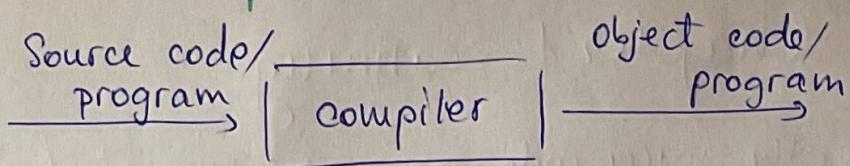


# Formal Languages and Compiler Design

## Course 1:

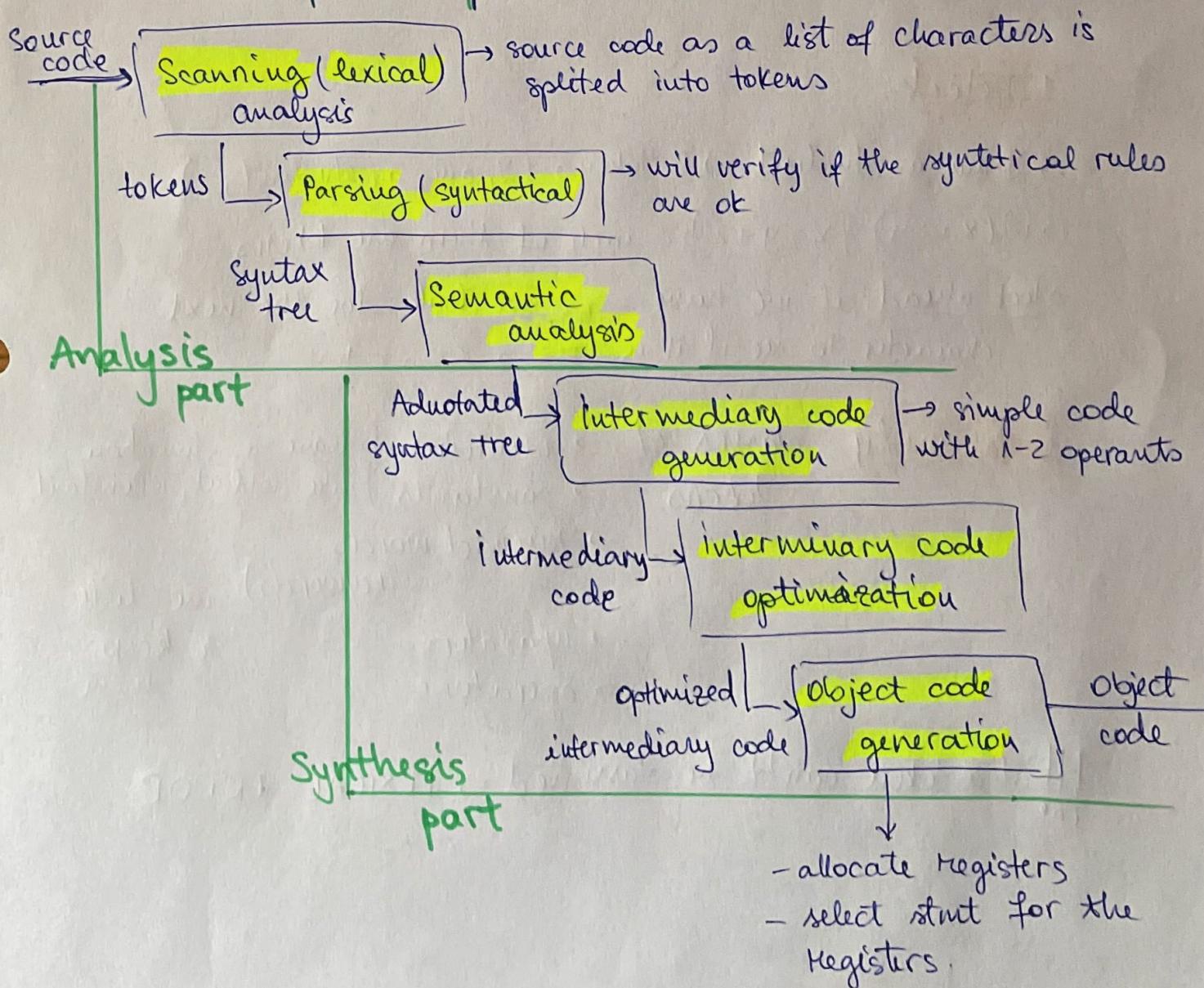
### What is a compiler?



↓ looks at the whole code

Interpreter: looks line by line and simulates execution

### Structure of a compiler:



Error handling: logical errors are usually not detected by the compiler

Symbol table management:

- it is important that the identifier is called m or a
- collects all symbolic names
- identifier = name to variables or functions
- symbolic names that are not identifiers id = exp.

## Chapter 1: Scanning

Def: Treats the source code as a sequence of characters, detect lexical tokens, classify and modify them

PIF - programme internal form

ST - symbol table

### 1. Detect:

- I am a student. ⇒ 5 tokens (I, am, a, student, .)
- if (x == y) { x = y + 2 } ← each of them are tokens
- look ahead: if we have =, < we look for the next character to see if it is ==, <=, =+

2. classify: classes of tokens

→ identifiers	} have a function like isident, isconstant (keywords)
→ constants	
→ reserved words	
→ separators	
→ operators	

use lists to keep them

!!! If a token can't be classified ⇒ LEXICAL ERROR

### 3. Modify:

- replace identifier with id
- identifier, constant  $\Rightarrow$  Symbol Table
- PIT  $\Rightarrow$  array of pairs  $\leftrightarrow$  pairs (token, position in ST)

### Symbol Table:

Def: Contains all information collected during compiling regarding the symbolic names from the source program  
↓  
identifiers, constants

#### Organization:

- \* unsorted table:  $O(n)$
- \* sorted table (alphabetic, numeric):  $O(\lg n)$
- \* Binary search tree (balanced):  $O(\lg n)$
- \* Hash table
  - open bucket: with linked lists  $O(1)$ .
  - close bucket: search for the first open position

For lab: hash function = sum of ASCII codes of char.

## Formal Languages Basic Notations

### Grammar:

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

- $N$  - set of nonterminal symbols,  $|N| < \infty$
- $\Sigma$  - set of terminal symbols (alphabet),  $|\Sigma| < \infty$
- $P$  - finite set of productions (rules)
- $S \in N$  - start symbol / axiom.

### Remarks :

1.  $(\alpha, \beta) \in P$  is a production denoted  $\alpha \rightarrow \beta$

2.  $A^* = \{a, aa, aaa, \dots\} \setminus \{a^0\}$

$A = \{a\}$

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

$a^0 = \epsilon$

3. Direct derivation :

$\alpha \Rightarrow \beta, \alpha, \beta \in (N \cup \Sigma)^*$  if  $\alpha = x_1 y_1$  and  $\beta = x_1 y_1$  and  $x \rightarrow y$   
 $\epsilon \in P$   
( $x$  is transformed into  $y$ ).

Definition: Language generated by a grammar

$G = (N, \Sigma, P, S)$  is :

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

### Remarks :

1. Operations defined for languages

$L_1 \cup L_2, L_1 \cap L_2, L_1 - L_2, \overline{L}, L^+ = \bigcup_{k>0} L^k$

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

2. Two grammars  $G_1$  and  $G_2$  are equivalent if they generate the same language  $L(G_1) = L(G_2)$

3. Context free grammar :  $A \rightarrow \alpha \in P$ , where  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$

4. Regular grammar :  $A \rightarrow aB \mid a \in P$

## Notations :

- $A, B, C, \dots$  - nonterminal symbols  $\rightarrow N$
- $s \in N$  - start symbol
- $a, b, c \in \Sigma$  - terminal symbols
- $\epsilon$  - empty word
- $x, y, z, w \in \Sigma^*$  - words
- $x, y, u, \dots \in (N \cup \Sigma)$  - grammar symbols (term or. nonterm.)

