

COMS30040

λ Types & -calculus

Lecture Notes for 2025/26

Steven Ramsay

School of Computer Science
University of Bristol

Contents

1	Proof: Propositional Reasoning	1
2	Proof: Quantifiers, Equality	9
3	Proof: Theories and Induction	15
4	Syntax: Terms and Binding	21
5	Syntax: Binding and Substitution	27
6	Semantics: One-Step Reduction	33
7	Semantics: Reduction and Confluence	39
8	Computability: Recursion	45
9	Computability: Effectiveness and Turing-Completeness	51
10	Computability: (Un)Decidability	59
11	Types: Judgements and the Type System	65
12	Types: Derivations and Typability	71
13	Types: Not Going Wrong	77
14	Types: The Curry-Howard Correspondence	83
15	Types: Constructive Proofs	89
<hr/>		
A	Proof Rules	93

1. Proof: Propositional Reasoning

In this first lecture, I want to introduce you to a way of thinking about mathematical proof as a kind of programming. The aim is to introduce you to proof (almost) from scratch, give you the rules you need to play the game and then look at examples of the rules in action.

1.1 Proof as a form of programming

Proofs are semi-formal arguments that make a case for the truth of some precise statement. The key word as far as proofs are concerned is *convincing*: a proof is a way of communicating your belief that something is true in a way that will convince another party.

Some proofs are more convincing than others. In mathematics, the statements that we want to give proofs for only have certain shapes — they are built from a few, standard logical operators. Over the years, mathematicians have developed equally standard approaches to dealing with the logical operators in proofs. By following these standard approaches, the benefit is that you will end up with a proof that will be convincing. Proofs of this kind are typically blocks of English text with a scattering of mathematical formulas.

It can be very helpful to think of the individual steps of a proof as being operations that manipulate the *proof state*. The proof state is a pair consisting of a set of *assumptions* and the current *goal*. The goal is the formula you are currently trying to prove and the assumptions are the resources that you have available in order to try to prove it. Each step of the proof (i.e. each application of one of the rules of natural deduction, signalled using an appropriate form of words) changes the proof state by either adding to the set of assumptions or changing the goal. For example, if we have assumptions A_1, \dots, A_k and we are aiming to prove $A \Rightarrow B$, then the magic words “assume A ”, which is a form of words used to signal a use of implication introduction, lead us to a proof state in which the assumptions are A_1, \dots, A_k, A and the goal is B . When you are reading a proof, it is very helpful to picture the proof state in your mind after each step of the proof.

The following example is a proof of a very simple fact from number theory: if two natural numbers n and m add up to zero, then m is not of the form $k + 1$ for any possible k .

Lemma 1.1. *For all $n, m, k \in \mathbb{N}$: if $n + m = 0$ then $m \neq k + 1$.*

Proof. Let n, m and k be natural numbers. Assume $n + m = 0$. To see that $m \neq k + 1$, we assume $m = k + 1$ and try to obtain a contradiction. Combining these assumptions we get $n + (k + 1) = 0$. By the associativity of addition this is the same as $(n + k) + 1 = 0$. Since $n + k$ is a natural number, we have $(n + k) + 1 \neq 0$. We have reached our contradiction. \square

Proof step	Assumptions	Goal
	$A,$ $B,$ C	$\forall nmk \in \mathbb{N}. n + m = 0 \Rightarrow m \neq k + 1$
Let n, m and k be natural numbers.	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N}$	$n + m = 0 \Rightarrow m \neq k + 1$
Assume $n + m = 0$.	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0$	$m \neq k + 1$
To see that $m \neq k + 1$, we assume $m = k + 1$ and try to obtain a contradiction.	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0,$ $m = k + 1$	\perp
Combining these we get $n + (k + 1) = 0$	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0,$ $m = k + 1,$ $n + (k + 1) = 0$	\perp
By the associativity of addition, this is the same as $(n + k) + 1 = 0$,	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0,$ $m = k + 1,$ $n + (k + 1) = 0,$ $(n + k) + 1 = 0$	\perp
Since $n + k$ is a natural number, we have $(n + k) + 1 \neq 0$.	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0,$ $m = k + 1,$ $n + (k + 1) = 0,$ $(n + k) + 1 = 0,$ $(n + k) + 1 \neq 0$	\perp
but this is a contradiction.	$A,$ $B,$ $C,$ $n, m, k \in \mathbb{N},$ $n + m = 0,$ $m = k + 1,$ $n + (k + 1) = 0,$ $(n + k) + 1 = 0,$ $(n + k) + 1 \neq 0,$ \perp	\perp

Table 1.1: Trace of the proof of $\forall nmk \in \mathbb{N}. n + m = 0 \Rightarrow m \neq k + 1$

Suppose someone gives you a piece of straightline C code (no loops) for computing some complicated function $f(x, y)$ of a pair of integers variables x and y , and asks you to check that it works. To understand what the code is doing, most likely you will start at the beginning of the code, read each statement in turn and keep track of the values of x and y in your head¹ as they are manipulated by the code. By the end, you should be able to tell whether or not the final value returned is $f(x, y)$ (for the original values of x and y) or not. Of course, to do this, you need to understand the “rules” of the C language, i.e. how the syntax fits together and how each statement affects the current state, i.e. the current values of x and y .

The same process is required to understand a proof. You will go through, statement by statement, keeping track of the proof state until, by the end, it should be clear whether or not the proof works, i.e. in the final proof state, the goal is something that is already known. Of course, to do this you will need to know the “rules” of mathematical proof. Giving you those rules is the purpose of the following section, but I think it is useful to go through this example first anyway so you can see what I mean.

One thing that we need to make clear before that, however, is what the starting point for the proof is.

In “real life” mathematicians are trying to prove statements that nobody has ever proven before (and may turn out to be false). To give themselves the best chance to prove their statement, they will employ all the mathematical knowledge that they have at their disposal: countless theorems already proven by the mathematicians that came before them, or that they proved during their career. In other words, they start the proof with a set of assumptions that is enormous (but they know each is justified by a proof, somewhere).

When learning mathematics, however, we are usually proving quite simple statements that someone else has already proven a long time ago. Hence, it can be a bit confusing to the student to know which already proven statements they are allowed to assume and which they are required to prove first. This will not be a problem during *Types and Lambda Calculus* because (I guess) you don’t know any basic theorems of this theory, so we will just agree that you can assume only those statements that we have proven so far during the course of the unit (and anything you know about sets and strings).

The example above, however, is in number theory, and you probably already have a good idea of what simple statements are true about natural numbers and which are not. So, for this example, I will be very clear about what we are allowing ourselves to assume at the start of the proof:

(A) for all $p \in \mathbb{N}$: $p + 1 \neq 0$

(B) for all p, q and r in \mathbb{N} : $p + (q + r) = (p + q) + r$

(C) for all $p, q \in \mathbb{N}$: $p + q \in \mathbb{N}$

Statement B is what is referred to in the proof as “associativity of addition”. Statements A and C are used in the last step, but are not referred to by any name. This probably means the author of the proof believes them to be so simple that they are not worth mentioning. Now let’s analyse the proof step by step starting in the proof state containing only these assumptions A , B and C . A trace of the

¹If you are not so experienced with C, or the code is very complicated, then it may be useful to write down the evolving values of x and y on a scrap of paper.

proof is given in Table 1.1, whose final two columns shows the proof state after executing each of the statements.

A few things about the general shape of this proof are worth observing.

- First, the proof appears to prove the statement, because in the final proof state the goal is already contained in the assumptions. However, we can't know that it is really a valid proof unless we check that each of the proof steps really transforms the proof state before-it into the proof state after-it, and this requires a knowledge of the rules.
- The set of assumptions grows larger over the course of the proof.
- The goal starts off quite large and complicated, but is progressively simplified over the first few proof steps and, once it is as simple as possible, it doesn't change from that point on.

Backward Reasoning Steps 1–3 in which the goal gets simpler are called *backwards reasoning*. It is quite typical that a proof starts by performing backward reasoning. Backwards reasoning steps are of the form “To prove a goal of shape A it is enough to prove the simpler goal(s) B_1, \dots, B_n ”. In the case of the example, the particular backwards steps that are employed only generate a single simple goal. For example, step 2 uses the backwards reasoning rule for implication, which says “to prove A implies B , it is enough to assume A is true and then try to prove B ”. Backwards reasoning is very easy because the backwards reasoning steps you are allowed to make mostly just depend upon the syntactic shape of the goal: is it a conjunction, an implication, a for all, etc. Often, but not always, backwards reasoning not only simplifies the goal, but also adds new assumptions (which is helpful for forwards reasoning).

Forwards Reasoning Steps 4–7 are *forwards reasoning* steps. Forwards reasoning consists of those proof steps that do not change the goal, but simply deduce more and more assumptions (facts) by combining existing ones. The majority of a non-trivial proof is spent on forwards reasoning. Forwards reasoning is a bit more difficult only because you can build up a large set of assumptions and spotting which ones can be fruitfully combined to deduce new ones *that are actually relevant to proving your goal* requires some insight.

1.2 Rules of the game

In Appendix A are listed all the rules that suffice for any mathematical proof in this unit, organised by logical connective. For each logical connective there is a list of the forwards and backwards rules for that connective: the backwards rules tell you how you can simplify a goal which is built from that connective at the outer level, the forwards rules tell you how you can deduce new assumptions from an existing one built from that connective at the outer level.

Each rule is given in form of a table entry of shape:

Assms	Goal

“blah blah”

Assms	Goal

This tells you the syntax of the proof rule, e.g. the words “blah blah”, and the semantics, i.e. how it changes the proof state. For a proof to be valid, the commands must hang together as a correct program, so you must ensure that you only use a rule when the proof state before has the correct shape. For example, the backwards reasoning rule for implication requires that the goal is literally an implication, i.e. has the shape $A \Rightarrow B$. It does not apply when the goal has shape, e.g. $\forall x \in X. A \Rightarrow B$, which is a for all (whose body is an implication), nor when the goal has shape $(A \Rightarrow B) \wedge C$, which is a conjunction (whose right conjunct is an implication). The exact form of words is not important, as long as it is clear from the context (i.e. the proof state) which rule are you trying to invoke.

Let’s take a look at some of the more interesting ones.

Implication

$A \Rightarrow B$	A implies B	If A then B .	Suppose A , then B .
-------------------	-----------------	-------------------	--------------------------

Backward Rule

To prove $A \Rightarrow B$ starting from some assumptions, it suffices to instead prove B with A added to the assumptions.

Assms	Goal
...	$A \Rightarrow B$

“Assume A . We prove B as follows...”

“To show $A \Rightarrow B$, we assume A”

“Assume A”

“Suppose A is true. ...”

Assms	Goal
..., A	B

Forwards Rule

If you already know $A \Rightarrow B$ and you already know A , then you know B .

Assms	Goal
$\dots, A \Rightarrow B, A$	C

"From $A \Rightarrow B$ and A we conclude B ."

"Applying $A \Rightarrow B$ to A we obtain B ."

"We have B by *modus ponens*."

Assms	Goal
$\dots, A \Rightarrow B, A, B$	C

Negation

$\neg A$ not A A is false A is absurd

Backwards Rule

To prove $\neg A$ starting from some assumptions, it suffices to derive a contradiction (give a proof of false) starting from the same assumptions with A added.

Assms	Goal
\dots	$\neg A$

"We assume A and try to obtain a contradiction."

"To show $\neg A$, we assume A ..."

"Assume A ."

Assms	Goal
\dots, A	\perp

Forwards Rule

If you know $\neg A$ and you also know A , then you know absurdity.

Assms	Goal
$\dots, \neg A, A$	C

“From A and $\neg A$ we obtain a contradiction.”

Assms	Goal
$\dots, \neg A, A, \perp$	C

1.3 Examples

The following is an example of pure logic. Let’s suppose we do not know any facts about logic already (i.e. our starting assumptions are empty).

Lemma 1.2. *A implies $\neg\neg A$*

Proof. Assume A . We wish to show $\neg\neg A$, so we assume $\neg A$ and try to obtain a contradiction. From A and $\neg A$ we obtain the desired contradiction. \square

Make sure you can follow which rule is being used in each step and how the proof state evolves. The first step in becoming proficient at proving is to at least be able to check that someone else’s proof is valid. The following table gives a listing of the proof state at the end of each step in the proof.

Proof step	Assumptions	Goal
		$A \Rightarrow \neg\neg A$
“Assume A ”	A	$\neg\neg A$
“... so we assume $\neg A$ and ... contradiction.”	$A, \neg A$	\perp
“From A and $\neg A$ we obtain...”	$A, \neg A, \perp$	\perp

In the following example, we collect quite a few assumptions from our initial backward reasoning. In such cases it can sometimes be helpful to label them as we go for when we want to combine them in new ways during forward reasoning.

Lemma 1.3. *if A implies B then $\neg B$ implies $\neg A$*

Proof. Assume $A \Rightarrow B$ (A1) and assume $\neg B$ (A2). We aim to prove $\neg A$, so we will assume A (A3) and try to obtain a contradiction. From (A1) and (A3), we obtain B (A4). From (A4) and (A2) we obtain the desired contradiction. \square

This is what is going on in my head when I read this proof back to myself. Let's again suppose that my starting knowledge is zero. Make sure that you can identify which backwards or forwards rule was used in each case.

Proof step	Assumptions	Goal
		$(A \Rightarrow B) \Rightarrow \neg B \Rightarrow A$
"Assume $A \Rightarrow B$ "	$A1:A \Rightarrow B$	$\neg B \Rightarrow \neg A$
"assume $\neg B$ "	$A1:A \Rightarrow B, A2:\neg B$	$\neg A$
"assume A ... try ... contradiction."	$A1:A \Rightarrow B, A2:\neg B, A3:A$	\perp
"From (A1) and (A3), we obtain B ."	$A1:A \Rightarrow B, A2:\neg B, A3:A, A4:B$	\perp
"From ... obtain ... contradiction."	$A1:A \Rightarrow B, A2:\neg B, A3:A, A4:B, \perp$	\perp

One important reason that it is helpful to have the proof state in your head is that, because the proof steps are just certain phrases in English, sometimes the proof state is needed for disambiguation. For example, it may only be possible to understand whether the word "assume" heralds a use of implication introduction or a use of negation introduction, based on goal at the point at which the word was used.

The next example also uses conjunction.

Conjunction

$A \wedge B$ A and B

Backwards Rule

To prove $A \wedge B$ from some assumptions, it suffices to give two separate proofs: one proof of A from the assumptions and one proof of B from the assumptions.

Assms	Goal
...	$A \wedge B$

"We prove A and B separately. To prove A ..."

"For A , we argue as follows ... For B , ..."

Assms	Goal	Assms	Goal
...	A	...	B

The proof splits into two at this point, so you have to keep in mind two states, although you can just work with one at a time.

Forwards Rule

If you know $A \wedge B$, then you can know A and you know B .

Assms	Goal
$\dots, A \wedge B$	C

"" (too trivial to mention)

Assms	Goal
$\dots, A \wedge B, A, B$	C

Lemma 1.4. $A \Rightarrow (B \Rightarrow C)$ implies $A \wedge B \Rightarrow C$

Proof. Assume $A \Rightarrow (B \Rightarrow C)$ (*). Assume A and B . From (*) and $A, B \Rightarrow C$ follows. From this and B, C follows. \square

Proof step	Assumptions	Goal
		$(A \Rightarrow (B \Rightarrow C)) \Rightarrow (A \wedge B) \Rightarrow C$
"Assume $A \Rightarrow (B \Rightarrow C)$ (*)"	*. $A \Rightarrow (B \Rightarrow C)$	$A \wedge B \Rightarrow C$
"Assume A and B "	*. $A \Rightarrow (B \Rightarrow C), A, B$	C
"From (*) and $A, B \Rightarrow C$ follows."	*. $A \Rightarrow (B \Rightarrow C), A, B, B \Rightarrow C$	C
"From this and B, C follows."	*. $A \Rightarrow (B \Rightarrow C), A, B, B \Rightarrow C, C$	C

In this example I omitted the step between assuming $A \wedge B$ and deducing A and B separately because there is never any reason not to immediately deduce the two conjuncts separately.

2. Proof: Quantifiers, Equality

In this second part of our crash course on proof, we introduce first-order reasoning with the rules for the quantifiers, equality and how to work with theories.

2.1 Quantifier rules

Universal quantification

$\forall x \in X. A$	$\forall x : X. A$	for all x in X, A	A holds of all x in X
----------------------	--------------------	-----------------------	-----------------------------

Backwards Rule

To give a proof of $\forall x \in X. A$ from some assumptions, it suffices to give a proof of A from the same assumptions with $x \in X$ added, as long as you have not made any prior assumptions about the name x (x doesn't appear in any of your existing assumptions).

Assms	Goal	with x not occurring free in the assms
...	$\forall x \in X. A$	

“Let $x \in X$. We show A .”

“Let x be an arbitrary member of X . We show A .”

“Suppose x is in X therefore A .”

Assms	Goal
..., $x \in X$	A

Forwards Rule

If you have know $\forall x \in X. A$ and you know $t \in X$, then you know A with all occurrences of x replaced by t .

Assms	Goal
$\dots \forall x \in X. A$	$t \in X$

"It follows that A holds of t "

Assms	Goal
$\dots, \forall x \in X. A, t \in X$	$A[t/x]$

Existential Quantification

$\exists x \in X. A$	$\exists x : X. A$	there exists x in X such that A	A holds of some x in X
----------------------	--------------------	---------------------------------------	------------------------------

Backwards Rule

To prove $\exists x : X. A$ from some assumptions, it suffices to find a witness t and give proofs that $t \in X$ and of A with every occurrence of x replaced by t , starting from the same assumptions.

Assms	Goal
\dots	$\exists x \in X. A$

"We show that A holds of t ."

"We take t as witness. To see $A[t/x]$..."

"We show that t is such an x ."

Assms	Goal	Assms	Goal
\dots	$t \in X$	\dots	$A[t/x]$

Here, $A[t/x]$ means A with all free occurrences of x replaced by t .

Forwards Rule

Assms	Goal
$\dots \exists x \in X. A$	C

with x not occurring free in the assms

“Let x be the witness...”

Assms	Goal
$\dots \exists x \in X. A \quad x \in X \quad A$	C

2.2 Playing the game

There are generally many ways that one can complete a proof because, in any given proof state, several forwards and backwards rules may be applicable. In many cases, it requires trial and error until you find the right strategy for a proof of a given statement. However:

An excellent rule of thumb is to:

1. Apply backwards reasoning until the goal is as simple as possible (doesn't involve any more logical connectives) and all the delicious extra assumptions have been extracted.
2. Rewrite the goal using any non-inductive/recursive definitions that you have assumed (see Theories and Equality in Section 2.3). This may give you even more opportunities to apply backwards reasoning steps.
3. Step back and consider what you now know to be true (your assumptions). Then apply forwards reasoning to deduce more and more of them until you have finally deduced the goal you are looking to prove.

This strategy will not work in all cases, but I would always use it as my first attempt at a proof. In this unit, it will serve you very well for almost every proof.

Here's a slightly fiddly example that demonstrates some of the quantifier rules. In it we use the notation $A[y/x]$ to mean the formula A but with any (free) occurrence of x replaced by y . For example, if A is $x = 0 \vee x = 1$, then $A[y/x]$ is $y = 0 \vee y = 1$. In the statement that follows, the point is that $\forall x \in X. A$ is really the same as $\forall y \in X. A[y/x]$, we have just renamed a bound variable to give it a more convenient name, y .

Lemma 2.1. *Suppose x does not occur free in B . Then $(\exists x \in X. A) \Rightarrow B$ iff $\forall x \in X. A \Rightarrow B$.*

Proof. Assume x does not occur free in B . We proceed to prove both directions separately:

- In the forward direction, suppose $(\exists x \in X. A) \Rightarrow B$ (*). Our goal is to show $\forall x \in X. A \Rightarrow B$, but let us rename the bound variable x to avoid confusion with the x in our assumptions. So we prove $\forall y \in X. A[y/x] \Rightarrow B[y/x]$. Then let $y \in X$ and assume $A[y/x]$. I claim that $\exists x \in X. A$. To see why, take the witness to be y . We already know that $y \in X$ and $A[y/x]$ so that proves the claim. It follows from this and (*) that, therefore, B . Our goal is to show $B[y/x]$, but we assumed that B does not contain any occurrence of x , so replacing all occurrences of x in B by y is just B , i.e. $B = B[y/x]$. Hence, we have proven the goal.
- In the backward direction, suppose $\forall x \in X. A \Rightarrow B$ (*). Then suppose $\exists x \in X. A$. Let the witness be some $y \in X$ (so we know $A[y/x]$). Then it follows from (*) that $A[y/x] \Rightarrow B[y/x]$. Since we have the former already, we obtain the latter $B[y/x]$. But, we started by assuming that x does not occur in B , so $B[y/x]$ is exactly B .

□

Here is the trace of the backward direction part of the proof, the initial goal is a bit long so let us shorten it to $G := (\forall x \in X. A \Rightarrow B) \Rightarrow ((\exists x \in X. A) \Rightarrow B)$.

Proof step	Assumptions	Goal
	$x \notin \text{FV}(B)$	G
"suppose ... (*)"	$x \notin \text{FV}(B), \forall x \in X. A \Rightarrow B$	$(\exists x \in X. A) \Rightarrow B$
"Then suppose ..."	$x \notin \text{FV}(B), \forall x \in X. A \Rightarrow B, \exists x \in X. A$	B
"Let the witness be..."	$x \notin \text{FV}(B), \forall x \in X. A \Rightarrow B, \exists x \in X. A, y \in X, A[y/x]$	B
"The it follows from (*)"	$x \notin \text{FV}(B), \dots, A[y/x], A[y/x] \Rightarrow B[y/x]$	B
"Since we have the former..."	$x \notin \text{FV}(B), \dots, A[y/x], A[y/x] \Rightarrow B[y/x], B[y/x]$	B

Bending the rules

There are many, many ways that people take short-cuts when writing proofs. Sometimes it is just laziness, but other times it makes the difference between a proof that is a reasonable size and a proof that is large beyond comprehension. Generally, it's ok to leave steps implicit as long as you believe that your reader can fill in the gaps: it should be possible, in principal, to derive all the explicit steps that you missed out without having to do too much search.

Sometimes backwards reasoning steps are left implicit. For example, when proving a goal of the form $\forall x \in X. A \Rightarrow B$, often the proof will launch straight into proving B and, as part of that proof, will use A and $x \in X$ as if they are already assumed.

Conversely, when forwards reasoning from an assumptions of shape $\forall x \in X. A \Rightarrow B$ and $A[t/x]$ and $t \in X$, often the proof will immediately deduce $B[t/x]$, performing the forwards rule for forall and the forwards rule for implication in one combined step.

I will try not to leave backwards steps implicit, because they are crucial to understanding where the proof is headed, which is important if you are just starting out. I will, however, often combine forwards steps together if I think it remains clear what is happening.

2.3 Equality

Most theories make use of some notion of equality over the inhabitants of the domain. The rules we can use for reasoning about equality are very straightforward:

At any time, you can assume $s = s$, for any term s , which is obviously true.

Assms	Goal
...	C

"" (too trivial to mention)

Assms	Goal
..., $s = s$	C

If you know an equation, you can use it to rewrite the goal.

Assms	Goal
..., $s = t$	$C[s/x]$

" $C[s/x]$ is just $C[t/x]$ "

" $C[s/x]$ is therefore the same as $C[t/x]$ "

Assms	Goal
..., $s = t$	$C[t/x]$

The use of $C[s/x]$ here is just a way of saying that you have some proposition which contains an occurrence of s and the corresponding $C[t/x]$ indicates replacing that occurrence of s by t .

If you know an equation, you can use it to rewrite an assumption.

Assms	Goal
..., $s = t, A[s/x]$	C

" $A[s/x]$ is just $A[t/x]$ "

" $A[s/x]$ is therefore the same as $A[t/x]$ "

Assms	Goal
..., $s = t, A[s/x], A[t/x]$	C

--

3. Proof: Theories and Induction

We have talked about the generic, purely logical aspects of proof but, in reality, most proofs argue something in a specific domain (or, in the logical jargon, a specific theory). When you are working with a particular theory, e.g. the theory of arithmetic over the natural numbers, the theory gives you more to work with: for example, the language of formulas is enlarged to include new constants, function symbols and relations like 0, + and < respectively. Along with those symbols come axioms or definitions, and you can make use of these as extra assumptions in your proofs.

A theory is just a collection of true sentences, typically organised around some theme, like a particular datatype and its operations “the theory of arithmetic” or a class of datatypes (an abstract datatype) “the theory of groups”. Theories bring with them definitions and existing theorems, which both contribute to your starting assumptions when you set out to prove something.

The following example uses the following recursive definition of multiplication of natural numbers:

$$\begin{aligned} 0 * p &= 0 \\ (q + 1) * p &= q * p + p \end{aligned}$$

If an equational definition of some function f is valid in a given theory, then we can just take the equations as extra assumptions when reasoning about f .

As usual, when we see a statement, such as the equations above, that contains free variables (here p and q), then the convention is to regard them as universally quantified, i.e:

$$\forall p \in \mathbb{N}. 0 * p = 0 \quad \text{and} \quad \forall p, q \in \mathbb{N}. (q + 1) * p = q * p + p$$

For this reason, it is extremely important to know your definitions well, because you will always be using them when you come to writing proofs. For the purpose of this example, we also allow ourselves to assume the commutativity of multiplication:

$$\forall p, q \in \mathbb{N}. p * q = q * p$$

Lemma 3.1. *for all $n, m \in \mathbb{N}$: if $n = 0$ or $m = 0$ it follows that $n * m = 0$*

Proof. Let $n, m \in \mathbb{N}$ and assume $n = 0 \vee m = 0$. We show that $n * m = 0$ by case analysis.

- When $n = 0$ then the goal $n * m = 0$ is $0 * m = 0$ which is true by the definition of multiplication.
- When $m = 0$ then the goal $n * m = 0$ is $n * 0 = 0$ which, by the commutativity of multiplication, is $0 * n = 0$ and this also true by definition.

