

Mathematical Structure for Programming

Conor McBride

May 8, 2024

Chapter 1

spuds and a lemon

Let's make some numbers! You'll need a bag of lemons, a sack of potatoes, and a load of sticks (e.g. cocktail sticks) which are pointy at both ends. (If you have suitable electrodes and wires, you can make batteries as well as number.)

There are two ways to get a number in your hand:

1. grab a lemon—it's a number all by itself;
2. if you already have a number in your hand, stab the thing in your hand with a stick (don't stab your hand), then press a potato onto the other end of the stick and hold onto that potato instead

So that means we can make (with your hand on the left)

lemon
spud-lemon
spud-spud-lemon
spud-spud-spud-lemon
⋮

(Incidentally, a crucial “let's pretend” in this game is that you can tell potatoes apart from lemons, but that blind to any difference between two potatoes or between two lemons.)

Now, if you only make numbers this way, starting only with a lemon and holding onto only the most recently joined potato, your numbers will all have an important property. **If you follow the sticks down the number, you will eventually get to the lemon.** No sneaky moves, like making a circular model from potatoes and sticks, or always adding on the next potato just in time.

If you're thinking “Those aren't numbers: that's just sculpture with groceries!”, then I have two things to tell you.

1. Our silly models with potatoes and lemons held together by sticks are as valid a representation of numbers as 5 or MMXXIV.

2. Count the spuds!

But how do you know you *can* count the spuds? How do you know that you won't just keep counting for ever? It's because you eventually get to the lemon.

1.1 addition as substitution

Adding these numbers is easy. If you're holding one of them in each hand, work your way along the left number until you reach the left lemon, then detach it, and stab the right number with the exposed stick, then work your way back out to the leftmost spud. In other words, *substitute* the right *number* for the left *lemon*

$$\begin{array}{r}
 \text{spud-spud-spud-spud-spud-lemon} \\
 + \\
 \text{spud-spud-spud-spud-spud-lemon} \\
 = \\
 \text{spud-spud-spud-spud-spud-spud-spud-spud-spud-spud-lemon}
 \end{array}$$

Note that the effectiveness of this procedure relies on the fact that you are guaranteed to find the left lemon.

It's also worth asking these three things:

1. What if the number in your left hand is a lemon? (You replace the lemon in your left hand with the number in your right, and the answer is exactly the right number.)
2. What if the number in your right hand is a lemon? (You replace the left lemon by the lemon in your right hand and the answer is indistinguishable from the left number you started with.)
3. What if you add three things? Does it matter whether you add left to middle first, or middle to right? (It doesn't matter. Both ways round, you end up with all the spuds you started with, in the same order (not that you can tell), and the rightmost lemon on the end.)

Chapter 2

seeing the trees

Computer programs are written in formal languages. You can't just say any old thing and expect a computer to make sense of it, even if these days they're quite willing to bullshit you that they understand. Computer programs tend to be written textually, as a sequence of *lines*, each of which is a sequence of *characters* (i.e., letters, digits, spaces, punctuation, emoji, etc), and that's not particularly helpful if you're trying to understand them—it's merely a convenient fit with our ancient methods for writing stuff down. Text doesn't have much obvious structure, but computer programs do. We need to learn to see, and to talk about how to see.

We have a language (the language of *grammars*) for talking about languages. Let me start with an example

$$\begin{array}{lcl} \langle number \rangle & ::= & \text{lemon} \\ & | & \text{spud-}\langle number \rangle \end{array}$$

How to read this? Where to begin? Which of these things is nearest to hand? You can't be expected to know until someone establishes the conventions.

The place to start is with $::=$, which you can pronounce 'is defined to be'. It's a variation on the theme of $=$, but with a particular spin. You're not being asked to agree with this equation, and there's no way you can check it. You're not supposed to think (or pretend) 'Oh I knew that!'. When I write this equation, I'm telling you something *new*, and asking you to put up with it for the time being.

