

## 2. Automata and Formal Languages

Automata and formal languages play an important role in projecting and realizing compilers. In the first section grammars and formal languages are defined. The different grammars and languages are discussed based on Chomsky hierarchy. In the second section we deal in detail with the finite automata and the languages accepted by them, while in the third section the pushdown automata and the corresponding accepted languages are discussed. Finally, references from a rich bibliography are given.

### 2.1. Languages and grammars

A finite and nonempty set of symbols is called an *alphabet*. The elements of an alphabet are *letters*, but sometimes are named also *symbols*. E.g. the set  $\Sigma = \{a, b, c, d, 0, 1, \sigma\}$  is an alphabet, with the letters  $a, b, c, d, 0, 1$  and  $\sigma$ .

With the letters of an alphabet words are composed. If  $a_1, a_2, \dots, a_n \in \Sigma, n \geq 0$ , then  $a_1 a_2 \dots a_n$  a  $\Sigma$  is a *word* over the alphabet  $\Sigma$  (the letters  $a_i$  are not necessary distinct). The number of letters of a word, with their multiplicities, constitutes the *length* of the word. If  $w = a_1 a_2 \dots a_n$ , then the length of  $w$  is  $|w| = n$ . If  $n = 0$ , then the word is an *empty word*, which will be denoted by  $\varepsilon$  (sometimes  $\lambda$  in other books). The set of words over the alphabet  $\Sigma$  will be denoted by  $\Sigma^*$ :

$$\Sigma^* = \{a_1 a_2 \dots a_n \mid a_1, a_2, \dots, a_n \in \Sigma, n \geq 0\} .$$

For the set of nonempty words over  $\Sigma$  the notation  $\Sigma^+ = \Sigma^* \setminus \{\varepsilon\}$  will be used. The set of words of length  $n$  over  $\Sigma$  will be denoted by  $\Sigma^n$ , and  $\Sigma^0 = \{\varepsilon\}$ . Then

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \dots \cup \Sigma^n \cup \dots \quad \text{and} \quad \Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \dots \cup \Sigma^n \cup \dots .$$

The words  $u = a_1 a_2 \dots a_m$  and  $v = b_1 b_2 \dots b_n$  are equal (i.e.  $u = v$ ), if  $m = n$  and  $a_i = b_i, i = 1, 2, \dots, n$ .

We define in  $\Sigma^*$  the binary operation called *concatenation*. The concatenation (or product) of the words  $u = a_1 a_2 \dots a_m$  and  $v = b_1 b_2 \dots b_n$  is the word  $uv = a_1 a_2 \dots a_m b_1 b_2 \dots b_n$ . It is clear that  $|uv| = |u| + |v|$ . This operation is associative but not commutative. Its neutral element is  $\varepsilon$ , because  $\varepsilon u = u\varepsilon = u$  for all  $u \in \Sigma^*$ .  $\Sigma^*$  with the concatenation is a monoid.

We introduce the power operation. If  $u \in \Sigma^*$ , then  $u^0 = \varepsilon$ , and  $u^n = u^{n-1}u$  for  $n \geq 1$ .

The **reversal** (or **mirror image**) of the word  $u = a_1a_2 \dots a_n$  is  $u^{-1} = a_n a_{n-1} \dots a_1$ . The reversal of  $u$  sometimes is denoted by  $u^R$  or  $\tilde{u}$ . It is clear that  $(u^{-1})^{-1} = u$  and  $(uv)^{-1} = v^{-1}u^{-1}$ .

Word  $v$  is a **prefix** of the word  $u$  if there exists a word  $z$  such that  $u = vz$ . If  $z \neq \varepsilon$  then  $v$  is a proper prefix of  $u$ . Similarly  $v$  is a **suffix** of  $u$  if there exists a word  $x$  such that  $u = xv$ . The proper suffix can also be defined. Word  $v$  is a **subword** of the word  $u$  if there are words  $p$  and  $q$  such that  $u = pvq$ . If  $pq \neq \varepsilon$  then  $v$  is a **proper subword**.

A subset  $L$  of  $\Sigma^*$  is called a **language** over the alphabet  $\Sigma$ . Sometimes this is called a **formal language** because the words are here considered without any meanings. Note that  $\emptyset$  is the empty language while  $\{\varepsilon\}$  is a language which contains the empty word.

### 2.1.1. Operations on languages

If  $L, L_1, L_2$  are languages over  $\Sigma$  we define the following operations

- **union**

$$L_1 \cup L_2 = \{u \in \Sigma^* \mid u \in L_1 \text{ or } u \in L_2\},$$

- **intersection**

$$L_1 \cap L_2 = \{u \in \Sigma^* \mid u \in L_1 \text{ and } u \in L_2\},$$

- **difference**

$$L_1 \setminus L_2 = \{u \in \Sigma^* \mid u \in L_1 \text{ and } u \notin L_2\},$$

- **complement**

$$\bar{L} = \Sigma^* \setminus L,$$

- **multiplication**

$$L_1 L_2 = \{uv \mid u \in L_1, v \in L_2\},$$

- **power**

$$L^0 = \{\varepsilon\}, \quad L^n = L^{n-1}L, \text{ if } n \geq 1,$$

- **iteration** or **star operation**

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L \cup L^2 \cup \dots \cup L^i \cup \dots,$$

- **mirror**

$$L^{-1} = \{u^{-1} \mid u \in L\}$$

We will use also the notation  $L^+$

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L \cup L^2 \cup \dots \cup L^i \cup \dots$$

The union, product and iteration are called **regular operations**.

### 2.1.2. Specifying languages

Languages can be specified in several ways. For example a language can be specified using

- 1) the enumeration of its words,
- 2) a property, such that all words of the language have this property but other word have not,
- 3) a grammar.

#### Specifying languages by listing their elements

For example the following are languages

$$L_1 = \{\varepsilon, 0, 1\},$$

$$L_2 = \{a, aa, aaa, ab, ba, aba\}.$$

Even if we cannot enumerate the elements of an infinite set infinite languages can be specified by enumeration if after enumerating the first some elements we can continue the enumeration using a rule. The following is such a language

$$L_3 = \{\varepsilon, ab, aabb, aaabbb, aaaabbbb, \dots\}.$$

### Specifying languages by properties

The following sets are languages

$$L_4 = \{a^n b^n \mid n = 0, 1, 2, \dots\},$$

$$L_5 = \{uu^{-1} \mid u \in \Sigma^*\},$$

$$L_6 = \{u \in \{a, b\}^* \mid n_a(u) = n_b(u)\},$$

where  $n_a(u)$  denotes the number of letters  $a$  in word  $u$  and  $n_b(u)$  the number of letters  $b$ .

### Specifying languages by grammars

Define the *generative grammar* or shortly the *grammar*.

**Definition 2.1** A *grammar* is an ordered quadruple  $G = (N, T, P, S)$ , where

- $N$  is the alphabet of **variables** (or **nonterminal symbols**),
- $T$  is the alphabet of **terminal symbols**, where  $N \cap T = \emptyset$ ,
- $P \subseteq (N \cup T)^* N (N \cup T)^* \times (N \cup T)^*$  is a finite set, that is  $P$  is the finite set of **productions** of the form  $(u, v)$ , where  $u, v \in (N \cup T)^*$  and  $u$  contains at least a nonterminal symbol,
- $S \in N$  is the **start symbol**.

*Remarks.* Instead of the notation  $(u, v)$  sometimes  $u \rightarrow v$  is used.

In the production  $u \rightarrow v$  or  $(u, v)$  word  $u$  is called the left-hand side of the production while  $v$  the right-hand side. If for a grammar there are more than one production with the same left-hand side, then these production

$$u \rightarrow v_1, u \rightarrow v_2, \dots, u \rightarrow v_r \quad \text{can be written as} \quad u \rightarrow v_1 \mid v_2 \mid \dots \mid v_r .$$

We define on the set  $(N \cup T)^*$  the relation called **direct derivation**

$$u \Longrightarrow v, \quad \text{if} \quad u = p_1 p p_2, \quad v = p_1 q p_2 \quad \text{and} \quad (p, q) \in P .$$

In fact we replace in  $u$  an appearance of the subword  $p$  by  $q$  and we get  $v$ . Another notations for the same relation can be  $\vdash$  or  $\models$ .

If we want to emphasize the used grammar  $G$ , then the notation  $\Longrightarrow$  can be replaced by  $\Longrightarrow_G$ . Relation  $\Longrightarrow$  is the reflexive and transitive closure of  $\Longrightarrow$ , while  $\Longrightarrow^+$  denotes its transitive closure. Relation  $\Longrightarrow^*$  is called a **derivation**.

From the definition of a reflexive and transitive relation we can deduce the following:  $u \Longrightarrow^* v$ , if there exist the words  $w_0, w_1, \dots, w_n \in (N \cup T)^*$ ,  $n \geq 0$  and  $u = w_0$ ,  $w_0 \Longrightarrow w_1$ ,  $w_1 \Longrightarrow w_2$ ,  $\dots, w_{n-1} \Longrightarrow w_n$ ,  $w_n = v$ . This can be written shortly  $u = w_0 \Longrightarrow w_1 \Longrightarrow w_2 \Longrightarrow \dots \Longrightarrow w_{n-1} \Longrightarrow w_n = v$ . If  $n = 0$  then  $u = v$ . The same way we can define the relation  $u \Longrightarrow^+ v$  except that  $n \geq 1$  always, so at least one direct derivation will be used.

**Definition 2.2** The **language generated** by grammar  $G = (N, T, P, S)$  is the set

$$L(G) = \{u \in T^* \mid S \xRightarrow{*} u\}.$$

So  $L(G)$  contains all words over the alphabet  $T$  which can be derived from the start symbol  $S$  using the productions from  $P$ .

**Example 2.1** Let  $G = (N, T, P, S)$  where

$$N = \{S\},$$

$$T = \{a, b\},$$

$$P = \{S \rightarrow aSb, S \rightarrow ab\}.$$

It is easy to see that  $L(G) = \{a^n b^n \mid n \geq 1\}$  because

$$S \xRightarrow{G} aSb \xRightarrow{G} a^2Sb^2 \xRightarrow{G} \dots \xRightarrow{G} a^{n-1}Sb^{n-1} \xRightarrow{G} a^n b^n,$$

where up to the last but one replacement the first production ( $S \rightarrow aSb$ ) was used, while at the last replacement the production  $S \rightarrow ab$ . This derivation can be written  $S \xRightarrow{*} a^n b^n$ . Therefore  $a^n b^n$  can be derived from  $S$  for all  $n$  and no other words can be derived from  $S$ .

**Definition 2.3** Two grammars  $G_1$  and  $G_2$  are **equivalent**, and this is denoted by  $G_1 \cong G_2$  if  $L(G_1) = L(G_2)$ .

**Example 2.2** The following two grammars are equivalent because both of them generate the language  $\{a^n b^n c^n \mid n \geq 1\}$ .

$G_1 = (N_1, T, P_1, S_1)$ , where

$$N_1 = \{S_1, X, Y\}, \quad T = \{a, b, c\},$$

$$P_1 = \{S_1 \rightarrow abc, S_1 \rightarrow aXbc, Xb \rightarrow bX, Xc \rightarrow Ybcc, bY \rightarrow Yb, aY \rightarrow aaX, aY \rightarrow aa\}.$$

$G_2 = (N_2, T, P_2, S_2)$ , where

$$N_2 = \{S_2, A, B, C\},$$

$$P_2 = \{S_2 \rightarrow aS_2BC, S_2 \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab, bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}.$$

First let us prove by mathematical induction that for  $n \geq 2$   $S_1 \xRightarrow{*}_{G_1} a^{n-1}Yb^n c^n$ . If  $n = 2$  then

$$S_1 \xRightarrow{G_1} aXbc \xRightarrow{G_1} abXc \xRightarrow{G_1} abYbcc \xRightarrow{G_1} aYb^2c^2.$$

The inductive hypothesis is  $S_1 \xRightarrow{*}_{G_1} a^{n-2}Yb^{n-1}c^{n-1}$ . We use production  $aY \rightarrow aaX$ , then  $(n-1)$  times production  $Xb \rightarrow bX$ , and then production  $Xc \rightarrow Ybcc$ , afterwards again  $(n-1)$  times production  $bY \rightarrow Yb$ . Therefore

$$S_1 \xRightarrow{G_1} a^{n-2}Yb^{n-1}c^{n-1} \xRightarrow{G_1} a^{n-1}Xb^{n-1}c^{n-1} \xRightarrow{G_1} a^{n-1}b^{n-1}Xc^{n-1} \xRightarrow{G_1} a^{n-1}b^{n-1}Ybc^n \xRightarrow{G_1} a^{n-1}Yb^n c^n.$$

If now we use production  $aY \rightarrow aa$  we get  $S_1 \xRightarrow{*}_{G_1} a^n b^n c^n$  for  $n \geq 2$ , but  $S_1 \xRightarrow{G_1} abc$  by the production  $S_1 \rightarrow abc$ , so  $a^n b^n c^n \in L(G_1)$  for any  $n \geq 1$ . We have to prove in addition that using the productions of the grammar we cannot derive only words of the form  $a^n b^n c^n$ . It is easy to see that a successful derivation (which ends in a word containing only terminals) can be obtained only in the presented way.

Similarly for  $n \geq 2$

$$S_2 \xRightarrow{G_2} aS_2BC \xRightarrow{G_2} a^{n-1}S_2(BC)^{n-1} \xRightarrow{G_2} a^n(BC)^n \xRightarrow{G_2} a^n B^n C^n$$

$$\xRightarrow{G_2} a^n b B^{n-1} C^n \xRightarrow{G_2}^* a^n b^n C^n \xRightarrow{G_2} a^n b^n c C^{n-1} \xRightarrow{G_2}^* a^n b^n c^n .$$

Here orderly were used the productions  $S_2 \rightarrow aS_2BC$  ( $n-1$  times),  $S_2 \rightarrow aBC$ ,  $CB \rightarrow BC$  ( $n-1$  times),  $aB \rightarrow ab$ ,  $bB \rightarrow bb$  ( $n-1$  times),  $bC \rightarrow bc$ ,  $cC \rightarrow cc$  ( $n-1$  times). But  $S_2 \xRightarrow{G_2} aBC \xRightarrow{G_2} abC \xRightarrow{G_2} abc$ , So  $S_2 \xRightarrow{G_2}^* a^n b^n c^n$ ,  $n \geq 1$ . It is also easy to see than other words cannot be derived using grammar  $G_2$ .

The grammars

$$G_3 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow \varepsilon\}, S) \text{ and}$$

$$G_4 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$$

are not equivalent because  $L(G_3) \setminus \{\varepsilon\} = L(G_4)$ .

**Theorem 2.4** *Not all languages can be generated by grammars.*

**Proof.** We encode grammars for the proof as words on the alphabet  $\{0, 1\}$ . For a given grammar  $G = (N, T, P, S)$  let  $N = \{S_1, S_2, \dots, S_n\}$ ,  $T = \{a_1, a_2, \dots, a_m\}$  and  $S = S_1$ . The encoding is the following:

$$\text{the code of } S_i \text{ is } \underbrace{10 \ 11 \ \dots \ 11 \ 01}_{i \text{ times}}, \quad \text{the code of } a_i \text{ is } \underbrace{100 \ 11 \ \dots \ 11 \ 001}_{i \text{ times}} .$$

In the code of the grammar the letters are separated by 000, the code of the arrow is 0000, and the productions are separated by 00000.

It is enough, of course, to encode the productions only. For example, consider the grammar

$$G = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S).$$

The code of  $S$  is 10101, the code of  $a$  is 1001001, the code of  $b$  is 10011001. The code of the grammar is

$$\underbrace{10101}_{S} \ 0000 \ \underbrace{1001001}_{a} \ 000 \ \underbrace{10101}_{S} \ 000 \ \underbrace{10011001}_{b} \ 00000 \ \underbrace{10101}_{S} \ 0000 \ \underbrace{1001001}_{a} \ 000 \ \underbrace{10011001}_{b} .$$

From this encoding results that the grammars with terminal alphabet  $T$  can be enumerated<sup>1</sup> as  $G_1, G_2, \dots, G_k, \dots$ , and the set of these grammars is a denumerable infinite set.

Consider now the set of all languages over  $T$  denoted by  $\mathcal{L}_T = \{L \mid L \subseteq T^*\}$ , that is  $\mathcal{L}_T = \mathcal{P}(T^*)$ . The set  $T^*$  is denumerable because its words can be ordered. Let this order  $s_0, s_1, s_2, \dots$ , where  $s_0 = \varepsilon$ . We associate to each language  $L \in \mathcal{L}_T$  an infinite binary sequence  $b_0, b_1, b_2, \dots$  the following way:

$$b_i = \begin{cases} 1, & \text{if } s_i \in L \\ 0, & \text{if } s_i \notin L \end{cases} \quad i = 0, 1, 2, \dots .$$

It is easy to see that the set of all such binary sequences is not denumerable, because each sequence can be considered as a positive number less than 1 using its binary representation

<sup>1</sup>Let us suppose that in the alphabet  $\{0, 1\}$  there is a linear order  $<$ , let us say  $0 < 1$ . The words which are codes of grammars can be enumerated by ordering them first after their lengths, and inside the equal length words, alphabetically, using the order of their letters. But we can use equally the lexicographic order, which means that  $u < v$  ( $u$  is before  $v$ ) if  $u$  is a proper prefix of  $v$  or there exists the decompositions  $u = xay$  and  $v = xby'$ , where  $x, y, y'$  are subwords,  $a$  and  $b$  letters with  $a < b$ .

(The decimal point is considered to be before the first digit). Conversely, to each positive number less than 1 in binary representation a binary sequence can be associated. So, the cardinality of the set of infinite binary sequences is equal to cardinality of interval  $[0, 1]$ , which is of continuum power. Therefore the set  $\mathcal{L}_T$  is of continuum cardinality. Now to each grammar with terminal alphabet  $T$  associate the corresponding generated language over  $T$ . Since the cardinality of the set of grammars is denumerable, there will exist a language from  $\mathcal{L}_T$ , without associated grammar, a language which cannot be generated by a grammar. ■

### 2.1.3. Chomsky hierarchy of grammars and languages

Putting some restrictions on the form of productions, four type of grammars can be distinguished.

**Definition 2.5** Define for a grammar  $G = (N, T, P, S)$  the following four types.

A grammar  $G$  is of type 0 (**phrase-structure grammar**) if there are no restrictions on productions.

A grammar  $G$  is of type 1 (**context-sensitive grammar**) if all of its productions are of the form  $\alpha A \gamma \rightarrow \alpha \beta \gamma$ , where  $A \in N$ ,  $\alpha, \gamma \in (N \cup T)^*$ ,  $\beta \in (N \cup T)^+$ . A production of the form  $S \rightarrow \varepsilon$  can also be accepted if the start symbol  $S$  does not occur in the right-hand side of any production.

A grammar  $G$  is of type 2 (**context-free grammar**) if all of its productions are of the form  $A \rightarrow \beta$ , where  $A \in N$ ,  $\beta \in (N \cup T)^+$ . A production of the form  $S \rightarrow \varepsilon$  can also be accepted if the start symbol  $S$  does not occur in the right-hand side of any production.

A grammar  $G$  is of type 3 (**regular grammar**) if its productions are of the form  $A \rightarrow aB$  or  $A \rightarrow a$ , where  $a \in T$  and  $A, B \in N$ . A production of the form  $S \rightarrow \varepsilon$  can also be accepted if the start symbol  $S$  does not occur in the right-hand side of any production.

If a grammar  $G$  is of type  $i$  then language  $L(G)$  is also of type  $i$ .

This classification was introduced by Noam Chomsky.

A language  $L$  is of type  $i$  ( $i = 0, 1, 2, 3$ ) if there exists a grammar  $G$  of type  $i$  which generates the language  $L$ , so  $L = L(G)$ .

Denote by  $\mathcal{L}_i$  ( $i = 0, 1, 2, 3$ ) the class of the languages of type  $i$ . Can be proved that

$$\mathcal{L}_0 \supset \mathcal{L}_1 \supset \mathcal{L}_2 \supset \mathcal{L}_3 .$$

By the definition of different type of languages, the inclusions ( $\supseteq$ ) are evident, but the strict inclusions ( $\supset$ ) must be proved.

**Example 2.3** We give an example for each type of context-sensitive, context-free and regular grammars.

*Context-sensitive grammar.*  $G_1 = (N_1, T_1, P_1, S_1)$ , where  $N_1 = \{S_1, A, B, C\}$ ,  $T_1 = \{a, 0, 1\}$ .

Elements of  $P_1$  are:

$$\begin{aligned} S_1 &\rightarrow ACA, \\ AC &\rightarrow AACA \mid ABa \mid AaB, \\ B &\rightarrow AB \mid A, \\ A &\rightarrow 0 \mid 1. \end{aligned}$$

Language  $L(G_1)$  contains words of the form  $uav$  with  $u, v \in \{0, 1\}^*$  and  $|u| \neq |v|$ .

*Context-free grammar.*  $G_2 = (N_2, T_2, P_2, S)$ , where  $N_2 = \{S, A, B\}$ ,  $T_2 = \{+, *, (, ), a\}$ .

Elements of  $P_2$  are:

$$\begin{aligned} S &\rightarrow S + A \mid A, \\ A &\rightarrow A * B \mid B, \\ B &\rightarrow (S) \mid a. \end{aligned}$$

Language  $L(G_2)$  contains algebraic expressions which can be correctly built using letter  $a$ , operators  $+$  and  $*$  and brackets.

*Regular grammar:*  $G_3 = (N_3, T_3, P_3, S_3)$ , where  $N_3 = \{S_3, A, B\}$ ,  $T_3 = \{a, b\}$ .

Elements of  $P_3$  are:

$$\begin{aligned} S_3 &\rightarrow aA \\ A &\rightarrow aB \mid a \\ B &\rightarrow aB \mid bB \mid a \mid b. \end{aligned}$$

Language  $L(G_3)$  contains words over the alphabet  $\{a, b\}$  with at least two letters  $a$  at the beginning.

It is easy to prove that any finite language is regular. The productions will be done to generate all words of the language. For example, if  $u = a_1a_2 \dots a_n$  is in the language, then we introduce the productions:  $S \rightarrow a_1A_1$ ,  $A_1 \rightarrow a_2A_2$ ,  $\dots$ ,  $A_{n-2} \rightarrow a_{n-1}A_{n-1}$ ,  $A_{n-1} \rightarrow a_n$ , where  $S$  is the start symbol of the language and  $A_1, \dots, A_{n-1}$  are distinct nonterminals. We define such productions for all words of the language using different nonterminals for different words, excepting the start symbol  $S$ . If the empty word is also an element of the language, then the production  $S \rightarrow \varepsilon$  is also considered.

The empty set is also a regular language, because the regular grammar  $G = (\{S\}, \{a\}, \{S \rightarrow aS\}, S)$  generates it.

### Eliminating unit productions

A production of the form  $A \rightarrow B$  is called a **unit production**, where  $A, B \in N$ . Unit productions can be eliminated from a grammar in such a way that the new grammar will be of the same type and equivalent to the first one.

Let  $G = (N, T, P, S)$  be a grammar with unit productions. Define an equivalent grammar  $G' = (N, T, P', S)$  without unit productions. The following algorithm will construct the equivalent grammar.

ELIMINATE-UNIT-PRODUCTIONS( $G, G'$ )

- 1 if the unit productions  $A \rightarrow B$  and  $B \rightarrow C$  are in  $P$  put also the unit production  $A \rightarrow C$  in  $P$  while  $P$  can be extended,
- 2 if the unit production  $A \rightarrow B$  and the production  $B \rightarrow \alpha$  ( $\alpha \notin N$ ) are in  $P$  put also the production  $A \rightarrow \alpha$  in  $P$ ,
- 3 let  $P'$  be the set of productions of  $P$  except unit productions.

Clearly,  $G$  and  $G'$  are equivalent. If  $G$  is of type  $i \in \{0, 1, 2, 3\}$  then  $G'$  is also of type  $i$ .

**Example 2.4** Use the above algorithm in the case of the grammar  $G = (\{S, A, B, C\}, \{a, b\}, P, S)$ , where  $P$  contains

$$\begin{array}{lllll} S \rightarrow A, & A \rightarrow B, & B \rightarrow C, & C \rightarrow B, & D \rightarrow C, \\ S \rightarrow B, & A \rightarrow D, & & C \rightarrow Aa, & \\ & A \rightarrow aB, & & & \\ & A \rightarrow b. & & & \end{array}$$

Using the first step of the algorithm, we get the following new unit productions:

$S \rightarrow D$	(because of $S \rightarrow A$ and $A \rightarrow D$ ),
$S \rightarrow C$	(because of $S \rightarrow B$ and $B \rightarrow C$ ),
$A \rightarrow C$	(because of $A \rightarrow B$ and $B \rightarrow C$ ),
$B \rightarrow B$	(because of $B \rightarrow C$ and $C \rightarrow B$ ),
$C \rightarrow C$	(because of $C \rightarrow B$ and $B \rightarrow C$ ),
$D \rightarrow B$	(because of $D \rightarrow C$ and $C \rightarrow B$ ).

In the second step of the algorithm will be considered only productions with  $A$  or  $C$  in the right-hand side, since productions  $A \rightarrow aB$ ,  $A \rightarrow b$  and  $C \rightarrow Aa$  can be used (the other productions are all unit productions). We get the following new productions:

$S \rightarrow aB$	(because of $S \rightarrow A$ and $A \rightarrow aB$ ),
$S \rightarrow b$	(because of $S \rightarrow A$ and $A \rightarrow b$ ),
$S \rightarrow Aa$	(because of $S \rightarrow C$ and $C \rightarrow Aa$ ),
$A \rightarrow Aa$	(because of $A \rightarrow C$ and $C \rightarrow Aa$ ),
$B \rightarrow Aa$	(because of $B \rightarrow C$ and $C \rightarrow Aa$ ).

The new grammar  $G' = (\{S, A, B, C\}, \{a, b\}, P', S)$  will have the productions:

$S \rightarrow b,$	$A \rightarrow b,$	$B \rightarrow Aa,$	$C \rightarrow Aa,$
$S \rightarrow aB,$	$A \rightarrow aB,$		
$S \rightarrow Aa$	$A \rightarrow Aa,$		

### Grammars in normal forms

A grammar is to be said a **grammar in normal form** if its productions have no terminal symbols in the left-hand side.

We need the following notions. For alphabets  $\Sigma_1$  and  $\Sigma_2$  a *homomorphism* is a function  $h : \Sigma_1^* \rightarrow \Sigma_2^*$  for which  $h(u_1u_2) = h(u_1)h(u_2)$ ,  $\forall u_1, u_2 \in \Sigma_1^*$ . It is easy to see that for arbitrary  $u = a_1a_2 \dots a_n \in \Sigma_1^*$  value  $h(u)$  is uniquely determined by the restriction of  $h$  on  $\Sigma_1$ , because  $h(u) = h(a_1)h(a_2) \dots h(a_n)$ .

If a homomorphism  $h$  is a bijection then  $h$  an *isomorphism*.

**Theorem 2.6** *To any grammar an equivalent grammar in normal form can be associated.*

#### Proof.

Grammars of type 2 and 3 have in left-hand side of any productions only a nonterminal, so they are in normal form. The proof has to be done for grammars of type 0 and 1 only.

Let  $G = (N, T, P, S)$  be the original grammar and we define the grammar in normal form as  $G' = (N', T, P', S)$ .

Let  $a_1, a_2, \dots, a_k$  be those terminal symbols which occur in the left-hand side of productions. We introduce the new nonterminals  $A_1, A_2, \dots, A_k$ . The following notation will be used:  $T_1 = \{a_1, a_2, \dots, a_k\}$ ,  $T_2 = T \setminus T_1$ ,  $N_1 = \{A_1, A_2, \dots, A_k\}$  and  $N' = N \cup N_1$ .