□

In this proof we also used a rule to do with disjunction – we worked forwards from an assumption that $n = 0 \vee m = 0$. The formulation of the rule is interesting:

Forwards Rule

If you have already know $A \vee B$ and, additionally, you can give the following two proofs:

- a proof of C starting from your assumptions and A
- a proof of C starting from your assumptions and B

Then you will also know C .

Assms	Goal
$\dots, A \vee B$	C

“We proceed by cases on $A \vee B$.

– Assume A ... Hence C .

– Assume B ... Hence C .”

“We analyse the two cases.”

“We proceed by case analysis on $A \vee B$.”

Assms	Goal	Assms	Goal
$\dots, A \vee B, A$	C	$\dots, A \vee B, B$	C

This forces the proof to split into two, and you’re not done until you have completed both of them.

So, this is how the proof state evolves, with G abbreviating the original statement to be proven and the three assumptions that we know about from the theory abbreviated by their numbers:

(i) $\forall p. 0 * p = 0$

(ii) $\forall pq. (q + 1) * p = q * p + p$

(iii) $\forall pq. p * q = q * p$

Proof step	Assumptions	Goal
	(i), (ii), (iii)	G
“Let $n, m \in \mathbb{N}$...”	(i), (ii), (iii), $n, m \in \mathbb{N}$	if $n = 0$ or $m = 0$, then $n * m = 0$
“and assume $n = 0 \vee m = 0$ ”	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0$	$n * m = 0$

“We show ... by case analysis”

"When $n = 0$ "	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0, n = 0$	$n * m = 0$
" $n * m = 0$ is $0 * m = 0$... true by defn"	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0, n = 0$	$0 * m = 0$

"When $m = 0$ "	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0, m = 0$	$n * m = 0$
" $n * m = 0$ is $n * 0 = 0$ "	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0, m = 0$	$n * 0 = 0$
"by commutativity is $0 * n = 0$, true by defn"	(i), (ii), (iii), $n, m \in \mathbb{N}, n = 0 \vee m = 0, m = 0$	$0 * n = 0$

3.1 Induction

Some theories even come with their own proof rules for reasoning about the datatypes involved. Induction over the natural numbers is an excellent example:

To prove $\forall n \in \mathbb{N}. A$ it suffices to prove $A[0/n]$ and, assuming $A[k/n]$, to prove $A[k + 1/n]$ for all k .

Assms	Goal
...	$\forall n \in \mathbb{N}. A$

"We prove the result by induction on $x \in \mathbb{N}$.

- When $x = 0$... hence A is true of 0 .
- When $x = k + 1$, assume A holds of k
... A is true of $k + 1$."

"By induction on $n \in \mathbb{N}$. We analyse the cases..."

Assms	Goal	Assms	Goal
...	$A[0/n]$..., $A[k/n]$	$A[k + 1/n]$

Recall that $A[0/n]$ means the formula A but where n is replaced by 0 . So, when we invoke induction on the natural numbers in order to prove a formula of the form $\forall n \in \mathbb{N}. A$, we are required to prove two new goals: A but with n being specifically 0 , and A with n being $k + 1$ but under the assumption, usually given a special name for some reason: *the induction hypothesis*, of A with k replacing n . So, in the sense above, natural number induction is an alternative introduction rule (backward reasoning) for those formulas of shape $\forall x : X. A$ in the specific case that X is \mathbb{N} .

We get an induction rule every time we introduce a datatype defined by an "inductive" definition.

Grammars

Every context-free grammar defines an inductive datatype. Here is a grammar for a language that some of you may recognise, the language of regular expressions:

$$R, S ::= \emptyset \mid \epsilon \mid a \mid R \cdot S \mid R + S \mid R^*$$

Regular expressions are a standard way of expressing pattern matching on strings when programming, but formally, we give a regular expression R semantics by describing the set of all strings $L(R)$ that R can match:

$$\begin{array}{ll} L(\emptyset) &= \emptyset \\ L(\epsilon) &= \{\epsilon\} \\ L(a) &= \{a\} \end{array} \qquad \begin{array}{ll} L(R + S) &= L(R) \cup L(S) \\ L(R \cdot S) &= \{uv \mid u \in L(R_1) \wedge v \in L(R_2)\} \\ L(R^*) &= \bigcup_{n \in \mathbb{N}} R^n \end{array}$$

In particular, the $*$ operation, R^* , matches any finite sequence of (strings that themselves each match) R , just as the similar notation we used for CFGs in PLC. For example, the regular expression $0^* \cdot 1 \cdot 0^* \cdot 1 \cdot 0^*$ matches those strings over $\{0, 1\}$ that contain exactly two 1s. The regular expression $(0 + 1)^* \cdot 1 \cdot (0 + 1)^* \cdot 1 \cdot (0 + 1)^*$ matches those strings that contain at least two 1, and so on.

We don't really care about regular expressions in this unit, but it is a good example. From the structure of the grammar, we get an induction proof rule. Since it is a form of induction, it can be used to prove formulas of the form "for all regular expressions R , A ", where A is some statement that mentions R (i.e. a property of R). A complete use of the rule consists of proving a number of subgoals, one for each of the different productions/alternatives in the grammar. In this case we have 6 alternatives and so there are 6 goals to prove. Each goal consists of proving the body of the original goal, the A , for each of the 6 different shapes of R . Where a production in the grammar is recursive, i.e. that alternative contains occurrences of the metavariable R or S , then we can assume that ϕ holds of those "smaller" occurrences when we proving the goal corresponding to that case. These are the induction hypotheses. Note we abbreviate "for all regular expressions R " as " $\forall R \in \text{REGEXP}$ ".

Assms	Goal
...	$\forall R \in \text{REGEXP}. A$

"We prove the result by induction on $R \in \text{REGEXP}$.

- When $R = \epsilon$... hence A is true of ϵ .
- When $R = \emptyset$... hence A is true of \emptyset .
- When $R = a$... hence A is true of a .
- When R is of shape $(R_1 + R_2)$, assume A holds of R_1 and R_2 separately ... A is true of $(R_1 + R_2)$.
- When R is of shape $(R_1 \cdot R_2)$, assume A holds of R_1 and R_2 separately ... A is true of $(R_1 + R_2)$.
- When R is of shape R_1^* , assume A holds of R_1 ... A is true of R_1^* ."

Assms	Goal	Assms	Goal
...	$A[\epsilon/R]$...	$A[\emptyset/R]$

Assms	Goal	Assms	Goal
...	$A[a/R]$..., $A[R_1/R], A[R_2/R]$	$A[(R_1 + R_2)/R]$

Assms	Goal	Assms	Goal
..., $A[R_1/R], A[R_2/R]$	$A[(R_1 \cdot R_2)/R]$..., $A[R_1/R]$,	$A[R_1^*/R]$

Then we can prove the following fact about star-free languages:

Lemma 3.2. *For each regular expression R , if R does not contain $*$, then $L(R)$ is finite.*

Proof. We proceed by induction on R .

- When $R = \epsilon$, (suppose R does not contain $*$) then, by definition, $L(R) = \{\epsilon\}$, which is finite.
- When $R = \emptyset$, (suppose R does not contain $*$) then, by definition, $L(R) = \emptyset$, which is finite.
- When $R = a$, (suppose R does not contain $*$) then, by definition, $L(R) = \{a\}$ which is finite.
- When R is of shape $(R_1 + R_2)$, proceed as follows. Assume the induction hypotheses:
 - (i) If R_1 does not contain $*$, then $L(R_1)$ is finite.
 - (ii) If R_2 does not contain $*$, then $L(R_2)$ is finite.

Now, suppose $(R_1 + R_2)$ does not contain $*$. It follows that neither R_1 nor R_2 can contain $*$. Hence, from induction hypothesis (i), we obtain that $L(R_1)$ is finite and, from induction hypothesis (ii), we obtain that $L(R_2)$ is finite. By definition, $L(R_1 + R_2)$ is the union of $L(R_1)$ and $L(R_2)$, and the union of two finite sets is finite.

- When R is of shape $(R_1 \cdot R_2)$, proceed as follows. Assume the induction hypotheses:
 - (i) If R_1 does not contain $*$, then $L(R_1)$ is finite.
 - (ii) If R_2 does not contain $*$, then $L(R_2)$ is finite.

Now, suppose $(R_1 R_2)$ does not contain $*$. It follows that neither R_1 nor R_2 can contain $*$. Hence, from induction hypothesis (i), we obtain that $L(R_1)$ is finite and, from induction hypothesis (ii), we obtain that $L(R_2)$ is finite. By definition, $L(R_1 R_2)$ is the set of all concatenated pairs of words uv with u from $L(R_1)$ and $v \in L(R_2)$. Since we now know $L(R_1)$ and $L(R_2)$ are finite, it follows that there are only finitely many such pairings.

- When R is of shape (R_1^*) , proceed as follows. Assume the induction hypothesis:
 - (i) If R_1 does not contain $*$ then $L(R_1)$ is finite.

Suppose R_1^* does not contain $*$. Clearly (i.e. by the basic properties of strings) R_1^* contains $*$. Hence, we obtained a contradiction and the subgoal follows trivially.

□

4. Syntax: Terms and Binding

In this unit we are going to study a kind of prototypical functional programming language, called PCF, or more fully *Programming Computable Functions*. It has a bit of a strange heritage, because it was first introduced by Gordon Plotkin in 1977 as a programming language alternative to Dana Scott's *Logic of Computable Functions*, but it has become one of the standard programming languages that theorists use to investigate the central questions of programming languages and type systems.

In this lecture, we are going to get some intuitions about PCF, then we will see the formal definition of the syntax.

4.1 Intuitions

When we talk about programming languages, it usually suffices to talk about them abstractly as a group of features – the kinds of expression that they support. In these terms we would characterise PCF as a programming language characterised by supporting: *arithmetic*, *conditionals*, *higher-order functions*, and *general recursion*. Now let's explore the language informally by taking a look at each of these features.

Arithmetic

Arithmetic is provided by built-in natural numbers. The numbers are given in unary notation: there is a constant for zero, Z , and a constant for the successor of any given natural number, S . These behave like constructors in Haskell, so we could imagine that they have been introduced via a datatype declaration like:

```
data Nat = Z | S Nat
```

But our language is so simple that we will not have the ability to program new datatypes like this, S and Z are baked-in. With them we can construct the natural number k by applying the successor to the zero k -times, e.g. the first three naturals are:

$$Z, \quad S\ Z, \quad S\ (S\ Z), \quad S\ (S\ (S\ Z))$$

There is also a built-in predecessor function pred which performs subtraction-by-1. When it's argument is zero it is defined so as to just return zero again:

$$\begin{aligned} \text{pred } Z &\triangleright^* Z \\ \text{pred } (S\ N) &\triangleright^* N \end{aligned}$$

Here you can see that applying a function f to an input x is written using juxtaposition $f\ x$, as in typical functional programming languages. You can also see the \triangleright^* notation we will use for making computational steps.

Addition, subtraction, multiplication and so on are not built-in to the language, but we can define addition by iterating the successor, subtraction by iterating the predecessor, multiplication by iterating addition and so on.

Conditionals

To keep our language as simple as possible, there are no other base types other than natural numbers. In particular, there are no built-in Booleans. Hence, the built-in conditional, `ifz`, is perhaps not what you expect:

$$\begin{aligned} \text{ifz } Z\ N\ P &\triangleright^* N \\ \text{ifz } (S\ M)\ N\ P &\triangleright^* P \end{aligned}$$

It takes three arguments as input, and returns the second N or third P depending on whether the first was zero or not. So we can think of the arguments as the guard, the then-branch and the else-branch respectively. When the guard is Z , it returns the then-branch N and when the guard is the successor of something, it returns the else-branch P .

Higher-order functions

Whilst most general purpose programming languages we meet will have the ability to perform arithmetic and do conditional branching, it is the support of higher-order functions that marks PCF as a *functional programming language*. Higher-order functions are written explicitly using λ -abstraction. A function that takes x as input and returns expression M (which will typically contain x) is written:

$$\lambda x. M$$

As remarked above, when we want to supply such a function with an input, say N , we write the two juxtaposed and say that the former is *applied* to the latter, and the result is M but with x everywhere replaced by N , which we write as $M[N/x]$.

$$(\lambda x. M)\ N \triangleright^* M[N/x]$$

In fact, as we shall see when we define the language precisely, abstraction and application are the *only* mechanism for building programs! This is what makes our language a λ -calculus.

For example, the function that takes a number x as input and returns $x + 2$ is written: $\lambda x. S\ (S\ x)$. When we supply it with one as an argument, i.e. $S\ Z$, then the result is three, i.e:

$$(\lambda x. S\ (S\ x))\ (S\ Z) \triangleright^* S\ (S\ (S\ Z))$$

where we have taken the expression that is returned $S\ (S\ x)$ and replaced every occurrence of x by the argument, which was $S\ Z$.

The function that takes a number x and always returns 0 is written: $\lambda x. Z$. When we apply it to the input 2, written $S (S Z)$, then the output will be Z , because we compute by replacing all occurrences of the formal parameter x in the return expression Z by the actual parameter 2:

$$(\lambda x. Z) (S (S Z)) \triangleright^* Z$$

In Haskell, λ is written in the best ascii approximation, which was deemed to be the backslash `\` and an arrow `→` is used to separate the input from the output instead of the full-stop, so the above application would be written:

$$(\backslash x \rightarrow S (S x)) (S Z)$$

We use iterated λ -abstractions to define functions of more than one argument. For example the function that takes a function f and a number x as input and returns the result of applying f to x twice is written:

$$\lambda f. \lambda x. f (f x)$$

We sometimes call such functions *curried*. Strictly speaking, this is the function that takes a function f as input and returns a function $\lambda x. f (f x)$, which itself takes a number x and returns $f (f x)$, but in practice we tend to think of it as a single function of two arguments f and x . When we apply this function to two arguments, say pred and the numeral for 3, $S (S (S Z))$, then we can compute as follows:

$$\begin{aligned} & (\lambda f. \lambda x. f (f x)) \text{pred} (S (S (S Z))) \\ & \triangleright^* (\lambda x. \text{pred} (\text{pred } x)) (S (S (S Z))) \\ & \triangleright^* \text{pred} (\text{pred} (S (S (S Z)))) \\ & \triangleright^* S Z \end{aligned}$$

General recursion

Recursion (i.e. the ability to compute something by potentially unbounded repetition) is given by a built-in *fixed-point combinator* called `fix`. It's behaviour on any input M is as follows:

$$\text{fix } M \triangleright^* M (\text{fix } M)$$

It seems a bit mysterious and, indeed, we will need to devote a lecture to a proper discussion of `fix`, but for now let us just remark that we can use it to define e.g. addition:

$$\underline{\text{add}} = \text{fix } (\lambda f x y. \text{ifz } x y (f (\text{pred } x) (S y)))$$

Let make a few points about this notation.

First, notice we have abbreviated the iterated abstraction $\lambda f. \lambda x. \lambda y. \dots$ as $\lambda f x y. \dots$, this is something we will commonly do. Second, notice that we have introduced a new name `add` and stated that it is identical to `fix` $(\lambda f x y. \text{ifz } x y (f (\text{pred } x) (S y)))$.

PCF doesn't have a built-in mechanism for making definitions, for defining new functions. However, although this would make writing programs extremely painful in practice, it is not really a hindrance to doing theory because we can just agree among ourselves that `add` is our abbreviation

for the term on the right-hand side. This means that whenever we use add, e.g. when we state an equivalence like:

$$\text{add } (S\ Z)\ Z \triangleright^* S\ Z$$

it is just an abbreviation for the expanded statement:

$$(\text{fix } (\lambda f\ x\ y. \text{ifz } x\ y\ (f\ (\text{pred } x)\ (S\ y))))\ (S\ Z)\ Z \triangleright^* S\ Z$$

This actually gives us the best of all worlds, because we can define abbreviations and build compact programs out of them, but when we analyse an arbitrary program we only have to consider a handful of primitive notions. Although our language is so small and simple, we will be adding everything the modern programmer could possibly need as abbreviations.

4.2 Syntax

We define the possible expressions of the programming language precisely, using a grammar. First, let us assume we have a countably infinite supply of program variable names. It actually won't matter what the variable names concretely look like, you could think of them as Haskell identifiers but we will never actually need to write them down. What we will do is use the letters towards the end of the alphabet (and variations): x, y, y', x_1, x_2 and so on to stand for arbitrary ones. Then program expressions, which we call *terms*, are built from variables and constants (the built-ins of PCF) using abstraction and application.

Definition 4.1. *The **terms** of PCF are those strings that can be generated by the following grammar:*

$$\begin{array}{ll} \text{(Terms)} & M, N ::= x \mid c \mid (\lambda x. M) \mid (MN) \\ \text{(Constants)} & c ::= \text{fix} \mid Z \mid S \mid \text{pred} \mid \text{ifz} \end{array}$$

We use M, N, P, Q and other capital letters from the end of the alphabet to stand for arbitrary terms. In this grammar, x is any *variable*, c is any *constant*, the form $\lambda x. M$ is called a *λ -abstraction* and the form $M\ N$ is called *application*. We write Λ for the set of all terms of PCF.

You have seen examples of all the built-in constants already. Some examples of terms are:

$$x \quad \text{ifz} \quad Z \quad (\lambda x. (S\ (S\ x))) \quad (\text{pred } (S\ Z)) \quad (\lambda f. (\lambda x. (f\ (f\ x))))$$

However, the following are not a terms: λ (only allowed as part of a λ -abstraction), $(\lambda(\text{pred } x). (S\ x))$ (only a variable is allowed between the λ and the $.$ in an abstraction).

An important fragment of PCF is the *pure λ -calculus*.

Definition 4.2. *The **pure terms** (of the pure, untyped λ -calculus) are those terms of PCF that don't contain any occurrences of constants. That is, they are expressed by the grammar:*

$$M, N ::= x \mid (\lambda x. M) \mid (MN)$$

Examples of pure terms are: $x, (x\ y), (\lambda x. (x\ y))$, and so on – they *only* involve variables, lambda abstraction and application! We will mostly concentrate our efforts on PCF, but occasionally we will discuss pure terms specifically.

Definition 4.3. A *subterm* of a term M is any substring of M that is itself a term, except for the variable that occurs in an abstraction between the λ and the dot if it does not occur outside of this position. This always includes M itself.

For example, $(\lambda x.(xy))$ has x as a subterm, but $(\lambda x.y)$ does not have x as a subterm. The subterms of $(\lambda x.(((\text{ifz } x) Z)(S Z)))$ are:

$(\lambda x.(((\text{ifz } x) Z)(S Z))), ((\text{ifz } x) Z)(S Z), ((\text{ifz } x) Z), (\text{ifz } x), \text{ifz}, x, Z, (S Z), S$

You may have noticed that we are using rather more parentheses than we did in the informal introduction. According to the grammar we need parentheses around every function application $(M N)$, this means that, for example, we need to write the term $\text{ifz } x \ y \ (f \ (\text{pred } x) \ (S \ y))$ which you saw earlier, more officially as $(((((\text{ifz } x) y) (f (\text{pred } x)) (S y))))$ – ouch!

Syntactic Sugar

We do need all these parentheses in order to specify a simple but unambiguous grammar. For example, without parentheses we wouldn't know whether a triple application $M N P$ should be parsed as $(M N) P$ – i.e. supplying P as input to the function returned by $M N$ – or as $M (N P)$ – i.e. supplying $N P$ as argument to the function M . However, it is rather inconvenient, so what we shall do is to agree some conventions between us for how to understand expressions when we don't put in quite so many parentheses.

Definition 4.4. We adopt the following conventions regarding the omission of parentheses and the handling of λ :

- We will omit any outermost parentheses. For example, whenever we write $M N$ to mean a term, the term that we mean is $(M N)$.
- In any term, we will assume that application associates to the left. For example, whenever we write $M N P$, the term that we mean is $((M N) P)$.
- In any term, we will assume that the body of an abstraction extends as far to the right as possible. For example, when we write $\lambda x.M N$, the term that we mean is $(\lambda x.(M N))$.
- In any term, iterated abstractions can be grouped. For example, when we write $\lambda x y.M$, the term that we mean is $(\lambda x.(\lambda y.M))$

It is still perfectly ok to write a term with all the parentheses, e.g. $((x \ y) \ z)$, but generally we will prefer to use as few parentheses as possible because it tends to be more readable. One exception to this is that we will always wrap a λ -abstraction in parentheses when it is an argument to some other function, for example, according to the agreed conventions we *can* write $(\lambda x.x) \ \lambda y.Z$ to mean $((\lambda x.x) (\lambda y.Z))$, but in practice I think everyone finds $(\lambda x.x) (\lambda y.Z)$ to be more readable.

While we are making life a bit more bearable, let us also agree some notation for numbers.

Definition 4.5. For each natural number n , the PCF numeral \underline{n} is the term defined recursively as follows:

$$\begin{aligned}\underline{0} &= Z \\ \underline{n+1} &= S \underline{n}\end{aligned}$$

For example, the numeral for three can be computed as follows:

$$\underline{3} = \underline{0 + 1 + 1 + 1} = S \underline{0 + 1 + 1} = S (S \underline{0 + 1}) = S (S (S \underline{0})) = S (S (S Z))$$

Now we can write some quite readable programs and know how to parse them:

<u>What we write</u>	<u>What we mean</u>
$x\ y\ (\text{pred } y)$	$((x\ y)\ (\text{pred } y))$
$\text{ifz } x\ \underline{2}\ \underline{3}$	$((\text{ifz } x)\ (S\ (S\ Z)))\ (S\ (S\ (S\ Z))))$
$\lambda f\ g\ x.\ f\ (g\ x)$	$(\lambda f.\ (\lambda g.\ (\lambda x.\ (f\ (g\ x)))))$
$\lambda x.\ x\ (\lambda y.\ y)$	$(\lambda x.\ (x\ (\lambda y.\ y)))$

5. Syntax: Binding and Substitution

In this lecture, we continue building up the syntax of the programming language by adding some additional structure that makes λ into a *binder*. Then we can properly define the most important operation associated with syntax, *substitution*.

5.1 Name Binding

Our intuition is that an abstraction $\lambda x.M$ is a function with formal parameter x that returns M . Intuitively, we think of the scope of x as being limited to M , so that if we see this term within something larger:

$$(\lambda x.M) Z \text{ (pred } x)$$

Then we expect the x occurring on the right to be totally unconnected to any x occurring in M . In particular, when we compute with this term, we are supplying Z as input to $\lambda x.M$, so every x in M will be replaced by Z , but not the x on the right:

$$(M[Z/x]) \text{ (pred } x)$$

Recall that our notation $M[N/x]$, which is only informal for now, is meant to mean “the term M but with every occurrence of x replaced by N ”.

In fact, in this larger term, the x does not appear to be connected to anything – there is no outer function that it is the parameter of. We are going to say that any x in M is *bound* by the λx and, conversely, that the x on the right, which is not bound by any lambda, is *free*.

Definition 5.1. The set of *free variables* of a term M is defined by recursion on the structure of M :

$$\begin{aligned} \text{FV}(x) &= \{x\} \\ \text{FV}(c) &= \emptyset \\ \text{FV}(PQ) &= \text{FV}(P) \cup \text{FV}(Q) \\ \text{FV}(\lambda x.N) &= \text{FV}(N) \setminus \{x\} \end{aligned}$$

A variable that occurs in a term and which is not free is said to be **bound**. A term M without free variables is said to be **closed** or a **combinator**. The set of all closed terms is written Λ^0 .

By contrast, a term M for which $FV(M) \neq \emptyset$ is said to be *open*. The term $(\lambda x. (\lambda y. (yx)))$ is closed, because:

$$\begin{aligned} FV(\lambda x. (\lambda y. (yx))) &= FV(\lambda y. (yx)) \setminus \{x\} \\ &= FV(yx) \setminus \{x, y\} \\ &= (FV(x) \cup FV(y)) \setminus \{x, y\} \\ &= (\{x\} \cup \{y\}) \setminus \{x, y\} \\ &= \emptyset \end{aligned}$$

We can think of an open term as being incomplete. Contrast x with $(\lambda x. x)$. We think of $(\lambda x. x)$ as being a function waiting for an input which, when supplied, will return it. Without more information we have no idea what the term x is supposed to do.

Now, since the scope of x in $\lambda x. M$ is completely delimited (i.e. we know that all the x bound by λx occur inside M), it doesn't really matter how we choose to name this formal parameter. For example, we could use y instead: $\lambda y. M[y/x]$ (as long as we rename x to y consistently everywhere inside M); $\lambda x. x$ is not different from $\lambda y. y$ in any essential way. This act of renaming is called a change of bound variable names, and it can be defined more precisely as follows.

Definition 5.2. Suppose $\lambda x. P$ is a term and $y \notin FV(P)$. Then the act of replacing $\lambda x. P$ by $\lambda y. P[y/x]$ is called a **change of bound variable name**. If two terms M and N can be made identical just by changing bound variable names to fresh names (not already used inside them), we say they are **α -equivalent**.

Intuitively, the terms $\lambda x. xy$ and $\lambda z. zy$ are α -equivalent because, if I pick a fresh variable y' to use instead of x in the body of the first term and use the same y' instead of z in the second term, the terms become equal: $\lambda y'. y'y$.

Most of the time it will be obvious when two terms are α -equivalent, such as the case with $\lambda x. x$ and $\lambda y. y$, or $\lambda xy. x$ and $\lambda yx. y$. However, it is absolutely crucial that we only allow renaming of *bound* variables in the definition. The term $\lambda x. xy$ is *not* α -equivalent with $\lambda x. xz$ (assuming $y \neq z$).

You should think of free variables as being part of the explicit interface to the term. For example, consider the closed term $\lambda yz. (\lambda x. xy)(zy)$. The body of this abstraction consists of two subterms $\lambda x. xy$ and zy , with the former applied to the latter. We can think of these two subterms interacting via the variable y , which is free in both of them (but bound in the term of which they are both a part). On the other hand, a variable that is bound in one of these subterms, like x in the first, cannot possibly interact with the other because its whole lifetime exists within the abstraction, which is completely enclosed by that subterm.

