{*\}$
 $FIRST()) = \{)\}$

Simboli
nonterminali

$FIRST(F) = \{id, (\}$
 $FIRST(T) = FIRST(E) = FIRST(F)$
 $FIRST(E') = \{+, \epsilon\}$
 $FIRST(T') = \{*, \epsilon\}$

Per quanto riguarda le stringhe nelle produzioni

Per esempio, $FIRST(TE') = FIRST(T)$

$FIRST(+TE') = \{+\}, \dots$

PSEUDO CODICE PER FIRST

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

FIRST PUÒ NON BASTARE

Il nostro obiettivo è stabilire una regola per costruire una tabella che, in base al simbolo successivo a quello letto è in grado di stabilire in maniera univoca quale regola usare

A questo proposito FIRST non basta.

Se infatti abbiamo delle regole che producono in uno o più passi la **parola vuota** ($A \rightarrow^* \epsilon$), non siamo in grado mediante la sola funzione first di conoscere qual è il prossimo simbolo prodotto, e non possiamo quindi fare la scelta in base al prossimo simbolo di input.

In questo caso occorre utilizzare l'insieme **FOLLOW**

FIRST PUÒ NON BASTARE

Il nostro obiettivo è stabilire una regola per costruire una tabella che, in base al simbolo successivo a quello letto è in grado di stabilire in maniera univoca quale regola usare

A questo proposito FIRST non basta. Per esempio, 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 $\text{FIRST}(\gamma_1) \cap \text{FIRST}(\gamma_2)$ è non vuota.

Allora la grammatica non può essere analizzata col parsing predittivo top-down:

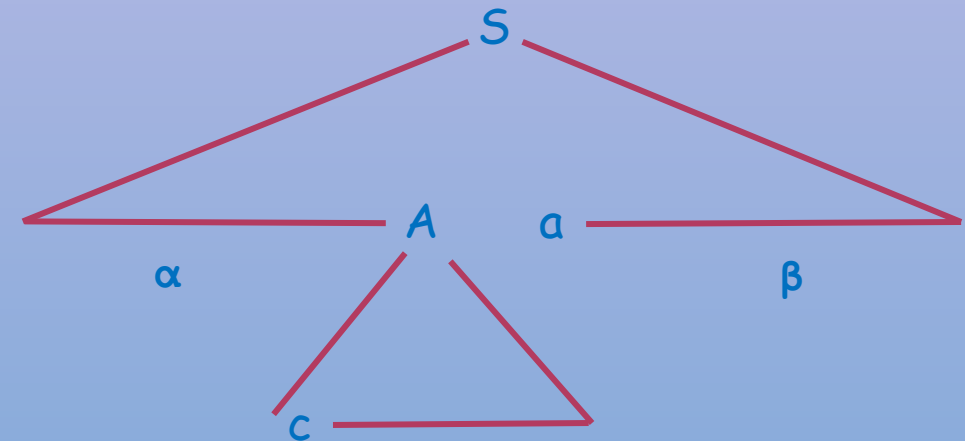
Infatti se $t \in \text{FIRST}(\gamma_1) \cap \text{FIRST}(\gamma_2)$ allora, il parsing discendente non saprà cosa fare quando il primo simbolo dell'input è t

In questo caso occorre utilizzare l'insieme **FOLLOW**

FOLLOW

FOLLOW(A) con **A** non terminale è l'insieme dei terminali che possono apparire immediatamente alla destra di **A** in una forma sentenziale, cioè l'insieme dei terminali **a** tali che esiste una derivazione della forma $X \rightarrow \alpha A a \beta$, con α e β generiche forme sentenziali. Tra **A** e **a** potrebbero anche esserci altri simboli che generano la parola vuota. Inoltre se **A** appare come ultimo simbolo allora **\$** appartiene a **FOLLOW(A)**