## Finite Automata :

Definition: A finite automata (FA) is a 5-tuple :

$$M = (Q, \Sigma, \delta, q_0, F)$$

where :

- $Q$  : finite set of states
- $\Sigma$  : finite alphabet
- $\delta$  : transition function :  $\delta: Q \times \Sigma \rightarrow P(Q)$
- $q_0$  : initial state  $q_0 \in Q$
- $F \subseteq Q$  - set of final states

## Remarks :

1. if  $|\delta(q, a)| \leq 1 \Rightarrow$  deterministic FA (DFA)
2. if  $|\delta(q, a)| > 1$  (more than a state obtained as a result)  $\Rightarrow$  nondeterministic FA (NFA)

Configuration :  $C = (q, x)$  where  $\rightarrow q$  state  
 $\rightarrow x$  unread sequence  
from input  $x \in \Sigma^*$

Initial configuration :  $(q_0, w)$  -  $w$  - whole sequence

Final configuration :  $(q_f, \epsilon)$ ,  $q_f \in F$ ,  $\epsilon$  - empty word

## Relations between configurations:

- $\vdash$  move / transition (simple, one step)  
 $(q, ax) \vdash (p, x) \text{ if } p = S(q, a)$
- $\vdash^k$   $k$  move = sequence of  $k$  simple transitions  $q_0 \vdash C_1 \dots \vdash C_k$
- $\vdash^*$  + move :  $C \vdash^* C'$  :  $\exists k \geq 0$  such that  $C \vdash^k C'$
- $\vdash^*$  \* move (star move) :  $C \vdash^* C'$ ,  $\exists k \geq 0$  such that  $C \vdash^k C'$

Definition: Language accepted by FA  $M = (Q, \Sigma, S, q_0, F)$   
 is  $L(M) = \{w \in \Sigma^* \mid (q_0, w) \vdash^* (q_f, \epsilon), q_f \in F\}$

## Regular Languages

### Regular grammars:

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

$\forall p \in P : A \rightarrow aB \text{ or } A \rightarrow b$

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

\*  $G$  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 right hand side of any other production

### Regular grammar:

\* right linear grammar  
 $A \rightarrow aB$  or  $A \rightarrow b$

\*  $\epsilon \in$  only for the starting symbol ( $s_0$ )

if  $\epsilon \in S_0 \Rightarrow S_0$  should not appear in the right hand side of any product.

### Example:

\*  $\begin{cases} S \rightarrow aA \\ A \rightarrow a \end{cases} \rightarrow \text{regular}$

\*  $\begin{cases} S \rightarrow aA \\ A \rightarrow aA \mid \epsilon \end{cases} \rightarrow \text{NOT regular}$

\*  $\begin{cases} S \rightarrow aS \mid aA \\ A \rightarrow bS \mid b \end{cases} \rightarrow \text{regular}$

\*  $\begin{cases} S \rightarrow aA \mid \epsilon \\ A \rightarrow aS \end{cases} \rightarrow \text{NOT regular}$

## Regular sets :

Definition: Let  $\Sigma$  be a finite alphabet. We define regular sets over  $\Sigma$  recursively in the following way :

1.  $\emptyset$  is a regular set over  $\Sigma$
2.  $\{\epsilon\}$  is a regular set over  $\Sigma$
3.  $\{a\}$  — u — u — u —,  $\forall a \in \Sigma$
4. If  $P, Q$  reg. sets  $\Rightarrow P \cup Q, PQ, P^*$  reg sets over  $\Sigma$
5. Nothing else is a regular set over  $\Sigma$

## Regular expressions :

Definition:  $\Sigma$ -finite alphabet. We define regular expressions over  $\Sigma$  recursively :

1.  $\emptyset$  is a regular expression denoting the reg set  $\emptyset$
2.  $\epsilon$  — u — u — u — u —  $\{\epsilon\}$
3.  $a$  — u — u — u —  $\{a\}$ ,  $\forall a \in \Sigma$
4. If  $p, q$  are reg. expr. denoting  $P, Q$  then  
 $p+q$  for  $P \cup Q$ ,  $pq$  for  $PQ$ ,  $p^* \rightarrow P^*$
5. Nothing else is a regular expression

## Remarks :

1.  $P^+ = PP^*$
2. Priority of operations :  $*$ , concat,  $+$  (from high).
3. Two regular expressions are equivalent iff they denote the same regular set

Example:  $(0+1)^*$  denotes  $\{\epsilon, 0, 1, 00, 01, 11, 10, \dots\}$   
 $0^*, 1^*$  denotes  $\{\epsilon, 0, 1, 01, 00, 11, \dots\}$

### Theorem :

A language is a regular set if and only if it is a linear language

!!! Use to find the

Reg Expression from a grammar !!!

$$G = (\{S, A, B\}, \{0, 1\}, P, S)$$

$$P: S \rightarrow 0A \mid 1B \mid \epsilon$$

$$A \rightarrow 0B \mid 1A$$

$$B \rightarrow 0S \mid 1$$

$$\left. \begin{aligned} x &= ax + b \\ \Rightarrow x &= a^*b \end{aligned} \right| (*)$$

$$\left. \begin{aligned} S &= 0A + 1B + \epsilon \\ A &= 0B + 1A \\ B &= 0S + 1 \end{aligned} \right\} \Rightarrow$$

$$A = 0(0S + 1) + 1A$$

$$A = 00S + 01 + 1A \text{ from } (*)$$

$$\Rightarrow A = 1^*(00S + 01)$$

$$S = 0A + 1B + \epsilon \rightarrow \text{Replace } A, B$$

$$S = 01^*(00S + 01) + 1(0S + 1) + \epsilon$$

$$S = \underbrace{(01^*00 + 10)}_a S + \underbrace{01^*01 + 11 + \epsilon}_b$$

$$\Rightarrow S = (01^*00 + 10)^* (01^*01 + 11 + \epsilon)$$

### Theorem :

A language is a regular set if and only if it is accepted by a FA

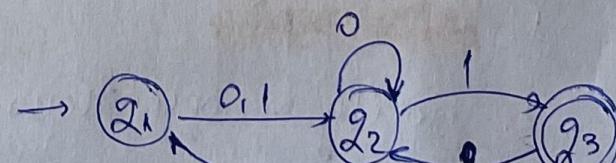
$$q_1 = q_3 0 + \epsilon \quad (1)$$

$$\left. \begin{aligned} q_2 &= q_1 0 + q_1 1 + q_2 0 + q_3 0 \quad (2) \\ q_3 &= q_2 1 \quad (3) \end{aligned} \right\}$$

$$3 \rightarrow 1 \Rightarrow q_1 = q_2 1 0 + \epsilon$$

$$2 \Rightarrow q_2 = (q_2 1 0 + \epsilon) 0 + (q_2 1 0 + \epsilon) 1 + q_2 0 + q_2 1 0$$

$$q_2 = q_2 1 0 0 + 0 + q_2 1 0 1 + 1 + q_2 0 + q_2 1 0$$



$$\left. \begin{aligned} x &= x a + b \\ \Rightarrow x &= b a^* \end{aligned} \right\}$$

Solution is the sum of the final states

$$q_2 = q_2(100 + 101 + 0 + 1) + 0 + 1$$

$$\Rightarrow q_2 = (0+1)(100+101+0+1)^*$$

$$\Rightarrow q_3 = q_2^* = (0+1)(100+101+0+1)^* 1.$$

Regular grammars:

Theorem 1: For every regular grammar  $G = (N, \Sigma, P, S)$  there exists a FA  $M = \{Q, \Sigma, \delta, q_0, F\}$  such that  $L(G) = L(M)$

Construct  $M$  based on  $G$

$$Q = N \cup \{k\}, k \notin N$$

$$q_0 = S$$

$$F = \{k\} \cup \{s \mid \text{if } s \rightarrow \epsilon \in P\}$$

$$\delta: \begin{aligned} \text{if } A \rightarrow aB \text{ then } \delta(A, a) &= B \\ \text{if } A \rightarrow a \text{ then } \delta(A, a) &= k \end{aligned}$$

Theorem 2: For any FA  $M = \{Q, \Sigma, \delta, q_0, F\}$  there exist a right linear grammar  $G = \{N, \Sigma, P, S\}$  such that  $L(G) = L(M)$

$$N = Q$$

$$S = q_0$$

$$P: \text{if } \delta(q_0, a) = p \text{ then } q \rightarrow ap \in P$$

$$\text{if } p \in F \text{ then } q \rightarrow a \in P$$

$$\text{if } q_0 \in F \text{ then } S \rightarrow \epsilon$$

Regular sets :

If  $L_1$  and  $L_2$  are right linear languages, then  
 $L_1 \cup L_2$ ,  $L_1 L_2$ ,  $L_1^*$  are right linear languages.

