long—to resist possible advances in the art of factoring. At the time of this writing (2009), RSA moduli were commonly in the range of 768 to 2048 bits. To create moduli of such sizes, we must be able to find large primes efficiently. Section 31.8 addresses this problem.

For efficiency, RSA is often used in a "hybrid" or "key-management" mode with fast non-public-key cryptosystems. With such a system, the encryption and decryption keys are identical. If Alice wishes to send a long message $M$ to Bob privately, she selects a random key $K$ for the fast non-public-key cryptosystem and encrypts $M$ using $K$, obtaining ciphertext $C$. Here, $C$ is as long as $M$, but $K$ is quite short. Then, she encrypts $K$ using Bob's public RSA key. Since $K$ is short, computing $P_B(K)$ is fast (much faster than computing $P_B(M)$). She then transmits $(C, P_B(K))$ to Bob, who decrypts $P_B(K)$ to obtain $K$ and then uses $K$ to decrypt $C$, obtaining $M$.

We can use a similar hybrid approach to make digital signatures efficiently. This approach combines RSA with a public *collision-resistant hash function $h$*—a function that is easy to compute but for which it is computationally infeasible to find two messages $M$ and $M'$ such that $h(M) = h(M')$. The value $h(M)$ is a short (say, 256-bit) "fingerprint" of the message $M$. If Alice wishes to sign a message $M$, she first applies $h$ to $M$ to obtain the fingerprint $h(M)$, which she then encrypts with her secret key. She sends $(M, S_A(h(M)))$ to Bob as her signed version of $M$. Bob can verify the signature by computing $h(M)$ and verifying that $P_A$ applied to $S_A(h(M))$ as received equals $h(M)$. Because no one can create two messages with the same fingerprint, it is computationally infeasible to alter a signed message and preserve the validity of the signature.

Finally, we note that the use of *certificates* makes distributing public keys much easier. For example, assume there is a "trusted authority" $T$ whose public key is known by everyone. Alice can obtain from $T$ a signed message (her certificate) stating that "Alice's public key is $P_A$." This certificate is "self-authenticating" since everyone knows $P_T$. Alice can include her certificate with her signed messages, so that the recipient has Alice's public key immediately available in order to verify her signature. Because her key was signed by $T$, the recipient knows that Alice's key is really Alice's.

### Exercises

***31.7-1***
Consider an RSA key set with $p = 11$, $q = 29$, $n = 319$, and $e = 3$. What value of $d$ should be used in the secret key? What is the encryption of the message $M = 100$?

**31.7-2**
Prove that if Alice's public exponent $e$ is 3 and an adversary obtains Alice's secret exponent $d$, where $0 < d < \phi(n)$, then the adversary can factor Alice's modulus $n$ in time polynomial in the number of bits in $n$. (Although you are not asked to prove it, you may be interested to know that this result remains true even if the condition $e = 3$ is removed. See Miller [255].)

**31.7-3** ★
Prove that RSA is multiplicative in the sense that

$$P_A(M_1)P_A(M_2) \equiv P_A(M_1 M_2) \pmod{n} .$$

Use this fact to prove that if an adversary had a procedure that could efficiently decrypt 1 percent of messages from $\mathbb{Z}_n$ encrypted with $P_A$, then he could employ a probabilistic algorithm to decrypt every message encrypted with $P_A$ with high probability.

## ★ 31.8 Primality testing

In this section, we consider the problem of finding large primes. We begin with a discussion of the density of primes, proceed to examine a plausible, but incomplete, approach to primality testing, and then present an effective randomized primality test due to Miller and Rabin.

### The density of prime numbers

For many applications, such as cryptography, we need to find large "random" primes. Fortunately, large primes are not too rare, so that it is feasible to test random integers of the appropriate size until we find a prime. The ***prime distribution function*** $\pi(n)$ specifies the number of primes that are less than or equal to $n$. For example, $\pi(10) = 4$, since there are 4 prime numbers less than or equal to 10, namely, 2, 3, 5, and 7. The prime number theorem gives a useful approximation to $\pi(n)$.

***Theorem 31.37 (Prime number theorem)***
$$\lim_{n\to\infty} \frac{\pi(n)}{n/\ln n} = 1 .$$
■

The approximation $n/\ln n$ gives reasonably accurate estimates of $\pi(n)$ even for small $n$. For example, it is off by less than 6% at $n = 10^9$, where $\pi(n) =$

50,847,534 and $n / \ln n \approx 48{,}254{,}942$. (To a number theorist, $10^9$ is a small number.)

We can view the process of randomly selecting an integer $n$ and determining whether it is prime as a Bernoulli trial (see Section C.4). By the prime number theorem, the probability of a success—that is, the probability that $n$ is prime—is approximately $1 / \ln n$. The geometric distribution tells us how many trials we need to obtain a success, and by equation (C.32), the expected number of trials is approximately $\ln n$. Thus, we would expect to examine approximately $\ln n$ integers chosen randomly near $n$ in order to find a prime that is of the same length as $n$. For example, we expect that finding a 1024-bit prime would require testing approximately $\ln 2^{1024} \approx 710$ randomly chosen 1024-bit numbers for primality. (Of course, we can cut this figure in half by choosing only odd integers.)

In the remainder of this section, we consider the problem of determining whether or not a large odd integer $n$ is prime. For notational convenience, we assume that $n$ has the prime factorization

$$n = p_1^{e_1} p_2^{e_2} \cdots p_r^{e_r} \ , \tag{31.39}$$

where $r \geq 1$, $p_1, p_2, \ldots, p_r$ are the prime factors of $n$, and $e_1, e_2, \ldots, e_r$ are positive integers. The integer $n$ is prime if and only if $r = 1$ and $e_1 = 1$.

One simple approach to the problem of testing for primality is **_trial division_**. We try dividing $n$ by each integer $2, 3, \ldots, \lfloor \sqrt{n} \rfloor$. (Again, we may skip even integers greater than 2.) It is easy to see that $n$ is prime if and only if none of the trial divisors divides $n$. Assuming that each trial division takes constant time, the worst-case running time is $\Theta(\sqrt{n})$, which is exponential in the length of $n$. (Recall that if $n$ is encoded in binary using $\beta$ bits, then $\beta = \lceil \lg(n + 1) \rceil$, and so $\sqrt{n} = \Theta(2^{\beta/2})$.) Thus, trial division works well only if $n$ is very small or happens to have a small prime factor. When it works, trial division has the advantage that it not only determines whether $n$ is prime or composite, but also determines one of $n$'s prime factors if $n$ is composite.

In this section, we are interested only in finding out whether a given number $n$ is prime; if $n$ is composite, we are not concerned with finding its prime factorization. As we shall see in Section 31.9, computing the prime factorization of a number is computationally expensive. It is perhaps surprising that it is much easier to tell whether or not a given number is prime than it is to determine the prime factorization of the number if it is not prime.

### Pseudoprimality testing

We now consider a method for primality testing that "almost works" and in fact is good enough for many practical applications. Later on, we shall present a re-

finement of this method that removes the small defect. Let $\mathbb{Z}_n^+$ denote the nonzero elements of $\mathbb{Z}_n$:

$$\mathbb{Z}_n^+ = \{1, 2, \ldots, n-1\} \ .$$

If $n$ is prime, then $\mathbb{Z}_n^+ = \mathbb{Z}_n^*$.

We say that $n$ is a ***base-a pseudoprime*** if $n$ is composite and

$$a^{n-1} \equiv 1 \pmod{n} \ . \tag{31.40}$$

Fermat's theorem (Theorem 31.31) implies that if $n$ is prime, then $n$ satisfies equation (31.40) for every $a$ in $\mathbb{Z}_n^+$. Thus, if we can find any $a \in \mathbb{Z}_n^+$ such that $n$ does *not* satisfy equation (31.40), then $n$ is certainly composite. Surprisingly, the converse *almost* holds, so that this criterion forms an almost perfect test for primality. We test to see whether $n$ satisfies equation (31.40) for $a = 2$. If not, we declare $n$ to be composite by returning COMPOSITE. Otherwise, we return PRIME, guessing that $n$ is prime (when, in fact, all we know is that $n$ is either prime or a base-2 pseudoprime).

The following procedure pretends in this manner to be checking the primality of $n$. It uses the procedure MODULAR-EXPONENTIATION from Section 31.6. We assume that the input $n$ is an odd integer greater than 2.

PSEUDOPRIME$(n)$

```
1  if MODULAR-EXPONENTIATION(2, n − 1, n) ≢ 1  (mod n)
2       return COMPOSITE        // definitely
3  else return PRIME            // we hope!
```

This procedure can make errors, but only of one type. That is, if it says that $n$ is composite, then it is always correct. If it says that $n$ is prime, however, then it makes an error only if $n$ is a base-2 pseudoprime.

How often does this procedure err? Surprisingly rarely. There are only 22 values of $n$ less than 10,000 for which it errs; the first four such values are 341, 561, 645, and 1105. We won't prove it, but the probability that this program makes an error on a randomly chosen $\beta$-bit number goes to zero as $\beta \to \infty$. Using more precise estimates due to Pomerance [279] of the number of base-2 pseudoprimes of a given size, we may estimate that a randomly chosen 512-bit number that is called prime by the above procedure has less than one chance in $10^{20}$ of being a base-2 pseudoprime, and a randomly chosen 1024-bit number that is called prime has less than one chance in $10^{41}$ of being a base-2 pseudoprime. So if you are merely trying to find a large prime for some application, for all practical purposes you almost never go wrong by choosing large numbers at random until one of them causes PSEUDOPRIME to return PRIME. But when the numbers being tested for primality are not randomly chosen, we need a better approach for testing primality.

As we shall see, a little more cleverness, and some randomization, will yield a primality-testing routine that works well on all inputs.

