

CIS511
Introduction to the Theory of Computation
Formal Languages and Automata
Models of Computation

Jean Gallier

February 25, 2009

Chapter 1

Basics of Formal Language Theory

1.1 Generalities, Motivations, Problems

In this part of the course we want to understand

- What is a language?
- How do we define a language?
- How do we manipulate languages, combine them?
- What is the complexity of a language?

Roughly, there are two dual views of languages:

- (A) The *recognition* point view.
- (B) The *generation* point of view.

No matter how we view a language, we are typically considering two things:

- (1) The *syntax*, i.e., what are the “legal” strings in that language (what are the “grammar rules”?).
- (2) The *semantics* of strings in the language, i.e., what is the *meaning* (or *interpretation*) of a string.

The semantics is usually a lot more interesting than the syntax but unfortunately much more difficult to deal with!

Therefore, sorry, we will only be dealing with syntax!

In (A), we typically assume some kind of “black box”, M , (an *automaton*) that takes a string, w , as input and returns two possible answers:

Yes, the string w is *accepted*, which means that w belongs to the language, L , that we are trying to define.

No, the string w is *rejected*, which means that w *does not* belong to the language, L .

Usually, the black box M gives a definite answer for every input after a finite number of steps, but not always.

For example, a Turing machine may go on computing forever and not give any answer for certain strings not in the language. This is an example of *undecidability*.

The black box may compute *deterministically* or *non-deterministically*, which means roughly that on input w , the machine M is allowed to try different computations and to ignore failing computations as long as there is some successful computation on input w .

This affects greatly the *complexity* of recognition, i.e., how many steps it takes to process w .

Sometimes, a nondeterministic version of an automaton turns out to be equivalent to the deterministic version (although, with different complexity).

This tends to happen for very restrictive models—where nondeterminism does not help, or for very powerful models—where again, nondeterminism does not help, but because the deterministic model is already very powerful!

We will investigate automata of increasing power of recognition:

- (1) Deterministic and nondeterministic finite automata (DFA's and NFA's, their power is the same).
- (2) Pushdown automata (PDA's) and deterministic pushdown automata (DPDA's), here $\text{PDA} > \text{DPDA}$.
- (3) Deterministic and nondeterministic Turing machines (their power is the same).
- (4) If time permits, we will also consider some restricted type of Turing machine known as LBA (linear bounded automaton).

In (B), we are interested in formalisms that specify a language in terms of *rules* that allow the generation of “legal” strings. The most common formalism is that of a formal *grammar*.

Remember:

- An automaton *recognizes* (or *accepts*) a language,
- a grammar *generates* a language.
- *gramm**a*r is spelled with an “a” (not with an “e”).
- The plural of automa*t**on* is automa*t**a* (not automa*t**ons*).

For “good” classes of grammars, it is possible to build an automaton, M_G , from the grammar, G , in the class, so that M_G recognizes the language, $L(G)$, generated by the grammar G .

However, grammars are nondeterministic in nature. Thus, even if we try to avoid nondeterministic automata, we usually can't escape having to deal with them.

We will investigate the following types of grammars (the so-called *Chomsky hierarchy*) and the corresponding families of languages:

- (1) Regular grammars (type 3-languages).
- (2) Context-free grammars (type 2-languages).
- (3) The recursively enumerable languages or r.e. sets (type 0-languages).
- (4) If time permit, context-sensitive languages (type 1-languages).

Miracle: The grammars of type (1), (2), (3), (4) correspond exactly to the automata of the corresponding type!

Furthermore, there are *algorithms* for converting grammars to the corresponding automata (and backward), although some of these algorithms are not practical.

Building an automaton from a grammar is an important practical problem in language processing. A lot is known for the regular and the context-free grammars, but there is still room for improvements and innovations!

There are other ways of defining families of languages, for example

Inductive closures.

In this style of definition, a collection of basic (atomic) languages is specified, some operations to combine languages are also specified, and the family of languages is defined as the smallest one containing the given atomic languages and closed under the operations.

Investigating closure properties (for example, union, intersection) is a way to assess how “robust” (or complex) a family of languages is.

Well, it is now time to be precise!

1.2 Alphabets, Strings, Languages

Our view of languages is that a language is a set of strings. In turn, a string is a finite sequence of letters from some alphabet. These concepts are defined rigorously as follows.

Definition 1.2.1 An *alphabet* Σ is any **finite** set.

We often write $\Sigma = \{a_1, \dots, a_k\}$. The a_i are called the *symbols* of the alphabet.

Examples:

$$\Sigma = \{a\}$$

$$\Sigma = \{a, b, c\}$$

$$\Sigma = \{0, 1\}$$

A string is a finite sequence of symbols. Technically, it is convenient to define strings as functions. For any integer $n \geq 1$, let

$$[n] = \{1, 2, \dots, n\},$$

and for $n = 0$, let

$$[0] = \emptyset.$$

Definition 1.2.2 Given an alphabet Σ , a *string over Σ (or simply a string) of length n* is any function

$$u: [n] \rightarrow \Sigma.$$

The integer n is the *length* of the string u , and it is denoted as $|u|$. When $n = 0$, the special string $u: [0] \rightarrow \Sigma$ of length 0 is called the *empty string, or null string*, and is denoted as ϵ .

Given a string $u: [n] \rightarrow \Sigma$ of length $n \geq 1$, $u(i)$ is the i -th letter in the string u . For simplicity of notation, we denote the string u as

$$u = u_1u_2 \dots u_n,$$

with each $u_i \in \Sigma$.

For example, if $\Sigma = \{a, b\}$ and $u: [3] \rightarrow \Sigma$ is defined such that $u(1) = a$, $u(2) = b$, and $u(3) = a$, we write

$$u = aba.$$

Strings of length 1 are functions $u: [1] \rightarrow \Sigma$ simply picking some element $u(1) = a_i$ in Σ . Thus, we will identify every symbol $a_i \in \Sigma$ with the corresponding string of length 1.