Define the isomorphism  $h : N \cup T \rightarrow N' \cup T_2$ , where

$$\begin{aligned} h(a_i) &= A_i, & \text{if } a_i \in T_1, \\ h(X) &= X, & \text{if } X \in N \cup T_2 \end{aligned}$$

Define the set  $P'$  of production as

$$P' = \{h(\alpha) \rightarrow h(\beta) \mid (\alpha \rightarrow \beta) \in P\} \cup \{A_i \rightarrow a_i \mid i = 1, 2, \dots, k\}$$

In this case  $\alpha \xrightarrow[G]{*} \beta$  if and only if  $h(\alpha) \xrightarrow[G']{*} h(\beta)$ . From this the theorem immediately



results because  $S \xRightarrow[G]{*} u \Leftrightarrow S = h(S) \xRightarrow[G']{*} h(u) = u$ . ■

**Example 2.5** Let  $G = (\{S, D, E\}, \{a, b, c, d, e\}, P, S)$ , where  $P$  contains

$$\begin{aligned} S &\rightarrow aebc \mid aDbc \\ Db &\rightarrow bD \\ Dc &\rightarrow Ebccd \\ bE &\rightarrow Eb \\ aE &\rightarrow aaD \mid aae \end{aligned}$$

In the left-hand side of productions the terminals  $a, b, c$  occur, therefore consider the new nonterminals  $A, B, C$ , and include in  $P'$  also the new productions  $A \rightarrow a, B \rightarrow b$  and  $C \rightarrow c$ .

Terminals  $a, b, c$  will be replaced by nonterminals  $A, B, C$  respectively, and we get the set  $P'$  as

$$\begin{aligned} S &\rightarrow AeBC \mid ADBC \\ DB &\rightarrow BD \\ DC &\rightarrow EBCCd \\ BE &\rightarrow EB \\ AE &\rightarrow AAD \mid AAe \\ A &\rightarrow a \\ B &\rightarrow b \\ C &\rightarrow c. \end{aligned}$$

Let us see what words can be generated by this grammars. It is easy to see that  $aebc \in L(G')$ , because  $S \Rightarrow AeBC \xRightarrow{*} aebc$ .

$S \Rightarrow ADBC \Rightarrow ABDC \Rightarrow ABEBCcd \Rightarrow AEBBCCd \Rightarrow AAeBBCCd \xRightarrow{*} aaebbccd$ , so  $aaebbccd \in L(G')$ .

We prove, using the mathematical induction, that  $S \xRightarrow{*} A^{n-1}EB^nC(Cd)^{n-1}$  for  $n \geq 2$ . For  $n = 2$  this is the case, as we have seen before. Continuing the derivation we get  $S \xRightarrow{*} A^{n-1}EB^nC(Cd)^{n-1} \Rightarrow A^{n-2}AADB^nC(Cd)^{n-1} \xRightarrow{*} A^nB^nDC(Cd)^{n-1} \Rightarrow A^nB^nEBCCd(Cd)^{n-1} \xRightarrow{*} A^nEB^{n+1}CCd(Cd)^{n-1} = A^nEB^{n+1}C(Cd)^n$ , and this is what we had to prove.

But  $S \xRightarrow{*} A^{n-1}EB^nC(Cd)^{n-1} \Rightarrow A^{n-2}AAeB^nC(Cd)^{n-1} \xRightarrow{*} a^n eb^n c(cd)^{n-1}$ . So  $a^n eb^n c(cd)^{n-1} \in L(G')$ ,  $n \geq 1$ . These words can be generated also in  $G$ .

#### 2.1.4. Extended grammars

In this subsection extended grammars of type 1, 2 and 3 will be presented.

**Extended grammar of type 1.** All productions are of the form  $\alpha \rightarrow \beta$ , where  $|\alpha| \leq |\beta|$ , excepted possibly the production  $S \rightarrow \varepsilon$ .

**Extended grammar of type 2.** All productions are of the form  $A \rightarrow \beta$ , where  $A \in N, \beta \in (N \cup T)^*$ .

**Extended grammar of type 3.** All productions are of the form  $A \rightarrow uB$  or  $A \rightarrow u$ , Where  $A, B \in N, u \in T^*$ .

**Theorem 2.7** *To any extended grammar an equivalent grammar of the same type can be associated.*

**Proof.** Denote by  $G_{ext}$  the extended grammar and by  $G$  the corresponding equivalent grammar of the same type.

*Type 1.* Define the productions of grammar  $G$  by rewriting the productions  $\alpha \rightarrow \beta$ , where  $|\alpha| \leq |\beta|$  of the extended grammar  $G_{ext}$  in the form  $\gamma_1\delta\gamma_2 \rightarrow \gamma_1\gamma\gamma_2$  allowed in the case of grammar  $G$  by the following way.

Let  $X_1X_2\dots X_m \rightarrow Y_1Y_2\dots Y_n$  ( $m \leq n$ ) be a production of  $G_{ext}$ , which is not in the required form. Add to the set of productions of  $G$  the following productions, where  $A_1, A_2, \dots, A_m$  are new nonterminals:

$$X_1X_2\dots X_m \rightarrow A_1X_2X_3\dots X_m$$

$$A_1X_2\dots X_m \rightarrow A_1A_2X_3\dots X_m$$

...

$$A_1A_2\dots A_{m-1}X_m \rightarrow A_1A_2\dots A_{m-1}A_m$$

$$A_1A_2\dots A_{m-1}A_m \rightarrow Y_1A_2\dots A_{m-1}A_m$$

$$Y_1A_2\dots A_{m-1}A_m \rightarrow Y_1Y_2\dots A_{m-1}A_m$$

...

$$Y_1Y_2\dots Y_{m-2}A_{m-1}A_m \rightarrow Y_1Y_2\dots Y_{m-2}Y_{m-1}A_m$$

$$Y_1Y_2\dots Y_{m-1}A_m \rightarrow Y_1Y_2\dots Y_{m-1}Y_mY_{m+1}\dots Y_n.$$

Furthermore, add to the set of productions of  $G$  without any modification the productions of  $G_{ext}$  which are of permitted form, i.e.  $\gamma_1\delta\gamma_2 \rightarrow \gamma_1\gamma\gamma_2$ .

Inclusion  $L(G_{ext}) \subseteq L(G)$  can be proved because each used production of  $G_{ext}$  in a derivation can be simulated by productions  $G$  obtained from it. Furthermore, since the productions of  $G$  can be used only in the prescribed order, we could not obtain other words, so  $L(G) \subseteq L(G_{ext})$  also is true.

*Type 2.* Let  $G_{ext} = (N, T, P, S)$ . Productions of form  $A \rightarrow \varepsilon$  have to be eliminated, only  $S \rightarrow \varepsilon$  can remain, if  $S$  doesn't occur in the right-hand side of productions. For this define the following sets:

$$U_0 = \{A \in N \mid (A \rightarrow \varepsilon) \in P\}$$

$$U_i = U_{i-1} \cup \{A \in N \mid (A \rightarrow w) \in P, w \in U_{i-1}^+\}.$$

Since for  $i \geq 1$  we have  $U_{i-1} \subseteq U_i$ ,  $U_i \subseteq N$  and  $N$  is a finite set, there must exist such a  $k$  for which  $U_{k-1} = U_k$ . Let us denote this set as  $U$ . It is easy to see that a nonterminal  $A$  is in  $U$  if and only if  $A \xRightarrow{*} \varepsilon$ . (In addition  $\varepsilon \in L(G_{ext})$  if and only if  $S \in U$ .)

We define the productions of  $G$  starting from the productions of  $G_{ext}$  in the following way. For each production  $A \rightarrow \alpha$  with  $\alpha \neq \varepsilon$  of  $G_{ext}$  add to the set of productions of  $G$  this one and all productions which can be obtained from it by eliminating from  $\alpha$  one or more nonterminals which are in  $U$ , but only in the case when the right-hand side does not become  $\varepsilon$ .

It is not difficult to see that this grammar  $G$  generates the same language as  $G_{ext}$  does, except the empty word  $\varepsilon$ . So, if  $\varepsilon \notin L(G_{ext})$  then the proof is finished. But if  $\varepsilon \in L(G_{ext})$ , then there are two cases. If the start symbol  $S$  does not occur in any right-hand side of productions, then by introducing the production  $S \rightarrow \varepsilon$ , grammar  $G$  will generate also the empty word. If  $S$  occurs in a production in the right-hand side, then we introduce a new start symbol  $S'$  and the new productions  $S' \rightarrow S$  and  $S' \rightarrow \varepsilon$ . Now the empty word  $\varepsilon$  can also be generated by grammar  $G$ .

*Type 3.* First we use for  $G_{ext}$  the procedure defined for grammars of type 2 to eliminate productions of the form  $A \rightarrow \varepsilon$ . From the obtained grammar we eliminate the unit productions using the algorithm ELIMINATE-UNIT-PRODUCTIONS (see page 82).

In the obtained grammar for each production  $A \rightarrow a_1a_2\dots a_nB$ , where  $B \in N \cup \{\varepsilon\}$ , add to the productions of  $G$  also the followings

$$\begin{aligned} A &\rightarrow a_1 A_1, \\ A_1 &\rightarrow a_2 A_2, \\ &\dots \\ A_{n-1} &\rightarrow a_n B, \end{aligned}$$

where  $A_1, A_2, \dots, A_{n-1}$  are new nonterminals. It is easy to prove that grammar  $G$  built in this way is equivalent to  $G_{ext}$ . ■

**Example 2.6** Let  $G_{ext} = (N, T, P, S)$  be an extended grammar of type 1, where  $N = \{S, B, C\}$ ,  $T = \{a, b, c\}$  and  $P$  contains the following productions:

$$\begin{array}{ll} S &\rightarrow aSBC \mid aBC & CB &\rightarrow BC \\ aB &\rightarrow ab & bB &\rightarrow bb \\ bC &\rightarrow bc & cC &\rightarrow cc. \end{array}$$

The only production which is not context-sensitive is  $CB \rightarrow BC$ . Using the method given in the proof, we introduce the productions:

$$\begin{aligned} CB &\rightarrow AB \\ AB &\rightarrow AD \\ AD &\rightarrow BD \\ BD &\rightarrow BC \end{aligned}$$

Now the grammar  $G = (\{S, A, B, C, D\}, \{a, b, c\}, P', S)$  is context-sensitive, where the elements of  $P'$  are

$$\begin{array}{ll} S &\rightarrow aSBC \mid aBC \\ CB &\rightarrow AB & aB &\rightarrow ab \\ AB &\rightarrow AD & bB &\rightarrow bb \\ AD &\rightarrow BD & bC &\rightarrow bc \\ BD &\rightarrow BC & cC &\rightarrow cc. \end{array}$$

It can be proved that  $L(G_{ext}) = L(G) = \{a^n b^p c^p \mid n \geq 1\}$ .

**Example 2.7** Let  $G_{ext} = (\{S, B, C\}, \{a, b, c\}, P, S)$  be an extended grammar of type 2, where  $P$  contains:

$$\begin{aligned} S &\rightarrow aSc \mid B \\ B &\rightarrow bB \mid C \\ C &\rightarrow Cc \mid \varepsilon. \end{aligned}$$

Then  $U_0 = \{C\}$ ,  $U_1 = \{B, C\}$ ,  $U_3 = \{S, B, C\} = U$ . The productions of the new grammar are:

$$\begin{aligned} S &\rightarrow aSc \mid ac \mid B \\ B &\rightarrow bB \mid b \mid C \\ C &\rightarrow Cc \mid c. \end{aligned}$$

The original grammar generates also the empty word and because  $S$  occurs in the right-hand side of a production, a new start symbol and two new productions will be defined:  $S' \rightarrow S, S' \rightarrow \varepsilon$ . The context-free grammar equivalent to the original grammar is  $G = (\{S', S, B, C\}, \{a, b, c\}, P', S')$  with the productions:

$$\begin{aligned} S' &\rightarrow S \mid \varepsilon \\ S &\rightarrow aSc \mid ac \mid B \\ B &\rightarrow bB \mid b \mid C \\ C &\rightarrow Cc \mid c. \end{aligned}$$

Both of these grammars generate language  $\{a^m b^n c^p \mid p \geq m \geq 0, n \geq 0\}$ .

**Example 2.8** Let  $G_{ext} = (\{S, A, B\}, \{a, b\}, P, S)$  be the extended grammar of type 3 under examination, where  $P$ :

$$\begin{aligned} S &\rightarrow abA \\ A &\rightarrow bB \\ B &\rightarrow S \mid \varepsilon. \end{aligned}$$

First, we eliminate production  $B \rightarrow \varepsilon$ . Since  $U_0 = U = \{B\}$ , the productions will be

$$\begin{aligned} S &\rightarrow abA \\ A &\rightarrow bB \mid b \\ B &\rightarrow S. \end{aligned}$$

The latter production (which a unit production) can also be eliminated, by replacing it with  $B \rightarrow abA$ . Productions  $S \rightarrow abA$  and  $B \rightarrow abA$  have to be transformed. Since, both productions have the same right-hand side, it is enough to introduce only one new nonterminal and to use the productions  $S \rightarrow aC$  and  $C \rightarrow bA$  instead of  $S \rightarrow abA$ . Production  $B \rightarrow abA$  will be replaced by  $B \rightarrow aC$ . The new grammar is  $G = (\{S, A, B, C\}, \{a, b\}, P', S)$ , where  $P'$ :

$$\begin{aligned} S &\rightarrow aC \\ A &\rightarrow bB \mid b \\ B &\rightarrow aC \\ C &\rightarrow bA. \end{aligned}$$

Can be proved that  $L(G_{ext}) = L(G) = \{(abb)^n \mid n \geq 1\}$ .

### 2.1.5. Closure properties in the Chomsky-classes

We will prove the following theorem, by which the Chomsky-classes of languages are closed under the regular operations, that is, the union and product of two languages of type  $i$  is also of type  $i$ , the iteration of a language of type  $i$  is also of type  $i$  ( $i = 0, 1, 2, 3$ ).

**Theorem 2.8** *The class  $\mathcal{L}_i$  ( $i = 0, 1, 2, 3$ ) of languages is closed under the regular operations.*

**Proof.** For the proof we will use extended grammars. Consider the extended grammars  $G_1 = (N_1, T_1, P_1, S_1)$  and  $G_2 = (N_2, T_2, P_2, S_2)$  of type  $i$  each. We can suppose that  $N_1 \cap N_2 = \emptyset$ .

*Union.* Let  $G_\cup = (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}, S)$ .

We will show that  $L(G_\cup) = L(G_1) \cup L(G_2)$ . If  $i = 0, 2, 3$  then from the assumption that  $G_1$  and  $G_2$  are of type  $i$  follows by definition that  $G_\cup$  also is of type  $i$ . If  $i = 1$  and one of the grammars generates the empty word, then we eliminate from  $G_\cup$  the corresponding production (possibly the both)  $S_k \rightarrow \varepsilon$  ( $k = 1, 2$ ) and replace it by production  $S \rightarrow \varepsilon$ .

*Product.* Let  $G_\times = (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}, S)$ .

We will show that  $L(G_\times) = L(G_1)L(G_2)$ . By definition, if  $i = 0, 2$  then  $G_\times$  will be of the same type. If  $i = 1$  and there is production  $S_1 \rightarrow \varepsilon$  in  $P_1$  but there is no production  $S_2 \rightarrow \varepsilon$  in  $P_2$  then production  $S_1 \rightarrow \varepsilon$  will be replaced by  $S \rightarrow S_2$ . We will proceed the same way in the symmetrical case. If there is in  $P_1$  production  $S_1 \rightarrow \varepsilon$  and in  $P_2$  production  $S_2 \rightarrow \varepsilon$  then they will be replaced by  $S \rightarrow \varepsilon$ .

In the case of regular grammars ( $i = 3$ ), because  $S \rightarrow S_1 S_2$  is not a regular production, we need to use another grammar  $G_\times = (N_1 \cup N_2, T_1 \cup T_2, P'_1 \cup P_2, S_1)$ , where the difference between  $P'_1$  and  $P_1$  lies in that instead of productions in the form  $A \rightarrow u, u \in T^*$  in  $P'_1$  will exist production of the form  $A \rightarrow uS_2$ .

*Iteration.* Let  $G_* = (N_1 \cup \{S\}, T_1, P, S)$ .

In the case of grammars of type 2 let  $P = P_1 \cup \{S \rightarrow S_1 S, S \rightarrow \varepsilon\}$ . Then  $G_*$  also is of type 2.

In the case of grammars of type 3, as in the case of product, we will change the productions, that is  $P = P'_1 \cup \{S \rightarrow S_1, S \rightarrow \varepsilon\}$ , where the difference between  $P'_1$  and  $P_1$  lies in that for each  $A \rightarrow u$  ( $u \in T^*$ ) will be replaced by  $A \rightarrow uS$ , and the others will be not changed. Then  $G_*$  also will be of type 3.

The productions given in the case of type 2 are not valid for  $i = 0, 1$ , because when applying production  $S \rightarrow S_1 S$  we can get the derivations of type  $S \xRightarrow{*} S_1 S_1$ ,  $S_1 \xRightarrow{*} \alpha_1 \beta_1$ ,  $S_1 \xRightarrow{*} \alpha_2 \beta_2$ , where  $\beta_1 \alpha_2$  can be a left-hand side of a production. In this case, replacing  $\beta_1 \alpha_2$  by its right-hand side in derivation  $S \xRightarrow{*} \alpha_1 \beta_1 \alpha_2 \beta_2$ , we can generate a word which is not in the iterated language. To avoid such situations, first let us assume that the language is in normal form, i.e. the left-hand side of productions does not contain terminals (see page 83), second we introduce a new nonterminal  $S'$ , so the set of nonterminals now is  $N_1 \cup \{S, S'\}$ , and the productions are the following:

$$P = P_1 \cup \{S \rightarrow \varepsilon, S \rightarrow S_1 S'\} \cup \{aS' \rightarrow aS \mid a \in T_1\}.$$

Now we can avoid situations in which the left-hand side of a production can extend over the limits of words in a derivation because of the iteration. The above derivations can be used only by beginning with  $S \xRightarrow{*} S_1 S'$  and getting derivation  $S \xRightarrow{*} \alpha_1 \beta_1 S'$ . Here we can not replace  $S'$  unless the last symbol in  $\beta_1$  is a terminal symbol, and only after using a production of the form  $aS' \rightarrow aS$ .

It is easy to show that  $L(G_*) = L(G_1)^*$  for each type. ■

### Exercises

**2.1-1** Give a grammar which generates language  $L = \{uu^{-1} \mid u \in \{a, b\}^*\}$  and determine its type.

**2.1-2** Let  $G = (N, T, P, S)$  be an extended context-free grammar, where

$$N = \{S, A, C, D\}, \quad T = \{a, b, c, d, e\},$$

$$P = \{S \rightarrow abCADE, C \rightarrow cC, C \rightarrow \varepsilon, D \rightarrow dD, D \rightarrow \varepsilon, A \rightarrow \varepsilon, A \rightarrow dDcCA\}.$$

Give an equivalent context-free grammar.

**2.1-3** Show that  $\Sigma^*$  and  $\Sigma^+$  are regular languages over arbitrary alphabet  $\Sigma$ .

**2.1-4** Give a grammar to generate language  $L = \{u \in \{0, 1\}^* \mid n_0(u) = n_1(u)\}$ , where  $n_0(u)$  represents the number of 0's in word  $u$  and  $n_1(u)$  the number of 1's.

**2.1-5** Give a grammar to generate all natural numbers.

**2.1-6** Give a grammar to generate the following languages, respectively:

$$L_1 = \{a^n b^m c^p \mid n \geq 1, m \geq 1, p \geq 1\},$$

$$L_2 = \{a^{2^n} \mid n \geq 1\},$$

$$L_3 = \{a^n b^m \mid n \geq 0, m \geq 0\},$$

$$L_4 = \{a^n b^m \mid n \geq m \geq 1\}.$$

**2.1-7** Let  $G = (N, T, P, S)$  be an extended grammar, where  $N = \{S, A, B, C\}$ ,  $T = \{a\}$  and  $P$  contains the productions:

$$S \rightarrow BAB, BA \rightarrow BC, CA \rightarrow AAC, CB \rightarrow AAB, A \rightarrow a, B \rightarrow \varepsilon.$$

Determine the type of this grammar. Give an equivalent, not extended grammar with the same type. What language it generates?

## 2.2. Finite automata and regular languages

Finite automata are computing models with input tape and a finite set of states (Fig. 2.1). Among the states some are called initial and some final. At the beginning the automaton read the first letter of the input word written on the input tape. Beginning with an initial state, the automaton read the letters of the input word one after another while change its

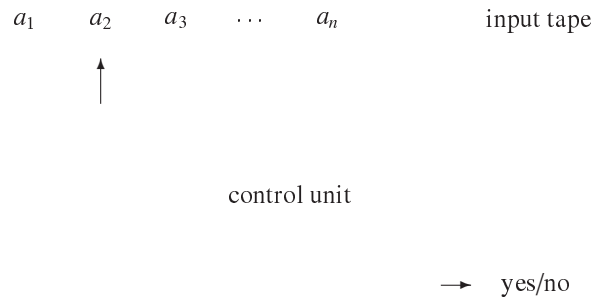


Figure 2.1. Finite automaton.

states, and when after reading the last input letter the current state is a final one, we say that the automaton accepts the given word. The set of words accepted by such an automaton is called the language accepted (recognized) by the automaton.

**Definition 2.9** A *nondeterministic finite automaton (NFA)* is a system  $A = (Q, \Sigma, E, I, F)$ , where

- $Q$  is a finite, nonempty set of *states*,
- $\Sigma$  is the *input alphabet*,
- $E$  is the set of *transitions* (or of edges), where  $E \subseteq Q \times \Sigma \times Q$ ,
- $I \subseteq Q$  is the set of *initial states*,
- $F \subseteq Q$  is the set of *final states*.

An NFA is in fact a directed, labelled graph, whose vertices are the states and there is a (directed) edge labelled with  $a$  from vertex  $p$  to vertex  $q$  if  $(p, a, q) \in E$ . Among vertices some are initial and some final states. Initial states are marked by a small arrow entering the corresponding vertex, while the final states are marked with double circles. If two vertices are joined by two edges with the same direction then these can be replaced by only one edge labelled with two letters. This graph can be called a transition graph.

**Example 2.9** Let  $A = (Q, \Sigma, E, I, F)$ , where  $Q = \{q_0, q_1, q_2\}$ ,  $\Sigma = \{0, 1, 2\}$ ,

$$E = \{(q_0, 0, q_0), (q_0, 1, q_1), (q_0, 2, q_2), \\ (q_1, 0, q_1), (q_1, 1, q_2), (q_1, 2, q_0), \\ (q_2, 0, q_2), (q_2, 1, q_0), (q_2, 2, q_1)\}$$

$$I = \{q_0\}, F = \{q_0\}.$$

The automaton can be seen in Fig. 2.2.

In the case of an edge  $(p, a, q)$  vertex  $p$  is the start-vertex,  $q$  the end-vertex and  $a$  the label. Now define the notion of the *walk* as in the case of graphs. A sequence

$$(q_0, a_1, q_1), (q_1, a_2, q_2), \dots, (q_{n-2}, a_{n-1}, q_{n-1}), (q_{n-1}, a_n, q_n)$$

of edges of a NFA is a walk with the label  $a_1 a_2 \dots a_n$ . If  $n = 0$  then  $q_0 = q_n$  and  $a_1 a_2 \dots a_n = \varepsilon$ . Such a walk is called an *empty walk*. For a walk the notation

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{n-1}} q_{n-1} \xrightarrow{a_n} q_n,$$

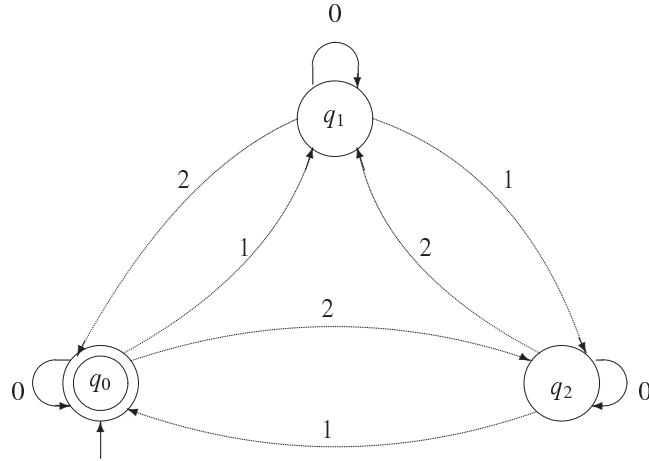


Figure 2.2. The finite automaton of Example 2.9..

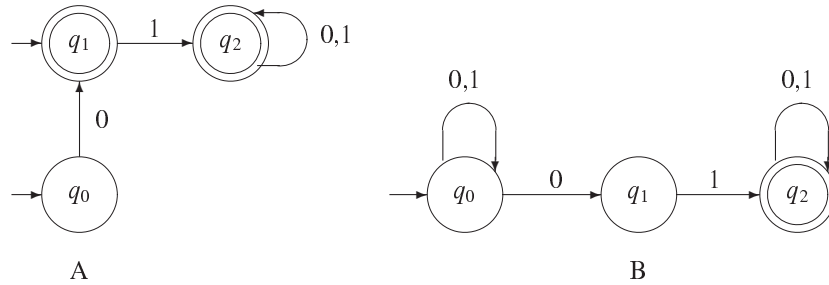


Figure 2.3. Nondeterministic finite automata.

will be used, or if  $w = a_1 a_2 \dots a_n$  then we write shortly  $q_0 \xrightarrow{w} q_n$ . Here  $q_0$  is the start-vertex and  $q_n$  the end-vertex of the walk. The states in a walk are not necessary distinct.

A walk is **productive** if its start-vertex is an initial state and its end-vertex is a final state. We say that an NFA **accepts** or **recognizes** a word if this word is the label of a productive walk. The empty word  $\varepsilon$  is accepted by an NFA if there is an empty productive walk, i.e. there is an initial state which is also a final state.

The set of words accepted by an NFA will be called the language accepted by this NFA. The **language accepted** or **recognized** by NFA A is

$$L(A) = \{w \in \Sigma^* \mid \exists p \in I, \exists q \in F, \exists p \xrightarrow{w} q\} .$$

The NFA  $A_1$  and  $A_2$  are **equivalent** if  $L(A_1) = L(A_2)$ .

Sometimes it is useful the following **transition function**:

$$\delta : Q \times \Sigma \rightarrow P(Q), \quad \delta(p, a) = \{q \in Q \mid (p, a, q) \in E\} .$$

$\delta$	0	1
$q_0$	$\{q_1\}$	$\emptyset$
$q_1$	$\emptyset$	$\{q_2\}$
$q_2$	$\{q_2\}$	$\{q_2\}$

A

$\delta$	0	1
$q_0$	$\{q_0, q_1\}$	$\{q_0\}$
$q_1$	$\emptyset$	$\{q_2\}$
$q_2$	$\{q_2\}$	$\{q_2\}$

B

**Figure 2.4.** Transition tables of the NFA in Fig. 2.3.

This function associate to a state  $p$  and input letter  $a$  the set of states in which the automaton can go if its current state is  $p$  and the head is on input letter  $a$ .

Denote by  $|H|$  the cardinal (the number of elements) of  $H$ .<sup>2</sup> An NFA is a **deterministic finite automaton** (DFA) if

$$|I| = 1 \text{ and } |\delta(q, a)| \leq 1, \forall q \in Q, \forall a \in \Sigma.$$

In Fig. 2.2 a DFA can be seen.

Condition  $|\delta(q, a)| \leq 1$  can be replaced by

$$(p, a, q) \in E, (p, a, r) \in E \implies q = r, \forall p, q, r \in Q, \forall a \in \Sigma.$$

If for a DFA  $|\delta(q, a)| = 1$  for each state  $q \in Q$  and for each letter  $a \in \Sigma$  then it is called a **complete DFA**.

