

Recursion

Recursion means to define something in terms of itself. The power of recursion lies in the possibility of defining an infinite set of objects by a finite statement. For instance, we can define a (rooted) *tree* as follows: a tree consists of a root node, which is connected to the roots of zero or more trees.

In computer science, recursion means a function that calls itself (directly, or indirectly). Recursion is a powerful and natural problem solving technique:

Humans solve difficult problems by dividing them into easier subproblems, and then solve these subproblems. Sometimes it turns out that the “easier subproblems” are actually the same problem we originally started with, but on a smaller or easier input. In this case, there is no separate function for the subproblem, and instead our main function calls itself *recursively* to solve the subproblem.

We can loosely define recursion like this:

- If the problem is small or simple enough, solve it directly (*base case*).
- Otherwise, divide the problem into one or more simpler instances of the same problem, solve these recursively, and combine the solutions.

Printing a number in any base

Let’s start with a simple example: How do you find the representation of the number 83790 in base 8? More generally, how do you find the base- b representation of a given number n ?

It is easy to find the *last digit* of this representation: it is $n \bmod b$ (`n % b` in Scala). The remaining digits are then the base- b representation of n/b .

The following method implements this recursive strategy:

```
def printInt(n : Int, b : Int)
{
  if (n >= b)
    printInt(n / b, b)
  print(n % b)
}
```

Note that there is a base case: If the number is less than b , we do not make a recursive call and actually print the number n directly.

To design a recursive function, do not try to mentally simulate how the program is executed. Instead, reason about your method by assuming that the recursive call works correctly, and then argue that the method itself will do the right thing.

If you find the recursive call to `printInt` confusing, it’s helpful to imagine that someone else is going to solve the simpler problems.

If your problem is to bake a cake, you would partition it into subproblems: break 5 eggs, warm some butter, measure flour, etc. You could give each of the subproblems to a different little helper (let’s assume there are some fairies in your dormitory that want to help you). When your problem is recursive, use the same thinking: Your *only* task is to *simplify* the original problem, or to solve it directly when simplification is either unnecessary or impossible. The *Recursion Fairy* will magically take care of the simpler subproblems.

In mathematics, recursion appears all the time, in definitions such as the definition of trees above (or, for instance, in the following definition of the natural numbers \mathbb{N} : *0 is a natural number; if n is a natural number, then $n + 1$ is a natural number*).

Mathematical induction is a form of recursion, and indeed, when we want to prove that a recursive method works correctly, we use a proof by induction. For example, to prove that our method `printInt` is correct, we would use induction on k , where k is the number of digits of n when written in base b (in other words, $k = \lceil \log_b(n + 1) \rceil$).

- Base case $k = 1$: If n has only one digit, then $0 \leq n < b$, and so we are in the base case of `printInt`, which prints one digit correctly.

- Inductive step $k > 1$: We assume that `printInt` works correctly for numbers with less than k digits. Consider now a number n with k digits in base b . Then n/b has $k - 1$ digits. By the inductive assumption, this means that the recursive call `printInt(n/b)` correctly prints n/b . The method then prints the last digit, and therefore the number n has been printed correctly.

In a recursive algorithm, there must be no infinite sequence of reductions to 'simpler' and 'simpler' subproblems. Eventually, the recursive reductions must stop with an elementary base case that can be solved by some other method; otherwise, the recursive algorithm will never terminate.

This means that, first, we have to remember to include a base case. Second, the recursive calls must in some sense be "easier" than the original call, so that progress is made and the base case is eventually reached.

Towers of Hanoi

The Towers of Hanoi puzzle was first published by the mathematician Fran cois Édouard Anatole Lucas in 1883, under the pseudonym 'N. Claus (de Siam)' (an anagram of 'Lucas dAmiens'). The following year, Henri de Parville described the puzzle with the following remarkable story:

In the great temple at Benares beneath the dome which marks the centre of the world, rests a brass plate in which are fixed three diamond needles, each a cubit high and as thick as the body of a bee. On one of these needles, at the creation, God placed sixty-four discs of pure gold, the largest disc resting on the brass plate, and the others getting smaller and smaller up to the top one. This is the Tower of Bramah. Day and night unceasingly the priests transfer the discs from one diamond needle to another according to the fixed and immutable laws of Bramah, which require that the priest on duty must not move more than one disc at a time and that he must place this disc on a needle so that there is no smaller disc below it. When the sixty-four discs shall have been thus transferred from the needle on which at the creation God placed them to one of the other needles, tower, temple, and Brahmins alike will crumble into dust, and with a thunderclap the world will vanish.

