

# Linguaggi Formali e Traduttori

## 4.1 Parsing top-down e grammatiche LL(1).

- Sommario
- Strategia per il parsing top-down
- Stringhe annullabili (NULL)
- Esempi di stringhe annullabili
- Inizi di una stringa (FIRST)
- Come calcolare FIRST
- Esempi di calcolo di FIRST
- Seguiti di una variabile (FOLLOW)
- Come calcolare FOLLOW
- Esempi di calcolo di FOLLOW
- Insiemi guida
- Grammatiche LL(1)
- Esempio: espressioni aritmetiche
- Esercizi

È proibito condividere e divulgare in qualsiasi forma i materiali didattici caricati sulla piattaforma e le lezioni svolte in videoconferenza: ogni azione che viola questa norma sarà denunciata agli organi di Ateneo e perseguita a termini di legge.

# Sommario

## Problema

- Data una grammatica  $G = (V, T, P, S)$  e una stringa  $w \in T^*$ , determinare se

$$S \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow w$$

o, equivalentemente, se esiste un albero sintattico di  $G$  con radice  $S$  e prodotto  $w$ .

- La **costruzione dell'automa** corrispondente a  $G$  produce un PDA non deterministico.
- Per alcune  $G$  sappiamo che **non è possibile trovare un DPDA**.

## In questa lezione

- Identifichiamo una famiglia di grammatiche libere per le quali è possibile costruire riconoscitori (parser) deterministici, cioè che non fanno uso di backtracking.
- Questi parser sono detti **top-down** perché costruiscono l'albero sintattico di  $w$  dalla radice (top) verso le foglie (down) o, equivalentemente, cercano una derivazione sinistra per  $w$ .

# Strategia per il parsing top-down

Data una grammatica  $G = (V, T, P, S)$  e una stringa  $w \in T^*$ , il parser cerca di ottenere una derivazione a sinistra  $S \Rightarrow_{lm}^* w$  in cui, al passo  $i$ , il parser sa che

$$S \Rightarrow_{lm}^* uA\beta$$

e deve stabilire se

$$uA\beta \Rightarrow_{lm}^* w$$

Ci sono due casi da considerare:

- Se  $u$  non è prefisso di  $w$ , allora il parser **rifiuta**  $w$ .
- Se  $w = uav$ , allora il parser deve **scegliere** una produzione per riscrivere  $A$

$$A \rightarrow \alpha_1 \mid \cdots \mid \alpha_n$$

e per farlo può usare  $a$  come “guida”, a patto che tale simbolo identifichi univocamente l' $\alpha_i$  tale che  $\alpha_i\beta \Rightarrow_{lm}^* av$ .

Per ogni produzione  $A \rightarrow \alpha_i$  occorre saper calcolare gli insiemi di simboli terminali che possono **iniziare** le stringhe derivate da  $\alpha_i\beta$  e richiedere che tali insiemi siano disgiunti.

# Stringhe annullabili (NULL)

## Definizione

Data una grammatica  $G = (V, T, P, S)$ , diciamo che  $\alpha \in (V \cup T)^*$  è **annullabile**, e scriviamo  $\text{NULL}(\alpha)$ , se e solo se  $\alpha \Rightarrow_G^* \varepsilon$ , ovvero se  $\alpha$  può essere riscritta nella stringa vuota.

## Come determinare se una stringa è annullabile

- (1) Se  $\text{NULL}(X_1), \dots, \text{NULL}(X_n)$ , allora  $\text{NULL}(X_1 \cdots X_n)$ .
- (2) Se esiste una produzione  $A \rightarrow \alpha \in P$  e  $\text{NULL}(\alpha)$ , allora  $\text{NULL}(A)$ .

## Note

- Come caso particolare di (1) quando  $n = 0$  abbiamo  $\text{NULL}(\varepsilon)$ .
- Combinando (1) e (2) abbiamo che  $A \rightarrow \varepsilon \in P$  implica  $\text{NULL}(A)$ .
- Una stringa che contiene simboli terminali non è mai annullabile.

