



UN COMPITO IMPORTANTE DI UN PARSER

Gestione degli errori sintattici

Giovedì 7 Novembre

GESTIONE DEGLI ERRORI IN UN PARSER

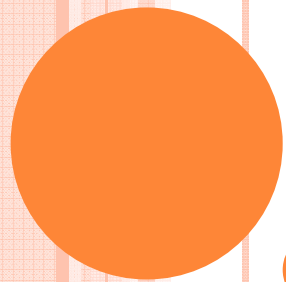
- Un parser deve essere in grado di *scoprire*, *diagnosticare* e *correggere* gli errori in maniera efficiente, per *riprendere* l'analisi e scoprire nuovi errori.
- Alcuni parser (LL e LR) hanno la proprietà "viable prefix": sono in grado di rilevare un errore non appena si presenta perché sono in grado di riconoscere i prefissi validi del linguaggio



STRATEGIE DI RIPARAZIONE

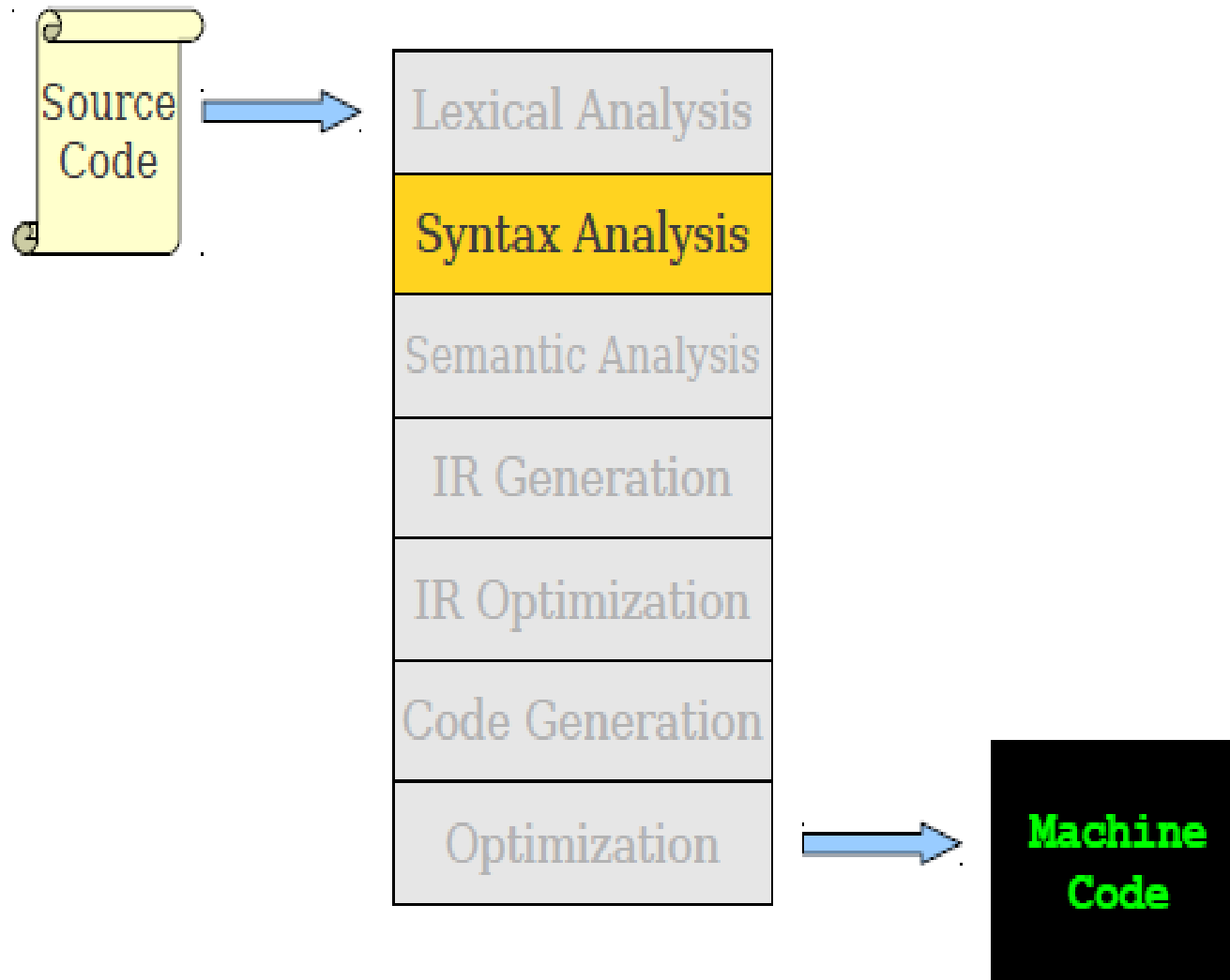
- “panic mode”: scoperto l'errore il parser riprende l'analisi in corrispondenza di alcuni token sincronizzanti predefiniti (es.: delimitatori begin end) *scartando* alcuni caratteri. Svantaggi: può essere *scartato* molto input.
- “phrase level”: correzioni locali ottenute inserendo, modificando, cancellando alcuni terminali per poter riprendere l'analisi (es.: ',' -> ';')
Svantaggi: difficoltà quando la distanza dall'errore è notevole.
- “error productions”: uso di produzioni che estendono la grammatica per generare gli errori più comuni. Metodo efficiente per la diagnostica.
- “global correction”: si cerca di “calcolare” la migliore correzione possibile alla derivazione errata (minimo costo di interventi per inserzioni/cancellazioni). Metodo globale poco usato in pratica, ma tecnica usata per ottimizzare la strategia “phrase level”.