Every DFA can be transformed in a complete DFA by introducing a new state, which can be called a snare state. Let  $A = (Q, \Sigma, E, \{q_0\}, F)$  be a DFA. An equivalent and complete DFA will be  $A' = (Q \cup \{s\}, \Sigma, E', \{q_0\}, F)$ , where  $s$  is the new state and  $E' = E \cup \{(p, a, s) \mid \delta(p, a) = \emptyset, p \in Q, a \in \Sigma\} \cup \{(s, a, s) \mid a \in \Sigma\}$ . It is easy to see that  $L(A) = L(A')$ .

Using the transition function we can easily define the transition table. The rows of this table are indexed by the elements of  $Q$ , its columns by the elements of  $\Sigma$ . At the intersection of row  $q \in Q$  and column  $a \in \Sigma$  we put  $\delta(q, a)$ . In the case of Fig. 2.2, the transition table is:

$\delta$	0	1	2
$q_0$	$\{q_0\}$	$\{q_1\}$	$\{q_2\}$
$q_1$	$\{q_1\}$	$\{q_2\}$	$\{q_0\}$
$q_2$	$\{q_2\}$	$\{q_0\}$	$\{q_1\}$

The NFA in 2.3 are not deterministic: the first (automaton A) has two initial states, the second (automaton B) has two transitions with 0 from state  $q_0$  (to states  $q_0$  and  $q_1$ ). The transition table of these two automata are in Fig. 2.4.  $L(A)$  is set of words over  $\Sigma = \{0, 1\}$  which do not begin with two zeroes (of course  $\varepsilon$  is in language),  $L(B)$  is the set of words which contain 01 as a subword.

### Eliminating inaccessible states

Let  $A = (Q, \Sigma, E, I, F)$  be a finite automaton. A state is *accessible* if it is on a walk which

<sup>2</sup>The same notation is used for the cardinal of a set and length of a word, but this is no matter of confusion because for word we use lowercase letters and for set capital letters. The only exception is  $\delta(q, a)$ , but this could not be confused with a word.



starts by an initial state. The following algorithm determines the inaccessible states building a sequence  $U_0, U_1, U_2, \dots$  of sets, where  $U_0$  is the set of initial states, and for any  $i \geq 1$   $U_i$  is the set of accessible states, which are at distance at most  $i$  from an initial state.

```

INACCESSIBLE-STATES(A, U)
   $U_0 \leftarrow I$ 
   $i \leftarrow 0$ 
  repeat
     $i \leftarrow i + 1$ 
    for all  $q \in U_{i-1}$  do
      for all  $a \in \Sigma$  do
         $U_i \leftarrow U_{i-1} \cup \delta(q, a)$ 
      endfor
    endfor
  until  $U_i = U_{i-1}$ 
   $U \leftarrow Q \setminus U_i$ 
  return  $U$ 

```

The inaccessible states of the automaton can be eliminated without changing the accepted language.

If  $|Q| = n$  and  $|\Sigma| = m$  then the running time of the algorithm (the number of steps) in the worst case is  $O(n^2m)$ , because the number of steps in the two embedded loops is at most  $nm$  and in the loop **repeat** at most  $n$ .

Set  $U$  has the property that  $L(A) \neq \emptyset$  if and only if  $U \cap F \neq \emptyset$ . The above algorithm can be extended by inserting the  $U \cap F \neq \emptyset$  condition to decide if language  $L(A)$  is or not empty.

### Eliminating nonproductive states

Let  $A = (Q, \Sigma, E, I, F)$  be a finite automaton. A state is *productive* if it is on a walk which ends in a terminal state. For finding the productive states the following algorithm uses the function  $\delta^{-1}$ :

$$\delta^{-1} : Q \times \Sigma \rightarrow \mathcal{P}(Q), \quad \delta^{-1}(p, a) = \{q \mid (q, a, p) \in E\}.$$

This function for a state  $p$  and a letter  $a$  gives the set of all states from which using this letter  $a$  the automaton can go into the state  $p$ .

```

NONPRODUCTIVE-STATES(A, V)
   $V_0 \leftarrow F$ 
   $i \leftarrow 0$ 
  repeat
     $i \leftarrow i + 1$ 
    for all  $p \in V_{i-1}$  do
      for all  $a \in \Sigma$  do
         $V_i \leftarrow V_{i-1} \cup \delta^{-1}(p, a)$ 
      endfor
    endfor
  until  $V_i = V_{i-1}$ 
   $V \leftarrow Q \setminus V_i$ 

```

**return**  $V$

The nonproductive states of the automaton can be eliminated without changing the accepted language.

If  $n$  is the number of states,  $m$  the number of letters in the alphabet, then the running time of the algorithm is also  $O(n^2m)$  as in the case of the algorithm `INACCESSIBLE-STATES`.

The set  $V$  given by the algorithm has the property that  $L(A) \neq \emptyset$  if and only if  $V \cap I \neq \emptyset$ . So, by a little modification it can be used to decide if language  $L(A)$  is or not empty.

### 2.2.1. Transforming nondeterministic finite automata in deterministic finite automata

As follows we will show that any NFA can be transformed in an equivalent DFA.

**Theorem 2.10** *For any NFA one may construct an equivalent DFA.*

**Proof.** Let  $A = (Q, \Sigma, E, I, F)$  be an NFA. Define a DFA  $\bar{A} = (\bar{Q}, \Sigma, \bar{E}, \bar{I}, \bar{F})$ , where

- $\bar{Q} = \mathcal{P}(Q) \setminus \emptyset$ ,
- edges of  $\bar{E}$  are those triplets  $(S, a, R)$  for which  $R, S \in \bar{Q}$  are not empty,  $a \in \Sigma$  and  $R = \bigcup_{p \in S} \delta(p, a)$ ,
- $\bar{I} = \{I\}$ ,
- $\bar{F} = \{S \subseteq Q \mid S \cap F \neq \emptyset\}$ .

We prove that  $L(A) = L(\bar{A})$ .

a) First prove that  $L(A) \subseteq L(\bar{A})$ . Let  $w = a_1 a_2 \dots a_k \in L(A)$ . Then there exists a walk

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{k-1}} q_{k-1} \xrightarrow{a_k} q_k, \quad q_0 \in I, \quad q_k \in F.$$

Using the transition function  $\bar{\delta}$  of NFA  $\bar{A}$  we construct the sets  $S_0 = \{q_0\}$ ,  $\bar{\delta}(S_0, a_1) = S_1$ ,  $\dots$ ,  $\bar{\delta}(S_{k-1}, a_k) = S_k$ . Then  $q_1 \in S_1, \dots, q_k \in S_k$  and since  $q_k \in F$  we get  $S_k \cap F \neq \emptyset$ , so  $S_k \in \bar{F}$ . Thus, there exists a walk

$$S_0 \xrightarrow{a_1} S_1 \xrightarrow{a_2} S_2 \xrightarrow{a_3} \dots \xrightarrow{a_{k-1}} S_{k-1} \xrightarrow{a_k} S_k, \quad S_0 \subseteq I, \quad S_k \in \bar{F}.$$

There are sets  $S'_0, \dots, S'_k$  for which  $S'_0 = I$ , and for  $i = 0, 1, \dots, k$  we have  $S_i \subseteq S'_i$ , and

$$S'_0 \xrightarrow{a_1} S'_1 \xrightarrow{a_2} S'_2 \xrightarrow{a_3} \dots \xrightarrow{a_{k-1}} S'_{k-1} \xrightarrow{a_k} S'_k$$

is a productive walk. Therefore  $w \in L(\bar{A})$ . That is  $L(A) \subseteq L(\bar{A})$ .

b) Now we show that  $L(\bar{A}) \subseteq L(A)$ . Let  $w = a_1 a_2 \dots a_k \in L(\bar{A})$ . Then there is a walk

$$\bar{q}_0 \xrightarrow{a_1} \bar{q}_1 \xrightarrow{a_2} \bar{q}_2 \xrightarrow{a_3} \dots \xrightarrow{a_{k-1}} \bar{q}_{k-1} \xrightarrow{a_k} \bar{q}_k, \quad \bar{q}_0 \in \bar{I}, \quad \bar{q}_k \in \bar{F}.$$

Using the definition of  $\bar{F}$  we have  $\bar{q}_k \cap F \neq \emptyset$ , i.e. there exists  $q_k \in \bar{q}_k \cap F$ , that is by the definitions of  $q_k \in F$  and  $\bar{q}_k$  there is  $q_{k-1}$  such that  $(q_{k-1}, a_k, q_k) \in E$ . Similarly, there are the states  $q_{k-2}, \dots, q_1, q_0$  such that  $(q_{k-2}, a_k, q_{k-1}) \in E, \dots, (q_0, a_1, q_1) \in E$ , where  $q_0 \in \bar{q}_0 = I$ , thus, there is a walk

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{k-1}} q_{k-1} \xrightarrow{a_k} q_k, \quad q_0 \in I, \quad q_k \in F,$$

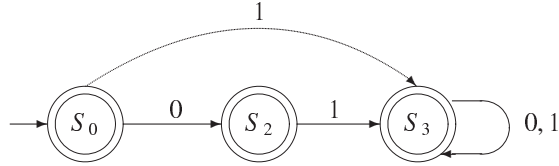


Figure 2.5. The equivalent DFA with NFA A in Fig. 2.3.

so  $L(\bar{A}) \subseteq L(A)$ . ■

In constructing DFA we can use the corresponding transition function  $\bar{\delta}$ :

$$\bar{\delta}(\bar{q}, a) = \left\{ \bigcup_{q \in \bar{q}} \delta(q, a) \right\}, \quad \forall \bar{q} \in \bar{Q}, \forall a \in \Sigma.$$

The empty set was excluded from the states, so we used here  $\emptyset$  instead of  $\{\emptyset\}$ .

**Example 2.10** Apply Theorem 2.10 to transform NFA A in Fig. 2.3. Introduce the following notation for the states of the DFA:

$$\begin{aligned} S_0 &:= \{q_0, q_1\}, & S_1 &:= \{q_0\}, & S_2 &:= \{q_1\}, & S_3 &:= \{q_2\}, \\ S_4 &:= \{q_0, q_2\}, & S_5 &:= \{q_1, q_2\}, & S_6 &:= \{q_0, q_1, q_2\}, \end{aligned}$$

where  $S_0$  is the initial state. Using the transition function we get the transition table:

$\bar{\delta}$	0	1
$S_0$	$\{S_2\}$	$\{S_3\}$
$S_1$	$\{S_2\}$	$\emptyset$
$S_2$	$\emptyset$	$\{S_3\}$
$S_3$	$\{S_3\}$	$\{S_3\}$
$S_4$	$\{S_5\}$	$\{S_3\}$
$S_5$	$\{S_3\}$	$\{S_3\}$
$S_6$	$\{S_5\}$	$\{S_3\}$

This automaton contains many inaccessible states. By algorithm INACCESSIBLE-STATES we determine the accessible states of DFA:

$$U_0 = \{S_0\}, \quad U_1 = \{S_0, S_2, S_3\}, \quad U_2 = \{S_0, S_2, S_3\} = U_1 = U.$$

Initial state  $S_0$  is also a final state. States  $S_2$  and  $S_3$  are final states. States  $S_1, S_4, S_5, S_6$  are inaccessible and can be removed from the DFA. The transition table of the resulted DFA is

$\bar{\delta}$	0	1
$S_0$	$\{S_2\}$	$\{S_3\}$
$S_2$	$\emptyset$	$\{S_3\}$
$S_3$	$\{S_3\}$	$\{S_3\}$

The corresponding transition graph is in Fig. 2.5.

The algorithm given in Theorem 2.10 can be simplified. It is not necessary to consider all subset of the set of states of NFA. The states of DFA  $\bar{A}$  can be obtained successively.

Begin with the state  $\bar{q}_0 = I$  and determine the states  $\bar{\delta}(\bar{q}_0, a)$  for all  $a \in \Sigma$ . For the newly obtained states we determine the states accessible from them. This can be continued until no new states arise.

In our previous example  $\bar{q}_0 := \{q_0, q_1\}$  is the initial state. From this we get

$$\begin{aligned} \bar{\delta}(\bar{q}_0, 0) &= \{\bar{q}_1\}, & \text{where } \bar{q}_1 &:= \{q_1\}, & \bar{\delta}(\bar{q}_0, 1) &= \{\bar{q}_2\}, & \text{where } \bar{q}_2 &:= \{q_2\}, \\ \bar{\delta}(\bar{q}_1, 0) &= \emptyset, & & & \bar{\delta}(\bar{q}_1, 1) &= \{\bar{q}_2\}, & & \\ \bar{\delta}(\bar{q}_2, 0) &= \{\bar{q}_2\}, & & & \bar{\delta}(\bar{q}_2, 1) &= \{\bar{q}_2\}. & & \end{aligned}$$

The transition table is

$\bar{\delta}$	0	1
$\bar{q}_0$	$\{\bar{q}_1\}$	$\{\bar{q}_2\}$
$\bar{q}_1$	$\emptyset$	$\{\bar{q}_2\}$
$\bar{q}_2$	$\{\bar{q}_2\}$	$\{\bar{q}_2\}$

which is the same (excepted the notation) as before.

The next algorithm will construct for an NFA  $A = (Q, \Sigma, E, I, F)$  the transition table  $M$  of the equivalent DFA  $\bar{A} = (\bar{Q}, \Sigma, \bar{E}, \bar{I}, \bar{F})$ , but without to determine the final states (which can easily be included). Value of  $\text{IsIn}(\bar{q}, \bar{Q})$  in the algorithm is true if state  $\bar{q}$  is already in  $\bar{Q}$  and is false otherwise. Let  $a_1, a_2, \dots, a_m$  be an ordered list of the letters of  $\Sigma$ .

**NFA-DFA**( $A, \bar{A}$ )

```

1   $\bar{q}_0 \leftarrow I$ 
2   $\bar{Q} \leftarrow \{\bar{q}_0\}$ 
3   $i \leftarrow 0$ 
4   $k \leftarrow 0$ 
5  repeat
6      for  $j = 1, 2, \dots, m$ 
7          do  $\bar{q} \leftarrow \bigcup_{p \in \bar{q}_i} \delta(p, a_j)$ 
8              if  $\bar{q} \neq \emptyset$ 
9                  then if  $\text{IsIn}(\bar{q}, \bar{Q})$ 
10                     then  $M[i, j] \leftarrow \{\bar{q}\}$ 
11                     else  $k \leftarrow k + 1$ 
12                          $\bar{q}_k \leftarrow \bar{q}$ 
13                          $M[i, j] \leftarrow \{\bar{q}_k\}$ 
14                          $\bar{Q} \leftarrow \bar{Q} \cup \{\bar{q}_k\}$ 
15                     else  $M[i, j] \leftarrow \emptyset$ 
16              $i \leftarrow i + 1$ 
17 until  $i = k + 1$ 
18 return transition table  $M$  of  $\bar{A}$ 

```

$\triangleright i$  counts the rows.  
 $\triangleright k$  counts the states.  
 $\triangleright j$  counts the columns.

Since loop **repeat** is executed as many times as the number of states of new automaton, in worst case the running time can be exponential, because, if the number of states in NFA is  $n$ , then DFA can have even  $2^n - 1$  states. (The number of subsets of a set of  $n$  elements is  $2^n$ , including the empty set.)

Theorem 2.10 will have it that to any NFA one may construct an equivalent DFA. Conversely, any DFA is also an NFA by definition. So, the nondeterministic finite automata

accepts the same class of languages as the deterministic finite automata.

### 2.2.2. Equivalence of deterministic finite automata

In this subsection we will use complete deterministic finite automata only. In this case  $\delta(q, a)$  has a single element. In formulae, sometimes, instead of set  $\delta(q, a)$  we will use its single element. We introduce for a set  $A = \{a\}$  the function  $elem(A)$  which give us the single element of set  $A$ , so  $elem(A) = a$ . Using walks which begin with the initial state and have the same label in two DFA's we can determine the equivalence of these DFA's. If only one of these walks ends in a final state, then they could not be equivalent.

Consider two DFA's over the same alphabet  $A = (Q, \Sigma, E, \{q_0\}, F)$  and  $A' = (Q', \Sigma, E', \{q'_0\}, F')$ . We are interested to determine if they are or not equivalent. We construct a table with elements of form  $(q, q')$ , where  $q \in Q$  and  $q' \in Q'$ . Beginning with the second column of the table, we associate a column to each letter of the alphabet  $\Sigma$ . If the first element of the  $i$ th row is  $(q, q')$  then at the cross of  $i$ th row and the column associated to letter  $a$  will be the pair  $(elem(\delta(q, a)), elem(\delta'(q', a)))$ .

	... a ...
...	...
$(q, q')$	$(elem(\delta(q, a)), elem(\delta'(q', a)))$
...	...

In the first column of the first row we put  $(q_0, q'_0)$  and complete the first row using the above method. If in the first row in any column there occur a pair of states from which one is a final state and the other not then the algorithm ends, the two automata are **not equivalent**. If there is no such a pair of states, every new pair is written in the first column. The algorithm continues with the next unfilled row. If no new pair of states occurs in the table and for each pair both of states are final or both are not, then the algorithm ends and the two DFA are **equivalent**.

DFA-EQUIVALENCE( $A, A'$ )

- 1 write in the first column of the first row the pair  $(q_0, q'_0)$
- 2  $i \leftarrow 0$
- 3 **repeat**
- 4      $i \leftarrow i + 1$
- 5     let  $(q, q')$  be the pair in the first column of the  $i$ th row
- 6     **for** all  $a \in \Sigma$
- 7         **do** write in the column associated to  $a$  in the  $i$ th row  
            the pair  $(elem(\delta(q, a)), elem(\delta'(q', a)))$
- 8         **if** one state in  $(elem(\delta(q, a)), elem(\delta'(q', a)))$  is final and the other not
- 9             **then return** NO
- 10         **else** write pair  $(elem(\delta(q, a)), elem(\delta'(q', a)))$  in the next empty row  
            of the first column, if not occurred already in the first column
- 11 **until** the first element of  $(i + 1)$ th row becomes empty
- 12 **return** YES

If  $|Q| = n$ ,  $|Q'| = n'$  and  $|\Sigma| = m$  then taking into account that in worst case loop **repeat** is executed  $nm'$  times, loop **for**  $m$  times, the running time of the algorithm in worst case will be  $O(nm'm)$ , or if  $n = n'$  then  $O(n^2m)$ .

Our algorithm was described to determine the equivalence of two complete DFA's. If we have to determine the equivalence of two NFA's, first we transform them into complete DFA's and after this we can apply the above algorithm.

**Example 2.11** Determine if the two DFA's in Fig. 2.6 are equivalent or not. The algorithm gives the table

	$a$	$b$
$(q_0, p_0)$	$(q_2, p_3)$	$(q_1, p_1)$
$(q_2, p_3)$	$(q_1, p_2)$	$(q_2, p_3)$
$(q_1, p_1)$	$(q_2, p_3)$	$(q_0, p_0)$
$(q_1, p_2)$	$(q_2, p_3)$	$(q_0, p_0)$

The two DFA's are equivalent because all possible pairs of states are considered and in every pair both

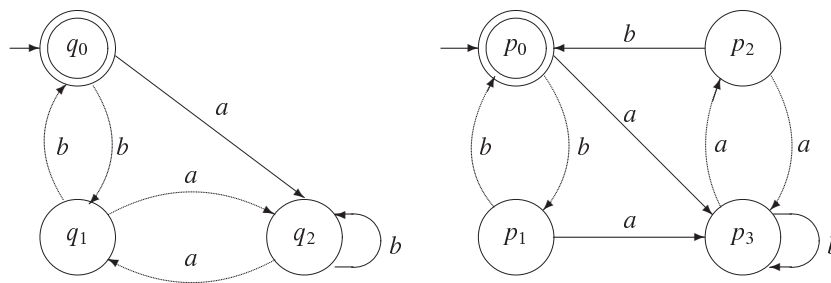


Figure 2.6. Equivalent DFA's (Example 2.11.).

states are final or both are not final.

**Example 2.12** The table of the two DFA's in Fig. 2.7 is:

	<i>a</i>	<i>b</i>
$(q_0, p_0)$	$(q_1, p_3)$	$(q_2, p_1)$
$(q_1, p_3)$	$(q_2, p_2)$	$(q_0, p_3)$
$(q_2, p_1)$		
$(q_2, p_2)$		

These two DFA's are not equivalent, because in the last column of the second row in the pair  $(q_0, p_3)$  the first state is final and the second not.

### 2.2.3. Equivalence of finite automata and regular languages

We have seen that NFA's accept the same class of languages as DFA's. The following theorem states that this class is that of regular languages.

**Theorem 2.11** *If  $L$  is a language accepted by a DFA, then one may construct a regular grammar which generates language  $L$ .*

**Proof.** Let  $A = (Q, \Sigma, E, \{q_0\}, F)$  be the DFA accepting language  $L$ , that is  $L = L(A)$ . Define the regular grammar  $G = (Q, \Sigma, P, q_0)$  with the productions:

- If  $(p, a, q) \in E$  for  $p, q \in Q$  and  $a \in \Sigma$ , then put production  $p \rightarrow aq$  in  $P$ .
- If  $(p, a, q) \in E$  and  $q \in F$ , then put also production  $p \rightarrow a$  in  $P$ .

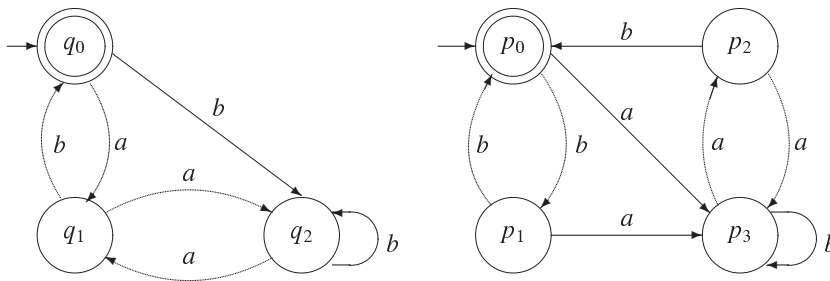
Prove that  $L(G) = L(A) \setminus \{\varepsilon\}$ .

Let  $u = a_1a_2 \dots a_n \in L(A)$  and  $u \neq \varepsilon$ . Thus, since  $A$  accepts word  $u$ , there is a walk

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{n-1}} q_{n-1} \xrightarrow{a_n} q_n, \quad q_n \in F.$$

Then there are in  $P$  the productions

$$q_0 \rightarrow a_1q_1, \quad q_1 \rightarrow a_2q_2, \quad \dots, \quad q_{n-2} \rightarrow a_{n-1}q_{n-1}, \quad q_{n-1} \rightarrow a_n$$



**Figure 2.7.** Non equivalent DFA's (Example 2.12).

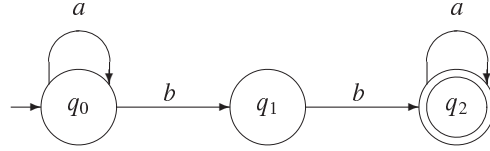


Figure 2.8. DFA of the Example 2.13..

(in the right-hand side of the last production  $q_n$  does not occur, because  $q_n \in F$ ), so there is the derivation

$$q_0 \Rightarrow a_1q_1 \Rightarrow a_1a_2q_2 \Rightarrow \dots \Rightarrow a_1a_2\dots a_{n-1}q_{n-1} \Rightarrow a_1a_2\dots a_n.$$

Therefore,  $u \in L(G)$ .

Conversely, let  $u = a_1a_2\dots a_n \in L(G)$  and  $u \neq \varepsilon$ . Then there exists a derivation

$$q_0 \Rightarrow a_1q_1 \Rightarrow a_1a_2q_2 \Rightarrow \dots \Rightarrow a_1a_2\dots a_{n-1}q_{n-1} \Rightarrow a_1a_2\dots a_n,$$

in which productions

$$q_0 \rightarrow a_1q_1, \quad q_1 \rightarrow a_2q_2, \quad \dots, \quad q_{n-2} \rightarrow a_{n-1}q_{n-1}, \quad q_{n-1} \rightarrow a_n$$

were used, which by definition means that in DFA A there is a walk

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{n-1}} q_{n-1} \xrightarrow{a_n} q_n,$$

and since  $q_n$  is a final state,  $u \in L(A) \setminus \{\varepsilon\}$ .

If the DFA accepts also the empty word  $\varepsilon$ , then in the above grammar we introduce a new start symbol  $q'_0$  instead of  $q_0$ , consider the new production  $q'_0 \rightarrow \varepsilon$  and for each production  $q_0 \rightarrow \alpha$  introduce also  $q'_0 \rightarrow \alpha$ . ■

**Example 2.13** Let  $A = (\{q_0, q_1, q_2\}, \{a, b\}, E, \{q_0\}, \{q_2\})$  be a DFA, where  $E = \{(q_0, a, q_0), (q_0, b, q_1), (q_1, b, q_2), (q_2, a, q_2)\}$ . The corresponding transition table is

$\delta$	a	b
$q_0$	$\{q_0\}$	$\{q_1\}$
$q_1$	$\emptyset$	$\{q_2\}$
$q_2$	$\{q_2\}$	$\emptyset$

The transition graph of A is in Fig. 2.8. By Theorem 2.11 we define regular grammar  $G = (\{q_0, q_1, q_2\}, \{a, b\}, P, q_0)$  with the productions in P

$$q_0 \rightarrow aq_0 \mid bq_1, \quad q_1 \rightarrow bq_2 \mid b, \quad q_2 \rightarrow aq_2 \mid a.$$

One may prove that  $L(A) = \{a^m bba^n \mid m \geq 0, n \geq 0\}$ .

The method described in the proof of Theorem 2.11 easily can be given as an algorithm. The productions of regular grammar  $G = (Q, \Sigma, P, q_0)$  obtained from the DFA  $A = (Q, \Sigma, E, \{q_0\}, F)$  can be determined by the following algorithm.



REGULAR-GRAMMAR-FROM-DFA(A, G)

```

1  P ← ∅
2  for all p ∈ Q
3      do for all a ∈ Σ
4          do for all q ∈ Q
5              do if (p, a, q) ∈ E
6                  then P ← P ∪ {p → aq}
7                      if q ∈ F
8                          then P ← P ∪ {p → a}
9  if q0 ∈ F
10 then P ← P ∪ {q0 → ε}

```

It is easy to see that the running time of the algorithm is  $\Theta(n^2m)$ , if the number of states is  $n$  and the number of letter in alphabet is  $m$ . In lines 2–4 we can consider only one loop, if we use the elements of  $E$ . Then the worst case running time is  $\Theta(p)$ , where  $p$  is the number of transitions of DFA. This is also  $O(n^2m)$ , since all transitions are possible. This algorithm is:

REGULAR-GRAMMAR-FROM-DFA'(A, G)

```

1  P ← ∅
2  for all (p, a, q) ∈ E
3      do P ← P ∪ {p → aq}
4          if q ∈ F
5              then P ← P ∪ {p → a}
6  if q0 ∈ F
7      then P ← P ∪ {q0 → ε}

```

**Theorem 2.12** *If  $L = L(G)$  is a regular language, then one may construct an NFA that accepts language  $L$ .*

**Proof.** Let  $G = (N, T, P, S)$  be the grammar which generates language  $L$ . Define NFA  $A = (Q, T, E, \{S\}, F)$ :

- $Q = N \cup \{Z\}$ , where  $Z \notin N \cup T$  (i.e.  $Z$  is a new symbol),
- For every production  $A \rightarrow aB$ , define transition  $(A, a, B)$  in  $E$ .
- For every production  $A \rightarrow a$ , define transition  $(A, a, Z)$  in  $E$ .
- $F = \begin{cases} \{Z\} & \text{if production } S \rightarrow \varepsilon \text{ does not occur in } G, \\ \{Z, S\} & \text{if production } S \rightarrow \varepsilon \text{ occurs in } G. \end{cases}$

Prove that  $L(G) = L(A)$ .