# Esempi di stringhe annullabili

$$A \rightarrow a \mid Bc$$

$$B \rightarrow \varepsilon \mid bB$$

$$C \rightarrow d \mid Cc \mid BB$$

# Esempi di stringhe annullabili

$$A \rightarrow a \mid Bc$$

$$B \rightarrow \varepsilon \mid bB$$

$$C \rightarrow d \mid Cc \mid BB$$

- Da  $\text{NULL}(\varepsilon)$  e dalla produzione  $B \rightarrow \varepsilon$  deduciamo  $\text{NULL}(B)$ .
- Da  $\text{NULL}(B)$  e dalla produzione  $C \rightarrow BB$  deduciamo  $\text{NULL}(C)$ .
- Da  $\text{NULL}(B)$  e  $\text{NULL}(C)$  deduciamo  $\text{NULL}(BC)$ .
- Da  $\neg\text{NULL}(a)$  e  $\neg\text{NULL}(Bc)$  deduciamo  $\neg\text{NULL}(A)$ .

# Inizi di una stringa (FIRST)

## Definizione

Data una grammatica  $G = (V, T, P, S)$  e una stringa  $\alpha \in (V \cup T)^*$ , indichiamo con  $\mathbf{FIRST}(\alpha)$  gli **inizi** di  $\alpha$ , ovvero l'insieme dei simboli terminali che possono trovarsi all'inizio delle stringhe derivate da  $\alpha$ . Formalmente:

$$\mathbf{FIRST}(\alpha) \stackrel{\text{def}}{=} \{a \in T \mid \alpha \Rightarrow_G^* a\beta\}$$

## Attenzione

Il libro di testo usa un'unica funzione  $\mathbf{FIRST}_{\text{libro}}$  che combina  $\mathbf{NULL}$  e  $\mathbf{FIRST}$  così:

$$\mathbf{FIRST}_{\text{libro}}(\alpha) = \begin{cases} \mathbf{FIRST}(\alpha) \cup \{\varepsilon\} & \text{se } \mathbf{NULL}(\alpha) \\ \mathbf{FIRST}(\alpha) & \text{altrimenti} \end{cases}$$

In pratica, l'approccio seguito dal libro ammette il simbolo speciale  $\varepsilon$  tra gli inizi di  $\alpha$  per indicare il fatto che  $\alpha$  è annullabile. Noi abbiamo definito un predicato  $\mathbf{NULL}(\alpha)$  apposito mentre  $\mathbf{FIRST}(\alpha)$  contiene solo simboli terminali.

# Come calcolare FIRST

È possibile calcolare  $\text{FIRST}(\alpha)$  per induzione su  $\alpha$ , usando le seguenti regole:

$$\begin{aligned}\text{FIRST}(\varepsilon) &= \emptyset \\ \text{FIRST}(a) &= \{a\} \\ \text{FIRST}(A) &= \bigcup_{A \rightarrow \alpha} \text{FIRST}(\alpha) \\ \text{FIRST}(X\alpha) &= \begin{cases} \text{FIRST}(X) \cup \text{FIRST}(\alpha) & \text{se } \text{NULL}(X) \\ \text{FIRST}(X) & \text{altrimenti} \end{cases}\end{aligned}$$

## Attenzione

Applicando le regole qui sopra, può capitare di arrivare a equazioni della forma

$$\text{FIRST}(A) = \text{FIRST}(A) \cup \mathcal{S}$$

dove  $\mathcal{S}$  è un insieme di terminali. Questa equazione si può semplificare a

$$\text{FIRST}(A) = \mathcal{S}$$

in quanto siamo interessati a ottenere il più piccolo insieme di terminali con la proprietà descritta nella [slide precedente](#).