PARSER DISCENDENTI

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.

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**
| **EXPR = EXPR**
| ...

STMTS

⇒ **STMT STMTS**

⇒ **EXPR; STMTS**

⇒ **EXPR = EXPR; STMTS**

⇒ **id = EXPR; STMTS**

⇒ **id = EXPR + EXPR; STMTS**

⇒ **id = id + EXPR; STMTS**

⇒ **id = id + constant; STMTS**

⇒ **id = id + constant;**



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

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

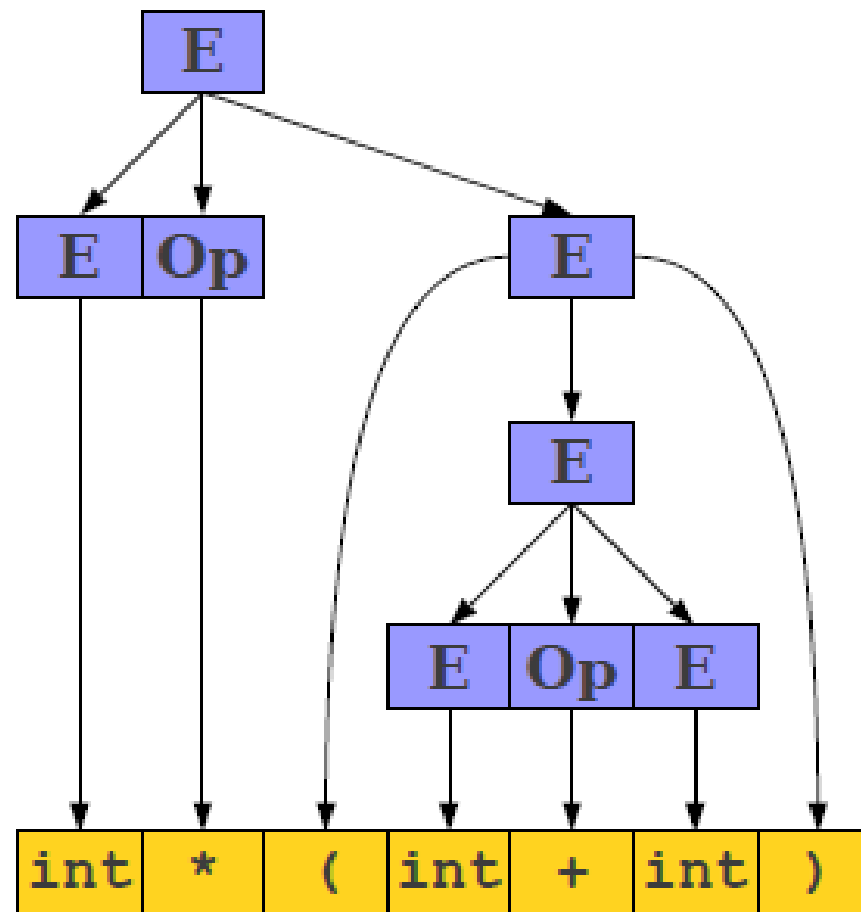
$\text{Op} \rightarrow + \mid - \mid * \mid /$

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?

