## FORMAL DEFINITION OF A FINITE AUTOMATON

In the preceding section, we used state diagrams to introduce finite automata. Now we define finite automata formally. Although state diagrams are easier to grasp intuitively, we need the formal definition, too, for two specific reasons.

First, a formal definition is precise. It resolves any uncertainties about what is allowed in a finite automaton. If you were uncertain about whether finite automata were allowed to have 0 accept states or whether they must have exactly one transition exiting every state for each possible input symbol, you could consult the formal definition and verify that the answer is yes in both cases. Second, a formal definition provides notation. Good notation helps you think and express your thoughts clearly.

The language of a formal definition is somewhat arcane, having some similarity to the language of a legal document. Both need to be precise, and every detail must be spelled out.

A finite automaton has several parts. It has a set of states and rules for going from one state to another, depending on the input symbol. It has an input alphabet that indicates the allowed input symbols. It has a start state and a set of accept states. The formal definition says that a finite automaton is a list of those five objects: set of states, input alphabet, rules for moving, start state, and accept states. In mathematical language, a list of five elements is often called a 5-tuple. Hence we define a finite automaton to be a 5-tuple consisting of these five parts.

We use something called a ***transition function***, frequently denoted $\delta$, to define the rules for moving. If the finite automaton has an arrow from a state $x$ to a state $y$ labeled with the input symbol 1, that means that if the automaton is in state $x$ when it reads a 1, it then moves to state $y$. We can indicate the same thing with the transition function by saying that $\delta(x, 1) = y$. This notation is a kind of mathematical shorthand. Putting it all together, we arrive at the formal definition of finite automata.

---

DEFINITION **1.5**

A ***finite automaton*** is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

1. $Q$ is a finite set called the ***states***,
2. $\Sigma$ is a finite set called the ***alphabet***,
3. $\delta\colon Q \times \Sigma \longrightarrow Q$ is the ***transition function***,[1]
4. $q_0 \in Q$ is the ***start state***, and
5. $F \subseteq Q$ is the ***set of accept states***.[2]

---

[1] Refer back to page 7 if you are uncertain about the meaning of $\delta\colon Q \times \Sigma \longrightarrow Q$.
[2] Accept states sometimes are called ***final states***.

EXAMPLE **1.15**

Consider a generalization of Example 1.13, using the same four-symbol alphabet $\Sigma$. For each $i \geq 1$ let $A_i$ be the language of all strings where the sum of the numbers is a multiple of $i$, except that the sum is reset to $0$ whenever the symbol $\langle\text{RESET}\rangle$ appears. For each $A_i$ we give a finite automaton $B_i$, recognizing $A_i$. We describe the machine $B_i$ formally as follows: $B_i = (Q_i, \Sigma, \delta_i, q_0, \{q_0\})$, where $Q_i$ is the set of $i$ states $\{q_0, q_1, q_2, \ldots, q_{i-1}\}$, and we design the transition function $\delta_i$ so that for each $j$, if $B_i$ is in $q_j$, the running sum is $j$, modulo $i$. For each $q_j$ let

$$\delta_i(q_j, 0) = q_j,$$
$$\delta_i(q_j, 1) = q_k, \text{where } k = j + 1 \text{ modulo } i,$$
$$\delta_i(q_j, 2) = q_k, \text{where } k = j + 2 \text{ modulo } i, \text{ and}$$
$$\delta_i(q_j, \langle\text{RESET}\rangle) = q_0.$$

## FORMAL DEFINITION OF COMPUTATION

So far we have described finite automata informally, using state diagrams, and with a formal definition, as a 5-tuple. The informal description is easier to grasp at first, but the formal definition is useful for making the notion precise, resolving any ambiguities that may have occurred in the informal description. Next we do the same for a finite automaton's computation. We already have an informal idea of the way it computes, and we now formalize it mathematically.

Let $M = (Q, \Sigma, \delta, q_0, F)$ be a finite automaton and let $w = w_1 w_2 \cdots w_n$ be a string where each $w_i$ is a member of the alphabet $\Sigma$. Then $M$ ***accepts*** $w$ if a sequence of states $r_0, r_1, \ldots, r_n$ in $Q$ exists with three conditions:

  **1.** $r_0 = q_0$,
  **2.** $\delta(r_i, w_{i+1}) = r_{i+1}$, for $i = 0, \ldots, n - 1$, and
  **3.** $r_n \in F$.

Condition 1 says that the machine starts in the start state. Condition 2 says that the machine goes from state to state according to the transition function. Condition 3 says that the machine accepts its input if it ends up in an accept state. We say that $M$ ***recognizes language*** $A$ if $A = \{w | M \text{ accepts } w\}$.

DEFINITION **1.16**

A language is called a ***regular language*** if some finite automaton recognizes it.