# Esempi di calcolo di FIRST

$$S \rightarrow Ac \mid Ba$$

$$A \rightarrow \varepsilon \mid a$$

$$B \rightarrow b$$

$$C \rightarrow a \mid Cb$$

$$D \rightarrow \varepsilon \mid d \mid Db$$

# Esempi di calcolo di FIRST

$$S \rightarrow Ac \mid Ba$$

$$A \rightarrow \varepsilon \mid a$$

$$B \rightarrow b$$

$$C \rightarrow a \mid Cb$$

$$D \rightarrow \varepsilon \mid d \mid Db$$

## Variabili annullabili

- $\text{NULL}(A)$
- $\text{NULL}(D)$

## Calcolo di FIRST di tutte le variabili

- $\text{FIRST}(B) = \text{FIRST}(b) = \{b\}$
- $\text{FIRST}(A) = \text{FIRST}(\varepsilon) \cup \text{FIRST}(a) = \{a\}$
- $\text{FIRST}(S) = \text{FIRST}(Ac) \cup \text{FIRST}(Ba) = \text{FIRST}(A) \cup \text{FIRST}(c) \cup \text{FIRST}(B) = \{a, b, c\}$
- $\text{FIRST}(C) = \text{FIRST}(a) \cup \text{FIRST}(Cb) = \{a\} \cup \text{FIRST}(C) = \{a\}$
- $\text{FIRST}(D) = \text{FIRST}(\varepsilon) \cup \text{FIRST}(d) \cup \text{FIRST}(Db) = \{d\} \cup \text{FIRST}(D) \cup \text{FIRST}(b) = \{b, d\}$

# Seguiti di una variabile (FOLLOW)

## Definizione

Data una grammatica  $G = (V, T, P, S)$  e una variabile  $A \in V$ , indichiamo con  $\text{FOLLOW}(A)$  i **seguiti** di  $A$ , ovvero l'insieme dei simboli terminali che possono seguire  $A$  in una forma sentenziale. Formalmente:

$$\text{FOLLOW}(A) \stackrel{\text{def}}{=} \{a \in T \mid S \Rightarrow_G^* \alpha A a \beta\}$$

## Attenzione

- Per convenzione aggiungeremo una sentinella  $\$$  ai seguiti del simbolo iniziale  $S$ .
- In questo modo il parser può capire quando è arrivato alla fine della stringa da riconoscere.

# Come calcolare FOLLOW

Il calcolo di **FOLLOW** si effettua in due fasi.

## Fase 1

In questa fase si annotano relazioni di appartenenza ed inclusione insiemistica secondo il seguente algoritmo:

- Annotare  $\$ \in \text{FOLLOW}(S)$ .
- Ripetere i passi seguenti per ogni produzione e per ogni variabile nel corpo di queste:
  1. Se  $A \rightarrow \alpha B \beta$ , allora annotare  $\text{FIRST}(\beta) \subseteq \text{FOLLOW}(B)$ .
  2. Se  $A \rightarrow \alpha B \beta$  e  $\text{NULL}(\beta)$ , allora annotare  $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$ .

Caso particolare di (2): se  $A \rightarrow \alpha B$ , allora annotare  $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$ .

## Fase 2

Si determinano i seguiti propagando i simboli terminali (e  $\$$ ) rispettando l'ordine delle inclusioni insiemistiche  $\subseteq$  che sono state annotate.

Per grammatiche complesse può essere utile fare una tabella con due colonne, l'elenco di tutte le variabili nella prima ed i seguiti corrispondenti alle variabili nella seconda.

# Esempi di calcolo di FOLLOW

$$S \rightarrow Ac \mid Ba$$

$$A \rightarrow \varepsilon \mid a$$

$$B \rightarrow b$$

$$C \rightarrow a \mid Cb$$

$$D \rightarrow \varepsilon \mid d \mid Db$$

# Esempi di calcolo di FOLLOW

$$S \rightarrow Ac \mid Ba$$

$$A \rightarrow \varepsilon \mid a$$

$$B \rightarrow b$$

$$C \rightarrow a \mid Cb$$

$$D \rightarrow \varepsilon \mid d \mid Db$$

## Fase 1

- $\$ \in \text{FOLLOW}(S)$
- $\text{FIRST}(c) \subseteq \text{FOLLOW}(A)$
- $\text{FIRST}(a) \subseteq \text{FOLLOW}(B)$
- $\text{FIRST}(b) \subseteq \text{FOLLOW}(C)$
- $\text{FIRST}(b) \subseteq \text{FOLLOW}(D)$

## Fase 2

$X$	$\text{FOLLOW}(X)$
$S$	$\{\$ \}$
$A$	$\{c\}$
$B$	$\{a\}$
$C$	$\{b\}$
$D$	$\{b\}$

