

Course 2

Chapter 1. Scanning

Definition = treats the source program as a sequence of characters, detect lexical tokens, classify and codify them

INPUT: source program
OUTPUT: PIF + ST

Algorithm Scanning v1
While (not (eof)) **do**
 detect (token) ;
 classify (token) ;
 codify (token) ;
End_while

Detect

I am a student.

- Separators => ***Remark 1) + 2)***

if (x==y) {x=y+2}

- Look-ahead => ***Remark 3)***

Codify

- Codification table
- Identifier, constant => Symbol Table (ST)
- PIF = Program Internal Form = array of pairs
- Token – replaced by pair (code, position in ST)

identifier, constant

Algorithm Scanning v2

```
While (not (eof)) do
    detect(token);
    if token is reserved word OR operator OR separator
        then genFIP(token, 0)
        else
            if token is identifier OR constant
                then index = pos(token, ST);
                    genFIP(token, index)
                else message "Lexical error"
            endif
        endif
    endwhile
```

Remarks:

- genPIF = adds a pair (token, position) to PIF
- Pos(token,ST) – searches *token* in symbol table ST; if found then return position; if not found insert in SR and return position
- Order of classification (reserved word, then identifier)
- If-then-else imbricate => detect error if a token cannot be classified
- **Also:** comments and white spaces are eliminated

Formal Languages

- *basic notions-*

Examples of languages

- natural (ex. English, Romanian)
- programming (ex. C,C++, Java, Python)
- formal

A formal language is a set

Ex.:

$$L = \{a^n b^n \mid n > 0\} \quad L = \{ab, aabb, aaabbb, \dots\}$$

$$L' = \{01^n \mid n \geq 0\} \quad L' = \{0, 01, 011, \dots\}$$

$$L'' = \{(01)^n \mid n \geq 0\} \quad L'' = \{\text{nothing}, 01, 0101, \dots\}$$

Example

a boy has a dog

$S \rightarrow PV$
 $P \rightarrow a N$
 $N \rightarrow boy \text{ or } N \rightarrow dog$
 $\quad \quad \quad (N \rightarrow boy \mid dog)$
 $V \rightarrow QC$
 $Q \rightarrow has$
 $C \rightarrow BN$
 $B \rightarrow a$

- $A \rightarrow \alpha = \text{rule}$
- $S, P, V, N, Q, C, B = \text{nonterminal symbols}$
- $a, boy, dog, has = \text{terminal symbols}$

Remarks

1. Sentence = word, sequence (contains only terminal symbols) ; denoted w.
2. $S \Rightarrow PV \Rightarrow a NV \Rightarrow a NQC \Rightarrow a N \text{ has } C$ - sentential form
In general : $w = a_1 a_2 \dots a_n$
3. The rule guarantees syntactical correctness, but not the semantical correctness (*A dog has a boy*)

Grammar

- **Definition:** A (formal) **grammar** is a 4-tuple: $G=(N,\Sigma,P,S)$

with the following meanings:

- N – set of nonterminal symbols and $|N| < \infty$
- Σ - set of terminal symbols (alphabet) and $|\Sigma| < \infty$
- P – finite set of productions (rules), with the property:
$$P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* X (N \cup \Sigma)^*$$
- $S \in N$ – start symbol /axiom

Remarks :

1. $(\alpha, \beta) \in P$ is a production denoted $\alpha \rightarrow \beta$
2. $N \cap \Sigma = \emptyset$

A^* = transitive and reflexive closure =

$$\{a, aa, aaa, \dots\} \cup \{a^0\}$$

$$A = \{a\}$$

$$A^+ = \{a, aa, aaa, \dots\}$$

$$X^0 = \epsilon$$

Binary relations defined on $(N \cup \Sigma)^*$

- **Direct derivation**

$\alpha \Rightarrow \beta , \alpha, \beta \in (N \cup \Sigma)^* \text{ if } \alpha=x_1xy_1 , \beta=x_1yy_1 \text{ and } x \rightarrow y \in P$
(x is transformed in y)