Now, it is often the case that we would really like to do a quick change of bound variable names in order to make life easier for ourselves. This happens when computing with terms and when doing hypothetical reasoning about them. In fact, because it happens so frequently, and because α -equivalent terms are essentially the same anyway, we will just build α -equivalence into our definition of identity for terms:

From now on we will consider α -equivalent terms to be identical.

The force of this statement is that, in the rest of this unit, we will not distinguish between e.g. $\lambda x. x$ and $\lambda y. y$, as far as we are concerned *they are the same string*. The situation is analogous to what you are used to with fractions and ratios. There are many different ways of writing down the same ratio $\frac{1}{2}$, $\frac{2}{4}$, $\frac{3}{6}$ and so on, but we consider them to be identical. And, as with fractions, we will feel

free to choose the most convenient representative at any given moment – for fractions this is usually the fraction that is in reduced form, where the numerator and denominator have no common factor – for terms this is usually the term in which all bound variables have different names to all the free variables.

Moreover, in proofs, if we have to reason about an arbitrary term M , we can assume that it is in this most convenient form (just as when we reason about an arbitrary fraction $\frac{m}{n}$ we can assume that m and n have no common factors). This is so helpful it gets a name:

The Variable Convention. If M_1, \dots, M_k occur within the same scope, then in these terms you may assume that all bound variables are chosen to be different from any others.

You will appreciate this convention more fully once we have done some proofs where it is useful.

5.2 Substitution

So far we have been working with an informal idea of substitution, which is that the term $M[N/x]$ is just the term M with all occurrences of x replaced by N . This simple understanding is perfectly fine for computing $M[N/x]$ (actually performing the substitution) most of the time, but there are a couple of cases where this substitution interacts with binding that require a more sophisticated understanding. It is a little bit subtle, but one thing to take away is that *these cases can always be avoided by renaming bound variables to ensure that no bound variable name clashes with any other*.

First, let's recap the simple case in which there are no bound variables. Here, to form $M[N/x]$ we just replace every occurrence of x in M by N :

$$(\text{ifz } x \text{ } 1 \text{ } (\text{pred } x))[\underline{2}/x] = \text{ifz } \underline{2} \text{ } 1 \text{ } (\text{pred } \underline{2})$$

If there are no occurrences of x , then by following this recipe one will just obtain M , i.e. $M[N/x] = M$:

$$(\text{ifz } y \text{ } 1 \text{ } (\text{pred } y))[\underline{2}/x] = \text{ifz } y \text{ } 1 \text{ } (\text{pred } y)$$

Avoid accidentally destroying existing binding structure

However, when the term M contains bound variables, blindly replacing in this way is not compatible with our assertion that terms are identical if they differ only in the names of their bound variables. For example, we consider $\lambda x. x$ and $\lambda y. y$ to be identical, they are the very same term (just as $\frac{2}{4}$ and $\frac{1}{2}$ are the very same rational number). However, if we try to substitute by blindly replacing for variables we meet an inconsistency:

$$\begin{aligned} (\lambda x. x)[\underline{1}/x] &\stackrel{?}{=} \lambda x. \underline{1} \\ (\lambda y. y)[\underline{1}/x] &\stackrel{?}{=} \lambda y. y \end{aligned}$$

If $\lambda x. x$ and $\lambda y. y$ are really the very same term, then it doesn't make sense that when we perform the same action on each of them we end up with different results – and the results *are* different, because the strings $\lambda x. \underline{1}$ and $\lambda y. y$ are not α -equivalent.

For this reason, we make the first refinement to our simple notion of substitution:

To calculate $M[N/x]$, replace all *free* occurrences of x by N .

So, in our notion of substitution, bound variables are not eligible for replacement. This makes sense, because you should think of the bound variables of M as being internal or private to M ; they are not part of its interface to the wider world. Thus, the second of the two equations above is correct, but not the first.

Avoid accidentally introducing new binding structure

Unfortunately, there is another, yet more subtle problem that can arise, which is called *variable capture*. This problem can happen when the term that is substituted-in, the N in $M[N/x]$, contains a free variable that clashes with the name of a bound variable in M . It seems a bit convoluted when I put it like that, so let's look at an example:

$$\begin{aligned} (\lambda x.z)[S\ x/z] &\stackrel{?}{=} \lambda x.S\ x \\ (\lambda y.z)[S\ x/z] &\stackrel{?}{=} \lambda y.S\ x \end{aligned}$$

As before, we are looking at *the same substitution* just written in two different ways: since $\lambda x.z$ and $\lambda y.z$ are identical terms, then when we perform the same replacement, namely replacing all *free* occurrences of z by $S\ x$ in each of them, we ought to end up with identical terms again. However, we don't: $\lambda x.S\ x$ and $\lambda y.S\ x$ are not two ways of writing the same term.

What has gone wrong this time is that, in calculating the substitution $M[N/x]$, a free variable in N got captured by a binder in M . Here, the free variable x in $S\ x$ ended up as a bound variable in the end result $\lambda x.S\ x$.

To avoid this problem, we simply disallow substitution in which capture might happen. So:

$$(\lambda y.z)[S\ x/z] = \lambda y.S\ x$$

is legitimate, because there are no name clashes between the free variables of $S\ x$ and the particular names of the bound variables $\lambda y.z$. However, $(\lambda x.z)[S\ x/z]$ is undefined, because a free variable would be captured as a result of replacement.

Capture-avoiding substitution

Consequently, this notion of substitution is often called *capture-avoiding*, since capture becomes impossible under the following definition.

Definition 5.3. We define *capture-avoiding substitution* of term N for variable x in term M , written $M[N/x]$, recursively on the structure of M :

$$\begin{aligned} c[N/x] &= c \\ y[N/x] &= y && \text{if } x \neq y \\ y[N/x] &= N && \text{if } x = y \\ (PQ)[N/x] &= P[N/x]Q[N/x] \\ (\lambda y.P)[N/x] &= \lambda y.P && \text{if } y = x \\ (\lambda y.P)[N/x] &= \lambda y.P[N/x] && \text{if } y \neq x \text{ and } y \notin \text{FV}(N) \end{aligned}$$

Notice the last two cases. The first says that we only substitute for free occurrences of x : if you attempt to substitute inside a term where x is bound, such as $\lambda x.P$, then there will be no free occurrences of x that can be replaced and so the result of the substitution $(\lambda x.P)[N/x]$ is just $\lambda x.P$. The second condition avoids the possibility of capture, you may only replace a free occurrence of x by N in $\lambda y.P$ if you know that y does not itself occur freely in N .

However, note that both of these conditions concern some bound variable y and so *one can always avoid any problems by first performing a renaming of bound variables* to change the bound variable y in $\lambda y.P$ for some other variable z that we haven't already used to obtain $\lambda z.P[z/y]$.

There is a potential for confusion due to the interaction between the notation for substitution $M[N/x]$ and our conventions for omitting parentheses when writing terms. So let's avoid it by agreeing that $MN[P/x]$ means M applied to $N[P/x]$ and $\lambda y.M[N/x]$ means an abstraction of y with body $M[N/x]$. Here are some examples of substitutions, make sure you follow the notation!

$$\begin{aligned}
(\lambda x.yx)[(\lambda z.z)/y] &= \lambda x.(yx)[\lambda z.z/y] \\
&= \lambda x.y[\lambda z.z/y]x[\lambda z.z/y] \\
&= \lambda x.(\lambda z.z)x \\
((\lambda x.yx)x)[y/x] &= (\lambda x.yx)[y/x]x[y/x] \\
&= (\lambda x.yx)y \\
(x(\lambda z.z)x)[\lambda z.z/x] &= x[\lambda z.z/x](\lambda z.z)[\lambda z.z/x]x[\lambda z.z/x] \\
&= (\lambda z.z)(\lambda z.z[\lambda z.z/x])(\lambda z.z) \\
&= (\lambda z.z)(\lambda z.z)(\lambda z.z) \\
(xx)[\underline{2}/y] &= x[\underline{2}/y]x[\underline{2}/y] \\
&= xx
\end{aligned}$$

On the other hand $(\lambda x.yx)[\lambda z.xz/y]$ and $(y(\lambda x.x))[\lambda z.x/y]$ are both undefined (even though, intuitively, the latter is not problematic).

6. Semantics: One-Step Reduction

Now we say what each basic kind of computational step looks like.

6.1 One-Step Reduction

To describe the behaviour of the built-in functions, we need to distinguish based on the shape of their input. For example, the predecessor function behaves differently whether it is given Z or a term starting with the successor $S\ N$ or if it is given something like $(\lambda x. x)$ that does not look like a numeral at all.

We do this by listing all the patterns of function applications that can possibly make a step, the *redexes*, and next to each we list its *contraction*, which is result of making the step.

Definition 6.1. We define the following redex/contraction pairs M / N . Any instance of the term on the left of the $/$ is called a **redex** and the term on the right its **contraction**.

$$\text{pred } Z / Z \tag{6.1}$$

$$\text{pred } (S\ N) / N \tag{6.2}$$

$$\text{ifz } Z\ N\ P / N \tag{6.3}$$

$$\text{ifz } (S\ M)\ N\ P / P \tag{6.4}$$

$$(\lambda x. M)\ N / M[N/x] \tag{6.5}$$

$$\text{fix } M / M\ (\text{fix } M) \tag{6.6}$$

The terminology *redex* is a portmanteau of “reducible expression” and this kind of computation is sometimes called “reduction” because, in some sense, the expression (term) somehow gets simple – it reduces in complexity – as a result of performing a step. However, one has to take this with a pinch of salt because there are lots of occasions in which performing a step of reduction actually creates a term which seems much *more* complex than the original!

The above listing of redexes and their contractions tells us what is involved in computing with terms that are themselves redexes, but, more generally, we want to compute with terms that merely contain redexes, e.g.

$$\text{ifz } (\text{pred } \underline{2})\ M\ N$$

This term is not itself a redex – it cannot make a step at the top-level – but it contains a redex, namely $\text{pred } \underline{2}$. If we were to contract this redex to give $\underline{1}$, we would then be in a position to compute the result of the conditional, because $\text{ifz } \underline{1}\ M\ N$ is a redex (it’s contraction is N).

The way we allow or disallow computational steps depending on where redexes occur in a given term gives rise to different kinds of reduction – *call-by-name*, *call-by-value*, *weak-head reduction* and so on. In this unit we take what is probably the most straightforward-to-use of these choices, the notion that is sometimes called *full β -reduction*. In full β -reduction, a term M makes a computational step by *contracting exactly one redex anywhere inside M* .

Definition 6.2 (Preliminary). *If N is obtained from M by contracting exactly one redex anywhere inside M , then we write $M \triangleright N$ and say that M has made **one reduction step** to N .*

We have marked this a preliminary definition because I am, at the moment, just going to rely on your intuitive understanding of strings to make sense of the phrase “anywhere inside”. This will be perfectly fine for using this definition to actually perform computations, but we will want to be a bit more concrete about what forms “exactly one redex anywhere inside” can take when we want to reason about arbitrary computation steps later on.

For example, $\text{ifz } (\text{pred } \underline{2}) \ M \ N \triangleright \text{ifz } \underline{1} \ M \ N$ by contracting the redex $\text{pred } \underline{2}$. Some more examples:

- $(\lambda x. x) \ \underline{2} \triangleright \underline{2}$
- $\text{ifz } ((\lambda x. x) \ \underline{2}) \ (\text{pred } \underline{2}) \ (S \ y) \triangleright \text{ifz } \underline{2} \ (\text{pred } \underline{2}) \ (S \ y)$
- $\text{ifz } ((\lambda x. x) \ \underline{2}) \ (\text{pred } \underline{2}) \ (S \ y) \triangleright \text{ifz } ((\lambda x. x) \ \underline{2}) \ \underline{1} \ (S \ y)$
- $(\lambda x y. x) \ (\lambda z. z) \ \underline{1} \triangleright (\lambda y. (\lambda z. z)) \ \underline{1}$

Note that there can be several possible reduction steps that we can perform, depending on which exactly one redex we choose (as in the second and third bullets). We do *not* have any of the following:

- $(\lambda x. x) \ \underline{2} \triangleright (\lambda x. x) \ \underline{2}$ — requires contracting zero redexes.
- $\text{ifz } ((\lambda x. x) \ \underline{2}) \ (\text{pred } \underline{2}) \ (S \ y) \not\triangleright \text{ifz } \underline{2} \ \underline{1} \ (S \ y)$ — requires contracting two redexes at once.
- $\text{ifz } ((\lambda x. x) \ \underline{2}) \ (\text{pred } \underline{2}) \ (S \ y) \triangleright S \ y$ — requires contracting two redexes in sequence.
- $(\lambda x y. x) \ (\lambda z. z) \ \underline{1} \triangleright (\lambda x y. x) \ \underline{1}$ — the subterm $((\lambda z. z) \ \underline{1})$ is a redex, but it does *not* occur inside this term! With the proper parens this term is $((\lambda x y. x) \ (\lambda z. z)) \ \underline{1}$, *not* $(\lambda x y. x) \ ((\lambda z. z) \ \underline{1})$.

However, as we start to look at reduction more closely, it will be helpful to make that a bit more concrete. For that we will use the notion of *one-hole context*.

Definition 6.3. *The **one-hole contexts** are those strings that can be generated by the following grammar:*

$$C[] ::= [] \mid M \ C[] \mid C[] \ N \mid \lambda x. C[]$$

Given a context $C[]$, we write $C[M]$ for the term obtained by replacing the hole $[]$ by M . We can also compose contexts by filling the hole in one context with another one: we write $C_1[C_2[]]$ for the context that arises by replacing the hole in $C_1[]$ with the context $C_2[]$.

Intuitively, the one-hole contexts are just terms that are missing exactly one subterm somewhere (which shows up as a single occurrence of the hole $[]$), so to make a one-hole context, you just need to think of a term and then replace some subterm by the hole. Examples of one-hole contexts are as follows:

$$[], \quad \lambda xyz. x [] (y z), \quad [] \underline{2}, \quad \text{ifz } [] (S x) (\text{pred } y), \quad \text{fix } (\lambda f x. S (f (\text{pred } [])))$$

The key operation associated with one-hole contexts is filling the hole. We can fill the hole in some context $C[]$ with a *term* N to obtain a *term* $C[N]$. For example, if we fill the hole of $\lambda xyz. x [] (y z)$ with z we obtain $\lambda xyz. x z (y z)$. Note, this operation does not respect binding, so we are happy to let variable capture happen (as in this case). Or, we could fill the same hole with another context, e.g. with $[] \underline{2}$, to obtain $\lambda xyz. x ([] \underline{2}) (y z)$.

We can use one-hole contexts to say “anywhere inside M ” as follows.

Definition 6.4. Then *one-step reduction relation*, written infix by \triangleright , is defined as $M \triangleright N$ just if there is a context $C[]$ and a redex/contraction pair P / Q such that $M = C[P]$ and $N = C[Q]$.

Let’s revisit our one-step reductions from earlier and justify them according to this definition:

- $(\lambda x. x) \underline{2} \triangleright \underline{2}$ – here $C[] = []$.
- $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright \text{ifz } \underline{2} (\text{pred } \underline{2}) (S y)$ – here $C[] = \text{ifz } [] (\text{pred } \underline{2}) (S y)$.
- $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright \text{ifz } ((\lambda x. x) \underline{2}) \underline{1} (S y)$ – here $C[] = \text{ifz } ((\lambda x. x) \underline{2}) [] (S y)$.
- $(\lambda xy. x) (\lambda z. z) \underline{1} \triangleright (\lambda y. (\lambda z. z)) \underline{1}$ – here $C[] = [] \underline{1}$.

6.2 Simple Programming

Now that we have a clearer idea of how terms compute, we can do some simple programming. I want to start with some very simple programs that will, nevertheless, come in handy to illustrate important aspects of this language later.

These first programs are extremely simple, but quite well known in the theory of λ -calculus, where they are just known simply as combinators. On the left are the combinators and how we shall abbreviate them. On the right an illustrative example of how they behave.

$\underline{\text{id}} = \lambda x. x$	$\underline{\text{id}} M \triangleright M$
$\underline{\text{const}} = \lambda xy. x$	$\underline{\text{const}} M N \triangleright (\lambda y. M) N \triangleright M$
$\underline{\text{sub}} = \lambda xyz. xz(yz)$	$\underline{\text{sub}} M N P \triangleright \dots \triangleright M P (N P)$
$\underline{\text{div}} = \text{fix } \underline{\text{id}}$	$\underline{\text{div}} \triangleright \underline{\text{id}} \underline{\text{div}} \triangleright \underline{\text{div}}$

The term $\underline{\text{id}}$ is sometimes called the *identity* combinator, for obvious reasons. The term $\underline{\text{const}}$ is sometimes called the *constant* combinator because, after accepting its first argument, say M , it then behaves like the “constantly M function” $\lambda y. M$ (i.e. that returns M no matter what its argument).

The term sub is sometimes called the *substitution* combinator because, historically, it has a role in simulating a kind of substitution if you don't have λ -abstraction. For us, it will not be so useful except when we need a handy example of a term that duplicates one of its arguments. Finally, the term div is known as *divergence* because it implements a diverging computation – one that unproductively runs forever.

Let me point out here that id, const etc are *merely abbreviations* (you can think of them as macros) that allow us humans to avoid having to write out enormous terms. That is, we are not adding anything new to the programming language, terms are still only built from either variables, abstractions, applications or the six constants. At any point we can eliminate all such abbreviations by simply replacing them by their definition.

Note: it *does not make any sense* to “define” a circular abbreviation like:

$$\text{fact} = \lambda x. \text{ifz } x \text{ } \underline{1} \text{ } (\text{mult } x \text{ } (\text{fact } (\text{pred } x))) \quad \text{Nonsense!}$$

(even assuming that we had defined mult). If we allowed something like this, then how could we make sense of a “term” like fact y ? If we replace the “abbreviation” fact by its definition, then we just end up with another “term”, $(\lambda x. \text{ifz } x \text{ } \underline{1} \text{ } (\text{mult } x \text{ } (\text{fact } (\text{pred } x)))) y$, involving fact. We could replace that occurrence by the definition of fact, but all this would do is create another one and there will never be an end to it. We will see in two lectures time that we can, in fact, define the factorial function and we will discuss recursive definition in general, and it will become clear then exactly what is going on.

Our version of PCF doesn't have tuple datatypes, but we can define them as a kind of abbreviation. We are going to use functions to represent tuples, which seems a bit strange at first but it works out quite nicely. We will use a technique called “Church Encoding”, named after the inventor of λ -calculus, Alonzo Church.

Let's start by discussing pairs. Recall that a pair (M, N) packages up two data items, which for us will be represented by terms, in this case M and N . Programming languages with pairs have a way to construct the pair (M, N) , given the items M and N , and they have a way to deconstruct the pair, e.g. projections fst and snd that behave like:

$$\text{fst } (M, N) \triangleright \dots \triangleright M \quad \text{snd } (M, N) \triangleright \dots \triangleright N$$

In other words, there is a series of reduction steps (computation steps) that take us from $\text{fst } (M N)$ to arrive at M , and similarly for the second projection.

In the Church encoding, a pair (M, N) will be represented by a function $\lambda p. p M N$. If you think about it, then, assuming M and N are in normal form (which we expect if they are meant to be data items), then $\lambda p. p M N$ is the simplest way to construct a closed term in normal form that involves M and N . The idea is that, if we supply this function $\lambda p. p M N$ with $\lambda xy. x$ as input then the result of that computation will be to extract M and, if we supply this function with $\lambda xy. y$ as input, then the result will be to project out N :

$$\begin{aligned} (\lambda p. p M N) (\lambda xy. x) &\triangleright (\lambda xy. x) M N \triangleright (\lambda y. M) N \triangleright M \\ (\lambda p. p M N) (\lambda xy. y) &\triangleright (\lambda xy. y) M N \triangleright (\lambda y. y) N \triangleright N \end{aligned}$$

So all we need to do is define the projections fst and snd so that they take a pair (encoded as a function) as input and pass to it $\lambda xy. x$ or $\lambda xy. y$ respectively:

$$\underline{\text{fst}} = \lambda p. p (\lambda xy. x) \quad \underline{\text{snd}} = \lambda p. p (\lambda xy. y)$$

To actually make a pair (M, N) , we just need a function that takes M and N and returns the encoding $\lambda p. p M N$:

$$\underline{\text{pair}} = \lambda m n p. p m n$$

Then we can check that everything works correctly by computing:

$$\begin{aligned} \underline{\text{fst}} (\underline{\text{pair}} M N) &\triangleright \dots \triangleright \underline{\text{fst}} (\lambda p. p M N) \triangleright \dots \triangleright M \\ \underline{\text{snd}} (\underline{\text{pair}} M N) &\triangleright \dots \triangleright \underline{\text{snd}} (\lambda p. p M N) \triangleright \dots \triangleright N \end{aligned}$$

So, indeed our encoding satisfies our specification of pairing in programming languages.

Let's extend this idea to tuples of every arity. We'll call the projection that extracts the i^{th} component of a k -ary tuple more generically $\underline{\text{proj}}_i^k$ and we'll call the function that constructs a k -ary tuple $\underline{\text{tuple}}^k$. So, for example, our $\underline{\text{fst}}$ arises as $\underline{\text{proj}}_1^2$, $\underline{\text{snd}}$ as $\underline{\text{proj}}_2^2$ and $\underline{\text{pair}}$ as $\underline{\text{tuple}}^2$.

Definition 6.5. For each natural number k and i , define k -ary pairing and projection functions as follows:

$$\begin{aligned} \underline{\text{proj}}_i^k &= \lambda p. p (\lambda x_1 x_2 \dots x_k. x_i) \\ \underline{\text{tuple}}^k &= \lambda x_1 x_2 \dots x_k p. p x_1 x_2 \dots x_k \end{aligned}$$

For ease of reading, let's agree to write each term of the form $\lambda p. p M_1 \dots M_k$ rather as $(M_1, \dots M_k)$.

7. Semantics: Reduction and Confluence

In this lecture we are going to look in more detail at reduction. In particular, we will see that reduction is *confluent*, which we can think of as a kind of eventual consistency. We will use these new tools in order to prove that certain programs are impossible to write.

7.1 Reduction

In the previous lecture we gave examples of reduction sequences, which were chains of steps one after another:

$$M_1 \triangleright M_2 \triangleright \cdots \triangleright M_{n-1} \triangleright M_n$$

Now we will formalise the idea that there exists such a reduction sequence leading from one term to another. In fact, a smooth way to do this is to introduce, for each natural number n , a relation $M \triangleright^n N$ to capture the idea that there exists a reduction sequence leading from M to N that comprises exactly n steps.

Definition 7.1. Define a family of *n-step reduction relations*, indexed by a natural number n and written $P \triangleright^n Q$, recursively as follows:

- $P \triangleright^0 Q$ just if $P = Q$.
- $P \triangleright^{k+1} Q$ just if there is some U such that $P \triangleright^k U$ and $U \triangleright Q$.

The *reduces-to* relation, written infix \triangleright^* , is defined by $M \triangleright^* N$ just if there is some n such that $M \triangleright^n N$. We can also write $M \triangleright^+ N$ to signify that $M \triangleright^n N$ with $n > 0$.

This definition captures our aim for the notation $M \triangleright^k N$ because if you unfold it for any specific k you will find it requires a k -step reduction sequence. For example, with $k = 3$:

$$\begin{aligned} & M \triangleright^3 N \\ & \text{just if there is some } U_2 \text{ such that } M \triangleright^2 U_2 \triangleright N \\ & \text{just if there is some } U_2, U_1 \text{ such that } M \triangleright^1 U_1 \triangleright U_2 \triangleright N \\ & \text{just if there is some } U_2, U_1, U_0 \text{ such that } M \triangleright^0 U_0 \triangleright U_1 \triangleright U_2 \triangleright N \\ & \text{just if there is some } U_2, U_1 \text{ such that } M \triangleright U_1 \triangleright U_2 \triangleright N \end{aligned}$$

So, for most purposes, one can think of the notation $M \triangleright^* N$ just to mean “there is a reduction sequence that leads from M to N ”. Note, this does not imply anything about N being terminal in

any way – there may still be reduction steps possible from N . The first few non-examples of one-step reductions from the last lecture are all examples of many-step reduction:

- $\text{id } \underline{2} \triangleright^* \text{id } \underline{2}$
– because there is a 0-step sequence: $\text{id } \underline{2} \triangleright^0 \text{id } \underline{2}$.
- $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright^* \text{ifz } \underline{2} \underline{1} (S y)$
– because there is e.g. a 2-step sequence: $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright^2 \text{ifz } \underline{2} \underline{1} (S y)$.
- $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright^* S y$
– because there is e.g. a 2-step sequence: $\text{ifz } ((\lambda x. x) \underline{2}) (\text{pred } \underline{2}) (S y) \triangleright^2 S y$.

The following terminology is not specific to PCF, but you will find it associated with any kind of computation by rewriting/reduction. It is very helpful for giving us a means to discuss reduction concisely.

Definition 7.2. *Some terminology:*

- If $M \triangleright^* N$ then the term N is said to be a **reduct** of M .
- If $M \triangleright^+ N$ then the term N is said to be a **proper reduct** of M .
- A term M without proper reduct, i.e. for which there is no N such that $M \triangleright N$, is a **normal form**.
- A term M that can reduce to normal form, i.e. for which there is a normal form N such that $M \triangleright^* N$ is said to **have a normal form** or be **normalisable**.
- A term M that has no infinite reduction sequences is said to be **strongly normalisable**.