Let  $u = a_1a_2 \dots a_n \in L(G)$ ,  $u \neq \varepsilon$ . Then there is in  $G$  a derivation of word  $u$ :

$$S \Rightarrow a_1A_1 \Rightarrow a_1a_2A_2 \Rightarrow \dots \Rightarrow a_1a_2 \dots a_{n-1}A_{n-1} \Rightarrow a_1a_2 \dots a_n.$$

This derivation is based on productions

$$S \rightarrow a_1A_1, A_1 \rightarrow a_2A_2, \dots, A_{n-2} \rightarrow a_{n-1}A_{n-1}, A_{n-1} \rightarrow a_n.$$

Then, by the definition of the transitions of NFA  $A$  there exists a walk

$$S \xrightarrow{a_1} A_1 \xrightarrow{a_2} A_2 \xrightarrow{a_3} \dots \xrightarrow{a_{n-1}} A_{n-1} \xrightarrow{a_n} Z, Z \in F.$$

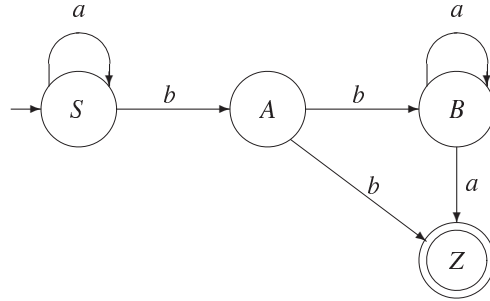


Figure 2.9. NFA associated to grammar in Example 2.14..

Thus,  $u \in L(A)$ . If  $\varepsilon \in L(G)$ , there is production  $S \rightarrow \varepsilon$ , but in this case the initial state is also a final one, so  $\varepsilon \in L(A)$ . Therefore,  $L(G) \subseteq L(A)$ .

Let now  $u = a_1 a_2 \dots a_n \in L(A)$ . Then there exists a walk

$$S \xrightarrow{a_1} A_1 \xrightarrow{a_2} A_2 \xrightarrow{a_3} \dots \xrightarrow{a_{n-1}} A_{n-1} \xrightarrow{a_n} Z, \quad Z \in F.$$

If  $u$  is the empty word, then instead of  $Z$  we have in the above formula  $S$ , which also is a final state. In other cases only  $Z$  can be as last symbol. Thus, in  $G$  there exist the productions

$$S \rightarrow a_1 A_1, \quad A_1 \rightarrow a_2 A_2, \quad \dots, \quad A_{n-2} \rightarrow a_{n-1} A_{n-1}, \quad A_{n-1} \rightarrow a_n,$$

and there is the derivation

$$S \Rightarrow a_1 A_1 \Rightarrow a_1 a_2 A_2 \Rightarrow \dots \Rightarrow a_1 a_2 \dots a_{n-1} A_{n-1} \Rightarrow a_1 a_2 \dots a_n,$$

thus,  $u \in L(G)$  and therefore  $L(A) \subseteq L(G)$ . ■

**Example 2.14** Let  $G = (\{S, A, B\}, \{a, b\}, \{S \rightarrow aS, S \rightarrow bA, A \rightarrow bB, A \rightarrow b, B \rightarrow aB, B \rightarrow a\}, S)$  be a regular grammar. The NFA associated is  $A = (\{S, A, B, Z\}, \{a, b\}, E, S, \{Z\})$ , where  $E = \{(S, a, S), (S, b, A), (A, b, B), (A, b, Z), (B, a, B), (B, a, Z)\}$ . The corresponding transition table is

$\delta$	$a$	$b$
$S$	$\{S\}$	$\{A\}$
$A$	$\emptyset$	$\{B, Z\}$
$B$	$\{B, Z\}$	$\emptyset$
$E$	$\emptyset$	$\emptyset$

The transition graph is in Fig. 2.9. This NFA can be simplified, states  $B$  and  $Z$  can be contracted in one final state.

Using the above theorem we define an algorithm which associate an NFA  $A = (Q, T, E, \{S\}, F)$  to a regular grammar  $G = (N, T, P, S)$

NFA-FROM-REGULAR-GRAMMAR( $G, A$ )

```

1   $E \leftarrow \emptyset$ 
2   $Q \leftarrow N \cup \{Z\}$ 
3  for all  $A \in N$ 
4      do for all  $a \in T$ 
5          do if  $(A \rightarrow a) \in P$ 
6              then  $E \leftarrow E \cup \{(A, a, Z)\}$ 
7              for all  $B \in N$ 
8                  do if  $(A \rightarrow aB) \in P$ 
9                      then  $E \leftarrow E \cup \{(A, a, B)\}$ 
10 if  $(S \rightarrow \varepsilon) \notin P$ 
11     then  $F \leftarrow \{Z\}$ 
12     else  $F \leftarrow \{Z, S\}$ 

```

As in the case of algorithm REGULAR-GRAMMAR-FROM-DFA, the running time is  $\Theta(n^2m)$ , where  $n$  is number of nonterminals and  $m$  the number of terminals. Loops in lines 3, 4 and 7 can be replaced by only one, which uses productions. The running time in this case is better and is equal to  $\Theta(p)$ , if  $p$  is the number of productions. This algorithm is:

NFA-FROM-REGULAR-GRAMMAR'( $G, A$ )

```

1   $E \leftarrow \emptyset$ 
2   $Q \leftarrow N \cup \{Z\}$ 
3  for all  $(A \rightarrow u) \in P$ 
4      do if  $u = a$ 
5          then  $E \leftarrow E \cup \{(A, a, Z)\}$ 
6          if  $u = aB$ 
7              then  $E \leftarrow E \cup \{(A, a, B)\}$ 
8  if  $(S \rightarrow \varepsilon) \notin P$ 
9     then  $F \leftarrow \{Z\}$ 
10    else  $F \leftarrow \{Z, S\}$ 

```

From theorems 2.10, 2.11 and 2.12 results that the class of regular languages coincides with the class of languages accepted by NFA's and also with class of languages accepted by DFA's. The result of these three theorems is illustrated in Fig. 2.10 and can be summarised also in the following theorem.

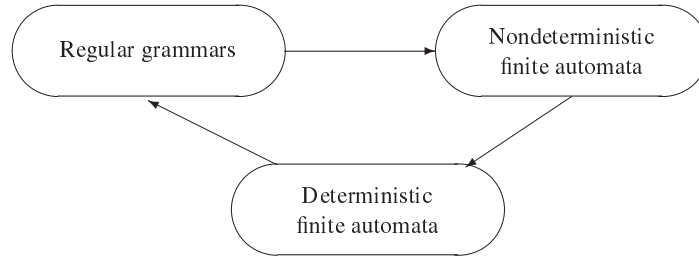
**Theorem 2.13** *The following three class of languages are the same:*

- the class of regular languages,
- the class of languages accepted by DFA's,
- the class of languages accepted by NFA's.

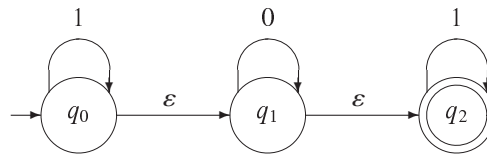
### Operation on regular languages

It is known (see Theorem 2.8) that the set  $\mathcal{L}_3$  of regular languages is closed under the regular operations, that is if  $L_1, L_2$  are regular languages, then languages  $L_1 \cup L_2$ ,  $L_1 L_2$  and  $L_1^*$  are also regular. For regular languages are true also the following statements.

*The complement of a regular language is also regular.* This is easy to prove using au-



**Figure 2.10.** Relations between regular grammars and finite automata. To any regular grammar one may construct an NFA which accepts the language generated by that grammar. Any NFA can be transformed in an equivalent DFA. To any DFA one may construct a regular grammar which generates the language accepted by that DFA.



**Figure 2.11.** Finite automata  $\epsilon$ -moves.

tomata. Let  $L$  be a regular language and let  $A = (Q, \Sigma, E, \{q_0\}, F)$  be a DFA which accepts language  $L$ . It is easy to see that the DFA  $\bar{A} = (Q, \Sigma, E, \{q_0\}, Q \setminus F)$  accepts language  $\bar{L}$ . So,  $\bar{L}$  is also regular.

*The intersection of two regular languages is also regular.* Since  $L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$ , the intersection is also regular.

*The difference of two regular languages is also regular.* Since  $L_1 \setminus L_2 = L_1 \cap \overline{L_2}$ , the difference is also regular.

### 2.2.4. Finite automata with $\epsilon$ -moves

A finite automaton with  $\epsilon$ -moves (FA with  $\epsilon$ -moves) extends NFA in such way that it may have transitions on the empty input  $\epsilon$ , i.e. it may change a state without reading any input symbol. In the case of a FA with  $\epsilon$ -moves  $A = (Q, \Sigma, E, I, F)$  for the set of transitions it is true that  $E \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$ .

The transition function of a FA with  $\epsilon$ -moves is:

$$\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow \mathcal{P}(Q), \quad \delta(p, a) = \{q \in Q \mid (p, a, q) \in E\}.$$

The FA with  $\epsilon$ -moves in Fig. 2.11 accepts words of form  $uvw$ , where  $u \in \{1\}^*$ ,  $v \in \{0\}^*$  and  $w \in \{1\}^*$ .

**Theorem 2.14** *To any FA with  $\epsilon$ -moves one may construct an equivalent NFA (without  $\epsilon$ -moves).*

Let  $A = (Q, \Sigma, E, I, F)$  be an FA with  $\epsilon$ -moves and we construct an equivalent NFA  $\bar{A} = (Q, \Sigma, \bar{E}, I, \bar{F})$ . The following algorithm determines sets  $\bar{F}$  and  $\bar{E}$ .

For a state  $q$  denote by  $\Lambda(q)$  the set of states (including even  $q$ ) in which one may go

from  $q$  using  $\varepsilon$ -moves only. This may be extended also to sets

$$\Lambda(S) = \bigcup_{q \in S} \Lambda(q), \quad \forall S \subseteq Q.$$

Clearly, for all  $q \in Q$  and  $S \subseteq Q$  both  $\Lambda(q)$  and  $\Lambda(S)$  may be computed. Suppose in the sequel that these are given.

The following algorithm determine the transitions using the transition function  $\bar{\delta}$ , which is defined in line 5.

If  $|Q| = n$  and  $|\Sigma| = m$ , then lines 2–6 show that the running time in worst case is  $O(n^2m)$ .

ELIMINATE-EPSILON-MOVES( $A, \bar{A}$ )

```

1  $\bar{F} \leftarrow F \cup \{q \in I \mid \Lambda(q) \cap F \neq \emptyset\}$ 
2 for all  $q \in Q$ 
3   do for all  $a \in \Sigma$ 
4     do  $\Delta \leftarrow \bigcup_{p \in \Lambda(q)} \delta(p, a)$ 
5        $\bar{\delta}(q, a) \leftarrow \Delta \cup \left( \bigcup_{p \in \Delta} \Lambda(p) \right)$ 
6  $\bar{E} \leftarrow \{(p, a, q), \mid p, q \in Q, a \in \Sigma, q \in \bar{\delta}(p, a)\}$ 

```

**Example 2.15** Consider the FA with  $\varepsilon$ -moves in Fig. 2.11. The corresponding transition table is:

$\delta$	0	1	$\varepsilon$
$q_0$	$\emptyset$	$\{q_0\}$	$\{q_1\}$
$q_1$	$\{q_1\}$	$\emptyset$	$\{q_2\}$
$q_2$	$\emptyset$	$\{q_2\}$	$\emptyset$

Apply algorithm ELIMINATE-EPSILON-MOVES.

$\Lambda(q_0) = \{q_0, q_1, q_2\}$ ,  $\Lambda(q_1) = \{q_1, q_2\}$ ,  $\Lambda(q_2) = \{q_2\}$

$\Lambda(I) = \Lambda(q_0)$ , and its intersection with  $F$  is not empty, thus  $\bar{F} = F \cup \{q_0\} = \{q_0, q_2\}$ .

$(q_0, 0)$ :

$$\Delta = \delta(q_0, 0) \cup \delta(q_1, 0) \cup \delta(q_2, 0) = \{q_1\}, \quad \{q_1\} \cup \Lambda(q_1) = \{q_1, q_2\}$$

$$\bar{\delta}(q_0, 0) = \{q_1, q_2\}.$$

$(q_0, 1)$ :

$$\Delta = \delta(q_0, 1) \cup \delta(q_1, 1) \cup \delta(q_2, 1) = \{q_0, q_2\}, \quad \{q_0, q_2\} \cup (\Lambda(q_0) \cup \Lambda(q_2)) = \{q_0, q_1, q_2\}$$

$$\bar{\delta}(q_0, 1) = \{q_0, q_1, q_2\}$$

$(q_1, 0)$ :

$$\Delta = \delta(q_1, 0) \cup \delta(q_2, 0) = \{q_1\}, \quad \{q_1\} \cup \Lambda(q_1) = \{q_1, q_2\}$$

$$\bar{\delta}(q_1, 0) = \{q_1, q_2\}$$

$(q_1, 1)$ :

$$\Delta = \delta(q_1, 1) \cup \delta(q_2, 1) = \{q_2\}, \quad \{q_2\} \cup \Lambda(q_2) = \{q_2\}$$

$$\bar{\delta}(q_1, 1) = \{q_2\}$$

$(q_2, 0)$ :  $\Delta = \delta(q_2, 0) = \emptyset$

$$\bar{\delta}(q_2, 0) = \emptyset$$

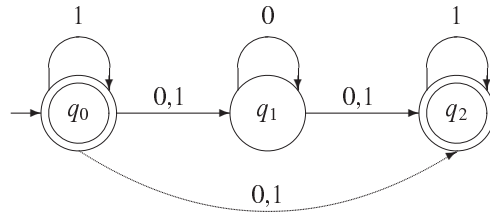


Figure 2.12. NFA equivalent to FA with  $\epsilon$ -moves given in Fig. 2.11.

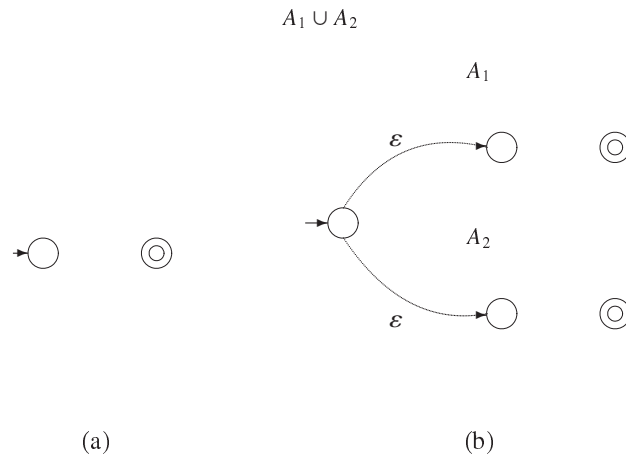


Figure 2.13. (a) Representation of an NFA. Initial states are represented by a circle with an arrow, final states by a double circle. (b) Union of two NFA's.

$(q_2, 1) :$   
 $\Delta = \delta(q_2, 1) = \{q_2\}, \quad \{q_2\} \cup \Lambda(q_2) = \{q_2\}$   
 $\bar{\delta}(q_2, 1) = \{q_2\}.$

The transition table of NFA  $\bar{A}$  is:

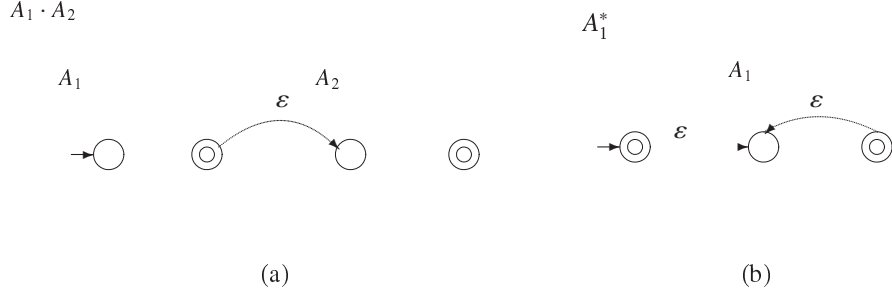
$\bar{\delta}$	0	1
$q_0$	$\{q_1, q_2\}$	$\{q_0, q_1, q_2\}$
$q_1$	$\{q_1, q_2\}$	$\{q_2\}$
$q_2$	$\emptyset$	$\{q_2\}$

and the transition graph is in Fig. 2.12.

Define regular operations on NFA: union, product and iteration. The result will be an FA with  $\epsilon$ -moves.

Operation will be given also by diagrams. An NFA is given as in Fig. 2.13(a). Initial states are represented by a circle with an arrow, final states by a double circle.

Let  $A_1 = (Q_1, \Sigma_1, E_1, I_1, F_1)$  and  $A_2 = (Q_2, \Sigma_2, E_2, I_2, F_2)$  be NFA. The result of any operation is a FA with  $\epsilon$ -moves  $A = (Q, \Sigma, E, I, F)$ . Suppose that  $Q_1 \cap Q_2 = \emptyset$  always. If



**Figure 2.14.** (a) Product of two FA. (b) Iteration of an FA.

not, we can rename the elements of any set of states.

*Union.*  $A = A_1 \cup A_2$ , where

$$Q = Q_1 \cup Q_2 \cup \{q_0\},$$

$$\Sigma = \Sigma_1 \cup \Sigma_2,$$

$$I = \{q_0\},$$

$$F = F_1 \cup F_2,$$

$$E = E_1 \cup E_2 \cup \bigcup_{q \in I_1 \cup I_2} \{(q_0, \varepsilon, q)\}.$$

For the result of the union see Fig. 2.13(b). The result is the same if instead of a single initial state we choose as set of initial states the union  $I_1 \cup I_2$ . In this case the result automaton will be without  $\varepsilon$ -moves. By the definition it is easy to see that  $L(A_1 \cup A_2) = L(A_1) \cup L(A_2)$ .

*Product.*  $A = A_1 \cdot A_2$ , where

$$Q = Q_1 \cup Q_2,$$

$$\Sigma = \Sigma_1 \cup \Sigma_2,$$

$$F = F_2,$$

$$I = I_1,$$

$$E = E_1 \cup E_2 \cup \bigcup_{\substack{p \in F_1 \\ q \in I_2}} \{(p, \varepsilon, q)\}$$

For the result automaton see Fig. 2.14(a). Here also  $L(A_1 \cdot A_2) = L(A_1)L(A_2)$ .

*Iteration.*  $A = A_1^*$ , where

$$Q = Q_1 \cup \{q_0\},$$

$$\Sigma = \Sigma_1,$$

$$F = F_1 \cup \{q_0\},$$

$$I = \{q_0\}$$

$$E = E_1 \cup \bigcup_{p \in I_1} \{(q_0, \varepsilon, p)\} \cup \bigcup_{\substack{q \in F_1 \\ p \in I_1}} \{(q, \varepsilon, p)\}.$$

The iteration of an FA can be seen in Fig. 2.14(b). For this operation it is also true that  $L(A_1^*) = (L(A_1))^*$ .

The definition of these tree operations proves again that regular languages are closed under the regular operations.

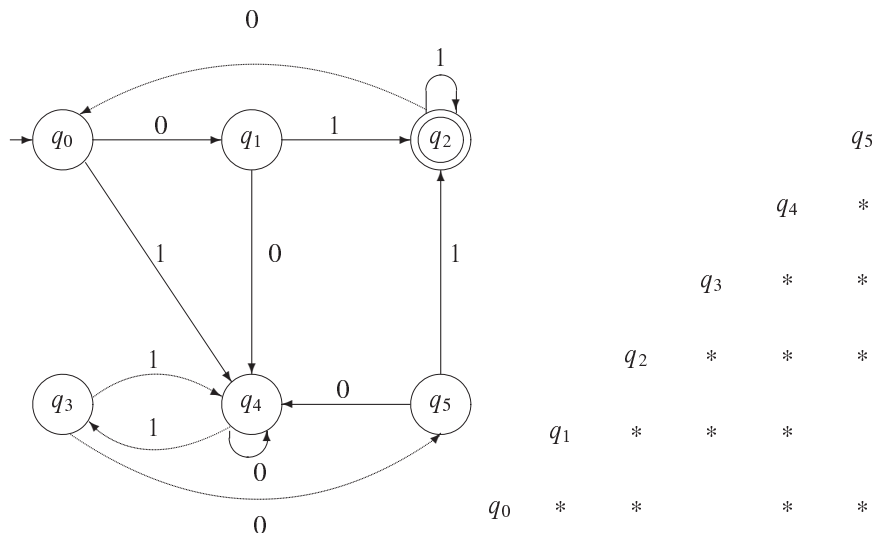


Figure 2.15. Minimization of DFA.

### 2.2.5. Minimization of finite automata

A DFA  $A = (Q, \Sigma, E, \{q_0\}, F)$  is called **minimum state automaton** if for any equivalent complete DFA  $A' = (Q', \Sigma, E', \{q'_0\}, F')$  it is true that  $|Q| \leq |Q'|$ . We give an algorithm which builds for any complete DFA an equivalent minimum state automaton.

States  $p$  and  $q$  of an DFA  $A = (Q, \Sigma, E, \{q_0\}, F)$  are **equivalent** if for arbitrary word  $u$  we reach from both either final or nonfinal states, that is

$$p \equiv q \text{ if for any word } u \in \Sigma^* \left\{ \begin{array}{l} p \xrightarrow{u} r, r \in F \text{ and } q \xrightarrow{u} s, s \in F \text{ or} \\ p \xrightarrow{u} r, r \notin F \text{ and } q \xrightarrow{u} s, s \notin F . \end{array} \right.$$

If two states are not equivalent, then they are distinguishable. In the following algorithm the distinguishable states will be marked by a star, and equivalent states will be merged. The algorithm will associate list of pair of states with some pair of states expecting a later marking by a star, that is if we mark a pair of states by a star, then all pairs on the associated list will be also marked by a star. The algorithm is given for DFA without inaccessible states. The used DFA is complete, so  $\delta(p, a)$  contains exact one element, function *elem* defined on page 96, which gives the unique element of the set, will be also used here.



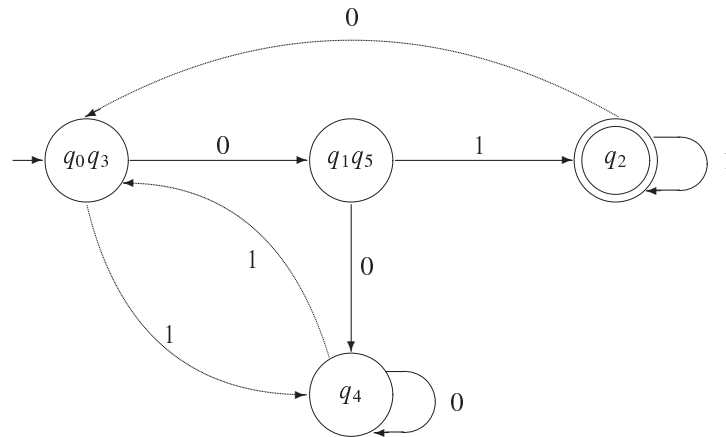


Figure 2.16. Minimum automaton equivalent with DFA in Fig. 2.15.

#### AUTOMATON-MINIMIZATION(A)

- 1 mark with a star all pairs of states  $\{p, q\}$  for which  $p \in F$  and  $q \notin F$  or  $p \notin F$  and  $q \in F$ ,
- 2 associate an empty list with each unmarked pair  $\{p, q\}$ ,
- 3 **for** all unmarked pair of states  $\{p, q\}$  and for all symbol  $a \in \Sigma$  examine pairs of states  $\{elem(\delta(p, a)), elem(\delta(q, a))\}$ ,  
**if** any of these pairs is marked,  
**then** mark also pair  $\{p, q\}$  with all the elements on the list before associated with pair  $\{p, q\}$ ,  
**else, if** all the above pairs are unmarked,  
**then** put pair  $\{p, q\}$  on each list associated with pairs  $\{elem(\delta(p, a)), elem(\delta(q, a))\}$ , unless  $\delta(p, a) = \delta(q, a)$ ,
- 4 merge all unmarked (equivalent) pairs.

After finishing the algorithm, if a cell of the table does not contain a star, then the states corresponding to its row and column index, are equivalent and may be merged. Merging states is continued until it is possible. We can say that the equivalence relation decomposes the set of states in equivalence classes, and the states in such a class may be all merged.

*Remark.* The above algorithm can be used also in the case of an DFA which is not complete, that is there are states for which does not exist transition. Then a pair  $\{\emptyset, \{q\}\}$  may occur, and if  $q$  is a final state, consider this pair marked.

**Example 2.16** Let be the DFA in Fig. 2.15. We will use a table for marking pairs with a star. Marking pair  $\{p, q\}$  means putting a star in the cell corresponding to row  $p$  and column  $q$  (or row  $q$  and column  $p$ ).

First we mark pairs  $\{q_2, q_0\}$ ,  $\{q_2, q_1\}$ ,  $\{q_2, q_3\}$ ,  $\{q_2, q_4\}$  and  $\{q_2, q_5\}$  (because  $q_2$  is the single final state). Then consider all unmarked pairs and examine them as the algorithm requires. Let us begin with pair  $\{q_0, q_1\}$ . Associate with it pairs  $\{elem(\delta(q_0, 0)), elem(\delta(q_1, 0))\}$ ,  $\{elem(\delta(q_0, 1)), elem(\delta(q_1, 1))\}$ ,

that is  $\{q_1, q_4\}$ ,  $\{q_4, q_2\}$ . Because pair  $\{q_4, q_2\}$  is already marked, mark also pair  $\{q_0, q_1\}$ .

In the case of pair  $\{q_0, q_3\}$  the new pairs are  $\{q_1, q_5\}$  and  $\{q_4, q_4\}$ . With pair  $\{q_1, q_5\}$  associate pair  $\{q_0, q_3\}$  on a list, that is

$$\{q_1, q_5\} \longrightarrow \{q_0, q_3\} .$$

Now continuing with  $\{q_1, q_5\}$  one obtain pairs  $\{q_4, q_4\}$  and  $\{q_2, q_2\}$ , with which nothing are associated by algorithm.

Continue with pair  $\{q_0, q_4\}$ . The associated pairs are  $\{q_1, q_4\}$  and  $\{q_4, q_3\}$ . None of them are marked, so associate with them on a list pair  $\{q_0, q_4\}$ , that is

$$\{q_1, q_4\} \longrightarrow \{q_0, q_4\}, \quad \{q_4, q_3\} \longrightarrow \{q_0, q_4\} .$$

Now continuing with  $\{q_1, q_4\}$  we get the pairs  $\{q_4, q_4\}$  and  $\{q_2, q_3\}$ , and because this latter is marked we mark pair  $\{q_1, q_4\}$  and also pair  $\{q_0, q_4\}$  associated to it on a list. Continuing we will get the table in Fig. 2.15, that is we get that  $q_0 \equiv q_3$  and  $q_1 \equiv q_5$ . After merging them we get an equivalent minimum state automaton (see Fig. 2.16).

### 2.2.6. Pumping lemma for regular languages

The following theorem, called *pumping lemma* for historical reasons, may be efficiently used to prove that a language is not regular. It is a sufficient condition for a regular language.