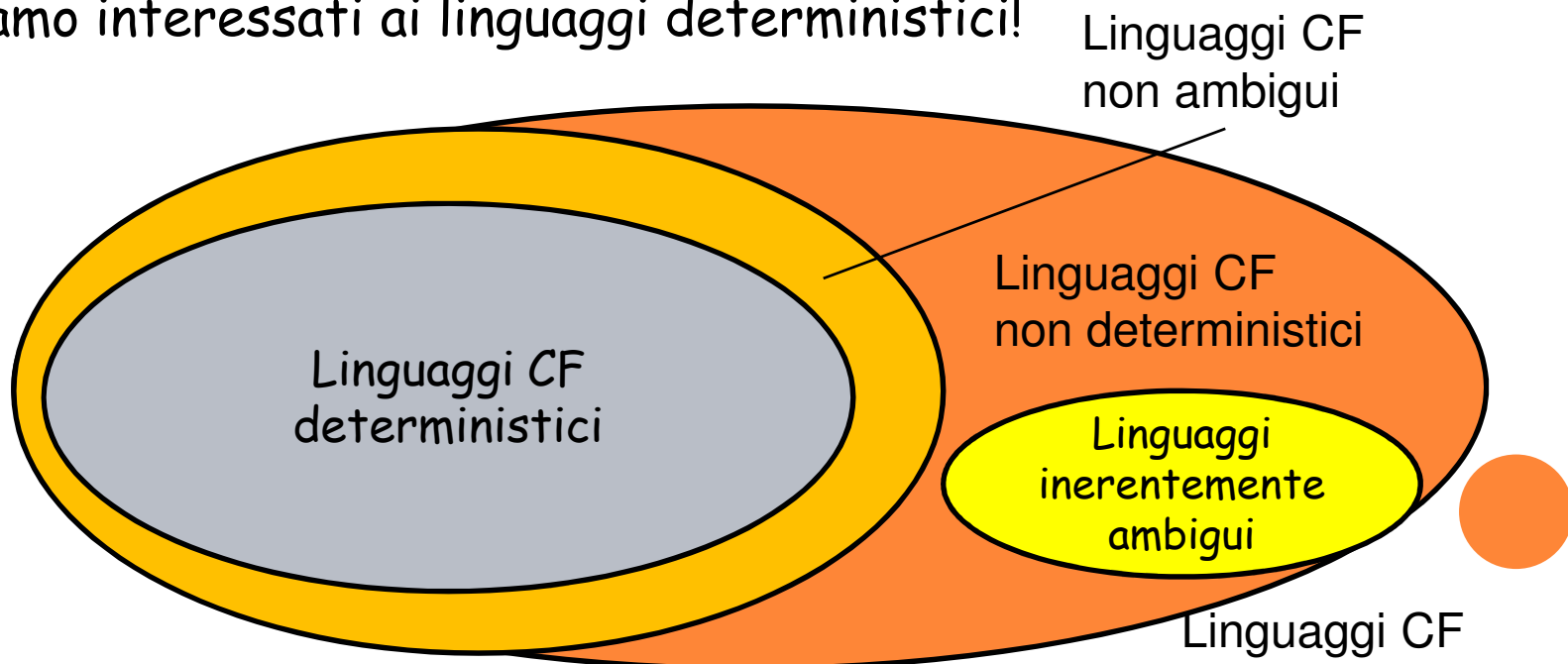
E
⇒ **E Op E**
⇒ **int Op E**
⇒ **int * E**
⇒ **int * (E)**
⇒ **int * (E Op E)**
⇒ **int * (int Op E)**
⇒ **int * (int + E)**
⇒ **int * (int + int)**



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


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.