# Insiemi guida

## Definizione

Data una grammatica  $G = (V, T, P, S)$  e una produzione  $A \rightarrow \alpha$ , indichiamo con  $\text{GUIDA}(A \rightarrow \alpha)$  l'**insieme guida** di  $A \rightarrow \alpha$ , ovvero l'insieme

$$\text{GUIDA}(A \rightarrow \alpha) \stackrel{\text{def}}{=} \begin{cases} \text{FIRST}(\alpha) \cup \text{FOLLOW}(A) & \text{se } \text{NULL}(\alpha) \\ \text{FIRST}(\alpha) & \text{altrimenti} \end{cases}$$

## Intuizione

Un parser predittivo che sceglie di riscrivere la variabile  $A$  usando la produzione  $A \rightarrow \alpha$  si aspetta di leggere nella stringa di input uno dei simboli nell'insieme guida di  $A \rightarrow \alpha$ .

Sono due i casi da considerare:

1. Il simbolo è uno degli inizi di  $\alpha$ , oppure
2.  $\alpha$  è annullabile ed il simbolo è uno dei seguiti di  $A$ .

# Grammatiche LL(1)

## Definizione

Diciamo che una grammatica  $G = (V, T, P, S)$  è LL(1) se, per ogni coppia di produzioni distinte  $A \rightarrow \alpha$  e  $A \rightarrow \beta$  in  $P$ , abbiamo che

$$\text{GUIDA}(A \rightarrow \alpha) \cap \text{GUIDA}(A \rightarrow \beta) = \emptyset$$

## Intuizione

Noto il simbolo da riscrivere  $A$ , note le produzioni  $A \rightarrow \beta_1 \mid \dots \mid \beta_n$  e noto il prossimo simbolo terminale  $a$  nella stringa da riconoscere, in una grammatica LL(1) esiste al massimo una produzione “giusta” tale che  $a \in \text{GUIDA}(A \rightarrow \beta_i)$  dunque il parser predittivo identifica univocamente la produzione  $A \rightarrow \beta_i$  a partire da  $a$ .

## Cosa c'è nel nome LL(1)

- L  $\rightarrow$  la stringa in input viene analizzata da sinistra (left) a destra;
- L  $\rightarrow$  il parser cerca di costruire una derivazione canonica sinistra (leftmost);
- 1  $\rightarrow$  il parser usa un solo simbolo terminale della stringa per scegliere la produzione.



# Esempio: espressioni aritmetiche

$$\begin{aligned}E &\rightarrow TE' \\E' &\rightarrow +TE' \mid \varepsilon \quad \text{NULL}(E') \\T &\rightarrow FT' \\T' &\rightarrow *FT' \mid \varepsilon \quad \text{NULL}(T') \\F &\rightarrow (E) \mid \text{id}\end{aligned}$$

- $\$ \in \text{FOLLOW}(E)$
- $\{+\} = \text{FIRST}(E') \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(E')$
- $\text{FOLLOW}(E') \subseteq \text{FOLLOW}(T)$
- $\{*\} = \text{FIRST}(T') \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(T')$
- $\text{FOLLOW}(T') \subseteq \text{FOLLOW}(F)$
- $\{)\} = \text{FIRST}()) \subseteq \text{FOLLOW}(E)$

# Esempio: espressioni aritmetiche

$$\begin{aligned}E &\rightarrow TE' \\E' &\rightarrow +TE' \mid \varepsilon \quad \text{NULL}(E') \\T &\rightarrow FT' \\T' &\rightarrow *FT' \mid \varepsilon \quad \text{NULL}(T') \\F &\rightarrow (E) \mid \text{id}\end{aligned}$$

$$\begin{aligned}\text{FIRST}(E) &= \text{FIRST}(T) = \{ (, \text{id} \} \\ \text{FIRST}(E') &= \{ + \} \\ \text{FIRST}(T) &= \text{FIRST}(F) = \{ (, \text{id} \} \\ \text{FIRST}(T') &= \{ * \} \\ \text{FIRST}(F) &= \{ (, \text{id} \}\end{aligned}$$

