

## 1 Modular Arithmetic

In several settings, such as error-correcting codes and cryptography, we sometimes wish to work over a smaller range of numbers. Modular arithmetic is useful in these settings, since it limits numbers to a predefined range  $\{0, 1, \dots, N - 1\}$ , and wraps around whenever you try to leave this range — like the hand of a clock (where  $N = 12$ ) or the days of the week (where  $N = 7$ ).

In this way, modular arithmetic is not unrelated to what you have seen when working multiplicatively on the unit circle in the complex plane in 16B. When multiplying numbers on the unit circle together, you never leave the unit circle. With modular arithmetic, we want a similarly bounded universe while supporting both multiplication and addition.

**Example: Calculating the time** When you calculate the time, you automatically use modular arithmetic. For example, if you are asked what time it will be 13 hours from 1 pm, you say 2 am rather than 14. Let's assume our clock displays 12 as 0. This is limiting numbers to a predefined range,  $\{0, 1, 2, \dots, 11\}$ . Whenever you add two numbers in this setting, you divide by 12 and provide the remainder as the answer.

If we wanted to know what the time would be 24 hours from 2 pm, the answer is easy. It would be 2 pm. This is true not just for 24 hours, but for any multiple of 12 hours (ignoring the detail of am/pm). What about 25 hours from 2 pm? Since the time 24 hours from 2 pm is still 2 pm, after 25 hours it would be 3 pm. Another way to say this is that we add 1 hour, which is the remainder when we divide 25 by 12.

This example shows that under certain circumstances it makes sense to do arithmetic within the confines of a particular number (12 in this example). That is, we only keep track of the remainder when we divide by 12, and when we need to add two numbers, instead we just add the remainders. This method is quite efficient in the sense of keeping intermediate values as small as possible, and we shall see in later lectures how useful it can be.

More generally, we can define  $x \bmod m$  (in words: “ $x$  modulo  $m$ ”) to be the remainder  $r$  when we divide  $x$  by  $m$ . I.e., if  $x \bmod m = r$ , then  $x = mq + r$  where  $0 \leq r \leq m - 1$  and  $q$  is an integer. Thus  $29 \bmod 12 = 5$  and  $13 \bmod 5 = 3$ .

### Computation

If we wish to calculate  $x + y \bmod m$ , we would first add  $x + y$  and then calculate the remainder when we divide the result by  $m$ . For example, if  $x = 14$  and  $y = 25$  and  $m = 12$ , we would compute the remainder when we divide  $x + y = 14 + 25 = 39$  by 12, and get the answer 3. Notice that we would get the same answer if we first computed  $2 = x \bmod 12$  and  $1 = y \bmod 12$  and added the results modulo 12 to get 3. The same holds for subtraction:  $x - y \bmod 12$  is  $-11 \bmod 12$ , which is 1. Again, we could have obtained this directly by simplifying first, i.e.,  $(x \bmod 12) - (y \bmod 12) = 2 - 1 = 1$ .

This idea saves us even more effort with multiplication: to compute  $xy \bmod 12$ , we could first compute  $xy = 14 \times 25 = 350$  and then compute the remainder when we divide by 12, which is 2. Notice that we get the same answer if we first compute  $2 = x \bmod 12$  and  $1 = y \bmod 12$  and simply multiply the results modulo 12.

More generally, while carrying out any sequence of additions, subtractions or multiplications mod  $m$ , we get the same answer if we reduce any intermediate results mod  $m$ . This can considerably simplify the calculations.

## Set Representation

There is an alternative view of modular arithmetic which helps understand all this better. For any integer  $m$ , we say that  $x$  and  $y$  are *congruent modulo  $m$*  if they differ by a multiple of  $m$  or, in symbols,

$$x \equiv y \pmod{m} \Leftrightarrow m \text{ divides } (x - y).$$

For example, 29 and 5 are congruent modulo 12 because 12 divides  $29 - 5$ . We can also write  $22 \equiv -2 \pmod{12}$ . Notice that  $x$  and  $y$  are congruent modulo  $m$  iff they have the same remainder modulo  $m$ .

What is the set of numbers that are congruent to 0 (mod 12)? These are all the multiples of 12:  $\{\dots, -36, -24, -12, 0, 12, 24, 36, \dots\}$ . What about the set of numbers that are congruent to 1 (mod 12)? These are all the numbers that give a remainder 1 when divided by 12:  $\{\dots, -35, -23, -11, 1, 13, 25, 37, \dots\}$ . Similarly the set of numbers congruent to 2 (mod 12) is  $\{\dots, -34, -22, -10, 2, 14, 26, 38, \dots\}$ . Notice in this way we get 12 such sets of integers, and every integer belongs to one and only one of these sets.

In general, if we work modulo  $m$ , then we get  $m$  such disjoint sets whose union is the set of all integers: these are often called *residue classes mod  $m$* . We can think of each set as represented by the unique element it contains in the range  $(0, \dots, m - 1)$ . The set represented by element  $i$  would be all numbers  $z$  such that  $z = mx + i$  for some integer  $x$ . Observe that all of these numbers have remainder  $i$  when divided by  $m$ ; they are therefore congruent modulo  $m$ .

