

10.3.2. Regular Expressions. Regular languages can be characterized as languages defined by *regular expressions*. Given an alphabet T , a regular expression over T is defined recursively as follows:

1. \emptyset , λ , and a are regular expressions for all $a \in T$.
2. If R and S are regular expressions over T the following expressions are also regular: (R) , $R + S$, $R \cdot S$, R^* .

In order to use fewer parentheses we assign those operations the following hierarchy (from do first to do last): $*$, \cdot , $+$. We may omit the dot: $\alpha \cdot \beta = \alpha\beta$.

Next we define recursively the language associated to a given regular expression:

$$\begin{aligned} L(\emptyset) &= \emptyset, \\ L(\lambda) &= \{\lambda\}, \\ L(a) &= \{a\} && \text{for each } a \in T, \\ L(R + S) &= L(R) \cup L(S), \\ L(R \cdot S) &= L(R)L(S) && \text{(language product),} \\ L(R^*) &= L(R)^* && \text{(language closure).} \end{aligned}$$

So, for instance, the expression a^*bb^* represents all strings of the form $a^n b^m$ with $n \geq 0$, $m > 0$, $a^*(b + c)$ is the set of strings consisting of any number of a 's followed by a b or a c , $a(a + b)^*b$ is the set of strings over $\{a, b\}$ than start with a and end with b , etc.

Another way of characterizing regular languages is as sets of strings recognized by finite-state automata, as we will see next. But first we need a generalization of the concept of finite-state automaton.

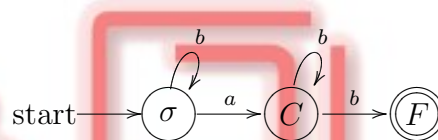
10.3.3. Nondeterministic Finite-State Automata. A *nondeterministic finite-state automaton* is a generalization of a finite-state automaton so that at each state there might be several possible choices for the “next state” instead of just one. Formally a nondeterministic finite-state automaton consists of

1. A finite set of *states* \mathcal{S} .
2. A finite set of *input symbols* \mathcal{J} .
3. A *next-state* or *transition function* $f : \mathcal{S} \times \mathcal{J} \rightarrow \mathcal{P}(\mathcal{S})$.
4. An *initial state* $\sigma \in \mathcal{S}$.

5. A subset \mathcal{F} of \mathcal{S} of *accepting* or *final states*.

We represent the automaton $A = (\mathcal{S}, \mathcal{J}, f, \sigma, \mathcal{F})$. We say that a nondeterministic finite-state automaton *accepts* or *recognizes* a given string of input symbols if in its transition diagram there is a path from the starting state to a final state with its edges labeled by the symbols of the given string. A path (which we can express as a sequence of states) whose edges are labeled with the symbols of a string is said to *represent* the given string.

Example: Consider the nondeterministic finite-state automaton defined by the following transition diagram:



This automaton recognizes precisely the strings of the form $b^n ab^m$, $n \geq 0$, $m > 0$. For instance the string $bbabb$ is represented by the path $(\sigma, \sigma, \sigma, C, C, F)$. Since that path ends in a final state, the string is recognized by the automaton.

Next we will see that there is a precise relation between regular grammars and nondeterministic finite-state automata.

Regular grammar associated to a nondeterministic finite-state automaton. Let A be a non-deterministic finite-state automaton given as a transition diagram. Let σ be the initial state. Let T be the set of inputs symbols, let N be the set of states, and $V = N \cup T$. Let P be the set of productions

$$S \rightarrow xS'$$

if there is an edge labeled x from S to S' and

$$S \rightarrow \lambda$$

if S is a final state. Let G be the regular grammar

$$G = (V, T, \sigma, P).$$

Then the set of strings recognized by A is precisely $L(G)$.