- $\$ \in \text{FOLLOW}(E)$
- $\{ + \} = \text{FIRST}(E') \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(E')$
- $\text{FOLLOW}(E') \subseteq \text{FOLLOW}(T)$
- $\{ * \} = \text{FIRST}(T') \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(T')$
- $\text{FOLLOW}(T') \subseteq \text{FOLLOW}(F)$
- $\{ ) \} = \text{FIRST}()) \subseteq \text{FOLLOW}(E)$

# Esempio: espressioni aritmetiche

$$\begin{aligned}E &\rightarrow TE' \\E' &\rightarrow +TE' \mid \varepsilon \quad \text{NULL}(E') \\T &\rightarrow FT' \\T' &\rightarrow *FT' \mid \varepsilon \quad \text{NULL}(T') \\F &\rightarrow (E) \mid \text{id}\end{aligned}$$

$$\begin{aligned}\text{FIRST}(E) &= \text{FIRST}(T) = \{ (, \text{id} \} \\ \text{FIRST}(E') &= \{ + \} \\ \text{FIRST}(T) &= \text{FIRST}(F) = \{ (, \text{id} \} \\ \text{FIRST}(T') &= \{ * \} \\ \text{FIRST}(F) &= \{ (, \text{id} \}\end{aligned}$$

- $\$ \in \text{FOLLOW}(E)$
- $\{ + \} = \text{FIRST}(E') \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(T)$
- $\text{FOLLOW}(E) \subseteq \text{FOLLOW}(E')$
- $\text{FOLLOW}(E') \subseteq \text{FOLLOW}(T)$
- $\{ * \} = \text{FIRST}(T') \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(F)$
- $\text{FOLLOW}(T) \subseteq \text{FOLLOW}(T')$
- $\text{FOLLOW}(T') \subseteq \text{FOLLOW}(F)$
- $\{ ) \} = \text{FIRST}()) \subseteq \text{FOLLOW}(E)$

$X$	$\text{FOLLOW}(X)$
$E$	$\$, )$
$E'$	$\$, )$
$T$	$\$, ), +$
$T'$	$\$, ), +$
$F$	$\$, ), +, *$

# Esercizi

1. Calcolare gli insiemi guida della grammatica nella [slide 14](#). La grammatica è LL(1)?
2. Calcolare gli insiemi guida della seguente grammatica e determinare se è LL(1).

$$A \rightarrow BC \mid D$$

$$B \rightarrow \varepsilon \mid a$$

$$C \rightarrow b \mid cCc$$

$$D \rightarrow \varepsilon \mid CD$$

3. Ripetere l'esercizio precedente per la grammatica

$$S \rightarrow \text{if } E \text{ then } SS' \text{ fi} \mid \text{skip}$$

$$S' \rightarrow \text{else } S \mid \varepsilon$$

$$E \rightarrow \text{true} \mid \text{false}$$

in cui  $S$ ,  $S'$  ed  $E$  sono variabili e **if**, **then**, ... sono terminali.

4. Ripetere l'esercizio precedente dopo aver rimosso il terminale **fi** dalla grammatica.

# Linguaggi Formali e Traduttori

## 4.2 Parsing ricorsivo discendente

- Sommario
- Struttura del parser ricorsivo
- Algoritmo di parsing ricorsivo
- Implementazione Java del parser (classe base)
- Esempio: parser per il linguaggio  $a^n b^n$
- Esercizi

È proibito condividere e divulgare in qualsiasi forma i materiali didattici caricati sulla piattaforma e le lezioni svolte in videoconferenza: ogni azione che viola questa norma sarà denunciata agli organi di Ateneo e perseguita a termini di legge.

# Sommario

## Problema

- Realizzazione pratica di un parser top-down.

## In questa lezione

- Studiamo una tecnica basata sulla **ricorsione** per la realizzazione pratica di un parser top-down.

# Struttura del parser ricorsivo

## Idea

Usare la pila del linguaggio di programmazione per “ricordare” il suffisso della forma sentenziale sinistra da riconoscere.

