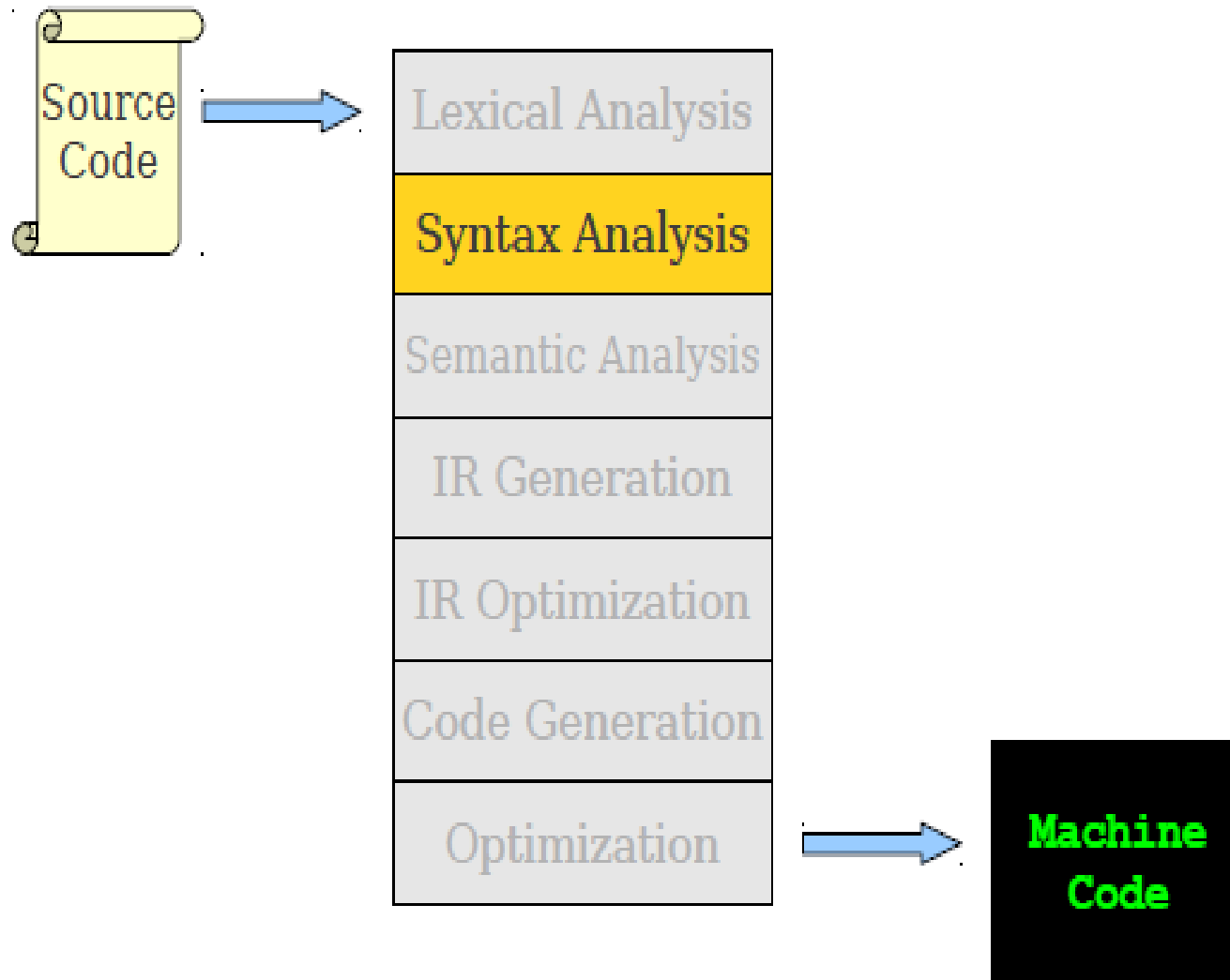


# Parser discendenti

Lunedì 12 Novembre

# Dove siamo?



# Input e output

- INPUT: Sequenza di token prodotti dall'analizzatore lessicale
- OUTPUT: Albero sintattico se la sequenza è generata dalla grammatica CF, altrimenti produce errori sintattici