The set of all strings over an alphabet Σ , including the empty string, is denoted as Σ^* .

Observe that when $\Sigma = \emptyset$, then

$$\emptyset^* = \{\epsilon\}.$$

When $\Sigma \neq \emptyset$, the set Σ^* is countably infinite. Later on, we will see ways of ordering and enumerating strings.

Strings can be juxtaposed, or concatenated.

Definition 1.2.3 Given an alphabet Σ , given any two strings $u: [m] \rightarrow \Sigma$ and $v: [n] \rightarrow \Sigma$, the *concatenation* $u \cdot v$ (also written uv) of u and v is the string $uv: [m + n] \rightarrow \Sigma$, defined such that

$$uv(i) = \begin{cases} u(i) & \text{if } 1 \leq i \leq m, \\ v(i - m) & \text{if } m + 1 \leq i \leq m + n. \end{cases}$$

In particular, $u\epsilon = \epsilon u = u$.

It is immediately verified that

$$u(vw) = (uv)w.$$

Thus, concatenation is a binary operation on Σ^* which is associative and has ϵ as an identity.

Note that generally, $uv \neq vu$, for example for $u = a$ and $v = b$.

Given a string $u \in \Sigma^*$ and $n \geq 0$, we define u^n as follows:

$$u^n = \begin{cases} \epsilon & \text{if } n = 0, \\ u^{n-1}u & \text{if } n \geq 1. \end{cases}$$

Clearly, $u^1 = u$, and it is an easy exercise to show that

$$u^n u = u u^n,$$

for all $n \geq 0$.

Definition 1.2.4 Given an alphabet Σ , given any two strings $u, v \in \Sigma^*$ we define the following notions as follows:

u is a prefix of v iff there is some $y \in \Sigma^*$ such that

$$v = uy.$$

u is a suffix of v iff there is some $x \in \Sigma^*$ such that

$$v = xu.$$

u is a substring of v iff there are some $x, y \in \Sigma^*$ such that

$$v = xuy.$$

We say that *u is a proper prefix (suffix, substring) of v* iff *u* is a prefix (suffix, substring) of *v* and $u \neq v$.

Recall that a partial ordering \leq on a set S is a binary relation $\leq \subseteq S \times S$ which is reflexive, transitive, and antisymmetric.

The concepts of prefix, suffix, and substring, define binary relations on Σ^* in the obvious way. It can be shown that these relations are partial orderings.

Another important ordering on strings is the lexicographic (or dictionary) ordering.

Definition 1.2.5 Given an alphabet $\Sigma = \{a_1, \dots, a_k\}$ assumed totally ordered such that $a_1 < a_2 < \dots < a_k$, given any two strings $u, v \in \Sigma^*$, we define the *lexicographic ordering* \preceq as follows:

$$u \preceq v \quad \left\{ \begin{array}{l} \text{if } v = uy, \text{ for some } y \in \Sigma^*, \text{ or} \\ \text{if } u = xa_iy, v = xa_jz, \\ \text{and } a_i < a_j, \text{ for some } x, y, z \in \Sigma^*. \end{array} \right.$$

It is fairly tedious to prove that the lexicographic ordering is in fact a partial ordering. In fact, it is a *total ordering*, which means that for any two strings $u, v \in \Sigma^*$, either $u \preceq v$, or $v \preceq u$.

The *reversal* w^R of a string w is defined inductively as follows:

$$\begin{aligned} \epsilon^R &= \epsilon, \\ (ua)^R &= au^R, \end{aligned}$$

where $a \in \Sigma$ and $u \in \Sigma^*$.

It can be shown that

$$(uv)^R = v^R u^R.$$

Thus,

$$(u_1 \dots u_n)^R = u_n^R \dots u_1^R,$$

and when $u_i \in \Sigma$, we have

$$(u_1 \dots u_n)^R = u_n \dots u_1.$$

We can now define languages.

Definition 1.2.6 Given an alphabet Σ , a *language over Σ* (or simply a *language*) is any subset L of Σ^* .

If $\Sigma \neq \emptyset$, there are uncountably many languages. We will try to single out countable “tractable” families of languages. We will begin with the family of *regular languages*, and then proceed to the *context-free languages*.

We now turn to operations on languages.

1.3 Operations on Languages

A way of building more complex languages from simpler ones is to combine them using various operations. First, we review the set-theoretic operations of union, intersection, and complementation.

Given some alphabet Σ , for any two languages L_1, L_2 over Σ , the *union* $L_1 \cup L_2$ of L_1 and L_2 is the language

$$L_1 \cup L_2 = \{w \in \Sigma^* \mid w \in L_1 \text{ or } w \in L_2\}.$$

The *intersection* $L_1 \cap L_2$ of L_1 and L_2 is the language

$$L_1 \cap L_2 = \{w \in \Sigma^* \mid w \in L_1 \text{ and } w \in L_2\}.$$

The *difference* $L_1 - L_2$ of L_1 and L_2 is the language

$$L_1 - L_2 = \{w \in \Sigma^* \mid w \in L_1 \text{ and } w \notin L_2\}.$$

The difference is also called the *relative complement*.

A special case of the difference is obtained when $L_1 = \Sigma^*$, in which case we define the *complement* \overline{L} of a language L as

$$\overline{L} = \{w \in \Sigma^* \mid w \notin L\}.$$

The above operations do not use the structure of strings. The following operations use concatenation.

Definition 1.3.1 Given an alphabet Σ , for any two languages L_1, L_2 over Σ , the *concatenation* L_1L_2 of L_1 and L_2 is the language

$$L_1L_2 = \{w \in \Sigma^* \mid \exists u \in L_1, \exists v \in L_2, w = uv\}.$$

For any language L , we define L^n as follows:

$$\begin{aligned} L^0 &= \{\epsilon\}, \\ L^{n+1} &= L^nL. \end{aligned}$$

The following properties are easily verified:

$$\begin{aligned}
 L\emptyset &= \emptyset, \\
 \emptyset L &= \emptyset, \\
 L\{\epsilon\} &= L, \\
 \{\epsilon\}L &= L, \\
 (L_1 \cup \{\epsilon\})L_2 &= L_1L_2 \cup L_2, \\
 L_1(L_2 \cup \{\epsilon\}) &= L_1L_2 \cup L_1, \\
 L^n L &= LL^n.
 \end{aligned}$$

In general, $L_1L_2 \neq L_2L_1$.