**Theorem 2.15** (pumping lemma). *For any regular language  $L$  there exists a natural number  $n \geq 1$  (depending only on  $L$ ), such that any word  $u$  of  $L$  with length at least  $n$  may be written as  $u = xyz$  such that*

- (1)  $|xy| \leq n$ ,
- (2)  $|y| \geq 1$ ,
- (3)  $xy^i z \in L$  for all  $i = 0, 1, 2, \dots$

**Proof.** If  $L$  is a regular language, then there is such an DFA which accepts  $L$  (by theorems 2.12 and 2.10). Let  $A = (Q, \Sigma, E, \{q_0\}, F)$  be this DFA, so  $L = L(A)$ . Let  $n$  be the number of its states, that is  $|Q| = n$ . Let  $u = a_1 a_2 \dots a_m \in L$  and  $m \geq n$ . Then, because the automaton accepts word  $u$ , there are states  $q_0, q_1, \dots, q_m$  and walk

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_{m-1}} q_{m-1} \xrightarrow{a_m} q_m, \quad q_m \in F.$$

Because the number of states is  $n$  and  $m \geq n$ , by the pigeonhole principle<sup>3</sup> states  $q_0, q_1, \dots, q_m$  can not all be distinct (see Fig. 2.17), there are at least two of them which are equal. Let  $q_j = q_k$ , where  $j < k$  and  $k$  is the least such index. Then  $j < k \leq n$ . Decompose word  $u$  as:

$$\begin{aligned} x &= a_1 a_2 \dots a_j \\ y &= a_{j+1} a_{j+2} \dots a_k \\ z &= a_{k+1} a_{k+2} \dots a_m. \end{aligned}$$

This decomposition immediately yields to  $|xy| \leq n$  and  $|y| \geq 1$ . We will prove that  $xy^i z \in L$  for any  $i$ .

Because  $u = xyz \in L$ , there exists an walk

$$q_0 \xrightarrow{x} q_j \xrightarrow{y} q_k \xrightarrow{z} q_m, \quad q_m \in F,$$

<sup>3</sup>*Pigeonhole principle:* If we have to put more than  $k$  objects into  $k$  boxes, then at least two boxes will contain at least two objects.

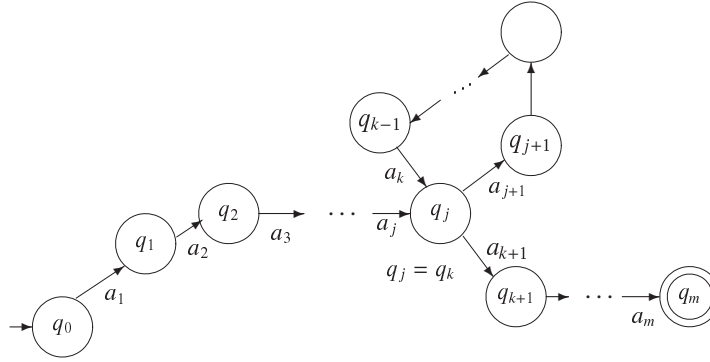


Figure 2.17. Sketch of DFA used in the proof of the pumping lemma.

and because of  $q_j = q_k$ , this may be written also as

$$q_0 \xrightarrow{x} q_j \xrightarrow{y} q_j \xrightarrow{z} q_m, q_m \in F .$$

From this walk  $q_j \xrightarrow{y} q_j$  can be omitted or can be inserted many times. So, there are the following walks:

$$q_0 \xrightarrow{x} q_j \xrightarrow{z} q_m, q_m \in F ,$$

$$q_0 \xrightarrow{x} q_j \xrightarrow{y} q_j \xrightarrow{y} \dots \xrightarrow{y} q_j \xrightarrow{z} q_m, q_m \in F .$$

Therefore  $xy^i z \in L$  for all  $i$ , and this proves the theorem. ■

**Example 2.17** We use the pumping lemma to show that  $L_1 = \{a^k b^k \mid k \geq 1\}$  is not regular. Assume that  $L_1$  is regular, and let  $n$  be the corresponding natural number given by the pumping lemma. Because the length of the word  $u = a^n b^n$  is  $2n$ , this word can be written as in the lemma. We prove that this leads to a contradiction. Let  $u = xyz$  be the decomposition as in the lemma. Then  $|xy| \leq n$ , so  $x$  and  $y$  can contain no other letters than  $a$ , and because we must have  $|y| \geq 1$ , word  $y$  contains at least one  $a$ . Then  $xy^i z$  for  $i \neq 1$  will contain a different number of  $a$ 's and  $b$ 's, therefore  $xy^i z \notin L_1$  for any  $i \neq 1$ . This is a contradiction with the third assertion of the lemma, this is why that assumption that  $L_1$  is regular, is false. Therefore  $L_1 \notin \mathcal{L}_3$ .

Because the context-free grammar  $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow ab, S \rightarrow aSb\}, S)$  generates language  $L_1$ , we have  $L_1 \in \mathcal{L}_2$ . From these two follow that  $\mathcal{L}_3 \subset \mathcal{L}_2$ .

**Example 2.18** We show that  $L_2 = \{u \in \{0, 1\}^* \mid n_0(u) = n_1(u)\}$  is not regular. ( $n_0(u)$  is the number of 0's in  $u$ , while  $n_1(u)$  the number of 1's).

We proceed as in the previous example using here word  $u = 0^n 1^n$ , where  $n$  is the natural number associated by lemma to language  $L_2$ .

**Example 2.19** We prove, using the pumping lemma, that  $L_3 = \{uu \mid u \in \{a, b\}^*\}$  is not a regular language. Let  $w = a^n b a^n b = xy z$  be, where  $n$  here is also the natural number associated to  $L_3$  by the pumping lemma. From  $|xy| \leq n$  we have that  $y$  contains no other letters than  $a$ , but it contains at least one. By lemma we have  $xz \in L_3$ , that is not possible. Therefore  $L_3$  is not regular.

Pumping lemma has several interesting consequences.

**Corollary 2.16** *Regular language  $L$  is not empty if and only if there exists a word  $u \in L$ ,  $|u| < n$ , where  $n$  is the natural number associated to  $L$  by the pumping lemma.*

**Proof.** The assertion in a direction is obvious: if there exists a word shorter than  $n$  in  $L$ , then  $L \neq \emptyset$ . Conversely, let  $L \neq \emptyset$  and let  $u$  be the shortest word in  $L$ . We show that  $|u| < n$ . If  $|u| \geq n$ , then we apply the pumping lemma, and give the decomposition  $u = xyz$ ,  $|y| > 1$  and  $xz \in L$ . This is a contradiction, because  $|xz| < |u|$  and  $u$  is the shortest word in  $L$ . Therefore  $|u| < n$ . ■

**Corollary 2.17** *There exists an algorithm that can decide if a regular language is or not empty.*

**Proof.** Assume that  $L = L(A)$ , where  $A = (Q, \Sigma, E, \{q_0\}, F)$  is a DFA. By consequence 2.16 and theorem 2.15 language  $L$  is not empty if and only if it contains a word shorter than  $n$ , where  $n$  is the number of states of automaton  $A$ . By this it is enough to decide that there is a word shorter than  $n$  which is accepted by automaton  $A$ . Because the number of words shorter than  $n$  is finite, the problem can be decided. ■

When we had given an algorithm for inaccessible states of a DFA, we remarked that the procedure can be used also to decide if the language accepted by that automaton is or not empty. Because finite automata accept regular languages, we can consider to have already two procedures to decide if a regular languages is or not empty. Moreover, we have a third procedure, if we take into account that the algorithm for finding productive states also can be used to decide on a regular language when it is empty.

**Corollary 2.18** *A regular language  $L$  is infinite if and only if there exists a word  $u \in L$  such that  $n \leq |u| < 2n$ , where  $n$  is the natural number associated to language  $L$ , given by the pumping lemma.*

**Proof.** If  $L$  is infinite, then it contains words longer than  $2n$ , and let  $u$  be the shortest word longer than  $2n$  in  $L$ . Because  $L$  is regular we can use the pumping lemma, so  $u = xyz$ , where  $|xy| \leq n$ , thus  $|y| \leq n$  is also true. By the lemma  $u' = xz \in L$ . But because  $|u'| < |u|$  and the shortest word in  $L$  longer than  $2n$  is  $u$ , we get  $|u'| < 2n$ . From  $|y| \leq n$  we get also  $|u'| \geq n$ .

Conversely, if there exists a word  $u \in L$  such that  $n \leq |u| < 2n$ , then using the pumping lemma, we obtain that  $u = xyz$ ,  $|y| \geq 1$  and  $xy^i z \in L$  for any  $i$ , therefore  $L$  is infinite. ■

Now, the question is: how can we apply the pumping lemma for a finite regular language, since by pumping words we get an infinite number of words? The number of states of a DFA accepting language  $L$  is greater than the length of the longest word in  $L$ . So, in  $L$  there is no word with length at least  $n$ , when  $n$  is the natural number associated to  $L$  by the pumping lemma. Therefore, no word in  $L$  can be decomposed in the form  $xyz$ , where  $|xyz| \geq n$ ,  $|xy| \leq n$ ,  $|y| \geq 1$ , and this is why we can not obtain an infinite number of words in  $L$ .

### 2.2.7. Regular expressions

In this subsection we introduce for any alphabet  $\Sigma$  the notion of regular expressions over  $\Sigma$  and the corresponding representing languages. A regular expression is a formula, and the

$$\begin{aligned}
x + y &\equiv y + x \\
(x + y) + z &\equiv x + (y + z) \\
(xy)z &\equiv x(yz) \\
(x + y)z &\equiv xz + yz \\
x(y + z) &\equiv xy + xz \\
(x + y)^* &\equiv (x^* + y)^* \equiv (x + y^*)^* \equiv (x^* + y^*)^* \\
(x + y)^* &\equiv (xy^*)^* \equiv (x^*y)^* \equiv (x^*y^*)^* \\
(x^*)^* &\equiv x^* \\
x^*x &\equiv xx^* \\
xx^* + \varepsilon &\equiv x^*
\end{aligned}$$

2.1. Table. Properties of regular expressions.

corresponding language is a language over  $\Sigma$ . For example, if  $\Sigma = \{a, b\}$ , then  $a^*$ ,  $b^*$ ,  $a^* + b^*$  are regular expressions over  $\Sigma$  which represent respectively languages  $\{a\}^*$ ,  $\{b\}^*$ ,  $\{a\}^* \cup \{b\}^*$ . The exact definition is the following.

**Definition 2.19** Define recursively a regular expression over  $\Sigma$  and the language it represent.

- $\emptyset$  is a regular expression representing the empty language.
- $\varepsilon$  is a regular expression representing language  $\{\varepsilon\}$ .
- If  $a \in \Sigma$ , then  $a$  is a regular expression representing language  $\{a\}$ .
- If  $x, y$  are regular expressions representing languages  $X$  and  $Y$  respectively, then  $(x + y)$ ,  $(xy)$ ,  $(x^*)$  are regular expressions representing languages  $X \cup Y$ ,  $XY$  and  $X^*$  respectively.

Regular expression over  $\Sigma$  can be obtained only by using the above rules a finite number of times.

Some brackets can be omitted in the regular expressions if taking into account the priority of operations (iteration, product, union) the corresponding languages are not affected. For example instead of  $((x^*)(x + y))$  we can consider  $x^*(x + y)$ .

Two regular expressions are **equivalent** if they represent the same language, that is  $x \equiv y$  if  $X = Y$ , where  $X$  and  $Y$  are the languages represented by regular expressions  $x$  and  $y$  respectively. Table 2.1 shows some equivalent expressions.

We show that to any finite language  $L$  can be associated a regular expression  $x$  which represent language  $L$ . If  $L = \emptyset$ , then  $x = \emptyset$ . If  $L = \{w_1, w_2, \dots, w_n\}$ , then  $x = x_1 + x_2 + \dots + x_n$ , where for any  $i = 1, 2, \dots, n$  expression  $x_i$  is a regular expression representing language  $\{w_i\}$ . This latter can be done by the following rule. If  $w_i = \varepsilon$ , then  $x_i = \varepsilon$ , else if  $w_i = a_1 a_2 \dots a_m$ , where  $m \geq 1$  depends on  $i$ , then  $x_i = a_1 a_2 \dots a_m$ , where the brackets are omitted.

We prove the theorem of Kleene which refers to the relationship between regular languages and regular expression.

**Theorem 2.20** (Kleene's theorem). Language  $L \subseteq \Sigma^*$  is regular if and only if there exists a regular expression over  $\Sigma$  representing language  $L$ .

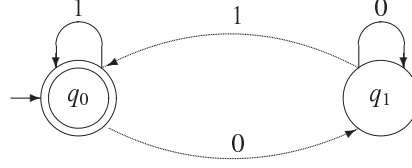


Figure 2.18. DFA from Example 2.20., to which regular expression is associated by Method 1.

**Proof.** First we prove that if  $x$  is a regular expression, then language  $L$  which represents  $x$  is also regular. The proof will be done by induction on the construction of expression.

If  $x = \emptyset$ ,  $x = \varepsilon$ ,  $x = a, \forall a \in \Sigma$ , then  $L = \emptyset$ ,  $L = \{\varepsilon\}$ ,  $L = \{a\}$  respectively. Since  $L$  is finite in all three cases, it is also regular.

If  $x = (x_1 + x_2)$ , then  $L = L_1 \cup L_2$ , where  $L_1$  and  $L_2$  are the languages which represent the regular expressions  $x_1$  and  $x_2$  respectively. By the induction hypothesis languages  $L_1$  and  $L_2$  are regular, so  $L$  is also regular because regular languages are closed on union. Cases  $x = (x_1x_2)$  and  $x = (x_1^*)$  can be proved by similar way.

Conversely, we prove that if  $L$  is a regular language, then a regular expression  $x$  can be associated to it, which represent exactly the language  $L$ . If  $L$  is regular, then there exists a DFA  $A = (Q, \Sigma, E, \{q_0\}, F)$  for which  $L = L(A)$ . Let  $q_0, q_1, \dots, q_n$  the states of the automaton  $A$ . Define languages  $R_{ij}^k$  for all  $-1 \leq k \leq n$  and  $0 \leq i, j \leq n$ .  $R_{ij}^k$  is the set of words, for which automaton  $A$  goes from state  $q_i$  to state  $q_j$  without using any state with index greater than  $k$ . Using transition graph we can say: a word is in  $R_{ij}^k$ , if from state  $q_i$  we arrive to state  $q_j$  following the edges of the graph, and concatenating the corresponding labels on edges we get exactly that word, not using any state  $q_{k+1}, \dots, q_n$ . Sets  $R_{ij}^k$  can be done also formally:

$$R_{ij}^{-1} = \{a \in \Sigma \mid (q_i, a, q_j) \in E\}, \text{ if } i \neq j,$$

$$R_{ii}^{-1} = \{a \in \Sigma \mid (q_i, a, q_i) \in E\} \cup \{\varepsilon\},$$

$$R_{ij}^k = R_{ij}^{k-1} \cup R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \text{ for all } i, j, k \in \{0, 1, \dots, n\}.$$

We can prove by induction that sets  $R_{ij}^k$  can be described by regular expressions. Indeed, if  $k = -1$ , then for all  $i$  and  $j$  languages  $R_{ij}^k$  are finite, so they can be expressed by regular expressions representing exactly these languages. Moreover, if for all  $i$  and  $j$  language  $R_{ij}^{k-1}$  can be expressed by regular expression, then language  $R_{ij}^k$  can be expressed also by regular expression, which can be corresponding constructed from regular expressions representing languages  $R_{ij}^{k-1}$ ,  $R_{ik}^{k-1}$ ,  $R_{kk}^{k-1}$  and  $R_{kj}^{k-1}$  respectively, using the above formula for  $R_{ij}^k$ .

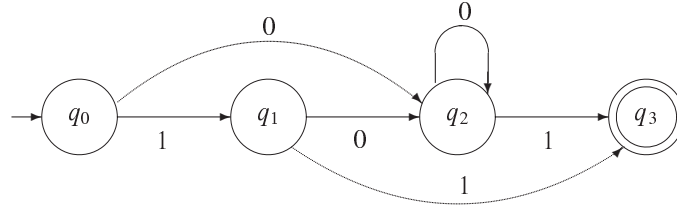
Finally, if  $F = \{q_{i_1}, q_{i_2}, \dots, q_{i_p}\}$  is the set of final states of the DFA  $A$ , then  $L = L(A) = R_{0i_1}^n \cup R_{0i_2}^n \cup \dots \cup R_{0i_p}^n$  can be expressed by a regular expression obtained from expressions representing languages  $R_{0i_1}^n, R_{0i_2}^n, \dots, R_{0i_p}^n$  using operation  $+$ . ■

Further on we give some procedures which associate DFA to regular expressions and conversely regular expression to DFA.

### Associating regular expressions to finite automata

We present here three methods, each of which associate to a DFA the corresponding regular expression.

*Method 1.* Using the result of the theorem of Kleene, we will construct the sets  $R_{ij}^k$ , and write a regular expression which represent the language  $L = R_{0i_1}^n \cup R_{0i_2}^n \cup \dots \cup R_{0i_p}^n$ , where



**Figure 2.19.** DFA in example 2.21, to which a regular expression is associated by Method 1. The computation are in the Table 2.2.

$F = \{q_{i_1}, q_{i_2}, \dots, q_{i_p}\}$  is the set of final states of the automaton.

**Example 2.20** Consider the DFA in Fig. 2.18.

$$L(A) = R_{00}^1 = R_{00}^0 \cup R_{01}^0 (R_{11}^0)^* R_{10}^0$$

$$R_{00}^0 : 1^* + \varepsilon \equiv 1^*$$

$$R_{01}^0 : 1^*0$$

$$R_{11}^0 : 11^*0 + \varepsilon + 0 \equiv (11^* + \varepsilon)0 + \varepsilon \equiv 1^*0 + \varepsilon$$

$$R_{10}^0 : 11^*$$

Then the regular expression corresponding to  $L(A)$  is  $1^* + 1^*0(1^*0 + \varepsilon)^*11^* \equiv 1^* + 1^*0(1^*0)^*11^*$ .

**Example 2.21** Find a regular expression associated to DFA in Fig. 2.19. The computations are in Table 2.2. The regular expression corresponding to  $R_{03}^3$  is  $11 + (0 + 10)^*1$ .

*Method 2.* Now we generalize the notion of finite automaton, considering words instead of letters as labels of edges. In such an automaton each walk determine a regular expression, which determine a regular language. The regular language accepted by a generalized finite automaton is the union of regular languages determined by the productive walks. It is easy to see that the generalized finite automata accept regular languages.

The advantage of generalized finite automata is that the number of its edges can be diminished by equivalent transformations, which do not change the accepted language, and leads to a graph with only one edge which label is exactly the accepted language.

The possible equivalent transformations can be seen in Fig. 2.20. If some of the vertices 1, 2, 4, 5 on the figure coincide, in the result they are merged, and a loop will arrive.

First, the automaton is transformed by corresponding  $\varepsilon$ -moves to have only one initial and one final state. Then, applying the equivalent transformations until the graph will have only one edge, we will obtain as the label of this edge the regular expression associated to the automaton.

**Example 2.22** In the case of Fig. 2.18 the result is obtained by steps illustrated in Fig. 2.21. This result is  $(1 + 00^*1)^*$ , which represents the same language as obtained by Method 1 (See example 2.20).

**Example 2.23** In the case of Fig. 2.19 is not necessary to introduce new initial and final state. The

	$k = -1$	$k = 0$	$k = 1$	$k = 2$	$k = 3$
$R_{00}^k$	$\varepsilon$	$\varepsilon$	$\varepsilon$	$\varepsilon$	
$R_{01}^k$	1	1	1	1	
$R_{02}^k$	0	0	$0 + 10$	$(0 + 10)0^*$	
$R_{03}^k$	$\emptyset$	$\emptyset$	11	$11 + (0 + 10)0^*1$	$11 + (0 + 10)0^*1$
$R_{11}^k$	$\varepsilon$	$\varepsilon$	$\varepsilon$	$\varepsilon$	
$R_{12}^k$	0	0	0	$00^*$	
$R_{13}^k$	1	1	1	$1 + 00^*1$	
$R_{22}^k$	$0 + \varepsilon$	$0 + \varepsilon$	$0 + \varepsilon$	$0^*$	
$R_{23}^k$	1	1	1	$0^*1$	
$R_{33}^k$	$\varepsilon$	$\varepsilon$	$\varepsilon$	$\varepsilon$	

**2.2. Table.** Determining a regular expression associated to DFA in Fig. 2.19 using sets  $R_{ij}^k$ .

steps of transformations can be seen in Fig. 2.22. The resulted regular expression can be written also as  $(0 + 10)0^*1 + 11$ , which is the same as obtained by the previous method.

*Method 3.* The third method for writing regular expressions associated to finite automata uses formal equations. A variable  $X$  is associated to each state of the automaton (to different states different variables). Associate to each state an equation which left side contains  $X$ , its right side contains sum of terms of form  $Ya$  or  $\varepsilon$ , where  $Y$  is a variable associated to a state, and  $a$  is its corresponding input symbol. If there is no incoming edge in the state corresponding to  $X$  then the right side of the equation with left side  $X$  contains  $\varepsilon$ , otherwise is the sum of all terms of the form  $Ya$  for which there is a transition labelled with letter  $a$  from state corresponding to  $Y$  to the state corresponding to  $X$ . If the state corresponding to



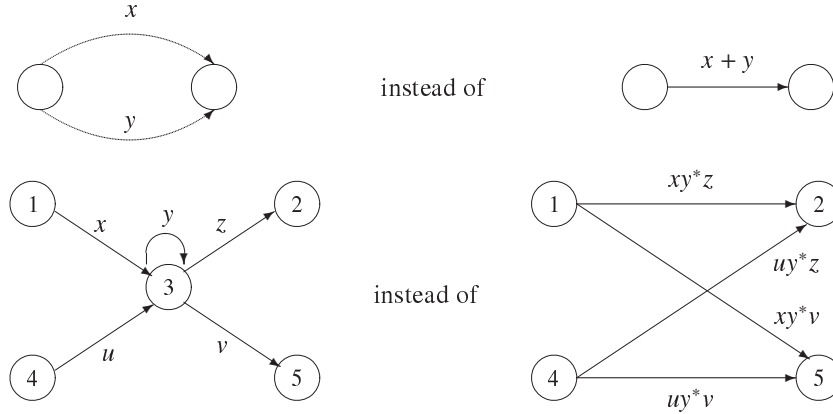


Figure 2.20. Possible equivalent transformations for finding regular expression associated to an automaton.

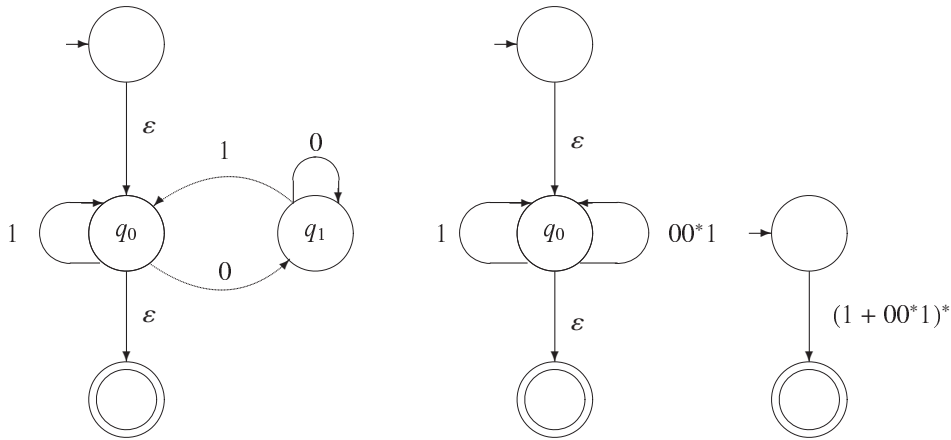


Figure 2.21. Transformation of the finite automaton in Fig. 2.18.

$X$  is also an initial and a final state, then on right side of the equation with the left side  $X$  will be also a term equal to  $\epsilon$ . For example in the case of Fig. 2.19 let these variable  $X, Y, Z, U$  corresponding to the states  $q_0, q_1, q_2, q_3$ . The corresponding equation are

$$\begin{aligned} X &= \epsilon \\ Y &= X1 \\ Z &= X0 + Y0 + Z0 \\ U &= Y1 + Z1. \end{aligned}$$

If an equation is of the form  $X = X\alpha + \beta$ , where  $\alpha, \beta$  are arbitrary words not containing  $X$ , then it is easy to see by a simple substitution that  $X = \beta\alpha^*$  is a solution of the equation.

Because these equations are linear, all of them can be written in the form  $X = X\alpha + \beta$  or  $X = X\alpha$ , where  $\alpha$  do not contain any variable. Substituting this in the other equations the number of remaining equations will be diminished by one. In such a way the system of equation can be solved for each variable.

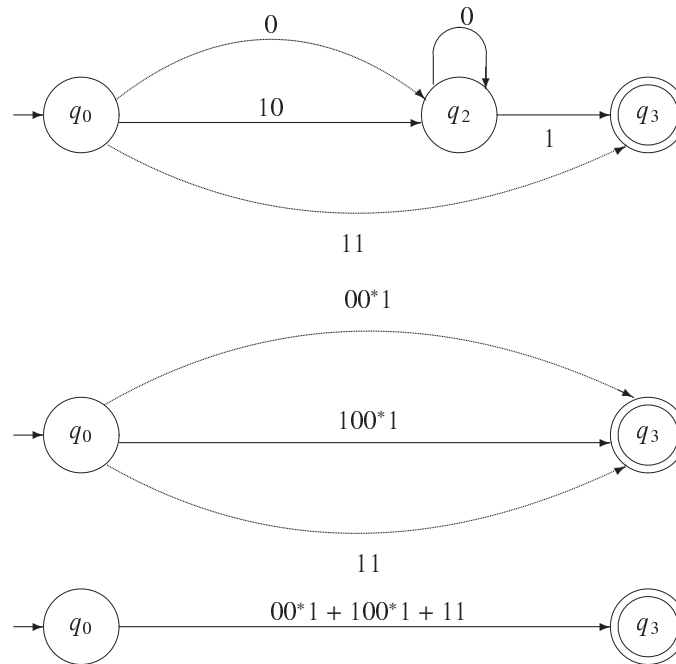


Figure 2.22. Steps of example 2.23..

The solution will be given by variables corresponding to final states summing the corresponding regular expressions.

In our example from the first equation we get  $Y = 1$ . From here  $Z = 0 + 10 + Z0$ , or  $Z = Z0 + (0 + 10)$ , and solving this we get  $Z = (0 + 10)0^*$ . Variable  $U$  can be obtained immediately and we obtain  $U = 11 + (0 + 10)0^*1$ .

Using this method in the case of Fig. 2.18, the following equations will be obtained

$$X = \varepsilon + X1 + Y1$$

$$Y = X0 + Y0$$

Therefore

$$X = \varepsilon + (X + Y)1$$

$$Y = (X + Y)0.$$

Adding the two equations we will obtain

$X + Y = \varepsilon + (X + Y)(0 + 1)$ , from where (considering  $\varepsilon$  as  $\beta$  and  $(0 + 1)$  as  $\alpha$ ) we get the result

$$X + Y = (0 + 1)^*.$$

From here the value of  $X$  after the substitution is

$$X = \varepsilon + (0 + 1)^*1,$$

which is equivalent to the expression obtained using the other methods.

### Associating finite automata to regular expressions

Associate to the regular expression  $r$  a generalized finite automaton:

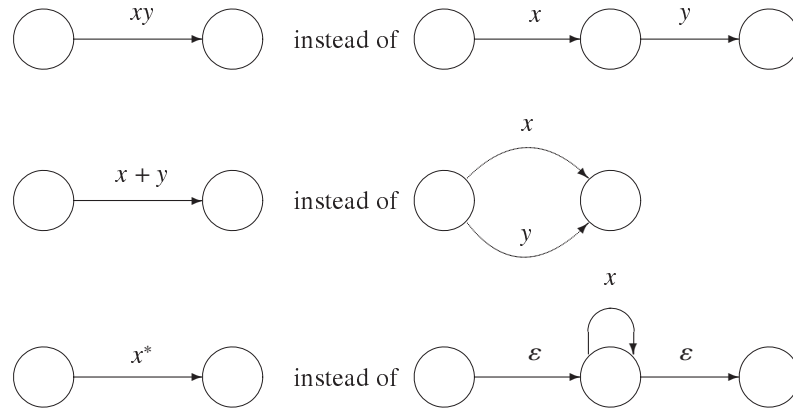
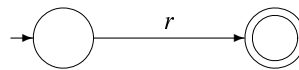


Figure 2.23. Possible transformations to obtain finite automaton associated to a regular expression.