# Esempio di CFG per un linguaggio di programmazione

```
BLOCK  → STMT  
        | { STMTS }  
  
STMTS → ε  
        | STMT STMTS  
  
STMT  → EXPR;  
        | if (EXPR) BLOCK  
        | while (EXPR) BLOCK  
        | do BLOCK while (EXPR);  
        | BLOCK  
        | ...  
  
EXPR   → identifier  
        | constant  
        | EXPR + EXPR  
        | EXPR - EXPR  
        | EXPR * EXPR  
        | ...
```

# Derivazione leftmost o rightmost

- Ad ogni passo si espande il simbolo non terminale più a sinistra, a differenza della rightmost in cui si espande il non terminale più a destra.

<b>BLOCK</b> → <b>STMT</b>   { <b>STMTS</b> }	
<b>STMTS</b> → $\epsilon$   <b>STMT STMTS</b>	<b>STMTS</b> ⇒ <b>STMT STMTS</b>
<b>STMT</b> → <b>EXPR;</b>   <b>if</b> ( <b>EXPR</b> ) <b>BLOCK</b>   <b>while</b> ( <b>EXPR</b> ) <b>BLOCK</b>   <b>do</b> <b>BLOCK</b> <b>while</b> ( <b>EXPR</b> );   <b>BLOCK</b>   ...	⇒ <b>EXPR; STMTS</b> ⇒ <b>EXPR = EXPR; STMTS</b> ⇒ <b>id = EXPR; STMTS</b> ⇒ <b>id = EXPR + EXPR; STMTS</b>
<b>EXPR</b> → <b>identifier</b>   <b>constant</b>   <b>EXPR + EXPR</b>   <b>EXPR - EXPR</b>   <b>EXPR * EXPR</b>   <b>EXPR = EXPR</b>   ...	⇒ <b>id = id + EXPR; STMTS</b> ⇒ <b>id = id + constant; STMTS</b> ⇒ <b>id = id + constant;</b>