## Elementi chiave

- Il parser ha una procedura per ogni variabile della grammatica.
- La procedura  $A$  nel parser riconosce le stringhe generate da  $A$  nella grammatica.
- La procedura  $A$  usa il simbolo corrente e gli insiemi guida, per scegliere la produzione  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  da usare per riscrivere  $A$ .
- Per ogni simbolo  $X$  trovato nel corpo della produzione scelta:
  - Se  $X$  è un simbolo terminale, il metodo controlla che il simbolo corrente sia proprio  $X$ . In tal caso, fa avanzare il lexer al simbolo successivo. In caso contrario, il metodo segnala un errore di sintassi.
  - Se  $X$  è una variabile, il metodo invoca la procedura  $X$ .

# Algoritmo di parsing ricorsivo

```
var  $w$  : string
```

```
var  $i$  : int
```

//  $w$  è la stringa da riconoscere con \$ in fondo  
//  $i$  è l'indice del prossimo simbolo di  $w$  da leggere

```
procedure match( $a$  : symbol)
```

```
  if  $w[i] = a$  then  $i \leftarrow i + 1$  else error
```

```
procedure parse( $v$  : string)
```

//  $v$  è la stringa da riconoscere

```
   $w \leftarrow v\$$ 
```

```
   $i \leftarrow 0$ 
```

```
   $S()$ 
```

```
  match( $\$$ )
```

//  $S$  è il simbolo iniziale della grammatica  
// controlla di aver letto tutta la stringa

```
procedure  $A()$ 
```

//  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  sono le produzioni per  $A$

```
  if  $w[i] \in \text{GUIDA}(A \rightarrow \alpha_1)$  then
```

```
    :
```

```
  else if  $w[i] \in \text{GUIDA}(A \rightarrow \alpha_k)$  then
```

```
    for  $X \in \alpha_k$  do
```

```
      if  $X$  è un terminale then match( $X$ ) else  $X()$ 
```

```
    :
```

```
  else error
```

//  $w[i]$  non è nell'insieme guida di nessuna produzione per  $A$



# Implementazione Java del parser (classe base)

```
public abstract class Parser {  
    private String w;           // stringa da riconoscere  
    private int i;              // indice del prossimo simbolo  
  
    protected char peek()      // legge il simbolo corrente  
    { return w.charAt(i); }  
  
    protected void match(char a) // controlla il simbolo corrente  
    { if (peek() == a) i++; else throw error(); }  
  
    public void parse(String v) { // avvia il parsing di v  
        w = v + "$";  
        i = 0;  
        S();  
        match('$');  
    }  
  
    protected abstract void S(); // simbolo iniziale della grammatica  
  
    protected SyntaxError error() { ... } // emette errore e interrompe  
}
```

- Per semplicità assumiamo che i simboli siano caratteri.

# Esempio: parser per il linguaggio $a^n b^n$

## Grammatica

$$S \rightarrow aSb \mid \varepsilon$$

## Insiemi guida

- $\text{GUIDA}(A \rightarrow aSb) = \{a\}$
- $\text{GUIDA}(A \rightarrow \varepsilon) = \{b, \$\}$

## Codice del parser

```
public class AnBn extends Parser {  
    protected void S() {  
        switch (peek()) {  
  
            case 'a': // S → aSb  
                match('a');  
                S();  
                match('b');  
                break;  
  
            case 'b': // S → ε  
            case '$':  
                break;  
  
            default:  
                throw error();  
        }  
    }  
}
```

# Esercizi

## Implementazione di parser

Implementare il parser ricorsivo discendente per le seguenti grammatiche:

1. La grammatica delle stringhe della forma  $wcw^R$  dove  $w \in \{0,1\}^*$ :
  - $S \rightarrow c \mid 0S0 \mid 1S1$
2. La grammatica delle stringhe di parentesi quadre bilanciate:
  - $S \rightarrow \varepsilon \mid [S]S$
3. La grammatica delle stringhe della forma  $a^n b^n c^m$ :
  - $S \rightarrow XC$
  - $X \rightarrow \varepsilon \mid aXb$
  - $C \rightarrow \varepsilon \mid cC$
4. La grammatica delle espressioni aritmetiche in forma prefissa:
  - $E \rightarrow 0 \mid 1 \mid \dots \mid 9 \mid +EE \mid *EE$
5. La grammatica delle **espressioni aritmetiche** in forma infissa.