We can understand the operations of addition, subtraction and multiplication in terms of these sets. When we add two numbers, say  $x \equiv 2 \pmod{12}$  and  $y \equiv 1 \pmod{12}$ , it does not matter which  $x$  and  $y$  we pick from the two sets, since the result is always an element of the set that contains 3. The same is true about subtraction and multiplication. It should now be clear that the elements of each set are interchangeable when computing modulo  $m$ , and this is why we can reduce any intermediate results modulo  $m$ .

Here is a more formal way of stating this observation:

**Theorem 7.1.** *If  $a \equiv c \pmod{m}$  and  $b \equiv d \pmod{m}$ , then  $a + b \equiv c + d \pmod{m}$  and  $a \cdot b \equiv c \cdot d \pmod{m}$ .*

*Proof.* We know that  $c = a + k \cdot m$  and  $d = b + \ell \cdot m$  for integers  $k, \ell$ , so  $c + d = a + k \cdot m + b + \ell \cdot m = a + b + (k + \ell) \cdot m$ , which means that  $a + b \equiv c + d \pmod{m}$ . The proof for multiplication is similar and left as an exercise.  $\square$

---

*Exercise.* Complete the proof of Theorem 7.1 for multiplication.

---

What this theorem tells us is that we can always reduce any arithmetic expression modulo  $m$  to a number in the range  $\{0, 1, \dots, m - 1\}$ . As an example, consider the expression  $(13 + 11) \cdot 18 \pmod{7}$ . Using the above theorem several times we can write:

$$\begin{aligned} (13 + 11) \cdot 18 &\equiv (6 + 4) \cdot 4 \pmod{7} \\ &= 10 \cdot 4 \pmod{7} \\ &\equiv 3 \cdot 4 \pmod{7} \\ &= 12 \pmod{7} \end{aligned}$$

$$\equiv 5 \pmod{7}. \tag{1}$$

In summary, we can always do basic arithmetic (multiplication, addition, subtraction) calculations modulo  $m$  by reducing intermediate results modulo  $m$ . (Note that we haven't mentioned division: much more on that later!)

## 2 Exponentiation

Another standard operation in arithmetic algorithms (used heavily, e.g., in primality testing and RSA) is raising one number to a power modulo another number. I.e., how do we compute  $x^y \pmod{m}$ , where  $x, y, m$  are natural numbers and  $m > 0$ ? A naïve approach would be to compute the sequence  $x \pmod{m}, x^2 \pmod{m}, x^3 \pmod{m}, \dots$  up to  $y$  terms, but this requires time exponential in the number of bits in  $y$ . We can do much better using the trick of *repeated squaring*:

```
algorithm mod-exp(x, y, m)
  if y = 0 then return(1)
  else
    z = mod-exp(x, y div 2, m)
    if y mod 2 = 0 then return(z * z mod m)
    else return(x * z * z mod m)
```

This algorithm uses the fact that any  $y > 0$  can be written as  $y = 2a$  or  $y = 2a + 1$ , where  $a = \lfloor \frac{y}{2} \rfloor$  (which we have written as  $y \text{ div } 2$  in the above pseudo-code), plus the facts

$$x^{2a} = (x^a)^2; \quad \text{and}$$

$$x^{2a+1} = x \cdot (x^a)^2.$$

---

*Exercise.* Use the above facts to prove by induction on  $y$  that the algorithm always returns the correct value.

---

What is its running time? The main task here, as is usual for recursive algorithms, is to figure out how many recursive calls are made. But we can see that the second argument,  $y$ , is being (integer) divided by 2 in each call, so the number of recursive calls is exactly equal to the number of bits,  $n$ , in  $y$ . (The same is true, up to a small constant factor, if we let  $n$  be the number of decimal digits in  $y$ .) Thus, if we charge only constant time for each arithmetic operation (`div`, `mod` etc.) then the running time of `mod-exp` is  $O(n)$ . Note that this is *very* efficient: it means that we can handle exponents with (at least) thousands of bits!

In a more realistic model (where we count the cost of operations at the bit level), we would need to look more carefully at the cost of each recursive call. Note first that the test on  $y$  in the `if`-statement just involves looking at the least significant bit of  $y$ , and the computation of  $\lfloor \frac{y}{2} \rfloor$  is just a shift in the bit representation. Hence each of these operations takes only constant time. The cost of each recursive call is therefore dominated by the `mod` operation<sup>1</sup> in the final result. A fuller analysis of such algorithms is performed in CS170.

---

<sup>1</sup>You may want to analyze grade-school long-division for binary numbers to understand how long a `mod` operation would take. Since all arithmetic is being done `mod m`, the cost of this operation depends only on the number of bits in  $m$  and  $x$  (and not on  $y$ ).

### 3 Bijections

Before talking about division, we are going to have to take a little detour to talk about inverses. From linear algebra and calculus, you already have a lot of intuition about when inverses do and do not exist. Recall that a square matrix has an inverse only if it does not have a non-trivial nullspace. Why? Because if it has a non-trivial nullspace, it maps lots of vectors to the zero vector and there is no way to recover the information that was lost. The concept of bijection just allows us to formalize the mathematical intuition that we already have.