So far, the operations that we have introduced, except complementation (since $\overline{L} = \Sigma^* - L$ is infinite if L is finite and Σ is nonempty), preserve the finiteness of languages. This is not the case for the next two operations.

Definition 1.3.2 Given an alphabet Σ , for any language L over Σ , the *Kleene *-closure* L^* of L is the language

$$L^* = \bigcup_{n \geq 0} L^n.$$

The *Kleene +-closure* L^+ of L is the language

$$L^+ = \bigcup_{n \geq 1} L^n.$$

Thus, L^* is the infinite union

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots \cup L^n \cup \dots,$$

and L^+ is the infinite union

$$L^+ = L^1 \cup L^2 \cup \dots \cup L^n \cup \dots$$

Since $L^1 = L$, both L^* and L^+ contain L .

In fact,

$$L^+ = \{w \in \Sigma^*, \exists n \geq 1, \\ \exists u_1 \in L \cdots \exists u_n \in L, w = u_1 \cdots u_n\},$$

and since $L^0 = \{\epsilon\}$,

$$L^* = \{\epsilon\} \cup \{w \in \Sigma^*, \exists n \geq 1, \\ \exists u_1 \in L \cdots \exists u_n \in L, w = u_1 \cdots u_n\}.$$

Thus, the language L^* always contains ϵ , and we have

$$L^* = L^+ \cup \{\epsilon\}.$$

However, if $\epsilon \notin L$, then $\epsilon \notin L^+$. The following is easily shown:

$$\begin{aligned}\emptyset^* &= \{\epsilon\}, \\ L^+ &= L^*L, \\ L^{**} &= L^*, \\ L^*L^* &= L^*.\end{aligned}$$

The Kleene closures have many other interesting properties.

Homomorphisms are also very useful.

Given two alphabets Σ, Δ , a *homomorphism* $h: \Sigma^* \rightarrow \Delta^*$ between Σ^* and Δ^* is a function $h: \Sigma^* \rightarrow \Delta^*$ such that

$$h(uv) = h(u)h(v)$$

for all $u, v \in \Sigma^*$.

Letting $u = v = \epsilon$, we get

$$h(\epsilon) = h(\epsilon)h(\epsilon),$$

which implies that (why?)

$$h(\epsilon) = \epsilon.$$

If $\Sigma = \{a_1, \dots, a_k\}$, it is easily seen that h is completely determined by $h(a_1), \dots, h(a_k)$ (why?)

Example: $\Sigma = \{a, b, c\}$, $\Delta = \{0, 1\}$, and

$$h(a) = 01, \quad h(b) = 011, \quad h(c) = 0111.$$

For example

$$h(abc) = 010110110111.$$

Given any language $L_1 \subseteq \Sigma^*$, we define the *image* $h(L_1)$ of L_1 as

$$h(L_1) = \{h(u) \in \Delta^* \mid u \in L_1\}.$$

Given any language $L_2 \subseteq \Delta^*$, we define the *inverse image* $h^{-1}(L_2)$ of L_2 as

$$h^{-1}(L_2) = \{u \in \Sigma^* \mid h(u) \in L_2\}.$$

We now turn to the first formalism for defining languages, Deterministic Finite Automata (DFA's)

Chapter 2

Regular Languages

2.1 Deterministic Finite Automata (DFA's)

First we define what DFA's are, and then we explain how they are used to accept or reject strings. Roughly speaking, a DFA is a finite transition graph whose edges are labeled with letters from an alphabet Σ .

The graph also satisfies certain properties that make it deterministic. Basically, this means that given any string w , starting from any node, there is a unique path in the graph “parsing” the string w .

Example 1. A DFA for the language

$$L_1 = \{ab\}^+ = \{ab\}^* \{ab\},$$

i.e.,

$$L_1 = \{ab, abab, ababab, \dots, (ab)^n, \dots\}.$$

Input alphabet: $\Sigma = \{a, b\}$.

State set $Q_1 = \{0, 1, 2, 3\}$.

Start state: 0.

Set of accepting states: $F_1 = \{2\}$.

Transition table (function) δ_1 :

	<i>a</i>	<i>b</i>
0	1	3
1	3	2
2	1	3
3	3	3

Note that state 3 is a *trap state* or *dead state*.

Example 2. A DFA for the language

$$L_2 = \{ab\}^* = L_1 \cup \{\epsilon\}$$

i.e.,

$$L_2 = \{\epsilon, ab, abab, ababab, \dots, (ab)^n, \dots\}.$$

Input alphabet: $\Sigma = \{a, b\}$.

State set $Q_2 = \{0, 1, 2\}$.

Start state: 0.

Set of accepting states: $F_2 = \{0\}$.

Transition table (function) δ_2 :

	<i>a</i>	<i>b</i>
0	1	2
1	2	0
2	2	2

State 2 is a *trap state* or *dead state*.

Example 3. A DFA for the language

$$L_3 = \{a, b\}^* \{abb\}.$$

Note that L_3 consists of all strings of a 's and b 's ending in abb .

Input alphabet: $\Sigma = \{a, b\}$.

State set $Q_3 = \{0, 1, 2, 3\}$.

Start state: 0.

Set of accepting states: $F_3 = \{3\}$.

Transition table (function) δ_3 :

	a	b
0	1	0
1	1	2
2	1	3
3	1	0

Is this a minimal DFA?