Unfortunately, we cannot entirely eliminate all the errors by simply checking equation (31.40) for a second base number, say $a = 3$, because there exist composite integers $n$, known as ***Carmichael numbers***, that satisfy equation (31.40) for *all* $a \in \mathbb{Z}_n^*$. (We note that equation (31.40) does fail when $\gcd(a, n) > 1$—that is, when $a \notin \mathbb{Z}_n^*$—but hoping to demonstrate that $n$ is composite by finding such an $a$ can be difficult if $n$ has only large prime factors.) The first three Carmichael numbers are 561, 1105, and 1729. Carmichael numbers are extremely rare; there are, for example, only 255 of them less than 100,000,000. Exercise 31.8-2 helps explain why they are so rare.

We next show how to improve our primality test so that it won't be fooled by Carmichael numbers.

### The Miller-Rabin randomized primality test

The Miller-Rabin primality test overcomes the problems of the simple test PSEU-DOPRIME with two modifications:

- It tries several randomly chosen base values $a$ instead of just one base value.

- While computing each modular exponentiation, it looks for a nontrivial square root of 1, modulo $n$, during the final set of squarings. If it finds one, it stops and returns COMPOSITE. Corollary 31.35 from Section 31.6 justifies detecting composites in this manner.

The pseudocode for the Miller-Rabin primality test follows. The input $n > 2$ is the odd number to be tested for primality, and $s$ is the number of randomly chosen base values from $\mathbb{Z}_n^+$ to be tried. The code uses the random-number generator RANDOM described on page 117: RANDOM$(1, n - 1)$ returns a randomly chosen integer $a$ satisfying $1 \le a \le n-1$. The code uses an auxiliary procedure WITNESS such that WITNESS$(a, n)$ is TRUE if and only if $a$ is a "witness" to the composite-ness of $n$—that is, if it is possible using $a$ to prove (in a manner that we shall see) that $n$ is composite. The test WITNESS$(a, n)$ is an extension of, but more effective than, the test

$$a^{n-1} \not\equiv 1 \pmod{n}$$

that formed the basis (using $a = 2$) for PSEUDOPRIME. We first present and justify the construction of WITNESS, and then we shall show how we use it in the Miller-Rabin primality test. Let $n - 1 = 2^t u$ where $t \ge 1$ and $u$ is odd; i.e., the binary representation of $n - 1$ is the binary representation of the odd integer $u$ followed by exactly $t$ zeros. Therefore, $a^{n-1} \equiv (a^u)^{2^t} \pmod{n}$, so that we can

compute $a^{n-1} \bmod n$ by first computing $a^u \bmod n$ and then squaring the result $t$ times successively.

WITNESS$(a, n)$

1  let $t$ and $u$ be such that $t \geq 1$, $u$ is odd, and $n - 1 = 2^t u$
2  $x_0$ = MODULAR-EXPONENTIATION$(a, u, n)$
3  **for** $i = 1$ **to** $t$
4      $x_i = x_{i-1}^2 \bmod n$
5      **if** $x_i == 1$ and $x_{i-1} \neq 1$ and $x_{i-1} \neq n - 1$
6          **return** TRUE
7  **if** $x_t \neq 1$
8      **return** TRUE
9  **return** FALSE

This pseudocode for WITNESS computes $a^{n-1} \bmod n$ by first computing the value $x_0 = a^u \bmod n$ in line 2 and then squaring the result $t$ times in a row in the **for** loop of lines 3–6. By induction on $i$, the sequence $x_0, x_1, \ldots, x_t$ of values computed satisfies the equation $x_i \equiv a^{2^i u} \pmod{n}$ for $i = 0, 1, \ldots, t$, so that in particular $x_t \equiv a^{n-1} \pmod{n}$. After line 4 performs a squaring step, however, the loop may terminate early if lines 5–6 detect that a nontrivial square root of 1 has just been discovered. (We shall explain these tests shortly.) If so, the algorithm stops and returns TRUE. Lines 7–8 return TRUE if the value computed for $x_t \equiv a^{n-1} \pmod{n}$ is not equal to 1, just as the PSEUDOPRIME procedure returns COMPOSITE in this case. Line 9 returns FALSE if we haven't returned TRUE in lines 6 or 8.

We now argue that if WITNESS$(a, n)$ returns TRUE, then we can construct a proof that $n$ is composite using $a$ as a witness.

If WITNESS returns TRUE from line 8, then it has discovered that $x_t = a^{n-1} \bmod n \neq 1$. If $n$ is prime, however, we have by Fermat's theorem (Theorem 31.31) that $a^{n-1} \equiv 1 \pmod{n}$ for all $a \in \mathbb{Z}_n^+$. Therefore, $n$ cannot be prime, and the equation $a^{n-1} \bmod n \neq 1$ proves this fact.

If WITNESS returns TRUE from line 6, then it has discovered that $x_{i-1}$ is a nontrivial square root of 1, modulo $n$, since we have that $x_{i-1} \not\equiv \pm 1 \pmod{n}$ yet $x_i \equiv x_{i-1}^2 \equiv 1 \pmod{n}$. Corollary 31.35 states that only if $n$ is composite can there exist a nontrivial square root of 1 modulo $n$, so that demonstrating that $x_{i-1}$ is a nontrivial square root of 1 modulo $n$ proves that $n$ is composite.

This completes our proof of the correctness of WITNESS. If we find that the call WITNESS$(a, n)$ returns TRUE, then $n$ is surely composite, and the witness $a$, along with the reason that the procedure returns TRUE (did it return from line 6 or from line 8?), provides a proof that $n$ is composite.

At this point, we briefly present an alternative description of the behavior of WITNESS as a function of the sequence $X = \langle x_0, x_1, \ldots, x_t \rangle$, which we shall find useful later on, when we analyze the efficiency of the Miller-Rabin primality test. Note that if $x_i = 1$ for some $0 \leq i < t$, WITNESS might not compute the rest of the sequence. If it were to do so, however, each value $x_{i+1}, x_{i+2}, \ldots, x_t$ would be 1, and we consider these positions in the sequence $X$ as being all 1s. We have four cases:

1. $X = \langle \ldots, d \rangle$, where $d \neq 1$: the sequence $X$ does not end in 1. Return TRUE in line 8; $a$ is a witness to the compositeness of $n$ (by Fermat's Theorem).

2. $X = \langle 1, 1, \ldots, 1 \rangle$: the sequence $X$ is all 1s. Return FALSE; $a$ is not a witness to the compositeness of $n$.

3. $X = \langle \ldots, -1, 1, \ldots, 1 \rangle$: the sequence $X$ ends in 1, and the last non-1 is equal to $-1$. Return FALSE; $a$ is not a witness to the compositeness of $n$.

4. $X = \langle \ldots, d, 1, \ldots, 1 \rangle$, where $d \neq \pm 1$: the sequence $X$ ends in 1, but the last non-1 is not $-1$. Return TRUE in line 6; $a$ is a witness to the compositeness of $n$, since $d$ is a nontrivial square root of 1.

We now examine the Miller-Rabin primality test based on the use of WITNESS. Again, we assume that $n$ is an odd integer greater than 2.

MILLER-RABIN$(n, s)$

```
1  for j = 1 to s
2      a = RANDOM(1, n − 1)
3      if WITNESS(a, n)
4          return COMPOSITE          // definitely
5  return PRIME                       // almost surely
```

The procedure MILLER-RABIN is a probabilistic search for a proof that $n$ is composite. The main loop (beginning on line 1) picks up to $s$ random values of $a$ from $\mathbb{Z}_n^+$ (line 2). If one of the $a$'s picked is a witness to the compositeness of $n$, then MILLER-RABIN returns COMPOSITE on line 4. Such a result is always correct, by the correctness of WITNESS. If MILLER-RABIN finds no witness in $s$ trials, then the procedure assumes that this is because no witnesses exist, and therefore it assumes that $n$ is prime. We shall see that this result is likely to be correct if $s$ is large enough, but that there is still a tiny chance that the procedure may be unlucky in its choice of $a$'s and that witnesses do exist even though none has been found.

To illustrate the operation of MILLER-RABIN, let $n$ be the Carmichael number 561, so that $n - 1 = 560 = 2^4 \cdot 35$, $t = 4$, and $u = 35$. If the procedure chooses $a = 7$ as a base, Figure 31.4 in Section 31.6 shows that WITNESS computes $x_0 \equiv a^{35} \equiv 241 \pmod{561}$ and thus computes the sequence

$X = \langle 241, 298, 166, 67, 1 \rangle$. Thus, WITNESS discovers a nontrivial square root of 1 in the last squaring step, since $a^{280} \equiv 67 \pmod{n}$ and $a^{560} \equiv 1 \pmod{n}$. Therefore, $a = 7$ is a witness to the compositeness of $n$, WITNESS$(7, n)$ returns TRUE, and MILLER-RABIN returns COMPOSITE.

If $n$ is a $\beta$-bit number, MILLER-RABIN requires $O(s\beta)$ arithmetic operations and $O(s\beta^3)$ bit operations, since it requires asymptotically no more work than $s$ modular exponentiations.

### Error rate of the Miller-Rabin primality test

If MILLER-RABIN returns PRIME, then there is a very slim chance that it has made an error. Unlike PSEUDOPRIME, however, the chance of error does not depend on $n$; there are no bad inputs for this procedure. Rather, it depends on the size of $s$ and the "luck of the draw" in choosing base values $a$. Moreover, since each test is more stringent than a simple check of equation (31.40), we can expect on general principles that the error rate should be small for randomly chosen integers $n$. The following theorem presents a more precise argument.

***Theorem 31.38***
If $n$ is an odd composite number, then the number of witnesses to the compositeness of $n$ is at least $(n-1)/2$.

***Proof***   The proof shows that the number of nonwitnesses is at most $(n-1)/2$, which implies the theorem.

We start by claiming that any nonwitness must be a member of $\mathbb{Z}_n^*$. Why? Consider any nonwitness $a$. It must satisfy $a^{n-1} \equiv 1 \pmod{n}$ or, equivalently, $a \cdot a^{n-2} \equiv 1 \pmod{n}$. Thus, the equation $ax \equiv 1 \pmod{n}$ has a solution, namely $a^{n-2}$. By Corollary 31.21, $\gcd(a, n) \mid 1$, which in turn implies that $\gcd(a, n) = 1$. Therefore, $a$ is a member of $\mathbb{Z}_n^*$; all nonwitnesses belong to $\mathbb{Z}_n^*$.