A function is a mapping from a set (called the *domain*) of inputs  $A$  to a set of outputs  $B$ : for input  $x \in A$ ,  $f(x)$  must be in the set  $B$ . To denote such a function, we write  $f : A \rightarrow B$ .

Consider the following examples of functions, where both functions map  $\{0, \dots, m-1\}$  to itself:

$$f(x) = x + 1 \pmod{m}$$

$$g(x) = 2x \pmod{m}$$

A *bijection* is a function for which every  $b \in B$  has a unique *pre-image*  $a \in A$  such that  $f(a) = b$ . Note that this consists of two conditions:

1.  $f$  is *onto*: every  $b \in B$  has a pre-image  $a \in A$ .
2.  $f$  is *one-to-one*: for all  $a, a' \in A$ , if  $f(a) = f(a')$  then  $a = a'$ .

Looking back at our examples, we can see that  $f$  is a bijection; the unique pre-image of  $y$  is  $y - 1$ . However,  $g$  is only a bijection if  $m$  is odd. Otherwise, it is neither one-to-one nor onto. The following lemma can be used to prove that a function is a bijection:

**Lemma:** For a finite set  $A$ ,  $f : A \rightarrow A$  is a bijection if there is an *inverse* function  $g : A \rightarrow A$  such that  $\forall x \in A$   $g(f(x)) = x$ .

*Proof.* If  $f(x) = f(x')$ , then  $x = g(f(x)) = g(f(x')) = x'$ . Therefore,  $f$  is one-to-one. Since  $f$  is one-to-one, there must be  $|A|$  elements in the range of  $f$ . This implies that  $f$  is also onto. (Can you see why? Prove this last fact for yourself. What kind of proof should you try?)  $\square$

The finiteness of the sets involved here make our life easier.

### 4 Inverses

We have so far discussed addition, multiplication and exponentiation. Subtraction is the inverse of addition and just requires us to notice that subtracting  $b$  modulo  $m$  is the same as adding  $-b \equiv m - b \pmod{m}$ .

What about division? This is a bit harder. Over the reals dividing by a number  $x$  is the same as multiplying by  $y = 1/x$ . Here  $y$  is that number such that  $x \cdot y = 1$ . Of course we have to be careful when  $x = 0$ , since such a  $y$  does not exist. Similarly, when we wish to divide by  $x \pmod{m}$ , we need to find  $y \pmod{m}$  such that  $x \cdot y \equiv 1 \pmod{m}$ ; then dividing by  $x$  modulo  $m$  will be the same as multiplying by  $y$  modulo  $m$ . Such a  $y$  is called the *multiplicative inverse* of  $x$  modulo  $m$ . In our present setting of modular arithmetic, can we be sure that  $x$  has an inverse mod  $m$ , and if so, is it unique (modulo  $m$ ) and can we compute it?

As a first example, take  $x = 8$  and  $m = 15$ . Then  $2x = 16 \equiv 1 \pmod{15}$ , so 2 is a multiplicative inverse of 8 mod 15. As a second example, take  $x = 12$  and  $m = 15$ . Then the sequence  $\{ax \pmod{m} : a = 1, 2, 3, \dots\}$  is

periodic, and takes on the values  $(12, 9, 6, 3, 0, 12, 9, 6, \dots)$ . [Exercise: check this!] Thus 12 *has no multiplicative inverse mod 15* since the number 1 never appears in this sequence.

This is the first warning sign that working in modular arithmetic might actually be very different from grade-school arithmetic. Two weird things are happening. First, no multiplicative inverse seems to exist for a number that isn't zero. (In normal arithmetic, the only number that has no inverse is zero.) Second, the "times table" for a number that isn't zero has zero showing up in it! So, e.g., 12 times 5 is equal to zero when we are considering numbers modulo 15. (In normal arithmetic, zero never shows up in the multiplication table for any number other than zero.)

So, when *does*  $x$  have a multiplicative inverse modulo  $m$ ? The answer is: if and only if the greatest common divisor of  $m$  and  $x$  is 1. Moreover, when the inverse exists it is unique. Recall that the *greatest common divisor* of two natural numbers  $x$  and  $y$ , denoted  $\gcd(x, y)$ , is the largest natural number that divides them both. For example,  $\gcd(30, 24) = 6$ . If  $\gcd(x, y)$  is 1, it means that  $x$  and  $y$  share no common factors (except 1). This is often expressed by saying that  $x$  and  $y$  are *relatively prime* or *coprime*.

**Theorem 7.2.** *Let  $m, x$  be positive integers such that  $\gcd(m, x) = 1$ . Then  $x$  has a multiplicative inverse modulo  $m$ , and it is unique (modulo  $m$ ).*

*Proof.* Consider the sequence of  $m$  numbers  $0, x, 2x, \dots, (m-1)x$ . We claim that these are all distinct modulo  $m$ . Since there are only  $m$  distinct values modulo  $m$ , it must then be the case that  $ax \equiv 1 \pmod{m}$  for exactly one  $a$  (modulo  $m$ ). This  $a$  is the unique multiplicative inverse of  $x$ .