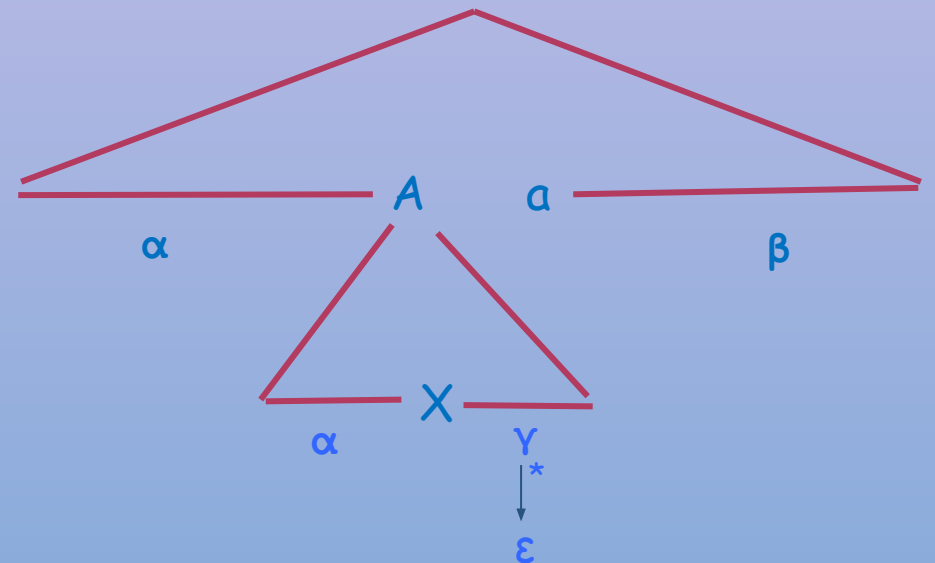
Nella figura a fianco, **c** appartiene a **FIRST(A)** e **a** appartiene a **FOLLOW(A)**



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 costruito 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** (γ) $\setminus \{\epsilon\}$ a **FOLLOW** (X);
3. Se c'è una produzione $A \rightarrow \alpha X$ oppure una produzione $A \rightarrow \alpha X \gamma$ per cui $\epsilon \in \text{FIRST}(\gamma)$, allora aggiungi **FOLLOW** (A) a **FOLLOW** (X).



FOLLOW - 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 " ϵ ";
3. FOLLOW(X) 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 α

PSEUDO CODICE PER 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

Nel caso della grammatica delle espressioni aritmetiche senza ricorsioni sinistre:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \mid id \end{aligned}$$
$$\begin{aligned} \text{FIRST}(id) &= \{id\} \\ \text{FIRST}('(') &= \{(\} \\ \text{FIRST}('+') &= \{+\} \\ \text{FIRST}('*') &= \{*\} \\ \text{FIRST}()) &= \{)\} \\ \text{FIRST}(T) &= \text{FIRST}(E) = \text{FIRST}(F) = \{id, '('\} \\ \text{FIRST}(E') &= \{+, \varepsilon\} \\ \text{FIRST}(T') &= \{*, \varepsilon\} \end{aligned}$$
$$\begin{aligned} \text{FOLLOW}(E) &= \{\$, \}\} \\ \text{FOLLOW}(E') &= \{\$, \}\} \\ \text{FOLLOW}(T) &= \{\$, \}, +\} \\ \text{FOLLOW}(T') &= \{\$, \}, +\} \\ \text{FOLLOW}(F) &= \{\$, \}, +, *\} \end{aligned}$$

COME COSTRUIRE LA TABELLA DI PARSING PREDITTIVO?

1. Per ogni regola $X \rightarrow \alpha$ della grammatica, e per ogni terminale a in $\text{FIRST}(\alpha)$, si aggiunga la regola $X \rightarrow \alpha$ in $M[X, a]$.
2. Per ogni regola $X \rightarrow \alpha$, se $\epsilon \in \text{FIRST}(\alpha)$ (come caso particolare la regola $X \rightarrow \epsilon$), per ogni $b \in \text{FOLLOW}(X)$, si inserisca nella casella $M(X, b)$ la regola $X \rightarrow \alpha$.
Se $\epsilon \in \text{FIRST}(\alpha)$ e $\$ \in \text{FOLLOW}(X)$, si inserisce la regola $X \rightarrow \alpha$ in $M(X, \$)$.