Proof:

$L_1, L_2$  - right linear lang  $\Rightarrow \exists G_1, G_2$  such that

$G_1 = (N_1, \Sigma_1, P_1, S_1)$  and  $L_1 = L(G_1)$  assume

$G_2 = (N_2, \Sigma_2, P_2, S_2)$  and  $L_2 = L(G_2)$ ,  $N_1 \cap N_2 = \emptyset$

\*  $L_1 \cup L_2 \rightarrow a+b$

$G_3 = (N_3, \Sigma_3, P_3, S_3)$

$N_3 = N_1 \cup N_2 \cup \{S_3\}$

$\Sigma_3 = \Sigma_1 \cup \Sigma_2$

$P_3 = P_1 \cup P_2 \cup \{S_3 \rightarrow \alpha_1 \mid S_1 \rightarrow \alpha_1 \in P_1\} \cup$   
 $\cup \{S_3 \rightarrow \alpha_2 \mid S_2 \rightarrow \alpha_2 \in P_2\}$

$\Rightarrow G_3$  - right linear language

and

$L(G_3) = L(G_1) \cup L(G_2)$

\*  $L_1 L_2$  (concatenation)  $\rightarrow ab$

$G_4 = (N_4, \Sigma_4, P_4, S_4)$  keep the second

$N_4 = N_1 \cup N_2$

And change

$\Sigma_4 = \Sigma_1 \cup \Sigma_2$

$S \rightarrow a$  with

$S_4 = S_1$

$S \rightarrow aS_1$ ,  $S_1$  - from  
the first

$P_4 : \{A \rightarrow aB \mid \text{if } A \rightarrow aB \in P_1\} \cup$

$\{A \rightarrow aS_2 \mid \text{if } A \rightarrow a \in P_1\} \cup$

$\{S_1 \rightarrow x \mid \text{if } S_1 \rightarrow \epsilon \text{ and } S_2 \rightarrow x\} \cup P_2$

$\Rightarrow G_4$  right linear language

if  $\epsilon \in \text{second}$

and

$L(G_4) = L(G_1) L(G_2)$

-10-

$\Rightarrow$  if 2-final st. of first  
2 will also be f.state

\*  $L_1^*$  // IDEA: Concatenate  $L_1$  with itself

$$G_5 = (N_5, \Sigma_1, P_5, S_5)$$

$$N_5 = N_1 \cup \{S_5\}$$

$$P_5 : P_1 \cup \{S_5 \rightarrow \epsilon\} \cup$$

$$\{S_5 \rightarrow d_1 \mid s_1 \rightarrow d_1 \in P_1\} \cup$$

$$\{A \rightarrow aS_1 \mid \text{if } A \rightarrow a \in P_1\}$$

Lemma:

If  $L_1$  and  $L_2$  are accepted by a FA then

$L_1 \cup L_2$ ,  $L_1 L_2$  and  $L_1^*$  are accepted by FA.

$$M_1 = (Q_1, \Sigma_1, \delta_1, q_{01}, F_1) \text{ such that } L_1 = L(M_1)$$

$$M_2 = (Q_2, \Sigma_2, \delta_2, q_{02}, F_2) \text{ such that } L_2 = L(M_2)$$

\*  $L(M_3) = L(M_1) \cup L(M_2) \rightarrow a+b$

$$M_3 = (Q_3, \Sigma_3, \delta_3, q_{03}, F_3)$$

$$Q_3 = Q_1 \cup Q_2 \cup \{q_{03}\}$$

$$\Sigma_3 = \Sigma_1 \cup \Sigma_2$$

$$F_3 = F_1 \cup F_2 \cup \{q_{03} \mid \text{if } q_{01} \in F_1 \text{ or } q_{02} \in F_2\}$$

$$\delta_3 = \delta_1 \cup \delta_2 \cup \{ \delta_3(q_{03}, a) = p \mid \exists \delta_1(q_{01}, a) = p \} \\ \cup \{ \delta_3(q_{03}, a) = p \mid \exists \delta_2(q_{02}, a) = p \}$$

\*  $L(M_3) = L(M_1) L(M_2) \rightarrow ab$

$$M_4 = (Q_4, \Sigma_4, \delta_4, q_{04}, F_4)$$

$$Q_4 = Q_1 \cup Q_2$$

$$q_{04} = q_{01}$$

$$F_4 = F_2 \cup \{q \in F_1 \mid \text{if } q_{02} \in F_2\}$$

$$\delta_4(q_1, a) = \delta_1(q_1, a) \text{ if } q_1 \in Q - F_1$$

$$\delta_1(q_1, a) \cup \delta_2(q_{02}, a) \text{ if } q_1 \in F_1$$

$$\delta_2(q_1, a) \text{ if } q_1 \in Q_2$$

Tip: • keep everything for the second language  
• for every single term, put  $S \rightarrow aS$

$$* L(M_3) = L(M_1)^*$$

$$M_5 = (Q_5, \Sigma_1, \delta_5, q_{05}, F_5)$$

$$Q_5 = Q_1$$

$$q_{05} = q_0$$

$$F_5 = F_1 \cup \{q_0\}$$

$$\delta_5(q, a) = \delta_1(q, a) \text{ if } q \in Q_1 - F_1$$

$$\delta_1(q, a) \cup \delta_1(q_0, a) \text{ if } q \in F_1$$

## Pumping Lemma :

!! To decide if a language is regular or not

Theorem (Bar-Hillel, Pumping lemma)

Let  $L$  be a regular language.  $\exists p \in \mathbb{N}$ , such that if  $w \in L$  with  $|w| > p$  then

$$w = xyz, \text{ where } 0 < |y| \leq p$$

and

$$xy^i z \in L, \forall i \geq 0$$

(Proof Course 5-4).

Example :  $L = \{0^n 1^n \mid n \geq 0\}$

Suppose  $L$  is regular  $\rightarrow w = xyz = 0^n 1^n$

$$xy^i z, \forall i \geq 0$$

Consider all possible decomposition  $\Rightarrow$

$$\text{I. } y = 0^k \Rightarrow xy^i z = 0^{n-k} 0^k 1^n$$

$$xy^i z = 0^{n-k} 0^{i k} 1^n \notin L$$

$$\text{II. } y = 1^k \Rightarrow xy^i z = 0^n 1^k 1^{n-k}$$

$$xy^i z = 0^n 1^{i k} 1^{n-k} \notin L$$

$$\text{III. } y = 0^k 1^l \Rightarrow xyz = 0^{n-k} 0^k 1^l 1^{n-l}$$

$$xy^iz = 0^{n-k} (0^k 1^l)^i 1^{n-l} \notin L$$

$$\text{IV. } y = 0^k 1^k \Rightarrow xyz = 0^{n-k} 0^k 1^k 1^{n-k}$$

$$xy^iz = 0^{n-k} (0^k 1^k)^i 1^{n-k} \notin L$$

$\Rightarrow L$  is not regular

## Context Free Grammars (CFG)

- Productions of the form:  $A \rightarrow \alpha, A \in N, \alpha \in (N \cup \Sigma)^*$
- More powerful
- Can model programming language

### Syntax Tree:

Definition: A syntax tree corresponding to a CFG

$G = (N, \Sigma, P, S)$  is a tree obtained in the following way:

1. Root is the starting symbol  $S$
2. Nodes  $\in N \cup \Sigma$ :
  - internal nodes  $\in N$
  - leaves  $\in \Sigma$
3. For a node  $A$  the descendants in order from left to right are  $x_1, x_2, \dots, x_n$  only if  $A \rightarrow x_1 x_2 \dots x_n \in P$ .

### Remarks:

\* Parse tree = syntax tree (result of parsing)

\* Derivation tree - condition 2.2 not satisfied

\* Abstract syntax tree (AST)  $\neq$  syntax tree

Property: in a CFG  $G = (N, \Sigma, P, S)$ ,  $w \in L(G)$  if and only if there exists a syntax tree with frontier  $w$ .

Example:

$$S \rightarrow aSbS \mid c$$

$$w = aacbcbc$$

Leftmost derivation:

$$\begin{aligned} S &\Rightarrow aSbS \Rightarrow aSbSbS \Rightarrow aacbSbS \Rightarrow aacbcbS \Rightarrow \\ &\Rightarrow aacbcbc \end{aligned}$$

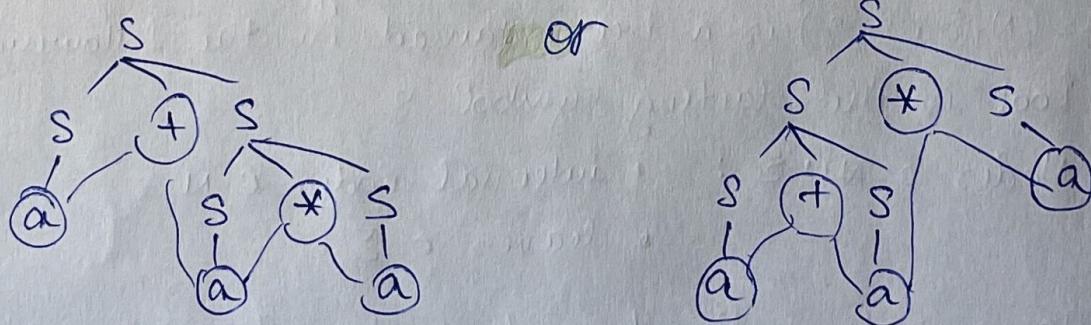
Rightmost derivation:

$$\begin{aligned} S &\Rightarrow aSbS \Rightarrow aSbc \Rightarrow aSbSbc \Rightarrow aSbcbC \Rightarrow \\ &\Rightarrow aacbcbc \end{aligned}$$

Definition: A cfg  $G = (N, \Sigma, P, S)$  is ambiguous if for a  $w \in L(G)$  there exists 2 distinct syntax tree with frontier  $w$ .

Example:  $S \rightarrow S+S \mid S * S \mid a$

$$w = a+a+a$$



Parsing modeled with ctg:

$$\text{cfg } G = (N, \Sigma, P, S)$$

- $N$  - nonterminal: syntact. constr: declaration, stmt, expression ...
- $\Sigma$  - terminals; elms of the lang: identifiers, const, reserved words, operators, separators
- $P$  - syntactical rules - expressed in BNF - simple trans
- $S$  - syntactical construct corrsp. to program

THEN

Program syntactical correct  $\Leftrightarrow w \in L(G)$ .