# Altro esempio di CFG

$E \rightarrow \text{int}$

$E \rightarrow E \text{ Op } E$

$E \rightarrow (E)$

$\text{Op} \rightarrow +$

$\text{Op} \rightarrow -$

$\text{Op} \rightarrow *$

$\text{Op} \rightarrow /$

# Esempio

$E \rightarrow \text{int} \mid E \text{ Op } E \mid (E)$

$\text{Op} \rightarrow + \mid - \mid * \mid /$

$E$

$\Rightarrow E \text{ Op } E$

$\Rightarrow \text{int Op } E$

$\Rightarrow \text{int} * E$

$\Rightarrow \text{int} * (E)$

$\Rightarrow \text{int} * (E \text{ Op } E)$

$\Rightarrow \text{int} * (\text{int Op } E)$

$\Rightarrow \text{int} * (\text{int} + E)$

$\Rightarrow \text{int} * (\text{int} + \text{int})$

$E$

$\Rightarrow E \text{ Op } E$

$\Rightarrow E \text{ Op } (E)$

$\Rightarrow E \text{ Op } (E \text{ Op } E)$

$\Rightarrow E \text{ Op } (E \text{ Op } \text{int})$

$\Rightarrow E \text{ Op } (E + \text{int})$

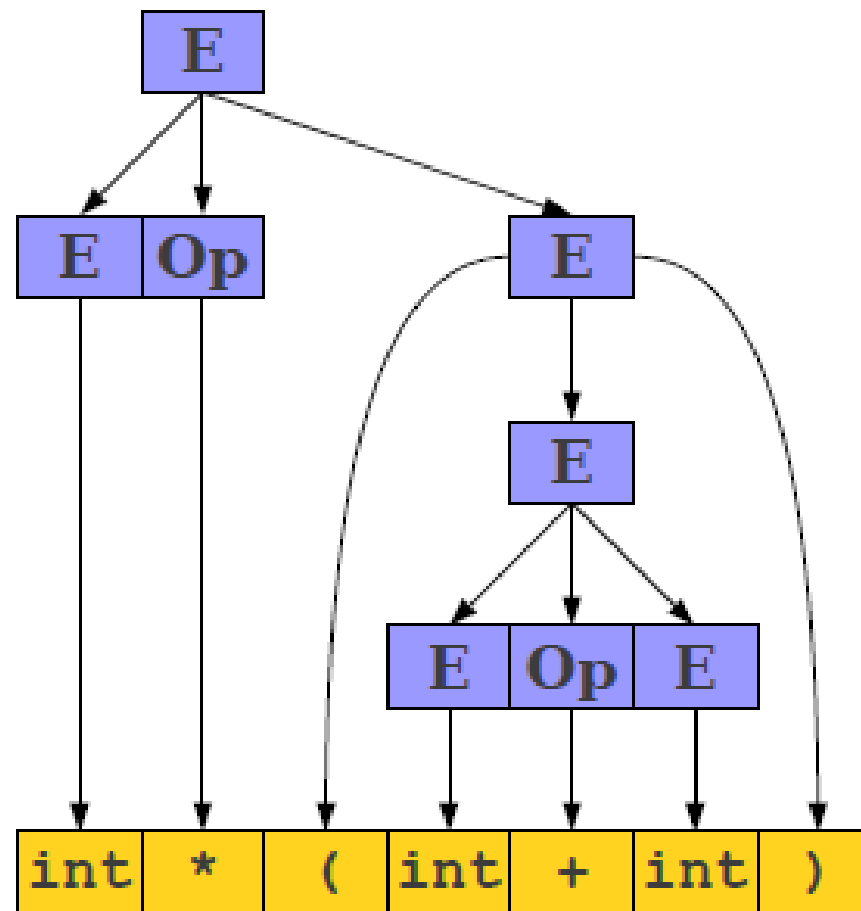
$\Rightarrow E \text{ Op } (\text{int} + \text{int})$

$\Rightarrow E * (\text{int} + \text{int})$

$\Rightarrow \text{int} * (\text{int} + \text{int})$

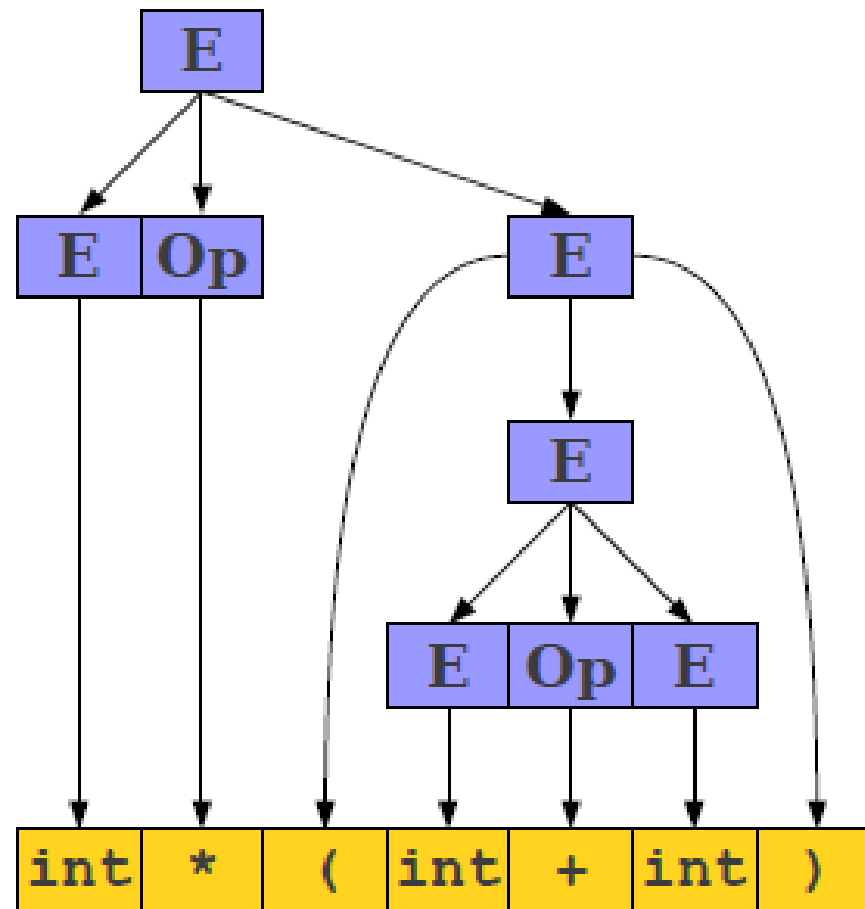
# Producono lo stesso syntax tree?