"any number" includes 0 as a possibility, the empty string $\varepsilon$ is always a member of $A^*$, no matter what $A$ is.

**EXAMPLE 1.24** ---------------------------------------------------------------

Let the alphabet $\Sigma$ be the standard 26 letters $\{\texttt{a}, \texttt{b}, \ldots, \texttt{z}\}$. If $A = \{\texttt{good}, \texttt{bad}\}$ and $B = \{\texttt{boy}, \texttt{girl}\}$, then

$A \cup B = \{\texttt{good}, \texttt{bad}, \texttt{boy}, \texttt{girl}\}$,

$A \circ B = \{\texttt{goodboy}, \texttt{goodgirl}, \texttt{badboy}, \texttt{badgirl}\}$, and

$A^* = \{\varepsilon, \texttt{good}, \texttt{bad}, \texttt{goodgood}, \texttt{goodbad}, \texttt{badgood}, \texttt{badbad},$
$\quad\quad \texttt{goodgoodgood}, \texttt{goodgoodbad}, \texttt{goodbadgood}, \texttt{goodbadbad}, \ldots\}$.

Let $\mathcal{N} = \{1, 2, 3, \ldots\}$ be the set of natural numbers. When we say that $\mathcal{N}$ is *closed under multiplication*, we mean that for any $x$ and $y$ in $\mathcal{N}$, the product $x \times y$ also is in $\mathcal{N}$. In contrast, $\mathcal{N}$ is not closed under division, as 1 and 2 are in $\mathcal{N}$ but $1/2$ is not. Generally speaking, a collection of objects is ***closed*** under some operation if applying that operation to members of the collection returns an object still in the collection. We show that the collection of regular languages is closed under all three of the regular operations. In Section 1.3, we show that these are useful tools for manipulating regular languages and understanding the power of finite automata. We begin with the union operation.

**THEOREM 1.25** ---------------------------------------------------------------

The class of regular languages is closed under the union operation.

In other words, if $A_1$ and $A_2$ are regular languages, so is $A_1 \cup A_2$.

**PROOF IDEA** We have regular languages $A_1$ and $A_2$ and want to show that $A_1 \cup A_2$ also is regular. Because $A_1$ and $A_2$ are regular, we know that some finite automaton $M_1$ recognizes $A_1$ and some finite automaton $M_2$ recognizes $A_2$. To prove that $A_1 \cup A_2$ is regular, we demonstrate a finite automaton, call it $M$, that recognizes $A_1 \cup A_2$.

This is a proof by construction. We construct $M$ from $M_1$ and $M_2$. Machine $M$ must accept its input exactly when either $M_1$ or $M_2$ would accept it in order to recognize the union language. It works by *simulating* both $M_1$ and $M_2$ and accepting if either of the simulations accept.

How can we make machine $M$ simulate $M_1$ and $M_2$? Perhaps it first simulates $M_1$ on the input and then simulates $M_2$ on the input. But we must be careful here! Once the symbols of the input have been read and used to simulate $M_1$, we can't "rewind the input tape" to try the simulation on $M_2$. We need another approach.

This concludes the construction of the finite automaton $M$ that recognizes the union of $A_1$ and $A_2$. This construction is fairly simple, and thus its correctness is evident from the strategy described in the proof idea. More complicated constructions require additional discussion to prove correctness. A formal correctness proof for a construction of this type usually proceeds by induction. For an example of a construction proved correct, see the proof of Theorem 1.54. Most of the constructions that you will encounter in this course are fairly simple and so do not require a formal correctness proof.

We have just shown that the union of two regular languages is regular, thereby proving that the class of regular languages is closed under the union operation. We now turn to the concatenation operation and attempt to show that the class of regular languages is closed under that operation, too.

**THEOREM   1.26**

The class of regular languages is closed under the concatenation operation.

In other words, if $A_1$ and $A_2$ are regular languages then so is $A_1 \circ A_2$.

To prove this theorem, let's try something along the lines of the proof of the union case. As before, we can start with finite automata $M_1$ and $M_2$ recognizing the regular languages $A_1$ and $A_2$. But now, instead of constructing automaton $M$ to accept its input if either $M_1$ or $M_2$ accept, it must accept if its input can be broken into two pieces, where $M_1$ accepts the first piece and $M_2$ accepts the second piece. The problem is that $M$ doesn't know where to break its input (i.e., where the first part ends and the second begins). To solve this problem, we introduce a new technique called nondeterminism.

# 1.2

## NONDETERMINISM

Nondeterminism is a useful concept that has had great impact on the theory of computation. So far in our discussion, every step of a computation follows in a unique way from the preceding step. When the machine is in a given state and reads the next input symbol, we know what the next state will be—it is determined. We call this *determinstic* computation. In a ***nondeterministic*** machine, several choices may exist for the next state at any point.