Note: every term in normal form has a normal form and is strongly normalising. Note also that it is possible to have a normal form and to have infinite reduction sequences at the same time. Some examples to illustrate:

- The reducts of $(\lambda x. (\lambda z. x) y y) (\lambda x. x y)$ are:
 - (i) $(\lambda x. (\lambda z. x) y y) (\lambda x. x y)$
 - (ii) $(\lambda x. x y) (\lambda x. x y)$
 - (iii) $(\lambda z x. x y) y y$
 - (iv) $(\lambda x. x y) y$
- Except (i), all those reducts are proper.
- Every numeral, e.g. $\underline{2}$, is in normal form and, therefore, also has a normal form and is strongly normalising.
- The term $(\lambda x. (\lambda z. x) y y) (\lambda x. x y)$ has a normal form, namely $y y$. Since there are no infinite reduction sequences starting from this term, it is strongly normalising.

- The term $\underline{\text{const } 1} \text{ div}$ has a normal form, namely $\underline{1}$, which can be obtained by, e.g.:

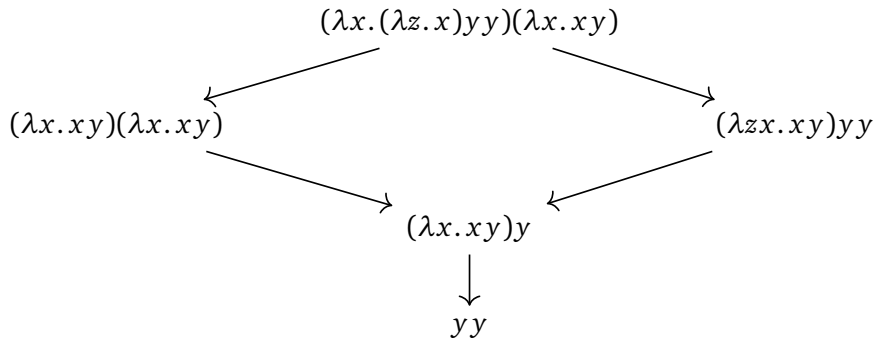
$$\underline{\text{const } 1} \text{ div} \triangleright (\lambda y. \underline{1}) \text{ div} \triangleright \underline{1}$$

However, it also has an infinite reduction sequence, e.g.:

$$\underline{\text{const } 1} \text{ div} \triangleright \underline{\text{const } 1} (\text{id div}) \triangleright \underline{\text{const } 1} (\text{id} (\text{id div})) \triangleright \underline{\text{const } 1} (\text{id} (\text{id} (\text{id div}))) \triangleright \dots$$

7.2 Confluence

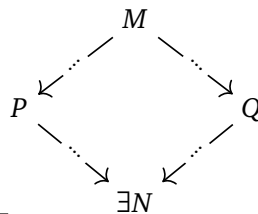
There is a sense in which one-step reduction is *nondeterministic*: as we saw in the earlier examples, for certain terms there is more than one possible reduct. This can be pictured very nicely if we draw the reduction graph of a term. The reduction graph of a term P is just a rooted directed graph. The vertices of the graph are P and all of the terms that can be reached by making one-step reductions from P , one-step reductions from those terms and so on. There is an edge between two of these vertices M and N just if $M \triangleright N$. For example, the reduction graph of $(\lambda x. (\lambda z. x)yy)(\lambda x. xy)$:



However, in this example, we can see that, although there are several possible ways to reduce $(\lambda x. (\lambda z. x)yy)(\lambda x. xy)$, they all eventually lead to yy . This is actually an important property of reduction, called *confluence*. We will state the theorem now and you will be expected to make use it from now on (i.e. you may assume it), but we will not prove it¹.

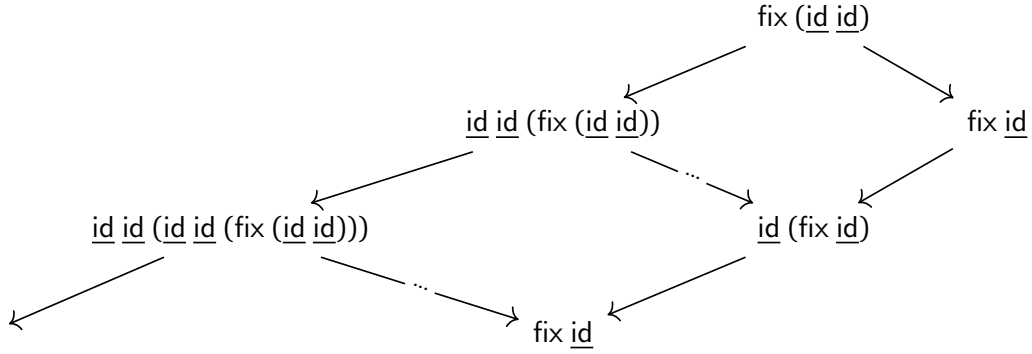
Theorem 7.1 (Confluence). *If $M \triangleright^* P$ and $M \triangleright^* Q$ then there exists a term N (not necessarily a normal form) such that $P \triangleright^* N$ and $Q \triangleright^* N$.*

This can be nicely illustrated by the following diagram, in which we use dots to indicate that more than one step may be involved:



¹Before the structure of TB1 changed so that we only had 7 weeks of teaching, we used to prove this - it is not difficult with the right definitions in place, but now there is not enough time.

It is a bit more nuanced than simply “if there are two different ways to reduce a term then they will always reach the same normal form” because it applies even when the term in question fails to have a normal form. Consider e.g. $\text{fix } (\underline{\text{id}} \ \underline{\text{id}})$. An illustration:



Use two fingers to trace out any two paths on this reduction graph for a while. Wherever each of your fingers got to, there is guaranteed to be a way to get to a common reduct.

There are many helpful consequences of confluence, but one particularly important one is that normal forms are unique when they exist.

Theorem 7.2 (Unique Normal Forms). *If $M \triangleright^* N$ and $M \triangleright^* P$ with N and P normal forms, then $N = P$.*

Proof. Suppose $M \triangleright^* N$ and $M \triangleright^* P$. Then, by confluence, there is a common reduct Q , i.e. $N \triangleright^* Q$ and $P \triangleright^* Q$. However, N is a normal form, so it must be that $N \triangleright^0 Q$, and P is a normal form, so $P \triangleright^0 Q$. Hence, by definition, $N = Q = P$. \square

Using the knowledge that terms have unique normal forms will actually allow us to exclude a variety of functions.

Theorem 7.3. *There is no PCF term IsNat such that, for all M :*

$$\text{IsNat } M \triangleright^* \begin{cases} \underline{1} & \text{if } M \text{ is a numeral} \\ \underline{0} & \text{otherwise} \end{cases}$$

Proof. We are trying to prove a statement of the form “not (there exists a term...)” so we assume there exists a term IsNat that satisfies the given specification, i.e. we assume:

- (i) for all M : if M is a numeral, then $\text{IsNat } M \triangleright^* \underline{1}$
- (ii) for all M : if M is not a numeral, then $\text{IsNat } M \triangleright^* \underline{0}$

Then we aim to find a contradiction. To that end, let us consider the behaviour of the program $\text{IsNat } ((\lambda x. x) \ \underline{2})$. Since $(\lambda x. x) \ \underline{2}$ is not a numeral, it follows from (ii) that:

$$\text{IsNat } ((\lambda x. x) \ \underline{2}) \triangleright^* \underline{0}$$

However, we can also construct an alternative reduction sequence, starting from the same term:

$$\text{IsNat } ((\lambda x. x) \ \underline{2}) \triangleright^* \text{IsNat } \underline{2}$$

Now we can continue this sequence using (i), because $\underline{2}$ is a numeral:

$$\text{IsNat } ((\lambda x. x) \underline{2}) \triangleright^* \text{IsNat } \underline{2} \triangleright^* \underline{1}$$

By the unique normal forms theorem, we obtain that $\underline{0} = \underline{1}$ which is clearly absurd (they are not identical strings modulo bound variable renaming). \square

8. Computability: Recursion

We usually use the term recursive function for those whose definition is given in terms of itself. The classic example is the factorial function on natural numbers:

$$\begin{aligned}\text{fact} &: \mathbb{N} \rightarrow \mathbb{N} \\ \text{fact}(n) &= \text{if } n = 0 \text{ then } 1 \text{ else } n \cdot \text{fact}(n - 1)\end{aligned}$$

The presence of the function's name, `fact`, on both sides of the `=` is the tell-tale sign that this is a recursive function. So, it seems that names are important for giving recursive function definitions. In λ -calculus there is no formal support for attaching names to functions — no syntax for naming that is built into the notion of reduction or conversion. So a question arises: can we even define recursive functions in the λ -calculus and, if not, how can we perform repetitive computations?

8.1 Recursive functions

To answer this question we have to understand more fully what we mean when we say that the above equation is a definition of the factorial function. As computer scientists, and especially as functional programmers, we are a bit spoiled by the very liberal notion of function definition given in languages like Haskell. In Haskell every equation of the form $f\ x = e$ (for some Haskell expression e), assuming it type-checks, is a valid function definition. Consequently, we may be led to believe that, in general, if we write a name f followed by a formal parameter x and then `=` followed by some expression e , we have defined a function.

However, this is not the case. Consider the following “definition” of a function `yuck` over the integers:

$$\begin{aligned}\text{yuck} &: \mathbb{Z} \rightarrow \mathbb{Z} \\ \text{yuck}(n) &= \text{yuck}(n) + 1\end{aligned}$$

Is this a valid function definition? The answer is no: we can prove it. Assume, for the purposes of obtaining a contradiction, that there is a function `yuck` on integers with the property that the following integer equation is true: $\text{yuck}(n) = \text{yuck}(n) + 1$. Then we can subtract $\text{yuck}(n)$ from both sides, to obtain the true equation $0 = 1$, which is absurd.

To understand this, we have to be clear on what we mean by the word function. Intuitively, we think of a function f as some black box which associates to each input x an output $f(x)$. This intuition is made precise through the *definition* of function in Zermelo Fraenkel (ZF) set theory, and this is usually the meaning of the word function in mathematics more generally. In ZF, a function f from a set A to a set B is a subset of $A \times B$ in which each element of A is associated with exactly one element of B .

According to this definition, a function is just a tabulation of its inputs and outputs, one can say that the notion of function and the notion of *graph of a function* are identified. So, strictly speaking, the factorial function is a certain set of pairs, starting off like:

$$\begin{array}{lcl} 0 & \mapsto & 1 \\ 1 & \mapsto & 1 \\ 2 & \mapsto & 2 \\ 3 & \mapsto & 6 \\ & \vdots & \end{array}$$

Since there is exactly one output for each input, the sentence “the output of fact corresponding to input 3” unambiguously refers to exactly one natural number, 6. However, it is a bit long-winded to say, so in 1734 Leonard Euler invented a more concise notation: `fact(3)`.

Unfortunately, the department’s budget does not stretch far enough to allow me to purchase enough paper to give a complete listing of the factorial function. So we have to look for some finite means by which we can *specify* this infinite set, unambiguously. This is purpose of the equation. The equation constitutes a *definition* of factorial in the following sense.

Definition 8.1. *We define fact as the unique function on natural numbers (i.e. subset of $\mathbb{N} \times \mathbb{N}$ with the property that exactly one output number is associated with each input number) that satisfies the equation (in parameter F):*

$$F(n) = \text{if } n = 0 \text{ then } 1 \text{ else } n \cdot F(n - 1)$$

So, in a sense, Definition 8.1 is the “official” definition of the factorial function, but, in practice we omit the parameter F and instead write the equation directly using `fact`, as we did at the start of this chapter. It should remind you of the way that we make inductive definitions of sets.

Of course, if we want to make such a definition then there is some obligation on us to be sure that the equation has a unique solution. However, we know a few syntactic tricks that make it often straightforward to tell. For example, if we do write the equation in the form $f(x) = e$ and there is a well-founded order on the domain and we only make recursive calls in e on strictly smaller values, then we can be sure that this equation specifies a genuine function.

We can now see what is the key mistake in our “definition” of `yuck`: there are *no* functions over integers that satisfy the equation (in parameter F):

$$F(n) = F(n) + 1$$

No wonder that by assuming that one exists, named `yuck`, we are led inevitably to a contradiction!

Before we leave this excursion into the world of recursive definitions outside of the λ -calculus, let me return briefly to Haskell. The equation:

$$\begin{array}{lcl} \text{yuck} & :: & \text{Int} \rightarrow \text{Int} \\ \text{yuck } n & = & \text{yuck } n + 1 \end{array}$$

has a solution in the space of Haskell programs. There is a Haskell program that satisfies these equations – the program that never terminates on any input!

8.2 Recursive functions in PCF

Now, let's turn to PCF. To understand when a PCF term satisfies an equational specification, we should first say what we mean by an equation between PCF programs.

When you learned about functional programming, you probably did some short proofs using equational reasoning, like:

$$\begin{aligned}\text{map } (+1) [3, 8, 6] &= 4 : \text{map } (+1) [8, 6] \\ &= 4 : 9 : \text{map } (+1) [6] \\ &= 4 : 9 : 7 : \text{map } (+1) [] \\ &= 4 : 9 : 7 : [] \\ &= [4, 9, 7]\end{aligned}$$

What do we really mean by $\text{map } (+1) [3, 8, 6] = 4 : \text{map } (+1) [8, 6]$? When I am doing arithmetic, I know what I mean by $3 + 1 = 4$, I mean that the left and right sides of the $=$ are literally the same natural number, they are identical: $3 + 1$ is a number and it is the very same number has 4. However, that's not what I mean in $\text{map } (+1) [3, 8, 6] = 4 : \text{map } (+1) [8, 6]$: it's plain to see that the left and right sides of this equation, which are Haskell expressions (i.e. strings), are *not* identical: the string on the left starts with 'm' whereas the one on the right starts with '4'.

You might complain at that and say that we are not talking about Haskell expressions as strings, but rather Haskell expressions as (compiled) programs, and that the left and right hand sides are identical programs. However, that is clearly untrue also, since the program on the left has more steps to evaluate than the one on the right before reaching a normal form.

What we need is some kind of "equality modulo computation".

We could perhaps make sense of these equations by justifying them by saying $\text{map } (+1) [3, 8, 6] = 4 : \text{map } (+1) [8, 6]$ because the left-hand side *evaluates to* the right-hand side. However, this too seems like a dishonest way to interpret an equation, because if

$$\text{map } (+1) [3, 8, 6] = 4 : \text{map } (+1) [8, 6]$$

then it should also be that

$$4 : \text{map } (+1) [8, 6] = \text{map } (+1) [3, 8, 6]$$

but surely it is not the case that $4 : \text{map } (+1) [8, 6]$ evaluates to $\text{map } (+1) [3, 8, 6]$.

A more satisfying approach is to say that two expressions are equal modulo computation when they *each evaluate to the same thing*. We say that they are *convertible*.

Definition 8.2. We write $M \approx N$ just if there is a term P such that $M \triangleright^* P$ and $N \triangleright^* P$. In this case we say that they are *convertible*.

For short, we will sometimes rephrase this as "convertible terms have a common reduct". The following are examples of convertible pairs of terms:

- $\text{ifz } \underline{0} (\text{pred } \underline{3}) \underline{1} \approx \underline{2}$ – the common reduct is $\underline{2}$
- $\text{ifz } x (\text{pred } \underline{3}) \underline{1} \approx \text{ifz } x \underline{2} (\text{pred } (\text{pred } \underline{3}))$ – a common reduct is $\text{ifz } x \underline{2} \underline{1}$.

In general, conversion is weaker than reduction, in the sense that $M \triangleright^* N$ implies $M \approx N$, but not necessarily vice-versa.

What is perhaps not so clear from this definition is that conversion behaves anything like an equality, that it is an equivalence. However, using confluence, it is quite straightforward to show.

Lemma 8.1. *Convertibility is a congruence, that is an equivalence relation that is compatible with contexts, i.e. all of the following hold:*

Reflexivity For all M : $M \approx M$.

Symmetry For all M, N : $M \approx N$ implies $N \approx M$.

Transitivity For all M, N and P : $M \approx P$ and $P \approx N$ implies $M \approx N$.

Compatibility For all M, N and $C[\]$: if $M \approx N$ then $C[M] \approx C[N]$.

It makes sense, for example, to ask: *is there a PCF-term M satisfying:*

$$Mx \approx x(Mx)$$

I hope you agree that, if there is such a term M then, we can think of it as a recursive function. In fact, there is a term satisfying this equation, fix , because $\text{fix } x \approx x(\text{fix } x)$.

Not all recursive specifications have solutions (terms that satisfy them) in PCF. For example, the following specification has no solution for M :

$$(\lambda x. \lambda y. y)M \approx (\lambda x. \lambda yz. y)M$$

How can you argue it? To show that it has *no* solution, we assume it has one, so suppose M is a solution. Then we have $(\lambda x. \lambda y. y)M \triangleright^* (\lambda x. \lambda yz. y)M$ is true of this M . But we also know that the following is true:

$$(\lambda x. \lambda y. y)M \triangleright^* \lambda y. y \text{ and } (\lambda x. \lambda yz. y)M \triangleright^* \lambda yz. y$$

However, this would mean that:

$$\begin{array}{ll} (\lambda x. \lambda y. y)M & \triangleright^* \lambda y. y \\ (\lambda x. \lambda y. y)M & \triangleright^* (\lambda x. \lambda yz. y)M \triangleright^* \lambda yz. y \end{array}$$

But, since these two are distinct normal forms, this contradicts the unique normal forms theorem.

Anyway, it is still a meaningful question to ask whether this recursive equation has a solution, but it is perhaps not a natural way to ask about recursive function definitions.

It might be more reasonable to restrict our attention to equations of the form:

$$Mx_1, \dots, x_n \approx N$$

for some fixed term N (which may, of course, contain occurrences of M). The equation that we started with, $Mx \approx x(Mx)$ has this shape, with N being $x(Mx)$. We will see that, in PCF, all such recursive equations have a solution for M .

Fixpoints

The fact that we can always find solutions to such equations follows from one important property enjoyed by PCF terms. Although very simple, some people would say that it is *the* fundamental property of untyped programs.

We start by transferring the concept of fixed point (or fixpoint) of a function to λ -terms. On mathematical (set-theoretic) functions, the concept is defined as follows. Let A be a set and f be a function from A to A . Then $a \in A$ is said to be a *fixed point* of f just if $f(a) = a$. For example, 0 is a fixed point of the “doubling” function $z \mapsto 2 \cdot z$ on natural numbers, because $2 \cdot 0 = 0$. Some functions on naturals do not possess any fixed points, for example, there is no natural z that is equal to its own successor $z + 1$, so the successor function has no fixed points. On the other hand, every natural is a fixed point of the identity function on naturals.

Let’s transfer this idea to PCF.

Definition 8.3. *A term N is said to be a **fixed point** of another term M just if $MN \approx N$.*

According to this definition there are analogues to the examples that we discussed above, but it’s not a perfect correspondence. For example, every natural number is a fixed point of the identity:

$$\text{id } \underline{n} \approx \underline{n}$$

It’s also the case that no *numeral* is the fixed point of the successor function.

$$S \underline{k} \not\approx \underline{k}$$

However, this implementation of the successor function *does* possess a fixed point! I can tell you that without even looking at its definition, because of the following theorem.

Theorem 8.1 (First Recursion Theorem). *Every term has a fixed point.*

Proof. We claim that $\text{fix } M$ is a fixed point of M . We calculate:

$$\text{fix } M \triangleright M (\text{fix } M)$$

Hence, $M (\text{fix } M) \approx \text{fix } M$. □

The proof of the First Recursion Theorem tells us more, not only does every term possess a fixed point, but the fixed point of M can be computed by a PCF program, $\text{fix } M$. The primitive fix is an example of a *fixed point combinator*, so called because it computes the fixed point of another term.

Fixpoints and recursive functions

We can use fixed point combinators to solve certain kinds of recursive equations, including *all* those of shape:

$$Mx_1 \cdots x_n \approx C[M]$$

i.e. where the right-hand side may contain M , as well as occurrences of x_1, \dots, x_n . Let’s suppose M does not occur in $C[\]$. There is a certain recipe that we can follow in order to find M .

First, we can use the compatibility property of reduction to reduce the problem to finding an M that satisfies the following equation instead:

$$M \approx \lambda x_1, \dots, x_n. C[M]$$

This is because we can go from the second equation to the first one: we just apply x_1, \dots, x_n in order to both sides – this amounts to using compatibility with $D[] = [] x_1 \cdots x_n$.

Next, we can “abstract out M ”, we replace all the occurrences of M in there by the same fresh variable. Let’s say that variable f does not occur anywhere else, then $C[f]$ is the result of replacing all the occurrences of M in $C[M]$ by f .

Then it is sufficient to find an M that satisfies:

$$M \approx (\lambda f. \lambda x_1, \dots, x_n. C[f]) M$$

We know about this shape. The First Recursion Theorem guarantees that there is a solution, and fix allow us to describe it directly. For example, $\text{fix } (\lambda f. \lambda x_1, \dots, x_n. C[f])$ is an M that satisfies this equation and hence also satisfies the very first equation.

9. Computability: Effectiveness and Turing-Completeness

The material in this chapter will not be assessed.

In this lecture we will show that in PCF one can program every function that is “intuitively” or “effectively” computable.

9.1 Computability

It is worth pausing to remind ourselves exactly what is meant by this, and to do that I want to go back to the origins of computability.

In the late 1920s and early 1930s a great unsolved problem preoccupied mathematical logicians. This problem is called the *entscheidungsproblem*, which is German for “decision problem”. Nowadays we are very familiar with decision problems, but back then the phrase referred to one in particular, which was the problem of determining the validity of a statement of first-order logic. The fact that there might be a step-by-step method, requiring no insight on the part of the practitioner, and by which one could determine the truth of any mathematical conjecture was an *entirely plausible* to the mathematicians at the turn of the 20th century.

The *entscheidungsproblem* was one pillar of the famous programme of David Hilbert. Hilbert really put meta-mathematics on the map. In the half-century before, mathematicians had become interested in putting more rigour into mathematics. On the one hand they were concerned by certain infinite constructions, like the limits of infinite series that arise in analysis, on the other there had been many important developments in mathematical logic which put a greater emphasis on axiomatisation. Hilbert was interested in analysing the axiomatisations themselves to understand when axioms could be said to be independent, when systems of axioms could be said to be consistent or complete and whether there was a procedure that could be followed in order to calculate the consequences of a set of axioms — hence the *entscheidungsproblem*.

If you believed that the decision problem has a positive solution, you need only propose an algorithm that will solve any instance. Even though the digital computer would not be invented for many years to come, logicians of the 1930s were quite able to recognise an algorithm when they saw one — (roughly speaking) a procedure of discrete steps, requiring no insight and which, if followed faithfully, would yield the desired result. However, if you believed that the decision problem has a

negative solution things are a little more difficult. To show that there can be *no* algorithm capable of solving the problem requires that you can make precise the class of all algorithms. What else can you do? You would need to give a mathematical definition of the intuitive notion of, as it was sometimes called at the time, *effective calculability*, which was thought of as meaning the quality of being solvable by an algorithm.

One of the key difficulties was to propose a precise definition of the set of computable functions that the mathematical community would agree captures the fuzzy idea of being “intuitively” or “effectively computable”. Whatever definition is proposed, someone may challenge and say “I think there are functions that can be computed by an algorithm that lie outside the set you defined... we just haven’t thought a specific example yet...”. The trouble is that it is not possible to *prove* that a definition of the set of computable functions completely captures the intuitive idea of computability, because only the former lives in the world of mathematics – the latter is more of a gut a feeling.

Anyway, the first serious proposal for a definition of effective calculability was given by Alonzo Church, using his recently invented pure, untyped λ -calculus. That is, untyped PCF, but *without any constants* including fix! We now agree that the pure, untyped λ -calculus *does* capture all the intuitively computable functions, but at the time Church proposed it, it was not generally accepted. However, soon after, Alan Turing proposed his own notion of (what is now know as) Turing machine. Turing’s great achievement was to give an informal, but very convincing argument as to why his theoretical machines – a kind of souped up finite automata – could do any computation that can ever be done. He also proved that his machines and the λ -calculus expressed exactly the same class of functions, so if you believe his argument then you must also believe that the pure, untyped λ -calculus captures computability as well.

The definition of computable functions that you saw last year was a bit more modern, going via *While*-programs, so let me first remind you of what a While program looks like:

(Expressions) $e ::= n \mid x \mid e + e \mid e - e \mid e * e$
 (Booleans) $b ::= e = e \mid e \leq e \mid \text{true} \mid \text{false} \mid ! b \mid b \wedge b$
 (Statements) $s ::= \text{skip} \mid x := e \mid \text{if } b \text{ then } s_1 \text{ else } s_2 \mid \text{while } b \text{ do } s \mid s_1 ; s_2$

In the while language you studied last year, n could range over any integer, but for simplicity we will only consider natural numbers $n \in \mathbb{N}$. Consequently, we will need to say what it means to subtract a larger number from a smaller one (which would normally go negative) and we will opt to truncate at 0, so that e.g. $4 - 6 = 0$.

Recall that the *state* of a While program is a map σ from program variables to their values. However, we are going to simplify this a bit in our presentation by assuming that all the variables in a given While program are numbered: i.e. they are called x_1, x_2, x_3 and so on for some initial segment of the natural numbers. Then we can represent a state as a tuple of values (n_1, \dots, n_k) rather than a map $\sigma(x_i) = n_i$. For example:

We will use: $(1, 6, 23, 4, 19)$ rather than $\begin{bmatrix} x_1 \mapsto 1 \\ x_2 \mapsto 6 \\ x_3 \mapsto 23 \\ x_4 \mapsto 4 \\ x_5 \mapsto 19 \end{bmatrix}$

I want to remind you of a couple of important points. First, recall that, although Expressions e and Booleans b cannot change the state of the program (they cannot embed an assignment or some such, as is possible in e.g. C), to understand their value ($\llbracket e \rrbracket$ and $\llbracket b \rrbracket$ respectively), we do need to supply the current state of the program, because they may contain variables:

$$\begin{aligned} \llbracket 3 * 4 \rrbracket(1, 6, 23, 4, 19) &= 12 & \llbracket \text{true} \wedge ! \text{false} \rrbracket(1, 6, 23, 4, 19) &= \text{true} \\ \llbracket 3 + (x_2 * x_4) \rrbracket(1, 6, 23, 4, 19) &= 27 & \llbracket (x_3 = 23) \wedge (x_4 \leq x_2) \rrbracket(1, 6, 23, 4, 19) &= \text{true} \\ \llbracket x_5 + (x_2 - x_4) \rrbracket(1, 6, 23, 4, 19) &= 21 & \llbracket !(x_1 + 3 \leq x_4) \rrbracket(1, 6, 23, 4, 19) &= \text{false} \end{aligned}$$

Whenever a While program s can start in a certain state σ , compute and then terminate in state some state σ' , we write $\langle s, \sigma \rangle \Downarrow \sigma'$. This relation is defined by the big step semantics of While.