Example: For the nondeterministic automaton defined above the corresponding grammar will be:

$T = \{a, b\}$, $N = \{\sigma, C, F\}$, with productions

$$\sigma \rightarrow b\sigma, \quad \sigma \rightarrow aC, \quad C \rightarrow bC, \quad C \rightarrow bF, \quad F \rightarrow \lambda.$$

The string $bbabb$ can be produced like this:

$$\sigma \Rightarrow b\sigma \Rightarrow bb\sigma \Rightarrow bbaC \Rightarrow bbabC \Rightarrow bbabbF \Rightarrow bbabb.$$

Nondeterministic finite-state automaton associated to a given regular grammar. Let $G = (V, T, \sigma, P)$ be a regular grammar. Let

$$\mathcal{J} = T.$$

$$\mathcal{S} = N \cup \{F\}, \text{ where } N = V - T, \text{ and } F \notin V.$$

$$f(S, x) = \{S' \mid S \rightarrow xS' \in P\} \cup \{F \mid S \rightarrow x \in P\}.$$

$$\mathcal{F} = \{F\} \cup \{S \mid S \rightarrow \lambda \in P\}.$$

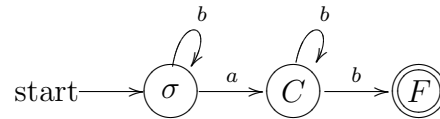
Then the nondeterministic finite-state automaton $A = (\mathcal{S}, \mathcal{J}, f, \sigma, \mathcal{F})$ recognizes precisely the strings in $L(G)$.

10.3.4. Relationships Between Regular Languages and Automata. In the previous section we saw that regular languages coincide with the languages recognized by nondeterministic finite-state automata. Here we will see that the term “nondeterministic” can be dropped, so that regular languages are precisely those recognized by (deterministic) finite-state automata. The idea is to show that given any nondeterministic finite-state automata it is possible to construct an equivalent deterministic finite-state automata recognizing exactly the same set of strings. The main result is the following:

Let $A = (\mathcal{S}, \mathcal{J}, f, \sigma, \mathcal{F})$ be a nondeterministic finite-state automaton. Then A is equivalent to the finite-state automaton $A' = (\mathcal{S}', \mathcal{J}', f', \sigma', \mathcal{F}')$, where

1. $\mathcal{S}' = \mathcal{P}(\mathcal{S})$.
2. $\mathcal{J}' = \mathcal{J}$.
3. $\sigma' = \{\sigma\}$.
4. $\mathcal{F}' = \{X \subseteq \mathcal{S} \mid X \cap \mathcal{F} \neq \emptyset\}$.
5. $f'(X, x) = \bigcup_{S \in X} f(S, x)$, $f'(\emptyset, x) = \emptyset$.

Example: Find a (deterministic) finite-state automaton A' equivalent to the following nondeterministic finite-state automaton A :



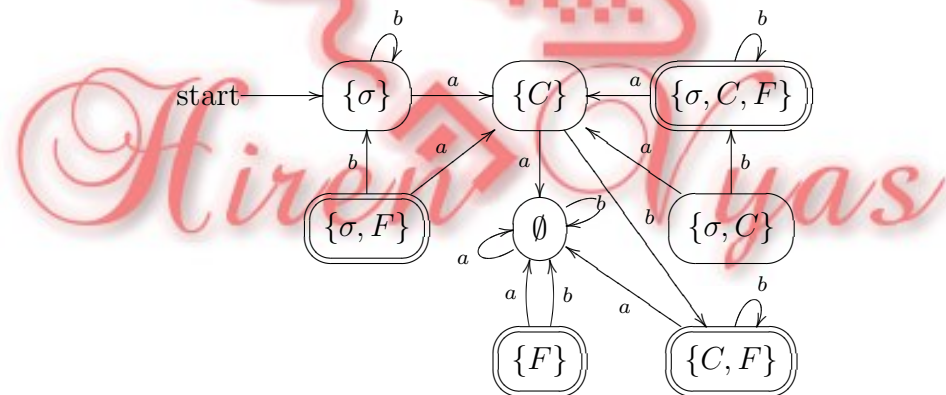
Answer: The set of input symbols is the same as that of the given automaton: $\mathcal{J}' = \mathcal{J} = \{a, b\}$. The set of states is the set of subsets of $\mathcal{S} = \{\sigma, C, F\}$, i.e.:

$$\mathcal{S}' = \{\emptyset, \{\sigma\}, \{C\}, \{F\}, \{\sigma, C\}, \{\sigma, F\}, \{C, F\}, \{\sigma, C, F\}\}.$$