**E**  
⇒ **E Op E**  
⇒ **int Op E**  
⇒ **int \* E**  
⇒ **int \* (E)**  
⇒ **int \* (E Op E)**  
⇒ **int \* (int Op E)**  
⇒ **int \* (int + E)**  
⇒ **int \* (int + int)**



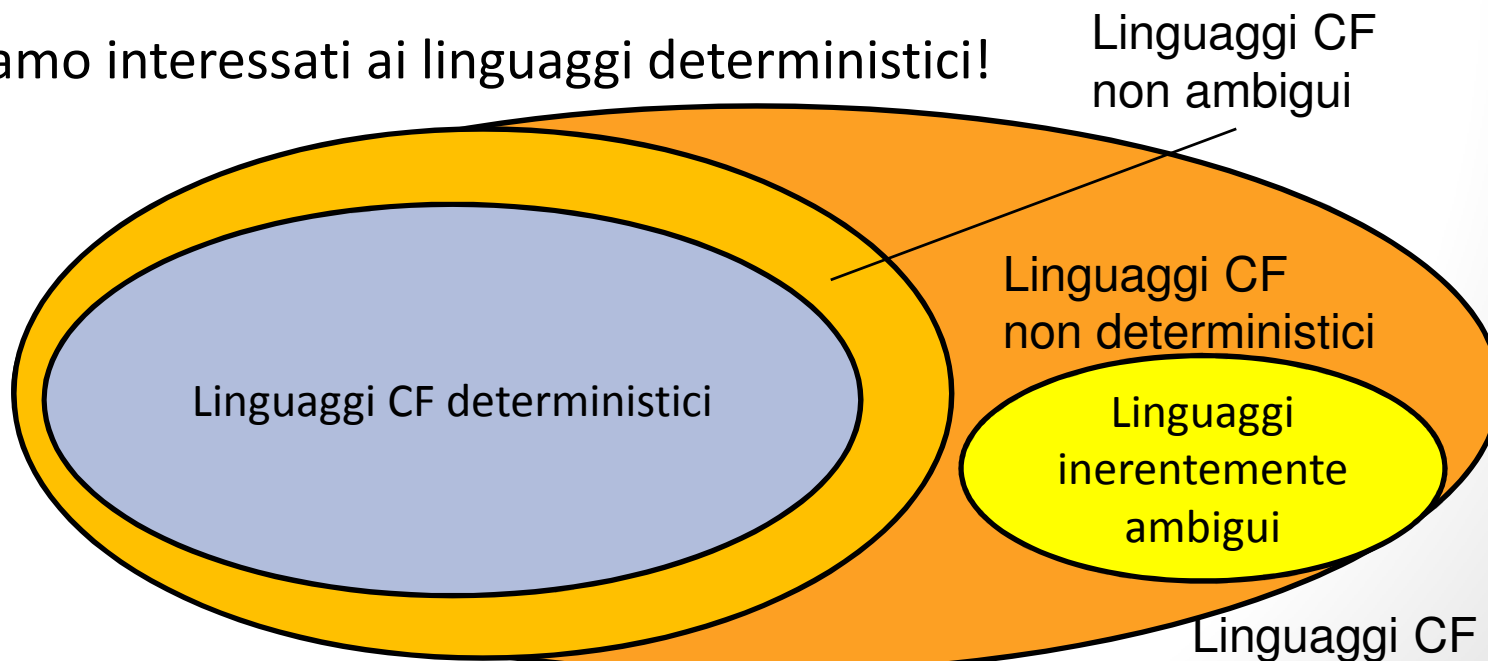


Si!

$$\begin{aligned} & E \\ \Rightarrow & E \text{ Op } E \\ \Rightarrow & E \text{ Op } (E) \\ \Rightarrow & E \text{ Op } (E \text{ Op } E) \\ \Rightarrow & E \text{ Op } (E \text{ Op } \text{int}) \\ \Rightarrow & E \text{ Op } (E + \text{int}) \\ \Rightarrow & E \text{ Op } (\text{int} + \text{int}) \\ \Rightarrow & E * (\text{int} + \text{int}) \\ \Rightarrow & \text{int} * (\text{int} + \text{int}) \end{aligned}$$


# Obiettivo del parser in un compilatore

- Costruire il syntax tree, ovvero quali produzioni vengono applicate piuttosto che l'ordine con cui si applicano.
- Se il linguaggio è non ambiguo, per ogni sequenza esiste un unico syntax tree
- Per l'insieme dei linguaggi non ambigui, il parser deve produrre un unico oggetto, ma potrebbe essere non deterministico.
- Siamo interessati ai linguaggi deterministici!



# Parser discendenti e ascendenti

- Considereremo due classi di parser:
  - Discendenti o TOP-DOWN: si costruisce la derivazione partendo dall'assioma; l'albero di derivazione si costruisce dalla radice alle foglie;
  - Ascendenti o BOTTOM-UP: si costruisce la derivazione ma nell'ordine riflesso, cioè l'albero si costruisce dalle foglie alla radice.