Nondeterminism is a generalization of determinism, so every deterministic finite automaton is automatically a nondeterministic finite automaton. As Figure 1.27 shows, nondeterministic finite automata may have additional features.
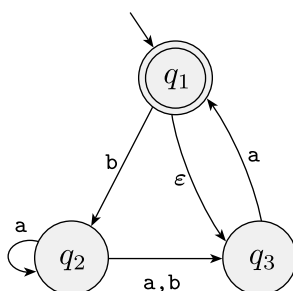
**FIGURE** **1.36**
The NFA $N_4$

## FORMAL DEFINITION OF A NONDETERMINISTIC FINITE AUTOMATON

The formal definition of a nondeterministic finite automaton is similar to that of a deterministic finite automaton. Both have states, an input alphabet, a transition function, a start state, and a collection of accept states. However, they differ in one essential way: in the type of transition function. In a DFA, the transition function takes a state and an input symbol and produces the next state. In an NFA, the transition function takes a state and an input symbol *or the empty string* and produces *the set of possible next states*. In order to write the formal definition, we need to set up some additional notation. For any set $Q$ we write $\mathcal{P}(Q)$ to be the collection of all subsets of $Q$. Here $\mathcal{P}(Q)$ is called the ***power set*** of $Q$. For any alphabet $\Sigma$ we write $\Sigma_\varepsilon$ to be $\Sigma \cup \{\varepsilon\}$. Now we can write the formal description of the type of the transition function in an NFA as $\delta \colon Q \times \Sigma_\varepsilon \longrightarrow \mathcal{P}(Q)$.

---

**DEFINITION** **1.37**

A ***nondeterministic finite automaton*** is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

**1.** $Q$ is a finite set of states,

**2.** $\Sigma$ is a finite alphabet,

**3.** $\delta \colon Q \times \Sigma_\varepsilon \longrightarrow \mathcal{P}(Q)$ is the transition function,

**4.** $q_0 \in Q$ is the start state, and

**5.** $F \subseteq Q$ is the set of accept states.

---

EXAMPLE **1.38**

Recall the NFA $N_1$:



The formal description of $N_1$ is $(Q, \Sigma, \delta, q_1, F)$, where

1. $Q = \{q_1, q_2, q_3, q_4\}$,
2. $\Sigma = \{0,1\}$,
3. $\delta$ is given as

|       | 0         | 1             | $\varepsilon$ |
|-------|-----------|---------------|---------------|
| $q_1$ | $\{q_1\}$ | $\{q_1, q_2\}$ | $\emptyset$   |
| $q_2$ | $\{q_3\}$ | $\emptyset$   | $\{q_3\}$     |
| $q_3$ | $\emptyset$ | $\{q_4\}$   | $\emptyset$   |
| $q_4$ | $\{q_4\}$ | $\{q_4\}$    | $\emptyset$,  |

4. $q_1$ is the start state, and
5. $F = \{q_4\}$.

The formal definition of computation for an NFA is similar to that for a DFA. Let $N = (Q, \Sigma, \delta, q_0, F)$ be an NFA and $w$ a string over the alphabet $\Sigma$. Then we say that $N$ ***accepts*** $w$ if we can write $w$ as $w = y_1 y_2 \cdots y_m$, where each $y_i$ is a member of $\Sigma_\varepsilon$ and a sequence of states $r_0, r_1, \ldots, r_m$ exists in $Q$ with three conditions:

1. $r_0 = q_0$,
2. $r_{i+1} \in \delta(r_i, y_{i+1})$, for $i = 0, \ldots, m - 1$, and
3. $r_m \in F$.

Condition 1 says that the machine starts out in the start state. Condition 2 says that state $r_{i+1}$ is one of the allowable next states when $N$ is in state $r_i$ and reading $y_{i+1}$. Observe that $\delta(r_i, y_{i+1})$ is the *set* of allowable next states and so we say that $r_{i+1}$ is a member of that set. Finally, condition 3 says that the machine accepts its input if the last state is an accept state.

## EQUIVALENCE OF NFAS AND DFAS

Deterministic and nondeterministic finite automata recognize the same class of languages. Such equivalence is both surprising and useful. It is surprising because NFAs appear to have more power than DFAs, so we might expect that NFAs recognize more languages. It is useful because describing an NFA for a given language sometimes is much easier than describing a DFA for that language.

Say that two machines are ***equivalent*** if they recognize the same language.

technique of nondeterminism. Reviewing the first proof, appearing on page 45, may be worthwhile to see how much easier and more intuitive the new proof is.

THEOREM    **1.45** ......................................................................................

The class of regular languages is closed under the union operation.

**PROOF IDEA**    We have regular languages $A_1$ and $A_2$ and want to prove that $A_1 \cup A_2$ is regular. The idea is to take two NFAs, $N_1$ and $N_2$ for $A_1$ and $A_2$, and combine them into one new NFA, $N$.