# Linguaggi Formali e Traduttori

## 4.3 Grammatiche fattorizzabili e ricorsive a sinistra

- Sommario
- Fattorizzazione
- Esempio di fattorizzazione
- Ricorsione immediata a sinistra
- Esempio di eliminazione della ricorsione
- Ricorsione a sinistra: caso generale
- Ricorsione indiretta a sinistra
- Eliminazione della ricorsione indiretta
- Esercizi

È proibito condividere e divulgare in qualsiasi forma i materiali didattici caricati sulla piattaforma e le lezioni svolte in videoconferenza: ogni azione che viola questa norma sarà denunciata agli organi di Ateneo e perseguita a termini di legge.

# Sommario

## Problema

- Molte grammatiche utili per descrivere linguaggi di programmazione non sono LL(1).
  1. Presenza di **produzioni fattorizzabili**

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

2. Presenza di **produzioni ricorsive a sinistra**

$$A \rightarrow A\alpha \mid \beta$$

3. Presenza di **ambiguità**

## In questa lezione

- Studiamo alcune tecniche per modificare produzioni fattorizzabili e ricorsive a sinistra senza cambiare il linguaggio generato dalla grammatica in modo da renderla – spesso, ma non sempre – LL(1).

# Fattorizzazione

## Problema

Data una grammatica con le produzioni

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

abbiamo

$$\text{GUIDA}(A \rightarrow \alpha\beta_1) \supseteq \text{FIRST}(\alpha) \quad \text{GUIDA}(A \rightarrow \alpha\beta_2) \supseteq \text{FIRST}(\alpha)$$

dunque

$$\text{GUIDA}(A \rightarrow \alpha\beta_1) \cap \text{GUIDA}(A \rightarrow \alpha\beta_2) \neq \emptyset$$

tranne nel caso degenere in cui  $\alpha$  genera solo  $\epsilon$ .

## Soluzione

**Fattorizzare** il prefisso comune  $\alpha$  introducendo una nuova variabile  $A'$ :

$$A \rightarrow \alpha A' \quad A' \rightarrow \beta_1 \mid \beta_2$$

# Esempio di fattorizzazione

La grammatica

- $S \rightarrow \text{if } E \text{ then } S \text{ else } S \text{ fi}$
- $S \rightarrow \text{if } E \text{ then } S \text{ fi}$
- $S \rightarrow a$
- $E \rightarrow b$

non è LL(1) infatti, dovendo espandere la variabile  $S$  quando il prossimo token nella stringa da riconoscere è **if**, il parser non saprebbe quale delle produzioni per  $S$  usare.

La grammatica è fattorizzabile nel modo seguente:

- $S \rightarrow \text{if } E \text{ then } S S'$
- $S' \rightarrow \text{else } S \text{ fi} \mid \text{fi}$
- $S \rightarrow a$
- $E \rightarrow b$

In particolare, gli insiemi guida delle produzioni della grammatica modificata sono ora disgiunti due a due.

# Ricorsione immediata a sinistra

## Problema

Una grammatica con le produzioni

$$A \rightarrow A\alpha \mid \beta$$

è detta **immediatamente ricorsiva a sinistra** in quanto la produzione  $A \rightarrow A\alpha$  ha  $A$  sia in testa che come primo simbolo del suo corpo. La grammatica non è LL(1):

$$\text{GUIDA}(A \rightarrow A\alpha) \supseteq \text{FIRST}(A) \supseteq \text{FIRST}(\beta) \quad \text{GUIDA}(A \rightarrow \beta) \supseteq \text{FIRST}(\beta)$$

## Osservazione

La grammatica genera stringhe  $\beta\alpha\alpha\cdots\alpha$  composte da una  $\beta$  seguita da zero o più  $\alpha$ .