To complete the proof, we show that not only are all nonwitnesses contained in $\mathbb{Z}_n^*$, they are all contained in a proper subgroup $B$ of $\mathbb{Z}_n^*$ (recall that we say $B$ is a *proper* subgroup of $\mathbb{Z}_n^*$ when $B$ is subgroup of $\mathbb{Z}_n^*$ but $B$ is not equal to $\mathbb{Z}_n^*$). By Corollary 31.16, we then have $|B| \leq |\mathbb{Z}_n^*|/2$. Since $|\mathbb{Z}_n^*| \leq n-1$, we obtain $|B| \leq (n-1)/2$. Therefore, the number of nonwitnesses is at most $(n-1)/2$, so that the number of witnesses must be at least $(n-1)/2$.

We now show how to find a proper subgroup $B$ of $\mathbb{Z}_n^*$ containing all of the nonwitnesses. We break the proof into two cases.

*Case 1:* There exists an $x \in \mathbb{Z}_n^*$ such that

$$x^{n-1} \not\equiv 1 \pmod{n} .$$

In other words, $n$ is not a Carmichael number. Because, as we noted earlier, Carmichael numbers are extremely rare, case 1 is the main case that arises "in practice" (e.g., when $n$ has been chosen randomly and is being tested for primality).

Let $B = \{b \in \mathbb{Z}_n^* : b^{n-1} \equiv 1 \pmod{n}\}$. Clearly, $B$ is nonempty, since $1 \in B$. Since $B$ is closed under multiplication modulo $n$, we have that $B$ is a subgroup of $\mathbb{Z}_n^*$ by Theorem 31.14. Note that every nonwitness belongs to $B$, since a nonwitness $a$ satisfies $a^{n-1} \equiv 1 \pmod{n}$. Since $x \in \mathbb{Z}_n^* - B$, we have that $B$ is a proper subgroup of $\mathbb{Z}_n^*$.

*Case 2:* For all $x \in \mathbb{Z}_n^*$,

$$x^{n-1} \equiv 1 \pmod{n} . \tag{31.41}$$

In other words, $n$ is a Carmichael number. This case is extremely rare in practice. However, the Miller-Rabin test (unlike a pseudo-primality test) can efficiently determine that Carmichael numbers are composite, as we now show.

In this case, $n$ cannot be a prime power. To see why, let us suppose to the contrary that $n = p^e$, where $p$ is a prime and $e > 1$. We derive a contradiction as follows. Since we assume that $n$ is odd, $p$ must also be odd. Theorem 31.32 implies that $\mathbb{Z}_n^*$ is a cyclic group: it contains a generator $g$ such that $\operatorname{ord}_n(g) = |\mathbb{Z}_n^*| = \phi(n) = p^e(1 - 1/p) = (p-1)p^{e-1}$. (The formula for $\phi(n)$ comes from equation (31.20).) By equation (31.41), we have $g^{n-1} \equiv 1 \pmod{n}$. Then the discrete logarithm theorem (Theorem 31.33, taking $y = 0$) implies that $n - 1 \equiv 0 \pmod{\phi(n)}$, or

$$(p-1)p^{e-1} \mid p^e - 1 .$$

This is a contradiction for $e > 1$, since $(p-1)p^{e-1}$ is divisible by the prime $p$ but $p^e - 1$ is not. Thus, $n$ is not a prime power.

Since the odd composite number $n$ is not a prime power, we decompose it into a product $n_1 n_2$, where $n_1$ and $n_2$ are odd numbers greater than 1 that are relatively prime to each other. (There may be several ways to decompose $n$, and it does not matter which one we choose. For example, if $n = p_1^{e_1} p_2^{e_2} \cdots p_r^{e_r}$, then we can choose $n_1 = p_1^{e_1}$ and $n_2 = p_2^{e_2} p_3^{e_3} \cdots p_r^{e_r}$.)

Recall that we define $t$ and $u$ so that $n - 1 = 2^t u$, where $t \geq 1$ and $u$ is odd, and that for an input $a$, the procedure WITNESS computes the sequence

$$X = \langle a^u, a^{2u}, a^{2^2 u}, \ldots, a^{2^t u} \rangle$$

(all computations are performed modulo $n$).

Let us call a pair $(v, j)$ of integers *acceptable* if $v \in \mathbb{Z}_n^*$, $j \in \{0, 1, \ldots, t\}$, and

$$v^{2^j u} \equiv -1 \pmod{n} .$$

Acceptable pairs certainly exist since $u$ is odd; we can choose $v = n - 1$ and $j = 0$, so that $(n-1, 0)$ is an acceptable pair. Now pick the largest possible $j$ such that there exists an acceptable pair $(v, j)$, and fix $v$ so that $(v, j)$ is an acceptable pair. Let

$$B = \{x \in \mathbb{Z}_n^* : x^{2^j u} \equiv \pm 1 \pmod{n}\} \,.$$

Since $B$ is closed under multiplication modulo $n$, it is a subgroup of $\mathbb{Z}_n^*$. By Theorem 31.15, therefore, $|B|$ divides $|\mathbb{Z}_n^*|$. Every nonwitness must be a member of $B$, since the sequence $X$ produced by a nonwitness must either be all 1s or else contain a $-1$ no later than the $j$th position, by the maximality of $j$. (If $(a, j')$ is acceptable, where $a$ is a nonwitness, we must have $j' \leq j$ by how we chose $j$.)

We now use the existence of $v$ to demonstrate that there exists a $w \in \mathbb{Z}_n^* - B$, and hence that $B$ is a proper subgroup of $\mathbb{Z}_n^*$. Since $v^{2^j u} \equiv -1 \pmod{n}$, we have $v^{2^j u} \equiv -1 \pmod{n_1}$ by Corollary 31.29 to the Chinese remainder theorem. By Corollary 31.28, there exists a $w$ simultaneously satisfying the equations

$$w \equiv v \pmod{n_1} \,,$$
$$w \equiv 1 \pmod{n_2} \,.$$

Therefore,

$$w^{2^j u} \equiv -1 \pmod{n_1} \,,$$
$$w^{2^j u} \equiv 1 \pmod{n_2} \,.$$

By Corollary 31.29, $w^{2^j u} \not\equiv 1 \pmod{n_1}$ implies $w^{2^j u} \not\equiv 1 \pmod{n}$, and $w^{2^j u} \not\equiv -1 \pmod{n_2}$ implies $w^{2^j u} \not\equiv -1 \pmod{n}$. Hence, we conclude that $w^{2^j u} \not\equiv \pm 1 \pmod{n}$, and so $w \notin B$.

It remains to show that $w \in \mathbb{Z}_n^*$, which we do by first working separately modulo $n_1$ and modulo $n_2$. Working modulo $n_1$, we observe that since $v \in \mathbb{Z}_n^*$, we have that $\gcd(v, n) = 1$, and so also $\gcd(v, n_1) = 1$; if $v$ does not have any common divisors with $n$, then it certainly does not have any common divisors with $n_1$. Since $w \equiv v \pmod{n_1}$, we see that $\gcd(w, n_1) = 1$. Working modulo $n_2$, we observe that $w \equiv 1 \pmod{n_2}$ implies $\gcd(w, n_2) = 1$. To combine these results, we use Theorem 31.6, which implies that $\gcd(w, n_1 n_2) = \gcd(w, n) = 1$. That is, $w \in \mathbb{Z}_n^*$.

Therefore $w \in \mathbb{Z}_n^* - B$, and we finish case 2 with the conclusion that $B$ is a proper subgroup of $\mathbb{Z}_n^*$.

In either case, we see that the number of witnesses to the compositeness of $n$ is at least $(n - 1)/2$. ∎

### Theorem 31.39
For any odd integer $n > 2$ and positive integer $s$, the probability that MILLER-RABIN$(n, s)$ errs is at most $2^{-s}$.

***Proof***    Using Theorem 31.38, we see that if $n$ is composite, then each execution of the **for** loop of lines 1–4 has a probability of at least $1/2$ of discovering a witness $x$ to the compositeness of $n$. MILLER-RABIN makes an error only if it is so unlucky as to miss discovering a witness to the compositeness of $n$ on each of the $s$ iterations of the main loop. The probability of such a sequence of misses is at most $2^{-s}$.    ∎

If $n$ is prime, MILLER-RABIN always reports PRIME, and if $n$ is composite, the chance that MILLER-RABIN reports PRIME is at most $2^{-s}$.

When applying MILLER-RABIN to a large randomly chosen integer $n$, however, we need to consider as well the prior probability that $n$ is prime, in order to correctly interpret MILLER-RABIN's result. Suppose that we fix a bit length $\beta$ and choose at random an integer $n$ of length $\beta$ bits to be tested for primality. Let $A$ denote the event that $n$ is prime. By the prime number theorem (Theorem 31.37), the probability that $n$ is prime is approximately

$$\begin{aligned} \Pr\{A\} &\approx 1/\ln n \\ &\approx 1.443/\beta \ . \end{aligned}$$

Now let $B$ denote the event that MILLER-RABIN returns PRIME. We have that $\Pr\{\overline{B} \mid A\} = 0$ (or equivalently, that $\Pr\{B \mid A\} = 1$) and $\Pr\{B \mid \overline{A}\} \le 2^{-s}$ (or equivalently, that $\Pr\{\overline{B} \mid \overline{A}\} > 1 - 2^{-s}$).