- **k derivation**

$\overset{k}{\alpha \Rightarrow \beta} , \alpha, \beta \in (N \cup \Sigma)^*$

sequence of k direct derivations $\alpha \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_{k-1} \Rightarrow \beta , \alpha, \alpha_1, \alpha_2, \dots, \alpha_{k-1}, \beta \in (N \cup \Sigma)^*$

- **+ derivation**

$\overset{+}{\alpha \Rightarrow \beta} \text{ if } \exists k > 0 \text{ such that } \overset{k}{\alpha \Rightarrow \beta}$ (there exists at least one direct derivation)

- *** derivation**

$\overset{*}{\alpha \Rightarrow \beta} \text{ if } \exists k \geq 0 \text{ such that } \overset{k}{\alpha \Rightarrow \beta}$ namely, $\overset{*}{\alpha \Rightarrow \beta} \Leftrightarrow \overset{+}{\alpha \Rightarrow \beta} \text{ OR } \overset{0}{\alpha \Rightarrow \beta} (\alpha = \beta)$

Definition: Language generated by a grammar $G=(N,\Sigma,P,S)$ is:

$$L(G) = \{w \in \Sigma^* \mid S \xrightarrow{*} w\}$$

Remarks:

1. $S \xrightarrow{*} \alpha, \alpha \in (N \cup \Sigma)^*$ = sentential form
 $S \xrightarrow{*} w, w \in \Sigma^*$ = word / sequence

2. Operations defined for languages (sets) :

$$L_1 \cup L_2, L_1 \cap L_2, L_1 - L_2, \overline{L} \text{ (complement)}, L^+ = \bigcup_{k>0} L^k, L^* = \bigcup_{k \geq 0} L^k$$

$$\text{Concatenation: } L = L_1 L_2 = \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$$

3. $|w|=0$ (empty word - denoted ϵ)

$$L_1 = \{a, b, aa\}$$

$$L_2 = \{c, d, cd\}$$

$$L_1 L_2 = \{ac, ad, acd, bc, bd, bcd, aac, aad, aacd\}$$

Definition: Two grammar G_1 and G_2 are equivalent if they generate the same language

$$L(G_1) = L(G_2)$$

Chomsky hierarchy(based on form $\alpha \rightarrow \beta \in P$)

- type 0 : no restriction
- type 1 : context dependent grammar ($x_1Ay_1 \rightarrow x_1\gamma y_1$)
- type 2 : context free grammar ($A \rightarrow \alpha \in P$,where $A \in N$ and $\alpha \in (N \cup \Sigma)^*$)
- type 3 : regular grammar ($A \rightarrow aB | a \in P$)

Remark :

$$\text{type 3} \subseteq \text{type 2} \subseteq \text{type 1} \subseteq \text{type 0}$$

Notations

- A, B, C, \dots – nonterminal symbols
- $S \in N$ – start symbol
- $a, b, c, \dots \in \Sigma$ – terminal symbol
- $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$ - sentential forms
- ε – empty word
- $x, y, z, w \in \Sigma^*$ - words
- $X, Y, U, \dots \in (N \cup \Sigma)$ – grammar symbols (nonterminal or terminal)

Regular grammars

- $G = (N, \Sigma, P, S)$ right linear grammar if

$\forall p \in P: A \rightarrow aB$ or $A \rightarrow \underline{a}$, where $A, B \in N$ and $a, b \in \Sigma$

- $G = (N, \Sigma, P, S)$ regular grammar if

- G is right linear grammar

and

- $A \rightarrow \epsilon \notin P$, with the exception that $S \rightarrow \epsilon \in P$, in which case S does not appear in the rhs (right hand side) of any other production

- $L(G) = \{w \in \Sigma^* \mid S \xrightarrow{*} w\}$ - right linear language

A- \rightarrow aA|a ok ✓
S- \rightarrow aA| ϵ and A- \rightarrow b ok ✓
S- \rightarrow aA| ϵ and A- \rightarrow ϵ NOT ok ✗
S- \rightarrow aA| ϵ and A- \rightarrow bS|a NOT ok ✗