## Soluzione

Introdurre una nuova variabile per spostare la ricorsione da sinistra a destra:

$$A \rightarrow \beta A' \quad A' \rightarrow \varepsilon \mid \alpha A'$$



# Esempio di eliminazione della ricorsione

La grammatica

- $E \rightarrow E + T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow n \mid (E)$

è immediatamente ricorsiva a sinistra nelle produzioni per  $E$  e per  $T$ . Ad esempio:

$$\text{GUIDA}(E \rightarrow E + T) = \text{FIRST}(E) = \{n, (\}$$

$$\text{GUIDA}(E \rightarrow T) = \text{FIRST}(T) = \{n, (\}$$

# Esempio di eliminazione della ricorsione

La grammatica

- $E \rightarrow E + T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow n \mid (E)$

è immediatamente ricorsiva a sinistra nelle produzioni per  $E$  e per  $T$ . Ad esempio:

$$\text{GUIDA}(E \rightarrow E + T) = \text{FIRST}(E) = \{n, (\}$$

$$\text{GUIDA}(E \rightarrow T) = \text{FIRST}(T) = \{n, (\}$$

Eliminando la ricorsione immediata a sinistra otteniamo la grammatica:

- $E \rightarrow TE'$
- $E' \rightarrow +TE' \mid \varepsilon$
- $T \rightarrow FT'$
- $T' \rightarrow *FT' \mid \varepsilon$
- $F \rightarrow n \mid (E)$

# Ricorsione a sinistra: caso generale

Una grammatica con le produzioni

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \cdots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$$

in cui nessun  $\beta_i$  inizia con  $A$ , genera stringhe della forma

$$\beta_i \alpha_{k_1} \alpha_{k_2} \cdots \alpha_{k_l}$$

L'eliminazione della ricorsione immediata a sinistra porta alla grammatica

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \cdots \mid \beta_n A' \\ A' &\rightarrow \varepsilon \mid \alpha_1 A' \mid \alpha_2 A' \mid \cdots \mid \alpha_m A' \end{aligned}$$

## Osservazioni

- In generale, l'eliminazione della ricorsione a sinistra non garantisce che la grammatica risultante sia LL(1).
- Ad esempio, nella grammatica qui sopra basta che uno degli  $\alpha_i$  sia annullabile per avere insiemi guida delle produzioni per  $A'$  non disgiunti.

# Ricorsione indiretta a sinistra

In alcune grammatiche alcune ricorsioni a sinistra sono “indirette”:

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow Ac \mid Sd \mid \varepsilon \end{aligned}$$

Tentando di eliminare la ricorsione a sinistra per le produzioni di  $A$  otteniamo

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow SdA' \mid A' \\ A' &\rightarrow \varepsilon \mid cA' \end{aligned}$$

ma la grammatica non è LL(1), infatti:

- $\text{GUIDA}(A \rightarrow SdA') \supseteq \text{FIRST}(S) \supseteq \text{FIRST}(A) \supseteq \text{FIRST}(A') \ni c$
- $\text{GUIDA}(A \rightarrow A') \supseteq \text{FIRST}(A') \ni c$

C'è una **ricorsione indiretta** a sinistra che riguarda la variabile  $A$

$$A \Rightarrow Sd \Rightarrow Aad$$

# Eliminazione della ricorsione indiretta

## Idea

Si può esporre la ricorsione indiretta facendo opportune riscritture di variabili.

## Algoritmo

1. Si impone un ordine (arbitrario) alle variabili della grammatica.
2. Considerando ogni variabile secondo l'ordine imposto, si elimina la ricorsione immediata per quella variabile e si riscrivono le occorrenze di quella variabile che compaiono nei corpi delle produzioni delle variabili seguenti.

## Esempio

$$\begin{array}{llll} S \rightarrow Aa \mid b & \Rightarrow & S \rightarrow Aa \mid b & \Rightarrow & S \rightarrow Aa \mid b \\ A \rightarrow Ac \mid Sd \mid \varepsilon & \Rightarrow & A \rightarrow Ac \mid Aad \mid bd \mid \varepsilon & \Rightarrow & A \rightarrow bdA' \mid A' \\ & & & & A' \rightarrow \varepsilon \mid cA' \mid adA' \end{array}$$

# Esercizi

1. Applicare le trasformazioni studiate in questa lezione alla grammatica non ambigua delle formule booleane per farla diventare LL(1).
2. Implementare il parser top-down per la grammatica ottenuta nell'esercizio precedente così ottenuta. Scegliere caratteri ASCII "normali" per rappresentare i connettivi logici, ad esempio & per  $\wedge$ , | per  $\vee$  e ~ per  $\neg$ .