# Esercizio

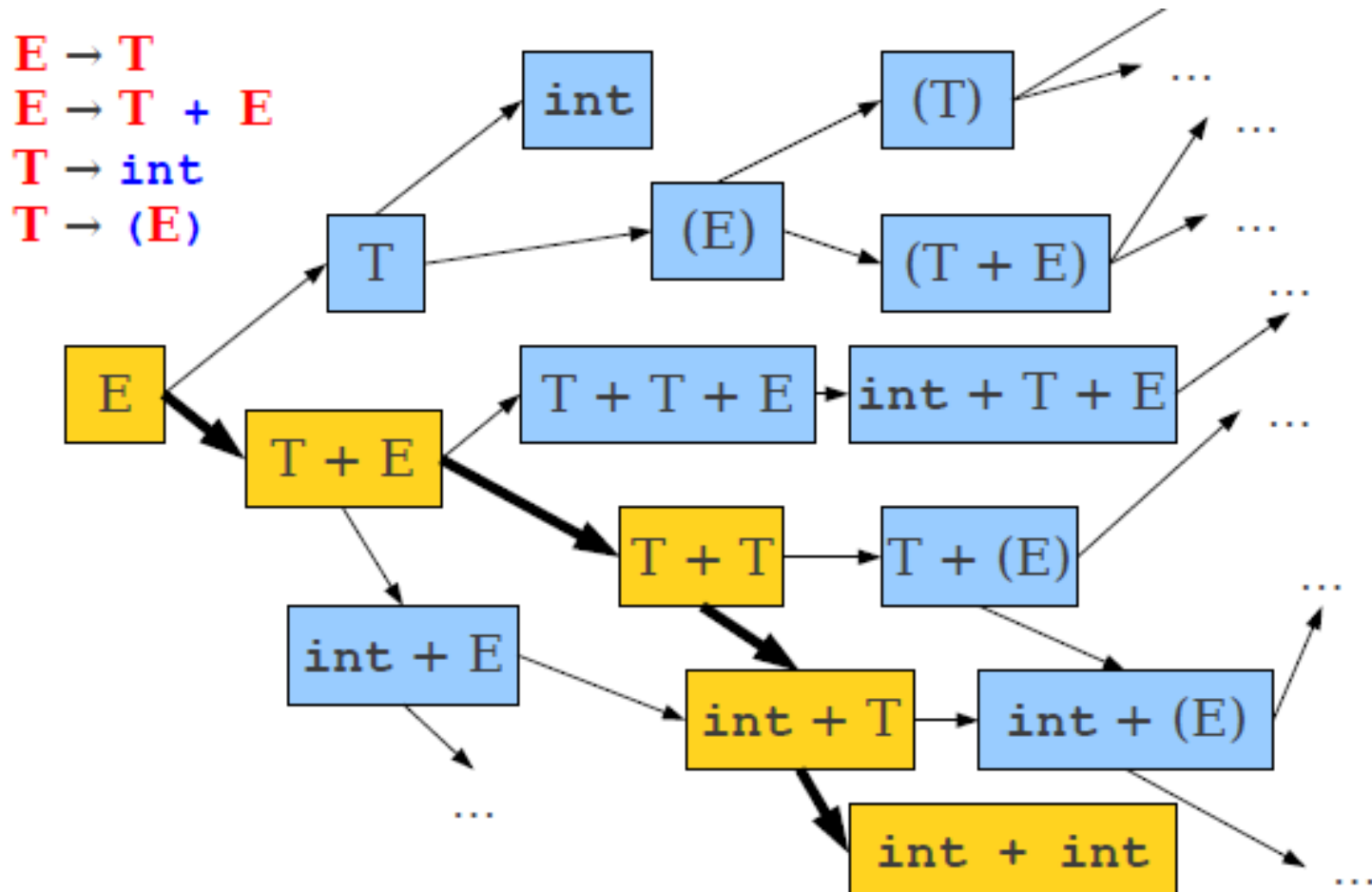
1.  $S \rightarrow aSAB$
2.  $S \rightarrow b$
3.  $A \rightarrow bA$
4.  $A \rightarrow a$
5.  $B \rightarrow cB$
6.  $B \rightarrow a$

La stringa  $a^2b^2a^4$  è generata dalla grammatica.  
Costruire l'albero sintattico.

# Parser top-down: problematiche

- I parser top-down iniziano il loro lavoro senza alcuna informazione iniziale.
- Cominciano con l'assioma, che va bene per tutti i programmi. Quale produzione applicare?
- Si può scommettere su una produzione, se il tentativo si rivela sbagliato si ritorna indietro e si ritenta (backtracking)
- Come si sceglie la produzione?

Il parsing top down è come  
ricercare un path in un grafo



# Due strategie di parser discendenti (o top down)

- PARSER a discesa ricorsiva (possono essere deterministici o non)
- PARSER LL(1) (parser deterministici)

# Parser a discesa ricorsiva

**L'idea: le regole grammaticali per un non-terminale A sono viste come costituenti una procedura che riconosce un A**

- la parte destra di ogni regola specifica la struttura del codice per questa procedura
- la sequenza di terminali e non terminali nelle varie regole corrisponde a un controllo che i terminali siano presenti nell'input e a invocazioni delle procedure dei simboli non terminali
- la presenza di diverse regole per A è modellata da **case** o **if**
- può richiedere backtracking (può richiedere di leggere più di una volta parte della stringa in ingresso, ovvero se l'applicazione di una produzione fallisce può tornare indietro).