After this, use the transformations in Fig. 2.23 step by step, until an automaton with labels equal to letters from  $\Sigma$  or  $\varepsilon$  will be obtained.

**Example 2.24** Get started from regular expression  $\varepsilon + (0+1)^*1$ . The steps of transformations are in Fig. 2.24(a)-(e). The last finite automaton (see Fig. 2.24(e)) can be done in a simpler form as can be seen in Fig. 2.24(f). After eliminating the  $\varepsilon$ -moves and transforming in a deterministic finite automaton the DFA in Fig. 2.25 will be obtained, which is equivalent to DFA in Fig. 2.18.

### Exercises

- 2.2-1** Give a DFA which accepts natural numbers divisible by 9.
- 2.2-2** Give a DFA which accepts the language containing all words formed by
- an even number of 0's and an even number of 1's,
  - an even number of 0's and an odd number of 1's,
  - an odd number of 0's and an even number of 1's,
  - an odd number of 0's and an odd number of 1's.
- 2.2-3** Give a DFA to accept respectively the following languages:
- $$L_1 = \{a^n b^m \mid n \geq 1, m \geq 0\}, \quad L_2 = \{a^n b^m \mid n \geq 1, m \geq 1\},$$
- $$L_3 = \{a^n b^m \mid n \geq 0, m \geq 0\}, \quad L_4 = \{a^n b^m \mid n \geq 0, m \geq 1\}.$$
- 2.2-4** Give an NFA which accepts words containing at least two 0's and any number of 1's. Give an equivalent DFA.
- 2.2-5** Minimize the DFA's in Fig. 2.26.
- 2.2-6** Show that the DFA in 2.27.(a) is a minimum state automaton.
- 2.2-7** Transform NFA in Fig. 2.27.(b) in a DFA, and after this minimize it.
- 2.2-8** Define finite automaton  $A_1$  which accepts all words of the form  $0(10)^n$  ( $n \geq 0$ ), and finite automaton  $A_2$  which accepts all words of the form  $1(01)^n$  ( $n \geq 0$ ). Define the union

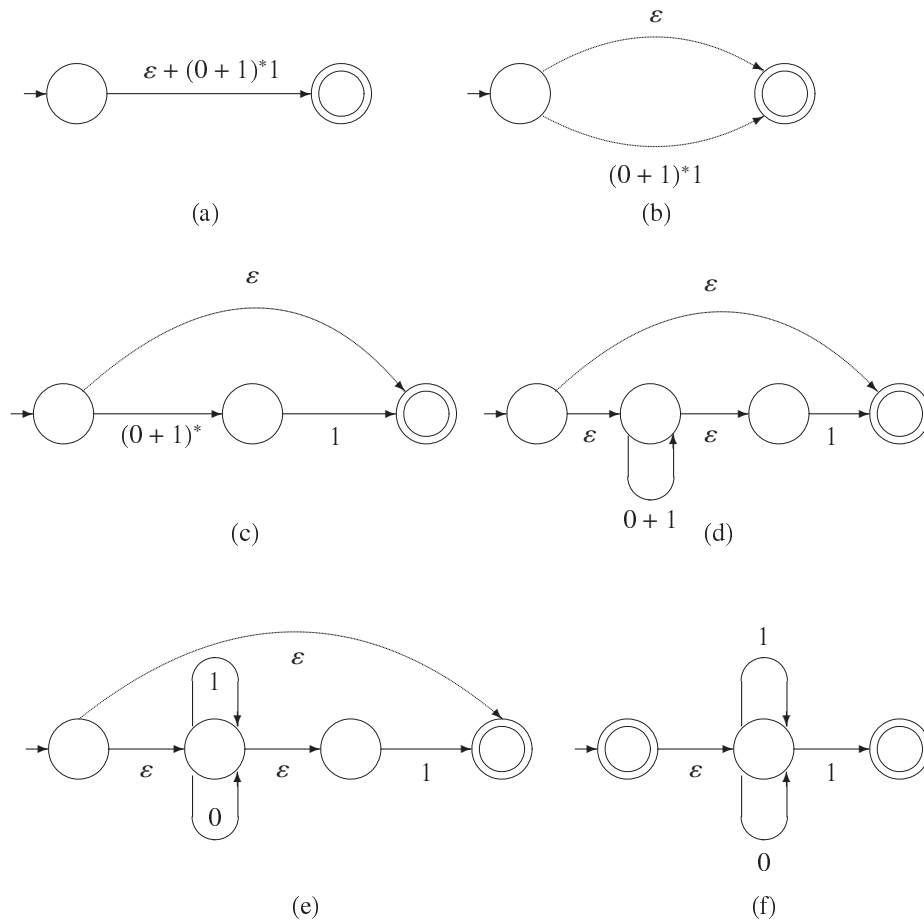


Figure 2.24. Associating finite automaton to regular expression  $\epsilon + (0 + 1)^*1$ .

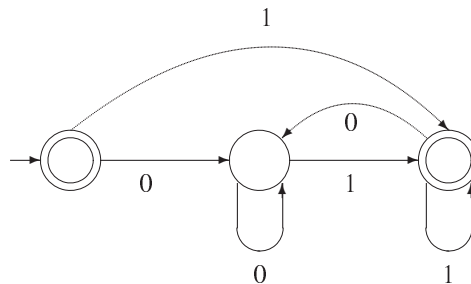


Figure 2.25. Finite automaton associated to regular expression  $\epsilon + (0 + 1)^*1$ .

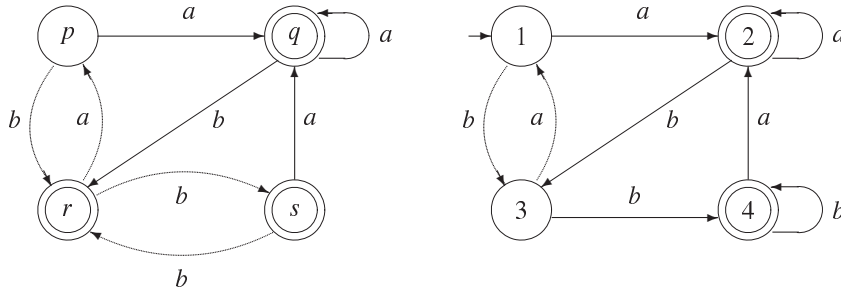


Figure 2.26. DFA's to minimize for exercise 2.2-5..

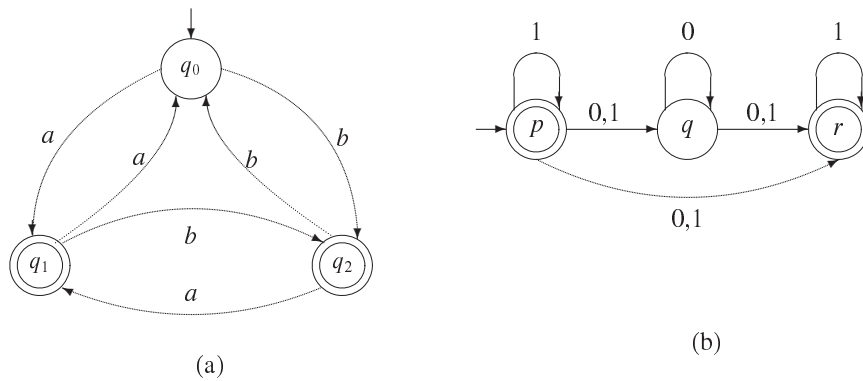


Figure 2.27. Finite automata for exercises 2.2-6. and 2.2-7..

automaton  $A_1 \cup A_2$ , and then eliminate the  $\epsilon$ -moves.

**2.2-9** Associate to DFA in Fig. 2.28 a regular expression.

**2.2-10** Associate to regular expression  $ab^*ba^* + b + ba^*a$  a DFA.

**2.2-11** Prove, using the pumping lemma, that none of the following languages are regular:

$$L_1 = \{a^n cb^n \mid n \geq 0\}, \quad L_2 = \{a^n b^n a^n \mid n \geq 0\}, \quad L_3 = \{a^p \mid p \text{ prim}\}.$$

**2.2-12** Prove that if  $L$  is a regular language, then  $\{u^{-1} \mid u \in L\}$  is also regular.

**2.2-13** Prove that if  $L \subseteq \Sigma^*$  is a regular language, then the following languages are also regular.

$$\text{pre}(L) = \{w \in \Sigma^* \mid \exists u \in \Sigma^*, wu \in L\}, \quad \text{suf}(L) = \{w \in \Sigma^* \mid \exists u \in \Sigma^*, uw \in L\}.$$

**2.2-14** Show that the following languages are all regular.

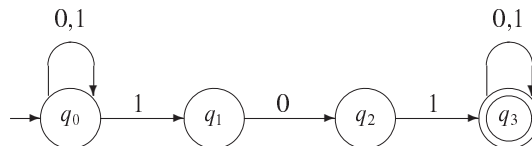


Figure 2.28. DFA for exercise 2.2-9..

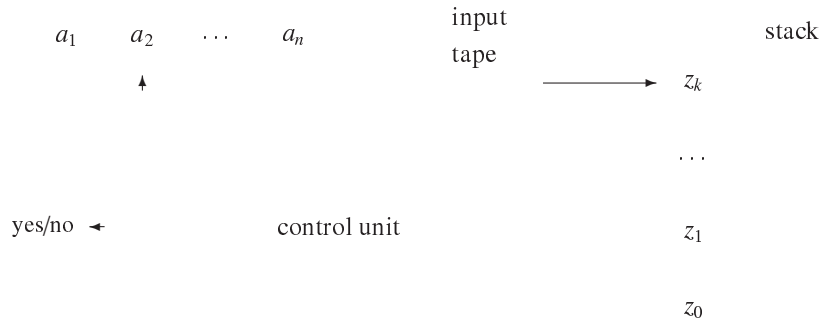


Figure 2.29. Pushdown automaton.

$$L_1 = \{ab^n cd^m \mid n > 0, m > 0\},$$

$$L_2 = \{(ab)^n \mid n \geq 0\},$$

$$L_3 = \{a^{kn} \mid n \geq 0, k \text{ constant}\}.$$

## 2.3. Pushdown automata and context-free languages

In this section we deal with the pushdown automata and the class of languages — the context-free languages — accepted by them.

As we have been seen in the section 2.1, a context-free grammar  $G = (N, T, P, S)$  is one with the productions of the form  $A \rightarrow \beta$ ,  $A \in N$ ,  $\beta \in (N \cup T)^+$ . The production  $S \rightarrow \varepsilon$  is also permitted if  $S$  does not appear in right hand side of any productions. Language  $L(G) = \{u \in T \mid S \xrightarrow{*}_G u\}$  is the context-free language generated by grammar  $G$ .

### 2.3.1. Pushdown automata

We have been seen that finite automata accept the class of regular languages. Now we get to know a new kind of automata, the so-called **pushdown automata**, which accept context-free languages. The pushdown automata differ from finite automata mainly in that to have the possibility to change states without reading any input symbol (i.e. to read the empty symbol) and possess a stack memory, which uses the so-called stack symbols (See Fig. 2.29).

The pushdown automaton get a word as input, start to function from an initial state having in the stack a special symbol, the initial stack symbol. While working, the pushdown automaton change its state based on current state, next input symbol (or empty word) and stack top symbol and replace the top symbol in the stack with a (possibly empty) word.

There are two type of acceptances. The pushdown automaton accepts a word by final state when after reading it the automaton enter a final state. The pushdown automaton accepts a word by empty stack when after reading it the automaton empties its stack. We show that these two acceptances are equivalent.

**Definition 2.21** A *nondeterministic pushdown automaton* is a system

$$V = (Q, \Sigma, W, E, q_0, z_0, F),$$

where

- $Q$  is the finite, non-empty set of states
- $\Sigma$  is the **input alphabet**,
- $W$  is the **stack alphabet**,
- $E \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times W \times W^* \times Q$  is the set of **transitions or edges**,
- $q_0 \in Q$  is the **initial state**,
- $z_0 \in W$  is the **start symbol of stack**,
- $F \subseteq Q$  is the set of **final states**.

A transition  $(p, a, z, w, q)$  means that if pushdown automaton  $V$  is in state  $p$ , reads from the input tape letter  $a$  (instead of input letter we can also consider the empty word  $\varepsilon$ ), and the top symbol in the stack is  $z$ , then the pushdown automaton enters state  $q$  and replaces in the stack  $z$  by word  $w$ . Writing word  $w$  in the stack is made by natural order (letters of word  $w$  will be put in the stack letter by letter from left to right). Instead of writing transition  $(p, a, z, w, q)$  we will use a more suggestive notation  $(p, (a, z/w), q)$ .

Here, as in the case of finite automata, we can define a transition function

$$\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times W \rightarrow \mathcal{P}(W^* \times Q),$$

which associate to current state, input letter and top letter in stack pairs of the form  $(w, q)$ , where  $w \in W^*$  is the word written in stack and  $q \in Q$  the new state.

Because the pushdown automaton is nondeterministic, we will have for the transition function

$\delta(q, a, z) = \{(w_1, p_1), \dots, (w_k, p_k)\}$  (if the pushdown automaton reads an input letter and moves to right), or

$\delta(q, \varepsilon, z) = \{(w_1, p_1), \dots, (w_k, p_k)\}$  (without move on the input tape).

A pushdown automaton is **deterministic**, if for any  $q \in Q$  and  $z \in W$  we have

- $|\delta(q, a, z)| \leq 1, \forall a \in \Sigma \cup \{\varepsilon\}$  and
- if  $\delta(q, \varepsilon, z) \neq \emptyset$ , then  $\delta(q, a, z) = \emptyset, \forall a \in \Sigma$ .

We can associate to any pushdown automaton a transition table, exactly as in the case of finite automata. The rows of this table are indexed by elements of  $Q$ , the columns by elements from  $\Sigma \cup \{\varepsilon\}$  and  $W$  (to each  $a \in \Sigma \cup \{\varepsilon\}$  and  $z \in W$  will correspond a column). At intersection of row corresponding to state  $q \in Q$  and column corresponding to  $a \in \Sigma \cup \{\varepsilon\}$  and  $z \in W$  we will have pairs  $(w_1, p_1), \dots, (w_k, p_k)$  if  $\delta(q, a, z) = \{(w_1, p_1), \dots, (w_k, p_k)\}$ .

The transition graph, in which the label of edge  $(p, q)$  will be  $(a, z/w)$  corresponding to transition  $(p, (a, z/w), q)$ , can be also defined.

**Example 2.25**  $V_1 = (\{q_0, q_1, q_2\}, \{a, b\}, \{z_0, z_1\}, E, q_0, z_0, \{q_0\})$ . Elements of  $E$  are:

$$\begin{array}{ll} (q_0, (a, z_0/z_0z_1), q_1) & \\ (q_1, (a, z_1/z_1z_1), q_1) & (q_1, (b, z_1/\varepsilon), q_2) \\ (q_2, (b, z_1/\varepsilon), q_2) & (q_2, (\varepsilon, z_0/\varepsilon), q_0) . \end{array}$$

The transition function:

$$\begin{array}{ll} \delta(q_0, a, z_0) = \{(z_0z_1, q_1)\} & \\ \delta(q_1, a, z_1) = \{(z_1z_1, q_1)\} & \delta(q_1, b, z_1) = \{(\varepsilon, q_2)\} \\ \delta(q_2, b, z_1) = \{(\varepsilon, q_2)\} & \delta(q_2, \varepsilon, z_0) = \{(\varepsilon, q_0)\} . \end{array}$$

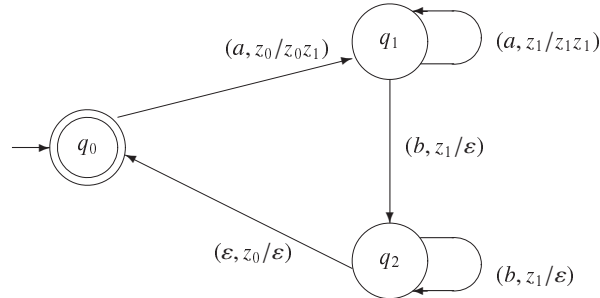


Figure 2.30. Example of pushdown automaton.

The transition table:

$\Sigma \cup \{\varepsilon\}$	$a$		$b$	$\varepsilon$
$W$	$z_0$	$z_1$	$z_1$	$z_0$
$q_0$	$(z_0z_1, q_1)$			
$q_1$		$(z_1z_1, q_1)$	$(\varepsilon, q_2)$	
$q_2$			$(\varepsilon, q_2)$	$(\varepsilon, q_0)$

Because for the transition function every set which is not empty contains only one element (e.g.  $\delta(q_0, a, z_0) = \{(z_0z_1, q_1)\}$ ), in the above table each cell contains only one element, and the set notation is not used. Generally, if a set has more than one element, then its elements are written one under other. The transition graph of this pushdown automaton is in Fig. 2.30.

The current state, the unread part of the input word and the content of stack constitutes a **configuration** of the pushdown automaton, i.e. for each  $q \in Q, u \in \Sigma^*$  and  $v \in W^*$  the triplet  $(q, u, v)$  can be a configuration.

If  $u = a_1a_2 \dots a_k$  and  $v = x_1x_2 \dots x_m$ , then the pushdown automaton can change its configuration in two ways:

- $(q, a_1a_2 \dots a_k, x_1x_2 \dots x_{m-1}x_m) \implies (p, a_2a_3 \dots a_k, x_1, x_2 \dots x_{m-1}w)$ ,  
if  $(q, (a_1, x_m/w), p) \in E$
- $(q, a_1a_2 \dots a_k, x_1x_2 \dots x_m) \implies (p, a_1a_2 \dots a_k, x_1, x_2 \dots x_{m-1}w)$ ,  
if  $(q, (\varepsilon, x_m/w), p) \in E$ .

The reflexive and transitive closure of the relation  $\implies$  will be denoted by  $\implies^*$ . Instead of using  $\implies$ , sometimes  $\vdash$  is considered.

How does work such a pushdown automaton? Getting started with the initial configuration  $(q_0, a_1a_2 \dots a_n, z_0)$  we will consider all possible next configurations, and after this the next configurations to these next configurations, and so on, until it is possible.



**Definition 2.22** Pushdown automaton  $V$  accepts (recognizes) word  $u$  **by final state** if there exist a sequence of configurations of  $V$  for which the following are true:

- the first element of the sequence is  $(q_0, u, z_0)$ ,
- there is a going from each element of the sequence to the next element, excepting the case when the sequence has only one element,
- the last element of the sequence is  $(p, \varepsilon, w)$ , where  $p \in F$  and  $w \in W^*$ .

Therefore pushdown automaton  $V$  accepts word  $u$  by final state, if and only if  $(q_0, u, z_0) \xRightarrow{*} (p, \varepsilon, w)$  for some  $w \in W^*$  and  $p \in F$ . The set of words accepted by final state by pushdown automaton  $V$  will be called the language accepted by  $V$  by final state and will be denoted by  $L(V)$ .

**Definition 2.23** Pushdown automaton  $V$  accepts (recognizes) word  $u$  **by empty stack** if there exist a sequence of configurations of  $V$  for which the following are true:

- the first element of the sequence is  $(q_0, u, z_0)$ ,
- there is a going from each element of the sequence to the next element,
- the last element of the sequence is  $(p, \varepsilon, \varepsilon)$  and  $p$  is an arbitrary state.

Therefore pushdown automaton  $V$  accepts a word  $u$  by empty stack if  $(q_0, u, z_0) \xRightarrow{*} (p, \varepsilon, \varepsilon)$  for some  $p \in Q$ . The set of words accepted by empty stack by pushdown automaton  $V$  will be called the language accepted by empty stack by  $V$  and will be denoted by  $L_\varepsilon(V)$ .

**Example 2.26** Pushdown automaton  $V_1$  of Example 2.25. accepts the language  $\{a^n b^n \mid n \geq 0\}$  by final state. Consider the derivation for words  $aaabbb$  and  $abab$ .

Word  $a^3 b^3$  is accepted by the considered pushdown automaton because

$$(q_0, aaabbb, z_0) \Rightarrow (q_1, aabbb, z_0 z_1) \Rightarrow (q_1, abbb, z_0 z_1 z_1) \Rightarrow (q_1, bbb, z_0 z_1 z_1 z_1)$$

$\Rightarrow (q_2, bb, z_0 z_1 z_1) \Rightarrow (q_2, b, z_0 z_1) \Rightarrow (q_2, \varepsilon, z_0) \Rightarrow (q_0, \varepsilon, \varepsilon)$  and because  $q_0$  is a final state the pushdown automaton accepts this word. But the stack being empty, it accepts this word also by empty stack.

Because the initial state is also a final state, the empty word is accepted by final state, but not by empty stack.

To show that word  $abab$  is not accepted, we need to study all possibilities. It is easy to see that in our case there is only a single possibility:

$$(q_0, abab, z_0) \Rightarrow (q_1, bab, z_0 z_1) \Rightarrow (q_2, ab, z_0) \Rightarrow (q_0, ab, \varepsilon),$$

but there is no further going, so word  $abab$  is not accepted.

**Example 2.27** The transition table of the pushdown automaton  $V_2 = (\{q_0, q_1\}, \{0, 1\}, \{z_0, z_1, z_2\}, E, q_0, z_0, \emptyset)$  is:

$\Sigma \cup \{\varepsilon\}$	0			1			$\varepsilon$
$W$	$z_0$	$z_1$	$z_2$	$z_0$	$z_1$	$z_2$	$z_0$
$q_0$	$(z_0 z_1, q_0)$	$(z_1 z_1, q_0)$ $(\varepsilon, q_1)$	$(z_2 z_1, q_0)$	$(z_0 z_2, q_0)$	$(z_1 z_2, q_0)$	$(z_2 z_2, q_0)$ $(\varepsilon, q_1)$	$(\varepsilon, q_1)$
$q_1$		$(\varepsilon, q_1)$				$(\varepsilon, q_1)$	$(\varepsilon, q_1)$

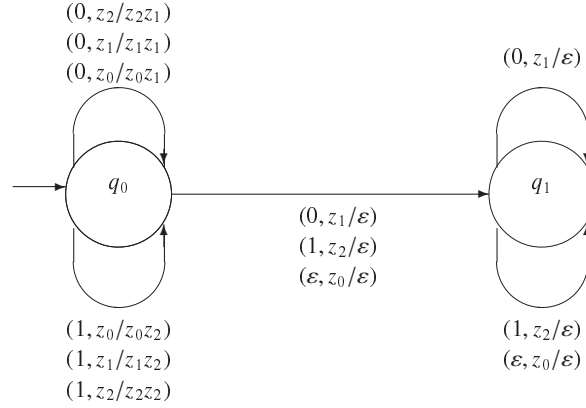


Figure 2.31. Transition graph of the Example 2.27..

The corresponding transition graph can be seen in Fig. 2.31. Pushdown automaton  $V_2$  accepts the language  $\{uu^{-1} \mid u \in \{0, 1\}^*\}$ . Because  $V_2$  is nondeterministic, all the configurations obtained from the initial configuration  $(q_0, u, z_0)$  can be illustrated by a **computation tree**. For example the computation tree associated to the initial configuration  $(q_0, 1001, z_0)$  can be seen in Fig. 2.32. From this computation tree we can observe that, because  $(q_1, \varepsilon, \varepsilon)$  is a leaf of the tree, pushdown automaton  $V_2$  accepts word 1001 by empty stack. The computation tree in Fig. 2.33 shows that pushdown automaton  $V_2$  does not accept word 101, because the configurations in leaves can not be continued and none of them has the form  $(q, \varepsilon, \varepsilon)$ .

**Theorem 2.24** *A language  $L$  is accepted by a nondeterministic pushdown automaton  $V_1$  by empty stack if and only if it can be accepted by a nondeterministic pushdown automaton  $V_2$  by final state.*

**Proof.** a) Let  $V_1 = (Q, \Sigma, W, E, q_0, z_0, \emptyset)$  be the pushdown automaton which accepts by empty stack language  $L$ . Define pushdown automaton  $V_2 = (Q \cup \{p_0, p\}, \Sigma, W \cup \{x\}, E', p_0, x, \{p\})$ , where  $p, p_0 \notin Q$ ,  $x \notin W$  and

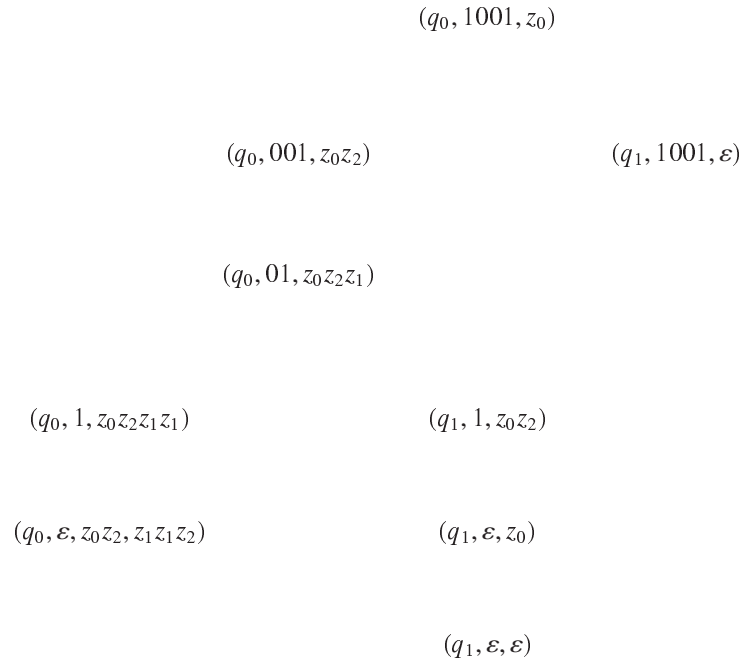
$$E' = E \cup \{(p_0, (\varepsilon, x/xz_0), q_0)\} \cup \{(q, (\varepsilon, x/\varepsilon), p) \mid q \in Q\}$$

Working of  $V_2$ : Pushdown automaton  $V_2$  with an  $\varepsilon$ -move first goes in the initial state of  $V_1$ , writing  $z_0$  (the initial stack symbol of  $V_1$ ) in the stack (beside  $x$ ). After this it is working as  $V_1$ . If  $V_1$  for a given word empties its stack, then  $V_2$  still has  $x$  in the stack, which can be deleted by  $V_2$  using an  $\varepsilon$ -move, while a final state will be reached.  $V_2$  can reach a final state only if  $V_1$  has emptied the stack.