Machine $N$ must accept its input if either $N_1$ or $N_2$ accepts this input. The new machine has a new start state that branches to the start states of the old machines with $\varepsilon$ arrows. In this way, the new machine nondeterministically guesses which of the two machines accepts the input. If one of them accepts the input, $N$ will accept it, too.

We represent this construction in the following figure. On the left, we indicate the start and accept states of machines $N_1$ and $N_2$ with large circles and some additional states with small circles. On the right, we show how to combine $N_1$ and $N_2$ into $N$ by adding additional transition arrows.



FIGURE    **1.46**
Construction of an NFA $N$ to recognize $A_1 \cup A_2$

**PROOF**

Let $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$ recognize $A_1$, and
    $N_2 = (Q_2, \Sigma, \delta_2, q_2, F_2)$ recognize $A_2$.

Construct $N = (Q, \Sigma, \delta, q_0, F)$ to recognize $A_1 \cup A_2$.

1. $Q = \{q_0\} \cup Q_1 \cup Q_2$.
   The states of $N$ are all the states of $N_1$ and $N_2$, with the addition of a new start state $q_0$.
2. The state $q_0$ is the start state of $N$.
3. The set of accept states $F = F_1 \cup F_2$.
   The accept states of $N$ are all the accept states of $N_1$ and $N_2$. That way, $N$ accepts if either $N_1$ accepts or $N_2$ accepts.
4. Define $\delta$ so that for any $q \in Q$ and any $a \in \Sigma_\varepsilon$,

$$\delta(q, a) = \begin{cases} \delta_1(q, a) & q \in Q_1 \\ \delta_2(q, a) & q \in Q_2 \\ \{q_1, q_2\} & q = q_0 \text{ and } a = \varepsilon \\ \emptyset & q = q_0 \text{ and } a \neq \varepsilon. \end{cases}$$

---

Now we can prove closure under concatenation. Recall that earlier, without nondeterminism, completing the proof would have been difficult.

THEOREM   **1.47**

The class of regular languages is closed under the concatenation operation.

**PROOF IDEA**   We have regular languages $A_1$ and $A_2$ and want to prove that $A_1 \circ A_2$ is regular. The idea is to take two NFAs, $N_1$ and $N_2$ for $A_1$ and $A_2$, and combine them into a new NFA $N$ as we did for the case of union, but this time in a different way, as shown in Figure 1.48.

Assign $N$'s start state to be the start state of $N_1$. The accept states of $N_1$ have additional $\varepsilon$ arrows that nondeterministically allow branching to $N_2$ whenever $N_1$ is in an accept state, signifying that it has found an initial piece of the input that constitutes a string in $A_1$. The accept states of $N$ are the accept states of $N_2$ only. Therefore, it accepts when the input can be split into two parts, the first accepted by $N_1$ and the second by $N_2$. We can think of $N$ as nondeterministically guessing where to make the split.

THEOREM **1.49** ........................................................................................................................

The class of regular languages is closed under the star operation.

**PROOF IDEA** We have a regular language $A_1$ and want to prove that $A_1^*$ also is regular. We take an NFA $N_1$ for $A_1$ and modify it to recognize $A_1^*$, as shown in the following figure. The resulting NFA $N$ will accept its input whenever it can be broken into several pieces and $N_1$ accepts each piece.

We can construct $N$ like $N_1$ with additional $\varepsilon$ arrows returning to the start state from the accept states. This way, when processing gets to the end of a piece that $N_1$ accepts, the machine $N$ has the option of jumping back to the start state to try to read another piece that $N_1$ accepts. In addition, we must modify $N$ so that it accepts $\varepsilon$, which always is a member of $A_1^*$. One (slightly bad) idea is simply to add the start state to the set of accept states. This approach certainly adds $\varepsilon$ to the recognized language, but it may also add other, undesired strings. Exercise 1.15 asks for an example of the failure of this idea. The way to fix it is to add a new start state, which also is an accept state, and which has an $\varepsilon$ arrow to the old start state. This solution has the desired effect of adding $\varepsilon$ to the language without adding anything else.



FIGURE **1.50**
Construction of $N$ to recognize $A^*$

**PROOF** Let $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$ recognize $A_1$.
Construct $N = (Q, \Sigma, \delta, q_0, F)$ to recognize $A_1^*$.

1. $Q = \{q_0\} \cup Q_1$.
   The states of $N$ are the states of $N_1$ plus a new start state.
2. The state $q_0$ is the new start state.
3. $F = \{q_0\} \cup F_1$.
   The accept states are the old accept states plus the new start state.

Seemingly, we are in danger of defining the notion of a regular expression in terms of itself. If true, we would have a ***circular definition***, which would be invalid. However, $R_1$ and $R_2$ always are smaller than $R$. Thus we actually are defining regular expressions in terms of smaller regular expressions and thereby avoiding circularity. A definition of this type is called an ***inductive definition***.