While Computability

We recall the notation $[x \mapsto n]$ for the state in which x is set to n and every other variable is set to 0 and we extend it to $[x_i \mapsto n_i \mid 1 \leq i \leq p]$, to stand for the state in which each variable x_i is set to n_i , for i between 1 and p , and all other variables are set to 0. We write \mathbb{N}^p for the set of all p -tuples (n_1, \dots, n_p) of natural numbers.

Definition 9.1. A While program s **computes** a partial function $f : \mathbb{N}^p \rightarrow \mathbb{N}$ on p -tuples of natural numbers (with respect to variables x_1, \dots, x_p) just if:

$$f(n_1, \dots, n_p) = m \quad \text{iff} \quad \langle s, [x_i \mapsto n_i \mid 1 \leq i \leq p] \rangle \Downarrow [x_1 \mapsto m]$$

We say that a partial function f on p -tuples of natural numbers is **While-computable** just if there is a While program s which computes it with respect to variables x_1, \dots, x_p .

Thesis 9.1 (Church-Turing : While). *The (intuitively) computable partial functions on natural numbers are exactly the While-computable functions (equivalently, the Turing-computable functions, the Haskell-definable functions etc).*

Now our aim is to show that every While program can be simulated faithfully by a term of PCF. Given a While program s over k variables, we will construct a PCF term $\langle s \rangle$ such that, for all possible states $\sigma = (n_1, \dots, n_k)$, $\langle s, \sigma \rangle \Rightarrow \langle s', \sigma' \rangle$ just if $\langle s \rangle \sigma \triangleright^* \langle s' \rangle \sigma'$. Note: here we are using σ as a While program state on the left of the “just if” and as a tuple of numerals (i.e. a PCF term) on the right.

Following the definition of the While syntax, we break our translation from While into PCF down into separate translations for Expressions, Booleans and Statements. Just as each of these syntactic classes requires a state in order to be evaluated correctly, the translation of an expression e , a Boolean b or a statement s will be a function $\langle e \rangle$, $\langle b \rangle$ and $\langle s \rangle$ respectively, that takes a state as input and delivers a state as output. In each case, the only difficulty is to thread the state of the computation through so that all the subterms have access to the values of the variables. Throughout we will assume we are dealing with a program of k variables.

Numeric Expressions

The translation of expressions is given recursively on the syntax:

$$\begin{aligned}
\llbracket n \rrbracket &= \text{const } \underline{n} \\
\llbracket x_i \rrbracket &= \text{proj}_i^k \\
\llbracket e_1 + e_2 \rrbracket &= \lambda z. \text{add } (\llbracket e_1 \rrbracket z) (\llbracket e_2 \rrbracket z) \\
\llbracket e_1 - e_2 \rrbracket &= \lambda z. \text{sub } (\llbracket e_1 \rrbracket z) (\llbracket e_2 \rrbracket z) \\
\llbracket e_1 * e_2 \rrbracket &= \lambda z. \text{mult } (\llbracket e_1 \rrbracket z) (\llbracket e_2 \rrbracket z)
\end{aligned}$$

Here the z parameter that shows up in the last three abstractions is meant to represent the state. This ensures that $\llbracket e \rrbracket(\sigma) = n$ iff $\llbracket e \rrbracket \sigma \triangleright^* \underline{n}$, in other words: a While expression e evaluates to n in a given state iff applying the PCF term $\llbracket e \rrbracket$ to the same tuple reduces to the numeral for n . For example we have:

$$\llbracket 3 + (x_2 * x_4) \rrbracket = \lambda z. \text{add } (\text{const } \underline{3} z) ((\lambda z. \text{mult } (\text{proj}_2^5 z) (\text{proj}_4^5 z)) z)$$

Then recall that $\llbracket 3 + (x_2 * x_4) \rrbracket(1, 6, 23, 4, 19) = 27$ and observe:

$$\begin{aligned}
&(\lambda z. \text{add } (\text{const } \underline{3} z) ((\lambda z. \text{mult } (\text{proj}_2^5 z) (\text{proj}_4^5 z)) z)) (1, 6, 23, 4, 19) \\
&\triangleright^* \text{add } (\text{const } \underline{3} (1, 6, 23, 4, 19)) ((\lambda z. \text{mult } (\text{proj}_2^5 z) (\text{proj}_4^5 z)) (1, 6, 23, 4, 19)) \\
&\triangleright^* \text{add } \underline{3} ((\lambda z. \text{mult } (\text{proj}_2^5 z) (\text{proj}_4^5 z)) (1, 6, 23, 4, 19)) \\
&\triangleright^* \text{add } \underline{3} (\text{mult } (\text{proj}_2^5 (1, 6, 23, 4, 19)) (\text{proj}_4^5 (1, 6, 23, 4, 19))) \\
&\triangleright^* \text{add } \underline{3} (\text{mult } \underline{6} \underline{4}) \\
&\triangleright^* \text{add } \underline{3} (\underline{24}) \\
&\triangleright^* \underline{27}
\end{aligned}$$

We don't have time to prove it here, but from now on we will assume that our encoding satisfies the following equivalence:

$$\llbracket e \rrbracket \sigma = n \quad \text{iff} \quad \llbracket e \rrbracket \sigma \triangleright^* \underline{n}$$

Boolean Expressions

The translation of Booleans is given recursively on the syntax:

$$\begin{aligned}
\llbracket e_1 = e_2 \rrbracket &= \lambda z. \text{eq } (\llbracket e_1 \rrbracket z) (\llbracket e_2 \rrbracket z) \\
\llbracket e_1 \leq e_2 \rrbracket &= \lambda z. \text{leq } (\llbracket e_1 \rrbracket z) (\llbracket e_2 \rrbracket z) \\
\llbracket \text{true} \rrbracket &= \text{const true} \\
\llbracket \text{false} \rrbracket &= \text{const false} \\
\llbracket ! b \rrbracket &= \lambda z. \text{not } (\llbracket b \rrbracket z) \\
\llbracket b_1 \wedge b_2 \rrbracket &= \lambda z. \text{and } (\llbracket b_1 \rrbracket z) (\llbracket b_2 \rrbracket z)
\end{aligned}$$

This ensures that, for all Boolean expressions b and all states σ : $\llbracket b \rrbracket(\sigma) = \text{true}$ iff $\langle b \rangle \sigma \triangleright^* \underline{\text{true}}$, and $\llbracket b \rrbracket(\sigma) = \text{false}$ iff $\langle b \rangle \sigma \triangleright^* \underline{\text{false}}$. For example, we have:

$$\langle \neg(x_1 + 3 \leq x_4) \rangle = \lambda z. \underline{\text{not}} ((\lambda z. \underline{\text{leq}} ((\lambda z. \underline{\text{add}} (\underline{\text{proj}}_1^5 z) (\underline{\text{const}} \underline{3} z)) z) (\underline{\text{proj}}_4^5 z)) z)$$

Then recall that $\llbracket \neg(x_1 + 3 \leq x_4) \rrbracket(1, 6, 23, 4, 19) = \text{false}$, and observe:

$$\begin{aligned} & (\lambda z. \underline{\text{not}} ((\lambda z. \underline{\text{leq}} ((\lambda z. \underline{\text{add}} (\underline{\text{proj}}_1^5 z) (\underline{\text{const}} \underline{3} z)) z) (\underline{\text{proj}}_4^5 z)) z)) (1, 6, 23, 4, 19) \\ & \triangleright^* \underline{\text{not}} ((\lambda z. \underline{\text{leq}} ((\lambda z. \underline{\text{add}} (\underline{\text{proj}}_1^5 z) (\underline{\text{const}} \underline{3} z)) z) (\underline{\text{proj}}_4^5 z)) (1, 6, 23, 4, 19)) \\ & \triangleright^* \underline{\text{not}} (\underline{\text{leq}} ((\lambda z. \underline{\text{add}} (\underline{\text{proj}}_1^5 z) (\underline{\text{const}} \underline{3} z)) (1, 6, 23, 4, 19)) (\underline{\text{proj}}_4^5 (1, 6, 23, 4, 19))) \\ & \triangleright^* \underline{\text{not}} (\underline{\text{leq}} ((\lambda z. \underline{\text{add}} (\underline{\text{proj}}_1^5 z) (\underline{\text{const}} \underline{3} z)) (1, 6, 23, 4, 19)) \underline{4}) \\ & \triangleright^* \underline{\text{not}} (\underline{\text{leq}} (\underline{\text{add}} (\underline{\text{proj}}_1^5 (1, 6, 23, 4, 19)) (\underline{\text{const}} \underline{3} (1, 6, 23, 4, 19))) \underline{4}) \\ & \triangleright^* \underline{\text{not}} (\underline{\text{leq}} (\underline{\text{add}} \underline{1} \underline{3}) \underline{4}) \\ & \triangleright^* \underline{\text{not}} (\underline{\text{leq}} \underline{4} \underline{4}) \\ & \triangleright^* \underline{\text{not}} \underline{\text{true}} \\ & \triangleright^* \underline{\text{false}} \end{aligned}$$

As above, we will not prove it, but from now on let us assume that our encoding satisfies the following equivalences:

$$\begin{array}{lll} \llbracket b \rrbracket \sigma = \text{true} & \text{iff} & \langle b \rangle \sigma \triangleright^* \underline{\text{true}} \\ \llbracket b \rrbracket \sigma = \text{false} & \text{iff} & \langle b \rangle \sigma \triangleright^* \underline{\text{false}} \end{array}$$

Statements

The translation of statements is given recursively on the syntax:

$$\begin{aligned} \langle \text{skip} \rangle &= \underline{\text{id}} \\ \langle x_i := e \rangle &= \lambda z. \underline{\text{proj}}_1^k z, \dots, \underline{\text{proj}}_{i-1}^k z, \langle e \rangle z, \underline{\text{proj}}_{i+1}^k z, \dots, \underline{\text{proj}}_k^k z \\ \langle \text{if } b \text{ then } s_1 \text{ else } s_2 \rangle &= \lambda z. \underline{\text{ifz}} (\underline{\text{not}} (\langle b \rangle z)) (\langle s_1 \rangle z) (\langle s_2 \rangle z) \\ \langle \text{while } b \text{ do } s \rangle &= \text{fix } (\lambda f z. \underline{\text{ifz}} (\underline{\text{not}} (\langle b \rangle z)) (f (\langle s \rangle z)) z) \\ \langle s_1 ; s_2 \rangle &= \lambda z. \langle s_2 \rangle (\langle s_1 \rangle z) \end{aligned}$$

9.2 PCF is Turing-Complete

Now we are in a position to show that every While program can be simulated by a PCF program. For this, we must remind ourselves of exactly how a While program computes. We will use the big-step

semantics:

$$\begin{array}{c}
\text{(BSkip)} \frac{}{\langle \text{skip}, \sigma \rangle \Downarrow \sigma} \quad \text{(BAss)} \frac{}{\langle v := a, \sigma \rangle \Downarrow \sigma[v \mapsto \llbracket a \rrbracket(\sigma)]} \\
\\
\text{(BSeq)} \frac{\langle s_1, \sigma \rangle \Downarrow \sigma' \quad \langle s_2, \sigma' \rangle \Downarrow \sigma''}{\langle s_1 ; s_2, \sigma \rangle \Downarrow \sigma''} \\
\\
\text{(BIf}_\top\text{)} \frac{\langle s_1, \sigma \rangle \Downarrow \sigma' \quad \llbracket b \rrbracket(\sigma) = \text{true}}{\langle \text{if } b \text{ then } s_1 \text{ else } s_2, \sigma \rangle \Downarrow \sigma'} \quad \text{(BIf}_\perp\text{)} \frac{\langle s_2, \sigma \rangle \Downarrow \sigma' \quad \llbracket b \rrbracket(\sigma) = \text{false}}{\langle \text{if } b \text{ then } s_1 \text{ else } s_2, \sigma \rangle \Downarrow \sigma'} \\
\\
\text{(BWhile}_\top\text{)} \frac{\langle s, \sigma \rangle \Downarrow \sigma' \quad \langle \text{while } b \text{ do } s, \sigma' \rangle \Downarrow \sigma'' \quad \llbracket b \rrbracket(\sigma) = \text{true}}{\langle \text{while } b \text{ do } s, \sigma \rangle \Downarrow \sigma''} \\
\\
\text{(BWhile}_\perp\text{)} \frac{}{\langle \text{while } b \text{ do } s, \sigma \rangle \Downarrow \sigma} \quad \llbracket b \rrbracket(\sigma) = \text{false}
\end{array}$$

Recall that, whenever a While program terminates, there is a proof tree, built from these rules that describes the computation. So, to prove that every (terminating) while computation can be simulated by the computation of a PCF term, we can use the proof tree associated with the While computation as a kind of template to guide us in constructing the corresponding PCF computation.

Lemma 9.1 (Simulation Lemma). *If $\langle s, \sigma \rangle \Downarrow \sigma'$ then $\llbracket s \rrbracket \sigma \triangleright^* \sigma'$*

Proof. By strong induction on the size of the proof tree witnessing $\langle s, \sigma \rangle \Downarrow \sigma'$. We proceed by analysing the possible cases for the conclusion of the proof tree:

- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (BSkip), it must be that $s = \text{skip}$ and $\sigma = \sigma'$. Since $\llbracket \text{skip} \rrbracket = \text{id}$, we obtain that $\llbracket \text{skip} \rrbracket \sigma \triangleright^* \sigma$ as required.
- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (BAss), it must be that s is of shape $x_i := e$ and σ' is σ but with the i^{th} component of the tuple replaced by $\llbracket e \rrbracket(\sigma)$. By definition:

$$\llbracket x_i := e \rrbracket = \lambda z. (\text{proj}_1^k z, \dots, \text{proj}_{i-1}^k z, \llbracket e \rrbracket z, \text{proj}_{i+1}^k z, \dots, \text{proj}_k^k z)$$

Now $\llbracket e \rrbracket \sigma \triangleright^* \llbracket \llbracket e \rrbracket(\sigma) \rrbracket$, so $(\text{proj}_1^k \sigma, \dots, \text{proj}_{i-1}^k \sigma, \llbracket e \rrbracket \sigma, \text{proj}_{i+1}^k \sigma, \dots, \text{proj}_k^k \sigma) \triangleright^* \sigma'$ and hence:

$$\llbracket x_i := e \rrbracket \sigma \triangleright^* \sigma'$$

- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (BSeq), it must be that s is of shape $s_1 ; s_2$ and there are (strictly smaller) subtrees concluding $\langle s_1, \sigma \rangle \Downarrow \sigma''$ and $\langle s_2, \sigma'' \rangle \Downarrow \sigma'$. Since the subtrees are smaller, we can invoke the induction hypothesis to obtain:

- $\llbracket s_1 \rrbracket \sigma \triangleright^* \sigma''$
- $\llbracket s_2 \rrbracket \sigma'' \triangleright^* \sigma'$

We can use these to construct a longer reduction sequence from $\langle s_1 ; s_2 \rangle \sigma$ to σ' :

$$\langle s_1 ; s_2 \rangle \sigma = (\lambda z. \langle s_2 \rangle (\langle s_1 \rangle z)) \sigma \triangleright \langle s_2 \rangle (\langle s_1 \rangle \sigma) \triangleright^* \langle s_2 \rangle \sigma'' \triangleright^* \sigma'$$

- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (Blf_\top) , it must be that s is of shape if b then s_1 else s_2 and $\llbracket b \rrbracket(\sigma) = \text{true}$. Therefore, we can assume that $\langle b \rangle \sigma \triangleright^* \underline{\text{true}}$ and hence that $\underline{\text{not}} (\langle b \rangle \sigma) \triangleright^* \underline{\text{false}} = 0$. Moreover, there must be a strictly smaller subtree, concluding $\langle s_1, \sigma \rangle \Downarrow \sigma'$. Therefore, we can invoke the induction hypothesis to obtain:

$$- \langle s_1 \rangle \sigma \triangleright^* \sigma'$$

We glue together all of these to construct the reduction sequence:

$$\begin{aligned} \langle \text{if } b \text{ then } s_1 \text{ else } s_2 \rangle \sigma &= \lambda z. \text{ifz } (\underline{\text{not}} (\langle b \rangle z)) (\langle s_1 \rangle z) (\langle s_2 \rangle z) \sigma \\ &\triangleright^* \text{ifz } 0 (\langle s_1 \rangle \sigma) (\langle s_2 \rangle \sigma) \\ &\triangleright^* \langle s_1 \rangle \sigma \\ &\triangleright^* \sigma' \end{aligned}$$

- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (Blf_\perp) , the proof is symmetrical to the previous case.
- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (BWhile_\top) , it must be that s is of shape while b do s' and $\llbracket b \rrbracket(\sigma) = \text{true}$ and there are two strictly smaller subtrees, concluding $\langle s', \sigma \rangle \Downarrow \sigma''$ and $\lambda \text{while } b \text{ do } s', \sigma''. \Downarrow \sigma'$. Therefore, $\underline{\text{not}} (\langle b \rangle \sigma) \triangleright^* 0$ and we can immediately invoke the induction hypothesis to obtain:

$$\begin{aligned} - \langle s' \rangle \sigma &\triangleright^* \sigma'' \\ - \langle \text{while } b \text{ do } s' \rangle \sigma'' &\triangleright^* \sigma' \end{aligned}$$

We glue all of these pieces of reduction sequence together to obtain:

$$\begin{aligned} \langle \text{while } b \text{ do } s' \rangle \sigma &= \text{fix } (\lambda f z. \text{ifz } (\underline{\text{not}} (\langle b \rangle z)) (f (\langle s' \rangle z)) z) \sigma \\ &\triangleright^* (\lambda z. \text{ifz } (\underline{\text{not}} (\langle b \rangle z)) (\langle \text{while } b \text{ do } s' \rangle (\langle s' \rangle z)) z) \sigma \\ &\triangleright^* \text{ifz } 0 (\langle \text{while } b \text{ do } s' \rangle (\langle s' \rangle \sigma)) \sigma \\ &\triangleright^* \langle \text{while } b \text{ do } s' \rangle (\langle s' \rangle \sigma) \\ &\triangleright^* \langle \text{while } b \text{ do } s' \rangle \sigma'' \\ &\triangleright^* \sigma' \end{aligned}$$

- When the proof tree for $\langle s, \sigma \rangle \Downarrow \sigma'$ is concluded by (BWhile_\perp) , it must be that s is of shape while b do s' and $\llbracket b \rrbracket(\sigma) = \text{false}$ and $\sigma' = \sigma$. Hence, $\underline{\text{not}} (\langle b \rangle \sigma) \triangleright^* \underline{1}$. This gives the following sequence:

$$\begin{aligned} \langle \text{while } b \text{ do } s' \rangle \sigma &= \text{fix } (\lambda f z. \text{ifz } (\underline{\text{not}} (\langle b \rangle z)) (f (\langle s' \rangle z)) z) \sigma \\ &\triangleright^* (\lambda z. \text{ifz } (\underline{\text{not}} (\langle b \rangle z)) (\langle \text{while } b \text{ do } s' \rangle (\langle s' \rangle z)) z) \sigma \\ &\triangleright^* \text{ifz } 1 (\langle \text{while } b \text{ do } s' \rangle (\langle s' \rangle \sigma)) \sigma \\ &\triangleright^* \sigma \end{aligned}$$

□

It is easy to make precise what it means for a PCF program to compute a function on natural numbers:

Definition 9.2. We say that a partial function $f : \mathbb{N}^p \rightarrow \mathbb{N}$ is **PCF-definable** just if there is a PCF term M such that, for all $n_1, \dots, n_p \in \mathbb{N}$:

- If f is defined at (n_1, \dots, n_p) , then $M \underline{n_1} \cdots \underline{n_p} \triangleright^* \underline{f(n_1, \dots, n_p)}$.
- Otherwise $M \underline{n_1} \cdots \underline{n_p}$ does not have a normal form.

Now, assuming the Church-Turing thesis, we can show that PCF can define all “intuitively computable” functions.

Theorem 9.1 (PCF is Turing-Complete). *Every computable function on natural numbers is PCF-definable.*

Proof. For lack of time, and because it is all we need in the following, we will only prove the case for total functions. Suppose $f : \mathbb{N}^p \rightarrow \mathbb{N}$ is a computable *total* function on p -tuples of natural numbers. Then, by the Church-Turing thesis, there is a While program s that computes f with respect to variables x_1, \dots, x_p . Let us suppose that this program contains $k \geq p$ variables. To show that f is also PCF-definable, we construct the following term:

$$\lambda x_1 \dots x_p. \underline{\text{proj}}_1^k (\llbracket s \rrbracket (x_1, \dots, x_p, \underline{0}, \dots, \underline{0}))$$

(where $(x_1, \dots, x_p, 0, \dots, 0)$ denotes the tuple with x_i in i^{th} position, from $i = 1$ to p and zeroes thereafter up to position k). We claim that this term PCF-defines f , so let $(n_1, \dots, n_p) \in \mathbb{N}^p$. Since f is total, f is defined at (n_1, \dots, n_p) and so, by definition, $\langle s, [x_i \mapsto n_i \mid 1 \leq i \leq p] \rangle \Rightarrow^* \langle \text{skip}, [x_1 \mapsto f(n_1, \dots, n_p)] \rangle$. It follows from the Simulation Lemma that, therefore $\llbracket s \rrbracket [x_i \mapsto n_i \mid 1 \leq i \leq p] \triangleright^* \llbracket \text{skip} \rrbracket [x_1 \mapsto f(n_1, \dots, n_p)]$. Hence, we can compute:

$$\begin{aligned} & (\lambda x_1 \dots x_p. \underline{\text{proj}}_1^k (\llbracket s \rrbracket (x_1, \dots, x_p, \underline{0}, \dots, \underline{0}))) \underline{n_1} \cdots \underline{n_p} \\ & \triangleright^* \underline{\text{proj}}_1^k (\llbracket s \rrbracket (\underline{n_1}, \dots, \underline{n_p}, \underline{0}, \dots, \underline{0})) \\ & \triangleright^* \underline{\text{proj}}_1^k (\llbracket \text{skip} \rrbracket (\underline{f(n_1, \dots, n_p)}, \underline{0}, \dots, \underline{0})) \\ & \triangleright^* \underline{\text{proj}}_1^k (\underline{f(n_1, \dots, n_p)}, \underline{0}, \dots, \underline{0}) \\ & \triangleright^* \underline{f(n_1, \dots, n_p)} \end{aligned}$$

□

10. Computability: (Un)Decidability

Hopefully, the previous lecture convinced you that PCF can implement all the computable functions. If so, this amounts to saying that you believe in the (generalised) Church-Turing thesis:

Thesis 10.1 (Church-Turing: PCF). *The (intuitively) computable partial functions on natural numbers are exactly the PCF-definable functions.*

Now we will use this to show that certain problems to do with PCF cannot be solved by a computable function, in fact, we will show that almost all properties of PCF programs are undecidable.

10.1 Decidable Properties of Terms

First, I want to talk a little bit about the fact that, throughout what we have done so far, we have assumed that the functions that we are interested in computing are functions on natural numbers. This is completely standard, and it makes sense if you think about it. Your experience of digital computers should suggest to you that all kinds of data — strings, arrays, trees, graphs etc — can all ultimately be represented using sequences of bits. Since we can think of a sequence of bits as a binary numeral, we can thus encode any kind of data as some natural number. Operations on data can then be implemented as operations on natural numbers, and to all intents and purposes basic arithmetic is all that is needed.

Consequently, formal systems (comprising syntax and rules) are especially powerful if they allow for the representation of numbers and the implementation of arithmetic. One way to see this precisely is to look at the consequences of encoding parts of the system within itself. This can typically be done because the a formal system is just a collection of strings and strings can be encoded as numbers.

This was the approach of Gödel for obtaining his famous incompleteness results about the theory of arithmetic in the early 1930s. In the theory of arithmetic, the central concept is provability. By encoding the formulas and proof system of arithmetic within the theory of arithmetic he was able to derive the conclusion that there are some formulas which, although true, are not provable. In the PCF the concept we are interested in is computability (stated precisely as PCF-definability), and the same technique will allow us to conclude that there are a large number of functions that are not computable.

The technique of encoding syntax by numbers takes Gödel's name. Here we are interested in encoding PCF-terms.

Definition 10.1. *A Gödel numbering is a pair of computable (partial) functions:*

- $\# : \Lambda \rightarrow \mathbb{N}$

- $\#^{-1} : \mathbb{N} \rightarrow \Lambda$

with the property that:

$$\#^{-1} (\# M) = M$$

Sometimes it is required that a Gödel encoding be a bijection, but we won't need that. There are actually lots of left-invertible injections $\#$ that we could define in order to satisfy the required properties and, for the rest of our development, it doesn't matter which one we choose. However, it is worth seeing one of them to make sure we all believe that it is possible. Let's consider an adaptation of Gödel's original.