But what is $\Pr\{A \mid B\}$, the probability that $n$ is prime, given that MILLER-RABIN has returned PRIME? By the alternate form of Bayes's theorem (equation (C.18)) we have

$$\begin{aligned} \Pr\{A \mid B\} &= \frac{\Pr\{A\}\Pr\{B \mid A\}}{\Pr\{A\}\Pr\{B \mid A\} + \Pr\{\overline{A}\}\Pr\{B \mid \overline{A}\}} \\ &\approx \frac{1}{1 + 2^{-s}(\ln n - 1)} \ . \end{aligned}$$

This probability does not exceed $1/2$ until $s$ exceeds $\lg(\ln n - 1)$. Intuitively, that many initial trials are needed just for the confidence derived from failing to find a witness to the compositeness of $n$ to overcome the prior bias in favor of $n$ being composite. For a number with $\beta = 1024$ bits, this initial testing requires about

$$\begin{aligned} \lg(\ln n - 1) &\approx \lg(\beta/1.443) \\ &\approx 9 \end{aligned}$$

trials. In any case, choosing $s = 50$ should suffice for almost any imaginable application.

In fact, the situation is much better. If we are trying to find large primes by applying MILLER-RABIN to large randomly chosen odd integers, then choosing a small value of $s$ (say 3) is very unlikely to lead to erroneous results, though

we won't prove it here. The reason is that for a randomly chosen odd composite integer $n$, the expected number of nonwitnesses to the compositeness of $n$ is likely to be very much smaller than $(n-1)/2$.

If the integer $n$ is not chosen randomly, however, the best that can be proven is that the number of nonwitnesses is at most $(n-1)/4$, using an improved version of Theorem 31.38. Furthermore, there do exist integers $n$ for which the number of nonwitnesses is $(n-1)/4$.

### Exercises

***31.8-1***
Prove that if an odd integer $n > 1$ is not a prime or a prime power, then there exists a nontrivial square root of 1 modulo $n$.

***31.8-2*** ★
It is possible to strengthen Euler's theorem slightly to the form

$$a^{\lambda(n)} \equiv 1 \pmod{n} \text{ for all } a \in \mathbb{Z}_n^* ,$$

where $n = p_1^{e_1} \cdots p_r^{e_r}$ and $\lambda(n)$ is defined by

$$\lambda(n) = \text{lcm}(\phi(p_1^{e_1}), \ldots, \phi(p_r^{e_r})) . \tag{31.42}$$

Prove that $\lambda(n) \mid \phi(n)$. A composite number $n$ is a Carmichael number if $\lambda(n) \mid n-1$. The smallest Carmichael number is $561 = 3 \cdot 11 \cdot 17$; here, $\lambda(n) = \text{lcm}(2, 10, 16) = 80$, which divides 560. Prove that Carmichael numbers must be both "square-free" (not divisible by the square of any prime) and the product of at least three primes. (For this reason, they are not very common.)

***31.8-3***
Prove that if $x$ is a nontrivial square root of 1, modulo $n$, then $\gcd(x-1, n)$ and $\gcd(x+1, n)$ are both nontrivial divisors of $n$.

## ★ 31.9 Integer factorization

Suppose we have an integer $n$ that we wish to ***factor***, that is, to decompose into a product of primes. The primality test of the preceding section may tell us that $n$ is composite, but it does not tell us the prime factors of $n$. Factoring a large integer $n$ seems to be much more difficult than simply determining whether $n$ is prime or composite. Even with today's supercomputers and the best algorithms to date, we cannot feasibly factor an arbitrary 1024-bit number.

**Pollard's rho heuristic**

Trial division by all integers up to $R$ is guaranteed to factor completely any number up to $R^2$. For the same amount of work, the following procedure, POLLARD-RHO, factors any number up to $R^4$ (unless we are unlucky). Since the procedure is only a heuristic, neither its running time nor its success is guaranteed, although the procedure is highly effective in practice. Another advantage of the POLLARD-RHO procedure is that it uses only a constant number of memory locations. (If you wanted to, you could easily implement POLLARD-RHO on a programmable pocket calculator to find factors of small numbers.)

POLLARD-RHO$(n)$

```
 1  i = 1
 2  x₁ = RANDOM(0, n − 1)
 3  y = x₁
 4  k = 2
 5  while TRUE
 6      i = i + 1
 7      xᵢ = (x²ᵢ₋₁ − 1) mod n
 8      d = gcd(y − xᵢ, n)
 9      if d ≠ 1 and d ≠ n
10          print d
11      if i == k
12          y = xᵢ
13          k = 2k
```

The procedure works as follows. Lines 1–2 initialize $i$ to 1 and $x_1$ to a randomly chosen value in $\mathbb{Z}_n$. The **while** loop beginning on line 5 iterates forever, searching for factors of $n$. During each iteration of the **while** loop, line 7 uses the recurrence

$$x_i = (x_{i-1}^2 - 1) \bmod n \tag{31.43}$$

to produce the next value of $x_i$ in the infinite sequence

$$x_1, x_2, x_3, x_4, \ldots , \tag{31.44}$$

with line 6 correspondingly incrementing $i$. The pseudocode is written using subscripted variables $x_i$ for clarity, but the program works the same if all of the subscripts are dropped, since only the most recent value of $x_i$ needs to be maintained. With this modification, the procedure uses only a constant number of memory locations.

Every so often, the program saves the most recently generated $x_i$ value in the variable $y$. Specifically, the values that are saved are the ones whose subscripts are powers of 2:

$x_1, x_2, x_4, x_8, x_{16}, \ldots$ .

Line 3 saves the value $x_1$, and line 12 saves $x_k$ whenever $i$ is equal to $k$. The variable $k$ is initialized to 2 in line 4, and line 13 doubles it whenever line 12 updates $y$. Therefore, $k$ follows the sequence $1, 2, 4, 8, \ldots$ and always gives the subscript of the next value $x_k$ to be saved in $y$.

Lines 8–10 try to find a factor of $n$, using the saved value of $y$ and the current value of $x_i$. Specifically, line 8 computes the greatest common divisor $d = \gcd(y - x_i, n)$. If line 9 finds $d$ to be a nontrivial divisor of $n$, then line 10 prints $d$.

This procedure for finding a factor may seem somewhat mysterious at first. Note, however, that POLLARD-RHO never prints an incorrect answer; any number it prints is a nontrivial divisor of $n$. POLLARD-RHO might not print anything at all, though; it comes with no guarantee that it will print any divisors. We shall see, however, that we have good reason to expect POLLARD-RHO to print a factor $p$ of $n$ after $\Theta(\sqrt{p})$ iterations of the **while** loop. Thus, if $n$ is composite, we can expect this procedure to discover enough divisors to factor $n$ completely after approximately $n^{1/4}$ updates, since every prime factor $p$ of $n$ except possibly the largest one is less than $\sqrt{n}$.
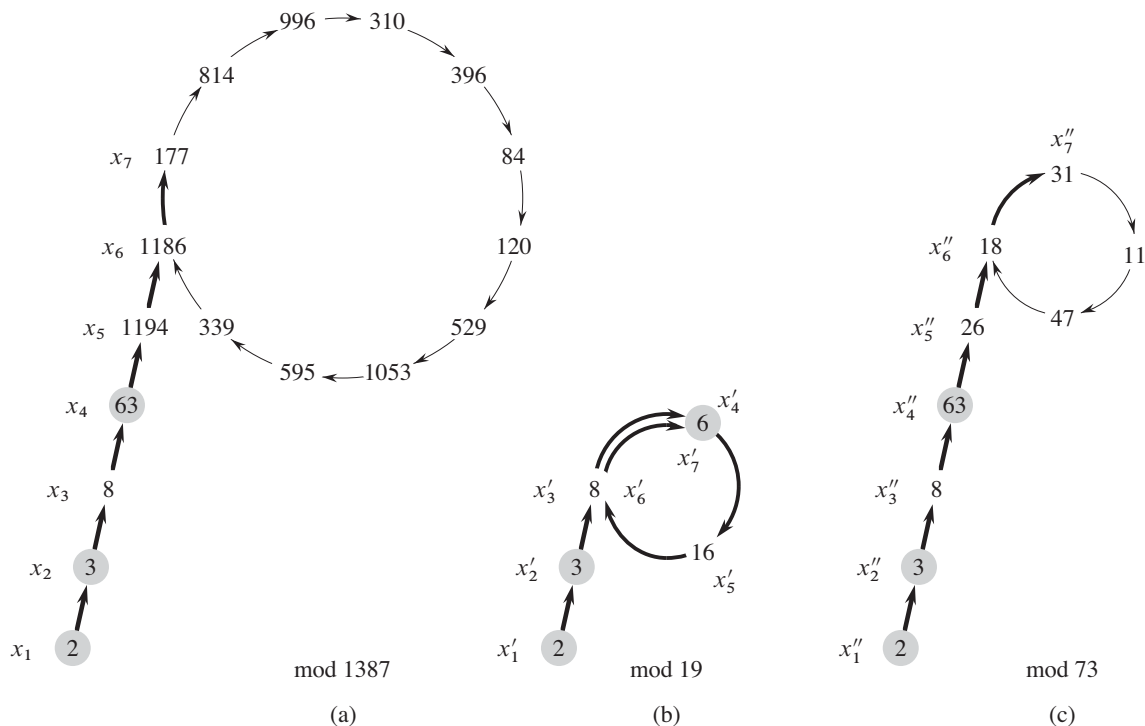
We begin our analysis of how this procedure behaves by studying how long it takes a random sequence modulo $n$ to repeat a value. Since $\mathbb{Z}_n$ is finite, and since each value in the sequence (31.44) depends only on the previous value, the sequence (31.44) eventually repeats itself. Once we reach an $x_i$ such that $x_i = x_j$ for some $j < i$, we are in a cycle, since $x_{i+1} = x_{j+1}$, $x_{i+2} = x_{j+2}$, and so on. The reason for the name "rho heuristic" is that, as Figure 31.7 shows, we can draw the sequence $x_1, x_2, \ldots, x_{j-1}$ as the "tail" of the rho and the cycle $x_j, x_{j+1}, \ldots, x_i$ as the "body" of the rho.