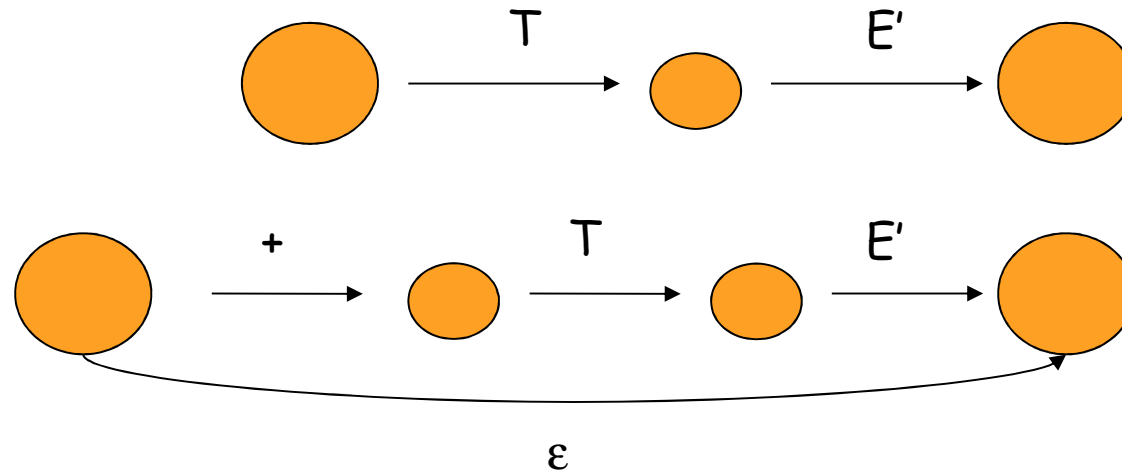


# Procedura tipica per un parser top-down a discesa ricorsiva

```
void A(){  
    scegli, per A, una produzione  $A \rightarrow X_1 X_2 \dots X_k$   
    for (da i a k) {  
        if ( $X_i$  è un non terminale)  
            richiama  $X_i()$ ;  
        else if ( $X_i$  è uguale al simbolo in input corrente)  
            procedi al simbolo successivo;  
        else si è verificato un errore;  
    }  
}
```

## Costruzione delle procedure ricorsive

$F \rightarrow (E) \text{ id}$   
 $E \rightarrow TE'$   
 $E' \rightarrow +TE'|\epsilon$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT'|\epsilon$



```
void E();
{
  T();
  E'();
}
```

```
void E'();
{
  If (tok==+) then
    {avanza input;
     T();
     E'();}
}
```

**esempio:**

**S** -> if **E** then **S** else **S**

**S** -> begin **S** **L**

**L** -> end

**L** -> ; **S** **L**

**S** -> print **E**

**E** -> num = num

token { IF, THEN,  
ELSE, BEGIN,  
PRINT, PUNTVIRG,  
NUM, EQ } ;

```
void S () { /* funzione per S */
```

```
if (tok == IF) then {  
  avanza(IF) ; E() ; avanza(THEN) ; S() ; avanza(ELSE) ; S() ;  
} else if (tok == BEGIN) {  
  avanza(BEGIN) ; S() ; L() ;  
} else if (tok == PRINT) {  
  avanza(PRINT) ; E() ;  
} else error() ;  
}
```

```
void L () { /* funzione per L */
```

```
if (tok == END) then  
  avanza(END) ;  
else if (tok == PUNTVIRG) {  
  avanza(PUNTVIRG) ; S() ; L() ;  
} else error() ;  
}
```

```
void E () { /* funzione per E */
```

```
  avanza(NUM) ; avanza(EQ) ; avanza(NUM) ;  
}
```

# Trasformazione della grammatica per l'analisi top-down

Due aspetti rendono una grammatica inadatta all'analisi top-down: la ricorsione sinistra e la presenza di prefissi comuni in più parti destre di regole associate allo stesso simbolo non terminale.