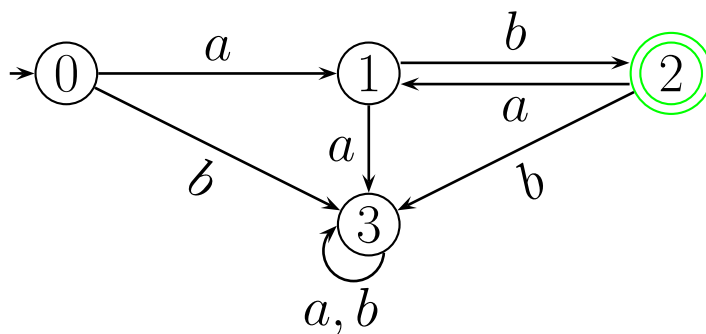


Figure 2.1: DFA for $\{ab\}^+$

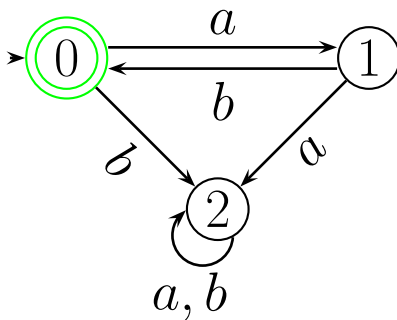


Figure 2.2: DFA for $\{ab\}^*$

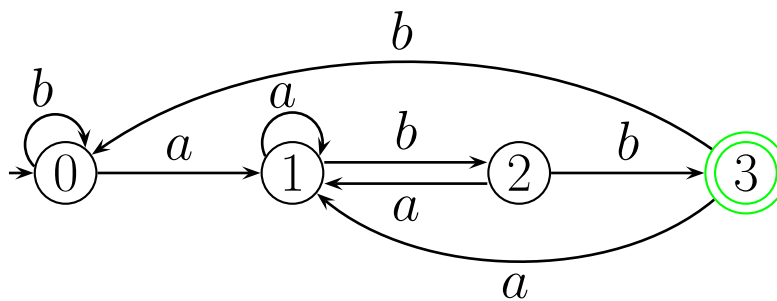


Figure 2.3: DFA for $\{a, b\}^*\{abb\}$

Definition 2.1.1 A *deterministic finite automaton (or DFA)* is a quintuple $D = (Q, \Sigma, \delta, q_0, F)$, where

- Σ is a finite *input alphabet*
- Q is a finite set of *states*;
- F is a subset of Q of *final (or accepting) states*;
- $q_0 \in Q$ is the *start state (or initial state)*;
- δ is the *transition function*, a function

$$\delta: Q \times \Sigma \rightarrow Q.$$

For any state $p \in Q$ and any input $a \in \Sigma$, the state $q = \delta(p, a)$ is uniquely determined. Thus, it is possible to define the state reached from a given state $p \in Q$ on input $w \in \Sigma^*$, following the path specified by w . Technically, this is done by defining the extended transition function $\delta^*: Q \times \Sigma^* \rightarrow Q$.

Definition 2.1.2 Given a DFA $D = (Q, \Sigma, \delta, q_0, F)$, the *extended transition function* $\delta^*: Q \times \Sigma^* \rightarrow Q$ is defined as follows:

$$\begin{aligned}\delta^*(p, \epsilon) &= p, \\ \delta^*(p, ua) &= \delta(\delta^*(p, u), a),\end{aligned}$$

where $a \in \Sigma$ and $u \in \Sigma^*$.

It is immediate that $\delta^*(p, a) = \delta(p, a)$ for $a \in \Sigma$. The meaning of $\delta^*(p, w)$ is that it is the state reached from state p following the path from p specified by w .

It is also easy to show that

$$\delta^*(p, uv) = \delta^*(\delta^*(p, u), v).$$

We can now define how a DFA accepts or rejects a string.

Definition 2.1.3 Given a DFA $D = (Q, \Sigma, \delta, q_0, F)$, the *language $L(D)$ accepted (or recognized) by D* is the language

$$L(D) = \{w \in \Sigma^* \mid \delta^*(q_0, w) \in F\}.$$

Thus, a string $w \in \Sigma^*$ is accepted iff the path from q_0 on input w ends in a final state.

We now come to the first of several equivalent definitions of the regular languages.

Regular Languages, Version 1

A language L is a *regular language* if it is accepted by some DFA.

Note that a regular language may be accepted by many different DFAs. Later on, we will investigate how to find minimal DFA's. (For a given regular language, L , a minimal DFA for L is a DFA with the smallest number of states among all DFA's accepting L . A minimal DFA for L must exist since every nonempty subset of natural numbers has a smallest element.)

In order to understand how complex the regular languages are, we will investigate the closure properties of the regular languages under union, intersection, complementation, concatenation, and Kleene $*$.

It turns out that the family of regular languages is closed under all these operations. For union, intersection, and complementation, we can use the cross-product construction which preserves determinism.

However, for concatenation and Kleene $*$, there does not appear to be any method involving DFA's only. The way to do it is to introduce nondeterministic finite automata (NFA's).

2.2 The “Cross-product” Construction

Let $\Sigma = \{a_1, \dots, a_m\}$ be an alphabet.

Given any two DFA's $D_1 = (Q_1, \Sigma, \delta_1, q_{0,1}, F_1)$ and $D_2 = (Q_2, \Sigma, \delta_2, q_{0,2}, F_2)$, there is a very useful construction for showing that the union, the intersection, or the relative complement of regular languages, is a regular language.