Let us consider the question of how long it takes for the sequence of $x_i$ to repeat. This information is not exactly what we need, but we shall see later how to modify the argument. For the purpose of this estimation, let us assume that the function

$$f_n(x) = (x^2 - 1) \bmod n$$

behaves like a "random" function. Of course, it is not really random, but this assumption yields results consistent with the observed behavior of POLLARD-RHO. We can then consider each $x_i$ to have been independently drawn from $\mathbb{Z}_n$ according to a uniform distribution on $\mathbb{Z}_n$. By the birthday-paradox analysis of Section 5.4.1, we expect $\Theta(\sqrt{n})$ steps to be taken before the sequence cycles.

Now for the required modification. Let $p$ be a nontrivial factor of $n$ such that $\gcd(p, n/p) = 1$. For example, if $n$ has the factorization $n = p_1^{e_1} p_2^{e_2} \cdots p_r^{e_r}$, then we may take $p$ to be $p_1^{e_1}$. (If $e_1 = 1$, then $p$ is just the smallest prime factor of $n$, a good example to keep in mind.)

**Figure 31.7** Pollard's rho heuristic.    **(a)** The values produced by the recurrence $x_{i+1} = (x_i^2 - 1) \bmod 1387$, starting with $x_1 = 2$. The prime factorization of 1387 is $19 \cdot 73$. The heavy arrows indicate the iteration steps that are executed before the factor 19 is discovered. The light arrows point to unreached values in the iteration, to illustrate the "rho" shape. The shaded values are the $y$ values stored by POLLARD-RHO. The factor 19 is discovered upon reaching $x_7 = 177$, when $\gcd(63 - 177, 1387) = 19$ is computed. The first $x$ value that would be repeated is 1186, but the factor 19 is discovered before this value is repeated. **(b)** The values produced by the same recurrence, modulo 19. Every value $x_i$ given in part (a) is equivalent, modulo 19, to the value $x_i'$ shown here. For example, both $x_4 = 63$ and $x_7 = 177$ are equivalent to 6, modulo 19. **(c)** The values produced by the same recurrence, modulo 73. Every value $x_i$ given in part (a) is equivalent, modulo 73, to the value $x_i''$ shown here. By the Chinese remainder theorem, each node in part (a) corresponds to a pair of nodes, one from part (b) and one from part (c).

The sequence $\langle x_i \rangle$ induces a corresponding sequence $\langle x_i' \rangle$ modulo $p$, where

$$x_i' = x_i \bmod p$$

for all $i$.

Furthermore, because $f_n$ is defined using only arithmetic operations (squaring and subtraction) modulo $n$, we can compute $x_{i+1}'$ from $x_i'$; the "modulo $p$" view of

the sequence is a smaller version of what is happening modulo $n$:

$$
\begin{aligned}
x'_{i+1} &= x_{i+1} \bmod p \\
&= f_n(x_i) \bmod p \\
&= ((x_i^2 - 1) \bmod n) \bmod p \\
&= (x_i^2 - 1) \bmod p &\text{(by Exercise 31.1-7)} \\
&= ((x_i \bmod p)^2 - 1) \bmod p \\
&= ((x'_i)^2 - 1) \bmod p \\
&= f_p(x'_i) \,.
\end{aligned}
$$

Thus, although we are not explicitly computing the sequence $\langle x'_i \rangle$, this sequence is well defined and obeys the same recurrence as the sequence $\langle x_i \rangle$.

Reasoning as before, we find that the expected number of steps before the sequence $\langle x'_i \rangle$ repeats is $\Theta(\sqrt{p})$. If $p$ is small compared to $n$, the sequence $\langle x'_i \rangle$ might repeat much more quickly than the sequence $\langle x_i \rangle$. Indeed, as parts (b) and (c) of Figure 31.7 show, the $\langle x'_i \rangle$ sequence repeats as soon as two elements of the sequence $\langle x_i \rangle$ are merely equivalent modulo $p$, rather than equivalent modulo $n$.

Let $t$ denote the index of the first repeated value in the $\langle x'_i \rangle$ sequence, and let $u > 0$ denote the length of the cycle that has been thereby produced. That is, $t$ and $u > 0$ are the smallest values such that $x'_{t+i} = x'_{t+u+i}$ for all $i \geq 0$. By the above arguments, the expected values of $t$ and $u$ are both $\Theta(\sqrt{p})$. Note that if $x'_{t+i} = x'_{t+u+i}$, then $p \mid (x_{t+u+i} - x_{t+i})$. Thus, $\gcd(x_{t+u+i} - x_{t+i}, n) > 1$.

Therefore, once POLLARD-RHO has saved as $y$ any value $x_k$ such that $k \geq t$, then $y \bmod p$ is always on the cycle modulo $p$. (If a new value is saved as $y$, that value is also on the cycle modulo $p$.) Eventually, $k$ is set to a value that is greater than $u$, and the procedure then makes an entire loop around the cycle modulo $p$ without changing the value of $y$. The procedure then discovers a factor of $n$ when $x_i$ "runs into" the previously stored value of $y$, modulo $p$, that is, when $x_i \equiv y \pmod{p}$.

Presumably, the factor found is the factor $p$, although it may occasionally happen that a multiple of $p$ is discovered. Since the expected values of both $t$ and $u$ are $\Theta(\sqrt{p})$, the expected number of steps required to produce the factor $p$ is $\Theta(\sqrt{p})$.

This algorithm might not perform quite as expected, for two reasons. First, the heuristic analysis of the running time is not rigorous, and it is possible that the cycle of values, modulo $p$, could be much larger than $\sqrt{p}$. In this case, the algorithm performs correctly but much more slowly than desired. In practice, this issue seems to be moot. Second, the divisors of $n$ produced by this algorithm might always be one of the trivial factors 1 or $n$. For example, suppose that $n = pq$, where $p$ and $q$ are prime. It can happen that the values of $t$ and $u$ for $p$ are identical with the values of $t$ and $u$ for $q$, and thus the factor $p$ is always revealed in the same gcd operation that reveals the factor $q$. Since both factors are revealed at the same

time, the trivial factor $pq = n$ is revealed, which is useless. Again, this problem seems to be insignificant in practice. If necessary, we can restart the heuristic with a different recurrence of the form $x_{i+1} = (x_i^2 - c) \bmod n$. (We should avoid the values $c = 0$ and $c = 2$ for reasons we will not go into here, but other values are fine.)

Of course, this analysis is heuristic and not rigorous, since the recurrence is not really "random." Nonetheless, the procedure performs well in practice, and it seems to be as efficient as this heuristic analysis indicates. It is the method of choice for finding small prime factors of a large number. To factor a $\beta$-bit composite number $n$ completely, we only need to find all prime factors less than $\lfloor n^{1/2} \rfloor$, and so we expect POLLARD-RHO to require at most $n^{1/4} = 2^{\beta/4}$ arithmetic operations and at most $n^{1/4}\beta^2 = 2^{\beta/4}\beta^2$ bit operations. POLLARD-RHO's ability to find a small factor $p$ of $n$ with an expected number $\Theta(\sqrt{p})$ of arithmetic operations is often its most appealing feature.

### Exercises

***31.9-1***
Referring to the execution history shown in Figure 31.7(a), when does POLLARD-RHO print the factor 73 of 1387?

***31.9-2***
Suppose that we are given a function $f : \mathbb{Z}_n \to \mathbb{Z}_n$ and an initial value $x_0 \in \mathbb{Z}_n$. Define $x_i = f(x_{i-1})$ for $i = 1, 2, \ldots$. Let $t$ and $u > 0$ be the smallest values such that $x_{t+i} = x_{t+u+i}$ for $i = 0, 1, \ldots$. In the terminology of Pollard's rho algorithm, $t$ is the length of the tail and $u$ is the length of the cycle of the rho. Give an efficient algorithm to determine $t$ and $u$ exactly, and analyze its running time.

***31.9-3***
How many steps would you expect POLLARD-RHO to require to discover a factor of the form $p^e$, where $p$ is prime and $e > 1$?

***31.9-4***  ★
One disadvantage of POLLARD-RHO as written is that it requires one gcd computation for each step of the recurrence. Instead, we could batch the gcd computations by accumulating the product of several $x_i$ values in a row and then using this product instead of $x_i$ in the gcd computation. Describe carefully how you would implement this idea, why it works, and what batch size you would pick as the most effective when working on a $\beta$-bit number $n$.

## Problems

### 31-1  *Binary gcd algorithm*
Most computers can perform the operations of subtraction, testing the parity (odd or even) of a binary integer, and halving more quickly than computing remainders. This problem investigates the ***binary gcd algorithm***, which avoids the remainder computations used in Euclid's algorithm.

***a.*** Prove that if $a$ and $b$ are both even, then $\gcd(a, b) = 2 \cdot \gcd(a/2, b/2)$.

***b.*** Prove that if $a$ is odd and $b$ is even, then $\gcd(a, b) = \gcd(a, b/2)$.

***c.*** Prove that if $a$ and $b$ are both odd, then $\gcd(a, b) = \gcd((a - b)/2, b)$.

***d.*** Design an efficient binary gcd algorithm for input integers $a$ and $b$, where $a \geq b$, that runs in $O(\lg a)$ time. Assume that each subtraction, parity test, and halving takes unit time.

### 31-2  *Analysis of bit operations in Euclid's algorithm*

***a.*** Consider the ordinary "paper and pencil" algorithm for long division: dividing $a$ by $b$, which yields a quotient $q$ and remainder $r$. Show that this method requires $O((1 + \lg q) \lg b)$ bit operations.

***b.*** Define $\mu(a, b) = (1 + \lg a)(1 + \lg b)$. Show that the number of bit operations performed by EUCLID in reducing the problem of computing $\gcd(a, b)$ to that of computing $\gcd(b, a \bmod b)$ is at most $c(\mu(a, b) - \mu(b, a \bmod b))$ for some sufficiently large constant $c > 0$.