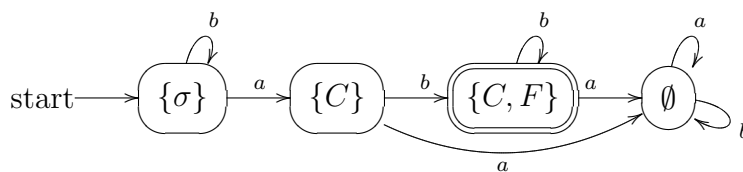
The starting state is $\{\sigma\}$. The final states of A' are the elements of \mathcal{S}' containing some final state of A :

$$\mathcal{F}' = \{\{F\}, \{\sigma, F\}, \{C, F\}, \{\sigma, C, F\}\}.$$

Then for each element X of \mathcal{S}' we draw an edge labeled x from X to $\bigcup_{S \in X} f(S, x)$ (and from \emptyset to \emptyset):



We notice that some states are unreachable from the starting state. After removing the unreachable states we get the following simplified version of the finite-state automaton:



So, once proved that every nondeterministic finite-state automaton is equivalent to some deterministic finite-state automaton, we obtain the main result of this section: A language L is regular if and only if there exists a finite-state automaton that recognizes precisely the strings in L .



APPENDIX A

A.1. Efficient Computation of Powers Modulo m

We illustrate an efficient method of computing powers modulo m with an example. Assume that we want to compute $3^{547} \pmod{10}$. First write 547 in base 2: 1000100011, hence $547 = 2^9 + 2^5 + 2 + 1 = ((2^4 + 1)2^4 + 1)2 + 1$, so: $3^{547} = ((3^{2^4} \cdot 3)^{2^4} \cdot 3)^2 \cdot 3$. Next we compute the expression beginning with the inner parenthesis, and reducing modulo 10 at each step: $3^2 = 9 \pmod{10}$, $3^{2^2} = 9^2 = 81 = 1 \pmod{10}$, $3^{2^3} = 1^2 = 1 \pmod{10}$, $3^{2^4} = 1^2 = 1 \pmod{10}$, $3^{2^4} \cdot 3 = 1 \cdot 3 = 3 \pmod{10}$, etc. At the end we find $3^{547} = 7 \pmod{10}$.

The algorithm in pseudocode would be like this:

```

1: procedure pow_mod(a,x,m) {computes a^x mod m}
2:   p := 1
3:   bx := binary_array(x) {x as a binary array}
4:   t := a mod m
5:   for k := 1 to length(bx)
6:     begin
7:       p := (p * p) mod m
8:       if bx[k] = 1 then
9:         p := (p * t) mod m
10:    end
11:   return p
12: end pow_mod

```

The following is a program in C implementing the algorithm:

```
int pow(int a, int x, int m) {
    int p = 1;
    int y = (1 << (8 * size of(int) - 2));

    a %= m;

    while (!(y & x)) y >>= 1;

    while (y) {
        p *= a;
        p %= m;
        if (x & y) {
            p *= a;
            p %= m;
        }
        y >>= 1;
    }
    return p;
}
```

The following is an alternative algorithm equivalent to running through the binary representation of the exponent from right to left instead of left to right:

```
1: procedure pow_mod(a,x,m) {computes  $a^x \bmod m$ }
2:   p := 1
3:   t := a mod m
4:   while x > 0
5:     begin
6:       if x is odd then
7:         p := (p * t) mod m
8:         t := (t * t) mod m
9:         x := floor(x/2)
10:    end
11:   return p
12: end pow_mod
```

A.2. Machines and Languages

A.2.1. Turing Machines. A *Turing machine* is a theoretical device intended to define rigorously the concept of *algorithm*. It consists of

1. An *infinite tape* made of a sequence of cells. Each cell may be empty or may contain a symbol from a given alphabet.
2. A *control unit* containing a finite set of instructions.
3. A *tape head* able to read and write (or delete) symbols from the tape.

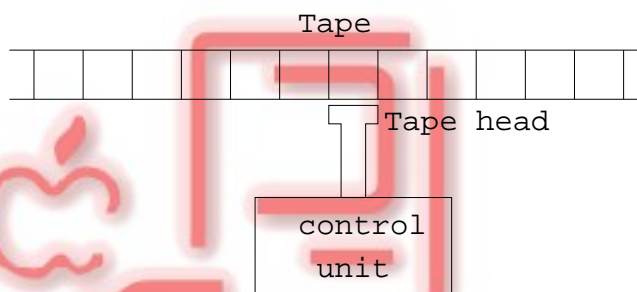


FIGURE A.1. Turing Machine.