# Equivalent transformation of cfg.

## Unproductive symbols:

Algorithm 1: Elimination of unproductive symbols

Idea: build  $N_0, N_1, \dots$  recursively (until saturation)

A nonterminal  $A$  is **unproductive** in a cfg if it does not generate any word:

$$\{w \mid A \xrightarrow{*} w \in \Sigma^*\} = \emptyset$$

Step 1:  $N_0 = \emptyset, i=1$

Step 2:  $N_i = N_{i-1} \cup \{A \mid A \rightarrow \alpha \in P, \alpha \in (N_{i-1} \cup \Sigma)^*\}$

Step 3: if  $N_i \neq N_{i-1}$  then  $i=i+1$ ; go to step 2  
else  $N' = N_i$

Step 4: if  $S \notin N'$  then  $L(G) = \emptyset$   
else  $P' = \{A \mid A \rightarrow \alpha \in P \text{ and } A \in N'\}$

Example:  $G = (\{S, A, B, C, D\}, \{a, b, c\}, P, S)$

$$\begin{aligned} P: \quad & S \rightarrow aA \mid \underline{ac} \\ & A \rightarrow AB \mid \underline{aB} \\ & B \rightarrow b \\ & C \rightarrow ac \mid CD \\ & D \rightarrow b \end{aligned}$$

$$N_0 = \emptyset$$

$$N_1 = N_0 \cup \{B, D\} = \{B, D\}$$

$$N_2 = N_1 \cup \emptyset = \{B, D\} \cup \emptyset$$

$$\begin{aligned} & \{A\} \\ & A \rightarrow AB \\ & \in \Sigma \in N_1 \end{aligned}$$

$$N_3 = \{B, A, D\} \cup \{S\}$$

$$N_4 = N_3 \cup \emptyset$$

$\Rightarrow C$  - unproductive

## Inaccessible symbols :

Algorithm 2: Elimination of inaccessible symbols

Step 1 :  $V_0 = \{S\}, i=1$

Step 2 :  $V_i = V_{i-1} \cup \{x \mid \exists A \rightarrow \alpha \in P, A \in V_{i-1}\}$

Step 3 : if  $V_i \neq V_{i-1}$  then  $i = i+1$ , go to step 2

else  $N' = N \cap V_i$

$\Sigma' = \Sigma \cap V_i$

$P' = \{A \rightarrow \alpha \mid A \rightarrow \alpha \in P, A \in N', \alpha \in (N \cup \Sigma)^*\}$

Example :

$S \rightarrow aA$

$A \rightarrow AB \mid aB$

$B \rightarrow b$

~~$D \rightarrow b$~~

eliminate D

$V_0 = \{S\}$

$V_1 = \{S, a, A\}$

$V_2 = \{S, a, A, B\}$

$V_3 = \{S, a, A, B, b\}$

$V_4 = \{S, a, A, B, b\}$

$\Rightarrow D$  inaccessible

## $\epsilon$ -productions :

Algorithm 3: Elimination of  $\epsilon$ -productions

Step 1: construct

$\bar{N} = \{A \mid A \in N, A \xrightarrow{*} \epsilon\}$

1.a.  $N_0 = \{A \mid A \rightarrow \epsilon \in P\}$   
 $i=1$

1.b.  $N_i = N_{i-1} \cup \{A \mid A \rightarrow \alpha \in P, \alpha \in \bar{N}_{i-1}^*\}$

1.c. if  $N_i \neq N_{i-1}$  then  $i = i+1 \Rightarrow$  go to step 1b  
else  $\underline{\bar{N} = N_i}$

A symbol  $x \in N \cup \Sigma$  is **inaccessible** if  $x$  does not appear in any sentential form :

$\# S \xrightarrow{*} \alpha, x \notin \alpha$

A cfg  $G = (N, \Sigma, P, S)$  is **without  $\epsilon$ -productions** if

1.  $P \not\ni A \rightarrow \epsilon$  ( $\epsilon$ -production)
- OR
2.  $\exists S \rightarrow \epsilon \text{ s.t. } S \notin \text{rhs}(p), \forall p \in P$

Step 2: Let  $P'$  - set of productions built:

2.a. if  $A \rightarrow \alpha_0 B_1 \alpha_1 B_2 \dots B_k \alpha_k \in P$ ,  $k \geq 0$

and for  $i = \overline{1, k}$   $B_i \in \bar{N}$  and  $\alpha_i \notin \bar{N}$ ,  $j = \overline{0, k}$

$\Rightarrow$  add to  $P'$  all prod of the form

$A \rightarrow \alpha_0 X_1 \alpha_1 X_2 \alpha_2 \dots X_k \alpha_k$

where  $X_i$  is  $B_i$  or  $\epsilon$  (not  $A \rightarrow \epsilon$ )

2b. if  $S \in N'$  then add  $S'$  to  $N'$  and  $S' \rightarrow S \mid \epsilon$  to  $P'$   
else  $N' = N$  and  $S' = S$

Example:

$S \rightarrow aA \mid aAbB$

$A \rightarrow aA \mid B$

$B \rightarrow bB \mid \epsilon$

Step 1:

$N_0 = \{B\}$

$N_1 = \{B, A\}$

$N_2 = \{B, A\} \Rightarrow \bar{N} = \{B, A\}$

Step 2:

$S \rightarrow aA \mid a$

$S \rightarrow aAbB \mid abB \mid aAb \mid ab$

$A \rightarrow aA \mid a \mid B$

$B \rightarrow bB \mid b$

single productions:

Algorithm 4: Elimination  
of single productions

A production of the form  
 $A \rightarrow B$  is called single  
production or renaming  
rule

For each  $A \in N$  build the set  $N_A = \{B \mid A \xrightarrow{*} B\}$

1.a.  $N_0 = \{A\}$   $i = 1$

1.b.  $N_i = N_{i-1} \cup \{C \mid B \rightarrow C \in P \text{ and } B \in N_{i-1}\}$

2.c. if  $N_i \neq N_{i-1}$  then  $i = i + 1$  goto 1b  
else  $N_A = N_i$

$P'$ : for all  $A \in N$  do

for all  $B \in N_A$  do

if  $B \rightarrow \alpha \in P$  and not "single" then

$A \rightarrow \alpha \in P'$

# Parsing

- Cfg  $G = (N, \Sigma, P, S)$  check if  $w \in L(G)$

- Construct parse tree

How? 1- Top down vs Bottom up

## 2. Recursive vs linear

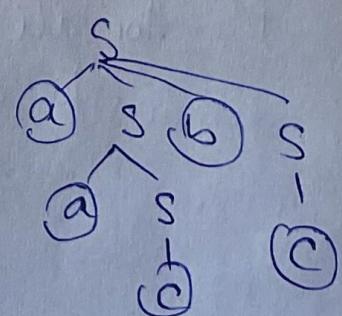
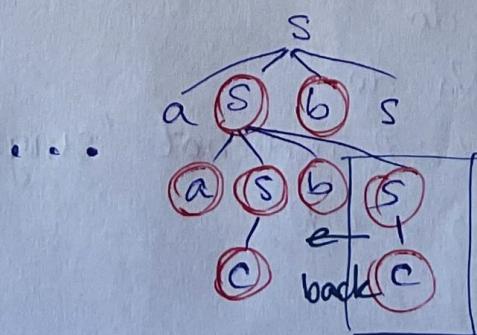
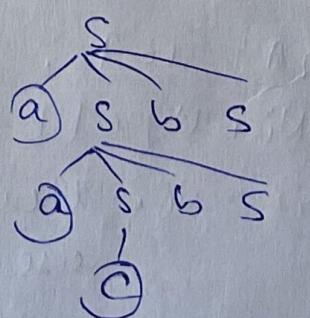
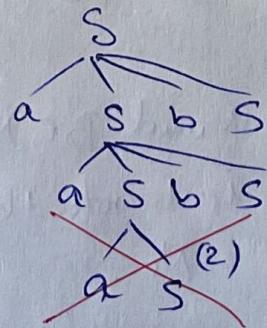
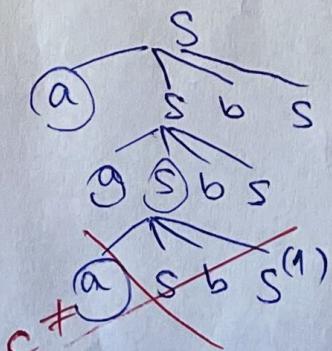
	Descendent	Ascendent
Recursive	Descendent recurs. parser	Ascendent recursive parser
Linear	$LL(k)$ : $LL(1)$	$LR(k)$ : $LR(0)$ , $SLR$ , $LR(1)$ , $LALR$

Result - parse tree representation

- **Arbitrary tree**: child nibbling repr (left child right sibl)
  - **Sequence of derivations**:  $S \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n = w$
  - **String of productions**: index of associated production  
(which production is used at each derivation step)