***c.*** Show that EUCLID$(a, b)$ requires $O(\mu(a, b))$ bit operations in general and $O(\beta^2)$ bit operations when applied to two $\beta$-bit inputs.

### 31-3  *Three algorithms for Fibonacci numbers*
This problem compares the efficiency of three methods for computing the $n$th Fibonacci number $F_n$, given $n$. Assume that the cost of adding, subtracting, or multiplying two numbers is $O(1)$, independent of the size of the numbers.

***a.*** Show that the running time of the straightforward recursive method for computing $F_n$ based on recurrence (3.22) is exponential in $n$. (See, for example, the FIB procedure on page 775.)

***b.*** Show how to compute $F_n$ in $O(n)$ time using memoization.

c. Show how to compute $F_n$ in $O(\lg n)$ time using only integer addition and multiplication. (*Hint:* Consider the matrix

$$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix}$$

and its powers.)

d. Assume now that adding two $\beta$-bit numbers takes $\Theta(\beta)$ time and that multiplying two $\beta$-bit numbers takes $\Theta(\beta^2)$ time. What is the running time of these three methods under this more reasonable cost measure for the elementary arithmetic operations?

### 31-4  *Quadratic residues*

Let $p$ be an odd prime. A number $a \in Z_p^*$ is a *quadratic residue* if the equation $x^2 = a \pmod{p}$ has a solution for the unknown $x$.

a. Show that there are exactly $(p-1)/2$ quadratic residues, modulo $p$.

b. If $p$ is prime, we define the *Legendre symbol* $\left(\frac{a}{p}\right)$, for $a \in Z_p^*$, to be 1 if $a$ is a quadratic residue modulo $p$ and $-1$ otherwise. Prove that if $a \in Z_p^*$, then

$$\left(\frac{a}{p}\right) \equiv a^{(p-1)/2} \pmod{p} .$$

Give an efficient algorithm that determines whether a given number $a$ is a quadratic residue modulo $p$. Analyze the efficiency of your algorithm.

c. Prove that if $p$ is a prime of the form $4k + 3$ and $a$ is a quadratic residue in $Z_p^*$, then $a^{k+1} \bmod p$ is a square root of $a$, modulo $p$. How much time is required to find the square root of a quadratic residue $a$ modulo $p$?

d. Describe an efficient randomized algorithm for finding a nonquadratic residue, modulo an arbitrary prime $p$, that is, a member of $Z_p^*$ that is not a quadratic residue. How many arithmetic operations does your algorithm require on average?

## Chapter notes

Niven and Zuckerman [265] provide an excellent introduction to elementary number theory. Knuth [210] contains a good discussion of algorithms for finding the

greatest common divisor, as well as other basic number-theoretic algorithms. Bach [30] and Riesel [295] provide more recent surveys of computational number theory. Dixon [91] gives an overview of factorization and primality testing. The conference proceedings edited by Pomerance [280] contains several excellent survey articles. More recently, Bach and Shallit [31] have provided an exceptional overview of the basics of computational number theory.

Knuth [210] discusses the origin of Euclid's algorithm. It appears in Book 7, Propositions 1 and 2, of the Greek mathematician Euclid's *Elements*, which was written around 300 B.C. Euclid's description may have been derived from an algorithm due to Eudoxus around 375 B.C. Euclid's algorithm may hold the honor of being the oldest nontrivial algorithm; it is rivaled only by an algorithm for multiplication known to the ancient Egyptians. Shallit [312] chronicles the history of the analysis of Euclid's algorithm.

Knuth attributes a special case of the Chinese remainder theorem (Theorem 31.27) to the Chinese mathematician Sun-Tsǔ, who lived sometime between 200 B.C. and A.D. 200—the date is quite uncertain. The same special case was given by the Greek mathematician Nichomachus around A.D. 100. It was generalized by Chhin Chiu-Shao in 1247. The Chinese remainder theorem was finally stated and proved in its full generality by L. Euler in 1734.

The randomized primality-testing algorithm presented here is due to Miller [255] and Rabin [289]; it is the fastest randomized primality-testing algorithm known, to within constant factors. The proof of Theorem 31.39 is a slight adaptation of one suggested by Bach [29]. A proof of a stronger result for MILLER-RABIN was given by Monier [258, 259]. For many years primality-testing was the classic example of a problem where randomization appeared to be necessary to obtain an efficient (polynomial-time) algorithm. In 2002, however, Agrawal, Kayal, and Saxema [4] surprised everyone with their deterministic polynomial-time primality-testing algorithm. Until then, the fastest deterministic primality testing algorithm known, due to Cohen and Lenstra [73], ran in time $(\lg n)^{O(\lg \lg \lg n)}$ on input $n$, which is just slightly superpolynomial. Nonetheless, for practical purposes randomized primality-testing algorithms remain more efficient and are preferred.

The problem of finding large "random" primes is nicely discussed in an article by Beauchemin, Brassard, Crépeau, Goutier, and Pomerance [36].

The concept of a public-key cryptosystem is due to Diffie and Hellman [87]. The RSA cryptosystem was proposed in 1977 by Rivest, Shamir, and Adleman [296]. Since then, the field of cryptography has blossomed. Our understanding of the RSA cryptosystem has deepened, and modern implementations use significant refinements of the basic techniques presented here. In addition, many new techniques have been developed for proving cryptosystems to be secure. For example, Goldwasser and Micali [142] show that randomization can be an effective tool in the design of secure public-key encryption schemes. For signature schemes,

Goldwasser, Micali, and Rivest [143] present a digital-signature scheme for which every conceivable type of forgery is provably as difficult as factoring. Menezes, van Oorschot, and Vanstone [254] provide an overview of applied cryptography.

The rho heuristic for integer factorization was invented by Pollard [277]. The version presented here is a variant proposed by Brent [56].

The best algorithms for factoring large numbers have a running time that grows roughly exponentially with the cube root of the length of the number $n$ to be factored. The general number-field sieve factoring algorithm (as developed by Buhler, Lenstra, and Pomerance [57] as an extension of the ideas in the number-field sieve factoring algorithm by Pollard [278] and Lenstra et al. [232] and refined by Coppersmith [77] and others) is perhaps the most efficient such algorithm in general for large inputs. Although it is difficult to give a rigorous analysis of this algorithm, under reasonable assumptions we can derive a running-time estimate of $L(1/3, n)^{1.902+o(1)}$, where $L(\alpha, n) = e^{(\ln n)^{\alpha}(\ln \ln n)^{1-\alpha}}$.
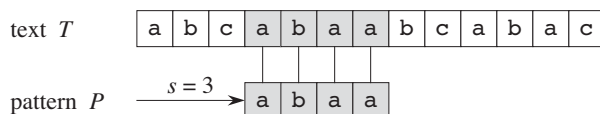
The elliptic-curve method due to Lenstra [233] may be more effective for some inputs than the number-field sieve method, since, like Pollard's rho method, it can find a small prime factor $p$ quite quickly. With this method, the time to find $p$ is estimated to be $L(1/2, p)^{\sqrt{2}+o(1)}$.

# 32    String Matching

Text-editing programs frequently need to find all occurrences of a pattern in the text. Typically, the text is a document being edited, and the pattern searched for is a particular word supplied by the user. Efficient algorithms for this problem—called "string matching"—can greatly aid the responsiveness of the text-editing program. Among their many other applications, string-matching algorithms search for particular patterns in DNA sequences. Internet search engines also use them to find Web pages relevant to queries.

We formalize the string-matching problem as follows. We assume that the text is an array $T[1 . . n]$ of length $n$ and that the pattern is an array $P[1 . . m]$ of length $m \leq n$. We further assume that the elements of $P$ and $T$ are characters drawn from a finite alphabet $\Sigma$. For example, we may have $\Sigma = \{0,1\}$ or $\Sigma = \{a, b, . . . , z\}$. The character arrays $P$ and $T$ are often called *strings* of characters.

Referring to Figure 32.1, we say that pattern $P$ *occurs with shift $s$* in text $T$ (or, equivalently, that pattern $P$ *occurs beginning at position $s + 1$* in text $T$) if $0 \leq s \leq n - m$ and $T[s + 1 . . s + m] = P[1 . . m]$ (that is, if $T[s + j] = P[j]$, for $1 \leq j \leq m$). If $P$ occurs with shift $s$ in $T$, then we call $s$ a *valid shift*; otherwise, we call $s$ an *invalid shift*. The *string-matching problem* is the problem of finding all valid shifts with which a given pattern $P$ occurs in a given text $T$.



**Figure 32.1**  An example of the string-matching problem, where we want to find all occurrences of the pattern $P = $ abaa in the text $T = $ abcabaabcabac. The pattern occurs only once in the text, at shift $s = 3$, which we call a valid shift. A vertical line connects each character of the pattern to its matching character in the text, and all matched characters are shaded.