Parentheses in an expression may be omitted. If they are, evaluation is done in the precedence order: star, then concatenation, then union.

For convenience, we let $R^+$ be shorthand for $RR^*$. In other words, whereas $R^*$ has all strings that are 0 or more concatenations of strings from $R$, the language $R^+$ has all strings that are 1 or more concatenations of strings from $R$. So $R^+ \cup \varepsilon = R^*$. In addition, we let $R^k$ be shorthand for the concatenation of $k$ $R$'s with each other.

When we want to distinguish between a regular expression $R$ and the language that it describes, we write $L(R)$ to be the language of $R$.

---

**EXAMPLE**    **1.53** ....................................................................................

In the following instances, we assume that the alphabet $\Sigma$ is $\{0,1\}$.

1. $0^*10^* = \{w|\ w \text{ contains a single } 1\}$.
2. $\Sigma^*1\Sigma^* = \{w|\ w \text{ has at least one } 1\}$.
3. $\Sigma^*001\Sigma^* = \{w|\ w \text{ contains the string } 001 \text{ as a substring}\}$.
4. $1^*(01^+)^* = \{w|\ \text{every } 0 \text{ in } w \text{ is followed by at least one } 1\}$.
5. $(\Sigma\Sigma)^* = \{w|\ w \text{ is a string of even length}\}$.[5]
6. $(\Sigma\Sigma\Sigma)^* = \{w|\ \text{the length of } w \text{ is a multiple of } 3\}$.
7. $01 \cup 10 = \{01, 10\}$.
8. $0\Sigma^*0 \cup 1\Sigma^*1 \cup 0 \cup 1 = \{w|\ w \text{ starts and ends with the same symbol}\}$.
9. $(0 \cup \varepsilon)1^* = 01^* \cup 1^*$.
   The expression $0 \cup \varepsilon$ describes the language $\{0, \varepsilon\}$, so the concatenation operation adds either $0$ or $\varepsilon$ before every string in $1^*$.
10. $(0 \cup \varepsilon)(1 \cup \varepsilon) = \{\varepsilon, 0, 1, 01\}$.
11. $1^*\emptyset = \emptyset$.
    Concatenating the empty set to any set yields the empty set.
12. $\emptyset^* = \{\varepsilon\}$.
    The star operation puts together any number of strings from the language to get a string in the result. If the language is empty, the star operation can put together 0 strings, giving only the empty string.

---

[5]The ***length*** of a string is the number of symbols that it contains.

We break this procedure into two parts, using a new type of finite automaton called a ***generalized nondeterministic finite automaton***, GNFA. First we show how to convert DFAs into GNFAs, and then GNFAs into regular expressions.

Generalized nondeterministic finite automata are simply nondeterministic finite automata wherein the transition arrows may have any regular expressions as labels, instead of only members of the alphabet or $\varepsilon$. The GNFA reads blocks of symbols from the input, not necessarily just one symbol at a time as in an ordinary NFA. The GNFA moves along a transition arrow connecting two states by reading a block of symbols from the input, which themselves constitute a string described by the regular expression on that arrow. A GNFA is nondeterministic and so may have several different ways to process the same input string. It accepts its input if its processing can cause the GNFA to be in an accept state at the end of the input. The following figure presents an example of a GNFA.



FIGURE   **1.61**
A generalized nondeterministic finite automaton

For convenience, we require that GNFAs always have a special form that meets the following conditions.

- The start state has transition arrows going to every other state but no arrows coming in from any other state.

- There is only a single accept state, and it has arrows coming in from every other state but no arrows going to any other state. Furthermore, the accept state is not the same as the start state.

- Except for the start and accept states, one arrow goes from every state to every other state and also from each state to itself.

We can easily convert a DFA into a GNFA in the special form. We simply add a new start state with an $\varepsilon$ arrow to the old start state and a new accept state with $\varepsilon$ arrows from the old accept states. If any arrows have multiple labels (or if there are multiple arrows going between the same two states in the same direction), we replace each with a single arrow whose label is the union of the previous labels. Finally, we add arrows labeled $\emptyset$ between states that had no arrows. This last step won't change the language recognized because a transition labeled with $\emptyset$ can never be used. From here on we assume that all GNFAs are in the special form.