The thing to the left of $::=$ is the name of a new *sort of thing*—by convention, the names of sorts of things are written in $\langle \cdot \rangle$ to distinguish them from *actual* things. The $\langle \cdot \rangle$ does the job of 'a' or 'an' in English—the indefinite article—we're talking about a generic class of things by imagining one, but not a special one. What I'm saying is 'I'm inventing a new sort of thing, namely a *number*, and I'm telling you of what a *number* may consist.'. Everything after $::=$ is the explanation of what I allow these *numbers* to be.

The next thing to pay attention to is the $|$ which you can read as 'or'. I'm giving you a choice of ways to make *numbers*.

The first choice I offer is ‘lemon’. See? In English prose I write it in quotation marks, because I mean the actual word ‘lemon’. It’s not a *variable* which could stand abstractly for a variety of things. It’s concretely what it says.

The second choice I offer is the concrete symbol ‘spud-’ which again (because no $\langle \cdot \rangle$) means only what it says, followed by $\langle number \rangle$. The latter indicates that any *number* you have made already can go in that place.

I didn’t offer any other choices. Implicit in this definition is that it is *exhaustive*. Spuds and lemons, that’s your lot! I’m saying “I’m inventing a new sort of thing, namely a *number*, which consists either of ‘lemon’ or of ‘spud-’ followed by another *number*.” and I’m also saying “You can’t possibly have *finished* making a *number* unless you eventually get to a ‘lemon’.”.

So what are these *numbers*?

lemon
spud-lemon
spud-spud-lemon
spud-spud-spud-lemon
⋮

But these examples and the ‘and so on’ indicated by ‘⋮’, don’t tell us precisely what *numbers* are, in general, whereas the grammar does.

The grammar does another useful thing, as well. It tells us how to perceive a *number* as a structure. We can draw a picture of how to see spud-spud-lemon as a *number*:

lemon
↑
spud- $\langle number \rangle$
↑
spud- $\langle number \rangle$
↑
 $\langle number \rangle$

I’ve drawn this picture, which is called a ‘parse tree’, with the bottom nearest to hand (they’re sometimes drawn the other way up), so you should read it bottom-to-top, like a detective, even though time went top-to-bottom¹. It records how the choices offered by the grammar can be employed to make sense of the text.

Old computer scientists *see* these trees when we *look* at text. We don’t even notice that we’re doing it, but we experience distress when it doesn’t just happen automatically and we look for someone else to blame. You are not born with this mode of perception, and it may take some practice to acquire it.

¹In the beginning was the lemon...

Chapter 3

squish a bunch of stuff

If you tap your foot five times, and then you tap your foot six times, you'll have tapped your foot eleven times. There's some kind of relationship between adding up numbers and repetitive action sequences.

If you have one bag of five potatoes and another of six, combining both bags into one will give you a bag of eleven potatoes. Bags of spuds don't behave this way because they were taught arithmetic at school. Arithmetic behaves this way because it helps us think about bags of spuds.

If you multiply two by five, that's ten; and two times six is twelve. Adding ten and twelve gives twenty-two, which is twice eleven. There's some kind of relationship between adding up numbers and adding up their doubles. (You may have heard this called "The Distributive Law", as if there were only one.)

If you raise two to the power five, that's thirty-two; and two to the six is sixty-four. I don't expect you to multiply numbers as large as those in your head, but you might add five and six to get eleven, then happen to know that two to the eleven is two thousand and forty-eight. (Having a head full of powers of two is an occupational hazard of programming computers.) Why does that trick work?

Exactly one of five and six is an odd number, and so is eleven. That's odd! It may seem trivial, but at least one of five and six is positive (as opposed to zero), and eleven is positive too. (The only way the sum of two counts can be zero is if they were both separately zero.)