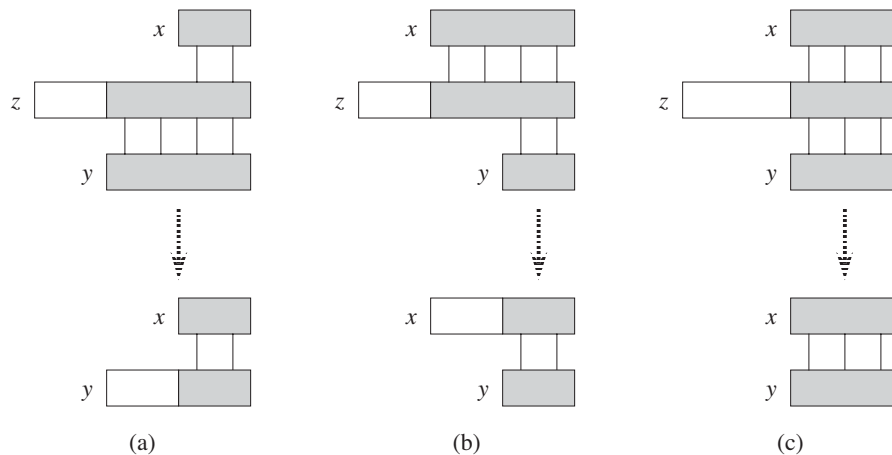| Algorithm | Preprocessing time | Matching time |
|---|---|---|
| Naive | 0 | $O((n-m+1)m)$ |
| Rabin-Karp | $\Theta(m)$ | $O((n-m+1)m)$ |
| Finite automaton | $O(m\,|\Sigma|)$ | $\Theta(n)$ |
| Knuth-Morris-Pratt | $\Theta(m)$ | $\Theta(n)$ |

**Figure 32.2**   The string-matching algorithms in this chapter and their preprocessing and matching times.

Except for the naive brute-force algorithm, which we review in Section 32.1, each string-matching algorithm in this chapter performs some preprocessing based on the pattern and then finds all valid shifts; we call this latter phase "matching." Figure 32.2 shows the preprocessing and matching times for each of the algorithms in this chapter. The total running time of each algorithm is the sum of the preprocessing and matching times. Section 32.2 presents an interesting string-matching algorithm, due to Rabin and Karp. Although the $\Theta((n-m+1)m)$ worst-case running time of this algorithm is no better than that of the naive method, it works much better on average and in practice. It also generalizes nicely to other pattern-matching problems. Section 32.3 then describes a string-matching algorithm that begins by constructing a finite automaton specifically designed to search for occurrences of the given pattern $P$ in a text. This algorithm takes $O(m\,|\Sigma|)$ preprocessing time, but only $\Theta(n)$ matching time. Section 32.4 presents the similar, but much cleverer, Knuth-Morris-Pratt (or KMP) algorithm; it has the same $\Theta(n)$ matching time, and it reduces the preprocessing time to only $\Theta(m)$.

### Notation and terminology

We denote by $\Sigma^*$ (read "sigma-star") the set of all finite-length strings formed using characters from the alphabet $\Sigma$. In this chapter, we consider only finite-length strings. The zero-length ***empty string***, denoted $\varepsilon$, also belongs to $\Sigma^*$. The length of a string $x$ is denoted $|x|$. The ***concatenation*** of two strings $x$ and $y$, denoted $xy$, has length $|x|+|y|$ and consists of the characters from $x$ followed by the characters from $y$.

We say that a string $w$ is a ***prefix*** of a string $x$, denoted $w \sqsubset x$, if $x = wy$ for some string $y \in \Sigma^*$. Note that if $w \sqsubset x$, then $|w| \le |x|$. Similarly, we say that a string $w$ is a ***suffix*** of a string $x$, denoted $w \sqsupset x$, if $x = yw$ for some $y \in \Sigma^*$. As with a prefix, $w \sqsupset x$ implies $|w| \le |x|$. For example, we have $\mathtt{ab} \sqsubset \mathtt{abcca}$ and $\mathtt{cca} \sqsupset \mathtt{abcca}$. The empty string $\varepsilon$ is both a suffix and a prefix of every string. For any strings $x$ and $y$ and any character $a$, we have $x \sqsupset y$ if and only if $xa \sqsupset ya$.

**Figure 32.3** A graphical proof of Lemma 32.1. We suppose that $x \sqsupset z$ and $y \sqsupset z$. The three parts of the figure illustrate the three cases of the lemma. Vertical lines connect matching regions (shown shaded) of the strings. **(a)** If $|x| \leq |y|$, then $x \sqsupset y$. **(b)** If $|x| \geq |y|$, then $y \sqsupset x$. **(c)** If $|x| = |y|$, then $x = y$.

Also note that $\sqsubset$ and $\sqsupset$ are transitive relations. The following lemma will be useful later.

*Lemma 32.1 (Overlapping-suffix lemma)*
Suppose that $x$, $y$, and $z$ are strings such that $x \sqsupset z$ and $y \sqsupset z$. If $|x| \leq |y|$, then $x \sqsupset y$. If $|x| \geq |y|$, then $y \sqsupset x$. If $|x| = |y|$, then $x = y$.

*Proof*   See Figure 32.3 for a graphical proof. ∎

For brevity of notation, we denote the $k$-character prefix $P[1 . . k]$ of the pattern $P[1 . . m]$ by $P_k$. Thus, $P_0 = \varepsilon$ and $P_m = P = P[1 . . m]$. Similarly, we denote the $k$-character prefix of the text $T$ by $T_k$. Using this notation, we can state the string-matching problem as that of finding all shifts $s$ in the range $0 \leq s \leq n - m$ such that $P \sqsupset T_{s+m}$.

In our pseudocode, we allow two equal-length strings to be compared for equality as a primitive operation. If the strings are compared from left to right and the comparison stops when a mismatch is discovered, we assume that the time taken by such a test is a linear function of the number of matching characters discovered. To be precise, the test "$x == y$" is assumed to take time $\Theta(t + 1)$, where $t$ is the length of the longest string $z$ such that $z \sqsubset x$ and $z \sqsubset y$. (We write $\Theta(t + 1)$ rather than $\Theta(t)$ to handle the case in which $t = 0$; the first characters compared do not match, but it takes a positive amount of time to perform this comparison.)

## 32.1   The naive string-matching algorithm

The naive algorithm finds all valid shifts using a loop that checks the condition $P[1 . . m] = T[s + 1 . . s + m]$ for each of the $n - m + 1$ possible values of $s$.
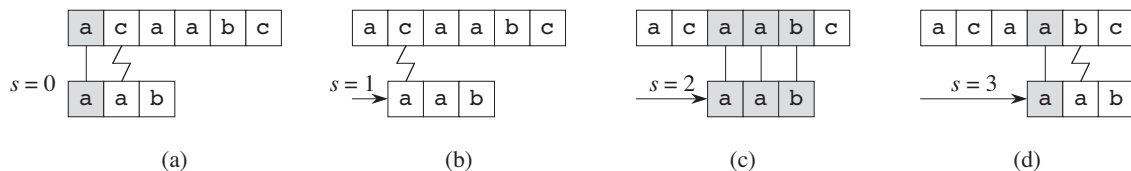
NAIVE-STRING-MATCHER$(T, P)$

```
1   n = T.length
2   m = P.length
3   for s = 0 to n − m
4       if P[1 . . m] == T[s + 1 . . s + m]
5           print "Pattern occurs with shift" s
```

Figure 32.4 portrays the naive string-matching procedure as sliding a "template" containing the pattern over the text, noting for which shifts all of the characters on the template equal the corresponding characters in the text. The **for** loop of lines 3–5 considers each possible shift explicitly. The test in line 4 determines whether the current shift is valid; this test implicitly loops to check corresponding character positions until all positions match successfully or a mismatch is found. Line 5 prints out each valid shift $s$.

Procedure NAIVE-STRING-MATCHER takes time $O((n - m + 1)m)$, and this bound is tight in the worst case. For example, consider the text string $a^n$ (a string of $n$ a's) and the pattern $a^m$. For each of the $n - m + 1$ possible values of the shift $s$, the implicit loop on line 4 to compare corresponding characters must execute $m$ times to validate the shift. The worst-case running time is thus $\Theta((n - m + 1)m)$, which is $\Theta(n^2)$ if $m = \lfloor n/2 \rfloor$. Because it requires no preprocessing, NAIVE-STRING-MATCHER's running time equals its matching time.

**Figure 32.4**   The operation of the naive string matcher for the pattern $P = \mathtt{aab}$ and the text $T = \mathtt{acaabc}$. We can imagine the pattern $P$ as a template that we slide next to the text. **(a)–(d)** The four successive alignments tried by the naive string matcher. In each part, vertical lines connect corresponding regions found to match (shown shaded), and a jagged line connects the first mismatched character found, if any. The algorithm finds one occurrence of the pattern, at shift $s = 2$, shown in part (c).

As we shall see, NAIVE-STRING-MATCHER is not an optimal procedure for this problem. Indeed, in this chapter we shall see that the Knuth-Morris-Pratt algorithm is much better in the worst case. The naive string-matcher is inefficient because it entirely ignores information gained about the text for one value of $s$ when it considers other values of $s$. Such information can be quite valuable, however. For example, if $P = $ aaab and we find that $s = 0$ is valid, then none of the shifts 1, 2, or 3 are valid, since $T[4] = $ b. In the following sections, we examine several ways to make effective use of this sort of information.

**Exercises**

***32.1-1***
Show the comparisons the naive string matcher makes for the pattern $P = $ 0001 in the text $T = $ 000010001010001.

***32.1-2***
Suppose that all characters in the pattern $P$ are different. Show how to accelerate NAIVE-STRING-MATCHER to run in time $O(n)$ on an $n$-character text $T$.

***32.1-3***
Suppose that pattern $P$ and text $T$ are *randomly* chosen strings of length $m$ and $n$, respectively, from the $d$-ary alphabet $\Sigma_d = \{0, 1, \ldots, d-1\}$, where $d \geq 2$. Show that the *expected* number of character-to-character comparisons made by the implicit loop in line 4 of the naive algorithm is

$$(n-m+1)\frac{1-d^{-m}}{1-d^{-1}} \leq 2(n-m+1)$$

over all executions of this loop. (Assume that the naive algorithm stops comparing characters for a given shift once it finds a mismatch or matches the entire pattern.) Thus, for randomly chosen strings, the naive algorithm is quite efficient.

***32.1-4***
Suppose we allow the pattern $P$ to contain occurrences of a ***gap character*** $\diamond$ that can match an *arbitrary* string of characters (even one of zero length). For example, the pattern ab$\diamond$ba$\diamond$c occurs in the text cabccbacbacab as

c ab cc ba cba c ab
   ab  $\diamond$  ba  $\diamond$  c

and as

c ab ccbac ba     c ab .
   ab   $\diamond$   ba $\diamond$ c

Note that the gap character may occur an arbitrary number of times in the pattern but not at all in the text. Give a polynomial-time algorithm to determine whether such a pattern $P$ occurs in a given text $T$, and analyze the running time of your algorithm.