Now we show how to convert a GNFA into a regular expression. Say that the GNFA has $k$ states. Then, because a GNFA must have a start and an accept state and they must be different from each other, we know that $k \geq 2$. If $k > 2$, we construct an equivalent GNFA with $k - 1$ states. This step can be repeated on the new GNFA until it is reduced to two states. If $k = 2$, the GNFA has a single arrow that goes from the start state to the accept state. The label of this arrow is the equivalent regular expression. For example, the stages in converting a DFA with three states to an equivalent regular expression are shown in the following figure.
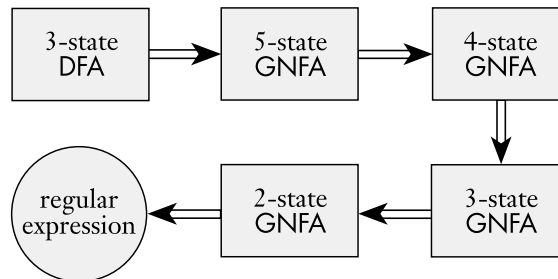


FIGURE   **1.62**
Typical stages in converting a DFA to a regular expression

The crucial step is constructing an equivalent GNFA with one fewer state when $k > 2$. We do so by selecting a state, ripping it out of the machine, and repairing the remainder so that the same language is still recognized. Any state will do, provided that it is not the start or accept state. We are guaranteed that such a state will exist because $k > 2$. Let's call the removed state $q_{rip}$.

After removing $q_{rip}$ we repair the machine by altering the regular expressions that label each of the remaining arrows. The new labels compensate for the absence of $q_{rip}$ by adding back the lost computations. The new label going from a state $q_i$ to a state $q_j$ is a regular expression that describes all strings that would

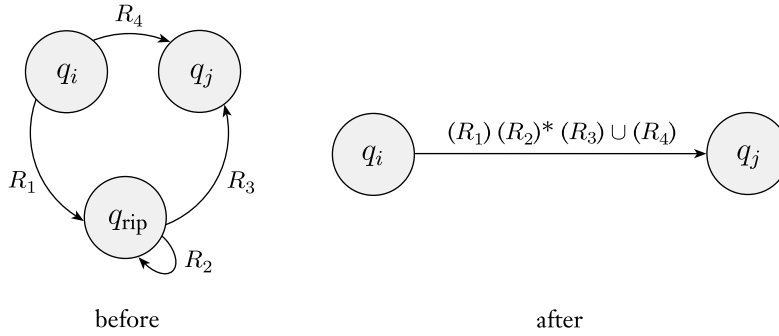take the machine from $q_i$ to $q_j$ either directly or via $q_{\text{rip}}$. We illustrate this approach in Figure 1.63.



before                                                    after

FIGURE  **1.63**
Constructing an equivalent GNFA with one fewer state

In the old machine, if

  **1.** $q_i$ goes to $q_{\text{rip}}$ with an arrow labeled $R_1$,
  **2.** $q_{\text{rip}}$ goes to itself with an arrow labeled $R_2$,
  **3.** $q_{\text{rip}}$ goes to $q_j$ with an arrow labeled $R_3$, and
  **4.** $q_i$ goes to $q_j$ with an arrow labeled $R_4$,

then in the new machine, the arrow from $q_i$ to $q_j$ gets the label

$$(R_1)(R_2)^*(R_3) \cup (R_4).$$

We make this change for each arrow going from any state $q_i$ to any state $q_j$, including the case where $q_i = q_j$. The new machine recognizes the original language.

**PROOF**   Let's now carry out this idea formally. First, to facilitate the proof, we formally define the new type of automaton introduced. A GNFA is similar to a nondeterministic finite automaton except for the transition function, which has the form

$$\delta \colon \big(Q - \{q_{\text{accept}}\}\big) \times \big(Q - \{q_{\text{start}}\}\big) \longrightarrow \mathcal{R}.$$

The symbol $\mathcal{R}$ is the collection of all regular expressions over the alphabet $\Sigma$, and $q_{\text{start}}$ and $q_{\text{accept}}$ are the start and accept states. If $\delta(q_i, q_j) = R$, the arrow from state $q_i$ to state $q_j$ has the regular expression $R$ as its label. The domain of the transition function is $\big(Q - \{q_{\text{accept}}\}\big) \times \big(Q - \{q_{\text{start}}\}\big)$ because an arrow connects every state to every other state, except that no arrows are coming from $q_{\text{accept}}$ or going to $q_{\text{start}}$.

---

**DEFINITION   1.64**

A ***generalized nondeterministic finite automaton*** is a 5-tuple, $(Q, \Sigma, \delta, q_{\text{start}}, q_{\text{accept}})$, where

1. $Q$ is the finite set of states,
2. $\Sigma$ is the input alphabet,
3. $\delta\colon (Q - \{q_{\text{accept}}\}) \times (Q - \{q_{\text{start}}\}) \longrightarrow \mathcal{R}$ is the transition function,
4. $q_{\text{start}}$ is the start state, and
5. $q_{\text{accept}}$ is the accept state.