Of course, being good computer scientists, we read this story and immediately substitute n for the hardwired constant sixty-four. How can we move a tower of n disks from one needle to another, using a third needles as an occasional placeholder, never placing any disk on top of a smaller disk? The trick to solving this puzzle

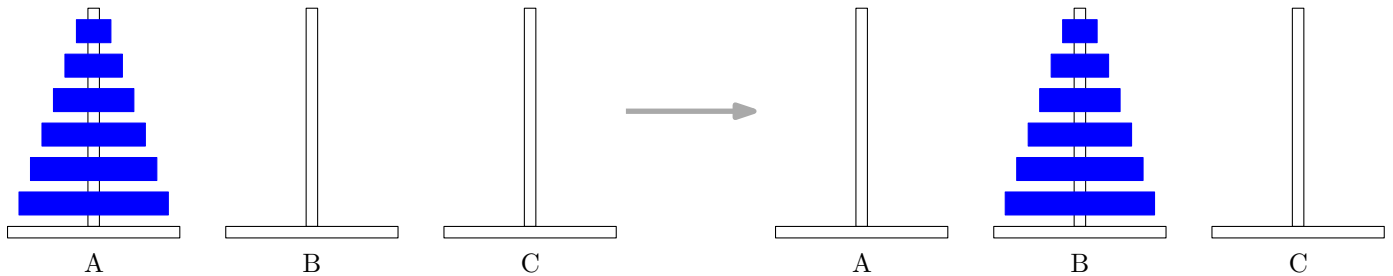


Figure 1: Towers of Hanoi

is to think recursively. Instead of trying to solve the entire puzzle all at once, let's concentrate on moving just the largest disk. We can't move it at the beginning, because all the other disks are covering it; we have to move those $n - 1$ disks to the third needle before we can move the n th disk. And then after we move the n th disk, we have to move those $n - 1$ disks back on top of it. So now all we have to figure out is how to...

STOP!! That's it! We're done! We've successfully solved the n -disk Tower of Hanoi problem by solving two instances of the $(n - 1)$ -disk Tower of Hanoi problem, which we can gleefully hand off to the Recursion Fairy (or, to carry the original story further, to the junior monks at the temple).

Our algorithm does make one subtle but important assumption: *There is a largest disk*. In other words, our recursive algorithm works for any $n \geq 1$, but it breaks down when $n = 0$. We must handle that base

case directly. Fortunately, the monks at Benares are quite adept at moving zero disks from one needle to another in no time at all, and we arrive at the following code:

```
def solveHanoi(n : Int, source : Char, destination : Char, spare : Char)
{
  if (n > 0) {
    solveHanoi(n-1, source, spare, destination)
    println("Move disc " + n + " from " + source + " to " + destination)
    solveHanoi(n-1, spare, destination, source)
  }
}

solveHanoi(n, 'A', 'B', 'C')
```

It is tempting to think about how all those smaller disks get moved—in other words, what happens when the recursion is unfolded—but it's not necessary. In fact, for more complicated problems, unfolding the recursive calls is merely *distracting*. Our only task is to reduce the problem to one or more simpler instances, or to solve the problem directly if such a reduction is impossible. Our algorithm is trivially correct when $n = 0$. For any $n \geq 1$, the Recursion Fairy correctly moves the top $n - 1$ disks, so our algorithm is clearly correct.

How many moves are needed to solve the problem for n disks? We can quickly check the answer for $n = 0, 1, 2, 3$, and obtain 0, 1, 3, 7 moves. We therefore guess that the number of moves for n disks must be $2^n - 1$, and prove this claim by mathematical induction:

- Base case ($n = 0$ disks): The number of moves is $0 = 2^0 - 1$, and so the claim is true.
- The inductive step: Let $n > 0$, and assume (*Induction hypothesis*) that the number of moves for $n - 1$ disks is $2^{n-1} - 1$. To move n disks, we first have to move the $n - 1$ top disks from A to C. By the induction hypothesis, this takes $2^{n-1} - 1$ moves. Then we have to move the largest disk (1 move), and finally we have to move the $n - 1$ top disks from C to B. This takes again $2^{n-1} - 1$ moves by the induction hypothesis. So the total number of moves is

$$(2^{n-1} - 1) + 1 + (2^{n-1} - 1) = 2^n - 1,$$

and we have proven the claim for n disks.

Fibonacci numbers

The Fibonacci numbers F_n are defined as follows:

$$\begin{aligned} F_0 &= 0 \\ F_1 &= 1 \\ F_n &= F_{n-1} + F_{n-2} \quad \text{for } n > 1 \end{aligned}$$

This immediately gives us a recursive method for computing Fibonacci numbers:

```
def fib(n : Int) : Long = {
  if (n == 0)      // base cases
    0
  else if (n == 1)
    1
  else
    fib(n - 1) + fib(n - 2)
}
```