Given any two languages L_1, L_2 over Σ , recall that

$$\begin{aligned} L_1 \cup L_2 &= \{w \in \Sigma^* \mid w \in L_1 \quad \text{or} \quad w \in L_2\}, \\ L_1 \cap L_2 &= \{w \in \Sigma^* \mid w \in L_1 \quad \text{and} \quad w \in L_2\}, \\ L_1 - L_2 &= \{w \in \Sigma^* \mid w \in L_1 \quad \text{and} \quad w \notin L_2\}. \end{aligned}$$

Let us first explain how to construct a DFA accepting the intersection $L_1 \cap L_2$. Let D_1 and D_2 be DFA's such that $L_1 = L(D_1)$ and $L_2 = L(D_2)$. The idea is to construct a DFA simulating D_1 and D_2 in parallel. This can be done by using states which are pairs $(p_1, p_2) \in Q_1 \times Q_2$. Thus, we define the DFA D as follows:

$$D = (Q_1 \times Q_2, \Sigma, \delta, (q_{0,1}, q_{0,2}), F_1 \times F_2),$$

where the transition function $\delta: (Q_1 \times Q_2) \times \Sigma \rightarrow Q_1 \times Q_2$ is defined as follows:

$$\delta((p_1, p_2), a) = (\delta_1(p_1, a), \delta_2(p_2, a)),$$

for all $p_1 \in Q_1$, $p_2 \in Q_2$, and $a \in \Sigma$.

Clearly, D is a DFA, since D_1 and D_2 are. Also, by the definition of δ , we have

$$\delta^*((p_1, p_2), w) = ((\delta_1^*(p_1, w), \delta_2^*(p_2, w))),$$

for all $p_1 \in Q_1$, $p_2 \in Q_2$, and $w \in \Sigma^*$.

Now, we have $w \in L(D_1) \cap L(D_2)$

- iff $w \in L(D_1)$ and $w \in L(D_2)$,
- iff $\delta_1^*(q_{0,1}, w) \in F_1$ and $\delta_2^*(q_{0,2}, w) \in F_2$,
- iff $(\delta_1^*(q_{0,1}, w), \delta_2^*(q_{0,2}, w)) \in F_1 \times F_2$,
- iff $\delta^*((q_{0,1}, q_{0,2}), w) \in F_1 \times F_2$,
- iff $w \in L(D)$.

Thus, $L(D) = L(D_1) \cap L(D_2)$.

We can now modify D very easily to accept $L(D_1) \cup L(D_2)$. We change the set of final states so that it becomes $(F_1 \times Q_2) \cup (Q_1 \times F_2)$.

Indeed, $w \in L(D_1) \cup L(D_2)$

- iff $w \in L(D_1)$ or $w \in L(D_2)$,
- iff $\delta_1^*(q_{0,1}, w) \in F_1$ or $\delta_2^*(q_{0,2}, w) \in F_2$,
- iff $(\delta_1^*(q_{0,1}, w), \delta_2^*(q_{0,2}, w)) \in (F_1 \times Q_2) \cup (Q_1 \times F_2)$,
- iff $\delta^*((q_{0,1}, q_{0,2}), w) \in (F_1 \times Q_2) \cup (Q_1 \times F_2)$,
- iff $w \in L(D)$.

Thus, $L(D) = L(D_1) \cup L(D_2)$.

We can also modify D very easily to accept $L(D_1) - L(D_2)$. We change the set of final states so that it becomes $F_1 \times (Q_2 - F_2)$. Indeed, $w \in L(D_1) - L(D_2)$

- iff $w \in L(D_1)$ and $w \notin L(D_2)$,
- iff $\delta_1^*(q_{0,1}, w) \in F_1$ and $\delta_2^*(q_{0,2}, w) \notin F_2$,
- iff $(\delta_1^*(q_{0,1}, w), \delta_2^*(q_{0,2}, w)) \in F_1 \times (Q_2 - F_2)$,
- iff $\delta^*((q_{0,1}, q_{0,2}), w) \in F_1 \times (Q_2 - F_2)$,
- iff $w \in L(D)$.

Thus, $L(D) = L(D_1) - L(D_2)$.

In all cases, if D_1 has n_1 states and D_2 has n_2 states, the DFA D has $n_1 n_2$ states.

2.3 Morphisms, F -Maps, B -Maps and Homomorphisms of DFA's

A map between DFA's is a certain kind of graph homomorphism. The following Definition is adapted from Eilenberg.

Definition 2.3.1 Given any two DFA's

$D_1 = (Q_1, \Sigma, \delta_1, q_{0,1}, F_1)$ and $D_2 = (Q_2, \Sigma, \delta_2, q_{0,2}, F_2)$,

a *morphism $h: D_1 \rightarrow D_2$ of DFA's* is a function

$h: Q_1 \rightarrow Q_2$ satisfying the following conditions:

(1)

$$h(\delta_1(p, a)) = \delta_2(h(p), a),$$

for all $p \in Q_1$ and all $a \in \Sigma$;

(2) $h(q_{0,1}) = q_{0,2}$.

An *F -map of DFA's*, for short, a *map*, is a morphism of DFA's $h: D_1 \rightarrow D_2$ that satisfies the condition

(3a) $h(F_1) \subseteq F_2$.

A *B -map of DFA's* is a morphism of DFA's

$h: D_1 \rightarrow D_2$ that satisfies the condition

(3b) $h^{-1}(F_2) \subseteq F_1$.

A *proper homomorphism of DFA's*, for short, a *homomorphism*, is an F -map of DFA's that is also a B -map of DFA's.

Now, for any function $f: X \rightarrow Y$ and any two subsets $A \subseteq X$ and $B \subseteq Y$,

$$f(A) \subseteq B \quad \text{iff} \quad A \subseteq f^{-1}(B).$$

Thus, (3a) & (3b) is equivalent to the condition

$$(3c) \quad h^{-1}(F_2) = F_1.$$

Note that the condition for being a proper homomorphism of DFA's is **not** equivalent to

$$h(F_1) = F_2.$$

It forces $h(F_1) = F_2 \cap h(Q_1)$, and furthermore, for every $p \in Q_1$, whenever $h(p) \in F_2$, then $p \in F_1$.

The reader should check that if $f: D_1 \rightarrow D_2$ and $g: D_2 \rightarrow D_3$ are morphisms (resp. F -map, resp. B -map), then $g \circ f: D_1 \rightarrow D_3$ is also a morphism (resp. F -map, resp. B -map).

