Concepts of Programming Languages

Lecture Notes: Week 6—Tail Recursion

Stefan Mitsch

School of Computing, DePaul University smitsch@depaul.edu

TAIL RECURSION

We well discuss tail recursion in terms of the details of the implementation of function calls. Remember that function calls store local data on the *stack* in the form of activation records (or stack frames). When we call a function, a new activation record is pushed onto the stack, on return we pop the topmost activation record and return to the caller, whose activation record becomes then active (new topmost on the stack). Since function calls and returns cannot be arbitrarily nested (a callee is called after its caller, and returns before it), a stack is the most natural data structure for storing activation records. Since memory in a computer is finite, however, there are limits to the size of the call stack, which can cause problems with deep recursion.

Let's consider a recursive definition of a count-down in C below:

```
int count_down (int x) {
     if (x == 0) {
2
       return o;
3
       else {
       return i + count_down (x - i);
5
6
   }
   int main (int argc, char ** argv) {
9
     long num = strtol (argv[1], NULL, 10);
ю
     count_down (num);
     return o;
13
```

For some arguments this code will succeed, for others it produces a segmentation fault when the computer runs out of stack size. Note that, in C, this behavior depends on the operating system and on the execution context: we won't necessarily find a limit, but will see behavior that some calls fail and then some with larger numbers succeed. The tendency is that the higher the argument the more often we will see the program fail. In Linux, we can modify some of those limits with shell arguments. In interpreted languages, such as Scheme, the same behavior occurs, but just like in C, it depends on the runtime environment. In Java, you may also have encountered a runtime exception of type StackOverflowError; this usually happens when we make a programming mistake in a recursive algorithm, but it can also happen in perfectly sensible code (for some definition of sensible) that just happens to exhaust the stack limits (e.g., code that has lots of backtracking options).

Intuitively, in the recursive definition of countDown in Scala below, each (i+...) represents a new activation record, because we must store the intermediate 1+ that still need to be summed up.

Let's take a more detailed look under the hood to find out why we have a call stack.

Tail-recursive calls have a special structure that has all recursive calls in tail position, which means that all recursive calls happen as the final step on the call stack. In assembly language, that corresponds to jumps at the very end of each recursive call, rather than a jump and then another operation after. This opens up the possibility for an optimization called *tail call optimization* that performs the work in a loop without creating new activation records on the stack. Such an optimization is vital when algorithms have potentially exponential work to perform, which can exhaust stack limits quickly. The optimization includes mutually recursive calls (f calls g, which calls f).

Tail-Recursive Implementations

The C code below computes the sum over a list of integers in a tail-recursive method.

```
typedef struct node node;
   struct node { int item; node *next; };
2
   int sum_aux (node *data, int result) {
     if (!data) {
6
        return result;
     } else {
        return sum aux (data -> next, result + data -> item);
9
   }
ю
11
   int sum (node *data) {
12
     return sum_aux (data, o);
   }
14
```

When inspecting the assembly language of that code, we can clearly see the loop behavior:

```
7 movq 8(%rdi), %rdi
8 testq %rdi, %rdi
9 jne .L9
10 .L7:
11 rep
12 ret
```

Lines 5–9 advance dereference the pointer, then advance it, then test for being null and jump back to Line 5 if not null. If null, the program continues into .L7 and returns the result.

Let us now inspect some behavior of non-tail-recursive functions and to practice rewriting of functions that are not tail-recursive into functions that are.

```
def longList(n: Int) =
if n==0 then List(1)

else
val sublist = longList(n-1)
sublist ::: sublist
end longList
```

First, let's make sense of the code: the function longList takes an argument n to produce a result list of length 2^n . It does so by doubling the length of the list on every n-decrement step using :::.

Exercise 1 (Tail-recursive or not?). Why is this function not tail-recursive?

Think for yourself before you read on. The recursive call to longList is not the very last step in the else-branch of longList.

Exercise 2 (Space complexity?). What is the space complexity of the function?

 \P Think for yourself before you read on. Exponential in n

Exercise 3 (Runtime complexity?). How many recursive calls are made?

 \P Think for yourself before you read on. Linear in n

Exercise 4 (Bad or ugly?). Is it necessary to store the created sublist in a variable or could we just directly make recursive calls as part of the last statement?

Think for yourself before you read on. If we're not reusing the result we end up with the Fibonacci disaster of exponential in n calls to longList instead of linear in n.

In order to make it tail-recursive, we follow an approach resembling forward vs. backward processing of folds, as seen in Lecture "Folds".

First, a non-tail-recursive function computing the sum of the elements in a list.

Exercise 5 (Make it tail-recursive). How do we need to change the code above to implement sum in a tail-recursive way?

Think for yourself before you read on. To make it tail-recursive, we need to first compute a partial result from the current element and the previous partial result. For this, we need to introduce an additional argument to pass in partial results.

4 Stefan Mitsch

```
def sum(xs: List[Int], result: Int = 0) = xs match
case Nil => result
case x:: tail => sum(tail, result+x)
end sum
```