To verify the above claim, suppose for contradiction that  $ax \equiv bx \pmod{m}$  for two distinct values  $a, b$  in the range  $0 \leq b \leq a \leq m-1$ . Then we would have  $(a-b)x \equiv 0 \pmod{m}$ , or equivalently,  $(a-b)x = km$  for some integer  $k$  (possibly zero or negative).

However,  $x$  and  $m$  are relatively prime, so  $x$  cannot share any factors with  $m$ . This implies that  $a-b$  must be an integer multiple of  $m$ . This is not possible, since  $a-b$  ranges between 1 and  $m-1$ .  $\square$

Actually it turns out that  $\gcd(m, x) = 1$  is also a *necessary* condition for the existence of an inverse: i.e., if  $\gcd(m, x) > 1$  then  $x$  has no multiplicative inverse modulo  $m$ .

---

*Exercise.* Verify this claim. [Hint: Assume for contradiction that  $x$  does have an inverse, say  $a$ . Then write down a (non-modular) equation that  $x, m$  and  $a$  must satisfy. Finally, derive a contradiction from the fact that  $x$  and  $m$  have a common factor  $d > 1$ .]

---

Since we know that multiplicative inverses are unique when  $\gcd(m, x) = 1$ , we shall write the inverse of  $x$  as  $x^{-1} \pmod{m}$ . Being able to compute the multiplicative inverse of a number is crucial to many applications, so ideally the algorithm used should be efficient. It turns out that we can use an extended version of Euclid's algorithm, which computes the gcd of two numbers, to compute the multiplicative inverse.

## 5 Computing Inverses: Euclid's Algorithm

Let us first discuss how computing the multiplicative inverse of  $x$  modulo  $m$  is related to finding  $\gcd(x, m)$ . For any pair of numbers  $x, y$ , suppose we could not only compute  $\gcd(x, y)$ , but also find integers  $a, b$  such that

$$d = \gcd(x, y) = ax + by. \tag{2}$$

(Note that this is not a modular equation; and the integers  $a, b$  could be zero or negative.) For example, we can write  $1 = \gcd(35, 12) = -1 \cdot 35 + 3 \cdot 12$ , so here  $a = -1$  and  $b = 3$  are possible values for  $a, b$ .

If we could do this then we'd be able to compute inverses, as follows. We first find integers  $a$  and  $b$  such that

$$1 = \gcd(m, x) = am + bx.$$

But this means that  $bx \equiv 1 \pmod{m}$ , so  $b$  is a multiplicative inverse of  $x$  modulo  $m$ . Reducing  $b$  modulo  $m$  gives us the unique inverse we are looking for. In the above example, we see that 3 is the multiplicative inverse of 12 mod 35. So, we have reduced the problem of computing inverses to that of finding integers  $a, b$  that satisfy equation (2). Remarkably, Euclid's algorithm for computing gcd's also allows us to find integers  $a$  and  $b$  as described above. So computing the multiplicative inverse of  $x$  modulo  $m$  is as simple as running Euclid's gcd algorithm on input  $x$  and  $m$ !

## Euclid's algorithm

If we wish to compute the gcd of two numbers  $x$  and  $y$ , how would we proceed? If  $x$  or  $y$  is 0, then computing the gcd is easy; it is simply the other number, since 0 is divisible by everything (although of course it divides nothing). The algorithm for other cases is ancient, and although associated with the name of Euclid, is almost certainly a folk algorithm invented by craftsmen (the engineers of their day) because of its intensely practical nature<sup>2</sup>.

This algorithm exists in cultures throughout the globe.

The algorithm for computing  $\gcd(x, y)$  uses the following theorem to eventually reduce to the case where one of the numbers is 0.

**Theorem 7.3.** *Let  $x \geq y > 0$ . Then  $\gcd(x, y) = \gcd(y, x \bmod y)$ .*

*Proof.* The theorem follows immediately from the fact that a number  $d$  is a common divisor of  $x$  and  $y$  if and only if  $d$  is a common divisor of  $y$  and  $x \bmod y$ . To see this, write  $x = qy + r$  where  $q$  is an integer and  $r = x \bmod y$ . Then, if  $d$  divides  $x$  and  $y$  then it also divides  $x$  and  $qy$ , and thus it also divides their difference  $r = x - qy$  (as we proved in Note 1). Conversely, if  $d$  divides  $y$  and  $r$  then it also divides  $qy$  and  $r$  and thus also their sum  $x = qy + r$ .  $\square$

Given this theorem, let's see how to compute  $\gcd(16, 10)$ :

$$\begin{aligned}\gcd(16, 10) &= \gcd(10, 6) \\ &= \gcd(6, 4) \\ &= \gcd(4, 2) \\ &= \gcd(2, 0) = 2\end{aligned}$$

In each line, we replace the pair of arguments  $(x, y)$  with  $(y, x \bmod y)$ , until the second argument becomes 0. At this point the gcd is just the first argument. By the theorem, each of these substitutions preserves the gcd.