Remark: In previous versions of these notes, an F -map was called simply a *map* and a B -map was called an F^{-1} -*map*. Over the years, the old terminology proved to be confusing. We hope the new one is less confusing!

Note that an F -map or a B -map is a special case of the concept of *simulation* of automata. A proper homomorphism is a special case of a *bisimulation*. Bisimulations play an important role in real-time systems and in concurrency theory.

The main motivation behind these definitions is that when there is an F -map $h: D_1 \rightarrow D_2$, somehow, D_2 simulates D_1 , and it turns out that $L(D_1) \subseteq L(D_2)$.

When there is a B -map $h: D_1 \rightarrow D_2$, somehow, D_1 simulates D_2 , and it turns out that $L(D_2) \subseteq L(D_1)$.

When there is a proper homomorphism $h: D_1 \rightarrow D_2$, somehow, D_1 bisimulates D_2 , and it turns out that $L(D_2) = L(D_1)$.

A DFA morphism (resp. F -map, resp. B -map), $f: D_1 \rightarrow D_2$, is an *isomorphism* iff there is a DFA morphism (resp. F -map, resp. B -map), $g: D_2 \rightarrow D_1$, so that

$$g \circ f = \text{id}_{D_1} \quad \text{and} \quad f \circ g = \text{id}_{D_2}.$$

The map g is unique and it is denoted f^{-1} . The reader should prove that if a DFA F -map is an isomorphism, then it is also a proper homomorphism and if a DFA B -map is an isomorphism, then it is also a proper homomorphism.

If $h: D_1 \rightarrow D_2$ is a morphism of DFA's, it is easily shown by induction on the length of w that

$$h(\delta_1^*(p, w)) = \delta_2^*(h(p), w),$$

for all $p \in Q_1$ and all $w \in \Sigma^*$.

As a consequence, we have the following Lemma:

Lemma 2.3.2 *If $h: D_1 \rightarrow D_2$ is an F -map of DFA's, then $L(D_1) \subseteq L(D_2)$. If $h: D_1 \rightarrow D_2$ is a B -map of DFA's, then $L(D_2) \subseteq L(D_1)$. Finally, if $h: D_1 \rightarrow D_2$ is a proper homomorphism of DFA's, then $L(D_1) = L(D_2)$.*

A DFA is *accessible, or trim*, if every state is reachable from the start state.

A morphism (resp. F -map, B -map) $h: D_1 \rightarrow D_2$ is *surjective* if $h(Q_1) = Q_2$.

It can be shown that if D_1 is trim, then there is a most one morphism $h: D_1 \rightarrow D_2$ (resp. F -map, B -map). If D_2 is also trim and we have a morphism $h: D_1 \rightarrow D_2$, then h is surjective.

It can also be shown that a minimal DFA D_L for L is characterized by the property that there is unique surjective proper homomorphism $h: D \rightarrow D_L$ from any trim DFA D accepting L to D_L .

Another useful notion is the notion of a congruence on a DFA.

Definition 2.3.3 Given any DFA

$D = (Q, \Sigma, \delta, q_0, F)$, a *congruence* \equiv on D is an equivalence relation \equiv on Q satisfying the following conditions:

- (1) If $p \equiv q$, then $\delta(p, a) \equiv \delta(q, a)$, for all $p, q \in Q$ and all $a \in \Sigma$.
- (2) If $p \equiv q$ and $p \in F$, then $q \in F$, for all $p, q \in Q$.

It can be shown that a proper homomorphism of DFA's $h: D_1 \rightarrow D_2$ induces a congruence \equiv_h on D_1 defined as follows:

$$p \equiv_h q \quad \text{iff} \quad h(p) = h(q).$$

Given a congruence \equiv on a DFA D , we can define the *quotient DFA* D / \equiv , and there is a surjective proper homomorphism $\pi: D \rightarrow D / \equiv$.

We will come back to this point when we study minimal DFA's.

2.4 Nondeterministic Finite Automata (NFA's)

NFA's are obtained from DFA's by allowing multiple transitions from a given state on a given input. This can be done by defining $\delta(p, a)$ as a **subset** of Q rather than a single state. It will also be convenient to allow transitions on input ϵ .

We let 2^Q denote the set of all subsets of Q , including the empty set. The set 2^Q is the *power set* of Q . We define NFA's as follows.

Example 4. A NFA for the language

$$L_3 = \{a, b\}^* \{abb\}.$$

Input alphabet: $\Sigma = \{a, b\}$.

State set $Q_4 = \{0, 1, 2, 3\}$.

Start state: 0.

Set of accepting states: $F_4 = \{3\}$.

Transition table δ_4 :

	a	b
0	$\{0, 1\}$	$\{0\}$
1	\emptyset	$\{2\}$
2	\emptyset	$\{3\}$
3	\emptyset	\emptyset

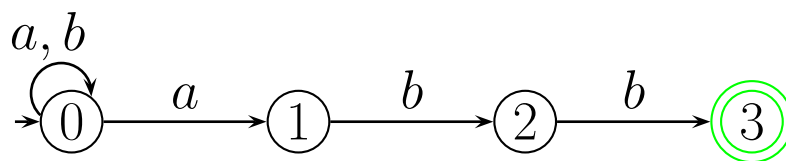


Figure 2.4: NFA for $\{a, b\}^* \{abb\}$

Example 5. Let $\Sigma = \{a_1, \dots, a_n\}$, let

$$L_n^i = \{w \in \Sigma^* \mid w \text{ contains an odd number of } a_i\text{'s}\},$$

and let

$$L_n = L_n^1 \cup L_n^2 \cup \dots \cup L_n^n.$$

The language L_n consists of those strings over Σ that contain an odd number of some letter $a_i \in \Sigma^*$.

Equivalently $\Sigma^* - L_n$ consists of those strings over Σ with an even number of *every* letter $a_i \in \Sigma^*$.

It can be shown that that every DFA accepting L_n has at least 2^n states.