My point is that

$$5 + 6 = 11$$

is not a random isolated truth. It has a relationship with a whole bunch of other truths about other senses of "combining stuff", which all go

Combining the fivey thing with the sixty thing gives you the eleveny thing.

and of course, five, six and eleven are overly concrete examples—what matters is that they are related by a truth about adding numbers, a truth which somehow

survives the translation to a truth about combining ny things, be they toe-taps, bags of spuds, or tests for oddness. Moreover, it wasn't only the things which mattered—it was also how they were combined. When we doubled we got a truth about adding, but when we two-to-the-powered we got a truth about multiplying, and when we tested for positivity we got a truth about “either or both”.

3.1 pictures of combining

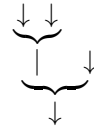
We can draw pictures of strategies for combining things like this

$$\begin{array}{c} 5 \quad 6 \\ \underbrace{\hspace{1cm}} \\ 11 \end{array}$$

Data flows from top to bottom. This component



has two inputs and one output. We can wire these components together to combine more things.



If we choose what sort of things we're combining (numbers, say) and how (adding, say), then by labelling the wires with data, we can make assertions about the results of combining things. I usually just write the data instead of the wire, so

$$\begin{array}{c} 5 \quad 6 \\ \underbrace{\hspace{1cm}} \\ 11 \quad 7 \\ \underbrace{\hspace{2cm}} \\ 18 \end{array} \quad \text{asserts} \quad 5 + 6 = 11 \quad \text{and} \quad 11 + 7 = 18$$

but it's allowed to label the same wire more than once as long as you label it with the same thing. Signals don't change value in the middle of a wire. The same calculation is indicated by

$$\begin{array}{c} 5 \quad 6 \quad 7 \\ \underbrace{\hspace{1cm}} \quad | \\ 11 \quad 7 \\ \underbrace{\hspace{2cm}} \\ 18 \\ | \\ 18 \end{array}$$

Now, some notions of combining things are equipped with “the sensible way to do nothing”. That’s given by a component with no inputs and one output

$$\begin{array}{c} \bullet \\ | \\ \downarrow \end{array} \quad \text{such that} \quad \begin{array}{c} \bullet \quad \cdot \\ \underbrace{\quad} \\ \cdot \end{array} = \begin{array}{c} \cdot \\ | \\ \cdot \end{array} = \begin{array}{c} \cdot \quad \bullet \\ \underbrace{\quad} \\ \cdot \end{array}$$

In other words, combining with nothing turns into wire: that the output must be the same as the input is what “doing nothing” means.

For adding numbers, 0 is the right way to do nothing,

$$\begin{array}{c} \bullet \\ | \\ 0 \end{array} \quad \text{so for any } n, \quad \begin{array}{c} \bullet \\ | \\ 0 \quad n \\ \underbrace{\quad} \\ n \end{array} = \begin{array}{c} n \\ | \\ n \end{array} = \begin{array}{c} n \quad 0 \\ \underbrace{\quad} \\ n \end{array}$$

It’s very useful to have a sensible way to do nothing. That way, we can look out of an aeroplane window and count the fish, or say how many spuds we have left once we’ve eaten them all. Our rules tell us that there’s only one sensible way to do nothing for any given way of combining. Imagine we had

$$\begin{array}{c} \bullet \\ | \\ a \end{array} \text{ and } \begin{array}{c} \bullet \\ | \\ b \end{array} \quad \text{then} \quad \begin{array}{c} b \\ | \\ c \end{array} = \begin{array}{c} \bullet \quad \bullet \\ | \quad | \\ a \quad b \\ \underbrace{\quad} \\ c \end{array} = \begin{array}{c} a \\ | \\ c \end{array} \quad \text{so} \quad a = c = b$$

That is, on the left, we do the left nothing to the right nothing, and on the right, we do the right nothing to the left nothing: our rules tell us we must end up with the same nothing. So if you see a “way of combining”, *look* for “the sensible way to do nothing”—you won’t find two, you might find one, and if there’s no sensible way to do nothing, that’s interesting, too.