- $A \rightarrow y\alpha^1 \mid \dots \mid y\alpha^n$

*Rimedio: fattorizzazione sinistra*

- $A \rightarrow Aa$  ( $A$  non terminale).

*Rimedio: eliminazione ricorsione sinistra*

*Può essere possibile adattare una grammatica in modo che sia applicabile un parser top-down.*

# Fattorizzazione sinistra

Idea: quando non è chiaro quale produzione usare per espandere un non terminale  $A$ , si possono riscrivere le produzioni in modo da “posticipare” la scelta, introducendo un non terminale supplementare.

ES.:

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \gamma$$

*fattorizzazione sinistra:*

$$\begin{aligned} A &\rightarrow \alpha A' \mid \gamma & A' \text{ è un nuovo simbolo non terminale} \\ A' &\rightarrow \beta_1 \mid \beta_2 \end{aligned}$$

$$\begin{aligned} \langle \text{stmt} \rangle &\rightarrow \mathbf{if} \langle \text{expr} \rangle \mathbf{then} \langle \text{stmt} \rangle \\ &\mid \mathbf{if} \langle \text{expr} \rangle \mathbf{then} \langle \text{stmt} \rangle \mathbf{else} \langle \text{stmt} \rangle \end{aligned}$$

*fattorizzazione sinistra:*

$$\begin{aligned} \langle \text{stmt} \rangle &\rightarrow \mathbf{if} \langle \text{expr} \rangle \mathbf{then} \langle \text{stmt} \rangle \langle S \rangle \\ \langle S \rangle &\rightarrow \epsilon \mid \mathbf{else} \langle \text{stmt} \rangle \end{aligned}$$

# Eliminazione ricorsione sinistra

Metodo: date le produzioni "ricorsive sinistre" e non, di un non terminale  $A$ :

$$A \rightarrow A \alpha_1 \mid \dots \mid A \alpha_n \mid \beta_1 \mid \dots \mid \beta_n$$

si sostituiscono con (eliminazione ricorsione sinistra immediata)

$$A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \varepsilon$$

Es. espressioni aritmetiche

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

$$E \rightarrow TE'$$

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

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

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

$$T \rightarrow FT'$$

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

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

NOTA: Non è detto che le ricorsioni sinistre siano solo immediate

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \varepsilon$$

# Algoritmo per eliminare le ricorsioni sinistre

- INPUT: Grammatica senza cicli né  $\epsilon$ -produzioni
- OUTPUT: Una grammatica equivalente senza ricorsioni sinistre (Tale grammatica potrebbe contenere  $\epsilon$ -produzioni)

Algoritmo:

1. Ordina i simboli non terminali  $A_1, A_2, \dots, A_n$ ;
2. For ( $i=1$  to  $n$ ) {  
    for ( $j=1$  to  $i-1$ ) {  
        sostituire ogni produzione  $A_i \rightarrow A_j \gamma$  con le produzioni  
         $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$  dove  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  sono tutte le  $A_j$   
        produzioni  
    }  
    eliminare le ricorsioni sinistre immediate  
};

Nell'esempio:  $S \rightarrow Aa \mid b$   
 $A \rightarrow Ac \mid Sd \mid \epsilon$

Si ottiene:  $S \rightarrow Aa \mid b$   
 $A \rightarrow bdA' \mid A'$   
 $A' \rightarrow cA' \mid ad A' \mid \epsilon$

# Parser predittivi

- Basandosi sull'input che resta da leggere, predice quale produzione usare senza fare uso di backtracking.
- Nella tecnica “a discesa ricorsiva” l'analisi sintattica viene effettuata attraverso una cascata di chiamate ricorsive.
- I parser discendenti deterministici possono anche essere guidati da una tabella. In tal caso si parla di parser predittivi.
- Nel parser LL(1), particolari parser predittivi, la pila delle chiamate ricorsive viene esplicitata nel parser, e quindi non si fa più uso di ricorsione



# 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