---

<sup>2</sup>This algorithm is used for figuring out a common unit of measurement for two lengths. You can imagine how this is extremely important for building something up from a scale model. Different lengths in a design can be expressed as integer multiples of a common length, and then a new measuring stick can be found for the scaled-up design. We will see how the algorithm itself can be executed without literacy or symbolic notation. It is fundamentally *physical* in its intuition and you should figure out how this can be executed using threads. In the homework, you will see how this algorithm reveals the secret hidden in plain sight within the Pentagram. The fact that Euclid's algorithm deals naturally with real numbers is also important in understanding the topic of continued fractions — a topic important in the understanding of approximations and numerical computing.

This algorithm can be written recursively as follows. The algorithm assumes that the inputs are natural numbers  $x, y$  satisfying  $x \geq y \geq 0$  and  $x > 0$ .

```
algorithm gcd(x, y)
  if y = 0 then return(x)
  else return(gcd(y, x mod y))
```

**Theorem 7.4.** *The algorithm above correctly computes the gcd of  $x$  and  $y$ .*

*Proof.* Correctness is proved by (strong) induction on  $y$ , the smaller of the two input numbers. For each  $y \geq 0$ , let  $P(y)$  denote the proposition that the algorithm correctly computes  $\text{gcd}(x, y)$  for all values of  $x$  such that  $x \geq y$  (and  $x > 0$ ). Certainly  $P(0)$  holds, since  $\text{gcd}(x, 0) = x$  and the algorithm correctly computes this in the `if`-clause. For the inductive step, we may assume that  $P(z)$  holds for all  $z < y$  (the inductive hypothesis); our task is to prove  $P(y)$ . The key observation here is that  $\text{gcd}(x, y) = \text{gcd}(y, x \bmod y)$  — that is, replacing  $x$  by  $x \bmod y$  does not change the gcd. This was proved in Theorem 7.3. Hence the `else`-clause of the algorithm will return the correct value provided the recursive call `gcd(y, x mod y)` correctly computes the value  $\text{gcd}(y, x \bmod y)$ . But since  $x \bmod y < y$ , we know this is true by the inductive hypothesis! This completes our verification of  $P(y)$ , and hence the induction proof.  $\square$

What is the running time of this algorithm? We shall see that, in terms of arithmetic operations on integers, it takes time  $O(n)$ , where  $n$  is the total number of bits in the input  $(x, y)$ . This is again very efficient. The argument for this fact will be similar to the one we used earlier for exponentiation, but slightly trickier: it is obvious that the arguments of the recursive calls become smaller and smaller (because  $y \leq x$  and  $x \bmod y < y$ ). The question is, how fast?

The key point we will prove is that, in the computation of  $\text{gcd}(x, y)$ , after *two* recursive calls the first (larger) argument is smaller than  $x$  by at least a factor of two (assuming  $x > 0$ ). (Note that we can't argue much about what happens in just one call.) There are two cases:

1.  $y \leq \frac{x}{2}$ . Then the first argument in the next recursive call,  $y$ , is already smaller than  $x$  by a factor of 2, and thus in the next recursive call it will be even smaller.
2.  $x \geq y > \frac{x}{2}$ . Then in two recursive calls the first argument will be  $x \bmod y$ , which is smaller than  $\frac{x}{2}$ .

So, in both cases the first argument decreases by a factor of at least two every two recursive calls. Thus after at most  $2n$  recursive calls, where  $n$  is the number of bits in  $x$ , the recursion must stop. (Note that the first argument is always a natural number.)

Note that the above argument only shows that the *number of recursive calls* in the computation is  $O(n)$ . We can make the same claim for the running time if we assume that each call only requires constant time. Since each call involves one integer comparison and one mod operation, it is reasonable to claim that its running time is constant. In a more realistic model of computation, however, we should really make the time for these operations depend on the size of the numbers involved. This will be discussed in CS170.

## Extended Euclid's algorithm

Recall that, in order to compute the multiplicative inverse, we need an algorithm which also returns integers  $a$  and  $b$  such that:

$$\text{gcd}(x, y) = ax + by. \quad (3)$$

Then, in particular, when  $\gcd(x,y) = 1$  we can deduce that  $b$  is an inverse of  $y \bmod x$ .

Now since this problem is a generalization of the basic gcd, it is perhaps not too surprising that we can solve it with a fairly straightforward extension of Euclid's algorithm.

The following recursive algorithm *extended-gcd* follows the same recursive structure as Euclid's original algorithm, but keeps track of the required coefficients  $a, b$  in equation (3) as the recursion unwinds. Specifically, the algorithm takes as input a pair of natural numbers  $x \geq y$  as in Euclid's algorithm, and returns a triple of integers  $(d, a, b)$  such that  $d = \gcd(x, y)$  and  $d = ax + by$ :

```
algorithm extended-gcd(x, y)
  if y = 0 then return(x, 1, 0)
  else
    (d, a, b) := extended-gcd(y, x mod y)
    return((d, b, a - (x div y) * b))
```

In this algorithm,  $x \text{ div } y$  denotes the usual truncated integer division, written mathematically as  $\lfloor x/y \rfloor$ .