# Descendent recursive parser

Example  $S \rightarrow aSbS \mid aS^2 \mid C^3$   $w = aacbcb$



## Formal Model:

\* Configuration:  $(s, i, \alpha, \beta)$

where :

- $s$  = state of the parsing, can be
  - $g$  - normal state
  - $b$  - back state
  - $f$  - final state  $\Leftrightarrow$  success  $w \in L(G)$
  - $e$  - error state  $\Leftrightarrow$  insuccess  $w \notin L(G)$
- $i$  - position of current symbol in input sequence :  
 $w = a_1 a_2 \dots a_n, i \in \{1, \dots, n+1\}$
- $\alpha$  = working stack, stores the way the parse is build
- $\beta$  = input stack, part of the tree to be build

\* Expand : when head of input stack ( $\beta$ ) is a nonterminal

$$(g, i, \alpha, A\beta) \vdash (g, i, \alpha A_1, A_1\beta)$$

where:  $A \rightarrow A_1 \mid A_2 \mid \dots$  productions correspond to  $A$   
 $A_1$  - first prod of  $A$

\* Advance : when head of input stack ( $\beta$ ) is a terminal =  
 = current symbol from input

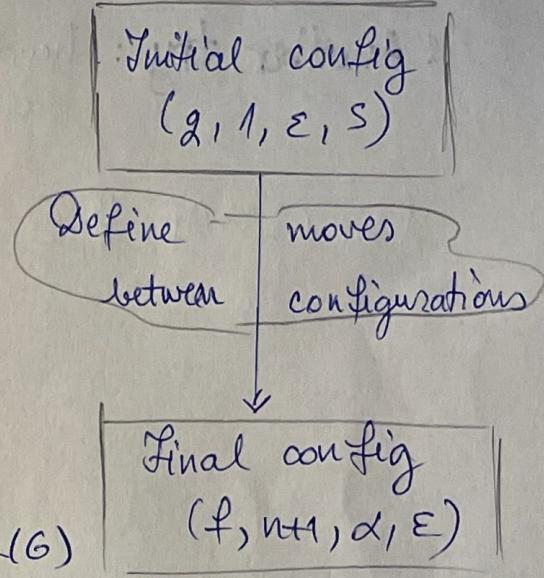
$$(g, i, \alpha, a_i\beta) \vdash (g, i+1, \alpha a_i, \beta)$$

\* Momentary insuccess : when head of the input stack  
 is a terminal  $\neq$  current symbol from input

$$(g, i, \alpha, a_i\beta) \vdash (b, i, \alpha, \beta)$$

\* Back : head of working stack ( $\alpha$ ) is a terminal

$$(b, i, \alpha a_i \beta) \vdash (b, i-1, \alpha, a_i \beta)$$



\* Another try: head of the working stack is a nonterminal

$(b, i, \alpha A_j, \gamma_j \beta) \leftarrow (g, i, \alpha A_{j+1}, \gamma_{j+1} \beta)$  if  $\exists A \rightarrow \gamma_{j+1}$   
 $(b, i, \alpha, \beta)$  otherwise with the exception  
 $(e, i, \alpha, \beta)$  if  $i=1, A=S$   
ERROR

\* Success:  $(g, m+1, \alpha, \epsilon) \leftarrow (f, m+1, \alpha, \epsilon)$

### Algorithm

$w \in L(G)$  - HOW

• Process  $\alpha$ :

- from left to right (reversed if stored as stack)
- skip terminal symbols
- Nonterminals - index of production

Example:  $\alpha = S_1 a S_2 a S_3 c b S_3 c$

When the algorithm never stops ???

$S \rightarrow S\alpha$  - expand indefinitely (left recursive)

# LL(1) Parser

$\text{First}_k \approx \text{first } k \text{ terminal symbols that can be generated from } \alpha$

Definition:

$$\text{FIRST}_k: (N \cup \Sigma)^* \rightarrow \mathcal{P}(\Sigma^k)$$

$$\text{FIRST}_k(\alpha) = \{ u \mid u \in \Sigma^k, \alpha \xrightarrow{*} ux, |u| = k \text{ and } \alpha \xrightarrow{*} u, |u| \leq k \}$$

Construct FIRST:

- If  $L_1, L_2$  are 2 languages over alphabet  $\Sigma$ , then

$$L_1 \oplus L_2 = \{ w \mid x \in L_1, y \in L_2, xy = w, |w| \leq 1 \text{ or} \\ \text{concatenation of} \\ \text{length 1.} \quad xy = w \text{, } |w| = 1 \}$$

$$\text{FIRST}(\alpha\beta) = \text{FIRST}(\alpha) \oplus \text{FIRST}(\beta)$$

$$\text{FIRST}(x_1 \dots x_n) = \text{FIRST}(x_1) \oplus \dots \oplus \text{FIRST}(x_n)$$

Example: Algorithm C6 slide 21 only terminals

$$E \rightarrow TE'$$

$$\text{First}(E) = \text{First}(T) = \{ (, \text{id}) \}$$

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

$$\text{First}(E') = \{ +, \epsilon \}$$

$$T \rightarrow FT'$$

$$\text{First}(T) = \text{First}(F) = \{ (, \text{id}) \}$$

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

$$\text{First}(T') = \{ *, \epsilon \}$$

$$F \rightarrow (E) \mid \text{id}$$

$$\text{First}(F) = \{ (, \text{id}) \}$$

FOLLOW:  $\text{FOLLOW}_k(A) = \text{next } k \text{ symbols generated after following } A$

$$\text{FOLLOW}: (N \cup \Sigma)^* \rightarrow \mathcal{P}(\Sigma)$$

$$\text{FOLLOW}(\beta) = \{ w \in \Sigma \mid s \xrightarrow{*} \alpha\beta w, w \in \text{FIRST}(\gamma) \}$$

Algorithm C6 slide 93

Example: Computing follow:

1. If A starting symbol, put  $\epsilon$  in FOLLOW(A)

2. Productions of the form  $B \rightarrow dA\beta$

$$\Rightarrow \text{FOLLOW}(A) = \text{FIRST}(\beta)$$

3. Productions of the form  $B \rightarrow dA$  or

$B \rightarrow dA\beta$  where  $\beta \Rightarrow^* \epsilon$

$$\text{Add } \text{FOLLOW}(A) = \text{FOLLOW}(B)$$

(a)  $E \rightarrow TE'$

$$\text{FOLLOW}(E) = \{ \epsilon \} \cup \{ \text{ } \} \xrightarrow{e \rightarrow 3}$$

(b)  $E' \rightarrow +TE' \mid \epsilon$

$$\text{FOLLOW}(E') = \text{FOLLOW}(E)$$

(c)  $T \rightarrow FT'$

$$\text{FOLLOW}(T) = \text{FIRST}(E') = \{ +, \epsilon \} \xrightarrow{b \rightarrow 2}$$

(d)  $T' \rightarrow *FT' \mid \epsilon$

$$\text{FOLLOW}(T') = \text{FOLLOW}(T)$$

(e)  $F \rightarrow (E) \mid \text{id}$

$$\text{FOLLOW}(F) = \text{FIRST}(T') = \{ *, \epsilon \} \xrightarrow{d \rightarrow 2}$$

$$E' \rightarrow +T \Rightarrow R_3 \Rightarrow \text{Follow}(T) = \text{Follow}(E')$$

$$T' \rightarrow *F \Rightarrow R_3 \Rightarrow \text{Follow}(F) = \text{Follow}(T')$$

$$\Rightarrow \text{Follow}(E) = \{ \epsilon, \} \}$$

$$\text{Follow}(E') = \{ \epsilon, \} \}$$

$$\text{Follow}(T) = \{ +, \epsilon \} \cup \text{Follow}(E') = \{ +, \epsilon, \} \}$$

$$\text{Follow}(T') = \{ +, \epsilon, \} \}$$

$$\text{Follow}(F) = \{ *, \epsilon \} \cup \text{Follow}(T') = \{ *, +, \epsilon, \} \}$$

Definition:

cfg is LL(k) if for any 2 leftmost derivations we have

1.  $S \xrightarrow{*} wA\alpha \Rightarrow w\beta\alpha \xrightarrow{*} w\gamma$

2.  $S \xrightarrow{*} wA\alpha \Rightarrow w\beta\alpha \xrightarrow{*} w\gamma$

such that  $\text{FIRST}_k(x) = \text{FIRST}_k(y) \Rightarrow \beta = \gamma$

### Theorem :

A grammar is  $LL(1)$  if and only if for any nonterm  $A$  with productions  $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$ ,

$\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$  and if  $\alpha_i \Rightarrow \epsilon$ ,

$\text{FIRST}(\alpha_i) \cap \text{FOLLOW}(A) = \emptyset$ ,  $\forall i, j = 1, n$ ,  $i \neq j$

### STEPS :

- 1) Construct  $\text{FIRST}$ ,  $\text{FOLLOW}$
- 2) Construct  $LL(1)$  parse table
- 3) Analyse sequence based on moves between configurations

