

# Analisi lessicale

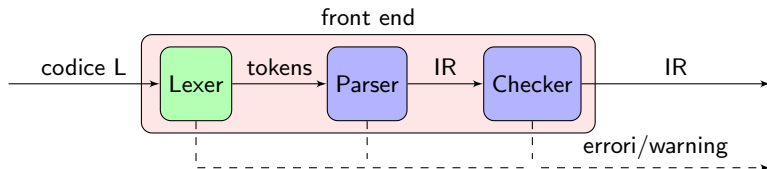
Enea Zaffanella

[enea.zaffanella@unipr.it](mailto:enea.zaffanella@unipr.it)

28 settembre 2020

Linguaggi, interpreti e compilatori  
Laurea Magistrale in Scienze informatiche

## Il lexer (analizzatore lessicale)



- Input: sequenza di caratteri
- Output: sequenza di *token*
- $token = \langle part\_of\_speech, lexeme \rangle$
- Esempi:  $\langle KWD, while \rangle$ ,  $\langle IDENT, somma \rangle$ ,  
 $\langle INT, 42 \rangle$ ,  $\langle FLOAT, 3.1415 \rangle$ ,  $\langle STR, "Hello" \rangle$ , ...

# Lexer: specifica vs implementazione

## Specifica

- Come definire in modo rigoroso quali sono i token validi?
- Linguaggio adeguato: **RE** (espressioni regolari)
- Comprensibile dal *progettista* (essere umano)

## Implementazione

- Fattore critico: efficienza
- Lexer detti anche “scanner”
- Riconoscitore: **DFSA** (automa a stati finiti deterministico)
- Spesso generato automaticamente partendo dalla specifica

# Richiami di nozioni base linguaggi formali

- $\Sigma$ : alfabeto (un insieme **finito** di simboli)
- $\Sigma^*$ : insieme di tutte le stringhe **finite** sull'alfabeto  $\Sigma$
- $\epsilon \in \Sigma^*$ : la stringa **vuota** (nessun simbolo, lunghezza 0)
- $L \subseteq \Sigma^*$ : linguaggio
  - $\emptyset \equiv \{\}$  è un linguaggio
  - $\{\epsilon\}$  è un linguaggio
  - $\Sigma$  è un linguaggio (abuso: simbolo  $\equiv$  stringa lunga 1)
  - $\Sigma^*$  è un linguaggio (infinito, se  $\Sigma \neq \emptyset$ )
- Problema di **decisione** per  $L$  (riconoscimento, appartenenza): trovare un **algoritmo** che, data una qualunque stringa  $s \in \Sigma^*$ , decida (in tempo finito) se vale  $s \in L$ .

# Le espressioni regolari: sintassi

## L'insieme $\text{RExpr}$ delle RE (espressioni regolari)

- $\epsilon \in \text{RExpr}$ ;  $a \in \text{RExpr}$ , per ogni  $a \in \Sigma$
- se  $e \in \text{RExpr}$ , allora:
  - $(e) \in \text{RExpr}$  (parentesi)
  - $e^* \in \text{RExpr}$  (chiusura di Kleene)
- se  $e_1 \in \text{RExpr}$  e  $e_2 \in \text{RExpr}$ , allora:
  - $e_1 e_2 \in \text{RExpr}$  (concatenazione)
  - $e_1 \mid e_2 \in \text{RExpr}$  (alternanza)

## Precedenza operatori

- la chiusura di Kleene ha la precedenza sulla concatenazione
- la concatenazione ha la precedenza sull'alternanza
- esempio:  $a \mid a^* b \equiv ((a) \mid ((a)^* b))$

# Le espressioni regolari: semantica

$L: \text{RExpr} \rightarrow \wp(\Sigma^*)$  (linguaggio denotato da una RE)

$$L(\epsilon) = \{\epsilon\};$$

$$L(a) = \{a\};$$

$$L(e_1 e_2) = L(e_1) L(e_2) = \{ s_1 s_2 \mid s_1 \in L(e_1), s_2 \in L(e_2) \};$$

$$L(e_1 \mid e_2) = L(e_1) \cup L(e_2);$$

$$L(e^*) = L(e)^* = \bigcup_{i=0}^{\infty} L(e)^i = \bigcup_{i=0}^{\infty} \overbrace{L(e) \cdots L(e)}^i;$$

$$L((e)) = L(e).$$

# Abbreviazioni RE

## Sintassi

- $e^+$  (chiusura positiva),  $e^?$  (opzionalità)
- $[a_1 \dots a_n]$ ,  $[a_1 - a_2]$  (classe/range di simboli)
- $[\hat{a}_1 \dots a_n]$ ,  $[\hat{a}_1 - a_2]$  (complemento di classe/range)