Here is the sequence of recursive calls (top row) along with the sequence of triples that they return (bottom row) for the same input  $(x, y) = (16, 10)$  as in our previous gcd example:

$$\begin{array}{ccccccccc} \text{e-gcd}(16, 10) & \longrightarrow & \text{e-gcd}(10, 6) & \longrightarrow & \text{e-gcd}(6, 4) & \longrightarrow & \text{e-gcd}(4, 2) & \longrightarrow & \text{e-gcd}(2, 0) \\ (2, 2, -3) & \longleftarrow & (2, -1, 2) & \longleftarrow & (2, 1, -1) & \longleftarrow & (2, 0, 1) & \longleftarrow & (2, 1, 0) \end{array}$$

*Exercise.* Check the above execution sequence.

The final triple returned is  $(d, a, b) = (2, 2, -3)$ , meaning that the gcd of  $(x, y) = (16, 10)$  is  $d = 2$ , and that it can be expressed as  $d = ax + by = 2 \cdot 16 - 3 \cdot 10$ , which we can easily see is correct. (In fact, the sequence of triples  $(d, a, b)$  returned each expresses  $d = 2$  as a linear combination of the corresponding inputs to that recursive call: so we can check that  $d = 1 \cdot 2 + 0 \cdot 0 = 0 \cdot 4 + 1 \cdot 2 = 1 \cdot 6 - 1 \cdot 4 = -1 \cdot 10 + 2 \cdot 6 = 2 \cdot 16 - 3 \cdot 10$ .) Since  $\gcd(16, 10) \neq 1$ , 10 doesn't have an inverse mod 16, so the coefficients  $a, b$  returned by the algorithm are not useful to us for computing inverses. However, the next exercise suggests an example where they allow us to read off the inverse, as described earlier.

*Exercise.* Hand-turn the algorithm yourself on the input  $(x, y) = (35, 12)$  and verify that it returns the triple  $(1, -1, 3)$ . Deduce that the inverse of 12 mod 35 is 3.

Let's now reverse-engineer the algorithm and understand why it works and why it was designed this way. In the base case ( $y = 0$ ), the algorithm returns the gcd value  $d = x$  as before, together with coefficients  $a = 1$  and  $b = 0$ ; clearly these satisfy  $ax + by = d$ , as required.

When  $y > 0$ , the algorithm first recursively computes values  $(d, a, b)$  such that  $d = \gcd(y, x \bmod y)$  and

$$d = ay + b(x \bmod y). \tag{4}$$

It then returns the triple  $(d, A, B)$ , where  $A = b$  and  $B = a - \lfloor x/y \rfloor b$ . Just as in our earlier analysis of the vanilla gcd algorithm, we know that the value  $d$  computed recursively will be equal to  $\gcd(x, y)$ . So the first component of the triple returned by the algorithm is correct.



What about the other two components,  $A$  and  $B$ ? From the specification of the algorithm, they should be integers that satisfy

$$d = Ax + By. \tag{5}$$

To figure out what  $A$  and  $B$  should be in terms of the previously returned values  $a$  and  $b$ , we can rearrange equation (4), as follows:

$$\begin{aligned} d &= ay + b(x \bmod y) \\ &= ay + b(x - \lfloor x/y \rfloor y) \\ &= bx + (a - \lfloor x/y \rfloor b)y. \end{aligned} \tag{6}$$

(In the second line here, we have used the fact that  $x \bmod y = x - \lfloor x/y \rfloor y$  — check this!) Now compare this last equation (6) with equation (5): comparing coefficients of  $x$  and  $y$ , we see that we need to take  $A = b$  and  $B = a - \lfloor x/y \rfloor b$ . But this is exactly what the last line of the algorithm does, and this is why it works!

---

*Exercise.* Turn the above argument into a formal proof by induction that the algorithm extended-gcd is correct.

---

Since the extended gcd algorithm has exactly the same recursive structure as the vanilla version, its running time will be the same up to constant factors (reflecting the increased time per recursive call). So once again the running time on  $n$ -bit numbers will be  $O(n)$  arithmetic operations. This means that we can find multiplicative inverses very efficiently.

## Division in modular arithmetic

Now that we know how to compute the inverse of  $x$  modulo  $m$  (assuming that  $x$  and  $m$  are coprime), how can we use it to do arithmetic? The simplest scenario is solving a modular equation such as the following:

$$8x \equiv 9 \pmod{15}. \tag{7}$$

To solve the analogous equation  $8x = 9$  over the rational numbers, we would multiply both sides by  $8^{-1}$  to get  $x = 9/8$ . Let's do the same thing in arithmetic mod 15. Recall that the inverse of  $8 \pmod{15}$  is 2 (since  $2 \cdot 8 = 16 \equiv 1 \pmod{15}$ ). Hence we can multiply both sides of equation (7) by  $8^{-1} \equiv 2$  to get

$$x \equiv 18 \equiv 3 \pmod{15}.$$

I.e., the solution to the modular equation (7) is  $x = 3$ , and this solution is unique modulo 15.