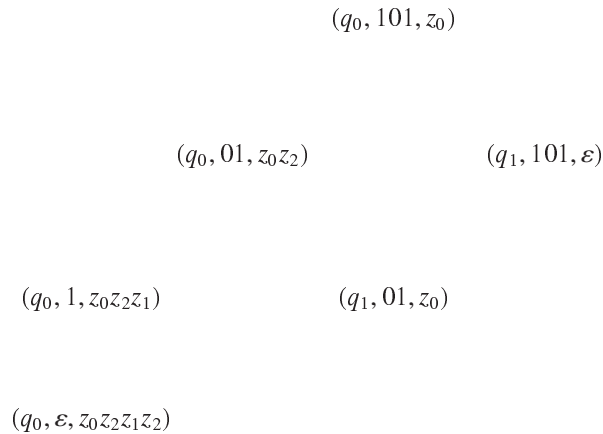
b) Let  $V_2 = (Q, \Sigma, W, E, q_0, z_0, F)$  be a pushdown automaton, which accepts language  $L$  by final state. Define pushdown automaton  $V_1 = (Q \cup \{p_0, p\}, \Sigma, W \cup \{x\}, E', p_0, x, \emptyset)$ , where  $p_0, p \notin Q$ ,  $x \notin W$  and

$$E' = E \cup \{(p_0, (\varepsilon, x/xz_0), q_0)\} \cup \{(q, (\varepsilon, z/\varepsilon), p) \mid q \in F, p \in Q, z \in W\} \\ \cup \{(p, (\varepsilon, z/\varepsilon), p) \mid p \in Q, z \in W \cup \{x\}\}$$

Working  $V_1$ : Pushdown automaton  $V_1$  with an  $\varepsilon$ -move writes in the stack beside  $x$  the initial



**Figure 2.32.** Computation tree to show acceptance of the word 1001 (see Example 2.27.).



**Figure 2.33.** Computation tree to show that the pushdown automaton in Example 2.27. does not accept word 101.

stack symbol  $z_0$  of  $V_2$ , then works as  $V_2$ , i.e reaches a final state for each accepted word. After this  $V_1$  empties the stack by an  $\varepsilon$ -move.  $V_1$  can empty the stack only if  $V_2$  goes in a final state. ■

The next two theorems prove that the class of languages accepted by nondeterministic pushdown automata is just the set of context-free languages.

**Theorem 2.25** *If  $G$  is a context-free grammar, then there exists such a nondeterministic pushdown automaton  $V$  which accepts  $L(G)$  by empty stack, i.e.  $L_\varepsilon(V) = L(G)$ .*

We outline the proof only. Let  $G = (N, T, P, S)$  be a context-free grammar. Define pushdown automaton  $V = (\{q\}, T, N \cup T, E, q, S, \emptyset)$ , where  $q \notin N \cup T$ , and the set  $E$  of transitions is:

- If there is in the set of productions of  $G$  a production of type  $A \rightarrow \alpha$ , then let put in  $E$  the transition  $(q, (\varepsilon, A/\alpha^{-1}), q)$ ,
- For any letter  $a \in T$  let put in  $E$  the transition  $(q, (a, a/\varepsilon), q)$ .

If there is a production  $S \rightarrow \alpha$  in  $G$ , the pushdown automaton put in the stack the mirror of  $\alpha$  with an  $\varepsilon$ -move. If the input letter coincides with that in the top of the stack, then the automaton deletes it from the stack. If in the top of the stack there is a nonterminal  $A$ , then the mirror of right-hand side of a production which has  $A$  in its left-hand side will be put in the stack. If after reading all letters of the input word, the stack will be empty, then the pushdown automaton recognized the input word.

The following algorithm builds for a context-free grammar  $G = (N, T, P, S)$  the pushdown automaton  $V = (\{q\}, T, N \cup T, E, q, S, \emptyset)$ , which accepts by empty stack the language generated by  $G$ .

FROM-CFG-TO-PUSHDOWN-AUTOMATON( $G, V$ )

- 1 **for** all production  $A \rightarrow \alpha$
- 2     **do** put in  $E$  the transition  $(q, (\varepsilon, A/\alpha^{-1}), q)$
- 3 **for** all terminal  $a \in T$
- 4     **do** put in  $E$  the transition  $(q, (a, a/\varepsilon), q)$

If  $G$  has  $n$  productions and  $m$  terminals, then the number of step of the algorithm is  $\Theta(n + m)$ .

**Example 2.28** Let  $G = (\{S, A\}, \{a, b\}, \{S \rightarrow \varepsilon, S \rightarrow ab, S \rightarrow aAb, A \rightarrow aAb, A \rightarrow ab\}, S)$ . Then  $V = (\{q\}, \{a, b\}, \{a, b, A, S\}, E, q, S, \emptyset)$ , with the following transition table.

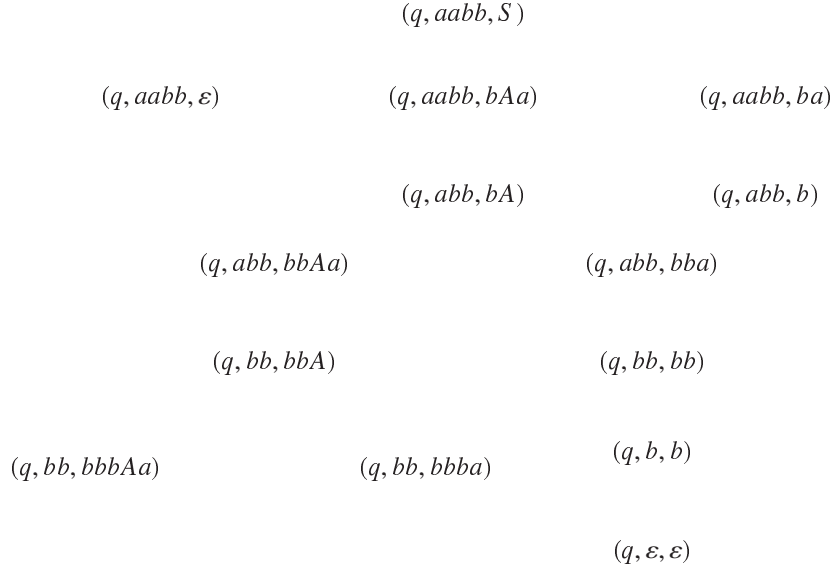
$\Sigma \cup \{\varepsilon\}$	$a$	$b$	$\varepsilon$	
$W$	$a$	$b$	$S$	$A$
$q$	$(\varepsilon, q)$	$(\varepsilon, q)$	$(\varepsilon, q)$ $(ba, q)$ $(bAa, q)$	$(bAa, q)$ $(ba, q)$

Let us see how pushdown automaton  $V$  accepts word  $aabb$ , which in grammar  $G$  can be derived in the following way:

$$S \implies aAb \implies aabb,$$

where productions  $S \rightarrow aAb$  and  $A \rightarrow ab$  were used. Word is accepted by empty stack (see Fig. 2.34).

**Theorem 2.26** *For a nondeterministic pushdown automaton  $V$  there exists always a*



**Figure 2.34.** Recognising a word by empty stack (see Example 2.28.).

context-free grammar  $G$  such that  $\mathbb{V}$  accepts language  $L(G)$  by empty stack, i.e.  $L_\varepsilon(\mathbb{V}) = L(G)$ .

Instead of a proof we will give a method to obtain grammar  $G$ . Let  $\mathbb{V} = (Q, \Sigma, W, E, q_0, z_0, \emptyset)$  be the nondeterministic pushdown automaton in question.

Then  $G = (N, T, P, S)$ , where

$N = \{S\} \cup \{S_{p,z,q} \mid p, q \in Q, z \in W\}$  and  $T = \Sigma$ .

Productions in  $P$  will be obtained as follows.

- For all state  $q$  put in  $P$  production  $S \rightarrow S_{q_0, z_0, q}$ .
- If  $(q, (a, z/z_k \dots z_2 z_1), p) \in E$ , where  $q \in Q$ ,  $z, z_1, z_2, \dots, z_k \in W$  ( $k \geq 1$ ) and  $a \in \Sigma \cup \{\varepsilon\}$ , put in  $P$  for all possible states  $p_1, p_2, \dots, p_k$  productions

$$S_{q,z,p_k} \rightarrow aS_{p,z_1,p_1} S_{p_1,z_2,p_2} \dots S_{p_{k-1},z_k,p_k}.$$

- If  $(q, (a, z/\varepsilon), p) \in E$ , where  $p, q \in Q, z \in W$ , and  $a \in \Sigma \cup \{\varepsilon\}$ , put in  $P$  production

$$S_{q,z,p} \rightarrow a.$$

The context-free grammar defined by this is an extended one, to which an equivalent context-free language can be associated. The proof of the theorem is based on the fact that to every sequence of configurations, by which the pushdown automaton  $\mathbb{V}$  accepts a word, we can associate a derivation in grammar  $G$ . This derivation generates just the word in question, because of productions of the form  $S_{q,z,p_k} \rightarrow aS_{p,z_1,p_1} S_{p_1,z_2,p_2} \dots S_{p_{k-1},z_k,p_k}$ , which were defined for all possible states  $p_1, p_2, \dots, p_k$ . In Example 2.27. we show how can be associated a derivation to a sequence of configurations. The pushdown automaton defined in the example recognizes word 00 by the sequence of configurations

$$(q_0, 00, z_0) \Rightarrow (q_0, 0, z_0 z_1) \Rightarrow (q_1, \varepsilon, z_0) \Rightarrow (q_1, \varepsilon, \varepsilon),$$

which sequence is based on the transitions

$$(q_0, (0, z_0/z_0 z_1), q_0),$$

$$(q_0, (0, z_1/\varepsilon), q_1),$$

$$(q_1, (\varepsilon, z_1/\varepsilon), q_1).$$

To these transitions, by the definition of grammar  $G$ , the following productions can be associated

- (1)  $S_{q_0, z_0, p_2} \rightarrow 0S_{q_0, z_1, p_1} S_{p_1, z_0, p_2}$  for all states  $p_1, p_2 \in Q$ ,
- (2)  $S_{q_0, z_1, q_1} \rightarrow 0$ ,
- (3)  $S_{q_1, z_0, q_1} \rightarrow \varepsilon$ .

Furthermore, for each state  $q$  productions  $S \rightarrow S_{q_0, z_0, q}$  were defined.

By the existence of production  $S \rightarrow S_{q_0, z_0, q}$  there exists the derivation  $S \Rightarrow S_{q_0, z_0, q}$ , where  $q$  can be chosen arbitrarily. Let choose in above production (1) state  $q$  to be equal to  $p_2$ . Then there exists also the derivation

$$S \Rightarrow S_{q_0, z_0, q} \Rightarrow 0S_{q_0, z_1, p_1} S_{p_1, z_0, q},$$

where  $p_1 \in Q$  can be chosen arbitrarily. If  $p_1 = q_1$ , then the derivation

$$S \Rightarrow S_{q_0, z_0, q} \Rightarrow 0S_{q_0, z_1, q_1} S_{q_1, z_0, q} \Rightarrow 00S_{q_1, z_0, q}$$

will result. Now let  $q$  equal to  $q_1$ , then

$$S \Rightarrow S_{q_0, z_0, q_1} \Rightarrow 0S_{q_0, z_1, q_1} S_{q_1, z_0, q_1} \Rightarrow 00S_{q_1, z_0, q_1} \Rightarrow 00,$$

which proves that word  $00$  can be derived used the above grammar.

The next algorithm builds for a pushdown automaton  $V = (Q, \Sigma, W, E, q_0, z_0, \emptyset)$  a context-free grammar  $G = (N, T, P, S)$ , which generates the language accepted by pushdown automaton  $V$  by empty stack.

FROM-PUSHDOWN-AUTOMATON-TO-CF-GRAMMAR( $V, G$ )

- 1 **for** all  $q \in Q$
- 2     **do** put in  $P$  production  $S \rightarrow S_{q_0, z_0, q}$
- 3 **for** all  $(q, (a, z/z_k \dots z_2 z_1), p) \in E$       $\triangleright q \in Q, z, z_1, z_2, \dots, z_k \in W (k \geq 1), a \in \Sigma \cup \{\varepsilon\}$
- 4     **do for** all states  $p_1, p_2, \dots, p_k$
- 5         **do** put in  $P$  productions  $S_{q, z, p_k} \rightarrow aS_{p, z_1, p_1} S_{p_1, z_2, p_2} \dots S_{p_{k-1}, z_k, p_k}$
- 6 **for** All  $(q(a, z/\varepsilon), p) \in E$       $\triangleright p, q \in Q, z \in W, a \in \Sigma \cup \{\varepsilon\}$
- 7     **do** put in  $P$  production  $S_{q, z, p} \rightarrow a$

If the automaton has  $n$  states and  $m$  productions, then the above algorithm executes at most  $n + mn + m$  steps, so in worst case the number of steps is  $O(nm)$ .

Finally, without proof, we mention that the class of languages accepted by deterministic pushdown automata is a proper subset of the class of languages accepted by nondeterministic pushdown automata. This points to the fact that pushdown automata behave differently as finite automata.

**Example 2.29** As an example, consider pushdown automaton  $V$  from the Example 2.28.:  $V = (\{q\}, \{a, b\}, \{a, b, A, S\}, E, q, S, \emptyset)$ . Grammar  $G$  is:

$$G = (\{S, S_a, S_b, S_s, S_A\}, \{a, b\}, P, S),$$

where for all  $z \in \{a, b, S, A\}$  instead of  $S_{q, z, q}$  we shortly used  $S_z$ . The transitions:

$$\begin{aligned} &(q, (a, a/\varepsilon), q), & &(q, (b, b/\varepsilon), q), \\ &(q, (\varepsilon, S/\varepsilon), q), & &(q, (\varepsilon, S/ba), q), & &(q, (\varepsilon, S/bAa), q), \\ &(q, (\varepsilon, A/ba), q), & &(q, (\varepsilon, A/bAa), q). \end{aligned}$$

Based on these, the following productions are defined:

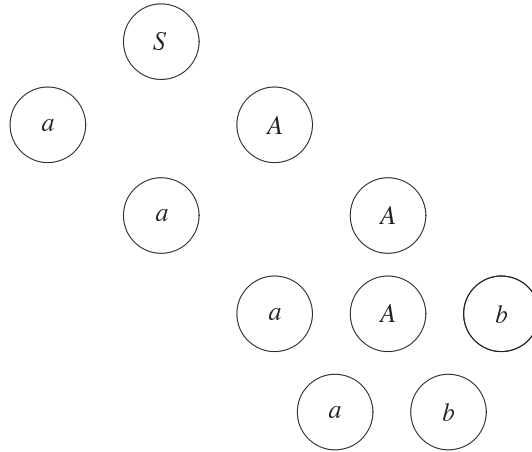


Figure 2.35. Derivation (or syntax) tree of word  $aaaabb$ .

$$\begin{aligned}
 S &\rightarrow S_S \\
 S_a &\rightarrow a \\
 S_b &\rightarrow b \\
 S_S &\rightarrow \varepsilon \mid S_a S_b \mid S_a S_A S_b \\
 S_A &\rightarrow S_a S_A S_b \mid S_a S_b.
 \end{aligned}$$

It is easy to see that  $S_S$  can be eliminated, and the productions will be:

$$\begin{aligned}
 S &\rightarrow \varepsilon \mid S_a S_b \mid S_a S_A S_b, \\
 S_A &\rightarrow S_a S_A S_b \mid S_a S_b, \\
 S_a &\rightarrow a, \quad S_b \rightarrow b,
 \end{aligned}$$

and these productions can be replaced:

$$\begin{aligned}
 S &\rightarrow \varepsilon \mid ab \mid aAb, \\
 A &\rightarrow aAb \mid ab.
 \end{aligned}$$

### 2.3.2. Context-free languages

Consider context-free grammar  $G = (N, T, P, S)$ . A **derivation tree** of  $G$  is a finite, ordered, labelled tree, which root is labelled by the the start symbol  $S$ , every interior vertex is labelled by a nonterminal and every leaf by a terminal. If an interior vertex labelled by a nonterminal  $A$  has  $k$  descendents, then in  $P$  there exists a production  $A \rightarrow a_1 a_2 \dots a_k$  such that the descendents are labelled by letters  $a_1, a_2, \dots, a_k$ . The **result** of a derivation tree is a word over  $T$ , which can be obtained by reading the labels of the leaves from left to right. Derivation tree is also called **syntax tree**.

Consider the context-free grammar  $G = (\{S, A\}, \{a, b\}, \{S \rightarrow aA, S \rightarrow a, S \rightarrow \varepsilon, A \rightarrow aA, A \rightarrow ab, A \rightarrow b\}, S)$ . It generates language  $L(G) = \{a^n b^m \mid n \geq m \geq 0\}$ . Derivation of word  $a^4 b^2 \in L(G)$  is:

$$S \Rightarrow aA \Rightarrow aaA \Rightarrow aaaAb \Rightarrow aaaabb.$$

In Fig. 2.35 this derivation can be seen, which result is  $aaaabb$ .

To every derivation we can associate a syntax tree. Conversely, to any syntax tree more

than one derivation can be associated. For example to syntax tree in Fig. 2.35 the derivation

$$S \Rightarrow aA \Rightarrow aaAb \Rightarrow aaaAb \Rightarrow aaaabb$$

also can be associated.

**Definition 2.27** Derivation  $\alpha_0 \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_n$  is a **leftmost derivation**, if for all  $i = 1, 2, \dots, n - 1$  there exist words  $u_i \in T^*$ ,  $\beta_i \in (N \cup T)^*$  and productions  $(A_i \rightarrow \gamma_i) \in P$ , for which we have  $\alpha_i = u_i A_i \beta_i$  and  $\alpha_{i+1} = u_i \gamma_i \beta_i$ .

Consider grammar:

$$G = (\{S, A\}, \{a, b, c\}, \{S \rightarrow bA, S \rightarrow bAS, S \rightarrow a, A \rightarrow cS, A \rightarrow a\}, S).$$

In this grammar word  $bcbaa$  has two different leftmost derivations:

$$S \Rightarrow bA \Rightarrow bcS \Rightarrow bcbAS \Rightarrow bcbaS \Rightarrow bcbaa,$$

$$S \Rightarrow bAS \Rightarrow bcSS \Rightarrow bcbAS \Rightarrow bcbaS \Rightarrow bcbaa.$$

**Definition 2.28** A context-free grammar  $G$  is **ambiguous** if in  $L(G)$  there exists a word with more than one leftmost derivation. Otherwise  $G$  is **unambiguous**.

The above grammar  $G$  is ambiguous, because word  $bcbaa$  has two different leftmost derivations. A language can be generated by more than one grammar, and between them can exist ambiguous and unambiguous too. A context-free language is **inherently ambiguous**, if there is no unambiguous grammar which generates it.

**Example 2.30** Examine the following two grammars.

Grammar  $G_1 = (\{S\}, \{a, +, *\}, \{S \rightarrow S + S, S \rightarrow S * S, S \rightarrow a\}, S)$  is ambiguous because

$$S \Rightarrow S + S \Rightarrow a + S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * S + S \Rightarrow a + a * a + S$$

$$\Rightarrow a + a * a + a \quad \text{and}$$

$$S \Rightarrow S * S \Rightarrow S + S * S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * S + S \Rightarrow a + a * a + S$$

$$\Rightarrow a + a * a + a.$$

Grammar  $G_2 = (\{S, A\}, \{a, *, +\}, \{S \rightarrow A + S \mid A, A \rightarrow A * A \mid a\}, S)$  is unambiguous.

Can be proved that  $L(G_1) = L(G_2)$ .

### 2.3.3. Pumping lemma for context-free languages

Like for regular languages there exists a pumping lemma also for context-free languages.

**Theorem 2.29** (pumping lemma). For any context-free language  $L$  there exists a natural number  $n$  (which depends only on  $L$ ), such that every word  $z$  of the language longer than  $n$  can be written in the form  $uvwxy$  and the following are true:

- (1)  $|w| \geq 1$ ,
- (2)  $|vx| \geq 1$ ,
- (3)  $|vwx| \leq n$ ,
- (4)  $uv^iwx^iy$  is also in  $L$  for all  $i \geq 0$ .

**Proof.** Let  $G = (N, T, P, S)$  be a grammar without unit productions, which generates language  $L$ . Let  $m = |N|$  be the number of nonterminals, and let  $\ell$  be the maximum of lengths of right-hand sides of productions, i.e.  $\ell = \max\{|\alpha| \mid \exists A \in N : (A \rightarrow \alpha) \in P\}$ . Let  $n = \ell^{m+1}$  and



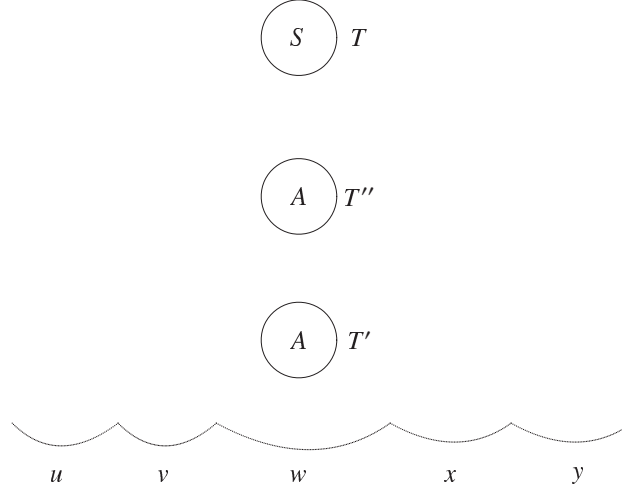


Figure 2.36. Decomposition of tree in the proof of pumping lemma.

$z \in L(G)$ , such that  $|z| > n$ . Then there exists a derivation tree  $T$  with the result  $z$ . Let  $h$  be the height of  $T$  (the maximum of path lengths from root to leaves). Because in  $T$  all interior vertices have at most  $\ell$  descendants,  $T$  has at most  $\ell^h$  leaves, i.e.  $|z| \leq \ell^h$ . On the other hand, because of  $|z| > \ell^{m+1}$ , we get that  $h > m + 1$ . From this follows that in derivation tree  $T$  there is a path from root to a leaf in which there are more than  $(m + 1)$  vertices. Consider such a path. Because in  $G$  the number of nonterminals is  $m$  and on this path vertices different from the leaf are labelled with nonterminals, by the pigeonhole principle, it must be a nonterminal on this path which occurs at least twice.

Let us denote by  $A$  the nonterminal being the first on this path from root to the leaf which firstly repeat. Denote by  $T'$  the subtree, which root is this occurrence of  $A$ . Similarly, denote by  $T''$  the subtree, which root is the second occurrence of  $A$  on this path. Let  $w$  be the result of the tree  $T'$ . Then the result of  $T''$  is in form  $vwx$ , while of  $T$  in  $uvwxy$ . Derivation tree  $T$  with this decomposition of  $z$  can be seen in Fig. 2.36. We show that this decomposition of  $z$  satisfies conditions (1)–(4) of lemma.

Because in  $P$  there are no  $\varepsilon$ -productions (except maybe the case  $S \rightarrow \varepsilon$ ), we have  $|w| \geq 1$ . Furthermore, because each interior vertex of the derivation tree has at least two descendants (namely there are no unit productions), also the root of  $T''$  has, hence  $|vx| \geq 1$ . Because  $A$  is the first repeated nonterminal on this path, the height of  $T''$  is at most  $m + 1$ , and from this  $|vwx| \leq \ell^{m+1} = n$  results.

After eliminating from  $T$  all vertices of  $T''$  excepting the root, the result of obtained tree is  $uAy$ , i.e.  $S \xRightarrow{*}_G uAy$ .

Similarly, after eliminating  $T'$  we get  $A \xRightarrow{*}_G vAx$ , and finally because of the definition of  $T'$  we get  $A \xRightarrow{*}_G w$ . Then  $S \xRightarrow{*}_G uAy$ ,  $A \xRightarrow{*}_G vAx$  and  $A \xRightarrow{*}_G w$ . Therefore  $S \xRightarrow{*}_G uAy \xRightarrow{*}_G uwy$  and  $S \xRightarrow{*}_G uAy \xRightarrow{*}_G uvAxy \xRightarrow{*}_G \dots \xRightarrow{*}_G uv^iAx^i y \xRightarrow{*}_G uv^iwx^i y$  for all

$i \geq 1$ . Therefore, for all  $i \geq 0$  we have  $S \xRightarrow{*} uv^iwx^iy$ , i.e. for all  $i \geq 0$   $uv^iwx^iy \in L(G)$ . ■

Now we present two consequences of the lemma.

**Corollary 2.30**  $\mathcal{L}_2 \subset \mathcal{L}_1$ .

**Proof.** This consequence states that there exists a context-sensitive language which is not context-free. To prove this it is sufficient to find a context-sensitive language for which the lemma is not true. Let this language be  $L = \{a^m b^m c^m \mid m \geq 1\}$ .

To show that this language is context-sensitive it is enough to give a convenient grammar. In Example 2.2. both grammars are extended context-sensitive, and we know that to each extended grammar of type  $i$  an equivalent grammar of the same type can be associated.

Let  $n$  be the natural number associated to  $L$  by lemma, and consider the word  $z = a^n b^n c^n$ . Because of  $|z| = 3n > n$ , if  $L$  is context-free  $z$  can be decomposed in  $z = uvwxy$  such that conditions (1)–(4) are true. We show that this leads us to a contradiction.

Firstly, we will show that word  $v$  and  $x$  can contain only one type of letters. Indeed if either  $v$  or  $x$  contain more than one type of letters, then in word  $uvvwxy$  the order of the letters will be not the order  $a, b, c$ , so  $uvvwxy \notin L(G)$ , which contradicts condition (4) of lemma.

If both  $v$  and  $x$  contain at most one type of letters, then in word  $uwy$  the number of different letters will be not the same, so  $uwy \notin L(G)$ . This also contradicts condition (4) in lemma. Therefore  $L$  is not context-free. ■

**Corollary 2.31** *The class of context-free languages is not closed under the intersection.*

**Proof.** We give two context-free languages which intersection is not context-free. Let  $N = \{S, A, B\}$ ,  $T = \{a, b, c\}$  and

$G_1 = (N, T, P_1, S)$  where  $P_1$  :

$S \rightarrow AB,$   
 $A \rightarrow aAb \mid ab,$   
 $B \rightarrow cB \mid c,$

and  $G_2 = (N, T, P_2, S)$ , where  $P_2$  :

$S \rightarrow AB,$   
 $A \rightarrow Aa \mid a,$   
 $B \rightarrow bBc \mid bc.$

Languages  $L(G_1) = \{a^n b^n c^m \mid n \geq 1, m \geq 1\}$  and  $L(G_2) = \{a^n b^m c^n \mid n \geq 1, m \geq 1\}$  are context-free. But

$$L(G_1) \cap L(G_2) = \{a^n b^n c^n \mid n \geq 1\}$$

is not context-free (see the proof of the Consequence 2.30). ■

### 2.3.4. Normal forms of the context-free languages

In the case of arbitrary grammars the normal form was defined (see page 83) as grammars with no terminals in the left-hand side of productions. The normal form in the case of the context-free languages will contains some restrictions on the right-hand sides of productions. Two normal forms (Chomsky and Greibach) will be discussed.

### Chomsky normal form

**Definition 2.32** A context-free grammar  $G = (N, T, P, S)$  is in Chomsky normal form, if all productions have form  $A \rightarrow a$  or  $A \rightarrow BC$ , where  $A, B, C \in N$ ,  $a \in T$ .

**Example 2.31** Grammar  $G = (\{S, A, B, C\}, \{a, b\}, \{S \rightarrow AB, S \rightarrow CB, C \rightarrow AS, A \rightarrow a, B \rightarrow b\}, S)$  is in Chomsky normal form and  $L(G) = \{a^n b^n \mid n \geq 1\}$ .

To each  $\varepsilon$ -free context-free language can be associated an equivalent grammar in Chomsky normal form. The next algorithm transforms an  $\varepsilon$ -free context-free grammar  $G = (N, T, P, S)$  in grammar  $G' = (N', T, P', S)$  which is in Chomsky normal form.