## Semantica

$$\begin{aligned}
 L(e^+) &= \cup_{i=1}^{\infty} L(e)^i; & L([a_1 \dots a_n]) &= \{a_1, \dots, a_n\}; \\
 L(e^?) &= \{\epsilon\} \cup L(e); & L([a_1 - a_2]) &= \{a \in \Sigma \mid a_1 \leq a \leq a_2\}; \\
 & & L([\hat{a}_1 \dots a_n]) &= \Sigma \setminus L([a_1 \dots a_n]).
 \end{aligned}$$

# I token dei linguaggi di programmazione

## Classificazione

- parole chiave
- identificatori
- costanti letterali (intere, floating point, stringa, ...)
- operatori (matematici, logici, ...)
- “punteggiatura” (parentesi, virgola, punto e virgola, ...)
- commenti (singola linea, multi linea)



## Esempi di specifica: keyword

### Parole chiave: facile

- `if, then, else, while, ...`
- si usa l'alternanza: `if | then | else | while ...`
- unica categoria sintattica: `KEYWORD`

### Case insensitive (e.g., SQL)

- `select, SELECT, SeLeCt, SElect, ...`
- Noioso e troppo verboso: si usano impostazioni *ad hoc*

## Esempi di specifica: identificatori

### Identificatori (1)

- $[a-zA-Z_][0-9a-zA-Z_]^*$

### Identificatori (2)

- $[a-zA-Z_]([0-9] \mid [a-zA-Z_])^*$

### Identificatori (3)

- $\text{DIGIT} = [0-9]$
- $\text{LETTER} = [a-zA-Z_]$
- $\{\text{LETTER}\}(\{\text{DIGIT}\} \mid \{\text{LETTER}\})^*$

## Esempi di specifica: costanti

### Costanti intere

- $\{\text{DIGIT}\}^+$
- Nota: accetta 000000, non accetta -1

### Costanti floating point

- $[+-]?[0-9]^+.[0-9]^*$
- occorre distinguere iterazione positiva dal carattere +?

### Costanti carattere

- $'[^']*'$
- come specifichiamo la costante carattere '?'

# Esempi di specifica: operatori e punteggiatura

Ogni lessema ha la sua categoria lessicale (singoletto)

lessema	categoria	lessema	categoria
(	OPEN_PAREN	)	CLOSE_PAREN
[	OPEN_BRACKET	]	CLOSE_BRACKET
{	OPEN_BRACE	}	CLOSE_BRACE
+	PLUS	-	MINUS
+=	PLUS_ASSIGN	-=	MINUS_ASSIGN
:	COLON	::	SCOPE
<	LESS_THAN	<<	SHIFT_LEFT
>	GREATER_THAN	>>	SHIFT_RIGHT
.	DOT	...	ELLIPSIS
...	...		

## Esempi di specifica: commenti

### Commento (singola linea) del C++

- `// [^\n]*\n`

### Commento (singola linea) di SQL

- `-- [^\n]*\n`

### Commento multilinea (C/C++/Java/SQL/...)

- `/\* ([^*] | \*+ [^/*] ) * \*+ /`
- Chiarissimo, vero?
- (Vedremo un modo migliore per ottenere lo stesso effetto)

## Risoluzione ambiguità

Più RE possono accettare (parti de-) lo stesso input

- Preferenza al **lessema più lungo**
  - `forwhile` è un unico IDENT  
(non sono le due KEYWORD `for` e `while`)
  - `>>` è un unico SHIFT\_RIGHT (non due GREATER\_THAN)

A parità di lunghezza?

- Si stabilisce un **ordine di priorità** tra le RE
  - `for` e `while` sono KEYWORD (non IDENT)

Nota: rispettare le regole può essere “noioso”

### Annidamento liste argomenti template nel C++

- Per lo standard C++03 e precedenti:
  - **corretto**: `std::vector<std::list<int> >`
  - **errore** (sintattico): `std::vector<std::list<int>>>`
- Per lo standard C++11 e successivi:
  - **corretto**: `std::vector<std::list<int> >`
  - **corretto**: `std::vector<std::list<int>>`

## Tra il dire e il fare ...

### Verso una esercitazione ...

- Goal: un lexer per una **porzione** (significativa) di un linguaggio di programmazione
- Focus: **specifica** (non implementazione)
- Prerequisito: sapere usare un **generatore** di lexer
- Lo strumento scelto: `flex`