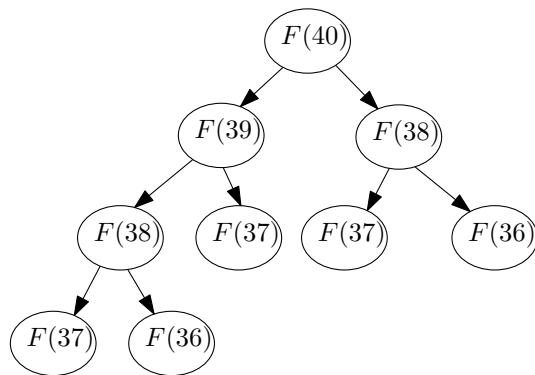


Figure 2: Recursive calls for computing `fib(40)`

It turns out that this method becomes where slow when $n \approx 40$. To see why, we draw a tree that shows all the calls and recursive calls to the function `fib`—see Fig. 2. As we can see in the tree, computing `fib(40)` means that we compute F_{38} two times, F_{37} three times, F_{36} five times—and so on, so that F_{41-k} is computed F_k times. Since the Fibonacci numbers grow exponentially, this procedure quickly becomes too slow. The problem in this example is that the recursive routine performs *redundant* calculations, that is, calculations that have already been done before. This example shows that recursion—blindly applied—is not always appropriate.

The solution is to either use an array to store all previous Fibonacci numbers (so no recursion is necessary), or to return both F_n and F_{n-1} from the function `fib(n)` (see example code `fibonacci2.scala` and `fibonacci3.scala`).

Stack frames

(We’ll discuss this a bit later in the course.)

The Scala runtime system (the Java virtual machine, or JVM) stores information in several separate memory areas, in particular the *heap* and the *runtime stack*.

As we discussed before, the *heap* stores *all objects*.

The *runtime stack* contains *stack frames* (also known as an *activation records*). Each stack frame stores the local variables of a method (which include the parameters of the method).

When a method is called, JVM creates the *stack frame* for the method, and pushes it onto the runtime stack. This method can call another method, which in turn calls another one, and so on. At any time, the stack frame at the top of the runtime stack corresponds to the currently executing method. When a method returns, its stack frame is popped from the runtime stack, and execution returns to the calling function.

A stack is the right data structure for the storage of stack frames, since at any time, only one method is actively executing, so we only need access to one stack frame. When this method returns, the method that called it becomes active again—which is exactly the one that becomes the top of the stack when we pop one stack frame.

Said differently, the starting and ending times of methods form a nicely nested pattern of opening and closing symbols.

When a method calls itself recursively, the runtime stack will hold multiple stack frames for the same method. Again, the top one is the one currently executing. For instance, calling `factorial(3)` for this method:

```

def factorial(n: Int): Long = {
  if (n <= 1)    // base case
    1
  else

```

```

    n * factorial(n - 1)
}

```

will lead to three stack frames for `factorial` being on the stack:

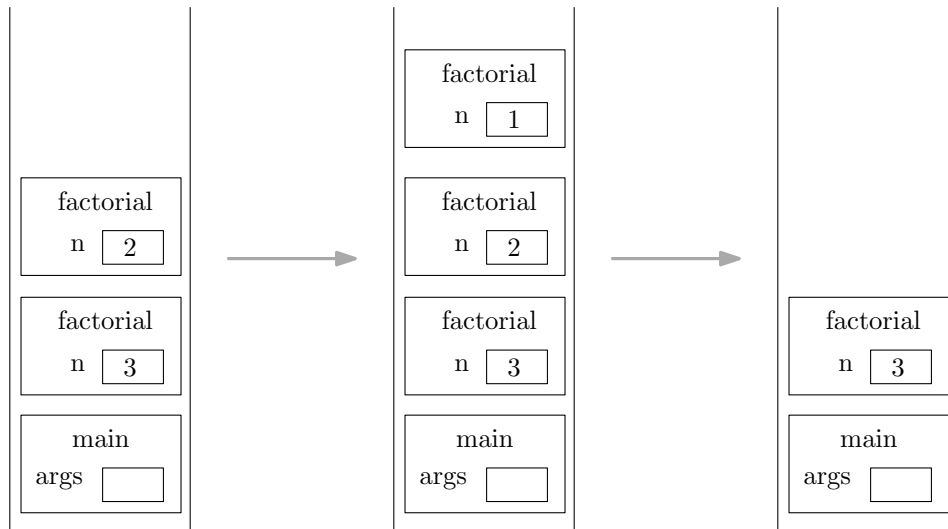


Figure 3: A call stack state: `factorial(3)`

When an exception is thrown, the runtime system uses the call stack to find the method where execution can continue. It pops stack frames until a stack frame is found that indicates that it catches the exception.

The `java.lang` package has a method `Thread.dumpStack()` that prints a list of the methods on the stack (but it doesn't print their local variables). This method can be convenient for debugging—for instance, when you're trying to figure out which method called another method with illegal parameters that made it crash. Another useful method is the `printStackTrace()` of `Exception` objects, which again allows you to find out what the contents of the runtime stack was when an exception happened.

Most operating systems give a program enough stack space for a few thousand stack frames. If you use a recursive procedure to walk through a million-node list, the system will try to create a million stack frames, and the stack will run out of space. The result is a run-time error. You should use iteration instead of recursion when the recursion is very deep.