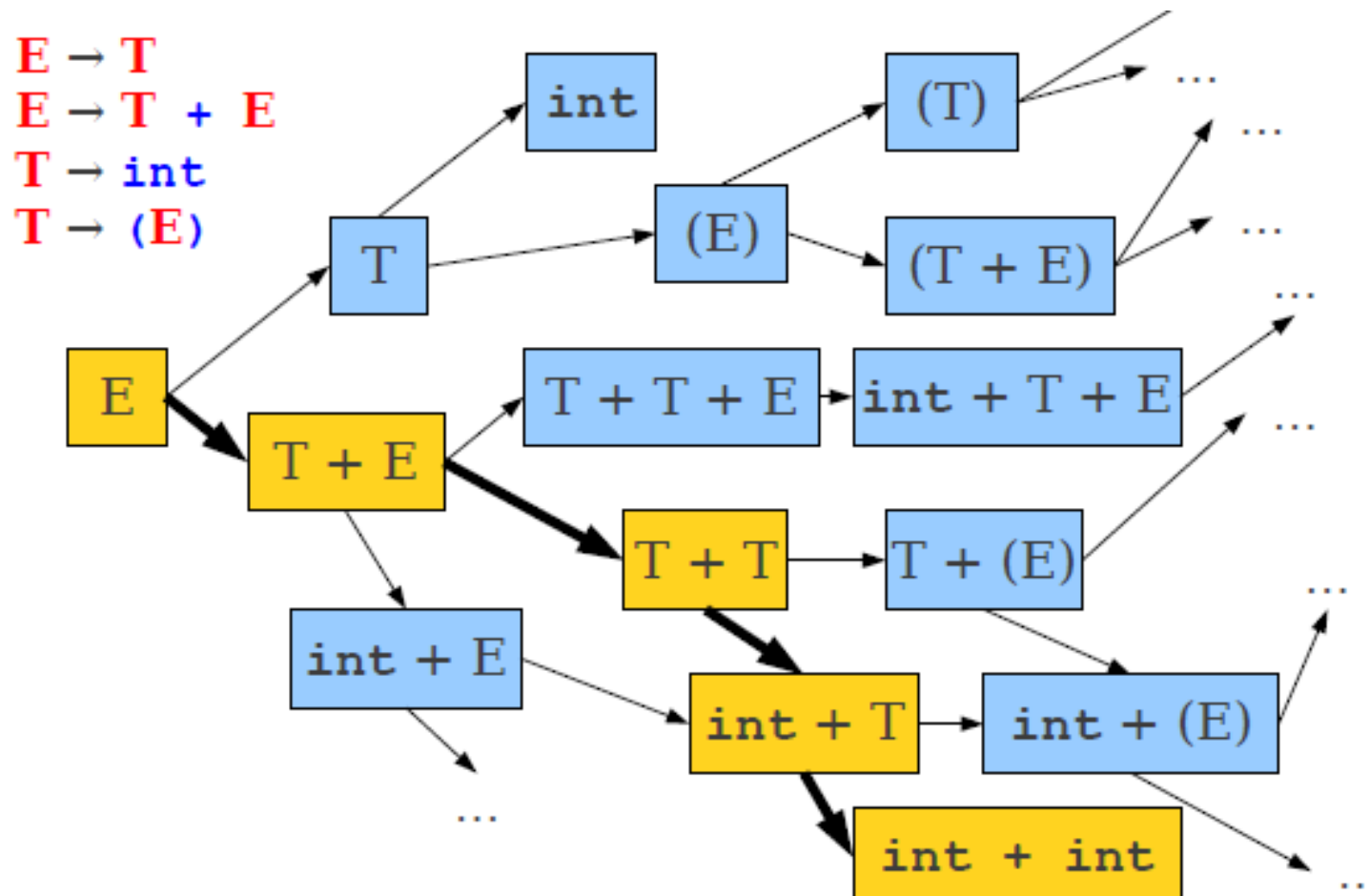


PARSER TOP-DOWN: PROBLEMATICHE

- 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).



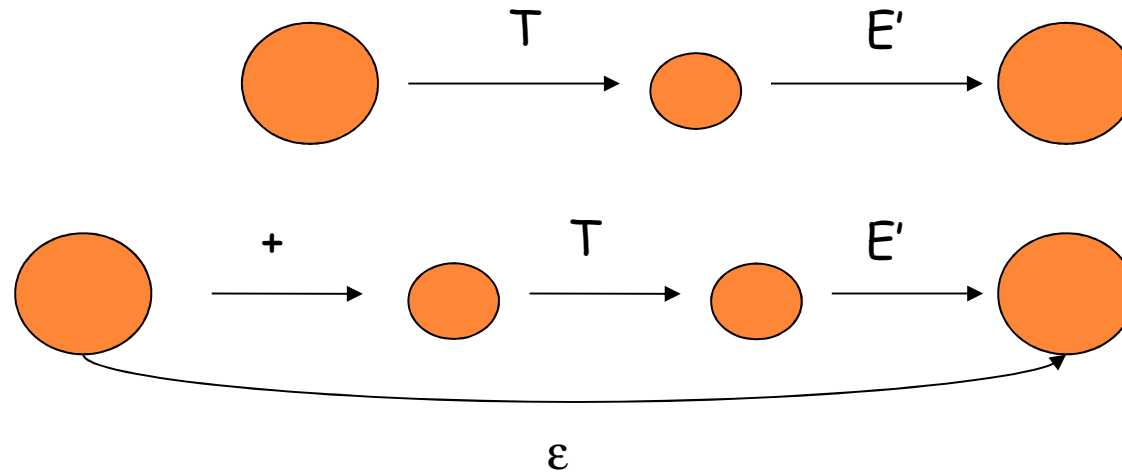
PROCEDURA TIPICA PER UN PARSER TOP-DOWN A DISCESA RICORSIVA

```
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

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



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

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

esempio:

$S \rightarrow \text{if } E \text{ then } S \text{ else } S$

$S \rightarrow \text{begin } S L$

$L \rightarrow \text{end}$

$L \rightarrow ; S L$

$S \rightarrow \text{print } E$

$E \rightarrow \text{num} = \text{num}$

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. Esiste un metodo più generale

$$S \rightarrow Aa \mid b$$

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



ALGORITMO PER ELIMINARE LE RICORSIONI SINISTRE

- INPUT: Grammatica senza cicli né ε -produzioni
- OUTPUT: Una grammatica equivalente senza ricorsioni sinistre (Tale grammatica potrebbe contenere ε -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: $S \rightarrow Aa \mid b$
 $A \rightarrow Ac \mid Sd \mid \varepsilon$

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



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