When we’re combining a bunch of stuff, we are sometimes in the lucky position that it’s the sequence of *the stuff* which determines the result of combining them, not the structure of pairwise combinations we choose. If we add up the list [5, 6, 7], we get 18, whether we calculate it this way

$$\begin{array}{c} \underbrace{5 \quad 6}_{11} \quad 7 \\ \underbrace{\quad}_{18} \end{array} \quad \text{or this way} \quad \begin{array}{c} \underbrace{6 \quad 7}_{13} \\ 5 \quad \underbrace{\quad}_{18} \end{array} \quad \text{or even this way} \quad \begin{array}{c} \bullet \\ | \\ 7 \quad 0 \\ \underbrace{\quad} \\ 6 \quad 7 \\ \underbrace{\quad} \\ 5 \quad 13 \\ \underbrace{\quad}_{18} \end{array}$$

In general, we want

$$\begin{array}{c} \cdot \quad \cdot \\ \underbrace{\quad} \\ | \\ \cdot \end{array} = \begin{array}{c} \cdot \quad \cdot \\ \underbrace{\quad} \\ | \\ \cdot \end{array}$$

so that we can regroup the calculation without reordering the inputs and be sure of getting the same output (even though the intermediate data change). That's why we don't bother writing brackets in $5+6+7$. You can combine data onto a "running total" one at a time, or you can give half the data to a friend and combine your outputs afterwards. Knowing what doesn't matter makes it far easier to gain confidence about what does!

Of course, I haven't really demonstrated that only the sequence of the inputs determines the outputs, not the structure of the combination tree. Let me sketch one way to see it. Let me say that one of our diagrams is *listy* if it is always overbalanced as far to the right as possible, and has a nothing in the top right corner, i.e.

- inputs occur only on the left of $\begin{array}{c} \cdot \\ \cdot \end{array}$
- only inputs occur on the left of $\begin{array}{c} \cdot \\ \cdot \end{array}$

Think carefully about the difference in meaning made by the difference in word order. The first condition rules out

$$\begin{array}{c} 6 \\ | \\ 6 \end{array} \quad \text{and} \quad \begin{array}{c} 5 \quad 6 \\ \underbrace{\hspace{1cm}} \\ 11 \end{array}$$

in both cases because 6 is somewhere it is not permitted. The second condition rules out

$$\begin{array}{c} \bullet \\ | \\ 0 \quad 6 \\ \underbrace{\hspace{1cm}} \\ 6 \end{array} \quad \text{and} \quad \begin{array}{c} 5 \quad 6 \\ \underbrace{\hspace{1cm}} \\ 11 \quad 7 \\ \underbrace{\hspace{1cm}} \\ 18 \end{array}$$

because some left things are not inputs.

Each of our sequences of inputs can be put into a listy diagram because

- if you have no inputs, only $\begin{array}{c} \bullet \\ | \\ \cdot \end{array}$ is listy;
- if you have a first input x , you must make $\begin{array}{c} x \quad L \\ \underbrace{\hspace{1cm}} \\ t \end{array}$ where L is the listy diagram made with all but the first input.

But that's not enough, because we need to know that *any* diagram can be transformed into its listy counterpart by using the three rules we gave ourselves

$$\begin{array}{c} \bullet \\ | \\ \cdot \end{array} = \begin{array}{c} \cdot \\ | \\ \cdot \end{array} = \begin{array}{c} \cdot \\ \bullet \\ \cdot \end{array} \quad \begin{array}{c} \cdot \\ \cdot \end{array} = \begin{array}{c} \cdot \\ \cdot \end{array} = \begin{array}{c} \cdot \\ \cdot \end{array}$$

Here's how:

- \uparrow is already listy;

- \cdot can be made listy like this $\begin{array}{c} \cdot \\ | \\ \cdot \end{array}$;

- if you have a diagram $\begin{array}{c} D_0 \quad D_1 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$, you can make it listy by first making D_0 into listy L_0 and then making listy D_1 into listy L_1 , so you have $\begin{array}{c} L_0 \quad L_1 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$. Now *rotate* the whole thing into being listy like this:

- if you have $\begin{array}{c} \uparrow L_1 \\ | \\ t \end{array}$, turn it into L_1 (whose ‘total’ must already be t);
- if you have $\begin{array}{c} x \quad L'_0 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$ rotate it to $\begin{array}{c} L'_0 \quad L_1 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$, then keep rotating on the right.