## 6 Chinese Remainder Theorem

It is worth stepping back for a moment and looking at what the EGCD revealed to us. It said that the GCD could be expressed as  $ax + by$  for two numbers  $x, y$ . To interpret this, we can imagine the number line, starting at zero and stretching out infinitely in both directions. Imagine that we are only allowed to take steps that are either  $x$  or  $y$  long. So, if  $x = 5$  and  $y = 7$ , then we can either move to the right or left by 5 units or 7 units. Suppose we start at zero, and want to know everywhere we can reach by taking a sequence of such moves.

Intuitively, if we can reach a number  $z$ , we can reach any multiple of  $z$  by simply repeating the steps it took to get to  $z$  over and over again. The fact that we can execute the steps of the Euclid's GCD algorithm tells

us that anything we can reach by taking steps of  $x$  and  $y$  must share all the common factors of  $x$  and  $y$ . This means that we can only reach any multiple of the GCD of  $x$  and  $y$ . The set of points that we can reach with such operations is called a “lattice” and this lattice-width interpretation of the GCD is interesting<sup>3</sup>.

When the GCD is 1, it means that we can reach all points on the integer lattice in this manner. Based on your experience with linear algebra, you should notice a very striking intellectual “rhyme” with the ideas of a basis and span. When their GCD is 1, it is as though the numbers  $x$  and  $y$  span all the integers<sup>4</sup>. The Chinese Remainder Theorem (CRT) can be interpreted as a way to make this interpretation even more striking.

Suppose we wanted to understand all the numbers mod  $pq$  where  $p$  and  $q$  are relatively prime to each other. If we had to arrange these numbers onto a sheet of paper, how would we do so? Going back to elementary school, it is natural to associate a product  $pq$  with a rectangle:  $p$  long on one side and  $q$  long on the other. So now, we know that we can place the  $pq$  numbers from 0 to  $pq - 1$  on this rectangle. But how? In what order? Given a number, how can you find its “x-coordinate” as something from  $0, 1, \dots, p - 1$  and its “y-coordinate” as something from  $0, 1, \dots, q - 1$ ? The natural first guess is to take a number  $z$  and just compute  $z \bmod p$  and  $z \bmod q$  to get two “coordinates” for  $z$ .

At this point, it is very useful to do a little exercise for yourself. Suppose  $p = 3$  and  $q = 5$  and just place all the numbers from 0 to 14 on this grid. You will see the coordinates as  $0 = (0, 0), 1 = (1, 1), 2 = (2, 2), 3 = (0, 3), 4 = (1, 4), 5 = (2, 0), 6 = (0, 1), 7 = (1, 2), 8 = (2, 3), 9 = (0, 4), 10 = (1, 0), 11 = (2, 1), 12 = (0, 2), 13 = (1, 3), 14 = (2, 4)$ . When writing them out, you will see that all the numbers lie on a diagonal line that wraps around the rectangle until it fills it. Notice that no two numbers from 0 to 14 have the same coordinates. Furthermore, notice that doing component-wise mod  $(3, 5)$  addition on the coordinates corresponds to doing mod 15 addition on the numbers themselves. Perhaps more interestingly, doing component-wise mod  $(3, 5)$  multiplication on the coordinates corresponds to doing mod 15 multiplication on the numbers themselves. (E.g.,  $3 * 4 = 12$  and  $(0, 3) * (1, 4) \equiv (0, 2)$ ). This means that operations can be equivalently performed component-wise in the tuple-representation.

Furthermore, we notice that there are two special tuples  $(1, 0) = 10$  and  $(0, 1) = 6$ . The corresponding numbers act like “orthonormal basis elements” (like the standard basis) do in linear algebra. They provide an easy way to map from coordinates back to numbers. So  $(a, b)$  in coordinates represents the same number as  $10a + 6b \bmod 15$ . For example,  $(2, 1) \rightarrow 20 + 6 = 26 \equiv 11 \pmod{15}$ . So, not only can we easily move from numbers to coordinates (by just taking mods), we can also easily move from coordinates to numbers (by using these special basis elements). Before we state the general form of the Chinese Remainder Theorem, it is useful to observe that the basis element 10 corresponding the first coordinate (obtained by modding by 3) is a multiple of the other modulus 5. This has to be true because its representation in coordinates is designed to have a zero in that other coordinate. Similarly, 6 corresponds to the second coordinate (obtained by modding by 5) and is a multiple of 3.

With this example in hand, we are ready to generalize and to state the result more formally.

**Chinese Remainder Theorem:** Let  $n_1, n_2, \dots, n_k$  be positive integers that are coprime to each other. Then, for any sequence of integers  $a_i$  there is a unique integer  $x$  between 0 and  $\prod_{i=1}^k n_i$  that satisfies the congru-

---

<sup>3</sup>This interpretation also makes short work of the classic family of puzzles of the form “you have a 5 oz cup and a 7 oz cup, an infinite reservoir of water, and a unlimited size mixing bowl. Can you manage to pour exactly  $z$  oz of water into a jar?” Do you see how such puzzles can be solved using EGCD?