However, there is an NFA with $2n + 1$ states accepting L_n (and even with $2n$ states!).

Definition 2.4.1 A *nondeterministic finite automaton (or NFA)* is a quintuple $N = (Q, \Sigma, \delta, q_0, F)$, where

- Σ is a finite *input alphabet*
- Q is a finite set of *states*;
- F is a subset of Q of *final (or accepting) states*;
- $q_0 \in Q$ is the *start state (or initial state)*;
- δ is the *transition function*, a function

$$\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q.$$

For any state $p \in Q$ and any input $a \in \Sigma \cup \{\epsilon\}$, the set of states $\delta(p, a)$ is uniquely determined. We write $q \in \delta(p, a)$.

Given an NFA $N = (Q, \Sigma, \delta, q_0, F)$, we would like to define the language accepted by N , and for this, we need to extend the transition function $\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$ to a function

$$\delta^*: Q \times \Sigma^* \rightarrow 2^Q.$$

The presence of ϵ -transitions (i.e., when $q \in \delta(p, \epsilon)$) causes technical problems, and to overcome these problems, we introduce the notion of ϵ -closure.

2.5 ϵ -Closure

Definition 2.5.1 Given an NFA $N = (Q, \Sigma, \delta, q_0, F)$ (with ϵ -transitions) for every state $p \in Q$, the *ϵ -closure of p* is set $\epsilon\text{-closure}(p)$ consisting of all states q such that there is a path from p to q whose spelling is ϵ . This means that either $q = p$, or that all the edges on the path from p to q have the label ϵ .

We can compute $\epsilon\text{-closure}(p)$ using a sequence of approximations as follows. Define the sequence of sets of states $(\epsilon\text{-clo}_i(p))_{i \geq 0}$ as follows:

$$\begin{aligned} \epsilon\text{-clo}_0(p) &= \{p\}, \\ \epsilon\text{-clo}_{i+1}(p) &= \epsilon\text{-clo}_i(p) \cup \\ &\quad \{q \in Q \mid \exists s \in \epsilon\text{-clo}_i(p), q \in \delta(s, \epsilon)\}. \end{aligned}$$

Since $\epsilon\text{-clo}_i(p) \subseteq \epsilon\text{-clo}_{i+1}(p)$, $\epsilon\text{-clo}_i(p) \subseteq Q$, for all $i \geq 0$, and Q is finite, there is a smallest i , say i_0 , such that

$$\epsilon\text{-clo}_{i_0}(p) = \epsilon\text{-clo}_{i_0+1}(p),$$

and it is immediately verified that

$$\epsilon\text{-closure}(p) = \epsilon\text{-clo}_{i_0}(p).$$

When N has no ϵ -transitions, i.e., when $\delta(p, \epsilon) = \emptyset$ for all $p \in Q$ (which means that δ can be viewed as a function $\delta: Q \times \Sigma \rightarrow 2^Q$), we have

$$\epsilon\text{-closure}(p) = \{p\}.$$

It should be noted that there are more efficient ways of computing $\epsilon\text{-closure}(p)$, for example, using a stack (basically, a kind of depth-first search).

We present such an algorithm below. It is assumed that the types *NFA* and *stack* are defined. If n is the number of states of an NFA N , we let

*eclo*type = **array**[1.. n] **of boolean**

```

function eclosure[N: NFA, p: integer]: eclotype;
  begin
    var eclo: eclotype, q, s: integer, st: stack;
    for each  $q \in \text{setstates}(N)$  do
      eclo[q] := false;
    endfor
    eclo[p] := true; st := empty;
    trans := deltatable(N);
    st := push(st, p);
    while  $st \neq \text{emptystack}$  do
      q = pop(st);
      for each  $s \in \text{trans}(q, \epsilon)$  do
        if eclo[s] = false then
          eclo[s] := true; st := push(st, s)
        endif
      endfor
    endwhile;
    eclosure := eclo
  end

```

This algorithm can be easily adapted to compute the set of states reachable from a given state p (in a DFA or an NFA).

Given a subset S of Q , we define ϵ -closure(S) as

$$\epsilon\text{-closure}(S) = \bigcup_{p \in S} \epsilon\text{-closure}(p).$$

When N has no ϵ -transitions, we have

$$\epsilon\text{-closure}(S) = S.$$

We are now ready to define the extension $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$ of the transition function $\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$.

2.6 Converting an NFA into a DFA

The intuition behind the definition of the extended transition function is that $\delta^*(p, w)$ is the set of all states reachable from p by a path whose spelling is w .

Definition 2.6.1 Given an NFA $N = (Q, \Sigma, \delta, q_0, F)$ (with ϵ -transitions), the *extended transition function* $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$ is defined as follows: for every $p \in Q$, every $u \in \Sigma^*$, and every $a \in \Sigma$,

$$\begin{aligned}\delta^*(p, \epsilon) &= \epsilon\text{-closure}(\{p\}), \\ \delta^*(p, ua) &= \epsilon\text{-closure}\left(\bigcup_{s \in \delta^*(p, u)} \delta(s, a)\right).\end{aligned}$$

The *language $L(N)$ accepted by an NFA N* is the set

$$L(N) = \{w \in \Sigma^* \mid \delta^*(q_0, w) \cap F \neq \emptyset\}.$$

We can also extend $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$ to a function

$$\widehat{\delta}: 2^Q \times \Sigma^* \rightarrow 2^Q$$

defined as follows: for every subset S of Q , for every $w \in \Sigma^*$,

$$\widehat{\delta}(S, w) = \bigcup_{p \in S} \delta^*(p, w).$$

Let \mathcal{Q} be the subset of 2^Q consisting of those subsets S of Q that are ϵ -closed, i.e., such that $S = \epsilon\text{-closure}(S)$. If we consider the restriction

$$\Delta: \mathcal{Q} \times \Sigma \rightarrow \mathcal{Q}$$

of $\widehat{\delta}: 2^Q \times \Sigma^* \rightarrow 2^Q$ to \mathcal{Q} and Σ , we observe that Δ is the transition function of a DFA. Indeed, this is the transition function of a DFA accepting $L(N)$. It is easy to show that Δ is defined directly as follows (on subsets S in \mathcal{Q}):

$$\Delta(S, a) = \epsilon\text{-closure}\left(\bigcup_{s \in S} \delta(s, a)\right).$$

Then, the DFA D is defined as follows:

$$D = (\mathcal{Q}, \Sigma, \Delta, \epsilon\text{-closure}(\{q_0\}), \mathcal{F}),$$

where $\mathcal{F} = \{S \in \mathcal{Q} \mid S \cap F \neq \emptyset\}$.