---

A GNFA accepts a string $w$ in $\Sigma^*$ if $w = w_1 w_2 \cdots w_k$, where each $w_i$ is in $\Sigma^*$ and a sequence of states $q_0, q_1, \ldots, q_k$ exists such that

1. $q_0 = q_{\text{start}}$ is the start state,
2. $q_k = q_{\text{accept}}$ is the accept state, and
3. for each $i$, we have $w_i \in L(R_i)$, where $R_i = \delta(q_{i-1}, q_i)$; in other words, $R_i$ is the expression on the arrow from $q_{i-1}$ to $q_i$.

Returning to the proof of Lemma 1.60, we let $M$ be the DFA for language $A$. Then we convert $M$ to a GNFA $G$ by adding a new start state and a new accept state and additional transition arrows as necessary. We use the procedure CONVERT($G$), which takes a GNFA and returns an equivalent regular expression. This procedure uses ***recursion***, which means that it calls itself. An infinite loop is avoided because the procedure calls itself only to process a GNFA that has one fewer state. The case where the GNFA has two states is handled without recursion.

CONVERT($G$):

1. Let $k$ be the number of states of $G$.
2. If $k = 2$, then $G$ must consist of a start state, an accept state, and a single arrow connecting them and labeled with a regular expression $R$.
   Return the expression $R$.
3. If $k > 2$, we select any state $q_{\text{rip}} \in Q$ different from $q_{\text{start}}$ and $q_{\text{accept}}$ and let $G'$ be the GNFA $(Q', \Sigma, \delta', q_{\text{start}}, q_{\text{accept}})$, where

$$Q' = Q - \{q_{\text{rip}}\},$$

   and for any $q_i \in Q' - \{q_{\text{accept}}\}$ and any $q_j \in Q' - \{q_{\text{start}}\}$, let

$$\delta'(q_i, q_j) = (R_1)(R_2)^*(R_3) \cup (R_4),$$

   for $R_1 = \delta(q_i, q_{\text{rip}})$, $R_2 = \delta(q_{\text{rip}}, q_{\text{rip}})$, $R_3 = \delta(q_{\text{rip}}, q_j)$, and $R_4 = \delta(q_i, q_j)$.
4. Compute CONVERT($G'$) and return this value.

# 1.4 ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪

## NONREGULAR LANGUAGES

To understand the power of finite automata, you must also understand their limitations. In this section, we show how to prove that certain languages cannot be recognized by any finite automaton.

Let's take the language $B = \{0^n 1^n | \; n \geq 0\}$. If we attempt to find a DFA that recognizes $B$, we discover that the machine seems to need to remember how many 0s have been seen so far as it reads the input. Because the number of 0s isn't limited, the machine will have to keep track of an unlimited number of possibilities. But it cannot do so with any finite number of states.

Next, we present a method for proving that languages such as $B$ are not regular. Doesn't the argument already given prove nonregularity because the number of 0s is unlimited? It does not. Just because the language appears to require unbounded memory doesn't mean that it is necessarily so. It does happen to be true for the language $B$; but other languages seem to require an unlimited number of possibilities, yet actually they are regular. For example, consider two languages over the alphabet $\Sigma = \{0,1\}$:

$C = \{w | \; w$ has an equal number of 0s and 1s$\}$, and

$D = \{w | \; w$ has an equal number of occurrences of 01 and 10 as substrings$\}$.

At first glance, a recognizing machine appears to need to count in each case, and therefore neither language appears to be regular. As expected, $C$ is not regular, but surprisingly $D$ is regular![6] Thus our intuition can sometimes lead us astray, which is why we need mathematical proofs for certainty. In this section, we show how to prove that certain languages are not regular.

## THE PUMPING LEMMA FOR REGULAR LANGUAGES

Our technique for proving nonregularity stems from a theorem about regular languages, traditionally called the ***pumping lemma***. This theorem states that all regular languages have a special property. If we can show that a language does not have this property, we are guaranteed that it is not regular. The property states that all strings in the language can be "pumped" if they are at least as long as a certain special value, called the ***pumping length***. That means each such string contains a section that can be repeated any number of times with the resulting string remaining in the language.

---

[6]See Problem 1.53.

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

**1.** for each $i \geq 0$, $xy^i z \in A$,

**2.** $|y| > 0$, and

**3.** $|xy| \leq p$.

Recall the notation where $|s|$ represents the length of string $s$, $y^i$ means that $i$ copies of $y$ are concatenated together, and $y^0$ equals $\varepsilon$.
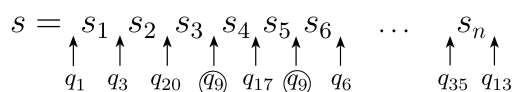
When $s$ is divided into $xyz$, either $x$ or $z$ may be $\varepsilon$, but condition 2 says that $y \neq \varepsilon$. Observe that without condition 2 the theorem would be trivially true. Condition 3 states that the pieces $x$ and $y$ together have length at most $p$. It is an extra technical condition that we occasionally find useful when proving certain languages to be nonregular. See Example 1.74 for an application of condition 3.