We start by making an arbitrary assignment code of natural numbers to the letters that make up the strings in the language Λ . For example maybe we decide to assign:

code('λ') = 1	code('S') = 5
code('(') = 2	code('pred') = 6
code(')') = 3	code('fix') = 7
code('.') = 4	code('ifz') = 8
	code('Z') = 9

and then enumerate the countably many variables as $\text{code}('x_0') = 10$, $\text{code}('x_1') = 11$, $\text{code}('x_2') = 12$ and so on. The key is then to map a whole string to a single number in such a way that we ensure that distinct strings get mapped to distinct numbers. The way Gödel did this was to use products of primes, with each prime raised to the power corresponding to the code of the letter:

$$\#("(\text{fix } (\lambda x_1. x_1))") = 2^2 \cdot 3^7 \cdot 5^2 \cdot 7^1 \cdot 11^{10} \cdot 13^4 \cdot 17^{10} \cdot 19^3 \cdot 23^3$$

It then follows from the fundamental theorem of arithmetic that different strings will give rise to different Gödel numbers.

Since we can represent numbers in PCF by our PCF numerals, let us investigate what we can say about PCF-terms computing with codings of PCF-terms. Our main aim is to be able to deduce some undecidability results about problems to do with PCF. For example, we will show that the problem of determining whether a PCF-term M normalises (has a normal form) is undecidable.

To that end, we first need to make precise what we mean by property of terms and what it means for such a property to be decidable. In mathematics, a property of a class of objects is often just represented as a the set of those objects which have the property. For example, the property of being even is just the set of even numbers. Then we can say that an object satisfies a property just if it is a member of the set. We will say that a property is decidable just when membership in the set is decidable.

Definition 10.2 (Decidability). *Let $\Phi \subseteq \Lambda$ be a property of PCF-terms. We say that Φ is **decidable** just if the characteristic function of Φ is computable, i.e. there is an algorithm for computing $\chi : \Lambda \rightarrow \mathbb{N}$ satisfying, for all terms $M \in \Lambda$:*

$$\chi(M) = \begin{cases} 1 & \text{if } M \in \Phi \\ 0 & \text{if } M \notin \Phi \end{cases}$$

To deduce that a particular problem χ is decidable, all we need is to convince ourselves that there is an algorithm which, for any term M , tells us faithfully if M satisfies the property (1) or not (0). For example:

Lemma 10.1. *The problem $\{M \mid M \text{ is an abstraction}\}$ is decidable.*

Proof. Given M , we can compute whether or not M is an abstraction by parsing M . This will correctly determine whether or not M is an abstraction for any term M . \square

In practice, a short sketch of an algorithm such as the one above always suffices as a proof.

10.2 Undecidability and the Second Recursion Theorem

To show that a problem is undecidable is, conceptually, more difficult. However, there is a very powerful theorem, independently discovered by Scott and Curry for the pure untyped λ -calculus, that allows us to conclude undecidability relatively straightforwardly, in many cases. We start by considering some computable functions on the codes of terms.

Lemma 10.2. *There exist terms app and gnum such that:*

$$\underline{\text{app}} \ \underline{\#M} \ \underline{\#N} \approx \underline{\#(MN)} \quad \text{and} \quad \underline{\text{gnum}} \ \underline{n} \approx \underline{\#n}$$

Proof. Given that $\#$ and $\#^{-1}$ are intuitively computable:

- The function $\text{app} : \mathbb{N}^2 \rightarrow \mathbb{N}$ satisfying $\text{app}(\#M, \#N) = \#(MN)$ is, intuitively, a computable function (I could come up with an algorithm to implement it).
- Likewise, the function $\text{gnum} : \mathbb{N} \rightarrow \mathbb{N}$ satisfying $\text{gnum}(n) = \#(\underline{n})$ is, intuitively, a computable function.

Hence, it follows from the Church-Turing thesis that there are (probably quite complicated) terms app and gnum to define them in PCF. \square

It is worth considering the two expressions M and $\#M$. *Both of these expressions are terms.* A useful way to think about the situation is that M is a program, $\#M$ is its source code and $\underline{\#M}$ is its source code represented as a datatype that we can compute with in PCF (i.e. a PCF numeral). Therefore, we can understand the combinator app as being the program that takes the source code of M and the source code of N and returns the source code for MN . If we specialise the equation for gnum so that n is, in fact, $\#M$, the code for some program, then we can understand gnum as being the program that takes the source code for some program M (represented as a PCF numeral) and returns the source code for that numeral.

The following important result is due to Kleene (for pure λ -calculus).

Theorem 10.1 (Second Recursion Theorem). *For all terms F , there exists a term M such that:*

$$F \ \underline{\#M} \approx M$$

Proof. Let F be an arbitrary term. Then define N as $\lambda x. F (\text{app } x (\text{gnum } x))$. Finally, take $N \# N$ to be the witness M : it is easily seen that this term satisfies $F \# M \approx M$. \square

First, pay special attention to the quantification here: it says that *for all* terms *there exists* an input which behaves in a certain way. The theorem is valuable because it says that, if you attempt to implement a program F to compute with the codes of terms, e.g. to answer 0 or 1 whether or not a given term satisfies some property, then no matter how you implement it (the theorem is *for all* F), there will always be some input $\#M$ about which you have no control over the corresponding output — it is always behaves like M , whether that is what you want or not!

10.3 The Scott-Curry Theorem

If you want to show that some set Φ is decidable then you must exhibit an F computing the characteristic function of Φ , that is it must be equal to $\underline{0}$ or $\underline{1}$ depending on whether the input is the code of a term M in Φ or not. This is quite difficult when the above theorem tells us there will always be some input $\#M$ over which you no control about the corresponding output.

In fact, for any interesting property, we can construct a tricky term that precisely demonstrates the difficulty. When we say “interesting” what we mean is that the property should be a *non-trivial* and *behavioural*.

Trivial properties are properties which are satisfied by no terms or by all terms. These kinds of properties are too easy to decide, because the characteristic function is constantly 0 or constantly 1 and these are obviously computable. Non-behavioural properties concern only the syntax of terms and not *what they compute*. For example, the property $\{M \in \Lambda \mid M \text{ is an abstraction}\}$ is not behavioural because there are terms like id const which do not satisfy the property and yet are convertible with a term const that does. Non-behavioural properties are often easy to compute because we just need to look at the syntax (but some non-behavioural properties are nevertheless undecidable).

However, if a property is both *non-trivial* and *behavioural* then it is undecidable.

Theorem 10.2 (Scott-Curry). *Let $\Phi \subseteq \Lambda$ satisfy the following properties:*

- Φ is *non-trivial*: $\emptyset \neq \Phi \neq \Lambda$
- Φ is *behavioural*: if $M \in \Phi$ and $M \approx N$, then $N \in \Phi$

Then it follows that Φ is undecidable.

Proof. Since Φ is non-trivial, there is some $P \in \Phi$ and some $Q \notin \Phi$. Suppose, for the purposes of obtaining a contradiction, that Φ is decidable. Then it follows that there is some term F that computes its characteristic function, i.e. satisfying:

$$F \# M \approx \begin{cases} \underline{1} & \text{if } M \in \Phi \\ \underline{0} & \text{if } M \notin \Phi \end{cases}$$

Define $G := \lambda x. \text{ifz } (F \ x) \ P \ Q$. It follows from the Second Recursion Theorem that there is some term N such that $G \# N = N$. One of $N \in \Phi$ or $N \notin \Phi$ is certainly true, so let us examine the two cases:

- If $N \in \Phi$, then by the properties of N, F and ifz :

$$N \approx G \# N \approx \text{ifz } (F \# N) P Q \approx Q$$

but we assumed that $N \in \Phi$ and Φ is closed under \approx , so we are forced to conclude that $Q \in \Phi$ which is a contradiction.

- If $N \notin \Phi$, then by the properties of N, F and ifz :

$$N \approx G \# N \approx \text{ifz } (F \# N) P Q \approx P$$

but we assumed that $N \notin \Phi$ and Φ is closed under \approx , so we are forced to conclude that $P \notin \Phi$ which is a contradiction.

□

The theorem is a direct analogue of Rice's Theorem for properties of Turing machines (you may have seen a version for While programs).

So which properties of PCF-terms are decidable? Essentially no interesting ones.

Corollary 10.1. *The problem $\Phi_{\text{norm}} := \{M \in \Lambda \mid M \text{ is normalising}\}$ is undecidable.*

Proof. The set Φ_{norm} is non-trivial because $\underline{\text{id}}$ is a member but $\underline{\text{div}}$ is not. Moreover, it is behavioural. To see this let M be normalising and $M \approx N$. Since M is normalising, there is a normal form P with $M \approx P$. Therefore, also $N \approx P$ and it follows that, therefore $N \triangleright^* P$. Consequently, N is normalising. Hence, undecidability of Φ_{norm} follows from the Scott-Curry theorem. □

11. Types: Judgements and the Type System

So far we have been considering the untyped syntax of PCF, in which it is perfectly possible to create strange looking terms like $\underline{1}$ `id` and compute with them, although the computation does not end well! In practice, we would like that such combinations of terms simply cannot arise, that they are forbidden, and this is where the type system comes in.

11.1 Intuitions

Typically, we are only interested in applications of a function to data that is within its domain and which yield data from its codomain — applications that are *well typed*. *Type theory* concerns the formalisation of this idea of type and rules for deducing well-typedness. In the type theory that we study in this second half of the unit, this is made precise by the *judgement*:

$$\Gamma \vdash M : A$$

which we can read as “Under assumptions Γ , M has type A ”. There are two particularly fruitful views of type theory, which help to explain the meaning of the words “has type” intuitively.

Types describe abstract properties

In this view, “type” and “property” are taken to be synonymous. Type theories are a syntax for reasoning about the properties of terms. For example, to say $M : A \rightarrow B$ is to assert that M has a certain property, namely: given any input satisfying the property A , its output will satisfy the property B . An example is the following:

$$\lambda x. x + 1 : \{n \in \mathbb{Z} \mid n \text{ is even}\} \rightarrow \{n \in \mathbb{Z} \mid n \text{ is odd}\}$$

The successor function has the property that, given an even integer, it will always return an odd integer. In this view it is typical that a term may enjoy the association of many different types. For example, the same term also satisfies the property $\{n \in \mathbb{Z} \mid n \geq 0\} \rightarrow \{n \in \mathbb{Z} \mid n \neq 0\}$.

This view finds important applications in the *analysis of programs*. The second property is exactly what you need to know about $\lambda x. x + 1$ in order to deduce that the following program does not raise an exception:

$$100 \div (\text{if } n \geq 0 \text{ then } (\lambda x. x + 1) \, n \text{ else } n)$$

Types prescribe abstract classifications

In this view, “type” and “classification” are taken as synonymous. Type theories are a syntax for enforcing the classification of terms. To say $M : A \rightarrow B$ is to assert that M belongs to a certain class, namely: the class of functions with domain A and codomain B . An example is the following:

$$\lambda x. x ++ [1] : [\text{int}] \rightarrow [\text{int}]$$

The append-one function is a function from list of integers to lists of integers. In this view it is typical that we think of the type as being an intrinsic part of the description of the function. For example, we consider $\lambda x. x ++ [1] : [\text{float}] \rightarrow [\text{float}]$ to be a *different function*, it’s just a coincidence that you happen to follow the same rule in order to compute its output.

Of course, these two views are far from mutually exclusive! On the one hand, it may only be sensible to infer properties of terms that are already classified, for example, the properties above only make sense for functions of type $\mathbb{Z} \rightarrow \mathbb{Z}$. On the other, a classification is a kind of very coarse property, for if M is classified as a function of type $[\text{int}] \rightarrow [\text{int}]$ then it certainly has the property that it maps inputs with the property of being a list of integers to outputs with the same property.

Just as there are an amazing variety of programming languages out there, there is an even greater variety of type systems. We will study the “simple type system” that is usually associated with PCF (here: *simple* has a technical connotation). Simple types go back all the way to the logicians of the early 20th Century, but they form the core part of any of the modern type systems that you will find attached to programming languages today.

11.2 Types and Environments

We start by saying exactly what the types of this system actually look like.

Definition 11.1. We assume a countable set of **type variables** \mathbb{A} , ranged over by a, b, c and other lowercase letters from the start of the alphabet. The **types**, written \mathbb{T} , are a set of strings A, B etc defined by the following grammar:

$$(\text{Types}) \quad A, B ::= \text{Nat} \mid a \mid (A \rightarrow B)$$

We immediately adopt the following conventions to make our life less tedious:

- We omit the outermost parenthesis when writing types.
- We assume that the arrow associates to the right: i.e. when we write $A_1 \rightarrow A_2 \rightarrow A_3$, the type that we mean is $(A_1 \rightarrow (A_2 \rightarrow A_3))$.

Here are some examples of types written with and without all their parentheses.

$$\begin{array}{ll} \text{Nat} \rightarrow \text{Nat} & (\text{Nat} \rightarrow \text{Nat}) \\ a \rightarrow b \rightarrow b & (a \rightarrow (b \rightarrow b)) \\ (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} & ((\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat}) \end{array}$$

We will consistently use A, B, C and other upper case letters at the start of the alphabet to refer to types generically.

Types that are of the form $A_1 \rightarrow A_2$ are variously called *arrow types* or *function types* and the idea is that if M has type $A_1 \rightarrow A_2$ then it represents a function that maps terms of type A_1 to terms of type A_2 . You can think of type variables as placeholders for any type. So, for example, from the fact that we will be able to assign the type $a \rightarrow a$ to the identity function $\text{id} : a \rightarrow a$ we can infer that we could assign *any* type of shape $A \rightarrow A$ to this function, for example, $\text{Nat} \rightarrow \text{Nat}$ or $(\text{Nat} \rightarrow \text{Nat}) \rightarrow (\text{Nat} \rightarrow \text{Nat})$.

The whole point of types is to classify terms, so we next need to say how that comes about. Let's first agree on how to write it.

Definition 11.2. A *type assignment* is a pair of a term M and a type A , written $M : A$. The term part of a type assignment is called the **subject** and the type part the **predicate**.

Clearly, not all type assignments $M : A$ should be allowed. For example, we would not expect $\underline{1}$ to have type $\text{Nat} \rightarrow \text{Nat}$. On the other hand, we would expect that $\lambda x. S\ x : \text{Nat} \rightarrow \text{Nat}$ or even $\lambda x. x : a \rightarrow a$.

The purpose of a *type system* is to provide a collection of rules in order to determine when $M : A$ is valid. In general M may contain free variables and then the type of M will usually depend on the type of the free variables. For example, term $\lambda x. y$ should be allowed to have a type of the form $A \rightarrow B$ where B is the type of y . So, this begs the question, what should be the type of a free variable?

We allow for *hypothetical* deductions. That is, we define when $M : A$ is valid only *with respect to* some given *assumptions* (also known as *hypotheses*) concerning the types of the free variables. These hypotheses are just type assignments of the form $x : B$, meaning that free variable x is assumed to have type B , which are collected together in a *type environment*, usually written Γ . We require that deductions can be made only under consistent assumptions.

Definition 11.3. A *type environment*, generically written Γ , is a finite set of type assignments of the form $x : A$ which is, moreover, **consistent**, in the sense that if $x : A \in \Gamma$ and $x : B \in \Gamma$, then $A = B$.

The **subjects** of Γ is the set, written $\text{dom } \Gamma$, of those term variables x for which there is some A such that $x : A \in \Gamma$.

A concrete example of a type environment is $\{x : a \rightarrow b, y : \text{Nat} \rightarrow \text{Nat}, z : \text{Nat}\}$ (whenever we write a type environment explicitly like this, we will assume that the *subjects* x , y and z are distinct variables).

11.3 Type System

The notion that a term may be assigned a type under some assumptions is made precise through the notion of *typing judgement*. A *type system* gives rules by which typing judgements can be deduced. We start by saying what should be the types of the built-in constants of the language.

Definition 11.4. Let \mathbb{C} be the following collection of type assignments:

$$\begin{aligned} &\{Z : \text{Nat}\} \\ &\cup \{S : \text{Nat} \rightarrow \text{Nat}\} \\ &\cup \{\text{pred} : \text{Nat} \rightarrow \text{Nat}\} \\ &\cup \{\text{ifz} : \text{Nat} \rightarrow A \rightarrow A \rightarrow A \mid A \in \mathbb{T}\} \\ &\cup \{\text{fix} : (A \rightarrow A) \rightarrow A \mid A \in \mathbb{T}\} \end{aligned}$$

You can see from this definition that we expect terms to be assigned several types, in general. For example, the ifz constant is assigned every type of the form $\text{Nat} \rightarrow A \rightarrow A \rightarrow A$ for some type A . This should ring true with your intuitions, because we expect the ifz constant to take three arguments: a number, the ‘then branch’ and the ‘else branch’ and it should return one of those branches. The branches can be numbers of type Nat : $\text{ifz } x \ 1 \ 2$; or they could be functions of type $\text{Nat} \rightarrow \text{Nat}$: $\text{ifz } x \ \text{id } S$; or of any other type A , as long as the type of the ‘then branch’ and the ‘else branch’ match. Recall that the idea of the fixed point combinator fix is that it takes a function as input and returns a fixed point of that function. This only makes sense for functions that have the same domain as their codomain because, if N is a fixed point of M then $M \ N \approx N$ – it returns something of the same type as that taken as input. Hence, fix should take a function of any type $A \rightarrow A$ (i.e. where the domain and codomain match) and return a fixed point, which is of course something of this type A .

Definition 11.5 (Type System). A **type judgement** is a triple consisting of a type environment Γ , a term M and a type A , which we write:

$$\Gamma \vdash M : A$$

Then the **type system** is the following collection of rules by which one can justify such type judgements using proof trees:

$$\begin{aligned} &x:A \in \Gamma \quad \frac{}{\Gamma \vdash x : A} \text{ (TVar)} \quad c:A \in \mathbb{C} \quad \frac{}{\Gamma \vdash c : A} \text{ (TCst)} \\ &\frac{\Gamma \vdash M : B \rightarrow A \quad \Gamma \vdash N : B}{\Gamma \vdash MN : A} \text{ (TApp)} \quad x \notin \text{dom } \Gamma \quad \frac{\Gamma \cup \{x : B\} \vdash M : A}{\Gamma \vdash \lambda x. M : B \rightarrow A} \text{ (TAbs)} \end{aligned}$$

In type theory, a proof tree justifying a type judgement is usually called a **type derivation**.

Since we are never in any danger of confusion, whenever we write out a type environment explicitly as part of a judgement, we will omit the braces part of the set notation. So, for example, we will write $x:\text{Nat} \rightarrow \text{Nat}, y : \text{Nat} \vdash xy : \text{Nat}$ instead of the more cumbersome $\{x:\text{Nat} \rightarrow \text{Nat}, y:\text{Nat}\} \vdash xy : \text{Nat}$. For the same reason, if we make no assumptions then we will simply write $\vdash M : A$ rather than $\emptyset \vdash M : A$.

The way to read a typing judgement $\Gamma \vdash M : A$ is “under assumptions Γ , M has type A ” and these rules allow us to conclude sound judgements – i.e. where that statement is true. Each rule should be read as:

“If any statement(s) above the line are true and any condition to the left is satisfied, then the statement below the line is also true”

The rule (TVar) allows us to make use of our assumptions, it says that, at any time (i.e. without premise), if we assumed $x:A$ then we can conclude $x : A$. The rule (TCst) plays a similar role for constants. It says that if $c:A$ is a primitive typing that we assigned to the constants in \mathbb{C} , then we can conclude exactly $c : A$. A difference between the assumptions on variables and the built-in assignment of types to the constants is that the former will typically change during a derivation, whereas the latter never changes.

The rule (TApp) allows us to make use of function types. It says that if we know M has a function type $B \rightarrow A$ and N has a matching type B then it follows that the application MN can be assigned the type according to what M is known to return, which is A .

The rule (TAbs) says that, in order to conclude that $\lambda x. M$ has type $B \rightarrow A$ from some assumptions Γ , I just have to show that I can conclude that M has type A (recall that M is what the function returns, and A is the return type of $B \rightarrow A$) under the additional assumption that x has type B .

The following is a proof tree showing that the function that adds two to its argument is of type $\text{Nat} \rightarrow \text{Nat}$. That is, that the following judgement is true: $\vdash \lambda x. S (S x) : \text{Nat} \rightarrow \text{Nat}$. We will take a look at proof trees in more detail in the next lecture.

$$\begin{array}{c}
 \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{ (TCst)} \quad \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{ (TCst)} \quad \frac{}{x:\text{Nat} \vdash x : \text{Nat}} \text{ (TVar)} \\
 \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{ (TCst)} \quad \frac{}{x:\text{Nat} \vdash S x : \text{Nat}} \text{ (TApp)} \quad \frac{}{x:\text{Nat} \vdash S (S x) : \text{Nat}} \text{ (TApp)} \\
 \frac{}{x:\text{Nat} \vdash S (S x) : \text{Nat}} \text{ (TApp)} \quad \frac{}{\vdash \lambda x. S (S x) : \text{Nat} \rightarrow \text{Nat}} \text{ (TAbs)}
 \end{array}$$

12. Types: Derivations and Typability

In this lecture we look at more detail at type derivations/proof trees, some nuances of the system and some of the key problems associated with typing.

12.1 Proof Trees

Let me remind you of the rules of the type system:

$$\begin{array}{c} x:A \in \Gamma \quad \frac{}{\Gamma \vdash x : A} \text{ (TVar)} \quad c:A \in \mathbb{C} \quad \frac{}{\Gamma \vdash c : A} \text{ (TCst)} \\[10pt] \frac{\Gamma \vdash M : B \rightarrow A \quad \Gamma \vdash N : B}{\Gamma \vdash MN : A} \text{ (TApp)} \quad x \notin \text{dom } \Gamma \quad \frac{\Gamma \cup \{x : B\} \vdash M : A}{\Gamma \vdash \lambda x. M : B \rightarrow A} \text{ (TAbs)} \end{array}$$

We can use the rules in order to construct *proof trees* (also called type derivations in this setting). A *proof tree* for this system is tree whose nodes are labelled by judgements, in such a way that, if a node is labelled $\Gamma \vdash M : A$ and its children are labelled $\Gamma_1 \vdash M_1 : A_1 \dots \Gamma_k \vdash M_k : A_k$, then it must be that:

$$\frac{\Gamma_1 \vdash M_1 : A_1 \quad \dots \quad \Gamma_k \vdash M_k : A_k}{\Gamma \vdash M : A}$$

is an instance of one of the given rules: (TVar), (TCst), (TAbs) or (TApp). An *instance* of a rule is just a copy of the rule but with the generic M, N, x, Γ and so on replaced *consistently* by particular objects like $(\lambda z. z), \{x:\text{Nat} \rightarrow \text{Nat}\}$ or y .

For example:

$$\frac{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat} \quad x:\text{Nat} \vdash x : \text{Nat}}{x:\text{Nat} \vdash S x : \text{Nat}} \text{ (TApp)}$$

is an instance of the (TApp) rule, with $M := S, N := x, \Gamma := \{x:\text{Nat}\}, A := \text{Nat}, B := \text{Nat}$. Similarly,

$$\frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{ (TCst)}$$

is an instance of the (TCst) rule, with $\Gamma := \{x:\text{Nat}\}, c = S$, and $A := \text{Nat} \rightarrow \text{Nat}$.

When we put rule instances together, in such a way that every branch of the tree is closed off by one of the axioms (those rules without premises) (TVar) or (TCst), then we obtain a proof tree.

The judgement at the root of the tree is then shown to be true. For example, the example from the previous lecture:

$$\begin{array}{c}
 \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{(TCst)} \quad \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{(TCst)} \quad \frac{}{x:\text{Nat} \vdash x : \text{Nat}} \text{(TVar)} \\
 \frac{}{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat}} \text{(TCst)} \quad \frac{x:\text{Nat} \vdash S : \text{Nat} \rightarrow \text{Nat} \quad x:\text{Nat} \vdash S x : \text{Nat}}{x:\text{Nat} \vdash S (S x) : \text{Nat}} \text{(TApp)} \\
 \frac{x:\text{Nat} \vdash S (S x) : \text{Nat}}{\vdash \lambda x. S (S x) : \text{Nat} \rightarrow \text{Nat}} \text{(TAbs)}
 \end{array}$$

When you want to construct a type derivation for a judgement, it is usually best to start from the bottom and build upwards. This type system has a particularly nice property that makes this approach work very well. In this system, there is exactly one rule for each possible shape of term. So given $\Gamma \vdash M : A$ there is only one rule that can be used to obtain this judgement in the conclusion, and it is completely determined by the shape of M . Type systems with this property are called *syntax-directed* because the syntax of the term directs you in building a derivation.

However, that is not to say that every type derivation is completely mechanical. The rule (TApp) requires some creativity. Imagine building a derivation that has $\Gamma \vdash MN : A$ in the conclusion. The shape MN determines that we must use the (TApp) rule, but to use this rule we have to somehow invent the type B . The type B is not necessarily mentioned in the conclusion (below the line) so it requires a little bit of thought (but usually not too much) to come up with a correct instantiation. Consider building the following derivation from the ground up, with $\Gamma = \{x : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat}\}$ (I have omitted the rule names so that it fits on the page, make sure you can reconstruct them):

$$\begin{array}{c}
 \frac{\Gamma, y:\text{Nat} \rightarrow \text{Nat} \vdash x : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat}}{\Gamma, y:\text{Nat} \rightarrow \text{Nat} \vdash xy : \text{Nat}} \quad \frac{\Gamma, y:\text{Nat} \rightarrow \text{Nat} \vdash y : \text{Nat} \rightarrow \text{Nat}}{\Gamma \vdash \lambda y. xy : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat}} \quad \frac{}{\Gamma, z:\text{Nat} \vdash z : \text{Nat}} \\
 \frac{\Gamma \vdash \lambda y. xy : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} \quad \Gamma \vdash \lambda z. z : \text{Nat} \rightarrow \text{Nat}}{\Gamma \vdash (\lambda y. xy)(\lambda z. z) : \text{Nat}}
 \end{array}$$

The tree is concluded using an instance of the rule (TApp) (this is the rule used at the root), with $B := \text{Nat} \rightarrow \text{Nat}$, but there is not much of a clue to why you should pick $B := \text{Nat} \rightarrow \text{Nat}$ if you only have the root judgement to go off.