### Step 2: Construct $LL(1)$ parse table

Add special char  $\$$  - marking for "empty stack"

#### Rules $LL(1)$ table :

1.  $M(A, a) = (\alpha, i)$   $\forall a \in \text{FIRST}(\alpha)$ ,  $a \neq \epsilon$

$A \rightarrow \alpha$  prod in  $P$  with nr  $i$

$M(A, b) = (\alpha, i)$  if  $\epsilon \in \text{FIRST}(\alpha)$ ,  $\forall b \in \text{FOLLOW}(A)$

$A \rightarrow \alpha$  production in  $P$  with nr  $i$

2.  $M(a, a) = \text{pop}$ ,  $\forall a \in \Sigma$

3.  $M($, $) = \text{acc}$

4.  $M(x, a) = \text{err}$

#### Remark :

A grammar is  $LL(1)$  if the  $LL(1)$  parse tree table does not contain conflicts - there exists at most one value in each cell of the table  $M(A, a)$

### Step 3: Define configurations and moves

Configuration:  $(\alpha, \beta, \pi)$   
where  $\alpha$  = input stack  
 $\beta$  = working stack  
 $\pi$  = output (result)

\* Initial configuration:  
 $(w\$, \$, \epsilon)$

\* Final configuration:  
 $(\$, \$, \pi)$

### Moves:

1. Push - put in stack

$(ux, A\alpha\$, \pi) \vdash (ux, \beta\alpha\$, \pi i)$   
if  $M(A, u) = (\beta, i)$   
(pop  $A$ , push symbols of  $\beta$ )

2. Pop - take off from stack (both stacks)

$(ux, \alpha\$ \pi) \vdash (x, \alpha\$, \pi)$  if  $M(\alpha, u) = \text{pop}$

3. Accept:  $(\$, \$, \pi) \vdash \text{acc}$

4. Error: otherwise

Algorithm: Input  $\rightarrow L(G)$  with no conflicts  
 $\rightarrow G$  - grammar  
 $\rightarrow$  sequence  $w = a_1 a_2 \dots a_n$

$\alpha = w\$$ ;  $\beta = \$\$$ ;  $\pi = \epsilon$

$go = \text{true}$

while  $go$  do:

if  $M(\text{head}(\beta), \text{head}(\alpha)) = (b, i)$  then

pop( $\beta$ ), push( $\beta, b$ ); addTo( $\pi, i$ )

else

if  $M(\text{head}(\beta), \text{head}(\alpha)) = \text{pop}$  then

pop( $\beta$ ); pop( $\alpha$ )

else

if  $M(\text{head}(\beta), \text{head}(\alpha)) = \text{acc}$  then

$go = \text{false}$ ;  $s = \text{acc}$

else

$go = \text{false}$ ,  $s = \text{err}$

print  $s$

# LR(0) Parser

L = left ( seq. is read from left to right)

R = right ( use rightmost derivations)

\* Enhanced grammar:

$$G = \{N, \Sigma, P, S\}$$

$$G' = \{N \cup \{S'\}, \Sigma, P \cup \{S' \rightarrow S\}, S'\}$$

!  $S'$  does not appear in any rhp. (right hand production)

EXAMPLE:  $G = \{S', S, A\}, \{a, b, c\}, P, S\}$

$$P: S' \rightarrow S$$

$$(1) S \rightarrow aA$$

$$(2) A \rightarrow bA$$

$$(3) A \rightarrow c$$

$$w = abbc$$

## I. Canonical collection of states:

Initialization:  $s_0 = \text{closure}(\{[S' \rightarrow .S]\}) = \{[S' \rightarrow .S]\} + \text{all productions of } S \text{ with } . \text{ in front of them.}$

### Rules:

CLOSURE: if  $.$  in front of nonterminal  $X \rightarrow$   
add all productions of  $X$  with  $.$  in front of  
their rhs

GOTO: search for all items with  $.$  in front of them,  
move  $.$  after and compute closure of (item  
with dot  $.$  after)

$$S_0 = \text{closure}(\{[S' \rightarrow .S]\}) = \{[S' \rightarrow .S], [S \rightarrow .aA]\}$$

$$S_1 = \text{goto}(S_0, S) = \text{closure}(\{[S' \rightarrow S.]\}) = \{[S' \rightarrow S.]\}$$

$$S_2 = \text{goto}(S_0, a) = \text{closure}(\{[S \rightarrow a.A]\}) = \{[S \rightarrow a.A], [A \rightarrow .bA], [A \rightarrow .c]\}$$

$$S_3 = \text{goto}(S_2, A) = \text{closure}(\{[S \rightarrow aA.]\}) = \{[S \rightarrow aA.]\}$$

$$S_4 = \text{goto}(S_2, b) = \text{closure}(\{[A \rightarrow b.A]\}) = \{[A \rightarrow b.A], [A \rightarrow .bA], [A \rightarrow .c]\}$$

$$S_5 = \text{goto}(S_2, c) = \text{closure}(\{[A \rightarrow c.]\}) = \{[A \rightarrow c.]\}$$

$$S_6 = \text{goto}(S_4, A) = \text{closure}(\{[A \rightarrow bA.]\}) = [A \rightarrow bA.]$$

$$\text{goto}(S_4, b) = \text{closure}(\{[A \rightarrow b.A]\}) = S_4$$

$$\text{goto}(S_4, c) = \text{closure}(\{[A \rightarrow c.]\}) = S_5$$

II Parsing table: reduce if . at the end.

\* Look at  $S_0$  it will be **shift**. and for the nonTerm.

with . in front look for  $\text{goto}(S_0, \text{nonTerm})$

and put the number of the state in the table

ex.  $\text{goto}(S_0, S) = S_1 \Rightarrow$  put 1 at **line 0**  
**column S**

\* If . at the end put **reduce + number of the production**

$S_5 = \{[A \rightarrow c.]\} \Rightarrow$  put reduce 3 at line 5

$\nearrow$   
(3)  $A \rightarrow C$

\* **Accept** if state =  $\{[S' \rightarrow S.]\}$

	ACTION	GOTO				
		a	b	c	s	A
0	shift	2			1	
1	accept.					
2	shift		4	5		3
3	reduce 1					
4	shift		4	5		6
5	reduce 3					
6.	reduce 2					

### III Parsing the sequence:

- ! workstack  $\rightarrow$  top at the right
- input stack  $\rightarrow$  top at the left
- output  $\rightarrow$  top at the left.

Always look at the top of the stacks

- \* Shift  $\rightarrow$  move the terminal from input stack to work stack and put the number of shift
- \* Reduce n  $\rightarrow$ 
  - \* look for production number n
  - \* take everything from rhs of prod n till the end of the workstack
  - \* put on the workstack the lhs of prod n and the number associate with (lhs, prev nr) from the table
  - \* add n to output band

work stack	input stack	output band
80	abbC \$	$\Sigma$
\$0a2	bbC \$	$\Sigma$
\$0a2b4	bc \$	$\Sigma$
\$0a2b4b4	c \$	$\Sigma$
\$0a2b4b4c5	\$	$\Sigma$
\$0a2b4b4A6	\$	3
\$0a2b4A6	\$	23
\$0a2A3	\$	223
\$0S1	\$	1223
accept	\$	1223

# SLR Parser

Example:  $G = (\{S', E, T\}, \{+, id, const, (, )\}, P, S')$

$$P: S' \rightarrow E$$

$$(1) E \rightarrow T$$

$$(2) E \rightarrow E + T$$

$$(3) T \rightarrow (E)$$

$$(4) T \rightarrow id$$

$$(5) T \rightarrow const$$

$$w = id + const$$

## I Canonical collection - SAME AS LR(0)

Rules:

- \* **CLOSURE:** If. in front of nonterminal  $X \Rightarrow$  add all productions of  $X$  with. in front of their rhs