**PROOF IDEA** Let $M = (Q, \Sigma, \delta, q_1, F)$ be a DFA that recognizes $A$. We assign the pumping length $p$ to be the number of states of $M$. We show that any string $s$ in $A$ of length at least $p$ may be broken into the three pieces $xyz$, satisfying our three conditions. What if no strings in $A$ are of length at least $p$? Then our task is even easier because the theorem becomes *vacuously* true: Obviously the three conditions hold for all strings of length at least $p$ if there aren't any such strings.

If $s$ in $A$ has length at least $p$, consider the sequence of states that $M$ goes through when computing with input $s$. It starts with $q_1$ the start state, then goes to, say, $q_3$, then, say, $q_{20}$, then $q_9$, and so on, until it reaches the end of $s$ in state $q_{13}$. With $s$ in $A$, we know that $M$ accepts $s$, so $q_{13}$ is an accept state.

If we let $n$ be the length of $s$, the sequence of states $q_1, q_3, q_{20}, q_9, \ldots, q_{13}$ has length $n + 1$. Because $n$ is at least $p$, we know that $n + 1$ is greater than $p$, the number of states of $M$. Therefore, the sequence must contain a repeated state. This result is an example of the ***pigeonhole principle***, a fancy name for the rather obvious fact that if $p$ pigeons are placed into fewer than $p$ holes, some hole has to have more than one pigeon in it.

The following figure shows the string $s$ and the sequence of states that $M$ goes through when processing $s$. State $q_9$ is the one that repeats.



$$s = s_1\ s_2\ s_3\ s_4\ s_5\ s_6\quad \cdots \quad s_n$$
$$q_1\quad q_3\quad q_{20}\quad \textcircled{$q_9$}\quad q_{17}\quad \textcircled{$q_9$}\quad q_6 \qquad q_{35}\quad q_{13}$$

FIGURE **1.71**
Example showing state $q_9$ repeating when $M$ reads $s$

We now divide $s$ into the three pieces $x$, $y$, and $z$. Piece $x$ is the part of $s$ appearing before $q_9$, piece $y$ is the part between the two appearances of $q_9$, and

To use the pumping lemma to prove that a language $B$ is not regular, first assume that $B$ is regular in order to obtain a contradiction. Then use the pumping lemma to guarantee the existence of a pumping length $p$ such that all strings of length $p$ or greater in $B$ can be pumped. Next, find a string $s$ in $B$ that has length $p$ or greater but that cannot be pumped. Finally, demonstrate that $s$ cannot be pumped by considering all ways of dividing $s$ into $x$, $y$, and $z$ (taking condition 3 of the pumping lemma into account if convenient) and, for each such division, finding a value $i$ where $xy^i z \notin B$. This final step often involves grouping the various ways of dividing $s$ into several cases and analyzing them individually. The existence of $s$ contradicts the pumping lemma if $B$ were regular. Hence $B$ cannot be regular.

Finding $s$ sometimes takes a bit of creative thinking. You may need to hunt through several candidates for $s$ before you discover one that works. Try members of $B$ that seem to exhibit the "essence" of $B$'s nonregularity. We further discuss the task of finding $s$ in some of the following examples.

EXAMPLE **1.73** ................................................................

Let $B$ be the language $\{0^n 1^n | n \geq 0\}$. We use the pumping lemma to prove that $B$ is not regular. The proof is by contradiction.

Assume to the contrary that $B$ is regular. Let $p$ be the pumping length given by the pumping lemma. Choose $s$ to be the string $0^p 1^p$. Because $s$ is a member of $B$ and $s$ has length more than $p$, the pumping lemma guarantees that $s$ can be split into three pieces, $s = xyz$, where for any $i \geq 0$ the string $xy^i z$ is in $B$. We consider three cases to show that this result is impossible.

1. The string $y$ consists only of 0s. In this case, the string $xyyz$ has more 0s than 1s and so is not a member of $B$, violating condition 1 of the pumping lemma. This case is a contradiction.

2. The string $y$ consists only of 1s. This case also gives a contradiction.

3. The string $y$ consists of both 0s and 1s. In this case, the string $xyyz$ may have the same number of 0s and 1s, but they will be out of order with some 1s before 0s. Hence it is not a member of $B$, which is a contradiction.

Thus a contradiction is unavoidable if we make the assumption that $B$ is regular, so $B$ is not regular. Note that we can simplify this argument by applying condition 3 of the pumping lemma to eliminate cases 2 and 3.

In this example, finding the string $s$ was easy because any string in $B$ of length $p$ or more would work. In the next two examples, some choices for $s$ do not work so additional care is required.

EXAMPLE **1.74** ................................................................

Let $C = \{w | \ w$ has an equal number of 0s and 1s$\}$. We use the pumping lemma to prove that $C$ is not regular. The proof is by contradiction.