12.2 Types involving type variables

So far we have been deriving judgements in which the types involved were all combinations of the base type Nat and the arrow (or function type). However, we can also prove judgements involving types with type variables. For example, let me abbreviate the environment $\{x:a \rightarrow b \rightarrow c, y:a \rightarrow b, z:a\}$ as Γ and omit the rule names in the following:

$$\begin{array}{c}
 \frac{\Gamma \vdash x : a \rightarrow b \rightarrow c \quad \Gamma \vdash z : a}{\Gamma \vdash xz : b \rightarrow c} \quad \frac{\Gamma \vdash y : a \rightarrow b \quad \Gamma \vdash z : a}{\Gamma \vdash yz : b} \\
 \frac{\Gamma \vdash xz : b \rightarrow c \quad \Gamma \vdash yz : b}{\Gamma \vdash xz(yz) : c} \\
 \frac{x : a \rightarrow b \rightarrow c, y : a \rightarrow b \vdash \lambda z. xz(yz) : a \rightarrow c}{x : a \rightarrow b \rightarrow c \vdash \lambda yz. xz(yz) : (a \rightarrow b) \rightarrow a \rightarrow c} \\
 \vdash \lambda xyz. xz(yz) : (a \rightarrow b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c
 \end{array}$$

What does it mean for a term like $\lambda xyz. x \ z \ (y \ z)$ to have a type of shape $(a \rightarrow b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c$? You can think of type variables as placeholders: it is possible to show that, if a provable judgement $\Gamma \vdash M : A$ involves type variables, then any consistent replacement of the type variables will also yield a provable judgement. As a simpler example, consider the proof:

$$\frac{\frac{}{x:a, y:b \vdash x:a} \text{ (TVar)}}{\frac{}{x:a \vdash \lambda y. x : b \rightarrow a} \text{ (TAbs)}} \text{ (TAbs)}$$

Any replacement of a and b by any other types will also yield a valid proof tree. E.g. $a := (\text{Nat} \rightarrow \text{Nat})$, $b := \text{Nat}$:

$$\frac{\frac{\frac{}{x:\text{Nat} \rightarrow \text{Nat}, y:\text{Nat} \vdash x : \text{Nat} \rightarrow \text{Nat}} \text{(TVar)}}{x:\text{Nat} \rightarrow \text{Nat} \vdash \lambda y. x : \text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat}} \text{(TAbs)}}{\vdash \lambda x y. x : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat}} \text{(TAbs)}$$

In fact, there is a sense in which the type $a \rightarrow b \rightarrow a$ tells you everything about the possible typings for the term $\lambda x y. x$.

No Polymorphism

However, this is an external property. The fact that I can assign a type like $a \rightarrow b \rightarrow a$ to $\lambda xy. x$ tells me, the human, that I will be able to construct a derivation of $\vdash \lambda xy. x : A \rightarrow B \rightarrow A$ for any particular types A and B . However, this is not a property that I can exploit within the system. This is because, unlike Haskell's type system, the simple type system we look at in this unit does not have *polymorphism*: that is, the ability to *use* a type involving type variables at different instances *within* a derivation. For example, consider the following derivation:

$$\frac{\frac{\frac{\vdots}{x:a \rightarrow a \vdash \text{ifz}(x \underline{0})(x S) : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} \rightarrow \text{Nat}}}{x:a \rightarrow a \vdash \text{ifz}(x \underline{0})(x S) \text{ pred} : \text{Nat} \rightarrow \text{Nat}}}{x:a \rightarrow a \vdash \text{ifz}(x \underline{0})(x S) : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} \rightarrow \text{Nat}} \quad \frac{x:a \rightarrow a \vdash S : \text{Nat} \rightarrow \text{Nat}}{x:a \rightarrow a \vdash \text{pred} : \text{Nat} \rightarrow \text{Nat}}$$

If you consider the invalid step $x:a \rightarrow a \not\vdash x : (\text{Nat} \rightarrow \text{Nat}) \rightarrow \text{Nat} \rightarrow \text{Nat}$, we know that x has type $a \rightarrow a$, so, in principle, whatever term x represents, we could type it by the replacement $a := \text{Nat} \rightarrow \text{Nat}$. However, there is no rule within the system that allows us to do this in the middle of the type derivation. When we have $x:a \rightarrow a$ as an assumption, the *only* thing we can do with it is to conclude $x : a \rightarrow a$.

This is in contrast to Haskell, which has a type system with a limited form of polymorphism. In Haskell, for example, supposing `id` is the identity function, which has Haskell type $a \rightarrow a$, and `succ` is the successor function and `pred` the predecessor with types $\text{Int} \rightarrow \text{Int}$, the expression:

if (id 0) then (id succ) else pred

is, in fact typable, because the Haskell polymorphic type system allows $\text{id} : a \rightarrow a$ to be used at any instance (here it is used at $\text{Int} \rightarrow \text{Int}$ on the left and $(\text{Int} \rightarrow \text{Int}) \rightarrow \text{Int} \rightarrow \text{Int}$ on the right).

12.3 Typability and other Problems

When we introduced the system, we said that not all type assignments $M : A$ should be allowed (to be derived).

Definition 12.1. We say that a closed term M is **typable** just if there is some type A such that $\vdash M : A$ is derivable in the type system.

It's possible to talk about the typability of open terms too, by extending this definition, but a little less neat, so we will stick to closed terms. To show that a closed term M is typable requires us to exhibit a type and to show that the corresponding judgement is derivable.

Not all closed terms are typable. For example, the term $\lambda x. xx$ is not typable. To prove it, we will rely on analysing the possible shape of *any* typing derivation that has $\lambda x. xx$ in the conclusion.

Lemma 12.1. The term $\lambda x. xx$ is untypable.

Proof. We suppose $\lambda x. xx$ is typable and try to obtain a contradiction. Since $\lambda x. xx$ is typable, there is a type A such that $\vdash \lambda x. xx : A$ is derivable. Since this judgement is derivable, and its subject is an abstraction, the derivation must conclude using (TAbs) at the root:

$$\frac{\vdots \quad x : B \vdash xx : C}{\vdash \lambda x. xx : A} \text{ (TAbs)}$$

but, for this to be a correct derivation, it must be that A has shape $B \rightarrow C$ for some types B and C . Now consider the next part of this derivation, rooted at $x : B \vdash xx : C$. This can only be justified using the (TApp) rule:

$$\frac{\frac{\vdots \quad x : B \vdash x : D \rightarrow C \quad \vdots \quad x : B \vdash x : D}{x : B \vdash xx : C} \text{ (TApp)}}{\vdash \lambda x. xx : B \rightarrow C} \text{ (TAbs)}$$

which means that there must be some type D such that the operator (on the left of the application) has type $D \rightarrow C$ and the operand (on the right of the application) has type D . In this case, both operator and operand are x . Both of the remaining judgements have x as the subject, and so must be justified using instances of (TVar):

$$\frac{\frac{\frac{}{x : B \vdash x : D \rightarrow C} \text{ (TVar)} \quad \frac{\frac{}{x : B \vdash x : D} \text{ (TVar)}}{x : B \vdash xx : C} \text{ (TApp)}}{\vdash \lambda x. xx : B \rightarrow C} \text{ (TAbs)}$$

For the left-hand instance of (TVar) to be a correct, it can only be that B and $D \rightarrow C$ are the same type. Similarly, from the right-hand instance of (TVar), it must be that B and D are the same type. Therefore, if this was a valid derivation, we would have that type $D = B = D \rightarrow C$. This is impossible! A type is just a string, and there is no finite string D that is equal to a string $D \rightarrow C$ that properly contains it. \square

These insights into what shape the various types have based on the form of a given judgement, which we used throughout the foregoing proof, are neatly summarised in the following *Inversion* theorem.

Theorem 12.1 (Inversion). *Suppose $\Gamma \vdash M : A$ (is derivable), then:*

- *If M is a variable x , then it must be that $x:A \in \Gamma$.*
- *If M is a constant c , then it must be that $c:A \in \mathbb{C}$.*
- *If M is an application PQ , then there is a type B such that $\Gamma \vdash P : B \rightarrow A$ and $\Gamma \vdash Q : B$.*
- *If M is an abstraction $\lambda x.P$, then there are types B and C such that $A = B \rightarrow C$, and $\Gamma, x : B \vdash P : C$.*

Proof. By inspection of the rules. □

I want to conclude this section by describing three of the main computational problems that are considered in connection to type systems. For simplicity, we give their definitions for the case of closed terms. The problems all extend to open terms, but the definitions are slightly messier.

- **Typability.** Given a closed term M , is M typable?
- **Type checking.** Given a closed term M and a type A , is $\vdash M : A$ derivable?
- **Type inference.** Given a closed term M , compute all types A such that $\vdash M : A$.

The three problems are all related, and it is likely that if one is decidable/computable for a given type system, then all of them are. For the simple type system that we study, all three are decidable/-computable.

13. Types: Not Going Wrong

Type systems constrain the programs that we can write. If we only allow ourselves well-typed programs, then we are excluding potentially useful terms, like $\lambda x. xx$. Many beginning programmers will decry typed programming for exactly this reason “why can’t I write the program the way I want to?”, “why does the type checker get in my way?”.

The reason we are happy to sacrifice expressivity - to forsake certain programs that would not be well-typed - is because terms that are well-typed behave well. By contrast, consider the following terms, which all crash at run-time (= their evaluation gets stuck without reaching a proper “value”):

$$\begin{array}{c} \text{pred } \underline{id} \\ \text{ifz } \underline{\text{const } 1} \ \underline{2} \\ (\lambda x. x \ \underline{4}) \ \underline{3} \end{array}$$

To be a bit more precise, we first need to say what we mean by “a (proper) value”.

Definition 13.1. A value, usually V, W , is any term of the following grammar:

$$V, W ::= \underline{n} \mid x \mid \lambda x. V \mid \text{pred} \mid S \mid \text{fix} \mid \text{ifz} \mid \text{ifz } \underline{n} \mid \text{ifz } \underline{n} \ W$$

A normalisable term M whose normal form is not a value, is said to go wrong.

The idea is that, if a term cannot reduce to a value, then either it diverges (reduces forever) or it goes wrong (i.e. would crash at runtime or throw an exception if implemented in a real system). So all of the terms $\text{pred } \underline{id}$, $\text{ifz } \underline{\text{const } 1} \ \underline{2}$, $(\lambda x. x \ \underline{4}) \ \underline{3}$ go wrong - in a sense the first two have already gone wrong because they are in normal form.

Thankfully, none of these terms are typable, so by restricting our attention to well-typed terms, we will have avoided accidentally writing one of them. For example:

Lemma 13.1. $\text{pred } \underline{id}$ is not typable.

Proof. We suppose $\text{pred } \underline{id}$ were typable and look to obtain a contradiction. According to the definition of typability, this means that there is some type A and a derivation of $\vdash \text{pred } \underline{id} : A$. By inversion (recall Theorem 12.1 of the previous chapter), it follows that there is some type B and derivations of:

(i) $\vdash \text{pred} : B \rightarrow A$.

(ii) $\vdash \underline{id} : B$.

However, by inversion on (i), it must be that $\text{pred} : B \rightarrow A \in \mathbb{C}$, from which we can deduce that both B and A must be Nat . But then (ii) is a derivation of $\vdash \text{id} : \text{Nat}$. Since $\text{id} = \lambda x. x$, by inversion on (ii), there are types B_1 and B_2 such that $\text{Nat} = B_1 \rightarrow B_2$. However, this is impossible, since these are obviously different strings. Hence, we have obtained our contradiction. \square

In fact, by restricting our attention to well-typed terms, we obtain a *guarantee*: no typable program will crash at run-time with a type error. In the case of PCF, we say: *well-typed programs don't go wrong*, and we can prove this as a theorem. By *programs*, we mean closed terms i.e. terms without any missing pieces.

Theorem 13.1 (Well-typed programs don't go wrong). *If $\vdash M : A$ then M does not go wrong.*

The proof consists of two parts:

- Showing that reduction steps preserve typing, this is called *preservation* or *subject reduction*.
- Showing that well-typed terms are either values or can make a step, this is called *progress*.

Preservation says that: if $\Gamma \vdash M : A$ and $M \triangleright^* N$, then also $\Gamma \vdash N : A$. If you think about it, this is a property that we rely on quite heavily in our typed programming. For example, if we know that our Haskell expression has type `String`, such as $\Gamma \vdash x ++ \text{"foo"} : \text{String}$, then this is telling us something about how $x ++ \text{"foo"}$, the subject, will evaluate, namely “if it terminates, this program expression will evaluate to a string”. This is good to know because it means that when I feed it to the function `words :: String → [String]` during evaluation, I can be sure that this combination will not crash. Imagine if this were not true, we could have that $\Gamma \vdash x ++ \text{"foo"} : \text{String}$ and yet $x ++ \text{"foo"}$ evaluated to an integer!

This property is called the *preservation of types under subject reduction* or more simply either *preservation* (particularly in literature on programming languages) or *subject reduction* (particularly in literature on pure λ -calculus).

Theorem 13.2 (Subject Reduction). *If $\Gamma \vdash M : A$ and $M \triangleright N$ then $\Gamma \vdash N : A$*

Of course, it will follow immediately that typing is also preserved under any chain of reduction steps:

$$M_0 \triangleright M_1 \triangleright M_2 \triangleright \cdots \triangleright M_{k-1} \triangleright M_k$$

whenever the first term M_0 in the chain is well-typed, simply by applying the theorem repeatedly to each next term in the sequence. Formally, this corresponds to an induction on the length of the reduction sequence.

Corollary 13.1. *If $\Gamma \vdash M : A$ and $M \triangleright^* N$ then $\Gamma \vdash N : A$*

Proof. We can prove $\Gamma \vdash M : A$ and $M \triangleright^n N$ implies $\Gamma \vdash N : A$, for all $n \in \mathbb{N}$, by induction on n . From which the given statement follows by definition. \square

The proof of preservation will be split into two parts. First we will show that the result holds for the particular case when M is a redex and N is its contraction. Then we will lift this to all terms by showing that the result still holds when you place the redex into any context.

Lemma 13.2. *If $\Gamma \vdash P : A$ and P / Q then $\Gamma \vdash Q : A$*

Proof. Suppose (i) $\Gamma \vdash P : A$ and (ii) P / Q . We analyse the cases on the possible pairs P and Q that satisfy (ii):

- In case P is of shape $\text{pred } Z$ and Q is Z , we have that (i) is $\Gamma \vdash \text{pred } Z : A$. By repeatedly applying inversion we can deduce that A must be Nat . Therefore, our goal is to show that $\Gamma \vdash Z : \text{Nat}$, and this follows immediately from (TCst).
- In case P is of shape $\text{pred } (S \ N)$ and Q is N , we have that (i) is $\Gamma \vdash \text{pred } (S \ N) : A$. By repeatedly applying inversion, we can deduce that A must be Nat and $\Gamma \vdash N : \text{Nat}$. This is exactly our goal.
- In case P is of shape $S \ V$ for some $V \notin \text{Val}_{\text{Nat}}$ and Q is wrong, we have that (i) is $\Gamma \vdash S \ V : A$. By repeatedly applying inversion, it follows that we must have a derivation of $\Gamma \vdash V : \text{Nat}$. However, since V is a non-natural number value, it is of one of the following shapes:

$$\lambda x.M \quad \text{pred} \quad S \quad \text{fix} \quad \text{ifz} \quad \text{ifz } M \quad \text{ifz } M \ N$$

If we painstakingly analyse each of these possibilities, e.g. $\Gamma \vdash \lambda x.M : \text{Nat}$, $\Gamma \vdash \text{pred} : \text{Nat}$, $\Gamma \vdash S : \text{Nat}$ and so on, using inversion we can show that each is not a valid typing. Hence, we have obtained a contradiction, and the desired conclusion follows trivially.

- The cases in which P is a redex headed by pred or ifz are similar to those above.
- The case in which P is a redex of shape $(\lambda x.M) \ N$ is more involved and you will prove it in one of this week's exercises.
- In case P is of shape $\text{fix } M$ and Q is M ($\text{fix } M$), we have that (i) is $\Gamma \vdash \text{fix } M : A$. Let us call the proof tree that witnesses this judgement Δ_1 . It follows from inversion that $\Gamma \vdash \text{fix} : (A \rightarrow A) \rightarrow A$ and, therefore, we have a derivation of $\Gamma \vdash M : A \rightarrow A$. Let us refer to this derivation, i.e. the proof tree rooted at $\Gamma \vdash M : A \rightarrow A$, as Δ_2 . The goal is to show that $\Gamma \vdash M (\text{fix } M) : A$, which we can do by combining our given trees Δ_1 and Δ_2 as follows:

$$\frac{\Delta_2 \quad \Delta_1}{\Gamma \vdash M (\text{fix } M) : A} \text{ (TApp)}$$

In all cases we can show that $\Gamma \vdash Q : A$, so we have completed the proof. □

Now we lift this result using one-hole contexts.

Lemma 13.3. *Suppose P / Q is a redex contraction pair. For all contexts $C[\]$, types A and environments Γ , if $\Gamma \vdash C[P] : A$ and $C[P] \triangleright C[Q]$ then $\Gamma \vdash C[Q] : A$.*

Proof. Suppose P / Q is a redex contraction pair. The proof is by induction on $C[\]$.

- If $C[\]$ is $[\]$, we reason as follows. Let A be a type, Γ an environment, and suppose $\Gamma \vdash P : A$ and $P \triangleright Q$. Then it follows from the previous lemma that $\Gamma \vdash Q : A$, which is exactly our goal.
- If $C[\]$ is of shape $D[\] \ N$, we may assume the induction hypothesis:

(IH) For all types A' and environments Γ' , if $\Gamma' \vdash D[P] : A'$ and $D[P] \triangleright D[Q]$ then $\Gamma' \vdash D[Q] : A'$

In this case $C[P] = D[P] N$ and $C[Q] = D[Q] N$. Let A be a type, Γ an environment, and suppose (i) $\Gamma \vdash D[P] N : A$ and (ii) $D[P] N \triangleright D[Q] N$. It follows from (i), by inversion, that there is a type B such that $(*) \Gamma \vdash D[P] : B \rightarrow A$ and $(**) \Gamma \vdash N : B$. It follows from (ii) it must be that $D[P] \triangleright D[Q]$. Hence, we can use the induction hypothesis with $A' := B \rightarrow A$ and obtain $(***) \Gamma \vdash D[Q] : B \rightarrow A$. Hence, it follows from (TApp) with premises $(***)$ and $(**)$ that $\Gamma \vdash D[Q] N : A$, which was exactly our goal.

- If $C[]$ is of shape $M D[]$, the argument is symmetrical to the previous case, so we will skip the detail.
- If $C[]$ is of shape $\lambda x. D[]$, we may assume the induction hypothesis:

(IH) For all types A' and environments Γ' , if $\Gamma' \vdash D[P] : A'$ and $D[P] \triangleright D[Q]$ then $\Gamma' \vdash D[Q] : A'$

In this case $C[P] = \lambda x. D[P]$ and $C[Q] = \lambda x. D[Q]$. We may assume by the variable convention that x does not occur outside of $D[P]$ and $D[Q]$. Let A be a type and Γ an environment, and suppose (i) $\Gamma \vdash \lambda x. D[P] : A$ and (ii) $\lambda x. D[P] \triangleright \lambda x. D[Q]$. It follows from (i) by inversion that there are types B_1 and B_2 such that $A = B_1 \rightarrow B_2$ and there is a derivation of $\Gamma, x : B_1 \vdash D[P] : B_2$. It follows from (ii) that $D[P] \triangleright D[Q]$. Therefore, we may use the induction hypothesis with $A' := B_2$ and $\Gamma' = \Gamma \cup \{x : B_1\}$ to obtain $\Gamma, x : B_1 \vdash D[Q] : B_2$. Hence, by (TAbs), $\Gamma \vdash \lambda x. D[Q] : B_1 \rightarrow B_2$ follows, and this is our goal.

□

From this, we immediately obtain Preservation:

Theorem 13.3 (Preservation). *If $\Gamma \vdash M : A$ and $M \triangleright N$ then $\Gamma \vdash N : A$*

Proof. Suppose (i) $\Gamma \vdash M : A$ and (ii) $M \triangleright N$. By definition of one-step reduction, there is a context $C[]$, and a redex/contraction pair P / Q such that $M = C[P]$ and $N = C[Q]$. Then, we see that (i) is $\Gamma \vdash C[P] : A$ and (ii) is $C[P] \triangleright C[Q]$ and so we can apply the previous lemma to obtain $\Gamma \vdash C[Q] : A$, and this is just $\Gamma \vdash N : A$. □

Progress, on the other hand, says that closed well-typed terms are either values (which is good, it means their execution terminated without going wrong) or they can still make a step (i.e. they are still executing).

Theorem 13.4 (Progress). *If $\vdash M : A$ then either M is a value or there is some N such that $M \triangleright N$.*

We don't have time to prove it, but if you take *Advanced Topics in Programming Languages* in year 4, then you will get plenty of practice at proving progress and preservation theorems.

And, from Preservation and Progress, we obtain Well Typed Programs Don't Go Wrong:

Theorem 13.5 (Well-Typed Programs Don't Go Wrong). *If $\vdash M : A$ then M does not go wrong.*

Proof. Suppose $\vdash M : A$. Then suppose for contradiction that M does go wrong, i.e. there is some reduction sequence $M \triangleright^* N$ with N a normal form which is not a value. Then it follows from Preservation that $(*) \vdash N : A$ too. However, by Progress, we have that therefore N is either a value or it can make a step. However, this contradicts the fact that N is a normal form which is not a value. \square

14. Types: The Curry-Howard Correspondence

In this lecture, we see the incredible connection of pure, typed λ -calculus and logic. We will only be considering the subset of *pure* terms – recall the definition given in Problem Sheet 5:

$$M, N ::= x \mid (MN) \mid (\lambda x. M)$$

There are no constants at all, no PCF numerals, no Z, S, pred, ifz, or fix. Although we can feel free to use the abbreviations you defined for them in Sheet 5.

14.1 Inhabitation

I want to start by considering a problem that is, in some sense, dual to the problem of finding the type of a term. The “elements” of a type (thinking of a type as something like a set of terms) are more usually called the *inhabitants*.

Definition 14.1. *The **pure-term inhabitation problem**, is the problem of, given a type A , determining if there a closed pure term M such that $\vdash M : A$. In such a case, M is said to be an **inhabitant** of A .*

Let’s look at some examples. Which of the following types are inhabited and which are not?

(T1) $a \rightarrow a$

(T2) $(a \rightarrow a \rightarrow b) \rightarrow a \rightarrow b$

(T3) $((a \rightarrow a) \rightarrow b) \rightarrow b$

(T4) $a \rightarrow b$ (with $a \neq b$)

(T5) $(a \rightarrow b) \rightarrow b$

The first three are inhabited, example inhabitants include:

$$\vdash \lambda x. x : a \rightarrow a$$

$$\vdash \lambda xy. xyy : (a \rightarrow a \rightarrow b) \rightarrow a \rightarrow b$$

$$\vdash \lambda x. x(\lambda y. y) : ((a \rightarrow a) \rightarrow b) \rightarrow b$$

The final two types, however, are not inhabited.

Now I would like you to consider the following formulas of propositional logic. Which of these formulas are valid, i.e. true for all possible instantiations of a and b ?

(F1) $a \Rightarrow a$

(F2) $(a \Rightarrow a \Rightarrow b) \Rightarrow a \Rightarrow b$

(F3) $((a \Rightarrow a) \Rightarrow b) \Rightarrow b$

(F4) $a \Rightarrow b$ (with $a \neq b$)

(F5) $(a \Rightarrow b) \Rightarrow b$

Obviously a implies a is true, no matter what a is. Formulas (F2) and (F3) are also valid and can be argued as follows.

Lemma 14.1. *Formulas $(a \Rightarrow a \Rightarrow b) \Rightarrow a \Rightarrow b$ and $((a \Rightarrow a) \Rightarrow b) \Rightarrow b$ are valid.*

Proof. For the first, assume $a \Rightarrow a \Rightarrow b$, call it (x) and then assume a (y). Applying (x) to (y) we have $a \Rightarrow b$ and applying this to (y) we obtain the truth of b .

For the second, assume $((a \Rightarrow a) \Rightarrow b)$, call it (x). We have $a \Rightarrow a$ already because if we assume a , then clearly a follows. Therefore, using (x) we can obtain b , as required. \square

On the other hand, the final two formulas are not valid. For example, $a \Rightarrow b$ is not true for $a := \top$ and $b := \perp$. Likewise, $(a \Rightarrow b) \Rightarrow b$ is not true for $a := \perp$ and $b := \top$.

If you look at the syntax, the only difference between (T1)–(T5) and (F1)–(F5) is the way I wrote the arrow, single-bodied in the former and double-bodied in the latter. Moreover, the types that were inhabited turned out to correspond to the formulas that were valid. Is this a coincidence?

Let's consider *why* we know that (T1)–(T3) are inhabited and *why* we know that (F1)–(F3) are valid: what is the evidence we are using to justify our conclusions in each case. For each type A , the evidence is that there is a closed term M and a type derivation for $\vdash M : A$. For each formula A , the evidence is that there is a proof of A . Let's put these side-by-side in the case of (T2) and (F2), using Γ as a short hand for $\{x : a \rightarrow a \rightarrow b, y : a\}$:

$\frac{\frac{\frac{\Gamma \vdash x : a \rightarrow a \rightarrow b}{\Gamma \vdash xy : a \rightarrow b} \quad \Gamma \vdash y : a}{\Gamma \vdash xyy : b} \quad \frac{x : a \rightarrow a \rightarrow b \vdash \lambda y. xyy : a \rightarrow b}{\vdash \lambda xy. xyy : (a \rightarrow a \rightarrow b) \rightarrow a \rightarrow b}$	<p><i>Proof.</i> Assume $a \Rightarrow a \Rightarrow b$, call it (x), and then assume a (y). Applying (x) to (y) we have $a \Rightarrow b$ and applying this to (y) we obtain b. \square</p>
---	---

If try to reconcile the two sides of this example, you will soon realise that the backwards rule for implication in the proof corresponds quite exactly to the (TAbs) rule. The (TAbs) rule says, reading bottom to top, that to show that a term $\lambda x. M$ has type $A \rightarrow B$ under hypotheses Γ , you should assume

that x has type A and go on to try to type M with B . The backwards rule for implication also says that, to prove $A \Rightarrow B$, you should assume A and try to prove B .