Le caselle non definite individuano un errore.

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

$A \rightarrow Ac$	Entrambe le regole contengono b nell'insieme first, dunque entrambe
$A \rightarrow b$	vanno nella stessa casella (A, b)

GRAMMATICA LL(1)

Definizione:

Una grammatica in cui ogni casella della tabella di parsing non contiene più di un elemento è detta **LL(1)**

Osservazione 1: per costruzione una grammatica **LL(1)** non è ambigua, né ricorsiva sinistra. Quindi in questi casi la grammatica non può essere **LL(1)**.

Osservazione 2: Per stabilire se una grammatica è **LL(1)** bisogna costruire la tabella e verificare che in ogni casella ci sia al più un elemento.

GRAMMATICA LL(1)

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

- α e β 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 α non deriva stringhe che cominciano con terminali che stanno in $FOLLOW(A)$. Analogamente per α .

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

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

PARSER LL(1)

E' sempre possibile realizzare un parser predittivo con una grammatica LL(1) perché ad ogni passo la produzione per un terminale può essere decisa, analizzando la tabella, in maniera unica analizzando solo il carattere di input successivo alla posizione della testina.

PARSER 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

La complessità dell'algoritmo del parser LL(1) è lineare nella lunghezza n della stringa sorgente, poiché viene consumato un carattere per volta.

ESERCIZI

Stabilire quali delle seguenti grammatiche sono LL(1)

1. $G: S \rightarrow aSb \mid \varepsilon$
2. $G: S \rightarrow +SS \mid *SS \mid id$
3. $G: S \rightarrow aSb \mid aSc \mid \varepsilon$
4. $G: S \rightarrow aSa \mid bSb \mid a \mid b$
5. $G: S \rightarrow iEtS \mid iEtEeS \mid a$
 $E \rightarrow b$

ESEMPIO

Costruire il parser **LL(1)** per la grammatica

$S \rightarrow (S)$

$S \rightarrow [S]$

$S \rightarrow \{ S \}$

$S \rightarrow \epsilon$

	First	FOLLOW
S	(,[,{,ε),,]}, \$

La tabella **LL(1)**:

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

Esercizio:

Effettuare il parsing della stringa $([\{\}])$ e il parsing di $([])$

ESERCIZIO 1

$G: S \rightarrow aSb \mid \epsilon$

$\text{FIRST}(S) = \{a, \epsilon\}$

$\text{FOLLOW}(S) = \{b, \$\}$

Tabella LL(1)

	a	b	\$
S	$S \rightarrow aSb$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$

ESERCIZIO 2

G :

$S \rightarrow +SS \mid *SS \mid id$

$FIRST(S) = \{+, *, id\}$

$FOLLOW(S) = \{+, *, id, \$\}$

	+	*	id	\$
S	$S \rightarrow +SS$	$S \rightarrow *SS$	$S \rightarrow id$	

ESERCIZIO 3

$G: S \rightarrow aSb | aSc | \epsilon$

Fattorizzazione
sinistra

$S \rightarrow aSA | \epsilon$
 $A \rightarrow b | c$

La grammatica presenta due regole con prefissi in comune.

$\text{FIRST}(S) = \{a, \epsilon\}$

$\text{FIRST}(A) = \{b, c\}$

$\text{FOLLOW}(S) = \{b, c, \$\}$

$\text{FOLLOW}(A) = \{b, c, \$\}$

	a	b	c	\$
S	$S \rightarrow aSA$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$
A		$A \rightarrow b$	$A \rightarrow c$	

ESERCIZIO 4

$G: S \rightarrow aSa \mid bSb \mid a \mid b$

Fattorizzazione
sinistra

$S \rightarrow aA \mid bB$
 $A \rightarrow Sa \mid \varepsilon$
 $B \rightarrow Sb \mid \varepsilon$