So we’ve used all and only the rules we asked for. The rotation process effectively pastes L_1 over the \uparrow in the top right corner of L_0 .

There are lots of things to say about this “explanation”, some good, some bad. Does it make sense to you? How would you check it?

How do you know, for example, that “then keep rotating on the right” actually achieves something other than endless rotation? There’s a sense that

rotation of $\begin{array}{c} L_0 \quad L_1 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$ works by looking at L_0 , and if $L_0 = \begin{array}{c} x \quad L'_0 \\ | \\ t \end{array}$, then we keep rotating

with $\begin{array}{c} L'_0 \quad L_1 \\ \underbrace{\hspace{1cm}} \\ t \end{array}$, where the new diagram on the left, L'_0 is smaller than the old diagram on the left L_0 . As long as we can’t keep cutting smaller diagrams out of bigger ones forever, we’ll be fine—sooner or later, we’ll hit \uparrow . But is that explanation of the explanation just more waffle, or is it something we can codify?

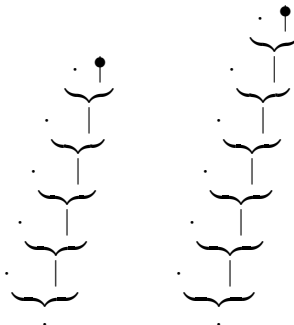
Another potential grumble is that the whole thing depends on inventing this notion of which diagrams are “listy”, and it looks like I just plucked that definition out of my arse. Where did it really come from?

As for the positives, for one thing, no numbers got added in the course of this narrative. We said which three rules we needed, and we deduced an “only the sequence of inputs determines the output” result for *any* notion of combining things which obeys those rules. If we want the same result for multiplication (where 1 is “the right way of doing nothing”), we know what we need to check.

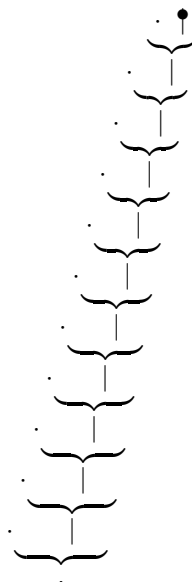
Moreover, our story about pictures of combining things involved inventing a way to combine *listy pictures* of combining *things*. To combine L_0 and L_1 , paste L_1 over the \uparrow in the top right corner of L_0 , or in other words, combine

the sequences by concatenating them—joining them up end-of-one-to-start-of-the-next in space. Here \uparrow is the “right way of doing nothing”, as it represents the empty sequence.

For example, here are the listy pictures for combining five things and six things, respectively:



and if you combine them, you get



which is, of course, the listy picture for combining eleven things, so perhaps numbers *did* get added, after all.

3.2 it’s called a “monoid”

While I’m trying to reduce the number of words by drawing more pictures, we are going to need a name for these “ways of combining a bunch of stuff”. When I do introduce terminology, I’ll try to make it the standard terminology. The

standard name is "monoid", which is Greek for "one-ish", as in "you can take any sequence of things and combine them into one thing". Let's review the definition.

A **monoid** is given by a type T , a value ε in T , and an operator \circ which takes two inputs from T and yields one output in T , satisfying the laws

$$\varepsilon \circ t = t \quad t \circ \varepsilon = t \quad (r \circ s) \circ t = r \circ (s \circ t)$$

We've been drawing

$$\begin{array}{c} \bullet \\ | \\ \varepsilon \end{array} \quad \begin{array}{c} s \quad t \\ \underbrace{\hspace{1cm}} \\ s \circ t \end{array}$$