\* **GOTO:** search for all items with . in front of them, compute closure (item with . AFTER)

$$S_0 = \text{closure}(\{[S' \rightarrow ., E]\}) = \{[S' \rightarrow ., E], [E \rightarrow ., T], [E \rightarrow ., E + T], [T \rightarrow ., (E)], [T \rightarrow ., id], [T \rightarrow ., const]\}$$

$$S_1 = \text{goto}(S_0, E) = \text{closure}(\{[S' \rightarrow E.], [E \rightarrow E. + T]\}) = \{[S' \rightarrow E.], [E \rightarrow E. + T]\}$$

$$S_2 = \text{goto}(S_0, T) = \text{closure}(\{[E \rightarrow T.]\}) = \{[E \rightarrow T.]\}$$

$$S_3 = \text{goto}(S_0, ( ) = \text{closure}(\{[T \rightarrow (.)]\}) = \{[T \rightarrow (.)], [E \rightarrow ., T] \\ [E \rightarrow ., E + T], [T \rightarrow ., (E)], [T \rightarrow ., id], [T \rightarrow ., const]\}$$

$$S_4 = \text{goto}(S_0, id) = \text{closure}(\{[T \rightarrow id.]\}) = \{[T \rightarrow id.]\}$$

$$S_5 = \text{goto}(S_0, const) = \text{closure}(\{[T \rightarrow const.]\}) = \{[T \rightarrow const.]\}$$

$$S_6 = \text{goto}(S_1, +) = \text{closure}(\{[E \rightarrow E + ., T]\}) = \{[E \rightarrow E + ., T], [T \rightarrow ., (E)], [T \rightarrow ., id], [T \rightarrow ., const]\}$$

$$S_7 = \text{goto}(S_3, E) = \text{closure}(\{[T \rightarrow (E.)], [E \rightarrow E. + T]\}) = \{[T \rightarrow (E.)], [E \rightarrow E. + T]\}$$

$$\text{goto}(s_3, T) = \text{closure}(\{[E \rightarrow T]\}) = \{[E \rightarrow T]\} = s_2$$

$$\text{goto}(s_3, ( ) ) = \text{closure}(\{[T \rightarrow (.)]\}) = s_3$$

$$\text{goto}(s_3, \text{id}) = \text{closure}(\{[T \rightarrow \text{id}]\}) = s_4$$

$$\text{goto}(s_3, \text{const}) = \text{closure}(\{[T \rightarrow \text{const}]\}) = s_5$$

$$s_8 = \text{goto}(s_6, T) = \text{closure}(\{[E \rightarrow E + T]\}) = \{[E \rightarrow E + T]\}$$

$$\text{goto}(s_6, ( ) ) = \text{closure}(\{[T \rightarrow (.)]\}) = s_3$$

$$\text{goto}(s_6, \text{id}) = \text{closure}(\{[T \rightarrow \text{id}]\}) = s_4$$

$$\text{goto}(s_6, \text{const}) = \text{closure}(\{[T \rightarrow \text{const}]\}) = s_5$$

$$g = \text{goto}(s_7, ( ) ) = \text{closure}(\{[T \rightarrow (.)]\}) = \{[T \rightarrow (.)]\}$$

$$\text{goto}(s_7, +) = \text{closure}(\{[E \rightarrow E + .]\}) = s_6$$

Compute FOLLOW rules :

1. If A starting symbol, put  $\epsilon$  in FOLLOW(A)

2. Productions of the form  $B \rightarrow \alpha A \beta$  stop

$$\Rightarrow \text{FOLLOW}(A) = \text{FIRST}(\beta)$$

3. Productions of the form  $B \rightarrow \alpha A$  or  
 $B \rightarrow \alpha A \beta$  with  $\beta \neq \epsilon$

$\Rightarrow$  Add to  $\text{FOLLOW}(A) \leftarrow \text{FOLLOW}(\beta)$

$$\text{FOLLOW}(E) = \{\epsilon, +, (\ )\}$$

$$\text{FOLLOW}(T) = \{\epsilon, +, (\ )\}$$

## II. Parsing table :: M2

ACTION - column for every terminal + \$

\* Shift : . in front of terminal  $\Rightarrow$  shift + nr of state from goto (current state, terminal)

\* Reduce : . at the end  $\Rightarrow$  reduce + nr of production to all elems from FOLLOW of nonterminal

\* Accept : only if . after starting symbol

GOTO : - column for every nonterminals

• in front of nonterminal  $x \Rightarrow$  put number of state that has goto (current state,  $x$ )

	ACTION							GOTO	
	+	(	)	id	const	\$	E	T	
0		Shift 3		Shift 4	Shift 5		1	2	
1	Shift 6					accept			
2	Reduce 1		Reduce 1			Reduce 1			
3		Shift 3		Shift 4	Shift 5		4	2	
4	Reduce 4		Reduce 4			Reduce 4			
5	Reduce 5		Reduce 5			Reduce 5			
6		Shift 3		Shift 4	Shift 5			8	
7	Shift 6		Shift 9						
8	Reduce 2		Reduce 2			Reduce 2			
9	Reduce 3		Reduce 3			Reduce 3			

### III Parsing the sequence: SAME AS LR(0)

\* **Shift** → move terminal from input stack to work stack and put number of shift

\* **Reduce n**: \* look for production nr n

\* take everything from rhs of prod n till the end of workstack

\* put on the work stack the lhs of prod n and the number associated with (lhs, previous nr) from the table

\* add n to output band

work stack	input stack	output band
\$ 0	id + const \$	$\Sigma$
\$ 0 id 4	+ const \$	$\Sigma$
\$ 0 T 2 3	+ const \$	4 ACTION
\$ 0 E 1 3	+ const \$	14
\$ 0 E 1 + 6	const \$	14
\$ 0 E 1 + 6 const 5	\$	14
\$ 0 E 1 + 6 T 8	\$	514
\$ 0 E 1	\$	2514
accept		
$E \xrightarrow{2} E + T \xrightarrow{5} E + \text{const} \xrightarrow{1} T + \text{const} \xrightarrow{4}$		
$\Rightarrow id + \text{const}$		

# Push-down automata (PDA)

Definition : A PDA is a 7-tuple  $M = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$

where :

- $Q$  - finite set of states
- $\Sigma$  - alphabet (input symbols)
- $\Gamma$  - stack alphabet
- $\delta : Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow P(Q \times \Gamma^*)$  - transition function
- $q_0 \in Q$  - initial state
- $z_0 \in \Gamma$  - initial stack symbol
- $F \subseteq Q$  - set of final states

\* A transition is determined by :

- current state
- current input symbol
- head of stack

\* Reading head  $\rightarrow$  input band :  $\begin{cases} \text{read one symbol} \\ \text{OR} \\ \text{no action} \end{cases}$

\* Stack :  $\rightarrow$  zero symbols  $\Rightarrow$  pop  
 $\rightarrow$  one symbol  $\Rightarrow$  push  
 $\rightarrow$  several symbols  $\Rightarrow$  repeat push

Configuration :  $(q, x, \alpha)$  where . PDA is in state  $q$   
• input band contains  $x$   
• head of stack is  $\alpha$

Initial configuration :  $(q_0, x, z_0)$

## Moves:

1. If  $\delta(q, a, \lambda) \ni (p, \gamma)$

$$\Rightarrow (q, aw, \lambda x) \xrightarrow{*} (p, w, \gamma x)$$

2. If  $\delta(q, \epsilon, \lambda) \ni (p, \gamma)$  ε-move

$$\Rightarrow (q, aw, \lambda x) \xrightarrow{*} (p, aw, \gamma x)$$

Language accepted by PDA :

- Empty stack-principle:

$$L_{\epsilon}(M) = \{ w \mid w \in \Sigma^*, (q_0, w, z_0) \xrightarrow{*} (q, \epsilon, \epsilon) \}$$

- Final state principle:

$$L_f(M) = \{ w \mid w \in \Sigma^*, (q_0, w, z_0) \xrightarrow{*} (q_f, \epsilon, \epsilon), q_f \in F \}$$

Example:

Construct the PDA for  $L = \{0^n 1^n \mid n \geq 1\}$

HOW?

For states:

- initial state:  $q_0$  - beginning and process only 0
- when first symbol 1 is found - move to new state  $\rightarrow q_1$
- final state:  $q_2$

For stack:

- initial symbol:  $\lambda_0$
- X - to count symbols:
  - \* when reading 0 - push X in stack
  - \* when reading 1 - pop X from stack

$$M = (\{q_0, q_1, q_2\}, \{0, 1\}, \{z_0, x, \delta\}, q_0, z_0, \{q_2\})$$

$$\delta(q_0, 0, z_0) = (q_0, X z_0)$$

$$\delta(q_0, 0, x) = (q_0, xz_0)$$

$$\delta(q_0, 1, x) = (q_1, \epsilon)$$

$$\delta(q_1, 1, x) = (q_1, \epsilon)$$

$$\delta(q_1, \epsilon, z_0) = (q_2, z_0) \text{ or } \delta(q_1, \epsilon, z_0) = (q_1, \epsilon)$$

$w = 0011$  - is it accepted by the language?

$$\begin{aligned}
 & (q_0, 0011, z_0) \xrightarrow{} (q_0, 011, X z_0) \xrightarrow{} (q_0, 11, XX z_0) \xrightarrow{\text{final state}} \\
 & \xrightarrow{} (q_1, 11, XX z_0) \xrightarrow{} (q_1, 1, X z_0) \xrightarrow{} (q_1, \epsilon, z_0) \xrightarrow{} (q_2, \epsilon, z_0) \\
 \Rightarrow 0011 \in L_f(M). & \qquad \qquad \qquad \xrightarrow{\text{empty stack}} \xrightarrow{} (q_1, \epsilon, \epsilon)
 \end{aligned}$$

### Properties:

**Theorem 1:** For any PDA  $M$ , there exists a PDA  $M'$  such that  $L_\epsilon(M) = L_f(M')$

**Theorem 2:** For any PDA  $M$ , there exists a context free grammar such that  $L_\epsilon(M) = L(G)$

**Theorem 3:** For any context free grammar there exists a PDA  $M$  such that  $L(G) = L_\epsilon(M)$ .

# Semantic analysis

- \* Parsing  $\Rightarrow$  syntax tree (ST)
  - \* Simplification  $\Rightarrow$  abstract ST (AST)
  - \* Annotated abstract ST (AAST)
    - ↳ attach semantic info in tree nodes