<sup>4</sup>And when the GCD is 2, we can reach all even numbers. The even numbers behave in a way analogous to a subspace in linear algebra.

ences:

$$x \equiv a_1 \pmod{n_1} \tag{8}$$

$$\vdots \tag{9}$$

$$x \equiv a_i \pmod{n_i} \tag{10}$$

$$\vdots \tag{11}$$

$$x \equiv a_k \pmod{n_k} \tag{12}$$

Moreover this integer  $x$  can be found:

$$x = \left( \sum_{i=1}^k a_i b_i \right) \bmod N \tag{13}$$

where  $N = \prod_{i=1}^k n_i$  and the “basis” numbers  $b_i$  are found using the formula  $b_i = \frac{N}{n_i} \left( \frac{N}{n_i} \right)^{-1}$  where  $\left( \frac{N}{n_i} \right)^{-1}$  denotes the multiplicative inverse  $\pmod{n_i}$  of the integer  $\frac{N}{n_i}$ .

**Proof:** The only question in being able to apply the formulas is to make sure that  $\left( \frac{N}{n_i} \right)^{-1}$  exists. To verify this, we first notice that  $\frac{N}{n_i} = \prod_{j \neq i} n_j$  is a nonzero integer that is coprime to  $n_i$  since by construction, they can share no common factors. So the multiplicative inverse exists. This means that the formula is indeed computable and because it involves modding by  $N$ , it clearly gives rise to an  $x$  between 0 and  $N - 1$ .

To see that this  $x$  solves the system of congruences, we need to take  $x \bmod n_i$  and see what happens. First notice that  $\frac{N}{n_r} = \prod_{j \neq r} n_j$  is congruent to 0 when we mod by  $n_i \neq n_r$ . This means that:

$$\begin{aligned} x \bmod n_i &= \left( \left( \sum_{i=1}^k a_i b_i \right) \bmod N \right) \bmod n_i \\ &= \left( \sum_{i=1}^k a_i b_i \right) \bmod n_i \\ &= a_i b_i \bmod n_i \\ &= a_i \left( \frac{N}{n_i} \left( \frac{N}{n_i} \right)^{-1} \right) \bmod n_i \\ &= a_i \bmod n_i \end{aligned}$$

where the last quality used the definition of multiplicative inverse and the second equality used the fact that modding by a product and then by one of terms in that product is the same as just modding by that single term.

The above establishes that  $x \equiv a_i \pmod{n_i}$  and so  $x$  does indeed solve the system of congruences. To see that it is unique, we have two arguments that we could use. The simplest argument is by counting. There are  $N = \prod_{i=1}^k n_i$  possible values for the  $(a_1, a_2, \dots, a_k)$  tuples and the  $N$  numbers from 0 to  $N - 1$  each land in exactly one of these. If two landed in one bin, then that means that another bin must be empty. But we can construct an  $x$  corresponding to that bin and so it cannot be empty. This means that there must be a bijection from the coordinate tuples  $(a_1, a_2, \dots, a_k)$  and the  $N$  numbers from 0 to  $N - 1$ .

Alternatively, suppose that some  $y$  also solves these congruences. Consider  $z = y - x$ . Clearly  $z \bmod n_i$  is zero for all the  $n_i$ . This means that  $z$  is a multiple of  $n_i$  for each  $i$  and since they are all coprime,  $z$  is a multiple of  $N$ , their product. But the difference of two numbers ranging from 0 to  $N - 1$  must have an absolute value of at most  $N - 1$ . This means that the only multiple of  $N$  that  $z$  can be is 0. This means that  $y = x$  and so indeed, the given solution is unique. ♠

The Chinese Remainder Theorem (CRT) is a very powerful tool since it lets us move between numbers and their coordinates for the purpose of doing computations. Although stated for moduli that are all coprime, it can be extended to moduli  $n_i$  that are not coprime. However, in those cases, one has to be more careful. First, the range of numbers that we are interested in now is the Least-Common-Multiple (LCM) of the  $n_i$  values. Second, we must beware of inconsistent congruences. For example, we cannot simultaneously be congruent to 1 (mod 2) and be congruent to 2 (mod 6). In general,  $a_i \equiv a_j \pmod{\gcd(n_i, n_j)}$  must hold for a pair of congruences to be consistent<sup>5</sup>. You might be tempted to just use the formulas above with  $N = LCM(n_1, n_2, \dots, n_k)$ , but that is not quite enough<sup>6</sup>.

The homework has problems that will help you discover for yourself how the CRT can be very useful in solving problems.

---

<sup>5</sup>Since we can just mod both sides of both congruences by the GCD of  $n_i$  and  $n_j$  to get a congruence mod the GCD. If these two disagree, then the system of equations is clearly inconsistent.

<sup>6</sup>Instead, you can proceed by turning all congruences into statements about remainders mod prime powers. For every congruence that involves a composite modulus, just replace it with the equivalent system of congruences in terms of the prime-power factors of the modulus. By the regular CRT, these are equivalent to the original congruence. Once this has been applied to all the congruences, you simply have to discard redundant information. The rule is simple: keep only the congruence involving the largest power of any given prime. All the congruences for smaller powers are redundant. At this point, you have expressed the original congruences into a set of canonical congruences in terms of the prime factorization of the LCM of the original moduli.