Each machine instruction contains the following five parts:

1. The current machine state.
2. A tape symbol read from the current tape cell.
3. A tape symbol to write into the current tape cell.
4. A direction for the tape head to move: L = 'move one cell to the left', R = 'move one cell to the right', S = 'stay in the current cell'.
5. The next machine state.

Turing machines are generalizations of finite-state automata. A finite-state automaton is just a Turing machine whose tape head moves always from left to right and never writes to the tape. The input of the finite-state automaton is presented as symbols written in the tape.

In general we make the following assumptions:

1. An input is represented on the tape by placing the letters of the strings in contiguous tape cells. All other cells contain the blank symbol, which we may denote λ .

2. The tape is initially positioned at the leftmost cell of the input string unless specified otherwise.
3. There is one *starting state*.
4. There is one *halt state*, which we denote by “Halt”.

The execution of a Turing machine stops when it enters the Halt state or when it enters a state for which there is no valid move. The output of the Turing machine is the contents of the tape when the machine stops.

We say that an input string is *accepted* by a Turing machine if the machine enters the Halt state. Otherwise the string is *rejected*. This can happen in two ways: by entering a state other than the Halt state from which there is no move, or by running forever (for instance executing an infinite loop).

If a Turing machine has at least two instructions with the same state and input letter, then the machine is *nondeterministic*. Otherwise it is *deterministic*.

Finite-State Automata. A finite-state automata can be interpreted as a Turing machine whose tape head moves only from left to right and never writes to the tape.

Pushdown Automata. A *pushdown automaton* is finite-state automaton with a stack, i.e., a storage structure in which symbols can be put and extracted from it by two operations: *push* (place on the top of the stack) and *pop* (take from the top of the stack)—consequently the last symbol put into the stack is the first symbol taken out. Additionally there is a third operation, *nop*, that leaves the stack intact. The next state function takes into account not only the current state and the symbol read from the input, but also the symbol at the top of the stack. After reading the next input symbol and the symbol at the top of the stack, the automaton executes a stack operation and goes to the next state. Initially there is a single symbol in the stack.

Linearly Bounded Automata. A *linearly bounded automaton* is a Turing machine whose tape is limited to the size of its input string plus two boundary cells that may not be changed.

Computable Functions. Consider a Turing machine T working on symbols from an alphabet of only one symbol $A = \{|\}$ (“stroke”). Let $f : \mathbb{N} \rightarrow \mathbb{N}$ the function defined so that $f(n) = m$ means that if the

initial input of T consists of a string of $n + 1$ strokes, the output of T is a string of $m + 1$ strokes. We say that f is *computed* by the Turing machine T . A *computable function* is a function computed by some Turing machine. A computable function $f(n)$ *halts* for a given value of its argument n if T with input $n + 1$ strokes halts. A computable function f is *total* if $f(n)$ halts for every n .

An *effective* enumeration of a set is a listing of its elements by an algorithm.

A.2.2. Hierarchy of Languages. Here we mention a hierarchy of languages that includes (and extends) Chomsky's classification, in increasing order of inclusion.

1. *Regular languages*. They are recognized by finite-state automata. *Example:* $\{a^m b^n \mid m, n = 1, 2, 3 \dots\}$.
2. *Deterministic context-free languages*, recognized by deterministic pushdown automata. *Example:* $\{a^n b^n \mid n = 1, 2, 3 \dots\}$.
3. *Context-free languages*, recognized by nondeterministic pushdown automata. *Example:* palindromes over $\{a, b\}$.
4. *Context-sensitive languages*, languages without λ recognized by linearly bounded automata. *Example:* $\{a^n b^n c^n \mid n = 1, 2, 3 \dots\}$
5. *Unrestricted or phrase-structure grammars*, recognized by Turing machines.
6. *Recursively enumerable languages*. A language is recursively enumerable if there is a Turing machine that outputs all the strings of the language. *Example:* $\{a^n \mid f_n(n) \text{ halts}\}$, where f_0, f_1, f_2, \dots is an effective enumeration of all computable functions.
7. *Nongrammatical languages*, languages that are not definable by any grammar and cannot be recognized by Turing machines. *Example:* $\{a^n \mid f_n \text{ is total}\}$.