## Example

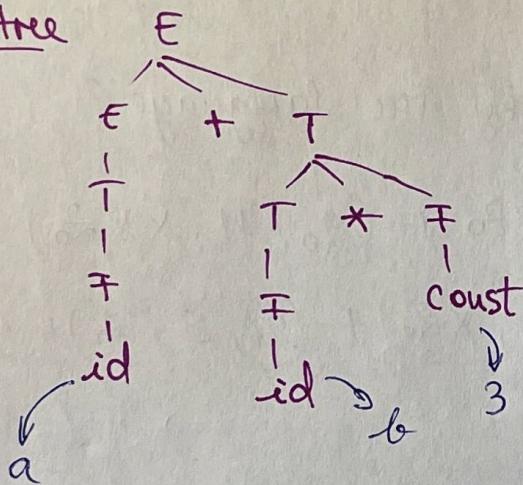
$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * \top \mid \top$$

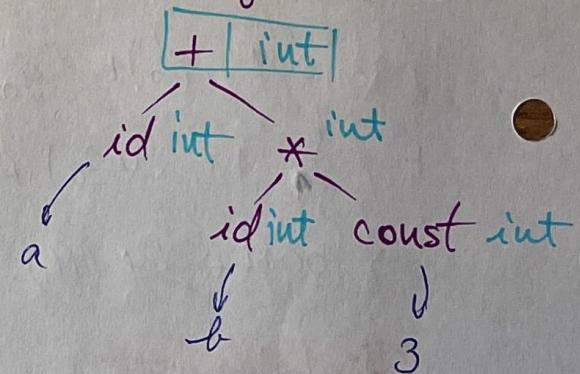
$\vdash \rightarrow (\vdash) \mid \text{id} \mid \text{const}$

parse : a + b \* 3

## Syntax tree



## Abstract syntax tree :



## Adnotated abstr. ST

-ex: type

- \* Attach meanings to syntactical constructions of a program
  - \* Constructions :
    - identifiers : values / how to be evaluated
    - statements : how to be executed
    - declaration : det. space to be allocated and location to be stored
  - \* Examples : type checking, verify properties.
  - \* How :
    - attribute grammars
    - manual methods

# Attribute grammar

- \* For nonterminal associate attributes (set of attributes)  
 $\uparrow x \in N \cup \Sigma : A(x)$
- \* For productions associate evaluation rules (set of rules)  
 $\uparrow p \in P : R(p)$

Definition:  $AG = (G, A, R)$  is called attribute grammar where

- \*  $G = (N, \Sigma, P, S)$  is a context free grammar
- \*  $A = \{A(x) \mid x \in N \cup \Sigma\}$  - finite set of attributes
- \*  $R = \{R(p) \mid p \in P\}$  - finite set of rules to compute/evaluate attributes

- \* Example → Compute the value of a nr ~~from~~ base 2, to base 10.

2. Write the rules:

$$G = \{N, B\}, \{0, 1\}, P, N\}$$

$$P: N \rightarrow NB$$

$$N \rightarrow B$$

$$B \rightarrow 0$$

$$B \rightarrow 1$$

$$N_1.V = 2 * N_2.V + B.V$$

$$N.V = B.V$$

$$B.V = 0$$

$$B.V = 1$$

- 1. Choose as attribute value of nr = v

→ Synthesized attributes:  $A(lhp)$  depends on rhp

→ Inherited attributes:  $A(rhp)$  depends on lhp.

- \* Special cases of AG :

- L-attribute grammars: for any node the depending attrs. are on the left; incorporated in top-down parser  $LL(1)$
- S-attribute grammars: synthesized attributes incorporated in bottom-up parser (LR).

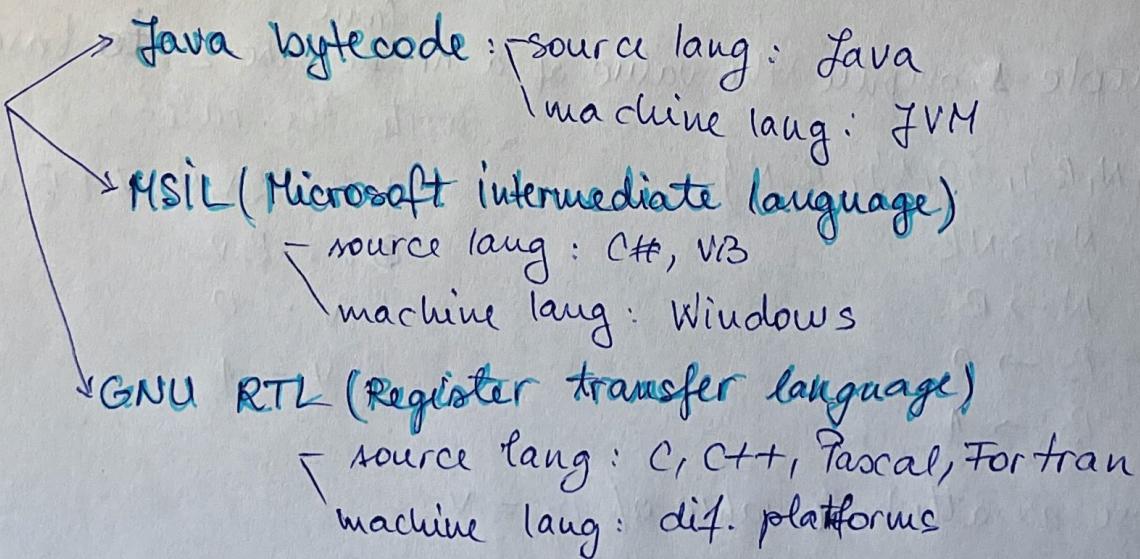
## Steps :

1. What? - decide what you want to compute (type, value)
2. Decide attributes:
  - \* how many
  - \* which attribute is defined for which symbol
3. Attach evaluation rules: for each production - which rule / rules.

More examples course 11 slide 23.

## Generate intermediary code

\* Forms of intermediary code:



\* Representations:

Annotated syntax tree: is generated in semantic analysis

Polish postfix form: operators appear in the order of execution

ex:  $\text{expr} = a + b * c$       PPF = abc\*+

$\text{expr} = a * b + c$       PPF = ab\*c+

3 address code

### 3 address code

Def: Sequence of simple format statements, close to object code, with the following general form:

$$\langle \text{result} \rangle = \langle \text{arg 1} \rangle \langle \text{op} \rangle \langle \text{arg 2} \rangle$$

Represented as:

\* **Quadruples**:  $\langle \text{op} \rangle \langle \text{arg 1} \rangle \langle \text{arg 2} \rangle \langle \text{result} \rangle$

\* **Triples**:  $\langle \text{op} \rangle \langle \text{arg 1} \rangle \langle \text{arg 2} \rangle$  (consider that the triple is storing the result)

\* **Indirected Triples**:

Example:  $b * b - 4 * a * c$ .

Quadruples:

op	arg 1	arg 2	rez
*	b	b	$t_1$
*	4	a	$t_2$
*	$t_2$	c	$t_3$
-	$t_1$	$t_3$	$t_4$

Triples

nr	op	arg 1	arg 2
(1)	*	b	b
(2)	*	4	a
(3)	*	(2)	c
(4)	-	(1)	(3)

### Optimize intermediary code

- Local optimizations:

- \* perform computation at compile time - constant values
- \* eliminate redundant computations
- \* eliminate inaccessible code

- Loop optimizations:

- \* factorization of loop invariants
- \* reduce power of operations

See examples  
course 12  
slide 11

# Generate object code

Def: Translate intermediary code statements into statements of object code (machine language)

- 2 aspects
- Register allocation: way in which variables are stored and manipulated
  - Instruction selection: way and order in which the intermediary code statements are mapped to machine instructions

## Computer with accumulator

- \* A stack machine consists of a stack (for storing and manipulating values) and 2 type of statements:
  - move and copy values in and from head of stack mem.
  - operations on stack head : operands are popped from stack, execute operation and then put the result in stack
- \* Accumulator : to execute operation
- \* Stack to store subexpressions and results.  
See example course 12 slide 17

## Computer with registers - registers + memory

- \* Instructions:
  - LOAD  $v, R$  : load value  $v$  in register  $R$
  - STORE  $R, v$  : put value  $v$  from register  $R$  in memory
  - ADD  $R_1, R_2$  : add value from  $R_1$  to the value from  $R_2$  and store the result in  $R_1$

### Remarks :

\* A register can be available or occupied  $\Rightarrow$

$\text{VAR}(R)$  = set of variables whose values are stored in register R

\* For every variable, the place (register, stack or memory) in which the current value of it exists  $\Rightarrow$

$\text{MEM}(x)$  = set of locations in which the value of variable x exists (will be stored in symbol table)

More info course 12 slide 20.