It is not difficult to show that $L(D) = L(N)$, that is, D is a DFA accepting $L(N)$. Thus, we have converted the NFA N into a DFA D (and gotten rid of ϵ -transitions).

Since DFA's are special NFA's, the subset construction shows that DFA's and NFA's accept *the same* family of languages, the *regular languages, version 1* (although not with the same complexity).

The states of the DFA D equivalent to N are ϵ -closed subsets of Q . For this reason, the above construction is often called the “subset construction”. It is due to Rabin and Scott. Although theoretically fine, the method may construct useless sets S that are not reachable from the start state $\epsilon\text{-closure}(\{q_0\})$. A more economical construction is given next.

An Algorithm to convert an NFA into a DFA: The “subset construction”

Given an input NFA $N = (Q, \Sigma, \delta, q_0, F)$, a DFA $D = (K, \Sigma, \Delta, S_0, \mathcal{F})$ is constructed. It is assumed that K is a linear array of sets of states $S \subseteq Q$, and Δ is a 2-dimensional array, where $\Delta[i, a]$ is the target state of the transition from $K[i] = S$ on input a , with $S \in K$, and $a \in \Sigma$.

$S_0 := \epsilon\text{-closure}(\{q_0\}); total := 1; K[1] := S_0;$

$marked := 0;$

while $marked < total$ **do**;

$marked := marked + 1; S := K[marked];$

for each $a \in \Sigma$ **do**

$U := \bigcup_{s \in S} \delta(s, a); T := \epsilon\text{-closure}(U);$

if $T \notin K$ **then**

$total := total + 1; K[total] := T$

endif;

$\Delta[marked, a] := T$

endfor

endwhile;

$\mathcal{F} := \{S \in K \mid S \cap F \neq \emptyset\}$

Let us illustrate the subset construction on the NFA of Example 4.

A NFA for the language

$$L_3 = \{a, b\}^* \{abb\}.$$

Transition table δ_4 :

	a	b
0	$\{0, 1\}$	$\{0\}$
1	\emptyset	$\{2\}$
2	\emptyset	$\{3\}$
3	\emptyset	\emptyset

Set of accepting states: $F_4 = \{3\}$.

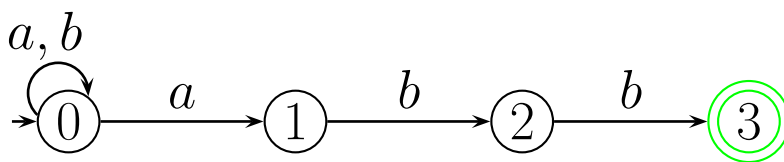


Figure 2.5: NFA for $\{a, b\}^* \{abb\}$

The pointer \Rightarrow corresponds to *marked* and the pointer \rightarrow to *total*.

Initial transition table Δ .

\Rightarrow	names	states	<i>a</i>	<i>b</i>
\rightarrow	<i>A</i>	{0}		

Just after entering the while loop

$\Rightarrow \rightarrow$	names	states	<i>a</i>	<i>b</i>
	<i>A</i>	{0}		

After the first round through the while loop.

\Rightarrow	names	states	<i>a</i>	<i>b</i>
	<i>A</i>	{0}	<i>B</i>	<i>A</i>
\rightarrow	<i>B</i>	{0, 1}		

After just reentering the while loop.

$\Rightarrow \rightarrow$	names	states	<i>a</i>	<i>b</i>
	<i>A</i>	{0}	<i>B</i>	<i>A</i>
	<i>B</i>	{0, 1}		

After the second round through the while loop.

\Rightarrow	names	states	<i>a</i>	<i>b</i>
	<i>A</i>	{0}	<i>B</i>	<i>A</i>
\rightarrow	<i>B</i>	{0, 1}	<i>B</i>	<i>C</i>
	<i>C</i>	{0, 2}		

After the third round through the while loop.

	names	states	a	b
	A	$\{0\}$	B	A
	B	$\{0, 1\}$	B	C
\Rightarrow	C	$\{0, 2\}$	B	D
\rightarrow	D	$\{0, 3\}$		

After the fourth round through the while loop.

	names	states	a	b
	A	$\{0\}$	B	A
	B	$\{0, 1\}$	B	C
	C	$\{0, 2\}$	B	D
$\Rightarrow \rightarrow$	D	$\{0, 3\}$	B	A

This is the DFA of Figure 2.3, except that in that example A, B, C, D are renamed $0, 1, 2, 3$

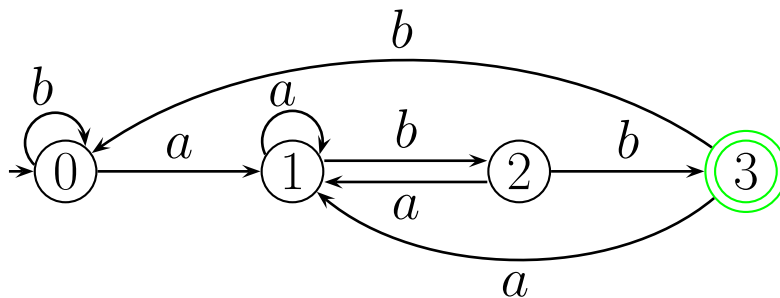


Figure 2.6: DFA for $\{a, b\}^*\{abb\}$