$\text{FIRST}(S) = \{a, b\}$
 $\text{FIRST}(A) = \{a, b, \varepsilon\}$
 $\text{FIRST}(B) = \{a, b, \varepsilon\}$

$\text{FOLLOW}(S) = \{a, b, \$\}$
 $\text{FOLLOW}(A) = \{a, b, \$\}$
 $\text{FOLLOW}(B) = \{a, b, \$\}$

	a	b	\$
S	$S \rightarrow aA$	$S \rightarrow bB$	
A	$A \rightarrow Sa$ $A \rightarrow \varepsilon$	$A \rightarrow Sa$ $A \rightarrow \varepsilon$	$A \rightarrow \varepsilon$
B	$A \rightarrow Sb$ $B \rightarrow \varepsilon$	$A \rightarrow Sb$ $B \rightarrow \varepsilon$	$B \rightarrow \varepsilon$

Possiamo dedurre che
non è una grammatica LL(1)

ESEMPIO

$G: S \rightarrow iEtS \mid iEtSeS \mid a$
 $E \rightarrow b$

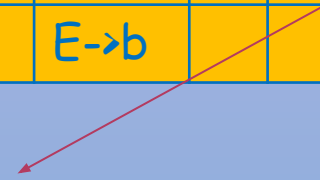


$S \rightarrow iEtSS' \mid a$
 $S' \rightarrow eS \mid \epsilon$
 $E \rightarrow b$

$\text{FIRST}(S) = \{i, a\}$
 $\text{FIRST}(S') = \{e, \epsilon\}$
 $\text{FIRST}(E) = \{b\}$

$\text{FOLLOW}(S) = \{e, \$\}$
 $\text{FOLLOW}(S') = \{e, \$\}$
 $\text{FOLLOW}(E) = \{t\}$

	i	a	b	t	e	\$
S	$S \rightarrow iEtSS'$	$S \rightarrow a$				
S'					$S' \rightarrow eS \quad S' \rightarrow \epsilon$	$S' \rightarrow \epsilon$
E			$E \rightarrow b$			



2 possibili scelte: la grammatica è ambigua

COME OTTENERE GRAMMATICHE LL(1)

- Verificare, se è possibile, che non sia ambigua. In caso contrario, **si tenti, se possibile, di rimuovere l'ambiguità;**
- Controllare che non presenti **ricorsioni sinistre**. In caso contrario **trasformarle in ricorsioni destre**.
- Se un simbolo non terminale ammette **più derivazioni con lo stesso prefisso** applicare **il metodo della fattorizzazione sinistra**.
- Dopo aver calcolato **First** e **Follow** per ogni simbolo, **costruire la tabella di parsing predittivo**. Se ogni casella della tabella ammette al più una derivazione allora la grammatica è LL(1)
- In alternativa, può essere necessario **allungare la lunghezza della prospezioni**. (Equivale a considerare parser LL(k), $k > 1$)

LIMITI DELLA FAMIGLIA LL(K)

Non tutti i linguaggi verificabili da parser deterministici sono generabili da grammatiche LL(k).

Per esempio:

$L = \{a^n b^n \mid n \geq 0\}$ linguaggio deterministico ma non LL(k)

È generato dalla grammatica

$S \rightarrow A \mid aS$

$A \rightarrow aAb \mid \epsilon$

Altro esempio:

$S \rightarrow R \mid (S)$

$R \rightarrow E=E$

$E \rightarrow a \mid (E+E)$

Definisce relazioni di uguaglianza tra espressioni aritmetiche additive. Una stringa che inizia con "(" può essere una relazione R parentesizzata oppure un'espressione.

Il linguaggio non è LL(k) ma è deterministico.

RELAZIONE FRA GRAMMATICHE $LL(K)$

E' possibile generalizzare la nozione di **FIRST** nel Parsing predittivo in modo da restituire i primi k token in input, e costruire una tabella predittiva in cui le righe sono i simboli nonterminali e le colonne sono sequenze di k terminali .

Le grammatiche non ambigue corrispondenti sono le $LL(K)$