#### CHOMSKY-NORMAL-FORM( $G, G'$ )

- 1  $N' \leftarrow N$
- 2 eliminate unit productions, and let  $P'$  the new set of productions (see algorithm ELIMINATE-UNIT-PRODUCTIONS on page 82)
- 3 in  $P'$  replace in each production with at least two letters in right-hand side all terminals  $a$  by a new nonterminal  $A$ , and add this nonterminal to  $N'$  and add production  $A \rightarrow a$  to  $P'$
- 4 replace all productions  $B \rightarrow A_1 A_2 \dots A_k$ , where  $k \geq 3$  and  $A_1, A_2, \dots, A_k \in N$ , by the following:

$$\begin{aligned} B &\rightarrow A_1 C_1, \\ C_1 &\rightarrow A_2 C_2, \\ &\dots \\ C_{k-3} &\rightarrow A_{k-2} C_{k-2}, \\ C_{k-2} &\rightarrow A_{k-1} A_k, \end{aligned}$$

where  $C_1, C_2, \dots, C_{k-2}$  are new nonterminals, and add them to  $N'$ .

**Example 2.32** Let  $G = (\{S, D\}, \{a, b, c\}, \{S \rightarrow aSc, S \rightarrow D, D \rightarrow bD, D \rightarrow b\}, S)$ . It is easy to see that  $L(G) = \{a^n b^m c^n \mid n \geq 0, m \geq 1\}$ . Steps of transformation to Chomsky normal form are the following:

Step 1:  $N' = \{S, D\}$

Step 2: After eliminating the unit production  $S \rightarrow D$  the productions are:

$$\begin{aligned} S &\rightarrow aSc \mid bD \mid b, \\ D &\rightarrow bD \mid b. \end{aligned}$$

Step 3: We introduce three new nonterminals because of the three terminals in productions. Let these be  $A, B, C$ . Then the production are:

$$\begin{aligned} S &\rightarrow ASC \mid BD \mid b, \\ D &\rightarrow BD \mid b, \\ A &\rightarrow a, \\ B &\rightarrow b, \\ C &\rightarrow c. \end{aligned}$$

Step 4: Only one new nonterminal (let this  $E$ ) must be introduced because of a single production with three letters in the right-hand side. Therefore  $N' = \{S, A, B, C, D, E\}$ , and the productions in  $P'$  are:

$$\begin{aligned} S &\rightarrow AE \mid BD \mid b, \\ D &\rightarrow BD \mid b, \end{aligned}$$

$$\begin{aligned} A &\rightarrow a, \\ B &\rightarrow b, \\ C &\rightarrow c, \\ E &\rightarrow SC. \end{aligned}$$

All these productions are in required form.

### Greibach normal form

**Definition 2.33** A context-free grammar  $G = (N, T, P, S)$  is in **Greibach normal form** if all production are in the form  $A \rightarrow aw$ , where  $A \in N$ ,  $a \in T$ ,  $w \in N^*$ .

**Example 2.33** Grammar  $G = (\{S, B\}, \{a, b\}, \{S \rightarrow aB, S \rightarrow aSB, B \rightarrow b\}, S)$  is in Greibach normal form and  $L(G) = \{a^n b^n \mid n \geq 1\}$ .

To each  $\varepsilon$ -free context-free grammar an equivalent grammar in Greibach normal form can be given. We give an algorithm which transforms a context-free grammar  $G = (N, T, P, S)$  in Chomsky normal form in a grammar  $G' = (N', T, P', S)$  in Greibach normal form.

First, we give an order of the nonterminals:  $A_1, A_2, \dots, A_n$ , where  $A_1$  is the start symbol. The algorithm will use the notations  $x \in N'^+$ ,  $\alpha \in TN'^* \cup N'^+$ .

GREIBACH-NORMAL-FORM( $G, G'$ )

```

1   $N' \leftarrow N$ 
2   $P' \leftarrow P$ 
3  for  $i \leftarrow 2$  to  $n$                                 ▷ Case  $A_i \rightarrow A_j x$ ,  $j < i$ 
4      do for  $j \leftarrow 1$  to  $i - 1$ 
5          do for all productions  $A_i \rightarrow A_j x$  and  $A_j \rightarrow \alpha$  (where  $\alpha$  has no  $A_j$  as first letter)
              in  $P'$  productions  $A_i \rightarrow \alpha x$ ,
              delete from  $P'$  productions  $A_i \rightarrow A_j x$ 
6      if there is a production  $A_i \rightarrow A_i x$                                 ▷ Case  $A_i \rightarrow A_i x$ 
7          then put in  $N'$  the new nonterminal  $B_i$ ,
              for all productions  $A_i \rightarrow A_i x$  put in  $P'$  productions  $B_i \rightarrow xB_i$  and  $B_i \rightarrow x$ ,
              delete from  $P'$  production  $A_i \rightarrow A_i x$ ,
              for all production  $A_i \rightarrow \alpha$  (where  $A_i$  is not the first letter of  $\alpha$ )
              put in  $P'$  production  $A_i \rightarrow \alpha B_i$ 
8      for  $i \leftarrow n - 1$  downto 1                                ▷ Case  $A_i \rightarrow A_j x$ ,  $j > i$ 
9          do for  $j \leftarrow i + 1$  to  $n$ 
10             do for all productions  $A_i \rightarrow A_j x$  and  $A_j \rightarrow \alpha$ 
                    put in  $P'$  production  $A_i \rightarrow \alpha x$  and
                    delete from  $P'$  productions  $A_i \rightarrow A_j x$ ,
11      for  $i \leftarrow 1$  to  $n$                                 ▷ Case  $B_i \rightarrow A_j x$ 
12          do for  $j \leftarrow 1$  to  $n$ 
13             do for all productions  $B_i \rightarrow A_j x$  and  $A_j \rightarrow \alpha$ 
                    put in  $P'$  production  $B_i \rightarrow \alpha x$  and
                    delete from  $P'$  productions  $B_i \rightarrow A_j x$ 

```

The algorithm first transform productions of the form  $A_i \rightarrow A_jx$ ,  $j < i$  such that  $A_i \rightarrow A_jx$ ,  $j \geq i$  or  $A_i \rightarrow \alpha$ , where this latter is in Greibach normal form. After this, introducing a new nonterminal, eliminate productions  $A_i \rightarrow A_ix$ , and using substitutions all production of the form  $A_i \rightarrow A_jx$ ,  $j > i$  and  $B_i \rightarrow A_jx$  will be transformed in Greibach normal form.

**Example 2.34** Transform productions in Chomsky normal form

$$\begin{aligned} A_1 &\rightarrow A_2A_3 \mid A_2A_4 \\ A_2 &\rightarrow A_2A_3 \mid a \\ A_3 &\rightarrow A_2A_4 \mid b \\ A_4 &\rightarrow c \end{aligned}$$

in Greibach normal form.

Steps of the algorithm:

3-5: Production  $A_3 \rightarrow A_2A_4$  must be transformed. For this production  $A_2 \rightarrow a$  is appropriate. Put  $A_3 \rightarrow aA_4$  in the set of productions and eliminate  $A_3 \rightarrow A_2A_4$ .

The productions will be:

$$\begin{aligned} A_1 &\rightarrow A_2A_3 \mid A_2A_4 \\ A_2 &\rightarrow A_2A_3 \mid a \\ A_3 &\rightarrow aA_4 \mid b \\ A_4 &\rightarrow c \end{aligned}$$

6-7: Elimination of production  $A_2 \rightarrow A_2A_3$  will be made using productions:

$$\begin{aligned} B_2 &\rightarrow A_3B_2 \\ B_2 &\rightarrow A_3 \\ A_2 &\rightarrow aB_2 \end{aligned}$$

Then, after steps 6-7. the productions will be:

$$\begin{aligned} A_1 &\rightarrow A_2A_3 \mid A_2A_4 \\ A_2 &\rightarrow aB_2 \mid a \\ A_3 &\rightarrow aA_4 \mid b \\ A_4 &\rightarrow c \\ B_2 &\rightarrow A_3B_2 \mid A_3 \end{aligned}$$

8-10: We make substitutions in productions with  $A_1$  in left-hand side. The results is:

$$A_1 \rightarrow aA_3 \mid aB_2A_3 \mid aA_4 \mid aB_2A_4$$

11-13: Similarly with productions with  $B_2$  in left-hand side:

$$B_2 \rightarrow aA_4B_2 \mid aA_3A_4B_2 \mid aA_4 \mid aA_3A_4$$

After the elimination in steps 8-13 of productions in which substitutions were made, the following productions, which are now in Greibach normal form, result:

$$\begin{aligned} A_1 &\rightarrow aA_3 \mid aB_2A_3 \mid aA_4 \mid aB_2A_4 \\ A_2 &\rightarrow aB_2 \mid a \\ A_3 &\rightarrow aA_4 \mid b \\ A_4 &\rightarrow c \\ B_2 &\rightarrow aA_4B_2 \mid aA_3A_4B_2 \mid aA_4 \mid aA_3A_4 \end{aligned}$$

**Example 2.35** Language

$$L = \{a^n b^k c^{n+k} \mid n \geq 0, k \geq 0, n+k > 0\}$$

can be generated by grammar

$$G = \{\{S, R\}, \{a, b, c\}, \{S \rightarrow aSc, S \rightarrow ac, S \rightarrow R, R \rightarrow bRc, R \rightarrow bc\}, S\}$$

First, will eliminate the single unit production, and after this we will give an equivalent grammar in Chomsky normal form, which will be transformed in Greibach normal form.

Productions after the elimination of production  $S \rightarrow R$ :

$$S \rightarrow aSc \mid ac \mid bRc \mid bc$$

$$R \rightarrow bRc \mid bc.$$

We introduce productions  $A \rightarrow a, B \rightarrow b, C \rightarrow c$ , and replace terminals by the corresponding nonterminals:

$$S \rightarrow ASC \mid AC \mid BRC \mid BC,$$

$$R \rightarrow BRC \mid BC,$$

$$A \rightarrow a, B \rightarrow b, C \rightarrow c.$$

After introducing two new nonterminals ( $D, E$ ):

$$S \rightarrow AD \mid AC \mid BE \mid BC,$$

$$D \rightarrow SC,$$

$$E \rightarrow RC,$$

$$R \rightarrow BE \mid BC,$$

$$A \rightarrow a, B \rightarrow b, C \rightarrow c.$$

This is now in Chomsky normal form. Replace the nonterminals to be letters  $A_i$  as in the algorithm.

Then, after applying the replacements

$S$  replaced by  $A_1$ ,  $A$  replaced by  $A_2$ ,  $B$  replaced by  $A_3$ ,  $C$  replaced by  $A_4$ ,  $D$  replaced by  $A_5$ ,

$E$  replaced by  $A_6$ ,  $R$  replaced by  $A_7$ ,

our grammar will have the productions:

$$A_1 \rightarrow A_2A_5 \mid A_2A_4 \mid A_3A_6 \mid A_3A_4,$$

$$A_2 \rightarrow a, A_3 \rightarrow b, A_4 \rightarrow c,$$

$$A_5 \rightarrow A_1A_4,$$

$$A_6 \rightarrow A_7A_4,$$

$$A_7 \rightarrow A_3A_6 \mid A_3A_4.$$

In steps 3–5 of the algorithm the new productions will occur:

$$A_5 \rightarrow A_2A_5A_4 \mid A_2A_4A_4 \mid A_3A_6A_4 \mid A_3A_4A_4 \text{ then}$$

$$A_5 \rightarrow aA_5A_4 \mid aA_4A_4 \mid bA_6A_4 \mid bA_4A_4$$

$$A_7 \rightarrow A_3A_6 \mid A_3A_4, \text{ then}$$

$$A_7 \rightarrow bA_6 \mid bA_4.$$

Therefore

$$A_1 \rightarrow A_2A_5 \mid A_2A_4 \mid A_3A_6 \mid A_3A_4,$$

$$A_2 \rightarrow a, A_3 \rightarrow b, A_4 \rightarrow c,$$

$$A_5 \rightarrow aA_5A_4 \mid aA_4A_4 \mid bA_6A_4 \mid bA_4A_4$$

$$A_6 \rightarrow A_7A_4,$$

$$A_7 \rightarrow bA_6 \mid bA_4.$$

Steps 6–7 will be skipped, because we have no left-recursive productions. In steps 8–10 after the appropriate substitutions we have:

$$A_1 \rightarrow aA_5 \mid aA_4 \mid bA_6 \mid bA_4,$$

$$A_2 \rightarrow a,$$

$$A_3 \rightarrow b,$$

$$A_4 \rightarrow c,$$

$$A_5 \rightarrow aA_5A_4 \mid aA_4A_4 \mid bA_6A_4 \mid bA_4A_4$$

$$A_6 \rightarrow bA_6A_4 \mid bA_4A_4,$$

$$A_7 \rightarrow bA_6 \mid bA_4.$$

## Exercises

**2.3-1** Give pushdown automata to accept the following languages:

$$L_1 = \{a^n cb^n \mid n \geq 0\},$$

$$L_2 = \{a^n b^{2^n} \mid n \geq 1\},$$

$$L_3 = \{a^{2n}b^n \mid n \geq 0\} \cup \{a^n b^{2n} \mid n \geq 0\},$$

**2.3-2** Give a context-free grammar to generate language  $L = \{a^n b^n c^m \mid n \geq 0, m \geq 0\}$ , and transform it in Chomsky and Greibach normal forms. Give a pushdown automaton which accepts  $L$ .

**2.3-3** What languages are generated by the following context-free grammars?

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow SSa, \rightarrow b\}, S), \quad G_2 = (\{S\}, \{a, b\}, \{S \rightarrow SaS, \rightarrow b\}, S)$$

**2.3-4** Give a context-free grammar to generate words with an equal number of letters  $a$  and  $b$ .

**2.3-5** Prove, using the pumping lemma, that a language which words contains an equal number of letters  $a$ ,  $b$  and  $c$  can not be context-free.

**2.3-6** Let the grammar  $G = (V, T, P, S)$ , where

$$V = \{S\},$$

$$T = \{if, then, else, a, c\},$$

$$P = \{S \rightarrow if\ a\ then\ S, S \rightarrow if\ a\ then\ S\ else\ S, S \rightarrow c\},$$

Show that word *if a then if a then c else c* has two different leftmost derivations.

**2.3-7** Prove that if  $L$  is context-free, then  $L^{-1} = \{u^{-1} \mid u \in L\}$  is also context-free.

## Problems

### 2-1. Linear grammars

A grammar  $G = (N, T, P, S)$  which has productions only in the form  $A \rightarrow u_1 B u_2$  or  $A \rightarrow u$ , where  $A, B \in N$ ,  $u, u_1, u_2 \in T^*$ , is called a **linear grammar**. If in a linear grammar all production are of the form  $A \rightarrow B u$  or  $A \rightarrow v$ , then it is called a left-linear grammar. Prove that the language generated by a left-linear grammar is regular.

### 2-2. Operator grammars

An  $\varepsilon$ -free context-free grammar is called **operator grammar** if in the right-hand side of productions there are no two successive nonterminals. Show that, for all  $\varepsilon$ -free context-free grammar an equivalent operator grammar can be built.

### 2-3. Complement of context-free languages

Prove that the class of context-free languages is not closed on complement.

## Chapter notes

In the definition of finite automata instead of transition function we have used the transition graph, which in many cases help us to give simpler proofs.

There exist a lot of classical books on automata and formal languages. We mention from these the following: two books of Aho and Ullman [1, 2] in 1972 and 1973, book of Gécseg and Peák [7] in 1972, two books of Salomaa [21, 22] in 1969 and 1973, a book of Hopcroft and Ullman [11] in 1979, a book of Harrison [9] in 1978, a book of Manna [17], which in 1981 was published also in Hungarian. We notice also a book of Sipser [25] in 1997 and a monograph of Rozenberg and Salomaa [20]. In a book of Lothaire (common name of French authors) [16] on combinatorics of words we can read on other types of automata. Paper of Giammarresi and Montalbano [8] generalise the notion of finite automata. A new

monograph is of Hopcroft, Motwani and Ullman [10]. In German we recommend the student book of Asteroth and Baier [3]. The concise description of the transformation in Greibach normal form is based on this book.

Other books in English: : [4, 5, 6, 12, 13, 14, 15, 18, 19, 23, 24, 26, 27].

At the end of the chapter on compilers another books on the subject are mentioned.



# Bibliography

- [1] A. V. [Aho](#), J. D. [Ullman](#). *The Theory of Parsing, Translation and Compiling Vol. I*. [Prentice-Hall](#), 1972. [138](#)
- [2] A. V. [Aho](#), J. D. [Ullman](#). *The Theory of Parsing, Translation and Compiling Vol. II*. [Prentice-Hall](#), 1973. [138](#)
- [3] A. Asteroth, C. Christel. *Theoretische Informatik*. [Pearson Studium](#), 2002. [139](#)
- [4] J. G. [Brookshear](#). *Theory of Computation – Formal Languages, Automata, and Complexity*. The [Benjamin/Cummings](#) Publishing Company, 1989. [139](#)
- [5] J. Carroll, D. [Long](#). *Theory of Finite Automata*. [Prentice Hall](#), 1989. [139](#)
- [6] D.-Z. Du, K-I. Ko. *Problem Solving in Automata, Languages, and Complexity*. John [Wiley](#) & Sons, 2001. [139](#)
- [7] F. Gécseg, I. Peák. *Algebraic Theory of Automata*. [Akadémiai Kiadó](#), 1972. [138](#)
- [8] D. [Giammarresi](#), R. Montalbano. Deterministic generalized automata. *Theoretical [Computer Science](#)*, 215(1–2):191–208, 1999. [138](#)
- [9] M. A. [Harrison](#). *Introduction to Formal Language Theory*. [Addison-Wesley](#), 1978. [138](#)
- [10] J. E. [Hopcroft](#), R. Motwani, J. D. [Ullman](#). *Introduction to Automata Theory, Languages, and Computation*. [Addison-Wesley](#), 2001 (in German: *Einführung in Automatentheorie, Formale Sprachen und Komplexitätstheorie*, [Pearson Studium](#), 2002). 2. edition. [139](#)
- [11] J. E. [Hopcroft](#), J. D. [Ullman](#). *Introduction to Automata Theory, Languages, and Computation*. [Addison-Wesley](#), 1979. [138](#)
- [12] D. Kelley. *Automata and Formal Languages*. [Prentice Hall](#), 1995. [139](#)
- [13] D. C. [Kozen](#). *Automata and Computability*. [Springer-Verlag](#), 1995. [139](#)
- [14] M. V. [Lawson](#). *Finite Automata*. [Chapman & Hall/CRC](#), 2004. [139](#)
- [15] P. Linz. *An Introduction to Formal Languages and Automata*. [Jones](#) and Barlett Publishers, 2001. [139](#)
- [16] M. Lothaire. *Algebraic Combinatorics on Words*. [Cambridge University Press](#), 2002. [138](#)
- [17] Z. [Manna](#). *Mathematical Theory of Computation*. [McGraw-Hill Book Co.](#), 1974. [138](#)
- [18] A. [Meduna](#). *Automata and Languages: Theory and Applications*. [Springer-Verlag](#), 2000. [139](#)
- [19] R. N. [Moll](#), M. A. Arbib, A. J. Kfoury. *An Introduction to Formal Language Theory*. [Springer-Verlag](#), 1988. [139](#)
- [20] G. [Rozenberg](#), A. [Salomaa](#). *Handbook of Formal Languages, Vols. I–III*. [Springer-Verlag](#), 1997. [138](#)
- [21] A. [Salomaa](#). *Theory of Automata*. [Pergamon Press](#), 1969. [138](#)
- [22] A. [Salomaa](#). *Formal Languages*. [Academic Press](#), 1987 (2. updated edition). [138](#)
- [23] M. Simon. *Automata Theory*. World [Scientific](#) Publishing Company, 1999. [139](#)
- [24] D. A. [Simovoci](#), R. L. Tenney. *Theory of Formal Languages with Applications*. World [Scientific](#) Publishing Company, 1999. [139](#)
- [25] M. Sipser. *Introduction to the Theory of Computation*. [PWS Publishing Company](#), 1997. [138](#)
- [26] H. [Straubing](#). *Finite Automata, Formal Logic, and Circuit Complexity*. [Birkhäuser](#), 1994. [139](#)
- [27] T. A. [Sudkamp](#). *Languages and Machines*. [Addison-Wesley](#), 1997. [139](#)

# Subject Index

## A, Á

alphabet, [76](#)  
ambiguous CF grammar, [131](#)  
DFA-EQUIVALENCE, [97](#)  
  equivalence, [96](#)  
  minimization, [107](#)  
  nondeterministic pushdown, [121](#)  
  pushdown, [121](#)  
AUTOMATON-MINIMIZATION, [108](#)  
automaton  
  finite, [88](#)  
  with  $\epsilon$ -moves, [103](#)

## C

Chomsky hierarchy, [81](#)  
CHOMSKY-NORMAL-FORM, [134](#)  
Chomsky normal form, [134](#)  
concatenation, [76](#)  
configuration, [123](#)  
context-free language  
  pumping lemma, [131](#)

## D

derivation, [78](#)  
derivation tree, [130](#)  
  result of, [130](#)  
deterministic finite automaton, [91](#)  
deterministic pushdown automata, [122](#)

## E, É

edge  
  of finite automaton, *see* transition  
  of pushdown automata, *see* transition  
ELIMINATE-EPSILON-MOVES, [104](#)  
equivalent expressions, [112](#)  
equivalent states, [107](#)  
extended grammar, [84](#)

## F

final state  
  of finite automaton, [89](#)  
  of pushdown automata, [122](#)  
finite automata, [88](#)

  minimization, [107](#)  
finite automaton  
  complete deterministic, [91](#)  
  deterministic, [91](#)  
  nondeterministic, [89](#)  
  with  $\epsilon$ -moves, [103](#)  
FROM-CFG-TO-PUSHDOWN-AUTOMATON, [127](#)

## G

grammar, [78](#)  
  context-free, [81](#), [121](#)  
  context-sensitive, [81](#)  
  generative, [78](#)  
  left-linear, [138](#)  
  linear, [138](#)  
  normal form, [83](#)  
  of type 0,1,2,3, [81](#)  
  operator-, [138](#)  
  phrase structure, [81](#)  
  regular, [81](#)  
GREIBACH-NORMAL-FORM, [135](#)  
Greibach normal form, [135](#)

## I, Í

INACCESSIBLE-STATES, [92](#)  
initial state  
  of finite automaton, [89](#)  
  of pushdown automata, [122](#)  
input alphabet  
  of finite automaton, [89](#)  
  of pushdown automaton, [122](#)

## K

Kleene's theorem, [112](#)

## L

language  
  complement, [77](#)  
  context-free, [81](#), [130](#)  
  context-sensitive, [81](#)  
  iteration, [77](#)  
  mirror, [77](#)  
  of type 0,1,2,3, [81](#)

- phrase-structure, [81](#)
  - power, [77](#)
  - regular, [81](#), [88](#)
  - star operation, [77](#)
  - language generated, [78](#)
  - languages
    - difference, [77](#)
    - intersection, [77](#)
    - multiplication, [77](#)
    - specifying them, [78](#)
    - union, [77](#)
  - leftmost derivation, [131](#)
- M**
- mirror image, [77](#)
- N**
- NFA-DFA, [95](#)
  - NFA-FROM-REGULAR-GRAMMAR, [102](#)
  - nondeterministic finite automaton, [89](#)
  - nondeterministic pushdown automaton, [121](#)
  - NONPRODUCTIVE-STATES, [92](#)
  - normal form
    - Chomsky, [134](#)
    - Greibach, [135](#)
- O, Ó**
- operations
    - on languages, [77](#)
    - on regular languages, [102](#)
- P**
- prefix, [77](#)
  - production, [78](#)
  - proper subword, [77](#)
  - pumping lemma
    - for context-free languages, [131](#)
    - for regular languages, [109](#)
  - pushdown automata, [121](#)
    - deterministic, [122](#)
  - FROM-PUSHDOWN-AUTOMATON-TO-CF-GRAMMAR, [129](#)
- R**
- regular
    - expression, [112](#)
    - language, [88](#)
    - operation, [77](#)
  - regular expression, [111](#)
  - REGULAR-GRAMMAR-FROM-DFA, [100](#)
  - reversal, [77](#)
- S**
- stack alphabet
    - of pushdown automaton, [122](#)
  - start symbol, [78](#)
    - of pushdown automata, [122](#)
  - state
    - inaccessible, [91](#)
    - nonproductive, [92](#)
    - of finite automaton, [89](#)
    - of pushdown automaton, [122](#)
  - subword, [77](#)
  - suffix, [77](#)
  - syntax tree, [130](#)
- T**
- terminal symbols, [78](#)
  - transition
    - of finite automaton, [89](#)
    - of pushdown automata, [122](#)
- U, Ú**
- ELIMINATE-UNIT-PRODUCTIONS, [82](#)
- V**
- variables, [78](#)
- W**
- word, [76](#)
    - power, [76](#)

# Name index

## A, Á

Aho, Alfred V., [138](#), [140](#)  
Arbib, M. A., [140](#)  
Asteroth, Alexander, [139](#), [140](#)

## B

Baier, Christel, [139](#), [140](#)  
Brookshear, J. G., [140](#)

## C

Carroll, J., [140](#)  
Chomsky, Noam, [76](#), [81](#)

## D

Du, Ding-Zhu, [140](#)

## G

Gécseg, Ferenc, [138](#), [140](#)  
Giammarresi, Dora, [138](#), [140](#)

## H

Harrison, Michael A., [138](#), [140](#)  
Hopcroft, John E., [139](#), [140](#)

## K

Kelley, D., [140](#)  
Kfoury, A. J., [140](#)  
Kleene, Stephen C., [112](#)  
Ko, Ker-I, [140](#)  
Kozen, D. C., [140](#)

## L

Lawson, M. V., [140](#)  
Linz, P., [140](#)  
Long, Darrell, [140](#)  
Lothaire, M., [138](#), [140](#)

## M

Manna, Zohar, [138](#), [140](#)  
Meduna, A., [140](#)  
Moll, R. N., [140](#)  
Montalbano, Rosa, [138](#), [140](#)  
Motwani, Rajeev, [139](#), [140](#)

## P

Peák, István (1936–1989), [138](#), [140](#)

## R

Rozenberg, Grzegorz, [138](#), [140](#)

## S

Salomaa, Arto, [138](#), [140](#)  
Simon, M., [140](#)  
Simovoci, Dan A., [140](#)  
Sipser, Michael, [138](#), [140](#)  
Straubing, H., [140](#)  
Sudkamp, Thomas A., [140](#)

## T

Tenney, R. L., [140](#)

## U, Ú

Ullman, Jeffrey David, [138–140](#)

# Contents

<b>2. Automata and Formal Languages (Kása Zoltán)</b>	<b>76</b>
2.1. Languages and grammars	76
2.1.1. Operations on languages	77
2.1.2. Specifying languages	77
Specifying languages by listing their elements	77
Specifying languages by properties	78
Specifying languages by grammars	78
2.1.3. Chomsky hierarchy of grammars and languages	81
Eliminating unit productions	82
Grammars in normal forms	83
2.1.4. Extended grammars	84
2.1.5. Closure properties in the Chomsky-classes	87
2.2. Finite automata and regular languages	88
Eliminating inaccessible states	91
Eliminating nonproductive states	92
2.2.1. Transforming nondeterministic finite automata in deterministic finite automata	93
2.2.2. Equivalence of deterministic finite automata	96
2.2.3. Equivalence of finite automata and regular languages	98
Operation on regular languages	102
2.2.4. Finite automata with $\epsilon$ -moves	103
2.2.5. Minimization of finite automata	107
2.2.6. Pumping lemma for regular languages	109
2.2.7. Regular expressions	111
Associating regular expressions to finite automata	113
Associating finite automata to regular expressions	117
2.3. Pushdown automata and context-free languages	121
2.3.1. Pushdown automata	121
2.3.2. Context-free languages	130
2.3.3. Pumping lemma for context-free languages	131
2.3.4. Normal forms of the context-free languages	133
Chomsky normal form	134
Greibach normal form	135

*Contents*

145

<b>Bibliography</b> . . . . .	<b>140</b>
<b>Subject Index</b> . . . . .	<b>141</b>
<b>Name index</b> . . . . .	<b>143</b>