## 32.2   The Rabin-Karp algorithm

Rabin and Karp proposed a string-matching algorithm that performs well in practice and that also generalizes to other algorithms for related problems, such as two-dimensional pattern matching. The Rabin-Karp algorithm uses $\Theta(m)$ preprocessing time, and its worst-case running time is $\Theta((n-m+1)m)$. Based on certain assumptions, however, its average-case running time is better.

This algorithm makes use of elementary number-theoretic notions such as the equivalence of two numbers modulo a third number. You might want to refer to Section 31.1 for the relevant definitions.

For expository purposes, let us assume that $\Sigma = \{0, 1, 2, \ldots, 9\}$, so that each character is a decimal digit. (In the general case, we can assume that each character is a digit in radix-$d$ notation, where $d = |\Sigma|$.) We can then view a string of $k$ consecutive characters as representing a length-$k$ decimal number. The character string 31415 thus corresponds to the decimal number 31,415. Because we interpret the input characters as both graphical symbols and digits, we find it convenient in this section to denote them as we would digits, in our standard text font.

Given a pattern $P[1 \ldots m]$, let $p$ denote its corresponding decimal value. In a similar manner, given a text $T[1 \ldots n]$, let $t_s$ denote the decimal value of the length-$m$ substring $T[s+1 \ldots s+m]$, for $s = 0, 1, \ldots, n-m$. Certainly, $t_s = p$ if and only if $T[s+1 \ldots s+m] = P[1 \ldots m]$; thus, $s$ is a valid shift if and only if $t_s = p$. If we could compute $p$ in time $\Theta(m)$ and all the $t_s$ values in a total of $\Theta(n-m+1)$ time,[1] then we could determine all valid shifts $s$ in time $\Theta(m) + \Theta(n-m+1) = \Theta(n)$ by comparing $p$ with each of the $t_s$ values. (For the moment, let's not worry about the possibility that $p$ and the $t_s$ values might be very large numbers.)

We can compute $p$ in time $\Theta(m)$ using Horner's rule (see Section 30.1):

$$p = P[m] + 10\,(P[m-1] + 10(P[m-2] + \cdots + 10(P[2] + 10P[1])\cdots))\,.$$

Similarly, we can compute $t_0$ from $T[1 \ldots m]$ in time $\Theta(m)$.

---

[1] We write $\Theta(n-m+1)$ instead of $\Theta(n-m)$ because $s$ takes on $n-m+1$ different values. The "+1" is significant in an asymptotic sense because when $m = n$, computing the lone $t_s$ value takes $\Theta(1)$ time, not $\Theta(0)$ time.

To compute the remaining values $t_1, t_2, \ldots, t_{n-m}$ in time $\Theta(n-m)$, we observe that we can compute $t_{s+1}$ from $t_s$ in constant time, since

$$t_{s+1} = 10(t_s - 10^{m-1} T[s+1]) + T[s+m+1] . \tag{32.1}$$

Subtracting $10^{m-1} T[s+1]$ removes the high-order digit from $t_s$, multiplying the result by 10 shifts the number left by one digit position, and adding $T[s+m+1]$ brings in the appropriate low-order digit. For example, if $m = 5$ and $t_s = 31415$, then we wish to remove the high-order digit $T[s+1] = 3$ and bring in the new low-order digit (suppose it is $T[s+5+1] = 2$) to obtain

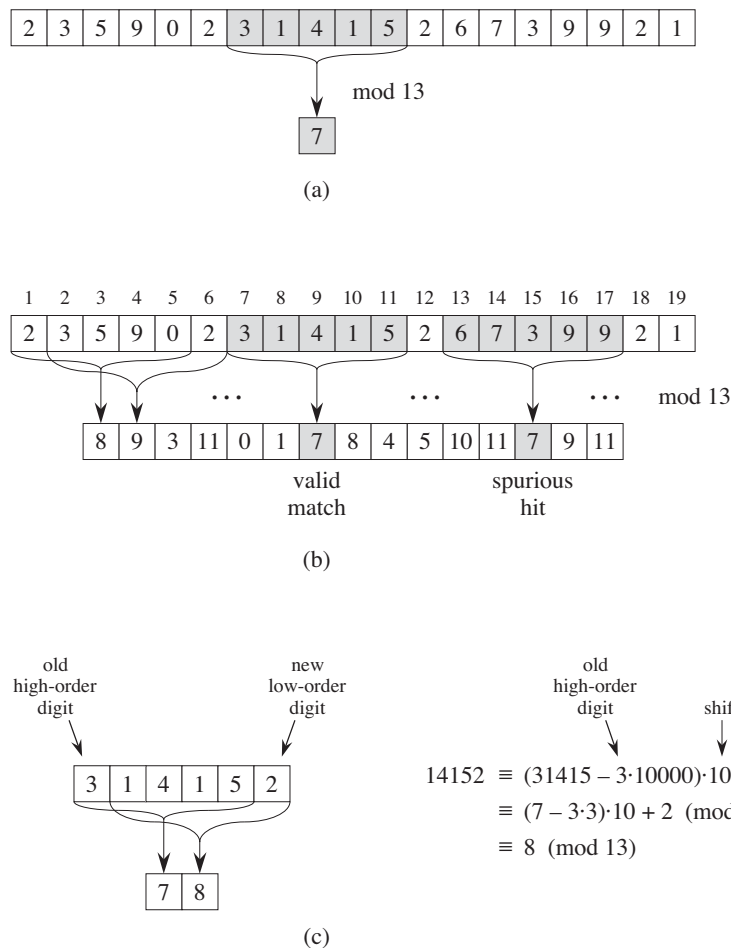$$\begin{aligned} t_{s+1} &= 10(31415 - 10000 \cdot 3) + 2 \\ &= 14152 . \end{aligned}$$

If we precompute the constant $10^{m-1}$ (which we can do in time $O(\lg m)$ using the techniques of Section 31.6, although for this application a straightforward $O(m)$-time method suffices), then each execution of equation (32.1) takes a constant number of arithmetic operations. Thus, we can compute $p$ in time $\Theta(m)$, and we can compute all of $t_0, t_1, \ldots, t_{n-m}$ in time $\Theta(n-m+1)$. Therefore, we can find all occurrences of the pattern $P[1 \mathinner{.\,.} m]$ in the text $T[1 \mathinner{.\,.} n]$ with $\Theta(m)$ preprocessing time and $\Theta(n-m+1)$ matching time.

Until now, we have intentionally overlooked one problem: $p$ and $t_s$ may be too large to work with conveniently. If $P$ contains $m$ characters, then we cannot reasonably assume that each arithmetic operation on $p$ (which is $m$ digits long) takes "constant time." Fortunately, we can solve this problem easily, as Figure 32.5 shows: compute $p$ and the $t_s$ values modulo a suitable modulus $q$. We can compute $p$ modulo $q$ in $\Theta(m)$ time and all the $t_s$ values modulo $q$ in $\Theta(n-m+1)$ time. If we choose the modulus $q$ as a prime such that $10q$ just fits within one computer word, then we can perform all the necessary computations with single-precision arithmetic. In general, with a $d$-ary alphabet $\{0, 1, \ldots, d-1\}$, we choose $q$ so that $dq$ fits within a computer word and adjust the recurrence equation (32.1) to work modulo $q$, so that it becomes

$$t_{s+1} = (d(t_s - T[s+1]h) + T[s+m+1]) \bmod q , \tag{32.2}$$

where $h \equiv d^{m-1} \pmod{q}$ is the value of the digit "1" in the high-order position of an $m$-digit text window.

The solution of working modulo $q$ is not perfect, however: $t_s \equiv p \pmod{q}$ does not imply that $t_s = p$. On the other hand, if $t_s \not\equiv p \pmod{q}$, then we definitely have that $t_s \neq p$, so that shift $s$ is invalid. We can thus use the test $t_s \equiv p \pmod{q}$ as a fast heuristic test to rule out invalid shifts $s$. Any shift $s$ for which $t_s \equiv p \pmod{q}$ must be tested further to see whether $s$ is really valid or we just have a *spurious hit*. This additional test explicitly checks the condition

| 2 | 3 | 5 | 9 | 0 | 2 | 3 | 1 | 4 | 1 | 5 | 2 | 6 | 7 | 3 | 9 | 9 | 2 | 1 |

mod 13

| 7 |

(a)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 2 | 3 | 5 | 9 | 0 | 2 | 3 | 1 | 4 | 1 | 5 | 2 | 6 | 7 | 3 | 9 | 9 | 2 | 1 |

···    ···    ···    mod 13

| 8 | 9 | 3 | 11 | 0 | 1 | 7 | 8 | 4 | 5 | 10 | 11 | 7 | 9 | 11 |

valid            spurious
match            hit

(b)

old                new                old                          new
high-order         low-order          high-order        shift      low-order
digit              digit              digit                         digit

| 3 | 1 | 4 | 1 | 5 | 2 |

| 7 | 8 |

$$14152 \equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13}$$
$$\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13}$$
$$\equiv 8 \pmod{13}$$

(c)

**Figure 32.5**   The Rabin-Karp algorithm. Each character is a decimal digit, and we compute values modulo 13. **(a)** A text string. A window of length 5 is shown shaded. The numerical value of the shaded number, computed modulo 13, yields the value 7. **(b)** The same text string with values computed modulo 13 for each possible position of a length-5 window. Assuming the pattern $P = 31415$, we look for windows whose value modulo 13 is 7, since $31415 \equiv 7 \pmod{13}$. The algorithm finds two such windows, shown shaded in the figure. The first, beginning at text position 7, is indeed an occurrence of the pattern, while the second, beginning at text position 13, is a spurious hit. **(c)** How to compute the value for a window in constant time, given the value for the previous window. The first window has value 31415. Dropping the high-order digit 3, shifting left (multiplying by 10), and then adding in the low-order digit 2 gives us the new value 14152. Because all computations are performed modulo 13, the value for the first window is 7, and the value for the new window is 8.