Similarly, the forwards rule for implication corresponds to the (TApp) rule. The (TApp) rule says, reading *top to bottom*, that if I already know that M has type $A \rightarrow B$ under some hypotheses Γ and I already know that N has type A under some hypotheses Γ , then I can conclude that MN has type B under the same hypotheses. The forwards rule for implication says as much: if I know $B \Rightarrow A$ and I know A , then I know B .

The (TVar) just corresponds to where we are recalling an assumption that we already made by using the name we gave it at the time (in this case, x or y).

Take a look at (T3) vs (F3), this time letting $\Gamma = \{x : (a \rightarrow a) \rightarrow b\}$:

$\frac{\frac{\Gamma \vdash x : (a \rightarrow a) \rightarrow b}{\Gamma \vdash x(\lambda y. y) : b} \quad \frac{\Gamma, y : a \vdash y : a}{\Gamma \vdash \lambda y. y : a \rightarrow a}}{\vdash \lambda x. x(\lambda y. y) : ((a \rightarrow a) \rightarrow b) \rightarrow b}$	<p><i>Proof.</i> Assume $((a \Rightarrow a) \Rightarrow b)$, call it (x). We have $a \Rightarrow a$ already because if we assume a, then clearly a follows. Therefore, using (x) we can obtain b, as required. \square</p>
---	--

Try to match up the backwards reasoning with uses of the (TAbs) rule and forwards reasoning with uses of the (TApp) rule.

14.2 C-H for the implicational fragment

I deliberately named all the assumptions in the forgoing proofs and I deliberately chose names that came from our set of term variables e.g. $x, y \in \mathbb{V}$. Let us say that a proof is *properly labelled* if all starting assumptions and those introduced by backwards rules are given a name from \mathbb{V} and names are given sensibly – i.e. we don't give the same name to two different formulas in the same proof. Obviously, it is not a real restriction to require proofs to be properly labelled, given a proof I can always do the labelling afterwards, but it is necessary to state a strong correspondence:

Theorem 14.1 (Curry-Howard). *If we consider properly labelled proofs in the implicational fragment of propositional logic (i.e. the only connective is implication), then the following is true:*

- From a constructive proof of A from starting assumptions Γ , one can extract a pure term M such that $\Gamma \vdash M : A$.
- From each pure term M such that $\Gamma \vdash M : A$, we can extract a constructive proof that A follows from assumptions Γ .

(This requires we accept that \rightarrow is an equally good way to write implication.)

We will come to the *constructive* requirement in a moment. First, let's consider another example. Suppose you give me a proof of the transitivity of implication:

$$(a \Rightarrow b) \Rightarrow (b \Rightarrow c) \Rightarrow a \Rightarrow c$$

Proof. Suppose $a \Rightarrow b$ (x) and $b \Rightarrow c$ (y). Then assume a (z). It follows from (x) and (z) that b holds, and it follows from that and (y) that c holds, as required. \square

I can translate each use of the backwards rule for implication into a use of (TAb), each use of the forwards rule for implication into a use of (TApp) and each recall of a labelled assumption into (TVar). This gives me a skeleton of a type derivation, in which I know which rules are used to conclude which types, but we don't yet have any of the terms in the interior:

$$\begin{array}{c}
\frac{}{\Gamma \vdash y : b \rightarrow c} \text{ (TVar)} \quad \frac{}{\Gamma \vdash x : a \rightarrow b} \text{ (TVar)} \quad \frac{}{\Gamma \vdash z : a} \text{ (TVar)} \\
\frac{}{\Gamma \vdash ? : b} \text{ (TApp)} \\
\frac{}{\Gamma \vdash ? : c} \text{ (TApp)} \\
\frac{}{x : a \rightarrow b, y : b \rightarrow c \vdash ? : a \rightarrow c} \text{ (TAbs)} \\
\frac{}{x : a \rightarrow b \vdash ? : (b \rightarrow c) \rightarrow a \rightarrow c} \text{ (TAbs)} \\
\frac{}{\vdash ? : (a \rightarrow b) \rightarrow (b \rightarrow c) \rightarrow a \rightarrow c} \text{ (TAbs)}
\end{array}$$

But, if I know all the types and which rules were used, then the terms are completely determined, I can just fill them in starting from the leaves and working my way down:

$$\begin{array}{c}
\frac{}{\Gamma \vdash y : b \rightarrow c} \text{ (TVar)} \quad \frac{}{\Gamma \vdash x : a \rightarrow b} \text{ (TVar)} \quad \frac{}{\Gamma \vdash z : a} \text{ (TVar)} \\
\frac{}{\Gamma \vdash xz : b} \text{ (TApp)} \\
\frac{}{\Gamma \vdash y(xz) : c} \text{ (TApp)} \\
\frac{}{x : a \rightarrow b, y : b \rightarrow c \vdash \lambda z. y(xz) : a \rightarrow c} \text{ (TAbs)} \\
\frac{}{x : a \rightarrow b \vdash \lambda yz. y(xz)(b \rightarrow c) \rightarrow a \rightarrow c} \text{ (TAbs)} \\
\frac{}{\vdash \lambda xyz. y(xz) : (a \rightarrow b) \rightarrow (b \rightarrow c) \rightarrow a \rightarrow c} \text{ (TAbs)}
\end{array}$$

The term that I obtain in this way is sometimes called the *proof term* of the formula. In this case, the proof term for the transitivity of implication is none other than the program that implements function composition:

$$N \circ M := (\lambda x y z. y(xz)) N M$$

There is a sense in which, if it is in normal form, a term of type A is a proof of A . For example, let's revisit (T2). The term here is $\lambda x. x(\lambda y. y)$ and the type is $(a \rightarrow a) \rightarrow b \rightarrow b$. Since this term is in normal form, the judgement $\vdash \lambda x. x(\lambda y. y) : (a \rightarrow a) \rightarrow b \rightarrow b$ completely determines the derivation tree:

$$\frac{\frac{\Gamma \vdash x : (a \rightarrow a) \rightarrow b}{\Gamma \vdash x(\lambda y. y) : b} \quad \frac{\Gamma, y : a \vdash y : a}{\Gamma \vdash \lambda y. y : a \rightarrow a}}{\vdash \lambda x. x(\lambda y. y) : ((a \rightarrow a) \rightarrow b) \rightarrow b}$$

There are no choices that can be made in how we construct this tree. But then the derivation tree determines a proof, by mapping each use of (TVar) to the recall of an assumption, each use of (TApp) to the forwards rule for implication and each use of (TAbs) to the backwards rule for implication. We obtain:

Proof. Assume $(a \Rightarrow a) \Rightarrow b$, call it (x). We have $a \Rightarrow a$ by the following reasoning: assume a and call it (y), then a follows immediately from (y). From $a \Rightarrow a$ and (x) we obtain b . \square

Notice how the subterm $\lambda y. y$ corresponds to a subproof of $a \Rightarrow a$. In a sense, $\lambda y. y$ is a proof of $a \Rightarrow a$.

Of course, the implicational fragment of propositional logic is not very interesting: we can hardly carry out many proofs in that system. Remarkably, however, the Curry-Howard correspondence carries over from this tiny λ -calculus and implicational logic, all the way up to first- and even higher-order logic, i.e. all the propositional connectives and all the quantifiers. In other words, all (constructive) mathematics corresponds to programming up λ -terms in such a way that they have the right type (of course, just as the first-order logic is larger than and more complicated than its implicational fragment, so the corresponding λ -calculi are larger and more complicated too). This is the so-called “proofs as programs” paradigm, which has given rise to new kinds of programming languages and new kinds of logic.

Types	\leftrightarrow	Formulas
Programs	\leftrightarrow	Proofs
Inhabitation	\leftrightarrow	Provability

15. Types: Constructive Proofs

The material in this chapter will not be assessed.

The idea of proof is that, if we have a valid proof of a proposition A , then we should be able to conclude that A is *true*. Do you believe in this? Do you believe in your proofs? Why should you believe that if you can write one of these little paragraphs of text with a scattering of mathematical formulas that therefore A is actually true? In other words, why should you believe that the rules of proof are *sound*?

Classically, we can justify the soundness of the proof rules by appealing to an independent definition of when formulas are *true*. We start from the assumption that every atomic proposition is either true or false, and then the truth of compound formulas is determined by the truth of their parts. For example, truth tables define the truth of the logical connectives in terms of their arguments. For universal quantification, $\forall x \in \mathbb{N}. A$ is true when $A[n/x]$ is true for all choices of natural number n . Such a definition of truth ensures that every closed formula, i.e. *sentence*, is either true or false under a given interpretation of the non-logical parts (like numbers). However, given any particular proposition C , we may not be able to tell from the definition which it is, because it may require verifying that *infinitely many* other formulas (e.g. the $A[n/x]$ when $C = \forall x \in \mathbb{N}. A$) are themselves true, and this is best attempted only if you are a deity.

According to this view, proof rules are a system that mere mortals can use to try to get a handle on the divine notion of *truth*. Human beings can write down proofs, which are finite objects and, as a result, conclude the truth of formulas — even universal quantifications which seem to require an infinite amount of checking.

A proof of A from starting assumptions Γ is a certificate guaranteeing that A is true whenever Γ is true.

This view of proof as being a syntactical certificate of truth allows us to justify why this particular collection of rules for deriving certain combinations of strings from other combinations of strings makes sense as a logic. For example, under this view, the forwards proof rule for implication is sound. This rule says, roughly: “if I have proven $A \Rightarrow B$ and I have proven A , then I have proven B ”. If we view “proven” as a way of saying “certified to be true using finite, mortal means”, then indeed, if I have certified $A \Rightarrow B$ to be true and I have certified A to be true, then by the definition of truth for

$A \Rightarrow B$, i.e.

A	B	$A \Rightarrow B$
\perp	\perp	\top
\perp	\top	\top
\top	\perp	\perp
\top	\top	\top

it can only be that B is true (i.e. we are in the last row). Similarly, the backwards proof rule for implication is sound. This rule says, roughly “to prove $A \Rightarrow B$, it suffices to assume that A is true and try to prove B ”. Indeed, if I want to certify that $A \Rightarrow B$ is true, then by the definition of truth, then it suffices to assume that A is true and try to certify that B is true as a consequence. This is because, if A is false, then $A \Rightarrow B$ is true by definition, so there is nothing to check.

15.1 Truth from Proof

Not all logicians are in agreement about this idea of truth at the foundation of logic and that proof is something we invent afterwards as a mortal means to discover it. For example, a problem with the foregoing account is that a proposition can be true even when we don’t have a certificate of it. We can even know that some proposition A is true without *ever* proving it *directly*, for example, if we merely have a certificate that $\neg A$ is false.

By and large, *intuitionists* believe that mathematics is sum total of the mental constructions that mathematicians devise and communicate to each other — there is no ideal or platonic truth behind the scenes that we just trying to get a glimpse of. In particular, intuitionists regard a proposition as true only if there is some concrete *evidence*. This gives an alternative interpretation of the proof, which does not presuppose a definition of truth.

A proof of A from starting assumptions Γ is a method for *constructing* evidence of A from evidence of Γ .

According to this view, proof is not about making arguments to certify truth, but instead about manipulating evidence. The idea is credited to Brouwer and Heyting and, independently, Kolmogorov. We start by explaining what kinds of evidence is needed in order to believe in any formula.

Each atomic proposition has its own notion of evidence, for example, the reading on a thermometer might be evidence for a basic proposition “the temperature is 23°C”. A proof tree using the rules (Var), (Abs) and (App) is evidence for the basic proposition $\lambda x y. yx \in \Lambda$.

We can understand more complex formulas in terms of the evidence required for their components.

- There can be no evidence for the truth of false.
- No evidence is needed for the (self-evident) truth of true.
- Evidence for $A \wedge B$ is a pair consisting of evidence for A and evidence for B .
- Evidence for $A \vee B$ is either a piece of evidence of A or a piece of evidence of B .
- Evidence for $A \Rightarrow B$ is a *procedure* for transforming evidence of A into evidence of B .

- Evidence of $\neg A$ is a procedure for transforming any evidence of A into evidence of false.
- Evidence of $\forall x \in X. A$ is a procedure for transforming any y and evidence of $y \in X$ into evidence of $A[y/x]$.
- Evidence of $\exists x \in X. A$ is a triple consisting of a y , evidence of $y \in X$ and evidence of $A[y/x]$.

A crucial criterion for intuitionists is that evidence of an existential statement actually constructs a witness. Indeed, such logics are sometimes called *constructive* because, in general, the evidence can be seen as some construction that witnesses the truth of the formula.

This gives us an alternative way to justify our proof rules. We can justify the forwards rule for implication as follows. If we have a proof of $A \Rightarrow B$ we interpret this to mean that we have, in our hand, a procedure for transforming evidence of A to evidence of B . So, if we also have evidence of A , we can apply our procedure to it in order to obtain some evidence of B .

These can be a bit of a mouthful to write this out in English, but luckily it can be precisely summarised by the rule:

$$\frac{\Gamma \vdash M : B \rightarrow A \quad \Gamma \vdash N : B}{\Gamma \vdash MN : A} \text{ (TApp)}$$

Read, from top to bottom: if I have evidence M of $B \rightarrow A$ and I have evidence N of B , then I can construct evidence of A by applying M to N , i.e. the evidence of A is MN .

Similarly, we can justify the backwards rule for implication as follows. If we have a proof of B starting from assumptions Γ including A , then this means that we have a method that can construct evidence for B from evidence for each of the propositions in Γ and evidence for A . Therefore, given evidence of each of the propositions in Γ , we have a way for converting any piece of evidence of the remaining proposition A into evidence for the proposition A , i.e. we simply supply that new piece of evidence along with our evidence for Γ to the method to obtain evidence of B . Hence we can package this up as a procedure evidencing $A \Rightarrow B$.

This is neatly summarised by the following rule:

$$x \notin \text{dom } \Gamma \quad \frac{\Gamma \cup \{x : A\} \vdash M : B}{\Gamma \vdash \lambda x. M : A \rightarrow B} \text{ (TAbs)}$$

This rule reads: if I can construct evidence of B from evidence of Γ and some evidence x of A , then from evidence of Γ only, I can construct a function that can already take any such x and return evidence of B . In other words, $\lambda x. M$ is the evidence by which we can believe that $A \rightarrow B$ is true.

We can justify e.g. the backwards rule for forall too. If we have a proof of A starting from assumptions Γ and that proof makes no assumption on say $x \in \mathbb{N}$, then, effectively, I have a procedure for converting arbitrary natural numbers n into evidence of $A[n/x]$: if you give me such a number, say $6 \in \mathbb{N}$, then I can substitute it everywhere into my proof of A to obtain a proof of $A[6/x]$, i.e. I have produced a new proof that is specific to the number 6.

Our type system is too simple to capture universal quantification on individuals, but more advanced type systems, such as the type system for the *Agda* programming language does so, and the terms of type $\forall x. A$ are functions.

The justification of the other *constructive* proof rules is similar. However, there are a few other rules of proof, such as *double negation elimination* and *the law of the excluded middle* that we did not cover in our crash-course on proofs. These rules are said to be *non-constructive*.

In the classical truth-theoretic reading of logic, every proposition is either true or false, although it is quite often the case that for any particular proposition, we mortals don't know which one it happens to be. A consequence of this is that, if I know that something is not false, it must be true. So, in this truth-theoretic reading, *reductio ad absurdum* (a.k.a the double negation rule) is justified: if I know that $\neg A$ is false, A must be true — in other words, to prove A it suffices to prove $\neg\neg A$.

This is quite different from the evidence-based approach to logic. If there can be no evidence for $\neg A$, does that mean that I have evidence for A ? For example, to have evidence of a disjunction $(\lambda x. x \in \Lambda) \vee (\vdash \lambda x. x : b)$ is to actually hold evidence for one of the disjuncts, i.e. a grammar derivation for $\lambda x. x \in \Lambda$ or a typing derivation for $\vdash \lambda x. x : b$. If I have evidence for $\neg(\neg((\lambda x. x \in \Lambda) \vee (\vdash \lambda x. x : b)))$ then I am holding a procedure that *would* transform any evidence of $\neg((\lambda x. x \in \Lambda) \vee (\vdash \lambda x. x : b))$ into evidence of false. That is to say, given a procedure that transforms evidence of one of the disjuncts into evidence of false as an input, I have in my hand a procedure that will produce evidence of false as output. Having such a procedure is clearly not the same as actually holding evidence of one of the disjuncts, which is a derivation. In particular, if I had evidence of one of the disjuncts then I would know that particular disjunct was true, but if I merely have a procedure for converting evidence of their negation into evidence of false, then I can't tell you anything about which of the disjuncts is actually true, which is the whole basis for believing that a disjunction is true in the first place.

Therefore, the double negation rule is rejected by intuitionists. They don't consider proofs that use this rule to be convincing. Since our Curry-Howard correspondence only applies to constructive proofs, types that correspond to formulas that are classically valid (i.e. true in the classical sense) but not intuitionistically valid (i.e. for which there is no constructive proof) are uninhabited. For example, Peirce's law $((a \Rightarrow b) \Rightarrow a) \Rightarrow a$ is an equivalent way to state double negation that only uses implication. Peirce's law is classically valid, i.e. true for all assignments of truth values to a and b , but the only proofs are non-constructive, for example, the following proof uses the double negation rule (proof by contradiction):

Proof. Assume $(a \Rightarrow b) \Rightarrow a$, call it (x). Our goal is to prove a , but we will proceed by contradiction, i.e. we will instead prove $\neg\neg a$. So assume $\neg a$ (y). We obtain $a \Rightarrow b$ by a subproof as follows: assume a (z), then by combining (y) with (z) we obtain \perp and from \perp anything follows, including b . Then, returning to the main proof, from $a \Rightarrow b$ and (x) we obtain a , but this contradicts (y). \square

Consequently, the correspondence does not apply and, indeed, we can show that the type $((a \rightarrow b) \rightarrow a) \rightarrow a$ is uninhabited in our type system.

A. Proof Rules

Implication

$A \Rightarrow B$	A implies B	If A then B .	Suppose A , then B .
-------------------	-----------------	-------------------	--------------------------

Backward Rule

To prove $A \Rightarrow B$ starting from some assumptions, it suffices to instead prove B with A added to the assumptions.

Assms	Goal
...	$A \Rightarrow B$

"Assume A . We prove B as follows..."

"To show $A \Rightarrow B$, we assume A"

"Assume A"

"Suppose A is true. ..."

Assms	Goal
..., A	B

Forwards Rule

If you already know $A \Rightarrow B$ and you already know A , then you know B .

Assms	Goal
$\dots, A \Rightarrow B, A$	C

"From $A \Rightarrow B$ and A we conclude B ."

"Applying $A \Rightarrow B$ to A we obtain B ."

"We have B by *modus ponens*."

Assms	Goal
$\dots, A \Rightarrow B, A, B$	C

Conjunction

$A \wedge B$ A and B

Backwards Rule

To prove $A \wedge B$ from some assumptions, it suffices to give two separate proofs: one proof of A from the assumptions and one proof of B from the assumptions.

Assms	Goal
\dots	$A \wedge B$

"We prove A and B separately. To prove A ..."

"For A , we argue as follows ... For B , ..."

Assms	Goal	Assms	Goal
\dots	A	\dots	B

The proof splits into two at this point, so you have to keep in mind two states, although you can just work with one at a time.

Forwards Rule

If you know $A \wedge B$, then you can know A and you know B .

Assms	Goal
$\dots, A \wedge B$	C

"" (too trivial to mention)

Assms	Goal
$\dots, A \wedge B, A, B$	C

Disjunction

$A \vee B$ A or B

Backwards Rules

To prove $A \vee B$ from some assumptions, it suffices to give a proof of A from those assumptions.

Assms	Goal
\dots	$A \vee B$

"To see $A \vee B$, observe that A is true since..."

"We prove A ."

"It suffices to prove A ."

Assms	Goal
\dots	A

To prove $A \vee B$ from some assumptions, it suffices to give a proof of B from those assumptions.

Assms	Goal
\dots	$A \vee B$

"To see $A \vee B$, observe that B is true since..."

"We prove B ."

"It suffices to prove B ."

Assms	Goal
\dots	B

Forwards Rule

If you have already know $A \vee B$ and, additionally, you can give the following two proofs:

- a proof of C starting from your assumptions and A
- a proof of C starting from your assumptions and B

Then you will also know C .

Assms	Goal
$\dots, A \vee B$	C

"We proceed by cases on $A \vee B$.

Assume A ... Hence C .

Assume B ... Hence C ."

"We analyse the two cases."

"We proceed by case analysis on $A \vee B$."

Assms	Goal	Assms	Goal
$\dots, A \vee B, A$	C	$\dots, A \vee B, B$	C

This forces the proof to split into two, and you're not done until you have completed both of them.

Negation

$\neg A$	not A	A is false	A is absurd
----------	---------	--------------	---------------

Backwards Rule

To prove $\neg A$ starting from some assumptions, it suffices to derive a contradiction (give a proof of false) starting from the same assumptions with A added.

Assms	Goal
\dots	$\neg A$

"We assume A and try to obtain a contradiction."

"To show $\neg A$, we assume A ..."

"Assume A ."

Assms	Goal
\dots, A	\perp

Forwards Rule

If you know $\neg A$ and you also know A , then you know absurdity.

Assms	Goal
$\dots, \neg A, A$	C

“From A and $\neg A$ we obtain a contradiction.”

Assms	Goal
$\dots, \neg A, A, \perp$	C

If and only if

$\Leftrightarrow \equiv \text{iff}$

Backwards Rule

An iff can be thought of as two implications, one in each direction. So to prove an iff it suffices to prove each separately.

Assms	Goal
\dots	$A \Leftrightarrow B$

“In the forwards direction...

... in the backwards direction...”

“We prove each direction separately...”

Assms	Goal	Assms	Goal
\dots	$A \Rightarrow B$	\dots	$B \Rightarrow A$

Forwards Rule

Conversely, to use one, first just extract the two directions.

Assms	Goal
$\dots, A \Leftrightarrow B$	C

"" (too trivial to mention)

Assms	Goal
$\dots, A \Leftrightarrow B, A \Rightarrow B, B \Rightarrow A$	C

False

\perp false absurdity

Forwards Rule

If you know absurdity, then you know everything.

Assms	Goal
\dots, \perp	C

"Hence we obtain the desired result."

"Ex falso quodlibet"

"A follows trivially"

Assms	Goal
\dots, \perp, A	C

Note, this A can be anything you like, including C .

Universal quantification

$\forall x \in X. A$	$\forall x : X. A$	for all x in X , A	A holds of all x in X
----------------------	--------------------	--------------------------	-----------------------------

Backwards Rule

To give a proof of $\forall x \in X. A$ from some assumptions, it suffices to give a proof of A from the same assumptions with $x \in X$ added, as long as you have not made any prior assumptions about the name x (x doesn't appear in any of your existing assumptions).

Assms	Goal	with x not occurring free in the assms
...	$\forall x \in X. A$	

"Let $x \in X$. We show A ."

"Let x be an arbitrary member of X . We show A ."

"Suppose x is in X therefore A ."

Assms	Goal
..., $x \in X$	A

Forwards Rule

If you have know $\forall x \in X. A$ and you know $t \in X$, then you know A with all occurrences of x replaced by t .

Assms	Goal
... $\forall x \in X. A$	$t \in X$

"It follows that A holds of t "

Assms	Goal
..., $\forall x \in X. A$	$t \in X, A[t/x]$

Here, $A[t/x]$ means A with all free occurrences of x replaced by t .

Existential Quantification

$\exists x \in X. A$	$\exists x : X. A$	there exists x in X such that A	A holds of some x in X
----------------------	--------------------	---------------------------------------	------------------------------

Backwards Rule

To prove $\exists x : X. A$ from some assumptions, it suffices to find a witness t and give proofs that $t \in X$ and of A with every occurrence of x replaced by t , starting from the same assumptions.

Assms	Goal
...	$\exists x \in X. A$

“We show that A holds of t .”

“We take t as witness. To see $A[t/x]$...”

“We show that t is such an x .”

Assms	Goal	Assms	Goal
...	$t \in X$...	$A[t/x]$

Here, $A[t/x]$ means A with all free occurrences of x replaced by t .

Forwards Rule

Assms	Goal
... $\exists x \in X. A$	C

with x not occurring free in the assms

“Let x be the witness...”

Assms	Goal
... $\exists x \in X. A$	C

Cut

At any point during a proof of some proposition C , you may launch into proving a subsidiary proposition A . Once proven, this proposition may be assumed when continuing the proof of C .

Assms	Goal
...	C

"We claim A . To prove A , proceed as follows... Now returning to the proof of C ..."

Assms	Goal	Assms	Goal
...	A	... A	C

Equality

At any time, you can assume $s = s$, for any term s , which is obviously true.

Assms	Goal
...	C

"" (too trivial to mention)

Assms	Goal
..., $s = s$	C

If you know an equation, you can use it to rewrite the goal.

Assms	Goal
..., $s = t$	$C[s/x]$

" $C[s/x]$ is just $C[t/x]$ "

" $C[s/x]$ is therefore the same as $C[t/x]$ "

Assms	Goal
..., $s = t$	$C[t/x]$

The use of $C[s/x]$ here is just a way of saying that you have some proposition which contains an occurrence of s and the corresponding $C[t/x]$ indicates replacing that occurrence of s by t .

If you know an equation, you can use it to rewrite an assumption.

Assms	Goal
$\dots, s = t, A[s/x]$	C

" $C[s/x]$ is just $C[t/x]$ "

" $C[s/x]$ is therefore the same as $C[t/x]$ "

Assms	Goal
$\dots, s = t, A[s/x], A[t/x]$	C